Perspective pubs.acs.org/jpr

50

51

## Next Steps on in Silico 2DE Analyses of Chromosome 18 <sup>2</sup> Proteoforms

<sup>3</sup> Stanislav N. Naryzhny,<sup>\*,†,‡</sup> Elena S. Zorina,<sup>†</sup> Arthur T. Kopylov,<sup>†</sup> Victor G. Zgoda,<sup>†</sup> Olga A. Kleyst,<sup>‡</sup> 4 and Alexander I. Archakov<sup>†</sup>

s <sup>†</sup>Institute of Biomedical Chemistry of Russian Academy of Medical Sciences, Pogodinskaya 10, Moscow 119121, Russia

6<sup>‡</sup>Petersburg Nuclear Physics Institute, National Research Center "Kurchatov Institute", Leningrad Region, Gatchina 188300,

Russia 7

22

Supporting Information 8

ABSTRACT: In the boundaries of the chromosome-centric Human Proteome 9 Project (c-HPP) to obtain information about proteoforms coded by chromosome 10 18, several cell lines (HepG2, glioblastoma, LEH), normal liver, and plasma were 11 analyzed. In our study, we have been using proteoform separation by two-dimensional 12 electrophoresis (2DE) (a sectional analysis) and a semivirtual 2DE with following 13 shotgun mass spectrometry using LC-ESI-MS/MS. Previously, we published a first 14 draft of this research, where only HepG2 cells were tested. Here, we present the 15 next step using more detailed analysis and more samples. Altogether, confident 16 17 (2 significant sequences minimum) information about proteoforms of 117 isoforms coded by 104 genes of chromosome 18 was obtained. The 3D-graphs showing 18 distribution of different proteoforms from the same gene in the 2D map were 19 generated. Additionally, a semivirtual 2DE approach has allowed for detecting 20 more proteoforms and estimating their pI more precisely. Data are available via 21 ProteomeXchange with identifier PXD010142.



KEYWORDS: proteome, inventory, proteoforms, chromosome 18, two-dimensional electrophoresis, mass spectrometry, 23 ESI LC-MS/MS, chromosome-centric 24

## 25 INTRODUCTION

26 The final goal of the "Human Proteome Project" (HPP) is 27 complete knowledge about all human proteins. Keeping in 28 mind the huge volume of information that needs to be dis-29 closed and efforts to be done, a special project was launched. 30 This project is chromosome-centric (c-HPP) to promote more  $_{31}$  effective collaborations<sup>1-3</sup> and to create a comprehensive knowl-32 edge base of all human proteins for use as a tool in different areas 33 such as education, science, or medicine that need this infor-34 mation. Russia is responsible for chromosome 18, which 35 harbors 275 protein coding genes. Since these proteins are 36 presented as multiple variants/proteoforms/protein species, we 37 expect to detect at least 3000 proteoforms generated from 38 chromosome 18. Usage of an interactive virtual 2DE map of 39 proteins coded by genes of chromosome 18 in combination 40 with experimental 2DE data will allow more effective execution 41 of the Russian part of c-HPP. As protein extracts of normal as 42 well as cancer cells are used in our study, information about 43 proteoform distribution (proteome) will allow better under-44 standing of protein changes connected to carcinogenesis. It will 45 allow the identification of characteristic quantitative and qual-46 itative differences between normal and cancer cells for use as 47 biomarkers or targets for therapy. Here, we represent our data 48 extracted from analysis of fibroblasts (LEH), HepG2, 49 glioblastoma, human liver, and plasma.

## EXPERIMENTAL SECTION

## **Chemicals and Materials**

All reagents used were obtained from Sigma-Aldrich (St. Louis, 52 MO, USA), unless another manufacturer is specified. The 53 remaining reagents were obtained from the following com- 54 panies: Thermo Fisher Scientific (Waltham, MA, USA): 55 dithiothreitol (DTT), protease inhibitor cocktail, penicillin, 56 streptomycin; GE Healthcare (Pittsburgh, PA, USA): IPG 57 DryStrip (gel strips), IPG-buffers, DryStrip-coating liquid, 58 Coomassie R350; Promega (Madison, WI, USA): Trypsin 59 Gold; Bio-Rad (Hercules, CA, USA) molecular weight markers 60 for protein electrophoresis; Biolot (Moscow, Russia): DMEM/ 61 F12 for cell growth, fetal calf serum, Trypsin-EDTA solution; 62 Orange Scientific (Braine-l'Alleud, Belgium): Carrel culture 63 flasks. 64

#### Cells

Human hepatocellular carcinoma (HepG2), lung embryonic 66 fibroblasts (LEH), and glioblastoma cells were cultured in 67 medium (DMEM/F12 supplemented with 10% fetal bovine 68 serum (FBS) and 100 U/mL penicillin/streptomycin) under 69

