

**Development of Blood Metabolomics for Biomarker Discovery Using Chemical Isotope Labeling and Liquid Chromatography-Mass Spectrometry**

by

Wei Han

A thesis submitted in partial fulfillment of the requirements for the degree of

Doctor of Philosophy

Department of Chemistry  
University of Alberta

© Wei Han, 2017

## **Abstract**

Blood-metabolite-based biomarkers are useful for the early-stage diagnosis, accurate prognostic prediction and personalized treatment of various diseases. Blood contains a massive number of metabolites that can potentially become biomarkers, but the metabolome coverage of current analytical techniques remains insufficient. Considering that traditional liquid chromatography-mass spectrometry (LC-MS) platforms are limited by the low metabolome coverage and quantification accuracy, our lab has developed the high-performance chemical isotope labeling (CIL) LC-MS platform which can significantly increase the metabolome coverage and efficiently overcome the detection variability. In general, blood biomarker discovery is susceptible to experimental interferences and biological variations. To improve the reliability of metabolomics discoveries, large sample sizes and time-resolved analysis are highly desirable.

Towards these challenges, the major part of my thesis focuses on assessing and minimizing the experimental and biological variations that could interfere with biomarker discoveries. With the CIL LC-MS platform, experimental variations have been largely overcome and biological variability carefully evaluated. The establishment of a serum metabolome database can facilitate future studies, and the high-coverage analysis of one microliter finger blood opens new vistas for biomarker discovery and environmental exposure assessment. The major motivation of this thesis is to consider the time factor in blood metabolomics, and the reported technical improvement has made time-dependent metabolomics possible at a minimal cost. Additionally, to demonstrate the benefits of

adding the time dimension to the study design, we report a cohort study for the diagnosis and prognosis of Parkinson's disease and an intervention study for the exposure assessment of DDT, which is a banned pesticide and an endocrine disruptor. Overall, this thesis work has demonstrated enhanced analytical performance for blood biomarker discovery, and high-quality biomarkers come from well-designed experiments, careful sample handling, high-performance analysis platforms, and a solid understanding of data analysis principles.

## Preface

In this research work, the collection of urine, venipuncture blood and finger sticking blood received research ethics approval from the University of Alberta Research Ethics Board.

A version of Chapter 2 was published as: Wei Han and Liang Li. "Matrix effect on chemical isotope labeling and its implication in metabolomic sample preparation for quantitative metabolomics." *Metabolomics* 11.6 (2015): 1733-1742. I was responsible for the design of experiment, data collection and analysis, as well as manuscript preparation. Dr. Liang Li supervised the project and edited the manuscript.

A version of Chapter 6 was published as: Wei Han, Shraddha Sapkota, Richard Camicioli, Roger A. Dixon, and Liang Li. "Profiling novel metabolic biomarkers for Parkinson's disease using in-depth metabolomic analysis." *Movement Disorders* (2017). I was responsible for sample preparation, analytical analysis, data interpretation and manuscript preparation. Dr. Shraddha Sapkota contacted the participants, collected the serum samples and helped with the revision of the manuscript. Dr. Richard Camicioli, Dr. Roger Dixon and Dr. Liang Li supervised the project and edited the manuscript.

A version of Chapter 7 was published as: Weifeng Shen, Wei Han, Yunong Li, Zhiqi Meng, Leiming Cai and Liang Li. "Development of chemical isotope labeling liquid chromatography mass spectrometry for silkworm hemolymph metabolomics." *Analytica chimica acta* 942 (2016): 1-11. I was responsible for the analytical analysis, data

interpretation, and part of the manuscript preparation. Dr. Zhiqi Meng prepared the silkworm model. Dr. Weifeng Shen raised the silkworms, collected the samples and prepared the draft of the paper. Yunong Li helped with the statistical analysis. Dr. Leiming Cai and Dr. Liang Li supervised the project. Dr. Liang Li polished the manuscript.

## **Acknowledgments**

This thesis would not have been possible without the support from many people, and it is my pleasure to convey my gratitude to them all.

First and foremost, I am glad to express my sincere gratitude and appreciation to my supervisor, Dr. Liang Li, for his invaluable guidance, advice, encouragement, and support throughout my graduate studies. His profound insights in the omics area have been and will continue to be the guide of my research career. I also want to extend my gratitude to Mrs. Li for revising the first chapter of my thesis.

Second, I would like to extend my appreciation to the members of my supervisory committee, Dr. John Vederas and Dr. Richard McCreery, for their insightful comments and suggestions on my annual reports and during my candidacy exam. Also, here I express my sincere gratitude to Dr. Arthur Mar and Dr. Anas El-Aneed for the time attending my oral examination and reviewing my thesis. I also want to thank Dr. John Klassen and Dr. Arthur Mar for attending my candidacy exam and providing me suggestions.

Furthermore, I want to thank my collaborators, without whom the application works in my thesis would not be possible. Dr. Roger Dixon, Dr. Richard Camicioli and Dr. Shraddha Sapkota helped me a lot with polishing the manuscript of the work on Parkinson's disease. And their deep and broad knowledge in the clinical area has greatly benefited my other

research works. I also really enjoyed the time working together with Dr. Weifeng Shen on the silkworm metabolome analysis.

My deep gratefulness also goes to all the previous and current members of the Li group. Importantly, I want to sincerely thank Dr. Tao Huan, Dr. Nan Wang, Dr. Difei Sun, Dr. Yiman Wu, Dr. Zhendong Li, Dr. Chiao-Li Tseng, Dr. Ruokun Zhou and Dr. Jun Peng for their training, advice and support at the beginning of my research work. I express special thanks to the people that I have worked with: Tran Tran, Yunong Li, Kevin Hooton, Dorothea Mung, Xiaohang Wang, Adriana Zardini Buzatto, Minglei Zhu and Xinyun Gu. I have also been benefited from the knowledge of the other group members, including but not limited to Jaspaul Tatlay, Xian Luo, Shuang Zhao, Hao Li and Wan Chan.

In addition, I would like to appreciate the support from all the staffs in the Department of Chemistry. I'm especially thankful to Mr. Gareth Lambkin of the Biological Services Facility for his kind assistance in the use of bio-hazard level 2 fume hood. I also want to acknowledge Dr. Randy Whittal, Ms. Jing Zheng and Mr. Bela Reiz of the Mass Spectrometry Facility for their professional support.

I am very grateful to the Alberta Doctoral Awards for Chinese Students (ADACS) Program and the Department of Chemistry at the University of Alberta for providing me the financial support during my graduate studies.

Finally and most importantly, I would never be where I am today without my family. Words cannot express my heartfelt gratitude to my beloved wife, Yanjun, for having done more than any faith could have done to encourage me to step forward, more than any fate could have done to make me happy, and more than any fairy could have done to show me a colorful life. And our parents, Mr. Yusheng Han, Ms. Yumei Chen, Mr. Zhong Guo, and Ms. Pingfeng Yan, as well as our grandparents, Mr. Yuansuo Yan and Ms. Ruiying Niu, deserve our highest appreciation for being our peerless listeners and incomparable rocks.

## Table of Contents

<b>Chapter 1 - Introduction</b> .....	1
<b>1.1 Introduction to Blood Biomarkers</b> .....	1
<b>1.2 Small-molecule Biomarkers and Metabolomics Analysis</b> .....	5
<b>1.3 Major Metabolomics Platforms for Blood Biomarker Discovery</b> .....	10
<b>1.4 Chemical Isotope Labeling in LC-MS-based Metabolomics</b> .....	13
<b>1.5 Blood Sample Handling Techniques</b> .....	17
<b>1.6 Statistical Analyses for Biomarker Discovery</b> .....	19
1.6.1 Uni-variate Analysis .....	20
1.6.2 Multi-variate Analysis .....	27
<b>1.7 Challenges and Solutions</b> .....	30
1.7.1 Statistical Over-fitting.....	31
1.7.2 Biological Variability .....	34
1.7.3 Study Design .....	36
<b>1.8 Overview of Thesis</b> .....	38
<b>Chapter 2 - Matrix Effect on Chemical Isotope Labeling and Its Implication in Metabolomic Sample Preparation for Quantitative Metabolomics</b> .....	41
<b>2.1 Introduction</b> .....	41
<b>2.2 Materials and Methods</b> .....	43

2.2.1 Chemicals and Reagents .....	43
2.2.2 Urine Sample Collection .....	43
2.2.3 Dansylation Labeling.....	43
2.2.4 LC-UV Quantification .....	44
2.2.5 LC-FTICR-MS.....	45
2.2.6 Data Analysis.....	45
<b>2.3 Results and Discussion .....</b>	<b>46</b>
2.3.1 Matrix Effect on CIL .....	46
2.3.2 Origin of the Matrix Effect.....	52
2.3.3 Minimizing the Matrix Effect on Metabolomic Profiling .....	56
<b>2.4 Conclusions .....</b>	<b>59</b>

**Chapter 3 - Development of a Human Serum Metabolome Database and Analysis of Metabolome Variations Using Isotope Labeling and High-resolution LC-MS.....60**

<b>3.1 Introduction.....</b>	<b>60</b>
<b>3.2 Materials and Methods .....</b>	<b>63</b>
3.2.1 Chemicals and Reagents .....	63
3.2.2 Serum Sample Collection and the Universal Serum Standard.....	63
3.2.3 Dansylation Labeling of Serum Sample .....	65
3.2.4 DMPA-labeling of Serum Sample.....	66
3.2.5 LC-UV Quantification and Pre-acquisition Sample Normalization .....	66
3.2.6 LC-FTICR-MS Analysis.....	67
3.2.7 Data Processing and Statistical Analysis .....	68
3.2.8 Metabolite Identification.....	69

<b>3.3 Results and Discussion</b> .....	70
3.3.1 Quality Control and Sample Normalization .....	70
3.3.2 Development of a Human Serum Metabolome Database .....	74
3.3.3 Variations Associated with Sex .....	79
3.3.4 Age Effects .....	88
3.3.5 BMI Effects .....	89
<b>3.4 Conclusions</b> .....	97

**Chapter 4 - High-coverage Metabolomics Analysis of One Microliter of Blood Using Chemical Isotope Labeling and High-resolution LC-MS..... 114**

<b>4.1 Introduction</b> .....	114
<b>4.2 Materials and Methods</b> .....	117
4.2.1 Chemicals and Reagents .....	117
4.2.2 Universal Serum Standard .....	117
4.2.3 Dansylation Labeling of Serum Sample .....	118
4.2.4 Finger Blood Sample Collection and Processing.....	118
4.2.5 Sample Quantification and Normalization .....	119
4.2.6 LC-FTICR-MS Analysis.....	120
4.2.7 Nano-LC-QTOF-MS.....	120
4.2.8 Data Processing and Metabolite Identification .....	121
<b>4.3 Results and Discussion</b> .....	121
4.3.1 Finger Blood Collection and Sample Preparation .....	121
4.3.2 Relative Quantification and Internal Standard .....	124
4.3.3 Optimization of the Labeling Method .....	126

4.3.4 Metabolome Coverage of the Optimized Method .....	130
4.3.5 Sample Normalization .....	132
4.3.6 Studying the Dietary Effect of Coffee with Finger Blood Analysis.....	134
4.3.7 Time-resolved Metabolic Analysis.....	140
<b>4.4 Conclusions</b> .....	<b>143</b>
<b>Chapter 5 - High-coverage Metabolomics Analysis of One Microliter Blood Using Two Isotope Labelings and High-resolution LC-MS</b> .....	<b>145</b>
<b>5.1 Introduction</b> .....	<b>145</b>
<b>5.2 Materials and Methods</b> .....	<b>146</b>
5.2.1 Chemicals and Reagents .....	146
5.2.2 Finger Blood Sample Collection .....	147
5.2.3 Dansyl-labeling.....	147
5.2.4 DMPA-labeling .....	148
5.2.5 LC-QTOF-MS.....	148
5.2.6 Data Processing and Statistical Analysis .....	149
5.2.7 Metabolite Identification.....	149
<b>5.3 Results and Discussion</b> .....	<b>150</b>
5.3.1 Metabolome Coverage of the Two Differential Isotope Labeling Methods.....	150
5.3.2 Time-resolved Metabolomics Analysis .....	152
5.3.3 Studying the Diet Effects of an Energy Drink with Two Labeling Methods .....	156
<b>5.4 Conclusions</b> .....	<b>160</b>

**Chapter 6 - Profiling Novel Metabolic Biomarkers for Parkinson’s Disease Using In-depth Metabolomic Analysis.....165**

**6.1 Introduction.....165**

**6.2 Methods .....167**

6.2.1 Participants.....167

6.2.2 Serum Samples and Dansylation LC-MS Metabolomic Profiling .....170

6.2.3 Data Processing and Statistical Analysis.....172

**6.3 Results .....172**

6.3.1 Submetabolome and Metabolite Identification .....172

6.3.2 Comparative Metabolome Analysis for PD biomarker Discovery .....173

6.3.3 Comparative Metabolome Analysis of PD with and without Incipient Dementia ....174

6.3.4 Common Discriminating Metabolites.....176

6.3.5 ROC Curves .....176

**6.4 Discussion.....183**

6.4.1 Catecholamine Metabolism.....185

6.4.2 Tryptophan Metabolism .....187

6.4.3 Caffeine Metabolism .....187

6.4.5 Oxidative Stress .....189

**6.5 Limitations.....190**

**6.6 Conclusions.....192**

**Chapter 7 - Development of Chemical Isotope Labeling Liquid Chromatography-Mass Spectrometry for Silkworm Hemolymph Metabolomics.....193**

<b>7.1 Introduction</b> .....	193
<b>7.2 Materials and methods</b> .....	195
7.2.1 Chemicals and Reagents .....	195
7.2.2 Silkworm Rearing and DDT Treatment .....	195
7.2.3 Hemolymph Collection and Preparation .....	196
7.2.4 Dansylation Labeling of Hemolymph.....	197
7.2.5 LC-UV Quantification .....	197
7.2.6 LC-MS.....	198
7.2.7 Data Analysis.....	199
<b>7.3 Results and Discussion</b> .....	200
7.3.1 CIL LC-MS Analysis of Labeled Hemolymph.....	200
7.3.2 Metabolome Profile of Silkworm Hemolymph.....	203
7.3.3 Metabolome Comparison of Silkworm Hemolymph with Different DDT Treatment	206
7.3.4 Significance of Candidate Metabolite Biomarkers .....	213
<b>7.4 Conclusions</b> .....	219
<b>Chapter 8 - Conclusions and Future Work</b> .....	222
<b>References</b> .....	229
<b>Appendix</b> .....	248

## List of Tables

<b>Table 2.1</b> Results of t-tests (p-values) showing the significance of ratio difference obtained from a given matrix group vs. the H <sub>2</sub> O/H <sub>2</sub> O group (the ratios for each metabolite are shown in Figure 2.4).....	50
<b>Table 2.2</b> Results of t-tests (p-values) showing the significance of ratio difference obtained from a given matrix vs. H <sub>2</sub> O (the ratio values for matrix solutions are shown in Figure 2.5B) .....	52
<b>Table 3.1</b> Description of the sample set. Number in the cell represents the number of individuals that meet the corresponding condition. BMI groups were defined as: Underweight (BMI < 18.5), Normal (18.5 < BMI < 24.9), and Overweight (BMI > 25.0) .....	64
<b>Table 3.2</b> 148 identified dansyl-labeled metabolites and 31 identified DMPA-labeled metabolites that have fold change (female/male) > 1.2 (or < 0.83) and q value < 0.05 for the difference between males and females. ....	98
<b>Table 3.3</b> 25 identified dansyl-labeled metabolites and 6 identified DMPA-labeled metabolites that have fold change (overweight/normal) > 1.2 (or < 0.83) and q value < 0.05 for the difference between the normal group and the overweight group. (Asterisk means the metabolite is also a significant metabolite for sex differences.) .....	110

<b>Table 3.4</b> 25 identified dansyl-labeled metabolites and one identified DMPA-labeled metabolite that have fold change (underweight/normal) > 1.2 (or < 0.83) and q value < 0.05 for the difference between the normal group and the underweight group. ....	112
<b>Table 4.1</b> List of 6 metabolites that have significant concentration changes among the 10 participants after coffee intake. ....	139
<b>Table 5.1</b> Average concentrations and RSDs of the 98 positively identified metabolites from one microliter of finger blood sample. ....	161
<b>Table 6.1</b> Baseline demographic and clinical characteristics. ....	170
<b>Table 6.2</b> List of 46 identified significant metabolites found in human serum samples that differentiate the PD group and the healthy control group. ....	180
<b>Table 6.3</b> List of 18 identified significant metabolites found in human serum samples that differentiate the PDND subgroup and the PDID subgroup. ....	182
<b>Table 7.1</b> List of positively identified metabolites showing significant concentration changes after DDT treatment. (Number in bold and italic means the fold change was $\geq 1.20$ or $\leq 0.83$ with p-value $\leq 0.05$ ) ....	210

## List of Figures

<b>Figure 1.1</b> Reaction schemes of (A) the dansyl-labeling and (B) the DMPA-labeling.....	16
<b>Figure 1.2</b> Box-and-whisker plot, showing the data distributions of the control group and the disease group.....	21
<b>Figure 1.3</b> Volcano plot, showing the significantly decreased metabolites (FC < 1.2, p-value < 0.05, in blue) and significantly increased metabolites (FC > 1.2, p-value < 0.05, in red).....	24
<b>Figure 1.4</b> (A) Schematic of the selection of criterion value and (B) the confusion matrix.....	26
<b>Figure 1.5</b> An ROC curve with AUC = 0.915, showing the optimal cut-off value (-0.109) with sensitivity of 83.7% and specificity of 84.5%.....	26
<b>Figure 1.6</b> Score plots of the same data set given by (A) PCA, (B) PLS-DA and (C) OPLS-DA.....	28
<b>Figure 1.7</b> (A) Permutation test accepts a valid model. (B) Permutation test rejects an over-fitted model.....	33

**Figure 2.1** LC-UV quantification of the total concentration of labeled metabolites in five mouse urine samples re-dissolved in water and PBS.....47

**Figure 2.2** (A) PLS-DA and (B) OPLS-DA score plots of dansylation LC-MS data obtained from five <sup>12</sup>C-dansylated mouse urine samples mixed with a <sup>13</sup>C-dansylated pooled sample. For each sample, three experimental replicates were performed. Label X/Y (X, Y=water or PBS) denotes a mixture of an individual urine re-dissolved in X and labeled with <sup>12</sup>C-dansylation and a pooled urine re-dissolved in Y and labeled with <sup>13</sup>C-dansylation.....48

**Figure 2.3** Distribution of the PBS/H<sub>2</sub>O ratios (i.e., peak pair ratio in PBS vs. peak pair ratio in water) for metabolites in urine #1.....49

**Figure 2.4** Relative intensities of four representative metabolites in urine #1 vs. a pooled urine determined from experimental triplicate analysis of four different mixtures as shown in Figure 2.2.....49

**Figure 2.5** (A) Standard addition curves for tyrosine in mouse urine labeled in water and PBS. (B) Comparison of relative intensities of tyrosine labeled in different matrices.....51

**Figure 2.6** Relative intensities of 16 amino acids as a function of NaCl concentration in the sample solution.....54

**Figure 2.7** Schematic of dansyl-labeling for amine in the presence of salts.....55

**Figure 2.8** (A) PLS-DA score plot for 5-fold diluted urine #1 labeled in water (red) and NaCl solutions (50 mM in orange, 150 mM in yellow, 250 mM in green, and 350 mM in blue). (B) PLS-DA score plot for the comparison of an undiluted urine sample labeled with dansylation (black) and 5-fold diluted urine samples labeled in water and NaCl solutions. The injection amount for all the samples based on LC-UV measurement was the same. (C) PLS-DA score plot for the comparison of 4-fold diluted urine #1 labeled at different concentrations of PBS solution (e.g., 25% PBS refers to 4-fold dilution of the PBS solution). For each sample, five experimental replicates were performed in (A) and (B), while three experimental replicates were performed in (C).....58

**Figure 3.1** Workflow of the non-targeted serum metabolome profiling using chemical isotope labeling and high-resolution LC-MS.....71

**Figure 3.2** (A) PCA score plot for dansylation LC-MS data obtained from 100 healthy subjects (in blue) and 20 QC runs (in red). (“PC” represents for “principal component” and the corresponding percentage is the percentage of the variance among all the data points that this principal component covers.) (B) PCA score plot for DMPA-labeled LC-MS data obtained from 100 healthy subjects (in blue) and 20 QC runs (in red).....72

<b>Figure 3.3</b> Distribution of the total metabolite concentration among the serum samples from 100 people.....	74
<b>Figure 3.4</b> (A) Injection optimization curve, showing the numbers of dansyl-labeled peak pairs when different amounts of the QC sample are analyzed by the LC-MS system. (B) Injection optimization curve, showing the numbers of DMPA-labeled peak pairs at different injection volumes of the QC sample.....	76
<b>Figure 3.5</b> (A) Distribution of relative standard deviations defining the inter-subject variability in the amine/phenol-containing metabolites. (B) Distribution of relative standard deviations defining the inter-subject variability in the carboxyl-containing metabolites.....	79
<b>Figure 3.6</b> (A) PLS-DA score plot for studying the statistical differences in the amine/phenol-submetbaolome between males and females. (B) Response permutation test result of the PLS-DA model for males vs. females (amine/phenol-submetablome).....	81
<b>Figure 3.7</b> (A) PLS-DA score plot for studying the statistical differences in the carboxyl-submetbaolome between males and females. (B) Response permutation test result of the PLS-DA model for males vs. females (carboxyl-submetablome).....	82

**Figure 3.8** (A) PLS-DA score plot of the combined amine/phenol-submetabolome and carboxyl-submetabolome, showing a separation between male and female subjects. (B) Response permutation test result of the PLS-DA model in Figure 3.8A. A total of 200 permutations were implemented.....83

**Figure 3.9** (A) Volcano plot for males vs. females (amine/phenol-submetabolome), showing 147 significantly decreased metabolites (fold change (female/male) < 0.83, p-value < 0.034) in blue and 59 significantly increased metabolites (fold change > 1.2, p-value < 0.034) in red. (B) Volcano plot for males vs. females (carboxyl-submetabolome), showing 7 significantly decreased metabolites (fold change < 0.83, p-value < 0.0062) in blue and 31 significantly increased metabolites (fold change > 1.2, p-value < 0.0062) in red.....85

**Figure 3.10** Box plots, demonstrating the metabolite concentration distributions in the male group and the female group for (A) tyrosine, (B) homovanillic acid, (C) hydroxyl-phenyllactic acid, and (D) cholic acid.....88

**Figure 3.11** (A) PLS-DA score plot for studying the statistical differences in the amine/phenol-submetbaolome between the young group (< 26 years old) and the old group (>= 26 years old). (B) PLS-DA score plot for studying the statistical differences in the carboxyl-submetbaolome between the young group (< 26 years old) and the old group (>= 26 years old).....89

**Figure 3.12** (A) PLS-DA score plot for studying the statistical differences in the amine/phenol-submetabolome between the normal group and the overweight group. (B) PLS-DA score plot for studying the statistical differences in the carboxyl-submetabolome between the normal group and the overweight group.....90

**Figure 3.13** (A) Response permutation test result of the PLS-DA model for overweight vs. normal (amine/phenol-submetabolome). (B) Response permutation test result of the PLS-DA model for overweight vs. normal (carboxyl-submetabolome).....91

**Figure 3.14** (A) Volcano plot for overweight vs. normal (amine/phenol-submetabolome), showing 18 significantly decreased metabolites (fold change (overweight/normal) < 0.83, p-value < 0.005) in blue and 22 significantly increased metabolites (fold change > 1.2, p-value < 0.005) in red. (B) Volcano plot for overweight vs. normal (carboxyl-submetabolome), showing 4 significantly decreased metabolites (fold change < 0.83, p-value < 0.00066) in blue and 2 significantly increased metabolites (fold change > 1.2, p-value < 0.00066) in red. (C) Volcano plot for underweight vs. normal (amine/phenol-submetabolome), showing 2 significantly decreased metabolites (fold change (underweight/normal) < 0.83, p-value < 0.0016) in blue and 27 significantly increased metabolites (fold change > 1.2, p-value < 0.0016) in red. (D) Volcano plot for underweight vs. normal (carboxyl-submetabolome), showing one significantly increased metabolite (fold change > 1.2, p-value < 0.000018) in red.....93

**Figure 3.15** Box plots, demonstrating the metabolite concentration distributions in the normal group and the overweight group for (A) taurine and (B) palmitic acid.....93

**Figure 3.16** Venn diagram for the volcano plot significant metabolites associated with sex, age and BMI. (All of these numbers refer to dansyl-labeled metabolites only.).....95

**Figure 3.17** (A) PLS-DA score plot for studying the statistical differences in the amine/phenol-submetabolome between the normal group and the underweight group. (B) PLS-DA score plot for studying the statistical differences in the carboxyl-submetabolome between the normal group and the underweight group.....95

**Figure 3.18** (A) Response permutation test result of the PLS-DA model for underweight vs. normal (amine/phenol-submetabolome). (B) Response permutation test result of the PLS-DA model for underweight vs. normal (carboxyl-submetabolome).....96

**Figure 4.1** Work flow of the finger blood analysis based on CIL-LC-MS.....126

**Figure 4.2** (A) Numbers of peak pairs detected with different protein precipitation solvents. (B) Numbers of peak pairs detected with different concentrations of dansyl chloride reagent. (C) Numbers of peak pairs detected with different reaction lengths and reaction temperatures.....129

**Figure 4.3** (A) Numbers of peak pairs detected with different starting volumes of the USS sample (The injection volume was fixed at 15  $\mu$ L). (B) Metabolome coverage with different starting sample volumes when the nano-LC-MS system is used.....131

**Figure 4.4** (A) PLS-DA score plot, showing that the clusters of data points before and after coffee intake overlap each other. The performance indicators,  $R^2$  and  $Q^2$ , are 0.995 and 0.754, respectively. (B) Permutation test result, confirming that the statistical separation between the two study groups is not valid.....137

**Figure 4.5** (A) Volcano plot, showing 6 significantly increased metabolites (in blue) after coffee intake. (B) Box plot of theophylline, showing the distributions of its blood concentration before and after coffee. (C) Box plot of catechol sulfate, showing the distributions of its blood concentration before and after coffee.....138

**Figure 4.6** Concentration-time curves of (A) theophylline, (B) catechol sulfate, (C) glutamic acid and (D) proline, illustrating their concentration changes during the 4-hour period after coffee intake. The blue curve represents data acquired on Day 1, the red curve represents Day 2 and the green curve represents Day 3.....142

**Figure 4.7** Concentration-time curves of (A) catechol sulfate, (B) glutamic acid and (C) proline, showing the day-to-day variations of their blood concentrations during a week..143

**Figure 5.1** Numbers of peak pairs detected from different volumes of USS sample by dansyl-labeling (in blue) or DMPA-labeling (in red).....151

**Figure 5.2** Distribution of relative standard deviations defining the day-to-day variability in (A) amine/phenol-containing metabolites and (B) carboxyl-containing metabolites...153

**Figure 5.3** The concentration-day curves of (A) alanine, (B) tyrosine, (C) taurine and (D) salicylic acid, showing variations in their concentrations during the week.....155

**Figure 5.4** The concentration-day curves of (A) heptanoic acid, (B) ursodeoxycholic acid, (C) perillic acid and (D) hexadecanedioic acid, showing variations in their concentrations during the week.....156

**Figure 5.5** Concentration-time curves of (A) theophylline, (B) taurine, lactic acid, (C) sorbic acid and (D) citric acid, showing the concentration changes after the energy drink intake.....158

**Figure 5.6** PLS-DA score plots, showing the statistical differences between samples collected at different time points in (A) amine/phenol-submetabolome and (B) carboxyl-submetabolome.....160

**Figure 6.1** (A) PLS-DA and (B) OPLS-DA score plots of dansylation LC-MS data obtained from 42 healthy controls (in green) and 43 PD patients (in red). (“PC” represents

for “principal component” and the corresponding percentage is the percentage of the variance among all the data points that this principal component covers.)  $R^2$  and  $Q^2$  values given by cross-validation are: 0.977 and 0.791 for PLS-DA; 0.974 and 0.866 for OPLS-DA. (C) Response permutation test result of the PLS-DA model in Figure 6.1A. (D) Volcano plot of the comparison between healthy control and PD showing 28 variables with  $FC > 1.2$ ,  $q < 0.1$  (in red) and 48 variables with  $FC < 0.83$ ,  $q < 0.1$  (in blue).....174

**Figure 6.2** (A) PLS-DA and (B) OPLS-DA score plots of dansylation LC-MS data obtained from 27 PD patients without dementia (PDND in red) and 16 PD patients with incipient dementia (PDID in blue).  $R^2$  and  $Q^2$  values given by cross-validation are: 0.974 and 0.866 for PLS-DA; 0.982 and 0.813 for OPLS-DA. (C) Response permutation test result of the PLS-DA model in Figure 6.2A. (D) Volcano plot of the comparison between the PDND subgroup and the PDID subgroup showing 21 variables with  $FC > 1.2$ ,  $q < 0.1$  (in red) and 15 variables with  $FC < 0.83$ ,  $q < 0.1$  (in blue).....175

**Figure 6.3** (A) The receiver operating characteristic curve generated by the random forest model using the following 5 metabolite biomarker candidates: vanillic acid, 3-hydroxykynurenine, isoleucyl-alanine, 5-acetylamino-6-amino-3-methyluracil, and theophylline. (B) The receiver operating characteristic curve generated by the random forest model using the following 8 metabolite biomarker candidates: His-Asn-Asp-Ser, 3, 4-dihydroxy-phenyl-lactone, desamino-tyrosine, hydroxy-isoleucine, alanylalanine, putrescine [-2H], purine [+O], and its riboside.....179

**Figure 6.4** Box plots of the relative concentrations of vanillylmandelic acid (A), vanillylmandelic acid-isomer (B), vanillic acid (C), and vanillic acid (D) in the control group, the Parkinson's disease with no dementia subgroup, and the PD patients with incipient dementia subgroup.....186

**Figure 6.5** Box plots of the relative concentrations of tryptophan (A), kynurenine (B), and 3-hydroxykynurenine (C) in the control group, the Parkinson's disease with no dementia subgroup, and the PD patients with incipient dementia subgroup.....187

**Figure 6.6** (A) A simplified schematic of the caffeine metabolism pathway. Box plots of the relative concentrations of theophylline (B), 5-acetylamino-6-amino-3-methyluracil (C) and xanthine (D) in the control group, the PDND subgroup and the PDID subgroup.....188

**Figure 6.7** Box plots of the relative concentrations of methionine sulfoxide (A), methionine sulfoxide-isomer (B) and citrulline (C) in the control group, the PDND subgroup and the PDID subgroup.....190

**Figure 7.1** Workflow for silkworm hemolymph metabolome profiling using CIL LC-MS.....200

**Figure 7.2** (A) Averaged total concentrations of labeled metabolites in five groups of hemolymph samples. (B) A representative base-peak ion chromatogram obtained from LC-FTICR-MS analysis of a  $^{12}\text{C}$ -/ $^{13}\text{C}$ -labeled hemolymph sample. (C) Expanded mass

spectrum showing peak pair of serine with the m/z 339.0957 peak from the <sup>12</sup>C-labeled serine in an individual sample and the m/z 341.1032 peak from the <sup>13</sup>C-labeled serine in the pooled sample. (D) The peak pair number detected and the percentage of common peak pairs as a function of the number of samples.....201

**Figure 7.3** The Venn diagram of the metabolite distribution for the five comparative groups.....204

**Figure 7.4** (A) PCA score plot for QC data (yellow), control group (pink) and the samples from different concentrations of DDT treatment (1 ppm in red, 0.1 ppm in green, 0.01 ppm in dark blue and 0.001 ppm in sky blue). (B) PLS-DA score plot for the control group (pink) and the samples from different concentrations of DDT treatment (1 ppm in red, 0.1 ppm in green, 0.01 ppm in dark blue and 0.001 ppm in sky blue).....207

**Figure 7.5** List of 40 significant metabolites with the highest PLS-DA VIP scores, showing the correlation between the metabolite concentration in each group and their DDT concentrations.....208

**Figure 7.6** Volcano plots (fold change  $\geq 1.2$  and  $p \leq 0.05$ ) for the binary comparison of (A) the 1 ppm DDT group vs. the control group, (B) the 0.1 ppm DDT group vs. the control group, (C) the 0.01 ppm DDT group vs. the control group, and (D) the 0.001 ppm DDT group vs. the control group.....209

**Figure 7.7** (A) Overview of metabolic pathway analysis. (B) Pathway of glycine, serine and threonine metabolism.....213

**Figure 7.8** Box plots of four metabolites showing their relative concentrations in the control group and four DDT treatment groups.....215

## List of Abbreviations

ANOVA	Analysis of Variance
AUC	Area Under the Curve
BMI	Body Mass Index
BPC	Base Peak Chromatogram
CID	Collision-Induced Dissociation
CIL	Chemical Isotope Labeling
CSF	Cerebrospinal Fluid
DBS	Dried Blood Spot
DDT	Dichlorodiphenyltrichloroethane
DMPA	p-Dimethylaminophenacyl
Dansyl	5-(Dimethylamino)naphthalene-1-sulfonyl
EDTA	Ethylenediaminetetraacetic Acid
EI	Electron Ionization
EML	Evidence-based Metabolome Library
ESI	Electrospray Ionization
FC	Fold Change
FDA	Food & Drug Administration
FT-ICR-MS	Fourier Transform-Ion Cyclotron Resonance-Mass Spectrometry
GC	Gas Chromatography
GC-MS	Gas Chromatography-Mass Spectrometry
HILIC	Hydrophilic Interaction Liquid Chromatography
HMDB	Human Metabolome Database
HPLC	High Performance Liquid Chromatography
LC-MS	Liquid Chromatography-Mass Spectrometry
LC-UV	Liquid Chromatography-Ultraviolet
m/z	Mass to Charge Ratio
MS	Mass Spectrometry
MS/MS	Tandem Mass Spectrometry
NMR	Nuclear Magnetic Resonance
OPLS-DA	Orthogonal Partial Least Squares-Discriminant Analysis
PBS	Phosphate Buffered Saline
PC	Principal Component
PCA	Principal Component Analysis
PD	Parkinson's disease
PDID	Parkinson's disease with incipient dementia
PLS-DA	Partial Least Squares-Discriminant Analysis
ppm	part(s) per million
QC	Quality Control
QTOF-MS	Quadrupole Time-Of-Flight Mass Spectrometry

ROC	Receiver Operating Characteristic
RPLC	Reversed Phase Liquid Chromatography
RSD	Relative Standard Derivation
RT	Retention Time
SVM	Support Vector Machine
UPLC	Ultra Performance Liquid Chromatography
USS	Universal Serum Standard
VIP	Variable Importance on the Projection

# Chapter 1

## Introduction

### 1.1 Introduction to Blood Biomarkers

A biomarker has been broadly defined as “a characteristic that is objectively measured and evaluated as an indicator of normal biological processes, pathogenic processes or pharmacological responses to a therapeutic intervention”.<sup>1</sup> This definition may include a thermometer readout, a urine test of a drug product, or whole-genome sequencing. During the past decade, biomarker-driven analyses of drug candidates have greatly boosted the decision-making process in the pharmaceutical drug development, significantly reducing the cost and risk of this expensive business.<sup>2-3</sup> Unlike pharmaceutical companies, most academic researchers study the disease biomarkers, which can help with the understanding, prediction, diagnosis, prognosis, and treatment assessment of a disease state.<sup>4</sup>

A diagnostic biomarker can differentiate patients with a specific disease from healthy people. For instance, hemoglobin A1c may serve as a diagnostic biomarker for type II diabetes.<sup>5</sup> This example elucidates an important point that a disease biomarker is not necessarily involved in the pathogenic mechanism of the disease. Nonetheless, many biomarkers are directly related to the onset of the disease, and studying the biological processes behind them can profoundly deepen our understanding of the disease

mechanisms. Importantly, diagnostic biomarkers can sometimes demonstrate detectable changes before the disease symptoms become noticeable, enabling early diagnosis of the disease. With preventive interventions, the onset and progression of the disease may be dramatically delayed. Considerable efforts have been made to discover biomarkers for risk assessment or early diagnosis of cancer.<sup>6-8</sup> In addition, a prognostic biomarker would indicate a likely outcome of the disease independent of the treatment.<sup>9</sup> Thus, a biomarker could play both diagnostic and prognostic roles, according to its extent of alteration. For example, biomarker candidates have been reported for both the diagnosis and cognitive impairment prediction of neurodegenerative diseases.<sup>10-11</sup>

Furthermore, biomarker-based analysis has opened new possibilities for personalized medical treatment. Since inter-patient variability in efficacy and toxicity has been observed for many medications,<sup>12</sup> biomarkers that can predict the drug responses may lead the way in personalized treatment design. At present, U.S. Food & Drug Administration (FDA) has listed 198 drugs with pharmacogenomic biomarkers, which can direct the optimization of drug dose and management of adverse effects.<sup>13-14</sup> The presence of certain biomarkers could determine the efficacy of the medication in a specific patient. While other biomarkers, involved in the metabolism of the drugs, could be used to vary the optimal dose among patients.

The heterogeneity in medication response has also led us to notice the inter-individual or inter-population differences in disease pathophysiology.<sup>15-16</sup> In addition to genetic factors,

the environment can also play a fundamental role in these differences. For example, Ng, Poon and their coworkers studied the cause of liver cancers and reported that in Taiwan, 78% of the patients showed significant signature of exposure to aristolochic acids, which are strong mutagens and mainly originate from various Chinese herbs, while in North America only 5% of the subjects demonstrated obvious signature of the aristolochic acids.<sup>17</sup> Moreover, studies on many diseases that were believed to be highly genetic have suggested that the gene-environment interaction is the true cause.<sup>18</sup> In recent years, the dramatically increasing prevalence of autism, which was thought to be a genetic disease, has convinced researchers to look into environmental exposures during pregnancy.<sup>19</sup> This example illustrates the importance of studying environmental factors for the understanding and treatment of diseases. The biomarkers for the assessment of environmental exposures are called exposure biomarkers, and they could be adopted in the early-stage diagnosis of various gene-environment-cause disorders, such as neurological diseases.<sup>4</sup>

Overall, disease-related biomarkers are commonly discovered by studying patients or animal models with a specific disease or environmental exposure. Before a biomarker is used in clinical practice, the workflow usually includes candidate discovery, biomarker validation and clinical assay development.<sup>20-21</sup> According to how they work in clinical diagnosis, biomarkers can be classified into three types. The first type refers to the biomarkers that only exist in the diseased group, due to abnormally altered metabolic pathways or a specific environmental exposure. This is an ideal case since any qualitative detection of the biomarker is a strong indicator of the disease state. Most of the second type biomarkers are present in both the healthy and diseased populations, but there are

significant concentration differences between the two groups. The second type of biomarkers requires accurate quantification results for the clinical applications. Thirdly, levels of a group of biomarkers are sometimes used to increase the diagnosis power. In this case, they are called a biomarker panel, which can be a linear combination of multiple compounds.<sup>22-23</sup>

Blood is the major biofluid of the human body, delivering nutrients and oxygen to tissues and transporting waste products away. Blood is also a complex treasury of blood cells, proteins, mineral ions, hormones and other small molecules, making itself an ideal material for biomarker studies. Unlike most of other biofluids (e.g., urine and sweat), which mainly contain the waste products of the body, blood can more accurately represent the ongoing biological processes in tissues. Importantly, compared with tissue samples or cerebrospinal fluid (CSF), the collection of blood is much less invasive.

Blood analysis has a long history in clinical analysis. Analyses of blood cells (e.g., a complete blood count<sup>24</sup>) or even the physical properties of blood (e.g., a hemorheology test<sup>25</sup>) have been used as clinical indicators. A large number of compounds in blood have been found to be biomarker candidates. The application of blood glucose or hemoglobin A1c for monitoring diabetes is a commonly seen example. More comprehensive blood-based biomarker analyses are also emerging for cancer screening. FDA has approved a handful of blood-based assays for the diagnosis or prognosis of cancer. For instance, researchers have found that several proteins in blood could improve the diagnostic power

for the detection of early-stage ovarian cancer. After seven years of assay development and validation, a biomarker panel of five proteins became the OVA1, an FDA-approved biomarker assay.<sup>26</sup> Since a large number of blood biomarker candidates have been reported, more and more validated blood-based biomarker assays will become commercialized in the future.

## **1.2 Small-molecule Biomarkers and Metabolomics Analysis**

Although many of the currently approved biomarker assays are based on RNAs or proteins, small-molecule biomarkers have attracted increasing attention. Small molecules are building blocks of larger biological components, function regulators of cells and messengers in signaling. Compared with the genome level or proteome level, small molecules are more time sensitive to the biological changes in the body. Also, because small molecules are directly involved in the body-environment interactions, they can provide firsthand information about the environmental stimulations. Most prescription drugs are small molecules, and personalized health management based on small molecules has never been far from our daily life, such as taking vitamin C for the prevention of scurvy and using glucosamine to delay the progression of osteoarthritis.<sup>27</sup>

Because of the importance of small molecules, metabolomics, which is the comprehensive and systematic analysis of small biological molecules in a given biological subject, has become a rapidly growing field in biomarker discovery. Metabolites are the intermediates and products of metabolism. In metabolomics, the term usually covers all the meaningful

small molecules (MW < 1,500 Da), including but not limited to endogenous nucleosides, oligonucleotides, amino acids, peptides, alcohols, amines, organic acids, ketones, aldehydes, lipids, steroids, as well as exogenous nutrients, food additives, toxins, and pollutants. The word “metabolome”, which refers to the complete set of metabolites synthesized by an organism, was first used by Olivier et al. in 1998.<sup>28</sup> As the downstream end of genome, transcriptome and proteome, the metabolome is at the frontier of system biology<sup>29</sup> and carries valuable information for viewing the whole picture of biological processes. Based on the hypothesis that certain metabolic changes may occur in the human body when a disease state develops, metabolites can potentially become powerful biomarkers for the early-stage diagnosis of diseases. Additionally, in some environmental studies, the total set of chemicals from environmental exposures is defined as exposome, and the corresponding assessment is called exposomics.<sup>30</sup>

Blood serum is a primary carrier of metabolites, transporting all the small molecules that are being secreted or excreted by different tissues in response to various physiological conditions.<sup>31</sup> Many blood biomarker candidates have been reported in the literature. For example, Sato and his coworkers have shown that the blood concentration of a small molecule, desmosterol, could be utilized for the diagnosis of Alzheimer’s disease.<sup>32</sup> Nishiumi et al. reported 18 blood metabolites to be biomarker candidates for the detection of pancreatic cancer.<sup>33</sup> Efforts have also been made to profile the whole blood metabolome. Lawton et al. applied multiple analytical methods to study the blood metabolome of 269 healthy adults and successfully quantified 300 metabolites.<sup>34</sup> However, this number is far lower than the total number of blood metabolites. The Human Serum Metabolome

Database lists 4,229 confirmed human serum metabolites.<sup>31</sup> In addition to the low-abundance endogenous metabolites that have not been detected, there are also numberless food/environment-originated compounds and their downstream metabolites existing in the blood, so the size of the blood metabolome remains unknown. With more sensitive analytical techniques in the future, more blood metabolites will continue to be discovered and contribute to the realm of biomarkers.

Generally speaking, metabolomics studies are inductive rather than deductive, expanding the limited biological understanding by digging into the huge amount of data. There are two routes of analyses for blood metabolome: targeted and non-targeted, depending at which stage the metabolite identification is performed. Targeted analysis focuses on a group of metabolites, which are usually associated with particular pathways or environmental exposures. After the samples are analyzed, the target metabolites are identified and quantified based on the raw data. The results are then used for biological interpretation. A targeted analysis is usually optimized to have high sensitivity and accuracy for the pre-defined group of metabolites.<sup>35-36</sup> The high-confidence information of metabolite identities and quantities enables accurate and clear interpretation of the result. However, targeted analyses cannot find new biomarkers or new pathways.

In contrast, a non-targeted analysis aims at detecting as many metabolites as possible in a given biological system. In this kind of studies, data processing is done before the identification of metabolites. Statistical tools are used to investigate the raw data and to

select the compounds that can significantly reflect the overall metabolome changes. After that, researchers identify these compounds with databases and try to interpret the biological implications. The non-targeted approach offers opportunities to discover more blood metabolites and metabolic pathways. However, non-target analysis is facing several technical challenges. First, unlike proteins which are composed of 20 amino acids and have similar chemical characteristics, metabolites encounter enormous diversity in physical and chemical properties, which has been the prime hurdle in detecting all the metabolites with a single platform. To maximize the metabolome coverage, many non-targeted blood metabolome profiling reports utilize a combination of multiple analysis techniques, such as gas chromatography-mass spectrometry (GC-MS) and liquid chromatography-mass spectrometry (LC-MS). Second, the concentrations of blood metabolites spread over a vast range of 9 orders of magnitude.<sup>37</sup> Technical improvements are required to correct the bias towards detection of high-abundance metabolites. At last, metabolite identification is the most challenging part. In LC-MS-based metabolomics, positive identification, which means metabolites are matched to known standards, is highly preferred. However, the number of standards available in a research laboratory is always limited, and it is not practical to use standards to confirm every single LC-MS peak, let alone a large portion of the blood metabolome are unstudied compounds without available standards. Therefore, putative identification is often used, complementary to the positive identification.

Putative identification is conducted based on an accurate mass match or MS/MS match to the metabolite databases. The Human Metabolome Database (HMDB)<sup>38</sup> is one of the most widely used metabolite databases. Until now, the HMDB database has been enriched to

114,100 metabolites, including endogenous metabolites as well as 2,800 drug metabolites, 3,670 toxins, and 28,000 food components. Created in 2004, METLIN<sup>39</sup> has grown to a popular database that includes 961,829 compounds, ranging from endogenous metabolites to small peptides, drugs and toxins. Among them, over 14,000 metabolites have been individually analyzed, and another 200,000 have *in silico* MS/MS data. In addition, MassBank<sup>40</sup> is another rapidly growing small-molecule database of mass spectral, especially ESI-MS<sup>2</sup> spectra. Considering a large number of the unknown metabolites are biologically modified products of the known metabolites, our lab has developed an Evidence-based Metabolome Library (EML), which employs 76 commonly encountered metabolic reactions for the *in silico* biotransformation of 8,021 known endogenous metabolites provided by HMDB, expanding the coverage to 375,809 one-reaction-derived metabolites and 10,583,901 two-reaction-derived metabolites.<sup>41</sup> To enable the MS/MS match, the predicted MS/MS spectra have also been simulated for the 8,021 known metabolites and 375,809 predicted metabolites.<sup>42</sup> Overall, taking advantage of the metabolite-specific databases and high-accuracy mass measurement enabled by high-resolution mass analyzers, we are able to generate putative identification results at acceptable confidence.

A number of non-targeted metabolome profiling works of blood samples have been published.<sup>34, 43-46</sup> Many of them adopted a combination of LC-MS and GC-MS to increase the metabolome coverage, but none of them reported more than 500 positively identified metabolites. Although some works provided 1,000 to 5,000 metabolite features (a chromatographic peak with a unique m/z), a large portion of them were not truly

metabolites, in most cases. More advanced analytical platforms with higher metabolome coverage and stronger identification power are needed for future studies in blood metabolome.

### **1.3 Major Metabolomics Platforms for Blood Biomarker Discovery**

As discussed before, most blood biomarkers are concentration-based, so quantitative analyses with high accuracy are always preferred in blood metabolomics. Because many blood metabolites are at very low concentrations, the analysis should also provide adequate sensitivity. Although most clinical blood collections collect at least 4.0 mL of blood, the samples are commonly used for multiple analyses or studies, and the amount available in one specific experiment is usually limited. Therefore, the required sample amount is also a factor for evaluating the performance of a blood analysis platform. Among a number of analytical methods, nuclear magnetic resonance spectroscopy (NMR) and chromatography-mass spectrometry are the major platforms for blood metabolomics study due to their relatively high metabolome coverage and abundant database resources.

The major advantage of NMR analysis is the ease of methodology.<sup>47</sup> In general, sample preparation just involves mixing the plasma or serum sample with a buffer and transferring the mixture into an NMR tube.<sup>48</sup> In NMR analysis, the signals correlate directly and linearly with compound abundance, and the reproducibility is very high.<sup>49</sup> The metabolite identification can be done with several commercialized spectra libraries. HMDB also provides NMR spectra for a part of its library. A non-negligible drawback of NMR is the

relatively low sensitivity. High-resolution NMR analysis requires a minimum metabolite concentration of 5  $\mu\text{M}$ ,<sup>50</sup> which is larger than the level of many blood metabolites. Subsequently, the required sample amount in NMR analysis is also relatively high, typically 200  $\mu\text{L}$  of serum or plasma.<sup>48</sup> Since the NMR analysis is not destructive to the samples, we can potentially reuse the sample in MS-based analyses,<sup>51</sup> despite the fact that the buffer may cause severe matrix effects.

Chromatography-MS-based methods are more sensitive than NMR by several orders of magnitude. Among the mass analyzers, quadrupoles and ion traps have been widely used in metabolomics applications but limited by mass accuracy and mass resolution in metabolite identification.<sup>52</sup> These are good options in targeted analyses which have internal standards for metabolite identification, and can attempt to reach a balance between the performance and ease of maintenance. Particularly, they offer excellent detection sensitivity in the MS/MS mode.<sup>53</sup> For non-targeted metabolome profiling, high-resolution mass analyzers, such as Quadrupole-time-of-flight-MS (Q-TOF-MS), Fourier transform-ion cyclotron resonance-MS (FTICR-MS) and Orbitrap-MS, are often used.<sup>54</sup> FTICR analyzers have the highest resolving power and mass accuracy of all mass analyzers.<sup>55</sup> At  $m/z$  400, a 21-Tesla FT-ICR-MS can reach the resolving power of more than 300,000 for a 0.76 s detection period.<sup>56</sup> However, the scan speed of FTICR-MS is relatively low, and the cost of the instrumentation can be very considerable.<sup>57</sup> Although the typical resolution for  $m/z$  400 on our Q-TOF-MS is 30,000, it provides higher scan speed and acceptable mass accuracy for identification. Importantly, our group previously found that by injecting a large but optimized amount of sample into a Q-TOF-MS that equipped with a high-

dynamic-range (HD) mode, the resulting metabolome coverage is higher than those on a non-HD Q-TOF platform or an FT-ICR platform.<sup>58</sup> Orbitrap-MS is a relatively new technology. FT-ICR and Orbitrap share a number of similar features, such as the high resolving power, while Orbitrap-MS has a much smaller size. However, the mass analyzing process in an Orbitrap-MS usually does not favor low-abundance species, which might limit its applications in biomarker discovery.

In order to reduce the complexity of the spectrum and thereby increase the metabolome coverage, a chromatographic separation is commonly used before the mass analyzer. The GC-MS platform using capillary columns has superior reproducibility and relatively low cost.<sup>59</sup> Taking advantage of the highly reproducible GC retention indices and electron ionization (EI) spectra, researchers have constructed comprehensive standard libraries compatible with data from different laboratories, regardless of the manufacturer of the instrument. In addition to the high-confidence identification, GC-MS requires a smaller volume of the serum sample, typically 10 to 50  $\mu\text{L}$ .<sup>44, 60-61</sup> However, the metabolome coverage of GC-MS is limited because the metabolites must be volatile or can be volatile after derivatization.<sup>62</sup> The targeted metabolites mainly have low molecular weights and low boiling points, such as alcohols and derivatized amino acids.

Being complementary to GC-MS, the LC-MS platform is capable of detecting metabolites with higher molecular weights or boiling points. In LC-MS applications, the reversed-phase-liquid-chromatography (RPLC) is used for moderately polar and non-polar

metabolites, and the hydrophilic-interaction-chromatography (HILIC) is used for highly polar compounds.<sup>63</sup> Unlike GC-MS, LC-MS employs a soft ionization method, the electrospray ionization (ESI), to generate ions,<sup>64</sup> and the MS/MS spectra can be obtained from collision-induced-dissociation (CID).<sup>65</sup> With the accurate mass calculated by the measured  $m/z$  of the  $(M+H)^+$  ion (positive mode), putative identification can be easily performed by searching through the metabolite databases. However, the retention time and CID spectra are not reproducible between different systems.<sup>62</sup> With 50 to 200  $\mu\text{L}$  of serum or plasma sample, LC-MS analysis can detect more than 1,000 metabolite features,<sup>66-68</sup> which are defined as chromatographic peaks with a specific retention time and a unique  $m/z$ . It has long been a problem to differentiate the weak signal of very low-abundant metabolites from the background noise. Also, a single metabolite may be detected in multiple forms, including adduct ions, in-source fragment ions, dimers, trimers, etc.<sup>62</sup> Therefore, a large portion of the metabolite features are not truly existing metabolites. Another disadvantage of the LC-MS platform is the variability in quantification due to ion suppression. Sample matrix or coeluting compounds can contribute to this effect and interfere with the quantification of metabolites.<sup>69</sup> When a large set of samples is being analyzed, the minor variability in sample matrix may significantly affect the quantification and mislead the data interpretation.

#### **1.4 Chemical Isotope Labeling in LC-MS-based Metabolomics**

As discussed above, even though the combination of GC-MS and LC-MS has been used in many metabolomics studies, the metabolome coverage remains low. The discovery of non-volatile, low-abundance blood biomarkers depends on the improvement of LC-MS

technique. To improve the quantification accuracy of LC-MS, we can introduce an internal standard, preferably an isotopic internal standard, into the sample to overcome the ion suppression effect via relative quantification. In non-targeted metabolomics, it is not possible to acquire the isotopic standard for every single metabolite. Alternatively, chemical derivatization to the metabolites with isotopic tag groups can realize the relative quantification for each metabolite, as long as it exists in both of the two comparative samples. More specifically, a metabolite in the sample being studied is derivatized by a labeling reagent, while in a reference sample, the same metabolite is labeled by the isotopic counterpart of the labeling reagent. After the labeled samples are mixed, the metabolite from the reference sample serves as the internal standard, and the measured concentration is relative to it. This process is called differential stable isotope labeling.

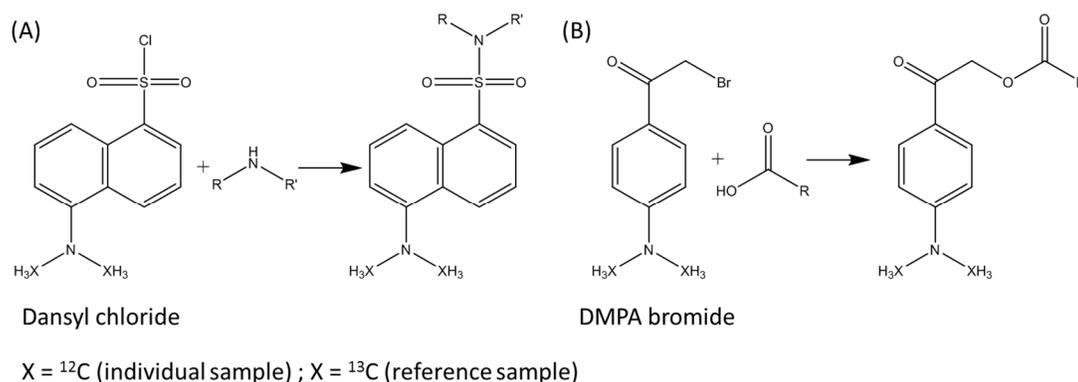
Previously, our group has reported a chemical isotope labeling (CIL) technique using dansyl chloride as the labeling reagent,<sup>70</sup> and applied this method to the LC-MS-based metabolomics. In this platform, a pooled sample works as the reference sample for all the individual samples. Each individual sample is labeled by the  $^{12}\text{C}$ -dansyl reagent, and the pooled sample is labeled by the  $^{13}\text{C}$ -dansyl reagent. Then the mixed sample is analyzed by LC-MS. For each metabolite, a peak pair is detected instead of a single mass peak. The light peak of the pair represents the  $^{12}\text{C}$ -labeled metabolite from the individual sample, and the heavy peak of the pair is the  $^{13}\text{C}$ -labeled metabolite from the pooled sample. To avoid the interferences from the natural isotopic peak of the  $^{12}\text{C}$ -reagent, the  $^{13}\text{C}$ -reagent has two  $^{13}\text{C}$  atoms, and the distance between the two peaks equals 2.00671 Da (Appendix Figure 2). We measure the relative concentration by calculating the intensity ratio of the two peaks

in a pair. Consequently, every metabolite has a corresponding internal reference to accurately measure its relative concentration. Although the quantification is relative, the information is adequate for metabolomics analysis to find the metabolites with significant changes. Absolute quantification of confirmed biomarker candidates can be conducted afterward.

Importantly, the chemical isotope labeling also benefits the LC-MS-based metabolomics in other ways. First, after being attached to the hydrophobic dansyl group, the polar metabolites become relatively non-polar, improving their separation on an RPLC column and reducing the ion suppression effect due to co-elution. Second, the non-polar dansyl group provides good surface activity during the ESI process,<sup>71</sup> and its tertiary amine group has excellent chargeability in the positive mode. Consequently, the detection sensitivity and metabolome coverage are significantly enhanced. Third, the addition of the dansyl group shifts low-mass metabolites to the higher mass region, which usually has cleaner background, so the signal-to-noise ratio is also improved.

The dansyl-labeling platform detects metabolites in forms of peak pairs instead of single mass peaks, making it easier to differentiate metabolites from the background noise peaks. Although sometimes there are still adduct ions, in-source fragment ions, and dimers, the IsoMS software,<sup>72</sup> which is designed to automatically pick the peak pairs, can filter out these interferences, as well as the constant background noises. Therefore, compared to the mass features in traditional LC-MS, the peak pairs are much more likely to be true

metabolites. Overall, the CIL LC-MS platform can significantly promote the blood metabolomics studies by quantifying a large number of high-confidence metabolites.



**Figure 1.1** Reaction schemes of (A) the dansyl-labeling and (B) the DMPA-labeling.

Although the dansyl-labeling is limited to the reaction with amine/phenol-containing metabolites, other CIL reagents have been developed to study more groups of metabolites and to expand the metabolome coverage. This strategy is called “Divide-and-Conquer”, which means dividing the whole blood metabolome into several submetabolomes, and then studying each of them with corresponding CIL methods. The p-dimethylaminophenacyl (DMPA) bromide reagent has been used to analyze the carboxyl-containing metabolites.<sup>73</sup> The reaction offers the same beneficial effects as dansyl-labeling does, and it has good selectivity against other functional groups. Recently, our group has reported another two labeling methods, focusing on the hydroxyl-submetabolome and carbonyl-submetabolome respectively.<sup>74-75</sup> With the four labeling methods, we are expecting more than 5,000 blood metabolites to be detected.

## 1.5 Blood Sample Handling Techniques

Despite the CIL LC-MS platform that can effectively overcome the variations during the detection, the sample preparation can also greatly affect the quantification results. For example, during the sample handling process, some impurities might be introduced to the sample, and at the same time, some metabolites might experience degradation. Sometimes other matrices are added to the sample, so the changes in matrix effect should be carefully evaluated.

Although there are a few studies on the blood cell metabolome,<sup>76</sup> most blood metabolomics studies remove blood cells and store blood samples in the form of serum or plasma. Venipuncture is usually performed in a clinic by a trained personnel, and the whole blood is collected into a specific blood collection tube according to the needs. For collecting serum samples, whole blood is allowed to naturally clot in a common plastic tube. After centrifuging, the supernatant above the clot is the serum. Some serum collection tubes are equipped with clot activators, which are usually silica beads. These tubes should be avoided as evidence has shown that the clot activators can cause interferences to the analysis (e.g., surface adsorption).<sup>77-78</sup> When whole blood is transferred into tubes coated with anticoagulants, the supernatant after centrifugation is called plasma. Commonly used anticoagulants include ethylenediaminetetraacetic acid (EDTA), citrate and heparin.<sup>79</sup> EDTA and citrate are very strong matrices that can significantly affect the quantification in the traditional LC-MS analysis, and heparin is relatively preferable.<sup>80</sup> Using the same

type of sample is recommended in a specific LC-MS-based study, and caution should be exercised when comparing metabolomics data obtained with different sample types.<sup>81</sup> Regarding the CIL LC-MS platform, we previously examined three types of plasma (EDTA, citrate, and heparin) and serum, and we proved that the differential isotope labeling could largely overcome the variability in metabolite detection and quantification.<sup>82</sup>

After collection, sample aliquots should be immediately frozen and stored at -80 °C. No detailed studies have assessed the stability of metabolites during storage, and the role of sample matrix (e.g., anticoagulant) is largely unknown. Nonetheless, it is recommended that the number of freeze-thaw cycles should be minimized and in one analysis all the samples should have experienced the same number of freeze-thaw cycles.

When venipuncture blood collection is not feasible or necessary, other blood collection approaches are used. For example, the measurement of blood glucose is often performed by the diabetes patients themselves with a finger stick and a blood glucose monitor.<sup>83</sup> Starting from the 1960s, dried blood spot (DBS) method has been applied to newborn screening.<sup>84</sup> The sample collection is done by a finger or heel prick, then typically 50 to 200 µL of whole blood is directly applied onto the sampling paper/card within a pre-marked circle. The sample is allowed to naturally dry at room temperature and then stored at room temperature or at -20 °C for many weeks, months or even years.<sup>85</sup> There are a few metabolomics studies based on DBS that have been reported.<sup>86-87</sup> However, the DBS is not

an ideal material for biomarker discovery due to the paper matrix, cross-contamination, metabolite degradation, etc.

Because of the increasing need for high-throughput metabolomics analysis and point-of-care diagnosis, microfluidic devices are rapidly emerging to increase the efficiency of sample preparation in metabolomics studies.<sup>88</sup> Even if the sample volume is very small, blood cell separation can be easily achieved in a microfluidic device.<sup>89</sup> More importantly, microchip-based LC separation has been added to mass analyzers,<sup>90</sup> and the microchip-LC-MS system has been widely applied to proteomics and other studies.<sup>91-92</sup> Compared to traditional methods, the microchip elutes a smaller volume into the mass analyzer, therefore generating a greater response.<sup>93</sup> With more microfluidic devices designed for blood metabolomics in the future, the metabolome coverage and analysis throughput will be significantly improved, making the metabolomics analysis feasible for personalized health monitoring or point-of-care diagnosis.

## **1.6 Statistical Analyses for Biomarker Discovery**

Although the metabolome coverage is very limited compared to the number of metabolites existing in blood, current metabolomics analyses are generating a massive amount of data. In a typical study with 200 subjects, detecting 1,000 metabolite features means 200,000 concentration values to be further analyzed. Therefore, both the discovery and application of biomarkers rely on statistics. The statistical analysis for biomarker discovery commonly follows two strategies: uni-variate analysis and multi-variate analysis. In these methods,

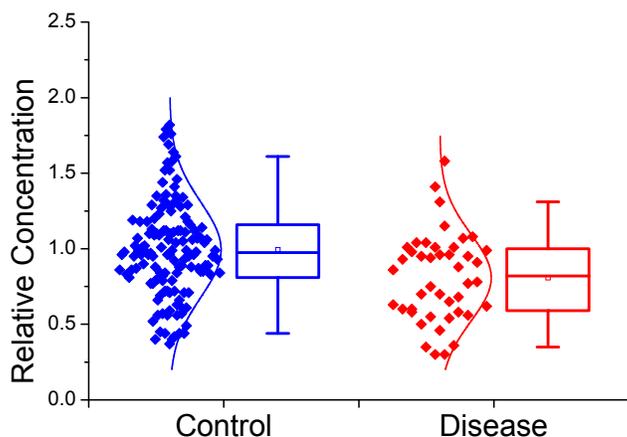
each metabolite (e.g., a mass feature or a peak pair) is considered as a variable, and each sample (e.g., a healthy participant or a patient) is regarded as an observation. The variable-observation matrix contains all the quantitative values we have measured. The uni-variate analysis studies one variable at a time, while the multi-variate analysis treats the matrix as a whole. The two approaches complement each other and usually used together in biomarker discovery studies.

### **1.6.1 Uni-variate Analysis**

In a biomarker discovery study, we usually have a control group and a disease group, and we assume the size of the control group is large enough so that these observations can statistically represent the distribution of the whole healthy population, and so is the disease group. The simplest way to visualize the concentration distribution of each metabolite is the box-and-whisker plot. A typical box plot, as shown in Figure 1.2, illustrates the 5th percentile (bottom whisker), 25th percentile (bottom of the box), median (middle line of the box), mean (small dot), 75th percentile (top of the box) and 95th percentile (top whisker). The difference between 25th percentile and 75 percentile is defined as the interquartile range (IQR), and sometimes (25th percentile – 1.5 IQR) and (75th percentile +1.5 IQR) are used to discover outliers. The box plot tells the difference between the two groups in an intuitive manner.

We also need statistical tools to study the difference numerically. For each variable (metabolite), we calculate the average concentration of the control group and the disease

group, respectively. Then the ratio of these two average values is defined as the fold change. The fold change tells the magnitude of the difference between the two groups, or theoretically the two populations. There is no standard cut-off value to determine a significant fold change. Most researchers set the cut-off according to their experimental design.



**Figure 1.2** Box-and-whisker plot, showing the data distributions of the control group and the disease group.

Although the fold change is a straightforward parameter, it does not consider the variability of data. Therefore, we need hypothesis testing to show the statistical significance of the inter-group difference. Welch's t-test and the one-way analysis of variance (one-way ANOVA) are the most used tools. Welch's t-test determines if the two populations have equal means by analyzing the two groups of data. The null hypothesis ( $H_0$ ) refers to the situation that the mean of the disease population equals to that of the healthy population. In other words, the disease does not cause any concentration change to this metabolite. The alternative hypothesis ( $H_1$ ) means that the means of the two populations are different,

indicating that the disease causes a noticeable change in the blood concentration of this metabolite. In this case, the metabolite becomes a biomarker candidate that can potentially be used to differentiate between patients and healthy people. If we choose to accept the alternative hypothesis when the null hypothesis is actually true, we would make a type I error (false positive). And if we mistakenly reject the alternative hypothesis, we would make a type II error (false negative). The hypothesis testing is designed to gauge the type I error.

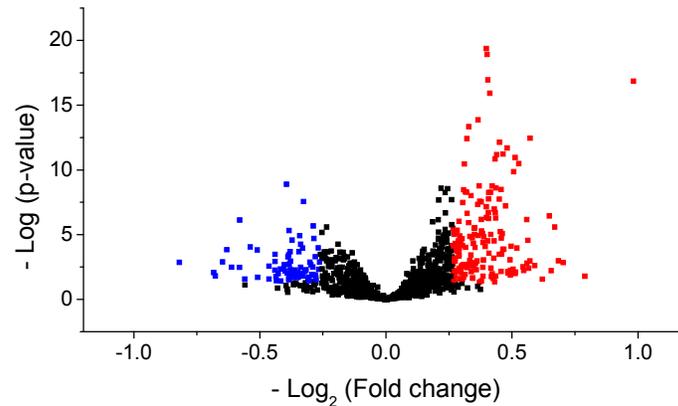
Instead of the t-value, we usually calculate the p-value in metabolomics analysis. The p-value is the probability of observing another set of data that is at least as extreme as the current observation, when the null hypothesis is true. Most researchers set a cut-off value of 0.05. When the p-value of a metabolite is smaller than 0.05, the difference between healthy people and patients is defined to be statistically significant. It is important to know that the p-value is not the chance of making a mistake by accepting the alternative hypothesis. Calculated based on the assumption that the null hypothesis is true, the p-value plays the role as a filter to select biomarker candidates but has no other clear statistical meanings. Furthermore, setting the cut-off at 0.05 is arbitrary. Some statisticians have criticised the misuse of the p-value in biological and medical sciences, and suggested lowering the cut-off to 0.005.<sup>94</sup> However, for biomarker discovery, we may care more about the type II error than the type I error. The uni-variate analysis is a preliminary screening of the potential biomarkers, and we do not want to miss any mistakenly. The false positives can be excluded in the following analyses or validation processes. To avoid increasing the chance of false negatives, we can keep using 0.05 as the cut-off. Another

concern is the normality of the data. Welch's t-test requires the data follow a normal distribution. If not, the Mann-Whitney u-test is supposed to be used.<sup>95</sup> Nonetheless, theoretical and data-based comparisons have elucidated that the performance difference between t-test and u-test is very minor.<sup>96-97</sup>

An alternative to the Welch's t-test is the Bayesian t-test. It compares the posterior odds of the two hypotheses and calculates the Bayesian factor  $\Omega = \Pr(H_0 | \text{data}) / \Pr(H_1 | \text{data})$ .<sup>98</sup> Unlike Welch's t-test, which does not tell the possibility that the alternative hypothesis is true, the Bayesian t-test quantitatively demonstrates which hypothesis is more likely to be true. If the Bayesian factor is 5,  $H_0$  is 5 times more probable than  $H_1$ , given the data. The Bayesian approach surpasses Welch's t-test in many aspects. However, since Welch's t-test has been widely used for many years, the application of Bayesian factor in metabolomics remains very scarce until now.

