



Review

Dual use of peptide mass spectra: Protein atlas and genome annotation



Justin W. Walley, Steven P. Briggs*

University of California San Diego, Section of Cell & Developmental Biology, La Jolla, CA 92093-0380, United States

ARTICLE INFO

Article history:

Received 26 August 2014

Received in revised form 3 February 2015

Accepted 24 February 2015

Keywords:

Proteogenomics

Proteomics

Atlas

Annotation

ABSTRACT

One of the objectives of genome science is the discovery and accurate annotation of all protein-coding genes. Proteogenomics has emerged as a methodology that provides orthogonal information to traditional forms of evidence used for genome annotation. By this method, peptides that are identified via tandem mass spectrometry are used to refine protein-coding gene models. Namely, these peptides are used to confirm the translation of predicted protein-coding genes, as evidence of novel genes or for correction of current gene models. Proteogenomics requires deep and broad sampling of the proteome in order to generate sufficient numbers of unique peptides. Therefore, we propose that proteogenomic projects are designed so that the generated peptides can also be used to create a comprehensive protein atlas that quantitatively catalogues protein abundance changes during development and in response to environmental stimulus.

© 2015 Published by Elsevier B.V. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

Contents

1. Introduction	21
2. Proteogenomic enabled annotation	22
3. Proteome sampling for proteogenomics	22
4. Proteome atlas	23
5. Perspective	23
Acknowledgments	23
References	23

1. Introduction

The primary goal of genome annotation efforts is the discovery and accurate annotation of all protein-coding genes. A complete and accurately annotated proteome provides the building blocks for hypothesis-driven research seeking to enhance our understanding of biology. Genome annotation is a complex process involving multiple integrated tools, which have been described in detail [1–5] and are beyond the scope of this review. Briefly, traditional methods of genome annotation rely on combining various forms of evidence. This includes *de novo* gene prediction, which utilizes only patterns in the genomic sequence to infer gene structure. Additionally, transcript sequences from cDNA libraries can be leveraged to enhance

gene prediction. Lastly, sequence conservation with related species can be incorporated into annotation pipelines. While DNA/RNA-based genome annotation approaches perform remarkably well, given the complexity of the challenge, they are currently unable to accurately predict all protein coding genes and their structure. Experimental evidence is required to determine if a transcript is translated and if the predicted protein sequence is correct.

The field of proteogenomics has emerged as a genome-wide method to improve genome annotations as well as to characterize the pattern of gene expression at the protein level. The concept of proteogenomics was introduced, by Jaffe and colleagues [6], as a method that utilizes peptides identified from their tandem mass spectra, for genome annotation (reviewed by [2,7–9]). Since its introduction, proteogenomics has successfully aided in the annotation of numerous prokaryotic and eukaryotic organisms. These studies have demonstrated that deep and broad sampling of the proteome is necessary, for proteogenomics, requiring the generation of hundreds of millions of mass spectra. Furthermore, protein

* Corresponding author.

E-mail address: sbriggs@ucsd.edu (S.P. Briggs).