Special Issue: Human Proteome Project 2018

Received: May 29, 2018 Published: September 21, 2018

70 standard conditions (5% CO<sub>2</sub>, 37 °C). To prepare cell samples 71 for protein extraction, the cells were detached with 0.25% 72 Trypsin-EDTA solution, washed 3 times with PBS, and treated 73 by lysis buffer.<sup>4,5</sup> Liver tissue samples of people who died as a 74 result of an accident were purchased from the ILSBioBiobank 75 within the framework of collaboration on the Chromosome-76 Centric Human Proteome Project (c-HPP). All ethical require-77 ments have been met. Extraction was performed by lysis buffer 78 after grinding the tissue in liquid nitrogen according to 2-DE 79 protocol described in the literature.<sup>6,7</sup>

## 80 Sample Preparation and Two-Dimensional Electrophoresis 81 (2DE)

<sup>82</sup> Samples were prepared as described previously.<sup>8,9</sup> Cells ( $\sim 10^7$ ) 83 containing 2 mg of protein were treated by 100  $\mu$ L of lysis 84 buffer (7 M urea, 2 M thiourea, 4% CHAPS, 1% dithiothreitol  $_{85}$  (DTT), 2% (v/v) ampholytes, pH 3–10, protease inhibitor 86 cocktail). The protein concentration in the sample was deter-87 mined by the method of Bradford.<sup>10</sup> 2DE was performed using 88 immobilized pH gradient (IPG)<sup>11</sup> for isoelectric focusing 89 (IEF). Proteins were separated by IEF using Immobiline 90 DryStrip 3-11 NL, 7 cm (GE Healthcare) following the manu-91 facturer's protocol. The samples in the lysis buffer were mixed 92 with rehydrating buffer (7 M urea, 2 M thiourea, 2% CHAPS, 93 0.3% DTT, 0.5% IPG (v/v) buffer, pH 3-11 NL, 0.001% 94 bromophenol blue) in final volume of 125  $\mu$ L (250–300  $\mu$ g of 95 protein)/strip. Strips (4 per experiment) were passively 96 rehydrated for 4 h at 4 °C. IEF was performed on an IPGphor (GE Healthcare), which was programmed as follows: first step, 97 98 300 V, 1 h; second step, gradient to 1000 V, 1 h, third step; 99 gradient to 5000 V, 1.5 h; fourth step, 5000 V, 1 h, temperature 100 20 °C and maintained at the voltage 500 V. After IEF, strips <sup>101</sup> were soaked 10 min in the equilibration solution (50 mM Tris, 102 pH 8.8, 6 M urea, 2% sodium dodecyl sulfate (SDS) and 30% 103 (v/v) glycerol, 1% DTT). This process was followed by 10 min 104 incubation in the equilibration solution containing iodoaceta-105 mide instead of DTT (50 mM Tris, pH 8.8, 6 M urea, 2% 106 sodium dodecyl sulfate (SDS) and 30% (v/v) glycerol, 5% 107 iodoacetamide). The strips were placed on the top of the 12% 108 polyacrylamide gel of the second direction and sealed with a 109 hot solution of 1 mL of 0.5% agarose prepared in electrode 110 buffer (25 mM Tris, pH 8.3, 200 mM glycine, and 0.1% SDS) 111 and electrophoresed to second direction under denaturing 112 conditions using the system Ettan DALTsix  $(240 \times 200 \times 1 \text{ mm})$ 113 GE Healthcare). Electrophoresis was carried out at room tem-114 perature at constant power 3.5 W per gel.<sup>12,13</sup> Gels were 115 stained with Coomassie Blue R350, scanned by ImageScanner 116 III (GE Healthcare) or Typhoon FLA 9500 (GE Healthcare), 117 and analyzed using ImageMaster 2D Platinum 7.0 (GE 118 Healthcare).

119 In the case of a semivirtual 2DE, the 24 cm IPG strip was cut 120 into 48 equal sections using scissors, and each section was 121 transferred to Eppendorf tube. For complete reduction, 300  $\mu$ L 122 of 3 mM DTT, 100 mM ammonium bicarbonate was added to 123 each sample and incubated at 50 °C for 15 min. For alkylation, 124 20 µL of 100 mM iodoacetamide (IAM) were added to the 125 same tube and incubated in the dark at room temperature for 126 15 min. For digestion, stock solution of trypsin (0.1 mg/mL) 127 was diluted 1:10 by 25 mM ammonium bicarbonate, and 100  $\mu$ L 128 of diluted trypsin was added into each tube. Samples were 129 incubated overnight for 4–24 h at 37 °C. Supernatants that 130 may contain peptides that have diffused out of the gel slices 131 were collected to new labeled 0.5 mL tubes. Peptides were

extracted by adding 150 µL of 60% acetonitrile, 0.1% 132 trifluoroacetic acid (TFA) to each tube containing gel slices. 133 Extracts were dried in Speed Vac, reconstituted in 20  $\mu$ L of 134 0.1% TFA, and analyzed by Orbitrap Q-Exactive Plus mass 135 spectrometer. Protein identification and relative quantification 136 were performed using Mascot 2.4.1 (Matrix Science) and emPAI. 137 A table with information about all detected protein proteoforms 138 was built. All proteins detected in the same section were given 139 the pI of this section. Accordingly, the same proteins detected 140 in different sections were considered as different proteoforms. 141 142

