

# In quest of SPARQLing lights in the dark human proteome

## Lydie Lane

CASSS Mass spectrometry symposium Long Beach, CA, Sept 28-30 2022





# Why is it important to characterize the human proteome?

- Proteins (not genes) are the active components of cells
- Defects in proteins can cause diseases (e.g. genetic diseases or cancers)
- Proteins can be drug targets
- Proteins can be biomarkers

# What is "the" human proteome?



### human genome

# ~ 20'000 protein coding-genes

nu

## human proteome ~ 5'000'000 different

proteoforms

#### post-translational modifications of proteins

(PTMs)

Ite, (3S)-3-hydroxyasparagine, 1'-histidyl-3'-tyrosine, 1nine, 3'-(S-cysteinyl)-tyrosine, 3-hydroxyproline, 34-carboxyglutamate, 4-hydroxyproline, 5-glutamyl, 5lroxylysine, 5-imidazolinone, ADP-ribosylasparagine,

ADP-ribosylcysteine, ADP-ribosylserine, Allysine, Arginine amide, Asparagine amide, Aspartate 1-(chondroitin 4-sulfate)-ester, Asymmetric dimethylarginine, Beta-decarboxylated aspartate, Cholesterol glycine ester, Citrulline, Cysteine methyl ester, Cysteine sulfenic acid, Cysteinyl-selenocysteine, Deamidated asparagine, Deamidated glutamine, Dimethylated arginine, Diphthamide, Disulfide bond, GPI-anchor amidated alanine, GPI-anchor amidated asparagine, GPI-anchor amidated aspartate, GPIanchor amidated cysteine, GPI-anchor amidated glycine, GPI-anchor amidated serine, Glutamic acid 1amide, Glutamine amide, Glycine amide, Glycyl adenylate, Glycyl lysine isopeptide, Hydroxyproline, Hydroxyproline, Hypusine, Isoglutamyl cysteine thioester, Isoglutamyl lysine isopeptide, Isoleucine amide, Leucine amide, Leucine methyl ester, Lysine amide, Lysine tyrosylquinone, Methionine amide, N,N,Ntrimethylalanine, N-acetylalanine, N-acetylaspartate, N-acetylcysteine, N-acetylglutamate, Nacetylglycine, N-acetylmethionine, N-acetylproline, N-acetylserine, N-acetylthreonine, N-acetylvaline, Nmyristoyl glycine, N-palmitoyl cysteine, N-palmitoyl glycine, N-pyruvate 2-iminyl-valine, N4,N4dimethylasparagine, N6,N6,N6-trimethyllysine, N6,N6-dimethyllysine, N6-(pyridoxal phosphate)lysine, N6-(retinylidene)lysine, N6-1-carboxyethyl lysine, N6-acetyllysine, N6-biotinyllysine, N6-carboxylysine, N6lipoyllysine, N6-methylated lysine, N6-methyllysine, N6-myristoyl lysine, Nitrated tyrosine, O-(pantetheine 4'-phosphoryl)serine, O-AMP-threonine, O-AMP-tyrosine, O-acetylserine, O-acetylthreonine, O-decanoyl serine, O-palmitoyl serine, Omega-N-methylarginine, Omega-N-methylated arginine, Omegahydroxyceramide glutamate ester, Phenylalanine amide, Phosphatidylethanolamine amidated glycine, Phosphohistidine, Phosphoserine, Phosphothreonine, Phosphothrosine, PolyADD-ribosyl alutamic acid

## Proline amide, Pyrrolidone cart FAD cysteine, S-Lysyl-methion cysteine, S-glutathionyl cystein + all the different glycosylation forms and the processing events !

cysteine, Sulfoserine, Sulfotyrosine, Symmetric aimetnyiarginine, Leie-Baipna-FAD nistiaine, Leiemethylhistidine, Thyroxine, Triiodothyronine, Tyrosine amide, Valine amide

alternative splicing of mRNA

2-5 fold increase

~ 50 to 100'000 transcripts (mRNAs)

+ all the individual variations (SNPs, disease mutations, etc)

# The main tasks of protein knowledgebases

- provide a reference set of protein sequences predicted from genome(s) analysis
- gather experimental validations for the predicted proteins
- keep information up to date about the predicted/confirmed functions of proteins and proteoforms in health and disease

# Curating the human proteome

## Known knowns

- Well characterized proteins
- Known PTMs and variants regulating function or localization of proteins

#### MANUAL CURATION OF PAPERS

#### Known unknowns

• Expected proteins

#### **GENOMIC/TRANSCRIPTOMIC DATA ANALYSIS**

• Uncharacterized proteins identified in human samples

#### **PROTEOMIC DATA ANALYSIS**

• Predicted PTMs

#### **AUTOMATIC PREDICTIONS**

• Genomic variants of unknown significance

**GENOMIC DATA ANALYSIS** 

# Tracking the known unknowns in the human proteome

- The Human Proteome Project (HPP) (<u>www.thehpp.org</u>), launched in 2010, is an international project organized by the Human Proteome Organization (HUPO) *that aims to revolutionize our understanding of the human proteome via a coordinated effort by many research laboratories around the world*.
- It is designed to map the entire human proteome in a systematic effort using currently available and emerging techniques. Completion of this project will enhance understanding of human biology at the cellular level and lay a foundation for development of diagnostic, prognostic, therapeutic, and preventive medical applications.
  - Goal 1: validate the existence of 1 protein product per gene
  - Goal 2: find at least one biological function for every gene

Human Proteome Project





## neXtProt, the human protein knowledge platform

- neXtProt (<u>www.nextprot.org</u>) was created in 2011 in order to organize human protein knowledge
- It has been chosen as the reference database for the HUPO Human Proteome Project in 2013
- neXtProt was the first resource to implement the HUPO Protein Standard Initiative (PSI) PEFF format, an enriched FASTA format allowing MS search engines and other tools to easily and consistently access sequence variations and PTM data.



#### neXtProt: Organizing Protein Knowledge in the Context of Human Proteome Projects

Pascale Gaudet,<sup>†</sup> Ghislaine Argoud-Puy,<sup>†</sup> Isabelle Cusin,<sup>†</sup> Paula Duek,<sup>†</sup> Olivier Evalet,<sup>†</sup> Alain Gateau,<sup>†</sup> Anne Gleizes,<sup>†</sup> Mario Pereira,<sup>†</sup> Monique Zahn-Zabal,<sup>†</sup> Catherine Zwahlen,<sup>†</sup> Amos Bairoch,<sup>†,‡</sup> and Lydie Lane<sup>\*,†,‡</sup>

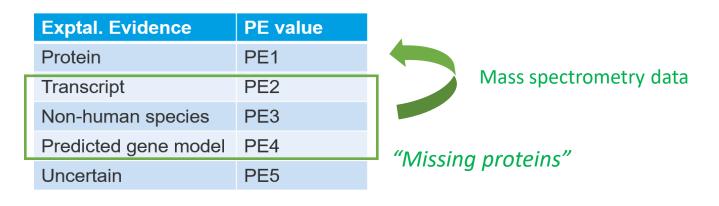
# neXtProt, the human protein knowledge platform

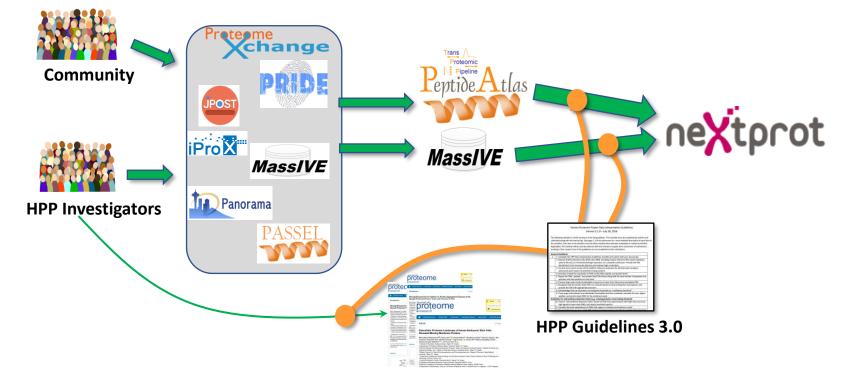
Huge amounts of data on the ~20,300 predicted protein-coding human genes constantly being generated by the life science community

[genomic variations in health and disease, mRNA expression in tissues, protein identification by antibody-based or mass spectrometry-based techniques, protein-protein interactions, PTMs, etc.]

