**Supplementary Materials for:**

# Simultaneous Improvement in the *Precision*, *Accuracy* and *Robustness* of Label-free Proteome Quantification by Optimizing Data Manipulation Chains

Jing TANG, Jianbo FU, Yunxia WANG, Yongchao LUO, Qingxia YANG, Bo LI, Gao TU, Jiajun HONG, Xuejiao CUI, Yuzong CHEN, Lixia YAO, Weiwei XUE and Feng ZHU*

[§] College of Pharmaceutical Sciences, Zhejiang University, Hangzhou 310058, China

[ξ] School of Pharmaceutical Sciences, Chongqing University, Chongqing 401331, China

[€] Department of Bioinformatics, Chongqing Medical University, Chongqing 400016, China

[ʃ] Department of Pharmacy, National University of Singapore, Singapore 117543, Singapore

[π] Department of Health Sciences Research, Mayo Clinic, Rochester MN 55905, United States

*To whom the correspondence should be addressed: Prof. Feng ZHU, College of Pharmaceutical Sciences, Zhejiang University, Hangzhou 310058, China. E-mail: zhufeng@zju.edu.cn; prof.zhufeng@gmail.com

**Running Title**: Simultaneous Improvement of LFQ Precision, Accuracy and Robustness

**Table S1.** A variety of popular quantification tools and subsequent manipulation methods for analyzing MS-based proteomic data together with the representative proteomics studies adopting each tool/method.

| Name of Tool / Method | Extensive Application of Each Tool / Method in Current Proteomic Researches |
|---|---|
| **(1) *Quantification Tools for Each Acquisition Technique*** | |

| | Name of Tool / Method | Extensive Application of Each Tool / Method in Current Proteomic Researches |
|---|---|---|
| **Peak Intensity** | MaxQuant | Applied to investigate the properties and function of the inner colonic mucus layer in mouse models of diet-induced and genetic obesity[1]. |
| | MFPaQ | Applied to evaluate a quantitative proteomic workflow for impoving reproduciblity and accuracy quantification[2]. |
| | OpenMS | Applied to deciphering physiological changes that mediate transition of Mycobacterium tuberculosis between replicating and nonreplicating states[3]. |
| | PEAKS | Employs de novo sequencing as a subroutine and exploits the de novo sequencing results to improve both the speed and accuracy of the database search[4]. |
| | Progenesis | Used to gain insight into mechanisms underpinning corticosteroid effects on neural stem cells[5]. |
| | Proteios SE | Applied to identify the variant protein in intestinal epithelial cells from healthy subjects (H) and Crohn's disease patients (CD)[6]. |
| | Scaffold | Used to investage interactions between NS4B involved in viral replication and immune evasion and human proteins[7]. |
| | Thermo Proteome Discoverer | Applied to analyze and compare the total proteome of aqueous humor) from patients with primary angle closure glaucoma, open angle glaucoma and age-related cataract[8]. |
| **Spectral Counting** | Abacus | Applied to extract and pre-process spectral count data for the label-free quantitative proteomic analysis[9]. |
| | Census | Enabling the large-scale differential proteome analysis in *Plasmodium falciparum* under drug treatment[10]. |
| | DTASelect | Used to facilitate to develop diagnostic biomarkers for subclinical IAI in amniotic fluid and blood of women with preterm labor[11]. |
| | IRMa-hEIDI | Applied to facilitate to obtain molecular and cellular cerebral imprints in the striatum of anesthetized monkeys[12]. |
| | MaxQuant | Applied to investigate the properties and function of the inner colonic mucus layer in mouse models of diet-induced and genetic obesity[1]. |
| | MFPaQ | Applied to evaluate a quantitative proteomic workflow for impoving reproduciblity and accuracy quantification[2]. |
| | ProteinProphet | Applied to identify secreted glycoproteins of human prostate and bladder stromal cells by comparative quantitative proteomics[13]. |
| **SWATH-MS** | DIA-Umpire | Applied to quantity and identify Host Cell Proteins (HCPs) of an IgG1 monoclonal antibody (mAb) sample[14]. |
| | OpenSWATH | Enabling to identify protein characterization of lung extracellular matrix for describing the specific matrisome remodeling mechanisms[15]. |
| | PeakView | Used to evaluate and identify optimal protein extraction method for proteomics analysis of green algae Chlorella vulgaris[16]. |
| | Skyline | Applied to explore the total proteome and glycoproteins of synovial fluid obtained from osteoarthritis patients[17]. |
| | Spectronaut | Applied to discovery the high throuput and accurate serum proteome profiling workflow[18]. |

| | **(2) Transformation Methods** | | |
|---|---|---|---|
| *Transformation* | 1 | **BOX** | Used to identify of novel biomarkers and the development of new therapeutic targets for seven important liver diseases[19]. |
| | 2 | **LOG** | Applied for identifying the new therapeutic targets of the treatment of early-stage hepatocellular carcinoma (HCC)[20]. |
| | 3 | **VSN** | Helping to address the accuracy and precision issues in the isobaric tags for relative and absolute quantification (iTRAQ)[21]. |

| | **(3) Pretreatment Methods** | | |
|---|---|---|---|
| *Centering* | 1 | **MEC** | Used for facilitating the improvement of the sensitivity of significance test in spectral counting-based comparative discovery proteomics[22]. |
| | 2 | **MDC** | Facilitating the normalization procedures in LC-MS proteomics experiments through dataset dependent ranking of normalization scaling factors[23]. |
| *Scaling* | 3 | **ATO** | Applied to discover the proteomic biomarkers for a variety of diseases, such as psoriasis and psoriasis arthritis[24]. |
| | 4 | **PAR** | Implemented into the proteomic experiments based on the LC-MS/MS with great potential to be applied to metaproteomic research[25]. |
| | 5 | **VAS** | Assessing the impact of delayed storage on the measured proteome and metabolome of human cerebrospinal fluid[26]. |
| | 6 | **RAN** | Manipulating the non-targeted ultra-high performance liquid chromatography tandem mass spectrometry (UHPLC-MS) proteomic/metabolomic data[27]. |
| *Normalization* | 7 | **CYC** | A frequently adopted normalization method in the quantitative label-free proteomics, systematically compared with other methods[28]. |
| | 8 | **EIG** | Enabling the normalization of peak intensities in bottom-up MS-based proteomics and label-free LC-MS based proteomics analysis[29]. |
| | 9 | **LIN** | Developed to normalize or scale the label-free relative quantification of the endogenous peptides[30]. |
| | 10 | **LOW** | Designed to normalize and statistical analyze the quantitative proteomics data generated by metabolic labeling[31]. |
| | 11 | **MEA** | Applied to the analysis by high throughput gel free quantitative proteomics and metaproteomic-related research[32]. |
| | 12 | **MED** | Used to achieve the reproducible and consistent quantification of the Saccharomyces cerevisiae proteome by SWATH-mass spectrometry[32]. |
| | 13 | **MAD** | Normalizing and improving the quality control procedure of the peptide-centric LC-MS proteomics data[33]. |
| | 14 | **PQN** | Normlizing and analyzing the protein profiles of antibody arrays based on a longitudinal twin model[34]. |
| | 15 | **QUA** | Facilitating the rapid mass spectrometric conversion of tissue biopsy samples into the permanent quantitative digital proteome maps[35]. |
| | 16 | **RLR** | Enabling the multidimensional normalization to minimize the plate effects of the suspension bead array data[36]. |
| | 17 | **TIC** | Applied to achieve SELDI-TOF-MS proteomic profiling of serum, urine, and amniotic fluid in neural tube defects[37]. |
| | 18 | **TMM** | Frequently used as a multi-model statistical approach for the proteomic spectral count quantitation[38]. |

| | | | |
|---|---|---|---|
| **(4)** *Methods for Missing-value Imputation* | | | |
| *Imputation* | 1 | **BAK** | Treating the missing values for multivariate statistical analysis of the gel-based proteomics data[39]. |
| | 2 | **BPC** | Used to improve the detection of differentially abundant proteins in current proteomic analysis[40]. |
| | 3 | **CEN** | An integrative imputation method based on multi-omics datasets, especially in proteomic analysis[41]. |
| | 4 | **KNN** | A frequently used imputation method in the quantitative label-free proteomics, systematically compared with other methods[42]. |
| | 5 | **LLS** | Facilitating the normalization of peak intensities in bottom-up MS-based proteomic analysis[29]. |
| | 6 | **SVD** | Helping the realizing the visualization, manipulating and quantitation of the isobaric tagged mass spectrometry based proteomic data[43]. |
| | 7 | **ZER** | Treating the missing values for multivariate statistical analysis of the gel-based proteomics data[39]. |

**Table S2.** The level of dependence of *precision* and *accuracy* on the selection of LFQs assessed based on the benchmark datasets of the second study in **Table 1**[12]. The studied LFQs were collectively defined by 4 quantification tools and 8 representative manipulation chains. As reported in original study[12], only the log *transformation* was adopted for analyses, which defined its manipulation chain as *LOG-[NON-NON-NON]-NON* (highlighted in grey background). The *precision* levels were defined by PMAD values (*superior*: <0.14, *good*: 0.14~0.3, *fair*: 0.3~0.7 and *poor*: >0.7), and the *accuracy* levels were defined by the absolute deviation from the expected abundance ratios (*high*: <0.05, *medium*: 0.05~0.1 and *low*: >0.1). Each manipulation method in a chain was abbreviated by a three-letter code which was systematically defined in **Table 1**.

| Representative Chains for Data Manipulation | Precision: reproducibility among technical replicates (PMAD) | | | | Accuracy: absolute deviation from expected abundance ratios | | | |
|---|---|---|---|---|---|---|---|---|
| | IRMa-hEIDI | MFPaQ | MaxQuant | Scaffold | IRMa-hEIDI | MFPaQ | MaxQuant | Scaffold |
| LOG-[NON-NON-NON]-NON | 2.02E-01 (Good) | 2.50E-01 (Good) | 2.56E-01 (Good) | 2.09E-01 (Good) | 1.42E-04 (High) | 5.07E-03 (High) | 2.80E-02 (High) | 5.15E-03 (High) |
| LOG-[MDC-PAR-MAD]-SVD | 6.81E-01 (Fair) | 7.27E-01 (Poor) | 7.56E-01 (Poor) | 6.47E-01 (Fair) | 1.34E-01 (Low) | 5.00E-02 (High) | 1.56E-01 (Low) | 1.06E-01 (Low) |
| LOG-[MDC-RAN-EIG]-KNN | 1.03E-17 (Superior) | 8.75E-18 (Superior) | 1.25E-17 (Superior) | 2.45E-17 (Superior) | 3.83E-02 (High) | 1.80E-02 (High) | 4.50E-02 (High) | 3.82E-02 (High) |
| BOX-[NON-VAS-EIG]-BAK | 1.14E+00 (Poor) | 1.83E+00 (Poor) | 2.82E+00 (Poor) | 8.73E-01 (Poor) | 4.16E-03 (High) | 3.58E-03 (High) | 7.24E-03 (High) | 4.15E-03 (High) |
| BOX-[MDC-PAR-TMM]-SVD | 2.63E-01 (Good) | 3.67E-01 (Fair) | 1.68E-01 (Good) | 7.99E-02 (Superior) | 1.20E-02 (High) | 1.74E-01 (Low) | 1.84E-01 (Low) | 7.04E-02 (Medium) |
| BOX-[MEC-PAR-TIC]-KNN | 3.60E-02 (Superior) | 2.33E-02 (Superior) | 1.46E-02 (Superior) | 3.04E-02 (Superior) | 4.99E-01 (Low) | 6.57E-01 (Low) | 8.77E-01 (Low) | 6.91E-01 (Low) |
| BOX-[MEC-RAN-MAD]-ZER | 3.97E-01 (Fair) | 4.78E-01 (Fair) | 5.03E-01 (Fair) | 4.07E-01 (Fair) | 2.18E-01 (Low) | 1.07E-02 (High) | 6.25E-02 (Medium) | 7.90E-04 (High) |
| BOX-[MEC-ATO-RLR]-BPC | 3.91E+00 (Poor) | 2.35E+00 (Poor) | 2.31E+00 (Poor) | 2.09E+00 (Poor) | 4.34E-01 (Low) | 1.37E+00 (Low) | 7.31E-01 (Low) | 6.79E-01 (Low) |

**Table S3.** The PMAD values of ten representative manipulation chains for *peak intensity* data that performed consistently better across five quantification tools than the manipulation chain adopted by the original study (*LOG-[NON-NON-MED]-NON*, highlighted in dark grey background).

| Manipulation Chains | | DecyderMS | MaxQuant | PEAKS | OpenMS | Sieve |
|---|---|---|---|---|---|---|
| | LOG-[NON-NON-MED]-NON | 5.56E-01 | 5.41E-01 | 6.43E-01 | 4.80E-01 | 4.71E-01 |
| 1 | BOX-[MEC-RAN-PQN]-KNN | 1.29E-01 | 1.24E-01 | 1.62E-01 | 1.47E-01 | 1.50E-01 |
| 2 | BOX-[MDC-RAN-QUA]-BAK | 1.90E-01 | 1.94E-01 | 2.09E-01 | 2.06E-01 | 1.38E-01 |
| 3 | BOX-[MEC-ATO-LOW]-SVD | 5.49E-02 | 6.22E-02 | 7.51E-02 | 5.76E-02 | 6.25E-02 |
| 4 | BOX-[MDC-RAN-TMM]-CEN | 1.48E-02 | 1.17E-02 | 3.02E-02 | 5.56E-03 | 8.15E-03 |
| 5 | BOX-[MDC-VAS-TIC]-BPC | 5.38E-04 | 3.61E-04 | 2.29E-04 | 2.74E-04 | 6.27E-04 |
| 6 | LOG-[MDC-VAS-TIC]-CEN | 8.39E-03 | 1.68E-03 | 2.82E-03 | 1.10E-02 | 6.79E-04 |
| 7 | LOG-[MDC-PAR-LOW]-SVD | 4.06E-02 | 2.39E-02 | 3.87E-02 | 2.48E-02 | 3.06E-02 |
| 8 | LOG-[MEC-ATO-LOW]-ZER | 4.23E-02 | 4.17E-02 | 6.01E-02 | 3.92E-02 | 5.22E-02 |
| 9 | LOG-[MEC-RAN-PQN]-BAK | 1.64E-01 | 1.94E-01 | 2.64E-01 | 4.63E-01 | 5.76E-02 |
| 10 | LOG-[MDC-PAR-TMM]-CEN | 4.87E-02 | 4.16E-02 | 5.56E-02 | 2.34E-02 | 6.39E-02 |

