

1 **Proteotranscriptomics assisted gene annotation and spatial proteomics of *Bombyx***

2 ***mori* BmN4 cell line**

3 Michal Levin\*, Marion Scheibe and Falk Butter\*

4 Institute of Molecular Biology (IMB), Ackermannweg 4, 55128 Mainz, Germany

5

6 \* To whom correspondence should be addressed:

7 Michal Levin or Falk Butter

8 Institute of Molecular Biology (IMB), Ackermannweg 4, 55128 Mainz, Germany

9 Phone: +49-6131-39-21573; +49-6131-39-21570

10 E-mail: [m.levin@imb.de](mailto:m.levin@imb.de); [f.butter@imb.de](mailto:f.butter@imb.de)

## Genome-free proteotranscriptomics assisted gene annotation

11 **Abstract**12 **Background**

13 The process of identifying all coding regions in a genome is crucial for any study at the level  
14 of molecular biology, ranging from single-gene cloning to genome-wide measurements using  
15 RNA-Seq or mass spectrometry. While satisfactory annotation has been made feasible for  
16 well-studied model organisms through great efforts of big consortia, for most systems this  
17 kind of data is either absent or not adequately precise.

18 **Results**

19 Combining in-depth transcriptome sequencing and high resolution mass spectrometry, we  
20 here use proteotranscriptomics to improve gene annotation of protein-coding genes in the  
21 *Bombyx mori* cell line BmN4 which is an increasingly used tool for the analysis of piRNA  
22 biogenesis and function. Using this approach we provide the exact coding sequence and  
23 evidence for more than 6,200 genes on the protein level. Furthermore using spatial  
24 proteomics, we establish the subcellular localization of thousands of these proteins. We show  
25 that our approach outperforms current *Bombyx mori* annotation attempts in terms of accuracy  
26 and coverage.

27 **Conclusions**

28 We show that proteotranscriptomics is an efficient, cost-effective and accurate approach to  
29 improve previous annotations or generate new gene models. As this technique is based on  
30 de-novo transcriptome assembly, it provides the possibility to study any species also in the  
31 absence of genome sequence information.

32

Genome-free proteotranscriptomics assisted gene annotation

### 33 **Keywords**

34 Proteotranscriptomics, mass spectrometry, gene assembly, gene annotation, spatial  
35 proteomics

### 36 **Background**

37 *Bombyx mori* was the first lepidopteran species whose draft genome was published in 2004  
38 [1,2]. In 2008, a more accurate genome assembly was generated by combining the raw data  
39 of these initial efforts within an international collaboration [3], and the results are available at  
40 SilkDB ([www.silkdb.org](http://www.silkdb.org)) and KAIKObase (sgp.dna.affrc.go.jp/index.html). However, for a large  
41 number of modern applications such as transcriptomic, epigenomic and proteomic studies,  
42 reverse genetic screens and genome editing tools such as TALEN and CRISPR/Cas9 the  
43 provided genome information is insufficient as this assembly contains numerous non-  
44 sequenced chromosome regions. Recently, parallel to our efforts to reannotate *Bombyx mori*  
45 using proteotranscriptomics, a new initiative provided an improved genome assembly [4]  
46 including revised gene predictions for 16,880 gene models. This new assembly has been  
47 made available at SilkBase and includes more genomic regions. However, the provided gene  
48 models are still based on automated gene prediction using limited full-length cDNA libraries,  
49 poly-A RNA-Seq data and previous *B. mori* NCBI annotations. These predictions are made  
50 with a mixture of data from various commercial and non-commercial strains of *Bombyx mori*,  
51 thus may not represent the genomic sequence of a single strain or its derived cell line due to  
52 intraspecies genetic variation. Furthermore, only very few predicted *Bombyx mori* genes have  
53 evidence at the protein level (roughly 150 genes in UniProt UP000005204). Hence, an  
54 improved strain-specific gene annotation would likely improve global analyses.

## Genome-free proteotranscriptomics assisted gene annotation

55 *Bombyx mori* is similar to humans in terms of sensitivities to pathogens and comparable  
56 effects of drugs. The advantages for research are the low cost of maintenance, little ethical  
57 constraints and no biohazard risks. Hence, it has been long recognized as an excellent  
58 system for drug screening and safety assessment [5], [6]. Furthermore, the BmN4 cell line of  
59 *Bombyx mori* [7] has been intensely used in studying many different biological aspects in the  
60 laboratory, including virus infection [8], [9] and germline piRNA biology [10]. Despite the  
61 widespread usage of this ovary-derived cell line, its exact genomic sequence and a tested  
62 gene structure model is still missing as this cell line does not originate from the genome  
63 sequenced *Bombyx mori* strains. Natural sequence variations will interfere with primer design  
64 for gene amplification, the design of CRISPR guides and limit matching of short sequencing  
65 reads in transcriptomics as well as peptide coverage in proteomic experiments.

66 By proteotranscriptomics, combining in-depth transcriptome sequencing and high-resolution  
67 mass spectrometry, we establish protein evidence for 6,273 genes. Using spatial proteomics,  
68 we additionally experimentally classify the localization of several thousand proteins using the  
69 recently published LOPIT-DC workflow [11], which utilizes differential ultracentrifugation  
70 following removal of unlysed cells to achieve enrichment of cellular compartments in different  
71 fraction. With this strategy, not only accurate gene models for the protein-coding genes, but  
72 also their subcellular localization is made available. Overall, we provide here a proof of  
73 concept for the generation of species-, strain- and cell line-specific gene annotation for  
74 protein-coding genes based on experimental evidence without the need of a sequenced  
75 genome.

Genome-free proteotranscriptomics assisted gene annotation

## 76 Results

### 77 Transcriptome assembly from RNA-Seq data

78 We generated 167.8 million paired reads (86.2 million reads from poly-A enriched and 81.6  
79 million from rRNA depleted total RNA samples) by Illumina paired-end sequencing of RNA  
80 prepared from the BmN4 cell line of *Bombyx mori*. To prove our concept of being able to  
81 produce a proper annotation without a genome and due to the unclear genetic background of  
82 the cell line, we applied a genome-free *de-novo* approach using the Trinity suite [12] (Fig 1A).  
83 After quality filtering, adapter trimming and erroneous k-mer removal almost 165 M paired  
84 reads and 158,589,380 bases were assembled into 186,401 ‘Trinity transcripts’ constituting  
85 120,287 distinct ‘Trinity genes’. The assembled raw transcriptome represents 98.32% of the  
86 input reads, which shows that the assembly is highly representative (Additional file 1: Table  
87 S1). The traditional N50 statistics describe the minimal transcript length of transcripts that are  
88 assembled from at least 50% of the reads. We found the N50 length for our assembly to be  
89 1553 bases. A better representation excludes lowly expressed transcripts as they might  
90 exhibit bigger biases. Hence, we investigated the N50 values across different expression level  
91 bins (ExN50) (Additional file 1: Fig S1). We found that the ExN50 peaks between the 80 and  
92 90 expression percentiles. Thus a better representation is the E90N50 statistic, which  
93 represents the minimal transcript length of transcripts in the 90<sup>th</sup> expression percentile that  
94 are assembled from at least 50% of the reads mapping to these transcripts. The E90N50  
95 transcript contig length is 2270 bases (Additional file 1: Table S2). We used TransRate [13] to  
96 validate the quality of the raw assembly (Additional file 1: Table S3). TransRate assesses the  
97 accuracy and completeness of a transcriptome assembly using only the input reads. It

## Genome-free proteotranscriptomics assisted gene annotation

98 proceeds by mapping the reads to the assembled contigs, analyzing the alignments,  
99 calculating metrics for each individual contig, integrating these contig-level metrics to provide  
100 a contig score, and then combining the accuracy of the assembly with the score of each  
101 contig to produce an overall assembly score. The crude overall and optimal Transrate  
102 assembly score is 0.31 and 0.41, respectively, of which both are in the 70<sup>th</sup> percentile range of  
103 Transrate assembly scores of 155 published *de-novo* assembled transcriptomes  
104 [13] (Additional file 1: Fig S2). The expression-level-weighted assembly score, which weights  
105 each constituent contig score by the relative abundance level of each contig raises to 0.54  
106 validating the high quality of the assembly and indicating that most of the low TransRate  
107 scores stem from contigs that are of relatively low abundance.