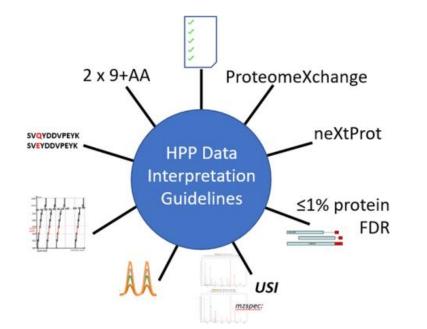


# HPP Goal 1: validate the existence of 1 protein product per gene





# Guidelines for Mass Spectrometry Data Interpretation



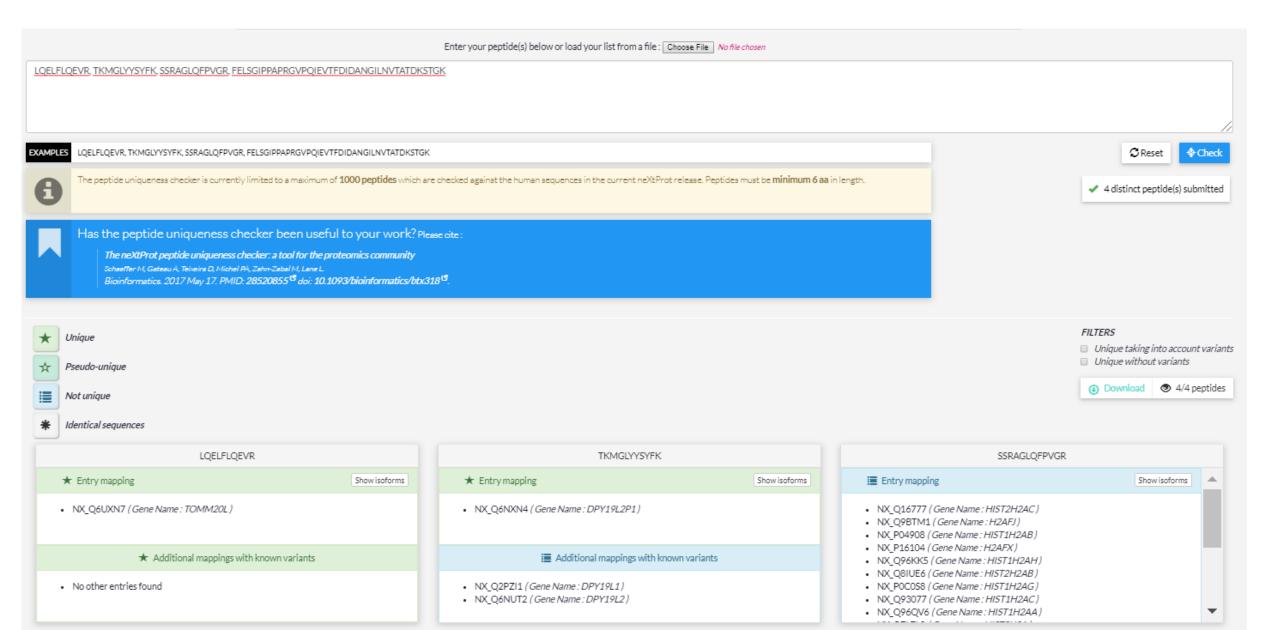
#### Deutsch et al, J. Proteome Res. 2019, 18, 12, 4108-4116

intormation (pizzzo for page iz fine zo, siz for supplementary table z, etc.)

nom	ation	(ארביבט וטו אמצי דב ווויד בט, דוב וטו געאאוידווידוומוץ נמטוב ב, בנכ.)						
Gene	eral gu	idelines for all manuscripts:						
٧	Loc	1. Complete this HPP MS Data Interpretation Guidelines checklist and submit with your manuscript.						
2. D	ata de	position guidelines						
		2a. Deposit all MS proteomics data to a ProteomeXchange repository as a "complete" submission.						
		2b. Include analysis reference files (search database, spectral library, transition list, etc.) in submission.						
		2c. Provide the PXD identifier(s) in the manuscript abstract.						
		2d. Provide the reviewer login credentials if the dataset is not yet public.						
		3. Use the most recent version of the neXtProt reference proteome for all informatics analyses,						
		particularly with respect to new PE1 protein detection claims.						
4. FI	DR-rel	ated guidelines						
		4a. Describe in detail the calculation of FDRs at the PSM, peptide, and protein levels.						
		4b. Report the PSM-, peptide-, and protein-level FDR values along with the total number of expected fals						
		positives at each level, using precision appropriate to the uncertainty in computed FDR.						
		4c. Present large-scale results thresholded at equal to or lower than 1% protein-level global FDR.						
		4d. If any large-scale datasets are individually thresholded and then combined, calculate the new, higher						
		peptide- and protein-level FDRs for the combined result.						
Guid	elines	for claims of new PE1 protein detections (i.e., presenting evidence to categorize a protein to PE1)						
		5a. If using DDA mass spectrometry for such claims, present high mass-accuracy, high signal-to-noise ratio						
		(SNR), and clearly annotated spectra. Scrutinize spectra for missing and extra peaks.						
		5b. Present high mass-accuracy, high-SNR, clearly annotated spectra of synthetic peptides that match the						
		spectra supporting the claims. Peptides from recombinant proteins are acceptable synthetics.						
		5c. Provide Universal Spectrum Identifiers (USIs) for all natural and synthetic peptide spectra that support						
		such claims, ideally as a supplementary data table.						
		6. If using SRM verification for such claims, present target traces alongside synthetic heavy-labeled peptid						
		traces, demonstrating co-elution and closely matching fragment mass intensity patterns.						
		7. If using DIA MS, then, if the data are analyzed with XICs, apply the above SRM guidelines (6); if the data						
		are analyzed by extracting deconvoluted spectra, apply the above DDA guidelines 5a-5c.						
		8. Even when very high confidence peptide identifications are demonstrated, consider alternate mappings						
		of the peptide to proteins other than the claimed one. Consider isobaric sequence/mass modification						
		variants, all known SAAVs, and unreported SAAVs.						
		9. Support such claims by two or more distinct uniquely-mapping, non-nested peptide sequences of lengtl						
		≥9 amino acids with the above evidence in the same paper. When 2 peptides overlap, the total extent						
		must be ≥18 amino acids. When weaker evidence is offered for such a claim, justify that other peptide						
		cannot be expected by any common digestion proteases. When 2 or more proteins are *exactly*						
		sequence identical (irrespective of SAAVs), peptides are considered uniquely mapping if they map only						
		to the group, and such proteins will have the same PE level and be counted.						
+ +		ments (use this snace and extra nages to explain any nonadherence $[NA/NC]$ in the above checklist):						

Author comments (use this space and extra pages to explain any nonadherence [NA/NC] in the above checklist):

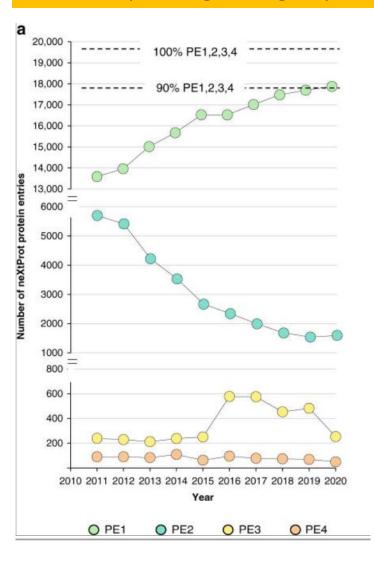
## The neXtProt peptide uniqueness checker



