



Japan Proteome Standard
Repository/Database

Proteomic Data Integration and Sharing by jPOST Repository/Database

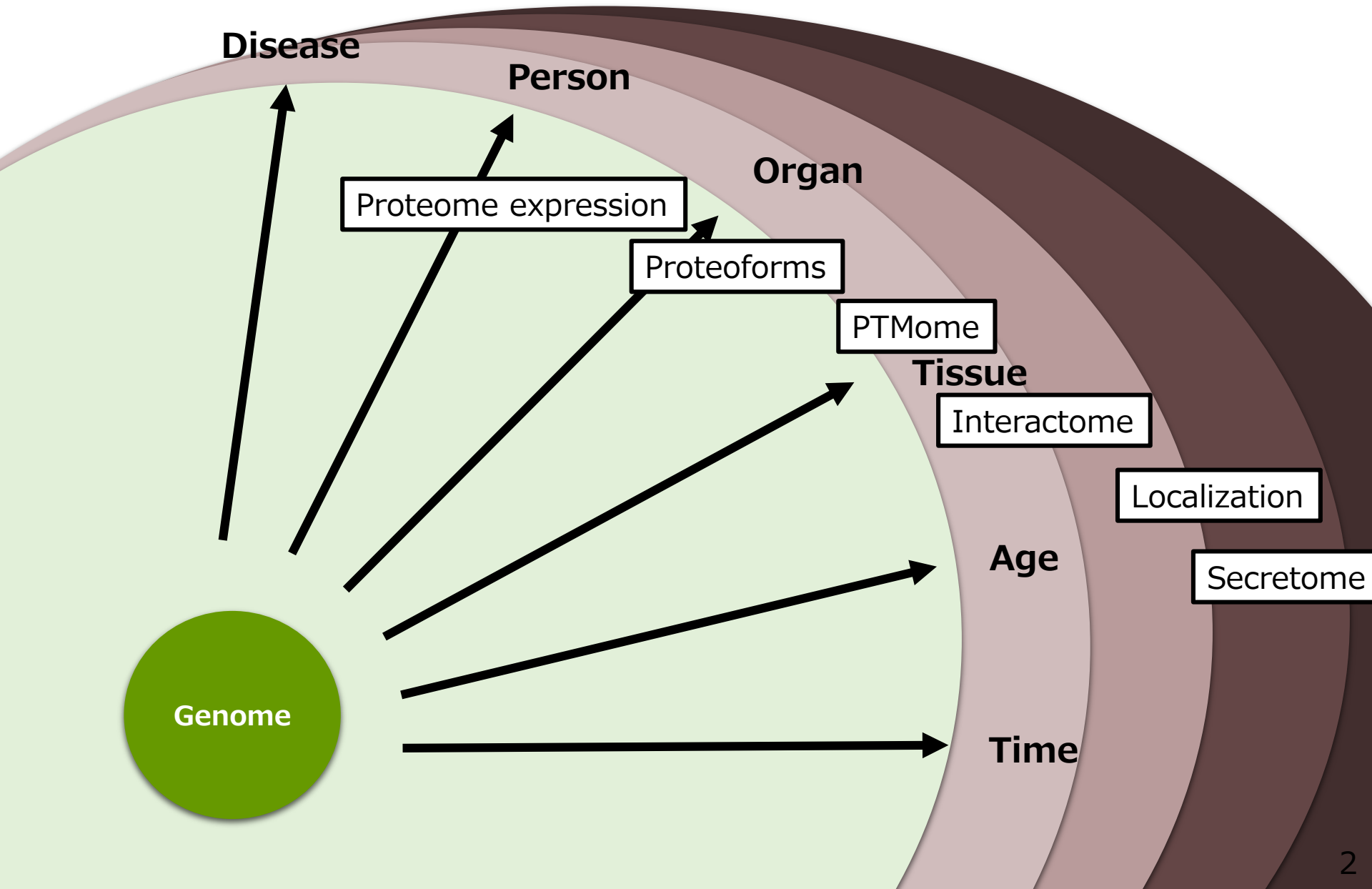
Yasushi Ishihama (Kyoto Univ) and jPOST project team

Y. Moriya¹, S. Kawano¹, S. Okuda², Y. Watanabe²,
M. Matsumoto³, T. Takami³, D. Kobayashi⁴,
N. Araki⁴, AC. Yoshizawa⁵, T. Tabata⁵,
M. Iwasaki⁵, S. Goto¹

1 DBCLS, 2 Niigata Univ, 3 Kyushu Univ, 4 Kumamoto Univ, 5 Kyoto Univ



Genome → Proteome → Function



Major Human Proteome DBs

UniProt/SwissProt, NCBIInr, ...

- **HPP-HUPO**
(Human Proteome Project)



Start in 2010,
international efforts

- **ProteomicsDB**
(LC/MS/MS)



Nature, 2014

- **Human Protein Atlas**
(Antibody-based)



Science, 2015

Nature. 2014 , DOI: [10.1038/nature13302](https://doi.org/10.1038/nature13302), PMID: 24870542

A draft map of the human proteome

[Min-Sik Kim](#); [Sneha M Pinto](#); [Derese Getnet](#); [Raja Nirujogi](#); [Srikanth S Manda](#); [Raghothama Chaerkady](#); [Anil K Madugundu](#); [Dhanashree S Kelkar](#); [Ruth Isserlin](#); [Shobhit Jain](#); [Joji K Thomas](#); [Babylakshmi Muthusamy](#); [Pamela Leal-Rojas](#); [Praveen Kumar](#); [Nandini A Sahasrabudde](#); [Lavanya Balakrishnan](#); [Jayshree Advani](#); [Bijesh George](#); [Santosh Renuse](#); [Lakshmi N Selvan](#); [Arun H Patil](#); [Vishalakshi Nanjappa](#); [Aneesha Radhakrishnan](#); [Samarjeet Prasad](#); [Tejaswini Subbannayya](#); [Rajesh Raju](#); [Manish Kumar](#); [Sreelakshmi K Sreenivasamurthy](#); [Arivusudar Marimuthu](#); [Gajanan J Sathe](#); [Sandip Chavan](#); [Keshava K Datta](#); [Yashwanth Subbannayya](#); [Apeksha Sahu](#); [Soujanya D Yelamanchi](#); [Savita Jayaram](#); [Pavithra Rajagopalan](#); [Jyoti Sharma](#); [Krishna R Murthy](#); [Nazia Syed](#); [Renu Goel](#); [Aafaque A Khan](#); [Sartaj Ahmad](#); [Gourav Dey](#); [Keshav Mudgal](#); [Aditi Chatterjee](#); [Tai-Chung Huang](#); [Jun Zhong](#); [Xinyan Wu](#); [Patrick G Shaw](#); ... (22 more)

The availability of human genome sequence has transformed biomedical research over the past decade. However, an equivalent map for the human proteome with direct measurements of proteins and peptides does not exist yet. Here we present a draft map of the human proteome using high-resolution Fourier-transform mass spectrometry. In-depth proteomic profiling of 30 histologically normal human samples, including 17 adult tissues, 7 fetal tissues and 6 purified primary haematopoietic cells, resulted in identification of proteins encoded by 17,294 genes accounting for approximately 84% of the total annotated protein-coding genes in humans. A unique and comprehensive strategy for proteogenomic analysis enabled us to discover a number of novel protein-coding regions, which includes translated pseudogenes, non-coding RNAs and upstream open reading frames. This large human proteome catalogue (available as an interactive web-based resource at <http://www.humanproteomemap.org>) will complement available human genome and transcriptome data to accelerate biomedical research in health and disease.

17,294 gene products

Nature. 2014 , DOI: [10.1038/nature13319](https://doi.org/10.1038/nature13319)

Mass-spectrometry-based draft of the human proteome

[Mathias Wilhelm](#); [Judith Schlegl](#); [Hannes Hahne](#); [Amin Moghaddas Gholami](#); [Marcus Lieberenz](#); [Mikhail M. Savitski](#); [Emanuel Ziegler](#); [Lars Butzmann](#); [Siegfried Gessulat](#); [Harald Marx](#); [Toby Mathieson](#); [Simone Lemeer](#); [Karsten Schnatbaum](#); [Ulf Reimer](#); [Holger Wenschuh](#); [Martin Mollenhauer](#); [Julia Slotta-Huspenina](#); [Joos-Hendrik Boese](#); [Marcus Bantscheff](#); [Anja Gerstmair](#); [Franz Faerber](#); [Bernhard Kuster](#)

Proteomes are characterized by large protein-abundance differences, cell-type- and time-dependent expression patterns and post-translational modifications, all of which carry biological information that is not accessible by genomics or transcriptomics. Here we present a mass-spectrometry-based draft of the human proteome and a public, high-performance, in-memory database for real-time analysis of terabytes of big data, called ProteomicsDB. The information assembled from human tissues, cell lines and body fluids enabled estimation of the size of the protein-coding genome, and identified organ-specific proteins and a large number of translated lincRNAs (long intergenic non-coding RNAs). Analysis of messenger RNA and protein-expression profiles of human tissues revealed conserved control of protein abundance, and integration of drug-sensitivity data enabled the identification of proteins predicting resistance or sensitivity. The proteome profiles also hold considerable promise for analysing the composition and stoichiometry of protein complexes. ProteomicsDB thus enables navigation of proteomes, provides biological insight and fosters the development of proteomic technology.

18,097 gene products 4

Analyzing the First Drafts of the Human Proteome

Iakes Ezkurdia,[†] Jesús Vázquez,[§] Alfonso Valencia,[‡] and Michael Tress^{*,‡}

[†]Unidad de Proteómica, [§]Laboratorio de Proteómica Cardiovascular, Centro Nacional de Investigaciones Cardiovasculares, Melchor Fernández Almagro, 3, Madrid 28029, Spain

[‡]Spanish National Cancer Research Centre (CNIO), Melchor Fernandez Almagro, 3, Madrid 28029, Spain

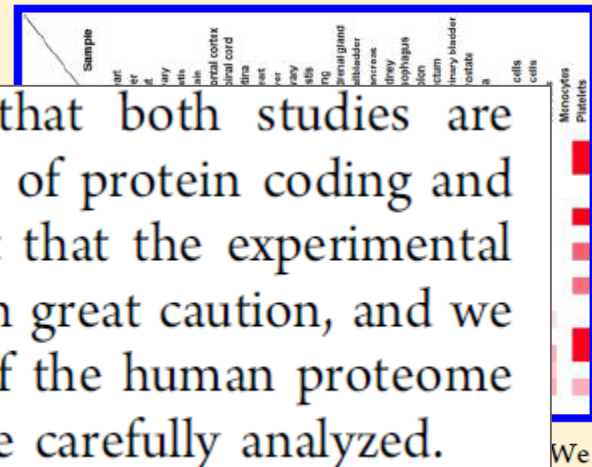
Supporting Information

ABSTRACT: This letter analyzes two large-scale proteomics studies published in the same issue of *Nature*. At the time of the release, both studies were

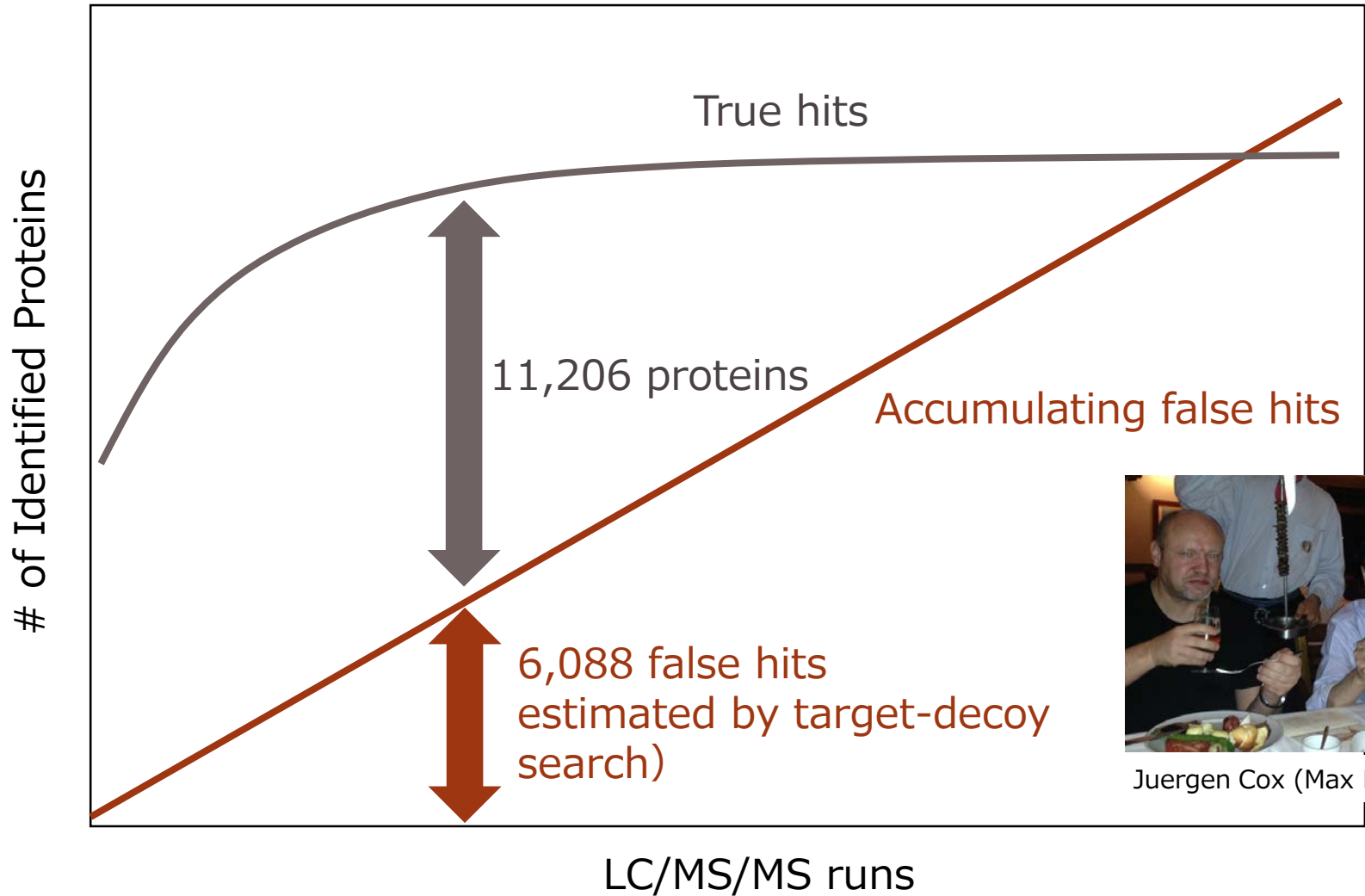
The results of our analysis show that both studies are substantially overestimating the number of protein coding and noncoding genes they find. We suggest that the experimental data from these two should be used with great caution, and we feel that these two unique draft maps of the human proteome should be put on hold until they can be carefully analyzed.

conclude that
KEYWORDS

These draft maps should be withdrawn!



