

Human Proteome Project Mass Spectrometry Data Interpretation Guidelines 2.1

DEUTSCH, Eric W, *et al.*

Abstract

Every data-rich community research effort requires a clear plan for ensuring the quality of the data interpretation and comparability of analyses. To address this need within the Human Proteome Project (HPP) of the Human Proteome Organization (HUPO), we have developed through broad consultation a set of mass spectrometry data interpretation guidelines that should be applied to all HPP data contributions. For submission of manuscripts reporting HPP protein identification results, the guidelines are presented as a one-page checklist containing fifteen essential points followed by two pages of expanded description of each. Here, we present an overview of the guidelines and provide an in-depth description of each of the fifteen elements to facilitate understanding of the intentions and rationale behind the guidelines, both for authors and for reviewers. Broadly, these guidelines provide specific directions regarding how HPP data are to be submitted to mass spectrometry data repositories, how error analysis should be presented, and how detection of novel proteins should be supported with additional confirmatory evidence. These [...]

Reference

DEUTSCH, Eric W, *et al.* Human Proteome Project Mass Spectrometry Data Interpretation Guidelines 2.1. *Journal of proteome research*, 2016, vol. 15, no. 11, p. 3961-3970

DOI : 10.1021/acs.jproteome.6b00392

PMID : 27490519

Available at:

<http://archive-ouverte.unige.ch/unige:86190>

Disclaimer: layout of this document may differ from the published version.



UNIVERSITÉ
DE GENÈVE

Human Proteome Project Mass Spectrometry Data Interpretation Guidelines 2.1

Eric W. Deutsch, Christopher M Overall, Jennifer E. Van Eyk, Mark S. Baker, Young-Ki Paik, Susan T. Weintraub, Lydie Lane, Lennart Martens, Yves Vandenbrouck, Ulrike Kusebauch, William S. Hancock, Henning Hermjakob, Ruedi Aebersold, Robert L. Moritz, and Gilbert S Omenn

J. Proteome Res., **Just Accepted Manuscript** • DOI: 10.1021/acs.jproteome.6b00392 • Publication Date (Web): 04 Aug 2016

Downloaded from <http://pubs.acs.org> on August 11, 2016

Just Accepted

“Just Accepted” manuscripts have been peer-reviewed and accepted for publication. They are posted online prior to technical editing, formatting for publication and author proofing. The American Chemical Society provides “Just Accepted” as a free service to the research community to expedite the dissemination of scientific material as soon as possible after acceptance. “Just Accepted” manuscripts appear in full in PDF format accompanied by an HTML abstract. “Just Accepted” manuscripts have been fully peer reviewed, but should not be considered the official version of record. They are accessible to all readers and citable by the Digital Object Identifier (DOI®). “Just Accepted” is an optional service offered to authors. Therefore, the “Just Accepted” Web site may not include all articles that will be published in the journal. After a manuscript is technically edited and formatted, it will be removed from the “Just Accepted” Web site and published as an ASAP article. Note that technical editing may introduce minor changes to the manuscript text and/or graphics which could affect content, and all legal disclaimers and ethical guidelines that apply to the journal pertain. ACS cannot be held responsible for errors or consequences arising from the use of information contained in these “Just Accepted” manuscripts.

Human Proteome Project Mass Spectrometry Data Interpretation

Guidelines 2.1

Eric W. Deutsch^{1,*}, Christopher M. Overall², Jennifer E. Van Eyk³, Mark S. Baker⁴, Young-Ki Paik⁵, Susan T. Weintraub⁶, Lydie Lane⁷, Lennart Martens^{8,9}, Yves Vandenbrouck¹⁰, Ulrike Kusebauch¹, William S. Hancock¹¹, Henning Hermjakob^{12,13}, Ruedi Aebersold^{14,15}, Robert L. Moritz¹, and Gilbert S. Omenn^{1,16}

¹ Institute for Systems Biology, Seattle, WA, USA

² Centre for Blood Research, Departments of Oral Biological & Medical Sciences, and Biochemistry & Molecular Biology, Faculty of Dentistry, University of British Columbia, Vancouver, Canada

³Advanced Clinical Biosystems Research Institute, Department of Medicine, Cedars Sinai Medical Center, Los Angeles, CA, USA

⁴ Department of Biomedical Sciences, Faculty of Medicine and Health Science, Macquarie University, NSW, Australia

⁵ Yonsei Proteome Research Center and Department of Biochemistry, Yonsei University, 50 Yonsei-ro, Sudaemoon-ku, Seoul, Korea

⁶ The University of Texas Health Science Center at San Antonio, San Antonio, Texas, USA

⁷ SIB Swiss Institute of Bioinformatics and Department of Human Protein Science, Faculty of medicine, University of Geneva, CMU, Michel Servet 1, 1211 Geneva 4, Switzerland

⁸ Department of Medical Protein Research, VIB, Ghent, Belgium

⁹ Department of Biochemistry, Ghent University, Ghent, Belgium

¹⁰ French Proteomics Infrastructure, Biosciences and Biotechnology Institute of Grenoble (BIG), Université Grenoble Alpes, CEA, INSERM, U1038, Grenoble, France.

¹¹ Department of Chemical Biology, Northeastern University, Boston, Massachusetts, USA

¹² European Molecular Biology Laboratory, European Bioinformatics Institute (EMBL-EBI), Wellcome Trust Genome Campus, Hinxton, Cambridge, CB10 1SD, UK

¹³ National Center for Protein Sciences, Beijing, China

¹⁴ Department of Biology, Institute of Molecular Systems Biology, ETH Zurich, Zurich, Switzerland

¹⁵ Faculty of Science, University of Zurich, 8006 Zurich, Switzerland

¹⁶ Departments of Computational Medicine & Bioinformatics, Internal Medicine, and Human Genetics and School of Public Health, University of Michigan, Ann Arbor, MI, 48109-2218, USA

*Address correspondence to: Eric W. Deutsch, Institute for Systems Biology, 401 Terry Ave N, Seattle, WA 98109, USA, Email: edeutsch@systemsbiology.org, Phone: 206-732-1200, Fax: 206-732-1299

Keywords: Guidelines, standards, Human Proteome Project, mass spectrometry, false-discovery rates, alternative protein matches

Abstract

Every data-rich community research effort requires a clear plan for ensuring the quality of the data interpretation and comparability of analyses. To address this need within the Human Proteome Project (HPP) of the Human Proteome Organization (HUPO), we have developed through broad consultation a set of mass spectrometry data interpretation guidelines that should be applied to all HPP data contributions. For submission of manuscripts reporting HPP protein identification results, the guidelines are presented as a one-page checklist containing fifteen essential points followed by two pages of expanded description of each. Here, we present an overview of the guidelines and provide an in-depth description of each of the fifteen elements to facilitate understanding of the intentions and rationale behind the guidelines, both for authors and for reviewers. Broadly, these guidelines provide specific directions regarding how HPP data are to be submitted to mass spectrometry data repositories, how error analysis should be presented, and how detection of novel proteins should be supported with additional confirmatory evidence. These guidelines, developed by the HPP community, are presented to the broader scientific community for further discussion.

Introduction

The flagship scientific project of the Human Proteome Organization (HUPO), known as the Human Proteome Project (HPP), is composed of 50 teams of scientists organized as the Chromosome-Centric HPP (C-HPP), the Biology and Disease-driven HPP (B/D-HPP), and the three resource pillars for Antibodies, Mass Spectrometry, and Knowledge Bases. The HPP is an international effort to advance the understanding of all aspects of the human proteome. Its initial primary aim is to develop a full “parts list” of proteins that are present in human cells, organs and biofluids. Beyond, the HPP aims to advance our understanding of protein interactions and functions in health and disease, and enable the widespread use of proteomics technologies through enhanced techniques and resources by the broader scientific community¹. One of the major goals for the C-HPP in establishing the full parts list is to obtain conclusive mass spectrometry (MS) evidence for what are termed “missing proteins”—the set of polypeptide sequences predicted to be translated from the genome and transcriptome, but for which there is not yet sufficient high-stringency evidence that such translation takes place²⁻⁴. The conclusive detection of these missing proteins, which are specified as having a PE (protein existence) designation of 2, 3, or 4 in the neXtProt⁵ knowledge base, as well as reported translation products from novel coding elements, requires compelling evidence. This includes an interpretation that clearly takes into account the inherent uncertainties currently found in high-throughput MS data acquisition techniques and sequence matching to still-evolving protein reference databases.

