1　Proteomics as a metrological tool to evaluate genome annotation accuracy

2　following *de novo* genome assembly: a case study using the Atlantic bottlenose

3　dolphin (*Tursiops truncatus*)

4

5　Benjamin A. Neely[1,*], Debra L. Ellisor[2], W. Clay Davis[1]

6

7

8　[1] *National Institute of Standards and Technology, Material Measurement Laboratory, Chemical Sciences*

9　*Division, Marine Biochemical Sciences Group, Hollings Marine Laboratory, 331 Fort Johnson Road,*

10　*Charleston, SC 29412, United States*

11　[2] *National Institute of Standards and Technology, Material Measurement Laboratory, Chemical Sciences*

12　*Division, Environmental Specimen Bank Group, Hollings Marine Laboratory, 331 Fort Johnson Road,*

13　*Charleston, SC 29412, United States*

14

15

16

17

18

19

20　* Corresponding Author: benjamin.neely@nist.gov

21　Keywords: *de novo* genome, genome accuracy, proteomics, Atlantic bottlenose dolphin (*Tursiops*

22　*truncatus*), marine mammal

23   **Abstract**

24   *Background*

25   The last decade has witnessed dramatic improvements in whole-genome sequencing capabilities coupled

26   to drastically decreased costs, leading to an inundation of high-quality *de novo* genomes. For this reason,

27   continued development of genome quality metrics is imperative. The current study utilized the recently

28   updated Atlantic bottlenose dolphin (*Tursiops truncatus*) genome and annotation to evaluate a

29   proteomics-based metric of genome accuracy.

30

31   *Results*

32   Proteomic analysis of six tissues provided experimental confirmation of 10 402 proteins from 4 711

33   protein groups, almost 1/3 of the possible predicted proteins in the genome. There was an increased

34   median molecular weight and number of identified peptides per protein using the current *T. truncatus*

35   annotation versus the previous annotation. Identification of larger proteins with more identified peptides

36   implied reduced database fragmentation and improved gene annotation accuracy. A metric is proposed,

37   $NP_{10}$, that attempts to capture this quality improvement. When using the new *T. truncatus* genome there

38   was a 21 % improvement in $NP_{10}$. This metric was further demonstrated by using a publicly available

39   proteomic data set to compare human genome annotations from 2004, 2013 and 2016, which had a 33 %

40   improvement in $NP_{10}$.

41

42   *Conclusions*

43   These results demonstrate that new whole-genome sequencing techniques can rapidly generate high

44   quality *de novo* genome assemblies and emphasizes the speed of advancing bioanalytical measurements

45   in a non-model organism. Moreover, proteomics may be a useful metrological tool to benchmark genome

46   accuracy, though there is a need for reference proteomic datasets to facilitate this utility in new *de novo*

47   and existing genomes.

Page 2

48 **Background**

49      Since 2007 there has been a rapid decrease in whole-genome sequencing costs coupled with

50 improved read lengths and development of long-range techniques such as synthetic long-reads and

51 mapping protocols. Concurrently, the access to high performance computing environments has improved

52 along with an endless supply of new genome assembly and annotation tools. With these new resources it

53 is now possible to rapidly generate high-quality *de novo* genomes for non-model organisms. Excellent

54 examples of this are two recently completed mammalian genomes (the domestic goat, *Capra hircus* [1, 2],

55 and the Hawaiian monk seal, *Neomonachus schauinslandi* [3]) that utilized a combination of approaches

56 including optical mapping, synthetic long reads, long read technology and chromatin interaction mapping

57 to generate highly contiguous (scaffold N50 > 29.5 Mbp) *de novo* genomes at a relatively low cost.

58 Overall, the result of these parallel advancements are numerous large-scale sequencing projects [4], the

59 most ambitious targeting approximately 9 000 eukaryotic species (Earth BioGenome Project). With the

60 forthcoming inundation of new high-quality *de novo* genomes, there is a continued need for improved

61 metrics to evaluate genome accuracy.

62      Genome assemblies and annotations are evaluated in terms of contiguity and completeness, both

63 indicators of genome accuracy. Measures of contiguity, such as scaffold N50 or N90 length, typically

64 correspond to the quality of the genome assembly [5]. Scaffold N50 or N90 length is similar to a median

65 or quantile scaffold length but is dependent on assembly size. Greater scaffold contiguity tends to result in

66 more protein-coding sequences and isoforms. For example, one of the initial finished human genome

67 assemblies from 2004 (NCBI Build 34) had a scaffold N50 of 27.2 Mbp and 27 180 protein-coding

68 sequences, which has since been improved to a scaffold N50 of 59.4 Mbp and 109 018 protein-coding

69 sequences (NCBI Release 108, March 2016). Gains can be even more pronounced in non-model

70 organisms with improved *de novo* genome assemblies. For example, the *Alligator mississippiensis*

71 (American alligator) genome recently improved from a scaffold N50 of 508 kbp to 10 Mbp using new

72 sequencing methods [6]. Similarly, the focus of this study, *Tursiops truncatus* (Atlantic bottlenose

73 dolphin), improved from a scaffold N50 of 116 kbp to 26.6 Mbp. Studies have shown that assembly

74   contiguity often corresponds to assembly quality [5] but does not necessarily correlate with genome

75   completeness and therefore accuracy [7]. One way to evaluate genome completeness is by using predicted

76   conserved gene products. First used in the Core Eukaryotic Genes Mapping Approach (CEGMA) [8, 9],

77   this concept has developed into Benchmarking Universal Single-Copy Orthologs (BUSCO), which is a

78   content-based quality assessment that uses universal single-copy markers to gauge genome completeness

79   [7]. It is evident that using many metrics to benchmark *de novo* genomes is essential to evaluating

80   genome quality. Given the orthogonal nature of proteomics and its dependence on accurately predicted

81   gene annotations, a quality metric based in this analytical domain may be advantageous.

82   Data-dependent acquisition bottom-up shotgun proteomics is one method to confirm gene

83   annotations by observing the predicted proteins using mass spectrometry. First, proteins are digested with

84   a known protease and the resulting peptides are fragmented within a mass spectrometer. Next, using an

85   accurate mass of the peptide and the resulting fragmentation pattern, search algorithms can

86   probabilistically identify peptides and then infer proteins in the search database. Alternatively, spectral

87   libraries directly match fragmentation patterns, though these initial assignments are typically made using

88   database-dependent approaches [10-12]. With the current generation of mass spectrometers, which have

89   high duty cycles with high mass accuracy and resolution, we may be approaching the era of being able to

90   infer the majority of proteins in a genome. For example, a recent proteomic analysis of HeLa tissue

91   accounted for 91.5 % of gene products measured in the same tissue by RNA-seq (12 209 protein coding

92   sequences versus 13 347 gene products) [13]. Since bottom-up shotgun proteomics relies completely on a

93   database for peptide identifications and protein inference, it may be possible that a high-quality mass

94   spectrometric dataset could be used to benchmark genome assembly and annotation quality.