30% false positives!



Juergen Cox (Max Planck)

Molecular & cellular proteomics : MCP. 2015 , DOI: [10.1074/mcp.M114.046995](https://doi.org/10.1074/mcp.M114.046995), PMID: 25987413

A scalable approach for protein false discovery rate estimation in large proteomic data sets.

Mikhail M Savitski; Mathias Wilhelm; Hannes Hahne; Bernhard Kuster; Marcus Bantscheff

Calculating the number of confidently identified proteins and estimating false discovery rate (FDR) is a challenge when analyzing very large proteomic datasets such as entire human proteomes. Biological and technical heterogeneity in proteomic experiments further add to the challenge and there are strong differences in opinion regarding the conceptual validity of a protein FDR and no consensus regarding the methodology for protein FDR determination. There are also limitations inherent to the widely used classic target-decoy strategy (TDS) that particularly show when analyzing very large data sets and that lead to a strong over-representation of decoy

identifications. In this study, we investigated the merits of the decoy-based protein FDR estimation approach taking advantage of a large-scale proteomic data collection comprised of ~19,000 LC-MS/MS runs deposited in ProteomicsDB (www.proteomicsdb.org). The "picked" protein FDR approach uses the same protein as a pair rather than as individual entities and a decoy sequence depending on which receives the highest score. The merits of this approach in combination with q-value based peptide scoring are independent of instrument and search engine-specific differences. The "picked" approach is the best when protein scoring was based on the best peptide score. The results demonstrate a stable number of true positive protein identifications over a range of data sets and demonstrate that this simple and unbiased strategy eliminates the commonly used, "classic" protein FDR approach that causes a significant decrease in protein identification in large data sets. The approach scales without losing performance, consistently increases the number of true positive protein identifications and is readily implemented in proteomics analysis software.

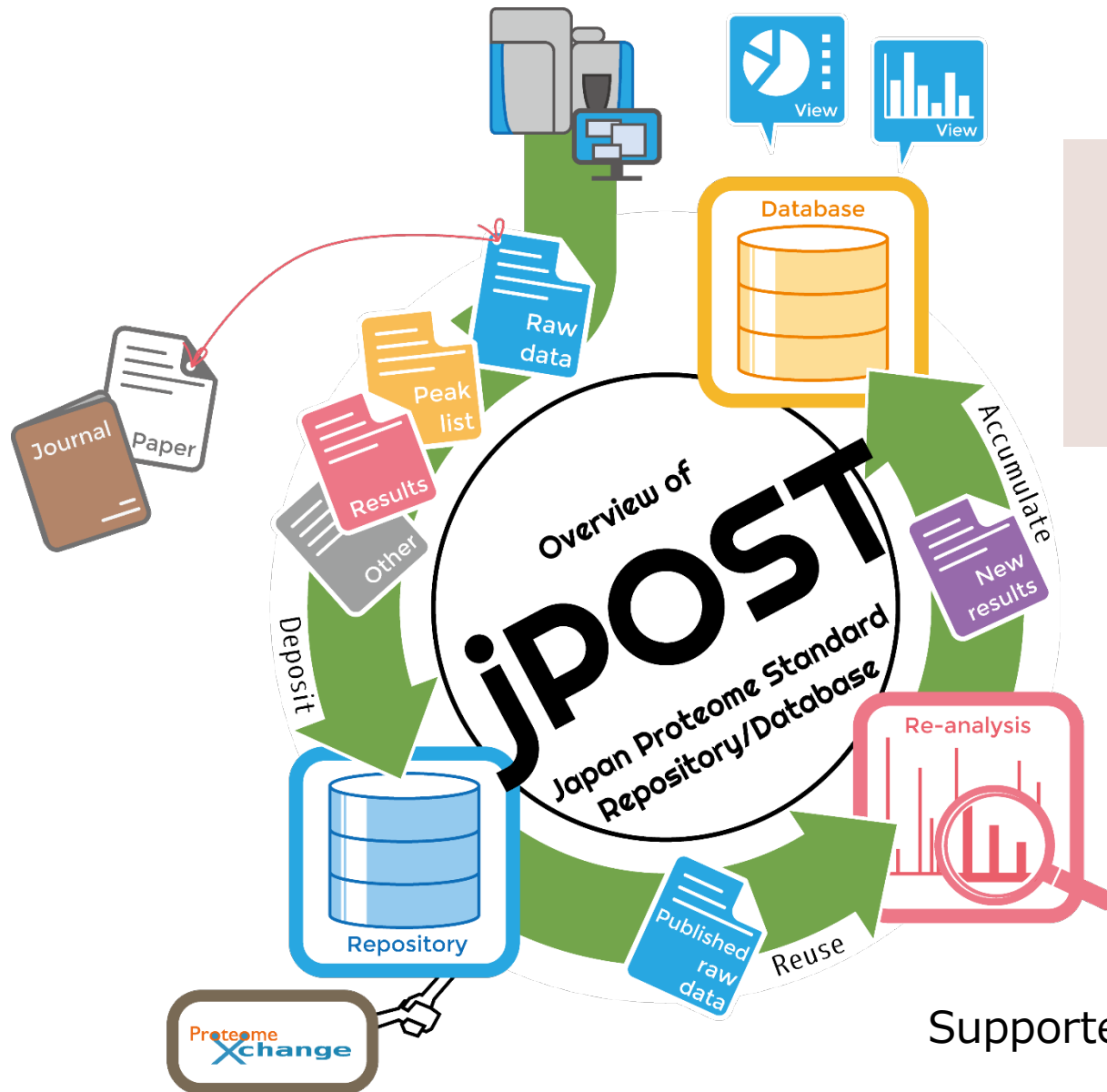
18,097 proteins (original)



for this difference (supplemental Figure 7). We next applied the described data analysis strategy to the subset of data stored in proteomicsDB corresponding to our earlier publication on a mass spectrometry based draft of the human proteome (9). Using the classic FDR strategy 14,035 proteins were observed at 1% protein FDR compared to 14,714 proteins using the picked strategy. Applying the picked strategy without any protein score threshold yielded 17,326 proteins of the target database at 11.3% protein FDR corresponding to 15,290 true positive protein identifications in the dataset. When analyzing the complete current content of proteomicsDB (including the data of the Pandey proteome (10) and a number of further datasets), the number of protein identifications at 1% FDR increased to 14,638 (classic) and 15,375 proteins (picked) respectively.

What is jPOST?

<http://jpostdb.org>



for Data Integration
& Sharing
in Life Science

Supported by NBDC-JST since 2015

← → 🏠 📄 jpostdb.org

🌐 アプリ ★ Bookmarks 🌐 Windows 🌐 リンクの変更 🌐 KUMail 📄 ログイン | JHUPO会員が 📄 薬学共用試験 📄 京大薬-cbt 📄 講習会

About Repository Database Workflow Contact 🔍



Japan Proteome Standard Repository/Database

Recent posts

other

jPOST joined to ProteomeXchange

🕒 2016-07-6 👤 jpost

We are pleased to announce that the jPOST repository has joined to ProteomeXchange consortium on July 6, 2016.

other

Server maintenance

🕒 2016-05-27 👤 jpost

jPOST repository server will be temporarily unavailable. May 29, 9:30 – 19:00 (UTC+9)

other

Announcement of jPOST repository

🕒 2016-05-2 👤 jpost

We are pleased to announce that our jPOST repository will be open on May 2, 2016. jPOST



The infographic illustrates a four-level workflow:

- Level 1 (Top):** 'Post raw data' (To repository) and 'Views' (To database). It shows a person at a computer with a 'Submission form data dashboard' and another person with 'Slices'.
- Level 2:** 'Register' (Repository) and 'Faceted search' (Globe). It shows a 'Cube' being 'Aggregate' into a 'Globe'.
- Level 3 (Bottom):** 'Reprocessing' (Results) and 'Pack into database'. It shows 'Results' being processed and then packed into a database.

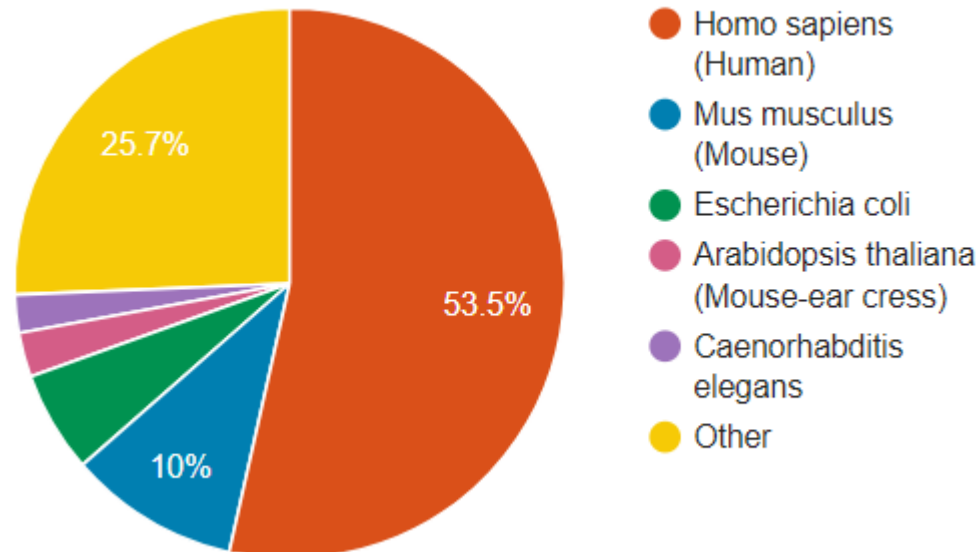
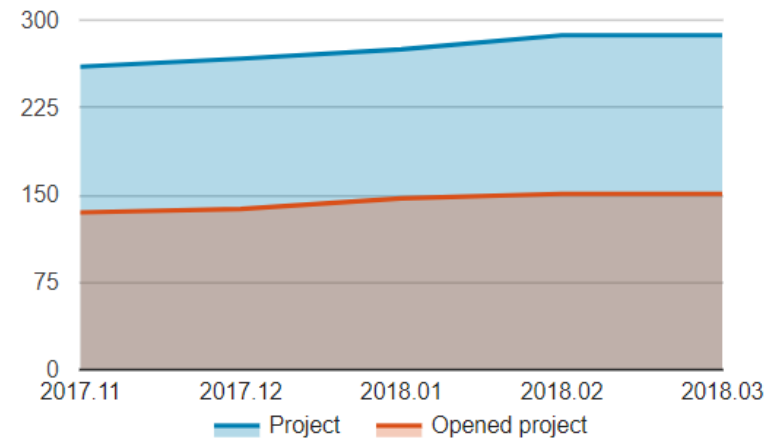
MS raw data with metadata are stored in ProteomeXchange formats.

Statistics

287 projects are registered. **151** are opened.

30559 files amount to **8.0 TB**.

43 species.



Molecular & Cellular Proteomics

Providing Access to Annotated Spectra

Currently, the guidelines for MCP require that annotated spectra be provided in these two cases:

- proteins identified on the basis of single peptide
- post-translationally modified peptides

The purpose of this document is to summarize in one location existing tools that an author may choose to utilize to convert different proteomic results formats from a variety of software tools into files that satisfy the MCP requirement for access to annotated spectra.

It is possible to make annotated spectra available from most search engines, although the options for how to do this differ between software. MCP requires the files required for annotated spectra to be stored in a public repository that is beyond the control of the authors, so a lab website is not a compliant location. It may be possible to submit them as supplementary files with the manuscript submission. However, these files are often large (>100 MB). If this is the case, there are a handful of public repositories that can be used to store these files and the authors just need to provide a link to the location where they have been uploaded at the time of manuscript submission. MCP does not officially endorse any one repository. However, repositories that are part of the proteomeXchange consortium (proteomexchange.org) are suitable choices.



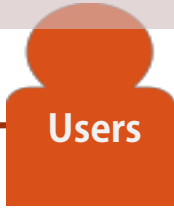
Instructions to Authors of *Journal of Proteome Research*

(Revised February 2018)

Important change about data deposition:

Author(s) are REQUIRED to deposit raw files and associated metadata in repositories such as ProteomeXchange (preferred) or other public repositories and to provide access to the information in the manuscript, including both the link as well as any necessary passwords (example shown below). Access to the information will be kept confidential while the manuscript is under review but will be open to the public upon publication. Please note: Providing this information on a link managed by the author(s) is not acceptable.