## Validation of "missing proteins" 2011-2021

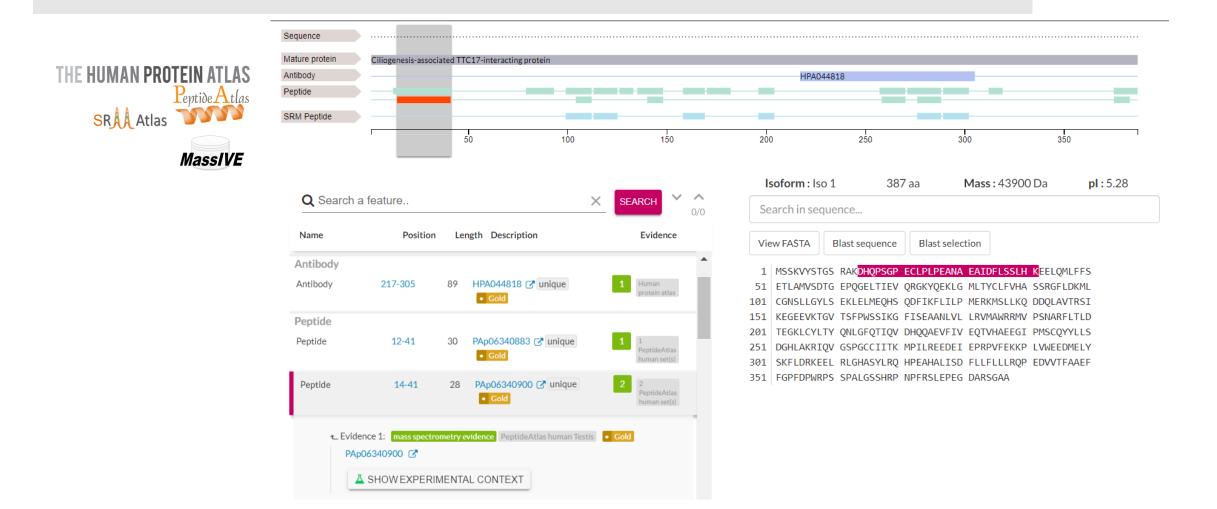
nature communications
PERSPECTIVE OPEN Mttps://doi.org/10.1038/s41467-020-19045-9 OPEN A high-stringency blueprint of the human proteome Subash Adhikari et al.#
Release 2021-02-15
Protein existence in neXtProt
Predicted (9) Uncertain (601) Inferred from homology (147) EVIdence at transcript level (1265)
Release statistics

#### 90.4% complete high-stringency human proteome



# Display of proteomics information in neXtProt

- ✓ Which (unique) MS/MS peptides from my protein have been found, and where ?
- ✓ What peptides can I use for targeted proteomics studies?



# Are there still missing proteins to validate?

1343 proteins are still "missing" (PE2-4) in neXtProt

- Difficult to detect with conventional mass spectrometry?
  - Low abundant proteins, or with a expression pattern restricted in time and/or space (incl. 392 olfactory/taste receptors)
  - Highly hydrophobic or hard to extract from cells?
  - Digestion with usual enzymes not leading to 2 unique peptides?
- Wrong predictions?

## **Potential solutions**

- Use different samples (selected using RNAseq data)
- Improve sample preparation
- Design targeted assays, use affinity-enrichment approaches
- Improve the sensitivity of the detection
- Use different digestion enzymes
- Adapt validation criteria?

## The neXtProt protein digestion tool

Enter neXtProt isoform accession number, i.e NX_P50222-1 NX_Q9HC47-1						
Max miscleavages O				Min peptide length 9	Max peptide length 40	RESET DIGEST
DIGESTED PEPTIDES FOR EACH PROTEASE :						
Copy CSV Excel Print						
Showing 1 to 27 of 27 entries 1 row selected						Search:
Protease name				It Peptide count	11 Unique peptide count	Į£
GLU_C_BICARBONATE				0	0	4
				0	0	
PEPSIN_PH_1_3				0	0	
PEPSIN_PH_GT_2				0	0	
PROTEINASE_K				0	0	
THERMOLYSIN				0	0	
TRYPSIN				0	0	
ASP_N				2	2	
CHYMOTRYPSIN_HIGH_SPEC				2	2	
GLU_C_PHOSPHATE				2	2	
Select a protease <b>Q</b> , GLU_C_PHOSPHATE						GET PEPTIDE
Copy CSV Excel Print Showing 1 to 2 of 2 entries						Search:
Peptide sequence	.↓↑ Length	↓↑ Missed cleavages	Position	L Unique without variants	1 Natural in neXtProt	Synthetic in neXtProt
Filter Q	32	37 0	1	Filter Q	Filter Q	Filter Q
MFVIISLHNCVVISFVLFLFGGNNFIQNFYLPQNYID		37	0 1-37	Yes	No	No
QFLLTSFPTFTSVGVLIVLVLCSAFLLLWQGE		32	0 38-69	Yes	No	No

#### Missing proteins



## Example: TMEM213 (2 tryptic peptides)

no enzyme generating at least 2 theoretical 9-50 aa peptides

at least 2 unique theoretical 9-50 aa peptides with another enzyme than trypsin

at least 2 unique theoretical 9-50 aa tryptic peptides

#### $TMEM213 \rightarrow Transmembrane protein 213$

Gene name : TMEM213

Entry whose protein(s) existence is based on evidence at transcript level.

Annotations in this section apply to all the isoforms if not specified otherwise.

**EXPRESSION** 

Expression in kidney and salivary gland. 💿 Gold

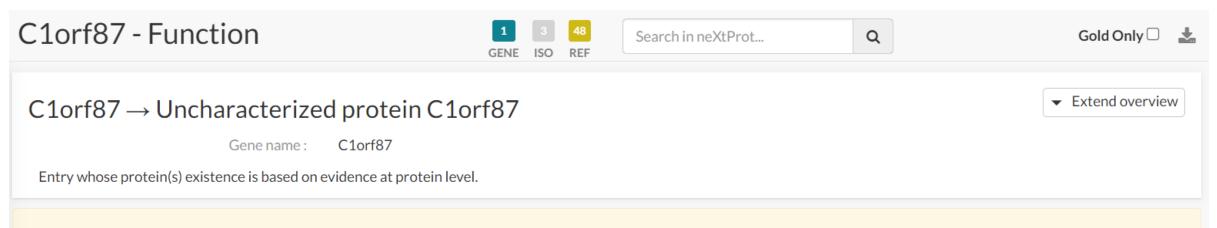
RNA tissue specificity: Group enriched (kidney, salivary gland). RNA tissue distribution: Detected in some. • Gold

Single cell type specificity: Cell type enhanced (Collecting duct cells, Distal tubular cells, Melanocytes). • Gold

# HPP Goal 2: "A Function for Every Protein"

- Most research on human genes only concentrates on approximately 2,000 of the 19,000 genes of the human genome

- In neXtProt 1191 proteins are validated at protein level but do not have any predicted or experimentally validated function



There is no function information available in neXtProt for NX\_Q8N0U7. Suggestions for updates and submission of functional predictions are welcome! Please contact us

# The CALIPHO laboratory (2009-2018)

Classical approach (hypothesis-driven, single gene)

With a focus on:

- Atypical (small) mitochondrial proteins
- Proteins with a typical "ciliary" profile
- Enzymatic pathway holes



Function prediction by data analysis

1 PhD or 1 post doc for 3-4 years for 1 gene

Gene knock-out in human cell lines or model organisms (zebrafish) Assay of the predicted function

# The CALIPHO laboratory (2009-2018)

## PLOS ONE

🔓 OPEN ACCESS 度 PEER-REVIEWED

RESEARCH ARTICLE

#### Functional Identification of APIP as Human mtnB, a Key Enzyme in the Methionine Salvage Pathway

Camille Mary 🔄, Paula Duek, Lisa Salleron, Petra Tienz, Dirk Bumann, Amos Bairoch, Lydie Lane

Published: December 28, 2012 • https://doi.org/10.1371/journal.pone.0052877



Contents lists available at ScienceDirect

Biochimica et Biophysica Acta

#### DERA is the human deoxyribose phosphate aldolase and is involved in stress response

Lisa Salleron <sup>a,\*</sup>, Giovanni Magistrelli <sup>b</sup>, Camille Mary <sup>a</sup>, Nicolas Fischer <sup>b</sup>, Amos Bairoch <sup>a,c</sup>, Lydie Lane <sup>a,c,\*\*</sup>

<sup>a</sup> Department of Human Protein Sciences, Faculty of Medicine, University of Geneva, Geneva, Switzerland <sup>b</sup> NovImmune SA, Plan-les-Ouates, Switzerland