108 To extract all potential protein-coding transcripts from the assembled contigs, we applied  
109 TransDecoder [14] and kept transcripts that comprise an open reading frame of at least 20  
110 amino acids. This filtering resulted in a list of 317,031 potential protein-coding open reading  
111 frames based on 95,817 individual genes. These potential protein-coding sequences are the  
112 basis of our further analysis. Using BUSCO [15], we detected that our assembled protein  
113 coding transcripts cover 94.8% of the arthropod BUSCO gene set (Fig 1B and Additional file  
114 1: Table S4). Currently, three main initiatives have provided *Bombyx mori* genome assembly  
115 and annotation. For precision estimation, we compared our data to the currently available  
116 annotations of the different *Bombyx mori* varieties from UniProt UP000005204 with 14,776  
117 gene models [16], NCBI Annotation Release 101 with 14,998 gene models, and SilkBase  
118 2017 with 16,880 gene models [4]. In general, correspondence between current annotations  
119 and *de-novo* predicted proteins is high, with the majority of transcripts sharing protein full  
120 sequence coverage (90-100%) to the respective UniProt, NCBI and SilkBase proteins (Fig

## Genome-free proteotranscriptomics assisted gene annotation

121 1C). When analyzing the genetic sequence variation between the BmN4 cell line and the  
122 transcripts from the NCBI and SilkBase annotation by mapping the RNA-Seq data to the  
123 respective CDS sequences of predicted gene models, we found on average exchange of 1 in  
124 129 bases for NCBI annotations (i.e. 126,300 changes in 16,393,027 bases) and 1 in 105  
125 bases for SilkBase annotations (i.e. 187,534 changes in 19,826,985 bases). The results of  
126 both comparisons are highly consistent. Approximately 75% of the detected changes are  
127 silent (synonymous) mutations while around 25% have missense (non-synonymous) and  
128 0.4% nonsense effects (Additional file 1: Table S5) . These results unveil an unexpected quite  
129 large variation between *Bombyx mori* strains and emphasize the importance of applying a  
130 genome-free approach to provide exact CDS sequences especially for molecular biology  
131 applications.

132 Functional annotation of TransDecoder predicted protein sequences was performed using  
133 Trinotate [17] including blastp searches against all model species Swissprot databases,  
134 HMMER to identify protein domains, signalP to predict signal peptides, tmHMM to predict  
135 transmembrane regions and RNAMMER to identify rRNA transcripts. Furthermore, we used  
136 deeploc to predict protein localizations from the respective protein sequences. All functional  
137 annotations are included in the annotation file (Additional file 2: Table S10).

138

139 **High resolution mass spectrometry data provides peptide evidence for protein coding  
140 transcripts**

141 In order to provide evidence for the protein coding capability of our predicted protein coding  
142 open reading frames (ORF), we performed mass spectrometry measurements of protein  
143 extracts from the BmN4 cell line. Using LOPIT-DC [11] as a strategy to fractionate our

## Genome-free proteotranscriptomics assisted gene annotation

sample, we aimed to increase protein detection depth, while also gaining cellular localization information for the detected protein sequences. Our full *Bombyx mori* proteome data set with 10 fractions contained 3,685,257 MS/MS spectra. Using the predicted open reading frames of the Trinity transcriptome assembly as search space for the mass spectrometry data, we noted a higher rate of identified MS/MS spectra compared to using the three currently available protein annotations from UniProt, NCBI and SilkBase (two-sided paired Wilcoxon signed rank test,  $p=3.7*10^{-8}$ ,  $p=4*10^{-8}$ , and  $p=3.7*10^{-8}$ , respectively (Fig 2A). Applying stringent filtering criteria to have at least 2 identified peptides (at least one of them being unique), we identified a total of 6,273 protein groups (fasta files of the CDS and protein sequences of these proteins are provided Material S1 and S2, respectively). This was 16%, 18% and 14% higher than for the currently available UniProt (5,254 identified protein groups), NCBI (5,125 identified protein groups) and SilkBase annotations (5,396 identified protein groups) (Fig 2B), emphasizing the power of strain specific gene sequences to increase proteome coverage and also validating the high quality of our genome-free *de-novo* transcriptome assembly. In order to investigate whether peptide identification could be hindered by provided protein annotations that include strain specific differences in protein sequences, we extracted peptide identification for relevant proteins from both Trinity and SilkBase annotations and chose as representatives those pairs that have an overlap of more than 80% in sequence but are not 100% identical in their protein sequence. We extracted all missense mutations that were identified for these annotation pairs and calculated for the respective locations the proportion of peptides that were not identified in the SilkBase search, although they could be identified in our *de-novo* Trinity annotation. We found that 88% of peptides assigned to predicted missense mutations

## Genome-free proteotranscriptomics assisted gene annotation

166 (2325 out of 2653 peptides in 1988 protein groups affected) indeed hamper peptide  
167 identification when using the SilkBase annotation as search base.

168 Detected protein groups show improved quality statistics when compared to the raw *in silico*  
169 predicted potential protein coding transcripts, e.g. better overall correspondence with current  
170 UniProt, NCBI and SilkBase annotations (Fig 2C, Additional file 1: Fig S3A), higher assembly  
171 scores (Fig 2D) and longer transcript lengths (Additional file 1: Fig S3B). We further observed  
172 that assembled contigs with high TransRate scores are indeed enriched with identification by  
173 mass spectrometry emphasizing the validity of the scoring approach used by TransRate,  
174 which is based on the raw read alignment features only (Additional file 1: Fig S3C).

175 Although the correspondence between our annotation and the annotation provided by NCBI  
176 and SilkBase is overall high (Fig 2C), there are still almost 20% of predicted coding  
177 sequences that correspond with less than 80% hit percentage coverage to the current  
178 annotations. We noted that some of the current SilkBase annotated transcripts are split into  
179 several (mostly two) separated genes in our annotation. This observation can have two main  
180 explanations: either the current gene annotations are interdependent (SilkBase includes NCBI  
181 annotations in the prediction process) and thus an erroneous annotation from earlier  
182 predictions could have been transferred to the newest SilkBase annotation, or, the separation  
183 of the non-corresponding annotations in our genome-free approach is wrong. To decide  
184 between these two possibilities, we checked the RNA-Seq read coverage across the gap  
185 between two separated proteins in our annotations that were suggested by SilkBase to be a  
186 single protein (941 pairs in total corresponding to 631 SilkBase genes). For this we  
187 investigated the RNA-Seq reads coverage for each of the relevant pairs of our Trinity  
188 predicted proteins and the respective genomic gap between these using the SilkBase

## Genome-free proteotranscriptomics assisted gene annotation

189 genome assembly. Overall we observe that 76% of the Trinity annotated splits (629 out of 826  
190 protein pairs) are well supported by clear gaps in the RNA-Seq raw data alignment at the split  
191 site (Fig S4 A and B). Only 5.5% of all Trinity predicted proteins (348 unique proteins in 195  
192 protein pairs out of 6,273 detected proteins) do not show an evident gap in read coverage and  
193 hence likely have been falsely split in the annotation process. Our set of identified ORFs also  
194 includes 188 predictions that have a less than 85% hit coverage with a SilkBase annotation  
195 entry. The length differences are shown in Additional file 1: Fig S5A. In order to investigate if  
196 the shorter ORFs are supported by read mapping data, we calculated the difference between  
197 the read mapping frequency in the ORF region with the coverage at the edges of the  
198 transcripts. 124 (66% of the short ORFs) identified proteins show an evident absence of reads  
199 at the edges of the transcripts, while the remaining 64 might have been falsely split as we  
200 could observe mapped reads adjacent to the edges (Additional file 1: Fig S4C). Based on  
201 these observations, we conclude that our method has a precision rate of at least 93.4%  
202 (5,861 out of 6,273) for assigning individual genes. The respective categorizations into “high  
203 correspondence”, “evidently split”, “probable false split”, “evidently shorter than SilkBase” and  
204 “probably falsely shorter than SilkBase” have been included in the annotation table for clarity  
205 (Additional file 2: Table S10). Furthermore, we found mass spectrometry evidence for 164  
206 predicted proteins that have been missed in any of the current annotations (marked with  
207 “newly annotated” in the annotation Additional file 2: Table S10, peptide evidence information  
208 for all novel proteins are provided in Additional file 4: Table S12). Another group of genes (513  
209 genes (8% of all ORFs detected by MS); marked with “longer than SilkBase”) are longer than  
210 annotated in the current SilkBase annotation, however differences are mostly neglectable

## Genome-free proteotranscriptomics assisted gene annotation

211 (Additional file 1: Fig S5B). We provide a website (<http://butterlab.org/bombyxviewer>) which  
212 incorporates data regarding all ORFs with mass spectrometric evidence including transcript  
213 and protein sequences and a genome viewer based on the SilkBase genome showing gene  
214 structure and individual RNA-Seq read mapping. Peptide evidence information and annotated  
215 MS-MS spectra of newly identified proteins are therewith also provided for download.

216

217 **Subcellular localization of proteins determined by LOPIT-DC**

218 To assign the *Bombyx mori* proteins to sub-cellular compartments, we performed label-free  
219 quantitative spatial proteomics adapting the recently released LOPIT-DC protocol [11]. Using  
220 differential centrifugation steps, we generated 10 subcellular fractions in independent  
221 quadruplicates. Fractionation replicates correlate very well (average Pearson correlation  
222 coefficient >0.97) (Fig 3A) and cluster together in the first two PCA dimensions (Fig 3B).

223 Within the gradient, fraction series [1-3], [4-7] and [8-9] are similar to each other, while fraction  
224 10 is most different from all others. This is consistent with fraction 10 constituting an acetone  
225 precipitation of all proteins that were not separated in the previous fractionation steps.

