

1 **A comprehensive mass spectral library for human**
2 **thyroid tissues**

3 **Authors**

4 Yaoting Sun^{1,2,3,4}, Lu Li^{2,3,4}, Weigang Ge⁵, Zhen Dong^{2,3,4}, Wei Liu⁵, Hao Chen⁵, Qi Xiao^{2,3,4}, Xue
5 Cai^{2,3,4}, Fangfei Zhang^{2,3,4}, Junhong Xiao⁶, Guangzhi Wang⁶, Yi He⁷, Oi Lian Kon⁸, N.
6 Gopalakrishna Iyer^{8,9}, Yongfu Zhao⁶, Tiannan Guo^{1,2,3,4}, *

7

8 ¹Zhejiang University, Hangzhou, China;

9 ²Key Laboratory of Structural Biology of Zhejiang Province, School of Life Sciences, Westlake
10 University, Hangzhou 310024, China;

11 ³Center for Infectious Disease Research, Westlake Laboratory of Life Sciences and Biomedicine,
12 Hangzhou 310024, China;

13 ⁴Institute of Basic Medical Sciences, Westlake Institute for Advanced Study, Hangzhou 310024, China;

14 ⁵Westlake Omics (Hangzhou) Biotechnology Co., Ltd., Hangzhou 310024, China;

15 ⁶Department of General Surgery, The Second Hospital of Dalian Medical University, Dalian, China;

16 ⁷Department of Urology, The Second Hospital of Dalian Medical University, Dalian, China;

17 ⁸Division of Medical Sciences, National Cancer Centre Singapore, Republic of Singapore;

18 ⁹Department of Head and Neck Surgery, National Cancer Centre Singapore, Republic of Singapore.

19 *Corresponding author(s): Tiannan Guo (guotiannan@westlake.edu.cn)

20

21 **Abstract**

22 Thyroid nodules occur in about 60% of the population. Current diagnostic strategies,
23 however, often fail at distinguishing malignant nodules before surgery, thus leading to
24 unnecessary, invasive treatments. As proteins are involved in all physio/pathological
25 processes, a proteome investigation of biopsied nodules may help correctly classify and
26 identify malignant nodules and discover therapeutic targets. Quantitative mass spectrometry
27 data-independent acquisition (DIA) enables highly reproducible and rapid throughput
28 investigation of proteomes. An exhaustive spectral library of thyroid nodules is essential for
29 DIA yet still unavailable. This study presents a comprehensive thyroid spectral library
30 covering five types of thyroid tissue: multinodular goiter, follicular adenoma, follicular and
31 papillary thyroid carcinoma, and normal thyroid tissue. Our library includes 925,330
32 transition groups, 157,548 peptide precursors, 121,960 peptides, 9941 protein groups, and
33 9826 proteins from proteotypic peptides. This library resource was evaluated using three
34 papillary thyroid carcinoma samples and their corresponding adjacent normal thyroid tissue,
35 leading to effective quantification of up to 7863 proteins from biopsy-level thyroid tissues.
36

37 **Background & Summary**

38 Thyroid nodules are common and, given the sensitivity of current diagnostic techniques, can
39 be detected in approximately 60% of the general population, especially in women^{1,2}. The
40 incidence of thyroid malignancy or thyroid carcinoma, has rapidly increased over the last
41 decades, although it is uncertain if this is a real increase or simply a result of widespread use
42 of screening ultrasonography^{3,4}. Most of these nodules are asymptomatic. Only 4-7% of
43 patients present with complaints attributed to thyroid nodules. Although ultrasonography and
44 ultrasound-guided fine-needle aspiration can help distinguish between benign and malignant
45 nodules, approximately 30% of thyroid nodules remain indeterminate by cytopathology and
46 require diagnostic surgery⁵, after which histopathology of surgical specimens provides a
47 definitive and complete diagnosis. More importantly, only 15% of indeterminate nodules
48 prove to be malignant. Because many benign nodules are clinically ambiguous and a source
49 of uncertainty, such patients often undergo unnecessary surgery. Nucleic acid-based
50 molecular tests, which require next-generation sequencing technology⁶, are currently used in
51 clinical practice to reduce overtreatment of thyroid nodules. However, the diagnostic
52 specificity of these tests remains modest at best (40-70%) for a myriad of reasons.
53

54 Unlike nucleic acids, proteins are directly involved in all life processes and determine cellular
55 and organismal phenotype. Proteins also have the potential to be critical biomarkers for
56 disease diagnosis and are themselves potential drug targets. For these reasons, there is
57 tremendous potential in exploring thyroid molecular pathology from a protein-based
58 perspective. Mass spectrometry (MS) -based proteomics has reached a high level of technical
59 and methodological development during the last decade. Data-independent acquisition (DIA),
60 in particular, enables comprehensive quantitation of peptides from complex compositions
61 with high reproducibility and throughput⁷. In the conventional data-dependent acquisition
62 (DDA) mode, only peptide precursors with high abundance in MS1 are fragmented. In DIA,