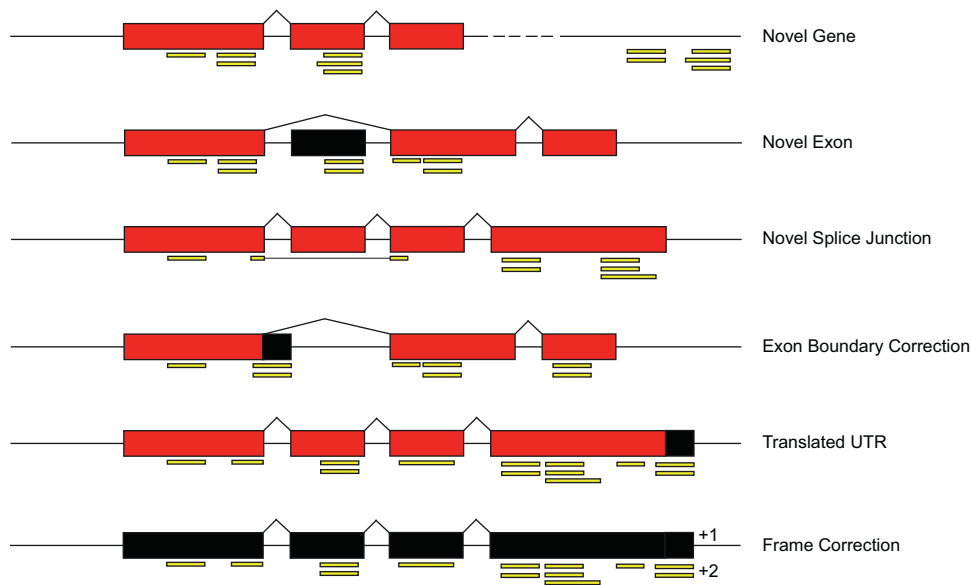


Fig. 1. Examples of gene model revision. Currently annotated exons are shown in red. Gene model revision suggested by novel peptides is depicted in black. Proteogenomically identified peptides are shown in yellow.

accumulation depends upon development and environmental conditions so spectra must be generated from a diverse set of samples to enable deep coverage of the proteome. Such broad sampling enables the additional use of the identified peptides for creation of a protein atlas that catalogs where, when, and how much of a given protein is present.

2. Proteogenomic enabled annotation

Proteogenomics provides a high-throughput method to incorporate protein level information into genome annotation. For this, tandem mass spectra are generated and then used to search genomic databases for peptide identification. The standard database utilized in proteogenomic pipelines is a six-frame translation of the genome [6]. Additionally, specialized types of databases such as an exon–splice graph, which is compact representation of predicted gene structures and splice junctions, have also been exploited [10]. The identified peptides fall into two categories. Namely, confirming peptides that match the current genome annotation and novel peptides, which do not (Fig. 1). It is important to emphasize that the confirming peptides represent critical events, as they directly confirm both the current structural annotation of a gene and demonstrate that the gene encodes a translated protein.

The novel peptides themselves can be further divided into two types of events. One category includes intergenic peptides, which map outside of known genes, and thus reveal the presence of novel genes. A second category is intragenic peptides that fall within a known locus, but do not match the currently annotated gene model. Intragenic peptides include those demonstrating the translation

of 5' or 3' untranslated regions (UTR), alternative start/stop sites, proteins out of frame, incorrect exon boundaries, novel exons or novel splice sites. While one may assume that the identification of these types of novel intergenic and intragenic peptides by proteogenomics to be rare, they are actually commonly found, even in well annotated model organisms (i.e. organisms that have been subjected to multiple rounds of genome annotation) (Table 1). This demonstrates that proteogenomics is a necessary addition to any comprehensive genome annotation effort.

3. Proteome sampling for proteogenomics

Deep and broad sampling of the proteome is necessary for comprehensive proteogenomic efforts. There are numerous strategies that have been developed for proteogenomic experiments to aid in maximizing the number of unique peptides identified by mass spectrometry [7,9,11]. Briefly, fractionation methods such as one-dimensional and two-dimensional gel electrophoresis, as well as gel-free chromatography based separations of proteins and peptides, aid in deep proteome sampling. Specialized sample preparations can also be used to sample subsets of the proteome such as phosphoproteins, basic proteins, small proteins, and N-terminal peptides [7,8,12–14]. Additionally, use of multiple proteases (examples include trypsin, chymotrypsin, Glu-C, and Lys-C) helps to increase the percentage of sequence covered for a given protein. Another consideration is that the proteome composition depends on both developmental and environmental factors. Thus, analyzing a diverse array of samples is critical for achieving comprehensive proteome coverage [12,13].

Table 1

Proteogenomic publications in plants. (If Novel Genes and Model Revision were not clearly identified all values went into the Model Revision Column.)

Organism	Peptides	Proteins	Novel peptides	Novel genes	Model revision	Citation
<i>Arabidopsis thaliana</i>	86,456	13,029	261	22	35	[28]
<i>Arabidopsis thaliana</i>	144,079	12,769	18,024	778	695	[13]
<i>Populus deltoides</i>	4943			56		[34]
<i>Chlamydomonas reinhardtii</i>	9336		932	3	65	[35]
<i>Oryza sativa</i>	15,121	5034	166		40	[36]
<i>Medicago truncatula</i>	78,647	9843	1568	32	293	[37]
<i>Zea mays</i>	225,166	14,615	24,782	165	1904	[38]
<i>Triticum aestivum</i>	203		17	5	8	[39]

Table 2

Plant protein atlas publications. To be considered a protein atlas publication we required the quantification of protein abundance for at least several thousand proteins across three or more cell-types and/or plant anatomical structures.

