# Optimization of Spectral Library Size

# Improves DIA-MS Proteome Coverage

Weigang Ge [1,2,3,4 #], Xiao Liang [1,2,3 #], Fangfei Zhang [1,2,3 #], Luang Xu [1,2,3], Nan Xiang [1,2,3], Rui Sun [1,2,3], Wei Liu [1,2,3] Zhangzhi Xue [1,2,3], Xiao Yi [1,2,3], Bo Wang [5], Jiang Zhu [6], Cong Lu [6], Xiaolu Zhan [7], Lirong Chen [8], Yan Wu [9,10], Zhiguo Zheng [11,12], Wangang Gong [11,12], Qijun Wu [13], Jiekai Yu [14], Zhaoming Ye [9,10], Xiaodong Teng [5], Shiang Huang [6], Shu Zheng [14], Tong Liu [7*], Chunhui Yuan [1,2,3*], Tiannan Guo [1,2,3*]

1, Key Laboratory of Structural Biology of Zhejiang Province, School of Life Sciences, Westlake University, Hangzhou, Zhejiang, China.

2, Westlake Laboratory of Life Sciences and Biomedicine, Hangzhou, Zhejiang, China

3, Institute of Basic Medical Sciences, Westlake Institute for Advanced Study, Hangzhou, Zhejiang, China

4, Westlake Omics (Hangzhou) Biotechnology Co.,Ltd. No.1, Yunmeng Road, Cloud Town, Xihu District, Hangzhou 310024, Zhejiang Province, China

5, Department of Pathology, The First Affiliated Hospital of College of Medicine, Zhejiang University, Hangzhou, Zhejiang Province, China

6, Center for Stem Cell Research and Application, Union Hospital, Tongji Medical College, Huazhong University of Science and Technology, Wuhan, Hubei Province, China

7, Harbin Medical University Cancer Hospital, Harbin 150081, China

8, Department of Pathology, The Second Affiliated Hospital of College of Medicine, Zhejiang University, Hangzhou 310009, Zhejiang Province, China

9, Department of Orthopaedics, The Second Affiliated Hospital of Zhejiang University School of Medicine, Hangzhou 310009, Zhejiang Province, China

10, Key Laboratory of Motor System Disease Research and Precision Therapy of Zhejiang Province, Hangzhou 310020, Zhejiang Province, China

11, The Cancer Hospital of the University of Chinese Academy of Sciences, Zhejiang Cancer Hospital, Hangzhou 310022, Zhejiang Province, China

12, Institute of Basic Medicine and Cancer (IBMC), Chinese Academy of Sciences, Hangzhou 310022, Zhejiang Province, China

13, Department of Clinical Epidemiology, Shengjing Hospital of China Medical University, Shenyang Province, China

14, Cancer Institute, Key Laboratory of Cancer Prevention and Intervention, Ministry of Education, The Second Affiliated Hospital, Zhejiang University School of Medicine, Hangzhou, Zhejiang, China

#, co-first authors

*, corresponding author. Tong Liu: liutong@hrbmu.edu.cn; Chunhui Yuan: yuanchunhui@westlake.edu.cn; Tiannan Guo: guotiannan@westlake.edu.cn

Running title: Optimization of Spectral Library Size for DIA

Keywords: Data-independent acquisition; protein identification; Pan-human library; Spectral library optimization.

2

44    **Abstract**

45        Efficient peptide and protein identification from data-independent acquisition  mass

46    spectrometric (DIA-MS) data typically rely on an experiment-specific spectral library with a

47    suitable size. Here, we report a computational strategy for optimizing the spectral library for a

48    specific DIA dataset based on a comprehensive spectral library, which is accomplished by *a*

49    *priori* analysis of the DIA dataset. This strategy achieved up to 44.7% increase in peptide

50    identification and 38.1% increase in protein identification in the test dataset of six colorectal

51    tumor samples compared with the comprehensive pan-human library strategy. We further applied

52    this strategy to 389 carcinoma samples from 15 tumor datasets and observed up to 39.2%

53    increase in peptide identification and 19.0% increase in protein identification. In summary, we

54    present a computational strategy for spectral library size optimization to achieve deeper

55    proteome coverage of DIA-MS data.

56

57    **Introduction**

58        Data-independent acquisition mass spectrometry (DIA-MS) based proteomics coupled with

59    targeted data analysis is playing an increasing role in biomedical studies (1), owing to its high

60    degree of reproducibility, quantitative accuracy, and high throughput (2, 3). Both spectral

61    library-free and library-based strategies are being applied to analyze DIA-MS data (4). While the

62    library-free strategies (5, 6) could identify peptides directly from DIA-MS itself without the

63    requirement of an external spectral library, the depth of proteomic coverage is limited at the

64    moment (7-9). The more widely adopted strategy is based on building a spectral library using the

65  corresponding data-dependent acquisition mass spectrometry (DDA-MS) datasets of the samples

66  of interest (10), or a pre-built library from public data repositories (11-14).

67      The size of the spectral library has a direct impact on the performance of DIA-MS data

68  analysis (15). A larger number of DDA-MS runs, particularly from fractionated samples, leads to

69  a more comprehensive spectral library enabling potential detection of a larger number of

70  peptides and proteins from the DIA-MS datasets (15). However, it also generates a larger search

71  space and reduces the statistical power to detect true positives (16, 17). Extra concerns are raised

72  where the proteins and peptides within the library may not be specific to a particular specimen,

73  potentially introducing more false positives (18). Other drawbacks include the prolonged

74  computational time which is approximately linearly correlated with the size of the library (19),

75  and distortion of retention time (RT) distribution for alignment (20).

76      The spectral library size could be optimized to improve DIA-MS performance. The Van

77  Eyk group have reported that applying a comprehensive fractionated library led to higher number

78  of protein/peptide identifications from DIA-MS datasets than un-fractionated libraries with

79  limited sizes (15). Similar results have been reported by Uszkoreit group, where they found

80  larger library led to higher peptide and protein identification but the increase was minimal when

81  the library is comprehensive enough (21). The combination of an in-house built library with

82  external libraries from public data improves DIA data analysis performance (17). Inclusion of

83  internal library extracted from DIA files also improved peptide and protein identification (9). On

84  the other hand, it has also been observed that libraries of very large size led to higher FDRs in

85  the DIA-MS analyses and hence compromises the identification results (17). It was further

86  demonstrated that, even within the same spectral library, controlling the confidence of peptide

87  identifications to exclude redundant peptides could improve peptide and protein identification

4

88    results (16). Although these studies have repetively reported the importance of the size of

89    spectral library size, a systematic evaluation and optimization of library size is still lacking.