The Cohen's effect size (Cohen's d-value) considers both the difference in mean and the variability. It is the ratio of the mean difference to the pooled standard deviation. When  $d > 0.8$ , it is called a "large effect".<sup>99</sup> More metabolomics studies chose to use the volcano plot (Figure 1.3) to visualize both the fold change and the p-value. In the volcano plot,  $-\log(p\text{-value})$  is plotted against  $\log_2(\text{fold change})$ , making a volcano-shaped scatter plot. Each point in the volcano plot represents a metabolite, and "significant metabolites" refers to those whose fold changes are larger than the cut-off and p-values smaller than the criterion. Again, there is no gold standard to set the cut-offs. As the fold change is also considered,

the false positive rate should be lower than using the p-value alone. We should keep in mind that the volcano plot serves as a screening method, and the significant metabolites need further validations.



**Figure 1.3** Volcano plot, showing the significantly decreased metabolites ( $FC < 1.2$ ,  $p\text{-value} < 0.05$ , in blue) and significantly increased metabolites ( $FC > 1.2$ ,  $p\text{-value} < 0.05$ , in red).

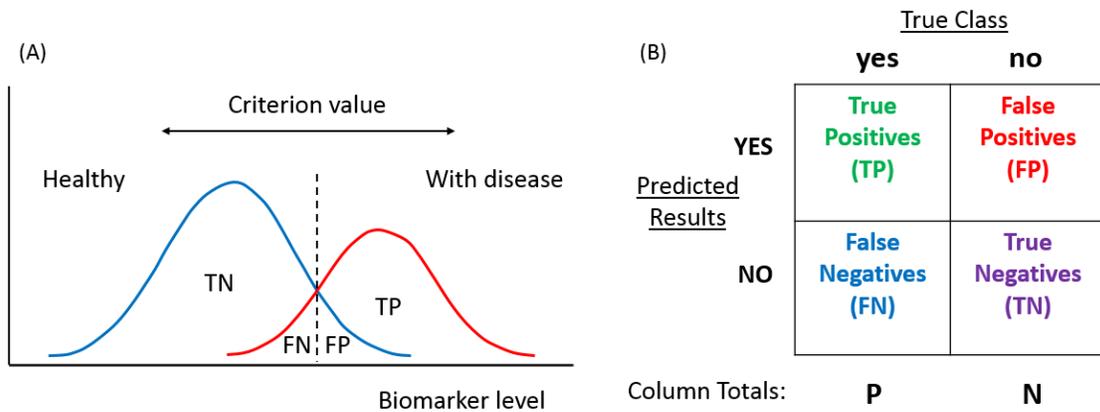
One-way ANOVA examines more than three independent groups of samples. The null hypothesis is that all the mean values are equal, and the alternative hypothesis is that at least one mean is statistically different. One-way ANOVA splits the total variance of the data into between-group variance and within-group variance, then calculates the ratio of these two. A large ratio means the alternative hypothesis is likely to be true. With the ratio, one-way ANOVA produces an F-test and also reports a p-value.

When there is a significant difference between the two populations (e.g., patients have significantly higher blood concentration) and this difference is validated by follow-up studies, the metabolite becomes a biomarker. In clinical diagnosis, when the blood

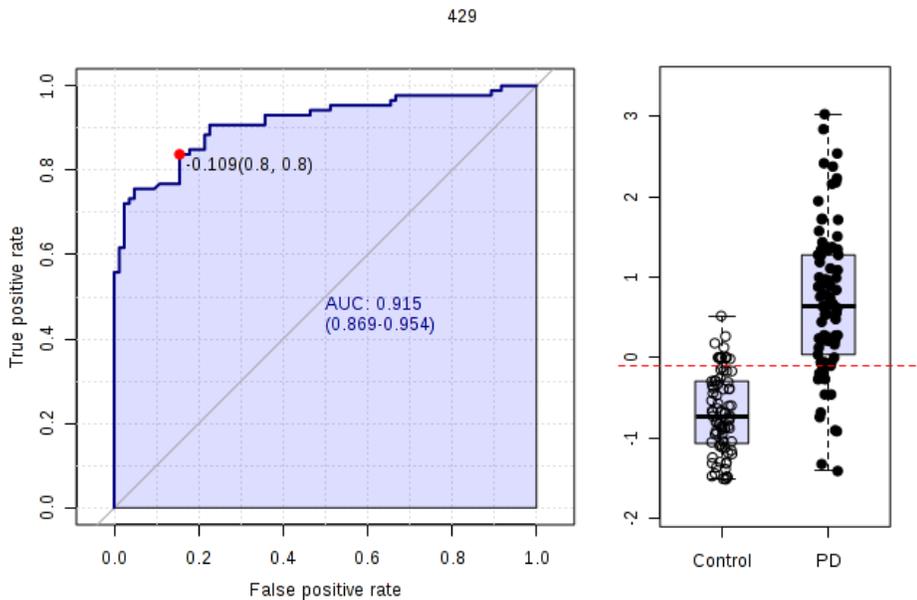
concentration of this biomarker is measured for a subject, we have to determine if the measured value belongs to the disease population. Logically, since we have the distributions of the general population and the disease population, we can do a one-sample t-test to prove the value is part of the disease population, and conduct another one-sample t-test to confirm that the value does not lie in the general distribution. A much easier and practical approach might be setting a criterion value, as shown in Figure 1.4. If the measured value exceeds the criterion, we can conclude that the subject has the disease. In the real-world, we rarely have a perfect separation between the two populations. Because of the overlap, we have to carefully adjust the criterion value to reach a balance between false positives and false negatives. In this case (disease > control), increasing the cut-off will decrease the number of false positives but increase the chance of false negatives.

A receiver operating characteristic curve (ROC curve) (Figure 1.5) is a graphical plot for evaluating the differentiating power of a binary classifier.<sup>100</sup> It examines a series of cut-off points, from low to high. For each criterion value, the true positive rate is calculated and named as sensitivity, and false positive rate is reported as (1-specificity). Then the ROC analysis draws a curve of sensitivity against (1-specificity). The classification power is evaluated by the Area-Under-the-Curve (AUC). AUC is a value between 0.5 and 1.0. 0.5 means no classification power at all, and 1.0 represents an excellent classifier. The AUC of a classifier is equivalent to the probability that the classifier ranks a randomly chosen positive instance above a randomly chosen negative one.<sup>101</sup> The uni-variate AUC values can be used for confirming and comparing the statistical significance of the significant metabolites output by the volcano plot. Several metabolites with medium AUCs can also

be combined to become a biomarker panel with higher AUC. The optimal cut-off point is the value when the Youden Index (sensitivity + specificity – 1) reaches the maximum. Researchers can report the sensitivity and specificity at the optimal point to demonstrate the performance of a biomarker.



**Figure 1.4** (A) Schematic of the selection of criterion value and (B) the confusion matrix.



**Figure 1.5** An ROC curve with AUC = 0.915, showing the optimal cut-off value (-0.109) with sensitivity of 83.7% and specificity of 84.5%.

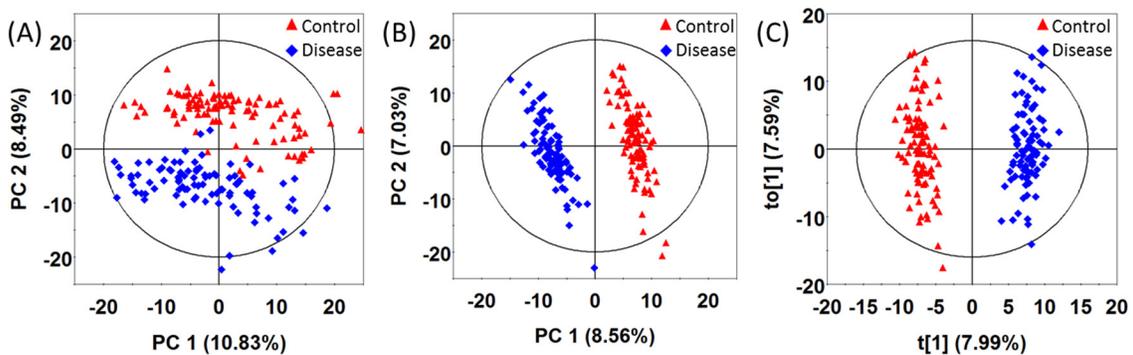
## 1.6.2 Multi-variate Analysis

The limitation of uni-variate analyses is that simply counting the number of significant metabolites cannot tell us how different two groups are. Multi-variate analyses, which treat the data matrix  $\mathbf{X}$  as a whole, can show us the inter-group differences in a broader perspective. Principal Component Analysis (PCA), Partial Least Squares-Discriminant Analysis (PLS-DA) and Orthogonal Partial Least Squares-Discriminant Analysis (OPLS-DA) are the most widely used multi-variate tools in metabolomics. All of them are projection methods, which means they project the high-dimensional data onto a 2D surface or 3D space, enabling visual presentation of the distribution of the data points. The axes of the 2D surface or 3D space are the principal components (PC).

PCA is an unbiased, high-confidence dimensionality reduction method. It finds the first principal component (PC1) by the linear combination of a set of variables. PC1 is supposed to account for as much of the variability in the data as possible. PC2, which is another linear combination of the variables, is orthogonal to PC1 and has the second largest variance. More PCs can be determined if needed. PCA decomposes data matrix  $\mathbf{X}$  into two orthogonal matrices ( $\mathbf{X} = \mathbf{V}^T \cdot \mathbf{U}$ ).<sup>102</sup> Matrix  $\mathbf{U}$  is called the scores matrix, which is a summary of the observations, and the loadings matrix  $\mathbf{V}$  is a summary of the variables. The 2D score plot projects all the observations onto the surface of PC1 and PC2. The distance between two data points is the variance, and when there is a statistically significant difference between two study groups, the inter-group variances should be more significant than the within-group variances. In other words, the observations in each group should cluster closely and there should be an obvious distance between the two clusters. PCA

provides a simple and graphical overview of the data without putting any extra assumptions into it. However, sometimes there is a true separation between the two groups, but the inter-group variation only accounts for a small portion of the total variability. In this case, PCA is not the best choice for studying the difference that we are interested in.

PLS-DA is a supervised method that considers the group assignment of the observations. It annotates the grouping information in numbers (e.g., 0 and 1) and then builds a linear regression model between the PCs and the group assignment. The chosen PCs should not only cover a relatively large portion of the total variance, but also satisfy a linear relationship to the observations. Although the variation coverage of PLS-DA PCs is usually lower than that of PCA PCs, PLS-DA provides a more focused view on the useful variations. OPLS-DA is an improved version of the PLS-DA with the stronger power to cope with unwanted variations. Fundamentally it has the same algorithm of PLS-DA and gives better-looking inter-group separation, which makes it easier to interpret the results.<sup>103</sup> OPLS-DA is more susceptible to over-fitting issues, and it can only be used after a significant separation is confirmed by PCA or PLS-DA.



**Figure 1.6** Score plots of the same data set given by (A) PCA, (B) PLS-DA and (C) OPLS-DA.

As mentioned before, multiple variables can be used to generate a multi-variate ROC curve. In this case, instead of a simple criterion value, we need a multi-variate classifier to determine the class of an observation. The development of the classifier is based on a part of the input data (as the training set), and then the other part of the data (as the testing set) is substituted into the classifier to generate the ROC curve. PLS-DA model can also be used as the classifier. However, as a regression model, PLS-DA is sensitive to missing values and the over-fitting problem if validations are not properly performed.

Several machine-learning methods can be applied for building the classifier. For example, Linear Support Vector Machine (SVM) maps the data into high-dimensional space that allows for the separation of two groups of samples into distinctive regions. It searches in the input high-dimensional space for an optimal plane that enables the maximization of the difference between the two groups. Subsequently, a new observation is classified based on which side it falls.<sup>104</sup> SVM outperforms PLS-DA because it is not influenced by the missing values and it can deal with non-linear models. However, it has no visualized scores and loadings, and the computational workload is burdensome.

The random forest method belongs to the family of classification trees.<sup>105</sup> It can easily handle missing values, outliers and relatively small sample sizes. A regular decision tree is the basic unit of the random forest. To generate an ROC curve, a multitude of decision trees are constructed using the training datasets. Then the testing datasets are classified

with these decision trees to generate the sensitivity and specificity values for plotting the ROC curve.

Generally speaking, a larger number of significant metabolites in a biomarker panel can always make better classification power. However, for clinical diagnostic purposes, we are trying to balance (1) achieving high sensitivity and specificity with (2) minimizing the required number of metabolites in the diagnosis panel. For real-world applications, a smaller set of biomarkers will be easier to quantify. We can choose the biomarkers according to the uni-variate results, or we may want to use the metabolites that have been positively identified or play roles in an essential metabolic pathway.

## **1.7 Challenges and Solutions**

In the past decades, metabolomics has discovered a large number of biomarker candidates, and the number is rapidly increasing. However, none of them has currently made the transition to routine use in clinical practice.<sup>106</sup> One possible reason is that, similar to protein biomarkers, the development and clinical implementation of a biomarker assay usually take several years. A number of companies are developing metabolite-based tests for the clinical diagnosis. For instance, in 2016, a commercialized assay based on three urine metabolites was put on the market for the diagnosis of adenomatous polyps.<sup>107</sup> In the future, more cooperative interactions between the researchers, the diagnostic industry as well as the clinicians will definitely speed up the application of metabolite biomarkers. Nonetheless,

metabolomics-based biomarker discovery is facing many challenges. Addressing these issues is a vital prerequisite for the wide application of metabolite biomarkers.

### **1.7.1 Statistical Over-fitting**

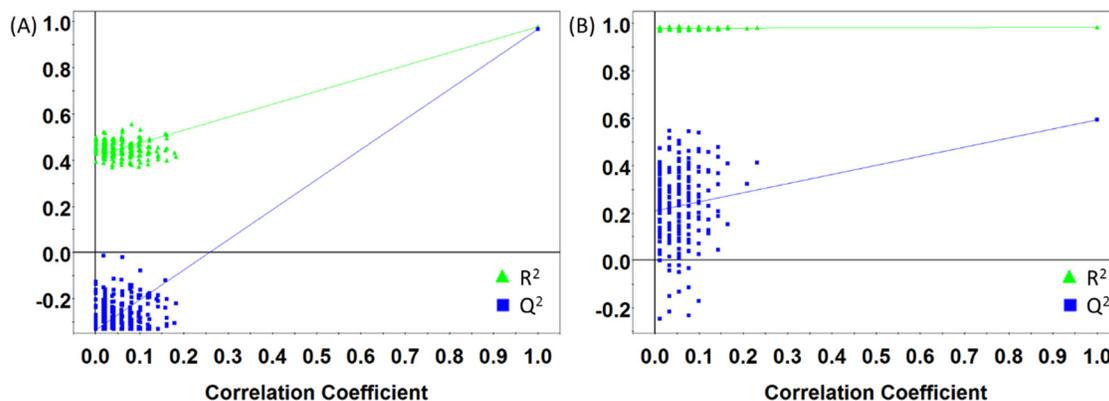
Most of the discovery and validation reports on biomarkers are based on assessing the statistical power. However, the metabolomics community has noticed that statistical significance does not always lead the way to biological meaningfulness.<sup>108</sup> Many factors may cause the false positive statistical outcome, such as the biological variability which will be discussed later. When the sample size is too small, or the sample collection is biased, the statistical analysis may not be able to accurately reflect the biological difference. As the metabolome coverage of the blood analysis is increasing, over-fitting is becoming a major problem, and fortunately, this issue can be largely overcome by careful statistical manipulation.

Mainly due to cost concerns, the number of observations in a biomarker discovery study is usually below 200. Therefore, we are having a much larger number of variables than observations. This may cause the over-fitting and significantly increase the chance of false positive results. In multi-variate analysis, when the sample sizes are relatively small, we can always see a nice separation on the score plot. Given a large set of variables, the algorithm can always find a combination that can coincidentally explain the pre-defined group separation. The over-fitting is monitored by internal validation methods, such as the cross-validation and the permutation test.

In PLS-DA, the cross-validation randomly picks 1/7th of the data as the testing set and leaves the other 6/7th as the training set. The training set is used for building a PLS model. The  $R^2$  value ( $1 - \text{residual sum of squares} / \text{sum of squares}$ ) can represent the quality of fitting. After that, the training set is substituted into the model, and the model predicts its observation values. By comparing the original values and the predicted values of the testing set, the algorithm gives the  $Q^2$  value ( $1 - \text{predicted residual sum of squares} / \text{sum of squares}$ ), which estimates the goodness of prediction. When over-fitting happens, the model might mistakenly give a high  $R^2$ , but the wrongly selected PCs cannot have high prediction power to reproduce the testing set, so a low  $Q^2$  value can indicate the over-fitting.

Most researchers accept that  $Q^2 > 0.5$  means an acceptable model and  $Q^2 > 0.9$  indicates an excellent model. However, the  $Q^2$  has no statistical significance to determine the proper cut-off. Therefore, we employ the permutation test to simulate the null distribution of the  $Q^2$  value. The group assignment of the observations is randomly permuted. Then for each permuted data set, we build a PLS-DA model and generate a pair of  $R^2$  and  $Q^2$ . The distribution of these new  $Q^2$  values is the null distribution when the null hypothesis is true. If the original  $Q^2$  is much larger than the permuted values, we can conclude that the  $Q^2$  has statistical significance and the PLS-DA model is valid. If the permuted data can generate very high  $Q^2$  values, we can know that the over-fitting is very severe, raising doubts about the separation in the original model. A response permutation plot, as shown in Figure 1.7, plots the original and permuted  $Q^2$  values against the correlation coefficient, and draws a straight line connecting the original value and the mean of the null distribution. A negative

intercept on the y-axis accepts the original PLS-DA model and a positive intercept indicates that the original PLS-DA model is over-fitted.



**Figure 1.7** (A) Permutation test accepts a valid model. (B) Permutation test rejects an over-fitted model.

Uni-variate analyses also need corrections due to the large number of variables. The multiple testing problem refers to the situation that when multiple comparisons are being interpreted simultaneously, the chance of false positives correspondingly increases. For example, if we have 100 variables, we will conduct 100 t-tests. With the significance level set at 0.05, even if the null hypothesis is true, the chance of having at least one significant result increases to  $1 - (1 - 0.05)^{100} = 99.4\%$ . Therefore, we will find at least one metabolite with statistical significance, but actually the observation is caused by random chance. With more variables, the multiple testing problem can become very significant. There are a number of tools for correcting the multiple testing problem. Bonferroni correction considers the most extreme case of over-fitting, and corrects the analysis by dividing the significance level by the number of variables. In the example above, the significance level should be set to  $0.05/100 = 0.0005$ .<sup>109</sup> Bonferroni correction is a bit too conservative, and

as discussed before, we do not want the primary screening of biomarker candidates to be too strict. More researchers prefer monitoring the false discovery rate (FDR) rather than the false positive rate. Storey et al. have developed a method to calculate the FDR-adjusted-p-value (q-value) based on the distribution of the p-value.<sup>110</sup> Setting the q-value threshold at 0.05, we are expecting that 5% of the selected significant variables are false positives, which will be excluded in the following analyses or the validation steps.

Overall, statistical over-fitting is mainly caused by the small or imbalanced sample sizes. When the variable-to-observation ratio is too big, it will be very hard to dig out the useful information from the countless false positive results. The ultimate solution is to increase the sample size. With the number of observations matched to or even larger than the number of variables, the reliability of the statistical analysis will be greatly improved.

### **1.7.2 Biological Variability**

Because of the inductive nature of metabolomics study and the lack of knowledge on the biological functions of the metabolites, some biological confounding factors may interfere with the statistical analysis. Humans are very diverse beings, and various genetic, environmental and lifestyle factors can cause variations in the blood metabolome. Blood metabolome variations due to sex, age, race, body weight and smoking have been reported.<sup>34, 45</sup> These variations may become confounding factors in biomarker discovery. For example, the onset rate of lymphedema in women is much higher than in men.<sup>111</sup> Therefore, in a metabolomics study, it is very likely that the majority of the disease group

are female. When the control group has more males, we will probably see some sex-dependent metabolites mistakenly recognized as the significant metabolites for the diagnosis of lymphedema. Logically, the disease group and the control group should have all these parameters matched. However, it is not very practical in clinical studies and there may be other interfering factors that we are not aware of. To overcome this issue, we need to perform metabolome profiling of large and diverse populations and develop metabolome databases for assessing all the possible confounding factors. For instance, with the knowledge of the sex-dependent metabolites, we are able to exclude these false positive findings from the lymphedema study.

In most metabolomics studies, the blood sample is only collected once from each subject. However, researchers have found significant time-of-day variations in the blood metabolome.<sup>43, 112</sup> Consequently, a single observation may not be able to represent the average level of the metabolite and may lead to false interpretations. Multiple blood collections during a period is a must for overcoming the within-individual variations and enabling a more comprehensive understanding of the individual blood metabolome. In addition, diet effects can significantly affect the blood metabolome as well, but how long the diet effect can last and whether the participants should fast overnight before giving blood remain unclear. Consecutive blood collections can also help with assessing the metabolic effects of dietary stimulation and other environmental exposures. However, performing venipuncture for many times within a day can be very invasive, and the cost will be too high. Less-invasive and cheaper blood analysis platforms are needed to achieve this goal.

Finally, due to the diluting effect of water and other reasons, the total metabolite concentration of blood may vary. This effect is very significant in urine metabolome studies, and urinary creatinine is widely used as the standard to normalize the samples.<sup>113</sup> Despite not being as significant as in urine, variations of total metabolite concentration in blood (about 20%) have also been reported.<sup>114</sup> According to our knowledge, most blood metabolomics studies only performed post-acquisition normalization (e.g., normalization based on the total useful signal). Pre-acquisition normalization, which is more accurate, has not been applied to the blood samples. And there is no assessment of the performance of post-acquisition normalization. Our group has previously developed an LC-UV based pre-acquisition normalization method, which measures the UV absorption of the dansyl-labeled amine/phenol-containing metabolites.<sup>115</sup> The pre-acquisition normalization can be performed for blood samples with our CIL LC-MS platform.

### **1.7.3 Study Design**

According to the above discussions, a large sample size is highly desirable in metabolomics studies, and an ideal study design should involve well-balanced, diverse and large study groups. In uni-variate analysis, the ideal sample size can be determined by the power analysis.<sup>116</sup> Unlike the t-test, the power analysis monitors the type II error. The power, defined as  $(1 - \text{false negative rate})$ , is the probability of making the correct decision if the alternative hypothesis is true. Given the desired effect size, significance level and power, power analysis determines the minimum sample size. Ferreira et al. have extended the

application of power analysis to multi-variate data.<sup>117</sup> Nonetheless, this method is designed to guide the study design according to previously obtained pilot data. In the real world, it is not very realistic to perform a pilot study to determine the sample size, and the number of samples is mainly determined by ethical and economical restrictions.<sup>97</sup> In a biomarker discovery study, we cannot pre-define the desired significance level and effect size, either. What we can do is to recruit as many samples as possible. Collecting the finger blood or heel blood is a cheaper and less-invasive alternative, however, with the small sample amount, the metabolome coverage and detectability of low-abundance metabolites should be carefully assessed. Since both increasing the sample size and performing the time-resolved analysis require more convenient blood sample collection methods, developing such a platform plays a central role in improving the reliability of blood metabolomics analysis.

Another important part of the study design is the validations. Most reported biomarker candidates do not have follow-up validation studies. Before a biomarker candidate goes into clinical use, multiple external validations should be conducted in independent, diverse, and large populations. This is not an easy task that can be achieved by one laboratory. Nonetheless, since many researchers are studying the same disease, sharing knowledge and collaboration among laboratories can largely promote the validation process. It is important to note that the external validation should be independent of the original study design. It is logically incorrect to add more observations into the study after the original analysis is done in order to improve the statistical performance, which is called “p-hacking”.<sup>118</sup> This is because when the null hypothesis is true, the relationship between the p-value and the

sample size is totally random. Therefore, even though the alternative hypothesis is false, adding more samples will always have a chance to improve the statistical performance.

Additionally, most biomarker discovery studies are case-control studies in epidemiology, in which there is a disease group and a control group. Adding the time dimension to the study design is also very beneficial. Despite the higher cost, cohort studies are more powerful than case-control studies. Cohort studies analyze the blood metabolome before the onset or progression of the disease, and examine the outcome after a period. Since the metabolic changes are measured before the outcome reveals, we have a temporal framework to assess the causality and therefore have more confidence in the biomarkers.<sup>119</sup> Similarly, intervention studies introduce an external stimulation to one of the two study groups and monitor the outcome. It is a useful tool for confirming the relationship between an environmental exposure and a disease, or assessing the performance of treatment.

## **1.8 Overview of Thesis**

My research started from developing and optimizing chemical isotope labeling (dansyl-labeling and DMPA-labeling) LC-MS platform for the analysis of blood metabolome. The first four chapters of my thesis focus on understanding and reducing the experimental and biological variations that interfere with biomarker discovery studies. The last two chapters are the applications of the CIL LC-MS platform for blood biomarker discovery and exposure assessment.

Extra matrices are commonly introduced to plasma or serum samples during the sample handling. Although the CIL LC-MS method can overcome the matrix effects during MS detection, the role of matrix effects during the chemical labeling process itself remains unclear. Chapter 2 assesses the matrix effects during the reaction and proposes a solution to minimize the interferences. To avoid confounding factors in biomarker discoveries, it is important to understand the metabolome variations among the general population due to genetic and environmental factors. In Chapter 3, we develop a serum metabolome database of 1,348 amine/phenol-containing metabolites and 1,065 carboxyl-containing metabolites, and by employing a universal serum standard, the information can be easily used in future discovery studies and clinical applications. The metabolome variations due to sex, age, and body weight are also assessed. Furthermore, a less-invasive and high-coverage metabolome profiling method plays a crucial role in increasing the sample size, overcoming the time-of-day variability and thereby improving the reliability of biomarker studies. Chapter 4 describes a dansyl-labeling LC-MS-based method that can accurately quantify 1,722 metabolites in one microliter of finger blood, opening the possibility of time-resolved metabolomics analysis. In Chapter 5, DMPA-labeling is added to the analysis of one microliter of whole blood, detecting more than 4,000 metabolites from the extremely low amount of sample. Additionally, the method is applied to the exposomics assessment of a dietary exposure.

In Chapter 6, a cohort study is designed to monitor the progression of Parkinson's disease (PD). Serum samples were collected three years before a part of the PD patients developed dementia. By analyzing the samples with CIL LC-MS, we report a 5-metabolite panel for

the diagnosis of PD and an 8-metabolite panel for predicting the onset of dementia. At last, animal models usually outperform humans in exposure assessment as animals are less diverse in terms of diet and living environment. In Chapter 7, we choose a silkworm model to study the metabolic changes in blood (hemolymph) under the exposure of an endocrine disruptor (dichlorodiphenyltrichloroethane (DDT)).

Overall, the major objective of my thesis work is to add the time factor into blood metabolomics, including both biomarker discovery and clinical diagnosis. We have established a metabolome database listing the common biological variations among healthy population. In the future, with our finger blood analysis technique, we will generate more knowledge about the time-dependent metabolic changes associated with diet effect or other environmental exposures. In biomarker discovery, this information can help us determine if the significant metabolome changes are truly originated from the disease state. And in clinical diagnosis, considering all these variations and standardizing the sample collection procedures will make the result more reliable.

## Chapter 2

### **Matrix Effect on Chemical Isotope Labeling and Its Implication in Metabolomic Sample Preparation for Quantitative Metabolomics**

#### **2.1 Introduction**

Chemical isotope labeling liquid chromatography mass spectrometry (CIL LC-MS) is an enabling analytical platform for generating comprehensive and quantitative metabolomic profiles for metabolomics research.<sup>70</sup> Using a proper labeling reagent to react with a class of metabolites (e.g., all amine-containing metabolites), a chemical-group-based submetabolome can be analyzed with improved LC separation and enhanced MS sensitivity to generate a comprehensive profile of the submetabolome.<sup>70, 73</sup> By combining the results of different submetabolomes generated using labeling reagents targeting different chemical groups, a large coverage of the entire metabolome may be achieved.<sup>120</sup> A growing number of CIL reagents have been developed for targeted metabolite analysis or group-based submetabolome profiling.<sup>121-131</sup>

For quantitative metabolomics, CIL LC-MS is performed using differential isotope labeling of individual samples (e.g., labeled with  $^{12}\text{C}$ -reagent) and their control (e.g., labeled with  $^{13}\text{C}$ -reagent), which overcomes the problems of matrix effect and ion suppression associated with MS detection.<sup>70, 73</sup> However, matrix compositions of individual samples such as the salt and buffer contents may be different from sample to sample or batch to batch. While it is important to control the sample collection and sample

preparation steps properly for quantitative metabolomics, differences in sample matrix are unavoidable due to inherent variations of sample matrix such as salt contents in a biofluid (e.g., urine and sweat) and logistical consideration in real world applications. An example of the latter is that samples may be collected at different centers or time under somewhat different conditions such as using different additives (buffers, EDTA, etc). By necessity, one may want to profile these samples for improving the overall performance of a metabolomics study (e.g., using samples from different centers after the initial work of disease biomarker discovery using a well-controlled sample set for biomarker validation).<sup>132</sup> In another situation, valuable samples that have been subjected to NMR analysis may be re-used for MS-based profiling.<sup>51</sup> The NMR samples with the addition of phosphate buffer (PB) or phosphate buffer saline (PBS) that are used for controlling pH and ionic strength to minimize the chemical shift changes<sup>133-136</sup> would obviously have different matrices from those of other untreated samples.

Blood metabolomics also encounters the matrix effect. Plasma samples are collected with the addition of an anticoagulant, such as EDTA, citrate and heparin. Serum samples are sometimes diluted by PBS solution due to specific experimental requirements. The matrix effect of these additives can potentially become confounding factors in blood metabolomics, and we should carefully assess the corresponding interferences during the CIL reaction. In this chapter, we report the presence of a matrix effect on chemical labeling in a dansylation isotope labeling LC-MS metabolomic profiling workflow. Since compared to blood, urine usually has a stronger matrix and more significant inter-individual variation, we use urine for the assessment of matrix effects. We illustrate that metabolomic profiles

of urine samples with and without the presence of high concentrations of NaCl, PB or PBS produced by dansylation LC-MS can be different. While matrix effect in LC-MS analysis can be overcome by differential CIL, our work points out the importance of using similar sample matrices for chemical labeling to maintain similar labeling efficiencies of individual metabolites in comparative metabolomics. We also demonstrate that for samples with varying concentrations of salts such as urine samples, simply diluting the samples to reduce the salt concentration prior to chemical labeling can overcome the matrix effect.

## **2.2 Materials and methods**

### **2.2.1 Chemicals and reagents**

All the chemicals and reagents, unless otherwise stated, were purchased from Sigma-Aldrich Canada (Markham, ON, Canada). For dansylation labeling reaction, the  $^{12}\text{C}$ -labeling reagent (dansyl chloride) was from Sigma-Aldrich and the  $^{13}\text{C}$ -labeling reagent were synthesized in our lab using the procedure published previously.<sup>70</sup> LC-MS grade water, methanol, and acetonitrile (ACN) were purchased from ThermoFisher Scientific (Nepean, ON, Canada).

### **2.2.2 Urine sample collection**

Urine samples were collected from five age-matched healthy mice. Equal volumes of the five individual samples were mixed together to make a pooled sample. After sample collection, the urine samples were immediately stored in  $-80\text{ }^{\circ}\text{C}$  freezer for further use.

### **2.2.3 Dansylation labeling**

The frozen urine samples were thawed in an ice-bath and then centrifuged at 14,000 rpm for 15 min. 12.5  $\mu\text{L}$  supernatant was taken into an Eppendorf tube and totally dried using a Speed Vac. The sample was re-dissolved to 50  $\mu\text{L}$  with water or a specific matrix solution. Then 25  $\mu\text{L}$  of 250 mM sodium carbonate/sodium bicarbonate buffer and 25  $\mu\text{L}$  of ACN were added into the sample. The solution was vortexed, spun down, and mixed with 50  $\mu\text{L}$  of freshly prepared  $^{12}\text{C}$ -dansyl chloride solution (18 mg/mL) (for light labeling) or  $^{13}\text{C}$ -dansyl chloride solution (18 mg/mL) (for heavy labeling). After 45 min incubation at 40  $^{\circ}\text{C}$ , 10  $\mu\text{L}$  of 250 mM NaOH was added to the reaction mixture to quench the excess dansyl chloride. The solution was then incubated at 40  $^{\circ}\text{C}$  for another 10 min. Finally, formic acid (425 mM) in 50/50 ACN/ $\text{H}_2\text{O}$  was added to consume excess NaOH and to make the solution acidic. The  $^{12}\text{C}$ - or  $^{13}\text{C}$ -labeled sample was centrifuged at 14,000 rpm for 10 min before injecting onto LC-UV for quantification. For LC-MS analysis, the  $^{12}\text{C}$ - and  $^{13}\text{C}$ -labeled samples were mixed in equal amounts based on the quantification results.

#### **2.2.4 LC-UV quantification**

A Waters ACQUITY UPLC system with a photodiode array (PDA) detector was used for the quantification of dansyl labeled metabolites for sample amount normalization as described earlier.<sup>115</sup> Briefly, 2  $\mu\text{L}$  of the labeled urine or amino acid solution was injected onto a Phenomenex Kinetex C18 column (2.1 mm  $\times$  5 cm, 1.7  $\mu\text{m}$  particle size) for a fast step-gradient run. Solvent A was 0.1% (v/v) formic acid in 5% (v/v) acetonitrile/water, and solvent B was 0.1% (v/v) formic acid in acetonitrile. The gradient started with 0% B for 1 min and was increased to 95% within 0.01 min and held at 95% B for 1 min to ensure complete elution of all labeled metabolites. The flow rate used was 0.45 mL/min. The peak

area, which can represent the total metabolite concentration in the sample, was integrated using the Empower software (6.00.2154.003).

### **2.2.5 LC-FTICR-MS**

An Agilent 1100 series binary system (Agilent, Palo Alto, CA) and an Agilent reversed-phase Eclipse plus C18 column (2.1 mm×100 mm, 1.8 µm particle size, 95 Å pore size) were used for LC-MS. LC solvent A was 0.1% (v/v) formic acid in 5% (v/v) ACN/H<sub>2</sub>O, and solvent B was 0.1% (v/v) formic acid in acetonitrile. The gradient elution profile was as follows: t = 0 min, 20% B; t = 3.5 min, 35% B; t = 18.0 min, 65% B; t = 24.0 min, 99% B; t = 32.0 min, 99% B. The flow rate was 180 µL/min. The flow from HPLC was split 1:2 and a 60 µL/min flow was loaded to the electrospray ionization (ESI) source of a Bruker 9.4 Tesla Apex-Qe Fourier transform ion-cyclotron resonance (FTICR) mass spectrometer (Bruker, Billerica, MA, USA), while the rest of the flow was delivered to waste. All MS spectra were obtained in the positive ion mode. To monitor the instrumental performance, a quality control sample (i.e., a differentially labeled urine sample) was injected every 10 to 12 sample injections.

### **2.2.6 Data analysis**

The <sup>12</sup>C/<sup>13</sup>C peak pairs were extracted by the IsoMS software reported.<sup>72</sup> IsoMS-Align was used to align the peak pair data from different samples by retention time and accurate mass.<sup>72</sup> The missing values were filled back by using the Zero-fill program.<sup>137</sup> Multivariate statistical analysis was carried out using SIMCA-P+ 12 (Umetrics AB, Umea, Sweden).

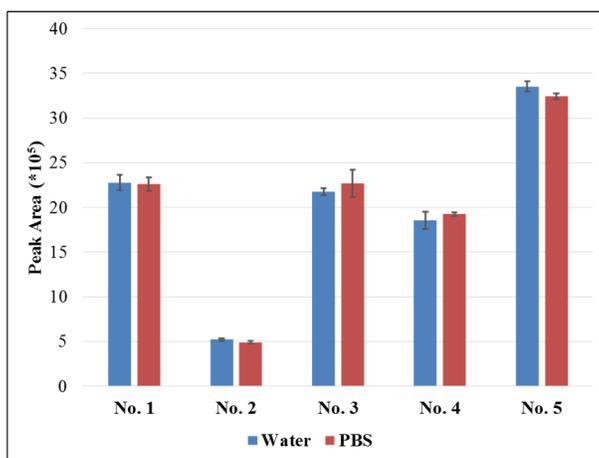
## 2.3 Results and discussion

### 2.3.1 Matrix effect on CIL

Differential CIL LC-MS provides relative quantification of metabolites in a sample vs. a control of similar type (e.g., an individual urine vs. a pooled urine) or absolute quantification of metabolites in a sample vs. a list of standards with known concentrations. When the sample matrices are not identical, matrix effect during the labeling process may occur, which could decrease the quantification accuracy. To examine the matrix effect on dansylation labeling which targets the amine/phenol submetabolome, we analyzed five individual mouse urine samples and a pooled sample in triplicates. For each sample, 12.5  $\mu\text{L}$  of sample were dried and then re-dissolved in 50  $\mu\text{L}$  of PBS or water. The individual samples were separately labeled by  $^{12}\text{C}$ -dansylation and the pooled sample was labeled by  $^{13}\text{C}$ -dansylation. Note that the concentration of dansyl chloride which was in excess was kept the same for labeling samples; we did not investigate how different concentrations of dansyl chloride affect the labeling efficiency of a sample containing different concentrations of salts or buffers. In CIL LC-MS, as long as the labeling efficiency is consistent for labeling the samples and the pool, relative quantification can be performed without much error. The total metabolite concentrations of the labeled samples were measured by LC-UV for sample amount normalization before mixing.

Figure 2.1 shows the LC-UV measurement results where the peak area was determined from a chromatographic peak of all the dansyl labeled metabolites eluted by using high organic solvent in a step-gradient LC chromatogram. The concentration of the 5 urine samples could differ by as much as 5-fold (sample #2 vs. sample #5). However, the total

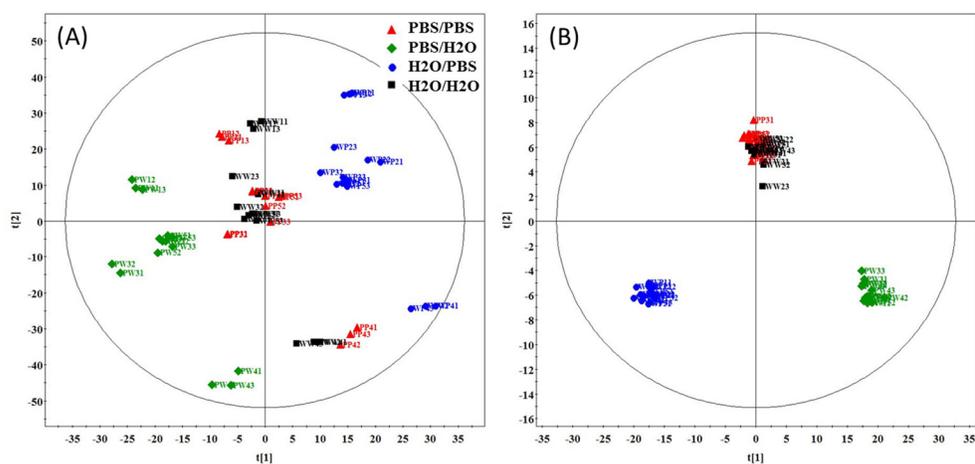
concentration of labeled metabolites is not affected by the presence or absence of PBS. Based on the LC-UV results alone, we could not detect any significant matrix effect on labeling. To normalize the sample amount for LC-MS analysis, equal amounts of a labeled sample and a labeled pooled sample were mixed. The same amount of the mixtures was injected into LC-MS for all samples. The intensity ratio of the  $^{12}\text{C}$ -/ $^{13}\text{C}$ -labeled peaks from a metabolite peak pair in mass spectra was measured and entered into a metabolite-intensity table for the samples and all the metabolite peak pairs. Since the same pooled sample was used as a reference for all the individual samples, the peak ratios found in the table for a given metabolite reflect the relative concentration differences of the metabolite in these samples. This is the basis of quantitative metabolomics using differential CIL LC-MS.



**Figure 2.1** LC-UV quantification of the total concentration of labeled metabolites in five mouse urine samples re-dissolved in water and PBS.

For the triplicate analyses of five urine samples, a total of 1662 peak pairs or putative metabolites were detected. We applied partial least squares discriminant analysis (PLS-DA) and orthogonal partial least squares discriminant analysis (OPLS-DA) to these data to

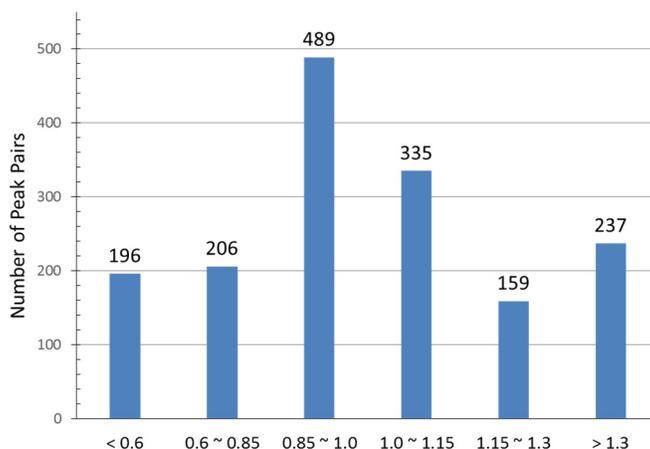
study the metabolomic changes. Figure 2.2 shows the score plots. As Figure 2.2 shows, there is a significant difference among the mixtures of urine re-dissolved in PBS labeled by  $^{12}\text{C}$ -dansylation and the pooled urine re-dissolved in  $\text{H}_2\text{O}$  labeled by  $^{13}\text{C}$ -dansylation (denoted as the PBS/ $\text{H}_2\text{O}$  group),  $^{12}\text{C}$ -urine in  $\text{H}_2\text{O}$  and  $^{13}\text{C}$ -pooled-urine in PBS (i.e.,  $\text{H}_2\text{O}/\text{PBS}$  group), and  $^{12}\text{C}$ -urine in  $\text{H}_2\text{O}$  and  $^{13}\text{C}$ -pooled-urine in  $\text{H}_2\text{O}$  (i.e.,  $\text{H}_2\text{O}/\text{H}_2\text{O}$  group). The PBS/ $\text{H}_2\text{O}$  group and the  $\text{H}_2\text{O}/\text{PBS}$  group are clearly separated from the  $\text{H}_2\text{O}/\text{H}_2\text{O}$  group.



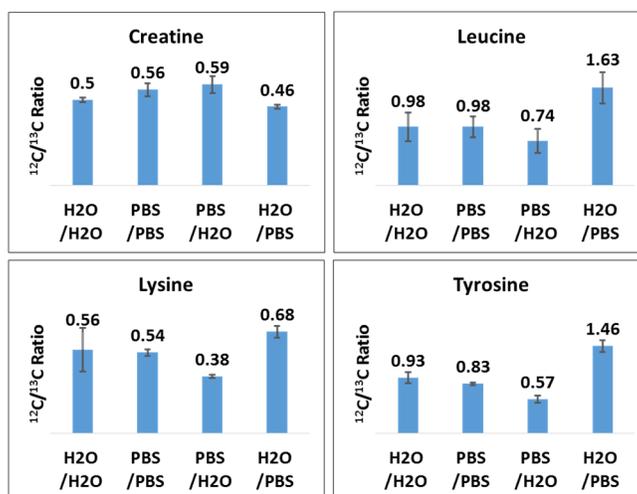
**Figure 2.2** (A) PLS-DA and (B) OPLS-DA score plots of dansylation LC-MS data obtained from five  $^{12}\text{C}$ -dansylated mouse urine samples mixed with a  $^{13}\text{C}$ -dansylated pooled sample. For each sample, three experimental replicates were performed. Label X/Y (X, Y=water or PBS) denotes a mixture of an individual urine re-dissolved in X and labeled with  $^{12}\text{C}$ -dansylation and a pooled urine re-dissolved in Y and labeled with  $^{13}\text{C}$ -dansylation.

At the individual metabolite ratio level, taking urine #1 as an example, out of 1622 peak pairs detected, 402 pairs in PBS have peak ratios decreased by more than 15%, compared to those in  $\text{H}_2\text{O}$ , while 396 pairs have ratios increased by more than 15% (Figure 2.3). Figure 2.4 shows four representative metabolites as an example to illustrate the matrix effect on labeling efficiency. Table 2.1 shows the p-values from t-test of the peak ratios

found in a given matrix group vs. the H<sub>2</sub>O/H<sub>2</sub>O group. The labeling efficiency of creatine was slightly increased in PBS, resulting in higher PBS/H<sub>2</sub>O ratio and lower H<sub>2</sub>O/PBS ratio. In contrast, labeling of leucine, lysine and tyrosine was suppressed by the PBS matrix. Because labeling efficiency of individual metabolites may increase or decrease, the total concentration of labeled metabolites was not affected by the presence of PBS as shown in Figure 2.1.



**Figure 2.3** Distribution of the PBS/H<sub>2</sub>O ratios (i.e., peak pair ratio in PBS vs. peak pair ratio in water) for metabolites in urine #1.



**Figure 2.4** Relative intensities of four representative metabolites in urine #1 vs. a pooled urine determined from experimental triplicate analysis of four different mixtures as shown in Figure 2.2.

**Table 2.1** Results of t-tests (p-values) showing the significance of ratio difference obtained from a given matrix group vs. the H<sub>2</sub>O/H<sub>2</sub>O group (the ratios for each metabolite are shown in Figure 2.4).

<b>Creatine</b>	PBS /PBS	PBS /H <sub>2</sub> O	H <sub>2</sub> O /PBS
H <sub>2</sub> O /H <sub>2</sub> O	0.051	0.035	0.012

<b>Leucine</b>	PBS /PBS	PBS /H <sub>2</sub> O	H <sub>2</sub> O /PBS
H <sub>2</sub> O /H <sub>2</sub> O	0.98	0.26	0.033

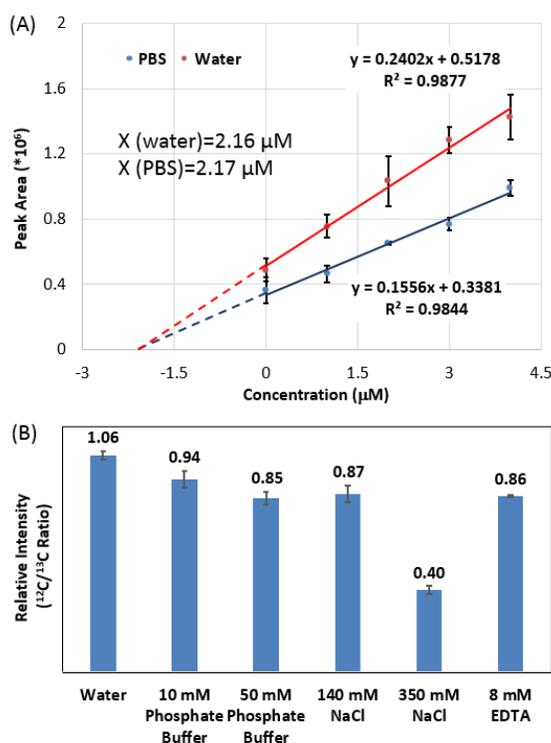
<b>Lysine</b>	PBS /PBS	PBS /H <sub>2</sub> O	H <sub>2</sub> O /PBS
H <sub>2</sub> O /H <sub>2</sub> O	0.76	0.094	0.25

<b>Tyrosine</b>	PBS /PBS	PBS /H <sub>2</sub> O	H <sub>2</sub> O /PBS
H <sub>2</sub> O /H <sub>2</sub> O	0.15	0.0050	0.00026

The results of PLS-DA and OPLS-DA plots and individual peak ratios clearly demonstrate that there was a matrix effect during the dansylation labeling process caused by the presence of PBS in a sample. Figure 2.2 also shows that the H<sub>2</sub>O/H<sub>2</sub>O group and the PBS/PBS group are overlapped on the score plots. The peak ratios of the four metabolites shown in Figure 2.4 are similar for the two groups. These results indicate that, while there was a matrix effect, the relative quantification results of metabolites was not affected if the individual samples and the pooled sample had the same or similar matrix.

To confirm the presence of matrix effect and investigate how it can influence absolute metabolite quantification, we chose tyrosine, which had a 39% peak ratio decrease in PBS, to generate standard addition curves. In this case, 12.5  $\mu$ L of urine #1 was dried and re-

dissolved to 50  $\mu\text{L}$  using either PBS or water spiked with 1  $\mu\text{M}$ , 2  $\mu\text{M}$ , 3  $\mu\text{M}$  and 4  $\mu\text{M}$  tyrosine standard solutions. The samples were separately labeled, followed by LC-MS analysis. The area of dansyl-tyrosine peak ( $m/z$  324.5953) was measured and plotted as a function of the concentration of spiked tyrosine. Figure 2.5A shows two standard addition curves for urine in PBS and water, respectively. There was a matrix effect on the labeling of tyrosine, causing a 35% decrease in labeling efficiency. However, the same absolute concentration (2.16-2.17  $\mu\text{M}$ ) was found from the two curves. It is clear that there was a matrix effect on metabolite quantification which could be overcome by the standard addition method.



**Figure 2.5** (A) Standard addition curves for tyrosine in mouse urine labeled in water and PBS. (B) Comparison of relative intensities of tyrosine labeled in different matrices.

**Table 2.2** Results of t-tests (p-values) showing the significance of ratio difference obtained from a given matrix vs. H<sub>2</sub>O (the ratio values for matrix solutions are shown in Figure 2.5B).

<b>Tyrosine</b>	10 mM Phosphate Buffer	50 mM Phosphate Buffer	140 mM NaCl	350 mM NaCl	8 mM EDTA
H <sub>2</sub> O	0.010	0.00053	0.0017	2.4*10 <sup>-6</sup>	8.8*10 <sup>-5</sup>

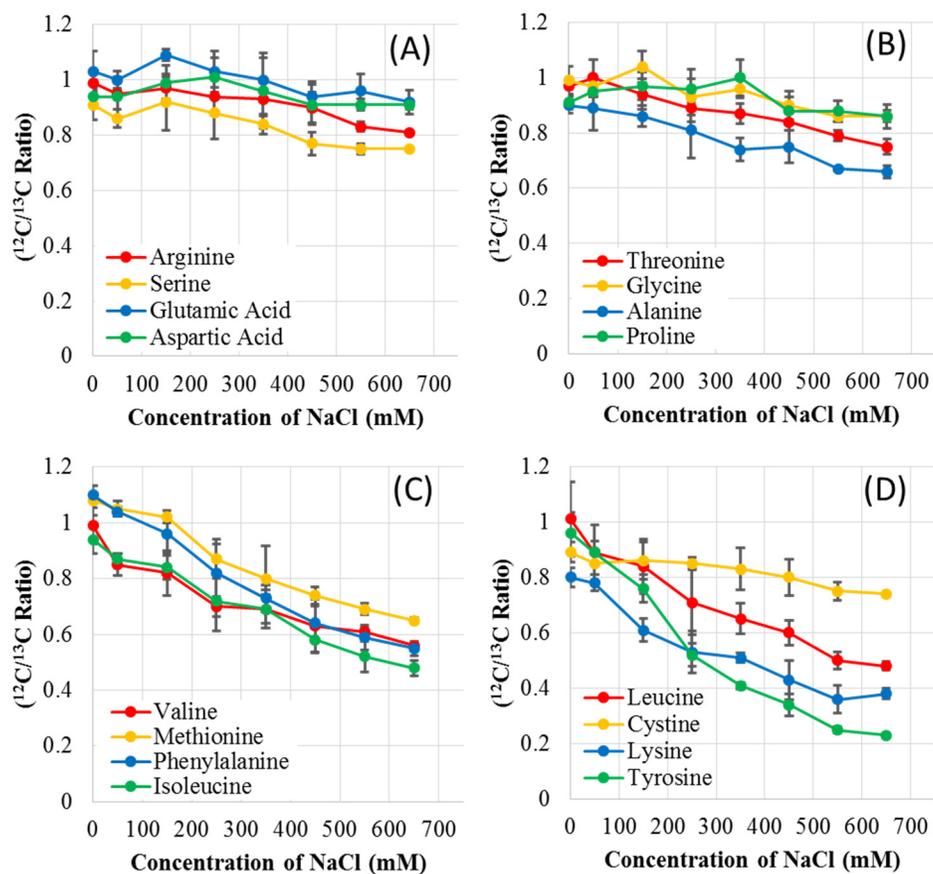
### 2.3.2 Origin of the matrix effect

To find a way to minimize the matrix effect on CIL, it is important to understand the origin and possible mechanism of this phenomenon. PBS solution contains 10 mM phosphate buffer (PB), 137 mM NaCl and 2.7 mM KCl. At first, we wanted to study PB and NaCl individually to see which one of them mostly contributes to the matrix effect in PBS. Because of the interest of performing MS analysis of the same sample already subjected to NMR measurement in our future research, we also examined the effect of a higher concentration PB, i.e., 50 mM. It should be noted that, using the tube-in-tube method,<sup>138-140</sup> no isotopic solvent needs to be added to the NMR sample and thus the sample, albeit containing high salt or buffer, can be readily transferred for MS analysis without any interference caused by the isotope solvent. Finally, 8 mM EDTA, which is approximately the final concentration of anticoagulant used in plasma, was also added to the list to examine its effect on labeling efficiency.

In our experiments, we dissolved tyrosine separately in water, 10 mM and 50 mM PB, 140 and 350 mM NaCl, and 8 mM EDTA, followed by labeling with <sup>12</sup>C-dansyl chloride and then 1:1 mixing with tyrosine in water labeled by <sup>13</sup>C-dansyl chloride. Note that dansyl chloride is highly soluble in acetonitrile or in a mixture of water and acetonitrile. As we

increased the salt or buffer concentration in the dansyl chloride solution, we did not observe any precipitation. Thus the presence of varying concentrations of salts or buffers in the sample did not change the solubility of dansyl chloride. The peak pair ratios of tyrosine were calculated after LC-MS analysis and the results are shown in Figure 2.5B. Table 2.2 shows the p-values from t-test of the peak ratios found in a given matrix vs. H<sub>2</sub>O. As Figure 2.5B shows, 140 mM NaCl has a larger effect on labeling efficiency, compared to 10 mM PB. The matrix effect of 50 mM PB is similar to that of 140 mM NaCl. More significantly, the tyrosine labeling efficiency in 350 mM NaCl solution is greatly reduced. In addition, 8 mM EDTA also causes matrix effect on tyrosine labeling to an extent similar to 50 mM PB or 140 mM NaCl.

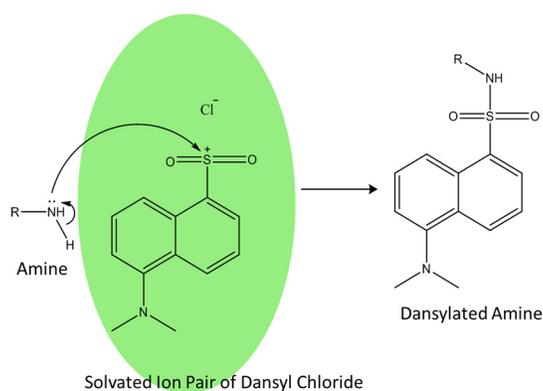
Although keeping the alkaline environment is crucial for the dansylation reaction,<sup>141</sup> no significant pH change of the reaction mixture was observed for these samples. Hence, the buffering property of PBS or PB was not the reason causing the matrix effect. Considering the results of the two NaCl samples (140 mM vs. 350 mM), we argue that increasing ionic strength in the labeling solution by high concentration of salts might be the main cause of the matrix effect. In the case of phosphate buffer, since HPO<sub>4</sub><sup>2-</sup> carries two charges, the ionic strength of PB should be higher than NaCl solution at the same concentration. However, the concentration of Na<sub>2</sub>HPO<sub>4</sub>/KH<sub>2</sub>PO<sub>4</sub> in PBS is much lower than NaCl. Therefore, the matrix effect of PBS could be considered mainly as the matrix effect of NaCl. The 8 mM EDTA solution only contains 16 mM Na<sup>+</sup>, but the matrix effect was also observed. This could be the effect of EDTA anion, which carries more charges than sodium cation.



**Figure 2.6** Relative intensities of 16 amino acids as a function of NaCl concentration in the sample solution.

To further examine how salt content or ionic strength can affect the labeling reaction, we dissolved a mixture of equal amounts of 16 amino acid standards in NaCl solution at a concentration ranging from 0 to 650 mM. When the salt concentration increased to 750 mM, acetonitrile and water could be separated into two layers by centrifuging. This is because the salt ions weaken the interaction between water molecules and acetonitrile molecules and change the solvent into an emulsion. This layer separation should be avoided for carrying out the labeling reaction. The amino acid mixture dissolved in different concentrations of NaCl was separately labeled by  $^{12}\text{C}$ -dansyl chloride and then 1:1 mixed

with a control standard labeled with  $^{13}\text{C}$ -dansyl chloride in water. The peak pair ratios were calculated and plotted against the concentration of NaCl and the results are shown in Figure 2.6. As Figure 2.6 shows, the 16 amino acids have different responses to the increasing salt concentration. The 8 relatively hydrophilic dansyl-amino acids (Figure 2.6A, B) have less effect by increasing NaCl concentration. However, the peak pair ratios of the relatively hydrophobic dansyl-amino acids except dansyl-cystine decrease significantly as the salt concentration increases (Figure 2.6C, D). For example, the peak pair ratio of dansyl-tyrosine in 650 mM NaCl is only about 24% of the ratio determined in water.



**Figure 2.7** Schematic of dansyl-labeling for amine in the presence of salts.

The fact that relatively hydrophobic amino acids are more sensitive to the salt matrix supports our hypothesis on the role of ionic strength on matrix effect. Carta and Tola reported the solubility of glycine, leucine, cystine and tyrosine in aqueous solution at different NaCl concentrations.<sup>142</sup> As the NaCl concentration increases, the solubility of cystine increases, while the solubility of leucine and tyrosine decreases. Glycine solubility is not affected by the salt. Our results of the matrix effect on the hydrophobic amino acids follow a similar trend. Dansylation reaction between dansyl chloride and an amine-

containing metabolite is initiated by nucleophilic attack of the amine to the dansyl sulfide. An intermediate or ion pair of dansyl and chloride is formed, followed by substitution of chloride by the amine (Figure 2.7). Ionic strength can influence the nucleophilic attack. Any ionic species from a matrix that surround the dansyl moiety may reduce or enhance the propensity of the amine to interact with dansyl to form a product, resulted in an decrease or increase in labeling efficiency. For amino acids with hydrophobic side chains, increasing ionic strength makes them less likely to interact with the dansyl moiety, which reduces the dansylation efficiency.

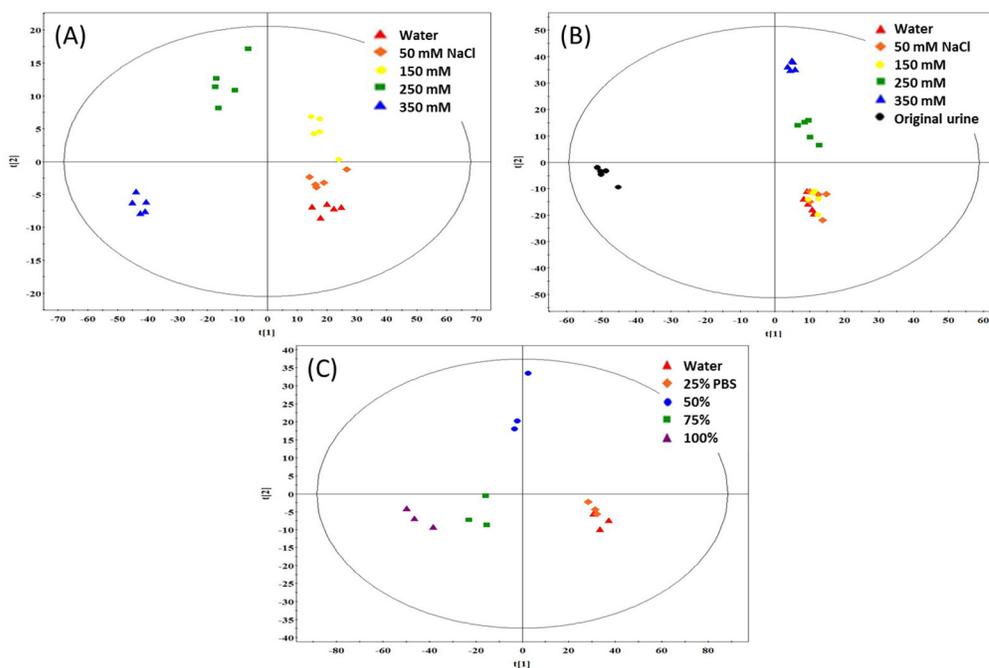
### **2.3.3 Minimizing the matrix effect on metabolomic profiling**

From the above results and discussion, it is apparent that high salts in various forms such as NaCl, PB or PBS can cause matrix effect on dansylation labeling. PB or PBS is introduced during sample preparation for various reasons. On the other hand, salts such as NaCl are inherently present at high concentrations in several biofluids such as urine and sweat. The salt concentration can vary from sample to sample. For example, in the study of urine sodium excretion trend between 1988 and 2010 among U.S. adults (age 20-59 y),<sup>143</sup> the urine sodium concentrations in spot urine specimens were reported to be 111 mM (mean) with 22.8 mM at 5th percentiles, 103 mM at 50th percentiles, and 221 mM at 95th percentiles for 1249 adults collected between 1988-1994. Slight increases in urine sodium concentrations were found for 1235 adults collected between 2003-2006 and 525 adults in 2010. The potassium concentrations in human urine are less than 1/3 of the sodium levels<sup>144</sup> and other metal ions concentrations are even lower.<sup>145</sup> For mouse urine samples, the

sodium concentration can range from ~40 mM to ~150 mM and the potassium concentrations can range from ~150 mM to ~500 mM, depending on diet and age.<sup>146</sup>

To examine how the salt contents can affect the metabolomic profiling, we selected one mouse urine and then diluted by 5-fold using different concentrations of NaCl solutions (50 mM, 150 mM, 250 mM and 350 mM). All these urine samples were labeled by <sup>12</sup>C-dansyl chloride and then mixed with a control sample diluted by water which was <sup>13</sup>C-dansylated. After LC-MS analysis, 1294 peak pairs or putative metabolites were detected. Figure 2.8A shows the PLS-DA score plot of the samples. There is a positive correlation between the salt concentration and the matrix effect. The higher the salt concentration, the further the data points are separated from the sample diluted by water. After adding the results from the un-diluted sample to the plot, Figure 2.8B shows the new plot. The samples diluted by water, 50 mM NaCl, and 150 mM NaCl cluster together. The samples diluted by 250 mM and 350 mM are clearly separated from the other diluted samples. The undiluted urine is significantly different from the 5-fold water diluted sample. These results indicate that the high salt concentration in the original urine had a strong matrix effect. After 5-fold dilution using water or low salt solution (50 mM), the matrix effect is greatly reduced or eliminated. Therefore, a dilute-then-label approach can be used to overcome the matrix effect on the CIL of biofluid samples. Figure 2.8A, B indicates that the presence of 50 mM NaCl in the urine samples does not cause significant matrix effect. This result is consistent with the data shown in Figure 2.6 where the peak ratios of most of the labeled metabolites do not differ significantly between the water and 50 mM NaCl samples.

The effects of varying concentrations of PBS were also examined. Similar trend to those of NaCl samples was observed for these samples (see Figure 2.8C). The samples labeled in 4-fold diluted PBS (2.5 mM phosphate buffer and 35 mM NaCl) are inseparable from the water-diluted samples as shown in Figure 2.8C. As the PBS concentration increases, the matrix effect becomes more significant. These results suggest that, for a NMR sample containing high concentration of PBS, we can simply dilute it to minimize the matrix effect on CIL LC-MS.



**Figure 2.8** (A) PLS-DA score plot for 5-fold diluted urine #1 labeled in water (red) and NaCl solutions (50 mM in orange, 150 mM in yellow, 250 mM in green, and 350 mM in blue). (B) PLS-DA score plot for the comparison of an undiluted urine sample labeled with dansylation (black) and 5-fold diluted urine samples labeled in water and NaCl solutions. The injection amount for all the samples based on LC-UV measurement was the same. (C) PLS-DA score plot for the comparison of 4-fold diluted urine #1 labeled at different concentrations of PBS solution (e.g., 25% PBS refers to 4-fold dilution of the PBS solution). For each sample, five experimental replicates were performed in (A) and (B), while three experimental replicates were performed in (C).

## 2.4 Conclusions

We demonstrated the presence of a matrix effect on chemical labeling that could affect the quantitative metabolomic profiling results in CIL LC-MS. Relative metabolite quantification in differential CIL LC-MS was not affected, if the same or similar matrix was present in comparative metabolomic samples. Because matrix effect is only present when the salt or buffer concentration is very high such as in some urine samples, the simplest way to minimize the influence of matrix effect on quantitative metabolomic profiling is to dilute all the samples by a specific factor (e.g., 4- or 5-fold for urine samples). Over-dilution is not recommended as it can increase the time required for concentrating the already labeled samples before LC-MS analysis. For different types of biofluids as well as different labeling chemistries, matrix effect on the chemical labeling process might be different and thus examining the matrix effect should be part of the protocol development for CIL LC-MS based metabolomics.

## Chapter 3

### **Development of a Human Serum Metabolome Database and Analysis of Metabolome Variations Using Isotope Labeling and High-resolution LC-MS**

#### **3.1 Introduction**

Human blood, as one of the most important biofluids, contains a treasury of known and unknown species that could reflect the ongoing physiologic state of all tissues,<sup>147</sup> and therefore it has long been used for clinical diagnosis, from the analysis of glutamic-pyruvic transaminase for assessment of hepatic diseases<sup>148</sup> to the application of blood protein biomarkers for tumor screening.<sup>149</sup> In addition to the large molecules, blood serum is a primary carrier of small-molecule metabolites, transporting all the small molecules that are being secreted or excreted by different tissues in response to various physiological conditions.<sup>31</sup> The regulation of these blood metabolites is not only governed by genetic effects,<sup>150</sup> but also determined by the interaction between human body and environmental factors.<sup>151</sup> Based on the hypothesis that when a specific disease state develops, certain physiological changes, as well as the resulting metabolic variations, may occur in the human body, blood metabolites can potentially become biomarkers for the early diagnosis of diseases.

Metabolomics, which is the high-throughput and systematic analysis of the small-molecule metabolites, has been emerging as a promising method for biomarker discovery.<sup>152</sup> Metabolomics analysis of serum has been used to discover metabolite biomarkers for the

diagnosis of various diseases, such as cobalamin deficiency,<sup>153</sup> colorectal cancer,<sup>154</sup> Alzheimer's disease<sup>155</sup> and Parkinson's disease.<sup>11</sup> To diagnose a disease with the quantitative analysis result of a biomarker from a subject, the concentration distribution among healthy people is a must for telling the analysis result is normal or not. However, humans are complex and diverse beings, and the diversity of genetic and environmental factors may cause huge variations in serum metabolome.<sup>62</sup> In order to accurately interpret metabolomics results, it is crucial to analyze the metabolomes of diverse populations that would allow us to determine if the presence, absence, over- or under-expression of specific metabolite(s) is a true representation of a disease state, and not due to other factors. Furthermore, biomarker discovery studies usually involve a comparison between the metabolic profiles of a control group and a disease group. Ideally, the populations to be compared are balanced for genetic and environmental factors (e.g., sex, age, race, and lifestyle), but this may not always be practical.<sup>34</sup> A possible situation can be that the average age of one study group is higher than that of the other. In this case, as metabolome variations associated with aging effects may cause confounding influences to the biomarker discovery, it is necessary to evaluate the metabolome variations due to aging within the healthy people. Logically, all the factors that cause metabolic heterogeneity among the general population should be carefully studied, and a non-targeted metabolomics profiling of a large, diverse and healthy population is an effective way to understand these metabolome variations and the underlying metabolic pathways.