<sup>c</sup> CALIPHO GroupSIB-Swiss Institute of Bioinformatics, Geneva, Switzerland

## PLOS ONE

🔓 OPEN ACCESS 💈 PEER-REVIEWED

RESEARCH ARTICLE

### C2orf62 and TTC17 Are Involved in Actin Organization and Ciliogenesis in Zebrafish and Human

Franck Bontems D, Richard J. Fish, Irene Borlat, Frédérique Lembo, Sophie Chocu, Frédéric Chalmel, Jean-Paul Borg, Charles Pineau, Marguerite Neerman-Arbez, Amos Bairoch, Lydie Lane D

Published: January 27, 2014 • https://doi.org/10.1371/journal.pone.0086476



C11orf83, a Mitochondrial Cardiolipin-Binding Protein Involved in  $bc_1$  Complex Assembly and Supercomplex Stabilization

Marjorie Desmurs,<sup>a</sup> Michelangelo Foti,<sup>b</sup> Etienne Raemy,<sup>c</sup> Frédéric Maxime Vaz,<sup>d</sup> Jean-Claude Martinou,<sup>c</sup> Amos Bairoch,<sup>a.e</sup> Lydie Lane<sup>a,e</sup>



ORIGINAL RESEARCH published: 16 March 2021 doi: 10.3389/fncel.2021.653075



#### ARTICLE https://doi.org/10.1038/s41467-020-14999-2 OPEN

## Mitochondrial peptide BRAWNIN is essential for vertebrate respiratory complex III assembly

Shan Zhang<sup>1</sup>, Boris Relijć<sup>6</sup>, <sup>2,12</sup>, Chao Liang<sup>6</sup>, <sup>12</sup>, Baptiste Kerouanton<sup>112</sup>, Joel Celio Francisco<sup>3</sup>, Jih Hou Peh<sup>1</sup>, Camille Mary<sup>6</sup>, <sup>4</sup>, Narendra Suhas Jagamathan<sup>6</sup>, <sup>5</sup>, Volodimir Olexiouk<sup>6</sup>, Claire Tang<sup>1</sup>, Gio Fidelito<sup>1</sup>, Srikanth Nama<sup>7</sup>, Ruey-Kuang Cheng<sup>8</sup>, Caroline Lei Wee<sup>6</sup>, Loo Chien Wang<sup>9</sup>, Paula Duek Roggli<sup>6</sup>, <sup>10</sup>, Prabha Sampath<sup>6</sup>, <sup>11</sup>, Lydie Lane<sup>6</sup>, Enrico Petretto<sup>5</sup>, Radoslaw M. Sobota<sup>6</sup>, Suresh Jesuthasan<sup>6</sup>, <sup>89</sup>, Lisa Tucker-Kellogg<sup>6</sup>, <sup>35</sup>, Bruno Reversade<sup>6</sup>, <sup>19</sup>, Gerben Menschaert<sup>6</sup>, Lei Sun<sup>6</sup>, David A. Stroud<sup>6</sup>, <sup>2</sup> & Lena Ho<sup>6</sup>, <sup>178</sup>

#### C21orf91 Regulates Oligodendroglial Precursor Cell Fate—A Switch in the Glial Lineage?

Laura Reiche<sup>1</sup>, Peter Göttle<sup>1</sup>, Lydie Lane<sup>2,3</sup>, Paula Duek<sup>2,3</sup>, Mina Park<sup>1</sup>, Kasum Azim<sup>1</sup>, Jana Schütte<sup>1</sup>, Anastasia Manousi<sup>1</sup>, Jessica Schira-Heinen<sup>1</sup> and Patrick Küry<sup>1+</sup>

## Lessons learned:

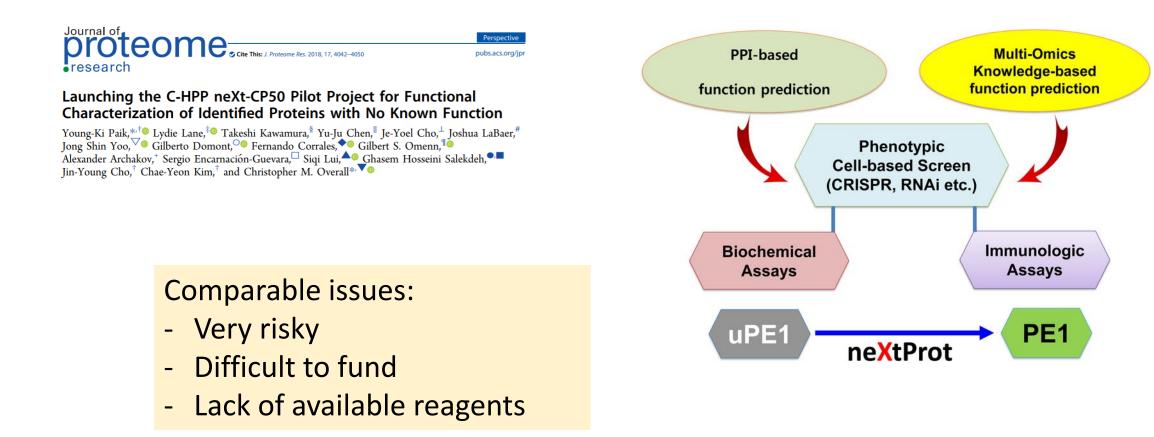
- Very risky

Check for updates

- Difficult to fund
- Lack of available reagents

# The HPP next-CP50 initiative (2018-2021)

On March 1, 2018, HUPO C-HPP announced launching the neXt-CP50 where CP stands for "characterization of proteins". This pilot project aims to **characterize the function of 50 identified proteins with unknown function** during ~3 years (2018-2021). This challenge is to test the feasibility of the functional characterization of large numbers of dark proteins



# The Human Proteome Grand Project 2022-2032

#### aims to:

- Build upon the large resources constructed through the significant efforts of the HPP
- Understand the Proteome in the context of Networks of Networks
- o Partner with complementary consortia to expand the reach of the HPP and increase its relevance
- Develop forward thinking strategies to provide functional information of each protein
- Be a focal point for translational researchers for mechanisms, diagnostics, and therapeutics to improve human health

### is:

- Open to all interested groups or individuals
- $\circ~$  Agnostic to disease and biological context
- $\circ$  will:

Utilize resources already familiar to the research groups

 $\circ~$  Advance the knowledge of systems being worked on through the access to resources to perturb these systems



## Network-based exploratory approach:

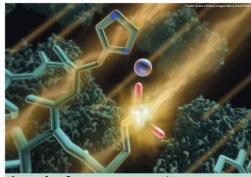


#### WHY TARGET 2035

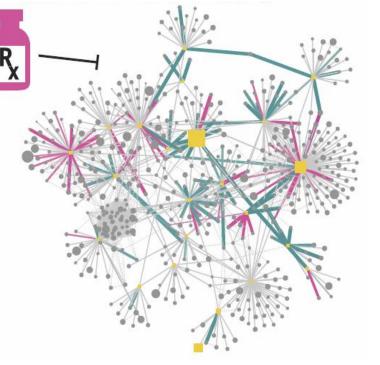
Almost 20 years after the human genome was sequenced, many of the genes linked to disease phenotypes or those associated with specific disease traits by genome-wide association studies remain severely understudied: the so-called 'dark genome'.

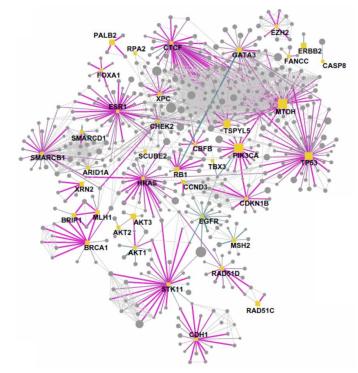
Generating chemical and biological reagents for these understudied proteins will enable research, lead to the validation of novel therapeutic targets, and the discovery of better medicines.

