

1 **Proteogenomic annotation of the Chinese hamster reveals extensive**
2 **novel translation events and endogenous retroviral elements**

3 Shangzhong Li^{1,2}, Seong Won Cha³, Kelly Hefner⁴, Deniz Baycin Hizal⁵, Michael Bowen⁵,
4 Raghobama Chaerkady⁵, Robert N. Cole⁶, Vijay Tejwani⁷, Prashant Kaushik⁸, Michael Henry⁸,
5 Paula Meleady⁸, Susan T. Sharfstein⁷, Michael J. Betenbaugh⁴, Vineet Bafna⁹, Nathan E.
6 Lewis^{1,2,10}

- 7 1. Department of Bioengineering, University of California, San Diego, La Jolla, CA, USA
8 2. Novo Nordisk Foundation Center for Biosustainability, University of California, San Diego,
9 La Jolla, CA, USA
10 3. Department of Electrical and Computer Engineering, University of California, San Diego,
11 La Jolla, CA, USA
12 4. Chemical and Biomolecular Engineering, Johns Hopkins University, Baltimore, Maryland,
13 USA
14 5. Antibody Discovery and Protein Engineering, MedImmune LLC, Gaithersburg, Maryland
15 USA
16 6. The Mass Spectrometry Core, Johns Hopkins School of Medicine, Baltimore, MD, 21205
17 7. Colleges of Nanoscale Science and Engineering, SUNY Polytechnic Institute, Albany,
18 NY, USA
19 8. National Institute for Cellular Biotechnology, Dublin City University, Dublin 9, Ireland
20 9. Department of Computer Science and Engineering, University of California, San Diego,
21 La Jolla, CA, USA
22 10. Department of Pediatrics, University of California, San Diego, La Jolla, CA USA
23

24

25

26

27

28

29

30

31

32

1 Abstract

2 A high quality genome annotation greatly facilitates successful cell line engineering. Standard
3 draft genome annotation pipelines are based largely on *de novo* gene prediction, homology, and
4 RNA-Seq data. However, draft annotations can suffer from incorrectly predictions of translated
5 sequence, incorrect splice isoforms and missing genes. Here we generated a draft annotation
6 for the newly assembled Chinese hamster genome and used RNA-Seq, proteomics, and Ribo-
7 Seq to experimentally annotate the genome. We identified 4,333 new proteins compared to the
8 hamster RefSeq protein annotation and 2,503 novel translational events (e.g., alternative
9 splices, mutations, novel splices). Finally, we used this pipeline to identify the source of
10 translated retroviruses contaminating recombinant products from Chinese hamster ovary (CHO)
11 cell lines, including 131 type-C retroviruses, thus enabling future efforts to eliminate retroviruses
12 by reducing the costs incurred with retroviral particle clearance. In summary, the improved
13 annotation provides a more accurate platform for guiding CHO cell line engineering, including
14 facilitating the interpretation of omics data, defining of cellular pathways, and engineering of
15 complex phenotypes.

16 **Keywords:** Chinese hamster, genome annotation, proteogenomics, endogenous retrovirus

17

18

19

20

21

1 Introduction

2 Chinese hamster ovary (CHO) cells are the primary workhorse for therapeutic protein
3 production(Golabgir et al., 2016). Sequencing and assembly of the CHO and Chinese hamster
4 genomes(Brinkrolf et al., 2013a; Lewis et al., 2013a; Xu et al., 2011) have enabled improvement
5 in protein production using genetic engineering and in cell line process optimization using omics
6 technologies(Kuo et al., 2017; Stolfa et al., 2018). A recent effort greatly improved the reference
7 Chinese hamster genome assembly by combining Pacific Biosciences Single Molecule Real
8 Time (SMRT) and short-read Illumina sequencing data, thus reducing the number of scaffolds
9 by 28-fold and filling 95% of the sequence gaps(Rupp et al., 2018). Despite these great
10 improvements in the assembly, the current genome annotation was based primarily on *ab initio*
11 prediction, protein homology, ESTs, and limited publicly available transcriptomic data. However,
12 these pipelines have difficulties in translation confirmation, splice form detection, and complete
13 novel gene identification(N. Castellana & Bafna, 2010). To improve cell line engineering
14 success, an accurate genome annotation is necessary.

15 Proteogenomics provides a way to address such challenges by integrating mass
16 spectrometry-based proteomics, RNA-Seq and genomic data. For example, peptides can be
17 identified by mapping tandem mass spectra to protein databases derived from RNA-Seq and
18 genome annotation. The peptides are then used to update the annotation with novel coding
19 regions and splice sites. Proteogenomics was first applied to *Mycoplasma pneumoniae*(Jaffe,
20 Berg, & Church, 2004) to identify new and extended open reading frames (ORFs) and remove
21 low quality gene models. It has also been applied to many eukaryotes including plants(N. E.
22 Castellana et al., 2008), yeast(Yagoub et al., 2015), and human(M.-S. Kim et al., 2014). In
23 addition to improving annotations, the proteomic data can also identify mutations (e.g., in
24 cancer(Woo et al., 2015) and post translational modifications(Cesnik, Shortreed, Sheynkman,
25 Frey, & Smith, 2016; Kaushik, Henry, Clynes, & Meleady, 2018)).

1 In addition to proteomics data, Ribo-Seq (data from sequencing ribosome-protected
2 coding reactions at single nucleotide resolution(Ingolia, Ghaemmaghami, Newman, &
3 Weissman, 2009)), provides a global view of actively translated mRNAs, and thus has been
4 utilized to predict ORFs and translation frames for proteins(Calviello & Ohler, 2017) and to
5 identify additional predicted proteins for proteogenomics(Crappé et al., 2015). Together
6 transcriptomic, Ribo-Seq, and proteomic data can be invaluable for refining annotation about
7 proteins.

8 To obtain a data-supported refinement of the Chinese hamster genome annotation, here
9 we integrated proteomics, RNA-Seq, and Ribo-Seq to verify coding regions, update gene
10 models, identify novel translated genes, and verify protein-coding variants in different CHO cell
11 lines (**Figure 1**). To further demonstrate the increased value of this resource, we investigated
12 the challenge associated with the Food and Drug Administration (FDA) requirement to ensure
13 that viral particles (particularly endogenous retroviral particles) are eliminated from the
14 therapeutic protein product, which contributes to the high costs in bioprocessing(Strauss et al.,
15 2009). Specifically, we identified all translated retrovirus particles in CHO cells, including
16 previously unannotated translated loci, thus providing potential knockout targets to increase
17 drug purity and reduce demands on viral clearance. This proteogenomic resource will be
18 invaluable for future efforts to study and engineer CHO cells for bioprocessing.

19 **Results**

20 **Draft annotation for the genome**

21 The CHO-K1(Xu et al., 2011) and Chinese hamster genome sequences(Brinkrolf et al.,
22 2013b; Lewis et al., 2013b) were originally assembled using short read (99bp) technology, and
23 therefore resulted in fragmented contigs and scaffolds. Thus, efforts to annotate the genomes
24 resulted in some errors in protein and gene models. The RefSeq pipeline has corrected some

1 such errors; however, the complete reassembly of the Chinese hamster genome(Rupp et al.,
2 2018) provides an opportunity to obtain a much improved annotation of coding regions in the
3 genes and their corresponding protein sequences. Therefore, we generated a new draft
4 annotation here (**Supplementary Figure S1**), including predicted protein sequences.

5 68 RNA-Seq samples were prepared from multiple CHO cell lines and hamster tissues
6 (**see Methods**). Transcripts were assembled for each sample separately using stringtie(Pertea
7 et al., 2015) and merged using stringtie-merge(Pertea et al., 2015), yielding 26,530 genes with
8 68,082 transcripts. Then 38,654 hamster RefSeq transcripts were mapped to the newly
9 assembled hamster reference genome using GMAP(Wu & Watanabe, 2005). Finally, the
10 RefSeq alignments and RNA-Seq assembled transcripts were merged to 86,790 transcripts and
11 grouped into 38,511 genes based on genomic locations using Program to Assemble Spliced
12 Alignments (PASA).

13 We then applied TransDecoder(Haas et al., 2013) to the 86,790 transcripts and
14 predicted 63,331 proteins with 47,829 unique protein sequences (the remaining are noncoding
15 transcripts). To annotate the function of the proteins, we aligned the protein sequences to the
16 hamster RefSeq and UniProt Swiss-Prot protein databases using BLASTP(Gish & States, 1993).
17 (**Supplementary Table S1**). We assigned UniProt gene names to the proteins in our draft
18 annotation except for those that only map to hamster RefSeq proteins. Furthermore, we
19 identified 4,640 non-coding transcripts by aligning the transcripts to the hamster RefSeq non-
20 coding transcripts using BLASTN.

21 **Proteogenomics helps identify novel proteins in the draft** 22 **annotation**

23 To quantify novel proteins predicted in the draft annotation compared to the hamster
24 Refseq proteins, we mapped 47,829 unique draft proteins to the RefSeq proteins using BLASTP.

1 We classified the mappings into 5 main categories: (1) 15,787 perfectly mapped proteins; (2)
2 7,483 proteins mapping perfectly on only one end between the draft and RefSeq sequences; (3)
3 11,780 high quality mapped proteins (over 90% percentage of identity (pident) and over 80%
4 percentage of length (plen) on both sides of homology protein mapping pair between draft and
5 RefSeq); (4) 6,688 high quality mapped proteins (over 90% pident and over 80% plen on either
6 side), but only mapping well on one end; (5) 5,820 low quality or non-mapping proteins. 289
7 proteins failed to map. We defined proteins that were not in category (1) as novel proteins.
8 Among the one-sided perfect mapping proteins, 3,336 proteins are shorter in the draft,
9 compared to RefSeq, while 4,147 proteins are longer than RefSeq. Interestingly, isoforms of
10 some of the former proteins map perfectly to RefSeq, which indicates the draft annotation
11 pipeline is sensitive to splice sites, which resulted in more isoforms being assembled than seen
12 in Refseq.

13 Next, we sought peptide support for the novel proteins in the draft annotation.
14 12,870,725 mass spectra were acquired and prepared from multiple CHO cell lines and hamster
15 tissues. We merged the draft and RefSeq proteins and extracted the unique protein sets as a
16 reference protein database. Then we used MS-GF+(S. Kim & Pevzner, 2014) to search the
17 peptide-spectrum matches (PSMs) with 1% FDR correction and identified 244,729 peptides.
18 Here we only consider proteins with at least two peptides and at least one of them mapping
19 uniquely. For each pair of homologous draft and RefSeq proteins, the draft protein was
20 considered as novel if it has extra peptide support compared to the corresponding RefSeq
21 protein. As a result, we identified 4,077 draft novel proteins, 3,890 of which have additional
22 unique peptide support not seen in the RefSeq sequence (**Figure 2A**). The remaining 187 novel
23 proteins have extra peptides mapping to multiple locations. (**Figure 2B**). The high quality protein
24 mappings category (>90% pident and >80% plen) has the most novel proteins. 5,608 draft
25 proteins have the same peptide support as similar RefSeq proteins, which may require
26 additional data to verify their novel features. The numbers of novel proteins in each mapping

1 category are depicted in **(Figure 2)**.