MS proteomics is a powerful technology that has enabled routine high-throughput identification and quantification of proteins in complex samples. There are several different MS techniques, including shotgun proteomics via data-dependent acquisition (DDA)^{6,7}, data-independent acquisition (DIA) (e.g., SWATH-MS⁸), and targeted proteomics via selected reaction monitoring (SRM; sometimes called multiple reaction monitoring, MRM). Each has different capabilities and strengths that can be brought to bear, depending on the goals of the analysis. Although many variations exist, a typical workflow involves

1
2
3 extracting and fractionating proteins from a sample, cleaving proteins into peptides using a protease such
4 as trypsin, fractionating the obtained peptides to reduce complexity through methods such as liquid
5 chromatography, and then introducing these fractionated peptides as charged ions into a mass
6 spectrometer, typically by coupling chromatography to an electrospray device. The resulting peptide ions
7 are subsequently fragmented in the instrument and spectral data of these fragments are recorded.
8
9

10
11
12 The data generated from the mass spectrometer are then subjected to extensive computational analysis to
13 determine which peptide ions likely yielded the observed fragment ions, along with confidence metrics
14 for identification and abundance measures⁹. There is a wide variety of informatics tools available for
15 these data analysis tasks, both commercial and free and/or open source¹⁰. However, most of these tools
16 are specific to only one type of MS technique. Confidence metrics reported by these tools are a crucial
17 component of the data analysis because different approaches, instruments and analysis parameters result
18 in different inherent uncertainties in data interpretation. These confidence metrics should be calculated at
19 the peptide-spectrum-match (PSM) level, the aggregated peptide level, and the aggregated protein level,
20 both at a global experiment level and individually. These confidence metrics must then be carefully
21 considered when performing downstream interpretation of the results and functional validation of missing
22 proteins.
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39

40
41 Every data-rich community research effort requires a clear plan to ensure that data are of high quality and
42 comparable between analyses. Over the years, several sets of guidelines have been developed in the field
43 of proteomics, including those from within HUPO. Each set of guidelines has been distinct in its focus
44 and goals; no single set of guidelines is applicable to all goals. The Minimum Information About a
45 Proteomics Experiment (MIAPE) guidelines¹¹ developed by the HUPO Proteomics Standards Initiative
46 (PSI)¹² focus specifically on the metadata annotation of experimental MS results. These metadata must
47 describe what was done to execute the experiment with sufficient detail that the results may be properly
48 interpreted or reproduced; MIAPE explicitly does not stipulate how an analysis is to be performed.
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 Several journals have developed their own guidelines, notably the *Journal of Proteome Research* (JPR)
4 (http://pubs.acs.org/paragonplus/submission/jprobs/jprobs_proteomics_guidelines.pdf), *Molecular and*
5
6 *Cellular Proteomics* (MCP)^{13,14}, and *Proteomics Clinical Applications*
7
8 ([http://onlinelibrary.wiley.com/journal/10.1002/%28ISSN%291862-](http://onlinelibrary.wiley.com/journal/10.1002/%28ISSN%291862-8354/homepage/ForAuthors.html#exp)
9
10 [8354/homepage/ForAuthors.html#exp](http://onlinelibrary.wiley.com/journal/10.1002/%28ISSN%291862-8354/homepage/ForAuthors.html#exp)). These guidelines specify what information must be included in a
11
12 submitted manuscript, as well as some basic expectations about how the acquired MS spectral data are
13
14 interpreted. Three tiers of guidelines for targeted quantitative workflows, depending on the purpose of the
15
16 assay, have been proposed by an NIH-NCI working group;¹⁵ another group has benchmarked and
17
18 proposed a set of guidelines specifically for proteogenomics efforts^{16,17}. These latter two sets provide
19
20 significant guidance on how an analysis should be performed, quite unlike the MIAPE approach, which
21
22 requires extensive disclosure on whatever analysis was performed to be fully compliant.
23
24
25
26
27
28

29 In 2012, the HPP posted an initial version of guidelines (available at <http://www.thehpp.org/guidelines>)
30
31 that focused primarily on ensuring that all data generated and published as part of the consortium effort
32
33 were deposited to one of the ProteomeXchange Consortium¹⁸ repositories for proteomics data, or another
34
35 suitable repository for other data types. A further requirement of the HPP version 1.0 guidelines required
36
37 an analysis threshold of no more than 1% false discovery rate (FDR) at the protein level. Despite having
38
39 served the HPP well for the past three years, it has been recognized that the pursuit of confident
40
41 identification of missing proteins, in particular, and also claims of novel translation products from long-
42
43 non-coding RNAs or pseudogenes, required an updated set of guidelines with more stringent criteria¹⁹.
44
45
46
47
48

49 A new set of guidelines, the HPP Mass Spectrometry Data Interpretation Guidelines Version 2.1, has,
50
51 therefore, been developed and discussed among the HPP community members to address the stringency
52
53 of data required to identify missing proteins or novel coding elements. These guidelines are intended to be
54
55 applied to identification results rather than quantitative results. The guidelines are presented as a one-page
56
57 checklist followed by two pages of expanded descriptions for each of the fifteen items in the checklist
58
59
60

1
2
3 (See Supplementary Material for this document). In this article we describe the development of the
4 guidelines, and we provide a deeper discussion on the reasoning behind these guidelines. We also provide
5 examples of common missteps seen in submitted manuscripts that prompted the development of the
6 guidelines. These guidelines have been adopted as a requirement for articles that will be published as part
7 of the HPP, and now are offered to the community for discussion and potential adoption elsewhere either
8 in whole or by incorporation into other guidelines.
9
10
11
12
13
14
15
16
17
18

19 **Development of the Guidelines**

20
21 A set of preliminary guidelines and discussion points was brought to the HPP Bioinformatics Workshop
22 at the 14th HUPO World Congress held in Vancouver, Canada, in September 2015. Each of the items was
23 discussed and additional input collected. The proposed guidelines were also extensively further debated at
24 the Bioinformatics Hub (Mohammed et al., in preparation) ([http://www.psidev.info/hupo2015-
25 bioinformatics-hub](http://www.psidev.info/hupo2015-bioinformatics-hub)) at the same Congress. This informal venue and the post-Congress HPP Workshop
26 enabled additional hours of discussion and refinement of the individual points. Following HUPO-2015,
27 the proposed elements were written into a draft guidelines document consisting of a one-page checklist
28 followed by two pages of additional detail about each of the checklist items. The document was circulated
29 among the HUPO and HPP leadership, and further edited and refined. The document was then approved
30 by the HPP and HUPO Executive Committees.
31
32
33
34
35
36
37
38
39
40
41
42
43
44

45 Version 2.0 of the guidelines was released on November 12, 2015 at <http://thehpp.org/guidelines> and at
46 www.c-hpp.org. After this release, minor clarifications to the wording were applied in versions 2.0.1 and
47 2.0.2. Further minor clarifications in the wording were applied during the preparation of this article,
48 resulting in version 2.0.3. These guidelines are in effect for all HPP papers submitted on or after
49 November 12, 2015, and are specifically applicable for all manuscripts for the JPR C-HPP 2016 Special
50 Issue. In response to the review of this article, guideline 2 was changed in a substantial manner, as
51
52
53
54
55
56
57
58
59
60

1
2
3 described below, and the guidelines document was updated to version 2.1.0 on July 6, 2016. Future
4 changes to the guidelines will be described at the same URLs above. Wording changes or clarifications
5 that do not change the intended interpretation of the guidelines will only invoke a version change in the
6 third digit (i.e. version x.x.1 to x.x.2). A substantive change to one or more guidelines will increment the
7 minor version digit (i.e. version x.0.x to x.1.0). A major rewrite would increment the major version digit
8 (i.e., version 2.x.x to 3.0.0). The process for making additional changes to these guidelines is as follows:
9 proposed changes are presented to the HPP Executive Committee for discussion and approval or
10 rejection. Approved changes trigger an increment in the version number as described above, and the
11 revised checklist and extended description is posted on the HPP web site, accompanied by
12 announcements to HPP participants and HUPO membership.
13
14
15
16
17
18
19
20
21
22
23
24
25
26

27 **Elaboration on the Guidelines**

28
29
30 As discussed above, page one of the version 2.1.0 guidelines (available at <http://thehpp.org/guidelines>) is
31 a checklist of fifteen items, with a check-box provided to signify compliance, and empty space for
32 explanation by authors of any elements of non-adherence to the guidelines. Each box should be checked,
33 marked as N/A (not applicable) for cases where the guideline is simply not relevant to a manuscript, or
34 marked with an asterisk (*) for non-compliance (NC), which must be explained and justified in the space
35 provided, extending to additional pages if necessary. In most instances, there should not be any NC
36 markings. However, if any specific non-compliance is well justified, or if scenarios that were unforeseen
37 by the drafters of these guidelines that prevent adherence do arise, reviewers or editors may make
38 exceptions.
39
40
41
42
43
44
45
46
47
48
49
50
51

52 The second part is an expanded description of each of the fifteen items. This section provides additional
53 points of clarification for each item and should be consulted by those not yet familiar with these
54 guidelines.
55
56
57
58
59
60

