1    **Accelerated Protein Biomarker Discovery from FFPE tissue samples**

2    **using Single-shot, Short Gradient Microflow SWATH MS**

3

4    Rui Sun [1,2]*, Christie Hunter [3]*[#], Chen Chen [4], Weigang Ge [1,2], Nick Morrice [3], Shuang
5    Liang[1,2], Chunhui Yuan[1,2], Qiushi Zhang [1,2], Xue Cai [1,2], Xiaoyan Yu [5], Lirong Chen [5],
6    Shaozheng Dai [6], Zhongzhi Luan [6], Ruedi Aebersold [7, 8], Yi Zhu [1,2#], Tiannan Guo [1,2 #]

7

8    1, Key Laboratory of Structural Biology of Zhejiang Province, School of Life Sciences,
9    Westlake University, 18 Shilongshan Road, Hangzhou 310024, Zhejiang Province, China

10   2, Institute of Basic Medical Sciences, Westlake Institute for Advanced Study, 18 Shilongshan
11   Road, Hangzhou 310024, Zhejiang Province, China

12   3, SCIEX, USA and UK (NM)

13   4, SCIEX, China

14   5, Department of Pathology, The Second Affiliated Hospital, Zhejiang University School of
15   Medicine, Hangzhou, Zhejiang, 310009, China

16   6, School of Computer Science and Engineering, Beihang University, Beijing, China

17   7, Department of Biology, Institute of Molecular Systems Biology, ETH Zurich, Switzerland

18   8, Faculty of Science, University of Zurich, Zurich, Switzerland

19   * co-first

20   # co-correspondence

21   **Emails:**

22   Christie Hunter: christie.hunter@sciex.com

23   Yi Zhu: zhuyi@westlake.edu.cn

24   Tiannan Guo: guotiannan@westlake.edu.cn

25
26

27 **ABSTRACT (no more than 4000 characters)**

28 We report and evaluated a microflow, single-shot, short gradient SWATH MS method

29 intended to accelerate the discovery and verification of protein biomarkers in clinical

30 specimens. The method uses 15-min gradient microflow-LC peptide separation, an optimized

31 SWATH MS window configuration and OpenSWATH software for data analysis.

32

33 We applied the method to a cohort 204 of FFPE prostate tissue samples from 58 prostate

34 cancer patients and 10 prostatic hyperplasia patients. Altogether we identified 27,976

35 proteotypic peptides and 4,043 SwissProt proteins from these 204 samples. Compared to a

36 reference SWATH method with 2-hour gradient the accelerated method consumed only 27%

37 instrument time, quantified 80% proteins and showed reduced batch effects. 3,800 proteins

38 were quantified by both methods in two different instruments with relatively high consistency

39 (r = 0.77). 75 proteins detected by the accelerated method with differential abundance

40 between clinical groups were selected for further validation. A shortlist of 134 selected

41 peptide precursors from the 75 proteins were analyzed using MRM-HR, exhibiting high

42 quantitative consistency with the 15-min SWATH method (r = 0.89) in the same sample set.

43 We further verified the capacity of these 75 proteins in separating benign and malignant

44 tissues (AUC = 0.99) in an independent prostate cancer cohort (n=154).

45

46 Overall our data show that the single-shot short gradient microflow-LC SWATH MS method

47 achieved about 4-fold acceleration of data acquisition with reduced batch effect and a

48 moderate level of protein attrition compared to a standard SWATH acquisition method.

49 Finally, the results showed comparable ability to separate clinical groups.

50

51

52 **Keywords**: SWATH MS; Data-independent Acquisition; FFPE; Biomarker; Prostate cancer

53

54    **INTRODUCTION**

55    A large number of clinical and pre-clinical research questions require biomarkers for the
56    classification of samples or phenotypes. Because they are thought to closely reflect the
57    biochemical state of samples, protein biomarkers are particularly valuable. Protein biomarkers
58    have been intensely sought to indicate disease type or stage, to report disease progression or
59    response or resistance to treatment. For the most part protein biomarker projects use mass
60    spectrometry as the base technique. In spite of enormous research efforts, the number of
61    protein biomarkers discovered by proteomic methods that have progressed to clinical utility
62    remains small (1-4).

63    Protein biomarker discovery and validation projects face significant technical and
64    logistical challenges, including the following: i) biological protein abundance variability.
65    Useful protein biomarkers will only be discovered if the variability within a population is
66    smaller than the variability between protein groups. In the context of a twin cohort study of
67    plasma proteins we have shown that the variability of proteins and the root cause for the
68    variability varies greatly in a human population and that particularly variable proteins are
69    unlikely to be selected as biomarkers (5). ii) confounding effects. Protein biomarker studies
70    suffer from a range of confounding effects, including batch effects of sample collection,
71    sample processing, data acquisition and data analysis. Batch effects are particularly severe
72    among different cohorts that might be required to validate results from a discovery cohort, iii)
73    sample availability. Frequently, sample cohorts of sufficient size and quality to generate
74    sufficient statistical power are not available and iv) technical limitations. Even if suitable
75    cohorts are available acquiring reproducible protein patterns by mass spectrometry from
76    extended cohorts has been costly and challenging. For example, typically, protein biomarkers
77    have multi-dimensional fractionation of the peptides generated from digested, tissue-extracted
78    proteins followed by the analysis of the fractions by shotgun MS analysis. Even if isotopic or
79    isobaric labeling methods increase the multiplexing capability of such analyses (6), the
80    general approach remains expensive and technically challenging (7-9). Overall, these
81    challenges convincingly support the need for the proteomic measurement of large sample
82    cohorts at moderate cost, limited batch effects and high degree of reproducibility. At present
83    state-of-the art, large scale clinical proteomic studies consist of 100 to 200 clinical samples
84    (9-11) and there are indication, *e.g* the lack of stability of discovered marker panels that
85    suggest that this number of samples is at the lower end of the required size range(12). Further,
86    these studies were for the most part carried out by highly specialized groups or consortia
87    using highly optimized analytical platforms. For many proteomic research groups that lack
88    the means to implement the involved consortia methods, meaningful protein biomarker
89    studies have therefore remained out of reach. Therefore, there is an urgent need for robust,
90    highly reproducible, high throughput methods that support large-scale biomarker studies at
91    moderate cost and with limited time consumption.