2 **Proteogenomics and Ribosome profiling identify** 3 **additional translational events**

4 In addition to verifying novel protein sequences, proteomics can also help identify other
5 translational events, e.g., novel splice sites, gene fusions, etc. **(Figure 3A)**. Thus, to obtain a
6 more comprehensive view of protein sequence verification and identification of other translation
7 events, we created 4 putative protein databases: (1) KnownDB-predicted from draft annotation,
8 (2) SnpDB-translated from RNA-Seq reads that have non-synonymous mutations and short
9 INDELs, (3) SpliceDB-translated from spliced RNA-Seq reads, and (4) SixframeDB-peptides
10 between stop codons in all 6 frames of the genome **(Figure 1)**. As previously
11 recommended(Woo et al., 2015), we performed multi-stage 1% **(Supplementary Figure S2)**
12 FDR correction for the databases sequentially. This pipeline **(Supplementary Figure S3)**
13 identified 3,656,801 (28%) significant PSMs resulting in 239,973 unique peptides mapping to
14 the KnownDB and 9,808 peptides mapping to the remaining databases **(Figure 3B)**. Among all
15 the peptides identified from KnownDB, 208,904 (87%) map to unique genomic locations.

16 We required each validated protein to have at least two peptides and at least one of
17 them map uniquely (i.e., to only one locus). Using this strategy, we verified 30,814 proteins,
18 which represent 64.4% of the sequences in the KnownDB.

19 After known protein sequence validation, we explored additional novel peptide events.
20 To guarantee high confidence of the events, we filtered out those with fewer than 3 RNA-Seq
21 supporting reads, and required each event to have at least one uniquely mapped peptide. Most
22 proteins are longer than 100 amino acids, and tryptic peptides sequenced are usually about 7-
23 32(Frank, 2009) amino acids long. Thus, if novel peptides are close and in the same
24 translational frame, they are likely to support the same protein. Therefore, we clustered novel
25 identified peptides that were in the same translational frame and fewer than 100 amino acids

1 away from each other to represent the same event. In total, we discovered 2,503 new
2 translational events, 86% of which are novel splice and nonsynonymous single nucleotide
3 polymorphisms (SNPs) (**Figure 3A**). Novel splice sites represent 47.3% of the total events,
4 covering 857 genes. We also identified 70 alternative splice events, 22 reverse strand
5 translation events, 9 novel ORFs (not predicted by the draft annotation) and 5 gene fusion
6 events.

7 Ribo-Seq offers orthogonal evidence to further support protein sequences, and can help
8 discover new proteins as well. Ribo-Seq and corresponding RNA-Seq data sets for an IgG-
9 producing CHO cell line were acquired at both exponential and stationary phase (Kallehauge et
10 al., 2017). We used RiboTaper (Calviello et al., 2016) to predict the translating ORFs under the
11 guidance of the draft annotation. 28,700 transcripts were predicted to encode proteins, with
12 24,709 (86%) having a single ORF (**Figure 3C**). Among these, 13,666 transcripts have the
13 same protein sequences as the draft annotation. In addition, 1,318 “non-coding transcripts” in
14 the draft annotation are predicted to encode proteins. The remaining Ribo-Seq predicted
15 sequences were classified into two groups: (1) in the same frame (8775) and (2) in different
16 frames (950) with the draft annotation (**Figure 3D, Supplementary Table S2**).

17 Protein predicted by Ribo-Seq can expand the databases for proteomics so that more
18 proteins can be verified. Here we used Ribo-Seq predicted proteins as a novel database and
19 ran the proteogenomics pipeline together with KnownDB. After filtering all peptides identified in
20 the previous proteogenomics pipeline, the Ribo-Seq data facilitated the identification of 2,581
21 new peptides. Here we require at least one unique peptide in an ORF to verify translation.
22 Among 1,628 non-protein-coding transcripts in the draft annotation, 218 were verified to encode
23 proteins by Ribo-Seq and proteomics. While Ribo-Seq enabled the successful identification of
24 many new peptides, including those in transcripts previously thought to be non-coding, it was
25 less successful at identifying the translation start sites. Supporting peptides were found for only
26 8 out of the 305 ORFs that were predicted to be longer than the draft annotation. In addition,

1 Ribo-Seq helped identify translation events in 5'UTR and 3'UTR regions in 234 genes,
2 consistent with previous reports in human (Ingolia et al., 2014). For more details, see
3 **Supplementary Table S2.**

4 **Proteomic based validation of SNPs and INDELs in CHO** 5 **cell lines and hamster tissues**

6 The peptides obtained from proteomic studies enabled the validation of genetic variants
7 in the various CHO cell lines and hamster tissues. To discover these mutations, we used our
8 SnpDB (**Figure 1**) as the novel database in the proteogenomics pipeline, which includes
9 peptides translated from all the RNA-Seq reads supporting SNPs and small INDELs.
10 Proteomics identified fewer mutations than RNA-Seq (**Supplementary Table S3**), mainly
11 because of its lower depth of coverage compared to RNA-Seq. Furthermore, mutated proteins
12 can be degraded, and therefore not detected. In total we identified 973 nonsynonymous SNPs,
13 located in 734 genes. Most genes have one SNP while there are 6 genes with more than 5
14 SNPs: GAPDH, GOLGB1, AHNAK2, PKM, MYH9 and EEF1A1. Surprisingly, only one protein
15 lost a stop codon (ribosome protein gene RPS23) while others change amino acids. 75% of the
16 SNPs are homozygous, which indicates CHO cell lines may have developed and retained those
17 mutations after long periods of evolution. Furthermore, the distributions of the 6 SNP types
18 showed that transitions occurred more frequently than transversions, and the proteomic data
19 captured a similar distribution of mutations as RNA-Seq data (**Figure 4A, Figure 4B**). We
20 identified 42 insertions and 6 deletions, located in 41 and 6 genes respectively. There are more
21 homogeneous frameshift INDELs than other types. The 8 genes harboring both SNPs and
22 insertions include AHNAK2, CALR, HNRNPUL1, HSP90B1, PLEKHG5, PTMA, RIF1, and VAT1,
23 while only SIK3 had both SNPs and deletions.

24 Next, we looked at the mutation distribution across the protein body. Since there are far
25 fewer INDELs than SNPs, focused on the distribution of SNPs. **Figure 4C** shows that SNPs are

1 distributed relatively evenly across the protein body. Interestingly, we also identified peptide-
2 supported mutations in 5' UTR regions, but the number is much smaller than for those in coding
3 regions.

4 Many different CHO cell lines (e.g., CHO-K1, CHO-S and DG44), have been used to
5 develop different recombinant protein-producing cells. Each cell line has a lengthy history of
6 mutation and selection during cell line development(Lewis et al., 2013a), and therefore can have
7 unique genomic variants(Lewis et al., 2013a; van Wijk et al., 2017). To check the variations
8 between these CHO cell lines, we compared their peptide-supported SNPs in the coding
9 regions. Most of the SNPs are shared among the cell lines, which means these variants have
10 been conserved during the long period of cell line development, either due to early mutations
11 obtained in CHO cells when derived in 1957 or genetic drift of the Chinese hamster colony since
12 then (**Figure 4D**).

13 **Proteomics elucidate translated retroviral elements in the** 14 **genome**

15 For decades, it has been known that CHO cells shed retroviral particles(Lieber,
16 Benveniste, Livingston, & Todaro, 1973); while these were shown to be non infectious(Anderson,
17 Low, Lie, Keller, & Dinowitz, 1991; Dinowitz et al., 1992), the safety concern has required
18 companies to filter out all such viral particles and conduct extensive testing to verify non-
19 infectivity. This adds a substantial cost to production. Viral particles have been isolated, but all
20 mRNAs that have been sequenced from these particles were all non-coding, in that they
21 contained many early stop codons. Thus, it remains unclear which loci encode the translated
22 viral particles. Here, we analyzed the transcriptomic and proteomic data to identify the loci of
23 expressed and translated retroviral particles, to enable further efforts to eliminate these particles
24 and reduce drug purification costs.

25 To identify translated endogenous retroviruses in CHO cells, we first extracted peptides

1 that support retroviral proteins from our draft annotation (**Figure 5A**). Since retroviruses can
2 have multiple similar copies across the whole genome, we consider peptides that can map to
3 multiple locations. We found 457 retroviral genes covered by 723 transcripts, 184 of these
4 genes (corresponding to 304 transcripts) have peptide support and 111 transcripts map well to
5 RefSeq (>60% pident, >80% plen) and UniProt (>50% pident, >80% plen). Infectious retroviral
6 DNA is flanked by two identical non-coding repeats called long terminal repeats (LTR)(Temin,
7 1982), which aid in retroviral mobility and integration into the host genome, along with regulation
8 of retroviral gene expression. In the hamster genome, we identified 3,324 LTR pairs in the
9 reference genome using LTRharvest(Ellinghaus, Kurtz, & Willhoeft, 2008) and found only 76
10 retroviral transcripts locate between those LTRs, 74 of which have peptide support. If one LTR
11 is disrupted, the other side can still effectively induce transcription. Thus, genes not flanked by
12 LTR pairs can still produce retroviral particles.

13 We next aimed to identify unannotated translated retroviral sites in the genome (**Figure**
14 **5A**). For this, we aligned all retroviral proteins from NCBI to the reference genome using
15 tblastn(Altschul, Gish, Miller, Myers, & Lipman, 1990). Then we overlapped the mapping sites
16 with the 5,104 novel peptides identified against novel protein databases from our
17 proteogenomics pipeline and obtained 41 novel retroviral sites, 1 of which localized between an
18 LTR pair and was covered by 5 peptides. The site showed homology to the gag proteins from
19 Gibbon ape leukemia virus and Spleen focus-forming virus. Finally, since CHO cells were
20 originally derived in 1957, we further checked if new infections may have emerged in CHO
21 (**Supplementary Figure S4**). For this analysis, we aligned all known retroviral proteins and
22 novel peptides from proteogenomics pipeline, using tblastn, to in-house Illumina-corrected
23 single molecule real time (SMRT) sequence data from CHO-S and the Chinese hamster(Rupp
24 et al., 2018). After filtering out the putative viral sites identified in hamster, we found no evidence
25 that wild type CHO-S cell lines have acquired any new retroviral sites.