**Table S4.** PMAD values of ten representative manipulation chains for *spectral counting* data that performed consistently better across four quantification tools than the manipulation chain adopted by the original study (*LOG-[NON-NON-NON]-NON*, highlighted in dark grey background).
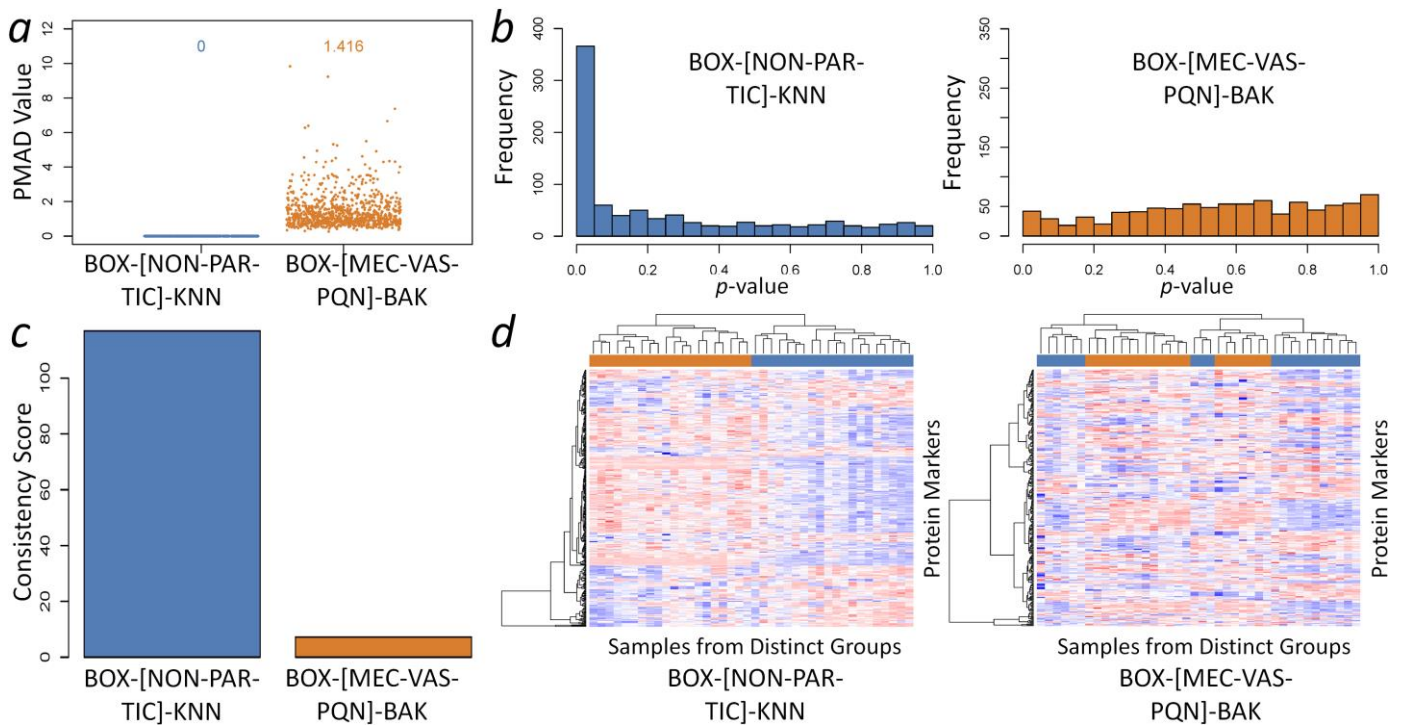
| Manipulation Chains | | IRMa-hEIDI | MFPaQ | MaxQuant | Scaffold |
|---|---|---|---|---|---|
| | LOG-[NON-NON-NON]-NON | 2.02E-01 | 2.50E-01 | 2.56E-01 | 2.09E-01 |
| 1 | BOX-[MEC-RAN-PQN]-CEN | 1.37E-01 | 1.64E-01 | 1.90E-01 | 1.37E-01 |
| 2 | BOX-[MEC-ATO-LOW]-KNN | 1.74E-02 | 2.75E-02 | 5.10E-02 | 1.82E-02 |
| 3 | BOX-[MEC-PAR-TIC]-BPC | 3.60E-02 | 2.33E-02 | 1.46E-02 | 3.04E-02 |
| 4 | BOX-[MDC-ATO-LOW]-KNN | 4.92E-02 | 6.03E-02 | 7.04E-02 | 4.62E-02 |
| 5 | BOX-[MDC-RAN-LOW]-SVD | 5.02E-03 | 2.04E-02 | 4.60E-03 | 4.24E-03 |
| 6 | LOG-[MEC-ATO-LOW]-KNN | 1.81E-02 | 2.68E-02 | 5.80E-02 | 1.45E-02 |
| 7 | LOG-[MEC-RAN-TIC]-CEN | 4.89E-02 | 1.10E-01 | 3.49E-02 | 5.32E-02 |
| 8 | LOG-[MDC-PAR-TIC]-BAK | 5.78E-02 | 2.84E-02 | 3.68E-02 | 6.29E-02 |
| 9 | LOG-[MDC-VAS-TIC]-SVD | 8.33E-03 | 8.05E-03 | 1.67E-02 | 1.34E-02 |
| 10 | LOG-[MDC-RAN-TMM]-ZER | 3.43E-02 | 5.25E-02 | 1.22E-01 | 3.44E-02 |

**Table S5**. Comparison among LFQ-related tools in terms of the acquisition technique, quantification tools, subsequent manipulation methods (*transformation*, *pretreatment* and *imputation*) and performance evaluation criteria. *Gmine* & *Perseus* integrated several manipulation methods in the quantification workflows, but no function of performance assessment was provided. *LFQbench* & *msCompare* were recognized as evaluating performances of 3~5 quantification tools, and *Normalyzer*, *SPANS* & *GiaPronto* were distinguished for being capable of assessing 1~8 *pretreatment* methods. Tools were ordered alphabetically. N.P.: not provided; PI: *peak intensity*; SC: *spectral counting*; SWATH: *SWATH-MS*.

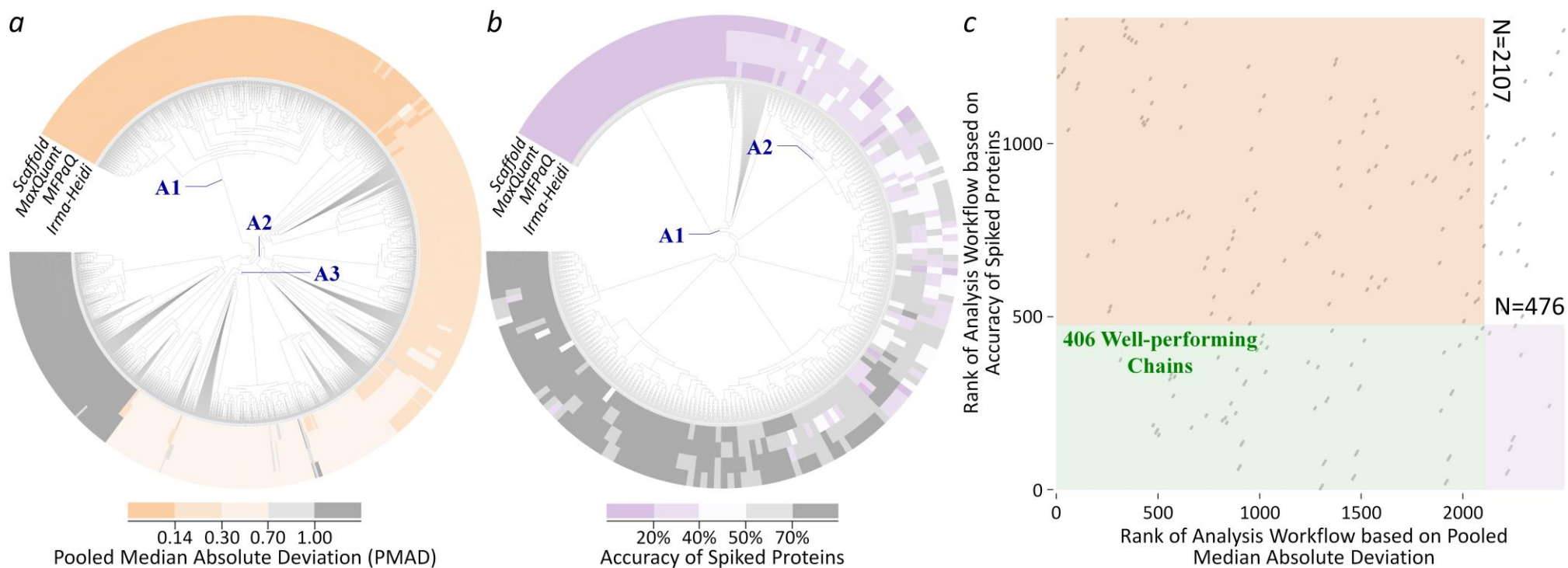| | Automatic Detection of the Output Format of Quantification Tools (Quantification Measurement) | No. of Quantification Tools | No. of Transformation Methods | No. of Pretreatment Methods | No. of Imputation Methods | Performance Evaluation | No. of Evaluation Criteria |
|---|---|---|---|---|---|---|---|
| **This Study** | **YES** (PI, SC, SWATH) | **18** | **4** | **18** | **7** | **YES** | **5** |
| GiaPronto | **YES** (PI, SC) | 1 | 1 | 1 | N.P. | **YES** | 1 |
| Gmine | **NO** | N.P. | 3 | 2 | N.P. | **NO** | N.P. |
| LFQbench | **NO** | 5 | N.P. | N.P. | N.P. | **YES** | 2 |
| msCompare | **NO** | 3 | N.P. | N.P. | N.P. | **YES** | 1 |
| Normalyzer | **NO** | N.P. | 1 | 8 | N.P. | **YES** | 1 |
| Perseus | **YES** (PI, SC) | 1 | 1 | 7 | 3 | **NO** | N.P. |
| SPANS | **NO** | N.P. | N.P. | 5 | N.P. | **YES** | 1 |