There are only a few comprehensive assessments of blood metabolome variations in the literature. Lawton et al. conducted LC-MS and GC-MS analyses to study more than 300

metabolites in human plasma samples provided by 259 participants, and confirmed metabolome variations due to sex, age and race differences.<sup>34</sup> Using LC-MS and GC-MS, Saito et al. determined the levels of 297 metabolites among 60 healthy Caucasian individuals and the result revealed inter-sex and inter-age differences.<sup>46</sup> Dunn and his coworkers analyzed more than 1,500 metabolite features which were detected by GC-MS and LC-MS. By studying a large UK population, they reported that variations in serum metabolome could be related to differences in gender, age, body mass index (BMI), blood pressure, and smoking.<sup>45</sup>

Although in these existing works the variations of hundreds of serum metabolites have been studied, there are at least 4,300 known human serum metabolites according to the Human Serum Metabolome Database,<sup>31</sup> which means a large portion of the metabolites have not been evaluated, let alone the possibility that the number of unknown metabolites is even larger. The detected metabolites in a non-targeted metabolome profiling are generally high-abundance. Serum biomarkers, on the other hand, can be in very low concentrations,<sup>156</sup> which require improved techniques with larger metabolome coverage.

The low metabolome coverage has always been a challenge to metabolomics analyses. To overcome this issue, a differential isotope labeling approach can be applied with improved LC separation and enhanced ESI ionization efficiency. Previously, we have developed a dansylation labeling approach for studying the amine/phenol-containing metabolites<sup>70</sup> and a p-dimethylaminophenacyl (DMPA) labeling for studying the carboxyl-containing metabolites.<sup>73</sup> In this chapter, we employ these two methods to achieve a high-coverage

metabolome profiling of serum samples from 100 healthy individuals, which can serve as a serum metabolome database for understanding the metabolic variations in the general population. With more than 2,400 metabolites detected, we also study the effects of sex, age and BMI. Importantly, we pooled these 100 healthy serum samples into a universal standard, which can be used as the internal reference in our metabolomics analysis. As a result, the findings of our work can be easily adopted to all the future studies that apply the same universal standard, regardless of the LC-MS platform.

## **3.2 Materials and methods**

### **3.2.1 Chemicals and reagents**

All the chemicals and reagents, unless otherwise stated, were purchased from Sigma-Aldrich Canada (Markham, ON, Canada). For chemical isotope labeling reactions, the  $^{12}\text{C}$ -labeling reagents (dansyl chloride and DMPA bromide) were purchased from Sigma-Aldrich, and the  $^{13}\text{C}$ -labeling reagents were synthesized in our lab using the procedures published previously.<sup>70, 73</sup> LC-MS grade water, methanol, and acetonitrile (ACN) were purchased from Thermo Fisher Scientific (Nepean, ON, Canada).

### **3.2.2 Serum sample collection and the universal serum standard**

A hundred healthy volunteers including 35 males and 65 females were recruited in Edmonton, Canada. At the time of sample collection, the eligible participants were between the ages of 19 and 39 years old, and all of them were non-smokers. They did not consume

any prescription medications, counter pharmaceuticals or natural supplements, either. The study was conducted in accordance with the codes of the University of Alberta’s Arts, Science, and Law Research Ethics Board, and all participants provided informed consent. Table 3.1 describes the detailed demographic information including gender, age, and BMI.

**Table 3.1** Description of the sample set. Number in the cell represents the number of individuals that meet the corresponding condition. BMI groups were defined as: Underweight (BMI < 18.5), Normal (18.5 < BMI < 24.9), and Overweight (BMI > 25.0).

Age	19-25			26-39			All ages		
Sex	Female	Male	Total	Female	Male	Total	Female	Male	Total
Underweight	3	0	3	6	0	6	9	0	9
Normal	27	13	40	18	11	29	45	24	69
Overweight	5	4	9	6	7	13	11	11	22
All BMIs	35	17	52	30	18	48	65	35	100

The participants were refrained from eating or drinking (except water) for at least 8 hours before giving blood. 10 mL of venipuncture blood was collected into a BD Vacutainer 10 mL serum collection tube. The raw blood was allowed to clot spontaneously at room temperature for one hour, and then centrifuged at 1,500 g for 15 min to separate the blood cells. The supernatant (serum) was divided into multiple 250 µL aliquots in 1.5 mL microcentrifuge tubes for analysis or storage in a -80 °C freezer.

A universal serum standard (USS) sample was made by mixing equal-volume aliquots from each individual sample, and finally the USS sample was stored in 1.5 mL micro-centrifuge tubes at -80 °C. In the future metabolomics studies by differential isotope labeling, the USS sample will be <sup>13</sup>C-labeled and serve as the internal reference.

### **3.2.3 Dansylation labeling of serum samples**

The frozen serum sample was thawed in an ice-bath and then centrifuged at 15,000 g for 15 min. In a microcentrifuge tube, 30 µL of supernatant was mixed with 90 µL of methanol. The mixture was then incubated at -20 °C for 2 hours before centrifuging at 15,000 g for 15 min to precipitate the proteins. 90 µL of clear supernatant was taken and dried using a Speed-Vac centrifugal evaporator. The sample was re-dissolved to 75 µL with 2:1 H<sub>2</sub>O/ACN. After that, 25 µL of 250 mM sodium carbonate/sodium bicarbonate buffer was added to the sample to make a basic environment which is optimal for the dansylation reaction. The solution was vortexed, spun down, and mixed with 50 µL of freshly prepared <sup>12</sup>C-DnsCl solution (20 mg/mL) (for light labeling) or <sup>13</sup>C-DnsCl solution (20 mg/mL) (for heavy labeling). After the sample was incubated at 40 °C for 45 min, 10 µL of 250 mM NaOH was added to quench the excess dansyl chloride. The solution was then incubated at 40 °C for another 10 min to allow the NaOH react will all the leftover amount of the labeling reagent. Finally, 50 µL of formic acid (425 mM) in 1:1 ACN/H<sub>2</sub>O was used to consume excess NaOH and to make the solution acidic. The labeled samples were sent to LC-MS analysis or stored at – 80 °C. For each subject, two aliquots of serum were labeled independently as experimental duplicates.

### **3.2.4 DMPA-labeling of serum sample**

After thawed in an ice-bath, the serum sample was centrifuged at 15,000 g for 15 min. In a microcentrifuge tube, 30  $\mu\text{L}$  of supernatant was mixed with 90  $\mu\text{L}$  of acetonitrile. The mixture was then stored at  $-20\text{ }^{\circ}\text{C}$  for 2 hours to precipitate the proteins. After this, the mixture was centrifuged at 15,000 g for 15 min. 90  $\mu\text{L}$  of supernatant was mixed with 20  $\mu\text{L}$  of 0.5 M triethanolamine, which worked as the catalyst, and 50  $\mu\text{L}$  of freshly prepared  $^{12}\text{C}$ -DmPA bromide solution (10 mg/mL) (for light labeling) or  $^{13}\text{C}$ -DmPA bromide solution (10 mg/mL) (for heavy labeling). In an incubator set at  $85\text{ }^{\circ}\text{C}$ , the reaction was allowed to proceed for one hour. Finally, the sample was cooled down in an ice bath. The labeled samples were ready for LC-MS analysis, as well as long-term storage at  $-80\text{ }^{\circ}\text{C}$ . For each subject, two aliquots of serum were labeled independently as experimental duplicates.

### **3.2.5 LC-UV quantification and pre-acquisition sample normalization**

Inter-individual variations in total metabolite amount must be minimized in order to accurately assess the concentration differences caused by the factors being studied. An LC-UV based method<sup>115</sup> was applied to determine the total concentration of dansylated amine/phenol-containing metabolites based on the UV absorption of the dansyl group. The experiment was performed with a Waters ACQUITY UPLC system UPLC (Waters, Milford, MA, USA) and a Phenomenex Kinetex C18 column (2.1 mm  $\times$  5 cm, 1.7  $\mu\text{m}$  particle size) (Phenomenex, Torrance, CA, USA). Two microliters of each dansyl-labeled sample were injected for a fast step-gradient run. Solvent A was 0.1% (v/v) formic acid in

5% (v/v) ACN/H<sub>2</sub>O, and solvent B was 0.1% (v/v) formic acid in ACN. Starting at 0% B for 1 min, the gradient was then increased to 95% B within 0.01 min and held at 95% B for 1 min to ensure complete elution of all labeled metabolites. The flow rate was 0.45 mL/min, and the total UV absorption of dansyl-labeled metabolites in the sample was measured by a photodiode array (PDA) detector. The peak area, which can represent the total metabolite concentration in the sample, was integrated using the Empower software (6.00.2154.003). According to the quantification results, the <sup>12</sup>C- and <sup>13</sup>C-labeled samples were mixed in equal amounts for the following LC-MS analysis. Based on the assumption that there is a linear relationship between the total amount of amine/phenol-containing metabolites and the total amount of all metabolites in the sample, the sample normalization of acid-labeling was also based on the LC-UV measurements of the dansyl-labeling.

### **3.2.6 LC-FTICR-MS analysis**

The LC-FTICR-MS analysis was performed using an Agilent 1100 series binary system (Agilent Palo Alto, CA) connected to a 9.4 T Apex-Qe FT-ICR-MS (Bruker, Billerica, MA). The MS data were acquired in the positive ion mode with an electrospray ionization (ESI) source. An Agilent reversed-phase Eclipse plus C18 column (2.1 mm×100 mm, 1.8 μm particle size, 95 Å pore size) was used for chromatographic separation. Solvent A was 0.1% (v/v) formic acid in 5% (v/v) ACN/H<sub>2</sub>O, and solvent B was 0.1% (v/v) formic acid in ACN.

The 32-min gradient for all dansyl-labeled samples was as follows: 0 min (20% B), 0-3.5 min (20-35% B), 3.5-18 min (45-65% B), 18-21 min (65-95% B), 21-24 min (95-99% B), and 24-32 min (99% B). The column was re-equilibrated with the initial mobile phase condition for 15 min before injecting the next sample. The flow rate was 180  $\mu\text{L}/\text{min}$ , and the injection volume was calculated according to the optimal injection amount.

The 40-min gradient for all DmPA-labeled samples was: 0 min (20% B), 0-9 min (20%-50% B), 9-22 min (50%-65% B), 22-26 min (65%-80% B), 26-28 min (80%-98% B) and 28-40 min (98% B). The column was re-equilibrated with the initial mobile phase condition for 15 min before injecting the next sample. The flow rate was 180  $\mu\text{L}/\text{min}$ , and the injection volume was calculated according to the optimal injection amount.

The representative LC-MS chromatograms are provided in Appendix Figure 1.

### **3.2.7 Data processing and statistical analysis**

With the  $^{12}\text{C}/^{13}\text{C}$ -mixed sample, metabolites are detected as  $^{12}\text{C}/^{13}\text{C}$ -peak pairs instead of metabolite features given by single peaks. For each metabolite, a  $^{12}\text{C}$ -labeled peak represents the individual sample, and a  $^{13}\text{C}$ -labeled peak represents the USS sample. Therefore, the intensity ratio of these two peaks is the relative concentration of the metabolite in the individual sample. The picking of peak pairs was done by our in-house developed IsoMS software.<sup>72</sup> The second step was aligning the detected metabolite

concentrations from the individual samples into a summarized data sheet. After the alignment, it is common to have missing values in the data set, which are due to the loss of weak signals during the data processing. We have developed a Zero-fill program to recover most of the missing information.<sup>137</sup> Since the IsoMS calculates the peak pair ratio by the intensity values at the highest chromatogram point, imperfect LC peak shapes may affect the accuracy of the relative quantification. To overcome this issue, we employed the IsoQuant software to re-calculate the peak pair ratios by the LC peak areas.<sup>157</sup> The final file with relative metabolite concentrations was exported to SIMCA-P+ 12.0 software (Umetrics, Umeå, Sweden) for multivariate statistical analysis.

### **3.2.8 Metabolite identification**

We performed three levels of metabolite identification: positive identification, putative identification with the Human Metabolome Database (HMDB) library,<sup>158</sup> and putative identification with the Evidence-based Metabolome Library (EML).<sup>41</sup> The positive identification of the dansyl-labeled metabolites was done by a Dansyl Library<sup>159</sup> which contains 315 metabolite standards. Meanwhile, the definitive identification of DMPA-labeled metabolites was conducted with a developing acid standard library, which currently has 187 carboxyl-containing metabolite standards. Putative identification refers to a structure in the database with its accurate mass matched to that of a detected peak pair or metabolite in the samples. These matched structures can be used as the starting point for future confirmation including the synthesis of standards to positively identify the putative matches. In our work, putative identification was done based on accurate mass matches to the metabolites in the HMDB library (8,021 known human endogenous metabolites) and

the EML library (375,809 predicted human metabolites with one biological reaction). The mass accuracy tolerance window was set at 0.008 Da for database search.

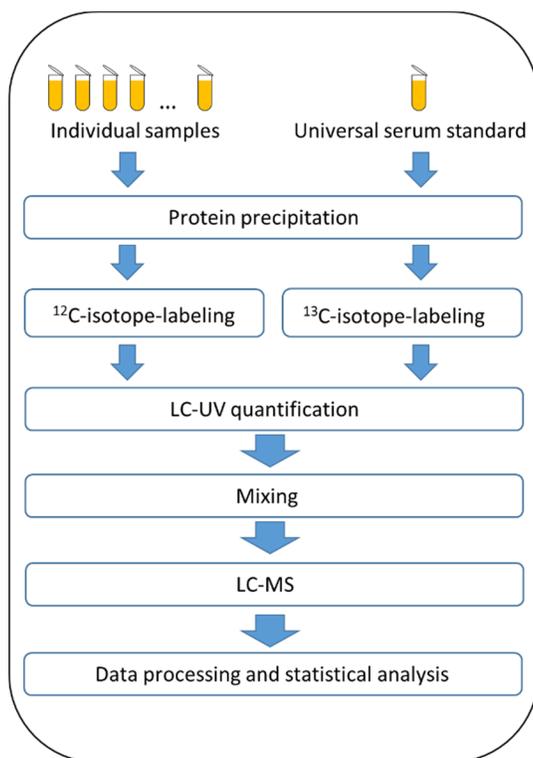
### **3.3 Results and Discussion**

#### **3.3.1 Quality control and sample normalization**

One of the major goals of our work in this chapter is the establishment of a high-coverage human serum metabolome database that can provide the average metabolite concentrations and variations among healthy people for future biomarker discovery studies. To achieve this goal, the metabolite quantification should be as accurate as possible. For large-scale metabolomics profiling works, especially those employing LC-MS, the issue of signal intensity drift over time is a major confounding factor which makes the metabolite quantification less accurate.<sup>62</sup> Also, as the serum is a complex biological mixture, the MS analysis may also suffer from the matrix effects.<sup>160</sup> Hence, the quality of the LC-MS data acquisition needs to be carefully monitored.

A major advantage of our differential isotope labeling methods is that the metabolite concentration is measured by the intensity ratio of a peak pair, instead of the absolute intensity of a single mass peak. According to the workflow shown in Figure 3.1, an individual sample is labeled by the <sup>12</sup>C-labeling reagent (<sup>12</sup>C-dansyl chloride or <sup>12</sup>C-DMPA bromide), and then mixed with the <sup>13</sup>C-labeled USS sample. Consequently, for each metabolite, a peak pair is detected in the LC-MS analysis. The light peak is from the <sup>12</sup>C-labeled metabolite in the individual sample, and the heavy peak is from the <sup>13</sup>C-labeled same metabolite in the USS sample. The intensity ratio of the two peaks can represent the

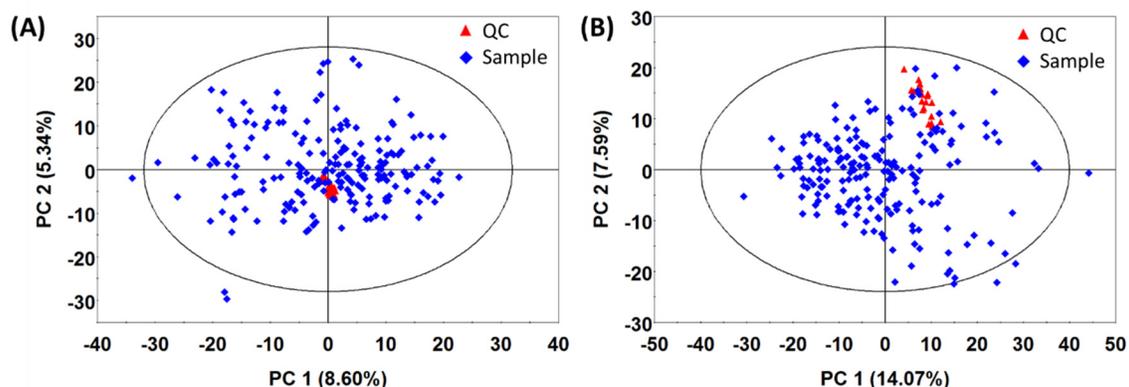
relative concentration of the metabolite in the individual sample. The absolute intensities of the light peak and the heavy peak may vary over time, but the ratio of them in a specific sample will remain constant. We have also confirmed that, even if the sample matrix is very extreme, the accuracy of relative quantification will not be affected.<sup>161</sup>



**Figure 3.1** Workflow of the non-targeted serum metabolome profiling using chemical isotope labeling and high-resolution LC-MS.

In our work, we applied a QC run after every tenth LC-MS analysis of samples. The QC sample was an equal-volume mixture of the <sup>12</sup>C-labeled and <sup>13</sup>C-labeled USS samples. The variations among QC runs represent the technical variations. If the technical variations are larger than the inter-subject variations, we will not be able to differentiate between the background noise and the metabolome changes that truly represent biological conditions. Therefore, the variations among QC runs must be smaller than the inter-subject variations.

As this is a typical multi-variate problem, we employed Principal Components Analysis (PCA) as an effective and straightforward way to study the variations. PCA linearly converts a combination of variables (metabolites) into a principal component. The first principal component (PC) accounts for as much of the variability in the data as possible. The second PC is orthogonal to the first and covers the second largest variance. The PCA score plots, as shown in Figure 3.2, project the high-dimensional data onto the 2-dimensional surface of PC 1 and PC 2, so that the inter-sample variances and inter-group variances can be visualized. In Figure 3.2A, the individual data points and the QC data points are plotted on the same surface. Unlike the individual data points, which randomly spread over the surface, the QC data points closely cluster, suggesting that technical variations are much smaller than the inter-subject variations. The comparison between DMPA-labeled individual samples and QCs in Figure 3.2B also confirms that our differential isotope labeling is a robust quantification technique and the quantification accuracy through our work is reliable.

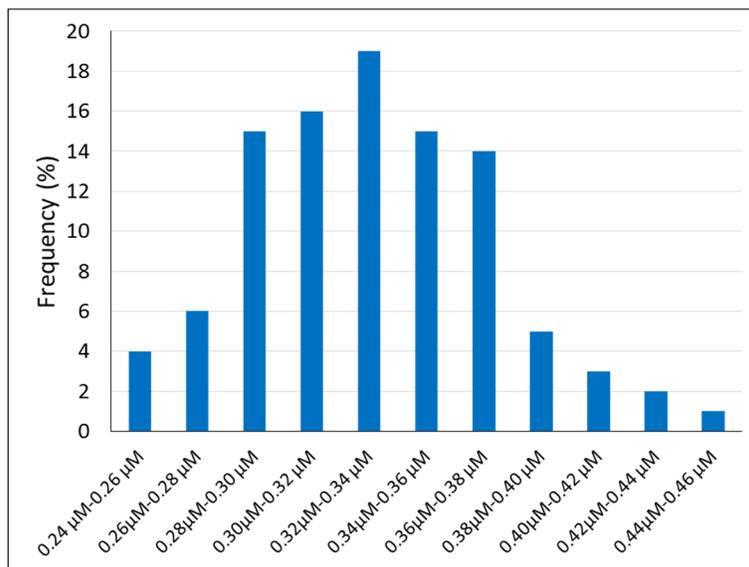


**Figure 3.2** (A) PCA score plot for dansylation LC-MS data obtained from 100 healthy subjects (in blue) and 20 QC runs (in red). (“PC” represents for “principal component” and the corresponding percentage is the percentage of the variance among all the data points that this principal component

covers.) (B) PCA score plot for DMPA-labeled LC-MS data obtained from 100 healthy subjects (in blue) and 20 QC runs (in red).

In addition to the instrumental drift, the inter-sample variations in the total metabolite concentration may also interfere with the metabolite quantification. For example, the concentration variability of urine due to the dilution effect of water can be very significant and the corresponding normalization is an indispensable step.<sup>162</sup> For serum metabolome, Roy et al. reported that the extent of concentration variation was between 0.8 and 1.2.<sup>114</sup> Although the variation of total metabolite concentration in serum is not as large as that in urine, it should be adjusted in order to acquire the high-quality data for the metabolome database.

In our work, we employed an LC-UV based method to quantify the total concentration of amine/phenol-containing metabolites in a dansyl-labeled sample.<sup>115</sup> Figure 3.3 shows the distribution of total metabolite concentration among serum samples from the 100 subjects. The average value is 0.34 mM, and the standard deviation is 0.04 mM, demonstrating a relatively narrow distribution. The most concentrated sample is 40% more concentrated than the average level, and the most diluted one's concentration is 29% lower than the average. According to the LC-UV quantification result, each dansyl-labeled individual sample was mixed with the USS sample at the same total metabolite amount. Based on the assumption that the total amount of the amine/phenol-submetabolome can proportionally represent the total amount of the whole metabolome, as well as other submetabolomes, we also used the LC-UV result to implement the normalization of DMPA-labeled samples.



**Figure 3.3** Distribution of the total metabolite concentration among the serum samples from 100 people.

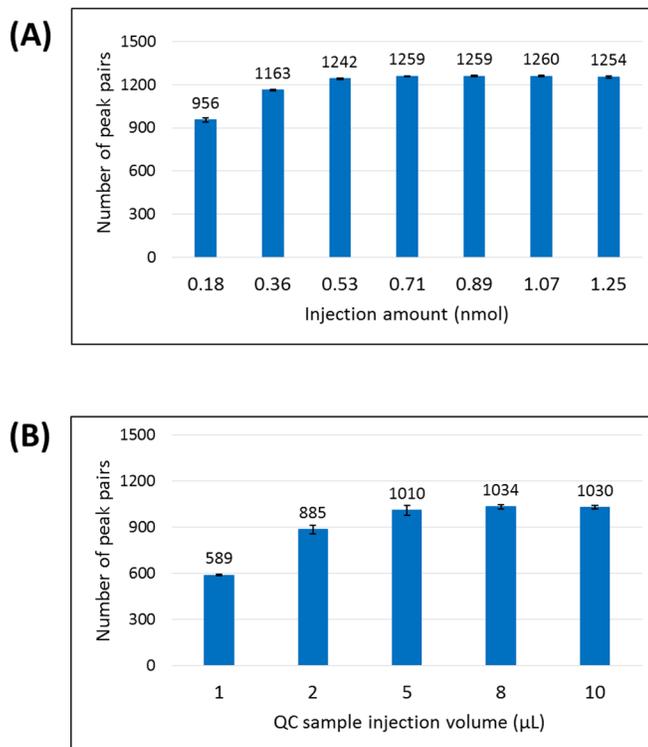
### 3.3.2 Development of a human serum metabolome database

A high-quality metabolome database should have as many high-confidence metabolites as possible. Since traditional LC-MS based profiling comes with background noise, protonated and deprotonated ions, adduct ions, fragment ions, dimers and trimers, the number of actually detected metabolites are always lower than the number of reported metabolite features.<sup>62</sup> Our method, on the other hand, gives a much more credible number of metabolites. This is because the derivatization can provide enhanced ionization efficiency of metabolites and thereby improve the MS signal by 10 to 1,000 fold. This not only allows us to achieve a much higher metabolome coverage, but also helps with differentiating the true metabolites from the background noises that cannot be labeled as peak pairs. Moreover, the addition of the labeling group shifts low-mass metabolites to the

higher mass region, which usually has cleaner background than the low-mass region. So the signal-to-noise ratio is also improved. Although sometimes there are still adduct ions, in-source fragment ions and dimers existing, the IsoMS software is designed to automatically filter out these interferences, as well as the background noises. Furthermore, we manually checked the output to make sure each reported peak pair is high-confidence. At last, the peak pairs with more than 50% missing values were excluded.

The sample injection amount was also optimized to obtain an optimal number of peak pairs. As shown in Figure 3.4, different amounts of the QC sample was injected into the LC-MS system. The number of detected peak pairs increased with larger injection amount, and we found the optimal injection amount when the number of peak pairs saturated. Consequently, 1 nmol dansyl-labeled QC or 5  $\mu$ L of DMPA-labeled QC was injected, and the injection volume of each individual sample was calculated to reach the same total metabolite amount. No carryovers were observed at the chosen injection amount. In our serum metabolome database, we report 1,348 amine/phenol-containing metabolites and 1,065 carboxyl-containing metabolites. Undoubtedly, a total number of 2,413 high-confidence metabolites can cover many more potential biomarkers, and the study of variations among normal population can be more broad and comprehensive with this data set. The DMPA-labeling focuses on relatively hydrophobic organic acids, and it has very good selectivity against other functional groups. Therefore, the overlap between the two labeling methods should be very small. Nonetheless, we note that there are a limited number of metabolites labeled by both reagents, and we will address this issue in the future by studying more standards and developing filtration algorithms. In our work, we study the amine/phenol-

submetabolome and carboxyl-submetabolome separately to avoid false findings due to the overlap.



**Figure 3.4** (A) Injection optimization curve, showing the numbers of dansyl-labeled peak pairs when different amounts of the QC sample are analyzed by the LC-MS system. (B) Injection optimization curve, showing the numbers of DMPA-labeled peak pairs at different injection volumes of the QC sample.

The database includes the retention time, labeled mass, accurate mass, and identification result of each peak pair. The average concentration and relative standard deviation (RSD) are also provided. Three levels of identification are annotated as “Library”, “HMDB” and “EML”. “Library” refers to the 56 dansylated metabolites which are positively identified with our standard library of 315 amine/phenol-containing compounds and the 32 DMPA-labeled metabolites which are confirmed by the standard library of 187 carboxyl-containing

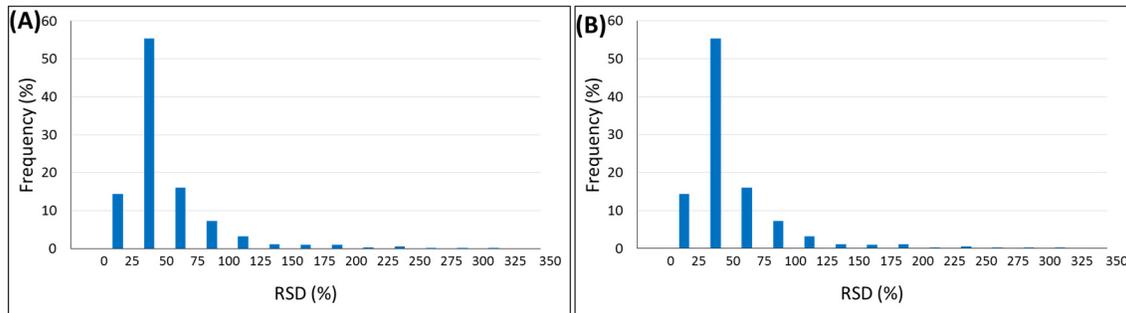
compounds. These 88 metabolites provide the highest confidence for further analyses. The standard libraries are being enriched, so in the future the structures of more metabolites in this database will be definitively verified.

For the rest of peak pairs, we conducted putative identification by matching the accurate mass to the HMDB library, which contains 8,021 known human endogenous metabolites.<sup>158</sup> 499 dansyl-labeled peak pairs and 298 DMPA-labeled peak pairs are putatively identified by this approach. With the high-resolution LC-MS platform, the data have very high mass accuracy and resolution. Also, unlike the other libraries which may include exogenous drugs and other artificial compounds, the HMDB library that we used focuses on known endogenous metabolites. As a result, the confidence of the identification should be highly acceptable. We note that this level of identification cannot be as accurate as the positive identification. In the future we will confirm the HMDB identification results with more standards or MS/MS spectra.

At last, one of the major challenges to metabolomics is that there are a large number of metabolites which remain unstudied and no standards are available for their identification. One possible origin of these unknown metabolites is the modification of primary metabolites. In the body, a known metabolite can be involved in various metabolic reactions in biological systems, producing different metabolic products.<sup>41</sup> In order to cover these modified metabolites, we have previously developed an EML library, which simulates 76 common metabolic reactions to the known metabolites of HMDB library and

enriches the library to 375,809 predicted human metabolites with one reaction. The EML library can identify 661 of the remaining dansyl-labeled metabolites and 658 of the remaining DMPA-labeled metabolites. With the suggested structures, it is easier to confirm the identities of them by MS/MS experiments in the future. Although the three levels of identification have covered 91% of the peak pairs, there are 132 dansyl-labeled and 77 DMPA-labeled peak pairs remaining unknown. With more metabolites and metabolic pathways discovered in the future, the structures of these metabolites will be eventually revealed.

Determining the inter-individual variation of each metabolite is also of great importance, as in a biomarker discover study, we need this information to evaluate how extreme an abnormal value is. The distributions of the RSDs are shown in Figure 3.5. For both the amine/phenol-submetabolome and the carboxyl-submetabolome, the majority of the RSDs are smaller than 100%. The RSDs of more than half of the metabolites are between 25% and 50%. We can conclude that the inter-individual variations for most metabolites are not extreme. However, these inter-individual variations can possibly be large enough to interfere with the application of biomarker candidates as many disease-dependent variations are not larger than 100%.



**Figure 3.5** (A) Distribution of relative standard deviations defining the inter-subject variability in the amine/phenol-containing metabolites. (B) Distribution of relative standard deviations defining the inter-subject variability in the carboxyl-containing metabolites.

In the future, we can employ our isotope labeling methods to any individual sample. As the labeled sample is mixed with the USS standard, we are able to collect the relative concentrations of the 2,413 metabolites. And for any of them, we have the RSD information to determine whether the concentration in the sample being studied is normal. This kind of experiment can be done in any laboratory with any LC-MS platform, and the relative quantification result will not be affected. The USS standard can also be applied to multiple biomarker discovery studies, enabling the inter-study comparisons.

### 3.3.3 Variations associated with sex

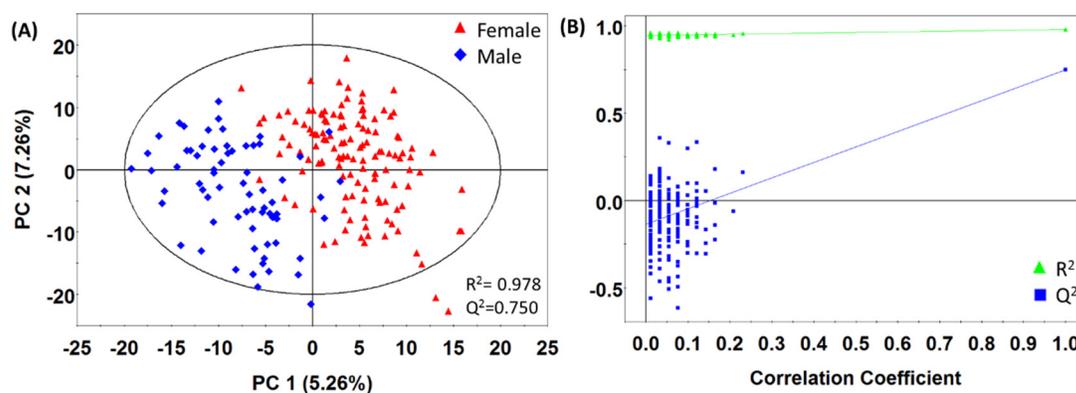
In addition to the concentration variability among the whole general population, the metabolome changes due to specific genetic or environmental factors should also be fully understood for the discovery and application of metabolite biomarkers. Sex, as a major genetic factor, has long been studied for the down-stream differences in proteome and metabolome.<sup>163-164</sup> According to Table 3.1, there are 65 female and 35 male subjects in our

study. Here we study the metabolome difference between males and females by another multi-variate statistical tool, the Partial Least Squares-Discriminant Analysis (PLS-DA). Unlike PCA, which is an unsupervised method, PLS-DA not only generates principal components, but also considers the group assignment and finds a linear regression model between the PCs and the grouping information. PLS-DA copes with the unwanted variances and focuses on the differences between study groups.

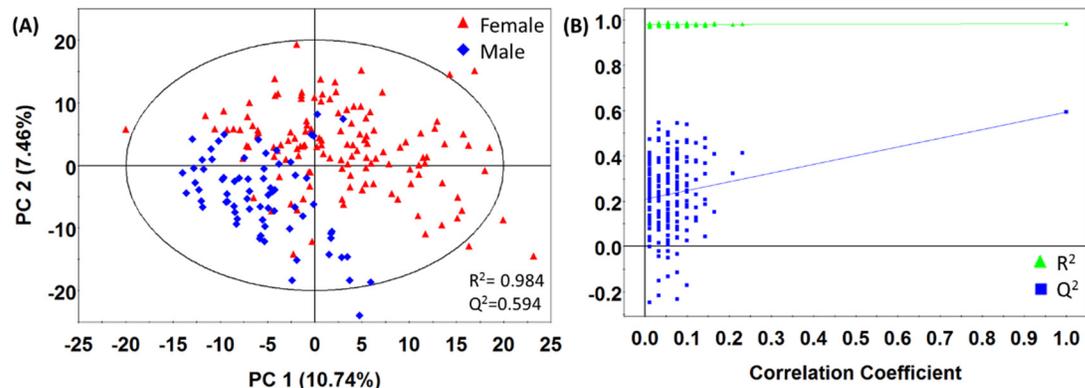
Figure 3.6A is the PLS-DA score plot for the separation between males and females. Clearly we can see that the two groups are separated in the direction of PC1, despite the fact that the inter-group variance is not much larger than the within-group variances, which are demonstrated in the direction of PC2. To evaluate the performance of a PLS-DA model, SIMCA-P outputs the  $R^2$  and  $Q^2$  values by the cross-validation process. When the data set is split into a training set and a testing set, the  $R^2$  value represents the quality of fitting for the building of the model, and the  $Q^2$  value demonstrates the power of the model to predict the testing set. Although the large number of variables in our study may always make the  $R^2$  relatively optimistic, the  $Q^2$  value can work as a major performance indicator of the separation model.

For the separation between males and females, the  $R^2$  is 0.978 and the  $Q^2$  is 0.750. It is generally accepted that  $Q^2 > 0.9$  means a perfect model and  $Q^2 > 0.5$  indicates an acceptable model, so our PLS-DA model can illustrate that there are statistically significant differences between the amine/phenol-submetabolomes of males and females. Since a

single  $Q^2$  value has no statistical significance, and the model can sometimes be over-fitted, a validation process is necessary to confirm the statistical findings. We used a permutation test to validate the PLS-DA models, and the result is shown in Figure 3.6B. The group assignment of the original data was randomized to generate multiple permuted data sets, and each of them outputs a pair of  $R^2$  and  $Q^2$ . The permuted data sets with low correlation efficient to the original data set should not be able to generate a good model and their  $Q^2$  values should be close to zero or even negative. On the permutation plot, if the y-axis intercept of the blue line is high, we will know that even totally randomized data can generate good models, raising doubts about the reliability of the original model. Usually, the intercept of the  $Q^2$  value should be negative. In Figure 3.6B, 200 permutation calculations were performed to the PLS-DA model described above, and the test successfully validated the model by showing a negative intercept.



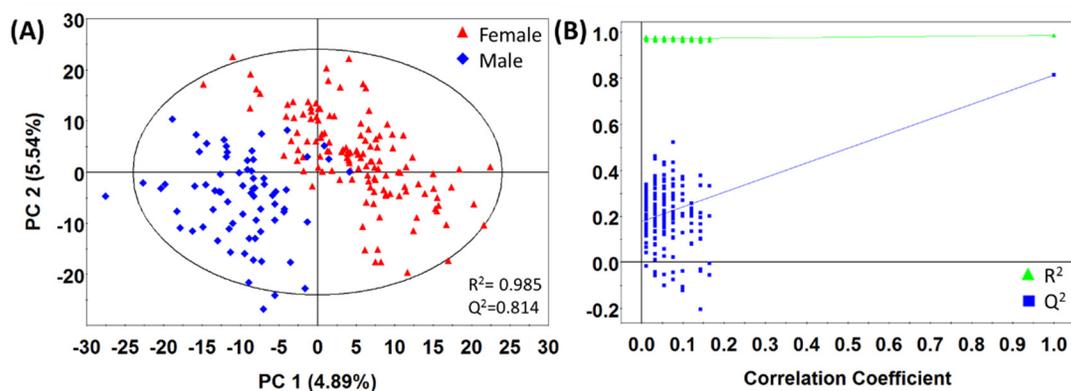
**Figure 3.6** (A) PLS-DA score plot for studying the statistical differences in the amine/phenol-submetbaolome between males and females. (B) Response permutation test result of the PLS-DA model for males vs. females (amine/phenol-submetablome).



**Figure 3.7** (A) PLS-DA score plot for studying the statistical differences in the carboxyl-submetabolome between males and females. (B) Response permutation test result of the PLS-DA model for males vs. females (carboxyl-submetabolome).

Figure 3.7A is the PLS-DA score plot for the differences in carboxyl-submetabolome between males and females. Unlike the amine/phenol-submetabolome, the separation between the two groups is not very significant. Interestingly, in Figure 3.6A there are two outliers of the male group positioning at the center of the female group, and in Figure 3.7A, the data points of these two outliers are lying at the same position, demonstrating consistency between the sub-metabolomes. The  $R^2$  value is 0.984, and the  $Q^2$  value is 0.594. Although the  $Q^2$  is larger than 0.5, when  $Q^2$  is substantially lower than  $R^2$ , the robustness of the model is questionable. The permutation plot in Figure 3.7B has a positive intercept, confirming that the model is not valid. We can conclude that although the amine/phenol-submetabolome has demonstrated clear sex-wise differences, the variations in the carboxyl-submetabolome is not significant.

We also examined the outcome of combining the data of the sub-metabolomes. The PLS-DA score plot is shown as Figure 3.8A. The distance between the two groups is clear, but not much more significant than that of the dansyl-labeled submetabolome. The  $Q^2$  value slightly increases to 0.814, but the permutation plot in Figure 3.8B rejects this model by showing a positive intercept. These results have illustrated that combining the data will not increase the discriminating power, and the large number of variables can increase the chance of over-fitting, which might account for the invalid permutation test.

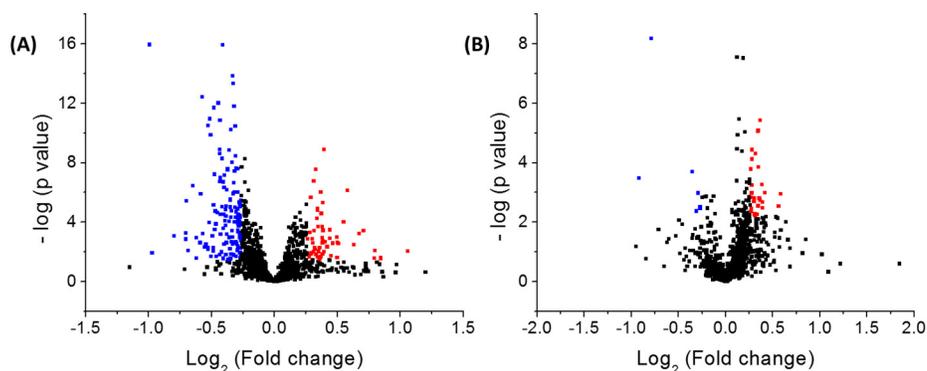


**Figure 3.8** (A) PLS-DA score plot of the combined amine/phenol-submetabolome and carboxyl-submetabolome, showing a separation between male and female subjects. (B) Response permutation test result of the PLS-DA model in Figure 3.8A. A total of 200 permutations were implemented.

To study the metabolites that have different concentration distributions in males and females, we did uni-variate analyses and made the volcano plot. A fold change was calculated as the ratio of the average concentration in the female group to that in the male group. For each metabolite, a t-test was employed to statistically study the concentration difference between the two groups. In order to overcome the multiple-testing problem, we

calculated the False-Discover-Rate-Adjusted p-value (q-value) for each metabolite. The calculation was performed by the “QVALUE” R package.<sup>110</sup> To control the false discovery rate, we let the q-value  $< 0.05$  and found the corresponding p-value threshold for selecting the statistically significant metabolites. A volcano plot is made by plotting  $-\log(\text{p-value})$  against  $\log_2(\text{fold change})$ , as shown in Figure 3.9A and 3.9B for the amine/phenol-submetabolome and carboxyl-submetabolome, respectively. If a metabolite's fold change is larger than 1.2 (or smaller than 0.83) and its p-value is smaller than the threshold that makes  $q < 0.05$ , we define it as a significant metabolite from the volcano plot. In Figure 3.9A, there are 147 significantly decreased metabolites (in blue) and 59 significantly increased metabolites (in red), making a total number of 206 dansyl-labeled metabolites having statistically significant sex differences out of the 1,348 metabolites in the whole list. In Figure 3.9B for the DMPA-labeled metabolites, there are seven significantly decreased ones, and 31 significantly increased ones. The number of significant metabolites is smaller because the statistical difference is not very significant. Furthermore, we note that statistical significance does not always lead to biological significance. Particularly, the p-value serves as a selection parameter but provides limited statistical meanings of the extent of variability for each metabolite. So in order to further improve the confidence of the findings, we also applied the uni-variate receiver operating characteristic (ROC) analysis, which illustrates the diagnostic ability of a binary classifier, to the significant metabolites in the volcano plots. For each metabolite, the ROC analysis reports an Area Under the Curve (AUC) as an indicator of the discriminative power. The AUC values are calculated with MetaboAnalyst 3.0.<sup>165</sup> In total, 148 dansyl-labeled and 31 DMPA-labeled significant

metabolites have been positively or putatively identified. Their identities, fold changes, q-values and AUCs are listed in Table 3.2.



**Figure 3.9** (A) Volcano plot for males vs. females (amine/phenol-submetabolome), showing 147 significantly decreased metabolites (fold change (female/male)  $< 0.83$ , p-value  $< 0.034$ ) in blue and 59 significantly increased metabolites (fold change  $> 1.2$ , p-value  $< 0.034$ ) in red. (B) Volcano plot for males vs. females (carboxyl-submetabolome), showing 7 significantly decreased metabolites (fold change  $< 0.83$ , p-value  $< 0.0062$ ) in blue and 31 significantly increased metabolites (fold change  $> 1.2$ , p-value  $< 0.0062$ ) in red.

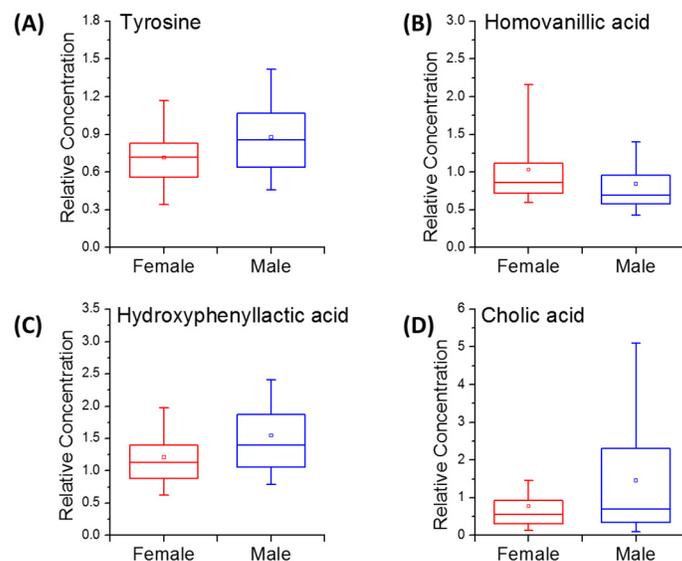
As shown in Table 3.2, 15 of the dansyl-labeled significant metabolites are positively identified: L-serine, homovanillic acid, L-aspartic acid, phenyl-alanyl-phenylalanine, L-aspartyl-L-phenylalanine, aminoadipic acid, leucyl-proline, L-glutamic acid, 3-hydroxymandelic acid, L-glutamic Acid [ $-\text{H}_2\text{O}$ ], L-tyrosine, L-leucine, L-methionine, trans-4-hydroxyl-L-proline and L-proline. In addition to the positively identified metabolites, 49 significant metabolites are putatively identified by the HMDB library and another 84 are putatively identified by the EML library. For the DMPA-labeled significant metabolites, five of them (hydroxyphenyllactic acid, 2-hydroxycaproic acid, 2-hydroxy-2-

methylbutyric acid, cholic acid and linoleic acid) are positively identified. Another five metabolites are HMDB-identified, and 21 metabolites are EML-identified.

Most of the positively identified dansyl-labeled compounds are amino acids or their derivatives, which agrees with the findings in many other metabolomics studies that amino acids are related to the sex-wise metabolome differences.<sup>34, 46, 166</sup> Importantly, sex differences in the regulation of amino acid metabolisms have also been reported,<sup>167-168</sup> underlying the sex differences in metabolome phenotype. For instance, we found that the average concentration of tyrosine in the male group was 24% higher than that in the female group. In Kawaguchi and his coworkers' report, they stated that testosterone, which is a sex hormone and a regulator of the synthetic enzyme of tyrosine, should be responsible for the significantly higher serum tyrosine levels in the male group, and the increased tyrosine level ulteriorly increased the insulin resistance.<sup>169</sup> This example tells us that a number of metabolic pathways can be different between sexes and studying these sex-wise metabolome differences can help us develop a deeper and more comprehensive understanding of the metabolic processes. Homovanillic acid has the second-highest univariate AUC (0.763) among the positively identified metabolites, demonstrating a high discriminating power for the sex-wise separation. Although the biological reason underlying the sex-dependent property of homovanillic acid is not clear, the concentration difference was also observed in other studies, showing the robustness of our method. For example, Koreen et al. studied the plasma homovanillic acid levels in 19 healthy volunteers, and reported that the concentration of the female group (7.0 +/- 2.3 ng/mL) was higher than that of the male group (6.3 + 1.9 ng/mL).<sup>170</sup> Homovanillic acid has been reported to be a

biomarker candidate for the diagnosis of Parkinson's disease<sup>23</sup> and also believed to be related to Alzheimer's disease.<sup>171</sup> In a biomarker discovery study, if the control group and the disease group are not sex-matched, the sex-wise difference of homovanillic acid will become a confounding factor. As the metabolome variability of sex is considerable, we suggest that in biomarker discovery studies the sex factor should always be carefully considered. The box plots showing the concentration distributions of tyrosine and homovanillic acid are shown as Figure 3.10A and 3.10B, respectively.

The DMPA-labeled significant metabolites are also important for revealing the biological differences between males and females. For example, hydroxyl-phenyllactic acid is a tyrosine metabolite, and its fold change (0.78) is close to that of tyrosine (0.82), suggesting the variability in the same metabolic pathway. Cholic acid is a primary bile acid, and our work shows that the serum concentration of cholic acid in men is almost twice as large as that in women, which is consistent with Fisher and Yousef's finding that the bile acid composition of human bile is sex-dependent.<sup>172</sup> Figure 3.10C and 3.10D are the box plots for the concentrations of hydroxyl-phenyllactic acid and cholic acid in males and females.

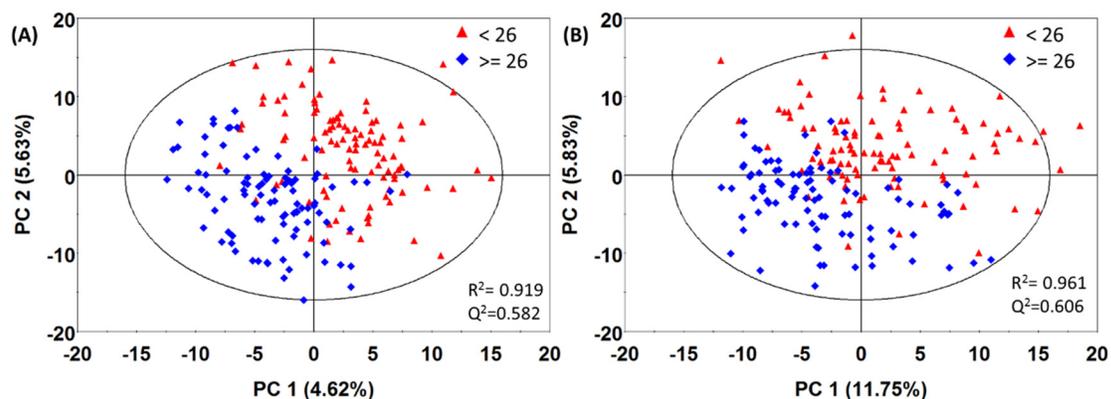


**Figure 3.10** Box plots, demonstrating the metabolite concentration distributions in the male group and the female group for (A) tyrosine, (B) homovanillic acid, (C) hydroxyl-phenyllactic acid, and (D) cholic acid.

### 3.3.4 Age effects

The aging process is believed to induce metabolome changes,<sup>173</sup> and the age-dependent metabolome variations have been reported in many reports.<sup>34, 174-175</sup> In our work, the participants aged 19 to 39 years old, with a median of 25 years old. So we split them into two groups: 52 participants younger than 26 years old, and 48 participants equal to or older than 26 years old. The PLS-DA score plots for the amine/phenol-submetabolome and carboxyl-submetabolome are shown as Figure 3.11A and Figure 3.11B, respectively. Neither of them demonstrates a clear separation between the two age groups. And as we can expect, both of the two  $Q^2$  values, 0.582 and 0.606, are very low. Only three dansyl-labeled metabolites are recognized as significant metabolites. One of them is EML-identified as Cystine [+CH<sub>2</sub>] and the other two are not identified. None of the DMPA-

labeled metabolites has q-value smaller than 0.05. The evidence leads us to believe that there are no age-related metabolome variations between our two age groups. Compared to the other works that always involve subjects older than 50 years, our study has a relatively young and narrow distribution of age. The aging effects can be significant among old people, but may not be noticeable among the relatively young population. We can conclude that, when all the subjects in a biomarker discovery study are younger than forty, the age factor can be safely ignored.

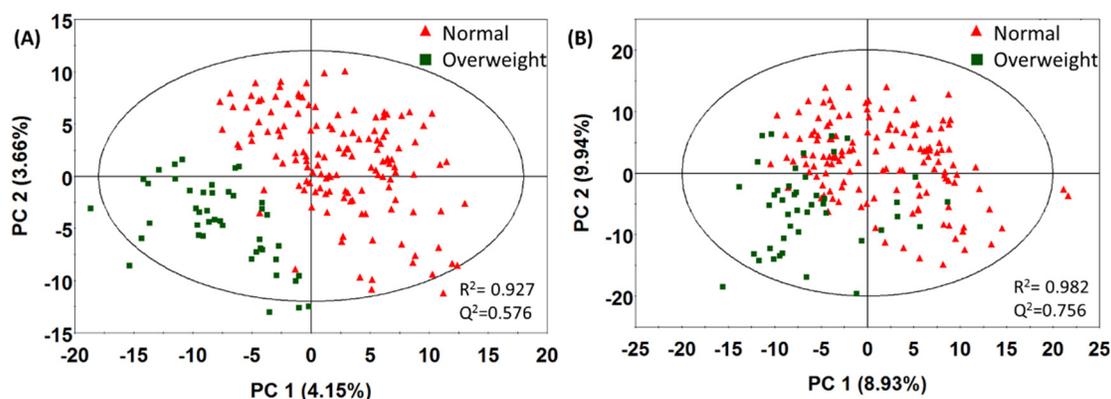


**Figure 3.11** (A) PLS-DA score plot for studying the statistical differences in the amine/phenol-submetabolome between the young group (< 26 years old) and the old group (≥ 26 years old). (B) PLS-DA score plot for studying the statistical differences in the carboxyl-submetabolome between the young group (< 26 years old) and the old group (≥ 26 years old).

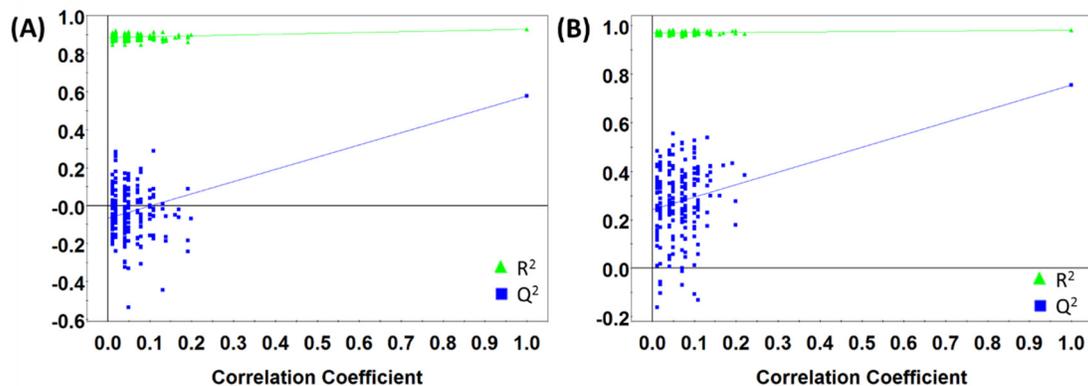
### 3.3.5 BMI effects

Obesity is a serious health problem that increases the risk of diabetes, cardiovascular disease, and cancer.<sup>176</sup> Metabolic abnormalities are believed to underlie this disorder and efforts have been made to find biomarkers for obesity.<sup>177</sup> Morio et al. have found that, in

comparison to normal people, the obese subjects demonstrated a different metabolic response to over-nutrition,<sup>178</sup> suggesting differences in metabolic pathways. As overweight and obesity are commonly seen conditions in the general population, it is worthwhile to study the related metabolic variations. Overweight can be defined as BMI between 25 and 30, and obesity refers to BMI larger than 30. In our work, the range of BMI is from 13.4 to 38.4, with a median of 22.5. For statistical analysis, we defined 9 subjects with BMI < 18.5 as the underweight group, 69 subjects with BMI between 18.5 and 25 as the normal group, and 22 subjects with BMI > 25 as the overweight group.



**Figure 3.12** (A) PLS-DA score plot for studying the statistical differences in the amine/phenol-submetabolome between the normal group and the overweight group. (B) PLS-DA score plot for studying the statistical differences in the carboxyl-submetabolome between the normal group and the overweight group.



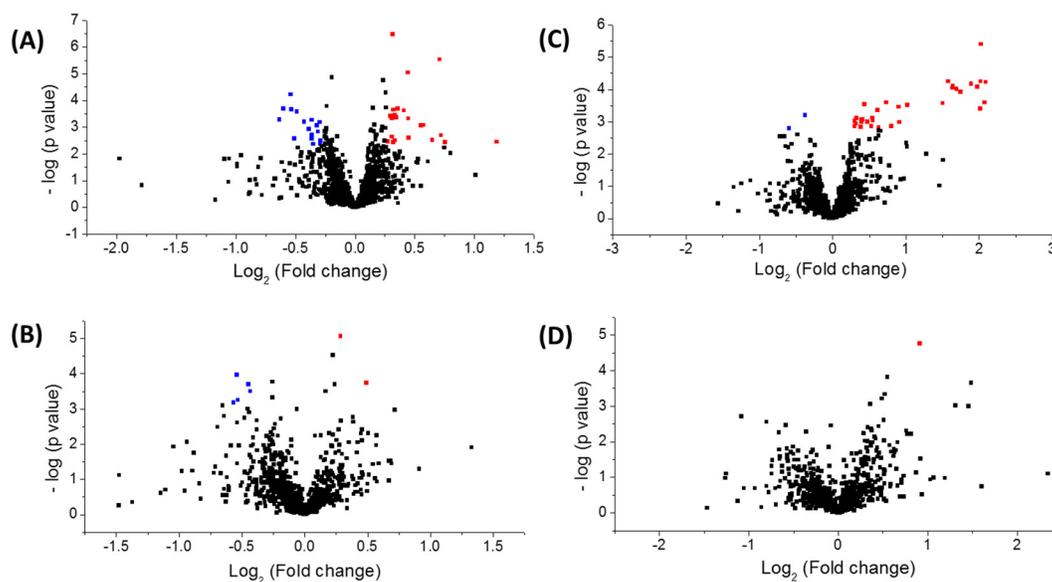
**Figure 3.13** (A) Response permutation test result of the PLS-DA model for overweight vs. normal (amine/phenol-submetabolome). (B) Response permutation test result of the PLS-DA model for overweight vs. normal (carboxyl-submetabolome).

The normal group has 35% of males and the overweight group has 50% of males. So the sex will not become a significant interfering factor. Figure 3.12A is the PLS-DA score plot for the dansyl-labeling data, showing a clear separation between the normal group and the overweight group. Although the  $Q^2$  (0.576) is relatively low, the model passed the permutation test, suggesting a minor but valid separation. For the DMPA-labeling data in Figure 3.12B, the two groups are visibly overlapping. Although the  $Q^2$  is higher than that of the amine/phenol-submetabolome, the model is rejected by the permutation test. The permutation test results are provided in Figure 3.13.

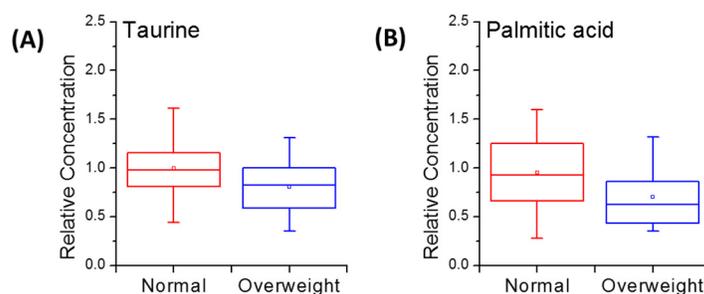
Figure 3.14A and 3.14B are the volcano plots for studying the overweight condition. The amine/phenol-submetabolome has 18 decreased (overweight/normal) metabolites and 22 increased metabolites. Meanwhile, the carboxyl-submetabolome still shows four decreased metabolites and two increased metabolites, despite the fact that the multi-variate difference

is not significant. Table 3.3 provides the information of 25 identified dansyl-labeled metabolites and six identified DMPA-labeled metabolites. Among them, taurine is the only positively identified one. Nonetheless, six are putatively identified by the HMDB library and another 24 are putatively identified by the EML library.

Taurine has a significantly lowered concentration in the overweight group, as shown by the box plot in Figure 3.15A. It has been proved that taurine produces a beneficial effect on lipid metabolism and helps with reducing body weight.<sup>179</sup> The concentration decrease of this beneficial regulator in the overweight subjects might be one of the reasons underlying the obese phenotype. Palmitic acid is an HMDB-identified metabolite, and its overweight/normal fold change is 0.74 (Figure 3.15B). It is one of the most common saturated fatty acids, and it also plays a role in the fat oxidation process.<sup>180</sup>

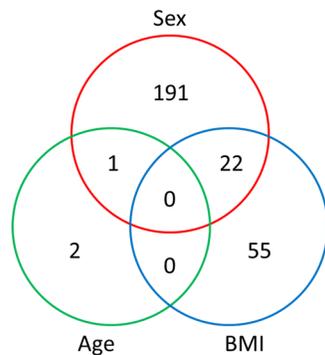


**Figure 3.14** (A) Volcano plot for overweight vs. normal (amine/phenol-submetabolome), showing 18 significantly decreased metabolites (fold change (overweight/normal) < 0.83, p-value < 0.005) in blue and 22 significantly increased metabolites (fold change > 1.2, p-value < 0.005) in red. (B) Volcano plot for overweight vs. normal (carboxyl-submetabolome), showing 4 significantly decreased metabolites (fold change < 0.83, p-value < 0.00066) in blue and 2 significantly increased metabolites (fold change > 1.2, p-value < 0.00066) in red. (C) Volcano plot for underweight vs. normal (amine/phenol-submetabolome), showing 2 significantly decreased metabolites (fold change (underweight/normal) < 0.83, p-value < 0.0016) in blue and 27 significantly increased metabolites (fold change > 1.2, p-value < 0.0016) in red. (D) Volcano plot for underweight vs. normal (carboxyl-submetabolome), showing one significantly increased metabolite (fold change > 1.2, p-value < 0.000018) in red.



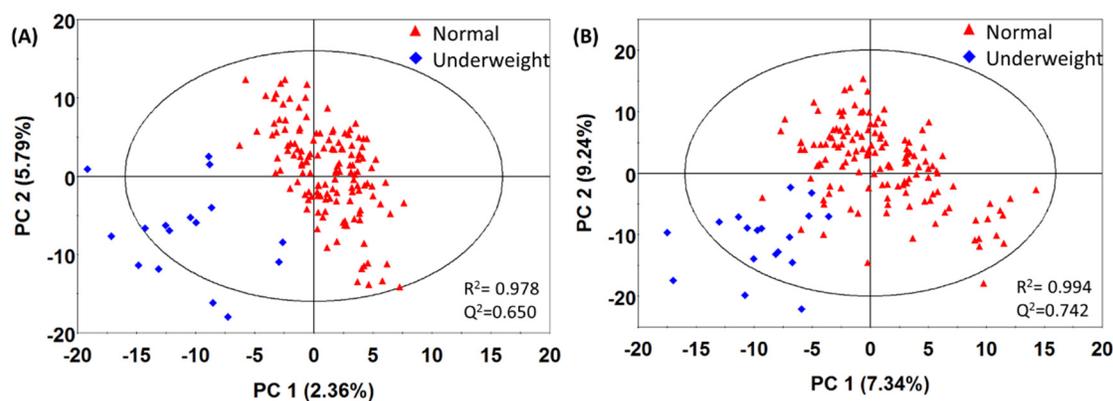
**Figure 3.15** Box plots, demonstrating the metabolite concentration distributions in the normal group and the overweight group for (A) taurine and (B) palmitic acid.

We also studied the metabolome variations due to the underweight condition. The underweight is a small group of nine subjects, and all of them are female, so the sex effect might be a confounding factor. The volcano plot analyses give 37 significant metabolites from the dansyl-labeling data and one significant metabolite from the DMPA-labeling data. Eight out of the 37 dansyl-labeled significant metabolites are also in the list of the sex effect and have similar fold changes in the two different comparisons. We excluded these eight metabolites from the following analysis to avoid interference of the sex effect. Moreover, to study the interactions among the three factors, we made the Venn diagram (Figure 3.16). For the amine/phenol-submetabolome, sex and age factors share one common significant metabolite, but the identity of this metabolite is unknown. 22 of the BMI-associated metabolites have significant sex interactions. Among them, eight are from the underweight group and have been excluded. Another 14 significant metabolites are shared by sex effect and overweight effect. These metabolites should be involved in the metabolic pathways that are not only regulated by sex-dependent genetic factors but also altered in the abnormal biological conditions of overweight. 11 of them have been putatively identified and highlighted with asterisks in Table 3.3. None of the DMPA-labeled metabolites demonstrated the inter-factor interactions.



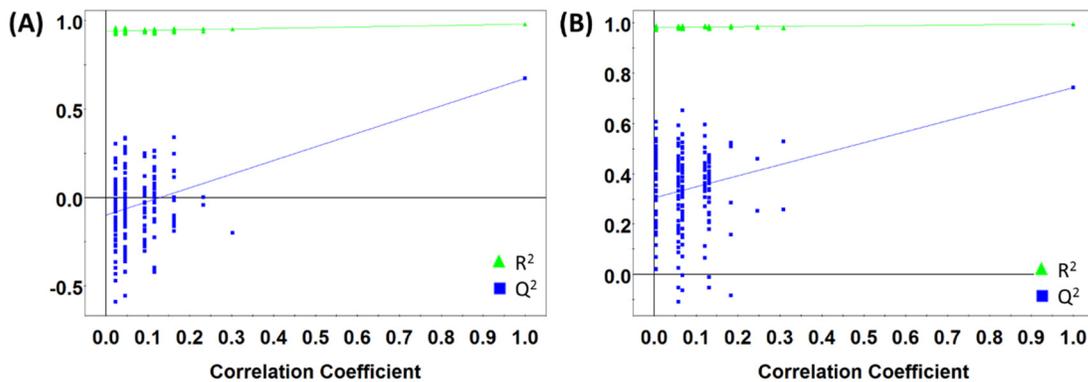
**Figure 3.16** Venn diagram for the volcano plot significant metabolites associated with sex, age and BMI. (All of these numbers refer to dansyl-labeled metabolites only.)

Figure 3.17A is the PLS-DA score plot for the underweight effect of the amine/phenol-submetabolome, showing a clear separation between the two groups. Although the  $Q^2$  (0.650) is relatively low, the model passed the permutation test (Figure 3.18A), suggesting an acceptable statistical difference. We can also find the group separation of the carboxyl-submetabolome in Figure 3.17B. However, the model failed the permutation test (Figure 3.18B).



**Figure 3.17** (A) PLS-DA score plot for studying the statistical differences in the amine/phenol-submetabolome between the normal group and the underweight group. (B) PLS-DA score plot for

studying the statistical differences in the carboxyl-submetabolome between the normal group and the underweight group.



**Figure 3.18** (A) Response permutation test result of the PLS-DA model for underweight vs. normal (amine/phenol-submetabolome). (B) Response permutation test result of the PLS-DA model for underweight vs. normal (carboxyl-submetabolome).

The volcano plot for the amine/phenol-submetabolome in Figure 3.14C has two decreased metabolites and 27 increased metabolites due to underweight. An isomer of methionine sulfoxide, an isomer of hypoxanthine and tryptophyl-phenylalanine are positively identified, with another 15 metabolites identified by HMDB and seven metabolites identified by EML library. Only one DMPA-labeled metabolite, which is putatively identified as glyoursodeoxycholic acid, is a statistically significant variable, as shown in Figure 3.14D. The information of these metabolites is given in Table 3.4. Although the metabolomics studies on underweight condition are scarce compared to those on obesity, our result has revealed that there should be abnormally regulated metabolic processes associated with the low body weight. With more positively identified metabolites and well-

studied pathways available in the future, the metabolic information about the BMI factor will be greatly enriched.

### **3.4 Conclusions**

Overall, we have successfully developed a high-coverage human serum metabolome database. In future metabolomics applications, with the use of differential isotope labelings and universal standard, the relative concentrations and statistical significances of these metabolites can be easily obtained from any individual sample, even with a very limited sample amount. Moreover, our work has demonstrated metabolic variations associated with sex and BMI. On the one hand, these variations are potential confounding factors in biomarker discovery and need to be carefully monitored. On the other hand, they can reveal the potential of metabolomics toward personalized health assessment. By studying these variations and the related metabolic pathways, we may develop a deeper understanding of the individual-dependent metabolic responses to environmental stimulations. Other lifestyle-related factors, such as coffee consumption and smoking, were strictly excluded in this study. However, they are very common variables accompanied with most clinical blood collections, so the metabolic impact of these factors should also be studied in the future.

Interestingly, the variability in the carboxyl-submetabolome is insignificant compared to that of the amine/phenol-submetabolome, suggesting differences in metabolic pathways for different groups of metabolites. In addition to our future plans of increasing the number of

positively identified metabolites in the current database, we will also expand the database with another two labeling techniques targeting at the hydroxyl-submetabolome<sup>74</sup> and carbonyl-submetabolome.<sup>75</sup> At last, this work is part of a collaborative project with Zhejiang University, Hangzhou, China, in the area of developing a Canadian and Chinese Metabolome Database (CCMD). By studying subjects living in different environments, it is possible to discover the metabolic responses to more environmental factors.

Note: In the work of this chapter, Tran Tran helped with the blood collection and Minglei Zhu helped with the DMPA-labeling of serum samples.

**Table 3.2** 148 identified dansyl-labeled metabolites and 31 identified DMPA-labeled metabolites that have fold change (female/male) > 1.2 (or < 0.83) and q value < 0.05 for the difference between males and females.