1
2
3
4
5 The final and third part is this article, which provides a full description and discussion of the reasoning
6 behind each guideline. The article should be read by all authors submitting an HPP manuscript, and by
7 those who still have questions about the guidelines that were not answered in the expanded description. It
8 is hoped that this three-tiered approach makes it as easy as possible to comply with the guidelines. Table
9
10 1 provides a list of the fifteen guidelines, each described in one or two sentences as listed in version 2.1.0
11 of the guidelines. Table 1 is provided here for ease of reading, but is not a substitute for the primary
12 checklist available in the Supplementary Material and periodically updated at <http://thehpp.org/guidelines>.
13
14
15
16
17
18
19
20
21
22

23 **Discussion of individual Guidelines**

24 25 26 27 **1. Complete this HPP Data Interpretation Guidelines checklist and submit with your** 28 **manuscript.** 29

30
31 For ease of use for completion and compliance checking, the checklist is presented as a one-page table.
32 This allows those familiar with the checklist to quickly assess compliance with each item and mark each
33 element appropriately. Each item in the checklist must be checked, marked as N/A, or marked with an
34 asterisk indicating non-compliance that must be further justified. Explanations for N/A entries or any
35 other variances marked with an asterisk must be provided in the Author Comments section. Submission
36 of a completed checklist was a requirement for initiation of peer review for the JPR 2016 HPP Special
37 Issue.
38
39
40
41
42
43
44
45
46
47
48

49 Please note that it will be common for manuscripts to have at least one N/A entry. For example, for an
50 SRM-only dataset, element 12 will likely be N/A, while for a dataset that does not include SRM data
51 element 13 will be N/A. Element 9 will be N/A if there is only a single dataset analyzed. Although full
52 compliance for all applicable elements is generally expected for manuscripts, in rare cases it may be
53 appropriate to allow particular non-compliance. If authors feel that compliance for a particular element is
54
55
56
57
58
59
60

1
2
3 applicable but not achievable, the element may be asterisked and explained. Reviewers and editors may
4 then consider whether the particular exception request is reasonable or should not be accepted. For
5 example, a reasonable exception for element 2 would be a meta-analysis of 1,000 datasets that are all
6 already publicly accessible in some form, although potentially not found in ProteomeXchange
7 repositories. Element 15 already has a potential exception described; some proteins are so short or their
8 sequences such that only one unique (i.e., proteotypic) peptide may be possible, even when considering
9 multiple enzymatic digests.
10
11
12
13
14
15
16
17
18
19

20
21 **2. Deposit all MS proteomics data (DDA, DIA, SRM), including analysis reference files**
22 **(search database, spectral library), to a ProteomeXchange repository as a complete**
23 **submission. Provide the PXD identifier(s) in the manuscript abstract and reviewer login**
24 **credentials.**
25
26
27
28

29 The 2012 HPP Guidelines were the first to require submission of data through one of the
30 ProteomeXchange consortium repositories¹⁸. At that time, this was only PRIDE²⁰⁻²² for shotgun data, and
31 PASSEL²³, a part of PeptideAtlas²⁴⁻²⁶, for SRM data. Compliance was not universally enforced, but data
32 were deposited to ProteomeXchange for most HPP Special Issue articles through 2015. Since the initial
33 guidelines were put into place, the MassIVE and jPOST repositories have joined ProteomeXchange; as a
34 result, there are now four repositories in the consortium. The new iProX repository has expressed interest
35 in joining the ProteomeXchange Consortium. Access to the raw data is essential for the standardized
36 reanalyses by the field. Indeed, reanalysis of MS data by PeptideAtlas and by GPMDB has greatly
37 advanced data quality and comparability of analysis in the field of proteomics and provided insights into
38 the metrics underlying these guidelines.
39
40
41
42
43
44
45
46
47
48
49
50

51
52
53 There are broadly two kinds of submissions supported by ProteomeXchange repositories: “partial” (also
54 called “unsupported”) and “complete” (also called “supported”). While both require the same amount of
55 information (metadata, raw data, and identification results), the key difference is that for a partial
56
57
58
59
60

1
2
3 submission, the receiving repository was not able to parse and fully load all of the data. In a complete
4
5 submission, the identification results and identified spectra are fully loaded and searchable via the
6
7 repository interface. The reason for this distinction is that there are many available software pipelines,
8
9 many of which do not use standardized or otherwise common output formats. The repositories may not be
10
11 able to support the parsing and loading of all possible formats on account of the limited resources
12
13 available to the repositories. Although complete submissions are most desirable, the partial submission
14
15 mechanism is supported by the repositories so that everyone can submit their data and results to
16
17 ProteomeXchange, even if the formats were not fully supported.
18
19
20

21
22
23 With these latest 2015/2016 guidelines, past requirements have been upgraded to that of mandatory
24
25 complete submission. This means that results must be submitted in a format that can be parsed by the
26
27 receiving repository, and some software tools may have to be excluded because their output cannot (yet)
28
29 be written or converted to a supported format. Although such a requirement was considered too
30
31 demanding in 2012 because the PSI mzIdentML format²⁷ was not universally supported by the
32
33 repositories at that time, it is now the case that mzIdentML is well supported and widely used. Although
34
35 not unanimous, the consensus opinion was that complete submission has been widely achievable for some
36
37 time, and workflows or tools that do not yet produce a suitable output need incentive to support complete
38
39 submission to ProteomeXchange. For such software that still does not permit complete submissions, we
40
41 hope that this raising of the bar will accelerate progress in this area. Complete submission is now clearly
42
43 presented as our long term goal. The HPP decided that a requirement for complete submission would be
44
45 an important component of such an effort, which has had broad support at the HUPO2015 Congress in
46
47 Vancouver and throughout the HPP community. If only a minority of submitters are unable to submit data
48
49 in a complete form, this should put pressure on developers of the software they use to develop solutions to
50
51 this problem. The HPP special issue editors are willing to be generous in allowing exceptions to this
52
53 guideline in the near term as we seek solutions for full compliance.
54
55
56
57
58
59
60

1
2
3 There was considerable debate about the details of the guideline about mandatory complete data
4 submission relative to timing—that is, whether complete data deposition should be required prior to
5 manuscript submission or not. It was generally agreed that submission after acceptance of a manuscript
6 was too late, as this does not give reviewers the opportunity to verify that the submission is appropriate
7 and matches all claims and descriptions found in the submitted manuscript. However, some felt that
8 deposition of datasets prior to initial manuscript submission would place undue burden on repositories,
9 since they are already operating with limited (human) resources, in that they would be forced to handle
10 unnecessary data submissions for any manuscripts destined to be returned without peer review or rejected,
11 leading to pollution of the databases. In versions 2.0.0 through 2.0.5, the current guideline was written so
12 that authors had the option of first submitting their manuscript and waiting to deposit their data until the
13 editors have signaled that the manuscript will be sent for review upon data deposition. In response to the
14 reviewer comments and additional debate by the HPP leadership, this guideline was changed as of version
15 2.1.0 to reflect the requirement that all data must be deposited prior to submission of the manuscript. This
16 policy is deemed simpler to implement and explain and desirable to expedite review.
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35

36 **3. Use the most recent version of the neXtProt reference proteome for all informatics**
37 **analyses, particularly with respect to potential missing proteins.**

38
39
40 The official reference knowledge base for the HPP is neXtProt⁵, and it is crucial that claims of detection
41 of missing proteins and possible translation products from other novel coding elements be compared with
42 the most current neXtProt release, rather than any earlier version. In some cases, this will require re-
43 processing data with an updated reference database prior to final submission of a manuscript. Manuscripts
44 that specially claim detection of missing proteins in their abstract, followed by comments in the
45 discussion section that some are no longer missing in the latest neXtProt release, which was permitted in
46 2015, are no longer acceptable. For submission to one of the HPP special issues, the call for papers will
47 denote which version of neXtProt, PeptideAtlas and other resources are to be used.
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 **4. Describe in detail the calculation of FDRs at the PSM, peptide, and protein levels.**
4

5 Each manuscript must describe the methods that were used to calculate the false discovery rates at the
6 PSM, distinct peptide, and protein levels. Importantly, no specific method is prescribed. Some methods or
7 tools can calculate all three levels at once, while in some cases multiple tools must be used. Use and
8 citation of existing tools is encouraged. It is not sufficient to state only “FDRs were calculated using tool
9 X.” Software versions, input parameters, apparent anomalies, input file formats, and output formats must
10 all be specified. Any variances or modifications to a previously published methodology should be
11 described. If custom, novel, or unpublished methods are used, they should be described in detail. If such
12 novel or unpublished methods are used, then the results should be compared in some way with results
13 from a more conventional analysis that has been previously published. Note that any assumptions should
14 be clearly stated. Be specific about the distinction between a global FDR (the fraction of incorrect entities
15 among all entities that pass the threshold) and a local FDR (the fraction of incorrect entities within a
16 subset of entities that share the same score, usually expressed for each entity or for the threshold score in
17 a list). The calculation at the peptide level may differentiate between different mass modifications, or
18 aggregate over multiple modifications, at the discretion of the authors.
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37