95   The purpose of the current study was two-fold: (i) provide detailed proteomic profiling of a

96   marine mammal and (ii) use this data to evaluate the new *T. truncatus* assembly and annotation. On

97   average over 4 800 proteins were identified in six different tissues, and when combined yielded 10 402

98   protein identifications. Although not an exhaustive proteomic dataset, it confirmed approximately 1/3 of

99   the predicted protein-coding genes. This dataset is an invaluable resource to support comparative

100 proteomics in diving mammals related to comparative evolution [14] and biomimicry [15] and

101 demonstrates the feasibility of accelerating cutting-edge bioanalytical approaches in non-model

102 organisms. Secondly, the new *de novo* assembly resulted in increased protein identifications but also a

103 decreased number of peptide identifications, despite more than a 200-fold improvement in scaffold N50

104 over the previous assembly. We investigated these differences at the peptide and protein level to identify

105 global trends and proposed a new measure of genome annotation quality, $NP_{10}$. This new measure was

106 further demonstrated by evaluating human genome improvements over the past decade using publicly

107 available proteomic data. Overall, these results highlight the improved annotation accuracy of the new *T.*

108 *truncatus* genome, the utility of proteomics as a metrological tool for evaluating genome annotation

109 quality, and emphasizes the need for reference proteomic datasets to facilitate metrology in new and

110 existing genomes.

111

112 **Results**

113 *Proteomic analysis of six tissues using NIST_Tur_tru v1*

114 The initial goal of this study was to advance metrological capabilities in *T. truncatus*. This was

115 accomplished by demonstrating proteomic measurements of six tissues from *T. truncatus*. On average, 2

116 199 protein groups and 4 888 proteins were identified in each tissue. The reason for performing proteomic

117 analysis on multiple tissue types was to capture more of the possible protein population. Although there

118 were 1 310 protein identifications shared across tissues, there was also diversity in protein identifications

119 between tissues with the brain and skin analyses having the most unique proteins (Figure 1). Proteomic

120 results for each tissue are available (Additional File Tables S6 –S11). It is interesting to note that the liver,

121 kidney and blubber came from the individual used for whole-genome sequencing. This dataset is

122 relatively diverse and provides experimental evidence for over 32 000 proteotypic peptides.

123

124 *Comparison of Ttru_1.4 and NIST_Tur_tru v1*

125   The second goal of the current study was to evaluate the new *T. truncatus de novo* genome assembly

126   (GCA_001922835.1) and annotation (NIST_Tur_tru v1). This genome assembly was generated in the fall

127   of 2016 using shotgun sequencing coupled to an *in vitro* histone ligation-based sequencing method (*i.e.*,

128   Chicago method) and proprietary assemblers described in detail by Putnam *et al.* [6]. This process

129   resulted in a genome assembly with a scaffold N50 of 26.6 Mbp. Of the 159 species with genomes

130   currently deposited on NCBI, 41 have scaffold N50 values greater than 26.6 Mbp. This level of contiguity

131   is becoming more commonplace with three marine mammal genomes released in 2017 with scaffold N50

132   greater than 19 Mbp (*T. truncatus*, *Neomonachus schauinslandi*, Hawaiian monk seal [3], and

133   *Delphinapterus leucas*, beluga whale [16]). For comparison, the prior NCBI *T. truncatus* annotation

134   (Ttru_1.4) was used. This assembly was a 2012 update [14] to the 2008 draft assembly based on Sanger

135   sequencing, Ttru_1.2 [17].

136        Both Ttru_1.4 and NIST_Tur_tru v1 are publicly available on NCBI and have been annotated

137   using NCBI's eukaryotic annotation pipeline and made available in RefSeq [18]. The current annotation

138   release, release 101 based on NIST_Tur_tru v1, has 24 026 genes and pseudogenes and 17 096 protein-

139   coding genes with 38 849 coding sequences. At the gene and transcript level, there were many changes

140   from Ttru_1.4 that are delineated based on alignment of genes and transcripts: identical, minor changes,

141   major changes, new, deprecated and other. These categories are defined and available through NCBI's

142   annotation report [19]. Briefly, 28 % of the prior genes and transcripts in Ttru_1.4 were deprecated, 72 %

143   had minor or major changes, and 21 % of the genes and transcripts in the NIST release are new.

144   Additionally, a small group of proteins have the prefix YP, which is not included in these NCBI

145   categories.

146        Tandem mass spectrometry data collected from all six tissues was searched against each release.

147   For both releases, almost 1/3 of the predicted protein-coding sequences were inferred by mass

148   spectrometry. Specifically the NIST assembly identified 32 582 peptide groups belonging to 10 402

149   proteins comprising 4 711 protein groups. The Ttru_1.4 assembly identified 33 738 peptide groups

150   belonging to 6 899 proteins comprising 5 292 protein groups. Many of the differences between the two

151 results were due to a loss of deprecated sequences and minor/major changes (Figure 2). Broadly, these

152 changes resulted in larger proteins with an increased median molecular weight and $NP_{10}$ molecular

153 weight.

154

155 *Confirming improvements in gene annotation*

156 There were 4,695 protein-coding sequences in the Ttru_1.4 annotation listed as partial, and one of the

157 main improvements in the new NIST annotation was that 86 % of these sequences were merged into

158 complete sequences. This offered an opportunity to evaluate the accuracy of these new assignments by

159 determining whether peptides identified by mass spectrometry supported the new complete sequences. Of

160 6 899 identified proteins using Ttru_1.4, 1 249 were partials. Of these 1 249 partial proteins identified

161 using Ttru_1.4, 534 had minor changes, 256 major, 450 were deprecated and 9 were other (defined simply

162 as other changes [19]). When this NIST annotation was used, 1 005 of these same 1 249 proteins were

163 identified, with 985 no longer being listed as partial. The median improvement within each protein was

164 two additional unique peptides and overall the median molecular weight improved 1.8-fold (Figure 3). Of

165 these 1 005 partial proteins identified using Ttru_1.4, when using the NIST annotation, 886 had increased

166 molecular weight and increased number of unique peptides.

167

168 *Comparing peptide identifications*

169 An unexpected result in the new annotation was that there were fewer peptide identifications. Given the

170 major changes between the two releases related to deprecated genes, new genes, and major changes, we

171 were interested in tracking these peptide level changes. Over 80 % of the peptide groups identified in

172 NIST annotation were also identified using the Ttru_1.4 annotation (Figure 4). The new peptide

173 identifications were linked to major and minor changes in genes with only 3.2 % due to new sequences.

174 As would be expected, many of the peptide groups not identified in the NIST annotation were deprecated

175 (41 %). Given that these 5 657 peptide groups lost using the NIST annotation were high-confidence

176   identifications, it may provide evidence for re-inclusion of these protein-coding sequences in future

177   annotation releases.