### ESI LC-MS/MS Analysis

All procedures were performed according to the protocol 143 described previously.<sup>13</sup> The gel was divided into 96 sections 144 with determined coordinates. Each section ( $\sim 0.7 \text{ cm}^2$ ) was cut, 145 shredded, and treated by trypsin according to the protocol for 146 single spots identification with proportionally increased 147 volumes of solutions. Tryptic peptides were eluted from the 148 gel with extraction solution (5% (v/v) ACN, 5% (v/v) formic 149 acid) and dried in a vacuum centrifuge. Peptides were 150 dissolved in 5% (v/v) formic acid. Tandem mass spectrometry 151 analysis was carried out in duplicate on an Orbitrap Q-Exactive 152 mass spectrometer (Thermo Scientific, USA). Mass spectra 153 were acquired in positive ion mode. High resolution data was 154 acquired in the Orbitrap analyzer with a resolution of 30 000 155 (m/z 400) for MS and 7500 (m/z 400) for MS/MS scans. 156

Identification of proteins was performed using Mascot 2.4.1 157 (Matrix Science, London, UK) by searching UniProt/Swiss- 158 Protein sequence database (October 2014, 20196 total 159 sequences). The following search parameters were used: trypsin 160 as the cutting enzyme, mass tolerance for the monoisotopic 161 peptide window was set to ±50 ppm, missed cleavages 1. Cyste- 162 ine carbamidomethyl was chosen as a fixed and oxidized methi- 163 onine as a variable modification. NeXtProt database was used 164 as a protein sequence database. For FDR assessment, a separate 165 decoy database was generated from the protein sequence 166 database. False positive rate of 1% was allowed for protein 167 identification. A minimum Mascot ion score of 30 was used for 168 accepting peptide MS/MS spectra. Data were also searched, 169 using the SearchGUI, an open-source graphical user inter- 170 face.<sup>14</sup> Two unique peptides per protein were required for all 171 protein identifications. Exponentially modified PAI (emPAI), 172 the exponential form of protein abundance index (PAI) 173 defined as the number of identified peptides divided by the 174 number of theoretically observable tryptic peptides for each 175 protein, was used to estimate protein abundance.<sup>1</sup> 176

## RESULTS AND DISCUSSION

## **High Yield Proteome Analysis**

177 178

We ran the same type of 2DE (7 cm  $\times$  8 cm) using different 179 human samples (HepG2, LEH, liver, glioblastoma, plasma) 180 (Figure S1). Each gel was stained and scanned, and the 181 produced experimental 2DE maps were calibrated according to 182 the position of the several known protein spots. Next, the gels 183 were identically divided into 96 sections, identified as 1-12 184 along the  $M_{\rm w}$  dimension and A–H along the pI dimension. All 185 these gel sections were chopped and treated with trypsin 186 according to protocol for mass spectrometry by LC-ESI-MS/ 187 MS. The tryptic peptides obtained from each 2DE section were 188 analyzed using an Orbitrap Q-Exactive mass spectrometer. 189 Finally, protein identification and relative quantification were 190 performed using Mascot 2.4.1 or SearchGUI and Peptide 191 Shaker. If the same protein (product of the same gene) was 192





Figure 1. continued



Figure 1. continued

Perspective





Figure 1. continued

Ε





Figure 1. continued



Figure 1. 3D graphs showing profiles of proteoforms of the same gene (isoform) in liver, HepG2, LEH, glioblastoma cells, or plasma. Proteins coded by chromosome 18 are shown. (A) AFG3L2 (Q9Y4W6-1), AFG3-like protein 2. (B) ATP5A1 (P25705-1), ATP synthase subunit alpha, mitochondrial. (C) CNDP2 (Q96KP4-1), cytosolic nonspecific dipeptidase. (D) TUBB6 (Q9BUF5-1), tubulin beta-6 chain. (E) MYL12B (O14950-1), myosin regulatory light chain 12B. (F) ACAA2 (P42765-1), 3-ketoacyl-CoA thiolase, mitochondrial isoform.