90         Here, we propose a two-step strategy called subLib to generate the experiment-specific

91    subset libraries using *a priori* analysis of the DIA data to improve the proteomic coverage. The

92    strategy to derive a subset library of optimal size was further applied to analyze the DIA data of

93    15 human tumors.

94

## Materials and Methods

### Colorectal cancer dataset

97         To evaluate our strategy, the DIA-MS datasets were collected from a colorectal cancer

98    proteomic project in our group (Xiang *et al.*, manuscript in preparation). Briefly, 286 FFPE

99    samples from 44 colorectal cancer patients were processed into peptides with a pressure cycling

100   technology (PCT)-based protocol as described in the previous study (22). They were subjected to

101   data acquisition on the nanoflow EASY-nLC™ 1200 System coupled with Q Exactive HF

102   hybrid Quadrupole-Orbitrap in DIA mode over a gradient of 60 min using 24 DIA windows

103   spanning from 400 Da to 1200 Da.

### Fifteen datasets of multiple tumor types

105        A total of 389 tumor tissue samples from 15 tumor types were collected. The gastric

106   carcinoma (n=30) and thyroid carcinoma (n=30) samples were collected from the First Affiliated

107   Hospital College of Medicine, Zhejiang University. The prostate carcinoma (n=30) and bone

108   carcinoma (n=30) samples were collected from the Second Affiliated Hospital College of

109     Medicine, Zhejiang University. The liver carcinoma (n=33) and leukemia (n=27) samples were

110     collected from Wuhan Union Hospital. The ovarian carcinoma (n=30) samples were collected

111     from Zhejiang Cancer Hospital. The cervical carcinoma (n=28) samples were collected from

112     Shengjing Hospital of China Medical University. The lung adenocarcinoma (n=32) , gallbladder

113     carcinoma (n=20), pancreatic adenocarcinoma (n=20), myosarcoma (n=19), clear cell renal cell

114     carcinoma (CCRCC, n=20), diffuse large B-cell lymphoma (DLBCL, n=19), and papillary

115     thyroid cancer (PTC, n=21) were collected from Harbin Medical University Cancer Hospital. All

116     samples were approved by the ethics committees of their respective hospitals. The tissue samples

117     were prepared with PCT-based tissue lysis and protein digestion protocol (22) and analyzed by

118     DIA-MS, as listed in Table S1. Ethics approvals for this study were obtained from the Ethics

119     Committee or Institutional Review.

120     **Proteomic data analysis workflow**

121     The raw DIA-MS data files were converted to mzXML format using the msConvert tool in

122     ProteomeWizard (23). The DIA-MS datasets were analyzed using the open-source software

123     OpenSWATH (version 2.4.0) (24) with the following criteria: common internal reference

124     peptides (CiRTs) of each tissue were applied respectively for retention time alignment; m/z

125     extraction window was set to 30 ppm, and RT extraction window was set between 200-800

126     seconds, depending on different gradients of the DIA-MS module (Table S1). PyProphet (version

127     2.1.3) (24) was used for statistical validation via setting the global cutoff of FDR as 0.01 at both

128     peptide and protein levels. Protein inference was performed as described previously (25). Unless

129     otherwise mentioned, the software parameters were kept the same for all the analyses in this

130     study.

131     **Subset library generation**

6

132     We proposed a two-step strategy to take a subset of the spectral library. Firstly, the public

133     library is taken to analyze the candidate DIA-MS dataset using the OpenSWATH workflow.

134     Different FDR cutoffs were set to generate a list of identification results. Afterwards, they were

135     matched against the public library to generate experiment-specific subset libraries.

136     In this study, we set the DIA Pan-Human Library (DPHL) (12) as the baseline library to

137     analyze the colorectal cancer dataset containing 284 DIA-MS data files. FDR cutoffs were set at

138     0.01, 0.02, 0.03, 0.04, 0.05, 0.06, 0.07, 0.08, 0.09, 0.1, 0.2, 0.3, 0.4, 0.5 and 0.6 (n=15), to

139     generate 15 identification results. After matching with DPHL, OpenSwathDecoyGenerator.exe in

140     OpenMS (version2.4.0) was applied to generate equal amount of decoys in mutated fashion. The

141     resultant subset library is a combination of DPHL subsets and decoys.

142

143     **Results and Discussions**

144     **Generation of the subset library by refining DPHL**

145     For data comprehensiveness and accessibility, DPHL built from 16 human tissue types

146     containing 359,627 peptide precursors and 14,782 protein groups was used as the baseline

147     spectral library. A DIA-MS dataset of 286 colorectal cancer sample cohort was analyzed to

148     derive the initial identifications. We set the FDR cut-off for peptide precursor and protein

149     identification to 0.01, 0.02, 0.03, 0.04, 0.05, 0.06, 0.07, 0.08, 0.09, 0.1, 0.2, 0.3, 0.4, 0.5 and 0.6

150     (a total of 15 tests), then retrieved the resultant subset libraries at each FDR cutoff. The four

151     representative DIA-MS data files (sample A1-A4) within the cohort and two external colorectal

152     cancer DIA-MS data files (sample B1 and B2) were taken to evaluate the identification

153     performance of each subset library (Figure 1A). The number of identified peptides shows a

7

154    generally decreasing trend as the FDR cutoff increases (Figure 1C), with the exceptions when

155    FDR increases from 0.01 to 0,02, and from 0.04 to 0.05. The number of identified proteins

156    increased as the FDR cut-off increased from 0.01 to 0.05, and gradually decreased afterward,

157    with a drastic decline when the cutoff was beyond 0.1 (Figure 1D). This is not unexpected since

158    the peptides identified with high FDR are more likely absent in the sample at the detection limit.