178

179   *Specific examples of annotation improvements*

180   The goal of evaluating differences at a broad level is to capture and describe relevant changes at the

181   granular level. At the peptide level, one the most striking improvements was related to titin, a major

182   component in muscle tissue. In Ttru_1.4, titin (XP_004322250.1) was a partial sequence of 2,167 amino

183   acids (241.7 kDa) and 60 unique peptides (40.2 %) were identified belonging to this sequence. In the

184   NIST annotation, the coding sequence for titin (XP_019787158.1) was 32 192 amino acids (3 812.8 kDa)

185   and 779 unique peptides (34.3 % coverage) were identified belonging to this sequence. This single

186   sequence improvement is responsible for many changes observed at the peptide level (Figure 4).

187        Almost 2 % of the identified proteins using the NIST annotation were considered new. One

188   important new protein of note is cystatin C (XP_019783122.1). This protein was not present in Ttru_1.4,

189   while using the NIST annotation the mass spectrometry data identified three unique peptides (41.3 %

190   coverage) belonging to the predicted 13.1 kDa protein. This protein has applications as a biomarker [20],

191   and with these proteomic results, it is possible to create SI traceable mass spectrometer-based assays

192   (similar to [21]). Another protein of note is serotransferrin (XP_019789750.1), which is 90 % identical

193   and 3.5 % longer than the entry in Ttru_1.4 annotation (XP_004329553.1). Most of these changes were

194   on the c-terminus section (from positions 537 to 634), which was supported by the proteomic data that

195   identified four peptides spanning this region. There were other slight changes to the sequence that resulted

196   in six more unique peptides identified in the improved serotransferrin, which supports the accuracy of the

197   new annotation. Overall, there are many changes related to the over 10 000 protein identifications and

198   many would be considered improvements as indicated by increased protein molecular weight and/or

199   greater peptide coverage. At a gene-by-gene level these results can be used to confirm and improve

200   annotations.

201

202    *Confirming quality metric in human annotations*

203    In order to gauge the broader applicability of using proteomics as a quality measure of genomic

204    annotations, we demonstrated $NP_{10}$ in a more mature genome with deeper proteomics. The recent work by

205    Bekker-Jensen *et al.* [13] is publicly available on ProteomeXchange [22, 23] and for this comparison the

206    data generated from a 39 fraction high pH pre-fractionation of a HeLa cell digest followed by LC-MS/MS

207    analysis was used for database searching. These data were searched against three human genome

208    annotations from 2004, 2013 and 2016, each with markedly increased scaffold N50 values and database

209    sizes (*i.e.*, number of coding-sequences; Table 1). The number of identified proteins was 13 341, 22 906,

210    and 48 019 proteins in Build 34, Release 105 and Release 108, respectively. The median molecular

211    weight improved 25 % (from 51.06 to 53.46 to 63.99 kDa, respectively) whereas the improvement in

212    $NP_{10}$ was more pronounced with a 33 % improvement (from 100.17 to 101.87 to 133.55 kDa,

213    respectively; Figure 5).

214

215    **Table 1. Descriptive statistics of human annotated databases and resulting proteomic**

216    **identifications.**

|  | Build 34 | Release 105 | Release 108 |
|---|---|---|---|
| **release date** | Feb 2004 | Jun 2013 | Mar 2016 |
| **scaffold N50** | 29.1 Mbp | 45.0 Mbp | 59.4 Mbp |
| **coding sequences** | 27 180 | 45 107 | 109 018 |
| **protein groups** | 9 762 | 10 059 | 10 219 |
| **proteins** | 13 341 | 22 906 | 48 019 |
| **peptide groups** | 175 895 | 184 580 | 184 806 |
| **peptide spectral matches** | 390 909 | 405 852 | 405 950 |

217

218

219

**Discussion**

220

221      Advances in bioanalytical platforms across domains (*i.e.*, genomics, transcriptomics, and

222   proteomics) are improving the accessibility of non-model organisms as viable research candidates. The

223   results of the current study provide secondary confirmation of 10 402 proteins from 4 711 protein groups

224   using a recently completed well-scaffolded high-coverage *T. truncatus* genome and shotgun proteomic

225   analysis of six different tissues. Previous proteomic studies of *T. truncatus* have identified less than 100

226   protein groups in serum [15, 21], while the most detailed published proteomic analysis of a marine

227   mammal identified 206 proteins in cerebrospinal fluid of *Zalophus californianus* (California sea lion)

228   [24]. Currently there are twelve marine mammal genomes that have been annotated by NCBI (of the 159

229   species with genomes currently deposited on NCBI), though only *T. truncatus* and *Z. californianus* have

230   published mass spectrometry based proteomic datasets. Work is underway to increase the number of

231   marine mammal genomes along with companion high-quality proteomic datasets and spectral libraries.

232   The results of the current study provide empirical confirmation of protein annotations, including

233   observable proteotypic peptides, which can be a resource for future targeted studies in *T. truncatus*. For

234   example, by improving the protein-coding sequence accuracy of serotransferrin in *T. truncatus,* future

235   studies can extrapolate metrological advances in human serotransferrin sialoforms [25] to *T. truncatus*

236   disease treatment [26]. Since the current results are not an exhaustive proteomic dataset, future studies

237   will utilize different solubilization techniques, proteases, and separation techniques to provide even

238   deeper proteome coverage (reviewed and demonstrated in the following [13, 27, 28]). Still, it is worth

239   noting that in single study using a simple experimental approach we have identified almost 1/3 of the

240   possible predicted proteins, emphasizing the ease of accomplishing bioanalytical advances in non-model

241   organisms using modern techniques.

242      In the current study, benchmark proteomic datasets were used to evaluate genome assembly and

243   annotation improvements in *T. truncatus* and *H. sapiens*. Typically, a reference database is used to

244   demonstrate proteomic improvements due to optimized protein extraction, solubilization and digestion,

245   peptide separation, mass spectrometer speed and mass accuracy, search algorithm performance and

246    database accuracy. In contrast, when the mass spectrometric data are held constant and instead the

247    database is varied, differences in proteomic results are indicative of database fragmentation and accuracy.

248    Proteomic analysis of multiple tissues allowed for greater protein diversity when evaluating *T. truncatus*,

249    though the publicly available human data performed exceptionally well despite using a single tissue since

250    it utilized highly optimized separation techniques. An optimum proteomic benchmark dataset would be

251    one that offers the possibility of the deepest proteome coverage. This would rely on using multiple

252    tissues, extraction protocols, enzymes and optimum separation techniques coupled to modern mass

253    spectrometers. These datasets could be developed in parallel to the exponential increase in *de novo*

254    genomes being released and annotated and would prove invaluable in exercises assessing assembly and

255    annotation performance (such as Assemblathon 2 [5]). Importantly, given the abundance and accessibility

256    of public proteomic data in this "Golden Age of Proteomics" (as coined by [29]) and modular open-access

257    proteogenomic pipelines such as Galaxy-P [30, 31], it would be possible to incorporate these reference

258    mass spectrometric datasets and proteomic derived quality metrics into genome assembly and annotation

259    pipelines.

260         In parallel to improvements in genome assembly contiguity and annotation accuracy, proteomic

261    results should have increased peptide numbers per protein, higher protein identifications due to isoform

262    resolution and improved coverage of higher molecular weight proteins due to better long-range accuracy.