**Figure S1.** Two representative manipulation chains performing well (*BOX-[NON-PAR-TIC]-KNN*) and poor (*BOX-[MEC-VAS-PQN]-BAK*) consistently across multiple criteria: (***a***) *precision*; (***b***) *differential abundance analysis*; (***c***) *robustness* and (***d***) *classification capacity*.
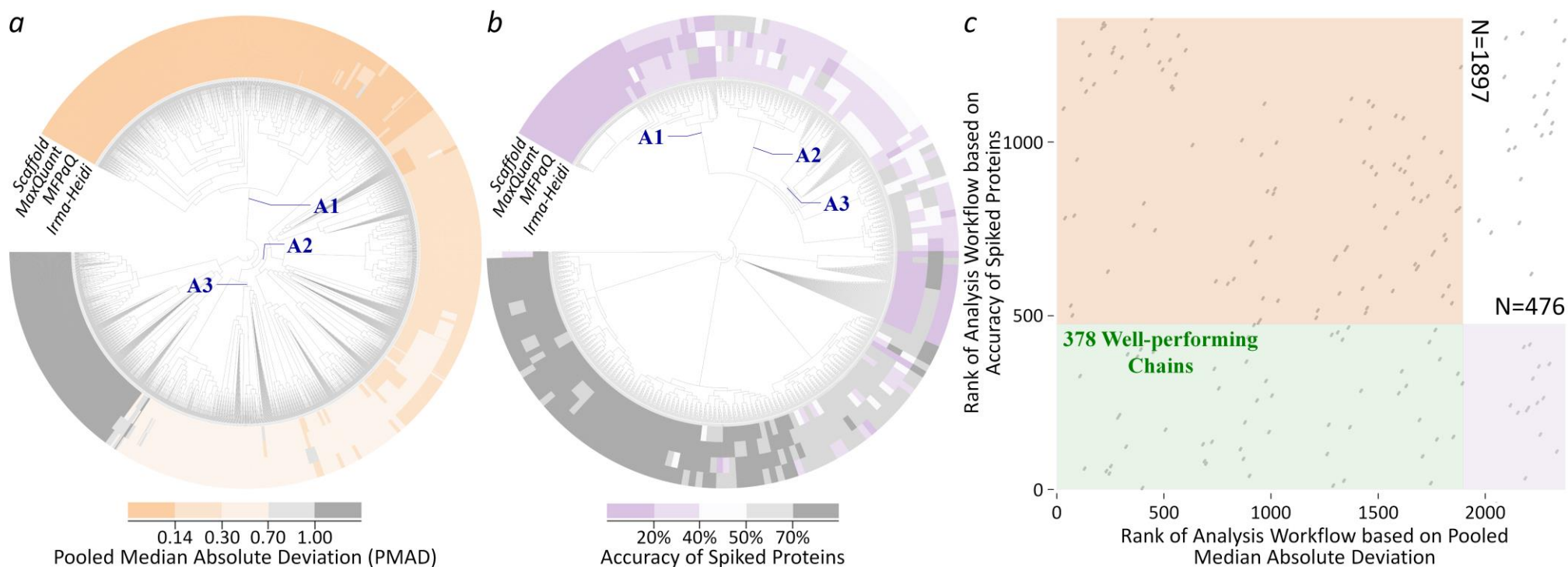
**Figure S2.** The strategy proposed in this study to discover manipulation chains with simultaneously improved *precision* and *accuracy* based on the benchmark from **Table 2** *study 2* of distinct concentrations (12.5 *vs* 25fmol/μg) of spiked UPS1 proteins. First, the clustering analyses among manipulation chains across four quantification tools were conducted for (***a***) *precision* and (***b***) *accuracy*. Second, 2D scatter plot (***c***) was drawn to provide the ranks of manipulation chains (represented by gray dots) collectively determined by *precision* (horizontal axis) and *accuracy* (vertical axis). The pink, violet and green areas in (***c***) indicated the chains of good *precision* (**A1**+**A2**+**A3** in (***a***)), good *accuracy* (**A1**+**A2** in (***b***)) and good-performance for both *precision* and *accuracy*, respectively. As a result, 728 chains (within the green region of (***c***)) were found to perform well under both criteria (*precision & accuracy*).
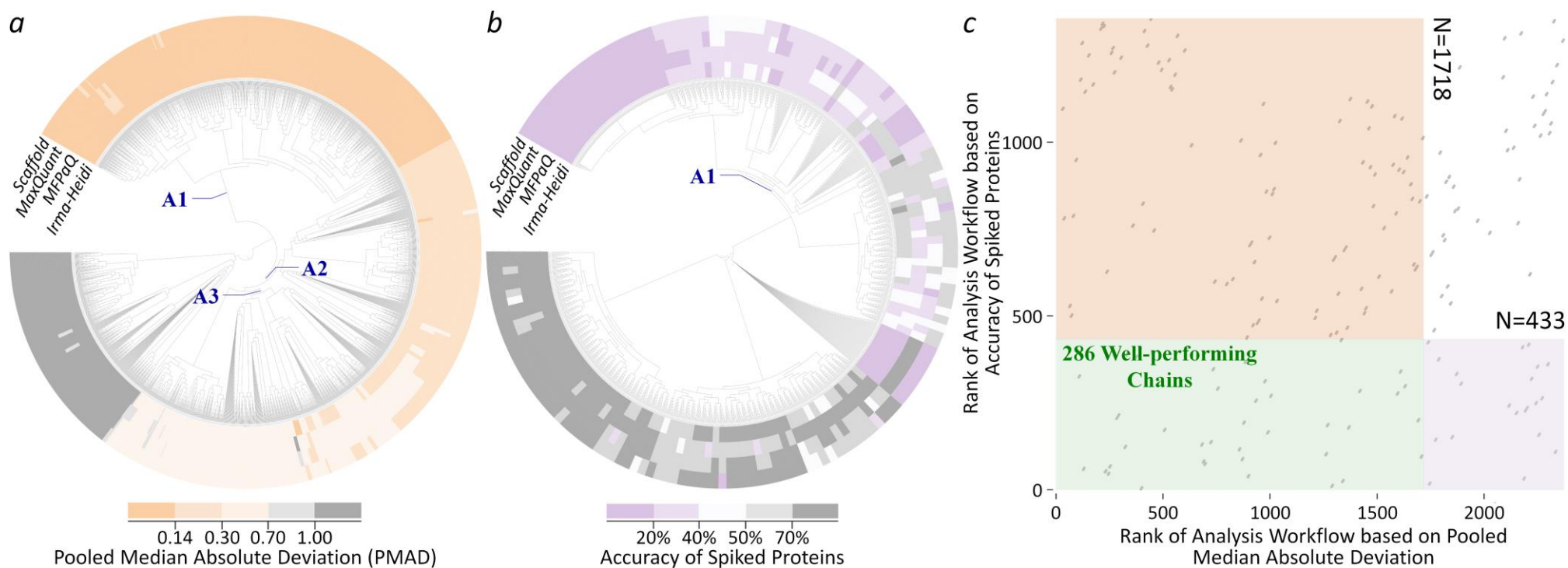
**Figure S3.** The strategy proposed in this study to discover manipulation chains with simultaneously improved *precision* and *accuracy* based on the benchmark from **Table 2** *study 2* of distinct concentrations (0.5 *vs* 5fmol/µg) of the spiked UPS1 proteins. First, the clustering analyses among manipulation chains across four quantification tools were conducted for (***a***) *precision* and (***b***) *accuracy*. Second, 2D scatter plot (***c***) was drawn to provide the ranks of manipulation chains (represented by gray dots) collectively determined by *precision* (horizontal axis) and *accuracy* (vertical axis). The pink, violet and green areas in (***c***) indicated the chains of good *precision* (**A1**+**A2**+**A3** in (***a***)), good *accuracy* (**A1**+**A2** in (***b***)) and good-performance for both *precision* and *accuracy*, respectively. As a result, 406 chains (within the green region of (***c***)) were found to perform well under both criteria (*precision* & *accuracy*).
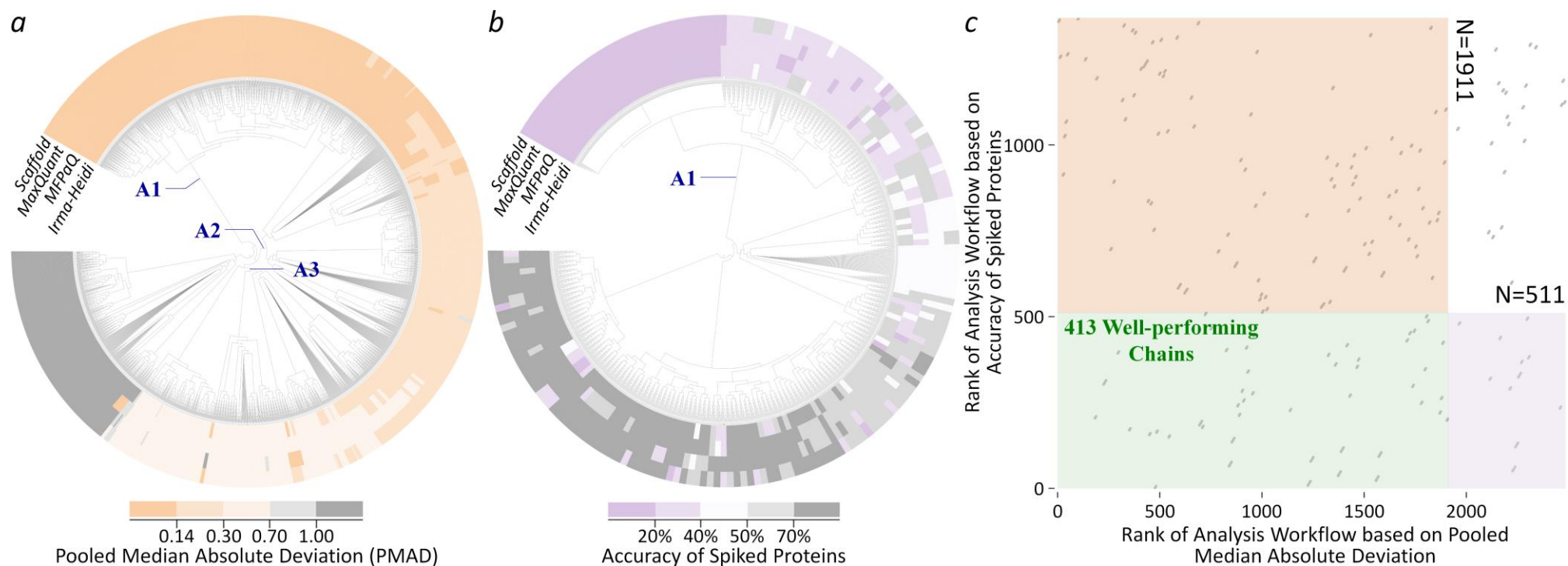
**Figure S4.** The strategy proposed in this study to discover manipulation chains with simultaneously improved *precision* and *accuracy* based on the benchmark from **Table 2** *study 2* of distinct concentrations (0.5 *vs* 12.5fmol/μg) of spiked UPS1 proteins. First, the clustering analyses among manipulation chains across four quantification tools were conducted for (***a***) *precision* and (***b***) *accuracy*. Second, 2D scatter plot (***c***) was drawn to provide the ranks of manipulation chains (represented by gray dots) collectively determined by *precision* (horizontal axis) and *accuracy* (vertical axis). The pink, violet and green areas in (***c***) indicated the chains of good *precision* (**A1**+**A2**+**A3** in (***a***)), good *accuracy* (**A1**+**A2**+**A3** in (***b***)) and good-performance for both *precision* & *accuracy*, respectively. As a result, 378 chains (within the green region of (***c***)) were found to perform well under both criteria (*precision* & *accuracy*).
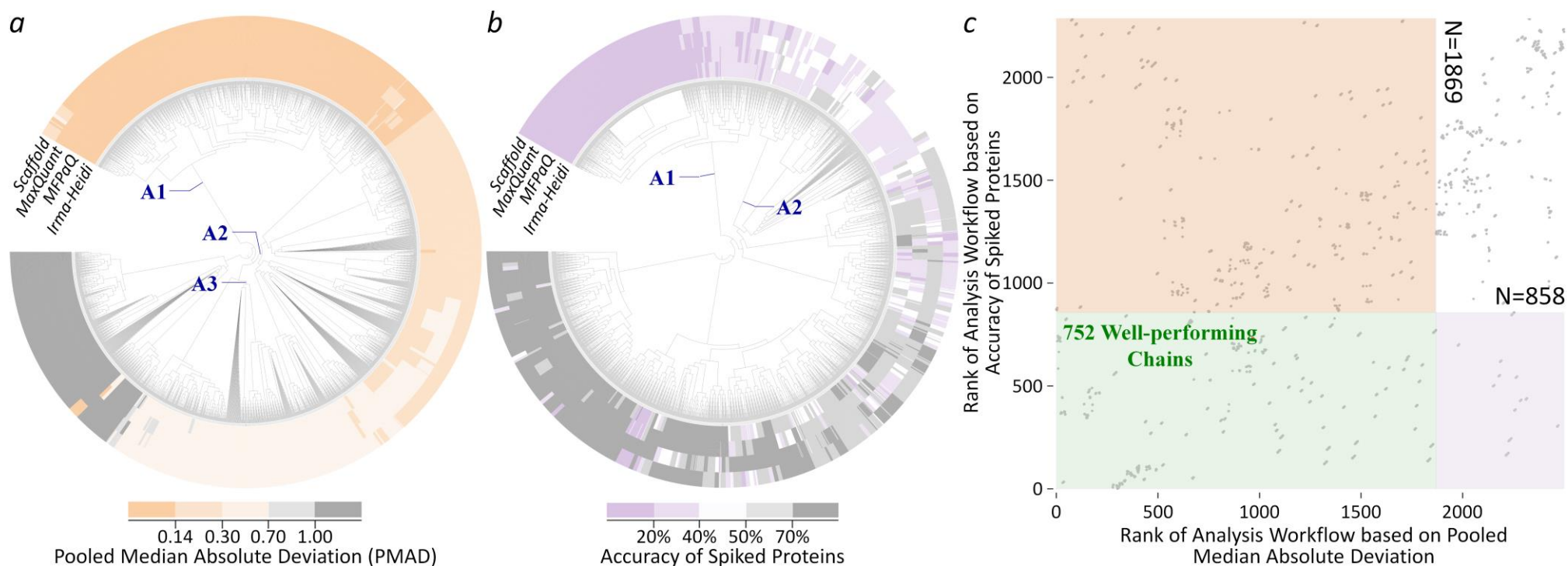
**Figure S5.** The strategy proposed in this study to discover manipulation chains with simultaneously improved *precision* and *accuracy* based on the benchmark from **Table 2** *study 2* of distinct concentrations (0.5 *vs* 25fmol/μg) of the spiked UPS1 protein. First, the clustering analyses among manipulation chains across four quantification tools were conducted for (***a***) *precision* and (***b***) *accuracy*. Second, 2D scatter plot (***c***) was drawn to provide the ranks of manipulation chains (represented by gray dots) collectively determined by *precision* (horizontal axis) and *accuracy* (vertical axis). The pink, violet and green areas in (***c***) indicated the chains of good *precision* (**A1**+**A2**+**A3** in (***a***)), good *accuracy* (**A1** in (***b***)) and good-performance for both *precision* and *accuracy*, respectively. All in all, 286 chains (within the green region of (***c***)) were found to perform well under both criteria (*precision & accuracy*).
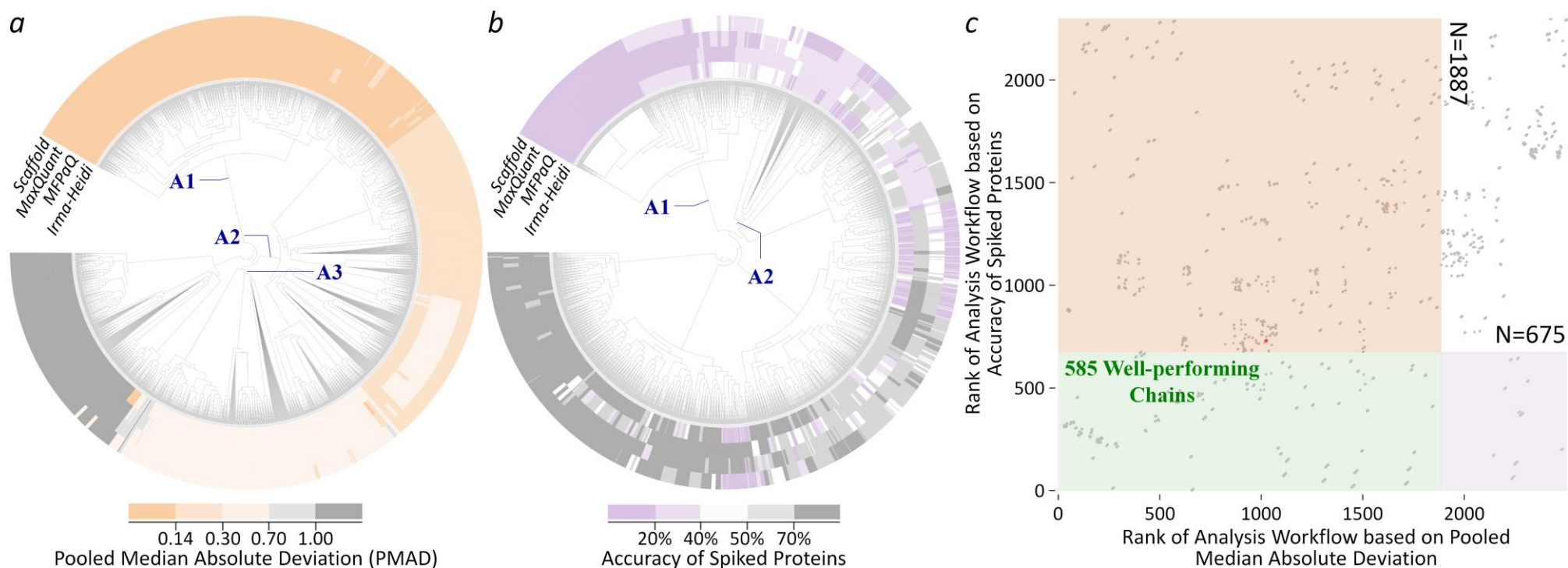
**Figure S6.** The strategy proposed in this study to discover manipulation chains with simultaneously improved *precision* and *accuracy* based on the benchmark from **Table 2** *study 2* of distinct concentrations (0.5 *vs* 50fmol/µg) of the spiked UPS1 protein. First, the clustering analyses among manipulation chains across four quantification tools were conducted for (***a***) *precision* and (***b***) *accuracy*. Second, 2D scatter plot (***c***) was drawn to provide the ranks of manipulation chains (represented by gray dots) collectively determined by *precision* (horizontal axis) and *accuracy* (vertical axis). The pink, violet and green areas in (***c***) indicated the chains of good *precision* (**A1**+**A2**+**A3** in (***a***)), good *accuracy* (**A1** in (***b***)) and good-performance for both *precision* and *accuracy*, respectively. All in all, 413 chains (within the green region of (***c***)) were found to perform well under both criteria (*precision* & *accuracy*).