159    As the library size increased,  the negative effect prevailed. The best result was obtained from the

160    library with a FDR cutoff of 0.05. The optimal library was composed of 85,655 peptide

161    precursors, 62,390 peptides, and 6,448 protein groups, leading to the identification of 29,979

162    peptide precursors and 4,418 protein groups, respectively. This optimized library led to 44.7%

163    and 38.1% increase of peptide precursors and protein groups, respectively, compared with the

164    results by the unfiltered DPHL (Figure S1). The subset library with the FDR cutoff of 0.05 was

165    the best subset library which was hence adopted for further evaluation. The DIA files used for

166    library size optimization from samples A1-A4 led to similar data to those from independent

167    samples (B1 and B2), suggesting that the library size optimization is generic and applicable to

168    DIA files of the same tissue type.

169    **Adding unidentified peptide procursors to the subset library sacrificed identification**

170        To check if unidentified peptide precursors in a spectral library would affect the DIA-MS

171    proteome coverage, we randomly generated nine sets of DPHL peptides that were excluded from

172    the subset library (defined as "unidentified peptides"), with precursor number equivalent to n%

173    of the subset library (n=10, 20, …, 90), and combined them with the subset library peptides

174    (Figure 1B). When applying the reconstructed spectral libraries to analyze the test DIA dataset, a

175    steady decrease of identified peptides and proteins was observed as more unidentified precursors

176    were included (Figure 1E, F), with the highest proteome coverage coming from the library with

177    no unidentified peptides, summing up to 29,712 peptide precursors and 4,433 protein groups.

178        We also replaced the unidentified peptides to *in silico* generated decoy peptides and

179    repeated the above analyses. Peptide/protein identifications decreased as the computational

180    peptide proportions increase from 0% to 60%. Further addition of decoys would, however, subtly

181    increase protein identifications (Figure 1G, H). The highest proteome coverage came from the

182    library with no decoy interferences, summing up to 19,322 peptide precursors and 3,461 protein

183    groups. We hence concluded that any false positive interference in the library would suppress the

184    peptide/protein identification.

**Adding subset library peptides to interferences improves identification**

186        We then conducted a backward analysis by adding increasing proportions of subset library

187    peptides to the unidentified peptides (Figure 1B). The spectral library composed by precursors of

188    unidentified peptides solely (n=0) could not identify any peptide or protein in the DIA-MS data.

189    The numbers of identification of peptides and proteins exhibited almost marked increase as n

190    increased (Figure 1I, K). Together with the above results, they validated the effectiveness of

191    setting FDR cutoff as 0.05 to eliminate false positive targets.

**Applying subLib to DIA-MS of 15 tumor sample types**

193        We named the library generation strategy "subLib" and further applied it to the fifteen DIA

194    datasets of different types of cancer samples, including bone, cervical, DLBCL, gallbladder,

195    gastric, leukemia, liver, lung, myosarcoma, ovarian, pancreatic, prostate, PTC, and CCRCC

196    (Figure 2A). Peptide/protein identifications using the subset library exceeded that from using

197    DPHL in most cases (Figure 3A), and over 99% of the protein identifications were overlapped in

198 every cancer type (Figure S2). We collectively found that the subLib strategy outperformed the

199 DPHL strategy in all cancer types, with the most prominent increase from PTC carcinoma

200 samples (19.02% increase in protein groups and 36.17% increase in peptide precursors, Figure

201 2B). Of note, the discrimination ability to separate the targets from decoys led to a marked

202 increase (Figure 2C), further validating that the subLib strategy can reduce false positives in

203 clinical proteomic data. Missing values were equivalent between DPHL and the subLib strategy

204 (Figure 3B), and the protein quantification results were in good accordance as well with Pearson

205 correlation ratios all over 0.92 across all the tumor tissue types (Figure 3B), suggesting that

206 decreasing library sizes by adjusting FDR values does not impair protein identification nor

207 quantification. Moreover, different tumor types could be well resolved using the thus generated

208 protein matrix (Figure 2D). These results indicate that this subLib strategy could be generically

209 used for DIA data generated from different samples.

**Concluding remarks**
210

211      In this study, we present a computational strategy to optimize library size for DIA data

212 analysis. In our DIA data of human tissue specimens, setting FDR to 0.05 enabled effective

213 spectral library subsetting. The application of this strategy to DIA data from 15 tumor types

214 further consoidated this conclusion. This subLib strategy reduced false positive identifications,

215 increased peptide and protein identifications, and generated protein data matrix quantitatively

216 comparable to the DIA analysis with unfiltered library. In conclusion, the subLib strategy for

217 DIA spectral library size optimization boosts proteome identifications of DIA-MS data.

218

**Acknowledgements**
219

226

**Conflict of interest statement**

228 The research group of Tiannan Guo is partly supported by Pressure Biosciences Inc, which

229 provided access to advanced sample preparation instrumentation. T.G is shareholder of Westlake

230 Omics Inc. W.G. is employee of Westlake Omics Inc. The remaining authors declare no

231 competing interests.