63 however, all precursors within a predefined range (also called window) of mass-to-charge
64 ratio (m/z) are fragmented by sequentially repeated cycling in windows, thus providing
65 detailed data without loss of any eluted peptides^{7,8}. Our group's established pressure cycling
66 technology (PCT)-based sample preparation methodology, coupled with DIA-MS, achieves
67 the acquisition of complete proteomic information in less than six hrs^{9,10}.
68 To optimize the efficiency of spectral identifications, DIA data analysis requires tissue- or
69 organism-specific spectral libraries^{8,11}. Although a pan-human library derived from healthy
70 subjects has already been established^{12,13}, this extensive and non-specific library could cause
71 inaccuracies during ion matching. In recent years, several novel software for DIA data
72 analysis, such as DIA-Umpire¹⁴, PECAN¹⁵, or DIA-NN¹⁶, no longer require spectral libraries.
73 However, this library-free mode should be applied with caution due to its lower sensitivity
74 and protein identification power compared to a library-based strategy¹⁷. A tissue-specific
75 library for thyroid nodules, both benign and malignant, as well as for healthy thyroid, would
76 thus provide an essential resource for the proteomic investigation of thyroid pathologies in a
77 high-throughput manner.

78
79 This study introduces a thyroid-specific spectral library to support protein identification and
80 quantification in thyroid nodules by DIA-MS (Figure 1). Five types of thyroid tissues were
81 collected, namely normal tissue, two types of benign nodules (multinodular goiter (MNG) and
82 follicular adenoma (FA), and two types of thyroid carcinomas (follicular thyroid carcinoma
83 (FTC) and papillary thyroid carcinoma (PTC). Normal thyroid and thyroid nodule tissues
84 were processed by PCT; extracted and desalted peptides were then combined into three
85 different pooled samples: (1) pooled sample containing all five types, (2) PTC pooled
86 samples, and (3) FA and FTC pooled sample. The pooled peptides were fractionated in two
87 ways, *i.e.* strong cation exchange (SCX) or high-pH reversed-phase chromatography, to
88 achieve higher peptide coverage. Peptide fractions were injected into HPLC-MS/MS with 60
89 min-gradient through DDA mode using Thermo Orbitrap Q ExactiveTM HF. 46 DDA files
90 were acquired in total. Our spectral library was built with Spectronaut 14.6 and included
91 925,330 transition groups, 157,548 precursors, 121,960 peptides, 9941 protein groups, and
92 9826 proteins from proteotypic peptides. We then validated this library by applying it to four
93 DIA datasets acquired with four different methods (Figure 1).

94

95 **Methods**

96 **Sample collection**

97 For the spectral library construction and testing, thyroid healthy and nodular samples were
98 collected, between 2011 and 2019, from two clinical centers in Singapore (Singapore General
99 Hospital) and China (The Second Hospital of Dalian Medical University). Ethical approval
100 was given by both hospitals. Tissue cores of 1 mm diameter (0.6-1.2 mg) were extracted from
101 the pathological regions of interest in formalin-fixed paraffin-embedded (FFPE) tissue blocks
102 demarcated by experienced histopathologists¹⁸. Four types of thyroid nodules (42 MNG, 49
103 FA, 33 FTC, and 54 PTC) and 10 normal thyroid tissues were used for building the library.
104 We also collected three paired PTC and corresponding tumor-adjacent tissues for validation
105 of the spectral library.

106 **Sample preparation assisted by PCT**

107 Samples were dewaxed, hydrated, and acidified using, in sequence, heptane, a decreasing
108 ethanol series (100%, 90%, and 75%), and formic acid. The samples were next kept under
109 basic hydrolysis conditions in Tris-HCl (100 mM, pH = 10) at 95 °C for 30 min, then
110 transferred into a solution containing 30 µL lysis buffer (6 M urea, 2 M thiourea), 5 µL tris(2-
111 carboxyethyl)phosphine (TECP, 10 mM), and 2.5 µL iodoacetamide (IAA) (40 mM). In PCT-
112 Micro tubes, samples were lysed, reduced, and hydroxylated at 30 °C using PCT (120 cycles,
113 45 Kpsi, 30 s on-time, 10 s off-time). Trypsin (enzyme-to-substrate ratio, 1:50; Hualishi
114 Scientific, China) and LysC (enzyme-to-substrate ratio, 1:40; Hualishi Scientific, China) were
115 then added, followed by PCT-assisted digestion (120 cycles, 20 Kpsi, 50 s on-time, 10 s off-
116 time). 1% trifluoroacetic acid (TFA) was added to terminate the digestion process. The
117 resulting peptides were desalted with 2% acetonitrile (ACN) and 0.1% TFA and reconstituted
118 with 2% ACN containing 0.1% formic acid. Peptide concentrations were measured by
119 Nanoscan (Analytic Jena, Germany) at A₂₈₀, and samples were stored at 4 °C for further
120 analysis. For sample testing, we used previously optimized methods^{10,19}. All the chemical
121 reagents were obtained from Sigma-Aldrich.