26 Mammalian retroviruses have been classified into different types based on the genomic

1 compositions. We mapped retroviral protein sequences identified from previous steps to UniProt
2 and used protein full names to discover 3 main retroviral types in hamster: type-A, type-B, and
3 type-C(Weiss, 1996). We found type-C retroviruses to be the most highly transcribed and
4 translated (**Supplementary Figure S5**). In addition, type-C viral particles have been identified in
5 CHO cell lines before and regulatory agencies now require the verification that products are
6 non-infectious type-C particles(Dinowitz et al., 1992; Lie et al., 1994). The vast majority of type-
7 C proteins had little or no peptide support, suggesting these are silenced or noncoding. Of the
8 131 type-C retroviral proteins with peptide support that we identified, most were gag and
9 envelope proteins (**Figure 5B**). Only 7 proteins had more than 20 peptides, and most proteins
10 had low coverage of supporting peptides (**Figure 5C**). **Figure 5D** shows an example of highly
11 translated envelope protein with 42 peptides, covering 32% of the coding sequence. Although it
12 has many secondary reads, it also has many uniquely mapping reads, which indicates this locus
13 is truly expressed. The proteins with high coverage and more supported peptides should be
14 prioritized in efforts to eliminate viral particle production. The genomic locations, RNA-Seq and
15 peptide coverage of all endogenous retroviral genes are provided in **Supplementary Table S4**.

16 Discussion and Conclusion

17 Here we presented the first proteogenomic reannotation of the Chinese hamster genome,
18 in which we utilized RNA-Seq, Ribo-Seq, and proteomics to improve the annotation. To identify
19 as many peptide-supported proteins as possible, we mapped spectra to a known protein
20 database from a draft annotation and several novel protein databases derived from different
21 data types. We identified 4,077 novel proteins in the draft annotation compared to the hamster
22 RefSeq protein database and 2,503 novel translational events and mutations in hamster and
23 CHO cell lines. Furthermore, we identified the potential sources of retroviral particles shed from
24 CHO cells, including 131 type-C retrovirus genes, 7 of which are supported by more than twenty
25 peptides.

1 Usually an annotation is required before running proteogenomics pipeline. The typical
2 first step of genome annotation is masking the repeats to avoid getting millions of seeds during
3 BLAST(Yandell & Ence, 2012). However masking the genome can hide important annotation
4 information, including common domains and retroviral elements. To avoid this loss of
5 information, we aligned assembled transcripts to the unmasked genome using gmap. Doing so
6 resulted in only 0.7% transcript mappings to be ambiguous (i.e., with >2 mappings). This
7 enabled the annotation and analysis of endogenous retroviral genes, which usually have
8 multiple similar copies. Masking repeats often removes such information(Slotkin & Keith Slotkin,
9 2018). Thus, future annotation efforts would benefit from the acquisition and alignment of high
10 quality transcripts or protein sequences to the genome of interest.

11 Different pipelines have been designed for genome annotation in higher
12 eukaryotes(Yandell & Ence, 2012), and most include *ab initio* prediction, mapping of
13 homologous protein sequences, and transcript assembly and alignment. As the throughput and
14 resolution of mass spectrometry techniques is increasing, more studies are integrating these
15 data type to refine annotations(Nesvizhskii, 2014). Here we discovered thousands of genes and
16 novel translational events using proteomics data. Despite their value, proteomics data can be
17 sparse since the chemical properties of some peptides are less compatible with the
18 experimental setup (e.g., due to hydrophobicity) or size of the peptides post
19 digestion(Chandramouli & Qian, 2009). Furthermore, many peptides have diverse post-
20 translational modifications, thus making it difficult to align peptides to predicted peptide
21 sequences from genomes. Thus, Ribo-Seq provides a complementary method to further
22 discover new genes or correct annotations(Ingolia, Brar, Rouskin, McGeachy, & Weissman,
23 2013). The Ribo-Seq datasets we used here are from the DG44 CHO cell line, and CHO cells
24 may only express half to two thirds of their genes(Rupp et al., 2014; Singh, Kildegaard, &
25 Andersen, 2018). Thus, further annotation efforts would benefit from the acquisition of Ribo-Seq
26 from many other hamster tissues and developmental stages. In summary, as more tools and

1 proteomic and Ribo-Seq data accumulates, these data types will become increasingly
2 integrated into standard pipelines for genome annotation.

3 Finally, the reannotation here provides an invaluable resource for the development of
4 improved CHO cell lines. While CHO cells provide several advantages as an expression host for
5 recombinant protein production, HCPs are continuously secreted(Hogwood, Bracewell, &
6 Smales, 2014; Kumar et al., 2015), thus impacting recombinant protein quality and safety. Thus,
7 expensive chromatographic columns, filtration systems, and infection and HCP assays are
8 required during downstream processing. This adds considerable cost to biopharmaceuticals.
9 One particular regulatory concern has been the endogenous retroviruses that are shed by CHO
10 cells(Lie et al., 1994). For these, assays were developed to quantify retroviral particles and
11 ensure infectious retroviral particles are not found in the drug product after extensive filtration(de
12 Wit, Fautz, & Xu, 2000). Here, we identified, among hundreds of retroviral genes in the hamster
13 genome, which ones are expressed and translated in several CHO cell lines. This information
14 will enable future efforts to remove these from CHO cells, thereby reducing burdens to
15 downstream processing. To further facilitate such efforts, many endogenous retroviral genes
16 show high levels of homology. In our work, this was manifested in the identification of many
17 retroviral RNA-Seq reads and tryptic peptides that map to multiple genomic loci. As many of
18 these also share DNA sequence, multiple viruses can be knocked out simultaneously by
19 targeting the conserved regions, as accomplished in the pig(Yang et al., 2015). Doing this in
20 CHO cells can reduce the costs by simplifying expensive purification steps where product can
21 also be lost, and also simplify viral testing steps for the final product.

22 In conclusion, our work provides a refined and more extensive annotation of the Chinese
23 hamster genome, which will enables more accurate CHO cell line engineering(Fischer, Handrick,
24 & Otte, 2015; Kuo et al., 2017; Lee, Grav, Lewis, & Faustrup Kildegaard, 2015; Richelle & Lewis,
25 2017)The improved annotation will also facilitate improved processing of omics data(Chen, Le,
26 & Goudar, 2017) as the gene models improve. Finally, a more complete list of all genes will

1 enable efforts to map out molecular pathways in CHO cells to enable systems approaches to
2 cell line development(Kuo et al., 2017).

3 **Methods**

4 **Proteomic sample preparation**

5 Proteomic data were acquired at two different locations using different protocols, thus
6 increasing the diversity in spectra used for annotation. These are referred to as batch 1 and
7 batch 2, as follow.

8 **Tissue Sample Collection for batch 1**

9 Chinese hamsters were generously provided by Dr. George Yerganian (Cytogen
10 Research, Roxbury, MA). Euthanization was performed by CO₂ and verified by puncture.
11 Harvested liver and ovary tissues were flash frozen on dry ice and stored at -80°C until analysis.

12 **Cell Culture Sample Collection for batch 1**

13 Suspension CHO cell lines (including CHO-S and CHO DG44) were grown in shake-
14 flask batch culture. CHO-S cells were cultured in CD-CHO medium supplemented with 8mM
15 glutamine (Thermo Fisher Scientific, Waltham, MA), and CHO DG44 cells were culture in DG44
16 medium supplemented with 2mM glutamine (Thermo Fisher Scientific, Waltham, MA). Samples
17 were collected on day 2 for exponential phase and day 4/5 for stationary phase. Cells were
18 incubated at 37°C, 8% CO₂, and 120RPM. For sample collection, approximately 3 million cells
19 were spun down, washed with PBS on ice, frozen rapidly on dry ice, and stored at -80°C until
20 analysis.

21 **Cell lysate and Tissue Sample preparation for batch 1**

22 Cell culture lysates and tissue samples were thawed on ice and suspended in 2%

1 sodium dodecyl sulfate (SDS) supplemented with 0.1mM phenylmethane sulfonyl fluoride
2 (PMSF) and 1mM ethylenediaminetetraacetic acid (EDTA), pH 7-8. Samples were lysed by
3 sonicating for 60 seconds at 20% amplitude followed by 90 seconds pause (for three cycles).
4 Protein concentration was measured with bicinchoninic acid (BCA) protein assay after briefly
5 spinning to remove cell debris. Three hundred micrograms of each sample was reduced in
6 10mM tris(2-carboxyethyl)phosphine (TCEP), pH 7-8, at 60°C for 1hr on a shaking platform.
7 After bringing each sample to room temperature, iodacetamide was added to alkylate the
8 sample to 17mM final concentration for 30 minutes. Next, samples were cleaned using 10kDa
9 filters to reduce the SDS concentration as suggested by the filter aided sample preparation
10 (FASP) protocol(Wiśniewski, Zougman, Nagaraj, & Mann, 2009). The samples were finally
11 digested using trypsin/LysC enzyme mix at an enzyme to substrate ratio of 1:10 (Promega
12 V507A, Madison, WI), overnight at 37°C on a shaking platform.

13 **Identification of Proteins by Mass Spectrometry for batch**

14 **1**

15 Digested peptides (100µg from each protein digest) were fractionated on a basic
16 reversed phase column (XBridge C18 Guard Column, Waters, Milford, MA). Fractions were
17 concatenated into 48 prior to second dimension LC and MS analysis. The use of fractionation
18 with equal peptides in each was designed to mimic biological replicates for each sample.
19 Tandem MS/MS analysis of the peptides was carried out on the LTQ Orbitrap Velos (Thermo
20 Fisher Scientific, Waltham, MA) MS interfaced to the Eksigent nanoflow liquid chromatography
21 system (Eksigent, Dublin, CA) with the Agilent 1100 auto sampler (Agilent Technologies, Santa
22 Clara, CA). Peptides were enriched on a 2cm trap column (YMC, Kyoto, Japan), fractionated on
23 Magic C18 AQ, 5µm, 100Å, 75µm x 15cm column (Bruker, Billerica, MA), and electrosprayed
24 through a 15µm emitter (SIS, Ringoes, NY). Reversed phase solvent gradient consisted of
25 solvent A (0.1% formic acid) with increasing levels of solvent B (0.1% formic acid, 90%

1 acetonitrile) over a period of 90 minutes. LTQ Orbitrap Velos parameters included 2.0kV spray
2 voltage, full MS survey scan range of 350-1800m/z, data dependent HCD MS/MS analysis of
3 top 10 precursors with minimum signal of 2000, isolation width of 1.9, 30s dynamic exclusion
4 limit and normalized collision energy of 35. Precursor and fragment ions were analyzed at
5 60000 and 7500 resolutions, respectively.