232

**Data Availability**

234 The raw data and peptide/protein matrixes were deposited in ProteomeXchange Consortium

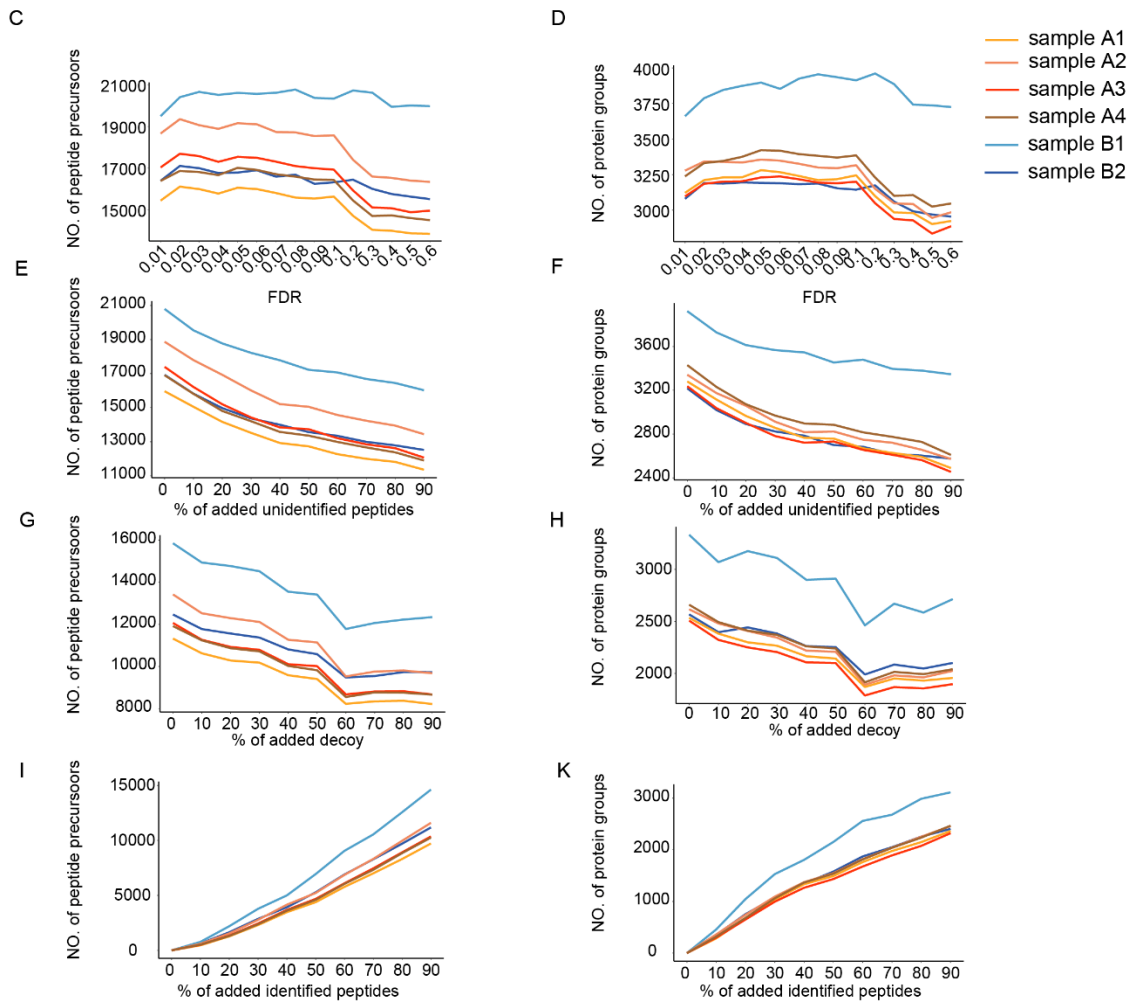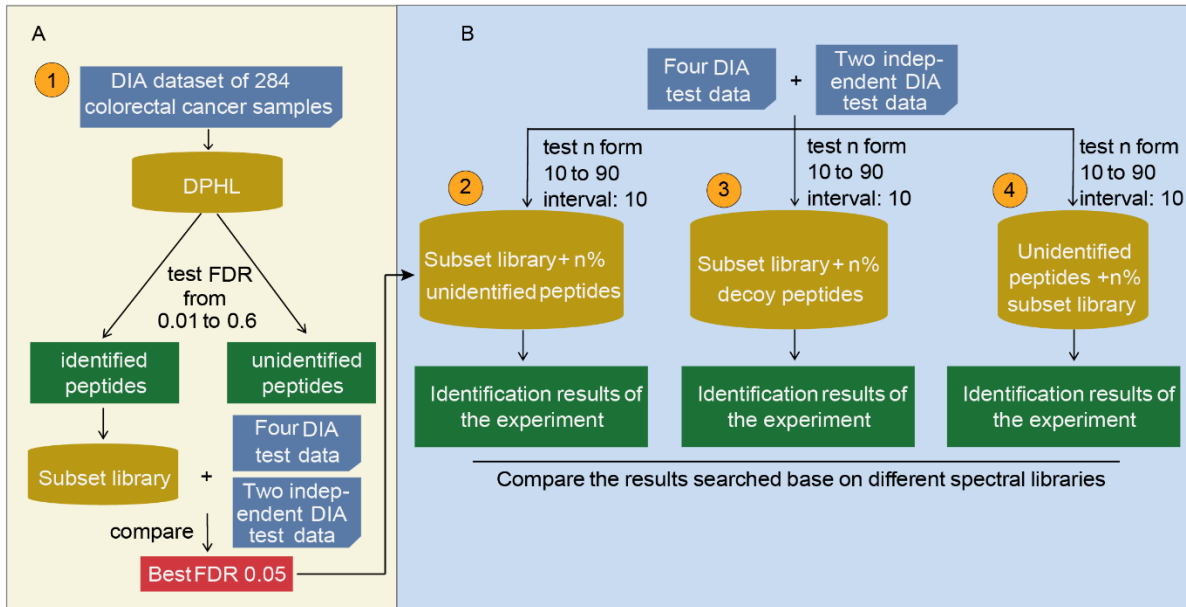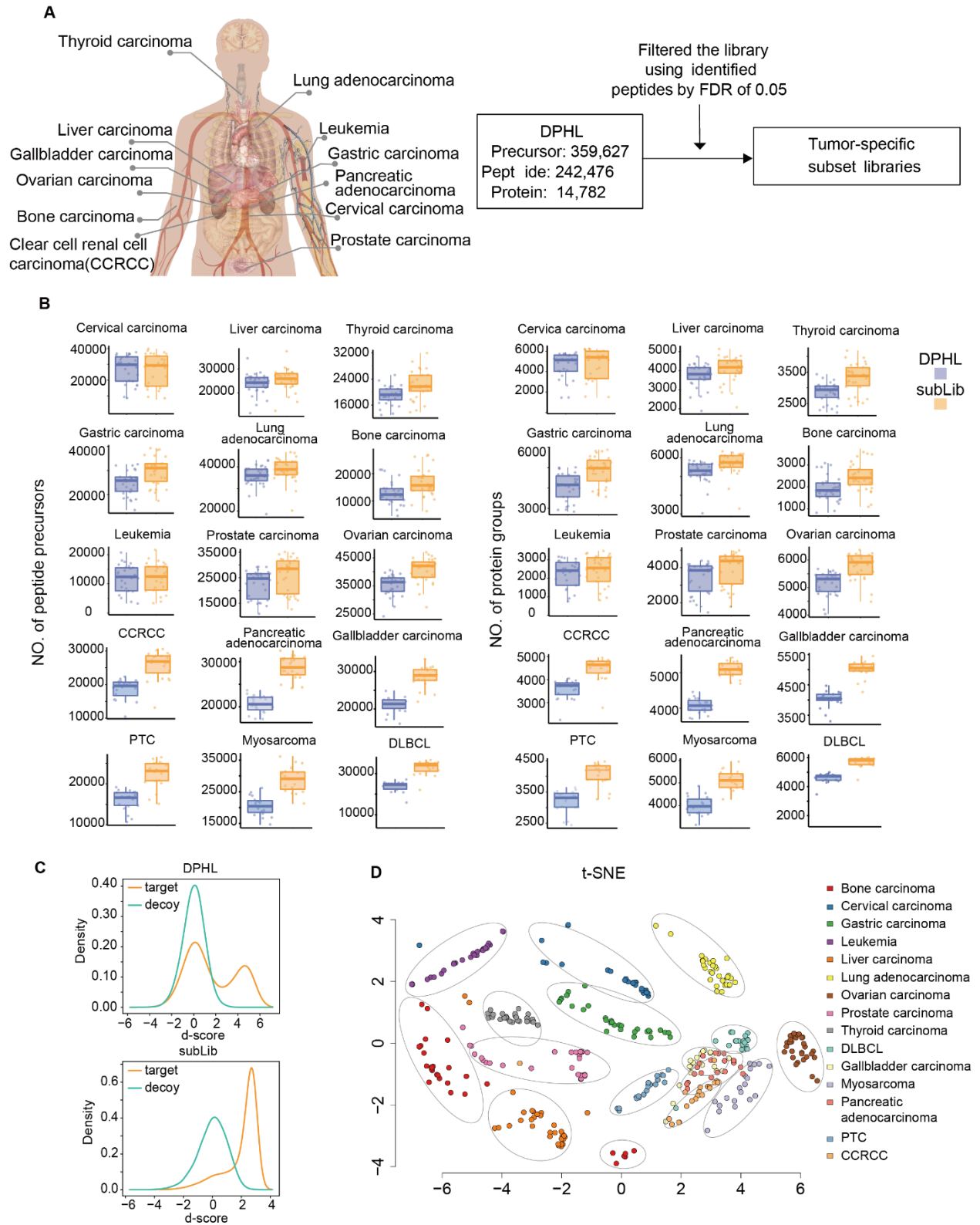235 (https://www.iprox.org/). Project ID: IPX0002439000 and IPX0001981000.