92    Sample throughput can be increased with short LC gradients for the separation of
93    peptides. Bekker-Jensen et.al have combined multiple dimension pre-fractionation with
94    relatively short LC gradient using shotgun proteomics to achieve deep proteome analysis (13);
95    however, this approach lacks reproducibility and it is still time-consuming for a large cohort
96    study. We and others have found that SWATH/DIA mass spectrometry (14) is a more suitable
97    acquisition method to classify samples in large sample cohorts(5, 15-17). SWATH/DIA is an
98    acquisition method for biomarker studies because it identifies and quantifies peptide
99    precursors via peak groups consisting of fragment ion chromatograms from highly convoluted
100    mass spectra (15) and thus obviates the need to isolate peptide precursors during acquisition.
101    This improves data completeness and enables efficient single-shot proteomic analysis. The
102    key to this MS technique is the ability to collect high-resolution MS/MS spectra at very high
103    acquisition rates, such that a wide mass range can be covered with a series of smaller Q1
104    isolation windows in an LC compatible cycle time. Thus the fast scanning rate of TripleTOF
105    system has been the key in enabling the shortening of LC gradients for analyzing complex
106    tissue proteomes, from 120 min (15) to 45 min (17) without strongly compromising proteome

107    depth, and has been increasingly applied to analyze various types of clinical samples
108    including plasma (5) and tumor tissues (15, 17, 18). A faster nano-LC and Orbitrap-based MS
109    method has been reported recently to allow analysis of plasma and cell line samples using a
110    21-min gradient (19). However, this method requires specialized LC system.

111         To further improve the robustness and throughput of the proteomic analysis of sizable
112    sample cohorts, the use of microflow chromatography is a promising option. An increasing
113    number of studies have demonstrated the applicability of microflow coupled with SWATH
114    MS (20-24). E.g. the Ralser group applied microflow-LC and SWATH to study yeast
115    proteome at a throughput of 60 samples per day (24).

116         Here, we established and optimized a 15-min gradient microflow LC SWATH method,
117    and rigorously examined its performance by analyzing 204 FFPE prostate tissue samples.
118    From the detected 4,043 proteins we prioritized 75 that were further verified with respect to
119    their ability to separate cancer from hyperplasia in an independent FFPE prostate tissue
120    sample cohort study by complementary methods. The results indicate the that short gradient
121    microflow-LC SWATH is a suitable and robust method for clinical protein biomarker studies.

122

123    **EXPERIMENTAL PROCEDURES**

124    **Standard protein digests**

125         Digests of proteins isolated from HEK 293 cell were prepared as previously described
126    (25) and provided by Dr Yansheng Liu from ETH Zurich (now in Yale University). Protein
127    digests from K562 cells were obtained from the SWATH Performance Kit (SCIEX). 10% (v/v)
128    iRT peptides (Biognosys, Switzerland) were spiked into peptide samples prior to MS analysis
129    for retention time calibration.

130    **PCa patient cohorts and formalin-fixed paraffin-embedded (FFPE) samples**

131         Two prostate cancer (PCa) sample cohorts termed PCZA and PCZB were used in this
132    study. The PCZA was acquired by the Second Affiliated Hospital College of Medicine,
133    Zhejiang University and consisted of 58 PCa patients and 10 benign prostatic hyperplasia
134    (BPH) patients. The PCZB cohort was acquired by the Second Affiliated Hospital College of
135    Medicine, Zhejiang University and consisted of 24 PCa patients and 30 BPH patients whose
136    benign and hyperplastic regions have been distinguished. All patients were recruited in 2017
137    and 2018. All cohorts were approved by the ethics committee of the respective hospitals for
138    the procedures of this study.

139         The two different cohort samples were handled by different pathology laboratories, fixed
140    and embedded by the respective staff. The samples were similarly processed and analyzed at
141    different time points. For the PCZA cohort, three biological replicates (size $1 \times 1 \times 5$ mm$^3$)
142    were collected and analyzed by SWATH MS and MRM-HR. For the PCZB cohort, two
143    biological replicates ($1.5 \times 1.5 \times 5$ mm$^3$) were analyzed by MRM-HR.

144    **Pressure cycling technology (PCT)-assisted peptide extraction from FFPE tissues**

145         About 0.5 mg of FFPE tissue was punched from the samples, weighed and processed for
146    each biological replicate via the FFPE-PCT workflow as described previously (26). Briefly,
147    the tissue punches were first dewaxed by incubating with 1 mL of heptane under gentle
148    vortexing at 600–800 rpm, followed by serial rehydration using 1 mL of 100%, 90%, and 75%
149    ethanol (General reagent, G73537B, Shanghai, China), respectively. The samples were further
150    incubated with 200 µL of 0.1% formic acid (FA) (Thermo Fisher Scientific, T-27563) at 30 °C
151    for 30 min for acidic hydrolysis. The tissue punches were then transferred into
152    PCT-MicroTubes (Pressure Biosciences Inc., Boston, MA, USA, MT-96) and briefly washed
153    with 100 µL of freshly prepared 0.1 M Tris-HCl (pH 10.0) to remove residual FA. Thereafter,
154    the tissues were incubated with 15 µL of freshly prepared 0.1 M Tris-HCl (pH 10.0) at 95 °C

155    for 30 m in with gentle vortexing at 600 rpm. Samples were immediately cooled to 4 °C after
156    basic hydrolysis.

157    Following the pretreatment described above, 25 µL of lysis buffer including 6 M urea
158    (Sigma, U1230), 2 M thiourea (Amresco, M226) in 100 mM ammonium bicarbonate (General
159    regent, G12990A, Shanghai, China), pH 8.5 was added to the PCT-MicroTubes containing
160    tissues. The tissue samples were further subjected to PCT-assisted tissue lysis and protein
161    digestion procedures using the Barocycler NEP2320-45K (Pressure Biosciences Inc., Boston,
162    MA, USA) as described previously (27). The PCT scheme for tissue lysis was set with each
163    cycle containing a period of 30 s of high pressure at 45 kpsi and 10 s at ambient pressure,
164    oscillating for 90 cycles at 30°C. Protein reduction and alkylation was performed at ambient
165    pressure by incubating protein extracts with 10 mM Tris(2-carboxyethyl) phosphine (TCEP)
166    (Sigma, C4706) and 20 mM iodoacetamide (IAA) (Sigma, I6125) in darkness at 25 °C for 30
167    min, with gentle vortexing at 600 rpm in a thermomixer. Then the proteins were digested with
168    MS grade Lys-C (Hualishi, Beijing, China, enzyme-to-substrate ratio, 1:40) using a PCT
169    scheme with 50 s of high pressure at 20 kpsi and 10 s of ambient pressure for each cycle,
170    oscillating for 45 cycles at 30 °C. Thereafter, the proteins were further digested with MS
171    grade trypsin (Hualishi, Beijing, China, enzyme-to-substrate ratio, 1:50) using a PCT scheme
172    with 50 s of high pressure at 20 kpsi and 10 s of ambient pressure in one cycle, oscillating for
173    90 cycles at 30 °C. Peptide digests were then acidified with 1% trifluoroacetic (TFA) (Thermo
174    Fisher Scientific, T/3258/PB05) to pH 2–3 and subjected to C18 desalting. iRT peptides were
175    spiked into peptide samples at a final concentration of 10% prior to MS analysis for RT
176    calibration.