**Figure S7.** The strategy proposed in this study to discover manipulation chains with simultaneously improved *precision* and *accuracy* based on the benchmark from **Table 2** *study 2* of distinct concentrations (5 *vs* 12.5fmol/μg) of the spiked UPS1 protein. First, the clustering analyses among manipulation chains across four quantification tools were conducted for (***a***) *precision* and (***b***) *accuracy*. Second, 2D scatter plot (***c***) was drawn to provide the ranks of manipulation chains (represented by gray dots) collectively determined by *precision* (horizontal axis) and *accuracy* (vertical axis). The pink, violet and green areas in (***c***) indicated the chains of good *precision* (**A1**+**A2**+**A3** in (***a***)), good *accuracy* (**A1**+**A2** in (***b***)) and good-performance for both *precision* and *accuracy*, respectively. As a result, 752 chains (within the green region of (***c***)) were found to perform well under both criteria (*precision* & *accuracy*).
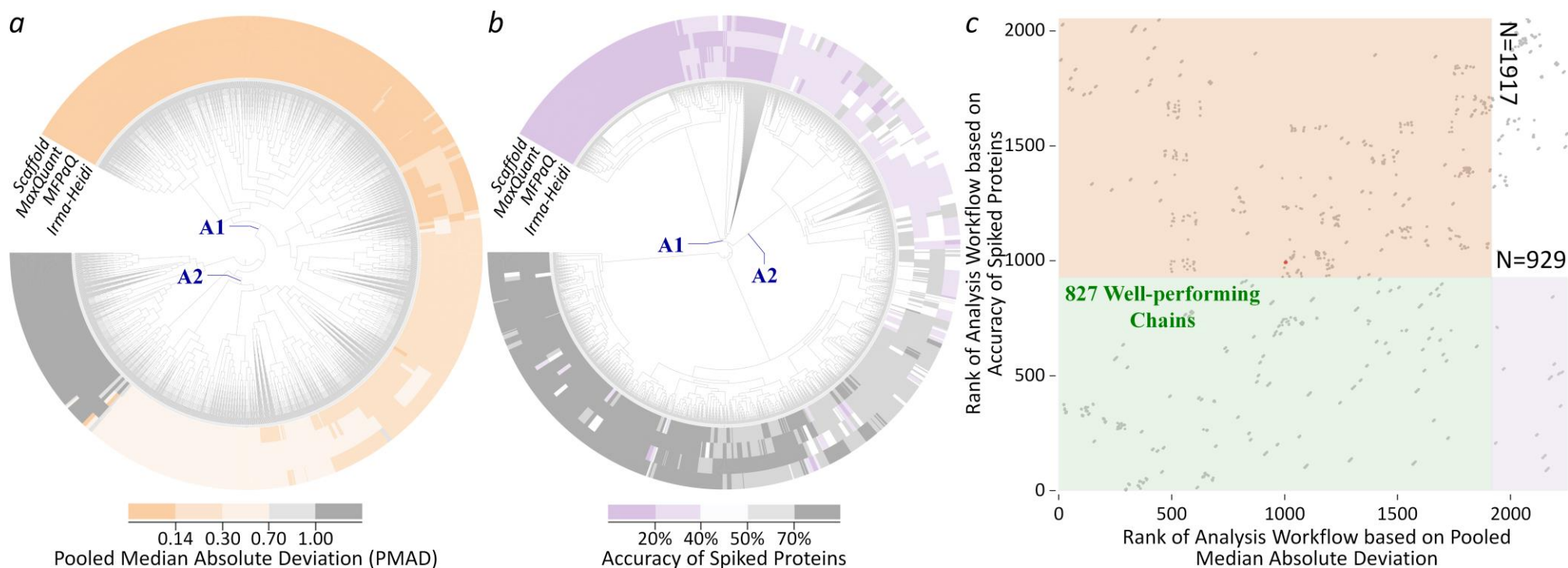
**Figure S8.** The strategy proposed in this study to discover manipulation chains with simultaneously improved *precision* and *accuracy* based on the benchmark from **Table 2** *study 2* of distinct concentrations (5 *vs* 25fmol/μg) of the spiked UPS1 proteins. First, the clustering analyses among manipulation chains across four quantification tools were conducted for (***a***) *precision* and (***b***) *accuracy*. Second, 2D scatter plot (***c***) was drawn to provide the ranks of manipulation chains (represented by gray dots) collectively determined by *precision* (horizontal axis) and *accuracy* (vertical axis). The pink, violet and green areas in (***c***) indicated the chains of good *precision* (**A1**+**A2**+**A3** in (***a***)), good *accuracy* (**A1**+**A2** in (***b***)) and good-performance for both *precision* and *accuracy*, respectively. As a result, 585 chains (within the green region of (***c***)) were found to perform well under both criteria (*precision & accuracy*).
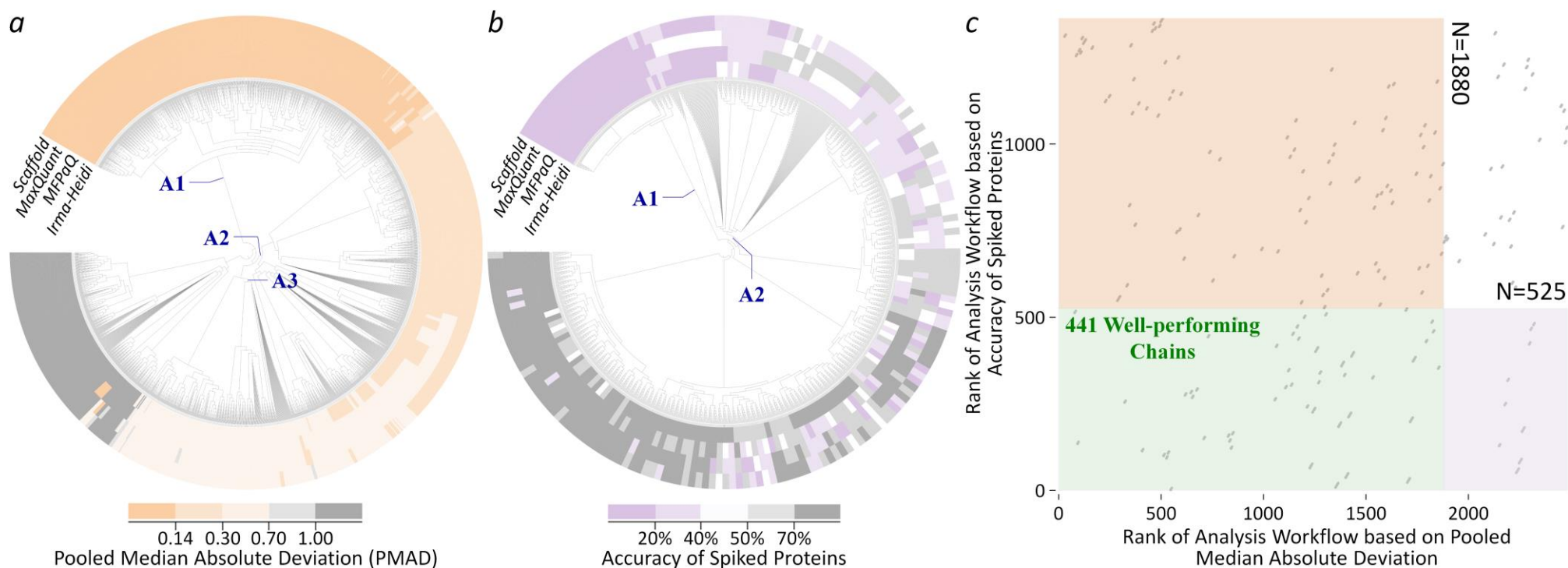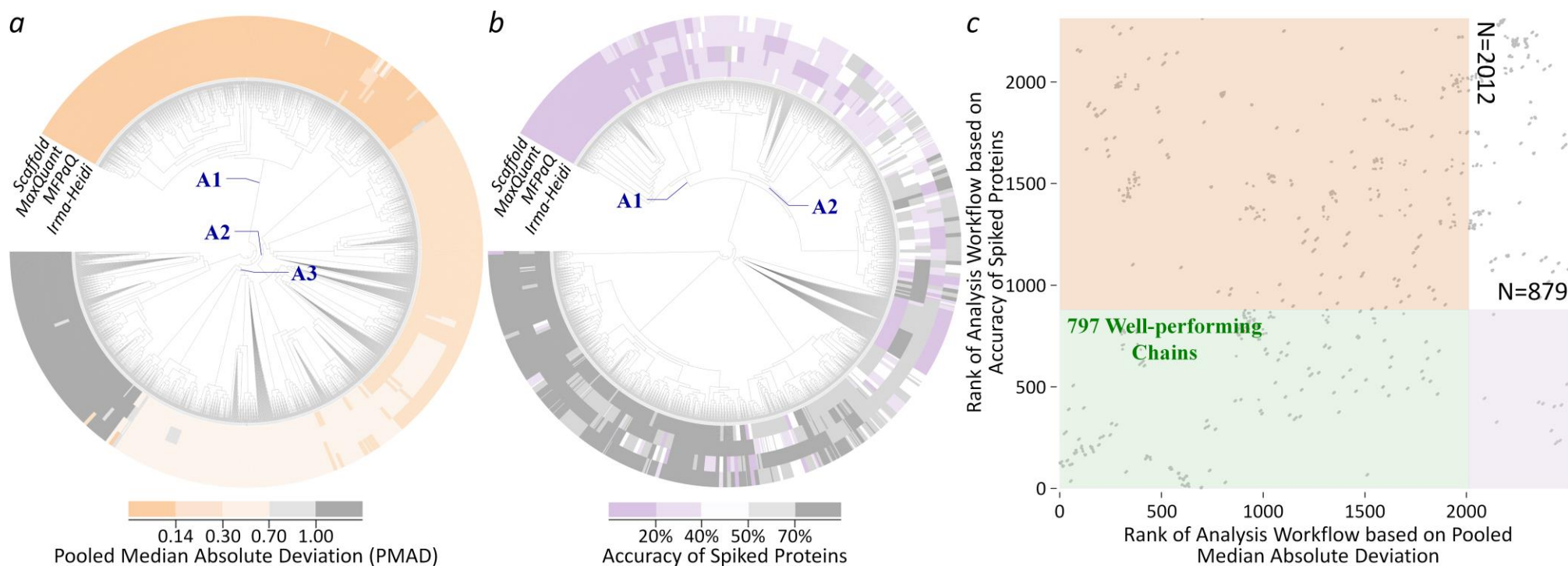
**Figure S9.** The strategy proposed in this study to discover manipulation chains with simultaneously improved *precision* and *accuracy* based on the benchmark from **Table 2** *study 2* of distinct concentrations (5 *vs* 50fmol/μg) of the spiked UPS1 proteins. First, the clustering analyses among manipulation chains across four quantification tools were conducted for (***a***) *precision* and (***b***) *accuracy*. Second, 2D scatter plot (***c***) was drawn to provide the ranks of manipulation chains (represented by gray dots) collectively determined by *precision* (horizontal axis) and *accuracy* (vertical axis). The pink, violet and green areas in (***c***) indicated the chains of good *precision* (**A1**+**A2** in (***a***)), good *accuracy* (**A1**+**A2** in (***b***)) and good-performance for both *precision* and *accuracy*, respectively. All in all, 827 chains (within the green region of (***c***)) were found to perform well under both criteria (*precision* & *accuracy*).

**Figure S10.** The strategy proposed in this study to discover manipulation chains with simultaneously improved *precision* & *accuracy* based on the benchmark from **Table 2** *study 2* of distinct concentrations (12.5 *vs* 50fmol/µg) of spiked UPS1 proteins. First, the clustering analyses among manipulation chains across four quantification tools were conducted for (***a***) *precision* and (***b***) *accuracy*. Second, 2D scatter plot (***c***) was drawn to provide the ranks of manipulation chains (represented by gray dots) collectively determined by *precision* (horizontal axis) and *accuracy* (vertical axis). The pink, violet and green areas in (***c***) indicated the chains of good *precision* (**A1**+**A2**+**A3** in (***a***)), good *accuracy* (**A1**+**A2** in (***b***)) and good-performance for both *precision* and *accuracy*, respectively. All in all, 441 chains (within the green region of (***c***)) were found to perform well under both criteria (*precision* & *accuracy*).
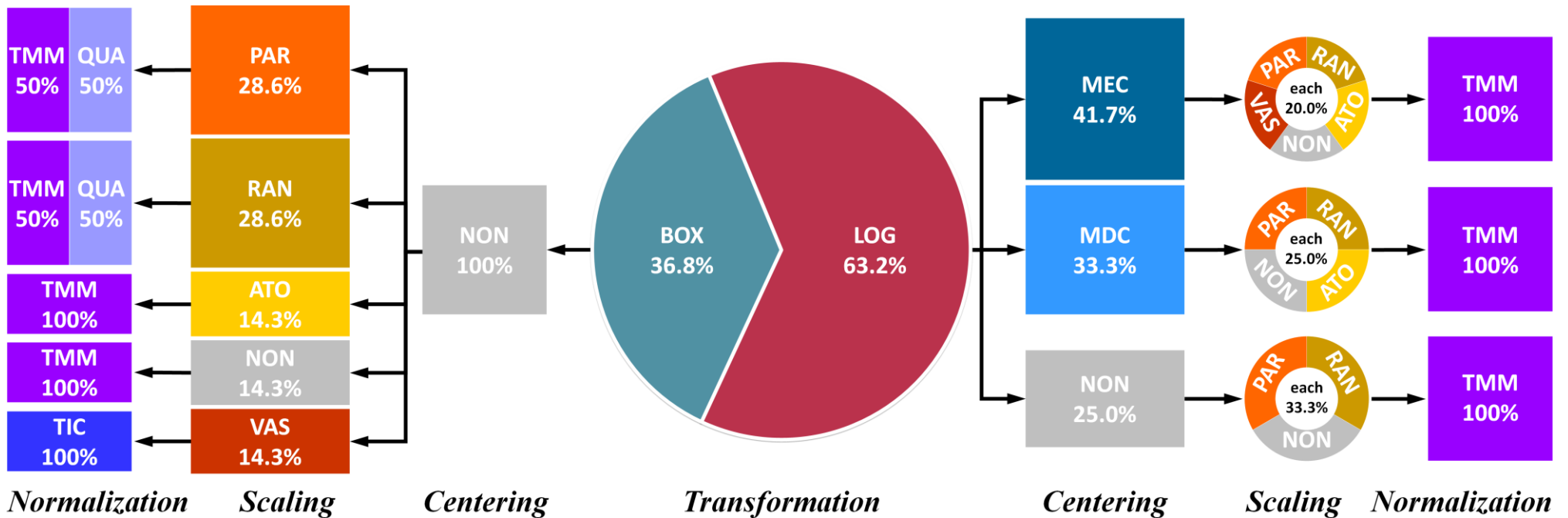
**Figure S11.** The strategy proposed in this study to discover manipulation chains with simultaneously improved *precision* & *accuracy* based on the benchmark from **Table 2** *study 2* of distinct concentrations (25 *vs* 50fmol/µg) of the spiked UPS1 protein. First, the clustering analyses among manipulation chains across four quantification tools were conducted for (***a***) *precision* and (***b***) *accuracy*. Second, 2D scatter plot (***c***) was drawn to provide the ranks of manipulation chains (represented by gray dots) collectively determined by *precision* (horizontal axis) and *accuracy* (vertical axis). The pink, violet and green areas in (***c***) indicated the chains of good *precision* (**A1**+**A2**+**A3** in (***a***)), good *accuracy* (**A1**+**A2** in (***b***)) and good-performance for both *precision* and *accuracy*, respectively. All in all, 797 chains (within the green region of (***c***)) were found to perform well under both criteria (*precision* & *accuracy*).

**Figure S12.** Distribution of manipulation methods in these 133 well-performing chains identified based on forty DDA-based benchmarks of **Table 2** *study 2*. Manipulation method was abbreviated by three-letter code which was defined in **Table 1**. For the seven *imputation* methods together with the non-*imputation* (NON), they demonstrated the exactly equal chances within the 133 well-performing chains, which showed that the selection of different *imputation* methods (even NON) had nothing to do with the performance under this circumstance. Therefore, the *imputation* methods were not displayed in this distribution.

# Supplementary Methods

## 1. Detailed Descriptions of Quantification Tools that Were Used in This Work

### 1.1 Quantification Tools for Pre-processing the Data Acquired Based on *SWATH-MS*

*DIA-UMPIRE*: A comprehensive computational workflow and open-source software for processing the data independent acquisition (DIA) MS-based proteomics[44]. It enables untargeted protein quantification based on SWATH-MS data obtained by Orbitrap family of mass spectrometers[45], and also enables targeted extraction of quantitative information based on peptides initially identified in only a subset of the samples, resulting in more consistent quantification across multiple samples[44]. It is used to identify similar number of peptide ions with better identification reproducibility between replicates and samples, than conventional data-dependent acquisition[46]. It has also been used to process untargeted data for identifying host cell protein[6] and export the peptide identification results of pseudo-MS2 spectra[47].

*OpenSWATH*: An open-source software that allows targeted analysis of DIA data based on SWATH-MS in automated, high-throughput fashion[48]. It is cross-platform software relying on open data formats, allowing it to analyze DIA data from multiple instrument vendors and is integrated and distributed with OpenMS[49]. It is widely applied to analyze the proteome of streptococcus pyogene[48], to estimate $q$-value of peptide and protein level[50]. Its generic utility for all types of modification and its scalability could enable confident quantification of the post-translational modifications in DIA-based large-scale studies[50].

*PeakView*: A commercial software (also name SWATH 2.0) which covers all major components of in-silico process in SWATH workflow from extended assay library building to final statistical analysis and reporting[51]. PeakView uses a set of processing settings to filter the ion library and determine which peptides or transitions should be adopted for proteome quantification[52], which is demonstrated to be a powerful strategy particularly for marker discovery and clinical applications[53]. It was used for N-linked glycoproteins enrichments prior to the tryptic digestion, library creation & analysis[54], evaluating the amount of sample needed for PCT-SWATH analysis[55] and selecting the best method for extracting green algae[16].

*Skyline*: An open source application for building selected reaction monitoring, multiple reaction monitoring, parallel reaction monitoring (targeted MS/MS), DIA/SWATH, targeted DDA of MS1 quantitative methods[16]. It was explicitly designed to accelerate targeted proteomics and foster broad sharing of the method and results across instrument platforms[56]. It has been used to peptides and transition selection for targeted experiments[57], the retention time determination for scheduled MS[58] and isolation window determination for DIA[59]

*Spectronaut*: For targeted analysis of DIA measurement based on SWATH-MS independent of MS vendor[60]. It demonstrates a powerful ability to peak picking and automatic interference correction by utilizing spectral libraries generated from the raw data acquired on multiple instrument platforms, and is specifically designed to support spectral-library-free workflow and targeted analysis of OMICs data by hyper reaction monitoring[61,

[62]. It is widely applied to the DIA-based quantitative proteome profiling[61], improved proteomic quantification by sequential window acquisition[62] and high-precision indexed retention time prediction in targeted DIA[60].

## 1.2 Quantification tool for Pre-processing the Data Acquired Based on *Peak Intensity*

*MaxQuant*: Integrated suite of algorithms for processing the high-resolution, quantitative mass-spectrometry data, which is one of the most frequently used platforms for analyzing the MS-based proteome data[63, 64]. It is widely used to analyze the tandem spectra generated by collision-induced dissociation (CID), higher-energy collisional dissociation (HECD) and electron transfer dissociation (ETD)[65]. *MaxQuant* is used for analyzing datasets derived from all major relative quantification techniques, including label-free quantification[64], MS1-level labeling readouts and isobaric MS2-level labeling readouts[66].

*MFPaQ*: A web-based tool which runs on a server on which Mascot Server 2.1 and Perl 5.8 must be installed. To perform quantification, the external module: *Extract Daemon* is developed to extract intensity values from the raw proteomic data[67]. One of its distinguished features lines in its quantification modules, which provides information on protein relative expression following the isotopic labeling and identification with the Mascot[2]. It has been used to quantify the membrane proteins from primary human endothelial cell[67], and SILAC-based proteomic profiling of the human MDA-MB-231 metastatic breast cancer cell line[68].

*OpenMS*: A robust, open-source, cross-platform software specifically designed for flexible and reproducible analysis of high-throughput MS data[66]. It uses the modern software engineering concepts with the emphasis on modularity, reusability and extensive testing using continuous integration, and implements common mass spectrometric data processing tasks through well-defined application programming interface and through the standardized open data format[52]. OpenMS is widely applied to the quantitative and variant enabled mapping of peptides to genome[69], the analysis of cerebrospinal fluid proteome in alzheimer's disease[70] and quantitative analysis of label-free LC-MS data for comparative candidate biomarker studies[71].