38 **5. Report the PSM-, peptide-, and protein-level FDR values along with the total number of**
39 **expected true positives and false positives at each level.**
40

41 Based on the methodology described, report global FDR values at each of the three levels. Unless unusual
42 methodology is employed, the PSM-level FDR should be lower than the peptide-level FDR, which should
43 be lower than the protein-level FDR. The larger the dataset, the more extreme these differences become.
44 In addition to the FDRs, report the total number of entities passing threshold at each level, and then also
45 state the expected or estimated number of incorrect entities passing threshold at each level. This means
46 that, in addition to stating that proteins are thresholded at a 1% global FDR, state that, for example, 5,000
47 proteins pass the threshold and, therefore, there are an estimated 50 incorrect identifications in the list.
48 Some software packages do not report all this information, or even FDRs at each of the three levels.
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 However, even a simple strategy of counting all the PSMs, distinct peptides, and proteins that pass
4 threshold, counting the corresponding decoys at each level that pass threshold, and entering those values
5 into a spreadsheet to calculate the decoy rates and presumed corresponding false discovery rates would be
6 sufficient.
7
8
9
10

11 12 13 14 **6. Present large-scale results thresholded at equal to or lower than 1% protein-level** 15 **global FDR.** 16 17

18 Although there is not universal agreement on what the best threshold is, and it may vary based on the
19 intent of the final protein list, the HPP has concluded that the baseline acceptable global FDR for a dataset
20 should be at most 1% at the protein level. Lower than 1% is strongly encouraged. As described above,
21 this will usually mean that the peptide-level and PSM-level FDRs will be far lower than 1% for large
22 datasets. We note that, for some datasets, the local FDR should be the factor that should be used to set the
23 threshold. Consider the extreme case where all identifications can be perfectly discriminated into correct
24 and incorrect populations; in order to achieve a 1% global FDR, one is forced to add known incorrect
25 identifications (with local FDR of 100%), which is clearly not an acceptable strategy. For some very high
26 quality datasets where discrimination is excellent, it may be best to apply a local FDR threshold of 10%
27 (where 1 in every 10 identifications near the threshold are incorrect), even though this may yield a global
28 FDR far lower than 1%.
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44

45 **7. Recognize that the protein-level FDR is an estimate based on several imperfect** 46 **assumptions, and present the FDR with appropriate precision.** 47 48

49 There are many different approaches to estimating FDRs. The most common is the target-decoy approach,
50 followed by a population modeling approach²⁸⁻³³. Both approaches make imperfect assumptions that
51 affect the accuracy of the results. Decoys are not representative of all kinds of false positives. For
52 example, identifications may be very nearly correct, but incorrect in one or two residues³⁴, and there tend
53 to be rather few decoys at very stringent thresholds, leading to problems with small-number statistics.
54
55
56
57
58
59
60

1
2
3 Consider a hypothetical dataset with 1,010 proteins that pass threshold, ten of which are decoys; one
4 might discard these ten decoys and presume there are another ten incorrect identifications among the
5 remaining 1,000, leading to a 1% FDR. However, the exact scores and occurrence of decoys depends on
6 many details of the exact decoy database used, and there could easily have been 9 or 11 decoys at the
7 same effective threshold. Such a change in a single decoy would then yield a calculated FDR of 0.9% or
8 1.1% for 9 and 11, respectively. Clearly the precision with which the true uncertainty is known when such
9 few decoys are present cannot be high. Model-based approaches may often fit well to the main part of the
10 population, but may fit less well at the very tail of the distribution where the stringent threshold lies,
11 leading to similar uncertainties. In addition, the model high-confidence tail can vary substantially
12 depending on the mathematical function used for the model. In summary, FDR values in a manuscript
13 should be quoted with appropriate precision; unjustified precision, *i.e.* more than two digits of precision,
14 should be avoided.
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30

31 **8. Acknowledge that not all proteins surviving the threshold are “confidently identified.”**

32
33 It is important that careful FDR estimation is not left behind during subsequent analysis of the protein
34 results. It is inappropriate to proceed with an analysis that treats all remarkable entries (e.g., missing
35 proteins) in the resulting list as “confidently identified” when errors are known to exist in the list. In fact,
36 the total number of remarkable identifications should be compared to the reported number of false
37 positives (guideline 5). For example, if one expects 30 incorrect identifications in a result (such as 1% of
38 3000 proteins), then a claim of the detection of 10 missing proteins should be treated with great caution.
39 The default hypothesis should be that these never-before-detected proteins (in mass spectrometry) are 10
40 of the expected 30 false positives. Orthogonal convincing evidence must be presented to rule out (or at
41 least significantly constrain) this default hypothesis. See guidelines 10, 11, and 14.
42
43
44
45
46
47
48
49
50
51
52
53
54

55 **9. If any large-scale datasets are individually thresholded and then combined, calculate** 56 **the new, higher peptide- and protein-level FDRs for the combined result.** 57 58 59 60

1
2
3 When several different proteomic (or MS) datasets are compared or combined in a manuscript, it is
4 important to be mindful that the combined results will have a different, usually higher FDR. Consider the
5 three cases in Figure 1. For example, in A, where there is no overlap in the correct proteins and no
6 overlap in the incorrect identifications, the combined FDR is truly the same as in the original datasets. In
7 case B, all of the correct identifications overlap, but the incorrect ones do not (because incorrect
8 identifications usually scatter over the proteome). The combined FDR is twice as high as the original. The
9 third case is a more real-world example where 50% of the correct identifications overlap, and none of the
10 incorrect ones does. The resulting FDR is ~1.5%, which is much larger than the original FDRs. Caution is
11 required with compendia of many experiments that have all been individually processed, thresholded, and
12 then combined, as the false discovery rates will inflate considerably. If all of the data discussed in a
13 manuscript are processed together with a single threshold, then this guideline will not be applicable.
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28

29 **10. Present “extraordinary detection claims” based on DDA mass spectrometry with high**
30 **mass-accuracy, high signal-to-noise ratio (SNR), and clearly annotated spectra.**
31

32 The concept of an “extraordinary detection claim” is purposely left somewhat vague in the guidelines.
33 Two obvious examples in this category are missing proteins (predicted proteins lacking PE=1 neXtProt
34 status) and novel coding elements (e.g., lncRNAs, novel exons, pseudogenes, or other sequences not
35 listed in neXtProt as entries with protein existence level 1 through 4). However, authors and reviewers
36 may consider other claims as extraordinary, such as a report of detection of a protein in a sample where
37 the protein would not likely be present and the transcript cannot be detected, such as an olfactory receptor
38 protein in liver.
39
40
41
42
43
44
45
46
47
48
49
50

51 Several journal guidelines already require annotated tandem mass spectra as supplementary material for
52 single-hit proteins. We have extended this requirement to all extraordinary detection claims, even when
53 supported by multiple peptides. Furthermore, the spectra must be of high signal-to-noise ratio, with a
54 recommendation that the highest 5% intensity peaks should have a signal at least 20 times those of the
55
56
57
58
59
60

1
2
3 lowest 5% intensity peaks, which are presumed to be mostly noise. Although low mass-accuracy (i.e., ion
4 trap) MS/MS spectra are still useful for many applications, MS/MS spectra supporting extraordinary
5 detection claims should be acquired in higher mass-accuracy (Fourier-transform, Orbitrap, TOF, Q-
6 Exactive, etc.) instruments.
7
8
9
10

11
12
13
14 **11. Consider alternative explanations of PSMs that appear to indicate extraordinary**
15 **results.**
16

17
18 In cases where a peptide identification corresponding to an extraordinary claim appears to have a well
19 annotated, high signal-to-noise ratio spectrum, consider whether a slightly different amino acid sequence
20 that can map to a different, common protein also could be a credible explanation. An example is the case
21 presented in Figure 5 of Deutsch *et al.*²⁶ of a spectrum in PeptideAtlas that appears to have excellent
22 coverage and, thus, a very high score for olfactory receptor 5A2 (Q8NGI9), with just a few missing and
23 unexplained peaks. However, careful scrutiny reveals a slightly different peptide sequence with an
24 unconsidered mass modification that yields an even better match, and also maps to a very commonly seen
25 protein (and peptide sequence), lactotransferrin (P02788). Such cases may be quite rare, but, among
26 millions of mass spectra, some of the ones that appear to implicate extraordinary results will be cases such
27 as this. This guideline does not require manual inspection of all spectra; rather it applies only to the
28 exceptional case of an extraordinary claim for a previously unreported protein match.
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43

44 **12. Present high mass-accuracy, high-SNR, clearly annotated spectra of synthetic**
45 **peptides that match the spectra supporting the extraordinary detection claims.**
46