**Optimization of microflow LC gradients coupled with SWATH MS**

178    During the optimization studies, 1 µg peptides were separated with different microflow
179    gradients and different SWATH MS parameters. Linear gradients of 3–35% acetonitrile (0.1%
180    formic acid) with durations of 5, 10, 20, 30, and 45 min were evaluated. The number of Q1
181    variable windows (40, 60, 100) and MS/MS accumulation times (15, 25 ms) constituted the
182    key parameters that were adjusted for the shorter gradients. The need for collision energy
183    spread with the optimized collision energy ramps was tested. Four replicates were performed
184    for each test, after which the data were processed with the PeakView® software with the
185    SWATH 2.0 MicroApp to evaluate the number of proteins and peptides quantified with FDR
186    < 1 % and CV < 20%. The optimized methods were then tested on multiple instruments with
187    different cell lysates to confirm the robustness of the method.

**SWATH MS acquisition**

189    Peptides were separated at a flow rate of 5 µL/min by a 15-min SWATH of 5–35% linear
190    LC gradient elution (buffer A: 2% ACN (Sigma, 34851), 0.1% formic acid; buffer B: 80%
191    ACN, 0.1% formic acid) on a column, 3 µm, ChromXP C18CL, 120 Å, 150 x 0.3 mm using
192    an Eksigent NanoLC$^{TM}$ 400 System coupled with a TripleTOF® 6600 system (SCIEX). The
193    DuoSpray Source was replumbed using the 25 µm ID hybrid electrodes to minimize
194    post-column dead volume. The applied SWATH method was composed of a 150 ms TOF MS
195    scan with m/z ranging from 350 to 1250 Da, followed by MS/MS scans performed on all
196    precursors (from 100 to 1500 Da) in a cyclic manner. A 100 variable Q1 isolation window
197    scheme was used in this study (Supplementary Table 1B). The accumulation time was set at
198    25 ms per isolation window, resulting in a total cycle time of 2.7 s.

199    We also included beta-galactosidase digest (β-gal) (SCIEX, 4465867) for mass and
200    retention time calibration which was analyzed every four injections. The target ion (m/z =
201    729.4) which is from a peptide precursor in the β-gal digest mixture was monitored under
202    high sensitivity mode. The RT, intensity, and m/z of targeted precursor and fragment ions
203    were respectively used for LC QC, the sensitivity test, and mass calibration separately.

**MRM-HR MS acquisition**

205        A time scheduled MRM-HR targeted quantification strategy was used to further validate
206  proteins observed to be differentially expressed proteins by SWATH MS as described above.
207  The same microflow LC approach was used for 15-min SWATH MS analysis. The TripleTOF
208  6600 mass spectrometer was operated in IDA mode for time-scheduling the MS/MS
209  acquisition for 134 peptides for the MRM-HR workflow. The method consisted of one 75 ms
210  TOF-MS scan for precursor ions with m/z ranging from 350 to 1250 Da, followed by MS/MS
211  scans for fragment ions with m/z ranging from 100 to 1500 Da, allowing for a maximum of
212  45 candidate ions being monitored per cycle (25 ms accumulation time, 50 ppm mass
213  tolerance, rolling collision energy, +2 to +5 charge states with intensity criteria above 2 000
214  000 cps to guarantee that no untargeted peptides should be acquired). The fragment ion
215  information including m/z and RT of a targeted precursor ion was confirmed by previous
216  SWATH results and was then added to the inclusion list for the targeted analysis. The intensity
217  threshold of targeted precursors in the inclusion list was set to 0 cps and the scheduling
218  window was 60 s. The targeted peptide sequences were the same as those found in the
219  previous SWATH MS analysis.

220        Targeted MRM-HR data were analyzed by 19.0.9.149 Skyline (28), which automatically
221  detected the extracted-ion chromatogram (XIC) from an LC run by matching the MS spectra
222  of the targeted ion against its spectral library generated from the IDA mode within a specific
223  mass tolerance window around its m/z. All peaks selected were checked manually after
224  automated peak detection using Skyline. Both MS1 and MS2 filtering were set as "TOF mass
225  analyzer" with a resolution power of 30 000 and 15 000, respectively, while the "Targeted"
226  acquisition method was defined in the MS/MS filtering.

227  **SWATH data analysis**

228        The optimization data for optimal LC gradients were processed using the SWATH 2.0
229  MicroApp in PeakView® software (SCIEX) using the pan-human library (29). RT calibration
230  was performed by first using iRT peptides with RT window at a 75 ppm XIC extraction width.
231  Replicate analysis was performed using the SWATH Replicate Analysis Template (SCIEX) to
232  determine the number of peptides and proteins quantified with FDR < 1% peptide and CV <
233  10 or 20%.

234        The data from prostate samples were processed using the OpenSWATH pipeline. Briefly,
235  SWATH raw data files were converted in profile mode to mzXML using msconvert and
236  analyzed using OpenSWATH (2.0.0) (30) as described previously (15). The retention time
237  extraction window was 600 s, while m/z extraction was performed with 0.03 Da tolerance. RT
238  was then calibrated using both iRT peptides. Peptide precursors were identified by
239  OpenSWATH and PyProphet (2.0.1) with d_score < 0.01 and FDR < 1%. For each protein,
240  the median MS2 intensity value of peptide precursor fragments which were detected to belong
241  to the protein was used to represent the protein abundance.

242

243  **RESULTS AND DISCUSSIONS**

244  **Establishment and optimization of the 15-min microflow SWATH MS method**

245        A HEK 293 cell lysate digest was used to establish and optimize the short microflow LC
246  gradient and SWATH acquisition schemes on TripleTOF 6600 systems. Specifically, we tested
247  the effects of LC gradient lengths of 5, 10, 20, and 45 min, and mass spectrometer parameters
248  including variable Q1 windows and accumulation time for MS2 (Supplementary Table 1). For
249  each injection 1 µg mass of total peptide was loaded onto a microflow column of 150 x 0.3
250  mm dimensions and analyzed under a range of conditions. To increase robustness of results,
251  four technical replicates of each condition were used. The acquired data were searched
252  consistently searched against PHL with the PeakView® software and the SWATH 2.0
253  MicroApp and the number of peptides and inferred proteins, as well as their intensities were
254  recorded. The data was processed as described in the methods section and evaluated

255     according to the number of proteins and peptides identified with FDR < 1% and quantified
256     with CV < 10% or CV < 20%, respectively. The whole dataset was acquired on two different
257     instruments. Supplementary Figure 1a shows that using shorter gradient methods generated
258     similar results between the two different 6600 instruments. Our data also showed that the 20
259     min microflow method detected 90% of the proteins quantified by 45 min method, while the
260     10 min LC method identified 70% of the proteins. With decreasing gradient length, the
261     number of identified proteins further decreased to 77% for a 10 min method to 53% for a 5
262     min method in their best condition (Supplementary Figure 1a).