122 **Strong cation exchange (SCX) fractionation of peptides**

123 Clean peptides were fractionated by 100 mg SCX solid-phase extraction (SPE) columns
124 (HyperSep™, Thermo Fisher Scientific) to enhance the peptide spectral information. 600 µg
125 of pooled peptides, including all five types of thyroid tissues (10 N, 42 MNG, 28 FA, 13 FTC,
126 38 PTC), were reconstituted in equilibration buffer (2.5 mM KH₂PO₄ / 25% ACN, pH = 3.0).
127 SCX columns were washed with MilliQ water and equilibration buffer. The pooled sample
128 was then loaded into a conditioned cartridge. Loaded columns were washed with six diluents
129 with different ratios of buffer A (10 mM KH₂PO₄ / 25% ACN, pH = 3.0) to buffer B (10 mM
130 KH₂PO₄ / 1 M KCl / 25% ACN, pH = 3.0) and increasing KCl concentration. The samples
131 were then split into six fractions and cleaned by C18 spin columns (The Nest Group, United
132 States).

133 **High-pH reversed-phase chromatography fractionation of peptides**

134 To further increase the peptide coverage in the spectral library, another fractionation method,
135 *i.e.* high-pH reversed-phase chromatography, was performed. Two pooled samples were
136 combined from 41 follicular thyroid neoplasms (21 FA and 20 FTC) and 16 PTC samples.
137 ~200 µg of each pooled sample was separated by Thermo DineX Ultramate 3000 with an
138 XBridge peptide BEH C18 column (4.6 mm X 250 mm, 5 µm, 1 / pkg) at 45 °C. The gradient
139 was 60 min long, with a flow rate of 1 mL/min, and the mobile phase consisting of buffer A
140 (ddH₂O water with 0.6% ammonia, pH = 10) and buffer B (98% ACN with 0.6% ammonia,
141 pH= 10). The gradient was from 5% to 35% buffer B in condition of pH 10.0 at a flow rate of
142 1 mL/min. 60 fractions were collected, separated by 1 min interval. The 60 fractions were
143 subsequently combined into 20 fractions for each pooled sample to build the library. The
144 resulting fractionated peptides were then resuspended into 20 µl buffer (2% of ACN, 0.1%
145 formic acid) for instrument injection.

146 **Data dependent acquisition (DDA)**

147 The fractionated peptides were separated by UltiMate™ 3000 RSLCnano System (Thermo
148 Fisher Scientific). The system was equipped with a 15 cm x 75 µm silica column custom
149 packed with 1.9 µm 100 Å C18-Aqua. The mobile phase comprised buffer A (2% ACN, 0.1%

150 formic acid) and buffer B (98% ACN, 0.1% formic acid). Peptides were separated on a 60 min
151 effective liquid chromatography (LC) buffer B gradient (3% to 28% at 300 nL/min). Ionized
152 peptides were transferred into a Q Exactive™ HF MS (Thermo Fisher Scientific). Full MS
153 scans were measured with an Orbitrap at a resolution of 60,000 full widths at half maximum
154 (FWHM) at 200 m/z covering 400 to 1200 m/z precursors, with automatic gain control (AGC)
155 target value of 3E6 charges and 80 ms maximum injection time (max IT). The top 20
156 precursor signals were chosen to be fragmented in a higher-energy collision (HCD) cell with
157 27% normalized collision energy and then transferred to an Orbitrap for MS/MS analysis at a
158 resolution of 30,000 FWHM and an AGC target value of 1E5. By using 60 min LC gradients,
159 we acquired a total of 46 DDA files (Details are listed in Supplementary Table 1).

160 **Spectral library construction based on DDA**

161 Spectronaut™ Pulsar X version 14.6 (Biognosys) was used to generate a spectral library
162 specific to the thyroid. All 46 DDA raw files were searched by Pulsar against a human Swiss-
163 Prot FASTA database (downloaded on 2020-01-22) which included 20,367 protein sequences
164 with FDR of 0.01. The enzyme setting for “trypsin/P” allowed no more than two missed
165 cleavages; cysteine carbamidomethyl was set as a fixed modification, and methionine
166 oxidation was set as a variable modification; mass tolerances were automatically determined,
167 while other settings were left to their default values.