226 Analyzing significant changes in pairwise comparisons of all fractions, efficient separation can  
227 be recapitulated, i.e. an increasing diversity of proteins can be observed throughout  
228 subsequent fractionation steps (Fig 3C). LFQ data and differential expression statistics across  
229 fractionation samples can be found in Additional file 3: Table S11.

230 As there are no experimentally validated marker proteins for specific cell compartments in  
231 *Bombyx mori*, we applied an unsupervised clustering approach to detect unique fractionation  
232 profiles. Unsupervised clustering of the normalized protein intensities of all proteins showing  
233 significant changes in at least one of the pairwise fractionation comparisons and appropriate

## Genome-free proteotranscriptomics assisted gene annotation

234 clustering assessment and filtering (see Methods and Additional file 1: Fig S6) revealed 8  
235 main clusters encompassing 3,942 protein groups as depicted in Fig 4A. The mean  
236 expression profiles across fractionation of the different clusters are depicted in Fig 4B (and in  
237 Additional file 1: Fig S7 for individual genes). In order to characterize the different clusters in  
238 terms of their potential subcellular localization or function, enrichment of specific categories  
239 as determined by signalP (prediction of signal peptides), tmHMM (prediction of  
240 transmembrane regions) and the Gene Ontology cellular-component (GO\_cc) annotation  
241 were calculated (Fig 4C, Additional file 1: Table S7). Clusters 1-4 represent proteins that show  
242 relatively high intensities in the early fractionation steps with diminished intensities in later  
243 steps. Generally, membrane associated proteins (framed with black box in Fig 4C) are highly  
244 enriched in clusters 2 and 3. To get a more specific insight into the subcellular localization, we  
245 subsequently checked for enrichment of the GO\_cc terms ‘lysosome’, ‘peroxisome’, ‘golgi  
246 apparatus’, ‘nucleus’, ‘chromatin’, ‘endoplasmic reticulum’ (ER), ‘mitochondrion’ and  
247 ‘ribosomes’, which were inferred by orthology to well-annotated model organisms in each  
248 cluster (framed with red box in Fig 4C). The most prominent enrichment was observed for  
249 cluster 1, which exhibits high levels in the first three fractions and low levels in later fractions.  
250 This cluster is highly enriched with mitochondrial genes (Fisher’s exact test,  $P = 10^{-197}$ , fold-  
251 enrichment = 35.56). Cluster 2 shows reduced intensity after fractionation step 4 and is  
252 enriched with peroxisome, ER and lysosome annotated proteins (Fisher’s exact test,  $P = 10^{-16}$   
253 ,  $P = 10^{-6}$  and  $10^{-5}$ , fold-enrichment = 24.14, 3.2 and 4, respectively). Cluster 3 has lower  
254 intensity starting at fractionation step 5 and represents mainly endoplasmatic reticulum (ER)  
255 proteins (Fisher’s exact test,  $P = 10^{-7}$ , fold-enrichment = 3.37). The profile of cluster 4 shows  
256 reduction of protein intensities after fractionation step 8 and contains proteins from Golgi

## Genome-free proteotranscriptomics assisted gene annotation

257 apparatus (Fisher's exact test,  $P = 10^{-5}$ , fold-enrichment = 2.9). The second highest  
258 enrichment could be observed for cluster 5, where measured protein intensities peak in  
259 fraction 7-9 and which is highly enriched with ribosomal 40s and 60s proteins (Fisher's exact  
260 test,  $P = 10^{-25}$ , fold-enrichment = 11.76). Cluster 6 encompasses proteins with low abundance  
261 in the initial fractionation steps, which increase until step 9 and decline to minimal levels in  
262 fraction 10. This cluster exhibits a mixed enrichment profile of nuclear and chromatin  
263 associated proteins, but also with ribosomal proteins (Fisher's exact test,  $P = 10^{-14}$ ,  $P = 10^{-3}$   
264 and  $P = 10^{-13}$ , fold-enrichment = 2.62, 2.97 and 6.92, respectively). Cluster 7 is the only  
265 cluster, that show exclusively high enrichment with nucleus associated proteins (Fisher's  
266 exact test,  $P = 10^{-19}$ , fold-enrichment = 3.18). Cluster 8, which shows abrupt protein intensity  
267 increase in the last two fractions couldn't be associated with any of the known localizations  
268 and hence was not further analyzed. Overall the cellular localization profiles of the different  
269 clusters correspond very well (average Spearman's correlation coefficient = 0.88) with those  
270 established in the LOPIT-DC method paper [11] where the assignment was established using  
271 experimentally validated marker genes (Additional file 1: Fig S8). This shows that the  
272 fractionation strategy is robust and widely applicable. Comparing the fractionation profile of  
273 each protein to each of the clusters enables localization prediction also for proteins that could  
274 not be annotated properly by previous *in silico* analysis. The resulting clusters allow to assign  
275 proteins especially to the following 6 compartments (in descending order of certainty):  
276 Mitochondria (cluster 1), Ribosome (cluster 5), Nucleus (cluster 7), Peroxisome (cluster 2),  
277 Endoplasmic reticulum (cluster 3) and Golgi apparatus (cluster 4). For each individual MS  
278 detected protein, we calculated the correlations between its expression profile across  
279 fractionations and the predicted median profile of the above depicted localization. These

## Genome-free proteotranscriptomics assisted gene annotation

280 correlation values are provided as confidence score for the localization probability of each  
281 protein. The information for each detected protein, including all relevant information such as  
282 transcripts type, length, score, annotations and localization categorization weights are  
283 provided in Additional file 2: Table S10).

## 284 Discussion

285 We here show that by combining comprehensive RNA-Seq data and high-resolution mass  
286 spectrometry data, we achieve a comparable and even slightly better annotation of protein-  
287 coding genes in *Bombyx mori* than previous efforts based on genome or transcriptome guided  
288 *de-novo* strategies. Even in the comparison of our genome-free to our own genome-guided  
289 assembly using the same raw RNA-Seq and mass spectrometry data a slightly better  
290 performance can be observed in the genome-free approach (see Additional file 1: Table S9).  
291 This fact emphasizes the importance of using genome-free approaches in conditions were  
292 provided genomes are suspected to stem from different genetic backgrounds as the  
293 measured system.

294 Our extensive comparison between our genome-free proteotranscriptomics annotation and  
295 the provided annotations from UniProt, NCBI and SilkBase showed that UniProt currently  
296 provides the annotation with the weakest performance. Gene annotations provided by NCBI  
297 and SilkBase are more comprehensive, however still do not report the full protein-coding  
298 potential demonstrated by the elevated percentage of identified MS-MS spectra in our tailor-  
299 made assembly (Fig 2A) and the identification of 164 new proteins in this study. Although the  
300 correspondence between our annotation and NCBI and SilkBase is overall high, there are still  
301 20% of coding sequences with less than 80% hit length, mostly attributed to split genes in our

## Genome-free proteotranscriptomics assisted gene annotation

302 assembly. Performing a detailed investigation of the gap region, we could provide evidence  
303 that in 80% of cases the split version in our assembly is supported, improving gene models  
304 for 451 genes.

305 Many studies have shown that small proteins ( $\leq 100$  amino acid residues) can be involved in  
306 important biological processes, including cell signaling, metabolism, and growth [18].  
307 However, they are underrepresented in many genome annotations as they are notoriously  
308 hard to predict because of their small ORF size [18]. To validate these small peptides, we kept  
309 all ORF with at least 20 amino acids. Indeed, we detected 308 small proteins (5% of all  
310 protein groups identified) with at least two peptides (one of them unique) and many of them  
311 reaching relatively high expression levels, however the overall expression levels are lower  
312 than for longer proteins (Additional file 1: Fig S9).

313 Additionally, we used the recently published LOPIT-DC approach [11] to provide experimental  
314 data for localization of our detected proteins. The high correspondence to the LOPIT-DC  
315 results, despite changing from TMT to LFQ and using cell lines from different species (human  
316 osteosarcoma U2OS vs. *B. mori* BmN4), indicate a universal applicability of the fractionation  
317 protocol and the resulting data. The resulting clusters allow assigning proteins especially to  
318 the following compartments: Mitochondria, Ribosome, Nucleus, Peroxisome, Endoplasmic  
319 reticulum and Golgi apparatus.

320 The current fractionation approach allowed us to detect peptide evidence for more than 6,200  
321 proteins. If more comprehensive databases are of interest, these limits may be overcome by  
322 using more diverse conditions or several different tissues for extraction of transcriptomic and  
323 proteomic data. While we here provide annotation for *Bombyx mori*, this approach is readily  
324 applicable to any species, including more complex organisms such as vertebrates and plants.

## Genome-free proteotranscriptomics assisted gene annotation

325 While possible parameters might need to be adapted such as even deeper RNA-Seq data,  
326 the high fraction of non-coding genome regions or highly repetitive structure that make  
327 genome assembly challenging can be disregarded.

328 Our developed proteotranscriptomics approach improves current gene annotations and  
329 provides the exact gene sequences for other applications such as gene amplifications via  
330 cDNA or planning CRISPR guides around the translation start site.