236

**References**

238 1. Yue, L., Zhang, F., Sun, R., Sun, Y., Yuan, C., Zhu, Y., and Guo, T. (2020) Generating Proteomic
239 Big Data for Precision Medicine. *PROTEOMICS* n/a, 1900358
240 2. Gillet, L. C., Navarro, P., Tate, S., Rost, H., Selevsek, N., Reiter, L., Bonner, R., and Aebersold, R.
241 (2012) Targeted data extraction of the MS/MS spectra generated by data-independent acquisition: a
242 new concept for consistent and accurate proteome analysis. *Mol Cell Proteomics* 11, O111 016717
243 3. Guo, T., Kouvonen, P., Koh, C. C., Gillet, L. C., Wolski, W. E., Rost, H. L., Rosenberger, G.,
244 Collins, B. C., Blum, L. C., Gillessen, S., Joerger, M., Jochum, W., and Aebersold, R. (2015) Rapid
245 mass spectrometric conversion of tissue biopsy samples into permanent quantitative digital proteome
246 maps. *Nat Med* 21, 407-413

247  4.  Zhang, F., Ge, W., Ruan, G., Cai, X., and Guo, T. (2020) Data-Independent Acquisition Mass
248      Spectrometry-Based Proteomics and Software Tools: A Glimpse in 2020. *PROTEOMICS* 20,
249      1900276
250  5.  Tsou, C. C., Avtonomov, D., Larsen, B., Tucholska, M., Choi, H., Gingras, A. C., and Nesvizhskii,
251      A. I. (2015) DIA-Umpire: comprehensive computational framework for data-independent acquisition
252      proteomics. *Nat Methods* 12, 258-+
253  6.  Demichev, V., Messner, C. B., Vernardis, S. I., Lilley, K. S., and Ralser, M. (2020) DIA-NN: neural
254      networks and interference correction enable deep proteome coverage in high throughput. *Nat
255      Methods* 17, 41-+
256  7.  Muntel, J., Gandhi, T., Verbeke, L., Bernhardt, O. M., Treiber, T., Bruderer, R., and Reiter, L. (2019)
257      Surpassing 10 000 identified and quantified proteins in a single run by optimizing current LC-MS
258      instrumentation and data analysis strategy. *Mol Omics* 15, 348-360
259  8.  Cho, K. C., Clark, D. J., Schnaubelt, M., Teo, G. C., Leprevost, F. D. V., Bocik, W., Boja, E. S.,
260      Hiltke, T., Nesvizhskii, A. I., and Zhang, H. (2020) Deep Proteomics Using Two Dimensional Data
261      Independent Acquisition Mass Spectrometry. *Anal Chem* 92, 4217-4225
262  9.  Zhong, C. Q., Wu, R., Chen, X., Wu, S., Shuai, J., and Han, J. (2020) Systematic Assessment of the
263      Effect of Internal Library in Targeted Analysis of SWATH-MS. *J Proteome Res* 19, 477-492
264  10. Schubert, O. T., Gillet, L. C., Collins, B. C., Navarro, P., Rosenberger, G., Wolski, W. E., Lam, H.,
265      Amodei, D., Mallick, P., MacLean, B., and Aebersold, R. (2015) Building high-quality assay libraries
266      for targeted analysis of SWATH MS data. *Nature Protocols* 10, 426-441
267  11. Rosenberger, G., Koh, C. C., Guo, T., Rost, H. L., Kouvonen, P., Collins, B. C., Heusel, M., Liu, Y.,
268      Caron, E., Vichalkovski, A., Faini, M., Schubert, O. T., Faridi, P., Ebhardt, H. A., Matondo, M.,
269      Lam, H., Bader, S. L., Campbell, D. S., Deutsch, E. W., Moritz, R. L., Tate, S., and Aebersold, R.
270      (2014) A repository of assays to quantify 10,000 human proteins by SWATH-MS. *Sci Data* 1,
271      140031
272  12. Zhu, T., Zhu, Y., Xuan, Y., Gao, H., Cai, X., Piersma, S. R., Pham, T. V., Schelfhorst, T., Haas, R. R.
273      G. D., Bijnsdorp, I. V., Sun, R., Yue, L., Ruan, G., Zhang, Q., Hu, M., Zhou, Y., Van Houdt, W. J.,
274      Lelarge, T. Y. S., Cloos, J., Wojtuszkiewicz, A., Koppers-Lalic, D., Böttger, F., Scheepbouwer, C.,
275      Brakenhoff, R. H., van Leenders, G. J. L. H., Ijzermans, J. N. M., Martens, J. W. M., Steenbergen, R.
276      D. M., Grieken, N. C., Selvarajan, S., Mantoo, S., Lee, S. S., Yeow, S. J. Y., Alkaff, S. M. F., Xiang,
277      N., Sun, Y., Yi, X., Dai, S., Liu, W., Lu, T., Wu, Z., Liang, X., Wang, M., Shao, Y., Zheng, X., Xu,
278      K., Yang, Q., Meng, Y., Lu, C., Zhu, J., Zheng, J. e., Wang, B., Lou, S., Dai, Y., Xu, C., Yu, C.,
279      Ying, H., Lim, T. K., Wu, J., Gao, X., Luan, Z., Teng, X., Wu, P., Huang, S. a., Tao, Z., Iyer, N. G.,
280      Zhou, S., Shao, W., Lam, H., Ma, D., Ji, J., Kon, O. L., Zheng, S., Aebersold, R., Jimenez, C. R., and
281      Guo, T. (2020) DPHL: A DIA Pan-human Protein Mass Spectrometry Library for Robust Biomarker
282      Discovery. *Genomics, Proteomics & Bioinformatics*
283  13. Blattmann, P., Stutz, V., Lizzo, G., Richard, J., Gut, P., and Aebersold, R. (2019) Generation of a
284      zebrafish SWATH-MS spectral library to quantify 10,000 proteins. *Scientific Data* 6, 190011
285  14. Zhang, H., Liu, P., Guo, T., Zhao, H., Bensaddek, D., Aebersold, R., and Xiong, L. (2019)
286      Arabidopsis proteome and the mass spectral assay library. *Sci Data* 6, 278
287  15. Parker, S. J., Venkatraman, V., and Van Eyk, J. E. (2016) Effect of peptide assay library size and
288      composition in targeted data-independent acquisition-MS analyses. *Proteomics* 16, 2221-2237
289  16. Zi, J., Zhang, S., Zhou, R., Zhou, B., Xu, S., Hou, G., Tan, F., Wen, B., Wang, Q., Lin, L., and Liu,
290      S. (2014) Expansion of the ion library for mining SWATH-MS data through fractionation
291      proteomics. *Anal Chem* 86, 7242-7246
292  17. Wu, J. X., Song, X., Pascovici, D., Zaw, T., Care, N., Krisp, C., and Molloy, M. P. (2016) SWATH
293      Mass Spectrometry Performance Using Extended Peptide MS/MS Assay Libraries. *Mol Cell
294      Proteomics* 15, 2501-2514
295  18. Rosenberger, G., Bludau, I., Schmitt, U., Heusel, M., Hunter, C. L., Liu, Y., MacCoss, M. J.,
296      MacLean, B. X., Nesvizhskii, A. I., Pedrioli, P. G. A., Reiter, L., Rost, H. L., Tate, S., Ting, Y. S.,