*PEAKS*: Software platform with complete solution for discovery proteomics, including protein identification and quantification, analysis of posttranslational modification and sequence variants, and peptide/protein de novo sequencing[4]. It relies on sophisticated dynamic programming algorithm to efficiently compute the best peptide sequence whose fragment ions can best interpret the peaks in the MS/MS spectrum. It is thus a useful tool for protein identification and quantification of known and unknown genomes[4]. PEAKS has matured into a comprehensive proteomic platform supporting the analysis of label-free and label-based data, and achieves significantly improved accuracy and sensitivity over other commonly applied software packages[72].

*Progenesis*: New generation of bioinformatics vehicle targeting small molecule analysis for proteomics that quantifies proteins based on peak intensity[73]. It allows full operator control over every processing step such as alignment of peptide ion signal landscapes and indeed individual peptide ion signal peaks[74]. It is used for protein label-free quantification and peak picking with the automatic sensitivity[75].

*Proteios SE*: A tool integrating protein identification search engine access into several proteomic workflows, both gel-based and liquid chromatography-based and allowing seamless combination of search result, protein inference, protein annotation and quantitation[76]. It is targeted large project of shared data, integrated sample tracking and aimed at becoming standard analysis platform in proteomics, whose major feature is automated linking of data from various proteomic pipelines. It has built-in support to some protein identification engines such as Mascot, X!Tandem, and combines search results from multiple engines, and automatically generates the protein identification reports containing information required for publication of proteomics results[77].

*Scaffold*: Commercial bioinformatics tool, which attempts to increase the confidence in protein identification reports through the use of several statistical methods. It supports a wide variety of search engines and uses a pipeline of several peptide and protein validation methods after an initial database-search analysis[78]. Scaffold has been widely used to the identification of proteome for new target to inhibit yellow fever virus replication[7], analysis of the follicle fluid proteome for preconception folic acids[79, 80].

*Thermo Proteome Discoverer*: Tool for workflow-driven data analysis in proteomics integrating all different steps in quantitative proteomic experiment (MS spectrum extraction, peptide identification and quantification) into the user-configurable, automated workflows[81]. It has a convenient graphical user interface in which users can load raw data directly from the instrument and explore and analyze it since it supports multiple sequence database search engines (Sequest HT, Mascot), spectral library searching, peptide spectrum-match validation (Percolator), as well as various quantification techniques, like isobaric mass tagging (e.g., iTRAQ, TMT) or SILAC[82]. It has been applied to iTRAQ-based quantitative analysis of protein mixture[83].

## 1.3 Quantification tool for Pre-processing the Data Acquired Based on *Spectral Counting*

*Abacus*: A stand-alone tool for extracting and processing *spectral count* data[9] aiming at streamlining analysis of *spectral count* data by providing automated solutions and extracting information from proteomics data for statistical analysis[4]. It has the disadvantage of losing information or attempting to apportion large number of spectra based on relatively small set of differentiating spectra[84]. It is compatible with popular trans-proteomic pipeline suite of tools and comes with a graphical user interface making it easy to interact with the program[9].

*Census*: A quantitative tool analyzing the high-throughput mass spectrometry data from shotgun proteomics experiments in efficient ways and various stable isotope labeling experiments (e.g., $^{15}$N, $^{18}$O, SILAC, iTRAQ and TMT) in addition to labeling-free experiment[85]. It is flexible in handling the most quantitative proteomics labeling strategies, as well as label-free experiment with multiple statistical algorithms to improve quality of results[86]. It is used to discover differential proteins in plasmodium falciparum patients under drug treatment[87].

*DTASelect*: A java tool used to organize, filter and interpret results generated by SEQUEST (one of the most widely used protein database searching programs for tandem mass spectrometry)[88]. It assembles protein-level signals from peptide data and focuses on peptides of interest by sweeping away the less likely identification[88].

It makes complex experiment feasible by streamlining data analysis for proteomics[89]. It can be applied for a proteogenomic study with a controlled false discovery rate[90] and palmitoylated protein identifications[91].

***IRMa-hEIDI***: IRMa toolbox provides an interactive application to assist in the validation of Mascot® search results, and allows an automatic filtering of Mascot identification result and manual confirmation or rejection of individual PSM (a match between a fragmentation mass spectrum and a peptide)[92]. Its main originality is to filter matches rather than identified proteins and its features are easy navigation within identification result and batch mode to automatically validate multiple identification results[92]. It is used to filter the *spectral count* workflows results with the compromise between sensitivity and false discovery rate[12].

***MFPaQ***: A software facilitating organization, mining and validation of Mascot results and offering different functionalities to work on validated protein lists, and data quantification by isotopic labeling methods or label free approaches[67]. MFPaQ extracts quantitative data from raw files obtained by nano-LC-MS/MS, calculates peptide ratios, and generates a non-redundant list of proteins identified in a multisearch experiment with their calculated averaged and normalized ratio[67]. It is used to large scale analysis of inflammatory endothelial cell[2], and the label-free quantification of cerebrospinal fluid by combining peptide ligand library treatment[93].

***ProteinProphet***: A statistical model designed for computing probabilities that proteins are present in a sample on the basis of peptides assigned to tandem mass (MS/MS) spectra acquired from proteolytic digest sample[94]. It allows the filtering of large-scale proteomic data with a predictable sensitivity and false positive discovery rates[94]. It was applied to discriminate true assignments of spectra to peptide sequences from false assignments, to assign probability for each identified peptide, and to compute sensitivity and error rates for the assignment of spectra to the sequences in each experiment[95]. It was used to infer the protein identification and to compute probabilities that a protein had been correctly identified, based on the available peptide sequence evidence[94].

## 2. Detailed Descriptions of Three *Transformation* Methods Applied in This Study

***Box-cox Transformation (BOX)***: A parametric power transformation technique in order to reduce anomalies such as non-additivity, non-normality and heteroscedasticity[96]. The method has been extensively studied, and an attempt is made to review its corresponding studies[96]. This method has been used to facilitate the discovery novel biomarkers and the development of new therapeutic target for seven important liver diseases based on the proteomic and transcriptomic data[19]. In this study, the parameter *lambda* ($\lambda$) of the BOX method was set to 0.3, and its algorithm was programed and implemented under the ***R*** environment (version 3.5.1).

***Log Transformation (LOG)***: For obtaining the more symmetric distribution prior to statistical analysis, LOG is carried out almost routinely[97]. It works for data where you can see that the residual get bigger for the bigger values of the dependent variable[97]. Such trends in the residuals occur frequently, because the errors or changes in the value of an outcome variable is often a percent of the value rather than an absolute value[97]. This method has been applied for identifying new therapeutic target of early-stage hepatocellular carcinoma[20]. In this study,

LOG was performed by <u>log2-scale</u>, and its algorithm was programed and implemented under **R** environment.

***Variance Stabilization Normalization (VSN)***: As a *transformation* method that integrated with *normalization* technique, the VSN was unique in determining data-dependent transformation parameters by having built-in *transformation*[98, 99]. It is one of the non-linear methods aiming at keeping variances constant over entire data range, therefore removing heteroscedasticity. For small protein abundances, it performs linear *transformation* behavior to make variance unchanged[99]. It was developed for label-free relative quantification of endogenous peptides[30]. In this study, the VSN method was implemented with the <u>*justvsn*</u> function using **R**/**Bioconductor** vsn-package and programed with the default parameter settings under **R** environment (version 3.5.1).

To facilitate the reproduction of the *transformation* methods discussed above, their ***source code*** can be readily downloaded from the official website (https://idrblab.org/anplea/) of the newly developed online tool.

## 3. Detailed Descriptions of Eighteen *Pretreatment* Methods Applied in This Study

### 3.1 Two *Centering* Methods

***Mean centering (MEC)***: Converting all the concentrations to fluctuations around zero instead of around the mean of the protein intensities. Hereby, it adjusts for differences in the offset between high and low abundant proteins. It is therefore used to focus on the fluctuating part of the data, and leaves only the relevant variation (being the variation between the samples) for analysis[100]. MEC has been applied for the improvement of the sensitivity of significance tests in *spectral counting*-based comparative discovery proteomics[22]. In this study, the algorithm of MEC was programed by integrating the basic <u>*mean*</u> function (<u>*mean value*</u>) in the **R**-statistical programming and implemented under the **R** environment (version 3.5.1).

***Median centering (MDC)***: Converting all the concentrations to the fluctuations around zero instead of around the median of the protein intensities. Hereby, it adjusts for differences in the offset between proteins of high and low abundances. It is thus applied to focus on the fluctuating part of the data, and leaves only the relevant variation (being the variation between samples) for analysis[100]. MDC facilitates the normalization procedures in LC-MS proteomics experiments through dataset dependent ranking of normalization scaling factors[23]. In this study, the algorithm of MDC was programed by integrating the basic <u>*median*</u> function (<u>*median value*</u>) in the **R**-statistical programming and implemented under the **R** environment (version 3.5.1).

To facilitate the reproduction of the *centering* methods discussed above, their ***source code*** can be downloaded from the official website (https://idrblab.org/anplea/) of the newly developed online tool.

### 3.2 Four *Scaling* Methods

***Auto Scaling (Unit Variance Scaling, ATO)***: One of the simplest methods for adjusting proteomics variances, scaling protein intensities based on the standard deviation of proteomic data[99]. It scales all protein intensities to unit variance, and all intensities are equally important and comparably scaled[101]. All data is analyzed based

on correlations and standard deviation of all intensities[99]. It has been used to identify proteomic markers for psoriasis and psoriasis arthritis[24] and normalize LC-MS proteomics based on scan-level data[102]. In this study, the algorithm of ATO was programed by integrating basic _sd_ function (_standard deviation_) in the **R**-statistical programming and implemented under the **R** environment (version 3.5.1).

***Pareto Scaling (PAR)***: Using square root of the standard deviation of the data as scaling factor[99], this method is able to reduce the weight of large fold changes in protein intensities, which is more significant than ATO[99]. The dominant weight of extremely large fold changes may still be unchanged[99]. Therefore, the disadvantage of PAR is the sensitivity to large fold changes[100]. The PAR was applied to normalize LC-MS proteomics data using scan-level information in the Gaussian process regression model[102]. In this study, the algorithm of PAR was programed by integrating two basic functions _sd_ and _sqrt_ (_standard deviation_ and _square root_) in the **R**-statistical programming and implemented under the **R** environment (version 3.5.1).

***Vast Scaling (VAS)***: An acronym of variable ability scaling and an extension of autoscaling that focuses on stable variables and uses standard deviation and the so-called coefficient of variation (cv) as scaling factors[100]. It was applied for investigating the feasibility of OMICs for immediate analysis of resection margins during breast cancer surgery[103]. In this study, the algorithm of VAS was programed by integrating two basic functions _var_ and _mean_ (_variance_ and _mean value_) in **R**-statistical program and implemented under **R** environment.

***Range Scaling (RAN)***: The measured intensity was divided by the range of the intensities over all samples[104]. Biological range is the difference between the minimal and maximal abundances reached by a certain protein in a set of experiments and RAN uses it as a scaling factor[100]. RAN was applied to a fuzzy C-means clustering for the chromatographic fingerprint analysis[105]. In this study, RAN's algorithm was programed by integrating two basic functions _max_ and _min_ (_maximum value_ and _minimum value_) in the **R**-statistical programming and implemented under the **R** environment (version 3.5.1).

To facilitate the reproduction of the _scaling_ methods shown above, the **source code** can be downloaded from the official website (https://idrblab.org/anplea/) of the newly developed online tool.

## 3.3 Twelve _Normalization_ Methods

***Cyclic Loess (Cyclic Locally Weighted Regression, CYC)***: Originated from the combination of MA-plot and logging Bland-Altman plot by assuming the existence of non-linear bias[99]. It can estimate regression surface using multivariate smoothing procedure[106], but is one of the most time-consuming methods[107]. CYC has been applied to the proteomic profiling in the context of common experimental designs[108]. In this study, the CYC method was implemented with the _CyclicLoess_ function using the limma-package in **R**/***Bioconductor***, which was then programed with the default parameter settings under **R** environment (version 3.5.1).

***EigenMS (EIG)***: Removing the biases of unknown complexity from LC/MS-based proteomics data. It allows for increased sensitivity in differential analysis, and aims at preserving the original difference while removing

the bias from the data[28]. It preserves true differences by estimating the treatment effects with ANOVA model, and has been used to profile MS-based quantitative label-free proteomics[109-111]. In this study, the EIG method was implemented using its **R**-codes available for downloading from the *Sourceforge-repositories*.

*Linear Baseline (Linear Baseline Scaling, LIN)*: Based on the assumption of the constant linear relationship between each feature of a given spectrum and the baseline[99], it maps each spectrum to the baseline[99]. Baseline is the median of each feature across all spectra and the scaling factor is computed as the ratio of mean intensity of the baseline to the mean intensity of each spectrum[99]. The intensities of all spectra are multiplied by their particular scaling factors[99], but this assumption of a linear correlation may be oversimplified[99]. In this study, the LIN method was programed by integrating two basic functions *median* and *mean* (*median value* and *mean value*) in the **R**-statistical programming and implemented under the **R** environment (version 3.5.1).

*Locally Weighted Scatterplot Smoothing (LOW)*: A *normalization* method assuming that the systematic bias is non-linearly dependent on the magnitude of peptide abundances[112]. This non-linearity potentially originates from the effects of the ion suppression on measured peptide abundances, or on measured peptide abundances approaching detector saturation or background[112]. This method has been applied to MS-based proteomics[110]. In this study, the LOW method was implemented with the *preprocess* function using the LPE-package in the **R** environment, which was then programed with the *LOWESS* parameter settings under **R** environment.

*Mean Normalization (MEA)*: Commonly used method to normalize data by the mean value of all signals to eliminate background effect[113]. Intensity of each protein in a given sample is adopted by the mean of intensity of all variables in samples[97]. To make the samples comparable, the mean of intensities for each experimental run is forced to equal to one another using this method[114]. Each sample is scaled such that the mean of protein abundances in a sample equals one[97]. This method has been applied in the profiling of urine peptidome[115]. In this study, the MEA method was implemented with the *Normalise* function using the metabolomics-package in **R**/**Bioconductor**, which was then programed with the *mean* parameter settings under **R** environment.

*Median normalization (MED)*: Based on assumption that the samples of a dataset are separated by a constant, it scales the samples so that they have the same medians[116]. As one of the commonly applied methods without the need for internal standards, it is more practical than sum *normalization* especially when several saturated abundances are associated with the factors of interests[116]. It has previously been used in MS-based label-free proteomics analysis for removing systematic biases associated with mass spectrometry[112]. In this study, MED method was implemented with the *Normalise* function using the metabolomics-package in **R**/**Bioconductor**, which was then programed with the *median* parameter settings under **R** environment.

*Median Absolute Deviation (MAD)*: A robust measure of how a set of univariate samples of quantitative data spreads out especially when data is unnormal. MAD takes the absolute deviations based on the median within a sample to normalize rather than directly uses median like MED[31]. It has been applied to improve the quality

control procedures of peptide-centric LC-MS proteomics data[31]. In this study, MAD method was programed by integrating two basic functions *median* and *mad* (*median value* and *median absolute deviation*) in the **R**-statistical programming and implemented under the **R** environment (version 3.5.1).

***Probabilistic Quotient Normalization (PQN)***: Based on an overall estimation on the most probable dilutions, it transforms the proteomics spectra[117]. It has been reported to be significantly robust and accurate comparing to the integral and vector length normalizations[117]. PQN performs an integral normalization of each spectrum, calculate the quotient between test and reference spectrum, then all variables of the test spectrum are divided by the median quotient[28]. PQN has been applied in MALDI-TOF mass spectrometry knowledge discovery[118]. In this study, the median values over all samples as the reference sample (to which all the other samples were normalized to) were first calculated, and then PQN method was implemented with the *median* function in the **R**-statistical programming and implemented under the **R** environment.

***Quantile (Quantile Normalization, QUA)***: Aiming at achieving the same distributions of protein abundances across all samples, quantile-quantile plot in this method is used to visualize distribution similarity[99]. Quantile is motivated by the idea that the distribution of two data vectors is the same if quantile-quantile plot is straight diagonal line[116]. QUA has been adopted for removing systematic bias associated with mass spectrometry and label-free proteomics[112]. In this study, QUA method was implemented with the *normalize.quantiles* function using the affy-package in **R**/***Bioconductor***, which was then programed with default parameter settings under **R** environment (version 3.5.1).