263     Next, we optimized the specific mass spec parameters including variable windows and
264     accumulation times to balance the width of the windows and scan times (Supplementary
265     Figure 1b). Typically, more variable windows led to more peptides and proteins quantified
266     robustly, but only up to a point where the MS/MS acquisition rates become too fast or the
267     cycle times too long, as evidenced in the 5min gradient optimization results. Thus, a higher
268     number of variable windows led to a higher number of peptide and protein identifications.
269     The optimal accumulation time was highly dependent on the LC time. Higher numbers of
270     acquisition windows necessitated shorter MS/MS accumulation times per precursor ion
271     window to maintain a cycle time that was compatible with the peak width generated by the
272     respective gradients. Considering the tradeoff between sample throughput and numbers of
273     proteins quantified, the gradient time from 10 min to 20 min is a better choice according to
274     the efficiency of peptides and proteins identification in unit of time (Supplementary Figure
275     1b). Therefore, we chose the 15 min gradient as the optimal LC condition (Supplementary
276     Figure 2).

277     **Application of short gradient microflow-SWATH to a PCa patient cohort**

278     We evaluated the performance of the optimized short gradient microflow LC SWATH
279     method on a set of prostate cancer (PCa) tissue samples named PCZA. The set consisted of
280     204 FFPE biospecimens collected from 58 PCa patients and 10 benign hyperplasia (BPH)
281     patients (Supplementary Table 2) for which clinical data were also available. The 204 samples
282     were randomly divided into seven batches and digested into peptides in barocyclers. Every
283     batch included a mouse liver sample as quality control (QC) for the PCT-assisted sample
284     preparation and a prostate tissue pool samples as the QC sample for SWATH MS
285     (Supplementary Figure 3).
286     We then subjected the resulting peptide samples to the15-min-SWATH method optimized
287     above (Figure 1a). The total sample set consisted of 58 PCa samples and 7 QC and reference
288     samples. The 204 samples were measured in 125.7 hrs (~5 days) and quantified 27,975
289     peptide precursors from 4,038 SwissProt proteins (without protein grouping) with 74.79%
290     missing value rate. On average, 5,615 peptide precursors from 1,018 proteins were quantified
291     for each sample. More peptides and proteins were quantified from tumor samples (5,861
292     peptide precursors from 1,078 proteins on average) than benign samples (3,988 peptide
293     precursors from 618 proteins on average). Totally 913 proteins were quantified in at least 50%
294     samples (Supplementary Table 3).
295     To allow a comparison of the accelerated short gradient method with a standard SWATH
296     MS method with respect to the number of proteins recorded and the respective clinically
297     relevant information content we re-acquired the whole sample set with a 120-min LC gradient
298     and 48 variable Q1 windows in a TripleTOF 5600+(26). These measurements consumed 467
299     hr (~20 days) and identified 38,338 peptide precursors from 5,059 SwissProt proteins with
300     61.86% missing value rate. On average, 10,751 peptide precursors from 1,921 proteins were
301     quantified for each sample. More peptides and proteins were quantified from tumor samples
302     (11,439 peptide precursors from 2,054 proteins on average) than benign samples (6,693
303     peptide precursors from 1192 proteins on average). Totally 1,914 proteins were quantified in
304     at least 50% samples (Supplementary Table 3). Compare to this 120-min method, the 15-min
305     method characterized about half of peptide precursors and proteins.
306

307       Overall, the data shows that the 15-min-SWATH coverage reached 50-80% of that
308 achieved by a standard method. In all samples, 3,800 proteins were quantified by both
309 methods. This result was generated at a 6-fold reduced acquisition time (time 125.7 hrs vs,
310 467 hrs) (Figure 1b) suggesting that clinical cohorts of significant size can be measured by the
311 accelerated method quickly, efficiently.

312 **Reproducibility and batch effect analysis**
313 We evaluated the reproducibility of the datasets produced by the 15-min gradient and the 120-
314 min gradient SWATH with respect to reproducibility and batch effect. We first assessed the
315 technical reproducibility by correlation between technical replicates for LC-MS. The
316 technical reproducibility of the data obtained by the 15-min SWATH method (r = 0.99) is
317 slightly higher than that from the 120-min SWATH method (r = 0.86) (Figure 2a). Thus, the
318 measured biological reproducibility is also slightly higher in the 15-min SWATH method
319 (Figure 2a). If we focused the analysis on the 3,800 proteins quantified by both methods, we
320 observed a high degree of similarity (r = 0.7681) between the methods (Figure 2b).
321       We next analyzed batch effects apparent in either dataset. Batch effects are an
322 unavoidable reality resulting from technical variation in multi-day MS analyses and are a
323 non-trivial complication for big cohort proteomics analysis. Several algorithms have been
324 developed to bioinformatically minimize the missing value rate, however, these imputation
325 approaches remain controversial (31). We evaluated the batch effect of the data acquired by
326 the 15-min SWATH, which is lower than that from the 120-min method (Figure 2c). Together,
327 the 15-min SWATH method improved quantitative reproducibility and reduced batch effect.
328

329 **Verification of differential expression proteins using MRM-HR**

330 On the path to clinical or preclinical use protein biomarkers detected by MS based cohort
331 studies face a number of verification and validation requirements. These include technical
332 verification of the abundance changes detected in the cohort study and validation in
333 independent sample cohorts.

334       To further validation the abundance changes detected in the SWATH data we selected a
335 panel of 75 proteins showing different abundance (absolute fold change larger than two and
336 adjusted p-value less than 0.05) between control and tumor tissue and measured their
337 respective intensities using the targeted MS method MRM-HR. The selected proteins were
338 associated with most strongly cancer dis regulated pathways and included 21 known
339 diagnosis biomarkers such as ACPP and FASN, and 10 drug targets (Supplementary Table
340 4A). The proteins were further annotated in IPA (Supplementary Table 4B) indicating that the
341 proteins suggested elevated cell migration, development and growth, and suppressed cell
342 death and survival (Supplementary Figure 5).