47
48 One method for increasing the confidence in the correctness of an identified peptide is to compare the
49 identification with a synthetic version of the peptide (i.e. same charge, same mass modifications, same
50 instrument fragmentation). The synthetic peptide fragment spectra should be shown alongside the
51 naturally-derived peptides, both with high spectrum quality and with similar peak intensity patterns
52 between the natural peptide and the synthetic peptide. A match in chromatographic elution time also is a
53
54
55
56
57
58
59
60

1
2
3 strong confirmation, but not sufficient, that the peptide is correctly identified. As in the JPR C-HPP 2015
4 special issue, the editors may allow stepwise presentation of “candidate missing protein identifications”,
5 followed by an explanation of how the candidate fared upon application of these more stringent
6 requirements. Such information may be a guide for others to seek more convincing evidence in the same
7 type of specimen or in another specimen guided by transcript expression data.
8
9
10
11
12

13
14
15
16
17 **13. If SRM verification for extraordinary detection claims is performed, present target**
18 **traces alongside synthetic heavy-labeled peptide traces, demonstrating co-elution and**
19 **very closely matching fragment mass intensity patterns.**
20
21

22 SRM can be a useful technology to confirm the unambiguous identification of peptides that appear to
23 support the extraordinary claim. Although its sensitivity can be better than conventional shotgun
24 technologies, it is not vastly better and, since fewer ions are often used as evidence, it is imperative that
25 SRM confirmation is performed with the use of spiked-in stable isotope-labeled synthetic peptides.
26 Maximal corresponding fragments (transitions) must be monitored for both heavy and light ions and of
27 predominantly higher mass transitions for better discrimination. The peak intensity order of those ions as
28 well as elution pattern must match with high similarity. Traces down at the detection limit are usually not
29 suitable, as the chance of spurious interferences is high at the detection limit. Furthermore, it is crucial to
30 exclude the possibility of light-peptide contamination in heavy-labeled spike-ins providing spurious
31 signal. For example, if a spike-in reference sample of heavy-labeled peptides contains 1% light peptide
32 contamination, then all samples analyzed with that reference will exhibit a false detection if the heavy-
33 labeled peptide signal is more than 100 times the level of detection. This should be prevented by spiking
34 in heavy-labeled peptides at a comparable abundance as the target peptide, or demonstrating that the
35 heavy-labeled reference has contamination much lower than a level at which putative target signals are
36 detected.
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 **14. Even when very high confidence peptide identifications are demonstrated, consider**
4 **alternate mappings of the peptides to proteins other than the claimed extraordinary**
5 **result. Consider isobaric sequence/mass modification variants, all known SAAVs, and**
6 **unreported SAAVs.**
7
8
9

10
11 Most of the earlier proteomic guidelines have been concerned with ensuring that peptide identifications
12 are of high quality. But even with nearly irrefutable evidence that a peptide identification is correct, the
13 peptide to protein mapping must also be considered very carefully. Clearly peptides that also map to a
14 common, well-observed protein cannot be held up as evidence in support of an extraordinary detection
15 claim, as the most likely explanation is that the peptide is derived from the common protein. Common
16 laboratory contaminant protein sequences should always be considered (e.g., the GPM distributes the very
17 comprehensive “cRAP” or “common Repository of Adventitious Proteins” set at
18 <http://www.thegpm.org/crap/>). Direct mapping is easy to determine, but it is also necessary to consider
19 alternative splice isoforms and single amino acid variants (SAAVs) in the mapping, as well. Substitutions
20 of I/L must be accounted for as these are isobaric and cannot be distinguished by current MS/MS
21 techniques used in mass spectrometers unless additional fragmentation routines are used^{35,36}. There are
22 other isobaric substitutions when one considers mass modifications. For example, deamidated N is
23 equivalent to D, and deamidated Q is equivalent to E. Note that there are many more substitutions that are
24 close but not exact, such as Q/K, that must be considered when analyzing low mass-accuracy spectra.
25 Low mass-accuracy data cannot easily distinguish between Q and K, F and oxidized M, and similar pairs,
26 which is another reason that guideline 11 excludes the use of low mass-accuracy ion trap spectra for
27 confirming evidence of extraordinary detection claims. As well, there is always the possibility that a
28 known or unknown PTM not taken into account during the search could lead alone or in combination with
29 misidentification to an incorrect match. A tool to assist with this analysis is available at neXtProt at
30 <https://search.nextprot.org/view/unicity-checker> and can be used to aid in compliance with this
31 guideline. For example, in PeptideAtlas peptide SITDVLSADDIAAALQECQDPDTFEPQK
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 appears to uniquely map to PE=5 protein Putative oncomodulin-2 (P0CE71) amongst the core
4
5 20,000 neXtProt predicted proteins. However, when SAAVs are considered, one finds that it also
6
7 maps to a known variant (dbSNP rs202012112) of PE=1 protein Oncomodulin-1 (P0CE72). This
8
9 peptide is therefore no longer uniquely mapping, and cannot be held up as protein evidence for
10
11 the existence of PE=5 protein P0CE71.
12
13

14
15
16
17 **15. Support extraordinary detection claims by two or more distinct uniquely-mapping,**
18
19 **non-nested peptide sequences of length ≥ 9 amino acids. When weaker evidence is**
20
21 **offered for detection of a previously unreported protein or a coding element proposed**
22
23 **translation product, justify that other peptides cannot be expected.**
24

25
26 As outlined above, it is clear that an apparently very high quality uniquely-mapping peptide identification
27
28 can still be incorrect as a protein match. In fact, in very large datasets using the thresholds advocated in
29
30 these guidelines, there will surely still be a few such cases. Therefore, in order to engender additional
31
32 confidence in extraordinary detection claims, we require the evidence of two distinct peptides of length 9
33
34 amino acids or more. Further, one of the peptides may not be fully nested within the other. Nested
35
36 peptides are not counted, because, while they increase the confidence of the sequence being accurately
37
38 identified, especially in the case of ragged peptides from termini, it does not generate additional
39
40 confidence in the uniqueness of the peptide-to-protein mapping.
41
42
43
44

45
46 Very short peptides usually map to many different proteins, and there are abundant examples in
47
48 PeptideAtlas where apparently “uniquely mapping” peptides can be better explained by mappings to
49
50 variants or nearly identical isobaric peptides for other proteins. This problem is so rampant with peptides
51
52 of length 6 or less, that they have long been completely discarded from PeptideAtlas and never shown. In
53
54 PeptideAtlas peptides of length 7 are retained and shown, but there are many cases where one cannot feel
55
56 confident that such short peptides are truly indicative of a protein detection alone. As a cautionary note,
57
58
59
60

1
2
3 there are also such cases for peptides of length 8, and we have therefore conservatively set a lower limit
4 of 9 amino acids for peptides that are needed to confer the canonical designation in PeptideAtlas and the
5 protein existence level 1 in neXtProt. It is useful to extend this same requirement for evidence of
6 extraordinary detections. If it is desirable to present evidence that does not meet these criteria (covered in
7 next paragraph), the implicated proteins may be offered as “candidate detections” to enable capture of this
8 information by other researchers and use in potential future experiments.
9
10
11
12
13
14
15
16
17

18 In some rare cases there are proteins that simply do not contain enough uniquely mapping peptides of
19 sufficient length to call a protein detected. For example, proteins with very few or excessive basic
20 residues produce only a few extremely long peptides, if any, on the one hand, and produce many
21 excessively short peptides on the other hand when trypsin is used as the cleavage reagent. The use of
22 other enzymes, e.g., GluC or chymotrypsin, or chemical cleavage reagents may provide additional
23 opportunities to detect a protein by generating different repertoires of peptides. It is still permissible to
24 present evidence that does not fully meet this guideline if there is a strong justification that additional
25 peptide evidence will be extraordinarily difficult to achieve. For example, if a single peptide or short
26 peptides are all that can be reasonably expected for a missing protein, even with the use of multiple
27 proteases, based on its sequence, and these are precisely the peptides that are observed, the community
28 and the neXtProt curators may be convinced to relax this guideline in such special cases.
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44

45 Discussion

46
47 Although there are already several sets of guidelines relevant to proteomics, there is minimal content
48 overlap. Most of the other guidelines focus on specific basic metadata that must be provided, while the
49 HPP guidelines focus mostly on addressing the need to provide well accepted evidence for analysis
50 quality and for new identifications whilst reducing the probability of false positive identifications that
51 seem to implicate proteins that were never before seen and are thus highly sought after as “missing
52
53
54
55
56
57
58
59
60

1
2
3 proteins". In this sense, these guidelines complement other data inclusion guidelines that may also apply.
4
5 For example, for manuscripts submitted to JPR, the journal data submission guidelines also apply. There
6
7 is minimal overlap between these two sets of guidelines; where these do overlap, complying with the HPP
8
9 guidelines should include compliance with the JPR guidelines.
10
11