6 **In-solution digestion of whole cell lysate for proteomics**

7 **batch 2**

8 A second batch of samples were prepared and analyzed using a different approach. For
9 these, 1mg of protein sample was transferred to a centrifuge tube, and all samples were
10 equalized to the same volume using the same lysis buffer. A fresh stock of 0.5M reducing agent
11 dithiothreitol (DTT) was prepared, and an appropriate volume of DTT was added to achieve a
12 final concentration of 5mM. Samples were incubated for 25 minutes at 56°C. Before alkylation
13 samples were cooled to room temperature and an appropriate volume of freshly prepared 0.5M
14 iodoacetamide was added to a final concentration of 14mM and incubated for 30 minutes at
15 room temperature. Untreated iodoacetamide was quenched by a second addition of 0.5M DTT
16 to make total concentration of DTT equal to 10mM and incubated for 15 min at room
17 temperature. The protein mixture was diluted 1:5 in 25 mM Tris-HCl, pH 8.2, to reduce the
18 concentration of urea to 1.6 M. A double digestion by trypsin was performed by adding trypsin to
19 1/50 enzyme: substrate ratio and incubated at 37°C. After 4 hours of primary incubation the
20 trypsin was topped up (enzyme: substrate ratio 1/100), and the protein mixture was left to digest
21 overnight at 37°C. After the overnight digestion, unused trypsin was quenched by adding TFA to
22 a final concentration of 0.4%.

23 **Peptide sample clean-up for proteomics batch 2**

24 Digested peptides were desalted and cleaned up using Sep-Pak c18 Vac cartridge,

1 200mg sorbent per cartridge, 55-500 μm Particle size (WAT054945) using negative pressure.
2 The C18 cartridge was washed and conditioned by using 9ml of ACN followed by 3ml of
3 50%ACN and 0.5% acetic acid. C18 resin was then equilibrated with 9ml of 0.1% TFA and
4 samples were loaded in 0.4%TFA. Loaded samples were desalted with 9ml 0.1%TFA. TFA was
5 removed with 1ml 0.5% acetic acid. Desalted peptides were eluted with 3ml of 50%ACN 0.5%
6 acetic acid. The eluted fraction was applied twice and collected in a 15ml conical tube. The
7 eluate was snap frozen in liquid nitrogen, and lyophilized overnight or until the white (sometimes
8 yellow) fluffy powder was observed. Dried peptides were stored at -20°C or otherwise dissolved
9 in the appropriate buffer for phosphopeptide enrichment.

10 **Draft genome annotation generation**

11 68 RNA-Seq samples from multiple CHO cell lines and hamster tissues were trimmed
12 and aligned to the newly assembled hamster reference genome using Trimmomatic 0.36(Bolger,
13 Lohse, & Usadel, 2014) and STAR 2.5.2(Dobin et al., 2013), respectively. The aligned reads
14 were assembled into transcripts using stringtie(Pertea et al., 2015) for each sample and then
15 the transcripts were consolidated into a union transcript set using stringtie-merge(Pertea et al.,
16 2015). To improve the transcript coverage, we also mapped the hamster RefSeq transcript
17 sequences to the newly assembled hamster genome using gmap(Wu & Watanabe, 2005) which
18 were then integrated with transcripts generated from stringtie-merge using the PASA
19 pipeline(Haas et al., 2003). Potential proteins encoded in the transcripts were predicted using
20 transdecoder(Haas et al., 2013). Finally the functions of predicted proteins were determined by
21 mapping them to the hamster RefSeq proteins and UniProt Swiss-Prot proteins(The UniProt
22 Consortium, 2018) using BLASTP(Gish & States, 1993). Proteins whose mapping lengths were
23 greater than 80% were considered and percentage identity of 60% when mapping to hamster
24 RefSeq and of 50% when mapping to UniProt were used as threshold to further classify proteins
25 into 4 categories with 1 to 4 indicating decreasing confidence scores as follows: 1. Pident and

1 plen are larger than the threshold and have the same gene name between hamster RefSeq and
2 UniProt. 2. Pident and plen are larger than the threshold and have different names between
3 hamster RefSeq and UniProt. 3. Pident and plen are less than the threshold and have the same
4 gene name between hamster RefSeq and UniProt. 4. Pident and plen are less than the
5 threshold and have different gene names between hamster RefSeq and UniProt. LncRNAs were
6 predicted by aligned transcripts to hamster lncRNAs using BLASTN with pident larger than 60%
7 and plen larger than 80%.

8 **Proteogenomics database construction**

9 We prepared 4 protein databases for mass spectrum matching: known protein database
10 (KnownDB), SNP database (SnpDB), splice database (SpliceDB) and six-frame translation of
11 the genome database (SixframeDB). The KnownDB includes protein sequences extracted from
12 draft annotation, the rest serve as novel protein databases. SnpDB was constructed by
13 translating RNA-Seq reads that have mutations(Woo, Cha, Na, et al., 2014). Mutations were
14 called using GATK3.7(Van der Auwera et al., 2013) and annotated using Annovar(Wang, Li, &
15 Hakonarson, 2010). The SpliceDB was constructed by translating all RNA-Seq reads that span
16 splice junctions(Woo, Cha, Merrihew, et al., 2014). The SixframeDB was derived from peptides
17 fragments between stop codons in all frames of the reference genome assembly.

18 Peptide identification The original MS/MS spectra were converted from RAW format to
19 Mascot Generic Format (MGF) using msconvert(Kessner, Chambers, Burke, Agus, & Mallick,
20 2008) and searched against each database independently using MSGFplus(S. Kim & Pevzner,
21 2014). Since different databases have different false discovery rates, it is recommended to
22 perform multistage FDR correction with 1% cut off for the databases(Woo et al., 2015), which
23 means the spectra failed to pass FDR correction were fed to the next database to correct again.
24 We corrected FDR for databases in the following order: KnownDB, SnpDB, SpliceDb,
25 SixframeDB.

1 **New translational event prediction**

2 Significant peptide-spectrum matches (PSM) against the novel databases were used to
3 discover new translation events using Enosi pipeline(Woo, Cha, Merrihew, et al., 2014). Briefly,
4 identified novel peptides were mapped to the novel databases to get loci relative to their
5 mapped proteins. Protein headers in the novel databases have loci relative to the reference
6 genome assembly. A custom python script was used to deduce peptide loci relative to the
7 reference genome assembly from those two loci. Then the loci were compared with the draft
8 annotation to decide event type. Peptides in the same translation frame that are less than 300
9 nucleotides (100 amino acids) away are grouped to represent the same event. For SNPs and
10 short INDELS, we filtered the false positives by variant calling using Illumina reads from
11 sequencing hamster genomic DNA against the reference genome assembly.

12 **Protein prediction using Ribo-Seq data**

13 We used previously published Ribo-Seq data of the CHO CS CS13-1.0 cell
14 line(Kallehauge et al., 2017). All Ribo-Seq and RNA-seq data of each biological sample were
15 trimmed and aligned to the reference genome assembly using Trimmomatic 0.36(Bolger et al.,
16 2014) and HISAT 2.2.1(D. Kim, Langmead, & Salzberg, 2015) respectively. The aligned bam
17 files were sorted and merged into one Ribo-Seq and one RNA-Seq bam files using Samtools
18 1.6(Li et al., 2009). The potential translated regions were predicted using RiboTaper(Calviello et
19 al., 2016), which takes advantage of coverage in coding regions and triplet periodicity of
20 ribosomal footprints. A custom python script was used to compare the translation prediction and
21 draft annotation.

22 We treated the predicted proteins from Ribo-Seq as a novel protein database and
23 combined it with the KnownDB. Then we mapped all the mass spectra to these two databases
24 using the proteogenomics pipeline. In this case, the identified novel peptides would be the
25 unique peptides from Ribo-Seq database.

1 **Retrovirus in the draft annotation**

2 All viral proteins in draft annotation were identified by mapping to UniProt and hamster
3 Refseq proteins using BLASTP and then the retroviral proteins were identify based on the full
4 gene names. Peptides supporting annotated retroviral proteins were identified by mapping the
5 known peptides to the KnownDB. LTRs were predicted using LTRharvest(Ellinghaus et al.,
6 2008). Retroviral proteins that locate between LTR were identified by overlapping LTR regions
7 with retroviral annotations using Bedtools2.27(Quinlan & Hall, 2010).

8 **New retrovirus discovery**

9 Since virus proteins lack introns, we filtered novel peptides by removing those with splice
10 sites. Retroviral proteins and filtered novel peptides were aligned to the reference genome
11 assembly using tblastn(Altschul et al., 1990) and mappings with more than 60 pident and 55
12 plen were considered for downstream analysis. Virus mapping and peptide mapping were
13 overlapped using Bedtools(Quinlan & Hall, 2010) to get the virus sites with peptide support,
14 which were then further overlapped with LTR regions to get virus sites with both LTR and
15 peptides support.

16 We decided the thresholds for virus tblastn mapping as follows (**Supplementary Figure**
17 **S6**). First, we started with low thresholds (30% for each). Then we assessed the overlap
18 between filtered mapping and the draft annotation to identify those associated with known virus
19 genes. If the mappings overlap with many non-viral genes, thresholds were incrementally
20 increased and overlapped with draft annotation again. This was repeated until the mappings
21 overlap with few non-viral genes and the number did not decrease anymore.

22 **Discovery of unique retroviruses in the CHO-S cell line**

23 The in-house CHO-S SMRT sequence data was used to check if CHO-S has unique
24 retroviral elements, compared to hamster. Firstly, we subsampled hamster SMRT reads to the

1 same depth as CHO-S SMRT data, and both datasets were corrected using Illumina paired-end
2 reads(Lewis et al., 2013a) through LoRDEC(Salmela & Rivals, 2014). Secondly, retroviral
3 proteins and filtered novel peptides from our proteogenomics pipeline were mapped to corrected
4 CHO-S and hamster SMRT reads using the same threshold as the previous tblastn mapping.
5 Thirdly, virus and peptide mappings were overlapped to get virus sites with peptide support for
6 CHO-S and hamster separately. Fourthly, we filtered CHO-S virus sites with hamster SMRT
7 virus sites and mapped the unique CHO-S sites to the reference genome assembly. The
8 mapped sites are the new retroviral sites and the unmapped sites are unique retroviral elements
9 in CHO-S.