***Robust Linear Regression (RLR)***: One robust measure is used for transference when you want to rescale one reference interval to another scale and it is robust against outliers in the data than linear regression using least squares estimation[28, 36]. In this study, the median values over all the samples as the reference sample (to which all the other samples were normalized to) were first calculated, and then RLR method was implemented with the *rlm* function using the Normalyzer-package in the **R** environment, which was then programed with default parameter settings.

***Total Ion Current (TIC)***: Summing all separate ion currents carried by the ions of different m/z, it contributes to complete mass spectrum or in specified m/z range of mass spectrum. The sum of all peak areas of peptides unique to a particular organism is called pTIC (proteome total ion current)[119]. This method has been used in MALDI-TOF and SELDI-TOF mass spectra proteomic profiling[120]. In this study, TIC method was programed by integrating the basic *sum* function (*summation*) in **R**-statistical programming and implemented under the **R** environment (version 3.5.1).

***Trimmed Mean of M Values (TMM)***: Estimating scale factors between samples that can be incorporated into current statistical methods for differential abundance analysis in proteomics, and removing the low-expressed proteins[121]. In this study, TMM method was implemented with the *tmm* function using the NOISeq-package

in *R*/*Bioconductor*, and was then programed with the default parameter settings under *R* environment.

To facilitate the reproduction of the *normalization* methods discussed above, their source code can be readily downloaded from the official website (https://idrblab.org/anplea/) of the newly developed online tool.

**4. Detailed Descriptions of Seven Missing-value *Imputation* Methods Applied in This Study**

***Background Imputation (BAK)***: Simulating the situation where protein values are missing because of having small concentration in samples and thus cannot be detected during MS runs[122]. Missing values were replaced with the lowest detected intensity value of dataset as a representative of background[122]. This method has been applied in popular proteomic software workflows for label-free proteome quantification and imputation[62]. In this study, the missing values were first identified, and then the BAK method was programed using the basic *min* function (*minimum value*) in the *R*-statistical programming and implemented under *R* environment.

***Bayesian Principal Component Imputation (BPCA)***: Reported as out-performing KNN and SVD[122]. One of its features allowing it better performance than these two is its capacity to auto-select the parameters used in the estimation[122]. BPCA produces improved estimation when the sample size is huge[122] and is applied to treat missing values for multivariate statistical analysis of gel-based proteomics data[39]. In this study, BPCA method was implemented with the *pca* function using the pcaMethods-package in *R*/*Bioconductor*, which was then programed with the parameter settings of *bpca method* and *nPcs equaling to 3* under *R* environment.

***Censored Imputation (CEN)***: Considered as being 'missing completely at random', no value is imputed for it if only single NA for a protein in a sample group was found[62]. If a protein contained more than one missing value in one sample group (consisting of technical replicates), they are considered missing due to being below detection capacity, and the lowest intensity value in the data set is imputed for them[62]. CEN has been used to improve detection of the differentially abundant proteins[40]. In this study, missing values were first identified and assessed. When the value of a particular protein is missing in multiple (>1) samples within single sample group, the CEN method was programed and applied using the basic *min* function (*minimum value*) in the *R*-statistical programming and implemented under *R* environment.

***K-nearest Neighbor Imputation (KNN)***: Aiming at identifying K proteins similar to the one of missing value, where the similarity is estimated by Euclidean distance measure, and the missing values are imputed with the values of weighted average from the neighboring proteins[122]. KNN tends to select proteins with an expression profile similar to the proteins of interest, and it outperforms BPCA and LLS in relatively small size datasets[122]. This method has been used in integrative analyses of multi-omics data[41]. In this study, the KNN method was implemented with the *impute.knn* function using the impute-package in *R*/*Bioconductor*, and then programed with the parameter settings of *k value equaling to 10* under *R* environment.

***Local Least Squares Imputation (LLS)***: This technique exploits the local similarity structures in the data, as well as the least squares optimization process[122]. It represents a protein of missing value as linear combination

29

of similar proteins[123]. Similar proteins are chosen by K-nearest neighbors that have large absolute values of *Pearson* correlation. As a nonparametric missing value estimation method, LLS was designed by introducing an automatic K-value estimator[123], and it is used in missing value imputation for proteomics data or any data that can be represented as matrix (e.g. NGS or microarray data)[42]. In this study, LLS method was implemented with the *llsImpute* function using the pcaMethods-package in **R**/**Bioconductor**, and then programed with the *k value equaling to 10* together with the default parameter settings under **R** environment.

***Singular Value Decomposition (SVD)***: Known as Karhunen-Loève expansion in the pattern recognition and as principal-component analyses in statistics[124], it is a linear transformation of protein abundance data[124]. In contrast to KNN imputation, SVD attempts to utilize global information in the entire matrix to predict missing values[125]. Its basic concept is to find the dominant component to summarize the entire matrix and then predict missing values in target protein by regressing against dominant components[125]. This method has been applied to enable greater accuracy and precision in quantitative comparison of the protein abundances[29]. In this study, the SVD method was implemented with the *pca* function using the pcaMethods-package in **R**/**Bioconductor**, and then programed with the parameter settings of *svdImpute method* under **R** environment.

***Zero Imputation (ZER)***: Deemed to the simplest imputation by replacing the missing values with zeros. This zero replacement method does not utilize any information about the data[125]. The integrity and usefulness of the data can be jeopardized by zero imputation since erroneous relationship among proteins can be artificially created[125]. ZER has been used in the analysis of quantitative proteomic experiment that use isobaric tagging[43]. In this study, the ZER method was programed in the **R**-statistical programming and implemented under the **R** environment (version 3.5.1).

To facilitate the reproduction of missing value *imputation* methods discussed above, their source code can be readily downloaded from the official website (https://idrblab.org/anplea/) of the newly developed online tool.

## 5. Detailed Descriptions of Criteria for Performance Evaluation Applied in This Work.

In total, 5 well-established criteria for a comprehensive evaluation on the performance of LFQs are provided, and each criterion is either quantitatively or qualitatively assessed by various metrics.

### 5.1. Precision Measuring the Reduction in Proteome Variation among Replicates

Different modes of acquisition, various kinds of software for pre-processing raw proteomics data, and diverse methods for data manipulation (such as *transformation*, *pretreatment* and value *imputation*) profoundly affect the *precision* of LFQ, which could be assessed by the pooled median absolute deviation (PMAD) of reported protein intensities among replicates[28]. In particular, the metric PMAD is designed to reflect LFQ's ability to reduce variations among replicates, and thus to enhance the technical reproducibility[126]. The lower value of PMAD denotes more thorough removal of experimentally induced noise and indicates better *precision*.

## 5.2. Classification Ability of Proteome Quantification between Distinct Sample Groups

An appropriate LFQ is expected to retain or even enlarge the difference in proteomic data between distinct sample groups[127]. A heatmap hierarchically clustering samples based on their protein intensities is therefore frequently used as effective metric to assess LFQ's classification ability[127]. Firstly, the total number of protein intensities in each sample is reduced using feature selection (*first*, the differential significance of each protein between distinct sample groups measured by adjusted *p*-value was calculated by ROTS package; *second*, the significant features (adjusted *p*-value <0.05) were selected for subsequent heatmap analyses). Then, proteins (rows) and samples (columns) are clustered based on their similarities in protein intensities. Detail procedure for assessing LFQ's classification ability can be found in the publication by Griffin NM, *et al*[127].

## 5.3. Differential Abundance Analysis in Proteomics Based on Reproducibility-optimization

To avoid overfitting/confounding, the distribution of *p*-values of protein intensities between distinct sample groups is explored[128]. Ideally, a uniform distribution for the bulk of non-differential proteins is expected with a peak in the [0.00, 0.05] interval corresponding to proteins with differential intensities[128]. In proteomics (and other OMICs) studies that explore the mechanism underling complex biological process, a limited number of proteins of differential abundance may resulted in false discovery[129]. Thus, the differential significance of protein intensities between distinct sample groups measured by *p*-values is first calculated by reproducibility-optimized test statistic (ROTS) package[130]. A skewed distribution of *p*-values may indicate overfitting and/or confounding[110].

## 5.4. Robustness among Different Sets of Protein Biomarkers Identified from Different Datasets

Consistency score is a popular criterion used to represent the robustness of protein marker identification[131], which is calculated to quantitatively measure the overlap of identified biomarkers among different partitions of a given dataset[132]. A higher consistency score represents the more robust results in marker identification[131]. Thus, the random sampling is first preformed within the quantified dataset to produce multiple sub-datasets. Then, each protein is ranked according to its significance measured by the *q*-value and absolute fold changes. Third, the top-ranked proteins in each sub-dataset are selected as biomarkers. Finally, a consistency score is calculated based on these markers using equation[132] as follow:

$$S = \sum_{i=2}^{C} \sum_{S \in I_i} 2^{i-2} \cdot n_S$$

where $C$ denoted the total number of sub-datasets, $I_i$ indicated a set of significant biomarkers containing the intersections of any $i$ sub-datasets, and $n_S$ referred to the number of markers in the intersection $S$.

## 5.5. Accuracy Assessing the Deviation of Spiked Proteins from Their Expected Abundance Ratio

Additional experimental data (e.g. spiked proteins) are frequently generated and used as references to validate

or adjust the performance of LFQ[61, 133], and expected log fold changes (logFCs) are known both for the spiked and background proteins (the expected LogFC for background proteins equals to zero)[62]. Herein, the logFCs of protein intensities (for both the spiked and background proteins) between distinct sample groups were first calculated, and the level of correspondence between the quantification and expected logFCs is then assessed using the mean squared error (MSE). The performance of LFQ can be reflected by how well the quantification logFCs corresponded to what are expected based on the references[62]. The preferred median values would be zero with minimized deviations.

# References

1   B. O. Schroeder; G. M. H. Birchenough; M. Stahlman; L. Arike; M. E. V. Johansson; G. C. Hansson; F. Backhed. Bifidobacteria or Fiber Protects against Diet-Induced Microbiota-Mediated Colonic Mucus Deterioration. Cell Host Microbe. 2018, 23(1): 27-40 e7

2   V. Gautier; E. Mouton-Barbosa; D. Bouyssie; N. Delcourt; M. Beau; J. P. Girard; C. Cayrol; O. Burlet-Schiltz; B. Monsarrat; A. Gonzalez de Peredo. Label-free quantification and shotgun analysis of complex proteomes by one-dimensional SDS-PAGE/NanoLC-MS: evaluation for the large scale analysis of inflammatory human endothelial cells. Mol Cell Proteomics. 2012, 11(8): 527-39

3   T. Cortes; O. T. Schubert; G. Rose; K. B. Arnvig; I. Comas; R. Aebersold; D. B. Young. Genome-wide mapping of transcriptional start sites defines an extensive leaderless transcriptome in Mycobacterium tuberculosis. Cell Rep. 2013, 5(4): 1121-31

4   B. Ma; K. Zhang; C. Hendrie; C. Liang; M. Li; A. Doherty-Kirby; G. Lajoie. PEAKS: powerful software for peptide de novo sequencing by tandem mass spectrometry. Rapid Commun Mass Spectrom. 2003, 17(20): 2337-42

5   R. S. Al-Mayyahi; L. D. Sterio; J. B. Connolly; C. F. Adams; W. A. Al-Tumah; J. Sen; R. D. Emes; S. R. Hart; D. M. Chari. A proteomic investigation into mechanisms underpinning corticosteroid effects on neural stem cells. Mol Cell Neurosci. 2018, 8630-40

6   P. Nanni; L. Mezzanotte; G. Roda; A. Caponi; F. Levander; P. James; A. Roda. Differential proteomic analysis of HT29 Cl.16E and intestinal epithelial cells by LC ESI/QTOF mass spectrometry. J Proteomics. 2009, 72(5): 865-73

7   A. Vidotto; A. T. Morais; M. R. Ribeiro; C. C. Pacca; A. C. Terzian; L. H. Gil; R. Mohana-Borges; P. Gallay; M. L. Nogueira. Systems Biology Reveals NS4B-Cyclophilin A Interaction: A New Target to Inhibit YFV Replication. J Proteome Res. 2017, 16(4): 1542-55

8   I. Kaur; J. Kaur; K. Sooraj; S. Goswami; R. Saxena; V. S. Chauhan; R. Sihota. Comparative evaluation of the aqueous humor proteome of primary angle closure and primary open angle glaucomas and age-related cataract eyes. Int Ophthalmol. 2018,

9   D. Fermin; V. Basrur; A. K. Yocum; A. I. Nesvizhskii. Abacus: a computational tool for extracting and pre-processing spectral count data for label-free quantitative proteomic analysis. Proteomics. 2011, 11(7): 1340-5

10  J. H. Prieto; S. Koncarevic; S. K. Park; J. Yates, 3rd; K. Becker. Large-scale differential proteome analysis in Plasmodium falciparum under drug treatment. PLoS One. 2008, 3(12): e4098

11  M. G. Gravett; M. J. Novy; R. G. Rosenfeld; A. P. Reddy; T. Jacob; M. Turner; A. McCormack; J. A. Lapidus; J. Hitti; D. A. Eschenbach; C. T. Roberts, Jr.; S. R. Nagalla. Diagnosis of intra-amniotic infection by proteomic profiling and identification of novel biomarkers. JAMA. 2004, 292(4): 462-9

12  C. Ramus; A. Hovasse; M. Marcellin; A. M. Hesse; E. Mouton-Barbosa; D. Bouyssie; S. Vaca; C. Carapito; K. Chaoui; C. Bruley; J. Garin; S. Cianferani; M. Ferro; A. V. Dorssaeler; O. Burlet-Schiltz; C. Schaeffer; Y. Coute; A. Gonzalez de Peredo. Spiked proteomic standard dataset for testing label-free quantitative software and statistical methods. Data Brief. 2016, 6286-94

13 Y. A. Goo; A. Y. Liu; S. Ryu; S. A. Shaffer; L. Malmstrom; L. Page; L. T. Nguyen; C. E. Doneanu; D. R. Goodlett. Identification of secreted glycoproteins of human prostate and bladder stromal cells by comparative quantitative proteomics. Prostate. 2009, 69(1): 49-61

14 S. Kreimer; Y. Gao; S. Ray; M. Jin; Z. Tan; N. A. Mussa; L. Tao; Z. Li; A. R. Ivanov; B. L. Karger. Host Cell Protein Profiling by Targeted and Untargeted Analysis of Data Independent Acquisition Mass Spectrometry Data with Parallel Reaction Monitoring Verification. Anal Chem. 2017, 89(10): 5294-302

15 E. Ahrman; O. Hallgren; L. Malmstrom; U. Hedstrom; A. Malmstrom; L. Bjermer; X. H. Zhou; G. Westergren-Thorsson; J. Malmstrom. Quantitative proteomic characterization of the lung extracellular matrix in chronic obstructive pulmonary disease and idiopathic pulmonary fibrosis. J Proteomics. 2018,

16 Y. Gao; T. K. Lim; Q. Lin; S. F. Li. Evaluation of sample extraction methods for proteomics analysis of green algae Chlorella vulgaris. Electrophoresis. 2016, 37(10): 1270-6

17 L. Balakrishnan; R. S. Nirujogi; S. Ahmad; M. Bhattacharjee; S. S. Manda; S. Renuse; D. S. Kelkar; Y. Subbannayya; R. Raju; R. Goel; J. K. Thomas; N. Kaur; M. Dhillon; S. G. Tankala; R. Jois; V. Vasdev; Y. Ramachandra; N. A. Sahasrabuddhe; K. Prasad Ts; S. Mohan; H. Gowda; S. Shankar; A. Pandey. Proteomic analysis of human osteoarthritis synovial fluid. Clin Proteomics. 2014, 11(1): 6

18 L. Lin; J. Zheng; Q. Yu; W. Chen; J. Xing; C. Chen; R. Tian. High throughput and accurate serum proteome profiling by integrated sample preparation technology and single-run data independent mass spectrometry analysis. J Proteomics. 2018, 1749-16

19 M. Kohl; D. A. Megger; M. Trippler; H. Meckel; M. Ahrens; T. Bracht; F. Weber; A. C. Hoffmann; H. A. Baba; B. Sitek; J. F. Schlaak; H. E. Meyer; C. Stephan; M. Eisenacher. A practical data processing workflow for multi-OMICS projects. Biochim Biophys Acta. 2014, 1844(1 Pt A): 52-62