Organism	Proteome coverage	Samples	Phosphorylation	Citation
<i>Arabidopsis thaliana</i>	13,029	Multiple developmental stages from roots, leaves, flowers and seeds	No	[28]
<i>Arabidopsis thaliana</i>	1995	Six root cell types	No	[29]
<i>Zea mays</i>	14,165	Aleurone/pericarp as well as multiple developmental stages of endosperm and embryo	Yes	[31]
<i>Populus tremula</i> × <i>alba</i>	7538	Leaf, root and stem	No	[30]

4. Proteome atlas

The extensive sampling required for a comprehensive proteogenomic project enables the dual use of the generated peptides for creation of a proteome atlas, which catalogues protein abundance throughout developmental time and/or in response to environmental stimulus. This type of catalogue is relatively common at the mRNA level, where extensive transcriptional atlases have been created for a range of plant species including *Arabidopsis thaliana* [15,16], barley [17], *Oryza sativa* [18,19], *Medicago truncatula* [20], *Glycine max* [21], *Solanum tuberosum* [22], *Zea mays* [23,24], *Rosa chinensis* [25], *Vitis vinifera* [26], and *Lotus japonicus* [27]. However, to our knowledge, there are only a handful of proteome atlas publications in plants, which we define as covering at least several thousand proteins from three or more cell-types and/or plant anatomical structures (Table 2) [28–31]. Well there are only a handful of proteome atlas publications there are several web-based resources including pep2pro [32] and MASCP Gator [33] that aggregate proteome datasets into a single information portal. Finally, an ideal comprehensive protein atlas would provide proteome-wide coverage and include multiple developmental stages, for each organ, as well as a range of environmental perturbations. While this is a daunting task, the ability to leverage the generated peptides for both proteogenomics, as well as building a protein atlas provides a considerable resource for the scientific community.

5. Perspective

Since its inception a decade ago proteogenomics has matured into a robust methodology, thanks in large part to rapid advances in mass spectrometry based proteomics. It is now possible to deeply sample the proteome identifying millions of mass spectra and hundreds of thousands of unique peptides. These unique peptides provide rich fodder not only for genome annotation but also for building protein atlases. Thus, in an ideal scenario all genome annotation pipelines would include proteogenomics and the proteogenomic component would be designed to enable the creation of a quantitative protein atlas.

Acknowledgments

This work was supported by National Science Foundation Grant 0924023 (to S.P.B.) and a National Institutes of Health National Research Service Award Postdoctoral Fellowship F32GM096707 (to J.W.W.).