10 **Type-C retrovirus detection in CHO cell lines**

11 The functions of all the identified retroviral proteins were determined by mapping the
12 protein sequences to UniProt using BLASTP. The full function descriptions of the proteins have
13 the organism resource. Therefore, the types of all the retroviruses were determined by manually
14 matching their full virus names to the types defined previously (see appendix “Retroviral
15 Taxonomy, Protein Structures, Sequences, and Genetic Maps” of this book(Coffin, Hughes, &
16 Varmus, 1997)). Peptide coverage of a protein equals to the number of amino acids covered by
17 peptides divided by the protein length. The RNA-Seq coverage along the protein body was
18 calculated using pysam(“Website,” n.d.).

19 **Accession**

20 RNA-Seq raw data: PRJNA504034; Proteomics raw data in MassIVE: doi:10.25345/C5M597,
21 Identified peptides in Synapse: doi:10.7303/syn17037373.

22 **Acknowledgements**

23 This work was supported by generous funding from the Novo Nordisk Foundation

1 provided to the Center for Biosustainability at the Technical University of Denmark (grant no.
2 NNF16CC0021858), and SL was supported with funding from the Frontiers of Innovation
3 Scholars Program at UCSD. VB was supported in part by grants from the NIH 1R01GM114362,
4 and P-41-RR24851.

5

6

7

8

9

10

11

12

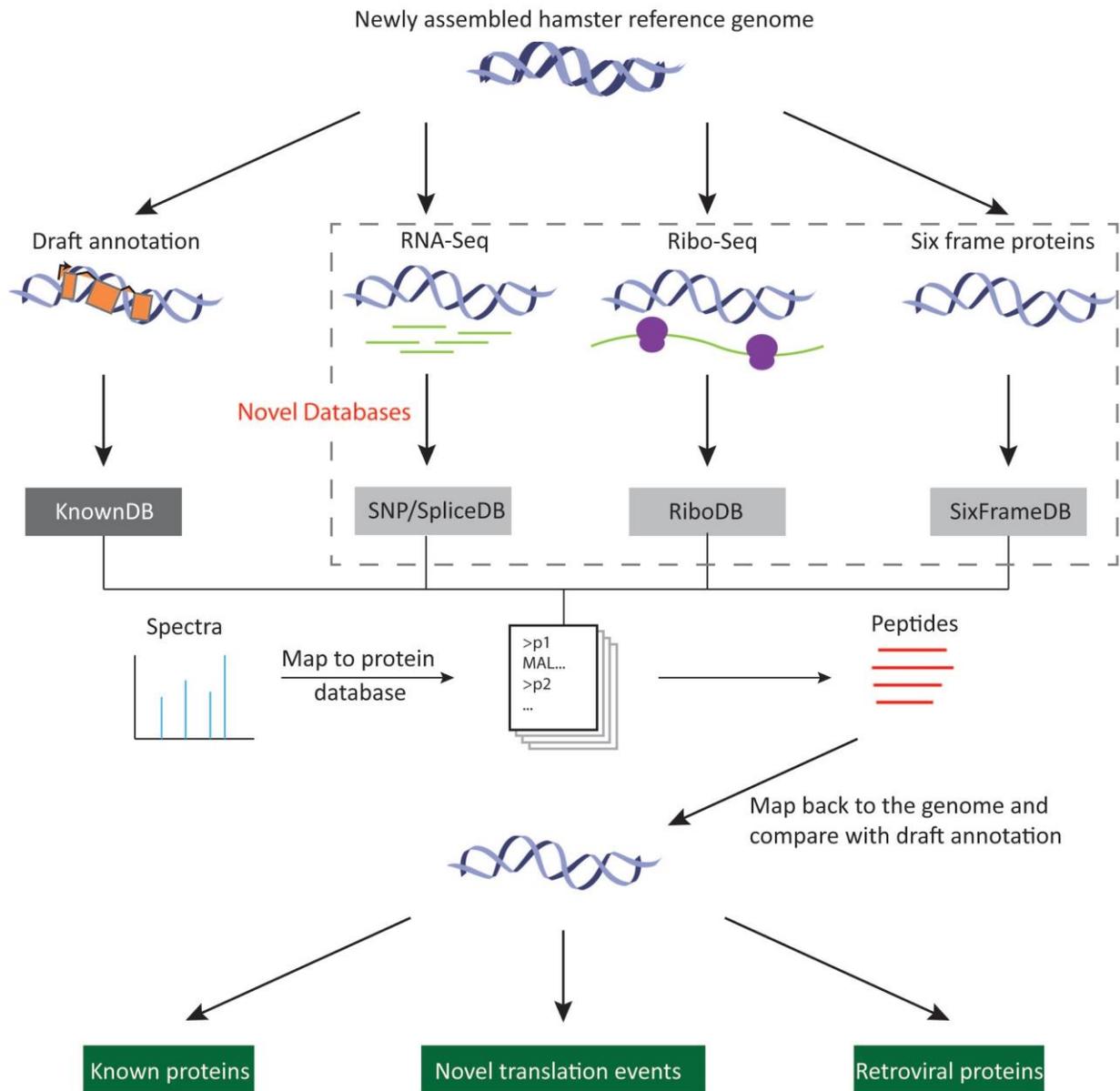
13

14

15

16

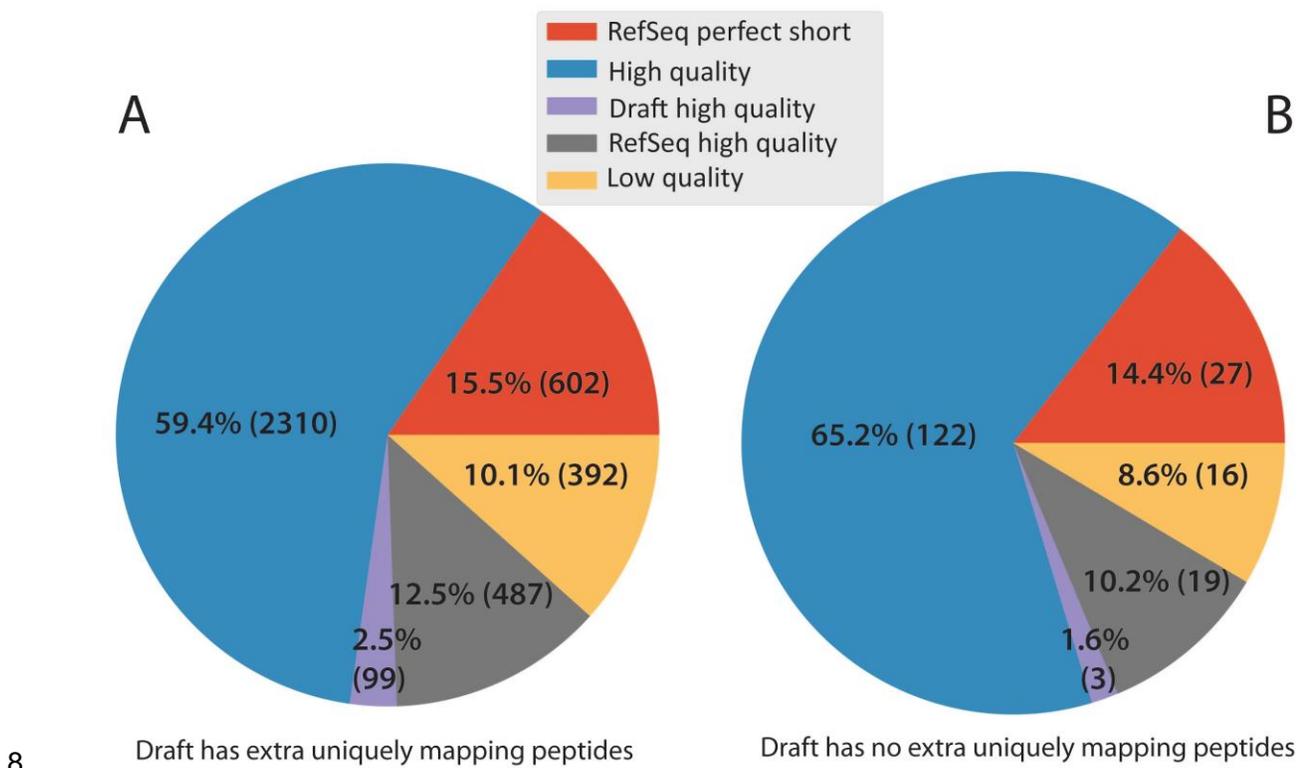
1 Figures:



2

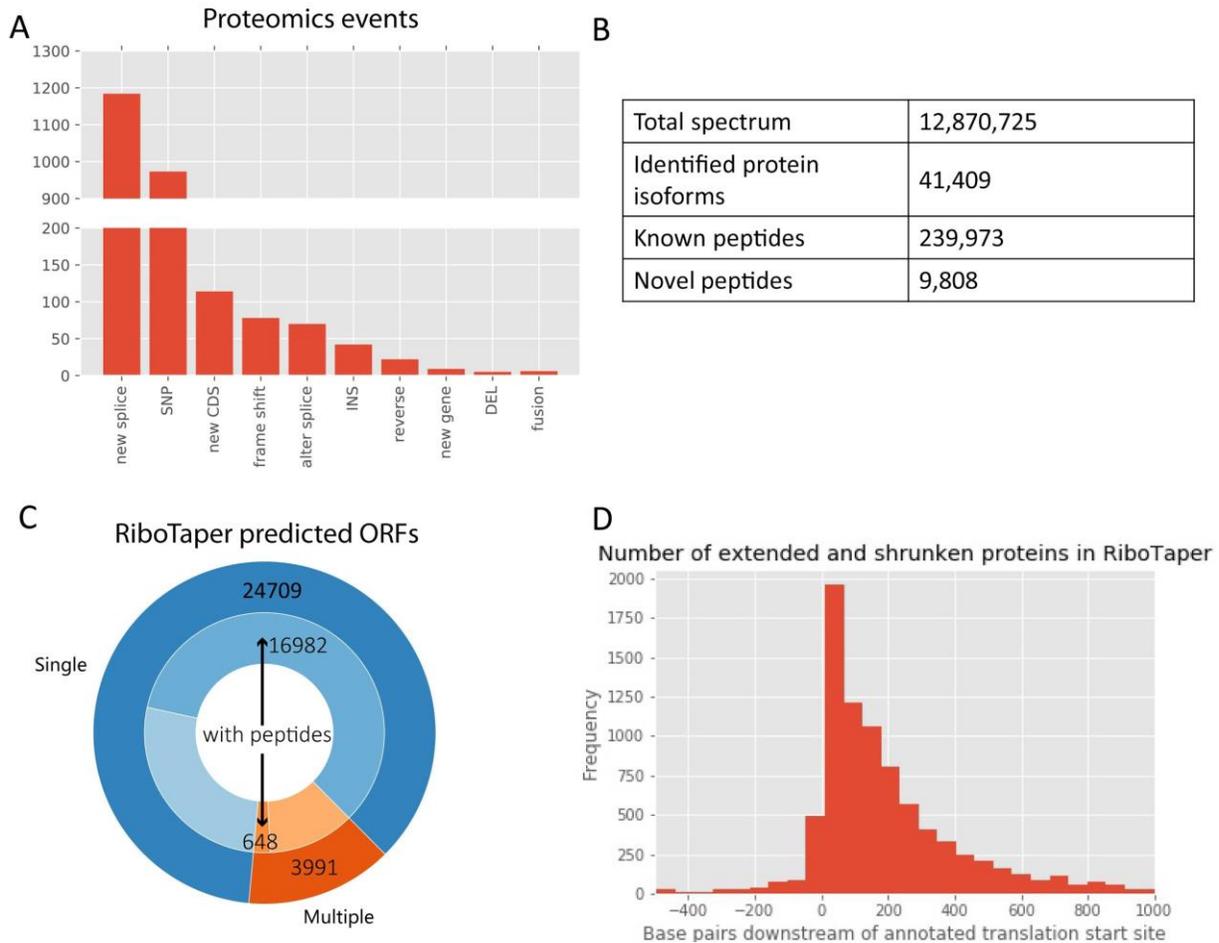
3 **Figure 1: Overview of the proteogenomic pipeline.** Multiple databases of putative protein
4 sequences were generated based on the newly assembled hamster genome (Rupp et al., 2018)
5 and additional data. The KnownDB contains protein sequences from our draft annotation
6 generated here. The SNP/SpliceDB was derived from RNA-Seq samples, and contains