193 identified in different sections, it was considered to exist as 194 different proteoforms. These forms have different  $pI/M_w$ 195 parameters (coordinates) and are in different areas (spots) of 196 the gel. After analysis of all sections and protein identification, 197 a table of proteins coded by chromosome 18 and detected in 198 extracts from several origins was generated (Table S1). Finally, 199 the 3D graphs showing patterns of proteoforms (profiles) of 200 the same isoform (unique UniProt number) coded by the 201 same gene of 18th chromosome were generated for all samples 202 analyzed (Figure S2). Some of these profiles are shown in 203 Figure 1A-F. These graphical images give us a general repre-204 sentation of these proteoform profiles in different cells or plasma. 205 The examples of genes presented are AFG3L2, ATP5A1, 206 CNDP2, TUBB6, MYL12B, ACAA2. Using virtual-experimen-207 tal approach we can actually try to explain these profiles and 208 proteoforms.<sup>13</sup> In these graphs, a master protein (product of a 209 nonmodified canonical sequence) is presented usually as a 210 most abundant proteoform. There are some exceptions of the 211 case, when polypeptide is processed into mature form by cleav-212 age of the signal sequence, for instance. It happens with 213 mitochondrial proteins that have a mitochondrial targeting 214 signal (1-43 AA), which should be usually cleaved once a 215 targeting is complete. In the case of chromosome 18, there are 216 several such proteins: AFG3-like protein 2 (AFG32 HUMAN), 217 ATP synthase subunit alpha (ATPA HUMAN), ferrochelatase 218 (HEMH HUMAN), (NADH dehydrogenase [ubiquinone] 219 flavoprotein 2 (NDUV2 HUMAN), 3-ketoacyl-CoA thiolase 220 (THIM HUMAN). Accordingly, in the case of mitochondrial 221 proteins, main proteoforms observed are not master but 222 mature forms. We often do not see a master protein in their 223 profiles at all (AFG32, HEMH, NDUV2) (Figure 1A) or see it 224 at a very low abundance (ATPA HUMAN) (Figure 1B). There 225 are some exceptions as well. For instance, mitochondrial transit 226 peptide (1-16) in 3-ketoacyl-CoA thiolase (THIM\_HUMAN) 227 is supposed to be not cleaved. Accordingly, in liver and HepG2 228 cells, we see that its main form has the same parameters 229 (pI/ $M_w$ , 8.32/42 000) as the master protein (Figure 1F). But 230 in LEH and glioblastoma cells, the most abundant peak of 231 THIM has more acidic pI, may be because of the signal 232 peptide cleavage (Figure 1F). An alternative reason for this

shift could be heavy acetylation or phosphorylation. In concor- 233 dance with this assumption, we can find in Proteoform Atlas 234 (http://atlas.topdownproteomics.org/) that it has already 235 detected 3 sites of THIM acetylation<sup>16</sup> 236

The sectional analysis, which we are performing, actually has <sup>237</sup> a low resolution (0.5–0.8 pH range per section), and single <sup>238</sup> PTMs (like acetylation, phosphorylation...) that produce shift <sup>239</sup> in the range of 0.05–0.5 pH unit/per PTM usually cannot be <sup>240</sup> differentiated here.<sup>17–19</sup> But in the case of multiple PTMs, the <sup>241</sup> corresponding proteoforms should be easily observed. Also, <sup>242</sup> some single PTMs like ubiquitination (~8.5 kDa, a single <sup>243</sup> shift), SUMOylation (~12 kDa, a single shift), glycosylation <sup>244</sup> (up to several kDa shift)<sup>20,21</sup> should be easily come out. These <sup>245</sup> could be reasons for the observation of proteoforms with  $M_w$  <sup>246</sup> bigger than theoretically predicted ones, like for serpin B3 <sup>247</sup> (SPB3\_HUMAN) or serpin B4 (SPB4\_HUMAN) in the case <sup>248</sup> of liver and HepG2 (Figure S2). <sup>249</sup>

Additionally, smaller proteoforms that are most likely 250 proteolytic products are often observed (ATPA, CNDP2, 251 DSC1) (Figure S2). One interesting example is AFG3-like protein 252 2 (AFG32\_HUMAN) (Figure 1A). The theoretical parameters 253  $pI/M_w$  of its mature mitochondrial form are 8.01/81184. But 254 in all observed samples (liver, HepG2, LEH, glioblastoma), we 255 see only acidic half-size proteoforms (pH 5–7,  $M_{\rm w} \sim 40\,000$ ) 256 (Figure 1A). The half-size proteoform with acidic pI  $\sim$  5.5 257 (C-terminus part) could be produced by cleavage of AFG32 by 258 endopeptidase somewhere in the middle of the chain. But in 259 this case, we should also see a N-terminus part with a basic pI 260 as well (pI  $\sim$  9.4). As we do not see this polypeptide, it could 261 be completely hydrolyzed by another proteinase. What is 262 interesting, a very similar situation is also observed with mito- 263 chondrial ferrochelatase (HEMH HUMAN). A mature form 264 of this protein should have  $pI/M_w$ : 8.61/42186.67. This 265 proteoform is observed only in liver cells, but in HepG2, LEH, 266 and glioblastoma cells, we see only half-size proteoform with pI 267  $5-6/M_{\rm w}\sim 21000$ . Taking together, we can assume that this 268 could be a result of an unknown type of processing for these 269 proteins. Complementing the theme of processing, very 270 interesting and unexpected results were obtained with human 271 plasma samples. Majority of proteins coded by chromosome 18 272



Perspective

Figure 2. A semivirtual 2DE of proteins coded by chromosome 18. (A) Liver. (B) HepG2, adapted from ref 31. (C) LEH. (D) Glioblastoma. (E) Plasma.