Label	Retention Time (min)	Detected m/z	Accurate Mass (Da)	ID Level	Compound Name	HMDB ID	Fold Change	q-value	Uni-variate AUC
Dansyl-36	4.02	391.0633	157.0050	EML	Alloxan [+NH]	HMDB 02818	1.26	2.67 E-02	0.740
Dansyl-56	4.63	373.0885	139.0301	EML	Taurine [+CH2]	HMDB 00251	0.62	1.85 E-02	0.550
Dansyl-62	4.78	531.1029	297.0446	HMDB	L-Cysteinylglycine disulfide	HMDB 00709	0.83	4.35 E-05	0.555
Dansyl-63	4.84	385.0883	151.0300	EML	Homocysteine [+O]	HMDB 00742	1.29	3.54 E-02	0.658
Dansyl-71	5.05	385.0341	150.9757	EML	Thiocysteine [-H2]	HMDB 03585	0.77	2.90 E-02	0.597

Dansyl-78	5.22	558.1586	324.1002	EML	L-Aspartyl-L-phenylalanine [+CO <sub>2</sub> ]	HMDB 00706	1.42	7.44 E-03	0.761
Dansyl-80	5.24	509.1695	275.1111	HMDB	Gamma-Glutamylglutamine	HMDB 11738	0.80	7.76 E-12	0.622
Dansyl-88	5.36	363.1008	129.0425	HMDB	Pyroglutamic acid	HMDB 00267	0.79	3.28 E-12	0.627
Dansyl-106	5.59	380.1637	146.1054	EML	Pipecolic acid [+NH <sub>3</sub> ]	HMDB 00070	0.73	3.37 E-02	0.576
Dansyl-132	6.25	452.1482	218.0898	HMDB	5-L-Glutamyl-L-alanine	HMDB 06248	0.80	1.47 E-05	0.551
Dansyl-141	6.41	423.1339	189.0756	HMDB	S-Prenyl-L-cysteine	HMDB 12286	1.20	2.76 E-02	0.678
Dansyl-144	6.43	367.0956	133.0373	HMDB	Iminodiacetate	HMDB 11753	1.29	2.11 E-04	0.769
Dansyl-152	6.55	339.0993	105.0410	Library	L-Serine	HMDB 00187	1.31	5.42 E-08	0.843
Dansyl-153	6.55	424.1132	190.0549	HMDB	L-beta-aspartyl-L-glycine	HMDB 11165	1.21	3.61 E-02	0.736
Dansyl-155	6.55	362.1905	128.1322	EML	Spermidine [-NH <sub>3</sub> ]	HMDB 01257	0.77	6.63 E-03	0.508
Dansyl-158	6.58	371.0861	137.0278	EML	Aminomalonic acid [+H <sub>2</sub> O]	HMDB 01147	1.25	6.95 E-07	0.826
Dansyl-163	6.61	410.1042	176.0459	HMDB	Ureidosuccinic acid	HMDB 00828	1.23	4.14 E-03	0.739
Dansyl-165	6.61	339.1376	105.0793	HMDB	Diethanolamine	HMDB 04437	1.22	2.94 E-05	0.812
Dansyl-168	6.65	505.2224	271.1640	EML	Homoanserine [+NH <sub>3</sub> ]	HMDB 05767	0.76	5.06 E-04	0.565
Dansyl-169	6.68	477.1432	243.0849	HMDB	Cytidine	HMDB 00089	1.41	3.41 E-03	0.746
Dansyl-172	6.83	477.1798	243.1215	EML	5-Methyldeoxycytidine [+H <sub>2</sub> ]	HMDB 02224	0.77	1.14 E-05	0.572

Dansyl-181	6.96	381.1108	147.0525	Library	L-Glutamic Acid	HMDB 00148	0.83	1.34 E-02	0.593
Dansyl-182	6.96	367.0924	133.0341	Library	L-Aspartic Acid	HMDB 00191	1.41	6.57 E-03	0.749
Dansyl-184	6.97	403.0897	169.0314	HMDB	2-Furoylglycine	HMDB 00439	0.83	2.28 E-02	0.545
Dansyl-186	6.98	365.1160	131.0577	Library	Trans-4-Hydroxyl-L-Proline	-	0.83	4.07 E-02	0.516
Dansyl-218	7.47	395.1272	161.0689	Library	Amino adipic acid	HMDB 00510	0.74	3.85 E-06	0.639
Dansyl-238	7.58	511.1745	277.1162	HMDB	Queuine	HMDB 01495	0.78	4.31 E-05	0.582
Dansyl-254	7.79	428.1366	194.0783	EML	N-Formyl-L-methionine [+NH3]	HMDB 01015	2.08	1.96 E-02	0.646
Dansyl-272	8.05	379.1321	145.0738	HMDB	(S)-5-Amino-3-oxohexanoate	HMDB 12131	1.73	4.16 E-02	0.654
Dansyl-275	8.10	480.1789	246.1206	HMDB	L-beta-aspartyl-L-leucine	HMDB 11166	0.80	1.54 E-10	0.637
Dansyl-277	8.10	512.1514	278.0931	EML	L-Aspartyl-L-phenylalanine [-H2]	HMDB 00706	0.75	4.14 E-14	0.685
Dansyl-294	8.31	452.1852	218.1269	EML	Diaminopimelic acid [+C2H4]	HMDB 01370	1.28	2.11 E-02	0.698
Dansyl-295	8.35	418.0816	184.0233	EML	Cysteic acid [+NH]	HMDB 02757	0.72	1.73 E-10	0.668
Dansyl-305	8.66	348.1381	114.0798	HMDB	3-Amino-2-piperidone	HMDB 00323	0.77	8.58 E-05	0.583
Dansyl-312	8.79	402.0862	168.0279	EML	Xanthine [+O]	HMDB 00292	0.70	6.73 E-09	0.686
Dansyl-321	8.86	431.1727	197.1144	EML	1-Methylhistidine [+C2H4]	HMDB 00001	0.70	3.48 E-02	0.581
Dansyl-324	8.94	420.1222	186.0639	EML	2-Furoylglycine [+NH3]	HMDB 00439	1.25	1.23 E-02	0.665
Dansyl-326	9.06	374.0808	140.0224	EML	Taurine [+NH]	HMDB 00251	1.36	3.77 E-03	0.757

Dansyl-349	9.33	502.1763	268.1180	EML	Carnosine [+C2H2O]	HMDB 00033	0.50	4.14 E-14	0.800
Dansyl-353	9.41	399.1011	165.0428	EML	2-Aminobenzoic acid [+CO]	HMDB 01123	1.46	5.77 E-04	0.751
Dansyl-354	9.49	420.1588	186.1004	EML	Glycylproline [+CH2]	HMDB 00721	0.83	5.16 E-03	0.527
Dansyl-356	9.53	478.2015	244.1432	EML	L-leucyl-L-proline [+O]	HMDB 11175	0.67	2.27 E-03	0.574
Dansyl-359	9.59	379.1329	145.0746	HMDB	Isobutyrylglycine	HMDB 00730	1.41	3.89 E-02	0.616
Dansyl-369	9.94	416.0648	182.0065	EML	Phosphoguanidinoacetate [-NH]	HMDB 03705	0.72	1.41 E-03	0.613
Dansyl-370	9.94	293.0951	59.0368	HMDB	N-Methylformamide	HMDB 01122	1.31	2.03 E-02	0.699
Dansyl-371	9.94	459.1690	225.1107	EML	Anserine [-NH]	HMDB 00194	0.72	1.36 E-06	0.635
Dansyl-372	9.95	478.1996	244.1413	EML	L-isoleucyl-L-proline [+O]	HMDB 11174	0.76	8.80 E-04	0.568
Dansyl-375	10.04	402.0861	168.0278	EML	Oxypurinol [+O]	HMDB 00786	0.77	2.79 E-04	0.583
Dansyl-382	10.22	494.1958	260.1375	HMDB	L-gamma-glutamyl-L-isoleucine	HMDB 11170	0.79	3.15 E-09	0.642
Dansyl-390	10.52	500.1845	266.1262	EML	Gamma-Glutamyltyrosine [+CO2]	HMDB 11741	1.36	3.19 E-02	0.695
Dansyl-395	10.60	402.1131	168.0547	EML	Imidazoleacetic acid [+C2H2O]	HMDB 02024	0.82	1.22 E-02	0.546
Dansyl-398	10.83	399.0674	165.0091	EML	2-Pyrrolidinone [+SO3]	HMDB 02039	0.82	2.35 E-04	0.519
Dansyl-403	11.00	390.1113	156.0529	HMDB	5-Hydroxymethyl-4-methyluracil	HMDB 00544	1.27	3.67 E-02	0.674
Dansyl-409	11.19	514.1634	280.1051	HMDB	L-beta-aspartyl-L-phenylalanine	HMDB 11167	1.23	2.39 E-02	0.653

Dansyl-411	11.25	363.1010	129.0427	Library	L-Glutamic Acid [-H2O]	HMDB 00148	0.83	3.48	0.573	E-03
Dansyl-419	11.44	480.1806	246.1223	HMDB	L-gamma-glutamyl-L-valine	HMDB 11172	1.32	6.02	0.681	E-03
Dansyl-429	11.67	393.1485	159.0902	HMDB	2-Methylbutyrylglycine	HMDB 00339	0.74	2.84	0.528	E-02
Dansyl-435	11.77	514.1632	280.1048	Library	L-Aspartyl-L-phenylalanine	HMDB 00706	1.23	2.35	0.653	E-02
Dansyl-436	11.82	438.0623	204.0039	EML	4-Hydroxybenzyl alcohol [+SO3]	HMDB 11724	1.80	4.02	0.643	E-02
Dansyl-441	11.97	420.1161	186.0578	EML	L-glycyl-L-hydroxyproline [-H2]	HMDB 11173	0.83	2.55	0.522	E-03
Dansyl-443	11.99	411.0908	354.0649	EML	5-amino-1-(5-phospho-D-ribose)imidazole-4-carboxylate [+NH]	HMDB 06273	0.73	1.85	0.615	E-04
Dansyl-444	11.99	349.1201	115.0617	Library	L-Proline	HMDB 00162	0.82	5.52	0.514	E-04
Dansyl-446	11.99	349.1383	115.0800	EML	Creatinine [+H2]	HMDB 00562	0.82	9.12	0.524	E-04
Dansyl-454	12.05	349.0732	115.0149	EML	Homocysteine thiolactone [-H2]	HMDB 02287	0.82	1.53	0.516	E-03
Dansyl-466	12.39	411.0911	354.0657	EML	Famotidine [+NH3]	HMDB 01919	0.72	7.78	0.625	E-05
Dansyl-480	12.64	383.1210	149.0627	HMDB	1-Methyladenine	HMDB 11599	0.79	2.72	0.538	E-02
Dansyl-488	12.71	383.1057	149.0474	Library	L-Methionine	HMDB 00696	0.83	7.99	0.521	E-04
Dansyl-491	12.74	411.0847	354.0528	EML	Phloretin [+HPO3]	HMDB 03306	0.64	6.42	0.684	E-06
Dansyl-500	13.15	418.1428	184.0845	EML	Pyridoxine [+NH]	HMDB 00239	0.82	8.78	0.517	E-04
Dansyl-501	13.16	400.1071	166.0488	HMDB	3-Methylxanthine	HMDB 01991	1.74	1.81	0.660	E-02

Dansyl-503	13.22	448.1897	214.1314	HMDB	Dethiobiotin	HMDB 03581	0.77	7.22 E-07	0.613
Dansyl-525	13.47	364.6244	261.1321	HMDB	Aspartyllysine	HMDB 04985	0.77	6.95 E-07	0.639
Dansyl-530	13.57	442.1081	208.0498	EML	L-Dopachrome [+NH]	HMDB 01430	0.82	1.60 E-03	0.575
Dansyl-549	14.14	371.6321	275.1475	HMDB	L-a-glutamyl-L-Lysine	HMDB 04207	0.81	1.98 E-09	0.581
Dansyl-563	14.49	462.2066	228.1483	Library	Leucyl-Proline	-	0.76	8.33 E-06	0.626
Dansyl-582	14.82	386.1165	304.1164	EML	D-Pantothenoyl-L-cysteine [-H2O]	HMDB 06834	0.83	3.20 E-04	0.548
Dansyl-586	14.89	395.6170	323.1174	EML	Bisdemethoxycurcumin [+NH]	HMDB 02114	0.76	1.36 E-04	0.599
Dansyl-591	14.91	402.0941	168.0358	Library	3-Hydroxymandelic acid	HMDB 00750	0.81	9.72 E-05	0.593
Dansyl-595	14.98	395.1220	322.1273	HMDB	D-Pantothenoyl-L-cysteine	HMDB 06834	0.83	5.29 E-04	0.558
Dansyl-597	15.01	365.1947	131.1364	HMDB	Norspermidine	HMDB 11634	0.83	6.75 E-04	0.540
Dansyl-610	15.05	436.1452	202.0869	EML	4-Aminophenol [+C4H3N3]	HMDB 01169	0.80	5.65 E-04	0.566
Dansyl-620	15.21	365.1500	131.0917	Library	L-leucine	HMDB 00687	0.80	1.76 E-05	0.556
Dansyl-622	15.24	397.1361	163.0778	EML	Aminoadipic acid [+H2]	HMDB 00510	0.83	1.30 E-02	0.548
Dansyl-624	15.30	410.1307	176.0724	EML	Aminoadipic acid [+NH]	HMDB 00510	0.83	2.40 E-02	0.543
Dansyl-663	15.96	462.1701	228.1117	HMDB	Prolylhydroxyproline	HMDB 06695	0.70	6.40 E-04	0.610
Dansyl-666	16.04	349.1580	115.0997	EML	Isovalerylglycine [-CO2]	HMDB 00678	0.71	1.25 E-02	0.592

Dansyl-673	16.14	612.2140	378.1557	EML	FMNH $[-\text{HPO}_3]$	HMDB011742B	0.75	9.61 E-05	0.622
Dansyl-677	16.18	305.0952	71.0369	HMDB	Acrylamide	HMDB04296	0.73	1.10 E-10	0.634
Dansyl-680	16.26	416.1157	182.0574	HMDB	Hydroxyphenyllactic acid	HMDB00755	0.76	1.51 E-06	0.594
Dansyl-681	16.30	307.0925	146.0683	HMDB	Alanylglycine	HMDB06899	0.72	7.36 E-04	0.580
Dansyl-684	16.37	377.1530	143.0947	EML	Pipecolic acid $[\text{+CH}_2]$	HMDB00070	1.25	1.56 E-02	0.685
Dansyl-703	16.95	496.1894	262.1311	HMDB	L-phenylalanyl-L-proline	HMDB11177	0.80	4.09 E-06	0.640
Dansyl-704	16.95	402.5971	337.0775	HMDB	2,8-Dihydroxyquinoline-beta-D-glucuronide	HMDB11658	0.82	1.64 E-03	0.525
Dansyl-706	17.16	331.1113	97.0530	EML	3-Amino-2-piperidone $[-\text{NH}_3]$	HMDB00323	0.81	5.90 E-03	0.543
Dansyl-738	17.85	355.6374	243.1582	EML	L-isoleucyl-L-proline $[\text{+NH}]$	HMDB11174	0.82	3.54 E-04	0.546
Dansyl-745	17.95	371.1069	137.0486	HMDB	2-Aminobenzoic acid	HMDB01123	1.30	3.52 E-03	0.731
Dansyl-750	18.03	348.6296	229.1425	EML	Gamma-Aminobutyryl-lysine $[-\text{H}_2]$	HMDB01959	0.76	2.29 E-04	0.578
Dansyl-760	18.15	356.6460	245.1754	EML	Gamma-Aminobutyryl-lysine $[\text{+CH}_2]$	HMDB01959	0.74	2.41 E-03	0.613
Dansyl-761	18.16	546.2048	312.1465	Library	Phenylalanylphenylalanine	HMDB13302	1.21	2.07 E-03	0.739
Dansyl-768	18.25	416.1159	182.0576	Library	Homovanillic acid	HMDB00118	1.22	2.57 E-02	0.763
Dansyl-784	18.79	355.6192	243.1217	EML	Porphobilinogen $[\text{+NH}_3]$	HMDB00245	0.65	3.50 E-03	0.571

Dansyl-787	18.88	460.1654	226.1071	HMDB	Carnosine	HMDB 00033	1.30	5.36 E-03	0.743
Dansyl-818	19.37	355.6372	243.1578	EML	Propranolol [-O]	HMDB 01849	0.82	1.19 E-04	0.557
Dansyl-851	20.10	370.1111	136.0528	EML	4-Hydroxybenzaldehyde [+CH2]	HMDB 11718	0.82	9.81 E-03	0.539
Dansyl-869	20.77	394.1804	160.1221	HMDB	Isoputrescine	HMDB 06009	0.69	1.17 E-02	0.589
Dansyl-877	20.99	368.1111	268.1056	EML	N-Acetylvani­lanine [+NH]	HMDB 11716	1.28	1.25 E-02	0.686
Dansyl-879	21.03	379.1684	145.1100	EML	L-Norleucine [+CH2]	HMDB 01645	0.82	1.91 E-03	0.532
Dansyl-883	21.07	370.1099	136.0515	EML	p-Hydroxyphenylacetic acid [-O]	HMDB 00020	0.82	4.84 E-04	0.543
Dansyl-884	21.08	510.2048	276.1465	EML	L-phenylalanyl-L-proline [+CH2]	HMDB 11177	1.34	1.11 E-02	0.649
Dansyl-894	21.32	362.1169	128.0586	HMDB	Dihydrothymine	HMDB 00079	0.79	3.41 E-03	0.538
Dansyl-903	21.54	431.5981	395.0797	EML	Quercetin [+C4H3N3]	HMDB 05794	1.55	8.26 E-03	0.709
Dansyl-908	21.73	402.0857	168.0274	EML	6,8-Dihydroxypurine [+O]	HMDB 01182	0.78	5.29 E-04	0.606
Dansyl-910	21.79	316.0922	164.0677	EML	L-Methionine [+NH]	HMDB 00696	0.83	1.37 E-02	0.520
Dansyl-911	21.80	544.1888	310.1305	EML	Olanzapine [-H2]	HMDB 05012	1.20	4.09 E-03	0.727
Dansyl-916	21.86	392.1644	158.1061	EML	Isoputrescine [-H2]	HMDB 06009	0.61	4.36 E-03	0.581
Dansyl-920	21.88	315.1158	162.1150	EML	Tryptamine [+H2]	HMDB 00303	0.69	7.87 E-03	0.579
Dansyl-921	21.88	314.1187	160.1208	HMDB	N(6)-Methyllysine	HMDB 02038	0.66	6.63 E-03	0.580

Dansyl-956	22.48	402.0863	168.0279	EML	3-Methyluric acid [-CH2]	HMDB 01970	0.75	3.85 E-06	0.629
Dansyl-965	22.60	403.6250	339.1333	EML	5-Hydroxytryptophol [+C6H10O5]	HMDB 01855	1.28	1.99 E-02	0.705
Dansyl-969	22.75	336.1211	204.1255	EML	Serotonin [+C2H4]	HMDB 00259	0.79	1.62 E-03	0.546
Dansyl-975	22.83	318.0724	168.0281	EML	1-Methyluric acid [-CH2]	HMDB 03099	0.74	9.81 E-10	0.663
Dansyl-991	23.29	325.0929	91.0346	EML	4-Hydroxy tolbutamide [+C2H4]	HMDB 06408	0.82	3.52 E-04	0.582
Dansyl-995	23.32	324.5927	181.0688	Library	L-Tyrosine	HMDB 00158	0.82	2.69 E-04	0.559
Dansyl-1003	23.37	401.0806	167.0222	HMDB	Homocysteinesulfinic acid	HMDB 06462	0.69	1.98 E-09	0.685
Dansyl-1007	23.38	337.1164	103.0580	EML	Ethanolamine [+C2H2O]	HMDB 00149	0.74	5.42 E-08	0.648
Dansyl-1008	23.40	319.0691	170.0216	HMDB	Gallic acid	HMDB 05807	0.75	1.78 E-07	0.606
Dansyl-1011	23.47	320.0999	172.0832	EML	Diaminopimelic acid [-H2O]	HMDB 01370	0.81	3.38 E-02	0.557
Dansyl-1035	23.99	359.6162	251.1158	EML	N-Acetyl-L-tyrosine [+C2H4]	HMDB 00866	0.78	3.93 E-05	0.591
Dansyl-1050	24.23	363.1734	129.1150	EML	2-Heptanone [+NH]	HMDB 03671	0.79	2.98 E-07	0.636
Dansyl-1051	24.25	284.1082	100.0998	EML	Cadaverine [-H2]	HMDB 02322	0.83	4.43 E-02	0.528
Dansyl-1057	24.43	342.1318	216.1470	EML	Gamma-Aminobutyryl-lysine [-NH]	HMDB 01959	0.75	4.34 E-02	0.547
Dansyl-1067	24.93	323.1024	178.0882	EML	5-Hydroxylysine [+O]	HMDB 00450	1.30	1.02 E-03	0.740
Dansyl-1068	24.93	322.1056	176.0945	HMDB	Canavanine	HMDB 02706	1.31	1.60 E-03	0.734

Dansyl-1069	24.95	318.0796	168.0426	HMDB	Homogentisic acid	HMDB 00130	0.67	1.85 E-05	0.597
Dansyl-1089	25.27	323.6064	179.0962	HMDB	2(N)-Methyl-norsalsolinol	HMDB 01189	1.36	1.50 E-03	0.726
Dansyl-1150	26.44	338.5839	209.0511	EML	1-Methylguanine [+CO2]	HMDB 03282	0.79	4.32 E-02	0.508
Dansyl-1180	26.93	432.6178	397.1189	EML	2'-Deoxysepiapterin [+C6H8O6]	HMDB 00389	0.65	4.24 E-02	0.516
Dansyl-1187	26.99	577.4068	343.3485	EML	Homophytanic acid [+NH3]	HMDB 02337	1.28	4.88 E-02	0.655
Dansyl-1188	27.00	327.1272	186.1378	EML	1-(3-Aminopropyl)-4-aminobutanal [+C2H2O]	HMDB 12135	0.77	4.99 E-02	0.591
Dansyl-1204	27.17	309.5878	151.0590	HMDB	Acetaminophen	HMDB 01859	0.61	2.20 E-03	0.571
Dansyl-1221	27.38	535.2666	301.2082	EML	Androstenedione [+NH]	HMDB 00053	0.82	1.23 E-02	0.513
Dansyl-1248	27.64	526.2275	292.1692	EML	Gingerol [-H2]	HMDB 05783	1.79	4.51 E-02	0.616
Dansyl-1282	28.40	377.6158	287.1150	HMDB	N-Ribosylhistidine	HMDB 02089	0.72	1.59 E-04	0.577
Dansyl-1288	28.50	496.1798	524.2429	HMDB	Phe Cys Gln Lys	-	0.83	5.91 E-05	0.531
Dansyl-1300	29.04	474.1549	240.0966	EML	Vanillic acid [+C2H4]	HMDB 00913	0.75	2.15 E-02	0.574
Dansyl-1302	29.07	400.0708	166.0125	EML	Uric acid [-H2]	HMDB 00289	0.79	7.84 E-03	0.541
Dansyl-1316	29.38	493.1362	259.0778	EML	Kinetin [+CO2]	HMDB 12245	0.51	2.35 E-02	0.501
Dansyl-1333	30.62	533.3398	299.2814	HMDB	Sphingosine	HMDB 00252	0.69	1.25 E-02	0.532
DMPA-14	6.33	502.1087	178.9323	EML	S-Carboxymethyl-L-cysteine	HMDB 29415	1.28	8.30 E-04	0.607

DMPA-44	7.26	896.3784	573.2019	EML	Aminopterin [+C5H3N5]	HMDB 01833	1.27	8.30 E-04	0.597
DMPA-45	7.30	488.1691	164.9926	EML	3-Sulfinylpyruvic acid [+CH2]	HMDB 01405	1.29	5.50 E-04	0.617
DMPA-55	7.58	504.2513	181.0748	HMDB	4-Hydroxy-4-(3-pyridyl)-butanoic acid	HMDB 01119	1.26	4.74 E-02	0.582
DMPA-92	9.18	491.2025	168.0261	EML	Butyric acid [+HPO3]	HMDB 00039	1.31	3.28 E-02	0.629
DMPA-98	9.26	344.1490	182.0571	Library	Hydroxyphenyllactic acid	HMDB 00755	0.78	8.61 E-03	0.733
DMPA-128	10.43	353.1712	191.0793	EML	Glutaryl-glycine [+H2]	HMDB 00590	1.21	1.73 E-02	0.602
DMPA-142	10.76	280.1538	118.0619	Library	2-Hydroxy-2-methylbutyric acid	HMDB 01987	0.81	4.47 E-02	0.690
DMPA-217	12.52	294.1696	132.0777	Library	2-Hydroxycaproic acid	HMDB 01624	0.83	3.99 E-02	0.723
DMPA-222	12.56	601.3375	278.1610	EML	Aspartylsine [+NH3]	HMDB 04985	1.24	4.81 E-02	0.579
DMPA-231	12.77	696.3856	373.2091	EML	6-Keto-decanoylcarnitine [+CO2]	HMDB 13202	1.21	4.27 E-02	0.569
DMPA-270	14.30	489.3312	327.2393	EML	13-L-Hydroperoxylinoleic acid [+NH]	HMDB 03871	0.83	4.27 E-02	0.688
DMPA-363	16.91	503.3471	341.2552	HMDB	trans-2-Dodecenoylcarnitine	HMDB 13326	0.82	2.63 E-02	0.760
DMPA-381	17.44	466.1582	304.0663	EML	Biotin sulfone [+CO]	HMDB 04818	1.23	3.21 E-02	0.597
DMPA-462	22.06	570.3778	408.2859	Library	Cholic acid	HMDB 00619	0.53	1.33 E-02	0.613
DMPA-510	24.78	421.1580	259.0661	EML	Phthalic acid [+C4H3N3]	HMDB 02107	1.23	4.74 E-02	0.568
DMPA-539	26.19	461.3726	299.2807	EML	Oleic acid [+NH3]	HMDB 00207	1.20	7.40 E-03	0.599

DMPA-599	28.57	489.4041	327.3122	EML	Phytanic acid [+NH]	HMDB 00801	1.25	2.73 E-03	0.616
DMPA-652	29.67	530.2788	368.1869	EML	Epsilon-(gamma-Glutamyl)-lysine [+C4H3N3]	HMDB 03869	1.21	2.40 E-03	0.593
DMPA-708	30.15	371.1242	209.0323	EML	5,6-Dihydroxyindole-2-carboxylic acid [+O]	HMDB 01253	1.24	4.74 E-02	0.585
DMPA-909	32.90	554.2848	392.1929	EML	13-L-Hydroperoxylinoleic acid [+SO3]	HMDB 03871	1.23	2.94 E-02	0.629
DMPA-928	33.65	381.1676	219.0757	EML	L-beta-aspartyl-L-threonine [-NH]	HMDB 11169	1.33	2.63 E-02	0.626
DMPA-947	34.38	428.3170	266.2251	EML	7,10-Hexadecadienoic acid [+CH2]	-	1.20	4.47 E-02	0.567
DMPA-959	34.65	442.3318	280.2399	Library	Linoleic acid	HMDB 00673	1.21	3.90 E-03	0.586
DMPA-968	34.88	416.3122	254.2203	HMDB	Hypogeic acid	HMDB 02186	1.50	2.69 E-02	0.624
DMPA-988	35.56	518.3640	356.2721	HMDB	Tetracosahexaenoic acid	HMDB 02007	1.28	2.94 E-02	0.566
DMPA-995	35.72	397.1442	235.0523	EML	5-Hydroxy-N-formylkynurenine [-NH3]	HMDB 04086	1.27	6.69 E-03	0.634
DMPA-1004	35.85	464.3128	302.2209	HMDB	Eicosapentaenoic acid	HMDB 01999	1.29	3.18 E-02	0.588
DMPA-1019	36.47	430.3329	268.2410	EML	Elaidic acid [-CH2]	HMDB 00573	1.24	4.74 E-02	0.578
DMPA-1033	37.31	456.3485	294.2566	EML	Linoleic acid [+CH2]	HMDB 00673	1.21	4.35 E-02	0.583
DMPA-1060	38.51	466.3194	304.2275	EML	Putreanine [+C7H13NO2]	HMDB 06078	1.31	1.71 E-02	0.602

**Table 3.3** 25 identified dansyl-labeled metabolites and 6 identified DMPA-labeled metabolites that have fold change (overweight/normal) > 1.2 (or < 0.83) and q value < 0.05 for the difference between the normal group and the overweight group. (Asterisk means the metabolite is also a significant metabolite for sex differences.)

Label	Retention time (min)	Detected m/z	Accurate mass (Da)	ID level	Compound Name	HMDB ID	Fold Change	q-value	Univariate ROC
Dansyl-15	2.66	521.0892	287.0308	EML	N-Acetylgalactosamine 6-sulfate [-CH <sub>2</sub> ]	-	1.27	1.74 E-02	0.599
Dansyl-54	4.61	359.0724	125.0141	Library	Taurine	HMDB 00251	0.81	2.03 E-02	0.642
Dansyl-62*	4.78	531.1029	297.0446	HMDB	L-Cysteinylglycine disulfide	HMDB 00709	1.24	2.56 E-04	0.639
Dansyl-281	8.18	542.1348	308.0765	EML	Prolylhydroxyproline [+HPO <sub>3</sub> ]	HMDB 06695	0.82	4.53 E-02	0.643
Dansyl-285	8.23	381.1423	147.0840	EML	5-Hydroxyhexanoic acid [+NH]	HMDB 00525	1.24	3.64 E-02	0.573
Dansyl-294*	8.31	452.1852	218.1269	EML	Diaminopimelic acid [+C <sub>2</sub> H <sub>4</sub> ]	HMDB 01370	0.78	4.53 E-02	0.585
Dansyl-295*	8.35	418.0816	184.0233	EML	Cysteic acid [+NH]	HMDB 02757	1.23	1.64 E-02	0.538
Dansyl-327	9.09	396.1357	324.1547	HMDB	Galactosylhydroxylysine	HMDB 00600	0.82	4.03 E-02	0.672
Dansyl-349*	9.33	502.1763	268.1180	EML	Carnosine [+C <sub>2</sub> H <sub>2</sub> O]	HMDB 00033	1.36	3.82 E-02	0.616
Dansyl-353*	9.41	399.1011	165.0428	EML	2-Aminobenzoic acid [+CO]	HMDB 01123	0.70	3.88 E-02	0.645
Dansyl-375*	10.04	402.0861	168.0278	EML	Oxypurinol [+O]	HMDB 00786	1.25	4.23 E-02	0.546
Dansyl-420	11.44	436.1900	202.1317	EML	Glycyl-L-leucine [+CH <sub>2</sub> ]	HMDB 00759	0.78	1.80 E-02	0.677
Dansyl-610*	15.05	436.1452	202.0869	EML	4-Aminophenol [+C <sub>4</sub> H <sub>3</sub> N <sub>3</sub> ]	HMDB 01169	1.24	1.75 E-02	0.582

Dansyl-720	17.49	478.2361	244.1778	EML	N(6)-(Octanoyl)lysine [- C2H4]	HMDB 11684	0.69	6.59 E-03	0.716
Dansyl-861	20.49	345.1269	111.0686	EML	4-Aminophenol [+H2]	HMDB 01169	1.69	4.24 E-02	0.575
Dansyl-863	20.61	405.5682	343.0198	EML	Epinephrine sulfate [+HPO3]	HMDB 01876	1.33	1.22 E-02	0.626
Dansyl-872	20.87	476.2206	242.1622	EML	L-isoleucyl-L-proline [+CH2]	HMDB 11174	0.78	3.88 E-02	0.689
Dansyl-888	21.18	324.5772	181.0377	EML	3-Hydroxyanthranilic acid [+CO]	HMDB 01476	0.64	1.80 E-02	0.693
Dansyl-933	22.13	575.2056	341.1473	EML	Tyramine glucuronide [+C2H4]	HMDB 10328	0.66	1.22 E-02	0.644
Dansyl- 1003*	23.37	401.0806	167.0222	HMDB	Homocysteinesulfinic acid	HMDB 06462	1.28	1.22 E-02	0.571
Dansyl- 1070	24.96	453.1077	438.0987	EML	2-Phenylaminoadenosine [+H2O]	HMDB 01069	1.47	2.35 E-02	0.612
Dansyl- 1150*	26.44	338.5839	209.0511	EML	1-Methylguanine [+CO2]	HMDB 03282	1.49	2.35 E-02	0.601
Dansyl- 1173	26.79	450.5982	433.0797	EML	Se- Adenosylselenomethioni ne [-CH2]	HMDB 11118	0.69	1.22 E-02	0.678
Dansyl- 1204*	27.17	309.5878	151.0590	HMDB	Acetaminophen	HMDB 01859	1.65	3.49 E-02	0.594
Dansyl- 1316*	29.38	493.1362	259.0778	EML	Kinetin [+CO2]	HMDB 12245	2.28	4.23 E-02	0.551
DMPA-89	9.13	435.1276	273.0357	EML	Phenylacetyl glycine [+SO3]	HMDB 00821	0.73	2.37 E-02	0.577
DMPA-121	10.29	450.1455	288.0536	HMDB	Orotidine	HMDB 00788	0.67	4.60 E-02	0.617
DMPA-662	29.79	483.1086	321.0167	EML	DOPA sulfate [+CO2]	HMDB 02028	1.41	2.37 E-02	0.658
DMPA-787	30.75	418.3373	256.2454	HMDB	Palmitic acid	HMDB 00220	0.74	2.95 E-02	0.599

DMPA-817	30.98	418.3345	256.2426	EML	Pentadecanoic acid [+CH2]	HMDB 00826	0.69	4.22 E-02	0.599
DMPA-861	31.83	545.3392	383.2473	EML	Linoleic acid [+C3H5NOS]	HMDB 00673	0.69	2.37 E-02	0.573

**Table 3.4** 25 identified dansyl-labeled metabolites and one identified DMPA-labeled metabolite that have fold change (underweight/normal) > 1.2 (or < 0.83) and q value < 0.05 for the difference between the normal group and the underweight group.

Label	Retention time (min)	Detected m/z	Accurate mass (Da)	ID level	Compound Name	HMDB ID	Fold Change	q-value	Uni-variate ROC
Dansyl-38	4.06	388.1072	154.0489	EML	1,3-Diaminopropane [+SO3]	HMDB 00002	3.23	1.25 E-02	0.616
Dansyl-118	5.94	321.0907	87.0324	HMDB	2-Aminoacrylic acid	HMDB 03609	1.74	4.61 E-02	0.501
Dansyl-154	6.55	399.0987	165.0404	Library	Methionine Sulfoxide - Isomer	-	1.46	4.06 E-02	0.531
Dansyl-190	7.05	466.1648	232.1065	HMDB	4-(Glutamylamino) butanoate	HMDB 12161	0.66	4.94 E-02	0.528
Dansyl-308	8.71	351.1005	117.0422	HMDB	L-2-Amino-3- oxobutanoic acid	HMDB 06454	4.07	4.52 E-03	0.571
Dansyl-348	9.32	279.0794	45.0211	HMDB	Formamide	HMDB 01536	1.23	4.06 E-02	0.552
Dansyl-386	10.35	528.1781	294.1198	HMDB	Glutamylphenylalanine	HMDB 00594	1.75	4.61 E-02	0.532
Dansyl-412	11.32	370.0973	136.0390	Library	Hypoxanthine - Isomer	-	2.98	1.25 E-02	0.571
Dansyl-458	12.15	454.1447	220.0864	EML	Canavanine [+CO2]	HMDB 02706	1.32	4.06 E-02	0.563
Dansyl-498	13.12	323.1375	89.0791	EML	Ethanolamine [+C3H4]	HMDB 00149	3.35	1.36 E-02	0.657

Dansyl -522	13.42	340.0672	106.0089	EML	3-Methylthiopropionic acid [-CH2]	HMDB 01527	1.66	2.34 E-02	0.513
Dansyl -560	14.37	363.1015	258.0864	HMDB	O-Desmethylangolensin	HMDB 04629	1.88	4.06 E-02	0.532
Dansyl -654	15.85	387.1007	153.0424	HMDB	3-Hydroxyanthranilic acid	HMDB 01476	3.93	1.25 E-02	0.595
Dansyl -725	17.60	585.2147	351.1564	Library	Tryptophyl- Phenylalanine	-	0.77	3.78 E-02	0.626
Dansyl -821	19.47	441.1472	207.0889	HMDB	3- Phenylpropionylglycine	HMDB 02042	1.23	4.61 E-02	0.634
Dansyl -985	23.21	311.0825	154.0484	EML	1,3-Diaminopropane [+SO3]	HMDB 00002	3.71	1.25 E-02	0.557
Dansyl -986	23.21	388.1072	154.0489	EML	6,8-Dihydropyrimidine [+H2]	HMDB 01182	3.12	1.25 E-02	0.614
Dansyl -1009	23.44	297.0851	126.0536	HMDB	2,4-Diamino-6- hydroxypyrimidine	HMDB 02128	4.21	2.34 E-02	0.525
Dansyl -1013	23.50	302.0774	136.0381	HMDB	Hypoxanthine	HMDB 00157	4.05	1.25 E-02	0.590
Dansyl -1073	25.02	711.1806	244.0713	HMDB	3,3',4'5- Tetrahydroxystilbene	HMDB 04215	1.30	4.69 E-02	0.542
Dansyl -1076	25.07	512.2095	278.1512	HMDB	Alpha-CEHC	HMDB 01518	2.83	2.34 E-02	0.524
Dansyl -1132	26.16	303.0751	138.0336	HMDB	Gentisate aldehyde	HMDB 04062	4.05	2.66 E-02	0.578
Dansyl -1149	26.43	302.0782	136.0399	EML	Purine [+O]	HMDB 01366	4.26	1.25 E-02	0.646
Dansyl -1178	26.90	648.1830	181.0736	HMDB	Beta-Tyrosine	HMDB 03831	1.31	4.06 E-02	0.543
Dansyl -1260	27.87	467.0908	233.0325	HMDB	Dopamine 4-sulfate	HMDB 04148	1.39	4.06 E-02	0.578
DMPA -415	19.17	611.4051	449.3132	HMDB	Glycoursodeoxycholic acid	HMDB 00708	1.88	1.64 E-02	0.607

## Chapter 4

### High-coverage Metabolomics Analysis of One Microliter of Blood Using Chemical Isotope Labeling and High-resolution LC-MS

#### 4.1 Introduction

Metabolomics, the comprehensive study of the complete set of small-molecule compounds within a biological sample, has become a powerful tool in biomarker discovery,<sup>35</sup> clinical diagnosis,<sup>37</sup> nutritional studies,<sup>181</sup> and toxicology applications.<sup>182</sup> In recent years, it has been increasingly accepted that environmental exposures, as well as gene-environment interactions, contribute to the development of many diseases.<sup>183-185</sup> When the body is under an environmental exposure, certain direct and indirect metabolic changes may occur, underlying the more complicated biological changes at the proteomics and genomics levels. Because metabolomics can monitor all metabolic variations caused by an environmental stimulation, it is also a promising technique for understanding the gene-environment-health relationship.<sup>186</sup>

Blood is a primary carrier of small molecules in the body, and according to the Human Serum Metabolome Database,<sup>31</sup> there are more than 4,600 already known small-molecule metabolites in the blood. As an important and easily accessible biological fluid, blood has been used in clinical tests for many years. And in recent years, blood metabolome analysis

has been used to study the metabolic changes caused by various environmental factors including diet,<sup>187</sup> smoking,<sup>188</sup> gut microflora,<sup>66</sup> medication<sup>189</sup> and air pollution.<sup>190</sup>

For metabolomics analyses, it is crucial to determine if the over- or under-expression of specific metabolites in blood is a true response to internal or external stimulations, and not due to the normal time-of-day or day-to-day variations. Using a non-targeted LC-MS approach to study 1,069 metabolites in human plasma, Ang et al. found that at least 19% of the metabolite features in their study exhibited significant 24-hour variations.<sup>115</sup> Therefore, if a blood sample is only collected once from each subject, the measured concentration of a specific metabolite may not be able to represent the average concentration of the metabolite within the day of sample collection. And if this metabolite is chosen as a biomarker candidate, the biological implications according to its concentration will become less reliable. This calls for multiple blood samples collected from the same subject at different time points to more accurately and comprehensively represent the blood metabolome. Furthermore, time-resolved metabolomics analysis is an important tool for the assessment of an environmental exposure and the resulting biological changes. For instance, Chourel et al. collected blood samples every 15 min from the participants during 90 min of ergometer-cycling to study the metabolic responses.<sup>191</sup> However, many metabolomics studies only have one blood sample from each subject under each experimental condition. This is partially because the commonly used venipuncture blood collection is very invasive and a trained phlebotomist is required.

Efforts have been made to find easier and less invasive alternatives to venipuncture blood collection. For example, starting from the 1960s, the dried blood spot (DBS) method has been used for clinical purposes, especially for newborn screening.<sup>84</sup> In addition to its advantages in sample collection and storage, DBS is believed to be less invasive than venipuncture, as the sample collection is done by a finger or heel prick. In recent years, other paper-based devices that perform both blood cell separation and targeted analysis of specific biomarkers have emerged, aiming to provide a low-cost and point-of care diagnosis platform.<sup>192-193</sup>

According to our knowledge, there hasn't been any non-targeted metabolic profiling investigation based on finger blood, as when the sample amount is too small, the metabolome coverage and detectability of low-abundant metabolites are usually limited. In this chapter, we apply our CIL-LC-MS method to the metabolic profiling of one microliter of finger blood. The blood sample is easily collected by a commercially available finger pricking device which is designed for diabetes patients to monitor the blood glucose level. With this technique, we achieve the high-coverage blood metabolome analysis of small sample amount, enabling time-resolved metabolomics analysis for various future studies.

## **4.2 Materials and methods**

### **4.2.1 Chemicals and reagents**

All the chemicals and reagents, unless otherwise stated, were purchased from Sigma-Aldrich Canada (Markham, ON, Canada). For dansylation labeling reaction, the  $^{12}\text{C}$ -labeling reagent (dansyl chloride) was from Sigma-Aldrich and the  $^{13}\text{C}$ -labeling reagent was synthesized in our lab using the procedure published previously.<sup>70</sup> LC-MS grade water, methanol, and acetonitrile (ACN) were purchased from Thermo Fisher Scientific (Nepean, ON, Canada).

### **4.2.2 Universal serum standard**

For making the universal serum standard (USS) sample, 10 mL blood samples were collected from 100 healthy volunteers in Edmonton, Canada. The participants are between the ages of 20 and 39 years old, and they were refrained from eating or drinking (except water) for at least 8 hours before giving blood. The University of Alberta health ethics review board approved this study, and all participants provided informed consent. The collected whole blood was allowed to clot by leaving it undisturbed at room temperature for 60 min. After clotting, the blood sample was centrifuged at 1,000 rpm for 15 min. An equal volume of supernatant was taken from each sample, and these aliquots were mixed to be the USS sample. The USS sample was finally stored in 1.5 mL micro-centrifuge tubes in  $-80\text{ }^{\circ}\text{C}$  freezer.

#### **4.2.3 Dansylation labeling of serum sample**

In a 0.2 mL PCR tube, 1  $\mu\text{L}$  of USS sample was mixed with 55  $\mu\text{L}$  of methanol to precipitate the proteins. After being incubated at  $-20\text{ }^{\circ}\text{C}$  for 1 hour, the mixture was centrifuged at 14,000 rpm for 15 min. 45  $\mu\text{L}$  supernatant was taken and totally dried using a Speed Vac. The sample was re-dissolved with 7.5  $\mu\text{L}$  of 85 mM sodium carbonate/sodium bicarbonate buffer, which makes a basic environment for the dansylation reaction. Then 7.5  $\mu\text{L}$  of freshly prepared  $^{12}\text{C}$ -dansyl chloride solution (20mg/mL) (for light labeling) or  $^{13}\text{C}$ -dansyl chloride solution (20mg/mL) (for heavy labeling) was added to the tube. The mixture was vortexed, spun down and then incubated at  $40\text{ }^{\circ}\text{C}$  for 45 min. To quench the excess dansyl chloride, 2  $\mu\text{L}$  of 250 mM NaOH solution was mixed with the reaction mixture before another 10 min incubation at  $40\text{ }^{\circ}\text{C}$ . Finally, 5  $\mu\text{L}$  of formic acid (425 mM) in 1:1 ACN/ $\text{H}_2\text{O}$  was added to consume excess NaOH and to make the solution acidic. After centrifuging at 14,000 rpm for 15 min, 15  $\mu\text{L}$  supernatants were taken from the  $^{12}\text{C}$ -labeled and  $^{13}\text{C}$ -labeled samples and mixed together for the following LC-MS analysis.

#### **4.2.4 Finger blood sample collection and processing**

The finger surface was cleaned with 70% isopropyl alcohol. Then the finger pricking was done by Bayer's Microlet 2 lancing device. When a blood drop formed, 1  $\mu\text{L}$  of whole blood sample was taken from the punching site with a micropipette, and immediately transferred to a 0.2 mL PCR tube which had 10  $\mu\text{L}$  PBS solution in it. The tube was vortexed gently to let the blood thoroughly mixed with the PBS solution, and then stayed at  $4\text{ }^{\circ}\text{C}$  for 0.5 hours. After that, the mixture was centrifuged at 1,000 rpm for 10 min and

8  $\mu\text{L}$  of supernatant was taken into another PCR tube. Finally, the supernatant was mixed with 48  $\mu\text{L}$  of methanol to precipitate the proteins. After centrifuging, 45  $\mu\text{L}$  of the supernatant was transferred to a new tube and dried. The labeling process is the same as the labeling of serum described above. Before the labeling reaction, the sample can also be stored in a  $-80\text{ }^{\circ}\text{C}$  freezer for future studies.

#### **4.2.5 Sample quantification and normalization**

For the pre-acquisition normalization of blood samples, we used an LC-UV based sample quantification method, which was previously developed in our lab.<sup>115</sup> With a Waters ACQUITY UPLC system UPLC (Waters, Milford, MA, USA), two microliters of each labeled sample was injected into a Phenomenex Kinetex C18 column (2.1 mm  $\times$  5 cm, 1.7  $\mu\text{m}$  particle size) (Phenomenex, Torrance, CA, USA) for a fast step-gradient run. Solvent A was 0.1% (v/v) formic acid in 5% (v/v) ACN/ $\text{H}_2\text{O}$ , and solvent B was 0.1% (v/v) formic acid in ACN. The gradient started with 0% B for 1 min and was increased to 95% B within 0.01 min and held at 95% B for 1 min to ensure complete elution of all labeled metabolites. The flow rate used was 0.45 mL/min, and the total UV absorption of dansyl-labeled metabolites in the sample was measured by a photodiode array (PDA) detector. For the post-acquisition normalization, the individual samples were normalized according to the total intensity of dansyl-labeled peaks.

#### 4.2.6 LC-FTICR-MS

An Agilent 1100 series binary system (Agilent, Palo Alto, CA, USA) and an Agilent reversed-phase Eclipse plus C18 column (2.1 mm×100 mm, 1.8 μm particle size, 95 Å pore size) were used for LC-MS. LC Solvent A was 0.1% (v/v) formic acid in 5% (v/v) ACN/H<sub>2</sub>O, and Solvent B was 0.1% (v/v) formic acid in ACN. The gradient elution profile was as follows: 0 min (20% B), 0-3.5 min (20-35% B), 3.5-18 min (35-65% B), 18-24 min (65-99% B), and 24-32 min (99% B). The flow rate was 180 μL/min, and the flow was loaded to the electrospray ionization (ESI) source of a Bruker 9.4 Tesla Apex-Qe Fourier transform ion-cyclotron resonance (FTICR) mass spectrometer (Bruker, Billerica, MA, USA). All MS spectra were obtained in the positive ion mode.

#### 4.2.7 Nano-LC-QTOF-MS

The nano-LC–MS experiments were performed on a Waters nanoAcquity connected to a Bruker Impact HD Quadrupole Time-of-flight (QTOF) mass spectrometer equipped with a nano ESI-source. LC separations were performed on an Acclaim PepMap RSLC C18 (75 μm × 150 mm, 2 μm) and Acclaim PepMap 100 trap column (75 μm × 20 mm, 3 μm) (Thermo Scientific, Sunnyvale, CA, USA). LC Solvent A was 0.1% (v/v) formic acid in water, and Solvent B was 0.1% (v/v) formic acid in acetonitrile. The 55 min gradient conditions were as follows: 0 min (15% B), 0-2.0 min (15% B), 2-4min (15-25% B), 4-34 min (15-60% B), 34-40 min (60-90% B), and 40-55 min (90 %B). The flow rate was 350 nL/min, and the injection volume was 5 μL. Prior to separating on the analytical column,

the sample is first pushed through the trap column by 98% mobile phase A at 7.0  $\mu\text{L}/\text{min}$  for 1.5 min. All MS spectra were obtained in the positive ion mode.

#### **4.2.8 Data processing and metabolite identification**

The  $^{12}\text{C}/^{13}\text{C}$ -peak pairs from each LC-MS run were extracted by the IsoMS software.<sup>72</sup> IsoMS-Align was used to align the peak pair data from different samples by retention time and accurate mass. The missing ratio values were filled back by the Zero-fill program.<sup>137</sup> IsoMS-Quant<sup>157</sup> was used to generate the final metabolite-intensity table, which was exported to SIMCA-P+12 (Umetrics AB, Umeå, Sweden) for statistical analysis. Metabolite identification was performed by matching mass and retention time to a dansyl standard library using DnsID.<sup>159</sup> Putative identification was done based on accurate mass matches to the metabolites in the human metabolome database (HMDB) (8,021 known human endogenous metabolites) and in the Evidence-based Metabolome Library (EML) (375,809 predicted human metabolites with one reaction) using MyCompoundID.<sup>42</sup> The mass accuracy tolerance window was set at 0.008 Da for database search.

### **4.3 Results and Discussion**

#### **4.3.1 Finger blood collection and sample preparation**

The starting point of our work is to determine the amount of blood we can collect each time. There is a common perception that the lower the collection volume is, the less invasive it is, and studies have shown that the finger sticking devices that give a larger sample volume usually create more pain.<sup>194-195</sup> Although many blood analyzing methods

based on finger blood have been developed, the sample collection amounts of those techniques have not been low enough to achieve a painless blood collection. For example, the typical sample collection amount for DBS is between 25  $\mu\text{L}$  and 100  $\mu\text{L}$ ,<sup>85</sup> and many of the fresh blood analyzing devices require at least 10  $\mu\text{L}$  of raw finger blood.<sup>193, 196</sup>

In our current work, we use a finger sticking device originally made for self-monitoring of blood glucose level. Various lancing devices have been made for the blood glucose meters, and the manufactures have put a lot of efforts into minimizing the required sample volume, which is usually one microliter nowadays. According to the result of Grady et al.'s study, the mean blood collection volume from 64 diabetes patients with their own lancing devices was 3.1  $\mu\text{L}$ , and 89% of the patients agreed that getting enough blood to fill a 1  $\mu\text{L}$ -testing strip was not painful.<sup>197</sup> The finger sticking device in our work creates blood drops of 2  $\mu\text{L}$  to 10  $\mu\text{L}$ , varying among individuals and sites of finger pricking. To avoid occasional difficulties in having enough sample, we have decided to collect a relatively small amount, one microliter of whole blood, for our metabolomics analysis. With this collection volume, the blood donors didn't experience any significant pain during serial blood collection. One major advantage of our dansyl-labeling method is that the MS signals of metabolites are improved by 10 to 1,000 fold, enabling high-coverage metabolomics analysis with a relatively small sample volume. Although the available blood volume is only one microliter, compared with other finger blood analyzing methods, our method can detect a larger number of metabolites.

After the finger pricking is done and a blood drop is formed, one microliter of finger blood is immediately collected by a disposable glass capillary micropipette and transferred into a 0.2 mL-PCR tube. For separating the blood cells, we prefer centrifugation than filtration, as the filter may absorb some metabolites and the resulting sample loss may greatly affect the detection of low-abundant metabolites. In addition to a laboratory centrifuge, which is the most convenient choice in a research lab, many point-of-care blood centrifuging devices have been reported.<sup>198-200</sup> For example, Haeberle et al. designed an on-disk device that can separate 2  $\mu$ L of plasma from 5  $\mu$ L of whole blood. We use a traditional centrifuge in our work, but in the future, we can switch to a smaller device and the sample collection can be done outside the lab. Also, we dilute the fresh blood with 10  $\mu$ L of PBS solution since the dilution makes it easier to separate the supernatant layer from the blood cells on the bottom. The processed sample can be considered as a PBS-diluted plasma, and can be long-termly stored in a -80 °C freezer for future analysis.

Since adding PBS makes the sample matrix more complicated, we note that it is very important to evaluate the matrix effects associated with the sample collection process. Blood metabolome analyses and further inter-study comparisons have long been accompanied by the concerns about matrix effects, as different types of blood samples (i.e., serum and plasma prepared using different anticoagulants) with vastly different matrix compositions are being used for different metabolomics analyses, and there is no universal standard for the sample collection protocol. To overcome the matrix effects, we have examined the results produced by dansyl-labeling LC-MS analysis of serum and plasma samples with different anticoagulants (EDTA, heparin and citrate). Among the four types

of samples, there is no significant difference in metabolite detectability and relative quantification precision.<sup>82</sup> We have also studied the matrix effect of PBS solution. The result showed that high concentrations of NaCl and phosphate buffer (>50 mM) or PBS could reduce or enhance the labeling efficiencies of metabolites. Nonetheless, by maintaining similar matrix contents in all the individual samples and the reference sample, accurate relative quantification of metabolites can be performed.<sup>161</sup> In our work, the concentrations of salts in the 1X PBS are not high enough to cause severe matrix effects, and more importantly, all of the collected blood samples have the same matrix, so the matrix effect will not significantly affect the analysis results.

#### **4.3.2 Relative quantification and internal standard**

For the metabolome analysis of a large set of samples using our CIL-LC-MS method, a <sup>13</sup>C-labeled sample is needed as the internal reference. As each blood sample is <sup>12</sup>C-labeled and mixed with the internal reference, the LC-MS analysis can generate the <sup>12</sup>C/<sup>13</sup>C-peak pairs of all the labeled metabolites. To achieve the relative quantification of these metabolites, the intensity ratios of the peak pairs are calculated. Each peak pair ratio represents the relative concentration of a specific metabolite in an individual sample, with respect to the internal reference. Since all the individual samples are mixed with the same internal reference, any change in the concentrations of metabolites can be easily detected.

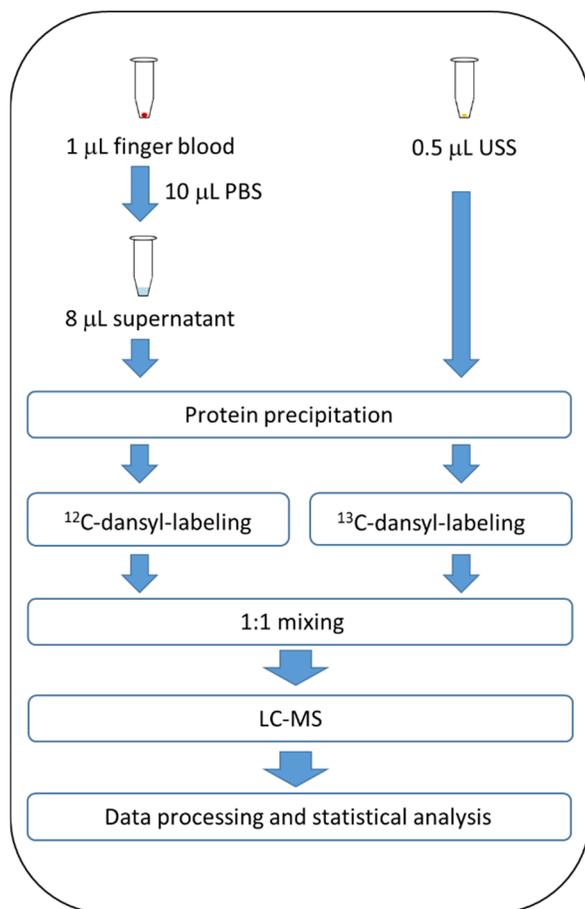
A pooled sample, which is usually made by mixing equal aliquots of all the individual samples, can be used as the internal reference. The pooled sample is easy to obtain, but it has some disadvantages. First, this method takes an aliquot from each sample to generate

the pooled sample, which may not be feasible if the available sample amount is limited. Only a very small volume of finger blood can be collected each time and the sample should be processed immediately after collection, so making a pooled sample can be very difficult. Moreover, when each metabolomics study has its own internal reference, results of different studies of the same sample type are not comparable. For a large-scale metabolomics study in which multiple sample sets are analyzed at different times or even different laboratories, a universal standard is preferred.

We have collected serum samples from 100 healthy blood donors and mixed the 100 samples to make a universal serum standard (USS). All the  $^{12}\text{C}$ -labeled individual samples are mixed with the  $^{13}\text{C}$ -labeled USS, so the measured metabolite concentrations are relative to the average level among the general population, which is represented by the USS sample. The isotopic labeled USS can be mass-produced in advance by making a large-scale reaction, and then stored for any future studies. This makes the experimental process easier than the pooled sample approach. Using the same reference sample, blood analysis results from different subjects and different time points can be readily compared.

Figure 4.1 summarizes the workflow of our metabolomics profiling of finger blood. After the collection of one microliter of whole blood, we perform protein precipitation using an organic solvent, and then the dried sample is labeled by  $^{12}\text{C}$ -dansyl chloride. Meanwhile, 0.5  $\mu\text{L}$  of the USS sample is labeled by  $^{13}\text{C}$ -dansyl chloride. We use 0.5  $\mu\text{L}$  of the USS sample because the total metabolite amounts in 0.5  $\mu\text{L}$  of serum and 1  $\mu\text{L}$  of whole blood are similar and the peak pair ratios of most metabolites will be close to 1. After equal

volume mixing of the  $^{12}\text{C}$ -labeled individual sample and the  $^{13}\text{C}$ -labeled internal standard, the mixed sample is analyzed by high-resolution LC-MS.



**Figure 4.1** Work flow of the finger blood analysis based on CIL-LC-MS.

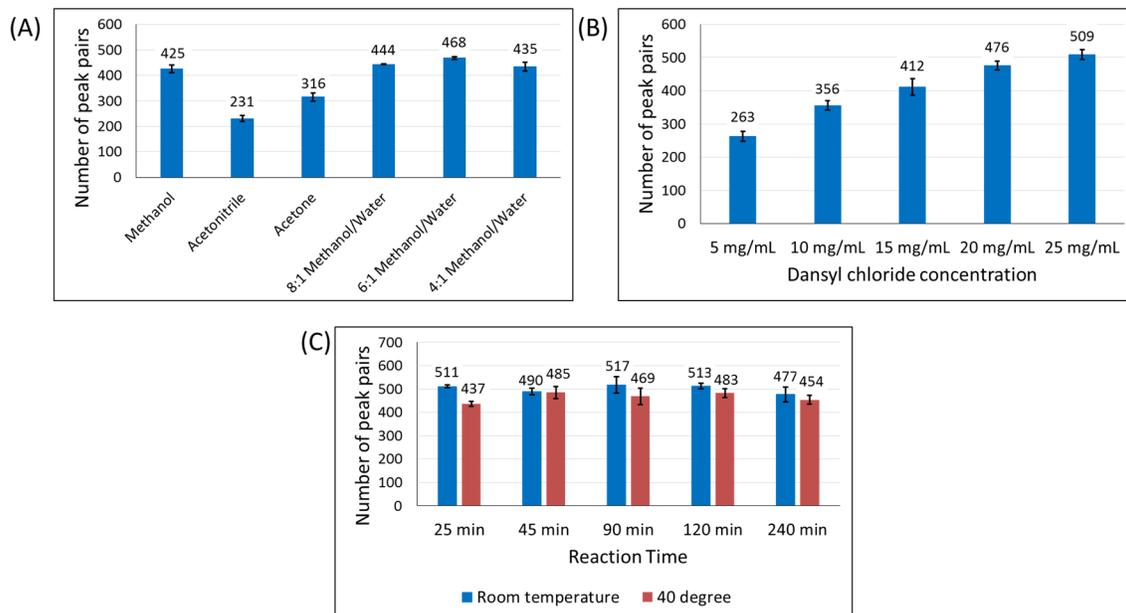
### 4.3.3 Optimization of the labeling method

Although we have previously developed a labeling protocol for the large-scale labeling of serum and plasma samples, the old protocol cannot be directly applied to the finger blood analysis as the sample amount is too small. In order to achieve a high metabolome coverage, we have optimized the small-volume protocol with one microliter of the USS sample. For each experimental condition, one microliter of the USS sample was  $^{12}\text{C}$ -labeled and another one microliter aliquot was  $^{13}\text{C}$ -labeled. The volume of the mixed sample was 30 μL, but

we injected only 3  $\mu\text{L}$  into the LC-MS system to avoid saturating the system and thereby making the difference between experimental conditions less significant. With a larger injection amount, more peak pairs will be detected.

The protein precipitation solvent was first optimized, and the results are shown in Figure 4.2A. For each organic solvent, 1  $\mu\text{L}$  of serum sample was mixed with 55  $\mu\text{L}$  of the solvent to precipitate the proteins. The labeling process and LC-MS analysis are the same for all the solvents. It is very clear that the methanol group gives more peak pairs than the acetone and acetonitrile groups. This is because most of the amine/phenol-containing metabolites are relatively hydrophilic and have better solubility in methanol, which is more hydrophilic than acetone and acetonitrile. Because of the same reason, adding a small portion of water into methanol can increase the number of detected metabolites. On the other hand, adding too much water decreases the protein precipitation performance, and the leftover proteins can consume labeling reagent during the reaction and thereby affect the detectability. In our experiment, we mixed methanol and water in ratios of 8:1, 6:1 and 4:1 (v/v), respectively. With more than 25% water in the solvent, no visible protein precipitates can be observed. The result shows that 6:1 methanol/water, which gives 468 peak pairs, is the optimal solvent for the protein precipitation. For the labeling of freshly collected finger blood sample, since the sample is already in 8  $\mu\text{L}$  of PBS solution, 48  $\mu\text{L}$  of pure methanol is added to the sample to precipitate the proteins and 45  $\mu\text{L}$  of the supernatant is taken for the following steps. And 6:1 methanol/PBS solution was used for the USS sample to maintain the same matrix.

In our previous dansyl-labeling protocol, in order to maximize the yield of the labeling reaction, an excess of labeling reagent was used. In this experiment, as the total amount of metabolites in the sample is smaller, we also tried to use a smaller amount of the reagent. We labeled one microliter of serum with dansyl chloride solutions at 5 mg/mL, 10 mg/mL, 15 mg/mL, 20 mg/mL and 25 mg/mL, respectively. When 25 mg/mL dansyl chloride was mixed with the buffer solution, a small amount of solid dansyl chloride formed on the bottom of the vial, so higher dansyl chloride concentrations were not studied. Figure 4.2B shows the number of detected peak pairs of these experimental conditions. Although the amount of dansyl chloride in 7.5  $\mu$ L of 5 mg/mL solution is much larger than the amount of amine/phenol-containing metabolites in one microliter of serum, there is a clear trend that a higher reagent concentration helps the labeling of metabolites. The concentrations of some low-abundant metabolites can be as low as several nano-molars during the reaction, and the main reaction has to compete with the side reaction between dansyl chloride and water, so an extremely high concentration of the labeling reagent is preferred to ensure the labeling of metabolites at low concentrations. As the difference between 20 mg/mL and 25 mg/mL is not significant ( $p = 0.08$ ), we choose to use 20 mg/mL as the optimal concentration of the labeling reagent.



**Figure 4.2** (A) Numbers of peak pairs detected with different protein precipitation solvents. (B) Numbers of peak pairs detected with different concentrations of dansyl chloride reagent. (C) Numbers of peak pairs detected with different reaction lengths and reaction temperatures.

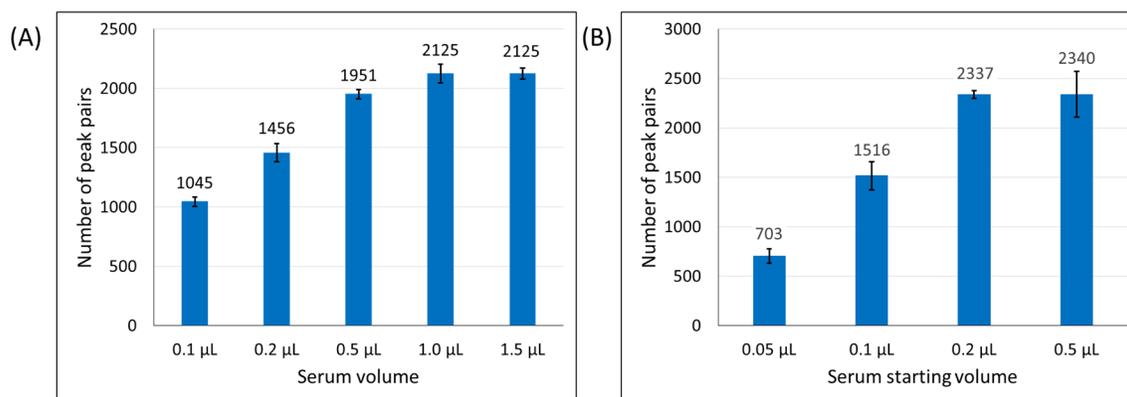
In fact, the danylation reaction can easily happen right after mixing the sample and the labeling reagent. Figure 4.2C shows the numbers of peak pairs detected with different reaction temperatures and reaction time lengths. The most important point is that, regardless of the reaction time, there is no significant difference between reactions at room temperature and at 40 °C. This characteristic of the reaction makes our technique possible to become a point-of-care analysis kit in the areas where an incubating oven is not easily accessible. Furthermore, there is no fundamental difference among different reaction time lengths, indicating that the main reaction completes within several minutes. In our work, 45 min is used to be consistent with our previous works. However, the reaction time can be shortened as needed in the future.

#### 4.3.4 Metabolome coverage of the optimized method

With the optimized labeling protocol, we investigated whether we can achieve an adequate metabolome coverage with only one microliter of the blood sample. We are also interested in the possibility to further decrease the required sample amount while maintaining a large number of metabolites that can be quantified. To perform a more complete and comprehensive whole metabolome analysis of blood, in the future, we will apply our other chemical isotope labeling methods for studying the carboxylic acids<sup>73</sup>, hydroxyl groups<sup>74</sup> and carbonyl groups<sup>75</sup>. In that case, the sample volume available for each analysis has to be minimized as we prefer not to increase the sample collection volume.

With the same labeling protocol, 0.1  $\mu\text{L}$ , 0.2  $\mu\text{L}$ , 0.5  $\mu\text{L}$ , 1.0  $\mu\text{L}$  and 1.5  $\mu\text{L}$  of the USS sample were processed and the same volume of  $^{12}\text{C}/^{13}\text{C}$ -mixed sample (15  $\mu\text{L}$ ) was injected into the LC-MS system. The numbers of detected peak pairs are shown in Figure 4.3A. There were 299 peak pairs detected in the method blank. The background peaks are from the side reactions between dansyl chloride and water, a minor amount of leftover methanol and the impurities in the plastic tubes. Most of the background peaks are very weak, and when a real sample is being analyzed, the signals of more abundant metabolites will suppress their MS signals. Also, our research group has developed software to recognize the background peak pairs and then subtract them from the following analysis. As our IsoMS software can automatically filter the adduct ions and method blank, we can confidently say that the number of peak pairs reflects the number of metabolites detected. With just 0.1  $\mu\text{L}$  of serum, our method detected 1,045 amine/phenol-containing metabolites, demonstrating superior metabolome coverage. And as expected, a larger starting volume

gives more peak pairs. The result of 1.0  $\mu\text{L}$  of serum has 2,125 detected metabolites, which is the highest number in this experiment. The same number of metabolites are detected from 1.5  $\mu\text{L}$  of sample, indicating that the LC-MS detection saturates when the injection amount is too large. Because approximately a half microliters of serum can be separated from one microliter of whole blood, it is important to gauge the metabolome coverage with one-half microliters of serum. The result shows that with 0.5  $\mu\text{L}$  of the USS sample, we can detect and quantify 1,951 metabolites, which is very close to the highest number given by 1.0  $\mu\text{L}$  of serum. Although 0.2  $\mu\text{L}$  serum gives a smaller peak pair number (1,456), it has demonstrated the significant technical advancement that more than a thousand metabolites can be quantified at the same time with less than 0.5  $\mu\text{L}$  of whole blood.