Manual curation of metadata



Users

Project ID

“Sample preset”のメタデータ

JPST000
203

JPST000
204

Species: Homo sapiens (Human)
Tissue: colorectal cancer cell

JPST000
205

JPST000
206

Species: Homo sapiens (Human)
Tissue: colon
Disease: carcinoma

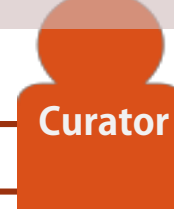
JPST000
207

JPST000
208

Species: Homo sapiens (Human)
Tissue: SW-480 cell, SW-620 cell
Disease: adenocarcinoma

JPST000
210

Species: Homo sapiens (Human)
Tissue: colon, colorectal cancer cell



Curator

Colorectal cancer tissue

JPST000
203

JPST000
204

JPST000
210

Species: Human
Sample Type: Tissues
Organ: Colorectum
Disease: Cancer
Disease name: Colorectal cancer

Colorectal cancer cell line

JPST000
205

JPST000
206

JPST000
207

JPST000
208

Species: Human
Sample Type: Cell line (HCT116)
Organ: Colorectum
Disease: Cancer
Disease name: Colorectal cancer

Species: Human
Sample Type: Cell line (SW480, SW620)
Organ: Colorectum
Disease: Cancer
Disease name: Colorectal cancer

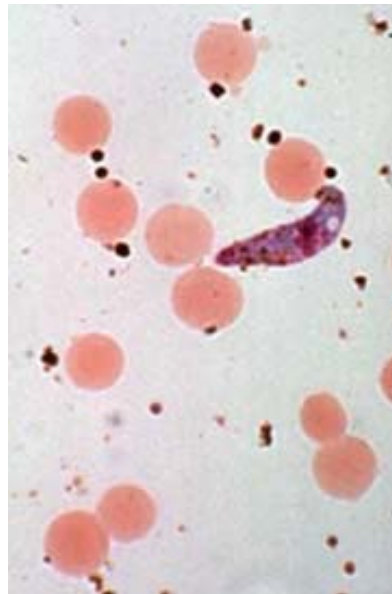
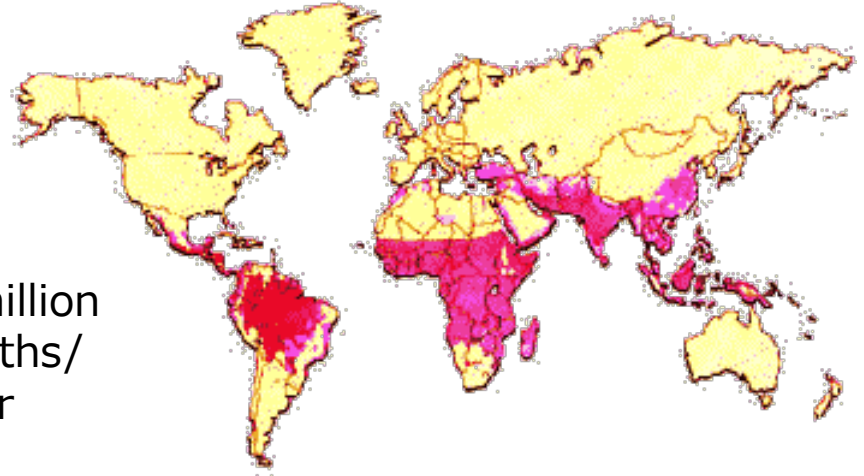
Proteomics community efforts against false hits

Malaria *Plasmodium falciparum* proteome by high-accuracy mass spectrometry in 2002

Hosts

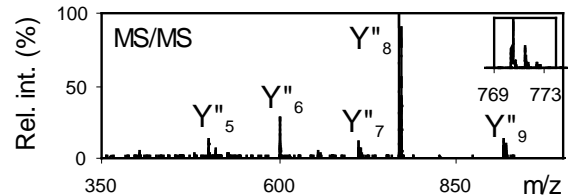
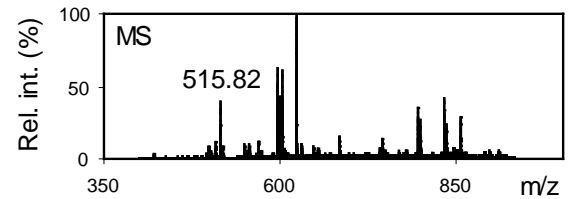
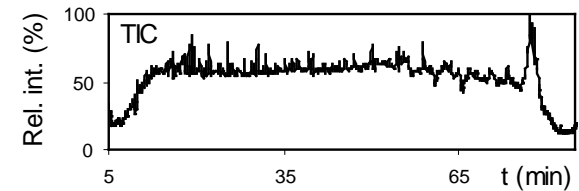
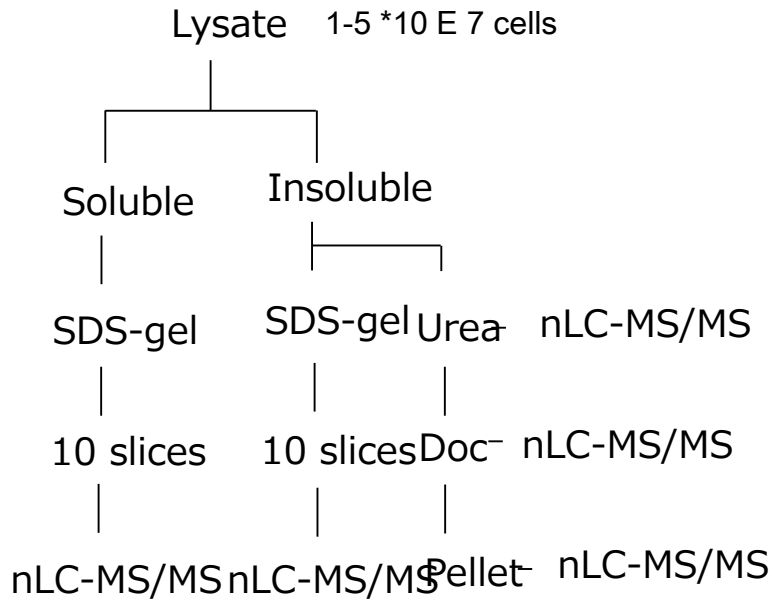
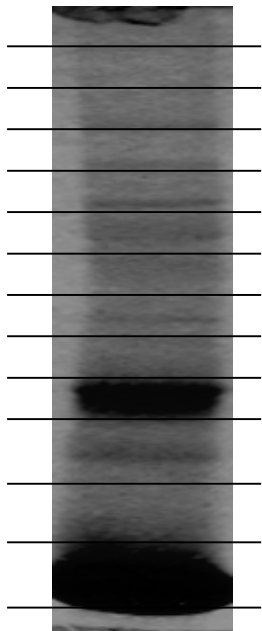
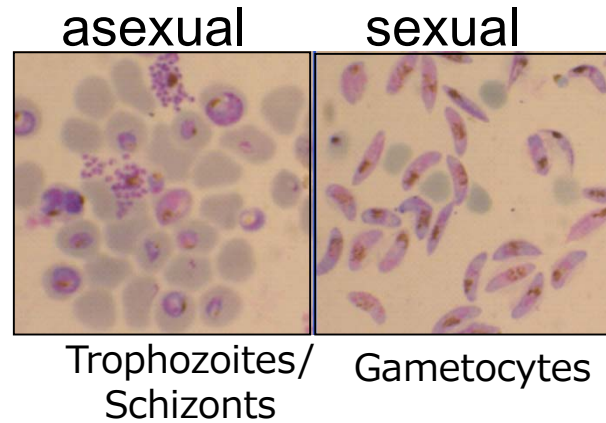
- Anopheles mosquito
- human

2 million
deaths/
year



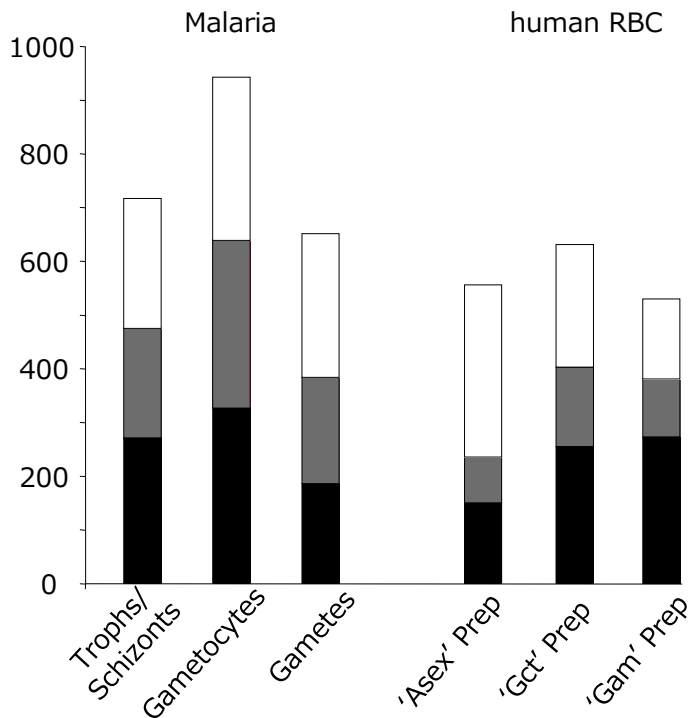
Proteomics of blood stages

infected RBCs

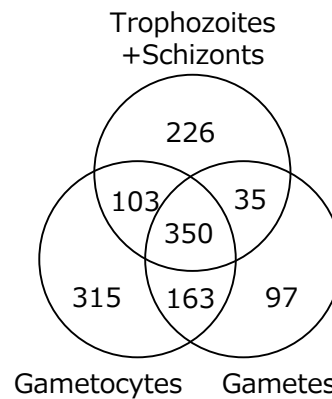


Identified proteins

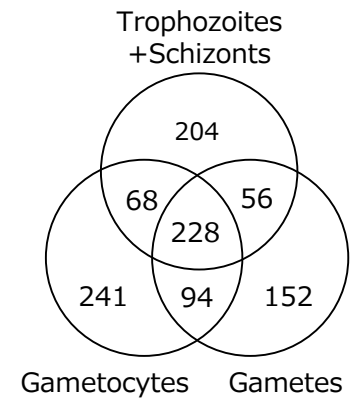
- soluble
- shared
- insoluble



Malaria proteins
Total unique: 1289



human proteins
Total unique: 1043



Most known merozoite surface proteins identified.

ca. 200 hypothetical proteins with transmembrane domains



New vaccine candidates

Nature: Malaria special issue in 2002

articles

A proteomic view of the *Plasmodium falciparum* life cycle

Laurence Florens*, Michael P. Washburn†, J. Dale Raine‡, Robert M. Anthony§, Munira Grainger||, J. David Haynes¶, J. Kathleen Moch§, Nemone Muster*, John B. Sacci§#, David L. Tabb*☆, Adam A. Witney§#, Dirk Wolters†#, Yimin Wu***, Malcolm J. Gardner††, Anthony A. Holder||, Robert E. Sinden‡, John R. Yates*† & Daniel J. Carucci§

* Department of Cell Biology, The Scripps Research Institute, SR-11, 10550 North Torrey Pines Road, La Jolla, California 92037, USA

† Department of Proteomics and Metabolomics, Torrey Mesa Research Institute, Syngenta Research & Technology, 3115 Merryfield Row, San Diego, California 92121-1125, USA

‡ Infection and Immunity Section, Department of Biological Sciences, Imperial College of Science, Technology & Medicine, Sir Alexander Fleming Building, South Kensington, London SW7 2AZ, UK

§ Naval Medical Research Center, Malaria Program (IDD), 503 Robert Grant Avenue, Room 3A40; and ¶ Department of Immunology, Walter Reed Army Institute of Research, Silver Spring, Maryland 20910-7500, USA

|| The Division of Parasitology, National Institute for Medical Research, The Ridgeway, Mill Hill, London NW7 1AA, UK

¶ Department of Genome Sciences, University of Washington, Seattle, Washington 98195, USA

** Malaria Research and Reference Reagent Resource Center, American Type Culture Collection, 10801 University Boulevard, Manassas, Virginia 20110-2209, USA

†† The Institute for Genomic Research, 9712 Medical Center Drive, Rockville, Maryland 20850, USA

The completion of the *Plasmodium falciparum* clone 3D7 genome provides a basis on which to conduct comparative proteomics studies of this human pathogen. Here, we applied a high-throughput proteomics approach to identify new potential drug and vaccine targets and to better understand the biology of this complex protozoan parasite. We characterized four stages of the parasite life cycle (sporozoites, merozoites, trophozoites and gametocytes) by multidimensional protein identification technology. Functional profiling of over 2,400 proteins agreed with the physiology of each stage. Unexpectedly, the antigenically variant proteins of *var* and *rif* genes, defined as molecules on the surface of infected erythrocytes, were also largely expressed in sporozoites. The detection of chromosomal clusters encoding co-expressed proteins suggested a potential mechanism for controlling gene expression.

Analysis of the *Plasmodium falciparum* proteome by high-accuracy mass spectrometry

Edwin Lasonder*†, Yasushi Ishihama*, Jens S. Andersen*, Adriaan M. W. Vermunt†, Arnab Pain‡, Robert W. Sauerwein§, Wijnand M. C. Eling§, Neil Hall‡, Andrew P. Waters||, Hendrik G. Stunnenberg† & Matthias Mann*

* Center for Experimental Bioinformatics, Department of Biochemistry and Molecular Biology, University of Southern Denmark, Campusvej 55, DK-5230 Odense M, Denmark

† Department of Molecular Biology, NCMLS, University of Nijmegen, Geert Grooteplein 26, 6525 GA Nijmegen, The Netherlands

‡ The Wellcome Trust Sanger Institute, The Wellcome Trust Genome Campus, Hinxton, Cambridge CB10 1SA, UK

§ Department of Medical Microbiology, NCMLS, University Medical Centre, PO Box 9101, 6500 HB Nijmegen, The Netherlands

|| Leiden Malaria Research Group, Department of Parasitology, Centre for Infectious Disease, Leiden University Medical Center, Albinusdreef 2, 2333 ZA Leiden, The Netherlands

The annotated genomes of organisms define a 'blueprint' of their possible gene products. Post-genome analyses attempt to confirm and modify the annotation and impose a sense of the spatial, temporal and developmental usage of genetic information by the

letters to nature

organism. Here we describe a large-scale, high-accuracy (average deviation less than 0.02 Da at 1,000 Da) mass spectrometric proteome analysis¹⁻³ of selected stages of the human malaria parasite *Plasmodium falciparum*. The analysis revealed 1,289 proteins of which 714 proteins were identified in asexual blood stages, 931 in gametocytes and 645 in gametes. The last two groups provide insights into the biology of the sexual stages of the parasite, and include conserved, stage-specific, secreted and membrane-associated proteins. A subset of these proteins contain domains that indicate a role in cell-cell interactions, and therefore can be evaluated as potential components of a malaria vaccine formulation. We also report a set of peptides with significant matches in the parasite genome but not in the protein set predicted by computational methods.