20 Y. Jiang; A. Sun; Y. Zhao; W. Ying; H. Sun; X. Yang; B. Xing; W. Sun; L. Ren; B. Hu; C. Li; L. Zhang; G. Qin; M. Zhang; N. Chen; M. Zhang; Y. Huang; J. Zhou; Y. Zhao; M. Liu; X. Zhu; Y. Qiu; Y. Sun; C. Huang; M. Yan; M. Wang; W. Liu; F. Tian; H. Xu; J. Zhou; Z. Wu; T. Shi; W. Zhu; J. Qin; L. Xie; J. Fan; X. Qian; F. He; C. Chinese Human Proteome Project. Proteomics identifies new therapeutic targets of early-stage hepatocellular carcinoma. Nature. 2019, 567(7747): 257-61

21 N. A. Karp; W. Huber; P. G. Sadowski; P. D. Charles; S. V. Hester; K. S. Lilley. Addressing accuracy and precision issues in iTRAQ quantitation. Mol Cell Proteomics. 2010, 9(9): 1885-97

22 J. Gregori; L. Villarreal; O. Mendez; A. Sanchez; J. Baselga; J. Villanueva. Batch effects correction improves the sensitivity of significance tests in spectral counting-based comparative discovery proteomics. J Proteomics. 2012, 75(13): 3938-51

23 B. J. Webb-Robertson; M. M. Matzke; J. M. Jacobs; J. G. Pounds; K. M. Waters. A statistical selection strategy for normalization procedures in LC-MS proteomics experiments through dataset-dependent ranking of normalization scaling factors. Proteomics. 2011, 11(24): 4736-41

24 J. Reindl; J. Pesek; T. Kruger; S. Wendler; S. Nemitz; P. Muckova; R. Buchler; S. Opitz; N. Krieg; J. Norgauer; H. Rhode. Proteomic biomarkers for psoriasis and psoriasis arthritis. J Proteomics. 2016, 14055-61

25  M. S. Bereman; R. Johnson; J. Bollinger; Y. Boss; N. Shulman; B. MacLean; A. N. Hoofnagle; M. J. MacCoss. Implementation of statistical process control for proteomic experiments via LC MS/MS. J Am Soc Mass Spectrom. 2014, 25(4): 581-7

26  T. Rosenling; M. P. Stoop; A. Smolinska; B. Muilwijk; L. Coulier; S. Shi; A. Dane; C. Christin; F. Suits; P. L. Horvatovich; S. S. Wijmenga; L. M. Buydens; R. Vreeken; T. Hankemeier; A. J. van Gool; T. M. Luider; R. Bischoff. The impact of delayed storage on the measured proteome and metabolome of human cerebrospinal fluid. Clin Chem. 2011, 57(12): 1703-11

27  R. Di Guida; J. Engel; J. W. Allwood; R. J. Weber; M. R. Jones; U. Sommer; M. R. Viant; W. B. Dunn. Non-targeted UHPLC-MS metabolomic data processing methods: a comparative investigation of normalisation, missing value imputation, transformation and scaling. Metabolomics. 2016, 1293

28  T. Valikangas; T. Suomi; L. L. Elo. A systematic evaluation of normalization methods in quantitative label-free proteomics. Brief Bioinform. 2016,

29  Y. V. Karpievitch; T. Taverner; J. N. Adkins; S. J. Callister; G. A. Anderson; R. D. Smith; A. R. Dabney. Normalization of peak intensities in bottom-up MS-based proteomics using singular value decomposition. Bioinformatics. 2009, 25(19): 2573-80

30  K. Kultima; A. Nilsson; B. Scholz; U. L. Rossbach; M. Falth; P. E. Andren. Development and evaluation of normalization methods for label-free relative quantification of endogenous peptides. Mol Cell Proteomics. 2009, 8(10): 2285-95

31  M. M. Matzke; K. M. Waters; T. O. Metz; J. M. Jacobs; A. C. Sims; R. S. Baric; J. G. Pounds; B. J. Webb-Robertson. Improved quality control processing of peptide-centric LC-MS proteomics data. Bioinformatics. 2011, 27(20): 2866-72

32  T. B. Bennike; K. Kastaniegaard; S. Padurariu; M. Gaihede; S. Birkelund; V. Andersen; A. Stensballe. Proteome stability analysis of snap frozen, RNAlater preserved, and formalin-fixed paraffin-embedded human colon mucosal biopsies. Data Brief. 2016, 6942-7

33  N. Selevsek; C. Y. Chang; L. C. Gillet; P. Navarro; O. M. Bernhardt; L. Reiter; L. Y. Cheng; O. Vitek; R. Aebersold. Reproducible and consistent quantification of the Saccharomyces cerevisiae proteome by SWATH-mass spectrometry. Mol Cell Proteomics. 2015, 14(3): 739-49

34  B. S. Kato; G. Nicholson; M. Neiman; M. Rantalainen; C. C. Holmes; A. Barrett; M. Uhlen; P. Nilsson; T. D. Spector; J. M. Schwenk. Variance decomposition of protein profiles from antibody arrays using a longitudinal twin model. Proteome Sci. 2011, 973

35  T. Guo; P. Kouvonen; C. C. Koh; L. C. Gillet; W. E. Wolski; H. L. Rost; G. Rosenberger; B. C. Collins; L. C. Blum; S. Gillessen; M. Joerger; W. Jochum; R. Aebersold. Rapid mass spectrometric conversion of tissue biopsy samples into permanent quantitative digital proteome maps. Nat Med. 2015, 21(4): 407-13

36  M. G. Hong; W. Lee; P. Nilsson; Y. Pawitan; J. M. Schwenk. Multidimensional Normalization to Minimize Plate Effects of Suspension Bead Array Data. J Proteome Res. 2016, 15(10): 3473-80

37  Z. Liu; Z. Yuan; Q. Zhao. SELDI-TOF-MS proteomic profiling of serum, urine, and amniotic fluid in neural tube defects. PLoS One. 2014, 9(7): e103276

38  O. E. Branson; M. A. Freitas. A multi-model statistical approach for proteomic spectral count quantitation. J Proteomics. 2016, 14423-32

39  R. Pedreschi; M. L. Hertog; S. C. Carpentier; J. Lammertyn; J. Robben; J. P. Noben; B. Panis; R. Swennen; B. M. Nicolai. Treatment of missing values for multivariate statistical analysis of gel-based proteomics data. Proteomics. 2008, 8(7): 1371-83

40  F. Koopmans; L. N. Cornelisse; T. Heskes; T. M. Dijkstra. Empirical Bayesian random censoring threshold model improves detection of differentially abundant proteins. J Proteome Res. 2014, 13(9): 3871-80

41  D. Lin; J. Zhang; J. Li; C. Xu; H. W. Deng; Y. P. Wang. An integrative imputation method based on multi-omics datasets. BMC Bioinformatics. 2016, 17247

42  W. S. Wu; M. J. Jhou. MVIAeval: a web tool for comprehensively evaluating the performance of a new missing value imputation algorithm. BMC Bioinformatics. 2017, 18(1): 31

43  L. Gatto; K. S. Lilley. MSnbase-an R/Bioconductor package for isobaric tagged mass spectrometry data visualization, processing and quantitation. Bioinformatics. 2012, 28(2): 288-9

44  C. C. Tsou; D. Avtonomov; B. Larsen; M. Tucholska; H. Choi; A. C. Gingras; A. I. Nesvizhskii. DIA-Umpire: comprehensive computational framework for data-independent acquisition proteomics. Nat Methods. 2015, 12(3): 258-64, 7 p following 64

45  C. C. Tsou; C. F. Tsai; G. C. Teo; Y. J. Chen; A. I. Nesvizhskii. Untargeted, spectral library-free analysis of data-independent acquisition proteomics data generated using Orbitrap mass spectrometers. Proteomics. 2016, 16(15-16): 2257-71

46  R. Bruderer; O. M. Bernhardt; T. Gandhi; S. M. Miladinovic; L. Y. Cheng; S. Messner; T. Ehrenberger; V. Zanotelli; Y. Butscheid; C. Escher; O. Vitek; O. Rinner; L. Reiter. Extending the limits of quantitative proteome profiling with data-independent acquisition and application to acetaminophen-treated three-dimensional liver microtissues. Mol Cell Proteomics. 2015, 14(5): 1400-10

47  L. Wu; S. Amon; H. Lam. A hybrid retention time alignment algorithm for SWATH-MS data. Proteomics. 2016, 16(15-16): 2272-83

48  H. L. Rost; G. Rosenberger; P. Navarro; L. Gillet; S. M. Miladinovic; O. T. Schubert; W. Wolski; B. C. Collins; J. Malmstrom; L. Malmstrom; R. Aebersold. OpenSWATH enables automated, targeted analysis of data-independent acquisition MS data. Nat Biotechnol. 2014, 32(3): 219-23

49  H. L. Rost; Y. Liu; G. D'Agostino; M. Zanella; P. Navarro; G. Rosenberger; B. C. Collins; L. Gillet; G. Testa; L. Malmstrom; R. Aebersold. TRIC: an automated alignment strategy for reproducible protein quantification in targeted proteomics. Nat Methods. 2016, 13(9): 777-83

50  G. Rosenberger; Y. Liu; H. L. Rost; C. Ludwig; A. Buil; A. Bensimon; M. Soste; T. D. Spector; E. T. Dermitzakis; B. C. Collins; L. Malmstrom; R. Aebersold. Inference and quantification of peptidoforms in large sample cohorts by SWATH-MS. Nat Biotechnol. 2017, 35(8): 781-88

51  J. X. Wu; X. Song; D. Pascovici; T. Zaw; N. Care; C. Krisp; M. P. Molloy. SWATH Mass Spectrometry Performance Using Extended Peptide MS/MS Assay Libraries. Mol Cell Proteomics. 2016, 15(7): 2501-14

52  H. L. Rost; T. Sachsenberg; S. Aiche; C. Bielow; H. Weisser; F. Aicheler; S. Andreotti; H. C. Ehrlich; P. Gutenbrunner; E. Kenar; X. Liang; S. Nahnsen; L. Nilse; J. Pfeuffer; G. Rosenberger; M. Rurik; U. Schmitt; J. Veit; M. Walzer; D. Wojnar; W. E. Wolski; O. Schilling; J. S. Choudhary; L. Malmstrom;

R. Aebersold; K. Reinert; O. Kohlbacher. OpenMS: a flexible open-source software platform for mass spectrometry data analysis. Nat Methods. 2016, 13(9): 741-8

53  S. I. Anjo; C. Santa; B. Manadas. SWATH-MS as a tool for biomarker discovery: From basic research to clinical applications. Proteomics. 2017, 17(3-4):

54  Y. Liu; R. Huttenhain; S. Surinova; L. C. Gillet; J. Mouritsen; R. Brunner; P. Navarro; R. Aebersold. Quantitative measurements of N-linked glycoproteins in human plasma by SWATH-MS. Proteomics. 2013, 13(8): 1247-56

55  S. Shao; T. Guo; C. C. Koh; S. Gillessen; M. Joerger; W. Jochum; R. Aebersold. Minimal sample requirement for highly multiplexed protein quantification in cell lines and tissues by PCT-SWATH mass spectrometry. Proteomics. 2015, 15(21): 3711-21

56  B. MacLean; D. M. Tomazela; N. Shulman; M. Chambers; G. L. Finney; B. Frewen; R. Kern; D. L. Tabb; D. C. Liebler; M. J. MacCoss. Skyline: an open source document editor for creating and analyzing targeted proteomics experiments. Bioinformatics. 2010, 26(7): 966-8

57  B. Schilling; M. J. Rardin; B. X. MacLean; A. M. Zawadzka; B. E. Frewen; M. P. Cusack; D. J. Sorensen; M. S. Bereman; E. Jing; C. C. Wu; E. Verdin; C. R. Kahn; M. J. Maccoss; B. W. Gibson. Platform-independent and label-free quantitation of proteomic data using MS1 extracted ion chromatograms in skyline: application to protein acetylation and phosphorylation. Mol Cell Proteomics. 2012, 11(5): 202-14

58  C. Escher; L. Reiter; B. MacLean; R. Ossola; F. Herzog; J. Chilton; M. J. MacCoss; O. Rinner. Using iRT, a normalized retention time for more targeted measurement of peptides. Proteomics. 2012, 12(8): 1111-21

59  J. Kang; J. Lu; X. Zhang. Metabolomics-based promising candidate biomarkers and pathways in Alzheimer's disease. Pharmazie. 2015, 70(5): 277-82

60  R. Bruderer; O. M. Bernhardt; T. Gandhi; L. Reiter. High-precision iRT prediction in the targeted analysis of data-independent acquisition and its impact on identification and quantitation. Proteomics. 2016, 16(15-16): 2246-56

61  P. Navarro; J. Kuharev; L. C. Gillet; O. M. Bernhardt; B. MacLean; H. L. Rost; S. A. Tate; C. C. Tsou; L. Reiter; U. Distler; G. Rosenberger; Y. Perez-Riverol; A. I. Nesvizhskii; R. Aebersold; S. Tenzer. A multicenter study benchmarks software tools for label-free proteome quantification. Nat Biotechnol. 2016, 34(11): 1130-36

62  T. Valikangas; T. Suomi; L. L. Elo. A comprehensive evaluation of popular proteomics software workflows for label-free proteome quantification and imputation. Brief Bioinform. 2017,

63  J. Cox; M. Mann. MaxQuant enables high peptide identification rates, individualized p.p.b.-range mass accuracies and proteome-wide protein quantification. Nat Biotechnol. 2008, 26(12): 1367-72

64  S. Tyanova; T. Temu; J. Cox. The MaxQuant computational platform for mass spectrometry-based shotgun proteomics. Nat Protoc. 2016, 11(12): 2301-19

65  S. Tyanova; T. Temu; A. Carlson; P. Sinitcyn; M. Mann; J. Cox. Visualization of LC-MS/MS proteomics data in MaxQuant. Proteomics. 2015, 15(8): 1453-6

66  H. Weisser; S. Nahnsen; J. Grossmann; L. Nilse; A. Quandt; H. Brauer; M. Sturm; E. Kenar; O. Kohlbacher; R. Aebersold; L. Malmstrom. An automated pipeline for high-throughput label-free quantitative proteomics. J Proteome Res. 2013, 12(4): 1628-44

67  D. Bouyssie; A. Gonzalez de Peredo; E. Mouton; R. Albigot; L. Roussel; N. Ortega; C. Cayrol; O. Burlet-Schiltz; J. P. Girard; B. Monsarrat. Mascot file parsing and quantification (MFPaQ), a new software to parse, validate, and quantify proteomics data generated by ICAT and SILAC mass spectrometric analyses: application to the proteomics study of membrane proteins from primary human endothelial cells. Mol Cell Proteomics. 2007, 6(9): 1621-37

68  E. Hoedt; K. Chaoui; I. Huvent; C. Mariller; B. Monsarrat; O. Burlet-Schiltz; A. Pierce. SILAC-based proteomic profiling of the human MDA-MB-231 metastatic breast cancer cell line in response to the two antitumoral lactoferrin isoforms: the secreted lactoferrin and the intracellular delta-lactoferrin. PLoS One. 2014, 9(8): e104563

69  C. N. Schlaffner; G. J. Pirklbauer; A. Bender; J. S. Choudhary. Fast, Quantitative and Variant Enabled Mapping of Peptides to Genomes. Cell Syst. 2017, 5(2): 152-56 e4

70  P. E. Khoonsari; A. Haggmark; M. Lonnberg; M. Mikus; L. Kilander; L. Lannfelt; J. Bergquist; M. Ingelsson; P. Nilsson; K. Kultima; G. Shevchenko. Analysis of the Cerebrospinal Fluid Proteome in Alzheimer's Disease. PLoS One. 2016, 11(3): e0150672

71  A. J. Weston; W. C. Dunlap; J. M. Shick; A. Klueter; K. Iglic; A. Vukelic; A. Starcevic; M. Ward; M. L. Wells; C. G. Trick; P. F. Long. A profile of an endosymbiont-enriched fraction of the coral Stylophora pistillata reveals proteins relevant to microbial-host interactions. Mol Cell Proteomics. 2012, 11(6): M111 015487