331 Importantly, the proposed annotation approach readily works without any genome reference  
332 and hence provides a precise, time- and cost-efficient method to construct annotations for  
333 protein-coding genes in any species where properly sequenced genomes are still out of  
334 reach.

335

## 336 **Conclusions**

337 Combining in-depth transcriptome sequencing and high resolution mass spectrometry, we  
338 here use proteotranscriptomics to improve gene annotation of protein-coding genes in the  
339 *Bombyx mori* cell line BmN4 which is an increasingly used tool for the analysis of piRNA  
340 biogenesis and function. Using this approach we provide the exact coding sequence and  
341 evidence for more than 6,200 genes on the protein level. We show that proteotranscriptomics  
342 is an efficient, cost-effective and accurate approach to improve previous annotations or  
343 generate new gene models. As this technique is based on de-novo transcriptome assembly, it  
344 provides the possibility to study any species also in the absence of genome sequence  
345 information.

346

## Genome-free proteotranscriptomics assisted gene annotation

347 **Methods**348 **Experimental design**

349 To build a genome-free proteotranscriptomics-based gene annotation we combined two types  
350 of data. First, RNA-Seq data of polyadenylated mRNA entities combining poly-A enriched and  
351 rRNA depleted samples from *Bombyx mori* BmN4 cell line was used to predict potential  
352 protein-coding genes. Secondly, MS/MS spectra data was used to find evidence for the  
353 predicted protein-coding genes. Profiling subcellular localization of proteins in *Bombyx mori*  
354 cells was performed using the LOPIT-DC strategy [11]. The procedures were performed with  
355 four biological replicates based on the high level of reproducibility between replicates  
356 (average Pearson correlation coefficient >0.97). Results are represented as averages of the  
357 biological replicates. We used Trinity for genome-free *de novo* RNA assembly and the  
358 MaxQuant data analysis platform [19] for quantitative proteomics analysis.

359

360 **Cell propagation and RNA extraction**

361 The *Bombyx mori* larval ovary-derived cell line BmN4 [7] was cultured in Insect media IPL-40  
362 (Pan Biotech) with 10% heat-inactivated FBS (Sigma) and 1x Penn-Strep (Sigma) at 27 °C.  
363 For RNA-Seq total RNA was extracted from 10 million cells using the RNAeasy Mini Kit  
364 (Qiagen) according to standard protocol. RNA integrity was tested by agarose gel  
365 electrophoresis and Bioanalyzer (RNA Nano Assay). RNA was quantified using Qubit.

366

367 **RNA-Seq measurements**

## Genome-free proteotranscriptomics assisted gene annotation

368 RNA-Seq libraries were prepared from total RNA using two different RNA enrichment  
369 protocols: 1. poly(A) purification using Illumina TruSeq stranded mRNA LT Sample Prep Kit  
370 following Illumina's standard protocol (Part # 15031047 Rev. E). [polyA-enriched] 2. depletion  
371 of ribosomal RNA using Illumina TruSeq stranded Total RNA LT Sample Prep Kit following  
372 Illumina's standard protocol (Part # 15031048 Rev. E) [ribo-minus].  
373 The libraries were prepared with a starting amount of 1000ng and amplified in 10 PCR cycles  
374 and profiled using a DNA 1000 Chip on a 2100 Bioanalyzer (Agilent technologies) and  
375 quantified using the Qubit dsDNA HS Assay Kit, in a Qubit 2.0 Fluorometer (Life  
376 technologies). The two libraries were pooled together in equimolar ratio and sequenced on 1  
377 NextSeq 500 Midoutput FC, PE for 2x 79 cycles plus 7 cycles for the index read. The  
378 measurements of polyA-enriched and ribo-minus RNA resulted in 86,178,436 and 81,597,503  
379 paired-end reads of length 80 bp, respectively. We assayed mycoplasma contamination by  
380 aligning all raw RNASeq forward reads to the genomes of *A.laidlawii*, *M.arginini*,  
381 *M.fermentans*, *M.hominis*, *M.hyorinis* and *M.orale*. The maximum percentage of uniquely  
382 mapped reads is below 0.00026 % and therefore a contamination can be excluded (see Table  
383 S8).

384

### 385 Transcriptome assembly

386 The two RNA-Seq datasets (polyA-enriched and ribo-minus RNA) were used in combination  
387 to assemble the transcriptome. First, both raw fastq files were cleaned from erroneous k-mers  
388 using Rcorrector [20] and the specialized scripts from TranscriptomeAssemblyTools  
389 (FilterUncorrectablePEfastq.py). Secondly, adapter sequences were removed using  
390 TrimGalore (a wrapper around Cutadapt [21] and FastQC [22]) and reads were filtered to

## Genome-free proteotranscriptomics assisted gene annotation

391 include only pairs consisting of proper pairs of minimum length of 36 nts each. These clean-  
392 up steps removed only 2% of the paired reads. The fastq files were then combined. For the  
393 genome-guided assembly raw RNA-Seq data was mapped to the *Bombyx mori* genome  
394 provided by SilkBase (<http://silkbbase.ab.a.u-tokyo.ac.jp/cgi-bin/index.cgi>) using STAR (version  
395 2.5.4b) [23]. The raw RNA-Seq or mapped data was used for a genome-free *de novo* or  
396 genome-guided assembly approach using the Trinity suite (Trinity version 2.4.0) [12] with the  
397 following parameter setting: [for genome-free: --seqType fq --SS\_lib\_type RF --min\_kmer\_cov  
398 1], [for genome-guided: Trinity --genome\_guided\_bam --genome\_guided\_max\_intron 30000 --  
399 genome\_guided\_min\_coverage 2]. The resulting Trinity fasta files were then further  
400 processed with TransDecoder version 5.4.0 [14] to predict potential protein coding transcripts  
401 using a length threshold of 20 amino acids. The resulting peptide fasta files were used as  
402 search space in subsequent steps for mass spectrometry data analysis.

403

404 **Quality check of transcriptome assembly**

405 The quality of the assembled transcriptome was assessed using several different state of the  
406 art approaches. These included general metrics of number of assembled transcripts, mean,  
407 median and Ex90N50 transcript lengths. The alignment rate of the raw reads to the assembly  
408 was calculated using Bowtie2 (version 2.3.4.3) [24] and dedicated scripts provided by Trinity  
409 (Trinity version 2.4.0) [12]. BUSCO (version 2.0) [15] with the arthropodae BUSCO database  
410 was used to assess the completeness of the assembly. Transrate scores and additional  
411 quality metrics were established using TransRate (version 1.0.3) [13]. Coherence with current  
412 annotations was measured using a combination of blastp (BLAST+ version 2.8.1) [25] and  
413 Trinity tools (Trinity version 2.4.0) [12]. For RNA-Seq coverage validations the combined

Genome-free proteotranscriptomics assisted gene annotation

414 cleaned RNA-seq data was mapped to the SilkBase genome assembly using STAR (version  
415 2.5.4b) [23]. Coverage per base was calculated using bedtools (version 2.26.0) [26] using the  
416 -pc option to also account for intronic alignment. Then using customized R (version 3.5.3) [27]  
417 scripts the average coverage per transcript or gap region was extracted (Fig S4).

418

419 **Annotation of identified transcripts**

420 Functional and domain annotations were produced using Trinotate (version 3.1.1)  
421 [17] combining the following applications: HMMER (version 3.2.1) [28] to identify protein  
422 domains, signalP (version 5.0) [29] to predict signal peptides, TMHMM (version 2.0c) [30] to  
423 predict transmembrane regions, RNAMMER (version1.2) [31] to identify rRNA transcripts in  
424 addition to infer Gene Ontology and KEGG terms from orthologies established by BLAST+  
425 (version 2.8.1 ) [25] with a swissprot database of all major model species. Further, localization  
426 predictions from protein sequences of the assembly were calculated using deeploc (version  
427 1.0) [32].

428

429 **Detection of variation level between NCBI and SilkBase protein coding sequences and**

430 **BmN4 RNASeq data**

431 In order to characterize the level of variation between the *Bombyx mori* coding sequences in  
432 the NCBI annotation and the BmN4 specific transcriptome, we first mapped the RNA-Seq  
433 data also used for transcriptome assembly to NCBI coding sequences using STAR aligner  
434 (version 2.5.4b) [23]. We then used GenomeAnalysisTK (version 3.8-0-ge9d806836) [33] to  
435 extract sequence variation information into a vcf file. SnpEff (version 4.4) [34] was used to  
436 annotate all sequence changes regarding their type (single nucleotide polymorphism (snp),

Genome-free proteotranscriptomics assisted gene annotation

437 deletion, insertion), their impact (low, moderate, high severeness) and their functional class  
438 (missense, nonsense, silent mutation). Results of this analysis are summarized in Table S5.