263    For instance, when evaluating the substantial reduction in partial sequences between Ttru_1.4 and

264    NIST_Tur_tru v1, there was an increase of 81 % in median molecular weight of these proteins that

265    coincided with more peptide identifications within these new complete sequences. The most drastic

266    example in this case study was titin, which went from 60 to 779 identified peptides with the addition of

267    over 32 000 amino acids to the previously partial sequence. This also emphasizes that greater numbers of

268    protein identifications does not imply higher quality since a more fragmented genome will give more

269    protein identifications. Instead, identification of larger proteins with more identified peptides is more

270    indicative of improved quality. The proposed metric, $NP_{10}$, attempts to capture this quality measure. One

271    issue is that the $NP_{10}$ may be glossing over how changes in spectral assignments to peptides with

Page 11

272 changing databases affect proteomic quality (such as false discovery rates). There is an opportunity to

273 develop a streamlined method to track MS/MS spectra assignments and quantify those changes with

274 database improvements in order to establish finer measures of search space effects on proteomic

275 performance. Overall, these results demonstrate that new whole-genome sequencing techniques can

276 provide high quality *de novo* genome assemblies and that proteomics is a useful metrological tool to

277 evaluate annotation and benchmark genome accuracy.

278

279 **Methods**

280 *Sample source and preparation*

281 Bottlenose dolphin tissues were collected from animals under appropriate permits (Additional File Table

282 S1) and stored at liquid nitrogen temperatures (-150 to -180 °C) until cryohomogenization in the National

283 Institute of Science and Technology's Marine Environmental Specimen Bank [32]. From the resulting

284 fine powder, 5 mg was subsampled and the proteins were extracted using RapiGest (Waters, Milford

285 MA). Briefly, 150 µL of 0.1 % (w/v) RapiGest (in 50 mM ammonium bicarbonate) was added, resulting

286 in a solution of 33 µg/µL tissue. The extraction mixture was shaken at 600 rpm for 25 min at room

287 temperature followed by removal of large debris using a benchtop microcentrifuge. From this solution, a

288 5 µL aliquot was removed and suspended in 35 µL of 0.1 % (w/v) RapiGest (in 50 mM ammonium

289 bicarbonate), followed by the addition of 40 uL of 50 mM ammonium bicarbonate. Next, the sample was

290 reduced with 10 µL of 45 mM dithiothreitol (DTT; final concentration of 5 mM) and incubated at 60 °C

291 for 30 min, then allowed to cool to room temperature. The mixture was alkylated using 3.75 µL of 375

292 mM iodoacetamide (Pierce, Thermo Scientific, Waltham, MA; final concentration of 15 mM) and

293 incubated in the dark at room temperature for 20 min. Prior to addition of trypsin, 100 µL of 50 mM

294 ammonium bicarbonate was added. A 3.3 µL aliquot of trypsin (MS-Grade; 1 µg/µl in 50 mM acetic acid)

295 was added (1:50 trypsin:protein) and samples were incubated overnight at 37 °C. The digestion was

296 halted and RapiGest cleaved with the addition of 100 µL of 3 % (v/v) trifluoroacetic acid (1% final

297    concentration) and incubated at 37 °C for 30 min before centrifugation and removal of the supernatant.

298    Samples were processed using Pierce C18 spin columns (8 mg of C18 resin; Thermo Scientific) according

299    to manufacturer's instructions. Each sample was processed in duplicate yielding at maximum of 60 µg

300    peptides. These solutions were evaporated to dryness in a vacufuge then reconstituted in 150 µL of 5 %

301    acetonitrile in water.

302

303    *Mass Spectrometry*

304    Samples were analyzed using an UltiMate 3000 Nano LC coupled to a Fusion Lumos mass spectrometer

305    (Thermo Fisher Scientific). Resulting peptide mixtures (10 µl) were loaded onto a PepMap 100 C18 trap

306    column (75 µm id x 2 cm length; Thermo Fisher Scientific) at 3 µL/min for 10 min with 2 % (v/v)

307    acetonitrile and 0.05 % (v/v) trifluoroacetic acid followed by separation on an Acclaim PepMap RSLC 2

308    µm C18 column (75µm id x 25 cm length; Thermo Fisher Scientific) at 40 °C. Peptides were separated

309    along a 130 min gradient of 5 % to 27.5 % mobile phase B [80 % (v/v) acetonitrile, 0.08 % (v/v) formic

310    acid] over 105 min followed by a ramp to 40 % mobile phase B over 15 min and lastly to 95 % mobile

311    phase B over 10 min at a flow rate of 300 nL/min. The mass spectrometer was operated in positive

312    polarity and data dependent mode (topN, 3 s cycle time) with a dynamic exclusion of 60 s (with 10 ppm

313    error). The RF lens was set at 30 %. Full scan resolution using the orbitrap was set at 120 000 and the

314    mass range was set to *m/z* 375 to1500. Full scan ion target value was 4.0e5 allowing a maximum injection

315    time of 50 ms. Monoisotopic peak determination was used, specifying peptides and an intensity threshold

316    of 1.0e4 was used for precursor selection. Data-dependent fragmentation was performed using higher-

317    energy collisional dissociation (HCD) at a normalized collision energy of 32 with quadrupole isolation at

318    *m/z* 0.7 width. The fragment scan resolution using the orbitrap was set at 30 000, *m/z* 110 as the first

319    mass, ion target value of 2.0e5 and a 60 ms maximum injection time.

320    *Protein Search parameters*

321  Resulting raw files from the analysis of six different *T. truncatus* tissues and raw files from a publicly

322  available 39 fraction HeLa experiment (ProteomeXchange Consortium [23] via the PRIDE partner

323  repository with the dataset identifier PXD004452) were processed and searched using Proteome

324  Discoverer (v.2.0.0.802). For *T. truncatus* analysis, Sequest HT and Mascot (v2.6.0; Matrix Science)

325  search algorithms were used, while only Sequest HT was used for human searches. For all searches, the

326  protein.faa fasta file was retrieved from NCBI RefSeq [18] via ftp [33]. For searches with the prior *T.*

327  *truncatus* annotation, GCF_000151865.2_Ttru_1.4 was used, while searches with the current *T. truncatus*

328  annotation, GCF_001922835.1_NIST_Tur_tru_v1 was used. These correspond to release 100 and 101 for

329  this organism on NCBI. The whole-genome sequencing projects can be found in GenBank [34] under

330  entries ABRN00000000.2 (Ttru_1.4) and MRVK00000000.1 (NIST_Tur_tru_v1). For the human

331  searches, the following were used: GCF_000001405.10_hg16_Build34.3 (Build 34),

332  GCF_000001405.25_GRCh37.p13 (Release 105) and GCF_000001405.33_GRCh38.p7 (Release 108).

333  The *T. truncatus* searches also used the common Repository of Adventitious Proteins database (cRAP;

334  2012.01.01; the Global Proteome Machine), though these sequences were removed from search results.