#### **NEWS & ANALYSIS**



A probe for every protein

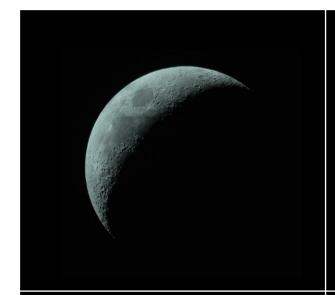




Define interactors and perturbers (variations amongst individuals, drugs...) to understand impact and function

**Comparative mapping** 

# The unknown unknowns in the human proteome



### Known unknowns

- Expected proteins
- Uncharacterized proteins identified in human samples
- Predicted PTMs
- Variants of unknown significance

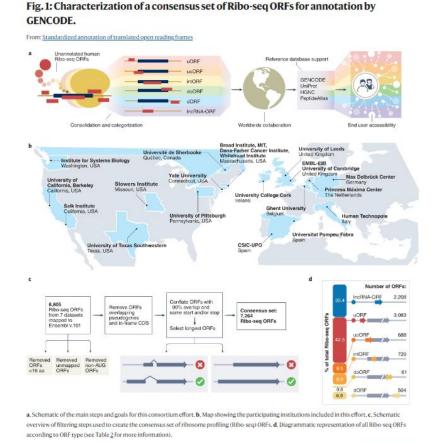
#### Unknown unknowns

- Are there coding elements in the genome outside the 1.2% predicted (e.g. smORFs and lncRNAs)?
- Novel functions to discover even for well-known proteins (one protein can have multiple functions!)
- Novel PTMs to discover
- · ...?

## How to track the unknown unknowns?

Back to article page >

 Validation of novel coding elements will require new reference databases and new guidelines



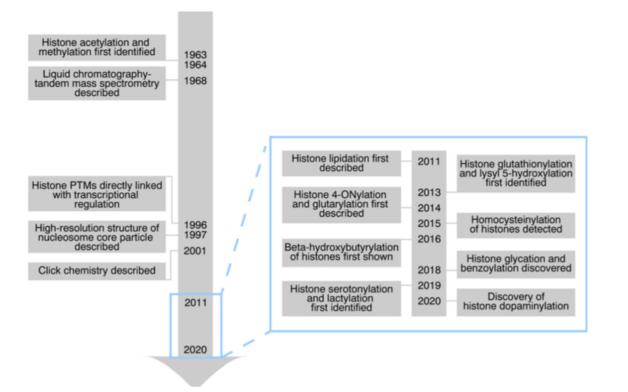
Example: A catalog of 7264 ORFs, some as short as 16 aa in length, with substantial evidence of ribosome translation activity based on Ribo-seq experiments currently analysed for MS evidence by PeptideAtlas

Mudge JM, Ruiz-Orera J, Prensner JR, Brunet MA, Calvet F, Jungreis I, Gonzalez JM, Magrane M, Martinez TF, Schulz JF, Yang YT, Albà MM, Aspden JL, Baranov PV, Bazzini AA, Bruford E, Martin MJ, Calviello L, Carvunis AR, Chen J, Couso JP, Deutsch EW, Flicek P, Frankish A, Gerstein M, Hubner N, Ingolia NT, Kellis M, Menschaert G, Moritz RL, Ohler U, Roucou X, Saghatelian A, Weissman JS, van Heesch S. Standardized annotation of translated open reading frames. Nat Biotechnol. 2022 Jul;40(7):994-999. doi: 10.1038/s41587-022-01369-0. PMID: 35831657.

Nature Biotechnology (Nat Biotechnol) | ISSN 1546-1696 (online) | ISSN 1087-0156 (print)

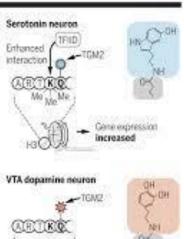
## How to track the unknown unknowns?

New PTMs and new functions are constantly discovered



Monoaminylation of histone H3

Serotonylated glutamine (Q) 5 in histone H3 during serotonergic neuron differentiation is catalyzed by transglutaminase 2 (TGM2) and associated with trimethylation of lysine 4 (K4me3). Serotonylation enhances general transcription factor ID (TFIID) binding to K4me3 and facilitates transcription. Oppaminylated Q5 in histone H3 is initially decreased during cocaine. withdrawal and then increased. This facilitates withdrawal-induced gene. expression alteration in ventral tegmental area (VTA) neurons and enhances dopaminergic neuron escitability and drug-seeking behavior in rats. A, alanine; R. arginine: T. threonine.



Gene expression

altered

Selected Milestones in Histone Post-Translational Modification (PTM) Discovery.

Chan JC, Maze I. Nothing Is Yet Set in (Hi)stone: Novel Post-Translational Modifications Regulating Chromatin Function. Trends in Biochemical Sciences. 2020 Oct;45(10):829-844. DOI: 10.1016/j.tibs.2020.05.009. PMID: 32498971; PMCID: PMC7502514.

## These newly discovered PTMs are now found in neXtProt



# Are there unknown knowns in the human proteome?

## Known knowns

- Well characterized proteins
- Known PTMs and variants regulating function or localization of proteins

## Unknown knowns

- Unannotated characterization data
- Unpublished/hidden characterization data
- Data from large scale studies that could be combined into functional knowledge

## Unannotated and unpublished data

- Curators @neXtProt @UniProt @GOA @KEGG @Reactome monitor the literature and update proteins' functional information every day, but there might be some gaps/delays!
- -> Don't hesitate to send update requests if you spot missing annotations

Unfortunately, negative results or results difficult to interpret are often unpublished

- In the last five years, we estimate that 4 5 papers describe newly characterized human proteins each month
- -> ~20 years to be completed! How can this be accelerated?

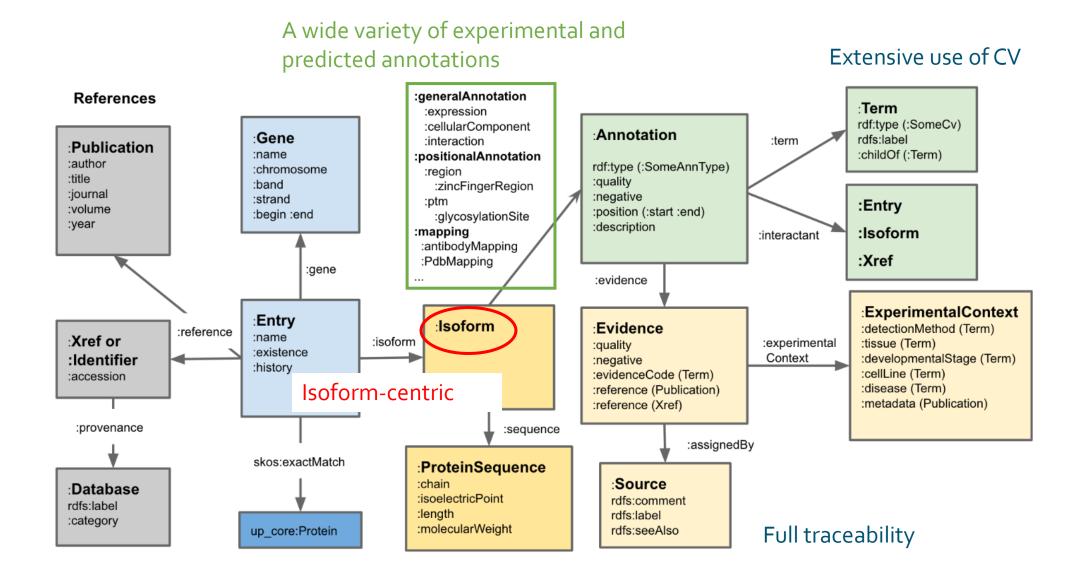
# Combining scattered data from large scale studies

Most proteins with no function annotated do have:

- Expression data
- Subcellular location information
- PPI data
- Orthologs in model organisms
- Predicted 3D structures

-> Can we retrieve this information from neXtProt, combine it, and deduce possible function(s)?