168 **Quantitative analysis of thyroid samples by data independent acquisition (DIA) and** 169 **PulseDIA**

170 Together with paired tissues adjacent to the tumor site, three PTC samples were prepared as
171 previously described. Proteomic data for these test thyroid samples was acquired by DIA or
172 PulseDIA, a gas phase fractionation method²⁰. For each run, the LC effective gradient was 45
173 min long, with 3% to 25% buffer B at 0.3 $\mu\text{L}/\text{min}$. MS1 was performed over an m/z range of
174 390-1010 for the DIA, and 390-1210 for the PulseDIA, with a resolution of 60,000 FWHM,
175 an AGC target of 3E6, and a max IT of 80 ms. MS2 was performed with a resolution of
176 30,000 FWHM, an AGC target of 1E6, and a max IT of 55 ms. For DIA, 24 isolation
177 windows were performed: 20 with 21 m/z wide windows, 2 with 41 m/z wide windows, and 2
178 with 61 m/z windows. For PulseDIA, five injections with 24 isolation windows per injection
179 were performed²⁰. DIA data were analyzed by Spectronaut™ version 14.6; all settings were
180 left to their default values.

181

182 **Data Records**

183 DDA raw files (Data Citation 1)

184 Spectral library files in the formats of xlsx, TSV and CSV (Data Citation 2)

185 DIA raw files (Data Citation 3)

186 **Technical Validation**

187 **Libraries evaluation**

188 Our thyroid-specific spectral library comprises 925,330 transition groups, 157,548 precursors,
189 121,960 peptides, 9941 protein groups, and 9826 proteins from proteotypic peptides. An
190 overview is provided in Table 1. Our library is, therefore, more comprehensive than currently
191 published data which include only 2682 proteins²¹.

192 To assess the quality of our spectral library, we first evaluated the composition and
193 distributions of precursors, peptides, and proteins. In our DIA-MS analysis, the precursor
194 mass range cover 400-1200 m/z , and approximately 82% of the precursors are between 400-
195 850 m/z (Figure 2A). Precursors primarily display two (53%) or three (37%) charges, and
196 their charge distributions are comparable to those of different spectral libraries (Figure 2B)²².
197 82% peptides are 8 to 20 amino acids long, with a median length of 14 amino acids,
198 consistently with the properties of trypsinized peptides (Figure 2C). We next focused on
199 peptide modifications. Oxidation on methionine, the most common modification in our
200 library, was detected in 22,853 peptides, 121,960 of the total peptides. Sample preparation
201 generated 2818 carbamidomethylated peptides at cysteine residues and 2231 N-terminal
202 acetylated ones (Figure 2D). A total of 7634 proteins were detected with at least three
203 proteotypic peptides, and the majority of proteins were found with more than ten (Figure 2E).
204 Additionally, fragments from y -ions were more frequently detected than those from b -ions
205 due to the collision mode. Our established spectral library achieved comprehensive peptide
206 and protein coverage with high quality.
207 We next used Gene Ontology to identify the main enriched protein categories within our
208 library. A total of 9,825 proteins were annotated by Ingenuity Pathway Analysis (IPA)
209 software: the enriched protein cellular locations (red words) and protein functions (black
210 words) are shown in Figure 2G. By matching our data to the kinase database KinMap²³, our
211 library was found to contain 340 kinases from 7 families, accounting for 63.4% (340/536) of
212 the entire kinase database (Figure 2H). These results demonstrate that our library provides a
213 valuable reference for the application of the DIA-MS method to human thyroid samples.

214

215 **Technical validation on four datasets**

216 To further validate our library, we analyzed three PTC samples, together with paired tissue
217 samples adjacent to the tumor site. Four datasets were then acquired with the following four
218 acquisition strategies: single-shot DIA (dataset 1), PulseDIA (dataset 2), pre-fraction DIA
219 (dataset 3), and a combination of pre-fraction and PulseDIA (dataset 4). All datasets were
220 subsequently analyzed using Spectronaut 14.6 and our thyroid nodule-specific spectral
221 library. The search results for the four datasets are shown in Figure 3. All three tumor tissues
222 expressed more proteins and peptides than the matched normal thyroid tissues (tumor-
223 adjacent tissues), and this was especially evident at the peptide level (Figure 3A, B). The
224 numbers of identified peptides and proteins using single-shot DIA were the lowest due to the
225 relatively short gradient and the highly abundant protein, thyroglobulin. PulseDIA and pre-
226 fraction DIA led to more identifications. PulseDIA identified more peptides than pre-fraction
227 DIA, but a comparable number of proteins. Finally, a combination of pre-fraction and
228 PulseDIA generated the best results at both peptide and protein levels: 65,544 peptides and
229 7863 proteins. These results showed that a longer gradient allows the detection of more
230 peptides and proteins.