273 are presented in plasma as proteoforms that are not master or 274 mature forms (Figure S2). Their profiles in plasma are very 275 different from the profiles in the cellular extracts (Figure 1B–E). 276 Most likely, these proteins are not bearing any function in 277 plasma and coming to the bloodstream as a byproduct from 278 different cells. Some of them could be a product of proteolysis 279 (smaller forms) but some (bigger forms) possibly are heavily 280 modified (Figure S2). An example of a protein coded by chromo-281 some 18 that is functional in plasma is beta-Al-His dipeptidase 282 (CNDP1 HUMAN), also known as serum carnosinase. This enzyme is abundant in blood serum, where it plays an impor- 283 tant role in hydrolysis of dietary carnosine.<sup>22,23</sup> Another protein, 284 transthyretin (TTHY\_HUMAN) is a serum and cerebrospinal 285 fluid protein that transports holo-retinol-binding protein and 286 thyroxine. Its serum concentration has been widely used to 287 assess clinical nutritional status.<sup>24</sup> 288

We should keep in mind that there is a possibility that some 289 of the observed truncated products can be produced during 290 2DE. 2DE is a long multistage procedure, and there is a chance 291 of extra modification of proteins during it. Previously, we had 292



Figure 3. continued



Figure 3. Comparison of a sectional analysis and a semivirtual 2DE for some proteins coded by chromosome 18. TXNL1 (thioredoxin-like protein), AFG32 (AFG3-like protein 2), SNAG (gamma-soluble NSF attachment protein), UBP14 (ubiquitin carboxyl-terminal hydrolase 14), TBB6 (tubulin beta-6 chain), ATPA (ATP synthase subunit alpha, mitochondrian).

293 been analyzing this situation on an example of PCNA 294 (proliferation cell nuclear antigen).<sup>25</sup> We found that its 295 proteasomal degradation can take place during 2DE (IEF) if 296 inhibitors (thiourea or proteasomal inhibitor) are absent.<sup>25,26</sup> 297 A positive side of this situation is that this proteasomal 298 degradation happens in vivo as well and may be relevant to the 299 physiological functions of proteins. Increased in vitro pro-300 teolysis that is detected using 2DE may be a kind of ampli-301 fication of the proteolysis that happens in vivo.<sup>25</sup> Though we 302 make every effort to exclude this situation using thiourea and 303 protease inhibitors cocktail, we cannot guarantee that it is not 304 completely happening in our experiments.<sup>26</sup> It seems that this 305 is a general situation with 2DE, which has manifested itself 306 after using sensitive large-scale mass spectrometry (ESI-LC-307 MS/MS).<sup>27</sup> As like with PCNA, a similar situation is observed 308 with ATPA (Figure 1B). It looks like the observed truncated 309 proteoforms of ATPA are produced by endopeptidase 310 (proteasome). According to detected peptides, we can even recon-311 struct ATPA proteoforms located in different sections. For 312 instance, ATPA proteoform located in section B12 highly 313 likely is AA59–168 polypeptide (theoretical  $pI/M_w$ : 4.57/ 314 11309), in section C12, AA2–168 ( $pI/M_w$ : 5.64/17253), in 315 section G7, N-terminus peptide AA24–317 ( $pI/M_{w}$ : 8.43/31495) 316 and C-terminus peptide AA253-554 (pI/M<sub>w</sub>: 7.93/32826), in  $_{317}$  section F9, N-terminus peptide AA46–260 (pI/M<sub>w</sub>: 7.73/22999) 318 and C-terminus peptide AA329- 554 ( $pI/M_w$  7.05/24468). To obtain more information about proteoforms, we also 320 performed a semivirtual 2DE, where high resolution of IEF was

implemented in a better degree than in the 2DE sectional anal- 321 ysis.<sup>28</sup> In the case of 24 cm strips with pH 3–11 (48 sections), 322 a pH range of each section taken for analysis is ~0.166 and 323 ~0.21 in the case of 18 cm strips (36 sections). Though 324 semivirtual 2DE does not give the value of experimental  $M_{\rm w}$  of 325 detected proteoforms, it allows to detect more of them and 326 measure their pI more precisely (Figure 2A-F, Table S2, 327 Figure S3A-F). As a complement, it is improving a resolution 328 of the previously performed sectional analysis. Let's see some 329 examples. For instance, in the case of AFG3-like protein 2 330 (AFG32 HUMAN), by 2DE sectional analysis of HepG2 331 proteins only one major peak in the section with coordinates 332 pI 5.11-5.80 and  $M_{\rm w}$  35 000-40 000 was detected. But using 333 semivirtual approach, seven proteoforms with pI around 334 5.8 were detected (Figure 3). For thioredoxin-like protein 335 1 (TXNL1 HUMAN), four proteoforms instead of one peak 336 around pI 5.0 were detected. In the case of tubulin beta-6 chain 337 (TBB6 HUMAN), 4 proteoforms around pI 5.0 were detected. 338 In the case of myosin regulatory light chain 12B (ML12B HU- 339 MAN), 6 proteoforms around pI 4.7 were detected. For 340 gamma-soluble NSF attachment protein isoform (SNAG HU- 341 MAN), two proteoforms around pI 5.5 were detected, and 342 for ubiquitin carboxyl-terminal hydrolase 14 isoform (UB- 343 P14 HUMAN), six proteoforms around pI 5.8 were detected 344 (Figure 3). In the case of mitochondrial ATP synthase subunit 345 alpha (ATPA\_HUMAN), the situation is more complicated. 346 A resolution of the semivirtual 2DE in pH direction is much 347 higher (its limit is 46 fractions against only 8 fractions in the 348