## www.nextprot.org RDF data model



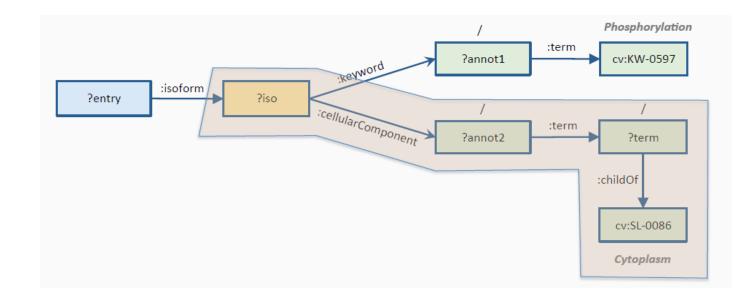
## RDF data can be queried using SPARQL queries

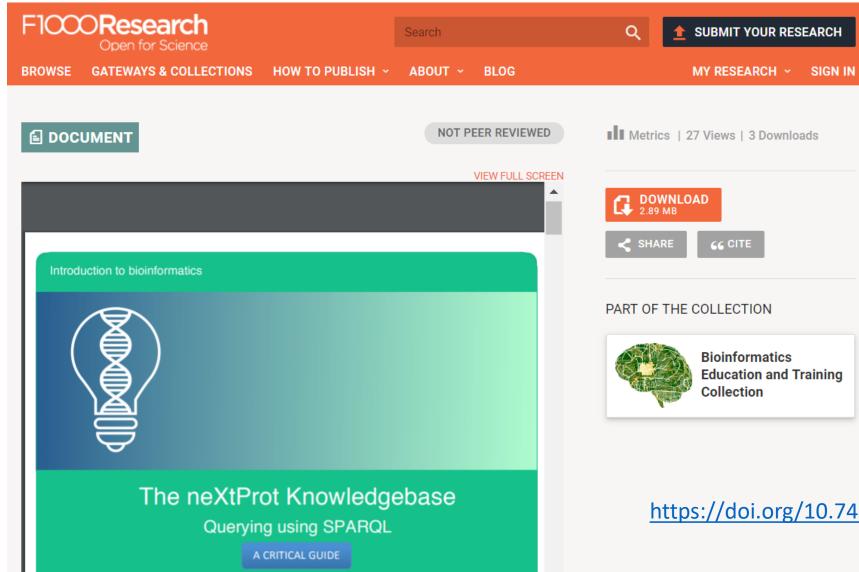
```
Simple search • Advanced search
  1 #Proteins phosphorylated and located in
                                                   +
                                                      X
                                                               Q Search
  2 #the cytoplasm
  3
    select distinct ?entry where {
  4
      ?entry :isoform ?iso.
  5
      ?iso :keyword / :term cv:KW-0597.
  6
      ?iso :cellularComponent /:term /:childOf cv:SL-0086.
  8
```

SPARQL queries consist of: 1) The query result variables – these will be shown when the query is executed; 2) the query pattern – query constructs defining the data selection, and generating any variables required; and 3) the optional query modifiers – used to modify the results.

SELECT ... query result variables ... WHERE {
... query pattern ...

... optional query modifiers ...





#### https://doi.org/10.7490/f1000research.1116829.1

## Almost 200 premade queries to explore the human proteome

.g:Sea	rch for MSH6 in proteins, Search for author Doolittle in publications, Search for liver in terms		
Deta	Tags - search in 145 queries (ex: liver)	New Quer	У
	Proteins phosphorylated and located in the cytoplasm	NXQ_00001	
	Proteins that are located in both the nucleus and in the cytoplasm	NXQ_00002	>
	Proteins with 7 transmembrane regions	NXQ_00003	>
	Proteins expressed in brain with IHC expression level: "high" but not expressed in testis	NXQ_00004	>
	Proteins located in mitochondrion and that lack a transit peptide	NXQ_00005	>
	Proteins whose genes are on chromosome 13 and are associated with a disease	NXQ_00006	>
	Proteins associated with diseases that are associated with cardiovascular aspects	NXQ_00007	>
	Proteins whose genes are less than 50000 bp away from the location of the gene coding for protein p53	NXQ_00008	>
	Proteins with 3 disulfide bonds and that are not annotated as hormones	NXQ_00009	•

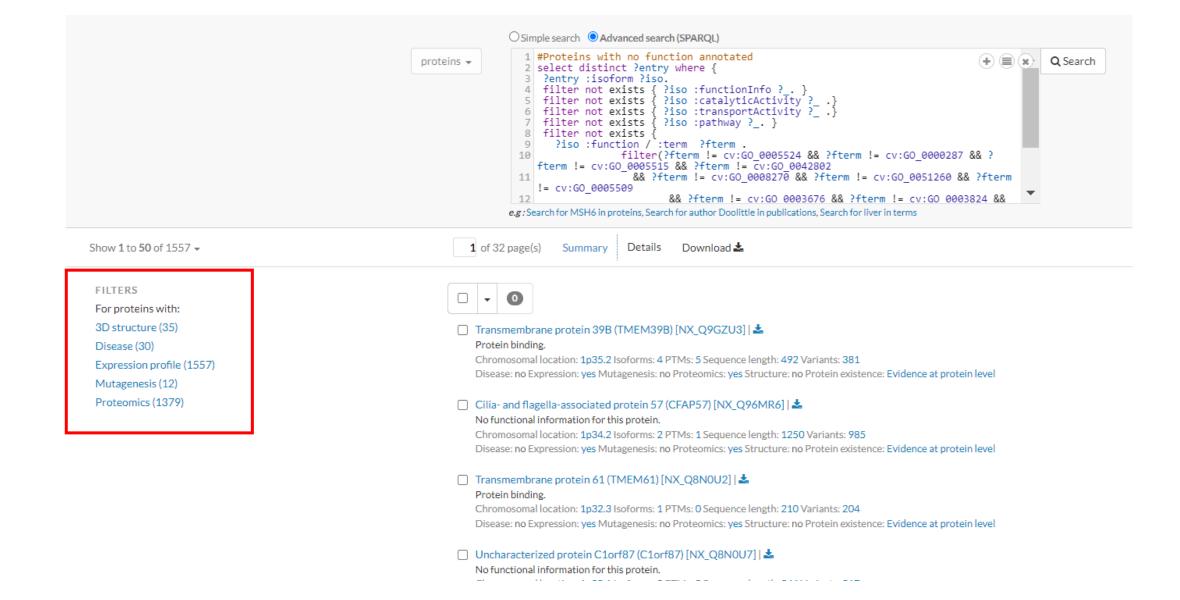
Missing proteins (NXQ\_00204) Proteins with no function (NXQ\_00022) Proteins with specific expression patterns (NXQ\_00004) Proteins associated with diseases (NXQ\_00007) Proteins with specific PTMs...

## A detailed help on neXtProt RDF entities is available

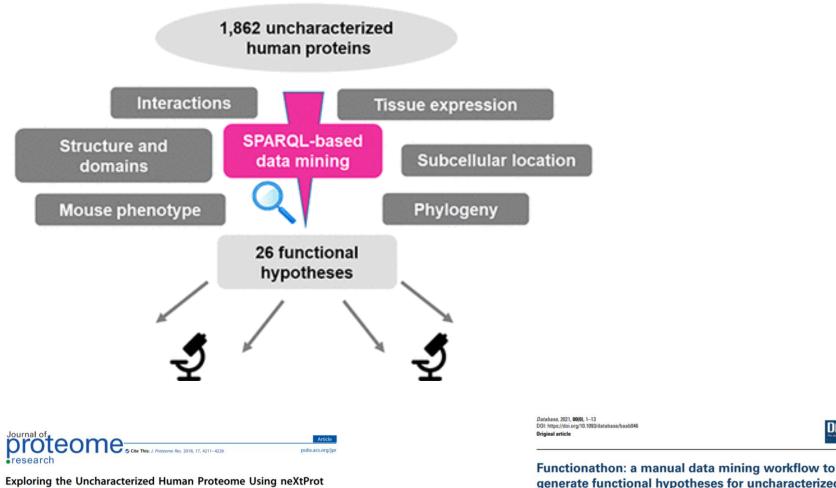
NonTerminalResidue (323)	PeptideMapping Values- 7727903				
NucleotidePhosphateBindingRegion (7109)					
ObservedExpression (5)	Peptide detected in a biological sample using mass spectrometry (MS).				
OmimCv (2634)	:Isoform :mapping	:PeptideMapping			
ORFName (2637)	:Isoform :peptideMapping	.герциемаррінд			
OrganelleCv (2)	:Isoform :positionalAnnotation	:end xsd:integer 7727903 Example: 1004			
Pathway (136835)		:entryAnnotationId xsd:string 7727903 Example: AN_NX_Q96Q:			
PdbMapping (120125)		:evidence :Evidence 63586484			
PeptideMapping (7727903)		:peptideName xsd:string 7727903 Example:NX_PEPT0088			
PeroxisomeTransitPeptide (4)		:peptideSource :Source 11636350 :peptideSource owl:NamedIndividual 11636350			
Person (4470928)		:peptideSource owl:Thing 11636350			
Pharmaceutical (89)		:peptideUniqueness xsd:string 7727903 NOT_UNIQUE V			
PhDependence (583)		:proteotypic xsd:boolean 7727903 false 🗸			
PhenotypicVariation (77044)		:quality :QualityQualifier 7727903 :quality owl:NamedIndividual 7727903			
Propeptide (1384)		:quality owl:Thing 7727903 :start xsd:integer 7727903 Example:1051			
ProteinExistence (5)		Example: 1051			
ProteinName (51868)	Evample				