439

#### 440 Subcellular fractionation

441 Subcelluar fractionation of *Bombyx mori* cells was performed in quadruplicates and based on  
442 the LOPIT-DC strategy [11] with some modifications. Per replicate 70 million BmN4 cells were  
443 collected and suspended in 4 ml lysis buffer (0.25 M sucrose, 10 mM HEPES pH7.5, 2mM  
444 EDTA, 2mM Magnesium acetate, Roche complete protease inhibitors). Cells were dounced in  
445 a 7 ml glass douncer with 50 strokes using a type B pestle. Samples were fractionated using  
446 centrifugation speed and times as indicated in the table below following the LOPIT-DC  
447 strategy [11]. For centrifugation steps a Heraeus Fresco 21 centrifuge (Thermo) or a Optima  
448 Max-XP benchtop centrifuge (Beckmann) with a TLA 100.3 rotor were used. Supernatant of  
449 the 9<sup>th</sup> fractionation step was precipitated using ice cold acetone and the pellet was  
450 resuspended in 50mM HEPES/KOH pH7.9 (represents fraction 10).

451

#### 452 Mass spectrometric sample preparation and measurement

453 50 µg protein from each fraction were loaded on a 4-12 % NOVEX Bis-Tris PAGE gel  
454 (Thermo) and separated for 7 min at 180 V in 1x MES buffer (Thermo). Proteins were fixated  
455 and stained with Coomassie. After destaining with water, in-gel digest preparation and MS  
456 stage tip purification were performed as previously described [35], [36]. Peptides were  
457 analyzed by nanoflow liquid chromatography on an EASY-nLC 1000 system (Thermo)  
458 coupled to a Q Exactive Plus mass spectrometer (Thermo). Peptides were separated on a  
459 C18-reversed-phase column packed with Reprosil aq1.9 (Dr. Maisch GmbH), directly

## Genome-free proteotranscriptomics assisted gene annotation

460 mounted on the electrospray ion source of the mass spectrometer. We used a ca. 200  
461 minutes gradient from 2% to 60% acetonitrile in 0.1% formic acid at a flow of 225 nL/min. The  
462 Q Exactive Plus was operated with a Top10 MS/MS spectra acquisition method per MS full  
463 scan.

**464 Protein Identification and Label Free Quantification of Protein Intensities**

465 MaxQuant (version 1.5.2.8)<sup>12</sup> was used for raw file peak extraction and protein identification  
466 against the following databases individually: UniProt UP000005204 (14,776 entries) [16],  
467 NCBI *Bombyx mori* Annotation Release 101 (14,998 entries), SilkBase 2017 (16,880 entries)  
468 [4] or our Trinity-based ORF library (317,031 entries). Protein quantification was performed  
469 with MaxQuant using the label free quantification (LFQ) algorithm [37]. The following  
470 parameters were applied: trypsin as cleaving enzyme; minimum peptide length of seven  
471 amino acids; maximal two missed cleavages; carbamidomethylation of cysteine as a fixed  
472 modification; N-terminal protein acetylation and oxidation of methionine as variable  
473 modifications. Peptide mass tolerance was set to 20 ppm and 7 ppm was used as MS/MS  
474 tolerance. Further settings were: “label-free quantification” with “FastLFQ” disabled, “match  
475 between runs” with time window 0.7 minutes for matching and 20 minutes for alignment;  
476 peptide and protein false discovery rates (FDR) were set to 0.01; common contaminants  
477 (standard MaxQuant contaminant list including trypsin, keratin etc.) were excluded. Detailed  
478 settings are available in the respective parameter files (uploaded to ProteomeXchange  
479 ([www.ProteomExchange.org](http://www.ProteomExchange.org)) via the PRIDE [38] partner repository with the dataset identifier  
480 PXD014626). MaxQuant LFQ data was further processed using in-house developed tools  
481 based on R (version 3.5.3) [27]. This included filtering out marked contaminants, reverse  
482 entries and proteins only identified by site. Protein groups with no unique and less than two

## Genome-free proteotranscriptomics assisted gene annotation

483 peptides were removed. Protein group averages were calculated from proteotypic peptides.  
484 Prior to imputation of missing LFQ values with a beta distribution ranging from 0.1 to 0.2  
485 percentile within each sample, the values were log2 transformed. For protein groups  
486 consisting of more than one Trinity annotation, we chose the longest as representative of the  
487 group for further analysis.

488

#### 489 **LFQ data analysis and unsupervised clustering**

490 For overall statistical quality calculations, we calculated proportions of assembly scores,  
491 transcript lengths and correspondences with UniProt, NCBI and SilkBase *Bombyx mori*  
492 annotations. For clustering purposes, we focused only on protein groups that have four  
493 measured LFQ levels in at least one of the fractions and showed significant enrichment in at  
494 least one of the pairwise fraction comparisons ending up in 5,058 out of 5,610 protein groups  
495 (90%). The mean LFQ data of the four biological replicates data was standardized protein-  
496 group-wise by adding the mean of all average fraction intensities and dividing by the standard  
497 deviation between all average fraction intensities of the respective protein group. An  
498 unsupervised machine learning approach was used to cluster all standardized profiles. We  
499 applied the kohonen R package [39] to build a self-organizing-map (SOM), i.e. to build an  
500 artificial neural network that is trained using unsupervised learning to produce a two-  
501 dimensional, approximated grouping of the input profiles. The standardized data was initially  
502 assigned to a SOM matrix of 12 basic clusters. After combining four clusters with high intra-  
503 cluster differences (mean differences above the 75%-ile intra-cluster distances of all protein  
504 groups; see Additional file 1: Fig S6) into one cluster with uncategorized profiles, eight distinct  
505 clusters representing similar profiles remained. These clusters were ordered according to

## Genome-free proteotranscriptomics assisted gene annotation

506 euclidean distances in the first two-component PCA space using the TSP (travelling salesman  
507 problem) R package (version 1.1-7) [40]. Using localization data retrieved through orthology  
508 and sequence screening approaches as described earlier, we checked for enrichment of  
509 categories relevant for cellular localization, namely TmHH (sequence based transmembrane  
510 region predictions), signalP (sequence based signal peptides identification), Gene Ontology  
511 associations based on blast orthology associations for cellular component annotations  
512 lysosome, peroxisome, golgi, nucleus, chromatin, endoplasmic reticulum (ER), mitochondrion  
513 and ribosome (based on 40S and 60S ribosomal proteins) annotations. Enrichment scores of  
514 each category and each fractionation profile cluster were calculated using Fisher's Exact test  
515 (one-sided test, alternate-hypothesis: cluster genes contain more genes belonging to the  
516 tested category than non-cluster genes). Correlation of the identified cluster profiles were  
517 compared to the profiles of the same cellular localization categories from the original LOPIT-  
518 DC paper<sup>11</sup> by calculating the spearman correlation between the median standardized LFQ  
519 profiles of our data and the standardized median TMT profiles of the TMT data from human  
520 osteosarcoma U-2 OS cell line across fractionations (normalized TMT data is provided in  
521 Supplementary data 1 of previous publication [11]). Respective profiles and correlations are  
522 shown in Additional file 1: Fig S8.

523

524 **Declarations**525 **Ethics approval and consent to participate**

526 Not applicable

527

## Genome-free proteotranscriptomics assisted gene annotation

528 **Consent for publication**

529 Not applicable

530

531 **Availability of data and materials**

532 RNA-Seq raw data has been submitted to the Sequence Read Archive (SRA) under SRA ID

533 SRR9685281 [Reviewer access: <https://dataview.ncbi.nlm.nih.gov/object/PRJNA554660?>

534 reviewer=e8g3blh527ts0crsu4bcvnmrlk].

535 The mass spectrometry proteomics data have been deposited to the ProteomeXchange

536 Consortium via the PRIDE partner repository [38] with the dataset identifier PXD014626.

537 [Reviewer access: Username: [reviewer88577@ebi.ac.uk](mailto:reviewer88577@ebi.ac.uk); Password: gU1zhhAS]

538 All CDS and protein sequences of the assembled and identified proteins are provided in

539 Additional file 5 and 6.

540

541 **Competing interests**

542 The authors declare that they have no competing interests.

543

544 **Funding**

545 The funder had no role in the design of the study, analysis, interpretation of data and in writing

546 the manuscript.

547

548 **Authors' contributions**

## Genome-free proteotranscriptomics assisted gene annotation

549 ML, MS and FB conceived the project. MS and FB designed the experiments. ML analyzed  
550 the data. MS conducted the experiments. ML and FB wrote the paper. All authors approved  
551 the final manuscript.

552

553 **Acknowledgements**

554 We are indebted to Walter Bronkhorst for advice on culturing BmN4 cells that were kindly  
555 provided by the Ketting group (IMB). The RNA sequencing was conducted at the IMB  
556 Genomics Core Facility. Parts of this research were conducted using the supercomputer  
557 Mogon and advisory services offered by Johannes Gutenberg University Mainz ([hpc.uni-  
558 mainz.de](http://hpc.uni-mainz.de)), which is a member of the AHRP (Alliance for High Performance Computing in  
559 Rhineland Palatinate, [www.ahrp.info](http://www.ahrp.info)) and the Gauss Alliance e.V.. We gratefully  
560 acknowledge the computing time granted on the supercomputer Mogon at Johannes  
561 Gutenberg University Mainz ([hpc.uni-mainz.de](http://hpc.uni-mainz.de)) and Christian Meesters for helpful instructions  
562 regarding the implementation of Trinity and Transdecoder on Mogon. We thank Mario Dejung  
563 from the Proteomics Core Facility at the Institute of Molecular Biology for setting up the web-  
564 server and the dedicated data viewer on the website <http://butterlab.org/bombyxviewer>.