343       For these measurements the MRM-HR method was optimized using a pooled prostate
344 sample to determine the best performing peptides from the selected proteins, and best target
345 fragment ions for quantitation. The information about proteins and peptides including the RTs
346 was imported into Skyline to build a spectral library. A total of 134 peptides for 75 proteins
347 were selected for targeted detection (Supplementary Table 2E). Time scheduling was used to
348 ensure at least eight data points were obtained across the LC peaks as well as an optimized
349 accumulation time of 25 ms for each peptide for high-quality quantitative data.

350       To confirm the quantitative accuracy of the 15-min SWATH data, we re-analyzed 99
351 samples in the PCZA cohort using the MRM-HR method. The 99 samples were randomly
352 allocated to five batches, each containing 20 samples and an extra MS QC sample which was
353 a pool of prostate tissue digests in PCZA. We firstly examined the reproducibility of XICs for
354 all peptides in MRM-HR assays. For the five pooled samples measured across five batches,
355 we found that 76.6% of precursors measured from the peptides were quantified with a CV
356 below 20%. The median CV was 13.4% (Supplementary Figure 6). Next the protein
357 fold-changes between tumor and normal samples were calculated to investigate the

358    correlation of 15-min SWATH with MRM-HR (Figure 3a).

359    We further quantified the expression levels of the 75-protein-panel in an independent
360    prostate cancer cohort, PCZB, containing 30 BPH and 24 PCa in duplicated biological
361    replicates using the same 15-min SWATH MRM-HR workflow (Supplementary Table 5). For
362    the six pooled samples measured across six batches, 75.6% of peptide precursors were
363    quantified with a CV below 20%. The median CV is 14.9% (Supplementary Figure 6).

364    To assess the power of the protein panel of differentially abundant proteins to separate
365    benign and malignant tissues, we assembled a random-forest model for the PCZA MRM-HR
366    dataset, and found an accuracy of 0.992 in this set (Supplementary Figure 7). Next, we tested
367    the ability of this panel to separate tumor from benign prostatic tissue samples in an
368    independent patient group, *i.e.* PCZB, including 24 PCa patients and 30 BPH patients. The
369    receiver operating curves (ROC) of the 75-protein-panel clearly distinguished PCa from BPH
370    patient groups (Figure 3b).

371    We then investigated in detail two proteins—PRDX3 (P30048) and COPA (P53621) which
372    were prioritized because of their role in TP53 oncogene regulation and as a potential drug
373    (decitabine) target (Supplementary Figure 8). The data show that these proteins significantly
374    up-regulated in tumor tissue from all three workflows, *i.e.*15-min-SWATH, and MRM-HR in
375    the PCZA cohort samples and MRM-HR in the PCZB cohort (Figure 3b). The ROC curve of
376    these two proteins from three different datasets distinguishing benign from malignant tissue
377    samples are shown in Supplementary Figure 9, with all of AUC over 0.78. Taken together, we
378    validated these dysregulated proteins quantification by SWATH showed higher reliability and
379    performed better prediction ability in different sample cohorts.

380

## Conclusion

382    In this study, we present a 15-min microflow-LC SWATH that supports the consistent
383    proteomic analysis of clinical (FFPE) samples at a throughput of ~50 samples per day
384    (excluding calibration and washing). The method is therefore well suited for the analysis of
385    large sample cohorts, even in a single investigator proteomic laboratory. The results show that
386    the presented method increases the throughput by ca six-fold compared to a conventional
387    SWATH MS method, at reduced batch effects and at an attrition of ca 20% of detected
388    proteins and increased missing value rate (~20% worse) in the prostate cancer cohort. For
389    individual samples, the number of detected proteins decreased by ~50%. The quantitative
390    accuracy of the short gradient method was comparable to that achieved by targeted
391    quantification using MRM-HR for shortlisted proteins. This work showed the potential of this
392    short gradient SWATH proteomics pipeline for accelerated discovery and verification of
393    protein biomarkers for precision medicine.

## Author Contributions

406

416

417 **ABBREVIATIONS**

418 AUC = area under the curve

419 BPH = benign prostatic hyperplasia

420 CV = coefficient of variation

421 DDA = data dependent acquisition

422 DIA = data independent acquisition

423 FA = formic acid

424 FDR = false discovery rate

425 FFPE = formalin fixed, paraffin embedded

426 IAA = iodoacetamide

427 IPA = ingenuity pathway analysis

428 LC = liquid chromatograph

429 MRM-HR = multiple reaction monitoring high-resolution

430 PCa = prostate cancer

431 PCT = pressure cycling technology

432 PRM = parallel reaction monitoring

433 QC = quality control

434 ROC = receiver operating characteristic

435 RF = random forest

436 RT = retention time

437 SWATH MS = sequential windowed acquisition of all theoretical fragment ion - mass spectra

438 TCEP = tris(2-carboxyethyl) phosphineTFA = trifluoroacetic

439 TMA = tissue microarray analysis

440 TOF = time of flight

441 XIC = extracted ion chromatogram

442

443 **References**

444 1. Olsen, M.; Ghannad, M.; Lok, C.; Bossuyt, P. M., Shortcomings in the evaluation of

445 biomarkers in ovarian cancer: a systematic review. *Clin Chem Lab Med* **2019**.

446 2. Rifai, N.; Gillette, M. A.; Carr, S. A., Protein biomarker discovery and validation: the long

447 and uncertain path to clinical utility. *Nat Biotechnol* **2006**, 24, (8), 971-83.

448    3.    Anderson, N. L.; Ptolemy, A. S.; Rifai, N., The riddle of protein diagnostics: future bleak or

449    bright? *Clin Chem* **2013,** 59, (1), 194-7.

450    4.    Frantzi, M.; Latosinska, A.; Kontostathi, G.; Mischak, H., Clinical Proteomics: Closing the

451    Gap from Discovery to Implementation. *Proteomics* **2018,** 18, (14), e1700463.

452    5.    Liu, Y.; Buil, A.; Collins, B. C.; Gillet, L. C.; Blum, L. C.; Cheng, L. Y.; Vitek, O.; Mouritsen,

453    J.; Lachance, G.; Spector, T. D.; Dermitzakis, E. T.; Aebersold, R., Quantitative variability of

454    342 plasma proteins in a human twin population. *Mol Syst Biol* **2015,** 11, (1), 786.

455    6.    Aebersold, R.; Mann, M., Mass-spectrometric exploration of proteome structure and

456    function. *Nature* **2016,** 537, (7620), 347-55.