https://snorql.nextprot.org/help/doc/introduction

## Proteins with no function annotated (NXQ\_00022)

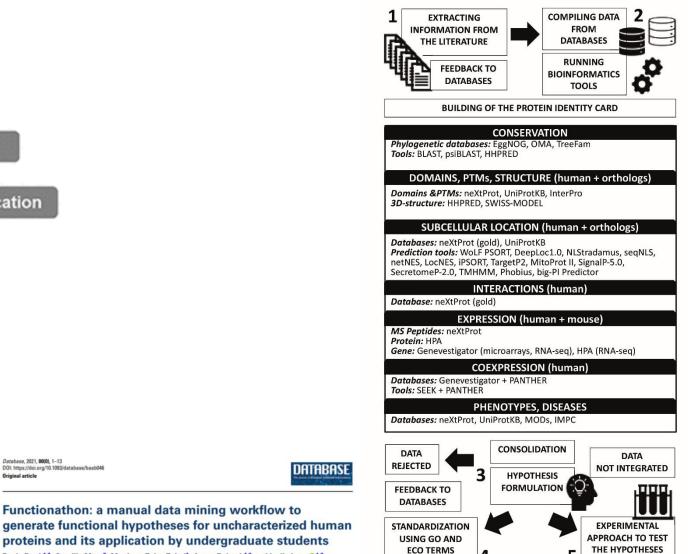


## Manual datamining on a selection of uncharacterized proteins



Paula Duek,<sup>†</sup> Alain Gateau,<sup>†</sup> Amos Bairoch,<sup>†,‡</sup> and Lydie Lane<sup>\*,†,‡</sup>

<sup>†</sup>CALIPHO Group, SIB-Swiss Institute of Bioinformatics, and <sup>‡</sup>Department of Microbiology and Molecular Medicine, Faculty of Medicine, University of Geneva, CMU, Michel-Servet 1, 1211 Geneva 4, Switzerland



Paula Duek<sup>1,2</sup>, Camille Mary<sup>2</sup>, Monique Zahn-Zabal<sup>1</sup>, Amos Bairoch<sup>1,2</sup> and Lydie Lane<sup>01,2,\*</sup>

# Standardization of the hypotheses as GO/ECO terms and integration in the neXtProt function prediction pages

ne <mark>X</mark> tprot	Tools → Portals → Download → Help → About → Contact	Login Expasy 🃩
NX_Q9BPX7	C7orf25 - Function predictions	Gold Only 🗆 🛛 🛓
Localization		
Sequence	$C7orf25 \rightarrow UPF0415$ protein C7orf25	✓ Extend overview
Proteomics	Gene name : C7orf25	
Structures	Family name : UPF0415	
Peptides	Entry whose protein(s) existence is based on evidence at protein level.	
Phenotypes		
Exons	Predictions in this section apply to all the isoforms if not specified otherwise.	▲ Hide evidences
Identifiers	GO MOLECULAR FUNCTION	1
REFERENCES	ribonuclease activity GO:0004540 Definition	L ev
Curated publications 36	Evidence 1: Match to InterPro member signature evidence used in manual assertion     Submitter:      10 0000-0002-5734-0298      10 0000-0002-0819-0473      10 0000-0002-9818-3030      2	
Additional publications	GO BIOLOGICAL PROCESS	
Patents 0	mRNA metabolic process GO:0016071 Definition	2
Submissions 1	<sup>t</sup> Evidence 1: Physical interaction evidence used in manual assertion	2 ev
Web resources	Submitter: (b) 0000-0002-5734-0298 C (b) 0000-0002-0819-0473 C (b) 0000-0002-9818-3030 C	
😁 COMMUNITY 🗸	<sup>t</sup> Evidence 2: Gene expression similarity evidence used in manual assertion	
C7orf25 Function predictions	Submitter: 🔞 0000-0002-5734-0298 🗷 🍈 0000-0002-0819-0473 🗹 🍺 0000-0002-9818-3030 🖉	
C7orf25 I-TASSER/COFACTOR	tRNA metabolic process GO:0006399 Definition	1
C7orf25 SAAVpedia	Evidence 1: Physical interaction evidence used in manual assertion	ev
C7orf25 3D structure	Submitter: (b) 0000-0002-5734-0298 C (b) 0000-0002-0819-0473 C (b) 0000-0002-9818-3030 C	

- We currently provide predictions for **204 uncharacterized proteins** in the neXtProt community pages
- Feel free to give us some feedback us and/or start experiments to test them!
- Publish your own hypotheses (derived from manual or automated tools)

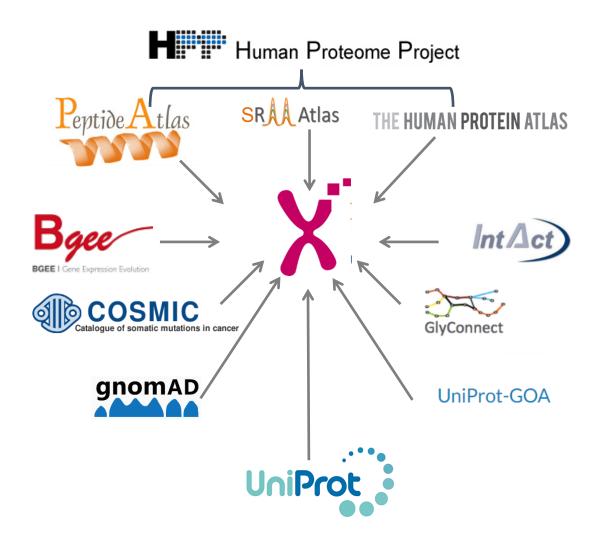
Please use the **Contact** option in neXtProt & provide:

- 1. Your **ORCID** *but you can choose not to have it displayed*
- 2. Protein **entry AC** *from query NXQ\_00022*
- **3. GO term** *describing molecular function* / *biological process*
- **4. ECO term** *describing the type of supporting evidence*
- **5. PubMed ID** for publication with the function prediction



https://www.nextprot.org/about/functional-proteome-project

## www.nextprot.org does not contain everything....



- ✓ MS-data
- ✓ Antibody-based data
- ✓ RNA-seq data
- ✓ PTMs
- ✓ Variants (genomic and somatic)
- ✓ PPIs
- ✓ Functional annotations
- Phylogenetic information
- Data on model organisms
- o Pharmacology data
- Clinical proteogenomic data
- $\circ~$  Structural data
- $\circ~$  Protein interactions with pathogens
- 0 ...