12
13
14 The overarching theme for the HPP guidelines is that careful control of false positives is crucial for
15
16 unambiguous protein identification when the goal of the work is to present claims of comprehensive
17
18 datasets of increasingly nearly complete proteomes and for the inclusion of never-before-confidently-
19
20 detected proteins. Unless the number of errors present in a final protein list is much smaller than the
21
22 number of claimed novel discoveries, confirmatory orthogonal evidence must be presented to demonstrate
23
24 that the novel claims are not merely one of the false positives.
25
26

27
28
29 There is one additional guideline that was considered and not included in the current release. There was
30
31 consideration for a requirement for a confirming detection in a second sample. At present, two peptides
32
33 from a single sample is all that is required. Majority consensus was that requiring detection of a protein
34
35 from at least two separate samples (i.e. biological replicates rather than technical replicates) was raising
36
37 the bar too high, and this guideline was not included. It may be considered for future guidelines upon
38
39 consultation with the proteomics community.
40
41
42

43
44 Another situation may arise that provides evidence for a missing protein, but which does not meet the
45
46 guidelines to its positive identification. For example, in PTM peptide enrichment studies, e.g., glyco and
47
48 phospho proteomics, and N and C terminomics³⁸, a single peptide from a protein is often identified with
49
50 high confidence. The point of the study may be to characterize the PTMs rather than to identify proteins,
51
52 but such studies also provide an orthogonal approach to provide proteomic evidence for proteins,
53
54 especially useful for missing proteins. For such studies with high quality spectra and peptides that meet
55
56 the spectral assignment and peptide identification guidelines otherwise, these peptides can be designated
57
58
59
60

1
2
3 as potentially having come from the missing protein. In any case specific caveats need to be stated, e.g.,
4 that peptide evidence was found for a missing protein or that “candidate missing proteins” were detected
5 by these high confidence peptides, but that further evidence is required for high confidence identification.
6
7 The hope is that these identifications can stimulate other groups to specifically seek further evidence of
8
9 the missing protein in that tissue, for example, or by using such approaches incorporated into broader
10
11 studies to identify recalcitrant missing proteins.
12
13
14
15
16
17

18 Although there has been extensive discussion and refinement of the guidelines, the first real test of the
19
20 guidelines has been this JPR 2016 HPP special issue. All submitted manuscripts were required to comply
21
22 with these guidelines. Completed checklists were submitted with the manuscripts. The special issue
23
24 editors agreed to perform a first pass of compliance checking before the manuscripts were sent out for
25
26 review. Reviewers were then asked to consider the guidelines as they review the manuscripts. Authors
27
28 generally complied with the guidelines, either upon submission or, in multiple cases, during revision. We
29
30 anticipate that JPR will consider adopting these guidelines for papers claiming identification of Missing
31
32 Proteins in regular journal issues. We encourage all journals, whether inside or outside the field of
33
34 proteomics, to consider and adopt these guidelines.
35
36
37
38
39

40 There is potential opportunity for integration with other guidelines, but this task will need effort from the
41
42 respective stakeholders. Many of the guidelines, including these, are directed to a specific purpose and
43
44 may not apply well in other experimental designs. The reasonable desire to have a single set of guidelines
45
46 might only result in a large and unwieldy document with many “if-then” sections for different strategies.
47
48 Despite these considerations, these HPP guidelines break new ground regarding the somewhat narrow
49
50 focus about claims of novel protein detection, and many of these individual guidelines may be suitable for
51
52 inclusion in more general fit-for-purpose guidelines.
53
54
55
56
57
58
59
60

Conclusion

We have presented the latest version (2.1.0) of the HPP MS Data Interpretation Guidelines. These guidelines expand substantially on the version 1.0 guidelines, which only required any kind of ProteomeXchange deposition and a 1% protein-level FDR threshold. For manuscript submission to the Journal of Proteome Research the primary guidelines comprise a one-page checklist followed by two pages of extended information. This article provides an in-depth history, reasoning, and expanded discussion of each of the guidelines so that the community may fully understand their intent and consider whether broader application to other projects is appropriate. The previous 2012 guidelines served the HPP well for three years. These guidelines will be further refined and expanded by the HPP as the field advances.

Acknowledgements

This work was funded in part by the National Institutes of Health through NIGMS grant R01GM087221 and NIBIB grant U54EB020406 (EWD) and NIEHS grant U54ES017885 (GSO). The authors have no conflicts of interest to declare.

Supporting Information

Supporting Information: Checklist document with extended description of the guidelines. Supplementary Document S1.

References

- (1) Legrain, P.; Aebersold, R.; Archakov, A.; Bairoch, A.; Bala, K.; Beretta, L.; Bergeron, J.; Borchers, C. H.; Corthals, G. L.; Costello, C. E.; et al. The human proteome project: current state and future direction. *Mol. Cell Proteomics* **2011**, *10* (7), M111.009993.

- 1
- 2
- 3
- 4 (2) Paik, Y.-K.; Jeong, S.-K.; Omenn, G. S.; Uhlen, M.; Hanash, S.; Cho, S. Y.; Lee, H.-J.;
5 Na, K.; Choi, E.-Y.; Yan, F.; et al. The Chromosome-Centric Human Proteome Project for
6 cataloging proteins encoded in the genome. *Nat. Biotechnol.* **2012**, *30* (3), 221–223.
- 7 (3) Paik, Y.-K.; Omenn, G. S.; Overall, C. M.; Deutsch, E. W.; Hancock, W. S. Recent
8 Advances in the Chromosome-Centric Human Proteome Project: Missing Proteins in the
9 Spot Light. *J. Proteome Res.* **2015**, *14* (9), 3409–3414.
- 10 (4) Horvatovich, P.; Lundberg, E. K.; Chen, Y.-J.; Sung, T.-Y.; He, F.; Nice, E. C.; Goode, R.
11 J.; Yu, S.; Ranganathan, S.; Baker, M. S.; et al. Quest for Missing Proteins: Update 2015
12 on Chromosome-Centric Human Proteome Project. *J. Proteome Res.* **2015**, *14* (9), 3415–
13 3431.
- 14 (5) Aebersold, R.; Mann, M. Mass spectrometry-based proteomics. *Nature* **2003**, *422* (6928),
15 198–207.
- 16 (6) Nilsson, T.; Mann, M.; Aebersold, R.; Yates, J. R.; Bairoch, A.; Bergeron, J. J. M. Mass
17 spectrometry in high-throughput proteomics: ready for the big time. *Nat. Methods* **2010**, *7*
18 (9), 681–685.
- 19 (7) Gillet, L. C.; Navarro, P.; Tate, S.; Röst, H.; Selevsek, N.; Reiter, L.; Bonner, R.;
20 Aebersold, R. Targeted data extraction of the MS/MS spectra generated by data-
21 independent acquisition: a new concept for consistent and accurate proteome analysis.
22 *Mol. Cell Proteomics* **2012**, *11* (6), O111.016717.
- 23 (8) Picotti, P.; Aebersold, R. Selected reaction monitoring-based proteomics: workflows,
24 potential, pitfalls and future directions. *Nat. Methods* **2012**, *9* (6), 555–566.
- 25 (9) Deutsch, E. W.; Lam, H.; Aebersold, R. Data analysis and bioinformatics tools for tandem
26 mass spectrometry in proteomics. *Physiol. Genomics* **2008**, *33* (1), 18–25.
- 27 (10) Nesvizhskii, A. I. A survey of computational methods and error rate estimation procedures
28 for peptide and protein identification in shotgun proteomics. *J Proteomics* **2010**, *73* (11),
29 2092–2123.
- 30 (11) Taylor, C. F.; Paton, N. W.; Lilley, K. S.; Binz, P.-A.; Julian, R. K.; Jones, A. R.; Zhu, W.;
31 Apweiler, R.; Aebersold, R.; Deutsch, E. W.; et al. The minimum information about a
32 proteomics experiment (MIAPE). *Nat. Biotechnol.* **2007**, *25* (8), 887–893.
- 33 (12) Deutsch, E. W.; Albar, J. P.; Binz, P.-A.; Eisenacher, M.; Jones, A. R.; Mayer, G.; Omenn,
34 G. S.; Orchard, S.; Vizcaíno, J. A.; Hermjakob, H. Development of data representation
35 standards by the human proteome organization proteomics standards initiative. *J Am Med*
36 *Inform Assoc* **2015**, *22* (3), 495–506.
- 37 (13) Bradshaw, R. A.; Burlingame, A. L.; Carr, S.; Aebersold, R. Reporting protein
38 identification data: the next generation of guidelines. *Mol. Cell Proteomics* **2006**, *5* (5),
39 787–788.
- 40 (14) Burlingame, A.; Carr, S. A.; Bradshaw, R. A.; Chalkley, R. J. On Credibility, Clarity, and
41 Compliance. *Mol. Cell Proteomics* **2015**, *14* (7), 1731–1733.
- 42 (15) Carr, S. A.; Abbatiello, S. E.; Ackermann, B. L.; Borchers, C.; Domon, B.; Deutsch, E.
43 W.; Grant, R. P.; Hoofnagle, A. N.; Hüttenhain, R.; Koomen, J. M.; et al. Targeted peptide
44 measurements in biology and medicine: best practices for mass spectrometry-based assay
45 development using a fit-for-purpose approach. *Mol. Cell Proteomics* **2014**, *13* (3), 907–
46 917.
- 47 (16) Nesvizhskii, A. I. Proteogenomics: concepts, applications and computational strategies.
48 *Nat. Methods* **2014**, *11* (11), 1114–1125.
- 49
- 50
- 51
- 52
- 53
- 54
- 55
- 56
- 57
- 58
- 59
- 60