457    7.    Sabrkhany, S.; Kuijpers, M. J. E.; Knol, J. C.; Olde Damink, S. W. M.; Dingemans, A. C.;

458    Verheul, H. M.; Piersma, S. R.; Pham, T. V.; Griffioen, A. W.; Oude Egbrink, M. G. A.; Jimenez,

459    C. R., Exploration of the platelet proteome in patients with early-stage cancer. *J Proteomics*

460    **2018,** 177, 65-74.

461    8.    Mun, D. G.; Bhin, J.; Kim, S.; Kim, H.; Jung, J. H.; Jung, Y.; Jang, Y. E.; Park, J. M.; Kim,

462    H.; Jung, Y.; Lee, H.; Bae, J.; Back, S.; Kim, S. J.; Kim, J.; Park, H.; Li, H.; Hwang, K. B.; Park,

463    Y. S.; Yook, J. H.; Kim, B. S.; Kwon, S. Y.; Ryu, S. W.; Park, D. Y.; Jeon, T. Y.; Kim, D. H.; Lee,

464    J. H.; Han, S. U.; Song, K. S.; Park, D.; Park, J. W.; Rodriguez, H.; Kim, J.; Lee, H.; Kim, K. P.;

465    Yang, E. G.; Kim, H. K.; Paek, E.; Lee, S.; Lee, S. W.; Hwang, D., Proteogenomic

466    Characterization of Human Early-Onset Gastric Cancer. *Cancer Cell* **2019,** 35, (1), 111-124

467    e10.

468    9.    Zhang, H.; Liu, T.; Zhang, Z.; Payne, S. H.; Zhang, B.; McDermott, J. E.; Zhou, J. Y.;

469    Petyuk, V. A.; Chen, L.; Ray, D.; Sun, S.; Yang, F.; Chen, L.; Wang, J.; Shah, P.; Cha, S. W.;

470    Aiyetan, P.; Woo, S.; Tian, Y.; Gritsenko, M. A.; Clauss, T. R.; Choi, C.; Monroe, M. E.;

471    Thomas, S.; Nie, S.; Wu, C.; Moore, R. J.; Yu, K. H.; Tabb, D. L.; Fenyo, D.; Bafna, V.; Wang,

472    Y.; Rodriguez, H.; Boja, E. S.; Hiltke, T.; Rivers, R. C.; Sokoll, L.; Zhu, H.; Shih, I. M.; Cope, L.;

473    Pandey, A.; Zhang, B.; Snyder, M. P.; Levine, D. A.; Smith, R. D.; Chan, D. W.; Rodland, K. D.;

474    Investigators, C., Integrated Proteogenomic Characterization of Human High-Grade Serous

475    Ovarian Cancer. *Cell* **2016,** 166, (3), 755-765.

476    10.  Vasaikar, S.; Huang, C.; Wang, X.; Petyuk, V. A.; Savage, S. R.; Wen, B.; Dou, Y.; Zhang,

477    Y.; Shi, Z.; Arshad, O. A.; Gritsenko, M. A.; Zimmerman, L. J.; McDermott, J. E.; Clauss, T. R.;

478    Moore, R. J.; Zhao, R.; Monroe, M. E.; Wang, Y. T.; Chambers, M. C.; Slebos, R. J. C.; Lau, K.

479    S.; Mo, Q.; Ding, L.; Ellis, M.; Thiagarajan, M.; Kinsinger, C. R.; Rodriguez, H.; Smith, R. D.;

480    Rodland, K. D.; Liebler, D. C.; Liu, T.; Zhang, B.; Clinical Proteomic Tumor Analysis, C.,

481    Proteogenomic Analysis of Human Colon Cancer Reveals New Therapeutic Opportunities.

482    *Cell* **2019,** 177, (4), 1035-1049 e19.

483    11.  Mertins, P.; Mani, D. R.; Ruggles, K. V.; Gillette, M. A.; Clauser, K. R.; Wang, P.; Wang,

484    X.; Qiao, J. W.; Cao, S.; Petralia, F.; Kawaler, E.; Mundt, F.; Krug, K.; Tu, Z.; Lei, J. T.; Gatza,

485    M. L.; Wilkerson, M.; Perou, C. M.; Yellapantula, V.; Huang, K. L.; Lin, C.; McLellan, M. D.; Yan,

486    P.; Davies, S. R.; Townsend, R. R.; Skates, S. J.; Wang, J.; Zhang, B.; Kinsinger, C. R.; Mesri,

487    M.; Rodriguez, H.; Ding, L.; Paulovich, A. G.; Fenyo, D.; Ellis, M. J.; Carr, S. A.; Nci, C.,

488    Proteogenomics connects somatic mutations to signalling in breast cancer. *Nature* **2016,** 534,

489    (7605), 55-62.

490    12. Thomas. S.; Friedrich, B., ; Schnaubelt M.; Chan. D.; Zhang H,; Aebersold. R.,

491    Orthogonal proteomic platforms and their implications for the stable classification of

492      high-grade serous ovarian cancer subtypes. *BioRxiv* **2019**.

493      13. Bekker-Jensen, D. B.; Kelstrup, C. D.; Batth, T. S.; Larsen, S. C.; Haldrup, C.; Bramsen, J.

494      B.; Sorensen, K. D.; Hoyer, S.; Orntoft, T. F.; Andersen, C. L.; Nielsen, M. L.; Olsen, J. V., An

495      Optimized Shotgun Strategy for the Rapid Generation of Comprehensive Human Proteomes.

496      *Cell Syst* **2017**, 4, (6), 587-599 e4.

497      14.  Gillet LC, N. P., Tate S, Röst H, Selevsek N, Reiter L, Bonner R, Aebersold R., Targeted

498      data  extraction  of  the  MS/MS  spectra  generated  by  data-independent  acquisition:  a  new

499      concept for consistent and accurate proteome analysis. *Mol Cell Proteomics.* **2012**.

500      15.  Guo, T.; Kouvonen, P.; Koh, C. C.; Gillet, L. C.; Wolski, W. E.; Rost, H. L.; Rosenberger,

501      G.; Collins, B. C.; Blum, L. C.; Gillessen, S.; Joerger, M.; Jochum, W.; Aebersold, R., Rapid

502      mass spectrometric conversion of tissue biopsy samples into permanent quantitative digital

503      proteome maps. *Nat Med* **2015**, 21, (4), 407-13.