335  The following search parameters were used for Mascot and Sequest: trypsin was specified as the

336  enzyme allowing for two mis-cleavages; carbamidomethyl (C) was fixed and acetylation (protein n-term),

337  deamidated (NQ), pyro-Glu (n-term Q), and oxidation (M) were variable modifications; 10 ppm precursor

338  mass tolerance and 0.02 Da fragment ion tolerance. Within Sequest, the peptide length was specified as a

339  minimum of six and maximum of 144 amino acids. Resulting peptide spectral matches were validated

340  using the percolator algorithm, based on q-values at a 1 % false discovery rate (FDR). The peptides that

341  were greater than six amino acids long were grouped into proteins according to the law of parsimony and

342  filtered to 1 % FDR and single peptide hits were allowed. Briefly, there may be more than one peptide

343  spectral match for a given peptide, which are then grouped to peptide groups. Protein inference is when

344  these peptide groups are assigned to proteins, but given similarity between some proteins (such as

345  isoforms or highly homologous sequences), peptides can match to more than one protein. For this reason,

346  protein families or protein groups are generated based on peptide overlap (and therefore sequence

347    overlap), which reduces inflation due to isoform identifications. For the described analyses, protein and

348    peptide groups are used and are available for each *T. truncatus* search in Additional File Tables S2 – S5.

349    Raw MS data and Mascot based search results for *T. truncatus*, as well as all fasta databases, have been

350    deposited to the ProteomeXchange Consortium [23] via the PRIDE partner repository with the dataset

351    identifier PXD008808 and 10.6019/PXD008808.

352

353    *Proteomic-based quality metric for annotation quality*

354    Evaluating proteomic results relies on qualifying how well a database explains the observed tandem mass

355    spectra: high numbers of protein identifications and percent identified spectra indicate good proteomic

356    performance. Another way of describing proteomic results is to plot the number of peptide identifications

357    versus protein molecular weight. A larger protein has potentially more peptide identifications but due to

358    solubilization and digestion effects (such as post-translational modifications and protein folding), larger

359    proteins do not always yield more unique peptides. For this reason, there is a somewhat Gaussian

360    distribution of peptide frequency around median protein molecular weight. This median can shift right

361    when the molecular weight of predicted protein-coding sequences increases and/or the number of

362    isoforms increases.

363        When evaluating and comparing *de novo* genome assemblies and annotations, the specific

364    question that proteomics can answer is the degree of database fragmentation and accuracy. If an

365    annotation improves partial coding sequences to complete protein-coding sequences with isoforms, then

366    there will be an increase in the molecular weight of identified proteins with more peptides assigned to

367    these longer sequences. By simply improving partial sequences there would be a shift to higher protein

368    molecular weight. One goal of the current study was to provide a more robust quality measure by

369    incorporating unique peptide counts (which corresponds to protein coverage) with the change of median

370    molecular weight of inferred proteins. The $NP_{10}$ is a proposed metric that first stratifies the results by

371    identifying the top decile (or $10^{th}$ 10-quantile) of proteins based on the number of peptides per protein and

372    then returns the median molecular weight of the resulting proteins (graphically demonstrated in

373     Additional File Figure S1). This metric is similar to simply calculating the median molecular weight of all

374     inferred proteins, but by removing protein identifications with relatively few peptide assignments, it

375     attempts to indicate accuracy of the improved/longer protein-coding sequences.

376

377

378     **Availability of supporting data**

379     The raw data and tissue specific search results along with all databases used are available at the

380     ProteomeXchange Consortium [23] via the PRIDE partner repository with the dataset identifier

381     PXD008808 and 10.6019/PXD008808. The proteomic data from Bekker-Jensen *et al.* [13] used for the

382     human comparison can be found at ProteomeXchange Consortium [23] via the PRIDE partner repository

383     with the dataset identifier PXD004452. Tabulated search results for combined analysis and for each tissue

384     can be found in Additional File Supplemental Tables S1-S11.

385           Additional File Figure S1. Graphical example of $NP_{10}$ calculation.

386           Additional File Table S1. Sample characteristics table.

387           Additional File Table S2. Protein Identifications using Ttru_1.4.

388           Additional File Table S3. Protein Identifications using NIST_Tur_tru v1.

389           Additional File Table S4. Peptide Group Identifications using Ttru_1.4.

390           Additional File Table S5. Peptide Group Identifications using NIST_Tur_tru v1.

391           Additional File Table S6. Protein Identifications in blubber tissue using NIST_Tur_tru v1.

392           Additional File Table S7. Protein Identifications in brain tissue using NIST_Tur_tru v1.

393           Additional File Table S8. Protein Identifications in kidney tissue using NIST_Tur_tru v1.

394           Additional File Table S9. Protein Identifications in liver tissue using NIST_Tur_tru v1.

395           Additional File Table S10. Protein Identifications in muscle tissue using NIST_Tur_tru v1.

396           Additional File Table S11. Protein Identifications in skin tissue using NIST_Tur_tru v1.

397

398

399 **Declarations**

400 **List of abbreviations**

401 BUSCO          Benchmarking Universal Single-Copy Orthologs

402 C              cysteine

403 CEGMA          Core Eukaryotic Genes Mapping Approach

404 Da             Dalton

405 FDR            false discovery rate

406 kbp            kilo base pairs

407 kDa            kilodaton

408 M              methionine

409 Mbp            mega base pairs

410 MW             molecular weight

411 N              asparagine

412 $NP_{10}$         proposed metric of the median molecular weight of proteins that had greater than or equal

413                unique peptides identified to the $10^{th}$-decile of unique peptides per protein; notation

414                derived from number of peptides in $10^{th}$-decile

415 NIST           National Institute of Standards and Technology

416 Q              glutamine

417

418 **Consent for publication**

419 Not applicable.

420

421 **Competing interests**

422 The authors declare they have no competing interests.

423

424

425 **Funding**

426 All authors were funded by the National Institute of Standards and Technology.

427

428 **Authors' contributions**

429 All authors helped conceived of the study, developed methodology and assisted in reviewing the

430 manuscript. DE and WD selected and processed samples for proteomic analysis and collected data. BN

431 analyzed the data and wrote the initial manuscript draft.

432

433 **Acknowledgements**

434 Specimens used for this study were collected by Wayne E. McFee (National Oceanic and Atmospheric

435 Administration, National Ocean Service, National Centers for Coastal Ocean Science) and William A.

436 McLellan (University of North Carolina, Wilmington) and provided by the National Marine Mammal

437 Tissue Bank, which is maintained as part of the Marine Environmental Specimen Bank at NIST and is

438 operated under the direction of NMFS and in collaboration with USGS, USFWS, BOEMRE (formerly

439 MMS), and NIST through the Marine Mammal Health and Stranding Response Program. All samples

440 were collected under approved permits issued to the MMHSRP (responsible party: Dr. Teri Rowles) and

441 all sampling protocols were reviewed and approved by a NOAA/NMFS *ad hoc* Institutional Animal Care

442 and Use Committee (IACUC). The authors wish to thank Michael G. Janech and staff at NIST for critical

443 feedback.