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
- (17) Ruggles, K. V.; Tang, Z.; Wang, X.; Grover, H.; Askenazi, M.; Teubl, J.; Cao, S.; McLellan, M. D.; Clauser, K. R.; Tabb, D. L.; et al. An Analysis of the Sensitivity of Proteogenomic Mapping of Somatic Mutations and Novel Splicing Events in Cancer. *Mol. Cell Proteomics* **2016**, *15* (3), 1060–1071.
 - (18) Vizcaíno, J. A.; Deutsch, E. W.; Wang, R.; Csordas, A.; Reisinger, F.; Ríos, D.; Dianes, J. A.; Sun, Z.; Farrah, T.; Bandeira, N.; et al. ProteomeXchange provides globally coordinated proteomics data submission and dissemination. *Nat. Biotechnol.* **2014**, *32* (3), 223–226.
 - (19) Omenn, G. S.; Lane, L.; Lundberg, E. K.; Beavis, R. C.; Nesvizhskii, A. I.; Deutsch, E. W. Metrics for the Human Proteome Project 2015: Progress on the Human Proteome and Guidelines for High-Confidence Protein Identification. *J. Proteome Res.* **2015**, *14* (9), 3452–3460.
 - (20) Martens, L.; Hermjakob, H.; Jones, P.; Adamski, M.; Taylor, C.; States, D.; Gevaert, K.; Vandekerckhove, J.; Apweiler, R. PRIDE: the proteomics identifications database. *Proteomics* **2005**, *5* (13), 3537–3545.
 - (21) Ternent, T.; Csordas, A.; Qi, D.; Gómez-Baena, G.; Beynon, R. J.; Jones, A. R.; Hermjakob, H.; Vizcaíno, J. A. How to submit MS proteomics data to ProteomeXchange via the PRIDE database. *Proteomics* **2014**, *14* (20), 2233–2241.
 - (22) Vizcaíno, J. A.; Csordas, A.; Del-Toro, N.; Dianes, J. A.; Griss, J.; Lavidas, I.; Mayer, G.; Perez-Riverol, Y.; Reisinger, F.; Ternent, T.; et al. 2016 update of the PRIDE database and its related tools. *Nucleic Acids Res.* **2016**, *44* (D1), D447–D456.
 - (23) Farrah, T.; Deutsch, E. W.; Kreisberg, R.; Sun, Z.; Campbell, D. S.; Mendoza, L.; Kusebauch, U.; Brusniak, M.-Y.; Hüttenhain, R.; Schiess, R.; et al. PASSEL: the PeptideAtlas SRMexperiment library. *Proteomics* **2012**, *12* (8), 1170–1175.
 - (24) Desiere, F.; Deutsch, E. W.; Nesvizhskii, A. I.; Mallick, P.; King, N. L.; Eng, J. K.; Aderem, A.; Boyle, R.; Brunner, E.; Donohoe, S.; et al. Integration with the human genome of peptide sequences obtained by high-throughput mass spectrometry. *Genome Biol.* **2005**, *6* (1), R9.
 - (25) Desiere, F.; Deutsch, E. W.; King, N. L.; Nesvizhskii, A. I.; Mallick, P.; Eng, J.; Chen, S.; Eddes, J.; Loevenich, S. N.; Aebersold, R. The PeptideAtlas project. *Nucleic Acids Res.* **2006**, *34* (Database issue), D655–D658.
 - (26) Deutsch, E. W.; Sun, Z.; Campbell, D.; Kusebauch, U.; Chu, C. S.; Mendoza, L.; Shteynberg, D.; Omenn, G. S.; Moritz, R. L. State of the Human Proteome in 2014/2015 As Viewed through PeptideAtlas: Enhancing Accuracy and Coverage through the AtlasProphet. *J. Proteome Res.* **2015**, *14* (9), 3461–3473.
 - (27) Jones, A. R.; Eisenacher, M.; Mayer, G.; Kohlbacher, O.; Siepen, J.; Hubbard, S. J.; Selley, J. N.; Searle, B. C.; Shofstahl, J.; Seymour, S. L.; et al. The mzIdentML data standard for mass spectrometry-based proteomics results. *Mol. Cell Proteomics* **2012**, *11* (7), M111.014381.
 - (28) Gaudet, P.; Michel, P.-A.; Zahn-Zabal, M.; Cusin, I.; Duek, P. D.; Evalet, O.; Gateau, A.; Gleizes, A.; Pereira, M.; Teixeira, D.; et al. The neXtProt knowledgebase on human proteins: current status. *Nucleic Acids Res.* **2015**, *43* (Database issue), D764–D770.
 - (29) Elias, J. E.; Gygi, S. P. Target-decoy search strategy for increased confidence in large-scale protein identifications by mass spectrometry. *Nat. Methods* **2007**, *4* (3), 207–214.
 - (30) Elias, J. E.; Gygi, S. P. Target-decoy search strategy for mass spectrometry-based proteomics. *Methods Mol. Biol.* **2010**, *604*, 55–71.

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
- (31) Choi, H.; Nesvizhskii, A. I. False discovery rates and related statistical concepts in mass spectrometry-based proteomics. *J. Proteome Res.* **2008**, *7* (1), 47–50.
- (32) Savitski, M. M.; Wilhelm, M.; Hahne, H.; Kuster, B.; Bantscheff, M. A Scalable Approach for Protein False Discovery Rate Estimation in Large Proteomic Data Sets. *Mol. Cell Proteomics* **2015**, *14* (9), 2394–2404.
- (33) Keller, A.; Nesvizhskii, A. I.; Kolker, E.; Aebersold, R. Empirical statistical model to estimate the accuracy of peptide identifications made by MS/MS and database search. *Anal. Chem.* **2002**, *74* (20), 5383–5392.
- (34) Nesvizhskii, A. I.; Keller, A.; Kolker, E.; Aebersold, R. A statistical model for identifying proteins by tandem mass spectrometry. *Anal. Chem.* **2003**, *75* (17), 4646–4658.
- (35) Colaert, N.; Degroeve, S.; Helsens, K.; Martens, L. Analysis of the resolution limitations of peptide identification algorithms. *J. Proteome Res.* **2011**, *10* (12), 5555–5561.
- (36) Armirotti, A.; Millo, E.; Damonte, G. How to discriminate between leucine and isoleucine by low energy ESI-TRAP MSⁿ. *J. Am. Soc. Mass Spectrom.* **2007**, *18* (1), 57–63.
- (37) Lebedev, A. T.; Damoc, E.; Makarov, A. A.; Samgina, T. Y. Discrimination of leucine and isoleucine in peptides sequencing with Orbitrap Fusion mass spectrometer. *Anal. Chem.* **2014**, *86* (14), 7017–7022.
- (38) Huesgen, P. F.; Lange, P. F.; Rogers, L. D.; Solis, N.; Eckhard, U.; Kleifeld, O.; Goulas, T.; Gomis-Rüth, F. X.; Overall, C. M. LysargiNase mirrors trypsin for protein C-terminal and methylation-site identification. *Nat. Methods* **2015**, *12* (1), 55–58.
- (39) Marino, G.; Eckhard, U.; Overall, C. M. Protein Termini and Their Modifications Revealed by Positional Proteomics. *ACS Chem. Biol.* **2015**, *10* (8), 1754–1764.