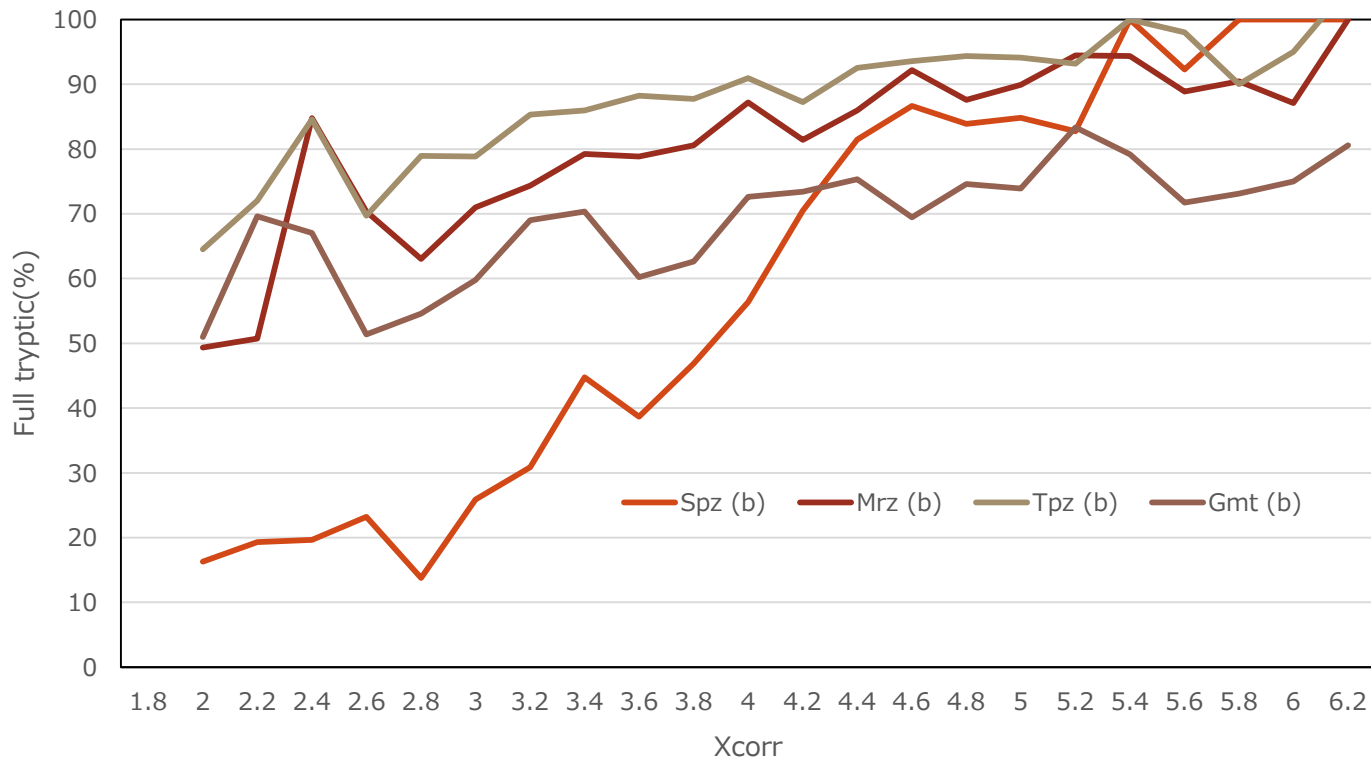
So sad..... 

MS/MS data set analysis

The SEQUEST algorithm was used to match MS/MS spectra to peptides in the sequence databases⁴¹. To account for carboxyamidomethylation, MS/MS data sets were searched with a relative molecular mass of 57,000 (M_r , 57K) added to the average molecular mass of cysteines. Peptide hits were filtered and sorted with DTASelect⁴². Spectra/peptide matches were only retained if they were at least half-tryptic (Lys or Arg at either end of the identified peptide) and with minimum cross-correlation scores (XCorr) of 1.8 for +1, 2.5 for +2, and 3.5 for +3 spectra and DeltaCn (top match's XCorr minus the second-best match's XCorr divided by the top match's XCorr) ≥ 0.08 . Peptide hits were deemed unambiguous only if they were not found in non-infected controls and were uniquely assigned to parasite proteins by searching against combined parasite–host databases. Finally, for low coverage loci, peptide/spectrum matches were visually assessed on two main criteria: any given MS/MS spectrum had to be clearly above the baseline noise, and both *b* and *y* ion series had to show continuity. The Contrast tool⁴² was used to compare and merge protein lists from replicate sample runs and to compare the proteomes established for the four stages.

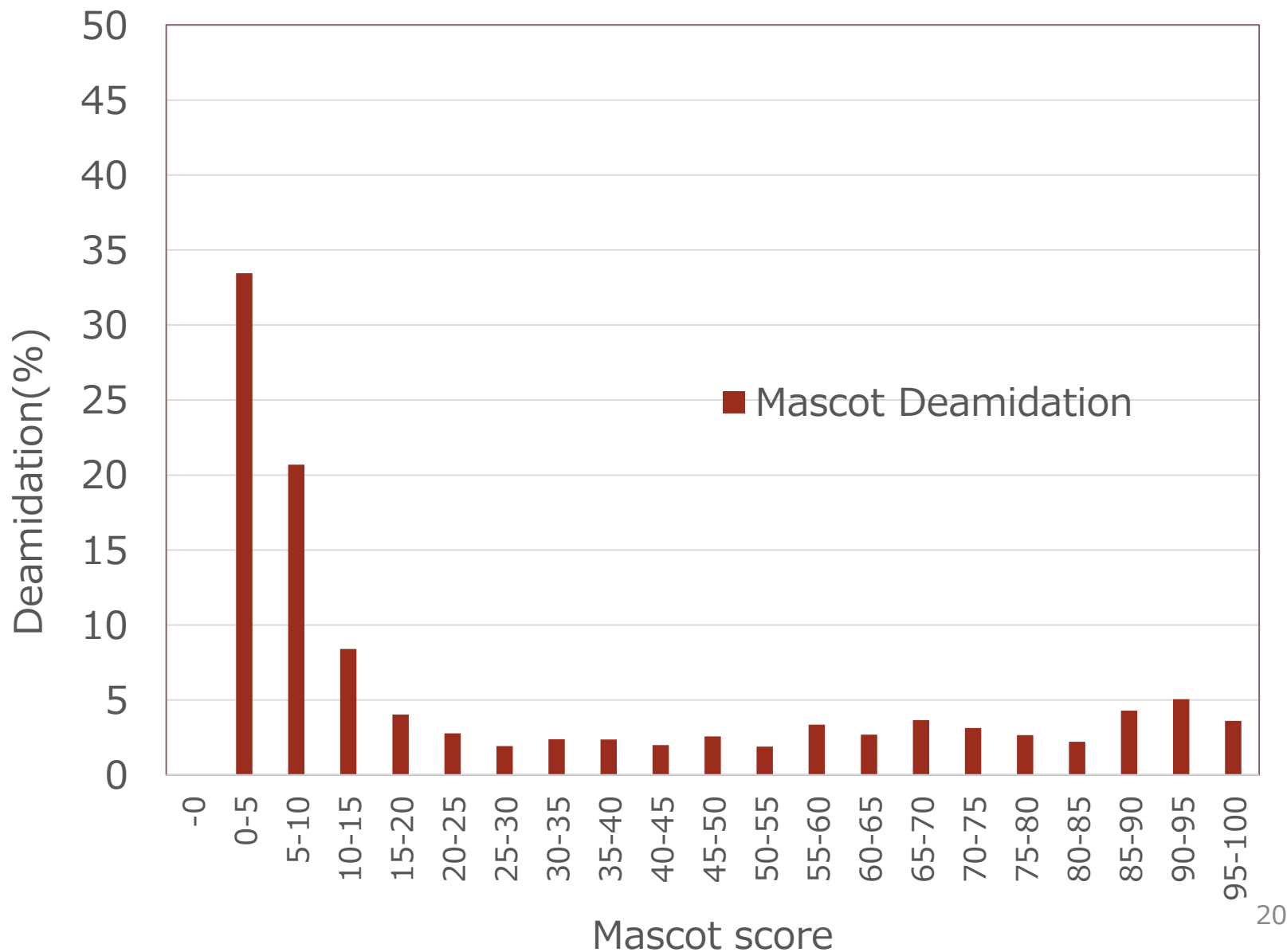
Umm, too wide search space??

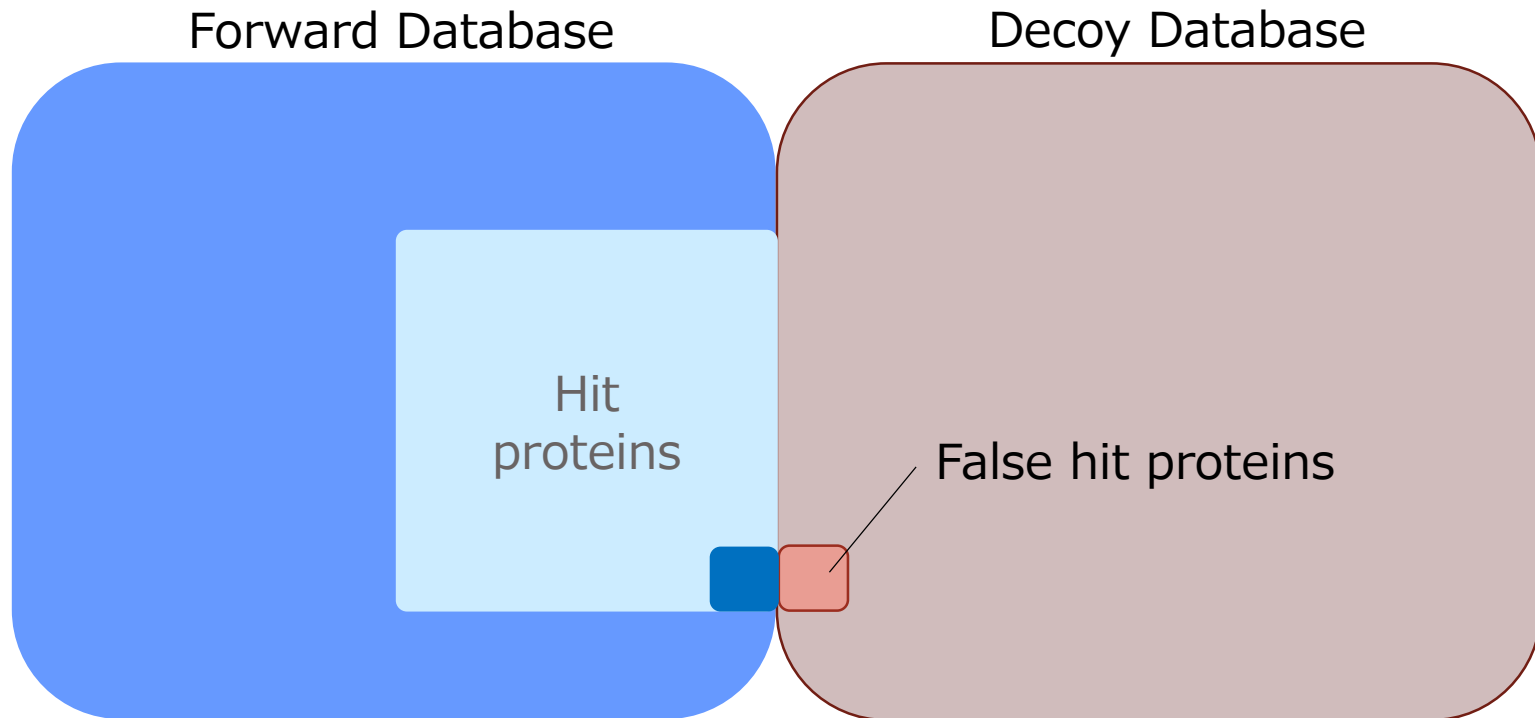
Score distribution of fully tryptic peptides (%)



This should be prohibited!!

Variable mod: De-amidation (N, Q)





$$\text{False Discovery Rate} = \frac{\text{False Positive}}{\text{False Positive} + \text{True Positive}}$$

Target-decoy search strategy for increased confidence in large-scale protein identifications by mass spectrometry

Joshua E Elias¹ & Steven P Gygi^{1,2}

NATURE METHODS | VOL.4 NO.3 | MARCH 2007 | 207

Rather than deciding exactly which peptide-spectral matches (PSMs) are correct or incorrect, the composite target-decoy database evaluates FP rates in large PSM populations. It permits estimation of the likelihood that a PSM is correct given that it came from a collection of PSMs with a measured FP rate. This is not to suggest that the search strategy removes all false identifications. Instead, the target-decoy approach allows the estimation of how many FP are associated with an entire data set.