References

- [1] V. Curwen, E. Eyra, T.D. Andrews, L. Clarke, E. Mongin, S.M.J. Searle, M. Clamp, The Ensembl automatic gene annotation system, *Genome Res.* 14 (2004) 942–950.
- [2] C. Ansong, S.O. Purvine, J.N. Adkins, M.S. Lipton, R.D. Smith, Proteogenomics: needs and roles to be filled by proteomics in genome annotation, *Brief. Funct. Genomics Proteomics* 7 (2008) 50–62.
- [3] M.R. Brent, Steady progress and recent breakthroughs in the accuracy of automated genome annotation, *Nat. Rev. Genet.* 9 (2008) 62–73.
- [4] C. Liang, L. Mao, D. Ware, L. Stein, Evidence-based gene predictions in plant genomes, *Genome Res.* 19 (2009) 1912–1923.
- [5] M. Yandell, D. Ence, A beginner's guide to eukaryotic genome annotation, *Nat. Rev. Genet.* 13 (2012) 329–342.
- [6] J.D. Jaffe, H.C. Berg, G.M. Church, Proteogenomic mapping as a complementary method to perform genome annotation, *Proteomics* 4 (2004) 59–77.
- [7] J. Armengaud, A perfect genome annotation is within reach with the proteomics and genomics alliance, *Curr. Opin. Microbiol.* 12 (2009) 292–300.
- [8] N. Castellana, V. Bafna, Proteogenomics to discover the full coding content of genomes: a computational perspective, *J. Proteomics* 73 (2010) 2124–2135.
- [9] S. Renue, R. Chaerkady, A. Pandey, Proteogenomics, *Proteomics* 11 (2011) 620–630.
- [10] S. Tanner, Z. Shen, J. Ng, L. Florea, R. Guigo, S.P. Briggs, V. Bafna, Improving gene annotation using peptide mass spectrometry, *Genome Res.* 17 (2007) 231–239.
- [11] K. Krug, S. Nahnsen, B. Macek, Mass spectrometry at the interface of proteomics and genomics, *Mol. Biosyst.* 7 (2011) 284–291.
- [12] E. Brunner, et al., A high-quality catalog of the *Drosophila melanogaster* proteome, *Nat. Biotechnol.* 25 (2007) 576–583.
- [13] N.E. Castellana, S.H. Payne, Z. Shen, M. Stanke, V. Bafna, S.P. Briggs, Discovery and revision of *Arabidopsis* genes by proteogenomics, *Proc. Natl. Acad. Sci. U. S. A.* 105 (2008) 21034–21038.
- [14] S. Gallien, E. Perrodou, C. Carapito, C. Deshayes, J.M. Reytrat, A. Van Dorsseleer, O. Poch, C. Schaeffer, O. Lecompte, Ortho-proteogenomics: multiple proteomes investigation through orthology and a new MS-based protocol, *Genome Res.* 19 (2009) 128–135.
- [15] P. Zimmermann, M. Hirsch-Hoffmann, L. Hennig, W. Gruissem, GENEVESTIGATOR. *Arabidopsis* microarray database and analysis toolbox, *Plant Physiol.* 136 (2004) 2621–2632.
- [16] T. Hruz, O. Laule, G. Szabo, F. Wessendorp, S. Bleuler, L. Oertle, P. Widmayer, W. Gruissem, P. Zimmermann, Genevestigator V3: a reference expression database for the meta-analysis of transcriptomes, *Adv. Bioinform.* 2008 (2008) 420747.
- [17] A. Druka, et al., An atlas of gene expression from seed to seed through barley development, *Funct. Integr. Genomics* 6 (2006) 202–211.
- [18] Y. Jiao, et al., A transcriptome atlas of rice cell types uncovers cellular, functional and developmental hierarchies, *Nat. Genet.* 41 (2009) 258–263.
- [19] L. Wang, W. Xie, Y. Chen, W. Tang, J. Yang, R. Ye, L. Liu, Y. Lin, C. Xu, J. Xiao, Q. Zhang, A dynamic gene expression atlas covering the entire life cycle of rice, *Plant J.* 61 (2010) 752–766.
- [20] V.A. Benedito, et al., A gene expression atlas of the model legume *Medicago truncatula*, *Plant J.* 55 (2008) 504–513.
- [21] M. Libault, A. Farmer, T. Joshi, K. Takahashi, R.J. Langley, L.D. Franklin, J. He, D. Xu, G. May, G. Stacey, An integrated transcriptome atlas of the crop model *Glycine max* and its use in comparative analyses in plants, *Plant J.* 63 (2010) 86–99.
- [22] A.N. Massa, K.L. Childs, H. Lin, G.J. Bryan, G. Giuliano, C.R. Buell, The transcriptome of the reference potato genome *Solanum tuberosum* Group Phureja clone DM1-3 516R44, *PLoS ONE* 6 (2011) e26801.
- [23] R.S. Sekhon, H. Lin, K.L. Childs, C.N. Hansey, C.R. Buell, N. de Leon, S.M. Kaeppler, Genome-wide atlas of transcription during maize development, *Plant J.* 66 (2011) 553–563.
- [24] R.S. Sekhon, R. Briskine, C.N. Hirsch, C.L. Myers, N.M. Springer, C.R. Buell, N. de Leon, S.M. Kaeppler, Maize gene atlas developed by RNA sequencing and comparative evaluation of transcriptomes based on RNA sequencing and microarrays, *PLOS ONE* 8 (2013) e61005.
- [25] A. Dubois, et al., Transcriptome database resource and gene expression atlas for the rose, *BMC Genomics* 13 (2012) 638.
- [26] M. Fasoli, et al., The grapevine expression atlas reveals a deep transcriptome shift driving the entire plant into a maturation program, *Plant Cell* 24 (2012) 3489–3505.
- [27] J. Verdier, I. Torres-Jerez, M. Wang, A. Andriankaja, S.N. Allen, J. He, Y. Tang, J.D. Murray, M.K. Udvardi, Establishment of the *Lotus japonicus* Gene Expression Atlas (LjGEA) and its use to explore legume seed maturation, *Plant J.* 74 (2013) 351–362.
- [28] K. Baerenfaller, J. Grossmann, M.A. Grobei, R. Hull, M. Hirsch-Hoffmann, S. Yalovsky, P. Zimmermann, U. Grossniklaus, W. Gruissem, S. Baginsky, Genome-scale proteomics reveals *Arabidopsis thaliana* gene models and proteome dynamics, *Science* 320 (2008) 938–941.
- [29] J.J. Petricka, M.A. Schauer, M. Megraw, N.W. Breakfield, J.W. Thompson, S. Georgiev, E.J. Soderblom, U. Ohler, M.A. Moseley, U. Grossniklaus, P.N. Benfey, The protein expression landscape of the *Arabidopsis* root, *Proc. Natl. Acad. Sci. U. S. A.* 109 (2012) 6811–6818.