231 We next calculated the coefficient of variation (CV) of peptides and proteins abundance to
232 evaluate the quality of these datasets. The median peptides CVs were less than 0.05 for all
233 datasets (Figure 3C). Similarly, the median proteins CVs were all less than 0.04 (Figure 3D).
234 These results indicate that all four datasets performed well as the quantifications had only

235 negligible differences. These results confirm that our spectral library as a valuable resource
236 provides a robust reference for proteomic exploration of thyroid disease.
237 Although five types of thyroid tissues and more than 10,000 proteins are in our spectral
238 library, some rare thyroid carcinomas such as anaplastic thyroid carcinoma and medullary
239 thyroid carcinoma were not included in the analysis. This could be addressed in the future
240 with the methodology adopted here. Targeted assays using parallel reaction monitoring
241 (PRM) and selected/multiple reaction monitoring S/MRM could also be developed based on
242 this DIA library¹³. In conclusion, our established DIA library offers a useful resource for
243 proteomic analysis of thyroid tissue specimens.
244

245 **References**

- 246 1 Burman, K. D. & Wartofsky, L. CLINICAL PRACTICE. Thyroid Nodules. *N Engl J*
247 *Med* **373**, 2347-2356, doi:10.1056/NEJMcp1415786 (2015).
- 248 2 Singh Ospina, N., Iniguez-Ariza, N. M. & Castro, M. R. Thyroid nodules: diagnostic
249 evaluation based on thyroid cancer risk assessment. *BMJ* **368**, l6670,
250 doi:10.1136/bmj.l6670 (2020).
- 251 3 Miranda-Filho, A. *et al.* Thyroid cancer incidence trends by histology in 25 countries:
252 a population-based study. *Lancet Diabetes Endocrinol* **9**, 225-234,
253 doi:10.1016/S2213-8587(21)00027-9 (2021).
- 254 4 Lim, H., Devesa, S. S., Sosa, J. A., Check, D. & Kitahara, C. M. Trends in Thyroid
255 Cancer Incidence and Mortality in the United States, 1974-2013. *JAMA* **317**, 1338-
256 1348, doi:10.1001/jama.2017.2719 (2017).
- 257 5 Fagin, J. A. & Wells, S. A., Jr. Biologic and Clinical Perspectives on Thyroid Cancer.
258 *N Engl J Med* **375**, 1054-1067, doi:10.1056/NEJMra1501993 (2016).
- 259 6 Wang, T. S. & Sosa, J. A. Thyroid surgery for differentiated thyroid cancer - recent
260 advances and future directions. *Nat Rev Endocrinol* **14**, 670-683,
261 doi:10.1038/s41574-018-0080-7 (2018).
- 262 7 Gillet, L. C. *et al.* Targeted data extraction of the MS/MS spectra generated by data-
263 independent acquisition: a new concept for consistent and accurate proteome analysis.
264 *Mol Cell Proteomics* **11**, O111.016717, doi:10.1074/mcp.O111.016717 (2012).
- 265 8 Zhang, F., Ge, W., Ruan, G., Cai, X. & Guo, T. Data-Independent Acquisition Mass
266 Spectrometry-Based Proteomics and Software Tools: A Glimpse in 2020. *Proteomics*
267 **20**, e1900276, doi:10.1002/pmic.201900276 (2020).
- 268 9 Guo, T. *et al.* Rapid mass spectrometric conversion of tissue biopsy samples into
269 permanent quantitative digital proteome maps. *Nat Med* **21**, 407-413,
270 doi:10.1038/nm.3807 (2015).
- 271 10 Gao, H. *et al.* Accelerated Lysis and Proteolytic Digestion of Biopsy-Level Fresh-
272 Frozen and FFPE Tissue Samples Using Pressure Cycling Technology. *J Proteome*
273 *Res* **19**, 1982-1990, doi:10.1021/acs.jproteome.9b00790 (2020).
- 274 11 Schubert, O. T. *et al.* Building high-quality assay libraries for targeted analysis of
275 SWATH MS data. *Nat Protoc* **10**, 426-441, doi:10.1038/nprot.2015.015 (2015).
- 276 12 Rosenberger, G. *et al.* A repository of assays to quantify 10,000 human proteins by
277 SWATH-MS. *Sci Data* **1**, 140031, doi:10.1038/sdata.2014.31 (2014).