Nature. 2014 , DOI: [10.1038/nature13302](https://doi.org/10.1038/nature13302), PMID: 24870542

A draft map of the human proteome

[Min-Sik Kim](#); [Sneha M Pinto](#); [Derese Getnet](#); [Raja Nirujogi](#); [Srikanth S Manda](#); [Raghothama Chaerkady](#); [Anil K Madugundu](#); [Dhanashree S Kelkar](#); [Ruth Isserlin](#); [Shobhit Jain](#); [Joji K Thomas](#); [Babylakshmi Muthusamy](#); [Pamela Leal-Rojas](#); [Praveen Kumar](#); [Nandini A Sahasrabudde](#); [Lavanya Balakrishnan](#); [Jayshree Advani](#); [Bijesh George](#); [Santosh Renuse](#); [Lakshmi N Selvan](#); [Arun H Patil](#); [Vishalakshi Nanjappa](#); [Aneesha Radhakrishnan](#); [Samarjeet Prasad](#); [Tejaswini Subbannayya](#); [Rajesh Raju](#); [Manish Kumar](#); [Sreelakshmi K Sreenivasamurthy](#); [Arivusudar Marimuthu](#); [Gajanan J Sathe](#); [Sandip Chavan](#); [Keshava K Datta](#); [Yashwanth Subbannayya](#); [Apeksha Sahu](#); [Soujanya D Yelamanchi](#); [Savita Jayaram](#); [Pavithra Rajagopalan](#); [Jyoti Sharma](#); [Krishna R Murthy](#); [Nazia Syed](#); [Renu Goel](#); [Aafaque A Khan](#); [Sartaj Ahmad](#); [Gourav Dey](#); [Keshav Mudgal](#); [Aditi Chatterjee](#); [Tai-Chung Huang](#); [Jun Zhong](#); [Xinyan Wu](#); [Patrick G Shaw](#); ... (22 more)



Juergen Cox (Max Planck)

The availability of human genome sequence has transformed biomedical research over the past decade. However, an equivalent map for the human proteome with direct measurements of proteins and peptides does not exist yet. Here we present a draft map of the human proteome using high-resolution Fourier-transform mass spectrometry. In-depth proteomic profiling of 30 histologically normal human samples, including 17 adult tissues, 7 fetal tissues and 6 purified primary haematopoietic cells, resulted in identification of proteins encoded by 17,294 genes accounting for approximately 84% of the total annotated protein-coding genes in humans. A unique and comprehensive strategy for

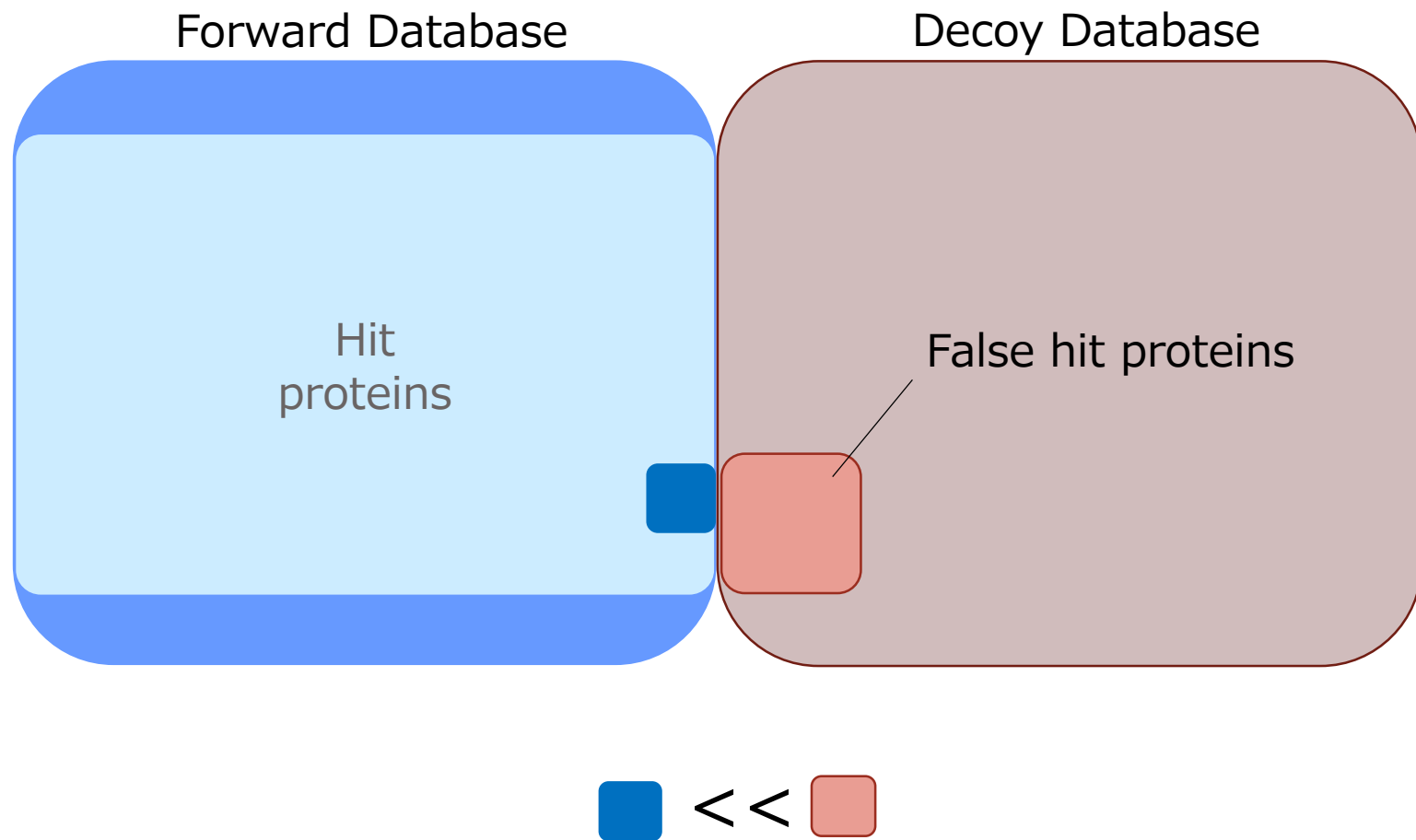
17,294 genes (84%)



Target-decoy search for all merged data
1% FDR at protein level

11,206 genes (57%)

Target-Decoy Approach for Ultra Large Datasets



Molecular & cellular proteomics : MCP. 2015 , DOI: [10.1074/mcp.M114.046995](https://doi.org/10.1074/mcp.M114.046995), PMID: 25987413

A scalable approach for protein false discovery rate estimation in large proteomic data sets.

Mikhail M Savitski; Mathias Wilhelm; Hannes Hahne; Bernhard Kuster; Marcus Bantscheff

Calculating the number of confidently identified proteins and estimating false discovery rate (FDR) is a challenge when analyzing very large proteomic datasets such as entire human proteomes. Biological and technical heterogeneity in proteomic experiments further add to the challenge and there are strong differences in opinion regarding the conceptual validity of a protein FDR and no consensus regarding the methodology for protein FDR determination. There are also limitations inherent to the widely used classic target-decoy strategy (TDS) that particularly show when analyzing very large data sets and that lead to a strong over-representation of decoy

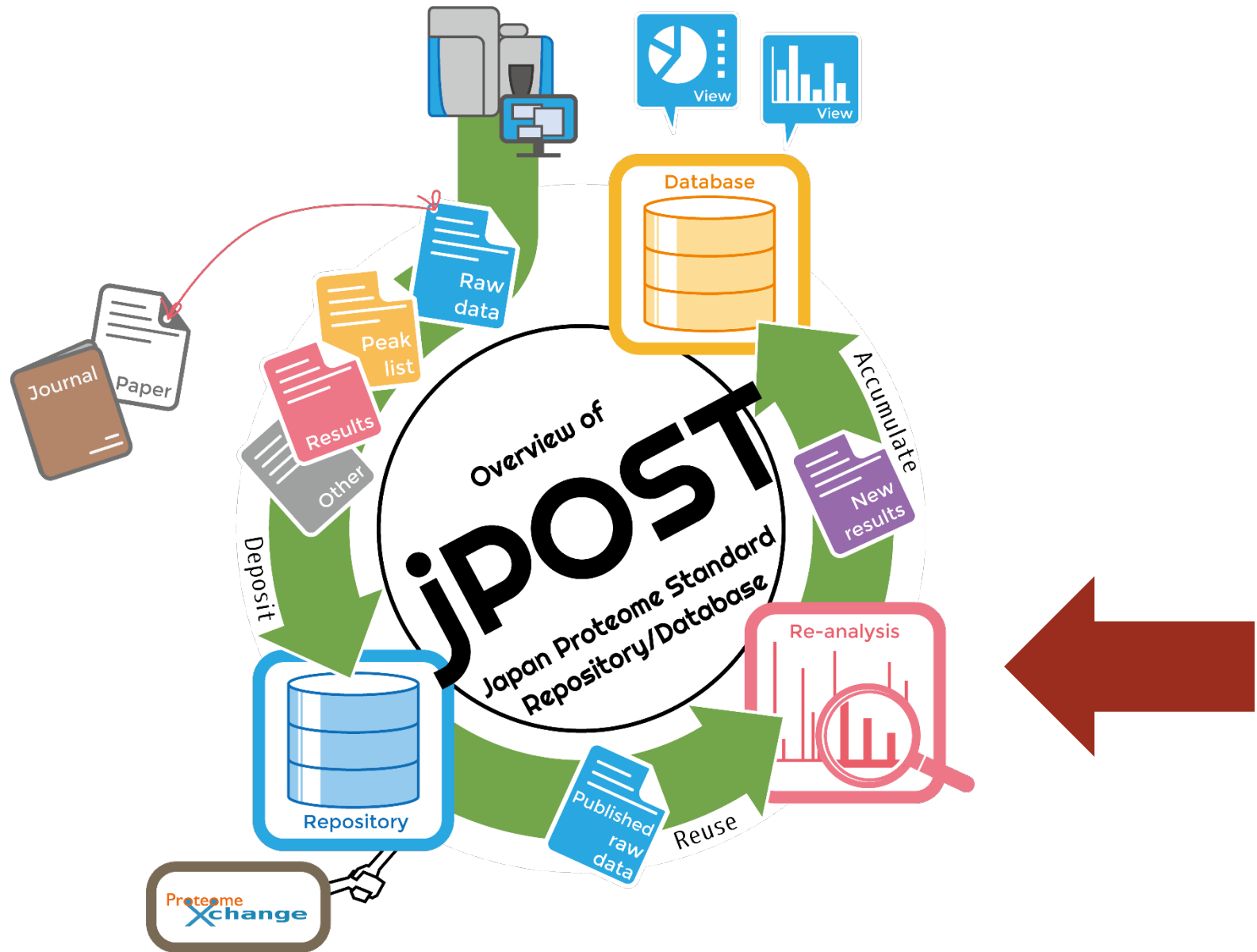
18,097 proteins (original)



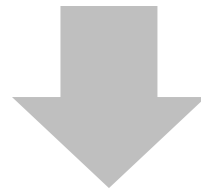
identifications. In this study, we investigated the merits of the decoy-based protein FDR estimation approach taking advantage of a large data collection comprised of ~19,000 LC-MS/MS runs deposited in proteomicsDB (www.proteomicsdb.org). The "picked" protein FDR approach uses the same protein as a pair rather than as individual entities and a decoy sequence depending on which receives the highest score. The merits of this approach in combination with q-value based peptide scoring are independent of instrument and search engine-specific differences. The "picked" approach is best when protein scoring was based on the best peptide score. It demonstrates a stable number of true positive protein identifications over a range of data sets. We demonstrate that this simple and unbiased strategy eliminates the commonly used, "classic" protein FDR approach that causes a significant decrease in protein identification in large data sets. The approach scales without losing performance, consistently increases the number of true positive protein identifications and is readily implemented in proteomics analysis software.

for this difference (supplemental Figure 7). We next applied the described data analysis strategy to the subset of data stored in proteomicsDB corresponding to our earlier publication on a mass spectrometry based draft of the human proteome (9). Using the classic FDR strategy 14,035 proteins were observed at 1% protein FDR compared to 14,714 proteins using the picked strategy. Applying the picked strategy without any protein score threshold yielded 17,326 proteins of the target database at 11.3% protein FDR corresponding to 15,290 true positive protein identifications in the dataset. When analyzing the complete current content of proteomicsDB (including the data of the Pandey proteome (10) and a number of further datasets), the number of protein identifications at 1% FDR increased to 14,638 (classic) and 15,375 proteins (picked) respectively.

jPOST re-analysis



How can we merge the results
from different sources?



jPOST score

- based on peak annotation in MSMS
- search engine independent
- MS instrument independent
- search DB independent
- can be used as universal threshold for peptide identification

Sequence Query

Introduction

The sequence query, in which one or more peptide molecular masses are combined with sequence, composition and fragment ion data, is potentially the most powerful search of all. The usual source of the sequence information is interpretation of an MS/MS spectrum. While it is very difficult to determine a complete and unambiguous peptide sequence from an MS/MS spectrum, it is often possible to find a series of peaks providing 3 or 4 residues of reliable sequence data.

This general approach was pioneered by Mann and co-workers at EMBL, who used the term "sequence tag" for the combination of a few residues of sequence data combined with molecular weight information [Mann, 1994]. They defined a sequence tag derived from an MS/MS spectrum as the mass of the precursor peptide, the mass of the first peak of the identified sequence ladder, a stretch of interpreted sequence, and the mass of the final peak₂₈ of the ladder.