- [30] P. Abraham, R.J. Giannone, R.M. Adams, U. Kalluri, G.A. Tuskan, R.L. Hettich, Putting the pieces together: high-performance LC-MS/MS provides network-, pathway-, and protein-level perspectives in *Populus*, *Mol. Cell. Proteomics* 12 (2013) 106–119.
- [31] J.W. Walley, Z. Shen, R. Sartor, K.J. Wu, J. Osborn, L.G. Smith, S.P. Briggs, Reconstruction of protein networks from an atlas of maize seed proteotypes, *Proc. Natl. Acad. Sci. U. S. A.* 110 (2013) E4808–E4817.
- [32] K. Baerenfaller, M. Hirsch-Hoffmann, J. Svozil, R. Hull, D. Russenberger, S. Bischof, Q. Lu, W. Gruissem, S. Baginsky, pep2pro: a new tool for comprehensive proteome data analysis to reveal information about organ-specific proteomes in *Arabidopsis thaliana*, *Integr. Biol.* 3 (3) (2011) 225–237.
- [33] H.J. Joshi, et al., MASP Gator: an aggregation portal for the visualization of *Arabidopsis* proteomics data, *Plant Physiol.* 155 (2011) 259–270.
- [34] X. Yang, et al., Discovery and annotation of small proteins using genomics proteomics, and computational approaches, *Genome Res.* 21 (2011) 634–641.
- [35] M. Specht, M. Stanke, M. Terashima, B. Naumann-Busch, I. Janßen, R. Höhner, E.F.Y. Hom, C. Liang, M. Hippler, Concerted action of the new Genomic Peptide Finder and AUGUSTUS allows for automated proteogenomic annotation of the *Chlamydomonas reinhardtii* genome, *Proteomics* 11 (2011) 1814–1823.
- [36] M. Helmy, M. Tomita, Y. Ishihama, OryzaPG-DB: rice proteome database based on shotgun proteogenomics, *BMC Plant Biol.* 11 (2011) 63.
- [37] J.D. Volkening, D.J. Bailey, C.M. Rose, P.A. Grimsrud, M. Howes-Podoll, M. Venkateshwaran, M.S. Westphall, J.-M. Ané, J.J. Coon, M.R. Sussman, A proteogenomic survey of the *Medicago truncatula* genome, *Mol. Cell. Proteomics* 11 (2012) 933–944.
- [38] N.E. Castellana, Z. Shen, Y. He, J.W. Walley, C.J. Cassidy, S.P. Briggs, V. Bafna, An automated proteogenomic method uses mass spectrometry to reveal novel genes in *Zea mays*, *Mol. Cell. Proteomics* 13 (2014) 157–167.
- [39] K.F.X. Mayer, et al., A chromosome-based draft sequence of the hexaploid bread wheat (*Triticum aestivum*) genome, *Science* 345 (2014) 1251788.