<sup>349</sup> semivirtual 2DE). But many proteoforms detected by the <sup>350</sup> sectional analysis are truncated products of ATPA and have the <sup>351</sup> similar pI (especially in the range 7.82–8.86) (Figure 3). <sup>352</sup> Accordingly, they are not separated by the semivirtual 2DE and <sup>353</sup> detected in the same fractions. This is the reason why the <sup>354</sup> amounts of detected by both methods proteoforms (44 and <sup>355</sup> 45) are similar. But the actual amount is higher. For instance, <sup>366</sup> instead of only 2 peaks ( $M_w$  6000–15 000) with acidic pH <sup>357</sup> (3.7–5.11) detected by the sectional analysis, the semivirtual <sup>358</sup> 2DE gives 13 proteoforms (Figure 3). It will rise the number of <sup>359</sup> proteoforms at least to 56. This more detailed set of proteo-<sup>360</sup> forms makes it possible to more accurately attribute to them <sup>361</sup> the information obtained about posttranslational modifications <sup>362</sup> (PTMs) (Table S2).

In summary, if the sectional analysis and the semivirtual 2DE act are used together, more detailed proteoform patterns can be obtained. First, information about proteoforms (emPAI and pI) act is extracted from the semivirtual 2DE. Next, additional information about the experimental  $M_w$  of proteoforms can be added from the sectional analysis.

## 369 CONCLUSION

370 This is our next step in the creation of a 2DE-based knowledge 371 database of proteins coded by chromosome 18 that was carried 372 out using small 2DE gels (7 cm  $\times$  8 cm). Such a scale and 373 division of the gels into 96 sections allowed to generate 374 proteoform profiles coded not only by the 18th chromosome 375 but by other chromosomes as well. A representation of 376 proteomes based on theoretical and experimental parameters  $_{377}$  (pI/ $M_{\rm w}$ ) allowed to get a clearer view of the general state at the 378 scale of complete proteomes and products of single genes. 379 In particular, using this way of representation, we have shown a 380 striking difference between specific profiles of proteoforms in 381 cellular and plasma proteomes in terms of PTMs. Addition of a 382 semivirtual 2DE allowed to produce more detailed profiles of 383 proteoforms coded by each gene, and a graphical representa-384 tion of the profiles gives a chance to get information about 385 profiles in a convenient way. Also, parameters  $(pI/M_w)$  of 386 detected proteoforms give us as a hint what kind of PTMs can 387 produce these proteoforms. But the next steps should be done 388 to disclose the fine details of proteoforms, including their 389 PTMs. At this moment, we can give only some information 390 about PTMs of proteins coded by chromosome 18 in HepG2 391 (Table S2). A deeper database search of our raw data for other 392 samples is in progress. Available in Proteoform Atlas, information 393 about already detected proteoforms also will be very helpful.

## 394 **ASSOCIATED CONTENT**

## 395 Supporting Information

396 The Supporting Information is available free of charge on the ACS 397 Publications website at DOI: 10.1021/acs.jproteome.8b00386.

- Figures S1 (adapted from ref 29) and S2 (PDF)
- 399 Table S1 (XLSX)
- 400 Table S2 (XLSX)

#### 401 **AUTHOR INFORMATION**

## 402 Corresponding Author

403 \*Tel: (+7) 9111764453. E-mail: snaryzhny@mail.ru.

404 ORCID 6

405 Stanislav N. Naryzhny: 0000-0002-4102-3423

406 Elena S. Zorina: 0000-0003-4456-6850

413

421

428

The authors declare no competing financial interest. 408 The mass spectrometry proteomics data have been deposited 409 to the ProteomeXchange Consortium via the PRIDE<sup>30</sup> partner 410 repository with the data set identifier PXD010142 and 411 10.6019/PXD010142. 412

## ACKNOWLEDGMENTS

This work was funded by a grant of RSF (Russian Science 414 Foundation) #14-25-00132. Mass-spectrometry measurements 415 were performed using the equipment of "Human Proteome" 416 Core Facilities of the Institute of Biomedical Chemistry 417 (Russia), which is supported by Ministry of Education and 418 Science of the Russian Federation (unique project ID 419 RFMEFI62117X0017). 420

ABBREVIATIONS

2DE, two-dimensional electrophoresis; ESI LC–MS/MS, 422 liquid chromatography-electrospray ionization-tandem mass 423 spectrometry; HCD, higher energy collisional dissociation; 424 ABC, ammonium bicarbonate; ACN, acetonitrile; PTM, post- 425 translation modification; emPAI, exponential modified form of 426 protein abundance index. 427