1 candidate mutated or novel spliced proteins compared to draft annotation. The RiboDB was
2 derived from predicted translated ORFs from Ribo-Seq and RNA-Seq. The SixFrameDB is
3 derived from the reference genome (Rupp et al., 2018). After database construction, mass
4 spectra were mapped against the protein databases using MSGF+ to identify the peptides. The
5 peptides were then mapped back to the genome and compared with the draft annotation to
6 verify translated known proteins, enumerate novel translation events and the identity of retroviral
7 proteins.



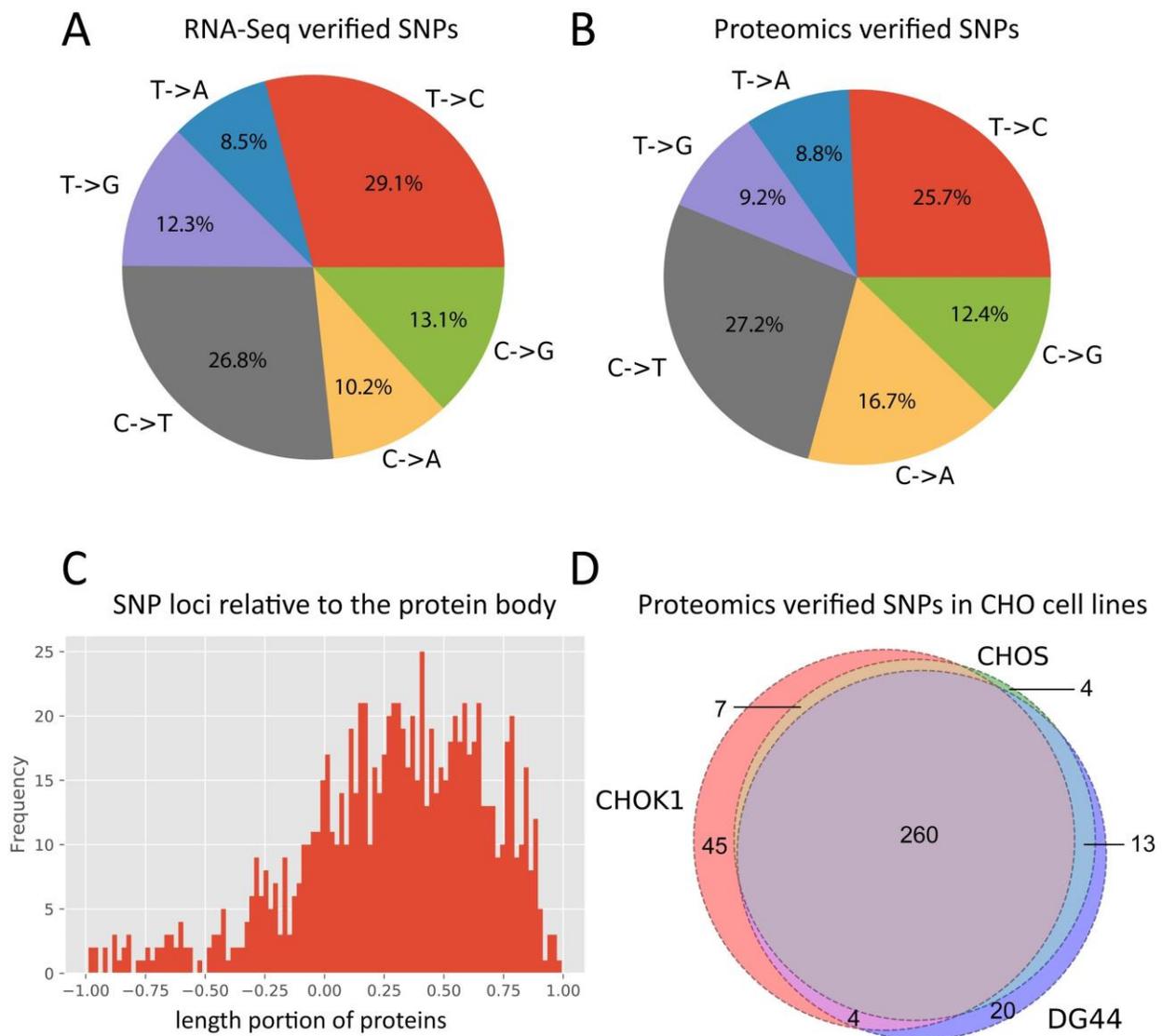
9 **Figure 2: Number of novel draft proteins verified by draft-only peptides in different**
10 **categories.** The draft annotation predicted thousands of novel protein sequences. **(A)** Of these
11 3,890 had uniquely mapped peptides supporting the novel protein sequences. **(B)** Only 187 did
12 not have extra peptide support from uniquely mapping peptides, and thousands provided
13 peptide support. **RefSeq perfect short:** RefSeq proteins map perfectly but are shorter than
14 draft proteins; **High quality:** high quality mapping proteins between draft and RefSeq; **Draft**

- 1 **high quality**: draft proteins map to RefSeq with high quality, but the reverse doesn't hold;
- 2 **RefSeq high quality**: RefSeq proteins map to draft with high quality, but the reverse doesn't
- 3 hold; **Low quality**: low quality mapping between draft and RefSeq.



- 4
- 5 **Figure 3: Proteogenomics and RiboTaper verified predicted protein sequences and**
- 6 **identified novel translation events. (A)** Numerous novel translational events were identified,
- 7 including novel splice sites that are not in the draft annotation file (new splice), non-synonymous
- 8 mutations (SNP), peptides that map to UTR regions or to transcripts with no CDS (new CDS),
- 9 alternative splice sites (alter splice), peptide mapping to reverse strand of reference CDS
- 10 (reverse), insertions (INS), peptide mapping to intergenic regions (new gene), deletion (DEL),
- 11 and gene fusions connecting two genes (fusion). **(B)** Statistics for the number of spectra,

1 peptides and protein isoforms identified in proteogenomics. **(C)** Number of ORFs identified
 2 using RiboTaper. **Outer circle:** Number of transcripts predicted with single ORF (blue) or
 3 multiple ORFs (orange). **Inner circle:** Number of transcripts with (darker blue and orange) or
 4 without (light blue and orange) peptide support. **(D)** Number of proteins that are shorter/longer
 5 than the draft annotation. Positive x axis means the RiboTaper proteins are shorter (i.e., start
 6 later) than the draft annotation.

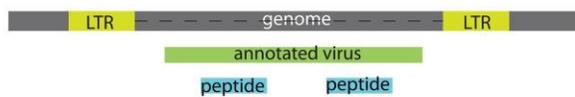


7
 8 **Figure 4: Hundreds of SNPs in hamster and different CHO cell lineages are validated. A**

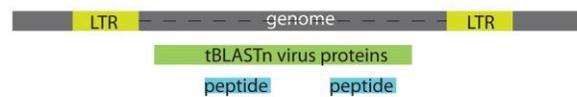
1 comparison of the **(A)** distribution of SNP types identified from RNA-Seq and **(B)** SNP types
 2 verified by proteomics validates the overall distribution of SNPs. **(C)** Peptide-validated non-
 3 synonymous SNPs are located throughout the protein bodies. The length of each protein is
 4 scaled to 1. 0 represents the start codon. SNPs that locate below 0 or above 1 represent
 5 peptide-supported SNPs in 5'-UTR and 3'-UTR regions, respectively. **(D)** Venn diagram of 353
 6 peptide-supported SNPs from CHO-K1, CHO-S and DG44 cell lines shows that most SNPs are
 7 shared across cell lines.