**Figure 4.3** (A) Numbers of peak pairs detected with different starting volumes of the USS sample (The injection volume was fixed at 15  $\mu\text{L}$ ). (B) Metabolome coverage with different starting sample volumes when the nano-LC-MS system is used.

Our CIL method is also compatible with nanoflow-LC-MS, which is a high-sensitivity platform for analyzing a small amount of samples.<sup>201</sup> A nano-LC-QTOF-MS system was used to test if any improvement of detectability can be achieved. Different volumes of the USS sample was labeled with the same protocol and the final  $^{12}\text{C}/^{13}\text{C}$ -mixed sample was

diluted by 5-fold with water. 5.0  $\mu\text{L}$  of the diluted sample was analyzed and the results are shown in Figure 4.3B. Starting with only 0.05  $\mu\text{L}$  of serum, which is equivalent to 0.1  $\mu\text{L}$  of whole blood, the platform can detect 703 metabolites. With 0.2  $\mu\text{L}$  of serum, 2,337 metabolites can be detected. This number is larger than that of 1.0  $\mu\text{L}$  serum on the micro-LC, demonstrating the high detection sensitivity of the nano-LC-MS platform. When it goes higher than 1.0  $\mu\text{L}$ , as the detection saturates, the number of detectable metabolites is not remarkably increased. Nonetheless, these results have proved that it's possible to split one microliter of whole blood into multiple aliquots for performing multiple metabolic analyses at the same time. In addition, since only a small fraction of the labeled product is injected into the nano-LC-MS system, the leftover amount can be used to run instrumental replicates or be stored for future studies.

#### **4.3.5 Sample normalization**

Sample normalization may also be a necessary step to minimize the inter-sample variations. The total concentration of metabolites in blood may differ among individuals, in other words, some people can have a bit more concentrated blood than the others. Since our major goal is to determine the concentration changes of individual metabolites in multiple comparable blood samples, the variations in total metabolite concentration, which can affect the relative quantification result, should be adjusted. We previously developed a pre-acquisition normalization approach, which is an LC-UV based method that determines the total concentration of dansyl-labeled metabolites by measuring the UV absorption of the dansyl group.<sup>115</sup> However, this approach is not applicable to the finger blood analysis due to the large amount of sample it consumes.

Fortunately, the inter-sample variations of total metabolite concentration among blood samples are relatively small, compared with those in other biofluids. As discussed in Chapter 3 and shown in Figure 3.3, the total concentrations of amine/phenol-containing metabolites in 100 dansyl-labeled serum samples were measured by the LC-UV method. The average value is 0.34 mM, and the standard deviation is 0.04 mM, showing the total metabolite concentrations of most blood samples are very close. The serum sample with highest total metabolite concentration is 40% more concentrated than the average level, and the most diluted one's is 29% lower than the average. In most cases, to be considered as a significant change in metabolite concentration, the fold change should be at least 1.5. Therefore, the small inter-sample variations, which can only make an error not greater than 40%, cannot greatly affect the metabolic analysis.

With the pre-acquisition normalization skipped to save the limited amount of sample, we did post-acquisition normalization, which is more convenient and easier to perform, to minimize the inter-sample variances. Warrack et al. have proposed a method called “MS total useful signal”, using the total intensity of components that are common to all samples. Compared with using the total ion intensity, the total useful signal excludes the possible interference from xenobiotics and artifacts and gives more accurate measurement results.

<sup>202</sup> In our study, sum of the peak intensities from all the peak pairs is considered as the total useful signal. We collected 120 finger blood samples from 10 subjects, and the relative standard deviation of the total <sup>12</sup>C-peak intensity among the 120 measured samples is 15.3 %, which confirms that variations in the total concentration of blood are small. As the

total metabolite amount of the individual sample and that of the reference sample are supposed to be the same, the ratio of the total  $^{12}\text{C}$ -peak intensity to the total  $^{13}\text{C}$ -peak intensity is used to normalize the measured peak pair ratios.

#### **4.3.6 Studying the dietary effect of coffee with finger blood analysis**

In order to demonstrate the performance of our finger blood metabolome analyzing method, we did a pilot study on the diet effect of coffee. It is not surprising that there is a close relationship between diet and human metabolomes. After a dietary intake, the blood metabolome may experience changes including the concentration increase of nutrients, metabolic processes of the compounds from the food and also other indirect effects. As discussed before, the dietary variations and the following metabolome variations can interfere with finding the true metabolome changes that are caused by the factors being studied. Walsh and her coworkers have studied the effect of dietary standardization on the metabolomic profiles of healthy humans,<sup>203</sup> and they found that urine metabolome was sensitive to the dietary intake, but the standardization of diet failed to make the plasma samples cluster closer in the multi-variate discriminative analysis. Still there was not enough evidence to prove that diet effect is insignificant in blood metabolome studies, since in Walsh's work, the subjects were asked to fast from midnight to the collection of biofluids. It is possible that certain blood metabolome changes occur after a dietary intake, and the blood metabolome can recover to a stable state within several hours or even a few minutes. Fasting is usually required for research purposes. However, there is no common agreement about how many hours the blood donors should fast before giving blood. In the real world, the blood banks often suggest a healthy meal before blood collection. If these

clinical samples are used for metabolomics research, the diet effect may become more significant. Therefore, it is very important to perform a time-resolved study for assessing the effect of dietary intake to the blood metabolome so that in future studies the sample collection time can also be carefully controlled. Moreover, as the compositions of foods are complicated and diverse, it might be difficult to monitor the metabolic changes when considering the blood metabolome as a whole. Nonetheless, a specific compound at a high concentration in the food should be detected in the blood, or at least its down-stream metabolites should be observed. Even if the dietary intake cannot make any global changes to the blood metabolome, when biomarkers being studied are also involved in the metabolism pathway of an ingested food, the diet effect can become a major issue.

Due to the considerations above, we decided to study the diet effect of coffee, a very common and relatively simple source of dietary stimulation. To make the model as simple as possible, no sugar or cream was added to the dark roast coffee used in the study. Coffee is known to have a large amount of caffeine, which is a central nervous system stimulant and an active small-molecule metabolite, so the metabolism of caffeine was monitored. To study the short-term metabolic effects, the interval between coffee intake and blood collection was set to two hours. Another major objective of this study is to demonstrate that the finger blood analysis has good metabolome coverage, repeatability and robustness comparable to those of the venous blood analysis.

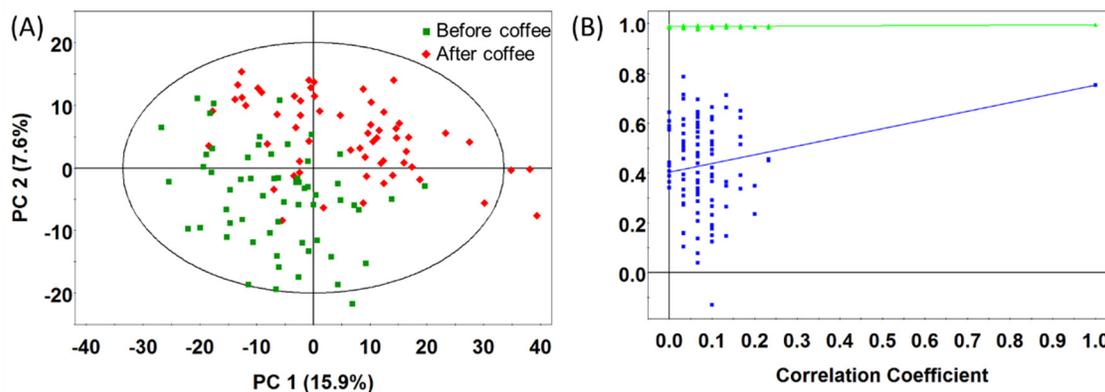
Ten participants were asked to refrain from consuming caffeine-containing food and drink for at least one day and to fast overnight before donating blood. Two aliquots of one-

microliter finger blood were collected as experimental duplicates. After that, the participants finished 350 mL of regular dark roast coffee (from Starbucks) within 20 minutes. Two hours after the first blood collection, we collected the second blood sample from each participant. This process was repeated for another two days, so totally 120 samples were collected and analyzed. Each individual sample was  $^{12}\text{C}$ -dansyl-labeled, and then mixed with equal volume of  $^{13}\text{C}$ -dansyl-labeled USS sample. The mixed samples were analyzed by the LC-FTICR-MS system. On average, 1,647 peak pairs were detected from each sample. Finally, the acquired peak pairs were aligned together, and their ratio values were normalized.

We applied the Partial Least Squares-Discriminant Analysis (PLS-DA) to reveal the statistical differences between the samples collected before and after coffee intake. To be included in the final list for statistical analysis, each peak pair should have valid ratio values in more than 80% of the samples of at least one study group (“before coffee” or “after coffee”), so that the percentage of missing values can be minimized and the statistical result can be more accurate. Finally, a total of 1,722 peak pairs and their ratio values were studied. Among them, 73 were positively identified by the Dansyl-library, 578 were putatively identified by the HMDB library and another 815 were putatively identified by the EML library.

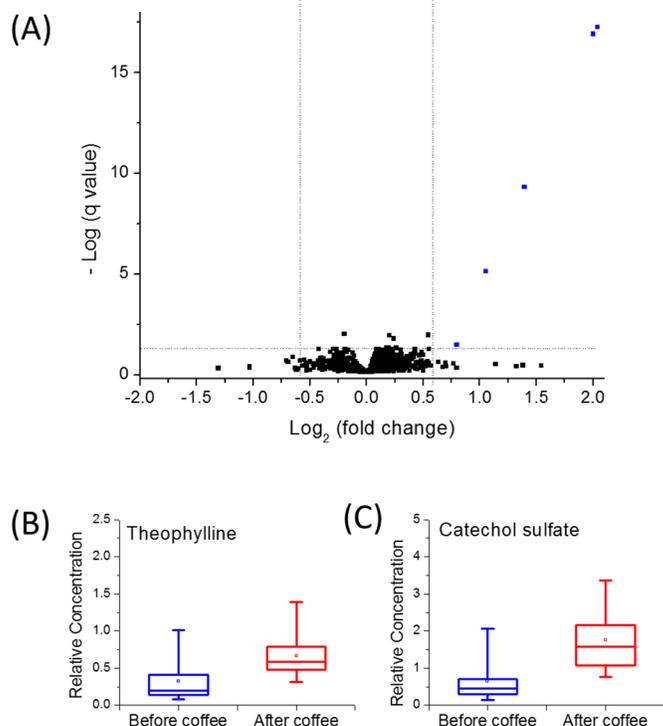
In Figure 4.4A, the score plot of the PLS-DA analysis, there is no obvious separation between the cluster of green dots (“before coffee”) and that of red dots (“after coffee”), indicating that there is no statistically significant difference between the two groups, which

is further confirmed by the permutation test result shown in Figure 4.4B. The performance of the discriminative model is gauged by two indicators:  $R^2$  and  $Q^2$ , which are 0.995 and 0.754, respectively. The permutation test result tells us that with totally randomized group assignment, the model may give similar  $R^2$  and  $Q^2$  values, which means the model is not valid.



**Figure 4.4** (A) PLS-DA score plot, showing that the clusters of data points before and after coffee intake overlap each other. The performance indicators,  $R^2$  and  $Q^2$ , are 0.995 and 0.754, respectively. (B) Permutation test result, confirming that the statistical separation between the two study groups is not valid.

We also used uni-variate analysis to study the changes of individual metabolites. For each metabolite, a fold change is calculated as the ratio of the average concentration in the “after coffee” group to that in the “before coffee” group, and an FDR-adjusted p-value (or q-value) is used to monitor the statistical significance. If the fold change of one metabolite is larger than 1.5 or smaller than 0.67, and the q-value is smaller than 0.05, we consider this metabolite is significantly changed. A volcano plot, which plots the q-value against the fold change, can visualize the results and it is shown as Figure 4.5A.



**Figure 4.5** (A) Volcano plot, showing 6 significantly increased metabolites (in blue) after coffee intake. (B) Box plot of theophylline, showing the distributions of its blood concentration before and after coffee. (C) Box plot of catechol sulfate, showing the distributions of its blood concentration before and after coffee.

The volcano plot agrees with the result of the PLS-DA analysis that no significant difference between the two study groups was observed. Only six metabolites experienced significant concentration change after coffee, and all of them had a concentration increase. The information of these six metabolites is given in Table 4.1. Theophylline is the only positively identified metabolite among them, and its fold change is 2.08. Figure 4.5B is the box plot that shows the distributions of theophylline's relative concentration before and after coffee intake. Caffeine plays a major role in the dietary exposure of coffee, but unfortunately, our dansyl-labeling method cannot label and detect caffeine itself. Nonetheless, theophylline is one of the major down-stream metabolites of caffeine, and the

concentration increase of theophylline can indirectly tell the absorption and metabolic degradation of caffeine. The other significant metabolites, except the one that cannot be identified, are putatively identified as prolyl-proline, catechol sulfate, xanthine [+C<sub>2</sub>H<sub>4</sub>] and oxypurinol [+C<sub>2</sub>H<sub>4</sub>]. Catechol is a known metabolite formed during coffee roasting<sup>204</sup> and by dansyl-labeling we detected the dansylated catechol in coffee. In blood, catechol is conjugated to sulfate, so an increased level of catechol sulfate was observed after coffee intake, as illustrated by the box plot in Figure 4.5C.

**Table 4.1** List of 6 metabolites that have significant concentration changes among the 10 participants after coffee intake.

Retention time (min)	Mass of dansylated metabolite (Da)	Mass of metabolite (Da)	Metabolite	HMDB ID	Identificaion type	Fold change	q value
7.14	446.1643	212.1060	L-prolyl-L-proline	HMDB11180	Putative	1.74	3.20E-02
8.91	424.0538	189.9955	Catechol sulfate	HMDB61713	Putative	2.63	4.67E-10
11.05	414.1243	180.0660	Xanthine + [C <sub>2</sub> H <sub>4</sub> ]	NA	Putative	2.63	4.67E-10
12.36	414.1245	180.0661	Oxypurinol + [C <sub>2</sub> H <sub>4</sub> ]	NA	Putative	4.11	5.77E-18
12.39	472.1784	238.1201	NA	NA	NA	4.00	1.22E-17
14.85	414.1249	180.0665	Theophylline	HMDB01889	Defenitive	2.08	7.55E-06

The result shows that for the amine/phenol-submetabolome, there is no statistically significant change that can be observed two hours after coffee intake. However, please note

that certain changes may happen at different time scales or at other submetabolomes. The purpose of this study is to demonstrate that the finger blood analysis can cover an adequate number of metabolites, sensitively detect concentration changes to individual metabolites and have good repeatability among experimental replicates. The results have proved that the finger blood can readily replace the venous blood in various metabolomics studies.

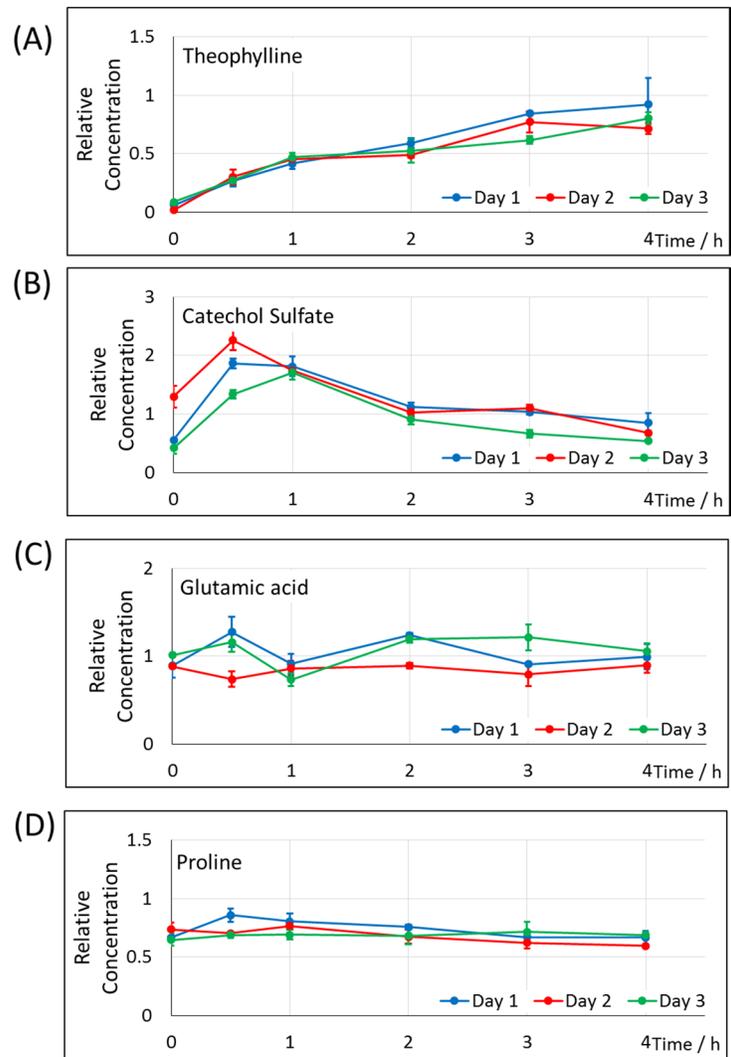
#### **4.3.7 Time-resolved metabolic analysis**

We also used our technique to perform a time-resolved monitoring of blood metabolome after coffee consumption. One subject fasted overnight and the finger blood samples were collected before drinking coffee. Thereafter, more samples were collected at 0.5, 1, 2, 3 and 4 hours after coffee intake. For each time point, experimental triplicates were analyzed, and the whole process was repeated for another two times named as “Day 2” and “Day 3”. The blood concentrations of theophylline and catechol sulfate, the significant metabolites confirmed in the first study, were monitored during the 4-hour period. In Figure 4.6A, the relative concentration of theophylline is plotted against the sample collection time, and the results from different days are shown in different colors. The subject was not a habitual coffee drinker, so the initial value before coffee was close to zero. Right after the coffee was consumed, the blood concentration of theophylline started to increase. At the 4-hour point, the concentration of theophylline reached a considerable level, which is many times higher than the starting point. This implies that if some metabolites of caffeine are considered as biomarkers, coffee intake before sample collection may interfere with the analysis and the sample collection time matters a lot. Importantly, the same changing trend was observed on three different days, indicating good biological and technical

reproducibility of the study. Catechol sulfate shown in Figure 4.6B reached its highest blood concentration between 0.5 and 1 hour after the coffee intake and then started to decrease to the original level. The metabolism of some metabolites can be as fast as catechol sulfate, or even faster, and this explains why no significant metabolome change was observed 2 hours after coffee. Overall, our technique has demonstrated the ability to sensitively and accurately monitor the blood concentration of an adequate number of metabolites based on frequent sample collections, which can be very aggressive in the conventional studies based on venous blood.

The concentration-time curves of glutamic acid and proline are also shown in Figure 4.6 as examples of biological variations in metabolites that are not directly associated with the coffee intake. Although the concentration of glutamic acid was fluctuating, it was never far away from the average level, showing the blood metabolome's ability to stabilize itself. The highest relative concentration of glutamic acid among all the data points is 1.27 and the lowest is 0.74. As the ratio of the highest to the lowest is 1.73, we can tell the variations in blood metabolome is a non-negligible factor in metabolomics studies and the time-resolved analyses can help us minimize the interference. The concentration of proline, as shown in Figure 4.6D, was more stable and the ratio of the highest to the lowest is 1.44.

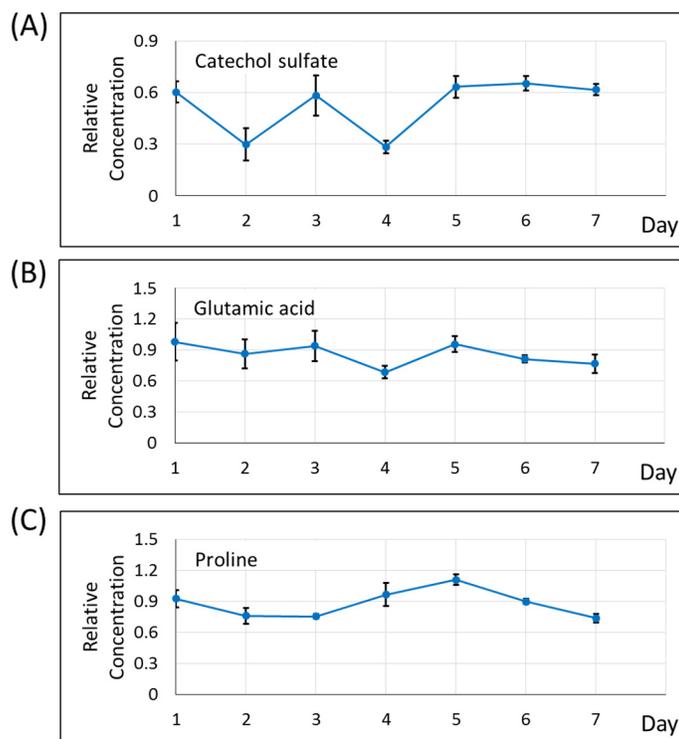
The day-to-day variations of these metabolites were also studied and the results are shown in Figure 4.7. The same subject donated blood in this study and the diet was not controlled. On each day during a week, finger blood samples were collected one hour after breakfast.



**Figure 4.6** Concentration-time curves of (A) theophylline, (B) catechol sulfate, (C) glutamic acid and (D) proline, illustrating their concentration changes during the 4-hour period after coffee intake. The blue curve represents data acquired on Day 1, the red curve represents Day 2 and the green curve represents Day 3.

The subject was not a habitual coffee drinker, so theophylline was not detected in his blood without any caffeine exposure. Catechol sulfate, which exists in many other kinds of food, was detected and it showed very significant concentration fluctuation, as shown in Figure 4.7A. In Figure 4.7B and 4.7C, glutamic acid and proline also demonstrated day-to-day variations in their blood concentrations. During the week, the ratio of the highest glutamic

acid concentration to the lowest is 1.43, and that of proline is 1.50. We can conclude from the results that if only two sample collection points are included in a metabolomics study, any fold change not higher than 1.50 should be carefully assessed as interference due to biological variations can be considerable.



**Figure 4.7** Concentration-time curves of (A) catechol sulfate, (B) glutamic acid and (C) proline, showing the day-to-day variations of their blood concentrations during a week.

#### 4.4 Conclusions

We have successfully developed a high coverage metabolomics analyzing method for one microliter or even lower volumes of finger blood based on CIL and high-resolution LC-MS. Our technique has demonstrated adequate detection sensitivity, repeatability and robustness, showing that the finger blood can be used to replace venous blood in various metabolomics studies without any significant drawbacks or limitations. Particularly, the

finger blood collection method that we used is much less invasive than venous blood collection and it can be done without any specially trained personnel. As more samples can be collected without a significant increase in financial and time costs, the accuracy and statistical power of measurements in metabolomics studies will be greatly improved. And more importantly, our method can readily achieve the frequent blood collection and time-resolved metabolomics analyses. As an example, we performed accurate and time-resolved monitoring of metabolite concentrations after a dietary exposure to coffee.

In the future, the method can be applied to study the metabolic consequences of various environmental exposures. Compared with the other time-resolved studies, our method has a major advantage that the relative quantification of a large number of metabolites can be done at the same time. In the future, with more peak pairs positively identified and some of them proved to be biomarkers, our technique will become a powerful tool for time-resolved studies of specific biomarkers during the progress of diseases. And such analyses will help with understanding the metabolic processes associated with a disease state. We note that there might be differences between venous blood and finger blood, however, since we mainly study the metabolite changes caused by a disease state or an environmental factor, and our pilot studies have demonstrated that the finger blood can sensitively reflect these changes, the relationship between venous blood metabolome and finger blood metabolome is not our major concern. If we want to compare the results of venous blood studies and finger blood studies in the future, the difference between the two materials will be further studied.

## Chapter 5

### High-coverage Metabolomics Analysis of One Microliter Blood Using Two Isotope Labelings and High-resolution LC-MS

#### 5.1 Introduction

Despite the fact that many reported biomarker candidates are amino acids or their derivatives, which belong to the amine and phenol categories, carboxylic acids have also demonstrated biological significance in various physiological processes. For example, it has been reported that short chain carboxylic acids, such as acetic acid, propionic acid, and butyric acid, can alter the function of cells.<sup>205-206</sup> Particularly in human blood, the short chain carboxylic acids play roles in the regulation of leukocyte function.<sup>207</sup> Furthermore, the deregulation of fatty acids has been discovered in hepatic diseases,<sup>208</sup> obesity,<sup>209</sup> and diabetes.<sup>210</sup>

Moreover, carboxylic acids have long been used as food additives and preservatives for extending the shelf life, though the antimicrobial mechanism is not fully understood.<sup>211-212</sup> In animal husbandry, carboxylic acids are also utilized as animal feed additives to improve the growth rate and feed conversion rate.<sup>213</sup> Consequently, the profile of carboxylic acids can become a powerful indicator for studying the metabolic interactions between human body and the environment, especially the diet effects. For instance, Hodson et al. demonstrated in their work that fatty acids could potentially be used as biomarkers of dietary intake.<sup>214</sup>

In chapter 5, we described the importance of time-resolved metabolomics analyses for assessing the metabolic outcomes of environmental exposures, and we also reported a high-coverage metabolome profiling method of one microliter of finger blood using dansylation isotope labeling and high-resolution LC-MS, enabling time-resolved metabolomics analyses with maximized metabolome coverage and minimally invasive sample collection. Since carboxylic acids are commonly seen species in the environment, especially in diet sources, time-resolved analyses of carboxylic acids will provide us abundant information about the body's physiological responses to environmental stimulations. To study the carboxyl-containing metabolites, we have previously developed a differential isotope labeling technique which uses p-dimethylaminophenacyl (DMPA) bromide as the reagent,<sup>73</sup> and this method has been used for profiling the carboxyl-submetabolome of venipuncture blood samples. In this chapter, we apply two isotope labeling methods, the dansylation labeling and the DMPA-labeling, to the metabolomics analysis of one microliter of finger blood. The metabolic responses of both the amine/phenol-submetabolome and the carboxyl-submetabolome to a dietary stimulation are also assessed.

## **5.2 Materials and methods**

### **5.2.1 Chemicals and reagents**

All the chemicals and reagents, unless otherwise stated, were purchased from Sigma-Aldrich Canada (Markham, ON, Canada). For chemical isotope labeling reactions, the <sup>12</sup>C-labeling reagents (dansyl chloride and DMPA bromide) were purchased from Sigma-Aldrich, and the <sup>13</sup>C-labeling reagents were synthesized in our lab using the procedures

published previously.<sup>70, 73</sup> LC-MS grade water, methanol, and acetonitrile (ACN) were purchased from Thermo Fisher Scientific (Nepean, ON, Canada).

### **5.2.2 Finger blood sample collection**

First, the finger surface was cleaned with 70% isopropyl alcohol. A Bayer's Microlet 2 lancing device was used for the finger pricking. After that, 1  $\mu\text{L}$  of whole blood sample was taken from the punching site with a micropipette, and immediately transferred to a 0.2 mL PCR tube to mix with 10  $\mu\text{L}$  of PBS solution. Then the mixture was centrifuged at 1,000 rpm for 10 min. One 5  $\mu\text{L}$  aliquot of supernatant was transferred into another PCR tube for the dansyl-labeling, and another 5  $\mu\text{L}$  aliquot was taken for the DMPA-labeling.

### **5.2.3 Dansyl-labeling**

5  $\mu\text{L}$  of the supernatant from the last step was mixed with 2.5  $\mu\text{L}$  of 250 mM sodium carbonate/sodium bicarbonate buffer and 7.5  $\mu\text{L}$  of freshly prepared  $^{12}\text{C}$ -DnsCl solution (20 mg/mL) (for light labeling). The solution was vortexed, spun down, and then incubated at 40 °C for 45min. After that, 2  $\mu\text{L}$  of 250 mM NaOH solution was added to quench the excess dansyl chloride. The solution was then incubated at 40 °C for another 10 min. Finally, 5  $\mu\text{L}$  of 425 mM formic acid in 1:1 ACN/H<sub>2</sub>O was added to consume excess NaOH and to make the solution acidic. For the internal reference, 0.25  $\mu\text{L}$  of the universal serum standard (USS) was diluted to 5  $\mu\text{L}$  with PBS, and then labeled by  $^{13}\text{C}$ -DnsCl with the same protocol.

#### 5.2.4 DMPA-labeling

5  $\mu\text{L}$  of the supernatant from the sample collection was mixed with 5  $\mu\text{L}$  of 0.2M triethanolamine and 15  $\mu\text{L}$  of freshly prepared  $^{12}\text{C}$ -DMPA bromide solution (7.5 mg/mL) (for light labeling). After the solution was vortexed and spun down, the reaction was allowed to proceed at 75  $^{\circ}\text{C}$  for 45 min. Finally, the excess amount of DMPA bromide was quenched by 15  $\mu\text{L}$  of 0.2M triglycine. For the internal reference, 0.25  $\mu\text{L}$  of the USS sample was diluted to 5  $\mu\text{L}$  with PBS, and then labeled by  $^{13}\text{C}$ -DMPA bromide with the same protocol.

#### 5.2.5 LC-QTOF-MS

For LC-QTOF-MS, an Agilent 1100 series binary system (Agilent, Palo Alto, CA) and an Agilent reversed-phase Eclipse plus C18 column (2.1 mm $\times$ 100 mm, 1.8  $\mu\text{m}$  particle size, 95 A pore size) were used. The flow was loaded to the electrospray ionization (ESI) source of a Bruker maXis impact high-resolution quadrupole time-of-flight (Q-TOF) mass spectrometer (Bruker, Billerica, MA). All MS spectra were obtained in the positive ion mode.

LC solvent A was 0.1% (v/v) formic acid in 5% (v/v) ACN/H<sub>2</sub>O, and solvent B was 0.1% (v/v) formic acid in ACN. The gradient elution profile for dansyl-labeled samples was as follows: 0 min (20% B), 0-3.5 min (20-35% B), 3.5-18 min (45-65% B), 18-21 min (65-95% B), 21-24 min (95-99% B), and 24-34 min (99% B). The column was re-equilibrated with the initial mobile phase condition for 15 min before injecting the next sample. The flow rate was 180  $\mu\text{L}/\text{min}$ , and the injection volume was 15.0  $\mu\text{L}$ .

The 42-min gradient for all DmPA-labeled samples is: 0 min (20% B), 0-9 min (20%-50% B), 9-22 min (50%-65% B), 22-26 min (65%-80% B), 26-29 min (80%-99% B) and 29-42 min (98% B). The column was re-equilibrated with the initial mobile phase condition for 15 min before injecting the next sample. The flow rate was 180  $\mu\text{L}/\text{min}$ , and the injection volume was 25.0  $\mu\text{L}$ .

Representative LC-MS chromatograms are provided in Appendix Figure 3.

### **5.2.6 Data processing and statistical analysis**

A software tool, IsoMS,<sup>72</sup> was used to process the raw data generated from multiple LC-MS runs by peak picking, peak pairing, peak-pair filtering and peak-pair intensity ratio calculation. The same peak pairs detected from multiple samples were then aligned together by IsoMS-Align. The missing ratio values were filled back by using the Zero-fill program.<sup>137</sup> Finally, IsoMS-Quant<sup>157</sup> was used to determine the chromatography-peak-intensity ratio of a  $^{12}\text{C}/^{13}\text{C}$ -pair. The final datasheet of metabolite concentrations was exported to MetaboAnalyst 3.0<sup>165</sup> for multivariate statistical analysis.

### **5.2.7 Metabolite identification**

Metabolite identification was performed based on mass and retention time match to a dansyl standard library<sup>159</sup> which contains 315 metabolite standards and a carboxylic acid standard library which includes 187 standards. Putative identification was done based on

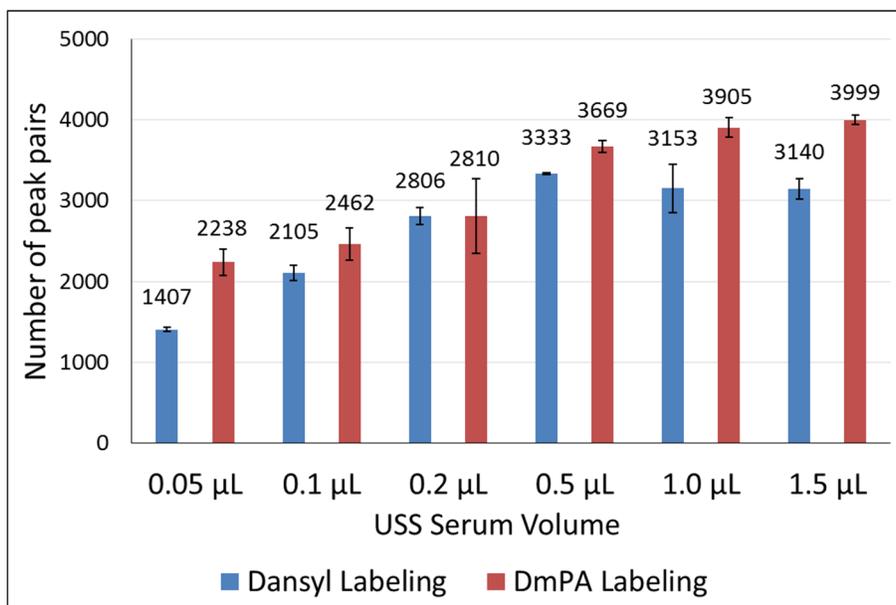
accurate mass matches to the metabolites in the human metabolome database (HMDB)<sup>158</sup> (8,021 known human endogenous metabolites) and in the Evidence-based Metabolome Library (EML)<sup>41</sup> (375,809 predicted human metabolites with one reaction) using MyCompoundID. The mass accuracy tolerance window was set at 0.008 Da for database search.

## **5.3 Results and Discussion**

### **5.3.1 Metabolome coverage of the two differential isotope labeling methods**

In this work, each microliter of raw blood is split into two equal aliquots. One of them is labeled by <sup>12</sup>C-dansyl chloride and then mixed with equal volume of <sup>13</sup>C-labeled USS sample. The other aliquot is labeled by <sup>12</sup>C-DMPA bromide before being mixed with the same volume of <sup>13</sup>C-labeled USS sample. For each labeling, the available sample amount is approximately equivalent to only 0.25  $\mu$ L of serum. Although the sample volume is extremely small, the isotope labeling techniques, accompanied with high-resolution LC-MS, can generate an adequate number of detected peak pairs. To examine the metabolome coverage of the two-labeling platform, different volumes (0.05  $\mu$ L, 0.2  $\mu$ L, 0.2  $\mu$ L, 0.5  $\mu$ L, 1.0  $\mu$ L, and 1.5  $\mu$ L) of the USS sample were labeled and <sup>12</sup>C/<sup>13</sup>C-mixed correspondingly. The injection volume of dansyl-labeled samples was 15  $\mu$ L, and that of DMPA-labeled samples was 25  $\mu$ L. The QTOF-MS was set to High-dynamic-range (HD) mode during the acquisition. Since the concentrations of different metabolites are very diverse, the HD mode can significantly increase the detected number of peak pairs.<sup>58</sup> The numbers of peak pairs acquired at different sample volumes are shown in Figure 5.1.

With just 0.05  $\mu\text{L}$  of serum, the dansyl-labeling provides 1,407 peak pairs and the DMPA-labeling gives 2,238 peak pairs, demonstrating the high detection sensitivity enabled by the differential isotope labeling. Both of the two techniques are able to detect more peak pairs as the starting volume increases. The dansyl-labeling reaches the highest number of peak pairs (3,333) at 0.5  $\mu\text{L}$  of USS serum, while the saturation point of the DMPA-labeling is at 1.0  $\mu\text{L}$ , having more than 3,900 peak pairs. Starting with 0.2  $\mu\text{L}$  of serum, whose total amount of metabolites is close to that in the 5  $\mu\text{L}$  diluted finger blood, the dansyl-labeling shows 2,806 peak pairs, which are only 16% fewer than the optimal number. Meanwhile, 2,810 of DMPA-labeled peak pairs are quantified, 28% lower than the highest number. Seeing that the analysis of 0.2  $\mu\text{L}$  of USS serum gives much higher metabolome coverage than traditional methods, we may conclude that our technique is qualified for monitoring the changes of amine, phenol and carboxyl submetabolomes by studying a minimal amount of finger blood.



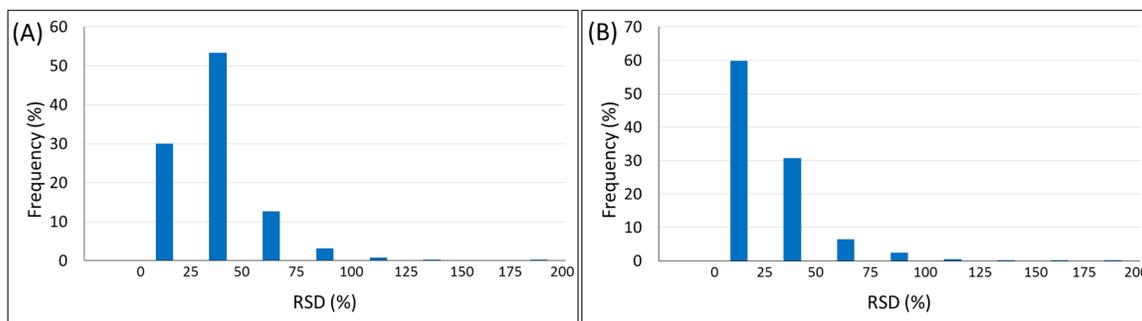
**Figure 5.1** Numbers of peak pairs detected from different volumes of USS sample by dansyl-labeling (in blue) or DMPA-labeling (in red).

### **5.3.2 Time-resolved metabolomics analysis for studying day-to-day metabolome variations**

The day-to-day metabolome variations within one subject are potentially confounding factors to biomarker discoveries and clinical applications. Notably, the diet is usually not strictly controlled in many studies, and therefore it could be a major environmental stimulation to cause day-to-day metabolic changes. To assess these variations, we collected finger blood samples from one male subject consecutively for seven days. During the study, the diet was not controlled, and on each day, finger blood samples (experimental triplicates) were collected one hour after breakfast. The individual samples were  $^{12}\text{C}$ -dansyl-labeled and  $^{12}\text{C}$ -DMPA-labeled before being mixed with the corresponding USS standards. As the same internal reference is applied to all samples, the concentrations of each metabolite over this period are directly comparable and can demonstrate the day-to-day variations.

After excluding the peak pairs with more than 50% missing values, we successfully detected 2,074 dansyl-labeled and 2,254 DMPA-labeled metabolites. We positively identified 72 dansyl-labeled metabolites and 26 DMPA-labeled metabolites using the standard libraries. The accurate masses, average concentrations and relative standard deviations (RSD) of these 98 metabolites are given in Table 5.1. Also, 490 dansyl-labeled metabolites and 501 DMPA-labeled metabolites are putatively identified by searching through the HMDB library. The EML library putatively identified 983 dansyl-labeled metabolites and 1237 DMPA-labeled metabolites.

For each metabolite, we calculated the RSD among the concentrations measured during the week. The distributions of the RSDs are shown in Figure 5.2. For the amine/phenol-containing metabolites (Figure 5.2A), the majority of the RSDs are below 100%, with a very small number of outliers, which mostly belong to the metabolites with a high percentage of missing values. The median of the RSDs is 32.52%, and more than half of them lie between 25% and 50%, showing that the day-to-day variations are existing but not very severe. Figure 5.2B illustrates that the RSDs of the carboxyl-containing metabolites are generally smaller than those of the dansyl-labeled metabolites. Except for a few outliers, most of the RSDs are also below 100% with a median value of 21.49%. And almost 60% of them are below 25%, suggesting that the day-to-day variations in the carboxyl-submetabolome are not very significant.

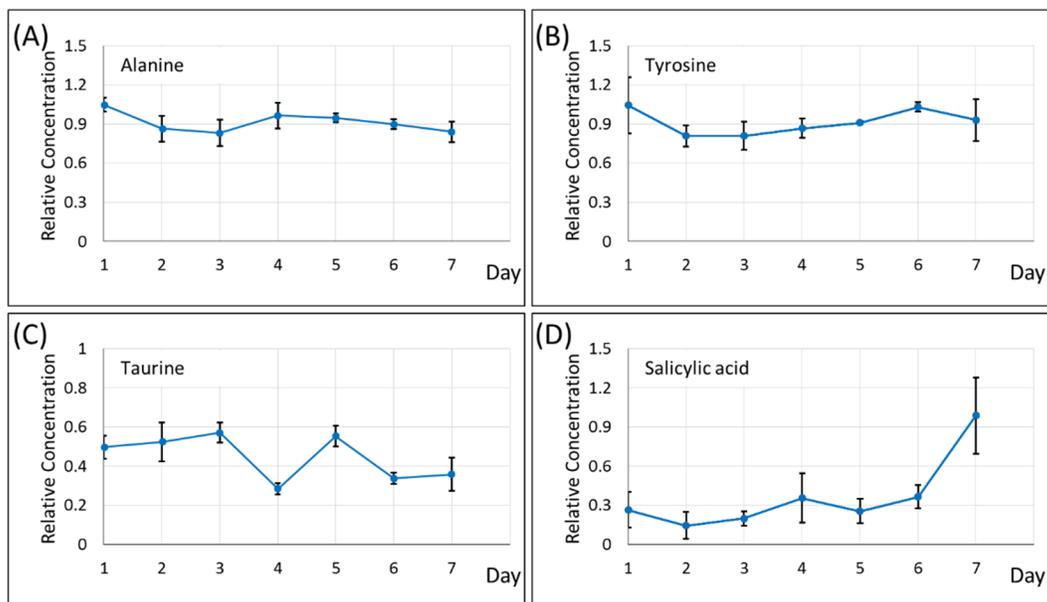


**Figure 5.2** Distribution of relative standard deviations defining the day-to-day variability in (A) amine/phenol-containing metabolites and (B) carboxyl-containing metabolites.

The concentration-day plots of four dansyl-labeled metabolites are shown in Figure 5.3 as examples. Among the positively identified amine/phenol-containing metabolites, alanine has the smallest RSD (11.07%). As shown in Figure 5.3A, the highest concentration (1.05) occurred on Day 1, and the lowest concentration (0.83) was on Day 3. Compared to the average concentration (0.92), the fold change is 1.14 for the highest value and is 0.90 for

the lowest level. As most biomarker candidates require fold change to be at least 1.2, the day-to-day variations in alanine concentration will not affect the biomarker studies. The concentration of tyrosine was also relatively stable, with RSD equal to 23.39%. Tyrosine is a proved indicator that can differentiate between males and females,<sup>169</sup> with the fold change of 0.82 (female/male). The fold change of the lowest concentration in Figure 5.3B is 0.95, indicating that the inter-sex variation is more significant than the within-individual variation. Taurine has multiple biological functions in the human body,<sup>215</sup> and it has been reported to be biomarker candidate for various diseases. For instance, Engelborghs et al. found that taurine levels in the cerebrospinal fluid of Parkinson's disease patients were significantly lower compared to the control group (fold change = 0.788).<sup>216</sup> Unlike the relatively stable alanine and tyrosine, taurine had a more fluctuating curve, as shown in Figure 5.3C. The lowest concentration (Day 4) is 38% lower than the average (fold change = 0.622). The evidence has led us to believe that the blood level of taurine is sensitive to multiple biological and environmental factors including the diet effects. As the fold change by chance is larger than that found in Parkinson's disease, if taurine is used as a biomarker, multiple sample collections over a period will be necessary to accurately profile the blood taurine level of one individual. Salicylic acid has one of the largest RSD (80.26%) among these metabolites. However, by looking at the plot in Figure 5.3D, we can find that the high RSD is mainly due to the concentration jump on Day 7. This metabolite is present in fruits and vegetables, and higher blood concentration can be found in vegetarians than non-vegetarians.<sup>217</sup> The origin of salicylic acid from food explains why the blood level of it soared 2.7-fold on Day 7. The discovery also confirms that diet effects contribute

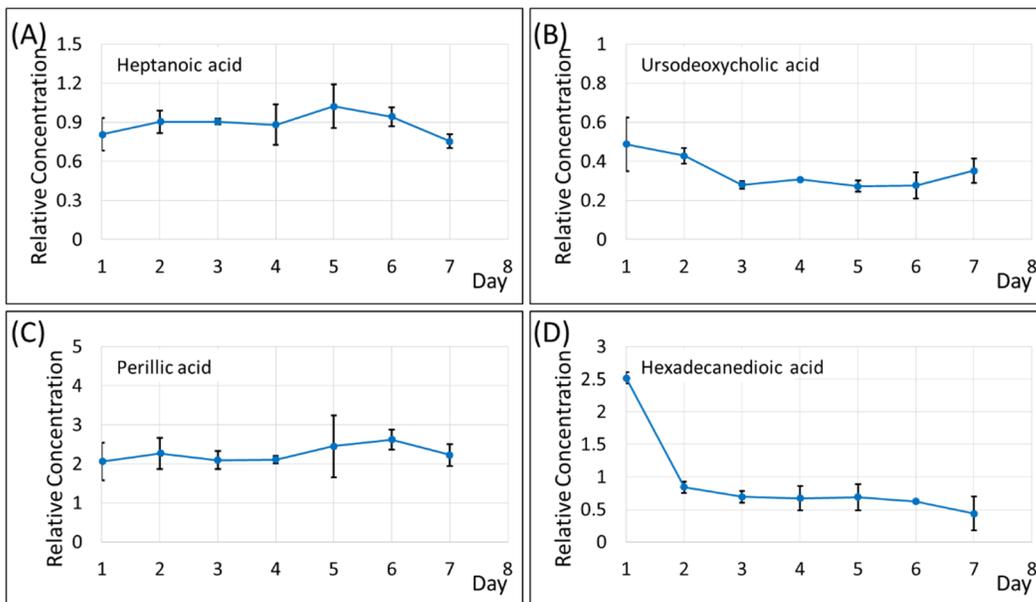
significantly to the variability of some metabolites, and our technique is an ideal platform for studying these exposomics effects.



**Figure 5.3** The concentration-day curves of (A) alanine, (B) tyrosine, (C) taurine and (D) salicylic acid, showing variations in their concentrations during the week.

For the DMPA-labeling, the concentration-day curves of heptanoic acid, ursodeoxycholic acid, perillic acid and hexadecanedioic acid are shown in Figure 5.4. Their RSDs are 14.71%, 30.24%, 23.41% and 69.36%, respectively. Heptanoic acid, a medium-chain fatty acid, demonstrated very stable blood concentration over the week. Ursodeoxycholic acid is an endogenous bile acid, and it protects against the membrane-damaging effects associated with hydrophobic bile acids.<sup>218</sup> The fold changes (compared with average) of its highest and lowest concentrations during the week are 1.39 and 0.79. This tells us that the within-individual changes can sometimes be larger than 1.2-fold, and increasing the fold change threshold to 1.5 can effectively cope with the variations of many metabolites. Perillic acid is an intermediate in the degradation pathway of limonene and pinene. It has

been discovered that perillic acid has protective functions against radiation<sup>219</sup> and cancer.<sup>220</sup> Its blood concentration also remained very stable during the study period. Hexadecanedioic acid is one of the carboxyl-containing metabolites that have prominent variations. Similar to the case of salicylic acid, its high RSD is mainly because of the remarkably high concentration on Day 1. Since hexadecanedioic acid has been found in some plants,<sup>221-222</sup> the abnormal concentration may also due to diet effects. These results have shown that our method has the adequate sensitivity and accuracy for the assessment of metabolic responses to dietary stimulations.

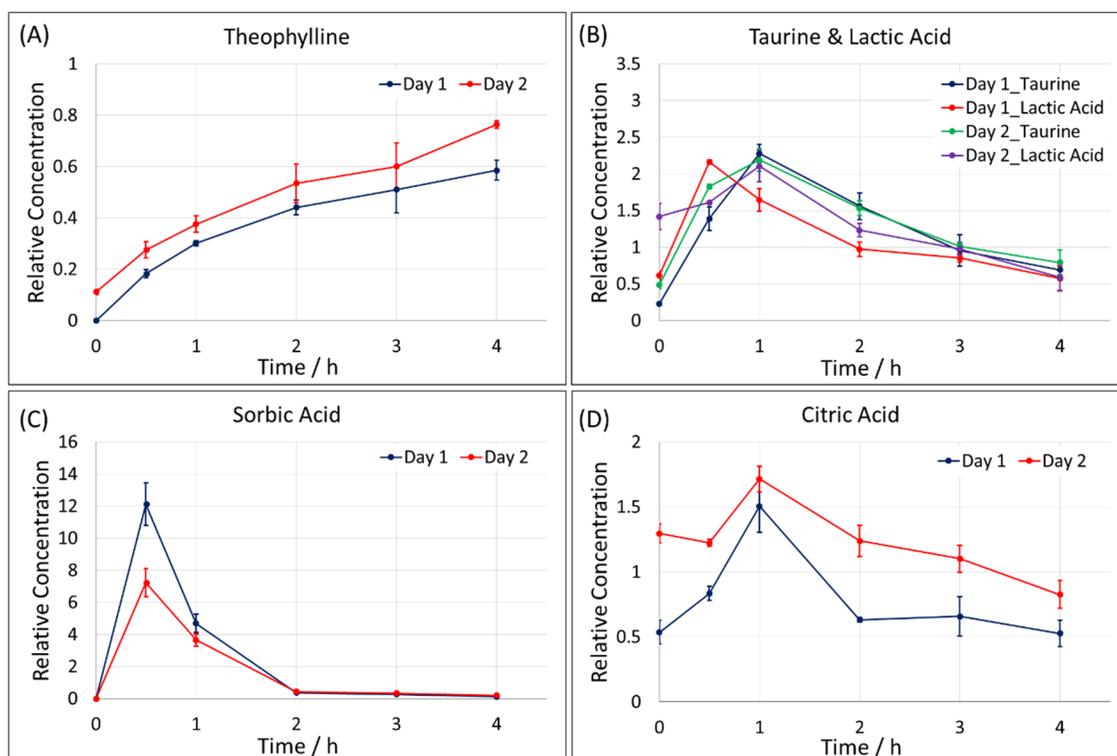


**Figure 5.4** The concentration-day curves of (A) heptanoic acid, (B) ursodeoxycholic acid, (C) perillic acid and (D) hexadecanedioic acid, showing variations in their concentrations during the week.

### 5.3.3 Studying the diet effects of an energy drink with two labeling methods

To study the diet effects more specifically within a shorter period, we chose the energy drink as the dietary stimulation. Compared with other food sources, the energy drink is a relatively simple mixture of a few compounds at very high concentrations. In addition to a large amount of sugars, an energy drink usually has caffeine and taurine as its active ingredients. Caffeine is widely used as a central nervous system stimulant, and taurine is believed to improve the force of skeletal muscles.<sup>223</sup> In a 16 oz. can of the energy drink (Monster Energy) we studied, there is approximately 160 mg of caffeine and 2000 mg of taurine, which can be a very strong environmental stimulation to the blood metabolome. Except for the active ingredients, the energy drink also has benzoic acid, citric acid and sorbic acid added as preservatives. It is also interesting to study the metabolism of these non-nutritional ingredients.

In this experiment, one subject fasted overnight before the experiment. The finger blood samples were collected in the early morning, labeled as time “0”. Then the subject finished a 16 oz. can of the energy drink within 10 min. More finger blood samples were collected at 0.5 hours, 1 hour, 2 hours, 3 hours and 4 hours after the energy drink intake. These samples were processed with the two chemical isotope labeling methods, and then mixed with the <sup>13</sup>C-labeled USS sample for LC-MS analysis. Forty-eight hours after the first experiment, the whole process was repeated as experimental duplicates, annotated as “Day 1” and “Day 2”. We studied the concentration changes of theophylline, taurine, lactic acid, sorbic acid and citric acid. For each metabolite, the concentrations are plotted against the sample collection time, and the curves are shown in Figure 5.5.

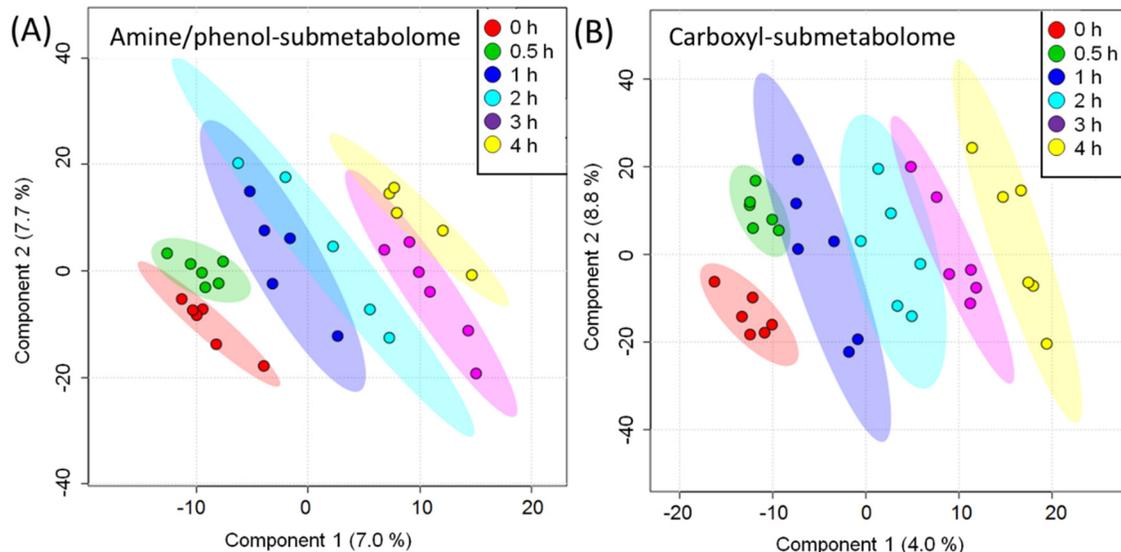


**Figure 5.5** Concentration-time curves of (A) theophylline, (B) taurine, lactic acid, (C) sorbic acid and (D) citric acid, showing the concentration changes after the energy drink intake.

Theophylline is one of the major metabolites of caffeine. As shown in Figure 5.5A, its relative concentration rose from zero to 0.59 during the 4-hour period. The increasing trend is consistent with the literature indicating that mean serum caffeine half-life for the healthy subjects is 5.7 hours.<sup>224</sup> On Day 2, despite a small leftover amount of theophylline as the baseline, the concentration change followed the same trend, proving the excellent reproducibility of this experiment. In Figure 5.5B, taurine demonstrated a shorter half-life, reaching the highest blood concentration at 1 hour and then starting to decrease to the initial level. Taurine plays a crucial role in the function of skeletal muscle,<sup>225</sup> and its concentration in skeletal muscle decreases after exercising.<sup>226</sup> Importantly, Manabe et al. studied rats and reported that the blood concentration increase of lactic acid after exercising became less significant with chronic treatment of taurine.<sup>227</sup> The accumulation of lactic acid in skeletal

muscle has long been noticed, though its relationship to the muscle fatigue is not clear yet.<sup>228</sup> It can be seen in Figure 5.5B that the blood levels of taurine and lactic acid changed in exactly the same manner. Our result suggests that even without the exercising and lactic acid accumulation in the skeletal muscles, taurine intake can prompt the skeletal muscle to release lactic acid and increase the blood lactic acid level. Sorbic acid is used as a preservative, and has very low level of toxicity.<sup>229</sup> In Figure 5.5C, its concentration significantly increased to the highest point within half an hour after the consumption of the energy drink. However, it had a remarkably high metabolism rate and decreased to a very low concentration within two hours. Another preservative, citric acid, only showed a minor concentration increase at 1 hour (Figure 5.5D). This can be expected since citric acid has multiple biological functions<sup>230-231</sup> and the energy drink was not the only source of it.

We also applied multi-variate analysis to study the dietary effects of the energy drink. Figure 5.6A and 5.6B are the PLS-DA score plots for the amine/phenol-submetabolome and carboxyl-submetabolome, respectively. On both of them, we can clearly see the separation from the 0-hour point to the 4-hour point. As the time increases, the distance to the initial state (0-hour) gets more significant. The  $Q^2$  value for the PLS-DA model of the dansyl-labeling data is 0.519, and the  $Q^2$  for the DMPA-labeling data is 0.574. They are relatively low because the energy drink only has a limited number of compounds and its impacts on the blood metabolome are limited in a limited number of pathways. Nonetheless, the changing trends successfully demonstrated the metabolome variations as responses to a dietary stimulation, and our method can effectively study this process.



**Figure 5.6** PLS-DA score plots, showing the statistical differences between samples collected at different time points in (A) amine/phenol-submetabolome and (B) carboxyl-submetabolome.

## 5.4 Conclusions

Overall, we have successfully applied the dansyl/DMPA-labeling LC-MS platform to the high-coverage metabolome analysis of one microliter of finger blood. Taking advantage of the high detection sensitivity enabled by the chemical isotope labeling, we can accurately quantify the blood concentrations of more than 4,000 metabolites from the extremely small amount of sample. The within-individual variations have been studied. Although these variations are not very significant, they should be considered in biomarker discovery studies, especially when the fold changes of biomarker candidates are not very large. In these applications, our method can become an ideal choice for the more accurate quantification of metabolites by taking the average of multiple measurements over a period. Furthermore, we have clearly demonstrated the metabolic responses to a dietary stimulation. Compared to traditional exposomics studies which focus on a limited number of target metabolites or target pathways, our method monitors how the whole metabolome

responds to environmental factors in a high-throughput and effective way. Future studies of these metabolic responses will significantly enrich our knowledge of the metabolic processes.

**Table 5.1** Average concentrations and RSDs of the 98 positively identified metabolites from one microliter of finger blood sample.

Retention time (min)	Detected m/z	Accurate mass (Da)	Compound name	HMDB ID	Average concentration	RSD (%)
2.16	375.0763	141.0179	O-Phosphoethanolamine	HMDB00224	2.36	40.01
2.48	403.1424	169.0840	3-methyl-histidine	HMDB00479	0.63	24.40
2.54	359.0736	125.0153	Taurine	HMDB00251	0.45	27.37
2.56	388.1063	154.0479	Hypoxanthine + H <sub>2</sub> O	HMDB00157	0.27	26.65
2.79	408.1694	174.1111	L-Arginine	HMDB00517	0.63	14.75
3.31	366.1117	132.0533	L-Asparagine	HMDB00168	0.66	15.33
3.38	422.1848	188.1265	Homo-L-arginine	HMDB00670	1.18	21.60
3.64	380.1280	146.0697	L-Glutamine	HMDB00641	0.46	14.46
3.83	409.1531	175.0948	Citrulline	HMDB00904	0.67	20.87
3.98	399.1028	165.0445	Methionine Sulfoxide	HMDB02005	0.80	34.20
4.13	307.1217	73.0634	Methylguanidine	HMDB01522	0.70	18.33
4.30	399.1033	165.0450	Methionine Sulfoxide - Isomer	HMDB02005 _2	0.82	36.33
4.47	339.1025	105.0441	L-Serine	HMDB00187	2.09	33.61
4.84	381.1118	147.0535	L-Glutamic Acid	HMDB00148	0.86	17.18
4.90	367.0966	133.0383	L-Aspartic Acid	HMDB00191	2.26	47.57
4.99	365.1162	131.0579	Trans-4-Hydroxyl-L- Proline	HMDB00725	0.35	19.60
5.44	353.1177	119.0594	L-Threonine	HMDB00167	0.83	24.36
5.49	395.1258	161.0675	Amino adipic acid	HMDB00510	0.57	16.57
5.70	295.1109	61.0525	Ethanolamine	HMDB00149	0.68	16.96

5.86	348.1016	114.0433	L-Asparagine - H2O	HMDB00168 _2	0.76	13.65
6.08	309.0909	75.0326	Glycine	HMDB00123	1.00	30.43
6.94	323.1068	89.0485	L-Alanine	HMDB00161	0.92	11.07
7.03	478.1264	244.0681	Uridine	HMDB00296	0.65	22.60
7.08	337.1216	103.0633	Gamma-Aminobutyric acid	HMDB00112	1.21	28.57
7.76	460.1165	226.0581	Uridine - H2O	HMDB00296 _2	0.66	21.67
7.92	337.1216	103.0633	3-Aminoisobutanoic acid	HMDB03911	1.58	19.64
8.00	386.0909	152.0326	Xanthine	HMDB00292	0.59	36.62
8.02	351.1371	117.0788	5-Aminopentanoic acid	HMDB03355	0.66	54.91
8.48	363.1019	129.0436	L-Glutamic Acid - H2O	HMDB00148 _2	1.03	17.33
8.62	337.1279	103.0695	D-Alpha-aminobutyric acid	HMDB00650	3.98	93.87
8.70	369.0930	135.0347	Methylcysteine	HMDB02108	0.81	24.38
9.29	349.1229	115.0646	L-Proline	HMDB00162	0.89	16.28
9.85	361.1321	127.0737	4-Guanidinobutanoic acid - H2O	HMDB03464 _2	0.76	21.04
9.90	383.1089	149.0506	L-Methionine	HMDB00696	0.75	43.39
9.92	351.1386	117.0803	L-Valine	HMDB00883	0.71	14.51
10.25	346.0851	112.0268	Uracil	HMDB00300	1.86	37.89
10.41	442.1448	208.0864	L-Kynurenine	HMDB00684	0.76	24.94
10.41	436.1888	202.1305	Alanyl-Leucine	HMDB28691	1.01	50.92
10.45	438.1480	204.0897	L-Tryptophan	HMDB00929	0.53	16.62
10.66	456.1576	222.0993	Glycyl-Phenylalanine	HMDB28848	0.54	34.78
11.71	399.1380	165.0797	L-Phenylalanine	HMDB00159	0.84	15.23
11.92	402.0982	168.0399	3-Hydroxymandelic acid	HMDB00750	0.78	28.94
11.96	365.1538	131.0955	L-Isoleucine	HMDB00172	0.79	22.11

12.11	363.1369	129.0786	L-Pipecolic acid	HMDB00716	0.44	13.20
12.22	365.1538	131.0955	L-leucine	HMDB00687	0.74	21.47
12.27	372.1006	138.0423	Urocanic acid	HMDB00301	4.93	34.89
12.99	354.0707	240.0247	L-Cystine	HMDB00192	2.87	15.79
13.12	551.2367	634.3567	Tryptophyl-Leucine	HMDB29087	1.79	54.25
13.14	416.1149	182.0566	Hydroxyphenyllactici acid	HMDB00755	0.63	18.42
14.12	372.0892	138.0308	Salicylic acid	HMDB01895	0.37	80.26
14.83	319.1108	85.0525	3-Aminoisobutanoic acid - H <sub>2</sub> O	HMDB03911 _2	1.33	18.27
14.87	385.1220	151.0637	Acetaminophen	HMDB01859	1.05	116.8 4
15.08	416.1152	182.0569	Homovanillic acid	HMDB00118	1.01	30.98
15.11	300.1037	132.0908	Ornithine	HMDB00214	1.66	41.23
15.18	386.1045	152.0462	3-Hydroxyphenylacetic acid	HMDB00440	0.93	18.53
15.46	386.1050	152.0467	3-Cresotinic acid	HMDB02390	0.97	25.07
15.93	402.0991	168.0408	Vanillic acid	HMDB00484	0.93	32.22
15.96	327.1146	93.0563	Aniline	HMDB03012	0.40	22.47
16.13	307.1114	146.1061	L-Lysine	HMDB00182	0.66	19.44
16.15	372.0895	138.0312	4-Hydroxybenzoic acid	HMDB00500	1.02	14.59
16.47	400.1205	166.0622	Desaminotyrosine	HMDB02199	0.56	77.54
16.83	398.1044	164.0461	m-Coumaric acid	HMDB01713	1.09	23.84
16.88	389.1276	155.0693	L-Histidine	HMDB00177	2.97	78.04
17.58	393.1848	159.1264	2-aminooctanoic acid	HMDB00991	0.82	31.94
17.70	395.1068	161.0485	Indole-3-carboxylic acid	HMDB03320	1.13	18.59
19.61	278.1078	88.0989	1,4-diaminobutane	HMDB01414	0.77	53.06
20.88	358.1108	124.0525	Guaiacol	HMDB01398	0.61	20.04
21.25	324.5960	181.0753	L-Tyrosine	HMDB00158	0.85	23.39
21.79	328.1004	94.0420	Phenol	HMDB00228	0.18	13.48

22.05	373.0839	139.0256	4-Nitrophenol	HMDB01232	0.17	37.00
23.21	342.1154	108.0570	o-Cresol	HMDB02055	0.06	40.80
25.10	289.0818	110.0470	pyrocatechol	HMDB00957	2.44	67.53
7.57	266.1386	104.0467	3-Hydroxybutyric acid	HMDB00357	0.42	35.40
7.98	266.1383	104.0464	Hydroxyisobutyric acid	HMDB00729	0.72	32.73
8.61	266.1380	104.0461	2-Hydroxybutyric acid	HMDB00008	1.19	20.39
11.65	300.1215	138.0296	3-Hydroxybenzoic acid	HMDB02466	0.83	18.65
15.24	298.1411	136.0492	Phenylacetic acid	HMDB00209	0.44	17.74
15.58	264.1586	102.0667	Isovaleric acid	HMDB00718	0.64	22.00
16.00	264.1591	102.0672	Valeric acid	HMDB00892	1.73	22.69
18.21	278.1732	116.0813	Isocaproic acid	HMDB00689	0.79	26.16
18.35	392.2407	230.1488	Dodecanedioic acid	HMDB00623	0.52	34.98
22.07	554.3847	392.2928	Ursodeoxycholic acid	HMDB00946	0.35	30.24
22.13	292.1893	130.0974	Heptanoic acid	HMDB00666	0.91	14.71
23.81	306.2039	144.1120	Valproic acid	HMDB01877	0.47	24.59
23.87	328.1894	166.0975	Perillic acid	HMDB04586	2.33	23.41
27.19	368.2192	206.1273	Ibuprofen	HMDB01925	0.19	56.10
28.96	554.3868	392.2949	Chenodeoxycholic acid	HMDB00518	0.43	29.83
29.19	553.3246	234.1625	Dodecanedioic acid	HMDB00623	0.87	28.91
31.03	476.3336	314.2417	Octadecanedioic acid	HMDB00782	0.47	50.94
31.18	360.2508	198.1589	5-Dodecenoic acid	HMDB00529	0.37	30.17
31.92	362.2673	200.1754	Dodecanoic acid	HMDB00638	0.37	30.77
32.33	444.3450	282.2531	Oleic acid	HMDB00207	1.46	19.82
32.43	305.1961	286.2084	Hexadecanedioic acid	HMDB00672	0.97	69.36
34.42	390.3002	228.2083	Myristic acid	HMDB00806	8.58	22.39
34.47	442.3334	280.2415	Linoleic acid	HMDB00673	0.26	54.57
34.86	404.3152	242.2233	Pentadecanoic acid	HMDB00826	0.10	33.50
34.91	468.3476	306.2557	Eicosatrienoic acid	HMDB02925	0.24	25.18
36.23	444.3481	282.2562	Vaccenic acid	HMDB03231	0.15	51.30

## Chapter 6

### Profiling Novel Metabolic Biomarkers for Parkinson's Disease Using In-depth Metabolomic Analysis

#### 6.1 Introduction

Parkinson's disease (PD) is a common progressive neurodegenerative disorder associated with the loss of dopaminergic neurons in the substantia nigra and production of Lewy bodies composed of  $\alpha$ -synuclein proteins.<sup>232</sup> Clinical PD diagnosis is based, in part, on impaired motor abilities such as bradykinesia, rigidity, tremor and postural instability. To date, no definitive single or set of biomarkers for PD have been discovered.<sup>233</sup> PD misdiagnosis rates were about 10% in 2001<sup>234</sup> and 6% in 2009<sup>235</sup>, rates which may depend on duration and stage of disease.<sup>236</sup> Therefore, the first goal of this research was to use systematic and unbiased metabolomics technology to detect a set of biomarkers that reflect disease pathways and might contribute to accurate discrimination. In addition to motoric impairment, up to 80% of PD patients eventually show cognitive impairment, including dementia, over the course of the disease, compromising quality of life and raising economic costs.<sup>237</sup> Therefore, the second goal of this research was to discover biomarkers that can be useful in discriminating PD patients who may remain dementia free (PDND) from those at risk for developing dementia. A better understanding of the metabolic pathways of PD patients with incipient dementia (PDID) can lead to improved disease monitoring and interventions targeted to patients according to dementia risk.