297     Collins, B. C., and Aebersold, R. (2017) Statistical control of peptide and protein error rates in large-
298     scale targeted data-independent acquisition analyses. *Nat Methods* 14, 921-927
299  19. Teleman, J., Hauri, S., and Malmstrom, J. (2017) Improvements in Mass Spectrometry Assay Library
300     Generation for Targeted Proteomics. *J Proteome Res* 16, 2384-2392
301  20. Escher, C., Reiter, L., MacLean, B., Ossola, R., Herzog, F., Chilton, J., MacCoss, M. J., and Rinner,
302     O. (2012) Using iRT, a normalized retention time for more targeted measurement of peptides.
303     *Proteomics* 12, 1111-1121
304  21. Barkovits, K., Pacharra, S., Pfeiffer, K., Steinbach, S., Eisenacher, M., Marcus, K., and Uszkoreit, J.
305     (2020) Reproducibility, Specificity and Accuracy of Relative Quantification Using Spectral Library-
306     based Data-independent Acquisition. *Molecular & Cellular Proteomics* 19, 181-197
307  22. Zhu, Y., and Guo, T. (2018) High-Throughput Proteomic Analysis of Fresh-Frozen Biopsy Tissue
308     Samples Using Pressure Cycling Technology Coupled with SWATH Mass Spectrometry. In: Sarwal,
309     M. M., and Sigdel, T. K., eds. *Tissue Proteomics: Methods and Protocols*, pp. 279-287, Springer
310     New York, New York, NY
311  23. Adusumilli, R., and Mallick, P. (2017) Data Conversion with ProteoWizard msConvert. *Methods Mol*
312     *Biol* 1550, 339-368
313  24. Rost, H. L., Rosenberger, G., Navarro, P., Gillet, L., Miladinovic, S. M., Schubert, O. T., Wolski, W.,
314     Collins, B. C., Malmstrom, J., Malmstrom, L., and Aebersold, R. (2014) OpenSWATH enables
315     automated, targeted analysis of data-independent acquisition MS data. *Nat Biotechnol* 32, 219-223
316  25. Gao, H., Zhang, F., Liang, S., Zhang, Q., Lyu, M., Qian, L., Liu, W., Ge, W., Chen, C., Yi, X., Zhu,
317     J., Lu, C., Sun, P., Liu, K., Zhu, Y., and Guo, T. (2020) Accelerated Lysis and Proteolytic Digestion
318     of Biopsy-level Fresh Frozen and FFPE Tissue Samples Using Pressure Cycling Technology. *J*
319     *Proteome Res*

320

14

321  **Figure 1. Optimizing DPHL in the DIA dataset of colorectal cancer.** (A) The workflow of

322  spectral library optimization. Step 1: Select the best FDR for refining the subset library from the

323  public DPHL library. The subset library refined from DPHL with FDR of 0.05 is considered as

324  the optimal subset library, which was used as a primary optimized subset spectral library in this

325  study. Step 2: Evaluate the performance of the spectral library consisting of the subset library

326  and n% unidentified peptides. Step 3: Evaluate the performance of spectral library consisting of

327  subset library and n% decoy peptides. Step 4:  Evaluate the performance of spectral library

328  consisting of unidentified peptides and n% peptides from the subset library. By comparing all the

329  identification results, the subset library refined from DPHL with FDR of 0.05 is the best

330  experiment-specific spectral library for DIA data analysis. The numbers of identified peptides

331  (C) and proteins (D) based on the subset libraries which were refined from DPHL at nine

332  different FDRs.  The numbers of identified peptides (E) and proteins (F) based on the spectral

333  libraries consisting of subset library and n% unidentified peptides. The numbers of identified

334  peptides (G) and proteins (H) based on the spectral libraries consisting of the subset library and

335  n% decoy peptides. The numbers of identified peptides (I) and proteins (K) based on the spectral