A

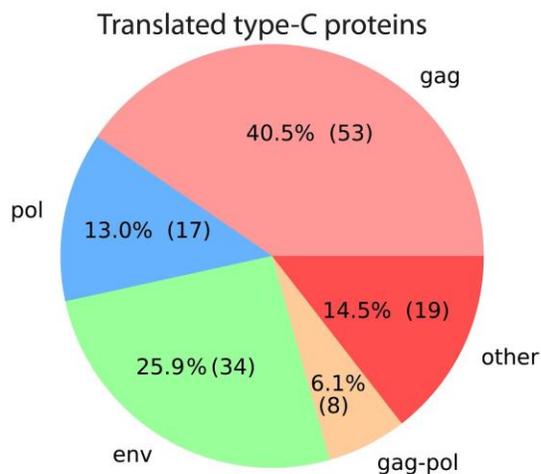
Endogenous retrovirus detection strategy 1



Endogenous retrovirus detection strategy 2

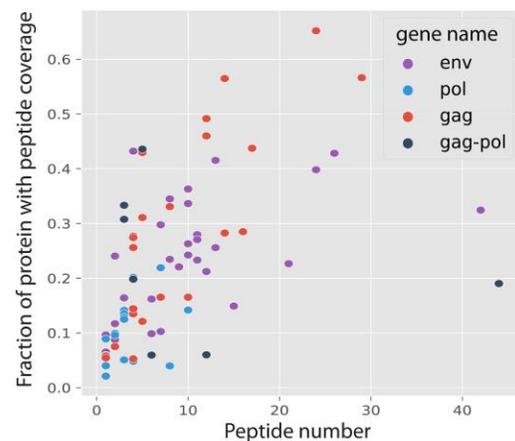


B



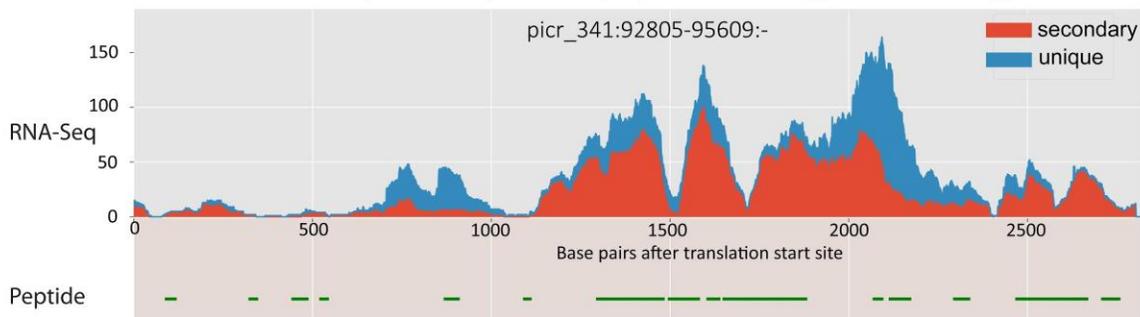
C

Peptide coverage of type-C retrovirus



D

RNA-Seq read depth and peptide coverage in an ENV gene



1 **Figure 5: A proteogenomic identification of the source of translated endogenous**
2 **retroviral particles shed from CHO cells. (A)** Two strategies were taken to identify translated
3 retroviral loci. In strategy 1, peptides were mapped to the annotated retroviral proteins. For
4 strategy 2, the sequences from the NCBI retroviral protein database were aligned to the
5 genome using BLASTP. Then we evaluated the overlap of these aligned peptides with the novel
6 peptides identified from the novel databases in our proteogenomics pipeline. **(B)** The strategies
7 recovered 131 type-C peptide-supported retroviral proteins in CHO cell lines (the “other”
8 category represents non-typical retroviral proteins, such as the p12 protein). **(C)** Peptide-
9 supported type-C virus proteins were analyzed to assess the portion of protein sequence
10 covered by peptides against peptide number. **(D)** Coverage of an envelope protein in reverse
11 strand. **uni**: reads map uniquely to the locus, **sec**: reads are secondary reads and map to
12 multiple loci.

13

14 **References**

- 15 Altschul, S. F., Gish, W., Miller, W., Myers, E. W., & Lipman, D. J. (1990). Basic local alignment
16 search tool. *Journal of Molecular Biology*, 215(3), 403–410.
- 17 Anderson, K. P., Low, M. A., Lie, Y. S., Keller, G. A., & Dinowitz, M. (1991). Endogenous origin
18 of defective retroviruslike particles from a recombinant Chinese hamster ovary cell line.
19 *Virology*, 181(1), 305–311.
- 20 Bolger, A. M., Lohse, M., & Usadel, B. (2014). Trimmomatic: a flexible trimmer for Illumina
21 sequence data. *Bioinformatics*, 30(15), 2114–2120.
- 22 Brinkrolf, K., Rupp, O., Laux, H., Kollin, F., Ernst, W., Linke, B., ... Borth, N. (2013a). Chinese
23 hamster genome sequenced from sorted chromosomes. *Nature Biotechnology*, 31(8), 694–
24 695.
- 25 Brinkrolf, K., Rupp, O., Laux, H., Kollin, F., Ernst, W., Linke, B., ... Borth, N. (2013b). Chinese

- 1 hamster genome sequenced from sorted chromosomes. *Nature Biotechnology*, 31(8), 694–
2 695.
- 3 Calviello, L., Mukherjee, N., Wyler, E., Zauber, H., Hirsekorn, A., Selbach, M., ... Ohler, U.
4 (2016). Detecting actively translated open reading frames in ribosome profiling data. *Nature*
5 *Methods*, 13(2), 165–170.
- 6 Calviello, L., & Ohler, U. (2017). Beyond Read-Counts: Ribo-seq Data Analysis to Understand
7 the Functions of the Transcriptome. *Trends in Genetics: TIG*, 33(10), 728–744.
- 8 Castellana, N., & Bafna, V. (2010). Proteogenomics to discover the full coding content of
9 genomes: a computational perspective. *Journal of Proteomics*, 73(11), 2124–2135.
- 10 Castellana, N. E., Payne, S. H., Shen, Z., Stanke, M., Bafna, V., & Briggs, S. P. (2008).
11 Discovery and revision of Arabidopsis genes by proteogenomics. *Proceedings of the*
12 *National Academy of Sciences of the United States of America*, 105(52), 21034–21038.
- 13 Cesnik, A. J., Shortreed, M. R., Sheynkman, G. M., Frey, B. L., & Smith, L. M. (2016). Human
14 Proteomic Variation Revealed by Combining RNA-Seq Proteogenomics and Global Post-
15 Translational Modification (G-PTM) Search Strategy. *Journal of Proteome Research*, 15(3),
16 800–808.
- 17 Chandramouli, K., & Qian, P.-Y. (2009). Proteomics: challenges, techniques and possibilities to
18 overcome biological sample complexity. *Human Genomics and Proteomics: HGP, 2009*.
19 <https://doi.org/10.4061/2009/239204>
- 20 Chen, C., Le, H., & Goudar, C. T. (2017). Evaluation of two public genome references for
21 chinese hamster ovary cells in the context of rna-seq based gene expression analysis.
22 *Biotechnology and Bioengineering*, 114(7), 1603–1613.
- 23 Coffin, J. M., Hughes, S. H., & Varmus, H. E. (1997). *Retroviruses*. Cold Spring Harbor.
- 24 Crappé, J., Ndah, E., Koch, A., Steyaert, S., Gawron, D., De Keulenaer, S., ... Menschaert, G.
25 (2015). PROTEOFORMER: deep proteome coverage through ribosome profiling and MS
26 integration. *Nucleic Acids Research*, 43(5), e29.

- 1 de Wit, C., Fautz, C., & Xu, Y. (2000). Real-time Quantitative PCR for Retrovirus-like Particle
2 Quantification in CHO Cell Culture. *Biologicals: Journal of the International Association of*
3 *Biological Standardization*, 28(3), 137–148.
- 4 Dinowitz, M., Lie, Y. S., Low, M. A., Lazar, R., Fautz, C., Potts, B., ... Anderson, K. (1992).
5 Recent studies on retrovirus-like particles in Chinese hamster ovary cells. *Developments in*
6 *Biological Standardization*, 76, 201–207.
- 7 Dobin, A., Davis, C. A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., ... Gingeras, T. R.
8 (2013). STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*, 29(1), 15–21.
- 9 Ellinghaus, D., Kurtz, S., & Willhoeft, U. (2008). LTRharvest, an efficient and flexible software
10 for de novo detection of LTR retrotransposons. *BMC Bioinformatics*, 9(1), 18.
- 11 Fischer, S., Handrick, R., & Otte, K. (2015). The art of CHO cell engineering: A comprehensive
12 retrospect and future perspectives. *Biotechnology Advances*, 33(8), 1878–1896.
- 13 Frank, A. M. (2009). A ranking-based scoring function for peptide-spectrum matches. *Journal of*
14 *Proteome Research*, 8(5), 2241–2252.
- 15 Gish, W., & States, D. J. (1993). Identification of protein coding regions by database similarity
16 search. *Nature Genetics*, 3(3), 266–272.
- 17 Golabgir, A., Gutierrez, J. M., Hefzi, H., Li, S., Palsson, B. O., Herwig, C., & Lewis, N. E. (2016).
18 Quantitative feature extraction from the Chinese hamster ovary bioprocess bibliome using a
19 novel meta-analysis workflow. *Biotechnology Advances*, 34(5), 621–633.
- 20 Haas, B. J., Delcher, A. L., Mount, S. M., Wortman, J. R., Smith, R. K., Jr, Hannick, L. I., ...
21 White, O. (2003). Improving the Arabidopsis genome annotation using maximal transcript
22 alignment assemblies. *Nucleic Acids Research*, 31(19), 5654–5666.
- 23 Haas, B. J., Papanicolaou, A., Yassour, M., Grabherr, M., Blood, P. D., Bowden, J., ... Regev,
24 A. (2013). De novo transcript sequence reconstruction from RNA-seq using the Trinity
25 platform for reference generation and analysis. *Nature Protocols*, 8(8), 1494–1512.
- 26 Hogwood, C. E. M., Bracewell, D. G., & Smales, C. M. (2014). Measurement and control of host

- 1 cell proteins (HCPs) in CHO cell bioprocesses. *Current Opinion in Biotechnology*, 30, 153–
2 160.
- 3 Ingolia, N. T., Brar, G. A., Rouskin, S., McGeachy, A. M., & Weissman, J. S. (2013). Genome-
4 wide annotation and quantitation of translation by ribosome profiling. *Current Protocols in*
5 *Molecular Biology / Edited by Frederick M. Ausubel ... [et Al.], Chapter 4*, Unit 4.18.
- 6 Ingolia, N. T., Brar, G. A., Stern-Ginossar, N., Harris, M. S., Talhouarne, G. J. S., Jackson, S.
7 E., ... Weissman, J. S. (2014). Ribosome profiling reveals pervasive translation outside of
8 annotated protein-coding genes. *Cell Reports*, 8(5), 1365–1379.
- 9 Ingolia, N. T., Ghaemmaghami, S., Newman, J. R. S., & Weissman, J. S. (2009). Genome-wide
10 analysis in vivo of translation with nucleotide resolution using ribosome profiling. *Science*,
11 324(5924), 218–223.
- 12 Jaffe, J. D., Berg, H. C., & Church, G. M. (2004). Proteogenomic mapping as a complementary
13 method to perform genome annotation. *Proteomics*, 4(1), 59–77.
- 14 Kallehauge, T. B., Li, S., Pedersen, L. E., Ha, T. K., Ley, D., Andersen, M. R., ... Lewis, N. E.
15 (2017). Ribosome profiling-guided depletion of an mRNA increases cell growth rate and
16 protein secretion. *Scientific Reports*, 7, 40388.
- 17 Kaushik, P., Henry, M., Clynes, M., & Meleady, P. (2018). The Expression Pattern of the
18 Phosphoproteome Is Significantly Changed During the Growth Phases of Recombinant
19 CHO Cell Culture. *Biotechnology Journal*, 13(10), e1700221.
- 20 Kessner, D., Chambers, M., Burke, R., Agus, D., & Mallick, P. (2008). ProteoWizard: open
21 source software for rapid proteomics tools development. *Bioinformatics*, 24(21), 2534–
22 2536.
- 23 Kim, D., Langmead, B., & Salzberg, S. L. (2015). HISAT: a fast spliced aligner with low memory
24 requirements. *Nature Methods*, 12(4), 357–360.
- 25 Kim, M.-S., Pinto, S. M., Getnet, D., Nirujogi, R. S., Manda, S. S., Chaerkady, R., ... Pandey, A.
26 (2014). A draft map of the human proteome. *Nature*, 509(7502), 575–581.