REFERENCES

(1) Legrain, P.; Aebersold, R.; Archakov, A.; Bairoch, A.; Bala, K.; 429 Beretta, L.; Bergeron, J.; Borchers, C. H.; Corthals, G. L.; Costello, C. 430 E.; et al. The Human Proteome Project : Current State and Future 431 Direction. *Mol. Cell. Proteomics* **2011**, *10*, 1–5. 432

(2) Paik, Y. K.; Jeong, S. K.; Omenn, G. S.; Uhlen, M.; Hanash, S.; 433 Cho, S. Y.; Lee, H. J.; Na, K.; Choi, E. Y.; Yan, F.; et al. The 434 Chromosome-Centric Human Proteome Project for Cataloging 435 Proteins Encoded in the Genome. *Nat. Biotechnol.* **2012**, *30*, 221–436 223. 437

(3) Paik, Y. K.; Omenn, G. S.; Uhlen, M.; Hanash, S.; Marko-Varga, 438 G.; Aebersold, R.; Bairoch, A.; Yamamoto, T.; Legrain, P.; Lee, H. J.; 439 et al. Standard Guidelines for the Chromosome-Centric Human 440 Proteome Project. *Journal of Proteome Research.* **2012**, *11*, 2005–2013. 441 (4) Naryzhny, S. N.; Maynskova, M. A.; Zgoda, V. G.; Ronzhina, N. 442 L.; Novikova, S. E.; Belyakova, N. V.; Kleyst, O. A.; Legina, O. K.; 443 Pantina, R. A. F. M. Proteomic Profiling of High-Grade Glioblastoma 444 Using Virtual-Experimental 2DE. *J. Proteomics Bioinf.* **2016**, *9* (6), 445 158–165.

(5) Naryzhny, S. N.; Lisitsa, A. V.; Zgoda, V. G.; Ponomarenko, E. 447 A.; Archakov, A. I. 2DE-Based Approach for Estimation of Number of 448 Protein Species in a Cell. *Electrophoresis* **2014**, 35 (6), 895–900. 449

(6) Zabel, C.; Klose, J. Protein Extraction for 2DE. *Methods Mol.* 450 *Biol.* **2009**, 519, 171–196. 451

(7) Naryzhny, S.; Maynskova, M.; Zgoda, V.; Archakov, A. Dataset 452 of Protein Species from Human Liver. *Data Br.* 2017, *12*, 584–588. 453

(8) Naryzhny, S.; Zgoda, V.; Kopylov, A.; Petrenko, E.; Archakov, A. 454 A Semi-Virtual Two Dimensional Gel Electrophoresis: IF–ESI LC– 455 MS/MS. *MethodsX* **2017**, *4*, 260. 456

(9) Naryzhny, S. N.; Maynskova, M. A. Proteomic Profiling of High- 457 Grade Glioblastoma Using Virtual Experimental 2DE. J. Proteomics 458 Bioinf. 2016, DOI: 10.4172/jpb.1000402. 459

(10) Bradford, M. M. A Rapid and Sensitive Method for the 460 Quantitation of Microgram Quantities Utilizing the Principle of. *Anal.* 461 *Biochem.* **1976**, *72*, 248–254. 462

(11) Gorg, A.; Postel, W.; Domscheit, A.; Gunther, S. Two- 463 Dimensional Electrophoresis with Immobilized pH Gradients of Leaf 464 Proteins from Barley (Hordeum Vulgare): Method, Reproducibility 465 and Genetic Aspects. *Electrophoresis* **1988**, *9* (11), 681–692. 466

(12) Naryzhny, S. N. Blue Dry Western: Simple, Economic, 467 Informative, and Fast Way of Immunodetection. *Anal. Biochem.* 468 **2009**, 392 (1), 90. 469

470 (13) Naryzhny, S. N.; Zgoda, V. G.; Maynskova, M. A.; Novikova, S.
471 E.; Ronzhina, N. L.; Vakhrushev, I. V.; Khryapova, E. V.; Lisitsa, A. V.;
472 Tikhonova, O. V.; Ponomarenko, E. A.; et al. Combination of Virtual
473 and Experimental 2DE Together with ESI LC-MS/MS Gives a

474 Clearer View about Proteomes of Human Cells and Plasma. 475 Electrophoresis **2016**, 37 (2), 302–309.

476 (14) Vaudel, M.; Barsnes, H.; Berven, F. S.; Sickmann, A.; Martens, 477 L. SearchGUI: An Open-Source Graphical User Interface for 478 Simultaneous OMSSA and X!Tandem Searches. *Proteomics* **2011**, 479 11 (5), 996–999.