336  libraries consisting of unidentified peptides and n% peptides from the subset library.
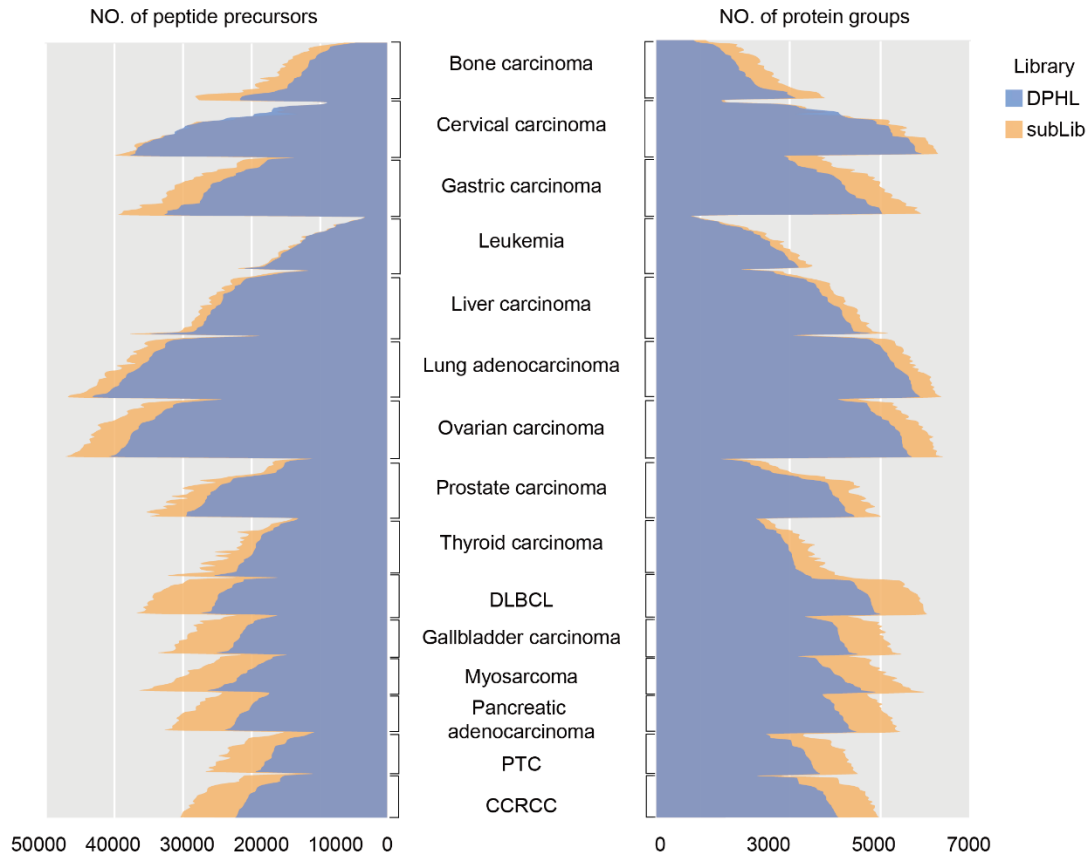
337

16

338    **Figure 2. Tumor-specific subset library improves the identifications compared with DPHL.**
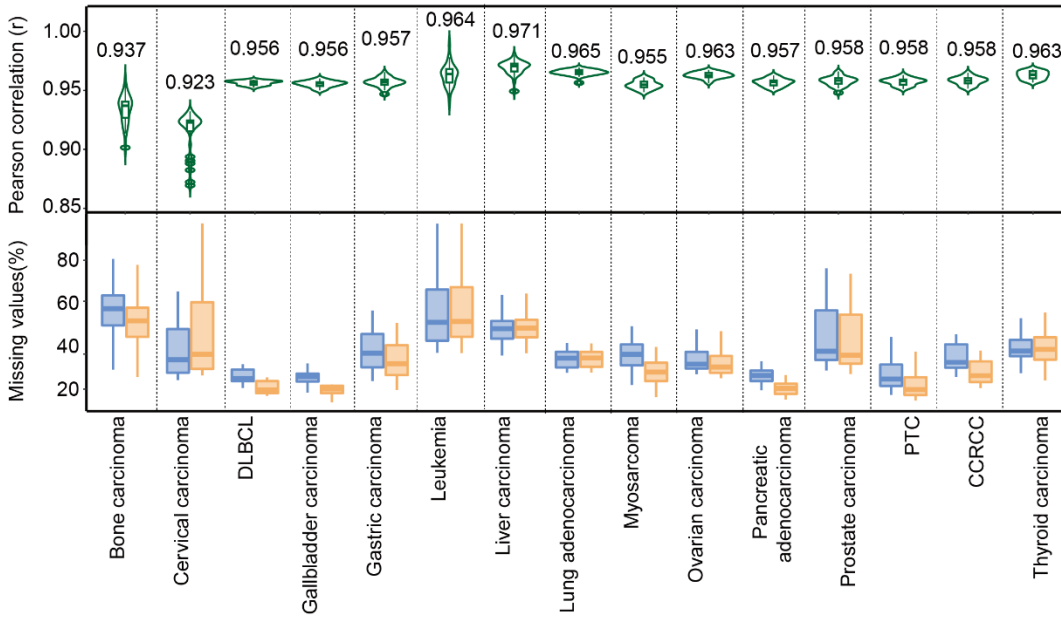
339    (A) The workflow of the subLib strategy. (B) The number of peptides and proteins identified

340    base on tumor-specific subLib and DPHL in 15 tumor types. (C) The distribution of

341    discrimination score (d-score) of the target and decoy of the subset library and DPHL. (D) The