- 1 Kim, S., & Pevzner, P. A. (2014). MS-GF+ makes progress towards a universal database
2 search tool for proteomics. *Nature Communications*, *5*, 5277.
- 3 Kumar, A., Baycin-Hizal, D., Wolozny, D., Pedersen, L. E., Lewis, N. E., Heffner, K., ...
4 Betenbaugh, M. J. (2015). Elucidation of the CHO Super-Ome (CHO-SO) by
5 Proteoinformatics. *Journal of Proteome Research*, *14*(11), 4687–4703.
- 6 Kuo, C.-C., Chiang, A. W., Shamie, I., Samoudi, M., Gutierrez, J. M., & Lewis, N. E. (2017). The
7 emerging role of systems biology for engineering protein production in CHO cells. *Current*
8 *Opinion in Biotechnology*, *51*, 64–69.
- 9 Lee, J. S., Grav, L. M., Lewis, N. E., & Fastrup Kildegaard, H. (2015). CRISPR/Cas9-mediated
10 genome engineering of CHO cell factories: Application and perspectives. *Biotechnology*
11 *Journal*, *10*(7), 979–994.
- 12 Lewis, N. E., Liu, X., Li, Y., Nagarajan, H., Yerganian, G., O'Brien, E., ... Palsson, B. O.
13 (2013a). Genomic landscapes of Chinese hamster ovary cell lines as revealed by the
14 *Cricetulus griseus* draft genome. *Nature Biotechnology*, *31*(8), 759–765.
- 15 Lewis, N. E., Liu, X., Li, Y., Nagarajan, H., Yerganian, G., O'Brien, E., ... Palsson, B. O.
16 (2013b). Genomic landscapes of Chinese hamster ovary cell lines as revealed by the
17 *Cricetulus griseus* draft genome. *Nature Biotechnology*, *31*(8), 759–765.
- 18 Lieber, M. M., Benveniste, R. E., Livingston, D. M., & Todaro, G. J. (1973). Mammalian cells in
19 culture frequently release type C viruses. *Science*, *182*(4107), 56–59.
- 20 Lie, Y. S., Penuel, E. M., Low, M. A., Nguyen, T. P., Mangahas, J. O., Anderson, K. P., &
21 Petropoulos, C. J. (1994). Chinese hamster ovary cells contain transcriptionally active full-
22 length type C proviruses. *Journal of Virology*, *68*(12), 7840–7849.
- 23 Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., ... 1000 Genome Project
24 Data Processing Subgroup. (2009). The Sequence Alignment/Map format and SAMtools.
25 *Bioinformatics*, *25*(16), 2078–2079.
- 26 Nesvizhskii, A. I. (2014). Proteogenomics: concepts, applications and computational strategies.

- 1 *Nature Methods*, 11(11), 1114–1125.
- 2 Perteua, M., Perteua, G. M., Antonescu, C. M., Chang, T.-C., Mendell, J. T., & Salzberg, S. L.
3 (2015). StringTie enables improved reconstruction of a transcriptome from RNA-seq reads.
4 *Nature Biotechnology*, 33(3), 290–295.
- 5 Quinlan, A. R., & Hall, I. M. (2010). BEDTools: a flexible suite of utilities for comparing genomic
6 features. *Bioinformatics*, 26(6), 841–842.
- 7 Richelle, A., & Lewis, N. E. (2017). Improvements in protein production in mammalian cells from
8 targeted metabolic engineering. *Current Opinion in Systems Biology*, 6, 1–6.
- 9 Rupp, O., Becker, J., Brinkrolf, K., Timmermann, C., Borth, N., Pühler, A., ... Goesmann, A.
10 (2014). Construction of a public CHO cell line transcript database using versatile
11 bioinformatics analysis pipelines. *PloS One*, 9(1), e85568.
- 12 Rupp, O., MacDonald, M. L., Li, S., Dhiman, H., Polson, S., Griep, S., ... Lee, K. H. (2018). A
13 reference genome of the Chinese hamster based on a hybrid assembly strategy.
14 *Biotechnology and Bioengineering*. <https://doi.org/10.1002/bit.26722>
- 15 Salmela, L., & Rivals, E. (2014). LoRDEC: accurate and efficient long read error correction.
16 *Bioinformatics*, 30(24), 3506–3514.
- 17 Singh, A., Kildegaard, H. F., & Andersen, M. R. (2018). An Online Compendium of CHO RNA-
18 Seq Data Allows Identification of CHO Cell Line-Specific Transcriptomic Signatures.
19 *Biotechnology Journal*, e1800070.
- 20 Slotkin, R. K., & Keith Slotkin, R. (2018). The case for not masking away repetitive DNA. *Mobile*
21 *DNA*, 9(1). <https://doi.org/10.1186/s13100-018-0120-9>
- 22 Stolfa, G., Smonskey, M. T., Boniface, R., Hachmann, A.-B., Gulde, P., Joshi, A. D., ...
23 Campbell, A. (2018). CHO-Omics Review: The Impact of Current and Emerging
24 Technologies on Chinese Hamster Ovary Based Bioproduction. *Biotechnology Journal*,
25 13(3), e1700227.
- 26 Strauss, D. M., Lute, S., Brorson, K., Blank, G. S., Chen, Q., & Yang, B. (2009). Removal of

- 1 endogenous retrovirus-like particles from CHO-cell derived products using Q sepharose
2 fast flow chromatography. *Biotechnology Progress*, 25(4), 1194–1197.
- 3 Temin, H. M. (1982). Function of the retrovirus long terminal repeat. *Cell*, 28(1), 3–5.
- 4 The UniProt Consortium. (2018). UniProt: the universal protein knowledgebase. *Nucleic Acids*
5 *Research*. <https://doi.org/10.1093/nar/gky092>
- 6 Van der Auwera, G. A., Carneiro, M. O., Hartl, C., Poplin, R., Del Angel, G., Levy-Moonshine,
7 A., ... DePristo, M. A. (2013). From FastQ data to high confidence variant calls: the
8 Genome Analysis Toolkit best practices pipeline. *Current Protocols in Bioinformatics /*
9 *Editorial Board, Andreas D. Baxevanis ... [et Al.]*, 43, 11.10.1–33.
- 10 van Wijk, X. M., Döhrmann, S., Hallström, B. M., Li, S., Voldborg, B. G., Meng, B. X., ... Esko, J.
11 D. (2017). Whole-Genome Sequencing of Invasion-Resistant Cells Identifies Laminin $\alpha 2$ as
12 a Host Factor for Bacterial Invasion. *mBio*, 8(1). <https://doi.org/10.1128/mBio.02128-16>
- 13 Wang, K., Li, M., & Hakonarson, H. (2010). ANNOVAR: functional annotation of genetic variants
14 from high-throughput sequencing data. *Nucleic Acids Research*, 38(16), e164.
- 15 Website. (n.d.). Retrieved July 23, 2018, from <https://github.com/pysam-developers/pysam>
- 16 Weiss, R. A. (1996). Retrovirus classification and cell interactions. *The Journal of Antimicrobial*
17 *Chemotherapy*, 37(suppl B), 1–11.
- 18 Wiśniewski, J. R., Zougman, A., Nagaraj, N., & Mann, M. (2009). Universal sample preparation
19 method for proteome analysis. *Nature Methods*, 6(5), 359–362.
- 20 Woo, S., Cha, S. W., Bonissone, S., Na, S., Tabb, D. L., Pevzner, P. A., & Bafna, V. (2015).
21 Advanced Proteogenomic Analysis Reveals Multiple Peptide Mutations and Complex
22 Immunoglobulin Peptides in Colon Cancer. *Journal of Proteome Research*, 14(9), 3555–
23 3567.
- 24 Woo, S., Cha, S. W., Merrihew, G., He, Y., Castellana, N., Guest, C., ... Bafna, V. (2014).
25 Proteogenomic database construction driven from large scale RNA-seq data. *Journal of*
26 *Proteome Research*, 13(1), 21–28.

- 1 Woo, S., Cha, S. W., Na, S., Guest, C., Liu, T., Smith, R. D., ... Bafna, V. (2014).
2 Proteogenomic strategies for identification of aberrant cancer peptides using large-scale
3 next-generation sequencing data. *Proteomics*, *14*(23-24), 2719–2730.
- 4 Wu, T. D., & Watanabe, C. K. (2005). GMAP: a genomic mapping and alignment program for
5 mRNA and EST sequences. *Bioinformatics*, *21*(9), 1859–1875.
- 6 Xu, X., Nagarajan, H., Lewis, N. E., Pan, S., Cai, Z., Liu, X., ... Wang, J. (2011). The genomic
7 sequence of the Chinese hamster ovary (CHO)-K1 cell line. *Nature Biotechnology*, *29*(8),
8 735–741.
- 9 Yagoub, D., Tay, A. P., Chen, Z., Hamey, J. J., Cai, C., Chia, S. Z., ... Wilkins, M. R. (2015).
10 Proteogenomic Discovery of a Small, Novel Protein in Yeast Reveals a Strategy for the
11 Detection of Unannotated Short Open Reading Frames. *Journal of Proteome Research*,
12 *14*(12), 5038–5047.
- 13 Yandell, M., & Ence, D. (2012). A beginner's guide to eukaryotic genome annotation. *Nature*
14 *Reviews. Genetics*, *13*(5), 329–342.
- 15 Yang, L., Güell, M., Niu, D., George, H., Lesha, E., Grishin, D., ... Church, G. (2015). Genome-
16 wide inactivation of porcine endogenous retroviruses (PERVs). *Science*, *350*(6264), 1101–
17 1104.
- 18
19
20
21
22
23
24
25

1

2

3

4

5

6

7