Metabolomics is an emerging field for biomarker discovery in human aging and neurodegenerative diseases.<sup>238</sup> A previous metabolomics study on cerebrospinal fluid (CSF) detected an increase in concentration of 3-hydroxykynurenine and decrease in concentration of glutathione in PD patients, suggesting the involvement of neurotoxicity and oxidative stress in PD pathogenesis.<sup>239</sup> Several potential biomarkers have also been identified in serum samples, which is less invasive to collect, including those involved in oxidative stress,<sup>240</sup> purine metabolism,<sup>241</sup> caffeine and xanthine metabolism.<sup>242</sup> However, because of the complexity of the metabolome, new technologies providing larger coverage and better quantitative capability will enable the discovery of more specific metabolic biomarkers, for both PD discrimination and early PDID sub-classification. Chemical isotope labeling liquid chromatography mass spectrometry (CIL LC-MS) that uses different labeling reagents to target chemical-group-based submetabolomes is a relatively new analytical platform for generating comprehensive and quantitative metabolomic profiles for biomarker research.<sup>70</sup> In the present study, we applied dansylation CIL LC-MS targeting the amine/phenol submetabolome to find the key metabolic differences in serum between two sets of groups. The overall study used a longitudinal design with baseline serum collection and two 18-month follow-ups. Using the baseline serum we performed two comparisons of metabolomics profiles: (1) PD patients and healthy controls and (2) PD patients who remained dementia-free for three years and PDID patients who developed dementia within this interval.

## 6.2 Methods

### 6.2.1 Participants

Clinically established PD patients ( $n = 52$ ) and age-and-sex matched healthy controls ( $n = 50$ ) between 64-84 years old volunteered for a 3-wave (18-month intervals) longitudinal study in Edmonton, Canada as previously detailed.<sup>243-244</sup> At baseline PD patients (1) met standard criteria for PD, (2) did not meet criteria for atypical parkinsonism, and (3) did not have unstable health conditions compromising survival. Patients who developed abnormal imaging such as a stroke or atypical features with follow up were also excluded. They were recruited from movement disorder clinics, the Parkinson's Society of Alberta, and from community neurologists. The control group was recruited by advertisement in seniors' centers and magazines, control and patient contacts, and general medicine clinics. The University of Alberta health ethics review board approved this study and all participants provided informed consent. For both groups, we excluded participants with baseline dementia, stroke, atypical parkinsonism, or attrition. In addition, one control participant was excluded as an outlier in the metabolomics analysis. Thus, there were 43 patients and 42 controls in the final groups.

We identified significant cognitive decline/dementia by caregiver and patient report of cognitive impairment in more than one cognitive domain that interfered with function. This was rated using the Clinical Dementia Rating (CDR), which was administered as a semi-structured interview. The investigators also used the Mini-Mental State Examination (MMSE), the Dementia Rating Scale (DRS), the Frontal Assessment Battery and the Short Orientation, Memory Concentration Test, which allowed assessment of a broad range of

cognitive function. This classification was highly correlated with cortical atrophy, suggesting validity. We did not use the Movement Disorder Society (MDS)-criteria, which were not available at the time of classification of subjects; however, we recently retrospectively examined all subjects using our CDR-based or DRS based classification and showed considerable overlap (57-78%) with classification using an independent neuropsychological battery (using cutoffs 1.5-2 SD suggested by MDS) (McDermott K, unpublished data, presented at the International Conference on Dementia, Banff, AB, May 2016). Our classification was more comprehensive taking into account all clinically available information.

Participants performed three waves of standardized assessments, including assessment for cognitive function and dementia. Of the 43 baseline PD patients 16 were diagnosed with dementia at wave 3. No further blood was taken and therefore the analyses (group comparison and prediction of dementia group at waver 3) were based on baseline blood work. Note that the exclusion described above was done before the metabolomic analysis. One outlier in the control group was excluded after the analysis, as its metabolomic data were too different from those of other participants. A possible reason of having this outlier could be contamination during sample collection.

We performed two pairwise metabolomics analyses. The first comparison evaluated the metabolomic profiles of the full available baseline groups, including 43 PD (no dementia at baseline) patients (M age = 70.71 years; sex = 44% female) and 42 controls (M age = 71.49 years; sex = 45% female) (see Table 6.1). The second comparison evaluated the

profiles for two PD subgroups, those who remained dementia-free at wave 3 ( $n = 27$ ; M age = 69.58 years) and those who were diagnosed with dementia at wave 3. The latter were classified post hoc as PDID at baseline ( $n = 16$ ; M age = 72.62 years). Duration of disease did not differ significantly between PDID ( $9.59 \pm 5.1$  years) and those remaining cognitively intact ( $7.75 \pm 4.1$ ,  $p=0.19$ ).

Baseline comparisons are shown in Table 6.1. PD patients did not differ from controls in age, education, sex distribution or cognitive status. The PDND and PDID subgroups were similar, differing slightly only on age and initial cognitive status. It is important to note that although there was minor statistical difference in the MMSE score, both means were above the impairment cut-off, which means there was no cognitive impairment in any of the patients at baseline. While they did not differ statistically on the key comparisons of levodopa equivalents and the Unified Parkinson's Disease Rating Scale (UPDRS), the average levodopa equivalent dose was higher in the PDID, despite slightly lower UPDRS part 3.

**Table 6.1** Baseline demographic and clinical characteristics.

	Control	PD	<i>p</i> -value	PDND	PDID	<i>p</i> -value
<i>N</i>	42	43	--	27	16	--
Age (years)	71.49 (5.01)	70.71 (4.14)	.434	69.58 (3.55)	72.62 (4.46)	.018
Education (years)	15.00 (3.42)	14.28 (2.98)	.303	14.74 (3.36)	13.50 (2.07)	.190
Sex (F/M)	19/23	19/24	.923	12/15	7/9	.966
MMSE	28.56 (1.48)	28.33 (1.67)	.503	28.85 (1.29)	27.36 (1.91)	.006
Folate (nmol/L)	879.81 (236.25)	842.56 (207.45)	.442	799.00 (185.45)	916.06 (227.41)	.073
Vitamin B12 (pmol/L)	393.79 (198.05)	293.26 (112.82)	.005	295.70 (92.16)	289.13 (144.52)	.856
Levodopa equivalents (mg)	N/A	644.00 (360.06)	--	611.83 (392.94)	703.76 (293.41)	.448
UPDRS part 3	N/A	16.12 (7.90)	--	16.67 (8.08)	15.19 (7.77)	.559

*Note.* PD, Parkinson's disease; PDND, Parkinson's disease no dementia; PDID, Parkinson's disease incipient dementia; MMSE, Mini Mental State Exam; UPDRS, Unified Parkinson's Disease Rating Scale. Standard deviations are in parentheses.

## 6.2.2 Serum samples and dansylation LC-MS metabolomic profiling

Serum was collected at baseline only from all participants and stored at -80 °C. To reveal small concentration variations of metabolites in comparative samples, we applied the CIL LC-MS technique to overcome potential inaccuracy due to matrix effects, ion suppression, or instrumental drift in MS detection. In our workflow, individual samples were labeled using <sup>12</sup>C-dansyl chloride and a pooled sample generated by mixing small aliquots of samples was labeled by <sup>13</sup>C-dansyl chloride. Each <sup>12</sup>C-labeled sample was mixed with an aliquot of <sup>13</sup>C-pooled sample, followed by LC-MS analysis of the mixture. All the labeled metabolites were detected as <sup>13</sup>C- and <sup>12</sup>C-peak pairs and the peak ratios were determined and used for quantitative metabolomic analysis of the individual samples.

We minimized variations in total sample amount in different samples in order to detect the individual metabolite concentration differences in comparative samples more accurately by performing sample normalization. Specifically, we applied an LC-UV method to determine the total concentration of dansyl-labeled metabolites based on the UV absorption of the dansyl group.<sup>115</sup> Before the LC-MS analysis, we mixed the <sup>12</sup>C-labeled individual sample with the same total amount of <sup>13</sup>C-labeled pool, according to the total concentrations of labeled metabolites, for sample amount normalization.

For LC-MS, an Agilent 1100 series binary system (Agilent, Palo Alto, CA) and an Agilent reversed-phase Eclipse plus C18 column (2.1 mm×100 mm, 1.8 μm particle size, 95 Å pore size) were used. LC solvent A was 0.1% (v/v) formic acid in 5% (v/v) ACN/H<sub>2</sub>O, and solvent B was 0.1% (v/v) formic acid in ACN. The gradient elution profile was as follows: t = 0 min, 20% B; t = 3.5 min, 35% B; t = 18.0 min, 65% B; t = 24.0 min, 99% B; t = 28.0 min, 99% B. The flow rate was 180 μL/min. After one injection, the column was re-equilibrated with the initial mobile phase conditions for 15 min before injecting the next sample. The flow was loaded to the electrospray ionization (ESI) source of a Bruker maXis impact high-resolution quadrupole time-of-flight (Q-TOF) mass spectrometer (Bruker, Billerica, MA). All MS spectra were obtained in the positive ion mode. According to the UV-quantification result, 1.5 nmol of each labeled and mixed samples were injected into the LC-MS system.

### **6.2.3 Data processing and statistical analysis**

The  $^{12}\text{C}$ -/ $^{13}\text{C}$ -peak pairs from each LC-MS run were extracted by the IsoMS software.<sup>72</sup> IsoMS-Align was used to align the peak pair data from different samples by retention time and accurate mass. The missing ratio values were filled back by using the Zero-fill program.<sup>137</sup> IsoMS-Quant<sup>157</sup> was used to generate the final metabolite-intensity table, which was exported to SIMCA-P+ 12 (Umetrics AB, Umeå, Sweden) for analysis. We followed the method described in the work of LeWitt et al.<sup>239</sup> for statistical analysis. To avoid over-fitting, we calculated the q-value for each p-value using QVALUE.<sup>110</sup> Metabolite identification was done using a DnsID standards library<sup>159</sup> for positive identification as well as the HMDB library and EML database for putative identification.<sup>41</sup>

## **6.3 Results**

### **6.3.1 Submetabolome and metabolite identification**

Dansylation labeling LC-MS targets the analysis of the amine/phenol submetabolome; many metabolomic pathways contain the amine- and phenol-containing metabolites. A total of 719 metabolites were commonly detected in 80% of the 85 samples. Among them, 66 metabolites were positively identified using an in-house developed dansyl standards library consisting of 273 compounds. For the remaining peak pairs, accurate mass search with a mass accuracy tolerance of 0.005 Da putatively identified 333 metabolites using the HMDB database and 282 metabolites in the EML database using MyCompoundID.<sup>41</sup> In total, 681 of the 719 metabolites (95%) were either definitely or putatively identified.

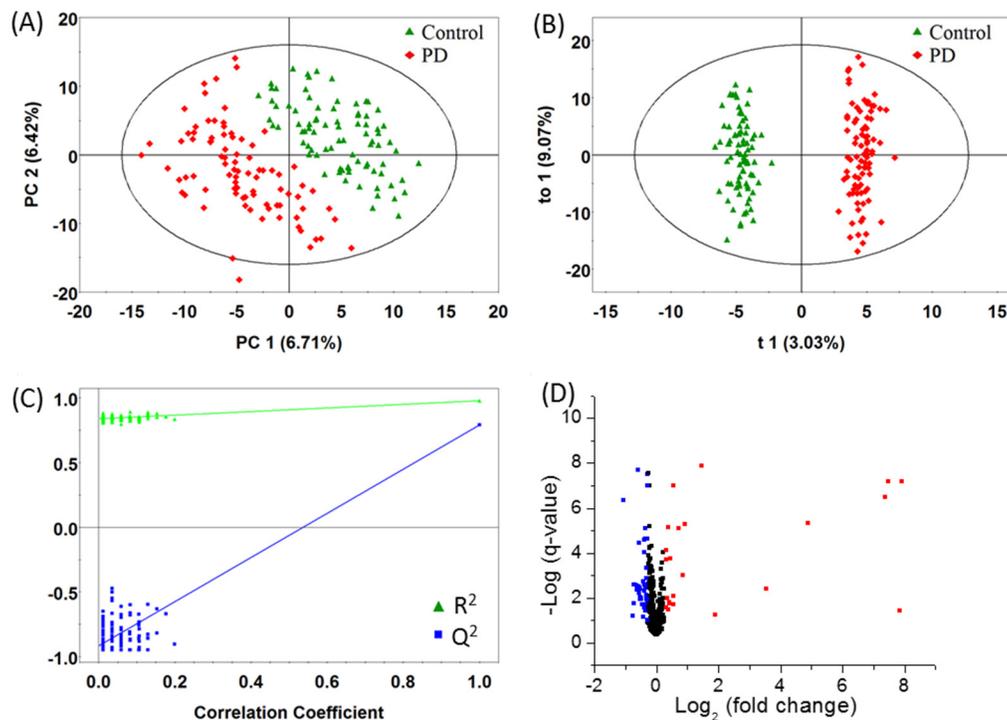
### 6.3.2 Comparative metabolome analysis for PD biomarker discovery

The Partial Least Squares-Discriminant Analysis (PLS-DA) and Orthogonal Partial Least Squares-Discriminant Analysis (OPLS-DA) score plots are shown in Figure 6.1. Figure 6.1A and Figure 6.1B show that there was a significant difference between the healthy controls (in green) and the PD patients (in red). This group separation was validated in a permutation test (Figure 6.1C). Note that in our data analysis the PLS and OPLS methods were used as multivariate calibration processes, and the resulting score plots helped with visualizing the inter-sample and inter-group variances. For the multivariate classifications (Control vs. PD, or PDND vs. PDID), we used random forest analyses. The model performance indicators (the  $R^2$  and  $Q^2$  values) are provided in the corresponding score plots.

The volcano plot shown in Figure 6.1D displays 28 metabolites with Fold Change (FC) > 1.2,  $q < 0.1$  (in red) and 48 metabolites with  $FC < 0.83$ ,  $q < 0.1$  (in blue). Table 6.2 lists the significant metabolites identified. Four significant metabolites, identified as 4-hydroxy-benzenepropanedioate, vanillylmandelic acid-isomer, alpha-methyldopa, and methylguanidine, were observed only in the PD group. The dansyl library identified citrulline, methionine sulfoxide, pantothenic acid, glycyl-valine, pipercolic acid, serotonin, vanillic acid, vanillylmandelic acid, theophylline, and hydroxykynurenine.

Among significant metabolites, there were several catecholamine metabolites. Some metabolites from the tryptophan pathway were also detected and quantified. Biopterin had an increased concentration in the PD group. Some significant metabolites in Table 6.2 were

related to oxidative stress, such as citrulline and methionine sulfoxide. Finally, caffeine metabolites were detected as altered in PD.



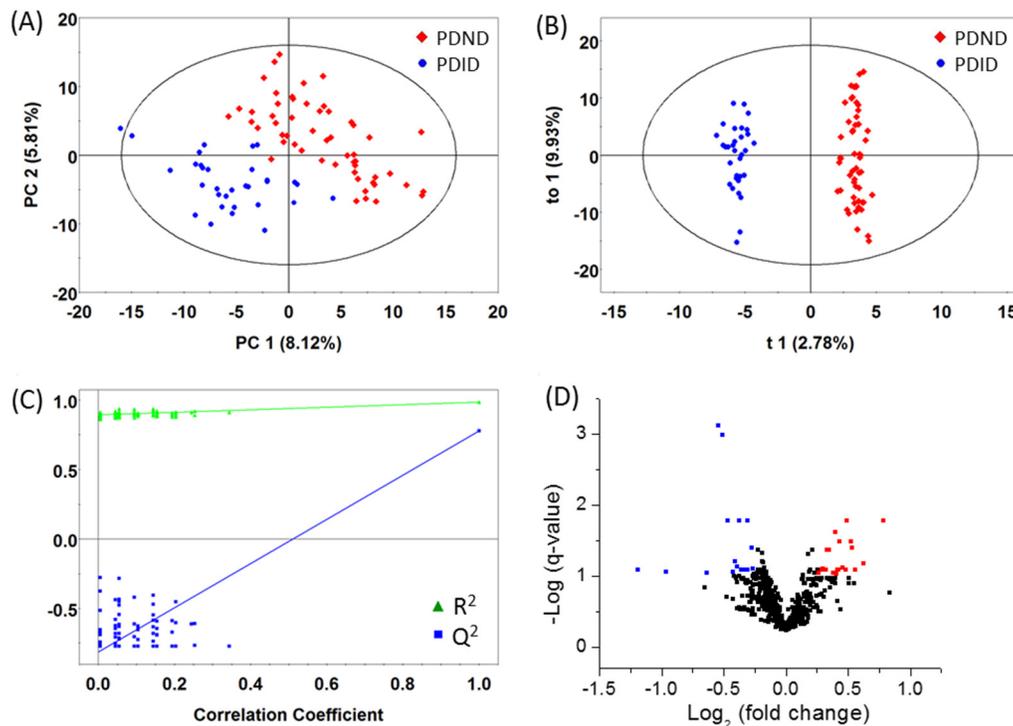
**Figure 6.1** (A) PLS-DA and (B) OPLS-DA score plots of dansylation LC-MS data obtained from 42 healthy controls (in green) and 43 PD patients (in red). (“PC” represents for “principal component” and the corresponding percentage is the percentage of the variance among all the data points that this principal component covers.)  $R^2$  and  $Q^2$  values given by cross-validation are: 0.977 and 0.791 for PLS-DA; 0.974 and 0.866 for OPLS-DA. (C) Response permutation test result of the PLS-DA model in Figure 6.1A. (D) Volcano plot of the comparison between healthy control and PD showing 28 variables with  $FC > 1.2$ ,  $q < 0.1$  (in red) and 48 variables with  $FC < 0.83$ ,  $q < 0.1$  (in blue).

### 6.3.3 Comparative metabolome analysis of PD with and without incipient dementia

Figure 6.2A and 6.2B show the PLS-DA and OPLS-DA score plots, respectively, for the comparison between PDND and PDID. The permutation test result is shown in Figure 6.2C.

These analyses indicate that the two groups were separated based on the metabolomic data set. Post-hoc comparisons showed the metabolic data discriminated subgroups of dementia severity but not age.

The volcano plot (Figure 6.2D) shows 21 metabolites with  $FC > 1.2$ ,  $q < 0.1$  (in red) and 15 metabolites with  $FC < 0.83$ ,  $q < 0.1$  (in blue). Among the 36 significant metabolites, 16 were identified (see Table 6.3 for the list). Among these 16 metabolites, two metabolites were definitely identified by the library and the others were putatively identified based on a database search. The two definitely identified metabolites are desaminotyrosine and 5-hydroxylysine.



**Figure 6.2** (A) PLS-DA and (B) OPLS-DA score plots of dansylation LC-MS data obtained from 27 PD patients without dementia (PDND in red) and 16 PD patients with incipient dementia (PDID

in blue).  $R^2$  and  $Q^2$  values given by cross-validation are: 0.974 and 0.866 for PLS-DA; 0.982 and 0.813 for OPLS-DA. (C) Response permutation test result of the PLS-DA model in Figure 6.2A. (D) Volcano plot of the comparison between the PDND subgroup and the PDID subgroup showing 21 variables with  $FC > 1.2$ ,  $q < 0.1$  (in red) and 15 variables with  $FC < 0.83$ ,  $q < 0.1$  (in blue).

### **6.3.4 Common discriminating metabolites**

Comparison of the metabolites listed in Table 6.2 and Table 6.3 indicate that there are five common significant metabolites. For example, 3, 4-dihydroxyphenylacetone, a significant metabolite in the PDND-PDID comparison (Table 6.3), also significantly differentiated PD patients from the healthy controls (Table 6.2). For this metabolite, the averaged peak pair ratio in the control group was  $0.006 \pm 0.012$ . In the PDND subgroup and the PDID subgroup, the averages were  $2.16 \pm 1.63$  and  $3.34 \pm 1.95$ , respectively.

### **6.3.5 ROC curves**

Receiver operating characteristic (ROC) curves are graphical plots that illustrate the performance of a binary classifier system as its discrimination threshold is varied. In our work, ROC analysis was used to show to diagnosis power of one or a group of metabolite candidates. There are several classification models available for building an ROC curve, such as the PLS, random forest, support vector machines, etc. Among them, we have chosen the random forest method for our analysis. Metaboanalyst 3.0 was used to generate the ROC curves for differentiating PD (baseline) from healthy controls and differentiating between PDND and PDID.

For building the ROC curves, we did not simply input all 709 variables into the model and then let it find the best five biomarker candidates to use in building the ROC curve. Instead, we manually selected the candidates based on standard procedures. To select the potential biomarkers for differentiating PD vs. control, 5 metabolites were selected according to the following criteria: (1) a large fold change, (2) a p-value of smaller than 0.05 and (3) metabolite identification with high confidence (positive IDs were preferred as they could be immediately used as biomarkers if they could be validated in future studies using large cohorts of samples).

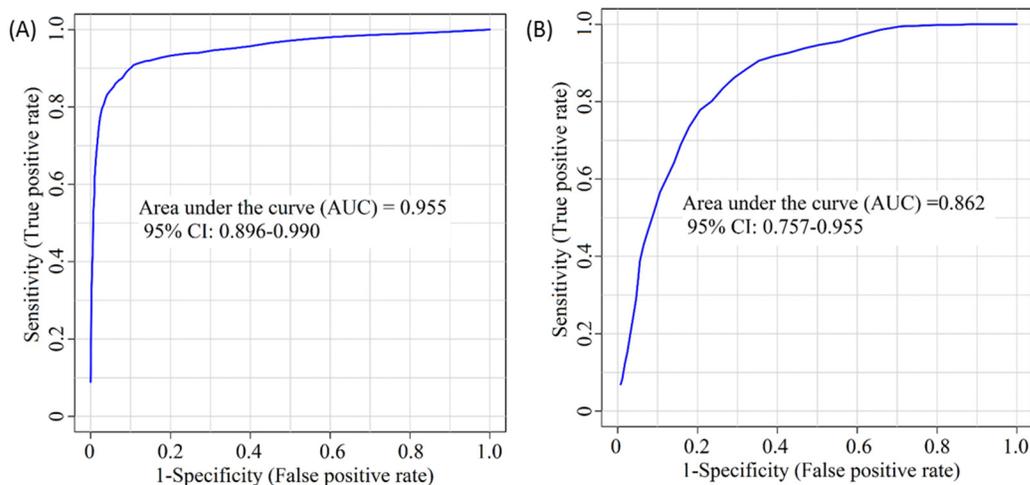
The metabolites were first ranked according to their fold change. There were seven putatively identified significant metabolites with fold changes of larger than 10, and each of these generated an ROC curve with very high AUC (the lowest one was 0.914). If these metabolites could be positively identified, they would work as very strong biomarkers for the diagnosis of PD. In order to use more positively identified metabolites to build the classification model, the positively identified metabolites were selected from the top 20 ranked metabolites (according to fold change). We excluded vanillylmandelic acid as it is potentially related to catecholamine metabolism, which can be affected by the DOPA medication. The three selected metabolites were vanillic acid, hydroxykynurenine and theophylline. Since the discrimination power of these three positively identified metabolites were not strong enough (AUC=0.944, but we wanted to go higher than 0.95), two putatively identified metabolites with large fold changes and high-confidence putative identifications, isoleucyl-alanine and 5-acetylamino-6-amino-3-methyluracil, were added to make a 5-metabolites biomarker panel. The ROC curve of this panel of 5 metabolites

produced an AUC value of 0.955. Adding more metabolites to this panel resulted in only minor increase in sensitivity and specificity. Considering that in real world clinical applications a panel of a few metabolites is preferred over a panel of many metabolites, we concluded that this panel of 5 metabolites was sufficient to illustrate the overall performance of the metabolic profiles for differentiating PD vs. Control.

We note that using CIL LC-MS there is no need to select the panel for targeted analysis in a validation study. CIL LC-MS is a quantitative method where thousands of metabolites can be analyzed in one LC-MS run. Thus, the selected biomarkers discussed in our work are only used to illustrate the separation performance. These biomarkers, along with all other labeled metabolites, will be monitored in our future validation studies. This approach will increase the likelihood of finding the high performance biomarkers which may be different from the initial biomarkers determined using a small cohort of discovery samples.

In summary, the classification model was built by the random forest method based on five metabolites: vanillic acid, 3-hydroxykynurenine, isoleucyl-alanine, 5-acetylamino-6-amino-3-methyluracil and theophylline. The AUC value for each of the metabolites separately was found to be 0.939, 0.781, 0.794, 0.730 and 0.714, respectively. Of interest, using vanillic acid alone we could achieve both sensitivity and specificity at 90.0%. The discriminating power was improved by combining these five metabolites into a biomarker panel. The corresponding ROC curve shown in Figure 6.3A produced an AUC value of 0.955, which is within the range of 0.896–0.990 at the 95% confidence interval.

Discrimination of PD from controls can be achieved at 87.5% sensitivity and 92.0% specificity. Using a permutation test, we did not find any over-fitting of the ROC results.



**Figure 6.3** (A) The receiver operating characteristic curve generated by the random forest model using the following 5 metabolite biomarker candidates: vanillic acid, 3-hydroxykynurenine, isoleucyl-alanine, 5-acetylamino-6-amino-3-methyluracil, and theophylline. (B) The receiver operating characteristic curve generated by the random forest model using the following 8 metabolite biomarker candidates: His-Asn-Asp-Ser, 3, 4-dihydroxy-phenylacetone, desaminotyrosine, hydroxy-isoleucine, alanylalanine, putrescine [-2H], purine [+O], and its riboside.

For differentiating the PDND and PDID subgroups at baseline, we found that the discrimination power of a panel based on the two definitely identified metabolites was not strong (AUC=0.673). The univariate AUC of 5-hydroxylysine is 0.659 and that of desaminotyrosine is 0.674. Thus, we selected the following eight putatively identified biomarker candidates with highly ranked independent AUCs to form a panel: His-Asn-Asp-Ser (AUC=0.597), 3, 4-dihydroxyphenylacetone (0.677), desaminotyrosine (0.640), hydroxy-isoleucine (0.610), alanyl-alanine (0.737), putrescine [-2H] (0.736), purine [+O] (0.627) and its riboside (0.597). As can be seen in Figure 6.3B, the ROC curve for this

biomarker panel produced an AUC value of 0.862, which is within the range of 0.757–0.955 at the 95% confidence interval. This panel can provide discrimination with sensitivity of 80.0% and specificity of 77.0%.

We have examined whether the use of any unidentified significant metabolites (there were 38 in total) could increase the differentiation power. We found that there was no significant increase in performance when one or more of these were used to replace the panels described above.

**Table 6.2** List of 46 identified significant metabolites found in human serum samples that differentiate the PD group and the healthy control group.

Retention time (min)	Mass of dansylated metabolite (Da)	Mass of metabolite (Da)	Metabolite	HMDB ID	ID	Fold change	q-value
2.65	432.1329	198.0746	5-Acetylamino-6-amino-3-methyluracil	HMDB04400	Putative	0.61	2.47E-03
3.20	502.1385	268.0802	Inosine	HMDB00195	Putative	0.74	6.86E-02
3.58	471.1422	237.0839	Biopterin	HMDB00468	Putative	1.36	1.74E-04
3.90	409.1534	175.0951	Citrulline	HMDB00904	Library	0.83	2.23E-05
4.08	399.1035	165.0452	Methionine Sulfoxide	HMDB02005	Library	1.21	2.55E-02
4.60	399.1037	165.0454	Methionine Sulfoxide - Isomer	HMDB02005	Library	1.24	2.68E-02
4.90	410.1372	176.0789	Ornithine [+CO <sub>2</sub> ]	HMDB00214	Putative	0.78	7.71E-06
6.81	415.1316	181.0733	L-Threo-3-Phenylserine	HMDB02184	Putative	0.76	9.15E-05
7.40	436.1894	202.1311	Isoleucyl-Alanine	HMDB28900	Putative	0.66	1.92E-08
7.63	309.1269	75.0685	1-Amino-propan-2-ol	HMDB12136	Putative	1.90	5.32E-06
7.75	465.1792	231.1209	Norphthalmic acid [-CO <sub>2</sub> ]	HMDB05766	Putative	1.34	1.51E-02
8.21	415.1314	181.0731	L-Threo-3-Phenylserine	HMDB02184	Putative	0.82	2.04E-02
8.56	453.1683	219.1100	Pantothenic acid	HMDB00210	Library	0.76	1.95E-02
9.34	408.1574	174.0991	Glycyl-Valine	HMDB28854	Library	1.45	1.91E-02

9.52	307.1116	73.0533	Aminoacetone	HMDB02134	Putative	1.48	1.02E-07
11.58	430.0946	196.0363	4-Hydroxy-benzenepropanedioate	HMDB59809	Putative	179.70	6.28E-08
11.77	400.1053	166.0470	Methylxanthine	HMDB10738	Putative	0.62	2.46E-03
12.02	400.1066	166.0483	3-Methylxanthine	HMDB01866	Putative	0.65	4.44E-03
12.90	382.5800	297.0434	L-Cysteinylglycine disulfide	HMDB00709	Putative	0.83	9.87E-08
13.27	432.1101	198.0518	Vanillylmandelic acid	HMDB00291	Library	3.80	2.07E-19
13.66	363.1361	129.0778	L-Pipecolic acid	HMDB00716	Library	0.80	1.93E-02
14.61	434.1739	200.1156	Glycylproline [+C2H4]	HMDB00721	Putative	0.82	1.18E-02
15.05	432.1101	198.0518	Isovanillylmandelic acid	NA	Putative	164.05	3.25E-07
15.87	414.1232	180.0649	Theophylline	HMDB01889	Library	0.48	4.42E-07
17.58	355.6188	243.1210	Aspartyllysine [-H2O]	HMDB04985	Putative	1.24	2.07E-04
17.77	372.0888	138.0305	4-Hydroxybenzoic acid	HMDB00500	Putative	0.80	7.15E-03
17.88	402.1001	168.0418	Vanillic acid	HMDB00484	Library	3.48	3.96E-20
20.43	414.1364	180.0781	p-Hydroxyphenylacetic acid [+C2H4]	HMDB00020	Putative	3.75	5.91E-02
20.81	333.1627	99.1044	Cyclohexylamine	HMDB31404	Putative	0.60	1.71E-02
20.82	539.3191	305.2608	Capsaicin	HMDB02227	Putative	1.65	7.89E-06
21.00	314.1180	160.1193	N(6)-Methyllysine	HMDB02038	Putative	0.65	2.91E-03
21.01	315.1180	162.1193	Tryptamine [+2H]	HMDB00303	Putative	0.68	6.05E-03
21.75	386.1039	152.0456	Vanillin	HMDB12308	Putative	1.27	1.01E-02
22.21	338.1026	208.0885	5-Hydroxyindoleacetic acid [+NH3]	HMDB00763	Putative	0.61	2.55E-19
22.28	679.1967	212.0873	Histidinyl-Glycine	HMDB28885	Putative	29.32	4.66E-06
23.42	346.0986	224.0805	Hydroxykynurenine	HMDB00732	Library	2.73	1.30E-08
23.75	328.0712	188.0259	L-Homocysteine sulfonic acid [+NH3]	HMDB02238	Putative	0.82	1.23E-02
24.23	310.0751	152.0336	6,8-Dihydroxypurine	HMDB01182	Putative	0.59	6.15E-02
24.38	302.0852	136.0538	Dopamine [-NH3]	HMDB00073	Putative	0.75	1.76E-03
24.38	342.1314	216.1461	Valyl-Valine	NA	Putative	0.82	9.95E-02
24.69	302.0841	136.0515	Dopamine [-NH3]	HMDB00074	Putative	0.83	1.70E-02
25.06	322.1050	176.0933	Serotonin	HMDB00259	Library	0.79	1.07E-02
27.48	304.4142	211.0677	a-Methyl dopa	HMDB11754	Putative	231.42	3.64E-02
27.54	304.0826	140.0485	Vanillic acid [-CO]	HMDB00484	Putative	11.77	3.81E-03
27.74	399.1134	165.0551	Methylguanine	HMDB03282	Putative	240.60	6.28E-08
27.76	317.0911	166.0655	3,4-Dihydroxy-phenylacetone	HMDB31132	Putative	208.33	1.38E-22

**Table 6.3** List of 18 identified significant metabolites found in human serum samples that differentiate the PDND subgroup and the PDID subgroup.

Retention time (min)	Mass of dansylated metabolite (Da)	Mass of metabolite (Da)	Metabolite	HMDB ID	ID	Fold change	q value
3.20	502.1385	268.0802	Riboside of Purine [+O]	NA	Putative	0.51	8.63E-02
10.38	381.1466	147.0883	Hydroxyisoleucine	NA	Putative	1.47	8.20E-02
11.58	430.0946	196.0363	4-Hydroxy-benzenepropanedioate	HMDB59809	Putative	1.33	9.22E-02
12.87	363.1360	129.0777	L-Proline [+CH2]	HMDB00612	Putative	0.81	8.17E-02
13.74	351.1000	117.0417	L-2-Amino-3-oxobutanoic acid	HMDB06454	Putative	1.35	3.32E-02
14.18	315.1080	162.0995	5-Hydroxylysine	HMDB00450	Library	1.23	7.94E-02
16.71	727.2133	493.1550	His-Asn-Asp-Ser	NA	Putative	1.72	1.64E-02
18.39	427.1304	193.0721	Phenylacetyl glycine	HMDB00821	Putative	1.27	4.33E-02
18.41	400.1196	166.0612	Desaminotyrosine	HMDB02199	Library	0.64	8.91E-02
19.09	370.0962	136.0379	Purine [+O]	HMDB01366	Putative	0.44	8.17E-02
20.83	314.1014	160.0862	Alanyl-alanine	NA	Putative	1.44	3.28E-02
21.43	321.1081	174.0995	Ornithine [+C2H2O]	HMDB00214	Putative	1.37	7.55E-02
21.43	322.1082	176.0997	Tryptamine [+O]	HMDB00303	Putative	1.33	8.20E-02
22.42	303.1154	69.0571	1-Pyrroline-2-carboxylic acid [-CO2]	HMDB06875	Putative	1.32	2.38E-02
22.43	277.1000	86.0834	Putrescine [-2H]	HMDB01414	Putative	1.41	1.64E-02
24.03	284.1077	100.0988	Cadaverine [-2H]	HMDB02322	Putative	1.20	8.17E-02
24.92	338.5921	209.0675	Hydroxyphenylacetyl glycine	HMDB00735	Putative	1.23	8.17E-02
27.76	317.0911	166.0655	3,4-Dihydroxyphenylacetone	HMDB31132	Putative	1.54	6.65E-02

## 6.4 Discussion

We performed two two-group metabolomic comparisons for identifying novel biomarkers of PD. First, using baseline serum samples and clinical characterizations and diagnoses, we compared a PD group with a comparable older adult control group (no PD, dementia, or impairment). Our results showed clear differentiation between groups. A panel of 5 metabolites in the ROC analysis gave an AUC of 0.955 with sensitivity of 87.5% and specificity of 93.0%. Second, we compared two subgroups of the initial PD group, again using baseline serum samples. PD patients who have yet to meet the criteria for a dementia diagnosis represent a detectably early and more advanced transitional phase of PD. Chia and colleagues observed metabolomic changes in the cerebrospinal fluid of PD patients already diagnosed with dementia and depression.<sup>245</sup> Notably, these subgroups were determined three years after the baseline clinical evaluation did not detect incipient dementia. The present task met the challenge of detecting biomarkers of PD dementia prior to clinical diagnosis. Specifically, using a panel of 8 metabolites in ROC analysis, we obtained an AUC of 0.862 with sensitivity of 80.0% and specificity of 77.0%. These results hold promise for PD discrimination and prognosis, as well as identification of pathways leading to dementia within PD patient groups that may assist in identifying targets for intervention. While the markers we have identified are not in presymptomatic or early patients, patients with established disease are at risk for cognitive decline and dementia, features which are critically relevant to clinical decisions regarding future planning and treatment options (such as having deep brain stimulation).

Although currently there is no metabolite biomarker has moved to the stage of clinical practice, many studies have reported promising biomarker candidates for the diagnosis of PD. For example, LeWitt et al. developed a 19-metabolite panel providing sensitivity of 83% and specificity of 91%.<sup>23</sup> Among these reports, our work also differentiated PD with very good diagnostic power. Considering the outstanding performance of these biomarker panels, in the future, we may combine the most significant and commonly existing metabolites from them to establish a more powerful biomarker panel with well-understood biological meanings. Importantly, our work contributes to the prediction of dementia during the progress of PD. According to our knowledge, there is no other cohort study that has achieved this goal. In clinical testing, screening examinations mentioned above are currently being used for the diagnosis of dementia. The AUC of the MMSE test remains low at 0.76, with sensitivity of 67% and specificity of 85%.<sup>246</sup> In the future, the application of biomarkers will definitely boost the diagnosis or prediction of dementia.

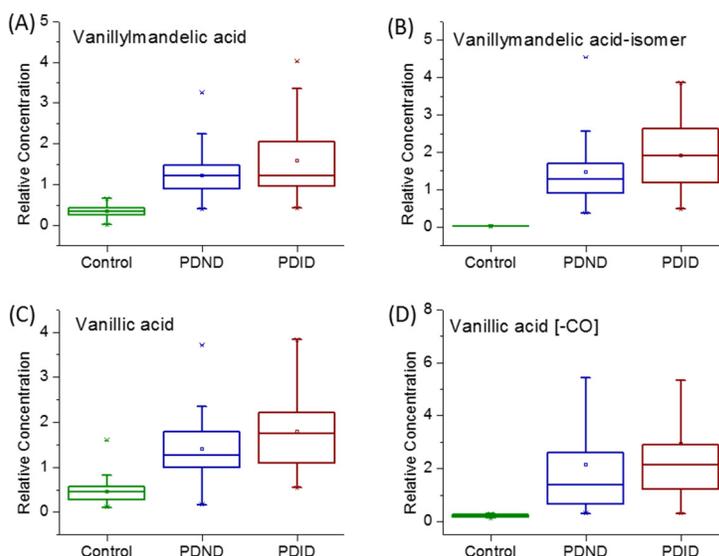
As Gerlach and colleagues suggested, biomarkers for PD should not only be linked to fundamental features of PD neuropathology, but also be correlated to the disease progression assessed by clinical rating scales.<sup>247</sup> Our approach targeted the amine/phenol-containing metabolites; other metabolites would require other labeling approaches. For the first set of results, the pathway analysis showed that catecholamine metabolism, tryptophan metabolism and caffeine metabolism were the most relevant metabolomics pathways found to be affected by PD. In addition, metabolites related to oxidative stress were also identified.

### 6.4.1 Catecholamine metabolism

The metabolites in the catecholamine pathway were excluded in order to avoid the interference of levodopa medication. However, vanillylmandelic acid (VMA), which is the end product of the catecholamine pathway, was retained as a significantly changed metabolite, because some studies showed that urinary excretions of VMA and its up-stream metabolites (epinephrine and normetanehrine) were not greatly affected by levodopa.<sup>248-249</sup> The fold change of VMA between the PD group and the control group was as large as 3.80, suggesting a significant metabolic change caused by PD rather than medication. In addition, the isomer of VMA was detected only in the PD group. Moreover, the vanillylmandelic acid-isomer showed very similar fold changes, although its origin and biological function is unclear.

This relationship between VMA and PD has not been previously reported; this change is likely related to a disorder of the catecholamine pathway, in which dopamine is produced. Vanillic acid is a food metabolite, which is often found in the urine of humans who have consumed coffee, tea and vanilla-flavored food.<sup>250</sup> It may have an increased concentration in PD via conversion from VMA or homovanillic acid. The derivative of vanillic acid, vanillic acid [-CO], was also a unique metabolite in the PD group. The relative concentrations of these metabolites in the control group, the PDND subgroup and the PDID subgroup are shown as box plots in Figure 6.4. Despite the fact that all four have an increased concentration in the serum of PD patients, the relative concentrations in the PDID subgroup were higher than those in the PDND subgroup. These results support our previous

discussion that a more advanced transitional phase of PD can be detected in PD patients biologically, even if not yet clinically.

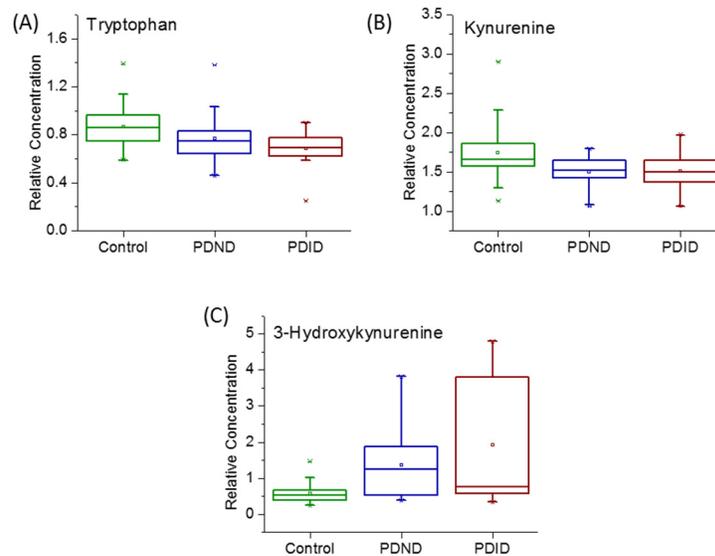


**Figure 6.4** Box plots of the relative concentrations of vanillylmandelic acid (A), vanillylmandelic acid-isomer (B), vanillic acid (C), and vanillic acid (D) in the control group, the Parkinson’s disease with no dementia subgroup, and the PD patients with incipient dementia subgroup.

Both methyl dopa and 3, 4-dihydroxyphenylacetone were increased in the PD samples. Methyl dopa is known to be an aromatic-amino-acid decarboxylase (AADC) inhibitor in animals and in humans.<sup>251</sup> Declining levels of AADC may contribute to decreasing effectiveness of L-dopa medication over time.<sup>252</sup> A metabolite of methyl dopa via AADC,<sup>253</sup> 3, 4-dihydroxyphenylacetone, was greatly increased in concentration in the PD group, but its role is unknown.

## 6.4.2 Tryptophan metabolism

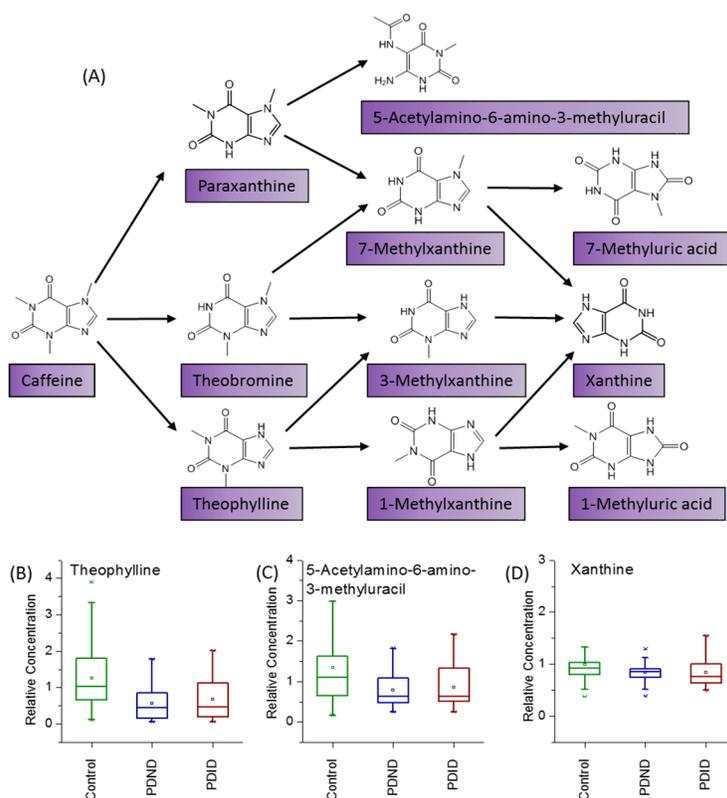
Although the concentrations of tryptophan and kynurenine did not change greatly, increased relative concentration of 3-hydroxykynurenine in the PD group was observed. Figure 6.5 shows the relative concentrations of tryptophan, kynurenine and 3-hydroxykynurenine in different groups. Metabolites of the kynurenine pathway are believed to play crucial roles in maintaining normal brain function.<sup>254</sup> Kynurenine is the down-stream metabolite of tryptophan that can be further converted to 3-hydroxykynurenine or kynurenic acid. 3-hydroxykynurenine is a neurotoxic metabolite which causes neuronal death.<sup>255</sup> However, kynurenic acid behaves as an endogenous neuroprotective agent.<sup>256</sup> Consistent with our results, LeWitt and colleagues reported that the CSF concentration of 3-hydroxykynurenine was increased in PD patients.<sup>239</sup>



**Figure 6.5** Box plots of the relative concentrations of tryptophan (A), kynurenine (B), and 3-hydroxykynurenine (C) in the control group, the Parkinson's disease with no dementia subgroup, and the PD patients with incipient dementia subgroup.

### 6.4.3 Caffeine metabolism

Although caffeine cannot be labeled by the dansylation reagent, some of its metabolites are labeled and were detected as significantly changed metabolites (Figure 6.6A). Theophylline can be labeled by the dansylation reagent. Paraxanthine cannot be labeled, but its down-stream metabolite, 5-acetylamino-6-amino-3-methyluracil, was detected as a significant metabolite. Although methylxanthine was detected and was in the VIP list (Table 6.2), we could not differentiate its three isomers without standards. Xanthine is the end product of caffeine metabolism. It was detected and identified by the dansyl standards library. Figures 6.6B-D show the relative concentrations of theophylline, 5-acetylamino-6-amino-3-methyluracil and xanthine in different groups. The concentrations of caffeine metabolites were lower in the PD group. Xanthine can also be converted from hypoxanthine and guanine in the purine pathway, and its concentration changed marginally.



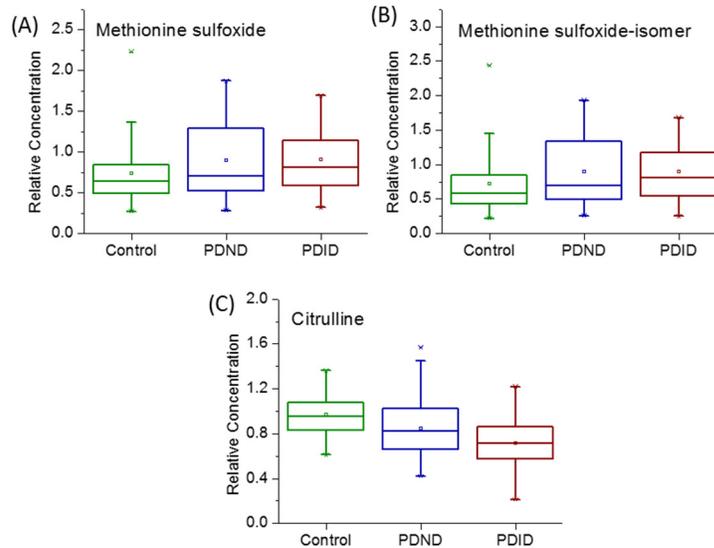
**Figure 6.6** (A) A simplified schematic of the caffeine metabolism pathway. Box plots of the relative concentrations of theophylline (B), 5-acetylamino-6-amino-3-methyluracil (C) and xanthine (D) in the control group, the PDND subgroup and the PDID subgroup.

It has been widely reported that coffee and tea consumption could protect against the risk of PD.<sup>257</sup> A possible reason for this phenomenon is that caffeine is an antagonist of the adenosine A2A receptor, which has a role in the regulation of glutamate and dopamine release.<sup>258</sup> Adenosine A2A antagonists can modify motor function and are being tested for the treatment of PD.<sup>259</sup> Caffeine metabolites, such as theophylline, may act as adenosine A2A receptor antagonists.<sup>260</sup> We did not determine caffeine intake in our participants; however, in recent work by Hatano and colleagues, serum levels of caffeine and caffeine metabolites were lower in PD patients than controls, even though there was no difference in caffeine consumption.<sup>242</sup>

#### **6.4.5 Oxidative stress**

Figure 6.7 shows the relative concentrations of methionine, an isomer of methionine and citrulline. According to our results, methionine sulfoxide and its isomer had increased concentrations in the PD group. Methionine sulfoxide is an oxidation product of methionine with reactive oxygen species so it is considered as a biomarker of oxidative stress.<sup>261</sup> It has also been reported that the oxidation of methionine residues could play an important role in the aggregation of normally soluble  $\alpha$ -synuclein in PD.<sup>262</sup> The average concentration of methionine sulfoxide in the two subgroups was similar, implying that the extent of peripheral oxidative stress does not increase dementia risk. On the other hand,

citrulline, an oxidization product of arginine, had a decreased concentration. It is reported to be an efficient hydroxyl radical scavenger.<sup>263</sup> Increased oxidative damage occurs in all human neurodegenerative diseases, including PD,<sup>264</sup> and some studies suggest that oxidative stress may be a factor in the loss of dopaminergic neurons.<sup>265</sup>



**Figure 6.7** Box plots of the relative concentrations of methionine sulfoxide (A), methionine sulfoxide-isomer (B) and citrulline (C) in the control group, the PDND subgroup and the PDID subgroup.

## 6.5 Limitations

First, our study examined a single cohort of subjects who were on treatment for PD at one center. Future studies should replicate and extend our results. One specific direction would be to examine untreated patients in order to estimate the potential confounding of medications. In our work, although not statistically significant, the PDID group did differ from the cognitively stable PD group in some characteristics, including baseline disease duration, UPDRS part 3, and Levodopa equivalents. We note, however, the important point

that the groups did slightly differ at baseline on global cognitive performance (MMSE), but no subject had functionally significant impairment (dementia). This highlights the importance of replicating and extending the present study with careful attention to matching on baseline characteristics.

Second, diet may affect the metabolic profiles. On diet effect, in our work, subjects came to the study in the morning after taking morning medication but no food or caffeine. They were given a breakfast after the samples were obtained. While it would have been ideal to take samples before medication this was a symptomatic group of patients and we felt it would be challenging (logistically and ethically) to withdraw patients from medication. Clearly this is something to consider in follow up validation studies.

Third, we collected blood only at baseline and thus intra-patient comparisons of metabolomic changes or correlations of changes with PD progress are not possible. We would recommend this follow-up blood collection for future research. We note that, in our work, our research goal was to compare PD with PDID, where the PD group was stable for the full period (and thus PDND) and the PDID is PDND at baseline but PDID 3 years later. This is a very interesting challenge, comparing two groups who, for all clinical purposes, were members of the same phenotype at baseline, with one developing dementia during the ensuing three years. We asked whether there would be detectable biomarker differences between these groups prior to their bifurcation into two subgroups. This attempt to discover early biomarkers of incipient dementia in PD patients was successful.

Fourth, the dansylation labeling mainly enhances the detection of the amine/phenol-containing metabolites, but it cannot detect other important metabolites in related pathways. Applying other labeling techniques to target different groups of submetabolomes will expand the overall metabolome coverage in the future.

Fifth, in our work, there was no external dataset for validation and the sample size was relatively small for validation work. In the future, more participants will be enrolled for the validation of the novel biomarker candidates comprising the panels observed in this study.

## **6.6 Conclusions**

Metabolomics analyses of serum are useful tools for identifying novel biomarker panels for both PD discrimination and the early detection of established PD patients who are at risk for transitioning to dementia. The significantly altered metabolites reported in our work can be used to differentiate (1) PD patients from healthy controls with high accuracy and (2) the stable PD with no dementia group from those with incipient dementia. Following further validation in larger cohorts, these metabolites could be used for both discrimination and establishing prognosis in PD. Follow-up studies can provide insights into potential pathways of PD neuropathology, including those associated with early discrimination and identification of risk for dementia.

## Chapter 7

### Development of Chemical Isotope Labeling Liquid Chromatography-Mass Spectrometry for Silkworm Hemolymph Metabolomics

#### 7.1 Introduction

Silkworm, *Bombyxmori*, has been an economically very important insect for over 5000 years, mainly for silk production. With recent advances in genetic engineering technology, silkworm may be potentially used to produce other functional proteins and biomaterials.<sup>266</sup> Because silkworm is very sensitive to pollutants such as pesticides,<sup>267</sup> heavy metals<sup>268</sup> and fluoride<sup>269</sup> as well as other chemicals such as pharmaceuticals,<sup>270</sup> it has been used as a target species in environmental and health safety evaluation. It has also been traditionally used as a model system for lepidopteran study.<sup>271</sup> Since the completion of silkworm genome sequencing,<sup>272-273</sup> functional genomic studies of silkworm on diverse areas of biological importance including developmental biology, reproduction and physiology have been extensively reported.<sup>271</sup> Many of these studies were focused on transcriptomic<sup>274-276</sup> or proteomic<sup>277-279</sup> investigation of silkworm. Very recently, research on using metabolomics to examine the metabolic changes induced by various stimulants or processes has been described.<sup>280-283</sup> Because metabolomics can provide complementary information to other Omics technologies, it is poised to play an increasingly important role in the future in large scale study of silkworm biology and related processes including developing genetically engineered silkworms.<sup>271</sup>

Metabolome profiling is a critical part of metabolomics studies of silkworm. Traditional methods including NMR, GC-MS and LC-MS have been used for metabolome analysis of silkworm hemolymph<sup>280-282, 284</sup> and larva brain,<sup>283</sup> but with limited metabolic coverage. Because of a small volume of sample available from each silkworm, generation of a metabolome profile with high coverage, which is often achieved by analyzing aliquots of the same sample using multiple techniques, is currently an analytical challenge. Mixing samples from a number of silkworms to form a pooled sample may increase the sample volume for analysis. However, this is not ideal to account for intra-group biological variations in individual silkworms to reveal inter-group metabolic differences, particularly if the changes are small. In this chapter, we report a sensitive method based on high-performance chemical isotope labeling (CIL) LC-MS<sup>141</sup> to perform in-depth submetabolome profiling of silkworm hemolymph. To demonstrate the utility and analytical performance of this method for silkworm metabolomics, we applied this method to examine the metabolomic changes in hemolymph samples collected from individual silkworms with and without the exposure of dichlorodiphenyltrichloroethane (DDT).

DDT was a popular organochlorine pesticide several decades ago, but is now regarded as an endocrine disruptor.<sup>285</sup> It can modulate the endocrine system through mimicking endogenous hormone action and can cause adverse effects in wildlife and human.<sup>286-287</sup> Although DDT had been banned since 1970, the negative effects will still exist for a long time because of the presence of residues in the environment and ecosystem. Silkworm should be particularly suitable for the evaluation of endocrine disrupting effects of exogenous chemicals such as DDT. It is known that the complete endocrine system of

silkworm consists of brain neurosecretory cells, suboesophageal ganglion, prothoracic glands, corpora allata, which can control the processes of growth, production, development and other aspects completely.<sup>288</sup> In addition, a wealth of background knowledge about genetics, physiology, biochemistry and genomics of silkworm<sup>271</sup> can provide us valuable information on endocrine disruption research. In our work, we applied CIL LC-MS metabolomics to generate metabolomic information in order to understand further how a simulant such as DDT affects silkworm growth as well as search for potential metabolite markers of DDT exposure. The latter is relevant to silk production in some part of the world where DDT residual levels in fields planted with mulberry trees could be still high.<sup>289-291</sup>

## **7.2 Materials and methods**

### **7.2.1 Chemicals and reagents**

All the chemicals and reagents, unless otherwise stated, were purchased from Sigma-Aldrich Canada. The pesticide DDT was purchased from AccuStandard USA. For dansylation labeling, the <sup>12</sup>C-labeling reagent (dansyl chloride) was purchased from Sigma-Aldrich and the <sup>13</sup>C-labeling reagent was synthesized according to the method published previously.<sup>292</sup> These reagents are also available from the University of Alberta (mcid.chem.ualberta.ca).

### **7.2.2 Silkworm rearing and DDT treatment**

The eggs of bivoltine hybrid “HuangKang 3” of silkworm were obtained from Sericultural Research Institute, Chinese Academy of Agricultural Sciences. The larvae were raised in incubator using a sterilized artificial diet developed at Zhejiang Academy of Agricultural Sciences, Hangzhou, China. The silkworms were raised intensively during the first instar with a condition of 29°C and 90% humidity, followed with 1°C temperature decrease and 5% humidity decrease in each instar. The door of incubator was kept open for 5 min to ventilate three times a day. Right after the first larval molted, size-matched larvae were selected and assigned randomly into the batches for the DDT treatment. Four concentrations of DDT (A=1.0 ppm, B=0.1 ppm, C=0.01 ppm, D=0.001 ppm) were used in this experiment. DDT was mixed with the diet. Meanwhile, the silkworms fed without DDT were considered as control. There were 3 replicate experiments for each DDT concentration as well as the control and 30 silkworms were used in each experiment. The diet was changed on alternate day from the first to the third instar, and every day in the fourth and fifth instar.

### **7.2.3 Hemolymph collection and preparation**

From mid to late fifth instar, silkworm started to prepare for cocoon spinning. Because this period is very crucial to silkworm development, we analyzed the metabolome of silkworm after DDT exposure at this point. Five out of thirty larvae were randomly selected from each replicate experiment on the third day of the fifth instar (i.e., after 12 days of DDT exposure). By cutting through the caudal horn, about 50  $\mu$ L of hemolymph was collected from each larva into a 1.5 mL Eppendorf tube. From 75 individual samples, 20  $\mu$ L of aliquot was taken from each sample and mixed with other aliquots to form a pooled sample

as control. The hemolymph sample was centrifuged for 10 min (14,000 rpm, 4°C) to precipitate the blood cells, and then 15 µL of the supernatant was taken into a new tube. After adding 45 µL of methanol, the sample was incubated in -20°C freezer for 2 h to precipitate the proteins. After centrifugation for 15 min (14,000 rpm, 4°C), 45 µL of the supernatant was taken out and dried using a Speed Vac at room temperature.

#### **7.2.4 Dansylation labeling of hemolymph**

For labeling, the dried sample was re-dissolved to 600 µL with 2:1 water/ACN and two 75 µL aliquots were taken for experimental duplicates. To each aliquot, 25 µL of 250 mM sodium carbonate/sodium bicarbonate buffer were added and the solution was vortexed, spun down, and mixed with 50 µL of freshly prepared <sup>12</sup>C-dansyl chloride solution (18 mg/mL in ACN) (for light labeling) or <sup>13</sup>C-dansyl chloride solution (18 mg/mL in ACN) (for heavy labeling). After 45 min incubation at 40°C, 10 µL of 250 mM NaOH was added to the reaction mixture to quench the excess dansyl chloride. The solution was then incubated at 40°C for another 10 min. Finally, 50 µL of formic acid (425 mM) in 1:1 ACN/H<sub>2</sub>O was added to consume excess NaOH and to make the solution acidic for analysis.

#### **7.2.5 LC-UV quantification**

The quenching reaction by adding NaOH after sample labeling converted the excess dansyl chloride to dansyl-OH which is very hydrophilic and can be washed away during the first step (aqueous solvent) in a step-gradient LC-UV analysis. The quenching reaction was actually very fast and complete and thus there was no leftover dansyl chloride in the labeled

sample from which the LC-UV measurement was performed. This was confirmed by targeted LC-MS analysis of dansyl chloride in a labeled sample; we could not detect any residual dansyl chloride. For LC-UV, a Waters ACQUITY UPLC system with a photodiode array (PDA) detector was used for the quantification of dansyl labeled metabolites for sample amount normalization as described earlier.<sup>115</sup> Briefly, 4  $\mu$ L of each labeled sample was injected onto a Phenomenex Kinetex C18 column (2.1 mm  $\times$  5 cm, 1.7  $\mu$ m particle size) for a fast step-gradient run. Solvent A was 0.1% (v/v) formic acid in 5% (v/v) ACN/H<sub>2</sub>O, and solvent B was 0.1% (v/v) formic acid in ACN. The gradient started with 0% B for 1 min and was increased to 95% within 0.01 min and held at 95% B for 1 min to ensure complete elution of all labeled metabolites. The flow rate used was 0.45 mL/min. The peak area related to the total labeled metabolite concentration in the sample was integrated using the Empower software (6.00.2154.003). Based on the quantification results, the <sup>12</sup>C-labeled sample and the <sup>13</sup>C-labeled pool were mixed in equal amounts.

### **7.2.6 LC-MS**

The HPLC system was an Agilent capillary 1100 binary system (Agilent, Palo Alto, CA). A reversed-phase Eclipse plus C18 column (2.1 mm $\times$ 100 mm, 1.8  $\mu$ m particle size, 95 Å pore size) was also purchased from Agilent. LC Solvent A was 0.1% (v/v) formic acid in 5% (v/v) acetonitrile, and Solvent B was 0.1% (v/v) formic acid in acetonitrile. The gradient elution profile was as follows: t = 0 min, 20% B; t = 3.5 min, 35% B; t = 18.0 min, 65% B; t = 24.0 min, 99% B; t = 32.0 min, 99% B. The flow rate was 180  $\mu$ L/min. The flow from HPLC was split 1:2 and a 60  $\mu$ L/min flow was loaded to the electrospray ionization (ESI) source of a Bruker 9.4 Tesla Apex-Qe Fourier transform ion-cyclotron

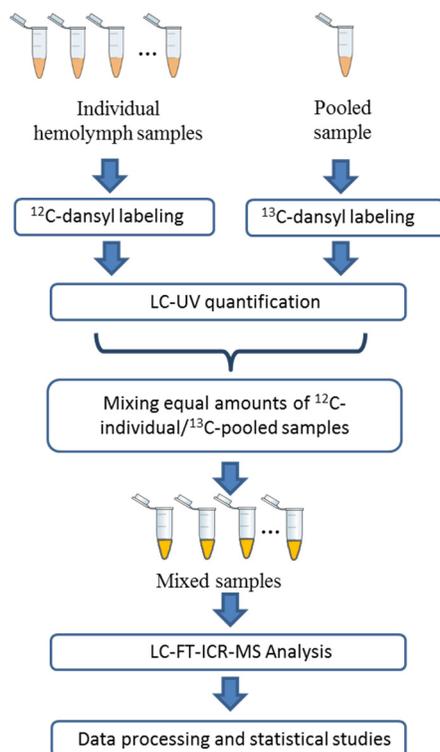
resonance (FTICR) mass spectrometer (Bruker, Billerica, MA, USA), while the rest of the flow was delivered to waste. All MS spectra were obtained in the positive ion mode. A quality control sample (QC) were injected several times to monitor the performance of LC-MS running during the whole experiment.

### **7.2.7 Data analysis**

The  $^{12}\text{C}$ -/ $^{13}\text{C}$ -peak pairs from each LC-MS run were extracted by the IsoMS software.<sup>72</sup> IsoMS-Align was used to align the peak pair data from different samples by retention time and accurate mass. The missing ratio values were filled back by using the Zero-fill program.<sup>137</sup> IsoMS-Quant<sup>157</sup> was used to generate the final metabolite-intensity table which was exported to MetaboAnalyst 3.0<sup>165</sup> for multivariate statistical analysis. Positive metabolite identification was done using a dansyl standards library consisting of 273 standards.<sup>159</sup> Putative metabolite identification was done by accurate mass matching of the experimental masses with those of the metabolites in the HMDB library and the EML database using MyCompoundID MS search.<sup>41</sup>

## 7.3 Results and discussion

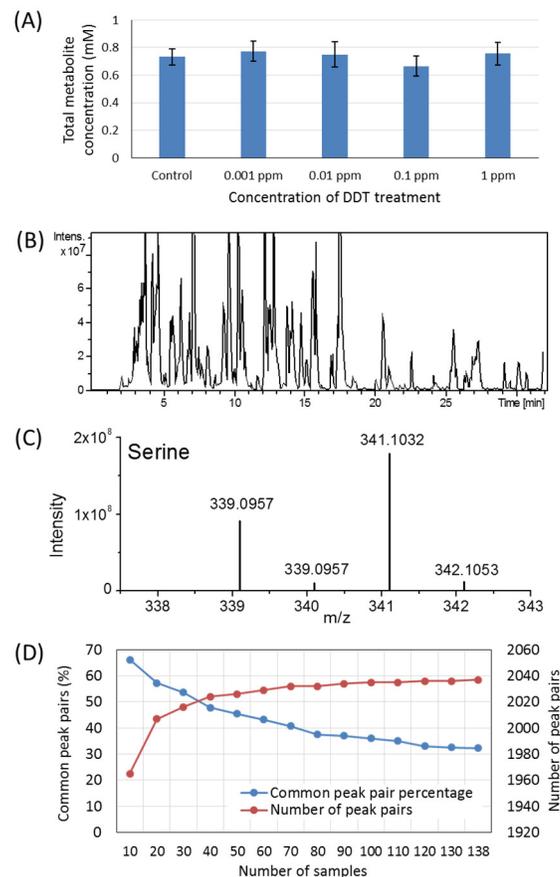
### 7.3.1 CIL LC-MS analysis of labeled hemolymph



**Figure 7.1** Workflow for silkworm hemolymph metabolome profiling using CIL LC-MS.

Figure 7.1 shows the workflow for silkworm hemolymph metabolome profiling using CIL LC-MS. One key step in the workflow is sample amount normalization. Variations in total sample amount in different samples can be greater than the analytical variation. This variation must be minimized in order to detect the concentration differences of individual metabolites caused by the DDT treatment. In our approach, we mixed the  $^{12}\text{C}$ -labeled individual sample with the same total amount of the  $^{13}\text{C}$ -labeled pooled sample, according to their total concentrations of labeled metabolites determined by the fast step-gradient LC-UV method. Figure 7.2A shows the average concentrations of labeled metabolites for the

five groups of samples. The concentrations of the non-treated samples are similar and within the range of 0.6 to 0.8 mM. The average concentrations of the other groups are similar to the control group. This result demonstrates that there is no significant change in total concentration of labeled metabolites caused by the DDT treatment. After applying sample normalization, sample amount variation in individual samples should not be an issue in our metabolomic dataset.



**Figure 7.2** (A) Averaged total concentrations of labeled metabolites in five groups of hemolymph samples. (B) A representative base-peak ion chromatogram obtained from LC-FTICR-MS analysis of a <sup>12</sup>C-/<sup>13</sup>C-labeled hemolymph sample. (C) Expanded mass spectrum showing peak pair of serine with the m/z 339.0957 peak from the <sup>12</sup>C-labeled serine in an individual sample and the m/z 341.1032 peak from the <sup>13</sup>C-labeled serine in the pooled sample. (D) The peak pair number detected and the percentage of common peak pairs as a function of the number of samples.

One benefit of knowing the concentration of labeled metabolites is that we can optimize the sample injection amount into LC-MS to ensure that a maximum number of peak pairs or metabolites are detected and the same sample amount is used for all the runs. By injecting varying amounts of a  $^{12}\text{C}$ -/ $^{13}\text{C}$ -labeled hemolymph sample while monitoring the number of peak pairs detected in LC-MS, it was found that the number of peak pairs reached the plateau at 8 nmol of injection. Subsequently, 8 nmol injection was used for all the sample runs. With this amount injection, no sample carryover from run to run was found and thus only a simple blank run was needed to re-equilibrate the column between sample runs.

Figure 7.2B shows a representative base-peak ion chromatogram obtained from LC-FTICR-MS analysis of a  $^{12}\text{C}$ -/ $^{13}\text{C}$ -labeled hemolymph sample. The advantage of using a rationally designed labeling reagent (e.g., dansyl) in CIL LC-MS is that metabolites with different polarity or even ionic species can be retained on a reversed-phase (RP) column after labeling. Figure 7.2B shows that many peaks are detected over the entire gradient elution window, indicating that labeled metabolites with a wide range of hydrophobicity are present in the sample. Figure 7.2C shows a typical mass spectrum covering the molecular ion region of a peak pair. The  $^{12}\text{C}$ -labeled peak is from the labeled serine in an individual sample and the  $^{13}\text{C}$ -labeled peak is from the pooled sample. Their peak ratio reflects the concentration difference of serine in the sample vs. the control. Since the same amount of the  $^{13}\text{C}$ -labeled pooled sample was spiked to all the  $^{12}\text{C}$ -labeled individual samples (see Figure 7.1), the ratio values of the peak pair determined from all the mass

spectra could be used to measure the concentration differences of the metabolite in these samples. In our work, after finding the peak pair, we actually reconstructed the extracted ion chromatograms of the  $^{12}\text{C}$ - and  $^{13}\text{C}$ -labeled peaks to calculate the chromatographic peak ratio and use it to measure the relative concentration more accurately and precisely, compared to using mass spectral peak intensity ratio.<sup>157</sup>

### **7.3.2 Metabolome profile of silkworm hemolymph**

In total, 150 samples were produced from duplicate experiments of 75 silkworm hemolymph samples. From the combined LC-MS results, we detected a total of 2,044 unique peak pairs with an average of  $1,467 \pm 41$  peak pairs from each individual sample. It was found that 6 duplicate-samples had significantly lower peak pair numbers than the other runs (i.e., less than 750 peaks); these were considered to be the outliers. After examining the possible causes of these outliers, it was deemed that vials of these samples might not be sealed properly and thus during shipment and storage the samples might be degraded. These samples were excluded from the final dataset and the remaining 138 analysis results were used for statistical studies. We note that spotting outliers is important in multiple sample analysis in metabolomics. Our workflow of using the  $^{13}\text{C}$ -labeled pool as the global standard for all the  $^{12}\text{C}$ -labeled individual samples and strict control of the same total amount used for sample normalization and sample injection allows us to find the outliers with ease.



our work. These data show that CIL LC-MS is a robust technique for quantifying a large number of metabolites that could be consistently detectable in the silkworm hemolymph samples.