PST-based jPOST SCORE



jPOST score

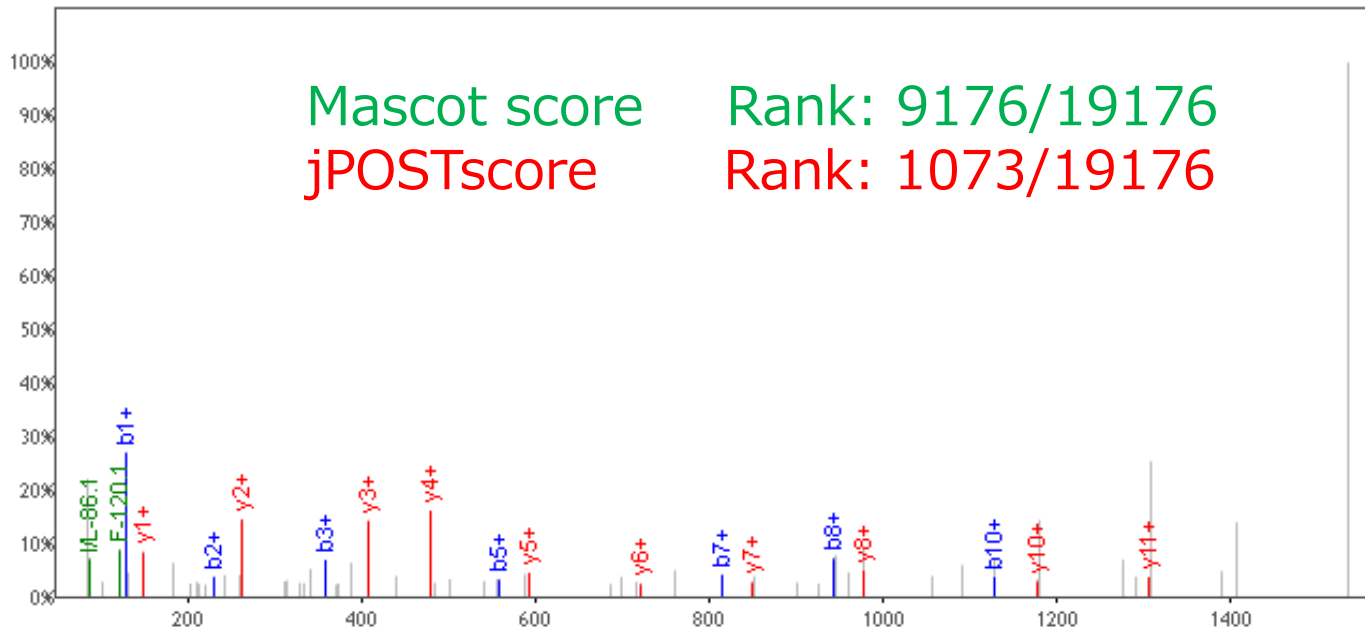
1. b, y-ion coverage
2. tag length
3. uncovered length

KVESLQEEIAFLK, MH+ 1533.8523, m/z 767.4298

File: 120201ry_aHDF1388-P9_1_3.wiff, PeptExpt: 2.0e-02, DeltaMass: -0.0045, Scan: 1.1.1257.7, Exp. m/z: 767.4275, Charge: 2

Mascot score
jPOSTscore

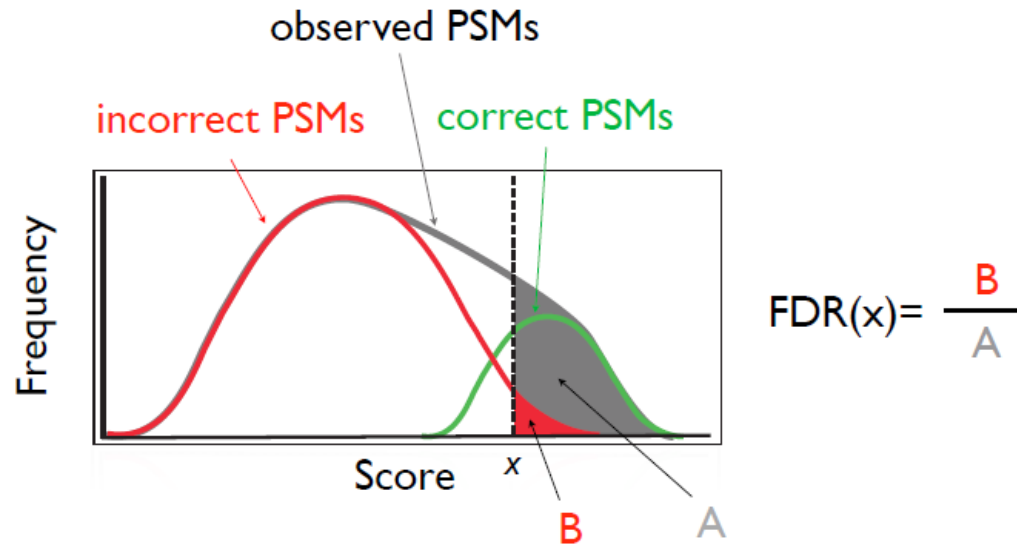
Rank: 9176/19176
Rank: 1073/19176



b+	#	Seq	#	y+
129.1022	1	K	13	
228.1707	2	V	12	1405.7573
357.2132	3	E	11	1306.6889
444.2453	4	S	10	1177.6463
557.3293	5	L	9	1090.6143
685.3879	6	Q	8	977.5302
814.4305	7	E	7	849.4716
943.4731	8	E	6	720.4291
1056.5572	9	I	5	591.3865
1127.5943	10	A	4	478.3024
1274.6627	11	F	3	407.2653
1387.7468	12	L	2	260.1969
	13	K	1	147.1128

[\[Click\]](#) to move table

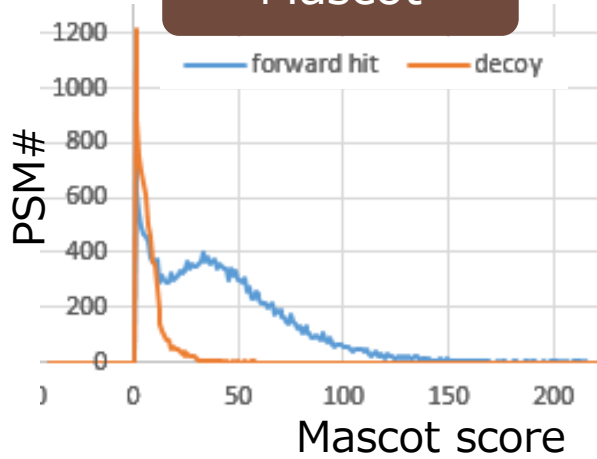
False Discovery Rate



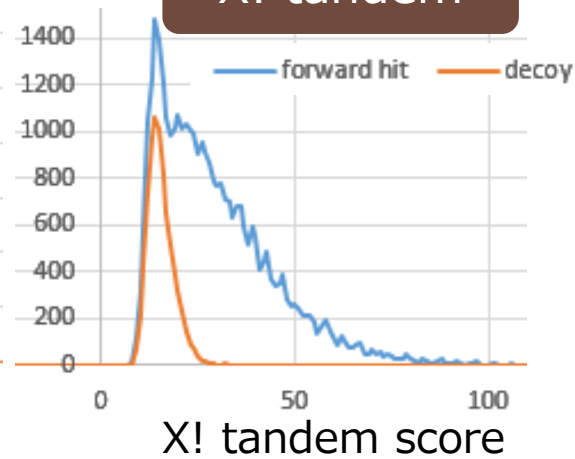
$FDR(x)$ is the expectation value of the fraction of detections above threshold x that are incorrect



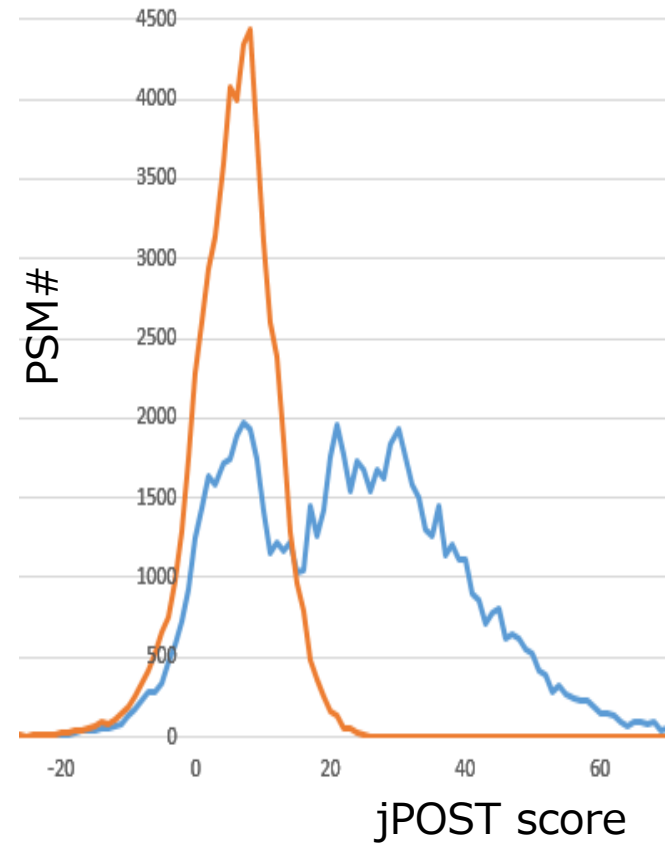
Mascot



X! tandem



Mascot + X!tandem + Comet + MQ

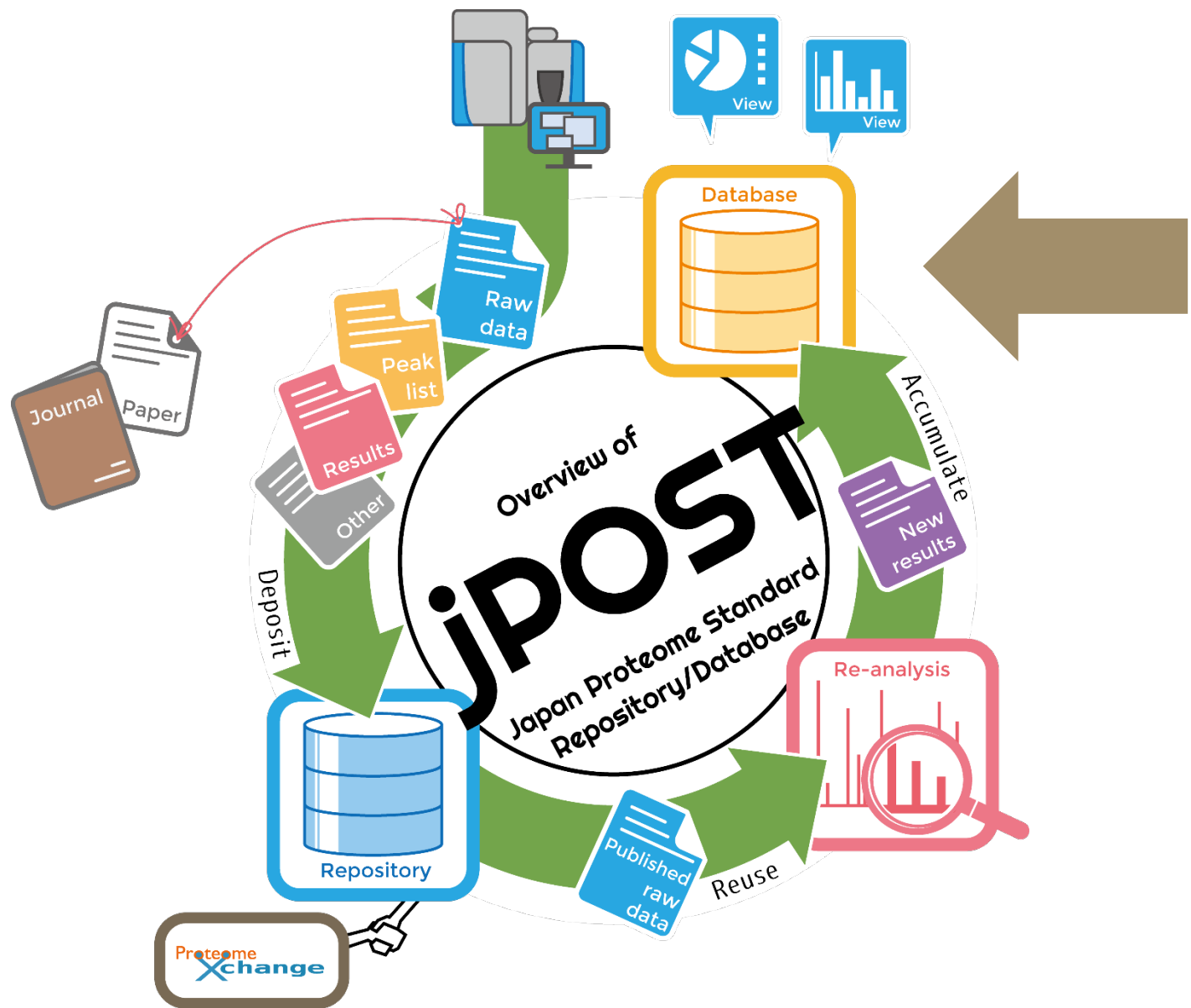


Dataset: PXD005159

Tryptic peptides from
human HeLa cells

by Thermo Q-Exactive

jPOST customizable database 'Slice'



jPOST slice database



Browser window showing the jPOST slice database interface. The URL is localhost/jPost-db/src/index.html. The navigation menu includes jPost, Search, Slices, and Compare.

Filters

Species: × Homo sapiens

Tissue: × colon, × colorectal cancer cell

Disease:

- neuroblastoma
- breast cancer
- carcinoma**
- adenocarcinoma

A hand cursor is pointing to the "carcinoma" option in the Disease filter.