- 278 13 Zhu, T. *et al.* DPHL: A DIA Pan-human Protein Mass Spectrometry Library for
279 Robust Biomarker Discovery. *Genomics Proteomics Bioinformatics* **18**, 104-119,
280 doi:10.1016/j.gpb.2019.11.008 (2020).
- 281 14 Tsou, C. C. *et al.* DIA-Umpire: comprehensive computational framework for data-
282 independent acquisition proteomics. *Nat Methods* **12**, 258-264, 257 p following 264,
283 doi:10.1038/nmeth.3255 (2015).
- 284 15 Ting, Y. S. *et al.* PECAN: library-free peptide detection for data-independent
285 acquisition tandem mass spectrometry data. *Nat Methods* **14**, 903-908,
286 doi:10.1038/nmeth.4390 (2017).
- 287 16 Demichev, V., Messner, C. B., Vernardis, S. I., Lilley, K. S. & Ralser, M. DIA-NN:
288 neural networks and interference correction enable deep proteome coverage in high
289 throughput. *Nat Methods* **17**, 41-44, doi:10.1038/s41592-019-0638-x (2020).
- 290 17 Blattmann, P. *et al.* Generation of a zebrafish SWATH-MS spectral library to quantify
291 10,000 proteins. *Sci Data* **6**, 190011, doi:10.1038/sdata.2019.11 (2019).
- 292 18 Sun, Y. *et al.* Protein Classifier for Thyroid Nodules Learned from Rapidly Acquired
293 Proteotypes. *medRxiv*, 2020.2004.2009.20059741, doi:10.1101/2020.04.09.20059741
294 (2020).
- 295 19 Zhu, Y. *et al.* High-throughput proteomic analysis of FFPE tissue samples facilitates
296 tumor stratification. *Mol Oncol*, doi:10.1002/1878-0261.12570 (2019).
- 297 20 Cai, X. *et al.* PulseDIA: Data-Independent Acquisition Mass Spectrometry Using
298 Multi-Injection Pulsed Gas-Phase Fractionation. *J Proteome Res*,
299 doi:10.1021/acs.jproteome.0c00381 (2020).
- 300 21 Martínez-Aguilar, J., Clifton-Bligh, R. & Molloy, M. P. Proteomics of thyroid
301 tumours provides new insights into their molecular composition and changes
302 associated with malignancy. *Sci Rep* **6**, 23660, doi:10.1038/srep23660 (2016).
- 303 22 Zhang, H. *et al.* Arabidopsis proteome and the mass spectral assay library. *Sci Data* **6**,
304 278, doi:10.1038/s41597-019-0294-0 (2019).
- 305 23 Eid, S., Turk, S., Volkamer, A., Rippmann, F. & Fulle, S. KinMap: a web-based tool
306 for interactive navigation through human kinome data. *BMC Bioinformatics* **18**, 16,
307 doi:10.1186/s12859-016-1433-7 (2017).

308 Acknowledgments

309 **Funding:** This work is supported by grants from National Key R&D Program of China (No.
310 2020YFE0202200), Zhejiang Provincial Natural Science Foundation for Distinguished Young
311 Scholars (LR19C050001), Hangzhou Agriculture and Society Advancement Program
312 (20190101A04), National Natural Science Foundation of China (81972492) and National
313 Science Fund for Young Scholars (21904107). Further support for this project was obtained
314 from the National Cancer Centre Research Fund (Peter Fu Program) and National Medical
315 Research Council Clinician-Scientist Award (NMRC/CSAINV/011/2016).

316 We thank for the for assistance in data storage, computation and peptide fractionation by the
317 Westlake University Supercomputer Center and the Mass Spectrometry & Metabolomics
318 Core Facility at the Center for Biomedical Research Core Facilities of Westlake University.
319

320 **Author contributions**

321 T.G., and Y.S. designed the project. N.G.I. and O.L.K. provided the Singapore set and Y.Z.,
322 G.W., Y.H., and J.X. collected the Chinese set. Y.S., L.L., W.L., Q.X and X.C. performed the
323 experiments. Y.S., L.L., W.G. and H.C. conducted proteomic data analysis. Y.S wrote the
324 manuscript, L.L., Z.D., and F.Z. revised the manuscript. T.G. supervised the project.
325