72  J. Zhang; L. Xin; B. Shan; W. Chen; M. Xie; D. Yuen; W. Zhang; Z. Zhang; G. A. Lajoie; B. Ma. PEAKS DB: de novo sequencing assisted database search for sensitive and accurate peptide identification. Mol Cell Proteomics. 2012, 11(4): M111 010587

73  J. Zhang; W. Yang; S. Li; S. Yao; P. Qi; Z. Yang; Z. Feng; J. Hou; L. Cai; M. Yang; W. Wu; D. A. Guo. An intelligentized strategy for endogenous small molecules characterization and quality evaluation of earthworm from two geographic origins by ultra-high performance HILIC/QTOF MS(E) and Progenesis QI. Anal Bioanal Chem. 2016, 408(14): 3881-90

74  M. R. Al Shweiki; S. Monchgesang; P. Majovsky; D. Thieme; D. Trutschel; W. Hoehenwarter. Assessment of Label-Free Quantification in Discovery Proteomics and Impact of Technological Factors and Natural Variability of Protein Abundance. J Proteome Res. 2017, 16(4): 1410-24

75  A. M. Almeida; P. Nanni; A. M. Ferreira; C. Fortes; J. Grossmann; R. J. B. Bessa; P. Costa. The longissimus thoracis muscle proteome in Alentejana bulls as affected by growth path. J Proteomics. 2017, 152206-15

76  P. Garden; R. Alm; J. Hakkinen. PROTEIOS: an open source proteomics initiative. Bioinformatics. 2005, 21(9): 2085-7

77  F. Levander; M. Krogh; K. Warell; P. Garden; P. James; J. Hakkinen. Automated reporting from gel-based proteomics experiments using the open source Proteios database application. Proteomics. 2007, 7(5): 668-74

78 M. C. Codrea; S. Nahnsen. Platforms and Pipelines for Proteomics Data Analysis and Management. Adv Exp Med Biol. 2016, 919203-15

79 J. M. Twigt; K. Bezstarosti; J. Demmers; J. Lindemans; J. S. Laven; R. P. Steegers-Theunissen. Preconception folic acid use influences the follicle fluid proteome. Eur J Clin Invest. 2015, 45(8): 833-41

80 J. Dorts; P. Kestemont; M. L. Thezenas; M. Raes; F. Silvestre. Effects of cadmium exposure on the gill proteome of Cottus gobio: modulatory effects of prior thermal acclimation. Aquat Toxicol. 2014, 15487-96

81 N. Colaert; H. Barsnes; M. Vaudel; K. Helsens; E. Timmerman; A. Sickmann; K. Gevaert; L. Martens. Thermo-msf-parser: an open source Java library to parse and visualize Thermo Proteome Discoverer msf files. J Proteome Res. 2011, 10(8): 3840-3

82 J. Veit; T. Sachsenberg; A. Chernev; F. Aicheler; H. Urlaub; O. Kohlbacher. LFQProfiler and RNP(xl): Open-Source Tools for Label-Free Quantification and Protein-RNA Cross-Linking Integrated into Proteome Discoverer. J Proteome Res. 2016, 15(9): 3441-8

83 J. Casado-Vela; M. J. Martinez-Esteso; E. Rodriguez; E. Borras; F. Elortza; R. Bru-Martinez. iTRAQ-based quantitative analysis of protein mixtures with large fold change and dynamic range. Proteomics. 2010, 10(2): 343-7

84 Y. Y. Chen; S. Dasari; Z. Q. Ma; L. J. Vega-Montoto; M. Li; D. L. Tabb. Refining comparative proteomics by spectral counting to account for shared peptides and multiple search engines. Anal Bioanal Chem. 2012, 404(4): 1115-25

85 S. K. Park; J. R. Yates, 3rd. Census for proteome quantification. Curr Protoc Bioinformatics. 2010, Chapter 13Unit 13 12 1-11

86 S. K. Park; J. D. Venable; T. Xu; J. R. Yates, 3rd. A quantitative analysis software tool for mass spectrometry-based proteomics. Nat Methods. 2008, 5(4): 319-22

87 J. H. Prieto; E. Fischer; S. Koncarevic; J. Yates; K. Becker. Large-scale differential proteome analysis in Plasmodium falciparum under drug treatment. Methods Mol Biol. 2015, 1201269-79

88 D. Cociorva; L. T. D; J. R. Yates. Validation of tandem mass spectrometry database search results using DTASelect. Curr Protoc Bioinformatics. 2007, Chapter 13Unit 13 4

89 D. L. Tabb; W. H. McDonald; J. R. Yates, 3rd. DTASelect and Contrast: tools for assembling and comparing protein identifications from shotgun proteomics. J Proteome Res. 2002, 1(1): 21-6

90 G. W. Park; H. Hwang; K. H. Kim; J. Y. Lee; H. K. Lee; J. Y. Park; E. S. Ji; S. R. Park; J. R. Yates, 3rd; K. H. Kwon; Y. M. Park; H. J. Lee; Y. K. Paik; J. Y. Kim; J. S. Yoo. Integrated Proteomic Pipeline Using Multiple Search Engines for a Proteogenomic Study with a Controlled Protein False Discovery Rate. J Proteome Res. 2016, 15(11): 4082-90

91 J. Wan; A. F. Roth; A. O. Bailey; N. G. Davis. Palmitoylated proteins: purification and identification. Nat Protoc. 2007, 2(7): 1573-84

92 V. Dupierris; C. Masselon; M. Court; S. Kieffer-Jaquinod; C. Bruley. A toolbox for validation of mass spectrometry peptides identification and generation of database: IRMa. Bioinformatics. 2009, 25(15): 1980-1

93  E. Mouton-Barbosa; F. Roux-Dalvai; D. Bouyssie; F. Berger; E. Schmidt; P. G. Righetti; L. Guerrier; E. Boschetti; O. Burlet-Schiltz; B. Monsarrat; A. Gonzalez de Peredo. In-depth exploration of cerebrospinal fluid by combining peptide ligand library treatment and label-free protein quantification. Mol Cell Proteomics. 2010, 9(5): 1006-21

94  A. I. Nesvizhskii; A. Keller; E. Kolker; R. Aebersold. A statistical model for identifying proteins by tandem mass spectrometry. Anal Chem. 2003, 75(17): 4646-58

95  A. Keller; A. I. Nesvizhskii; E. Kolker; R. Aebersold. Empirical statistical model to estimate the accuracy of peptide identifications made by MS/MS and database search. Anal Chem. 2002, 74(20): 5383-92

96  R. M. Sakia. The Box-Cox Transformation Technique - a Review. Journal Of the Royal Statistical Society Series D-the Statistician. 1992, 41(2): 169-78

97  A. M. De Livera; D. A. Dias; D. De Souza; T. Rupasinghe; J. Pyke; D. Tull; U. Roessner; M. McConville; T. P. Speed. Normalizing and integrating metabolomics data. Anal Chem. 2012, 84(24): 10768-76

98  W. Huber; A. von Heydebreck; H. Sultmann; A. Poustka; M. Vingron. Variance stabilization applied to microarray data calibration and to the quantification of differential expression. Bioinformatics. 2002, 18 Suppl 1S96-104

99  S. M. Kohl; M. S. Klein; J. Hochrein; P. J. Oefner; R. Spang; W. Gronwald. State-of-the art data normalization methods improve NMR-based metabolomic analysis. Metabolomics. 2012, 8(Suppl 1): 146-60

100 R. A. van den Berg; H. C. Hoefsloot; J. A. Westerhuis; A. K. Smilde; M. J. van der Werf. Centering, scaling, and transformations: improving the biological information content of metabolomics data. BMC Genomics. 2006, 7142

101 P. S. Gromski; Y. Xu; K. A. Hollywood; M. L. Turner; R. Goodacre. The influence of scaling metabolomics data on model classification accuracy. Metabolomics. 2015, 11(3): 684-95

102 M. R. Nezami Ranjbar; Y. Zhao; M. G. Tadesse; Y. Wang; H. W. Ressom. Gaussian process regression model for normalization of LC-MS data using scan-level information. Proteome Sci. 2013, 11(Suppl 1): S13

103 T. F. Bathen; B. Geurts; B. Sitter; H. E. Fjosne; S. Lundgren; L. M. Buydens; I. S. Gribbestad; G. Postma; G. F. Giskeodegard. Feasibility of MR metabolomics for immediate analysis of resection margins during breast cancer surgery. PLoS One. 2013, 8(4): e61578

104 A. K. Smilde; M. J. van der Werf; S. Bijlsma; B. J. van der Werff-van der Vat; R. H. Jellema. Fusion of mass spectrometry-based metabolomics data. Anal Chem. 2005, 77(20): 6729-36

105 H. Parastar; A. Bazrafshan. Fuzzy C-means clustering for chromatographic fingerprints analysis: A gas chromatography-mass spectrometry case study. J Chromatogr A. 2016, 1438236-43

106 B. J. Webb-Robertson; Y. M. Kim; E. M. Zink; K. A. Hallaian; Q. Zhang; R. Madupu; K. M. Waters; T. O. Metz. A Statistical Analysis of the Effects of Urease Pre-treatment on the Measurement of the Urinary Metabolome by Gas Chromatography-Mass Spectrometry. Metabolomics. 2014, 10(5): 897-908

107 K. V. Ballman; D. E. Grill; A. L. Oberg; T. M. Therneau. Faster cyclic loess: normalizing RNA arrays via linear models. Bioinformatics. 2004, 20(16): 2778-86

108 A. J. Keeping; R. A. Collins. Data variance and statistical significance in 2D-gel electrophoresis and DIGE experiments: comparison of the effects of normalization methods. J Proteome Res. 2011, 10(3): 1353-60

109 Y. V. Karpievitch; S. B. Nikolic; R. Wilson; J. E. Sharman; L. M. Edwards. Metabolomics data normalization with EigenMS. PLoS One. 2014, 9(12): e116221

110 Y. V. Karpievitch; A. R. Dabney; R. D. Smith. Normalization and missing value imputation for label-free LC-MS analysis. BMC Bioinformatics. 2012, 13 Suppl 16S5

111 Y. V. Karpievitch; A. D. Polpitiya; G. A. Anderson; R. D. Smith; A. R. Dabney. Liquid Chromatography Mass Spectrometry-Based Proteomics: Biological and Technological Aspects. Ann Appl Stat. 2010, 4(4): 1797-823

112 S. J. Callister; R. C. Barry; J. N. Adkins; E. T. Johnson; W. J. Qian; B. J. Webb-Robertson; R. D. Smith; M. S. Lipton. Normalization approaches for removing systematic biases associated with mass spectrometry and label-free proteomics. J Proteome Res. 2006, 5(2): 277-86

113 V. Andjelkovic; R. Thompson. Changes in gene expression in maize kernel in response to water and salt stress. Plant Cell Rep. 2006, 25(1): 71-9

114 B. A. Ejigu; D. Valkenborg; G. Baggerman; M. Vanaerschot; E. Witters; J. C. Dujardin; T. Burzykowski; M. Berg. Evaluation of normalization methods to pave the way towards large-scale LC-MS-based metabolomics profiling experiments. OMICS. 2013, 17(9): 473-85

115 A. Padoan; D. Basso; M. La Malfa; C. F. Zambon; P. Aiyetan; H. Zhang; A. Di Chiara; G. Pavanello; R. Bellocco; D. W. Chan; M. Plebani. Reproducibility in urine peptidome profiling using MALDI-TOF. Proteomics. 2015, 15(9): 1476-85

116 B. M. Bolstad; R. A. Irizarry; M. Astrand; T. P. Speed. A comparison of normalization methods for high density oligonucleotide array data based on variance and bias. Bioinformatics. 2003, 19(2): 185-93

117 F. Dieterle; A. Ross; G. Schlotterbeck; H. Senn. Probabilistic quotient normalization as robust method to account for dilution of complex biological mixtures. Application in $^1$H NMR metabonomics. Anal. Chem. 2006, 78(13): 4281-90

118 H. Lopez-Fernandez; H. M. Santos; J. L. Capelo; F. Fdez-Riverola; D. Glez-Pena; M. Reboiro-Jato. Mass-Up: an all-in-one open software application for MALDI-TOF mass spectrometry knowledge discovery. BMC Bioinformatics. 2015, 16318

119 M. Gaspari; L. Chiesa; A. Nicastri; C. Gabriele; V. Harper; D. Britti; G. Cuda; A. Procopio. Proteome Speciation by Mass Spectrometry: Characterization of Composite Protein Mixtures in Milk Replacers. Anal Chem. 2016, 88(23): 11568-74

120 S. P. Borgaonkar; H. Hocker; H. Shin; M. K. Markey. Comparison of normalization methods for the identification of biomarkers using MALDI-TOF and SELDI-TOF mass spectra. OMICS. 2010, 14(1): 115-26

121 Y. Lin; K. Golovnina; Z. X. Chen; H. N. Lee; Y. L. Negron; H. Sultana; B. Oliver; S. T. Harbison. Comparison of normalization and differential expression analyses using RNA-Seq data from 726 individual Drosophila melanogaster. BMC Genomics. 2016, 1728

122 L. E. Chai; C. K. Law; M. S. Mohamad; C. K. Chong; Y. W. Choon; S. Deris; R. M. Illias. Investigating the effects of imputation methods for modelling gene networks using a dynamic bayesian network from gene expression data. Malays J Med Sci. 2014, 21(2): 20-7

123 H. Kim; G. H. Golub; H. Park. Missing value estimation for DNA microarray gene expression data: local least squares imputation. Bioinformatics. 2005, 21(2): 187-98

124 O. Alter; P. O. Brown; D. Botstein. Singular value decomposition for genome-wide expression data processing and modeling. Proc Natl Acad Sci U S A. 2000, 97(18): 10101-6

125 X. Gan; A. W. Liew; H. Yan. Microarray missing data imputation based on a set theoretic framework and biological knowledge. Nucleic Acids Res. 2006, 34(5): 1608-19

126 A. Chawade; E. Alexandersson; F. Levander. Normalyzer: a tool for rapid evaluation of normalization methods for omics data sets. J Proteome Res. 2014, 13(6): 3114-20

127 N. M. Griffin; J. Yu; F. Long; P. Oh; S. Shore; Y. Li; J. A. Koziol; J. E. Schnitzer. Label-free, normalized quantification of complex mass spectrometry data for proteomic analysis. Nat Biotechnol. 2010, 28(1): 83-9

128 D. Risso; J. Ngai; T. P. Speed; S. Dudoit. Normalization of RNA-seq data using factor analysis of control genes or samples. Nat Biotechnol. 2014, 32(9): 896-902

129 B. J. Blaise. Data-driven sample size determination for metabolic phenotyping studies. Anal Chem. 2013, 85(19): 8943-50

130 A. Pursiheimo; A. P. Vehmas; S. Afzal; T. Suomi; T. Chand; L. Strauss; M. Poutanen; A. Rokka; G. L. Corthals; L. L. Elo. Optimization of Statistical Methods Impact on Quantitative Proteomics Data. J Proteome Res. 2015, 14(10): 4118-26

131 V. B. Dubinkina; A. V. Tyakht; V. Y. Odintsova; K. S. Yarygin; B. A. Kovarsky; A. V. Pavlenko; D. S. Ischenko; A. S. Popenko; D. G. Alexeev; A. Y. Taraskina; R. F. Nasyrova; E. M. Krupitsky; N. V. Shalikiani; I. G. Bakulin; P. L. Shcherbakov; L. O. Skorodumova; A. K. Larin; E. S. Kostryukova; R. A. Abdulkhakov; S. R. Abdulkhakov; S. Y. Malanin; R. K. Ismagilova; T. V. Grigoryeva; E. N. Ilina; V. M. Govorun. Links of gut microbiota composition with alcohol dependence syndrome and alcoholic liver disease. Microbiome. 2017, 5(1): 141

132 X. Wang; E. J. Gardiner; M. J. Cairns. Optimal consistency in microRNA expression analysis using reference-gene-based normalization. Mol Biosyst. 2015, 11(5): 1235-40

133 J. Kuharev; P. Navarro; U. Distler; O. Jahn; S. Tenzer. In-depth evaluation of software tools for data-independent acquisition based label-free quantification. Proteomics. 2015, 15(18): 3140-51