Search boxes for filtering datasets by the experimental metadata from jPOST database

iPS **colon** +

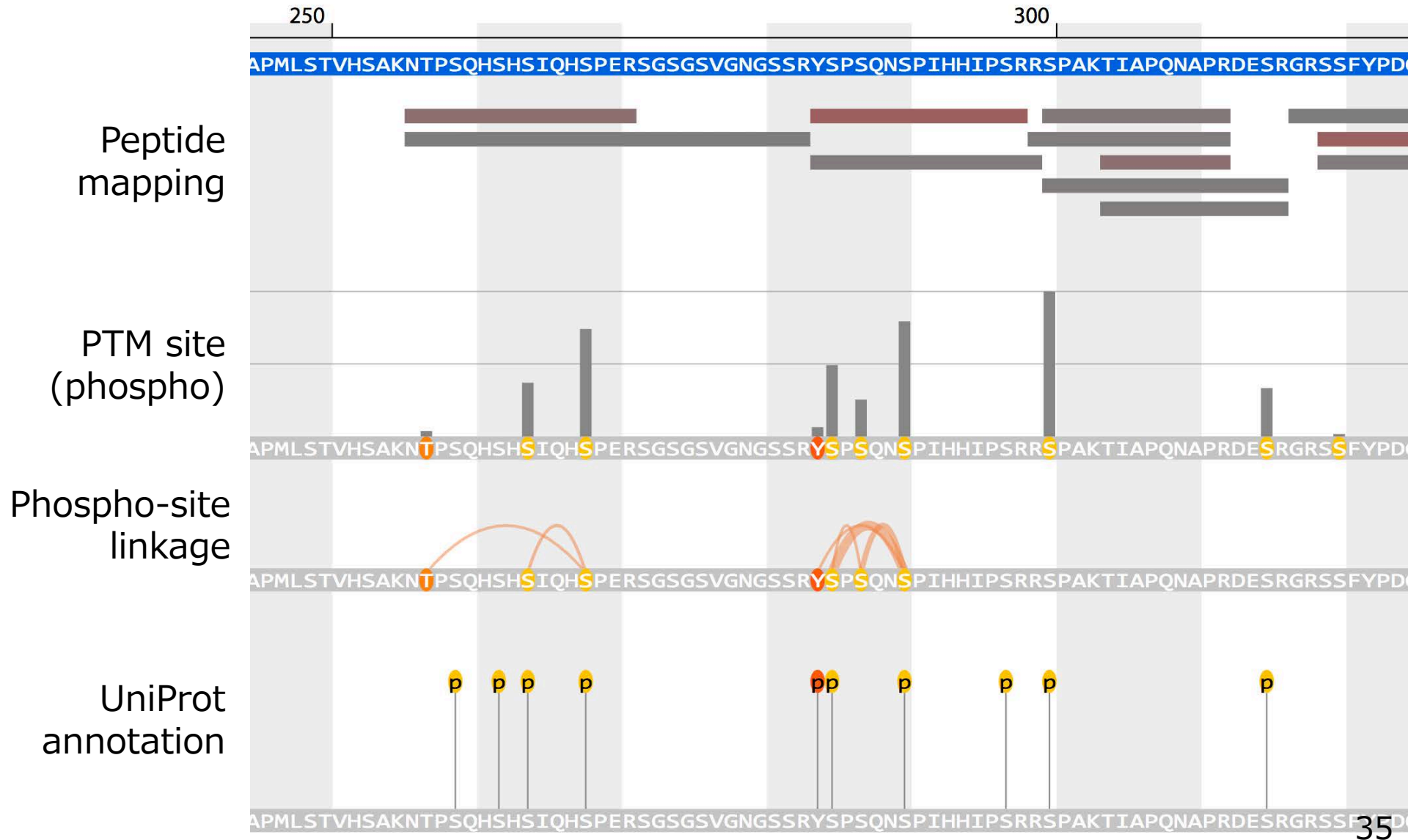
Dataset Protein

Selected dataset list
called ' **slice** '

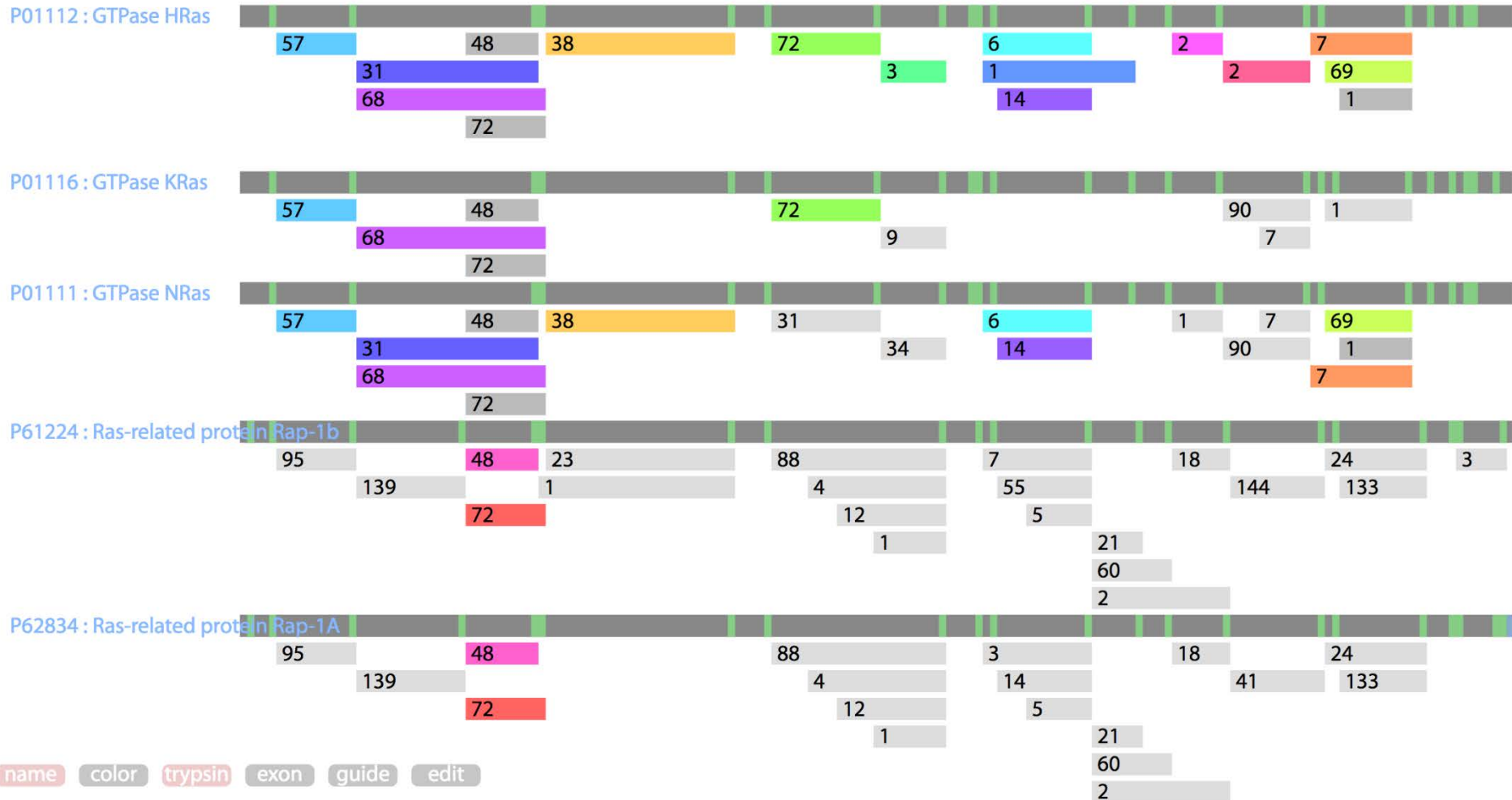
Showing 1 to 6 of 6 entries

ID ▲	Project ID ◆	Project Title ◆	Project Date ◆
DS203_1	JPST000203	Quantitative proteomics of colorectal cancer tissues	2016-10-18
DS204_1	JPST000204	Quantitative phosphoproteomics of colorectal cancer tissues	2016-10-18
DS205_1	JPST000205	Proteomic data of HCT116 cells	2016-10-18
DS206_1	JPST000206	Phosphoproteomic data of HCT116 cells	2016-10-18
DS210_1	JPST000210	Phosphoproteomics data of colon tissues (tumor and non-tumor)	2016-10-18
DS210_2	JPST000210	Phosphoproteomics data of colon tissues (tumor and non-tumor)	2016-10-18

Protein browser



Proteoform browser shows peptide sharing



KEGG pathway mapping with absolute quantitative value

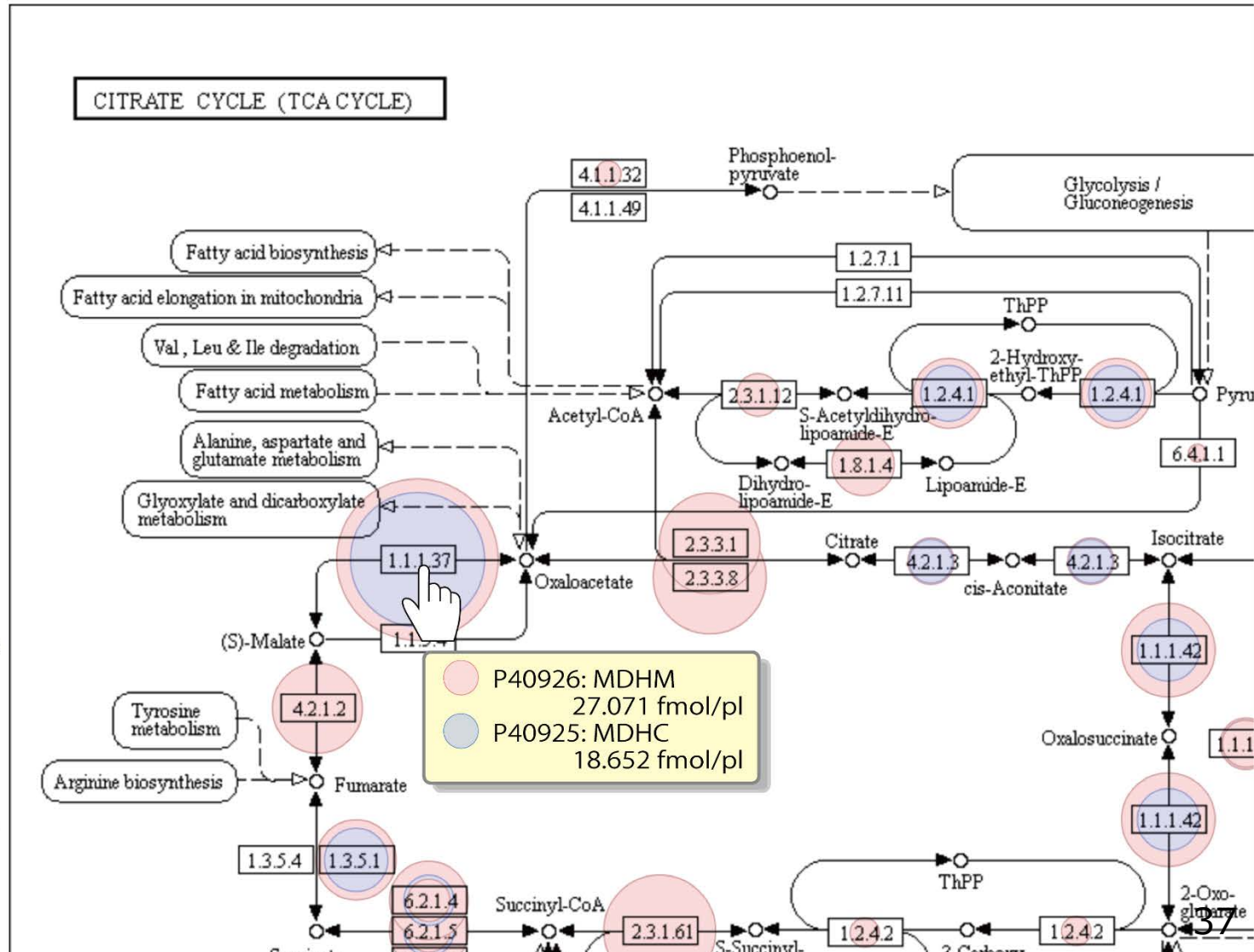
Metabolism

Carbohydrate metabolism

- Glycolysis / Gluconeogenesis
- Amino sugar and nucleotide sugar metabolism
- Pyruvate metabolism
- Inositol phosphate metabolism
- Citrate cycle (TCA cycle)
- Propanoate metabolism
- Fructose and mannose metabolism
- Pentose phosphate pathway
- Glyoxylate and dicarboxylate metabolism
- Galactose metabolism
- Starch and sucrose metabolism
- Butanoate metabolism
- Pentose and glucuronate interconversions
- Ascorbate and aldarate metabolism

Amino acid metabolism

- Valine, leucine and isoleucine degradation
- Cysteine and methionine metabolism
- Lysine degradation
- Arginine and proline metabolism
- Glycine, serine and threonine metabolism
- Alanine, aspartate and glutamate metabolism
- Tryptophan metabolism
- Arginine biosynthesis



Missing protein search

using latest **nextprot** & peptide uniqueness checker

chromosome:

protein evidence:

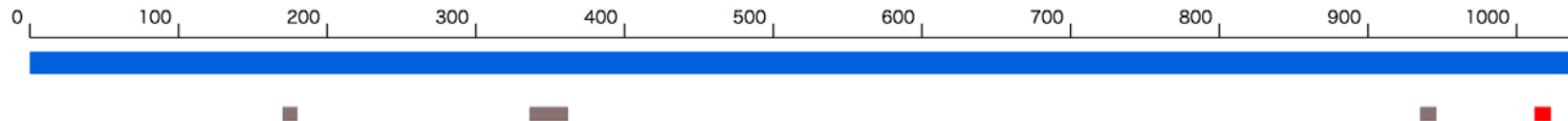
peptide length:

number of peptide:

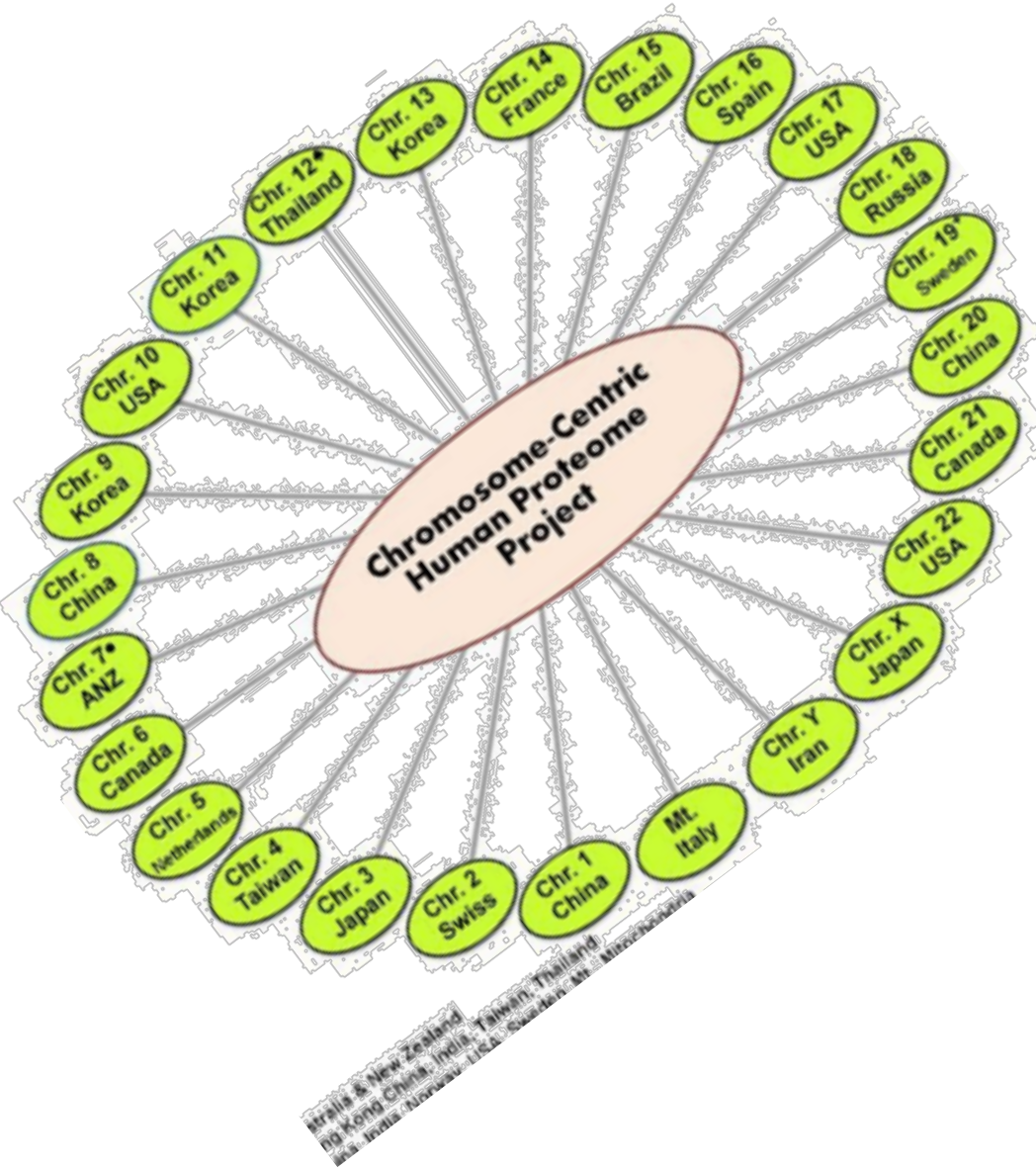
unique peptide:

PE	Chromosome	UniProt	Gene Symbol	Name	#Peptide	#Unique Peptide
3	X	A6NGH7	CC160_HUMAN	Coiled-coil domain-containing protein 160	2	2
2	X	Q5HY64	FA47C_HUMAN	Putative protein FAM47C	4	3
2	X	Q5HYW3	RGAG4_HUMAN	Retrotransposon gag domain-containing protein 4	3	3
2	X	Q6PI77	BHLH9_HUMAN	Protein BHLH9	3	3
2	X	Q8IZF6	AGRG4_HUMAN	Adhesion G-protein coupled receptor G4	2	2
2	X	Q8N7E2	ZN645_HUMAN	E3 ubiquitin-protein ligase ZNF645	2	2

1 - 6 / 6 1



jPOST meets C-HPP



Chr No.	Leader
Chr. 1	Ping Xu
Chr. 2	Lydie Lane
Chr. 3	Takashi Kawamura
Chr. 4	Yu Ju Chen
Chr. 5	Peter Horvatovich
Chr. 6	Christoph Borchers
Chr. 7	Edouard Nice
Chr. 8	Pengyuan Yang
Chr. 9	Je-Yoel Cho
Chr.10	Joshua Labaer
Chr.11	Jong Shin Yoo
Chr.12	Ravi Sirdeshmukh
Chr.13	Young-Ki Paik
Chr.14	Charles Pineau
Chr.15	Gilberto B. Domont
Chr.16	Fernando Corrales
Chr.17	Gilbert S. Omenn
Chr.18	Alexander Archakov
Chr.19	György Marko-Varga
Chr.20	Siqi Liu
Chr.21	Albert Sickmann
Chr.22	Akhileshi Pandey
Chr. X	Yasushi Ishihama
Chr. Y	Ghasem Hosseini Salekdeh
Mitochondria	Andrea Urbani

Hunting Missing Proteins using jPOST



Rapid and Deep Profiling of Human Induced Pluripotent Stem Cell Proteome by One-shot NanoLC–MS/MS Analysis with Meter-scale Monolithic Silica Columns

Ryota Yamana,^{†,‡} Mio Iwasaki,^{†,‡} Masaki Wakabayashi,[†] Masato Nakagawa,[§] Shinya Yamanaka,[§] and Yasushi Ishihama^{*,†}

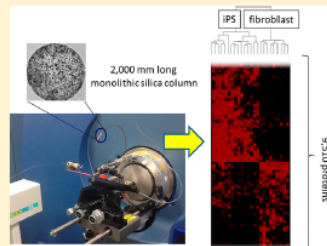
[†]Department of Molecular & Cellular BioAnalysis, Graduate School of Pharmaceutical Sciences, Kyoto University, Sakyo-ku, Kyoto 606-8501, Japan

[§]Center for iPS Cell Research and Application, Kyoto University, Sakyo-ku, Kyoto 606-8507, Japan

Supporting Information

ABSTRACT: Proteome analyses of human induced pluripotent stem cells (iPSC) were carried out on a liquid chromatography–tandem mass spectrometry system using meter-scale monolithic silica-C18 capillary columns without prefractionation. Tryptic peptides from five different iPSC lysates and three different fibroblast lysates (4 μg each) were directly injected onto a 200 cm long, 100 μm i.d. monolithic silica-C18 column and an 8-h gradient was applied at 500 nL/min at less than 20 MPa. We identified 98 977 nonredundant tryptic peptides from 9510 proteins (corresponding to 8712 genes), including low-abundance protein groups (such as 329 protein kinases) from triplicate measurements within 10 days. The obtained proteome profiles of the eight cell lysates were categorized into two groups, iPSC and fibroblast, by hierarchical cluster analysis. Further quantitative analysis based on an exponentially modified protein abundance index approach combined with UniProt keyword enrichment analysis revealed that the iPSC group contains more “transcription regulation”-related proteins, while the fibroblast group contained more “transport”-related proteins. Our results indicate that this simplified one-shot proteomics approach with long monolithic columns is advantageous for rapid, deep, sensitive, and reproducible proteome analysis.

KEYWORDS: *shotgun proteomics, monolithic silica column, iPS cell, one-shot proteomics*

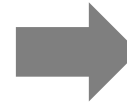


Special Issue: Chromosome-centric Human Proteome Project

Received: September 2, 2012

Published: December 4, 2012

dx.doi.org/10.1021/pr300837u | *J. Proteome Res.* 2013, 12, 214–221



Data list

Free word | **Ontology keyword**

human iPS

Project type

All Mass spectrometry Gel electrophoresis Antibody

Search | Reset



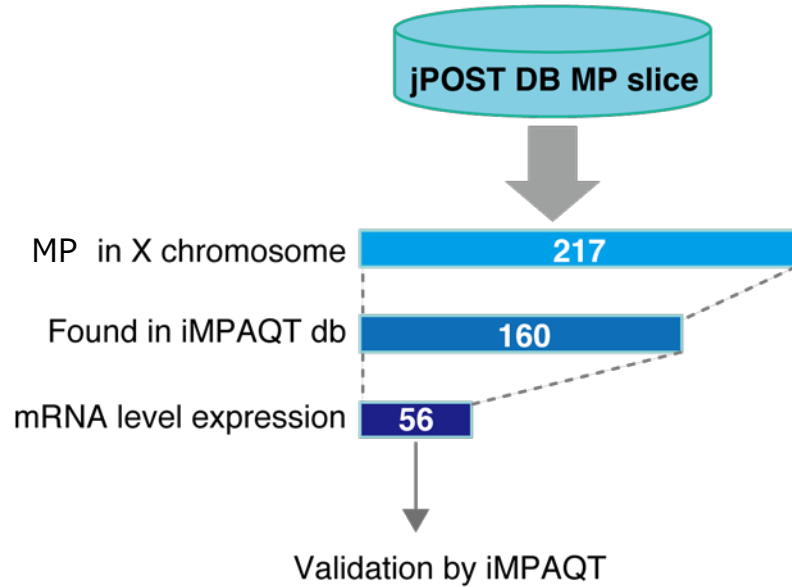
JPST000081	PXD004615	Human iPS cell_201B7-P32	Proteome analyses of human induced pluripotent ste ...	Complete	Ya Isl Ky ur
JPST000082	PXD004616	Human iPS cell_32R1-P32	Proteome analyses of human induced pluripotent ste ...	Complete	Ya Isl Ky ur
JPST000083	PXD004617	Human iPS cell_414C2-P43	Proteome analyses of human induced pluripotent ste ...	Complete	Ya Isl Ky ur
JPST000085	PXD004618	Human iPS cell_585A1-P55	Proteome analyses of human induced pluripotent ste ...	Complete	Ya Isl Ky ur
JPST000086	PXD004619	Human iPS 606A1-P46	Proteome analyses of human induced pluripotent ste ...	Complete	Ya Isl Ky ur
JPST000087	PXD004620	Human Fibroblast cell_aHDF1388-P9	Proteome analyses of human fibroblast cell line (a ...	Complete	Ya Isl Ky

Missing protein MRM transitions

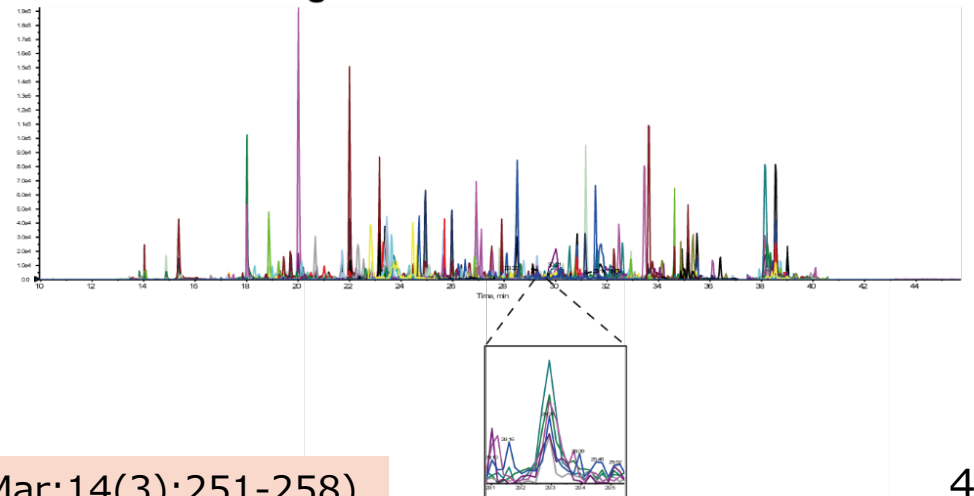
Gene ID list (Total number: 71)

Check all / Uncheck all

Gene ID	Symbol	Description	RefSeq	UniProt	PL1	Protein number (Accession/Protein/Conf/Natural)	HLA	WT	MISSING	TS	TSM	TSM	IMPQ
1266	CSNA	caudal type homeobox 4	NP_001338.1	G15627	FL051056	3 4	0.00	0.00	0.00	0.00	0.00	0.00	
1521	DMF2	dyadaplin-related protein 2	NP_001330.2	G12429	FL323310	25 30	1	0.00	0.00	2.00	1.10	2.47	
2556	GABRA3	gamma-aminobutyric acid (GABA) A receptor, alpha 3	NP_000705.1	P24953	FL222266	19 32	2.00	0.00	0.00	0.00	0.00	0.00	
2553	GPR	gastrin-releasing peptide receptor	NP_003305.1	P30550	FL261126	4 17	0.00	4.01	2.26	0.30	0.53		
4673	NAP1L3	nucleosome assembly protein 1-like 3	NP_003329.2	G20957	FL261126	8 31	0.00	0.75	0.11	0.00	0.00		
5628	PNOC	protein rich Gc (G-carboxypeptidase and) 1	NP_001189911.1	G14608	FL322256	5 13	1	6.77	4.06	3.10	4.00	3.42	
6120	NLGN3	neuroligin 3	NP_000955.1	P38255	Q91208010	1 2	26.17	1076.24	1262.41	1961.20	1466.72	1466.0000	
7526	ZNF780	Zinc Finger protein 780	NP_000962.2	P38215	Q91208099	31 40	1.10	2.41	1.77	1.61	1.63		
7712	ZNF157	Zinc Finger protein 157	NP_002407.2	P31786	Q91208080	8 44	0.00	0.00	0.00	0.10	0.00		
8223	SLC38A3	solute carrier family 38 (sodium/bic acid cotransporter family), member 3	NP_002807.1	P59131	FL380936	5 21	17.34	12.77	11.44	8.18	10.90		
8862	APLN	apelin	NP_005126.3	Q90411	FL385286	0 4	0.40	1.30	1.27	1.86	1.05		
9016	SLC35A14	solute carrier family 35 (inositol/serine carrier, member 14)	NP_012722.1	D62228	FL356116	2 4	3.81	5.47	3.79	5.55	6.73		
10060	CYSLTR1	cysteinyl leukotriene receptor 1	NP_000606.1	Q20213	FL262106	6 17	0.00	0.00	0.00	0.00	0.00		
10068	ZNF275	Zinc Finger protein 275	NP_001327962.2	P316728	FL216728	10 25	2.00	5.11	4.03	2.42	3.25		
12670	RRWD4	RNAi-templated, WW-domain 4	NP_026226.2	Q28959	FL222106	3 6	0.00	1.29	1.15	0.10	0.00		
20075	EGP16	EGF-like domain, multiple 6	NP_004322.2	Q91008	FL382306	14 40	0.00	0.00	0.00	0.00	0.00		
35290	ILKAPL2	integrin-like receptor accessory protein-like 2	NP_009112.1	Q28950	FL324356	16 21	0.00	0.00	0.00	0.00	0.00		



MRM chromatogram



Conclusions

- The jPOST repository, re-analysis and database have been successfully developed.
 - The jPOST repository is a part of ProteomeXchange.
 - Re-analysis is based on jPOST score, independent of search engines.
 - The jPOST team is involved in HPP, missing protein “next 50 challenge” projects.
- The jPOST scoring could be extended to proteomic analysis with wider search space such as proteogenomics and metaproteomics.