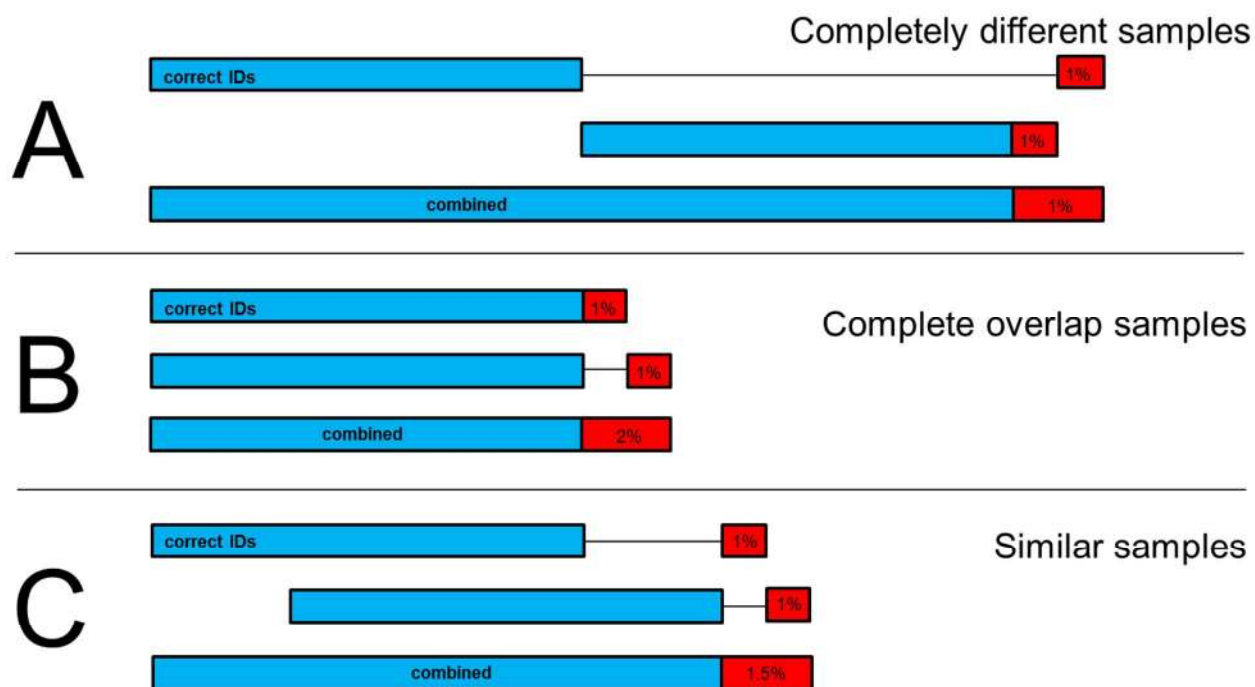
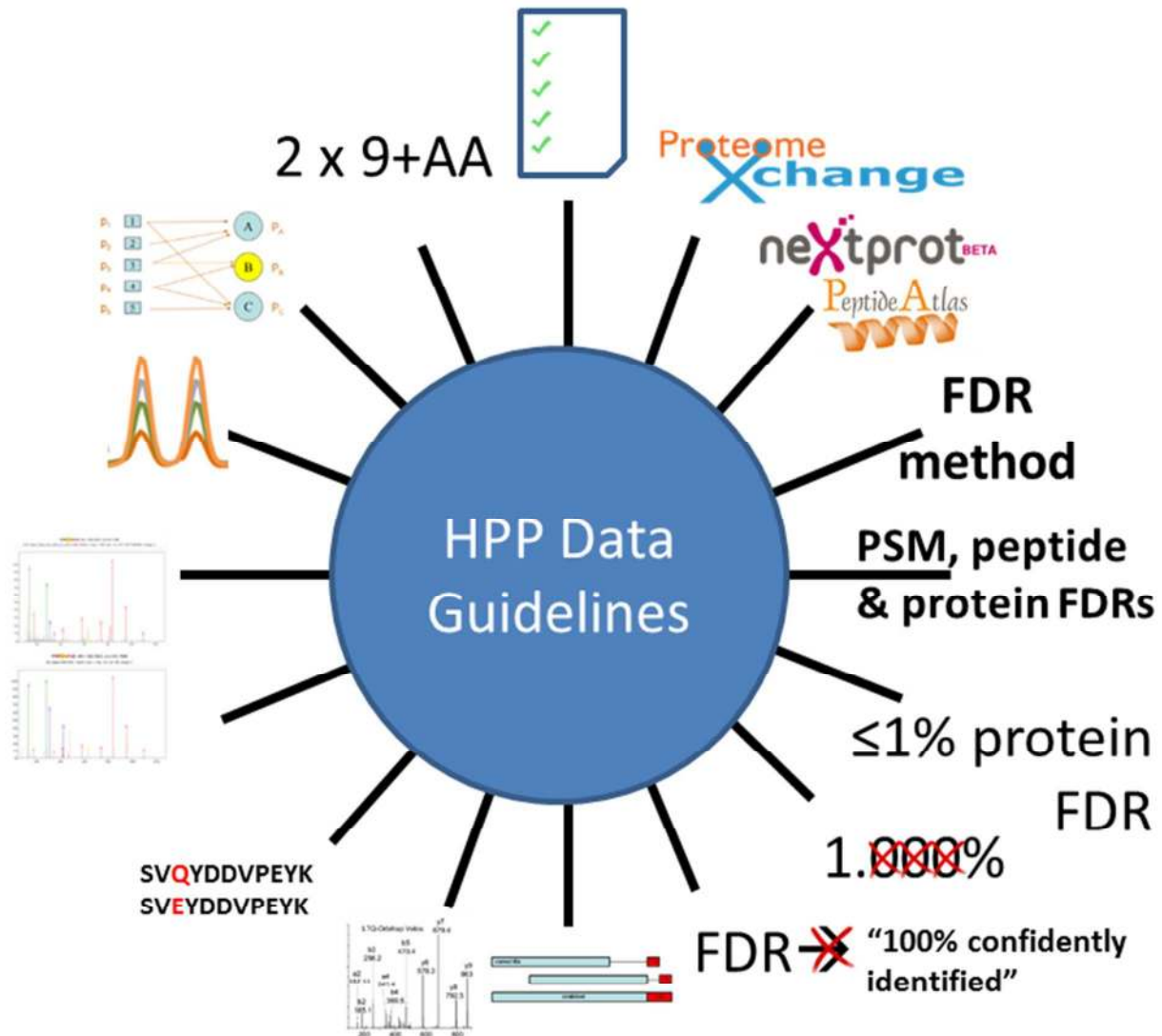


Figure 1. Three scenarios for how false discovery rates combine. True positives are shown in blue, and false positives in red. The red false positive boxes are depicted approximately 10 times larger than they should be for enhanced readability. A) If none of the true positives and 1% false positives overlap, then the final FDR does not expand. B) If all of the true positives overlap, but none of the false positives overlap (because they are false and random), then the final FDR is double the original rates. C) In a real-world scenario where the intersection of true positives overlaps by 50% and the false positives do not overlap, the combined FDR is 1.5%. This effect compounds as more datasets are merged.

Table 1. Checklist of the HUPO-2015 Human Proteome Project Data Interpretation Guidelines version 2.1.0.

General Guidelines:	
√	1. Complete this HPP Data Interpretation Guidelines checklist and submit with your manuscript.
	2. Deposit all MS proteomics data (DDA, DIA, SRM), including analysis reference files (search database, spectral library), to a ProteomeXchange repository as a complete submission. Provide the PXD identifier(s) in the manuscript abstract and reviewer login credentials.
	3. Use the most recent version of the neXtProt reference proteome for all informatics analyses, particularly with respect to potential missing proteins.
	4. Describe in detail the calculation of FDRs at the PSM, peptide, and protein levels.
	5. Report the PSM-, peptide-, and protein-level FDR values along with the total number of expected true positives and false positives at each level.
	6. Present large-scale results thresholded at equal to or lower than 1% protein-level global FDR.
	7. Recognize that the protein-level FDR is an estimate based on several imperfect assumptions, and present the FDR with appropriate precision.
	8. Acknowledge that not all proteins surviving the threshold are “confidently identified”.
	9. If any large-scale datasets are individually thresholded and then combined, calculate the new, higher peptide- and protein-level FDRs for the combined result.
Guidelines for extraordinary detection claims (e.g., missing proteins, novel coding elements)	
	10. Present “extraordinary detection claims” based on DDA mass spectrometry with high mass-accuracy, high signal-to-noise ratio (SNR), and clearly annotated spectra.
	11. Consider alternate explanations of PSMs that appear to indicate extraordinary results.
	12. Present high mass-accuracy, high-SNR, clearly annotated spectra of synthetic peptides that match the spectra supporting the extraordinary detection claims.
	13. If SRM verification for extraordinary detection claims is performed, present target traces alongside synthetic heavy-labeled peptide traces, demonstrating co-elution and very closely matching fragment mass intensity patterns.
	14. Even when very high confidence peptide identifications are demonstrated, consider alternate mappings of the peptides to proteins other than the claimed extraordinary result. Consider isobaric sequence/mass modification variants, all known SAAVs, and unreported SAAVs.
	15. Support extraordinary detection claims by two or more distinct uniquely-mapping, non-nested peptide sequences of length ≥ 9 amino acids. When weaker evidence is offered for a previously unreported protein or a coding element proposed translation product, justify that other peptides cannot be expected.



TOC Figure