444

445 **Disclaimer**

446 Identification of certain commercial equipment, instruments, software or materials does not imply

447 recommendation or endorsement by the National Institute of Standards and Technology, nor does it imply

448 that the products identified are necessarily the best available for the purpose.

449

450

451 **Figure Legends**

452 **Figure 1. Overlap and unique protein identifications by *T. truncatus* tissue.** Proteins unique to each

453 tissue and shared by all tissues are shown along with the total number of proteins identified in each

454 analysis.

455

456 **Figure 2. Descriptive statistics of identified proteins using different annotations.** The $NP_{10}$ molecular

457 weight improved 21.3 % from 67.59 kDa to 81.99 kDa (indicated by the red dotted line) along with an

458 improvement in median molecular weight of inferred proteins across genes with minor and major

459 changes. (note: these axes have been truncated for illustration and do not show all data points.)

460

461 **Figure 3. Confirming improved annotation of former partial proteins.** Proteins that were partial in

462 the Ttru_1.4 annotation were improved in the NIST annotation, and there was mass spectrometric

463 evidence to support the accuracy of these improvements corresponding to increased peptide

464 identifications and median molecular weight (the latter indicated by the red dotted line; note: these axes

465 have been truncated for illustration and do not show all data points.)

466

467 **Figure 4. Source of peptide identification differences using the two assemblies.** There was strong

468 overlap of identified peptides using the two assemblies with over 80 % overlap. The sources of the

469 differences were largely comprised of deprecated proteins in Ttru_1.4 (41 % of the 5 657) and

470 minor/major changes in NIST_Tur_tru_v1 (96 % of the 4 768).

471

472 **Figure 5. Similar trends with improved human assemblies.** As the contiguity of the human genome

473 has improved, there is a shift upward and to the right indicating annotations are more accurate (increased

474 coverage) and complete (increased molecular weight). The $NP_{10}$ improved 33 % and is indicated by the

475 red dotted line (note: these axes have been truncated for illustration and do not show all data points).
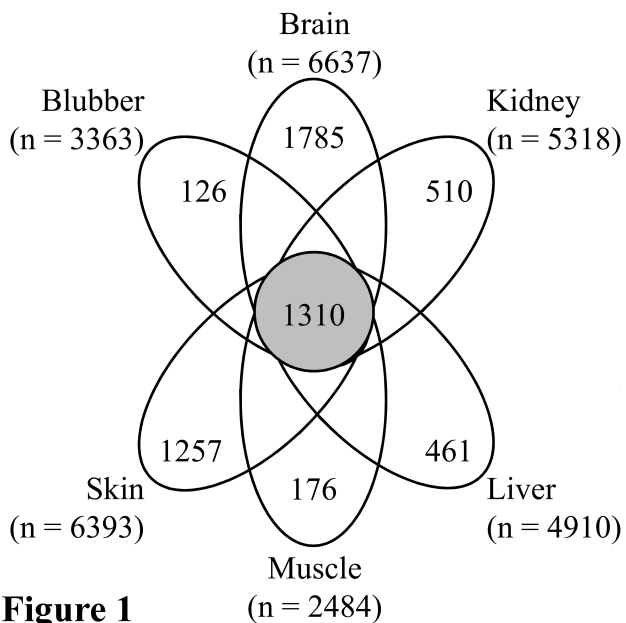
476

Page 19

477 **References**

478 1. Bickhart DM, Rosen BD, Koren S, Sayre BL, Hastie AR, Chan S, et al. Single-molecule

479 sequencing and chromatin conformation capture enable de novo reference assembly of the

480 domestic goat genome. Nature genetics. 2017;49 4:643-50. doi:10.1038/ng.3802.

481 2. Worley KC. A golden goat genome. Nature genetics. 2017;49 4:485-6. doi:10.1038/ng.3824.

482 3. Mohr DW, Naguib A, Weisenfeld N, Kumar V, Shah P, Church DM, et al. Improved de novo

483 Genome Assembly: Linked-Read Sequencing Combined with Optical Mapping Produce a High

484 Quality Mammalian Genome at Relatively Low Cost. bioRxiv. 2017:128348.

485 4. Richards S. It's More Than Stamp Collecting: How Genome Sequencing Can Unify Biological

486 Research. Trends in genetics : TIG. 2015;31 7:411-21. doi:10.1016/j.tig.2015.04.007.

487 5. Bradnam KR, Fass JN, Alexandrov A, Baranay P, Bechner M, Birol I, et al. Assemblathon 2:

488 evaluating de novo methods of genome assembly in three vertebrate species. GigaScience. 2013;2

489 1:10. doi:10.1186/2047-217x-2-10.

490 6. Putnam NH, O'Connell BL, Stites JC, Rice BJ, Blanchette M, Calef R, et al. Chromosome-scale

491 shotgun assembly using an in vitro method for long-range linkage. Genome research. 2016;26

492 3:342-50. doi:10.1101/gr.193474.115.

493 7. Simao FA, Waterhouse RM, Ioannidis P, Kriventseva EV and Zdobnov EM. BUSCO: assessing

494 genome assembly and annotation completeness with single-copy orthologs. Bioinformatics

495 (Oxford, England). 2015;31 19:3210-2. doi:10.1093/bioinformatics/btv351.

496 8. Parra G, Bradnam K and Korf I. CEGMA: a pipeline to accurately annotate core genes in

497 eukaryotic genomes. Bioinformatics (Oxford, England). 2007;23 9:1061-7.

498 doi:10.1093/bioinformatics/btm071.

499 9. Parra G, Bradnam K, Ning Z, Keane T and Korf I. Assessing the gene space in draft genomes.

500 Nucleic Acids Res. 2009;37 1:289-97. doi:10.1093/nar/gkn916.
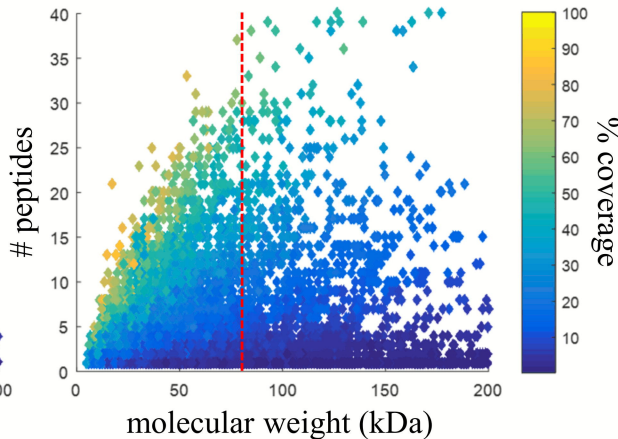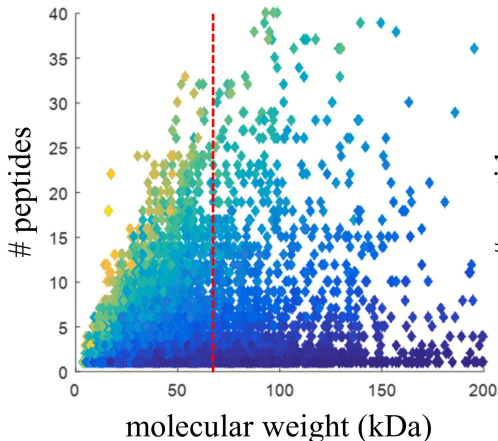
501    10.    Lam H, Deutsch EW, Eddes JS, Eng JK, Stein SE and Aebersold R. Building consensus spectral

502           libraries for peptide identification in proteomics. Nature methods. 2008;5 10:873-5.