565

566 **References**

1. Mita K, Kasahara M, Sasaki S, Nagayasu Y, Yamada T, Kanamori H, et al. The Genome Sequence of Silkworm, *Bombyx mori*. *DNA Res.* 2004;11: 27–35. doi:10.1093/dnares/11.1.27
2. Xia Q, Zhou Z, Lu C, Cheng D, Dai F, Li B, et al. A draft sequence for the genome of the domesticated silkworm (*Bombyx mori*). *Science*. 2004;306: 1937–1940. doi:10.1126/science.1102210

## Genome-free proteotranscriptomics assisted gene annotation

3. The International Silkworm Genome Consortium. The genome of a lepidopteran model insect, the silkworm *Bombyx mori*. *Insect Biochem Mol Biol*. 2008;38: 1036–1045. doi:10.1016/J.IBMB.2008.11.004
4. Kawamoto M, Jouraku A, Toyoda A, Yokoi K, Minakuchi Y, Katsuma S, et al. High-quality genome assembly of the silkworm, *Bombyx mori*. *Insect Biochem Mol Biol*. 2019;107: 53–62. doi:10.1016/J.IBMB.2019.02.002
5. Nwibo DD, Hamamoto H, Matsumoto Y, Kaito C, Sekimizu K. Current use of silkworm larvae (*Bombyx mori*) as an animal model in pharmaco-medical research. *Drug Discov Ther*. 2015;9: 133–135. doi:10.5582/ddt.2015.01026
6. Abdelli N, Peng L, Keping C. Silkworm, *Bombyx mori*, as an alternative model organism in toxicological research. *Environ Sci Pollut Res*. 2018;25: 35048–35054. doi:10.1007/s11356-018-3442-8
7. Grace TD. Establishment of a line of cells from the silkworm *Bombyx mori*. *Nature*. 1967;216: 613.
8. Katsuma S, Kawamoto M, Shoji K, Aizawa T, Kiuchi T, Izumi N, et al. Transcriptome profiling reveals infection strategy of an insect maculavirus. *DNA Res Int J Rapid Publ Rep Genes Genomes*. 2018;25: 277. doi:10.1093/dnares/dsx056
9. Tsukui K, Yagisawa C, Fujimoto S, Ogawa M, Kokusho R, Nozawa M, et al. Infectious Virions of *Bombyx Mori* Latent Virus Are Incorporated into *Bombyx Mori* Nucleopolyhedrovirus Occlusion Bodies. *Viruses*. 2019;11: 316. doi:10.3390/v11040316
10. Kawaoka S, Hayashi N, Suzuki Y, Abe H, Sugano S, Tomari Y, et al. The *Bombyx* ovary-derived cell line endogenously expresses PIWI/PIWI-interacting RNA complexes. *RNA*. 2009;15: 1258–64. doi:10.1261/rna.1452209
11. Geladaki A, Kočevar Britovšek N, Breckels LM, Smith TS, Vennard OL, Mulvey CM, et al. Combining LOPIT with differential ultracentrifugation for high-resolution spatial proteomics. *Nat Commun*. 2019;10: 331. doi:10.1038/s41467-018-08191-w
12. Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, et al. Full-length transcriptome assembly from RNA-Seq data without a reference genome. *Nat Biotechnol*. 2011;29: 644–652. doi:10.1038/nbt.1883
13. Smith-Unna R, Boursnell C, Patro R, Hibberd JM, Kelly S. TransRate: reference-free quality assessment of de novo transcriptome assemblies. *Genome Res*. 2016;26: 1134–44. doi:10.1101/gr.196469.115
14. Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, et al. De novo transcript sequence reconstruction from RNA-seq using the Trinity platform for reference generation and analysis. *Nat Protoc*. 2013;8: 1494–1512. doi:10.1038/nprot.2013.084

## Genome-free proteotranscriptomics assisted gene annotation

15. Waterhouse RM, Seppey M, Simão FA, Manni M, Ioannidis P, Klioutchnikov G, et al. BUSCO Applications from Quality Assessments to Gene Prediction and Phylogenomics. *Mol Biol Evol.* 2018;35: 543–548. doi:10.1093/molbev/msx319
16. Duan J, Li R, Cheng D, Fan W, Zha X, Cheng T, et al. SilkDB v2.0: a platform for silkworm (*Bombyx mori*) genome biology. *Nucleic Acids Res.* 2010;38: D453–6. doi:10.1093/nar/gkp801
17. Bryant DM, Johnson K, DiTommaso T, Tickle T, Couger MB, Payzin-Dogru D, et al. A Tissue-Mapped Axolotl De Novo Transcriptome Enables Identification of Limb Regeneration Factors. *Cell Rep.* 2017;18: 762–776. doi:10.1016/j.celrep.2016.12.063
18. Scheidler CM, Kick LM, Schneider S. Ribosomal Peptides and Small Proteins on the Rise. *ChemBioChem.* 2019;20: cbic.201800715. doi:10.1002/cbic.201800715
19. Cox J, Mann M. MaxQuant enables high peptide identification rates, individualized p.p.b.-range mass accuracies and proteome-wide protein quantification. *Nat Biotechnol.* 2008;26: 1367–1372. doi:10.1038/nbt.1511
20. Song L, Florea L. Rcorrector: efficient and accurate error correction for Illumina RNA-seq reads. *GigaScience.* 2015;4: 48. doi:10.1186/s13742-015-0089-y
21. Martin M. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet.journal.* 2011;17: 10. doi:10.14806/ej.17.1.200
22. Andrews S. FastQC: a quality control tool for high throughput sequence data. Babraham Bioinformatics. 2010 [cited 10 Jul 2019]. Available: <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>
23. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics.* 2013;29: 15–21. doi:10.1093/bioinformatics/bts635
24. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat Methods.* 2012;9: 357–359. doi:10.1038/nmeth.1923
25. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+: architecture and applications. *BMC Bioinformatics.* 2009;10: 421. doi:10.1186/1471-2105-10-421
26. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics.* 2010;26: 841–842. doi:10.1093/bioinformatics/btq033
27. R Core Team. R: A language and environment for statistical computing (Version 3.5.3)[Computer software]. Vienna, Austria: R Foundation for Statistical Computing. R Foundation for Statistical Computing; 2019.
28. Eddy SR. Accelerated Profile HMM Searches. Pearson WR, editor. *PLoS Comput Biol.* 2011;7: e1002195. doi:10.1371/journal.pcbi.1002195