326 **Competing interests**

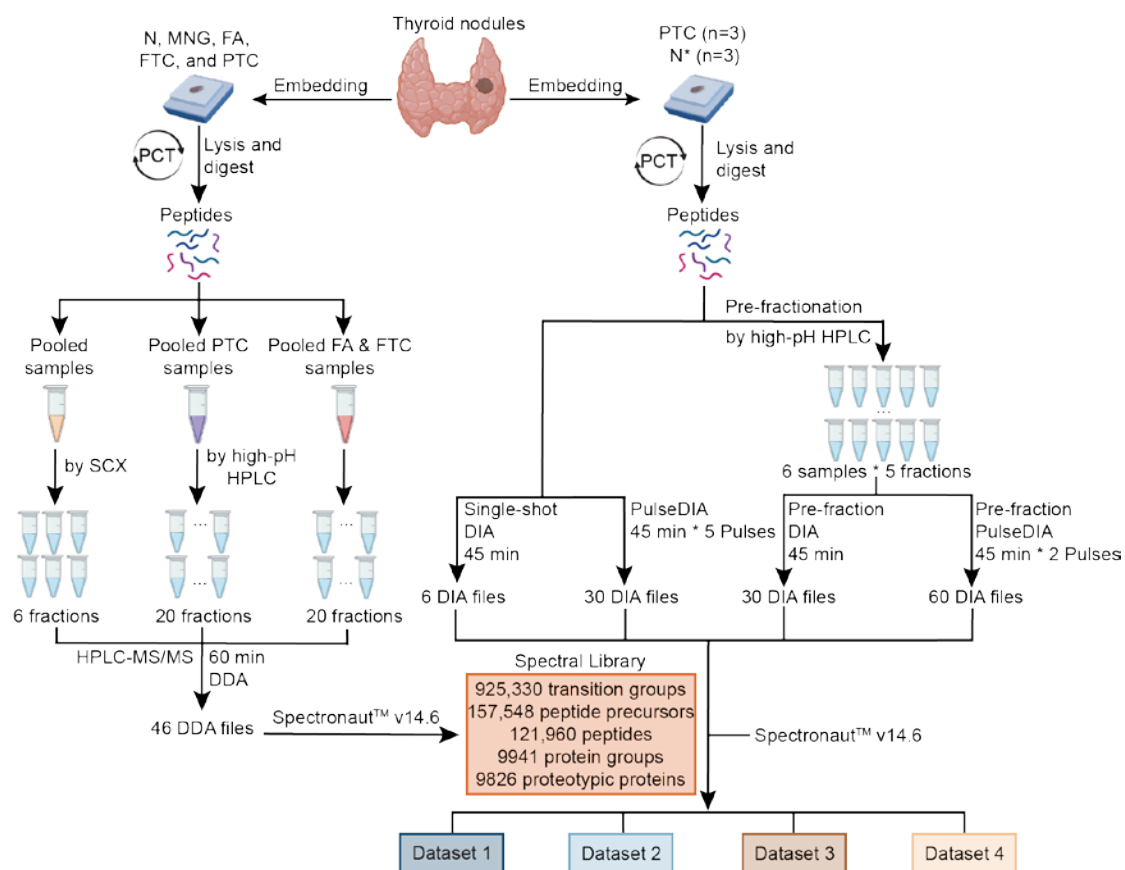
327 The T.G. group is supported by Pressure Biosciences Inc, which provides sample preparation
328 instrumentation including Barocycler and Barozyme. T.G. is a shareholder of Westlake
329 Omics Inc. W.G., W.L. and H.C. are employees of Westlake Omics Inc. The other authors
330 declare no competing interests in this paper.
331

332 **Tables and Figures**

333 **Table 1. Statistics of the thyroid-specific spectral library**

	Library
Transition groups	925,330
Peptide precursors	157,548
Peptides	121,960
Protein groups	9941
Proteotypic proteins	9826