503           doi:10.1038/nmeth.1254.

504    11.    Burke MC, Mirokhin YA, Tchekhovskoi DV, Markey SP, Heidbrink Thompson J, Larkin C, et

505           al. The Hybrid Search: A Mass Spectral Library Search Method for Discovery of Modifications

506           in Proteomics. Journal of proteome research. 2017;16 5:1924-35.

507    12.    Zhang Z, Burke M, Mirokhin YA, Tchekhovskoi DV, Markey SP, Yu W, et al. Reverse and

508           Random Decoy Methods for False Discovery Rate Estimation in High Mass Accuracy Peptide

509           Spectral Library Searches. Journal of proteome research. 2018;

510           doi:10.1021/acs.jproteome.7b00614.

511    13.    Bekker-Jensen DB, Kelstrup CD, Batth TS, Larsen SC, Haldrup C, Bramsen JB, et al. An

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

513           systems. 2017;4 6:587-99.e4. doi:10.1016/j.cels.2017.05.009.

514    14.    Foote AD, Liu Y, Thomas GW, Vinar T, Alfoldi J, Deng J, et al. Convergent evolution of the

515           genomes of marine mammals. Nature genetics. 2015;47 3:272-5. doi:10.1038/ng.3198.

516    15.    Sobolesky P, Parry C, Boxall B, Wells R, Venn-Watson S and Janech MG. Proteomic Analysis of

517           Non-depleted Serum Proteins from Bottlenose Dolphins Uncovers a High Vanin-1 Phenotype.

518           Scientific reports. 2016;6:33879. doi:10.1038/srep33879.

519    16.    Jones SJ, Taylor GA, Chan S, Warren RL, Hammond SA, Bilobram S, et al. The Genome of the

520           Beluga Whale (*Delphinapterus leucas*). Genes. 2017;8 12:378. doi:10.3390/genes8120378.

521    17.    Lindblad-Toh K, Garber M, Zuk O, Lin MF, Parker BJ, Washietl S, et al. A high-resolution map

522           of human evolutionary constraint using 29 mammals. Nature. 2011;478 7370:476-82.

523           doi:10.1038/nature10530.

524    18.    Pruitt KD, Brown GR, Hiatt SM, Thibaud-Nissen F, Astashyn A, Ermolaeva O, et al. RefSeq: an

525           update on mammalian reference sequences. Nucleic Acids Res. 2014;42 Database issue:D756-63.

526           doi:10.1093/nar/gkt1114.

527    19.    NCBI *Tursiops truncatus* Annotation Release 101 Annotation Report.

528          https://www.ncbi.nlm.nih.gov/genome/annotation_euk/Tursiops_truncatus/101/. Accessed 13

529          January 2017.

530    20.    Shamsi A and Bano B. Journey of cystatins from being mere thiol protease inhibitors to at heart

531          of many pathological conditions. International journal of biological macromolecules.

532          2017;102:674-93. doi:10.1016/j.ijbiomac.2017.04.071.

533    21.    Neely BA, Carlin KP, Arthur JM, McFee WE and Janech MG. Ratiometric Measurements of

534          Adiponectin by Mass Spectrometry in Bottlenose Dolphins (*Tursiops truncatus*) with Iron

535          Overload Reveal an Association with Insulin Resistance and Glucagon. Frontiers in

536          endocrinology. 2013;4:132. doi:10.3389/fendo.2013.00132.

537    22.    Deutsch EW, Csordas A, Sun Z, Jarnuczak A, Perez-Riverol Y, Ternent T, et al. The

538          ProteomeXchange consortium in 2017: supporting the cultural change in proteomics public data

539          deposition. Nucleic Acids Research. 2017;45 D1:D1100-D6. doi:10.1093/nar/gkw936.

540    23.    Vizcaíno JA, Deutsch EW, Wang R, Csordas A, Reisinger F, Ríos D, et al. ProteomeXchange

541          provides globally co-ordinated proteomics data submission and dissemination. Nature

542          biotechnology. 2014;32 3:223-6. doi:10.1038/nbt.2839.

543    24.    Neely BA, Soper JL, Gulland FM, Bell PD, Kindy M, Arthur JM, et al. Proteomic analysis of

544          cerebrospinal fluid in California sea lions (*Zalophus californianus*) with domoic acid toxicosis

545          identifies proteins associated with neurodegeneration. Proteomics. 2015;15 23-24:4051-63.

546          doi:10.1002/pmic.201500167.

547    25.    Ordonez YN, Anton RF and Davis WC. Quantification of total serum transferrin and transferrin

548          sialoforms in human serum; an alternative method for the determination of carbohydrate-deficient

549          transferrin in clinical samples. Analytical Methods. 2014;6 12:3967-74.

550          doi:10.1039/C4AY00159A.

551  26.  Johnson SP, Venn-Watson SK, Cassle SE, Smith CR, Jensen ED and Ridgway SH. Use of

552       phlebotomy treatment in Atlantic bottlenose dolphins with iron overload. Journal of the American

553       Veterinary Medical Association. 2009;235 2:194-200. doi:10.2460/javma.235.2.194.

554  27.  Leon IR, Schwammle V, Jensen ON and Sprenger RR. Quantitative assessment of in-solution

555       digestion efficiency identifies optimal protocols for unbiased protein analysis. Molecular &

556       cellular proteomics : MCP. 2013;12 10:2992-3005. doi:10.1074/mcp.M112.025585.

557  28.  Bryk AH and Wisniewski JR. Quantitative Analysis of Human Red Blood Cell Proteome. Journal

558       of proteome research. 2017;  doi:10.1021/acs.jproteome.7b00025.

559  29.  Martens L and Vizcaino JA. A Golden Age for Working with Public Proteomics Data. Trends in

560       biochemical sciences. 2017;42 5:333-41. doi:10.1016/j.tibs.2017.01.001.

561  30.  Jagtap PD, Johnson JE, Onsongo G, Sadler FW, Murray K, Wang Y, et al. Flexible and

562       accessible workflows for improved proteogenomic analysis using the Galaxy framework. Journal

563       of proteome research. 2014;13 12:5898-908. doi:10.1021/pr500812t.

564  31.  Sheynkman GM, Johnson JE, Jagtap PD, Shortreed MR, Onsongo G, Frey BL, et al. Using

565       Galaxy-P to leverage RNA-Seq for the discovery of novel protein variations. BMC genomics.

566       2014;15:703. doi:10.1186/1471-2164-15-703.