## Genome-free proteotranscriptomics assisted gene annotation

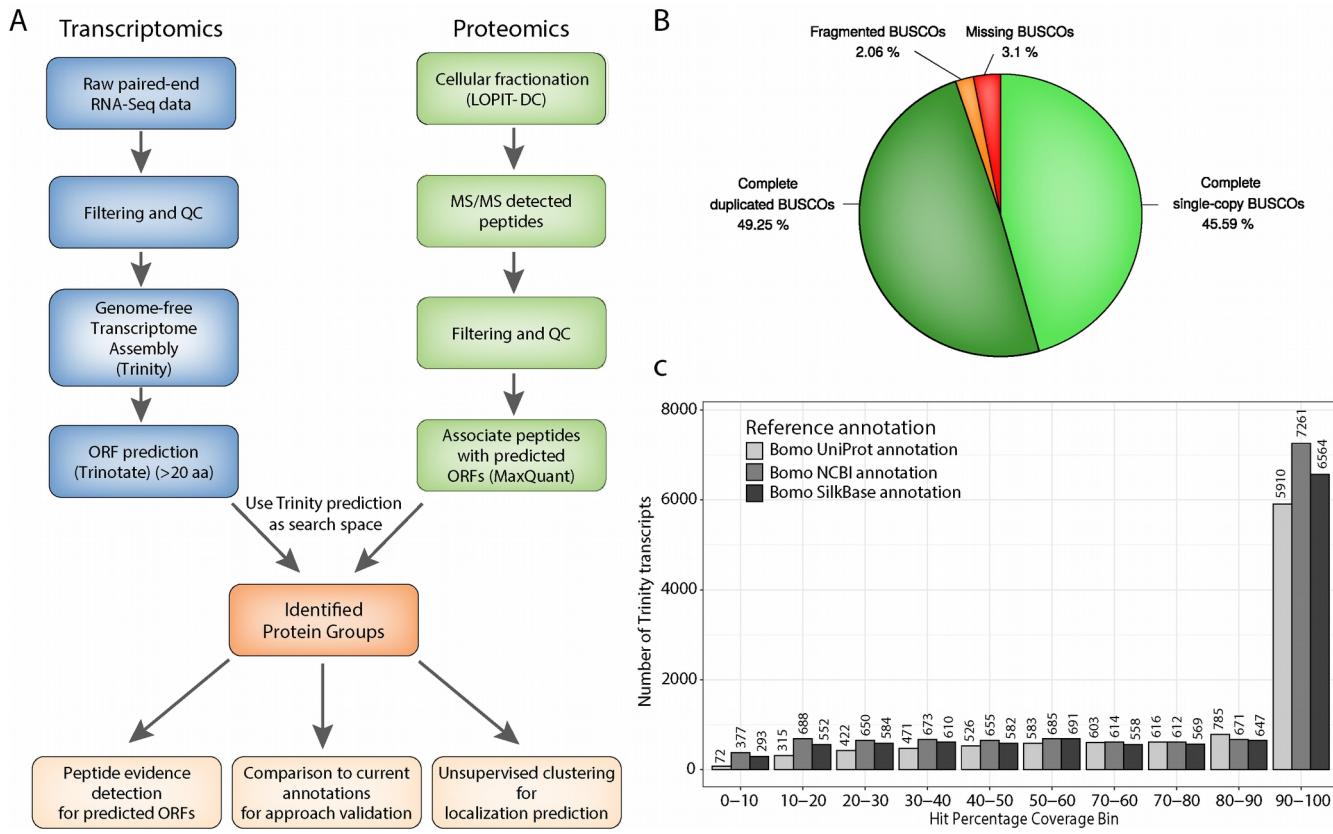
29. Almagro Armenteros JJ, Tsirigos KD, Sønderby CK, Petersen TN, Winther O, Brunak S, et al. SignalP 5.0 improves signal peptide predictions using deep neural networks. *Nat Biotechnol.* 2019;37: 420–423. doi:10.1038/s41587-019-0036-z
30. Krogh A, Larsson B, von Heijne G, Sonnhammer ELL. Predicting transmembrane protein topology with a hidden markov model: application to complete genomes. *J Mol Biol.* 2001;305: 567–580. doi:10.1006/JMBI.2000.4315
31. Lagesen K, Hallin P, Rødland EA, Stærfeldt H-H, Rognes T, Ussery DW. RNAmmer: consistent and rapid annotation of ribosomal RNA genes. *Nucleic Acids Res.* 2007;35: 3100–3108. doi:10.1093/nar/gkm160
32. Almagro Armenteros JJ, Sønderby CK, Sønderby SK, Nielsen H, Winther O. DeepLoc: prediction of protein subcellular localization using deep learning. Hancock J, editor. *Bioinformatics.* 2017;33: 3387–3395. doi:10.1093/bioinformatics/btx431
33. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* 2010;20: 1297–303. doi:10.1101/gr.107524.110
34. Cingolani P, Platts A, Wang LL, Coon M, Nguyen T, Wang L, et al. A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff. *Fly (Austin).* 2012;6: 80–92. doi:10.4161/fly.19695
35. Shevchenko A, Tomas H, Havli J, Olsen J V, Mann M. In-gel digestion for mass spectrometric characterization of proteins and proteomes. *Nat Protoc.* 2006;1: 2856–2860. doi:10.1038/nprot.2006.468
36. Rappaport J, Mann M, Ishihama Y. Protocol for micro-purification, enrichment, pre-fractionation and storage of peptides for proteomics using StageTips. *Nat Protoc.* 2007;2: 1896–1906. doi:10.1038/nprot.2007.261
37. Cox J, Hein MY, Luber CA, Paron I, Nagaraj N, Mann M. Accurate Proteome-wide Label-free Quantification by Delayed Normalization and Maximal Peptide Ratio Extraction, Termed MaxLFQ. *Mol Cell Proteomics.* 2014;13: 2513–2526. doi:10.1074/mcp.m113.031591
38. Perez-Riverol Y, Csordas A, Bai J, Bernal-Llinares M, Hewapathirana S, Kundu DJ, et al. The PRIDE database and related tools and resources in 2019: improving support for quantification data. *Nucleic Acids Res.* 2019;47: D442–D450. doi:10.1093/nar/gky1106
39. Wehrens R, Kruisselbrink J. Flexible Self-Organizing Maps in kohonen 3.0. *J Stat Softw Artic.* 2018;87: 1–18. doi:10.18637/jss.v087.i07
40. Hahsler M, Hornik K. **TSP** - Infrastructure for the Traveling Salesperson Problem. *J Stat Softw.* 2007;23: 1–21. doi:10.18637/jss.v023.i02

## Genome-free proteotranscriptomics assisted gene annotation

567

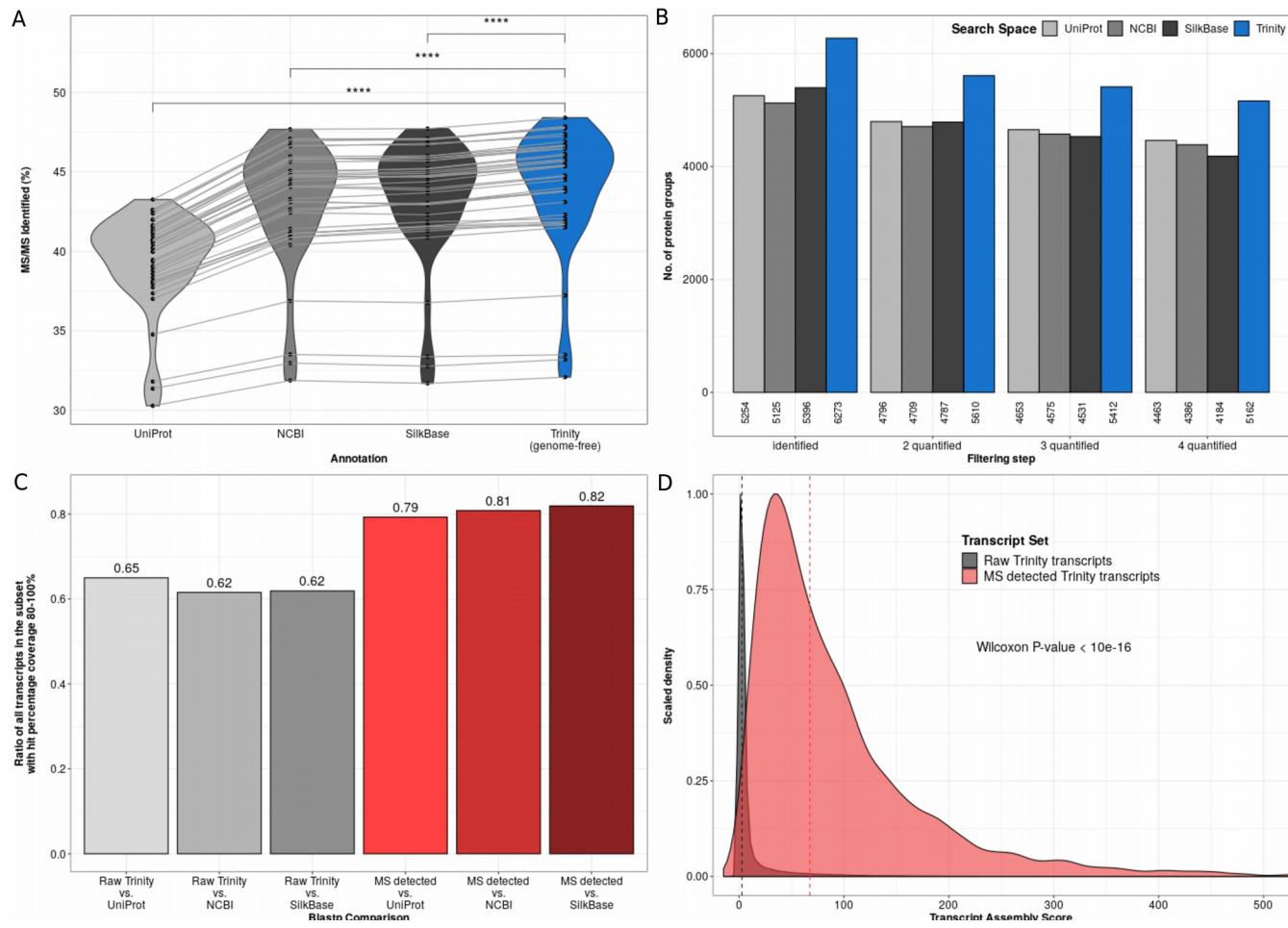
29

## Genome-free proteotranscriptomics assisted gene annotation

568 **FIGURES**

569 **Figure 1. Genome-free transcriptome assembly approach and assessment of**  
570 **annotation quality. A.** Overview of the proteotranscriptomics annotation approach. **B.** Pie-  
571 **chart of BUSCO analysis based on the BUSCO arthropoda gene set. C.** Barplot summarizing  
572 **the results of a full-length transcript comparison between the genome-free Trinity assembly to**  
573 **currently available annotations from UniProt, NCBI and SilkBase.**