480 (15) Ishihama, Y.; Oda, Y.; Tabata, T.; Sato, T.; Nagasu, T.; 481 Rappsilber, J.; Mann, M. Exponentially Modified Protein Abundance 482 Index (emPAI) for Estimation of Absolute Protein Amount in 483 Proteomics by the Number of Sequenced Peptides per Protein. *Mol.* 484 *Cell. Proteomics* **2005**, *4* (9), 1265–1272.

(16) Tran, J. C.; Zamdborg, L.; Ahlf, D. R.; Lee, J. E.; Catherman, A.
D.; Durbin, K. R.; Tipton, J. D.; Vellaichamy, A.; Kellie, J. F.; Li, M.;
et al. Mapping Intact Protein Isoforms in Discovery Mode Using Topdown Proteomics. *Nature* 2011, 480 (7376), 254–258.

(17) Naryzhny, S. N.; Lee, H. The Post-Translational Modifications
of Proliferating Cell Nuclear Antigen: Acetylation, Not Phosphorylation, Plays an Important Role in the Regulation of Its Function. *J.*Biol. Chem. 2004, 279 (19), 20194–20199.

493 (18) Halligan, B. D.; Ruotti, V.; Jin, W.; Laffoon, S.; Twigger, S. N.; 494 Dratz, E. A. ProMoST (Protein Modification Screening Tool): A

495 Web-Based Tool for Mapping Protein Modifications on Two-496 Dimensional Gels. *Nucleic Acids Res.* 2004, 32, W638.

497 (19) Halligan, B. D. ProMoST: A Tool for Calculating the pI and 498 Molecular Mass of Phosphorylated and Modified Proteins on Two-499 Dimensional Gels. *Methods Mol. Biol.* **2009**, *527*, 283.

500 (20) Xu, G.; Jaffrey, S. R. Proteomic Identification of Protein 501 Ubiquitination Events. *Biotechnol. Genet. Eng. Rev.* **2013**, *29* (1), 73– 502 109.

503 (21) Hendriks, I. A.; Vertegaal, A. C. O. A Comprehensive 504 Compilation of SUMO Proteomics. *Nat. Rev. Mol. Cell Biol.* **2016**, 505 17 (9), 581–595.

506 (22) Bellia, F.; Vecchio, G.; Rizzarelli, E. Carnosinases, Their 507 Substrates and Diseases. *Molecules* **2014**, *19*, 2299–2329.

508 (23) Peters, V.; Zschocke, J.; Schmitt, C. P. Carnosinase, Diabetes 509 Mellitus and the Potential Relevance of Carnosinase Deficiency. *J.* 

510 Inherited Metab. Dis. 2018, 41, 39–47.
511 (24) Buxbaum, J. N.; Reixach, N. Transthyretin: The Servant of

512 Many Masters. Cell. Mol. Life Sci. 2009, 66, 3095–3101.

513 (25) Naryzhny, S. N.; Lee, H. Observation of Multiple Isoforms and

514 Specific Proteolysis Patterns of Proliferating Cell Nuclear Antigen in 515 the Context of Cell Cycle Compartments and Sample Preparations. 516 *Proteomics* **2003**, *3*, 930–936.

(26) Castellanos-Serra, L.; Paz-Lago, D. Inhibition of Unwanted
 Proteolysis during Sample Preparation: Evaluation of Its Efficiency in
 Challenge Experiments. *Electrophoresis* 2002, 23 (11), 1745–1753.

520 (27) Thiede, B.; Koehler, C. J.; Strozynski, M.; Treumann, a.; Stein, 521 R.; Zimny-Arndt, U.; Schmid, M.; Jungblut, P. R. Protein Species 522 High Resolution Quantitative Proteomics of HeLa Cells Using 523 SILAC-2-DE-nanoLC/LTQ-Orbitrap Mass Spectrometry. *Mol. Cell.* 

524 Proteomics 2013, 12 (2), 529–538. 525 (28) Naryzhny, S.; Zgoda, V.; Kopylov, A.; Petrenko, E.; Archakov,

526 A. A Semi-Virtual Two Dimensional Gel Electrophoresis: IF-ESI 527 LC-MS/MS. *MethodsX* 2017, 4, 260-264.

528 (29) Naryzhny, S. N.; Maynskova, M. A.; Zgoda, V. G.; Ronzhina, N. 529 L.; Kleyst, O. A.; Vakhrushev, I. V.; Archakov, A. I. Virtual-530 Experimental 2DE Approach in Chromosome-Centric Human 531 Proteome Project. J. Proteome Res. **2016**, *15*, 525.

(30) Vizcaíno, J. A.; Csordas, A.; Del-Toro, N.; Dianes, J. A.; Griss,
J.; Lavidas, I.; Mayer, G.; Perez-Riverol, Y.; Reisinger, F.; Ternent, T.;
et al. 2016 Update of the PRIDE Database and Its Related Tools. *Nucleic Acids Res.* 2016, 44 (D1), D447–D456.

(31) Naryzhny, S. Inventory of Proteoforms as a Current Challenge
 of Proteomics: Some Technical Aspects. J. Proteomics 2018,
 538 DOI: 10.1016/j.jprot.2018.05.008.