567  32.  Pugh RS, Becker PR, Porter BJ, Ellisor MB, Moors AJ and Wise SA. Design and Applications of

568       the National Institute of Standards and Technology's (NIST's) Environmental Specimen Banking

569       Programs. Cell Preservation Technology. 2008;6 1:59-72. doi:10.1089/cpt.2007.0517.

570  33.  NCBI RefSeq. ftp://ftp.ncbi.nih.gov/genomes/refseq/. Accessed 28 June 2017.

571  34.  Clark K, Karsch-Mizrachi I, Lipman DJ, Ostell J and Sayers EW. GenBank. Nucleic Acids Res.
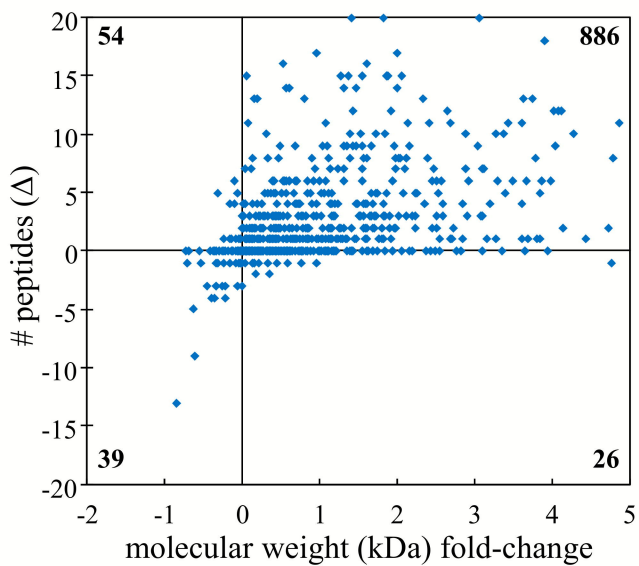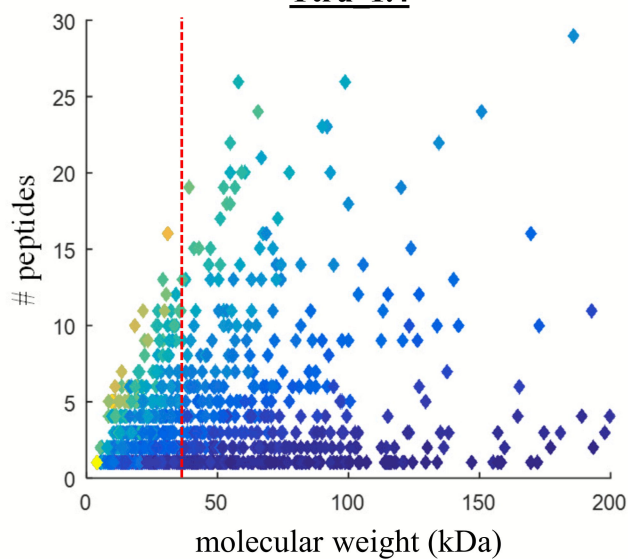
572       2016;44 D1:D67-72. doi:10.1093/nar/gkv1276.

573

**Figure 1**

| NCBI RefSeq category | Relative % | | Median MW (kDa) | |
|---|---|---|---|---|
| | Ttru 1.4 | NIST | Ttru 1.4 | NIST (Δ) |
| new | NA | 1.8 % | NA | 20.42 |
| minor | 55.5 % | 65.5 % | 47.24 | 59.07 (+25.0 %) |
| major | 24.8 % | 31.5 % | 36.97 | 41.08 (+11.1 %) |
| identical | 0.6 % | 0.4 % | 34.65 | 33.59 (-3.1 %) |
| other | 0.9 % | 0.7 % | 66.67 | 40.28 (-39.6 %) |
| deprecated | 18.0 % | NA | 28.66 | NA |
| YP-prefix | 0.2 % | 0.1 % | 35.76 | 35.76 (0.0 %) |
| | | total | 41.07 | 52.26 (+27.3 %) |

**Figure 2**

**Figure 3**

**Figure 4**

Ttru_1.4
(n = 28 293)

NIST_Tur_tru v1
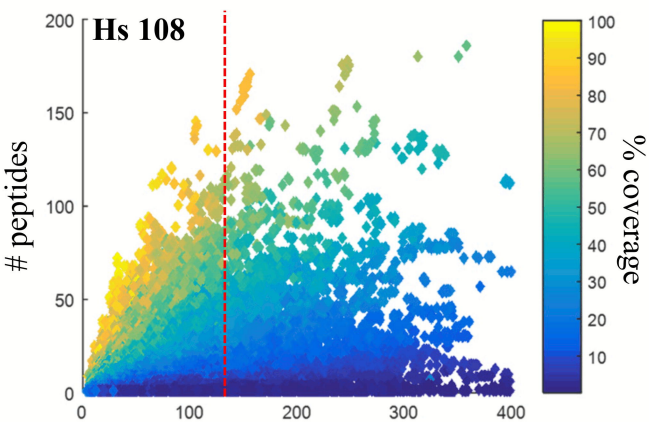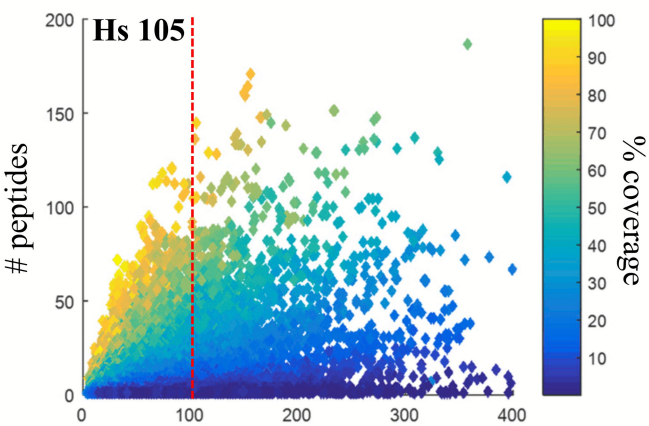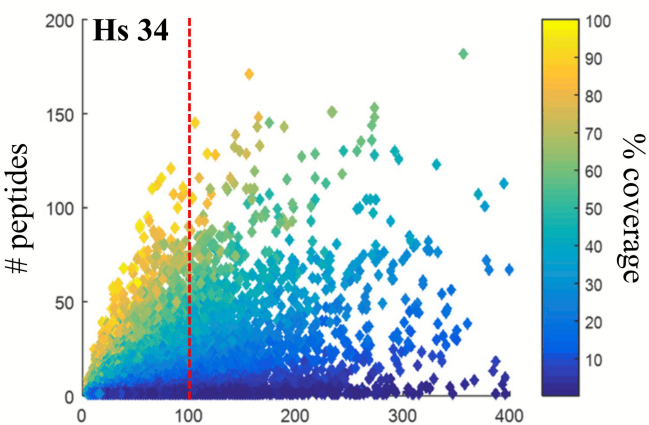(n = 27 404)

5657

22 636

4768

minor  major
other  deprecated

minor  major
other  new

**Figure 5** molecular weight (kDa)