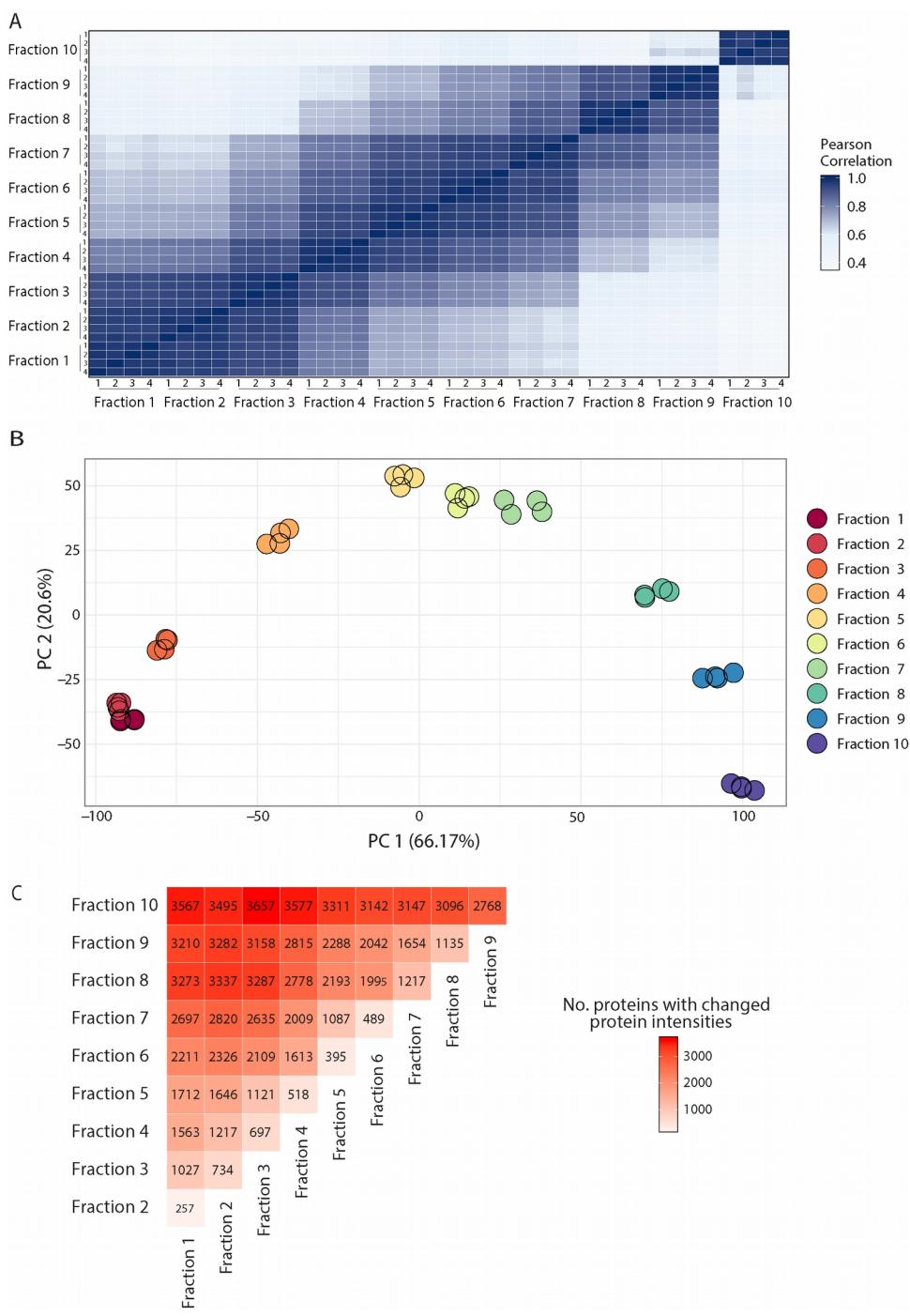
## Genome-free proteotranscriptomics assisted gene annotation



575 **Figure 2. High resolution mass spectrometry provides evidence for superior genome-**  
 576 **free annotation.** **A.** Violin plots show distribution of identified MS/MS spectra (in percent) for  
 577 each database used. With identical raw proteomic data the genome-free Trinity annotation  
 578 shows significantly higher identified tandem MS spectra percentages than the three currently  
 579 available annotations from UniProt, NCBI and SilkBase. Grey lines connect percentages  
 580 stemming from the same MS run. \*\*\*\* indicates two-sided paired Wilcoxon signed rank test p-  
 581 values below 0.0001. **B.** Barplot showing number of protein groups identified after different  
 582 filtering steps with UniProt, NCBI, SilkBase and genome-free Trinity annotation. The Trinity  
 583 annotation shows higher numbers of identified protein groups for identification and  
 584 quantification (including replicates). **C.** Barplot of the ratio of transcripts with a hit percentage  
 585 coverage of more than 80% when compared to current *Bombyx mori* annotations. Grey bars  
 586 include all Trinity annotated transcripts and red bars represent transcripts that have peptide  
 587 evidences detected by MS. **D.** Scaled density plot showing distribution of transcript assembly  
 588 scores of all Trinity annotated transcripts (gray) and transcripts with peptide evidences  
 589 detected by MS (red). Dashed vertical lines indicate the median assembly score of each  
 590 subset.

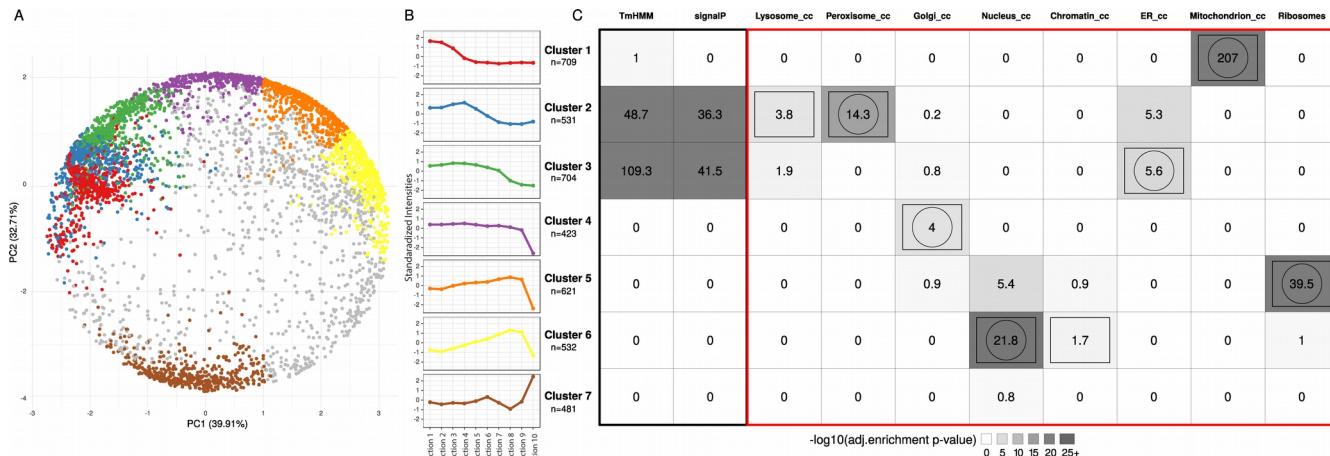
## Genome-free proteotranscriptomics assisted gene annotation

591

592 **Figure 3: Spatial proteomics unveils subcellular localization of *Bombyx mori* proteins.**

593 **A.** Heatmap of all pairwise comparisons (pearsons correlations) shows high concordance  
 594 between replicates and relatedness of adjacent fractions. **B.** Principal component plot based  
 595 on the 1000 most dynamic protein groups demonstrates differences of fractions and similarity  
 596 of replicates (same color) cluster together and (except for fraction 1 and 2) farther  
 597 away from the other fractions. **C.** Summary heatmap of number of proteins with significantly  
 598 different protein intensities across fractions.

## Genome-free proteotranscriptomics assisted gene annotation



599 **Figure 4: Unsupervised clustering of fractionation mass spectrometry data.**

600 **A.** Principal component plot of components 1 and 2 calculated from standardized average  
 601 protein intensities. Standardized protein intensities across fractionation samples were  
 602 clustered using SOM (self-organizing maps). Proteins in PCA space are colored according to  
 603 their assignment to one of eight distinct clusters (see color code in b). 795 proteins could not  
 604 be associated with any of the clusters (gray colored dots). **B.** Line plots of the median  
 605 standardized protein intensities across fractionation steps of the different clusters assigned by  
 606 SOM (see a). **C.** Heatmap summarizing enrichment analyses of cellular localization  
 607 annotations in the respective clusters. Color darkness corresponds to levels of enrichment (-  
 608 log<sub>10</sub> of adjusted P-values). Rectangles and circles indicate highest enrichment for the  
 609 corresponding localization and gene clusters, respectively. General and more specific  
 610 localization categories are framed in black and red, respectively.

## Genome-free proteotranscriptomics assisted gene annotation

611 **Supplementary information**612 **Additional file 1:** Supplementary Figures S1-S9 and Tables S1-S9613 **Additional file 2:** Table S10. Functional annotation and characteristics of ORFs614 **Additional file 3:** Table S11. Label-free quantification data615 **Additional file 4:** Table S12. Novel protein evidence data616 **Additional file 5:** BOMOPTA annotation CDS sequences617 **Additional file 6:** BOMOPTA annotation protein sequences