504      16.  Bouchal,  P.;  Schubert,  O.  T.;  Faktor,  J.;  Capkova,  L.;  Imrichova,  H.;  Zoufalova,  K.;

505      Paralova,  V.;  Hrstka,  R.;  Liu,  Y.;  Ebhardt,  H.  A.;  Budinska,  E.;  Nenutil,  R.;  Aebersold,  R.,

506      Breast Cancer Classification Based on Proteotypes Obtained by SWATH Mass Spectrometry.

507      *Cell Rep* **2019**, 28, (3), 832-843 e7.

508      17. Zhu, Y.; Zhu, J.; Lu, C.; Zhang, Q.; Xie, W.; Sun, P.; Dong, X.; Yue, L.; Sun, Y.; Yi, X.; Zhu,

509      T.; Ruan, G.; Aebersold, R.; Huang, S.; Guo, T., Identification of Protein Abundance Changes

510      in Hepatocellular Carcinoma Tissues Using PCT-SWATH. *Proteomics Clin Appl* **2019**, 13, (1),

511      e1700179.

512      18.  Guo, T.; Li, L.; Zhong, Q.; Rupp, N. J.; Charmpi, K.; Wong, C. E.; Wagner, U.; Rueschoff,

513      J. H.; Jochum, W.; Fankhauser, C. D.; Saba, K.; Poyet, C.; Wild, P. J.; Aebersold, R.; Beyer, A.,

514    Multi-region proteome analysis quantifies spatial heterogeneity of prostate tissue biomarkers.

515    *Life Sci Alliance* **2018,** 1, (2).

516    19.  Bache N1, G. P., 3, Bekker-Jensen DB3, Hoerning O1, Falkenby L1, Treit PV2, Doll S2,

517    Paron I2, Müller JB2, Meier F2, Olsen JV3, Vorm O1, Mann M, A Novel LC System Embeds

518    Analytes in Pre-formed Gradients for Rapid, Ultra-robust Proteomics. *Mol Cell Proteomics.*

519    **2018.**

520    20.  Shi, J.; Wang, X.; Lyu, L.; Jiang, H.; Zhu, H. J., Comparison of protein expression

521    between human livers and the hepatic cell lines HepG2, Hep3B, and Huh7 using SWATH and

522    MRM-HR proteomics: Focusing on drug-metabolizing enzymes. *Drug Metab Pharmacokinet*

523    **2018,** 33, (2), 133-140.

524    21.  He, B.; Shi, J.; Wang, X.; Jiang, H.; Zhu, H. J., Label-free absolute protein quantification

525    with data-independent acquisition. *J Proteomics* **2019,** 200, 51-59.

526    22.  Colgrave, M. L.; Byrne, K.; Blundell, M.; Heidelberger, S.; Lane, C. S.; Tanner, G. J.;

527    Howitt, C. A., Comparing Multiple Reaction Monitoring and Sequential Window Acquisition of

528    All Theoretical Mass Spectra for the Relative Quantification of Barley Gluten in Selectively

529    Bred Barley Lines. *Anal Chem* **2016,** 88, (18), 9127-35.

530    23.  Le Duff, M.; Gouju, J.; Jonchere, B.; Guillon, J.; Toutain, B.; Boissard, A.; Henry, C.;

531    Guette, C.; Lelievre, E.; Coqueret, O., Regulation of senescence escape by the

532    cdk4-EZH2-AP2M1 pathway in response to chemotherapy. *Cell Death Dis* **2018,** 9, (2), 199.

533    24.  Vowinckel, J.; Zelezniak, A.; Bruderer, R.; Mulleder, M.; Reiter, L.; Ralser, M.,

534    Cost-effective generation of precise label-free quantitative proteomes in high-throughput by

535    microLC and data-independent acquisition. *Sci Rep* **2018,** 8, (1), 4346.

536    25.   Liu, Y.; Mi, Y.; Mueller, T.; Kreibich, S.; Williams, E. G.; Van Drogen, A.; Borel, C.; Frank,

537    M.; Germain, P. L.; Bludau, I.; Mehnert, M.; Seifert, M.; Emmenlauer, M.; Sorg, I.; Bezrukov, F.;

538    Bena, F. S.; Zhou, H.; Dehio, C.; Testa, G.; Saez-Rodriguez, J.; Antonarakis, S. E.; Hardt, W.

539    D.; Aebersold, R., Multi-omic measurements of heterogeneity in HeLa cells across laboratories.

540    *Nat Biotechnol* **2019,** 37, (3), 314-322.

541    26.   Yi Zhu 1, 3*,Tobias Weiss 4*, Qiushi Zhang 1,2, Rui Sun 1,2, Bo Wang 5, Zhicheng Wu

542    1,2, Qing Zhong 6,7, Xiao Yi 1,2 , Huanhuan Gao 1,2, Xue Cai 1,2, Guan Ruan 1,2, Tiansheng

543    Zhu 1,2, Chao Xu  , Sai Lou 9, Xiaoyan Yu 10, Ludovic Gillet 3, Peter Blattmann 3, Karim Saba

544    11, Christian D.Fankhauser 11, Michael B. Schmid 11, Dorothea Rutishauser 6, Jelena

545    Ljubicic 6, Ailsa , Christiansen 6, Christine Fritz 6, Niels J. Rupp 6, Cedric Poyet 11, Elisabeth

546    Rushing 12, Michael Weller 4, Patrick Roth 4, Eugenia Haralambieva 6, Silvia Hofer 13, Chen

547    Chen 14, Wolfram Jochum 15, Xiaofei Gao 1,2, Xiaodong Teng 5, Lirong Chen 10, Peter J.

548    Wild 6,16#, Ruedi Aebersold 3,17# , Tiannan Guo, High-throughput proteomic analysis of

549    FFPE tissue samples facilitates tumor stratification. *biorxiv* **2019**.

550    27.   Zhu, Y.; Guo, T., High-Throughput Proteomic Analysis of Fresh-Frozen Biopsy Tissue

551    Samples Using Pressure Cycling Technology Coupled with SWATH Mass Spectrometry.

552    *Methods Mol Biol* **2018,** 1788, 279-287.

553    28.   MacLean, B.; Tomazela, D. M.; Shulman, N.; Chambers, M.; Finney, G. L.; Frewen, B.;

554    Kern, R.; Tabb, D. L.; Liebler, D. C.; MacCoss, M. J., Skyline: an open source document editor

555    for creating and analyzing targeted proteomics experiments. *Bioinformatics* **2010,** 26, (7),

556    966-8.

557    29.   Rosenberger, G.; Koh, C. C.; Guo, T.; Rost, H. L.; Kouvonen, P.; Collins, B. C.; Heusel,

558    M.; Liu, Y.; Caron, E.; Vichalkovski, A.; Faini, M.; Schubert, O. T.; Faridi, P.; Ebhardt, H. A.;

559    Matondo, M.; Lam, H.; Bader, S. L.; Campbell, D. S.; Deutsch, E. W.; Moritz, R. L.; Tate, S.;

560    Aebersold, R., A repository of assays to quantify 10,000 human proteins by SWATH-MS. *Sci*

561    *Data* **2014,** 1, 140031.