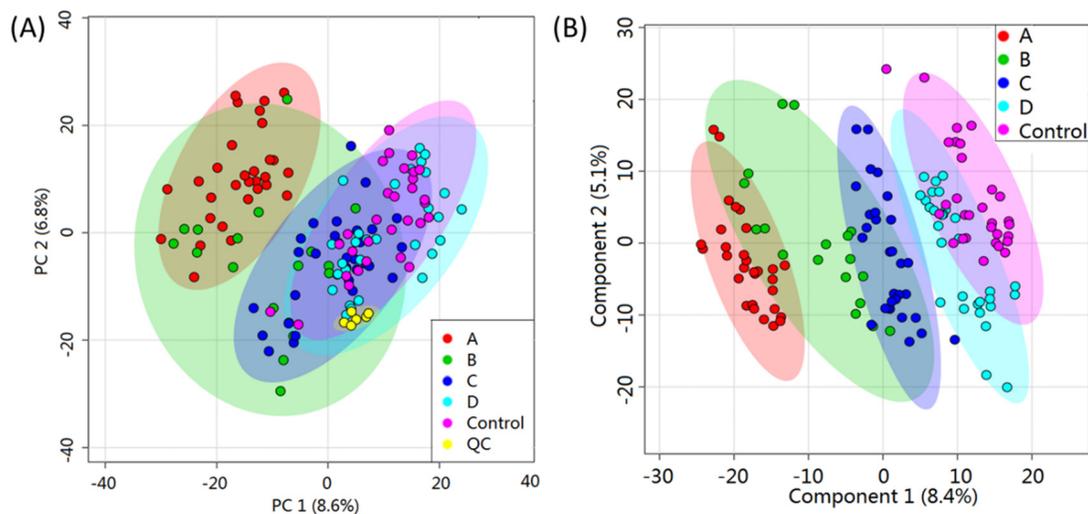
Metabolite identification for the 2044 peak pairs was carried out using two levels of database search. Positive metabolite identification was performed based on mass and retention time match (e.g., M-RT search) to the dansyl standard library containing 273 unique amines/phenols using DnsID.<sup>159</sup> We did not use MS/MS, as we have already shown in a previous paper<sup>159</sup> that accurate mass and retention time matches are already sufficient for identifying a dansyl labeled metabolite positively without the need of using MS/MS. This is because each labeled metabolite standard in the dansyl standard library has a unique combination of accurate mass and retention time. For each dansyl standard, we do have the MS/MS spectrum. However, if retention time is available for matching, this information can be used to replace MS/MS spectral match. This process of positive metabolite identification works in the same manner as we spike a standard into a real sample for accurate mass and retention time comparison for positive identification. In our previous paper,<sup>159</sup> we have shown that M-RT search gave 105 matches from a human urine sample, and manual inspection with MS/MS data did not find any mistake in the match result. Overall, according to our experiences in working with many different samples, M-RT search using the dansyl library is a reliable approach for positive metabolite identification. In total, 65 metabolites were identified by using the dansyl library searching in this study. In addition, using MyCompoundID MS search, 338 and 1,471 peak pairs were putatively identified based on accurate mass match against the HMDB library and the EML library,

respectively. It should be noted that while we designated this accurate mass match as a putative identification, the match is highly speculative and requires further information such as retention time or MS/MS comparison to that of a standard for more confident identification.

### **7.3.3 Metabolome comparison of silkworm hemolymph with different DDT treatment**

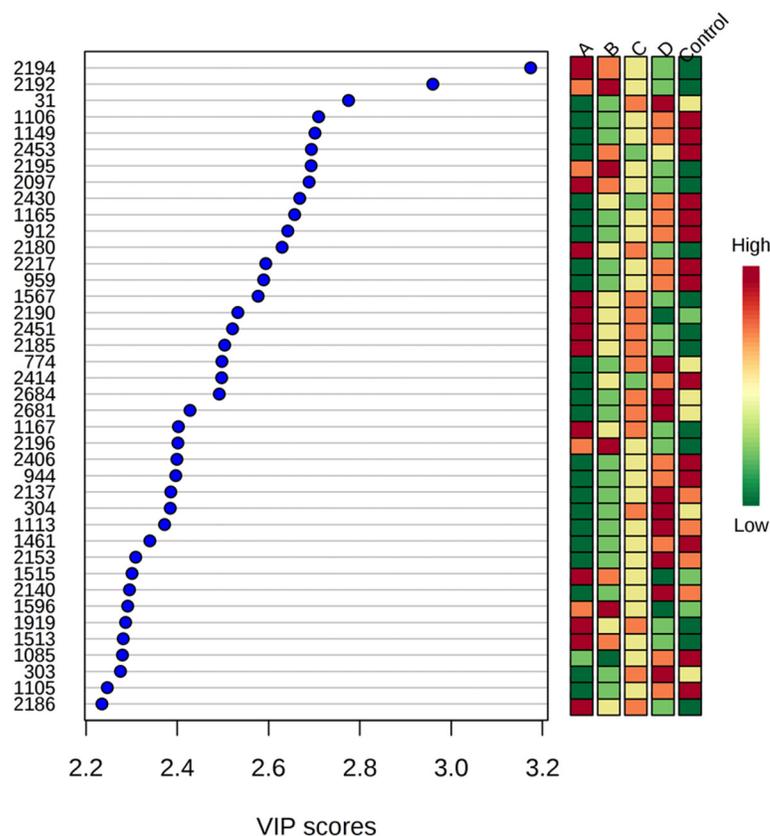
We used PCA, an unsupervised chemometric method, to provide an overall view of the whole data sets in order to determine if there are any clustering, trends or outliers. Figure 7.4A shows the PCA score plot of the 138 sample runs and 6 QC runs. Group A (red dots), B (green dots), C (dark blue dots) and D (sky blue dots) represent the samples from different concentrations of DDT treatment. The QC data (yellow dots) cluster together, indicating that excellent analytical reproducibility was achieved in LC-MS data acquisition. The PCA analysis shows some separations among the groups. Group A treated with the highest DDT concentration is largely separated from the control group, while Groups C and D with lower DDT concentrations are only slightly separated from the control.

The supervised method, PLS-DA, was used for further analysis of the five groups. Figure 7.4B shows the PLS-DA score plot of the five groups. Group separation is observed ( $R^2Y=0.999$ ,  $Q^2=0.942$ ), although there are a few overlapping data points between the adjacent groups. Figure 7.4B also shows a clear trend of increasing separation between the treated and control groups as the DDT concentration increases.



**Figure 7.4** (A) PCA score plot for QC data (yellow), control group (pink) and the samples from different concentrations of DDT treatment (1 ppm in red, 0.1 ppm in green, 0.01 ppm in dark blue and 0.001 ppm in sky blue). (B) PLS-DA score plot for the control group (pink) and the samples from different concentrations of DDT treatment (1 ppm in red, 0.1 ppm in green, 0.01 ppm in dark blue and 0.001 ppm in sky blue).

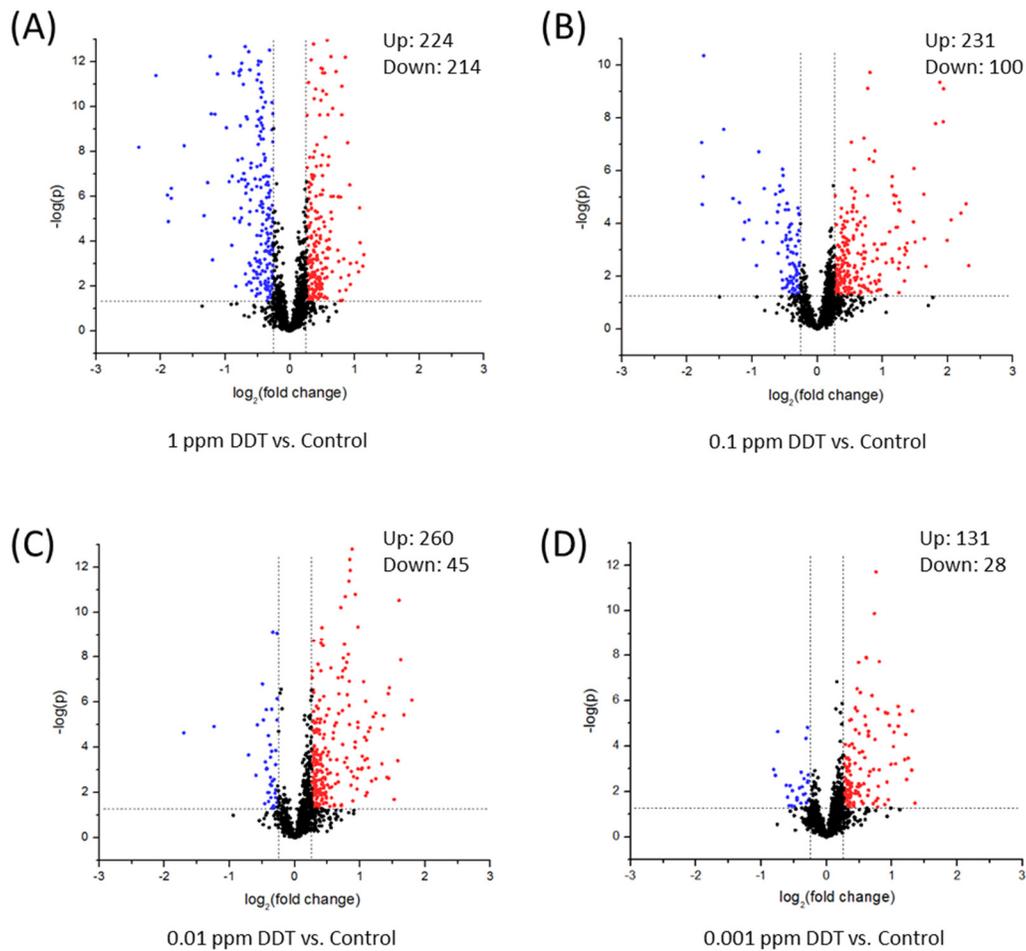
Figure 7.5 shows the study of any correlation between the metabolite concentration in each group and their DDT concentrations for the top 40 metabolites based on PLS-DA VIP scores. Different trends in concentration change are observed. For example, metabolites #2194 and #2097 increase their concentration with the increase of DDT concentration, while metabolites #2217 and #959 decrease their levels with the increase of DDT concentration. Most of the metabolites do not show strong correlations with the DDT concentration change. Thus, the metabolome of silkworm hemolymph was affected through various ways of individual metabolite level changes under different levels of DDT exposure.



**Figure 7.5** List of 40 significant metabolites with the highest PLS-DA VIP scores, showing the correlation between the metabolite concentration in each group and their DDT concentrations.

To reveal the changes of individual metabolites upon DDT exposure, binary comparison of the metabolome of treated group vs. control was carried out using volcano plots. A metabolite with a concentration change of  $\geq 1.2$ -fold and a p-value  $\leq 0.05$  was deemed to have a significant change. In this study, un-adjusted p-value was used but the volcano plot is just a relatively loose screening of potential biomarkers and we will mainly focus on the biological meanings of these candidates with pathway analysis. The numbers of up-regulated metabolites and down-regulated metabolites were found to be 224 and 214 in the comparison of the 1 ppm DDT group vs. the control group (Figure 7.6A), 231 and 100 in

the 0.1 ppm DDT group vs. the control group (Figure 7.6B), 260 and 45 in the 0.01 ppm DDT group vs. the control group (Figure 7.6C), 131 and 28 in the 0.001 ppm DDT group vs. the control group (Figure 7.6D). There is a trend of increasing the number of down-regulated metabolites as the DDT concentration increases.



**Figure 7.6** Volcano plots (fold change  $\geq 1.2$  and  $p \leq 0.05$ ) for the binary comparison of (A) the 1 ppm DDT group vs. the control group, (B) the 0.1 ppm DDT group vs. the control group, (C) the 0.01 ppm DDT group vs. the control group, and (D) the 0.001 ppm DDT group vs. the control group.

On the list of the significant metabolites found from the binary comparisons that may potentially serve as biomarkers related to DDT exposure, 33 of them could be positively identified (see Table 7.1) and many with large fold changes (fold change of  $\geq 2$  with a p-value  $\leq 0.05$ ) could be matched to some metabolite structures in HMDB or EML libraries. The positively identified metabolites are involved in various metabolic pathways. For example, tryptophan and alanine are involved in glycolysis and Krebs cycle process. Methyl-histidine and imidazoleacetic acid are involved in histidine metabolism process. Citrulline and ornithine are involved in arginine and proline metabolism process. Cystathionine is involved in serine metabolism process. These metabolic processes are all crucial to silkworm development as discussed below.

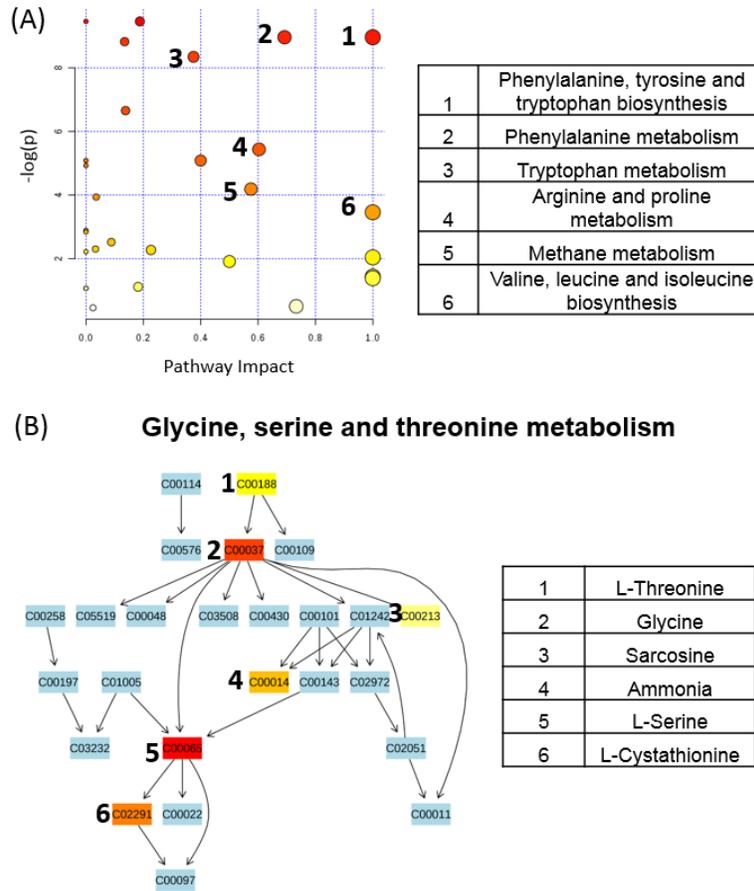
**Table 7.1** List of positively identified metabolites showing significant concentration changes after DDT treatment. (Number in bold and italic means the fold change was  $\geq 1.20$  or  $\leq 0.83$  with p-value  $\leq 0.05$ )

Ret. Time (min)	Accurate Mass	HMDB ID	Compound Name	Fold Change (vs. Control) and p Value (in brackets)			
				1 ppm A	0.1 ppm B	0.01 ppm C	0.001 ppm D
15.65	104.0564	HMDB02006	2,3-Diaminopropionic acid	<b><i>1.57</i></b> <b><i>(0.0022)</i></b>	1.29 (0.12)	1.33 (0.066)	1.11 (0.56)
13.69	222.0652	HMDB00099	L-Cystathionine-Isomer	<b><i>1.49</i></b> <b><i>(0.00021)</i></b>	<b><i>1.46</i></b> <b><i>(2.4E-05)</i></b>	1.22 (0.071)	1.15 (0.24)
12.74	165.0775	HMDB00159	L-Phenylalanine	<b><i>1.39</i></b> <b><i>(0.0046)</i></b>	1.01 (0.91)	<b><i>1.32</i></b> <b><i>(0.012)</i></b>	1.14 (0.23)
22.65	181.0708	HMDB00158	L-Tyrosine	<b><i>1.34</i></b> <b><i>(0.00065)</i></b>	1.10 (0.51)	<b><i>0.61</i></b> <b><i>(0.00022)</i></b>	0.96 (0.69)
11.44	204.0883	HMDB00929	L-Tryptophan	<b><i>1.30</i></b>	1.09	<b><i>1.26</i></b>	1.14

				<b>9.7E-09</b>	(0.25)	<b>1.5E-06</b>	(0.0019)
5.79	119.0576	HMDB00167	L-Threonine	<b>1.27</b> <b>(0.0024)</b>	1.25 (0.11)	<b>1.41</b> <b>(0.00095)</b>	<b>1.36</b> <b>(2.2E-06)</b>
4.40	105.0416	HMDB00187	L-Serine	<b>1.26</b> <b>(7.1E-05)</b>	1.05 (0.56)	1.12 (0.10)	1.04 (0.53)
4.52	130.1199	HMDB01432	Agmatine	<b>1.21</b> <b>(0.0027)</b>	<b>1.36</b> <b>(0.00037)</b>	<b>1.26</b> <b>(0.0050)</b>	0.98 (0.82)
7.57	89.0440	HMDB00161	L-Alanine	<b>1.21</b> <b>(0.033)</b>	1.11 (0.42)	1.15 (0.22)	1.18 (0.16)
13.06	131.0933	HMDB00172	L-Isoleucine	<b>1.21</b> <b>(0.031)</b>	1.27 (0.17)	<b>1.26</b> <b>(0.0095)</b>	1.17 (0.12)
15.80	118.0711	HMDB02362	2,4-Diaminobutyric acid	<b>1.21</b> <b>(8.0E-06)</b>	1.11 (0.025)	1.14 (0.00094)	1.00 (0.88)
13.34	222.0638	HMDB00099	L-Cystathionine	<b>1.20</b> <b>(0.015)</b>	<b>1.23</b> <b>(0.029)</b>	0.97 (0.63)	0.95 (0.54)
2.17	169.0841	HMDB00001	1-Methylhistidine	1.19 (0.028)	<b>1.50</b> <b>(6.0E-06)</b>	1.13 (0.28)	<b>1.23</b> <b>(0.048)</b>
18.09	155.0681	HMDB00177	L-Histidine	1.19 (0.00054)	1.16 (0.019)	<b>1.25</b> <b>(8.0E-06)</b>	1.09 (0.061)
8.67	103.0619	HMDB03911	3-Aminoisobutanoic acid	1.12 (0.0015)	1.11 (0.078)	<b>1.20</b> <b>(8.3E-08)</b>	1.12 (0.00086)
14.11	240.0192	HMDB00192	L-Cystine	1.12 (0.50)	1.02 (0.94)	<b>1.52</b> <b>(0.034)</b>	<b>1.72</b> <b>(0.020)</b>
3.32	146.0690	HMDB00641	L-Glutamine	1.11 (0.061)	1.25 (0.080)	<b>1.23</b> <b>(0.0056)</b>	<b>1.22</b> <b>(0.020)</b>
11.12	126.0414	HMDB02024	Imidazoleacetic acid	1.10 (0.10)	1.16 (0.017)	<b>1.24</b> <b>(0.0079)</b>	1.12 (0.051)
3.05	202.1415	HMDB03334	Symmetric dimethylarginine	1.07 (0.45)	<b>1.61</b> <b>(0.023)</b>	1.55 (0.061)	1.06 (0.65)
8.87	142.0362	HMDB00469	5-Hydroxymethyluracil	1.06 (0.19)	1.09 (0.38)	<b>1.21</b> <b>(0.0011)</b>	1.06 (0.18)
8.95	152.0320	HMDB00292	Xanthine	1.05 (0.60)	1.03 (0.71)	<b>1.38</b> <b>(0.0023)</b>	0.97 (0.69)

2.44	174.1100	HMDB00517	L-Arginine	1.00 (0.91)	<b>1.24</b> <b>(0.0051)</b>	1.12 (0.12)	1.16 (0.046)
9.39	240.0725	HMDB00884	Ribothymidine- H <sub>2</sub> O	1.01 (0.84)	<b>1.20</b> <b>(0.0046)</b>	1.15 (0.044)	0.94 (0.39)
5.82	17.0253	HMDB00051	Ammonia	1.00 (0.95)	<b>1.24</b> <b>(0.023)</b>	<b>1.75</b> <b>(4.7E-07)</b>	<b>1.27</b> <b>(0.0050)</b>
17.57	138.0305	HMDB00500	4-Hydroxybenzoic acid	0.99 (0.83)	<b>0.78</b> <b>(0.012)</b>	0.97 (0.83)	0.99 (0.79)
10.89	149.0506	HMDB00696	L-Methionine	0.98 (0.92)	<b>2.09</b> <b>(0.050)</b>	1.46 (0.27)	0.87 (0.52)
21.64	122.0349	HMDB00750	3-Hydroxymandelic acid - COOH	0.94 (0.18)	<b>0.75</b> <b>(0.0034)</b>	<b>0.77</b> <b>(0.0021)</b>	1.01 (0.91)
4.20	165.0462	HMDB02005	Methionine Sulfoxide - Isomer	0.93 (0.17)	<b>0.82</b> <b>(0.023)</b>	0.97 (0.63)	0.93 (0.21)
6.00	61.0479	HMDB00149	Ethanolamine	0.91 (0.31)	<b>1.47</b> <b>(0.00038)</b>	<b>1.50</b> <b>(0.018)</b>	1.16 (0.29)
3.00	59.0470	HMDB01842	Guanidine	<b>0.66</b> <b>(0.0071)</b>	1.13 (0.44)	0.74 (0.078)	<b>0.57</b> <b>(0.0011)</b>
7.53	147.0540	HMDB02393	N-methyl-D-aspartic acid	<b>0.66</b> <b>(0.0012)</b>	<b>0.73</b> <b>(0.015)</b>	0.95 (0.69)	<b>1.27</b> <b>(0.028)</b>
3.74	175.0936	HMDB00904	Citrulline	<b>0.56</b> <b>(0.010)</b>	2.08 (0.23)	1.64 (0.16)	0.68 (0.072)
16.58	132.0877	HMDB00214	Ornithine	<b>1.59</b> <b>(9.8E-07)</b>	1.22 (0.11)	<b>1.21</b> <b>(0.039)</b>	<b>1.24</b> <b>(0.0033)</b>

### 7.3.4 Significance of candidate metabolite biomarkers



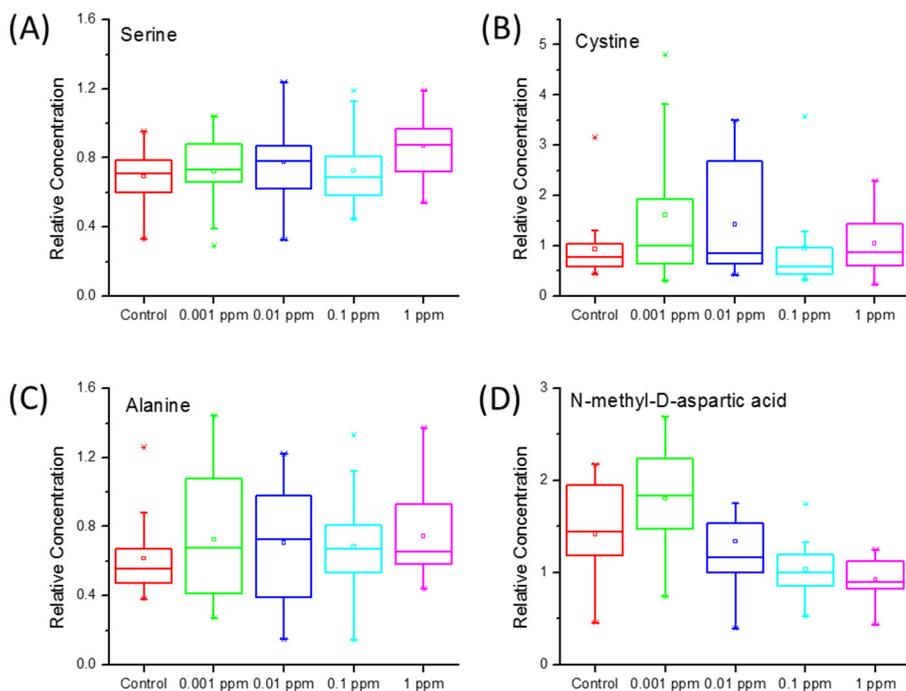
**Figure 7.7** (A) Overview of metabolic pathway analysis. (B) Pathway of glycine, serine and threonine metabolism.

The 65 definitively identified significant metabolites were exported into Metaboanalyst 3.0 for pathway analysis.<sup>165</sup> These compounds were matched to an insect pathway library which contained 79 pathways found from the fruit fly, and their concentration changes from the control to group A were also used to calculate their importance factors. A pathway impact and a p-value were calculated for each pathway, depending on how many hits it had and the importance factors of the hit compounds. The result is shown as Figure 7.7A in

which  $-\log(p\text{-value})$  is plotted against the pathway impact. The pathway on the top right corner has the highest pathway impact and statistical significance, and the mostly affected pathways include: “phenylalanine, tyrosine and tryptophan biosynthesis”; “phenylalanine metabolism”; “tryptophan metabolism”; “arginine and proline metabolism”; “methane metabolism”; and “valine, leucine and isoleucine biosynthesis”.

Among the list of other affected pathways from pathway analysis, we selected “glycine, serine and threonine metabolism” as an example (Figure 7.7B) for in-depth discussion, as this metabolic pathway has a direct relevance to silk production and a larger number of metabolites in this pathway were positively identified in our work. On the pathway schematic, light blue means the metabolite is not in the list of the 65 positively identified metabolites but is used as background for enrichment analysis. The six metabolites with other colors (varying from yellow to red) are the positively identified metabolites, including threonine, glycine, sarcosine, ammonia, serine and cystathionine. A red color indicates the corresponding metabolite has a more significant change between the control and group A comparing to a yellow-colored metabolite. For example, sarcosine is not a significantly changed metabolite. Both serine and its downstream metabolite, cystathionine, are important metabolites and serine has a larger fold change than cystathionine. With the help of the pathway analysis, some metabolite changes were thought to be likely due to endocrine system disrupted in silkworm by the DDT treatment, and the biological processes underlie these changes were also studied (see below). In addition to the comparison between the control and group A, the trend of the metabolite concentration change among the four different groups can also help us understand the DDT’s interference

on biological processes of silkworm. Figure 7.8 shows the metabolite concentrations in the control and the four study groups of four example metabolites. For example, the concentration of serine shows an increasing trend with the concentration of DDT, which may reveal a positive relationship between the DDT exposure and serine up-regulation.



**Figure 7.8** Box plots of four metabolites showing their relative concentrations in the control group and four DDT treatment groups.

Serine is derived from glycine and crucial in maintaining health of the neuron.<sup>293</sup> Alanine is one of the most important amino acids released by muscle.<sup>294</sup> Through pathway searching, we found that serine is included in both the “glycine, serine and threonine” and “methane” metabolism pathways. It is illustrated that serine is a crucial metabolite in

silkworm development. Serine, alanine and glycine are the three most important amino acids for silk protein synthesis in silkworm.<sup>295</sup> With 1 ppm DDT feed to silkworm from start of second instar to cocoon, 9% decrease of the weight of cocoon shell was observed in our previous test (unpublished). In the present study, we discovered that serine and alanine were both up-regulated in hemolymph of silkworm after 1 ppm DDT treatment (Figure 7.8). It might imply that serine and alanine had not fully been utilized in silk protein production. Since the process of silk protein synthesis is controlled by endogenous hormones, particularly juvenile hormone (JH) and molting hormone (MH),<sup>296</sup> it is possible that DDT disturbed the balance of JH and MH, thereby hindering the silk protein synthesis.

Methionine is an essential amino acid for growth and development of animals. It is required for cysteine synthesis, and the sulfur atom from methionine is transferred to cysteine.<sup>297</sup> Cystine is formed by linking two cysteine residues via a disulfide bond (cys-S-S-cys) between the -SH groups. Thus, methionine, cysteine and cystine have a close relationship in the “cysteine and methionine metabolism” pathway. In this study, we discovered that methionine and cystine were up-regulated in different DDT concentrations, which would hint that “cysteine and methionine” metabolism pathway was disturbed. Meanwhile, cysteine is very important to the structure of silkworm prothoracicotropic hormone (PTTH), which needs one inter-chain disulfide linkage (Cys-Cys) and three intra-chain disulfide linkages for manifesting biological activity.<sup>298</sup> PTTH is a neuropeptide hormone, and can stimulate prothoracic gland to produce MH.<sup>296</sup> The up-regulating of methionine and cystine in silkworm hemolymph suggests that the process of larva molting might be affected.

Tryptophan is an essential amino acid, and may contribute to mental retardation when it is in excess in the blood. Assessments of tryptophan deficiency had been done by studying excretion of tryptophan metabolites in the urine or blood from people who had nervous system disease.<sup>299-300</sup> Tryptophan is also the precursor of serotonin, a neurotransmitter and neurohormone found in many animals.<sup>301-303</sup> Serotonin had been proved to have antagonist roles in female reproduction and can stimulate reproduction in male of *M. rosenbergii*.<sup>304</sup> In our work, the concentration of tryptophan was up-regulated in silkworm hemolymph under 1 and 0.01 ppm DDT group. The disorder of the tryptophan metabolism might induce the change of serotonin, which could affect the reproduction of silkworm moth. However, in our work, we could not detect serotonin directly, probably due to very low level of serotonin present in silkworm hemolymph during the 5th instar. We believe that tryptophan may be a potential biomarker for endocrine disruptor evaluation using silkworm as a target insect.

Asymmetric dimethylarginine (ADMA) is an endogenous inhibitor of nitric oxide synthase (NOS). Symmetric dimethylarginine (SDMA) is a stereoisomer produced alongside ADMA, and has been considered as a risk factor for cardiovascular diseases. SDMA does not inhibit NOS activity directly, but might have an indirect effect by limiting cellular uptake of L-arginine.<sup>305</sup> Some researches had showed that the SDMA levels of plasma were raised in women with polycystic ovary syndrome<sup>306</sup> and people with hyperthyroidism.<sup>307</sup> The healthy postmenopausal women always have an insufficient level of estrogens. Verhoeven et al. found that SDMA concentration of plasma was reduced and arginine was transiently decreased after oral 17  $\beta$ -estradiol combined with norethisterone.<sup>308</sup> So far,

there is no correlated report about the change of SDMA in insect. In this study, the SDMA and arginine levels of silkworm hemolymph increased 1.61 and 1.24 fold, respectively, in the 0.1 ppm DDT treatment group. DDT is a pesticide that has homologous estrogen biologic character. The changing trends of SDMA and arginine in our experiment were contrary to the results of Verhoeven's report. We propose the reason of difference is that the changes of SDMA and arginine could show difference between the male and female silkworm larvae after DDT treatment. Thus, in the future work, SDMA, ADMA and arginine will need to be measured using male and female silkworm larvae, respectively. The results of SDMA changing may imply that the balance of hormone in silkworm was broken by DDT treatment.

N-Methyl-D-aspartic acid (NMDA) receptor plays significant roles during the development of nervous system in animal. It involves in many physiological processes, such as apoptosis, learning and memory.<sup>309</sup> NMDA acts as a specific agonist at the NMDA receptor.<sup>310</sup> In our study, NMDA had 1.27-fold increase in the silkworm hemolymph at 0.001 ppm DDT compared to the control (Figure 7.8D). The excess of NMDA might bind to more NMDA receptor, which would cause interference on the nervous system of silkworm and hormones synthesis. However, at 0.1 ppm and 1 ppm of DDT groups, NMDA had 0.73-fold and 0.66-fold decreases conversely, which would lead the NMDA receptor excessive and make the silkworm under an excited state. These results implied that different doses of DDT affected the silkworm development through different ways.

L-tyrosine can be changed into L-dopa with tyrosine hydroxylase catalysis. And L-dopa is the precursor for dopamine and melanin.<sup>280</sup> In lepidopteran insect, melanization is controlled by melanization and reddish coloration hormone (MRCH), which is an important neurohormone.<sup>311-312</sup> In our study, DDT exposure increased the tyrosine level (1.34 fold) at 1 ppm and decreased the level (0.61 fold) at 0.01 ppm, which might hint that the process of pigmentation was disturbed and the balance of endocrine hormone was broken by DDT feeding. This is an interesting finding, as if we could prove the relevance between the epidermal color and DDT concentrations, the color of larva might become a potential indicator for endocrine disruption evaluation using silkworm as a target.

These potential biomarkers discussed above demonstrated that our metabolomic analysis method could be used to discover untargeted metabolic biomarkers. At present, there are only limited amount of metabolomic data available on silkworm. The biomarkers found in this study need to be further verified to confirm the relationship with endocrine disruption effects by other methods such as determination hormone level of silkworm after DDT treatment.

#### **7.4 Conclusions**

In this study, we have developed a chemical isotope labeling LC-MS method for silkworm hemolymph metabolomics and apply this method to study the metabolomic changes in silkworms with and without different levels of DDT exposure. Using this method, a total of 2,044 peak pairs have been detected in 138 samples of five groups. By searching the

HMDB database and the EML library based on accurate mass match, 338 and 1471 metabolites have been putatively identified, respectively. Based on the mass and retention time match in the dansyl standard library, 65 unique amines/phenols were positively identified. Among them, 33 metabolites have  $\geq 1.20$ -fold or  $\leq 0.83$ -fold in one or more groups with p-value  $\leq 0.05$ .

Hemolymph metabolomic study on DDT treated silkworms showed that metabolomic profiles could differentiate the different DDT treatment groups from the control group. Several metabolite candidate biomarkers identified in this study had been detected in silkworm hemolymph in previous work, such as serine, alanine and arginine. These amines are very important to silkworm growth, and have been found in silkworm hemolymph throughout the whole stages of development. Many potentially new biomarkers were found in this work. For example, based on the functional analysis, we predict that the changes of methionine, cysteine and tyrosine could be due to the endocrine system disrupted by the pesticide DDT. We expect that other chemicals besides DDT could illicit similar responses. In the future, we will further explore the relationships among different pesticides or chemicals, biomarkers and endocrine disrupting effects and mechanism of endocrine disruption.

The present study focused on the use of dansylation labeling for the detection of compounds containing primary, secondary amines and phenol groups. However, the CIL LC-MS workflow described herein should be applicable to other submetabolome profiling

using other labeling reagents. In our future work, we plan to screen more biomarkers by using other labeling methods. In addition, we plan to verify the authenticity of these biomarkers using biochemical or molecular biology methods.

## Chapter 8

### Conclusions and Future Work

Blood-metabolite-based biomarkers are promising tools for the early-stage diagnosis, accurate prognostic prediction and personalized treatment of various diseases. LC-MS-based metabolomics profiling provides a sensitive and robust methodology for biomarker discovery and delineates the underlying metabolic pathways. Considering that traditional LC-MS platforms are limited by the low metabolome coverage and quantification accuracy, our group has developed the CIL LC-MS methods which can significantly increase the metabolome coverage and effectively overcome the detection variability. Nonetheless, blood metabolomics remains susceptible to biological variations and experimental interferences during the sample preparation. Large sample sizes and time-resolved studies are highly desirable in order to improve the reliability of findings. Adding the time dimension to the study design is also preferable, so cohort and intervention studies are recommended. The first part of my thesis focuses on assessing and minimizing the variability in blood metabolomics studies. The second part of my thesis describes the application of CIL LC-MS method to a cohort study and an intervention study.

Chapter 1 provides an overview of the basic concepts of the biomarker, metabolomics, CIL LC-MS platform, blood sample handling, statistical analyses and the challenges in blood biomarker discovery. Particularly, possible solutions to the challenges are proposed, leading the storyline of this thesis.

Matrix effect from various constituents in biological samples can reduce the accuracy of quantitative metabolomics. CIL LC-MS can overcome the matrix effect on MS detection based on measuring the intensity ratios of metabolite peak pairs detected in a mixture of a light-isotope labeled sample, and a heavy isotope labeled reference sample. However, the chemical labeling process itself may encounter matrix effect which can influence the overall quantitative results. Chapter 2 reports the effects of salts and buffers commonly present in metabolomic samples on dansylation labeling. It is shown that high concentrations of NaCl and phosphate buffer (>50 mM) or PBS can reduce or enhance the labeling efficiencies of metabolites. By maintaining similar matrix contents in an individual sample versus a reference sample, relative quantification of metabolites can be performed without compromising the metabolomic profiling results. For samples containing varying amounts of high salts such as urine, we demonstrate that the matrix effect can be largely overcome by diluting the original sample before dansylation labeling (e.g., fourfold dilution for urine).

As serum metabolomics is widely used for biomarker discovery, it is crucial to characterize the “normal” metabolite concentrations and inter-subject variations in the general population, as well as the potential confounding factors. In Chapter 3, non-targeted metabolome profiling was performed on serum samples from 100 healthy subjects, using two differential isotope labeling methods and high-resolution LC-MS. A high-coverage serum metabolome database including 1,348 amine/phenol-containing metabolites and 1,065 carboxyl-containing metabolites has been developed, providing the relative concentrations and inter-individual variations. In addition, this study demonstrates the

impact of sex and body mass index on human serum metabolome, indicating that these factors should be carefully assessed in metabolomics studies.

Blood is an important biofluid for metabolomics study and clinical diagnosis. Compared with venous blood, capillary blood collected by finger pricking is a more convenient and less invasive alternative. However, the detectability of low-abundant compounds is usually limited with these low-volume samples. Therefore, in Chapter 4, we developed a CIL method accompanied with high-resolution LC-MS to achieve metabolomics analyses with superior sensitivity and high metabolome coverage. By collecting only one microliter of finger blood from the participants, we detected and quantified 1,722 amine/phenol-containing metabolites. Among them, 73 have been positively identified and another 1,393 have been putatively identified. An even smaller sample volume also works without significant decrease of metabolome coverage, opening the possibility of analyzing other groups of metabolites without increasing the sample collection amount. Using the dietary exposure to coffee as an example, we have demonstrated that our technique has excellent sensitivity, repeatability and robustness for both biomarker discovery and time-resolved exposomics studies. With our method, finger blood can replace venous blood for metabolomics studies, making low-cost and point-of-care metabolic analyses feasible in the future.

Carboxylic acids have also demonstrated biological significance in various physiological processes. And since carboxylic acids are widely used as food additives, they can become

ideal biomarkers of dietary intake. In Chapter 5, the finger blood analysis method is further improved, and from only one microliter of whole blood we successfully detected 2,074 dansyl-labeled and 2,254 DMPA-labeled metabolites. Time-of-day variations of metabolites are studied. And the method is used to assess the metabolic response to energy drink consumption. We have proved that our method has the adequate sensitivity and accuracy for exposure studies.

In Chapter 6, we profiled the amine/phenol submetabolome to determine potential metabolite biomarkers associated with Parkinson's disease (PD) and PD with incipient dementia. It is a cohort study design that at baseline of a 3-wave (18-month intervals) longitudinal study, serum samples were collected from 42 healthy controls and 43 PD patients. By wave 3 (year 3) 16 PD patients were diagnosed with dementia and were classified as PD with incipient dementia at baseline. Metabolomic analyses detected 719 common metabolites in 80% of the samples. Some were significantly altered in a pairwise comparison of different groups (fold-change of  $>1.2$  or  $< 0.83$  with  $q < 0.1$ ). We discriminated PD and controls by using a 5-metabolite panel, vanillic acid, 3-hydroxykynurenine, isoleucyl-alanine, 5-acetylamino-6-amino-3-methyluracil, and theophylline. The Receiver Operating Characteristic curve produced an Area-Under-the-Curve value of 0.955 with 87.5% sensitivity and 93.0% specificity. In comparing PD with no dementia with PD with incipient dementia we used an 8-metabolite panel, His-Asn-Asp-Ser, 3, 4-dihydroxy-phenylacetone, desaminotyrosine, hydroxy-isoleucine, alanyl-alanine, putrescine [-2H], purine [+O] and its riboside. This produced an Area-Under-the-Curve value of 0.862 with 80.0% sensitivity and 77.0% specificity. The significantly

altered metabolites can be used to differentiate (1) PD patients from healthy controls with high accuracy and (2) the stable PD with no dementia group from those with incipient dementia. Following further validation in larger cohorts, these metabolites could be used for both clinical diagnosis and prognosis of PD.

Silkworm (*Bombyxmori*) is a very useful target insect for evaluation of endocrine disruptor chemicals (EDCs) due to mature breeding techniques, complete endocrine system and broad basic knowledge on developmental biology. Comparative metabolomics of silkworms with and without EDC exposure offers another dimension of studying EDCs. In Chapter 7, we report a workflow on metabolomic profiling of silkworm hemolymph based on CIL LC-MS and demonstrate its application in studying the metabolic changes associated with the pesticide DDT exposure in silkworm. Hemolymph samples were taken from mature silkworms after growing on a diet that contained DDT at four different concentrations (1, 0.1, 0.01, 0.001 ppm) as well as on diet without DDT as controls. They were subjected to the differential  $^{12}\text{C}$ -/ $^{13}\text{C}$ -dansyl labeling of the amine/phenol submetabolome, LC-UV quantification of the total amount of labeled metabolites for sample normalization, and LC-MS detection and relative quantification of individual metabolites in comparative samples. The total concentration of labeled metabolites did not show any significant change between four DDT-treatment groups and one control group. Multivariate statistical analysis of the metabolome data set showed that there was a distinct metabolomic separation between the five groups. Out of the 2044 detected peak pairs, 338 and 1471 metabolites have been putatively identified against the HMDB database and the EML library, respectively. 65 metabolites were identified by the dansyl library searching

based on the accurate mass and retention time. Among the 65 identified metabolites, 33 positive metabolites had changes of greater than 1.20-fold or less than 0.83-fold in one or more groups with p-value of smaller than 0.05. Several useful biomarkers including serine, methionine, tryptophan, asymmetric dimethylarginine, N-Methyl-D-aspartic and tyrosine were identified. The changes of these biomarkers were likely due to the disruption of the endocrine system of silkworm by DDT. Our work illustrates that the method of CIL LC-MS is useful to generate quantitative submetabolome profiles from a small volume of silkworm hemolymph with much higher coverage than conventional LC-MS methods, thereby facilitating the discovery of potential metabolite biomarkers related to EDC or other chemical exposure.

Overall, we have successfully developed a high-coverage CIL-LC-MS platform for profiling the amine, phenol and carboxyl submetabolomes of blood. Experimental variations have been largely overcome and biological variations have been carefully assessed. We have established a metabolome database listing the common biological variations among healthy population. With the figure blood analysis, we can conduct more time-dependent studies on the diet effect and other environmental exposures, and establish databases of these time-dependent metabolic changes. In biomarker discovery, when researchers focus on a relatively small list of biomarker candidates, which demonstrate significant changes in a comparison study, they can use this information of biological variations to determine if the change of a specific biomarker candidate is truly due to the disease. And in clinical diagnosis using biomarkers, this information can help us determine the criterion range. For example, setting the fold change at 1.5 for classification may work

well for a biomarker with very stable blood concentration, while for another biomarker with significant time-dependent variation following the circadian rhythm, we may need to increase the criterion to a 2-fold change. Moreover, the time-dependent metabolomics can tell us the length of diet effect for each food metabolite and whether the subjects need to fast before the testing of specific biomarkers. Importantly, we will be able to determine how long the subjects should fast or if any computational algorithms can cope with the diet effect, which is very useful in clinical practices.

The first work to be done in future is the enrichment of the serum metabolome database. The profiles of hydroxyl and carbonyl submetabolomes will be added. The overlap region of the four submetabolomes should be evaluated, and these metabolites can potentially be used to monitor the performance and consistency among the four labeling reactions. The dependent pairs of metabolites should also be noticed by algorithms and further studies. Second, the metabolite identification should be improved. We will obtain more standards to enlarge our standard libraries. The putatively identified metabolites in the serum metabolome database will also be further confirmed by MS/MS spectra. Third, for the finger blood analysis, we will try to simplify the sample handling process and to make the point-of-care analysis feasible. Microfluidic chip is a promising choice for the separation of blood cells and proteins. At last, with the less-invasive blood analysis method, larger sample sizes become more realistic to obtain, and the biomarker candidates will be validated for future clinical applications. In summary, qualified biomarkers come from well-designed experiments, careful sample handling, stable analysis platforms, and a solid understanding of data analysis principles.

## References

1. Naylor, S., Biomarkers: current perspectives and future prospects. *Expert Review of Molecular Diagnostics* **2003**, *3* (5), 525-529.
2. Frank, R.; Hargreaves, R., Clinical biomarkers in drug discovery and development. *Nature Reviews Drug Discovery* **2003**, *2* (7), 566-580.
3. Colburn, W. A., Biomarkers in drug discovery and development: from target identification through drug marketing. *The Journal of Clinical Pharmacology* **2003**, *43* (4), 329-341.
4. Mayeux, R., Biomarkers: potential uses and limitations. *NeuroRx* **2004**, *1* (2), 182-188.
5. Bunn, H. F.; Gabbay, K. H.; Gallop, P. M., The glycosylation of hemoglobin: relevance to diabetes mellitus. *Science* **1978**, *200* (4337), 21-27.
6. Hagmar, L.; Bonassi, S.; Strömberg, U.; Brøgger, A.; Knudsen, L. E.; Norppa, H.; Reuterwall, C., Chromosomal aberrations in lymphocytes predict human cancer: a report from the European Study Group on Cytogenetic Biomarkers and Health (ESCH). *Cancer research* **1998**, *58* (18), 4117-4121.
7. Belinsky, S. A.; Nikula, K. J.; Palmisano, W. A.; Michels, R.; Saccomanno, G.; Gabrielson, E.; Baylin, S. B.; Herman, J. G., Aberrant methylation of p16INK4a is an early event in lung cancer and a potential biomarker for early diagnosis. *Proceedings of the National Academy of Sciences* **1998**, *95* (20), 11891-11896.
8. Laxman, B.; Morris, D. S.; Yu, J.; Siddiqui, J.; Cao, J.; Mehra, R.; Lonigro, R. J.; Tsodikov, A.; Wei, J. T.; Tomlins, S. A., A first-generation multiplex biomarker analysis of urine for the early detection of prostate cancer. *Cancer research* **2008**, *68* (3), 645-649.
9. Ballman, K. V., Biomarker: predictive or prognostic? *Journal of Clinical Oncology* **2015**, *33* (33), 3968-3971.
10. Tarawneh, R.; D'angelo, G.; Macy, E.; Xiong, C.; Carter, D.; Cairns, N. J.; Fagan, A. M.; Head, D.; Mintun, M. A.; Ladenson, J. H., Visinin - like protein - 1: Diagnostic and prognostic biomarker in Alzheimer disease. *Annals of neurology* **2011**, *70* (2), 274-285.
11. Han, W.; Sapkota, S.; Camicioli, R.; Dixon, R. A.; Li, L., Profiling novel metabolic biomarkers for Parkinson's disease using in - depth metabolomic analysis. *Movement Disorders* **2017**.
12. Evans, W. E.; Johnson, J. A., Pharmacogenomics: the inherited basis for interindividual differences in drug response. *Annual review of genomics and human genetics* **2001**, *2* (1), 9-39.
13. Frueh, F. W.; Amur, S.; Mummaneni, P.; Epstein, R. S.; Aubert, R. E.; DeLuca, T. M.; Verbrugge, R. R.; Burckart, G. J.; Lesko, L. J., Pharmacogenomic biomarker information in drug labels approved by the United States food and drug administration: prevalence of related drug use. *Pharmacotherapy: The Journal of Human Pharmacology and Drug Therapy* **2008**, *28* (8), 992-998.
14. Sim, S. C.; Ingelman-Sundberg, M., Pharmacogenomic biomarkers: new tools in current and future drug therapy. *Trends in Pharmacological Sciences* **2011**, *32* (2), 72-81.
15. Barbour, A. G.; Heiland, R. A.; Howe, T. R., Heterogeneity of major proteins in Lyme disease borreliae: a molecular analysis of North American and European isolates. *Journal of Infectious Diseases* **1985**, *152* (3), 478-484.

16. Venook, A. P.; Papandreou, C.; Furuse, J.; de Guevara, L. L., The incidence and epidemiology of hepatocellular carcinoma: a global and regional perspective. *The oncologist* **2010**, *15* (Supplement 4), 5-13.
17. Ng, A. W.; Poon, S. L.; Huang, M. N.; Lim, J. Q.; Boot, A.; Yu, W.; Suzuki, Y.; Thangaraju, S.; Ng, C. C.; Tan, P., Aristolochic acids and their derivatives are widely implicated in liver cancers in Taiwan and throughout Asia. *Science Translational Medicine* **2017**, *9* (412), eaan6446.
18. Hunter, D. J., Gene–environment interactions in human diseases. *Nature Reviews Genetics* **2005**, *6* (4), 287-298.
19. London, E., The environment as an etiologic factor in autism: a new direction for research. *Environmental Health Perspectives* **2000**, *108* (Suppl 3), 401.
20. Rifai, N.; Gillette, M. A.; Carr, S. A., Protein biomarker discovery and validation: the long and uncertain path to clinical utility. *Nature biotechnology* **2006**, *24* (8), 971-983.
21. Ilyin, S. E.; Belkowski, S. M.; Plata-Salamán, C. R., Biomarker discovery and validation: technologies and integrative approaches. *Trends in biotechnology* **2004**, *22* (8), 411-416.
22. Younossi, Z. M.; Jarrar, M.; Nugent, C.; Randhawa, M.; Afendy, M.; Stepanova, M.; Rafiq, N.; Goodman, Z.; Chandhoke, V.; Baranova, A., A novel diagnostic biomarker panel for obesity-related nonalcoholic steatohepatitis (NASH). *Obesity surgery* **2008**, *18* (11), 1430-1437.
23. LeWitt, P.; Schultz, L.; Auinger, P.; Lu, M.; Investigators, P. S. G. D., CSF xanthine, homovanillic acid, and their ratio as biomarkers of Parkinson's disease. *Brain research* **2011**, *1408*, 88-97.
24. Spell, D. W.; Jones, D. V.; Harper, W. F.; Bessman, J. D., The value of a complete blood count in predicting cancer of the colon. *Cancer detection and prevention* **2004**, *28* (1), 37-42.
25. Ernst, E.; Weihmayr, T.; Schmid, M.; Baumann, M.; Matral, A., Cardiovascular risk factors and hemorheology: physical fitness, stress and obesity. *Atherosclerosis* **1986**, *59* (3), 263-269.
26. Fung, E. T., A recipe for proteomics diagnostic test development: the OVA1 test, from biomarker discovery to FDA clearance. *Clinical chemistry* **2010**, *56* (2), 327-329.
27. Reginster, J. Y.; Deroisy, R.; Rovati, L. C.; Lee, R. L.; Lejeune, E.; Bruyere, O.; Giacovelli, G.; Henrotin, Y.; Dacre, J. E.; Gossett, C., Long-term effects of glucosamine sulphate on osteoarthritis progression: a randomised, placebo-controlled clinical trial. *The Lancet* **2001**, *357* (9252), 251-256.
28. Oliver, S. G.; Winson, M. K.; Kell, D. B.; Baganz, F., Systematic functional analysis of the yeast genome. *Trends in biotechnology* **1998**, *16* (9), 373-378.
29. Weckwerth, W., Metabolomics in systems biology. *Annual review of plant biology* **2003**, *54* (1), 669-689.
30. Vineis, P.; Chadeau-Hyam, M.; Gmuender, H.; Gulliver, J.; Herceg, Z.; Kleinjans, J.; Kogevinas, M.; Kyrtopoulos, S.; Nieuwenhuijsen, M.; Phillips, D., The exposome in practice: design of the EXPOsOMICS project. *International journal of hygiene and environmental health* **2017**, *220* (2), 142-151.
31. Psychogios, N.; Hau, D. D.; Peng, J.; Guo, A. C.; Mandal, R.; Bouatra, S.; Sinelnikov, I.; Krishnamurthy, R.; Eisner, R.; Gautam, B., The human serum metabolome. *PloS one* **2011**, *6* (2), e16957.
32. Sato, Y.; Suzuki, I.; Nakamura, T.; Bernier, F.; Aoshima, K.; Oda, Y., Identification of a new plasma biomarker of Alzheimer's disease using metabolomics technology. *Journal of lipid research* **2012**, *53* (3), 567-576.
33. Nishiumi, S.; Shinohara, M.; Ikeda, A.; Yoshie, T.; Hatano, N.; Kakuyama, S.; Mizuno, S.; Sanuki, T.; Kutsumi, H.; Fukusaki, E., Serum metabolomics as a novel diagnostic approach for pancreatic cancer. *Metabolomics* **2010**, *6* (4), 518-528.

34. Lawton, K. A.; Berger, A.; Mitchell, M.; Milgram, K. E.; Evans, A. M.; Guo, L.; Hanson, R. W.; Kalhan, S. C.; Ryals, J. A.; Milburn, M. V., Analysis of the adult human plasma metabolome. **2008**.
35. Griffiths, W. J.; Koal, T.; Wang, Y.; Kohl, M.; Enot, D. P.; Deigner, H. P., Targeted metabolomics for biomarker discovery. *Angewandte Chemie International Edition* **2010**, *49* (32), 5426-5445.
36. Roberts, L. D.; Souza, A. L.; Gerszten, R. E.; Clish, C. B., Targeted metabolomics. *Current protocols in molecular biology* **2012**, 30.2. 1-30.2. 24.
37. Becker, S.; Kortz, L.; Helmschrodt, C.; Thiery, J.; Ceglarek, U., LC-MS-based metabolomics in the clinical laboratory. *Journal of Chromatography B* **2012**, *883*, 68-75.
38. Wishart, D. S.; Jewison, T.; Guo, A. C.; Wilson, M.; Knox, C.; Liu, Y.; Djoumbou, Y.; Mandal, R.; Aziat, F.; Dong, E., HMDB 3.0—the human metabolome database in 2013. *Nucleic acids research* **2012**, *41* (D1), D801-D807.
39. Smith, C. A.; O'Maille, G.; Want, E. J.; Qin, C.; Trauger, S. A.; Brandon, T. R.; Custodio, D. E.; Abagyan, R.; Siuzdak, G., METLIN: a metabolite mass spectral database. *Therapeutic drug monitoring* **2005**, *27* (6), 747-751.
40. Horai, H.; Arita, M.; Kanaya, S.; Nihei, Y.; Ikeda, T.; Suwa, K.; Ojima, Y.; Tanaka, K.; Tanaka, S.; Aoshima, K., MassBank: a public repository for sharing mass spectral data for life sciences. *Journal of mass spectrometry* **2010**, *45* (7), 703-714.
41. Li, L.; Li, R.; Zhou, J.; Zuniga, A.; Stanislaus, A. E.; Wu, Y.; Huan, T.; Zheng, J.; Shi, Y.; Wishart, D. S., MyCompoundID: using an evidence-based metabolome library for metabolite identification. *Analytical chemistry* **2013**, *85* (6), 3401-3408.
42. Huan, T.; Tang, C.; Li, R.; Shi, Y.; Lin, G.; Li, L., MyCompoundID MS/MS Search: Metabolite identification using a library of predicted fragment-ion-spectra of 383,830 possible human metabolites. *Analytical chemistry* **2015**, *87* (20), 10619-10626.
43. Ang, J. E.; Revell, V.; Mann, A.; Mäntele, S.; Otway, D. T.; Johnston, J. D.; Thumser, A. E.; Skene, D. J.; Raynaud, F., Identification of human plasma metabolites exhibiting time-of-day variation using an untargeted liquid chromatography-mass spectrometry metabolomic approach. *Chronobiology international* **2012**, *29* (7), 868-881.
44. Ciborowski, M.; Lipska, A.; Godzien, J.; Ferrarini, A.; Korsak, J.; Radziwon, P.; Tomasiak, M.; Barbas, C., Combination of LC-MS-and GC-MS-based metabolomics to study the effect of ozonated autohemotherapy on human blood. *Journal of proteome research* **2012**, *11* (12), 6231-6241.
45. Dunn, W. B.; Lin, W.; Broadhurst, D.; Begley, P.; Brown, M.; Zelena, E.; Vaughan, A. A.; Halsall, A.; Harding, N.; Knowles, J. D., Molecular phenotyping of a UK population: defining the human serum metabolome. *Metabolomics* **2015**, *11* (1), 9-26.
46. Saito, K.; Maekawa, K.; Pappan, K. L.; Urata, M.; Ishikawa, M.; Kumagai, Y.; Saito, Y., Differences in metabolite profiles between blood matrices, ages, and sexes among Caucasian individuals and their inter-individual variations. *Metabolomics* **2014**, *10* (3), 0.
47. Powers, R., NMR metabolomics and drug discovery. *Magnetic Resonance in Chemistry* **2009**, *47* (S1).
48. Beckonert, O.; Keun, H. C.; Ebbels, T. M.; Bundy, J.; Holmes, E.; Lindon, J. C.; Nicholson, J. K., Metabolic profiling, metabolomic and metabonomic procedures for NMR spectroscopy of urine, plasma, serum and tissue extracts. *Nature protocols* **2007**, *2* (11), 2692-2703.
49. Choi, Y. H.; Kim, H. K.; Linthorst, H. J.; Hollander, J. G.; Lefeber, A. W.; Erkelens, C.; Nuzillard, J.-M.; Verpoorte, R., NMR Metabolomics to Revisit the Tobacco Mosaic Virus Infection in *Nicotiana tabacum* Leaves. *Journal of Natural Products* **2006**, *69* (5), 742-748.

50. Eisenmann, P.; Ehlers, M.; Weinert, C. H.; Tzvetkova, P.; Silber, M.; Rist, M. J.; Luy, B.; Muhle-Goll, C., Untargeted NMR spectroscopic analysis of the metabolic variety of new apple cultivars. *Metabolites* **2016**, *6* (3), 29.
51. Beltran, A.; Suarez, M.; Rodríguez, M. A.; Vinaixa, M.; Samino, S.; Arola, L.; Correig, X.; Yanes, O., Assessment of compatibility between extraction methods for NMR-and LC/MS-based metabolomics. *Analytical chemistry* **2012**, *84* (14), 5838-5844.
52. Glish, G. L.; Burinsky, D. J., Hybrid mass spectrometers for tandem mass spectrometry. *Journal of the American Society for Mass Spectrometry* **2008**, *19* (2), 161-172.
53. Gummer, J.; Banazis, M.; Maker, G.; Solomon, P.; Oliver, R.; Trengove, R., Use of mass spectrometry for metabolite profiling and metabolomics. *Australian Biochemist* **2009**, *40* (3), 5-8.
54. Lu, W.; Bennett, B. D.; Rabinowitz, J. D., Analytical strategies for LC-MS-based targeted metabolomics. *Journal of Chromatography B* **2008**, *871* (2), 236-242.
55. Brown, S. C.; Kruppa, G.; Dasseux, J. L., Metabolomics applications of FT - ICR mass spectrometry. *Mass Spectrometry Reviews* **2005**, *24* (2), 223-231.
56. Hendrickson, C. L.; Quinn, J. P.; Kaiser, N. K.; Smith, D. F.; Blakney, G. T.; Chen, T.; Marshall, A. G.; Weisbrod, C. R.; Beu, S. C., 21 Tesla Fourier transform ion cyclotron resonance mass spectrometer: a national resource for ultrahigh resolution mass analysis. *Journal of the American Society for Mass Spectrometry* **2015**, *26* (9), 1626-1632.
57. Lei, Z.; Huhman, D. V.; Sumner, L. W., Mass spectrometry strategies in metabolomics. *Journal of Biological Chemistry* **2011**, *286* (29), 25435-25442.
58. Zhou, R.; Li, L., Effects of sample injection amount and time-of-flight mass spectrometric detection dynamic range on metabolome analysis by high-performance chemical isotope labeling LC-MS. *Journal of proteomics* **2015**, *118*, 130-139.
59. Kopka, J., Current challenges and developments in GC-MS based metabolite profiling technology. *Journal of biotechnology* **2006**, *124* (1), 312-322.
60. Elian, A. A., GC-MS determination of gamma-hydroxybutyric acid (GHB) in blood. *Forensic science international* **2001**, *122* (1), 43-47.
61. Nishiumi, S.; Kobayashi, T.; Ikeda, A.; Yoshie, T.; Kibi, M.; Izumi, Y.; Okuno, T.; Hayashi, N.; Kawano, S.; Takenawa, T., A novel serum metabolomics-based diagnostic approach for colorectal cancer. *PLoS one* **2012**, *7* (7), e40459.
62. Dunn, W. B.; Broadhurst, D.; Begley, P.; Zelena, E.; Francis-McIntyre, S.; Anderson, N.; Brown, M.; Knowles, J. D.; Halsall, A.; Haselden, J. N., Procedures for large-scale metabolic profiling of serum and plasma using gas chromatography and liquid chromatography coupled to mass spectrometry. *Nature protocols* **2011**, *6* (7), 1060-1083.
63. Chen, J.; Wang, W.; Lv, S.; Yin, P.; Zhao, X.; Lu, X.; Zhang, F.; Xu, G., Metabonomics study of liver cancer based on ultra performance liquid chromatography coupled to mass spectrometry with HILIC and RPLC separations. *Analytica Chimica Acta* **2009**, *650* (1), 3-9.
64. Milne, S.; Ivanova, P.; Forrester, J.; Brown, H. A., Lipidomics: an analysis of cellular lipids by ESI-MS. *Methods* **2006**, *39* (2), 92-103.
65. Weinmann, W.; Lehmann, N.; Müller, C.; Wiedemann, A.; Svoboda, M., Identification of lorazepam and sildenafil as examples for the application of LC/ion-spray-MS and MS-MS with mass spectra library searching in forensic toxicology. *Forensic science international* **2000**, *113* (1), 339-344.
66. Wikoff, W. R.; Anfora, A. T.; Liu, J.; Schultz, P. G.; Lesley, S. A.; Peters, E. C.; Siuzdak, G., Metabolomics analysis reveals large effects of gut microflora on mammalian blood metabolites. *Proceedings of the national academy of sciences* **2009**, *106* (10), 3698-3703.

67. Lin, L.; Huang, Z.; Gao, Y.; Yan, X.; Xing, J.; Hang, W., LC-MS based serum metabonomic analysis for renal cell carcinoma diagnosis, staging, and biomarker discovery. *Journal of proteome research* **2011**, *10* (3), 1396-1405.
68. Xiao, J. F.; Varghese, R. S.; Zhou, B.; Ranjbar, M. R. N.; Zhao, Y.; Tsai, T.-H.; Di Poto, C.; Wang, J.; Goerlitz, D.; Luo, Y., LC-MS based serum metabolomics for identification of hepatocellular carcinoma biomarkers in Egyptian cohort. *Journal of proteome research* **2012**, *11* (12), 5914.
69. Annesley, T. M., Ion suppression in mass spectrometry. *Clinical chemistry* **2003**, *49* (7), 1041-1044.
70. Guo, K.; Li, L., Differential <sup>12</sup>C-/<sup>13</sup>C-isotope dansylation labeling and fast liquid chromatography/mass spectrometry for absolute and relative quantification of the metabolome. *Analytical chemistry* **2009**, *81* (10), 3919-3932.
71. Cech, N. B.; Enke, C. G., Practical implications of some recent studies in electrospray ionization fundamentals. *Mass Spectrometry Reviews* **2001**, *20* (6), 362-387.
72. Zhou, R.; Tseng, C.-L.; Huan, T.; Li, L., IsoMS: automated processing of LC-MS data generated by a chemical isotope labeling metabolomics platform. *Analytical chemistry* **2014**, *86* (10), 4675-4679.
73. Guo, K.; Li, L., High-Performance Isotope Labeling for Profiling Carboxylic Acid-Containing Metabolites in Biofluids by Mass Spectrometry. *Analytical Chemistry* **2010**, *82* (21), 8789-8793.
74. Zhao, S.; Luo, X.; Li, L., Chemical Isotope Labeling LC-MS for High Coverage and Quantitative Profiling of the Hydroxyl Submetabolome in Metabolomics. *Analytical chemistry* **2016**, *88* (21), 10617-10623.
75. Zhao, S.; Dawe, M.; Guo, K.; Li, L., Development of High-Performance Chemical Isotope Labeling LC-MS for Profiling the Carbonyl Submetabolome. *Analytical Chemistry* **2017**.
76. Gevi, F.; D'Alessandro, A.; Rinalducci, S.; Zolla, L., Alterations of red blood cell metabolome during cold liquid storage of erythrocyte concentrates in CPD-SAGM. *Journal of proteomics* **2012**, *76*, 168-180.
77. Beyne, P.; Vigier, J.-P.; Bourgoin, P.; Vidaud, M., Comparison of single and repeat centrifugation of blood specimens collected in BD evacuated blood collection tubes containing a clot activator for cardiac troponin I assay on the ACCESS analyzer. *Clinical chemistry* **2000**, *46* (11), 1869-1870.
78. Sampson, M.; Ruddel, M.; Albright, S.; Elin, R. J., Positive interference in lithium determinations from clot activator in collection container. *Clinical chemistry* **1997**, *43* (4), 675-679.
79. Tuck, M. K.; Chan, D. W.; Chia, D.; Godwin, A. K.; Grizzle, W. E.; Krueger, K. E.; Rom, W.; Sanda, M.; Sorbara, L.; Stass, S., Standard operating procedures for serum and plasma collection: early detection research network consensus statement standard operating procedure integration working group. *Journal of proteome research* **2009**, *8* (1), 113.
80. Barri, T.; Dragsted, L. O., UPLC-ESI-QTOF/MS and multivariate data analysis for blood plasma and serum metabolomics: effect of experimental artefacts and anticoagulant. *Analytica chimica acta* **2013**, *768*, 118-128.
81. Gonzalez-Covarrubias, V.; Dane, A.; Hankemeier, T.; Vreeken, R. J., The influence of citrate, EDTA, and heparin anticoagulants to human plasma LC-MS lipidomic profiling. *Metabolomics* **2013**, *9* (2), 337-348.
82. Chen, D.; Han, W.; Su, X.; Li, L.; Li, L., Overcoming Sample Matrix Effect in Quantitative Blood Metabolomics Using Chemical Isotope Labeling Liquid Chromatography Mass Spectrometry. *Anal. Chem* **2017**, *89* (17), 9424-9431.