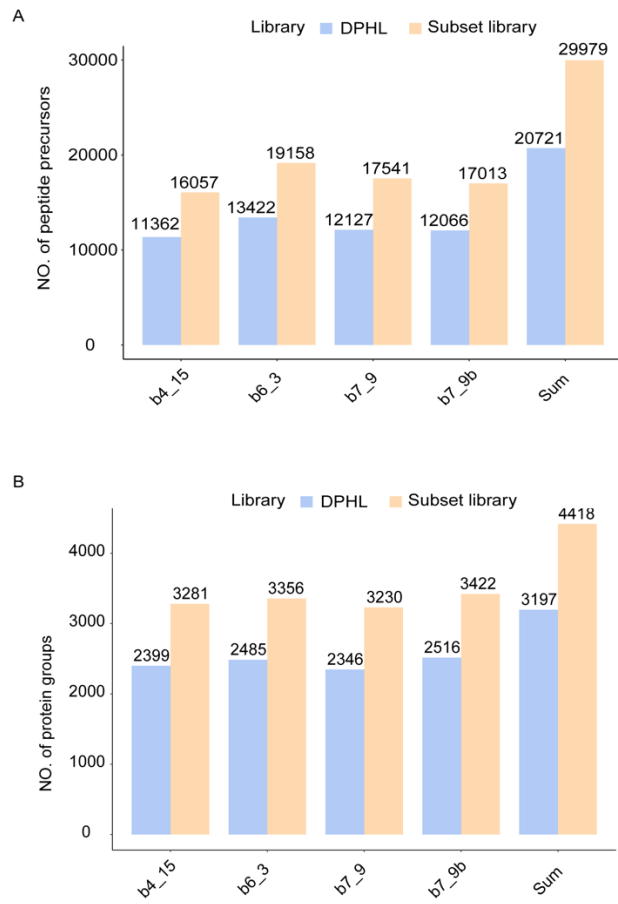342    tSNE plot shows the samples are well resolved by tissue type.
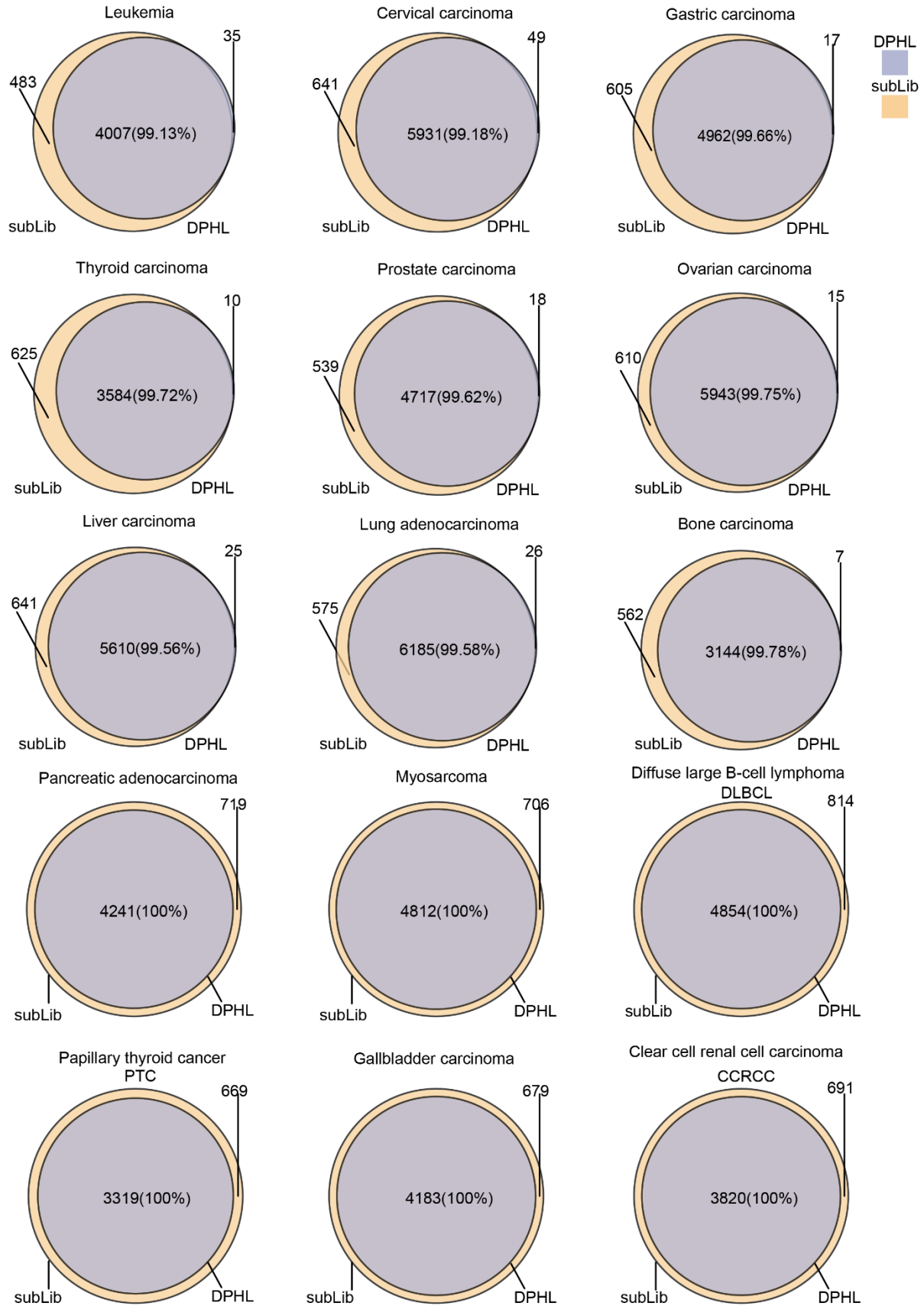
343

A



B



344

18

345 **Figure 3. Peptide precursor and protein identification using the optimized subset library**

346 **and DPHL.** (A) The number of peptide precursors and protein groups identified using the

347 optimized subset library and DPHL for each sample of every tumor type. subLib, the optimized

348 subset library. Protein identifications were shown on the right, and peptide precursor

349 identifications were shown on the left. (B) The correlation values on the protein level between

350 identification results of the optimized subset library and DPHL. The percentages of protein

351 missing values identified base on DPHL and the optimized library of each tumor type.

352

353

**Figure S1. Identification results of the four representative DIA-MS data in the colorectal cancer cohort.**

356

Leukemia — 483 / 4007(99.13%) / 35

Cervical carcinoma — 641 / 5931(99.18%) / 49

Gastric carcinoma — 605 / 4962(99.66%) / 17

DPHL
subLib

Thyroid carcinoma — 625 / 3584(99.72%) / 10

Prostate carcinoma — 539 / 4717(99.62%) / 18

Ovarian carcinoma — 610 / 5943(99.75%) / 15

Liver carcinoma — 641 / 5610(99.56%) / 25

Lung adenocarcinoma — 575 / 6185(99.58%) / 26

Bone carcinoma — 562 / 3144(99.78%) / 7

Pancreatic adenocarcinoma — 4241(100%) / 719

Myosarcoma — 4812(100%) / 706

Diffuse large B-cell lymphoma DLBCL — 4854(100%) / 814

Papillary thyroid cancer PTC — 3319(100%) / 669

Gallbladder carcinoma — 4183(100%) / 679

Clear cell renal cell carcinoma CCRCC — 3820(100%) / 691

357

358 **Figure S2. Venn diagrams showing overlap of protein identifications between the optimized**

359 **subset library and DPHL.**

360