334



335

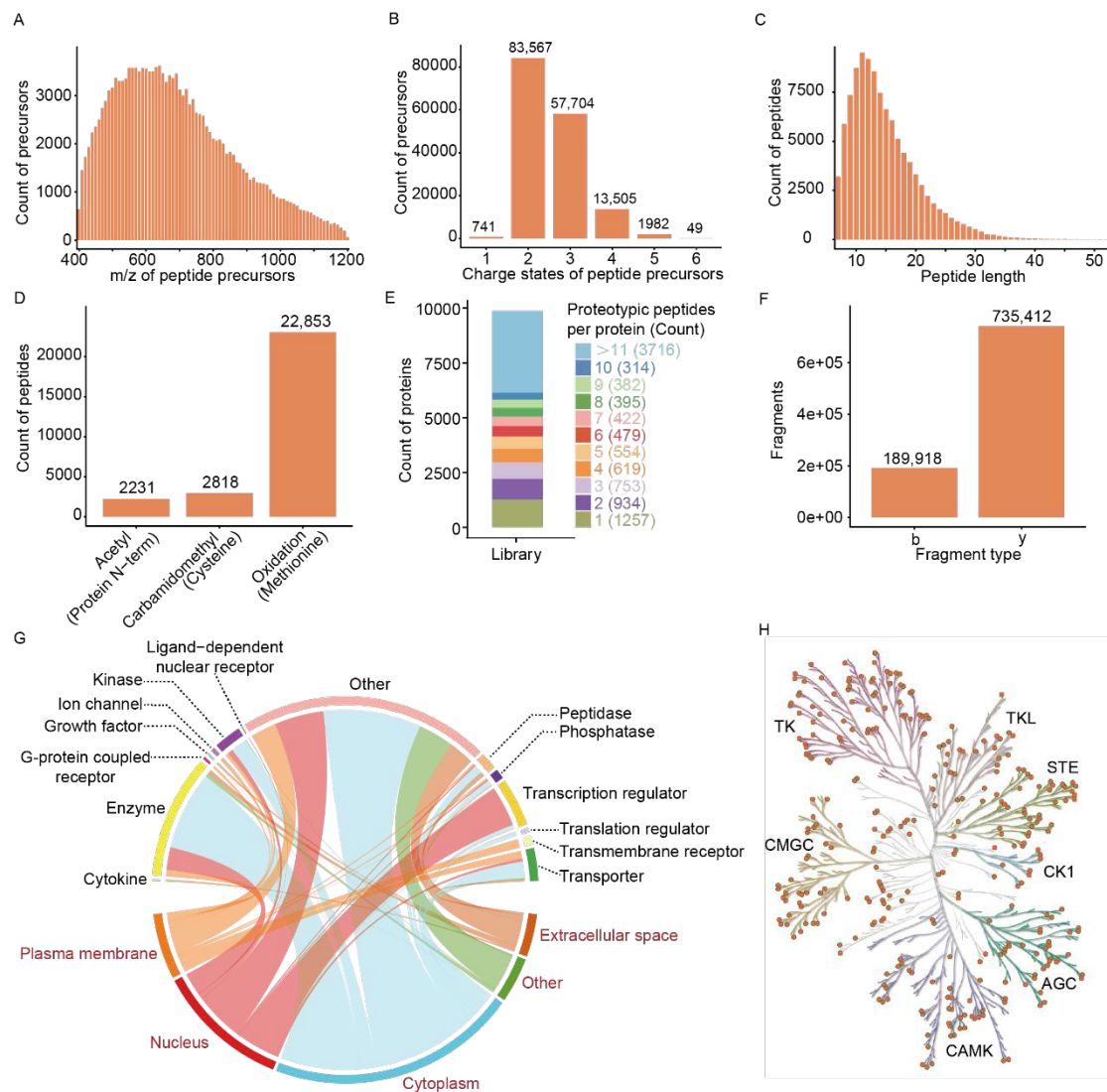
336

Figure 1. Workflow for generating a comprehensive thyroid-specific spectral library (left) and for its validation (right).

337

338

339



340

341

342 **Figure 2. Characterization and statistics of the thyroid-specific spectral library**

343 (A) Distribution of peptide precursor m/z . (B) Counts of different precursor charge states. (C)

344 Distribution of identified peptides lengths. (D) Modified peptides numbers and distribution of

345 three modifications. (E) Numbers of proteotypic peptides for each protein and their

346 corresponding ratios and counts. (F) Ion counts of each fragment type. (G) Proteins were

347 annotated according to two classification systems, subcellular location (words in red) and

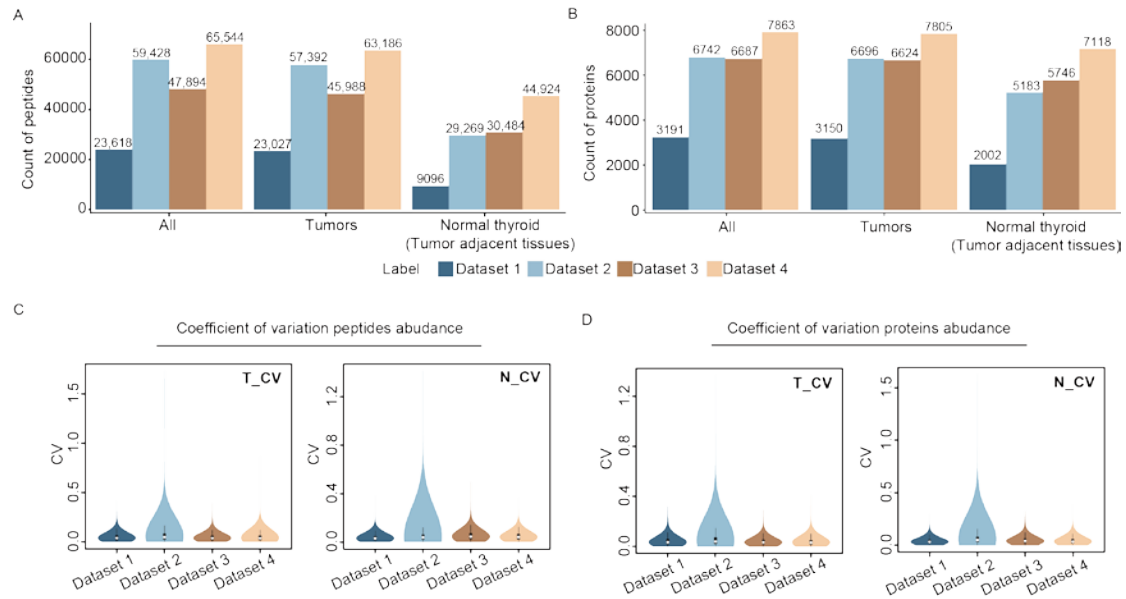
348 function type (words in black). Each curve represents one protein, linking the protein function

349 type with the corresponding subcellular location. (H) A total of 340 kinases (orange dots),

350 belonging to seven families (highlighted by the different tree colors) were identified in our

351 library.

352



353

354 **Figure 3. Results from a technical validation of our thyroid-specific spectral library.**

355 Four datasets were acquired with single-shot DIA (dataset 1), PulseDIA (dataset 2), pre-

356 fraction DIA (dataset 3), and a combination of pre-fraction and PulseDIA (dataset 4).

357 Identified peptides (A) and proteins (B) obtained by searching against our thyroid specific

358 spectral library. Coefficient of variation of peptides (C) and proteins (D) abundance in tumors

359 (T_CV) and their adjacent normal tissues (N_CV).

360

361