83. McGarraugh, G.; Schwartz, S.; Weinstein, R., Glucose Measurements Using Blood Extracted from the Forearm and the Finger. *TheraSense, Inc* **2001**, 16.
84. Guthrie, R.; Susi, A., A simple phenylalanine method for detecting phenylketonuria in large populations of newborn infants. *Pediatrics* **1963**, 32 (3), 338-343.
85. Li, W.; Tse, F. L., Dried blood spot sampling in combination with LC - MS/MS for quantitative analysis of small molecules. *Biomedical Chromatography* **2010**, 24 (1), 49-65.
86. Zukunft, S.; Sorgenfrei, M.; Prehn, C.; Möller, G.; Adamski, J., Targeted metabolomics of dried blood spot extracts. *Chromatographia* **2013**, 76 (19-20), 1295-1305.
87. Kong, S. T.; Lin, H.-S.; Ching, J.; Ho, P. C., Evaluation of dried blood spots as sample matrix for gas chromatography/mass spectrometry based metabolomic profiling. *Analytical chemistry* **2011**, 83 (11), 4314-4318.
88. Kraly, J. R.; Holcomb, R. E.; Guan, Q.; Henry, C. S., Microfluidic applications in metabolomics and metabolic profiling. *Analytica chimica acta* **2009**, 653 (1), 23-35.
89. Chen, X.; Liu, C. C.; Li, H., Microfluidic chip for blood cell separation and collection based on crossflow filtration. *Sensors and Actuators B: Chemical* **2008**, 130 (1), 216-221.
90. Ehlert, S.; Tallarek, U., High-pressure liquid chromatography in lab-on-a-chip devices. *Analytical and bioanalytical chemistry* **2007**, 388 (3), 517-520.
91. Fortier, M.-H.; Bonneil, E.; Goodley, P.; Thibault, P., Integrated microfluidic device for mass spectrometry-based proteomics and its application to biomarker discovery programs. *Analytical Chemistry* **2005**, 77 (6), 1631-1640.
92. Chu, C. S.; Niñonuevo, M. R.; Clowers, B. H.; Perkins, P. D.; An, H. J.; Yin, H.; Killeen, K.; Miyamoto, S.; Grimm, R.; Lebrilla, C. B., Profile of native N - linked glycan structures from human serum using high performance liquid chromatography on a microfluidic chip and time - of - flight mass spectrometry. *Proteomics* **2009**, 9 (7), 1939-1951.
93. Rainville, P., Microfluidic LC-MS for analysis of small-volume biofluid samples: where we have been and where we need to go. *Bioanalysis* **2011**, 3 (1), 1-3.
94. Benjamin, D. J.; Berger, J. O.; Johannesson, M.; Nosek, B. A.; Wagenmakers, E.-J.; Berk, R.; Bollen, K. A.; Brembs, B.; Brown, L.; Camerer, C., Redefine statistical significance. *Nature Human Behaviour* **2017**.
95. MacFarland, T. W.; Yates, J. M., Mann-Whitney U Test. In *Introduction to Nonparametric Statistics for the Biological Sciences Using R*, Springer: 2016; pp 103-132.
96. Zimmerman, D. W.; Zumbo, B. D., Rank transformations and the power of the Student t test and Welch t'test for non-normal populations with unequal variances. *Canadian Journal of Experimental Psychology/Revue canadienne de psychologie expérimentale* **1993**, 47 (3), 523.
97. Vinaixa, M.; Samino, S.; Saez, I.; Duran, J.; Guinovart, J. J.; Yanes, O., A guideline to univariate statistical analysis for LC/MS-based untargeted metabolomics-derived data. *Metabolites* **2012**, 2 (4), 775-795.
98. Rouder, J. N.; Speckman, P. L.; Sun, D.; Morey, R. D.; Iverson, G., Bayesian t tests for accepting and rejecting the null hypothesis. *Psychonomic bulletin & review* **2009**, 16 (2), 225-237.
99. Sullivan, G. M.; Feinn, R., Using effect size—or why the P value is not enough. *Journal of graduate medical education* **2012**, 4 (3), 279-282.
100. Zweig, M. H.; Campbell, G., Receiver-operating characteristic (ROC) plots: a fundamental evaluation tool in clinical medicine. *Clinical chemistry* **1993**, 39 (4), 561-577.
101. Gu, Q.; Cai, Z.; Zhu, L.; Huang, B. In *Data mining on imbalanced data sets*, Advanced Computer Theory and Engineering, 2008. ICACTE'08. International Conference on, IEEE: 2008; pp 1020-1024.

102. Shlens, J., A tutorial on principal component analysis. *arXiv preprint arXiv:1404.1100* **2014**.
103. Bylesjö, M.; Rantalainen, M.; Cloarec, O.; Nicholson, J. K.; Holmes, E.; Trygg, J., OPLS discriminant analysis: combining the strengths of PLS - DA and SIMCA classification. *Journal of Chemometrics* **2006**, *20* (8 - 10), 341-351.
104. Gromski, P. S.; Muhamadali, H.; Ellis, D. I.; Xu, Y.; Correa, E.; Turner, M. L.; Goodacre, R., A tutorial review: Metabolomics and partial least squares-discriminant analysis—a marriage of convenience or a shotgun wedding. *Analytica chimica acta* **2015**, *879*, 10-23.
105. Liaw, A.; Wiener, M., Classification and regression by randomForest. *R news* **2002**, *2* (3), 18-22.
106. Mamas, M.; Dunn, W. B.; Neyses, L.; Goodacre, R., The role of metabolites and metabolomics in clinically applicable biomarkers of disease. *Archives of toxicology* **2011**, *85* (1), 5-17.
107. Deng, L.; Fang, H.; Tso, V. K.; Sun, Y.; Foshaug, R. R.; Krahn, S. C.; Zhang, F.; Yan, Y.; Xu, H.; Chang, D., Clinical validation of a novel urine-based metabolomic test for the detection of colonic polyps on Chinese population. *International journal of colorectal disease* **2017**, *32* (5), 741-743.
108. Martinez-Abraín, A., Statistical significance and biological relevance: A call for a more cautious interpretation of results in ecology. *acta oecologica* **2008**, *34* (1), 9-11.
109. Armstrong, R. A., When to use the Bonferroni correction. *Ophthalmic and Physiological Optics* **2014**, *34* (5), 502-508.
110. Dabney, A.; Storey, J. D.; Warnes, G., qvalue: Q-value estimation for false discovery rate control. *R package version* **2010**, *1* (0).
111. Cheville, A. L., Current and future trends in lymphedema management: implications for women's health. *Physical medicine and rehabilitation clinics of North America* **2007**, *18* (3), 539-553.
112. Parlee, S. D.; Ernst, M. C.; Muruganandan, S.; Sinal, C. J.; Goralski, K. B., Serum chemerin levels vary with time of day and are modified by obesity and tumor necrosis factor- $\alpha$ . *Endocrinology* **2010**, *151* (6), 2590-2602.
113. Zhou, H.; Yuen, P. S.; Pisitkun, T.; Gonzales, P. A.; Yasuda, H.; Dear, J. W.; Gross, P.; Knepper, M. A.; Star, R. A., Collection, storage, preservation, and normalization of human urinary exosomes for biomarker discovery. *Kidney international* **2006**, *69* (8), 1471-1476.
114. Roy, S. M.; Anderle, M.; Lin, H.; Becker, C. H., Differential expression profiling of serum proteins and metabolites for biomarker discovery. *International Journal of Mass Spectrometry* **2004**, *238* (2), 163-171.
115. Wu, Y.; Li, L., Determination of total concentration of chemically labeled metabolites as a means of metabolome sample normalization and sample loading optimization in mass spectrometry-based metabolomics. *Analytical chemistry* **2012**, *84* (24), 10723-10731.
116. Cohen, J., Statistical power analysis. *Current directions in psychological science* **1992**, *1* (3), 98-101.
117. Ferreira, J. A.; Zwinderman, A. H., Approximate power and sample size calculations with the Benjamini-Hochberg method. *The International Journal of Biostatistics* **2006**, *2* (1).
118. Head, M. L.; Holman, L.; Lanfear, R.; Kahn, A. T.; Jennions, M. D., The extent and consequences of p-hacking in science. *PLoS biology* **2015**, *13* (3), e1002106.
119. Song, J. W.; Chung, K. C., Observational studies: cohort and case-control studies. *Plastic and reconstructive surgery* **2010**, *126* (6), 2234.

120. Fu, F. F.; Cheng, V. W. T.; Wu, Y. M.; Tang, Y. A.; Weiner, J. H.; Li, L., Comparative Proteomic and Metabolomic Analysis of *Staphylococcus warneri* SG1 Cultured in the Presence and Absence of Butanol. *Journal of Proteome Research* **2013**, *12* (10), 4478-4489.
121. Tang, Z. M.; Guengerich, F. P., Dansylation of Unactivated Alcohols for Improved Mass Spectral Sensitivity and Application to Analysis of Cytochrome P450 Oxidation Products in Tissue Extracts. *Analytical Chemistry* **2010**, *82* (18), 7706-7712.
122. Yuan, W.; Zhang, J. X.; Li, S. W.; Edwards, J. L., Amine Metabolomics of Hyperglycemic Endothelial Cells using Capillary LC-MS with Isobaric Tagging. *Journal of Proteome Research* **2011**, *10* (11), 5242-5250.
123. Toyo'oka, T., LC-MS determination of bioactive molecules based upon stable isotope-coded derivatization method. *Journal of Pharmaceutical and Biomedical Analysis* **2012**, *69*, 174-84.
124. Dai, W. D.; Huang, Q.; Yin, P. Y.; Li, J.; Zhou, J.; Kong, H. W.; Zhao, C. X.; Lu, X.; Xu, G. W., Comprehensive and Highly Sensitive Urinary Steroid Hormone Profiling Method Based on Stable Isotope-Labeling Liquid Chromatography Mass Spectrometry. *Analytical Chemistry* **2012**, *84* (23), 10245-10251.
125. Tayyari, F.; Gowda, G. A. N.; Gu, H. W.; Raftery, D., N-15-Cholamine-A Smart Isotope Tag for Combining NMR- and MS-Based Metabolite Profiling. *Analytical Chemistry* **2013**, *85* (18), 8715-8721.
126. Leng, J. P.; Wang, H. Y.; Zhang, L.; Zhang, J.; Wang, H.; Guo, Y. L., A highly sensitive isotope-coded derivatization method and its application for the mass spectrometric analysis of analytes containing the carboxyl group. *Analytica Chimica Acta* **2013**, *758*, 114-121.
127. Bueschl, C.; Krska, R.; Kluger, B.; Schuhmacher, R., Isotopic labeling-assisted metabolomics using LC-MS. *Analytical and Bioanalytical Chemistry* **2013**, *405* (1), 27-33.
128. Bruheim, P.; Kvitvang, H. F. N.; Villas-Boas, S. G., Stable isotope coded derivatizing reagents as internal standards in metabolite profiling. *Journal of Chromatography A* **2013**, *1296*, 196-203.
129. Ulbrich, A.; Bailey, D. J.; Westphall, M. S.; Coon, J. J., Organic Acid Quantitation by NeuCode Methylamidation. *Analytical Chemistry* **2014**, *86* (9), 4402-4408.
130. Liu, P.; Huang, Y. Q.; Cai, W. J.; Yuan, B. F.; Feng, Y. Q., Profiling of Thiol-Containing Compounds by Stable Isotope Labeling Double Precursor Ion Scan Mass Spectrometry. *Analytical Chemistry* **2014**, *86* (19), 9765-9773.
131. Hao, L.; Zhong, X. F.; Greer, T.; Ye, H.; Li, L. J., Relative quantification of amine-containing metabolites using isobaric N,N-dimethyl leucine (DiLeu) reagents via LC-ESI-MS/MS and CE-ESI-MS/MS. *Analyst* **2015**, *140* (2), 467-475.
132. Dane, A. D.; Hendriks, M.; Reijmers, T. H.; Harms, A. C.; Troost, J.; Vreeken, R. J.; Boomsma, D. I.; van Duijn, C. M.; Slagboom, E. P.; Hankemeier, T., Integrating Metabolomics Profiling Measurements Across Multiple Biobanks. *Analytical Chemistry* **2014**, *86* (9), 4110-4114.
133. Lauridsen, M.; Hansen, S. H.; Jaroszewski, J. W.; Cornett, C., Human urine as test material in H-1 NMR-based metabolomics: Recommendations for sample preparation and storage. *Analytical Chemistry* **2007**, *79* (3), 1181-1186.
134. Asiago, V. M.; Gowda, G. A. N.; Zhang, S.; Shanaiah, N.; Clark, J.; Raftery, D., Use of EDTA to minimize ionic strength dependent frequency shifts in the (1)H NMR spectra of urine. *Metabolomics* **2008**, *4* (4), 328-336.
135. Xiao, C. N.; Hao, F. H.; Qin, X. R.; Wang, Y. L.; Tang, H. R., An optimized buffer system for NMR-based urinary metabolomics with effective pH control, chemical shift consistency and dilution minimization. *Analyst* **2009**, *134* (5), 916-925.

136. Akira, K.; Hichiya, H.; Shuden, M.; Morita, M.; Mitome, H., Sample preparation method to minimize chemical shift variability for NMR-based urinary metabonomics of genetically hypertensive rats. *Journal of Pharmaceutical and Biomedical Analysis* **2012**, *66*, 339-344.
137. Huan, T.; Li, L., Counting Missing Values in a Metabolite-Intensity Data Set for Measuring the Analytical Performance of a Metabolomics Platform. *Analytical Chemistry* **2015**, *87* (2), 1306-1313.
138. Tukiainen, T.; Tynkkynen, T.; Makinen, V. P.; Jylanki, P.; Kangas, A.; Hokkanen, J.; Vehtari, A.; Grohn, O.; Hallikainen, M.; Soininen, H.; Kivipelto, M.; Groop, P. H.; Kaski, K.; Laatikainen, R.; Soininen, P.; Pirttila, T.; Ala-Korpela, M., A multi-metabolite analysis of serum by H-1 NMR spectroscopy: Early systemic signs of Alzheimer's disease. *Biochemical and Biophysical Research Communications* **2008**, *375* (3), 356-361.
139. Salek, R.; Cheng, K. K.; Griffin, J., THE STUDY OF MAMMALIAN METABOLISM THROUGH NMR-BASED METABOLOMICS. In *Methods in Enzymology, Vol 500: Methods in Systems Biology*, Jameson, D.; Verma, M.; Westerhoff, H. V., Eds. Elsevier Academic Press Inc: San Diego, 2011; Vol. 500, pp 337-351.
140. Glaves, J. P.; Li, M. X.; Mercier, P.; Fahlman, R. P.; Sykes, B. D., High-throughput, multi-platform metabolomics on very small volumes: H-1 NMR metabolite identification in an unadulterated tube-in-tube system. *Metabolomics* **2014**, *10* (6), 1145-1151.
141. Guo, K.; Li, L., Differential 12 C-/ 13 C-isotope Dansylation Labeling and Fast Liquid Chromatography/Mass Spectrometry for Absolute and Relative Quantification of the Metabolome. *Analytical Chemistry* **2009**, *81* (10), 3919-3932.
142. Carta, R.; Tola, G., Solubilities of L-cystine, L-tyrosine, L-leucine, and glycine in aqueous solutions at various pHs and NaCl concentrations. *Journal of Chemical and Engineering Data* **1996**, *41* (3), 414-417.
143. Pfeiffer, C. M.; Hughes, J. P.; Cogswell, M. E.; Burt, V. L.; Lacher, D. A.; LaVoie, D. J.; Rabinowitz, D. J.; Johnson, C. L.; Pirkle, J. L., Urine Sodium Excretion Increased Slightly among US Adults between 1988 and 2010. *Journal of Nutrition* **2014**, *144* (5), 698-705.
144. Donfrancesco, C.; Ippolito, R.; Lo Noce, C.; Palmieri, L.; Iacone, R.; Russo, O.; Vanuzzo, D.; Galletti, F.; Galeone, D.; Giampaoli, S.; Strazzullo, P., Excess dietary sodium and inadequate potassium intake in Italy: Results of the MINISAL study. *Nutrition Metabolism and Cardiovascular Diseases* **2013**, *23* (9), 850-856.
145. Sieniawska, C. E.; Jung, L. C.; Olufadi, R.; Walker, V., Twenty-four-hour urinary trace element excretion: reference intervals and interpretive issues. *Annals of Clinical Biochemistry* **2012**, *49*, 341-351.
146. Schmidt, K.; Ripper, M.; Tegtmeier, I.; Humberg, E.; Sterner, C.; Reichold, M.; Warth, R.; Bandulik, S., Dynamics of Renal Electrolyte Excretion in Growing Mice. *Nephron Physiology* **2013**, *124* (3-4), 7-13.
147. Liotta, L. A.; Ferrari, M.; Petricoin, E., Clinical proteomics: written in blood. *Nature* **2003**, *425* (6961), 905-905.
148. Wróblewski, F.; Ladue, J. S., Serum glutamic pyruvic transaminase (SGP-T) in hepatic disease: A preliminary report. *Annals of Internal Medicine* **1956**, *45* (5), 801-811.
149. Mor, G.; Visintin, I.; Lai, Y.; Zhao, H.; Schwartz, P.; Rutherford, T.; Yue, L.; Bray-Ward, P.; Ward, D. C., Serum protein markers for early detection of ovarian cancer. *Proceedings of the National Academy of Sciences of the United States of America* **2005**, *102* (21), 7677-7682.
150. Shin, S.-Y.; Fauman, E. B.; Petersen, A.-K.; Krumsiek, J.; Santos, R.; Huang, J.; Arnold, M.; Erte, I.; Forgetta, V.; Yang, T.-P., An atlas of genetic influences on human blood metabolites. *Nature genetics* **2014**, *46* (6), 543-550.

151. Bundy, J. G.; Davey, M. P.; Viant, M. R., Environmental metabolomics: a critical review and future perspectives. *Metabolomics* **2009**, *5* (1), 3.
152. Serkova, N. J.; Standiford, T. J.; Stringer, K. A., The emerging field of quantitative blood metabolomics for biomarker discovery in critical illnesses. *American journal of respiratory and critical care medicine* **2011**, *184* (6), 647-655.
153. Lindenbaum, J.; Savage, D. G.; Stabler, S. P.; Allen, R. H., Diagnosis of cobalamin deficiency: II. Relative sensitivities of serum cobalamin, methylmalonic acid, and total homocysteine concentrations. *American journal of hematology* **1990**, *34* (2), 99-107.
154. Tan, B.; Qiu, Y.; Zou, X.; Chen, T.; Xie, G.; Cheng, Y.; Dong, T.; Zhao, L.; Feng, B.; Hu, X., Metabonomics identifies serum metabolite markers of colorectal cancer. *Journal of proteome research* **2013**, *12* (6), 3000-3009.
155. Schwarz, M. J.; Guillemin, G. J.; Teipel, S. J.; Buerger, K.; Hampel, H., Increased 3-hydroxykynurenine serum concentrations differentiate Alzheimer's disease patients from controls. *European archives of psychiatry and clinical neuroscience* **2013**, *263* (4), 345-352.
156. Zerwekh, J. E., Blood biomarkers of vitamin D status. *The American journal of clinical nutrition* **2008**, *87* (4), 1087S-1091S.
157. Huan, T.; Li, L., Quantitative Metabolome Analysis Based on Chromatographic Peak Reconstruction in Chemical Isotope Labeling Liquid Chromatography Mass Spectrometry. *Analytical Chemistry* **2015**, *87* (14), 7011-7016.
158. Wishart, D. S.; Tzur, D.; Knox, C.; Eisner, R.; Guo, A. C.; Young, N.; Cheng, D.; Jewell, K.; Arndt, D.; Sawhney, S., HMDB: the human metabolome database. *Nucleic acids research* **2007**, *35* (suppl\_1), D521-D526.
159. Huan, T.; Wu, Y.; Tang, C.; Lin, G.; Li, L., DnsID in MyCompoundID for Rapid Identification of Dansylated Amine- and Phenol-Containing Metabolites in LC-MS-Based Metabolomics. *Analytical chemistry* **2015**, *87* (19), 9838-9845.
160. Matuszewski, B.; Constanzer, M.; Chavez-Eng, C., Matrix effect in quantitative LC/MS/MS analyses of biological fluids: a method for determination of finasteride in human plasma at picogram per milliliter concentrations. *Analytical Chemistry* **1998**, *70* (5), 882-889.
161. Han, W.; Li, L., Matrix effect on chemical isotope labeling and its implication in metabolomic sample preparation for quantitative metabolomics. *Metabolomics* **2015**, *11* (6), 1733-1742.
162. Barr, D. B.; Wilder, L. C.; Caudill, S. P.; Gonzalez, A. J.; Needham, L. L.; Pirkle, J. L., Urinary creatinine concentrations in the US population: implications for urinary biologic monitoring measurements. *Environmental health perspectives* **2005**, *113* (2), 192.
163. Khan, S. M.; Franke-Fayard, B.; Mair, G. R.; Lasonder, E.; Janse, C. J.; Mann, M.; Waters, A. P., Proteome analysis of separated male and female gametocytes reveals novel sex-specific Plasmodium biology. *Cell* **2005**, *121* (5), 675-687.
164. Markle, J. G.; Frank, D. N.; Mortin-Toth, S.; Robertson, C. E.; Feazel, L. M.; Rolle-Kampczyk, U.; von Bergen, M.; McCoy, K. D.; Macpherson, A. J.; Danska, J. S., Sex differences in the gut microbiome drive hormone-dependent regulation of autoimmunity. *Science* **2013**, *339* (6123), 1084-1088.
165. Xia, J.; Wishart, D. S., Using metaboanalyst 3.0 for comprehensive metabolomics data analysis. *Current Protocols in Bioinformatics* **2016**, *14.10*, 1-14.10. 91.
166. Thévenot, E. A.; Roux, A.; Xu, Y.; Ezan, E.; Junot, C., Analysis of the human adult urinary metabolome variations with age, body mass index, and gender by implementing a comprehensive workflow for univariate and OPLS statistical analyses. *Journal of proteome research* **2015**, *14* (8), 3322-3335.

167. Kobayashi, R.; Shimomura, Y.; Murakami, T.; Nakai, N.; Fujitsuka, N.; Otsuka, M.; Arakawa, N.; POPOV, M. K.; HARRIS, A. R., Gender difference in regulation of branched-chain amino acid catabolism. *Biochemical Journal* **1997**, *327* (2), 449-453.
168. Lamont, L. S.; McCullough, A. J.; Kalhan, S. C., Gender differences in the regulation of amino acid metabolism. *Journal of Applied Physiology* **2003**, *95* (3), 1259-1265.
169. Kawaguchi, T.; Nagao, Y.; Matsuoka, H.; Ide, T.; Sata, M., Branched-chain amino acid-enriched supplementation improves insulin resistance in patients with chronic liver disease. *International journal of molecular medicine* **2008**, *22* (1), 105-112.
170. Koreen, A. R.; Lieberman, J.; Alvir, J.; Mayerhoff, D.; Loebel, A.; Chakos, M.; Amin, F.; Cooper, T., Plasma homovanillic acid levels in first-episode schizophrenia: psychopathology and treatment response. *Archives of general psychiatry* **1994**, *51* (2), 132-138.
171. Bareggi, S. R.; Franceschi, M.; Bonini, L.; Zecca, L.; Smirne, S., Decreased CSF concentrations of homovanillic acid and  $\gamma$ -aminobutyric acid in Alzheimer's disease: Age-or disease-related modifications? *Archives of Neurology* **1982**, *39* (11), 709-712.
172. Fisher, M.; Yousef, I., Sex differences in the bile acid composition of human bile: studies in patients with and without gallstones. *Canadian Medical Association Journal* **1973**, *109* (3), 190.
173. Soltow, Q. A.; Jones, D. P.; Promislow, D. E., A network perspective on metabolism and aging. *Integrative and comparative biology* **2010**, *50* (5), 844-854.
174. Yu, Z.; Zhai, G.; Singmann, P.; He, Y.; Xu, T.; Prehn, C.; Römisch - Margl, W.; Lattka, E.; Gieger, C.; Soranzo, N., Human serum metabolic profiles are age dependent. *Aging cell* **2012**, *11* (6), 960-967.
175. Avanesov, A. S.; Ma, S.; Pierce, K. A.; Yim, S. H.; Lee, B. C.; Clish, C. B.; Gladyshev, V. N., Age-and diet-associated metabolome remodeling characterizes the aging process driven by damage accumulation. *Elife* **2014**, *3*, e02077.
176. Xie, B.; Waters, M. J.; Schirra, H. J., Investigating potential mechanisms of obesity by metabolomics. *BioMed Research International* **2012**, *2012*.
177. Rauschert, S.; Uhl, O.; Koletzko, B.; Hellmuth, C., Metabolomic biomarkers for obesity in humans: a short review. *Annals of Nutrition and Metabolism* **2014**, *64* (3-4), 314-324.
178. Morio, B.; Comte, B.; Martin, J.-F.; Chanseau, E.; Alligier, M.; Junot, C.; Lyan, B.; Boirie, Y.; Vidal, H.; Laville, M., Metabolomics reveals differential metabolic adjustments of normal and overweight subjects during overfeeding. *Metabolomics* **2015**, *11* (4), 920-938.
179. Zhang, M.; Bi, L.; Fang, J.; Su, X.; Da, G.; Kuwamori, T.; Kagamimori, S., Beneficial effects of taurine on serum lipids in overweight or obese non-diabetic subjects. *Amino acids* **2004**, *26* (3), 267-271.
180. Joshi - Barve, S.; Barve, S. S.; Amancherla, K.; Gobejishvili, L.; Hill, D.; Cave, M.; Hote, P.; McClain, C. J., Palmitic acid induces production of proinflammatory cytokine interleukin - 8 from hepatocytes. *Hepatology* **2007**, *46* (3), 823-830.
181. Jones, D. P.; Park, Y.; Ziegler, T. R., Nutritional metabolomics: progress in addressing complexity in diet and health. *Annual review of nutrition* **2012**, *32*, 183-202.
182. Robertson, D. G.; Watkins, P. B.; Reilly, M. D., Metabolomics in toxicology: preclinical and clinical applications. *Toxicological Sciences* **2010**, *120* (suppl\_1), S146-S170.
183. Rappaport, S. M.; Smith, M. T., Environment and Disease Risks. *Science* **2010**, *330* (6003), 460-461.
184. Hunter, D. J., Gene-environment interactions in human diseases. *Nature reviews. Genetics* **2005**, *6* (4), 287.

185. Calne, D.; McGeer, E.; Eisen, A.; Spencer, P., Alzheimer's disease, Parkinson's disease, and motoneuron disease: abiotropic interaction between ageing and environment? *The Lancet* **1986**, *328* (8515), 1067-1070.
186. Collino, S.; Martin, F. P. J.; Rezzi, S., Clinical metabolomics paves the way towards future healthcare strategies. *British journal of clinical pharmacology* **2013**, *75* (3), 619-629.
187. Kim, H.-J.; Kim, J. H.; Noh, S.; Hur, H. J.; Sung, M. J.; Hwang, J.-T.; Park, J. H.; Yang, H. J.; Kim, M.-S.; Kwon, D. Y., Metabolomic analysis of livers and serum from high-fat diet induced obese mice. *Journal of proteome research* **2010**, *10* (2), 722-731.
188. Gu, F.; Derkach, A.; Freedman, N. D.; Landi, M. T.; Albanes, D.; Weinstein, S. J.; Mondul, A. M.; Matthews, C. E.; Guertin, K. A.; Xiao, Q., Cigarette smoking behaviour and blood metabolomics. *International journal of epidemiology* **2015**, *45* (5), 1421-1432.
189. Fannin, R. D.; Russo, M.; O'connell, T. M.; Gerrish, K.; Winnike, J. H.; Macdonald, J.; Newton, J.; Malik, S.; Sieber, S. O.; Parker, J., Acetaminophen dosing of humans results in blood transcriptome and metabolome changes consistent with impaired oxidative phosphorylation. *Hepatology* **2010**, *51* (1), 227-236.
190. Vlaanderen, J.; Janssen, N.; Hoek, G.; Keski-Rahkonen, P.; Barupal, D.; Cassee, F.; Gosens, I.; Strak, M.; Steenhof, M.; Lan, Q., The impact of ambient air pollution on the human blood metabolome. *Environmental Research* **2017**, *156*, 341-348.
191. Chorell, E.; Moritz, T.; Branth, S.; Antti, H.; Svensson, M. B., Predictive metabolomics evaluation of nutrition-modulated metabolic stress responses in human blood serum during the early recovery phase of strenuous physical exercise. *Journal of proteome research* **2009**, *8* (6), 2966-2977.
192. Kim, J.-H.; Woenker, T.; Adamec, J.; Regnier, F. E., Simple, miniaturized blood plasma extraction method. *Analytical chemistry* **2013**, *85* (23), 11501-11508.
193. Vella, S. J.; Beattie, P.; Cademartiri, R.; Laromaine, A.; Martinez, A. W.; Phillips, S. T.; Mirica, K. A.; Whitesides, G. M., Measuring markers of liver function using a micropatterned paper device designed for blood from a fingerstick. *Analytical chemistry* **2012**, *84* (6), 2883-2891.
194. Fruhstorfer, H.; Schmelzeisen-Redeker, G. t.; Weiss, T., Capillary blood sampling: relation between lancet diameter, lancing pain and blood volume. *European Journal of Pain* **1999**, *3* (3), 283-286.
195. Fruhstorfer, H., Capillary blood sampling: the pain of single - use lancing devices. *European Journal of Pain* **2000**, *4* (3), 301-305.
196. Fan, R.; Vermesh, O.; Srivastava, A.; Yen, B. K.; Qin, L.; Ahmad, H.; Kwong, G. A.; Liu, C.-C.; Gould, J.; Hood, L., Integrated barcode chips for rapid, multiplexed analysis of proteins in microliter quantities of blood. *Nature biotechnology* **2008**, *26* (12), 1373-1378.
197. Grady, M.; Pineau, M.; Pynes, M. K.; Katz, L. B.; Ginsberg, B., A clinical evaluation of routine blood sampling practices in patients with diabetes: impact on fingerstick blood volume and pain. *Journal of diabetes science and technology* **2014**, *8* (4), 691-698.
198. Haerberle, S.; Brenner, T.; Zengerle, R.; Ducreé, J., Centrifugal extraction of plasma from whole blood on a rotating disk. *Lab on a Chip* **2006**, *6* (6), 776-781.
199. Wong, A. P.; Gupta, M.; Shevkoplyas, S. S.; Whitesides, G. M., Egg beater as centrifuge: isolating human blood plasma from whole blood in resource-poor settings. *Lab on a Chip* **2008**, *8* (12), 2032-2037.
200. Bhamla, M. S.; Benson, B.; Chai, C.; Katsikis, G.; Johri, A.; Prakash, M., Hand-powered ultralow-cost paper centrifuge. *Nature Biomedical Engineering* **2017**, *1*, 0009.
201. Li, Z.; Tatlay, J.; Li, L., Nanoflow LC-MS for High-Performance Chemical Isotope Labeling Quantitative Metabolomics. *Analytical chemistry* **2015**, *87* (22), 11468-11474.

202. Warrack, B. M.; Hnatyshyn, S.; Ott, K.-H.; Reily, M. D.; Sanders, M.; Zhang, H.; Drexler, D. M., Normalization strategies for metabonomic analysis of urine samples. *Journal of Chromatography B* **2009**, *877* (5), 547-552.
203. Walsh, M. C.; Brennan, L.; Malthouse, J. P. G.; Roche, H. M.; Gibney, M. J., Effect of acute dietary standardization on the urinary, plasma, and salivary metabolomic profiles of healthy humans. *The American journal of clinical nutrition* **2006**, *84* (3), 531-539.
204. Guertin, K. A.; Loftfield, E.; Boca, S. M.; Sampson, J. N.; Moore, S. C.; Xiao, Q.; Huang, W.-Y.; Xiong, X.; Freedman, N. D.; Cross, A. J., Serum biomarkers of habitual coffee consumption may provide insight into the mechanism underlying the association between coffee consumption and colorectal cancer. *The American journal of clinical nutrition* **2015**, *101* (5), 1000-1011.
205. Le Poul, E.; Loison, C.; Struyf, S.; Springael, J.-Y.; Lannoy, V.; Decobecq, M.-E.; Brezillon, S.; Dupriez, V.; Vassart, G.; Van Damme, J., Functional characterization of human receptors for short chain fatty acids and their role in polymorphonuclear cell activation. *Journal of Biological Chemistry* **2003**, *278* (28), 25481-25489.
206. Tedelind, S.; Westberg, F.; Kjerrulf, M.; Vidal, A., Anti-inflammatory properties of the short-chain fatty acids acetate and propionate: a study with relevance to inflammatory bowel disease. *World journal of gastroenterology: WJG* **2007**, *13* (20), 2826.
207. Nilsson, N. E.; Kotarsky, K.; Owman, C.; Olde, B., Identification of a free fatty acid receptor, FFA 2 R, expressed on leukocytes and activated by short-chain fatty acids. *Biochemical and biophysical research communications* **2003**, *303* (4), 1047-1052.
208. Zhou, L.; Wang, Q.; Yin, P.; Xing, W.; Wu, Z.; Chen, S.; Lu, X.; Zhang, Y.; Lin, X.; Xu, G., Serum metabolomics reveals the deregulation of fatty acids metabolism in hepatocellular carcinoma and chronic liver diseases. *Analytical and bioanalytical chemistry* **2012**, *403* (1), 203-213.
209. Jensen, M. D.; Haymond, M. W.; Rizza, R. A.; Cryer, P. E.; Miles, J., Influence of body fat distribution on free fatty acid metabolism in obesity. *Journal of Clinical Investigation* **1989**, *83* (4), 1168.
210. Yang, J.; Xu, G.; Hong, Q.; Liebich, H. M.; Lutz, K.; Schmülling, R.-M.; Wahl, H. G., Discrimination of Type 2 diabetic patients from healthy controls by using metabolomics method based on their serum fatty acid profiles. *Journal of Chromatography B* **2004**, *813* (1), 53-58.
211. Cherrington, C.; Hinton, M.; Mead, G.; Chopra, I., Organic acids: chemistry, antibacterial activity and practical applications. *Advances in microbial physiology* **1991**, *32*, 87-108.
212. Ricke, S., Perspectives on the use of organic acids and short chain fatty acids as antimicrobials. *Poultry science* **2003**, *82* (4), 632-639.
213. Roth, F.; Kirchgessner, M., Organic acids as feed additives for young pigs: Nutritional and gastrointestinal effects. *J. Anim. Feed Sci* **1998**, *7* (Suppl 1), 25-33.
214. Hodson, L.; Skeaff, C. M.; Fielding, B. A., Fatty acid composition of adipose tissue and blood in humans and its use as a biomarker of dietary intake. *Progress in lipid research* **2008**, *47* (5), 348-380.
215. Chesney, R., Taurine: its biological role and clinical implications. *Advances in pediatrics* **1984**, *32*, 1-42.
216. Engelborghs, S.; Marescau, B.; De Deyn, P., Amino acids and biogenic amines in cerebrospinal fluid of patients with Parkinson's disease. *Neurochemical research* **2003**, *28* (8), 1145-1150.
217. Blacklock, C.; Lawrence, J.; Wiles, D.; Malcolm, E.; Gibson, I.; Kelly, C.; Paterson, J., Salicylic acid in the serum of subjects not taking aspirin. Comparison of salicylic acid

- concentrations in the serum of vegetarians, non-vegetarians, and patients taking low dose aspirin. *Journal of clinical pathology* **2001**, *54* (7), 553-555.
218. Rodrigues, C.; Fan, G.; Ma, X.; Kren, B. T.; Steer, C. J., A novel role for ursodeoxycholic acid in inhibiting apoptosis by modulating mitochondrial membrane perturbation. *Journal of Clinical Investigation* **1998**, *101* (12), 2790.
219. Pratheeshkumar, P.; Raphael, T.; Kuttan, G., Protective Role of Perillic Acid Against Radiation- Induced Oxidative Stress, Cytokine Profile, DNA Damage, and Intestinal Toxicity in Mice. *Journal of Environmental Pathology, Toxicology and Oncology* **2010**, *29* (3).
220. Yeruva, L.; Pierre, K. J.; Elegbede, A.; Wang, R. C.; Carper, S. W., Perillyl alcohol and perillic acid induced cell cycle arrest and apoptosis in non small cell lung cancer cells. *Cancer letters* **2007**, *257* (2), 216-226.
221. Rahuman, A. A.; Gopalakrishnan, G.; Ghouse, B. S.; Arumugam, S.; Himalayan, B., Effect of *Feronia limonia* on mosquito larvae. *Fitoterapia* **2000**, *71* (5), 553-555.
222. Rodriguez-Garcia, I.; Guil-Guerrero, J. L., Evaluation of the antioxidant activity of three microalgal species for use as dietary supplements and in the preservation of foods. *Food chemistry* **2008**, *108* (3), 1023-1026.
223. Hamilton, E.; Berg, H.; Easton, C.; Bakker, A. J., The effect of taurine depletion on the contractile properties and fatigue in fast-twitch skeletal muscle of the mouse. *Amino acids* **2006**, *31* (3), 273-278.
224. Statland, B. E.; Demas, T. J., Serum caffeine half-lives: healthy subjects vs. patients having alcoholic hepatic disease. *American journal of clinical pathology* **1980**, *73* (3), 390-393.
225. Warskulat, U.; Flögel, U.; Jacoby, C.; Hartwig, H.-G.; Thewissen, M.; Merx, M. W.; Molojavyi, A.; Heller-Stilb, B.; Schrader, J.; Häussinger, D., Taurine transporter knockout depletes muscle taurine levels and results in severe skeletal muscle impairment but leaves cardiac function uncompromised. *The FASEB journal* **2004**, *18* (3), 577-579.
226. Yatabe, Y.; Miyakawa, S.; Ohmori, H.; Adachi, H. M. T., Effects of taurine administration on exercise. In *Taurine 7*, Springer: 2009; pp 245-252.
227. Manabe, S.; Kuroda, I.; Okada, K.; Morishima, M.; Okamoto, M.; Harada, N.; Takahashi, A.; Sakai, K.; Nakaya, Y., Decreased blood levels of lactic acid and urinary excretion of 3-methylhistidine after exercise by chronic taurine treatment in rats. *Journal of nutritional science and vitaminology* **2003**, *49* (6), 375-380.
228. Cairns, S. P., Lactic acid and exercise performance. *Sports Medicine* **2006**, *36* (4), 279-291.
229. Walker, R., Toxicology of sorbic acid and sorbates. *Food Additives & Contaminants* **1990**, *7* (5), 671-676.
230. He, W.; Miao, F. J.-P.; Lin, D. C.-H.; Schwandner, R. T.; Wang, Z.; Gao, J.; Chen, J.-L.; Tian, H.; Ling, L., Citric acid cycle intermediates as ligands for orphan G-protein-coupled receptors. *Nature* **2004**, *429* (6988), 188-194.
231. Krebs, H. A.; Johnson, W. A., The role of citric acid in intermediate metabolism in animal tissues. *Enzymologia* **1937**, *4*, 148-156.
232. Berg, D.; Postuma, R. B.; Bloem, B.; Chan, P.; Dubois, B.; Gasser, T.; Goetz, C. G.; Halliday, G. M.; Hardy, J.; Lang, A. E., *Annals of Neurology Mov. Disord.* **2014**, *29* (4), 454-462.
233. Miller, D. B.; O'Callaghan, J. P., Biomarkers of Parkinson's disease: present and future. *Metabolism*. **2015**, *64* (3), S40-S46.
234. Hughes, A. J.; Daniel, S. E.; Lees, A. J., Improved accuracy of clinical diagnosis of Lewy body Parkinson's disease. *Neurology* **2001**, *57* (8), 1497-1499.

235. Newman, E. J.; Breen, K.; Patterson, J.; Hadley, D. M.; Grosset, K. A.; Grosset, D. G., Accuracy of Parkinson's disease diagnosis in 610 general practice patients in the West of Scotland. *Mov. Disord.* **2009**, *24* (16), 2379-2385.
236. Adler, C. H.; Beach, T. G.; Hentz, J. G.; Shill, H. A.; Caviness, J. N.; Driver-Dunckley, E.; Sabbagh, M. N.; Sue, L. I.; Jacobson, S. A.; Belden, C. M., Low clinical diagnostic accuracy of early vs advanced Parkinson disease Clinicopathologic study. *Neurology* **2014**, *83* (5), 406-412.
237. Svenningsson, P.; Westman, E.; Ballard, C.; Aarsland, D., Cognitive impairment in patients with Parkinson's disease: diagnosis, biomarkers, and treatment. *The Lancet Neurology* **2012**, *11* (8), 697-707.
238. Xia, J.; Broadhurst, D. I.; Wilson, M.; Wishart, D. S., Translational biomarker discovery in clinical metabolomics: an introductory tutorial. *Metabolomics* **2013**, *9* (2), 280-299.
239. LeWitt, P. A.; Li, J.; Lu, M.; Beach, T. G.; Adler, C. H.; Guo, L., 3 - hydroxykynurenine and other Parkinson's disease biomarkers discovered by metabolomic analysis. *Movement Disorders* **2013**, *28* (12), 1653-1660.
240. Bogdanov, M.; Matson, W. R.; Wang, L.; Matson, T.; Saunders-Pullman, R.; Bressman, S. S.; Beal, M. F., Metabolomic profiling to develop blood biomarkers for Parkinson's disease. *Brain* **2008**, *131* (2), 389-396.
241. Johansen, K. K.; Wang, L.; Aasly, J. O.; White, L. R.; Matson, W. R.; Henchcliffe, C.; Beal, M. F.; Bogdanov, M., Metabolomic profiling in LRRK2-related Parkinson's disease. *PLoS One* **2009**, *4* (10), e7551.
242. Hatano, T.; Saiki, S.; Okuzumi, A.; Mohny, R. P.; Hattori, N., Identification of novel biomarkers for Parkinson's disease by metabolomic technologies. *J. Neurol. Neurosurg. Psychiatry* **2015**, jnnp-2014-309676.
243. Camicioli, R.; Sabino, J.; Gee, M.; Bouchard, T.; Fisher, N.; Hanstock, C.; Emery, D.; Martin, W., Ventricular dilatation and brain atrophy in patients with Parkinson's disease with incipient dementia. *Mov. Disord.* **2011**, *26* (8), 1443-1450.
244. Sapkota, S.; Gee, M.; Sabino, J.; Emery, D.; Camicioli, R., Association of homocysteine with ventricular dilatation and brain atrophy in Parkinson's disease. *Mov. Disord.* **2014**, *29* (3), 368-374.
245. Chia, L.-G.; Cheng, L.-J.; Chuo, L.-J.; Cheng, F.-C.; Cu, J.-S., Studies of dementia, depression, electrophysiology and cerebrospinal fluid monoamine metabolites in patients with Parkinson's disease. *J. Neurol. Sci.* **1995**, *133* (1), 73-78.
246. Hoops, S.; Nazem, S.; Siderowf, A.; Duda, J.; Xie, S.; Stern, M.; Weintraub, D., Validity of the MoCA and MMSE in the detection of MCI and dementia in Parkinson disease. *Neurology* **2009**, *73* (21), 1738-1745.
247. Gerlach, M.; Maetzler, W.; Broich, K.; Hampel, H.; Rems, L.; Reum, T.; Riederer, P.; Stöfler, A.; Streffer, J.; Berg, D., Biomarker candidates of neurodegeneration in Parkinson's disease for the evaluation of disease-modifying therapeutics. *J. Neural Transm.* **2012**, *119* (1), 39-52.
248. Hinterberger, H.; Andrews, C. J., Catecholamine metabolism during oral administration of Levodopa: effects of the medication in Parkinson's disease. *Arch. Neurol.* **1972**, *26* (3), 245-252.
249. Factor, S. A.; Schneider, A. S., Peripheral catecholamine output in Parkinson's disease: effects of drug treatment. *Exp. Neurol.* **1995**, *131* (1), 64-68.
250. Wishart, D. S.; Tzur, D.; Knox, C.; Eisner, R.; Guo, A. C.; Young, N.; Cheng, D.; Jewell, K.; Arndt, D.; Sawhney, S., HMDB: the human metabolome database. *Nucleic Acids Res.* **2007**, *35* (suppl 1), D521-D526.

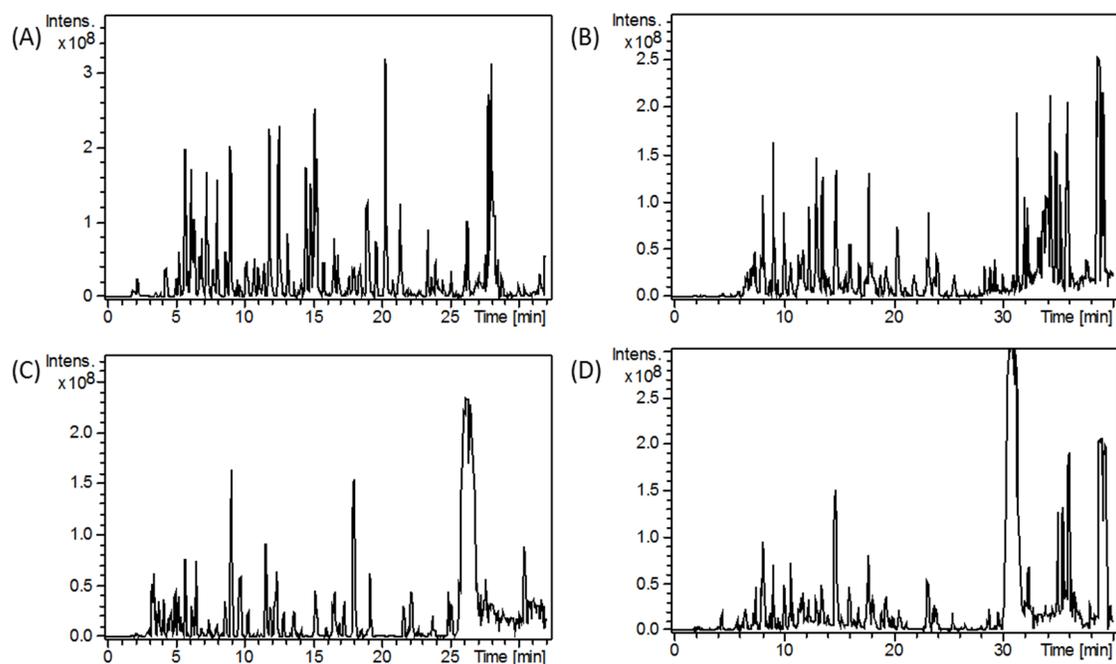
251. Culvenor, A. J.; Jarrott, B., Reduction of Aromatic L-Amino Acid Decarboxylase Protein in Rats after Chronic Administration of Alpha-Methyl dopa. *Mol. Pharmacol.* **1979**, *15* (1), 86-98.
252. Coune, P. G.; Schneider, B. L.; Aebischer, P., Parkinson's disease: gene therapies. *Cold Spring Harbor perspectives in medicine* **2012**, *2* (4), a009431.
253. Bertoldi, M.; Dominici, P.; Moore, P. S.; Maras, B.; Voltattorni, C. B., Reaction of dopa decarboxylase with  $\alpha$ -methyl dopa leads to an oxidative deamination producing 3, 4-dihydroxyphenylacetone, an active site directed affinity label. *Biochemistry* **1998**, *37* (18), 6552-6561.
254. Vamos, E.; Pardutz, A.; Klivenyi, P.; Toldi, J.; Vecsei, L., The role of kynurenines in disorders of the central nervous system: possibilities for neuroprotection. *J. Neurol. Sci.* **2009**, *283* (1), 21-27.
255. Chiarugi, A.; Meli, E.; Moroni, F., Similarities and differences in the neuronal death processes activated by 3OH - kynurenine and quinolinic acid. *J. Neurochem.* **2001**, *77* (5), 1310-1318.
256. Hartai, Z.; Klivenyi, P.; Janaky, T.; Penke, B.; Dux, L.; Vecsei, L., Kynurenine metabolism in plasma and in red blood cells in Parkinson's disease. *J. Neurol. Sci.* **2005**, *239* (1), 31-35.
257. Hu, G.; Bidel, S.; Jousilahti, P.; Antikainen, R.; Tuomilehto, J., Coffee and tea consumption and the risk of Parkinson's disease. *Mov. Disord.* **2007**, *22* (15), 2242-2248.
258. Schwarzschild, M. A.; Xu, K.; Oztas, E.; Petzer, J. P.; Castagnoli, K.; Castagnoli, N.; Chen, J.-F., Neuroprotection by caffeine and more specific A2A receptor antagonists in animal models of Parkinson's disease. *Neurology* **2003**, *61* (11 suppl 6), S55-S61.
259. Morelli, M.; Di Paolo, T.; Wardas, J.; Calon, F.; Xiao, D.; Schwarzschild, M. A., Role of adenosine A 2A receptors in parkinsonian motor impairment and L-DOPA-induced motor complications. *Prog. Neurobiol.* **2007**, *83* (5), 293-309.
260. Yasui, K.; Agematsu, K.; Shinozaki, K.; Hokibara, S.; Nagumo, H.; Nakazawa, T.; Komiyama, A., Theophylline induces neutrophil apoptosis through adenosine A2A receptor antagonism. *J. Leukocyte Biol.* **2000**, *67* (4), 529-535.
261. Mashima, R.; Nakanishi-Ueda, T.; Yamamoto, Y., Simultaneous determination of methionine sulfoxide and methionine in blood plasma using gas chromatography-mass spectrometry. *Anal. Biochem.* **2003**, *313* (1), 28-33.
262. Glaser, C. B.; Yamin, G.; Uversky, V. N.; Fink, A. L., Methionine oxidation,  $\alpha$ -synuclein and Parkinson's disease. *Biochim. Biophys. Acta, Proteins Proteomics* **2005**, *1703* (2), 157-169.
263. Akashi, K.; Miyake, C.; Yokota, A., Citrulline, a novel compatible solute in drought-tolerant wild watermelon leaves, is an efficient hydroxyl radical scavenger. *FEBS Lett.* **2001**, *508* (3), 438-442.
264. Halliwell, B., Oxidative stress and neurodegeneration: where are we now? *J. Neurochem.* **2006**, *97* (6), 1634-1658.
265. Zhou, C.; Huang, Y.; Przedborski, S., Oxidative stress in Parkinson's disease. *Ann. N. Y. Acad. Sci.* **2008**, *1147* (1), 93-104.
266. Xu, H. F.; O'Brochta, D. A., Advanced technologies for genetically manipulating the silkworm *Bombyx mori*, a model Lepidopteran insect. *Proceedings of the Royal Society B-Biological Sciences* **2015**, *282* (1810), 20150487.
267. Chi, Y.; Qiao, K.; Jiang, H.; Lin, R.; Wang, K., Comparison of Two Acute Toxicity Test Methods for the Silkworm (Lepidoptera: Bombycidae). *Journal of Economic Entomology* **2015**, *108* (1), 145-149.
268. Zhou, L.; Chen, X.; Shao, Z.; Zhou, P.; Knight, D. P.; Vollrath, F., Copper in the silk formation process of *Bombyx mori* silkworm. *FEBS Letters* **2003**, *554* (3), 337-341.

269. Wang, J.-X.; Bian, Y.-M., Fluoride effects on the mulberry-silkworm system. *Environmental Pollution* **1988**, *52* (1), 11-18.
270. Sekimizu, N.; Paudel, A.; Hamamoto, H., Animal welfare and use of silkworm as a model animal. *Drug discoveries & therapeutics* **2012**, *6* (4), 226-9.
271. Xia, Q. Y.; Li, S.; Feng, Q. L., Advances in Silkworm Studies Accelerated by the Genome Sequencing of *Bombyx mori*. In *Annual Review of Entomology, Vol 59, 2014*, Berenbaum, M. R., Ed. Annual Reviews: Palo Alto, 2014; Vol. 59, pp 513-536.
272. Xia, Q. Y.; Zhou, Z. Y.; Lu, C.; Cheng, D. J.; Dai, F. Y.; Li, B.; Zhao, P.; Zha, X. F.; Cheng, T. C.; Chai, C. L.; Pan, G. Q.; Xu, J. S.; Liu, C.; Lin, Y.; Qian, J. F.; Hou, Y.; Wu, Z. L.; Li, G. R.; Pan, M. H.; Li, C. F.; Shen, Y. H.; Lan, X. Q.; Yuan, L. W.; Li, T.; Xu, H. F.; Yang, G. W.; Wan, Y. J.; Zhu, Y.; Yu, M. D.; Shen, W. D.; Wu, D. Y.; Xiang, Z. H.; Yu, J.; Wang, J.; Li, R. Q.; Shi, J. P.; Li, H.; Li, G. Y.; Su, J. N.; Wang, X. L.; Li, G. Q.; Zhang, Z. J.; Wu, Q. F.; Li, J.; Zhang, Q. P.; Wei, N.; Xu, J. Z.; Sun, H. B.; Dong, L.; Liu, D. Y.; Zhao, S. L.; Zhao, X. L.; Meng, Q. S.; Lan, F. D.; Huang, X. G.; Li, Y. Z.; Fang, L.; Li, D. W.; Sun, Y. Q.; Zhang, Z. P.; Yang, Z.; Huang, Y. Q.; Xi, Y.; Qi, Q. H.; He, D. D.; Huang, H. Y.; Zhang, X. W.; Wang, Z. Q.; Li, W. J.; Cao, Y. Z.; Yu, Y. P.; Yu, H.; Li, J. H.; Ye, J. H.; Chen, H.; Zhou, Y.; Liu, B.; Ye, J.; Ji, H.; Li, S. T.; Ni, P. X.; Zhang, J. G.; Zhang, Y.; Zheng, H. K.; Mao, B. Y.; Wang, W.; Ye, C.; Li, S. G.; Wong, G. K. S.; Yang, H. M., A draft sequence for the genome of the domesticated silkworm (*Bombyx mori*). *Science* **2004**, *306* (5703), 1937-1940.
273. International Silkworm Genome Consortium, The genome of a lepidopteran model insect, the silkworm *Bombyx mori*. *Insect Biochemistry and Molecular Biology* **2008**, *38* (12), 1036-45.
274. Ou, J.; Deng, H.-M.; Zheng, S.-C.; Huang, L.-H.; Feng, Q.-L.; Liu, L., Transcriptomic analysis of developmental features of *Bombyx mori* wing disc during metamorphosis. *BMC Genomics* **2014**, *15* (1), 820.
275. Wang, H.; Fang, Y.; Wang, L.; Zhu, W.; Ji, H.; Wang, H.; Xu, S.; Sima, Y., Transcriptome analysis of the *Bombyx mori* fat body after constant high temperature treatment shows differences between the sexes. *Molecular Biology Reports* **2014**, *41* (9), 6039-6049.
276. Shi, X. F.; Li, Y. N.; Yi, Y. Z.; Xiao, X. G.; Zhang, Z. F., Identification and Characterization of 30 K Protein Genes Found in *Bombyx mori* (Lepidoptera: Bombycidae) Transcriptome. *Journal of Insect Science* **2015**, *15* (1), 71-71.
277. Chen, J.-E.; Li, J.-Y.; You, Z.-Y.; Liu, L.-L.; Liang, J.-S.; Ma, Y.-Y.; Chen, M.; Zhang, H.-R.; Jiang, Z.-D.; Zhong, B.-X., Proteome Analysis of Silkworm, *Bombyx mori*, Larval Gonads: Characterization of Proteins Involved in Sexual Dimorphism and Gametogenesis. *Journal of Proteome Research* **2013**, *12* (6), 2422-2438.
278. Hu, X.; Zhu, M.; Wang, S.; Zhu, L.; Xue, R.; Cao, G.; Gong, C., Proteomics analysis of digestive juice from silkworm during *Bombyx mori* nucleopolyhedrovirus infection. *Proteomics* **2015**, *15* (15), 2691-2700.
279. Wang, G.-B.; Zheng, Q.; Shen, Y.-W.; Wu, X.-F., Shotgun proteomic analysis of *Bombyx mori* brain: emphasis on regulation of behavior and development of the nervous system. *Insect Science* **2016**, *23* (1), 15-27.
280. Yin, W. M.; Xu, X.; He, Y.; Wei, G. B.; Sima, Y. H.; Shi-Qing, X., Metabonomic Analysis of *Bombyx mori* (Heterocera: Bombycidae) Treated With Acetaminophen. *Journal of Insect Science* **2014**, *14* (1), 225-225.
281. Zhou, L.; Li, H.; Hao, F.; Li, N.; Liu, X.; Wang, G.; Wang, Y.; Tang, H., Developmental Changes for the Hemolymph Metabolome of Silkworm (*Bombyx mori* L.). *Journal of Proteome Research* **2015**, *14* (5), 2331-2347.

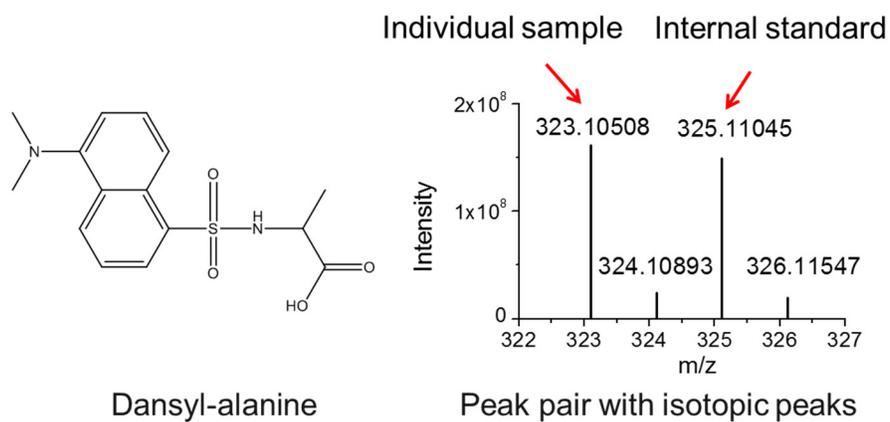
282. Chen, Q. M.; Liu, X. Y.; Zhao, P.; Sun, Y. H.; Zhao, X. J.; Xiong, Y.; Xu, G. W.; Xia, Q. Y., GC/MS-based metabolomic studies reveal key roles of glycine in regulating silk synthesis in silkworm, *Bombyx mori*. *Insect Biochemistry and Molecular Biology* **2015**, *57*, 41-50.
283. Li, Y.; Wang, X.; Hou, Y.; Zhou, X.; Chen, Q.; Guo, C.; Xia, Q.; Zhang, Y.; Zhao, P., Integrative Proteomics and Metabolomics Analysis of Insect Larva Brain: Novel Insights into the Molecular Mechanism of Insect Wandering Behavior. *Journal of Proteome Research* **2016**, *15* (1), 193-204.
284. Chikayama, E.; Suto, M.; Nishihara, T.; Shinozaki, K.; Hirayama, T.; Kikuchi, J., Systematic NMR Analysis of Stable Isotope Labeled Metabolite Mixtures in Plant and Animal Systems: Coarse Grained Views of Metabolic Pathways. *Plos One* **2008**, *3* (11), e3805.
285. de Cock, M.; van de Bor, M., Obesogenic effects of endocrine disruptors, what do we know from animal and human studies? *Environment International* **2014**, *70*, 15-24.
286. Fry, D.; Toone, C., DDT-induced feminization of gull embryos. *Science* **1981**, *213* (4510), 922-924.
287. Ayotte, P.; Giroux, S.; Dewailly, A. r.; Hernandez Avila, M.; Farias, P.; Danis, R.; Villanueva Daz, C., DDT Spraying for Malaria Control and Reproductive Function in Mexican Men. *Epidemiology* **2001**, *12* (3), 366-367.
288. Li, J. Y.; Chen, X.; Fan, W.; Moghaddam, S. H. H.; Chen, M.; Zhou, Z. H.; Yang, H. J.; Chen, J. E.; Zhong, B. X., Proteomic and Bioinformatic Analysis on Endocrine Organs of Domesticated Silkworm, *Bombyx mori* L. for a Comprehensive Understanding of Their Roles and Relations. *Journal of Proteome Research* **2009**, *8* (6), 2620-2632.
289. Bai, J. H.; Lu, Q. Q.; Zhao, Q. Q.; Wang, J. J.; Gao, Z. Q.; Zhang, G. L., Organochlorine pesticides (OCPs) in wetland soils under different land uses along a 100-year chronosequence of reclamation in a Chinese estuary. *Scientific Reports* **2015**, *5*, 10.
290. Liu, M. X.; Yang, Y. Y.; Yun, X. Y.; Zhang, M. M.; Wang, J., Occurrence and assessment of organochlorine pesticides in the agricultural topsoil of Three Gorges Dam region, China. *Environmental Earth Sciences* **2015**, *74* (6), 5001-5008.
291. Yi, Z. G.; Guo, P. P.; Zheng, L. L.; Huang, X. R.; Bi, J. Q., Distribution of HCHs and DDTs in the soil-plant system in tea gardens in Fujian, a major tea-producing province in China. *Agriculture Ecosystems & Environment* **2013**, *171*, 19-24.
292. Guo, K.; Li, L., Differential C-12/C-13-Isotope Dansylation Labeling and Fast Liquid Chromatography/Mass Spectrometry for Absolute and Relative Quantification of the Metabolome. *Analytical Chemistry* **2009**, *81* (10), 3919-3932.
293. Furuya, S.; Tabata, T.; Mitoma, J.; Yamada, K.; Yamasaki, M.; Makino, A.; Yamamoto, T.; Watanabe, M.; Kano, M.; Hirabayashi, Y., L-Serine and glycine serve as major astroglia-derived trophic factors for cerebellar Purkinje neurons. *Proceedings of the National Academy of Sciences of the United States of America* **2000**, *97* (21), 11528-11533.
294. Wagenmakers, A. J. M., Protein and Amino Acid Metabolism in Human Muscle. In *Advances in Experimental Medicine and Biology*, Springer Science + Business Media: 1998; pp 307-319.
295. Viswanathan, S.; Dignam, S. S.; Dignam, J. D., Control of the levels of alanyl-, glycy-, and seryl-tRNA synthetases in the silkgland of *Bombyx mori*. *Developmental Biology* **1988**, *129* (2), 350-357.
296. Lu, X. M., *Silkworm hormones and its control to development, in Research and Development of High-tech for Sericulture*. China Agricultural University Press, pp161-195: Beijing, China, 2012.
297. Stipanuk, M. H., Metabolism of Sulfur-Containing Amino Acids. *Annu. Rev. Nutr.* **1986**, *6* (1), 179-209.

298. Nagata, S.; Kataoka, H.; Suzuki, A. I., Silk Moth Neuropeptide Hormones: Prothoracicotropic Hormone and Others. *Annals of the New York Academy of Sciences* **2005**, *1040* (1), 38-52.
299. Ogawa, S.; Fujii, T.; Koga, N.; Hori, H.; Teraishi, T.; Hattori, K.; Noda, T.; Higuchi, T.; Motohashi, N.; Kunugi, H., Plasma L-Tryptophan Concentration in Major Depressive Disorder. *J. Clin. Psychiatry* **2014**, *75* (09), e906-e915.
300. Nishijo, M.; Tai, P. T.; Anh, N. T. N.; Nghi, T. N.; Nakagawa, H.; Van Luong, H.; Anh, T. H.; Morikawa, Y.; Waseda, T.; Kido, T.; Nishijo, H., Urinary Amino Acid Alterations in 3-Year-Old Children with Neurodevelopmental Effects due to Perinatal Dioxin Exposure in Vietnam: A Nested Case-Control Study for Neurobiomarker Discovery. *Plos One* **2015**, *10* (1), e0116778.
301. Pascucci, T.; Ventura, R.; Puglisi-Allegra, S.; Cabib, S., Deficits in brain serotonin synthesis in a genetic mouse model of phenylketonuria. *NeuroReport* **2002**, *13* (18), 2561-2564.
302. Prasad, P.; Ogawa, S.; Parhar, I. S., Role of serotonin in fish reproduction. *Front. Neurosci.* **2015**, *9*, doi: 10.3389/fnins.2015.00195.
303. Valim, V. r.; Natour, J.; Xiao, Y.; Pereira, A. o. F. A.; Lopes, B. B. d. C.; Pollak, D. F.; Zandonade, E.; Russell, I. J., Effects of physical exercise on serum levels of serotonin and its metabolite in fibromyalgia: a randomized pilot study. *Revista Brasileira de Reumatologia* **2013**, *53* (6), 538-541.
304. Siangcham, T.; Tinikul, Y.; Poljaroen, J.; Sroyraya, M.; Changklungmoa, N.; Phoungpetchara, I.; Kankuan, W.; Sumpownon, C.; Wanichanon, C.; Hanna, P. J.; Sobhon, P., The effects of serotonin, dopamine, gonadotropin-releasing hormones, and corazonin, on the androgenic gland of the giant freshwater prawn, *Macrobrachium rosenbergii*. *General and Comparative Endocrinology* **2013**, *193*, 10-18.
305. Mommersteeg, P. M. C.; Schoemaker, R. G.; Eisel, U. L. M.; Garrelds, I. M.; Schalkwijk, C. G.; Kop, W. J., Nitric Oxide Dysregulation in Patients With Heart Failure. *Psychosomatic Medicine* **2015**, *77* (3), 292-302.
306. Lakhani, K.; Kay, A. R.; Leiper, J.; Barry, J. A.; Hardiman, P. J., Symmetric dimethylarginine (SDMA) is raised in women with polycystic ovary syndrome: A pilot study. *Journal of Obstetrics and Gynaecology* **2011**, *31* (5), 417-419.
307. Arikan, E.; Karadag, C. H.; Guldiken, S., Asymmetric dimethylarginine levels in thyroid diseases. *Journal of Endocrinological Investigation* **2007**, *30* (3), 186-191.
308. Verhoeven, M. O.; Hemelaar, M.; Teerlink, T.; Kenemans, P.; van der Mooren, M. J., Effects of intranasal versus oral hormone therapy on asymmetric dimethylarginine in healthy postmenopausal women: A randomized study. *Atherosclerosis* **2007**, *195* (1), 181-188.
309. Yang, H.-Y.; Liu, Y.; Xie, J.-C.; Liu, N.-N.; Tian, X., Effects of repetitive transcranial magnetic stimulation on synaptic plasticity and apoptosis in vascular dementia rats. *Behavioural Brain Research* **2015**, *281*, 149-155.
310. Cline, H. T.; Constantine-Paton, M., NMDA receptor antagonists disrupt the retinotectal topographic map. *Neuron* **1989**, *3* (4), 413-426.
311. Matsumoto, S.; Isogai, A.; Suzuki, A., N-terminal amino acid sequence of an insect neurohormone, melanization and reddish coloration hormone (MRCH): heterogeneity and sequence homology with human insulin-like growth factor II. *FEBS Letters* **1985**, *189* (1), 115-118.
312. Matsumoto, S.; Isogai, A.; Suzuki, A.; Ogura, N.; Sonobe, H., Purification and properties of the melanization and reddish colouration hormone (MRCH) in the armyworm, *Leucania separata* (Lepidoptera). *Insect Biochemistry* **1981**, *11* (6), 725-733.

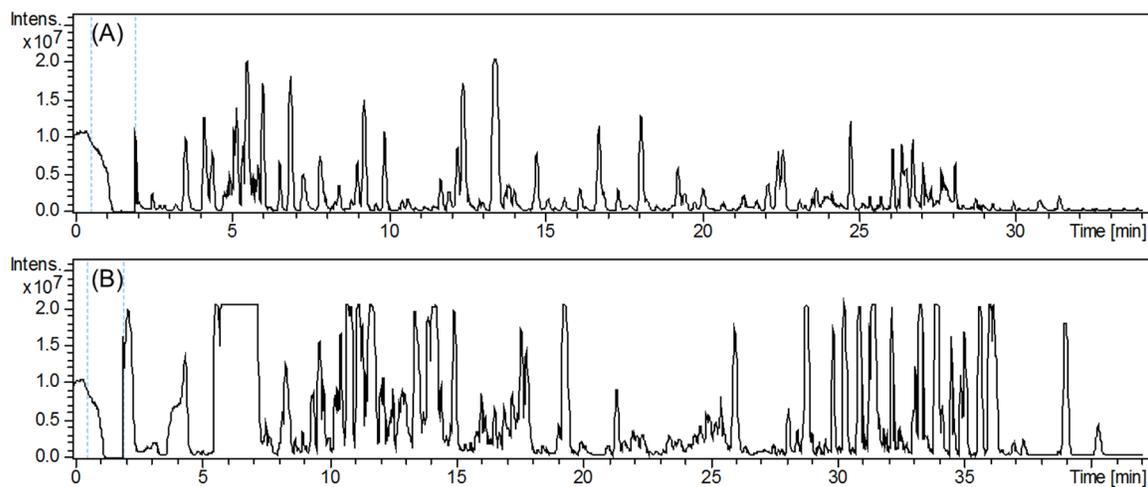
## Appendix



**Appendix Figure 1.** Representative LC-MS chromatograms from (A) a dansyl-labeled human serum sample, (B) a DMPA-labeled human serum sample, (C) a dansyl-labeled human plasma sample, and (D) a DmPA-labeled human plasma sample (acquired from LC-FT-ICR-MS, with an ESI source at positive mode).



**Appendix Figure 2.** Mass spectrum of a representative peak pair (dansyl-labeled alanine).



**Appendix Figure 3.** (A) A representative base peak chromatogram of dansyl-labeled 0.5  $\mu$ L finger blood, and (B) a representative base peak chromatogram of DMPA-labeled 0.5  $\mu$ L finger blood. (acquired by LC-Q-TOF-MS)