562    30.   Rost, H. L.; Rosenberger, G.; Navarro, P.; Gillet, L.; Miladinovic, S. M.; Schubert, O. T.;

563    Wolski, W.; Collins, B. C.; Malmstrom, J.; Malmstrom, L.; Aebersold, R., OpenSWATH enables

564    automated, targeted analysis of data-independent acquisition MS data. *Nat Biotechnol* **2014,**
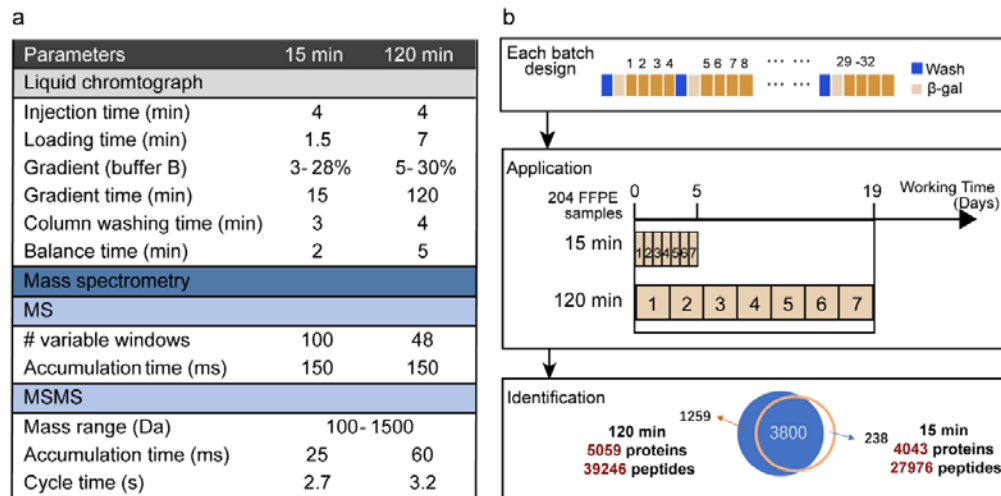
565    32, (3), 219-23.

566    31.   Goh, W. W. B.; Wong, L., Advanced bioinformatics methods for practical applications in

567    proteomics. *Brief Bioinform* **2019,** 20, (1), 347-355.
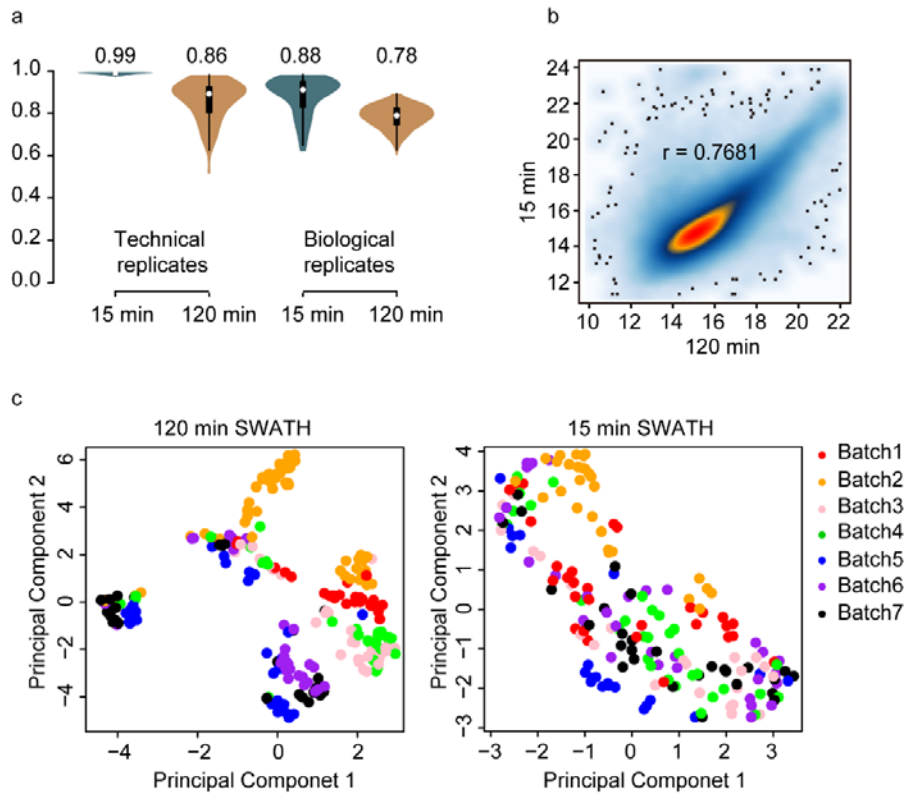
568

569

570     Figure 1. The short-gradient SWATH method and application in the PCZA PCa cohort. (a)
571     Comparison of parameters between the 15-min SWATH and 120-min SWATH methods. (b)
572     The workflow of the 15-min-SWATH and conventional SWATH for the PCZA cohort. We
573     designed seven randomly shuffled batches with a column washing run and a calibration (β-gal)
574     run inserted every four samples.



575

576

577    Figure 2. The reproducibility of the short gradient SWATH method in the PCZA PCa cohort.
578    (a) Violin plots show the technical replicates and biological replicates in the two methods. (b)
579    Pearson correlation of log2-scaled protein intensity values obtained from 3,800 proteins that
580    were quantified by both methods. (c) PCA analysis of all samples quantified by the 120-min
581    method (left) and the 15-min method (right).



582

583

584 Figure 3. Verification of proteomic data using MRM-HR. (a) Pearson correlation coefficient
585 between the 15-min SWATH and MRM-HR datasets based on the $\log_2(T/N)$ of protein
586 expression in PCZA. (b) The ROC curves of protein quantification from MRM-HR to predict
587 the tumor and normal tissues with the random forest algorithm in PCZB (T: PCa, N: BPH, H:
588 hyperplasia in BPH patients, B: benign in BPH patients). (c) MRM-HR validation of potential
589 diagnostic proteins using the PCZA and PCZB. PRDX3 (peptide: +2 DYGVLLEGSGLALR),
590 COPA (peptide: +2 DVAVMQLR). The left panel shows the fragment ion extracted-ion
591 chromatograms (XICs) for the peptide from each protein. The right panel of boxplots shows
592 the peptides quantified in the different data sets.



593