Extending neXtProt contents and empowering network-based approaches by federating other life science SPARQL endpoints

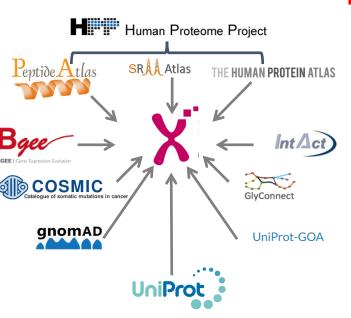
## Pharmacology/Clinical data



Systems biology



WIKIPATHWAYS Pathways for the People





Structural biology

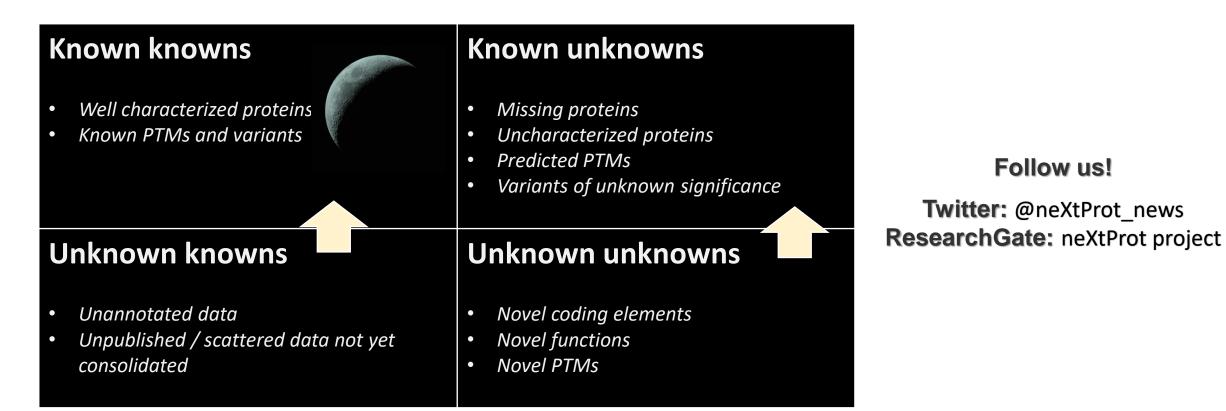
#### 15 federated queries already available

federated query- Filler sparql examples
NXQ_00139 - Protein kinases which are high-confidence drug targets according to CHEMBL drug engine (rederated every (function) evaluativ) function
NXQ_00140 - Proteins that interact with viral proteins federated query PPI tutorial
NXQ_00246 - Proteins which are enzymes catalyzing a reaction involving lipids enzyme (rederated query) function (tutorial)
NXQ_00253 - Human pathways in which at least one protein is mitochondrial GOLD [federated query] pathway (quality (snorql-only) subcellular location [tutorial]
NXQ_00254 - Proteins with associated pathways in WikiPathways federated query pathway (storial)
NXQ_00258 - Proteins involved in diseases due to intronic variants with one selected publication disease federated query publication snorq-only tutorial variant
NXQ_00259 - Proteins involved in diseases with clinical manifestations that include long organs disease (disease) (extended query (Latorial)
NXQ_00260 - Proteins with high-frequency missense variants involved in bacterial infection, with dbSNP identifiers and position on the canonical isoform disease (store) (stor
NXQ_00266 - Proteins binding estradiol and/or similar molecules (substructure search with SMILES) and their associated GO_MF terms (rederated cuery) function (similarity) small molecule interaction) sorrel-only (substructure search with SMILES) and their associated GO_MF terms
NXQ_00267 - Proteins binding estradiol and/or similar molecules (similarity search with SMILES), and their associated GO_MF terms federated query function similarity small molecule interaction snore-lonly tutorial
NXQ_00269 - Proteins with associated cancer pathways in WikiPathways (via Disease Ontology classification)           Tederated query         pathway         snord-only         tutorial
NXQ_00270 - Proteins belonging to Rett syndrome pathways, and their subcellular locations (GOLD) disease: (Ederated overy: (astronically: (astronical): (ast
NXQ_00272 - Proteins involved in coronaviruses/SARS-CoV-2 pathways with associated medical information disease. (Indexated query: (pathway: prorelently: (tutorial)
NXQ_00273 - Proteins involved in coronaviruses/SARS-CoV-2 pathways that are expressed in lung according to RNA-seq analysis and detected at high levels by IHC in at least one lung cell type expression (referated every) (pathway) (promionly) (tutorial)
NXQ_00276 - Diseases/phenotypes associated with coding variants and associated publications for a given gene disease [federated query] morei-only [statrial] variant

#### Example: Proteins binding estradiol and/or similar molecules (federated query with IDSM/SACHEM)

<pre>PREFIX sachem: <http: bioinfo.uochb.cas.cz="" rdf="" sachem#="" v1.0=""> PREFIX idsm: <https: endpoint="" idsm.elixir-czech.cz="" sparql=""></https:> PREFIX chembl: <http: chembl#="" rdf.ebi.ac.uk="" terms=""></http:></http:></pre>	NXQ_00139 - Protein kinases which are high-confidence drug targets according to CHEMBL drug enzyme federated query function quality tutorial	
<pre>ELECT distinct ?entry (group_concat(distinct str(?gomflab); SEPARATOR = ",") as SERVICE <https: endpoint="" idsm="" idsm.elixir-czech.cz="" sparql=""> { SERVICE <https: chembl="" endpoint="" idsm.elixir-czech.cz="" sparql=""> { ?compound sachem:similarCompoundSearch [ sachem:query "CC12CCC3C(C1CC)</https:></https:></pre>	NXQ_00140 - Proteins that interact with viral proteins federated query PPI tutorial	
<pre>}</pre>	e enetein tangat	NXQ_00246 - Proteins which are enzymes catalyzing a reaction involving lipids enzyme federated query function tutorial
<pre>?TARGET chembl:hasTargetComponent ?COMPONENT. ?COMPONENT chembl:targetCmptXref ?UNIPROT. filter(contains(str(?UNIPROT),"uniprot")) }</pre>	NXQ_00253 - Human pathways in which at least one protein is mitochondrial GOLD federated query pathway quality snorql-only subcellular location tutorial	
?entry`skos:exactMatch ?UNIPROT. ?entry :isoform ?iso. ?iso :goMolecularFunction / :term ?gomf . ?gomf rdfs:label ?gomflab .		NXQ_00254 - Proteins with associated pathways in WikiPathways federated query pathway snorql-only tutorial
roup by ?entry		
oteins binding estradiol and/or similar molecules (similarity search with SMILES), and t	eir associated GO_MF terms	NXQ_00258 - Proteins involved in diseases due to intronic variants with one selected publication disease federated query publication snorql-only tutorial variant
		NXQ_00259 - Proteins involved in diseases with clinical manifestations that include long organs disease federated query tutorial
		NXQ_00260 - Proteins with high-frequency missense variants involved in bacterial infection, with dbSNP identifiers and position on the canonical isoform
		disease federated query isoforms sequence snorql-only tutorial variant
html 🗸 Go Reset	eg. peroxisome, liver Q	NXQ_00266 - Proteins binding estradiol and/or similar molecules (substructure search with SMILES) and their associated GO_MF terms
und a bug? Improve this query!		federated query function similarity small molecule interaction snorql-only tutorial
Query time is 10.989[s] for 240 rows		NXQ_00267 - Proteins binding estradiol and/or similar molecules (similarity search with SMILES), and their associated GO_MF terms
entry gomfx	· · · · · · · · · · · · · · · · · · ·	federated query function similarity small molecule interaction snorql-only tutorial
ntry:NX_Q9HAW8 (neXtProt link) "enzyme binding,glucuronosyltransferase activity,protein activity,protein kinase C binding,retinoic acid binding"	eterodimerization activity,protein homodimerization	NXQ_00269 - Proteins with associated cancer pathways in WikiPathways (via Disease Ontology classification)
ntry:NX_P00918 (neXtProt link) arylesterase activity,carbonate dehydratase activity,hydr	-lyase activity,protein binding,zinc ion binding"	federated query pathway snorql-only tutorial
entry:NX P20813 (neXtProt link) "anandamide 11.12 epoxidase activity.anandamide 14.15	poxidase activity, anandamide 8,9 epoxidase activity, arachidonic	

# Let's complete the functional proteome together!



#### Interested in

- Integrating your data in neXtProt ?
- Federating your data with neXtProt data using semantic technologies ?
- Building personal, customized SPARQL queries for your own research ?

Let's discuss! lydie.lane@sib.swiss support@nextprot.org

# Special thanks



The neXtProt current team members

Paula Duek Pierre-André Michel Valentine Rech de Laval Kasun Samarasinghe Monique Zahn Amos Bairoch

...and all the past members!



Young-Ki Paik Chris Overall Charles Pineau

Gil Omenn

**Eric Deutsch** 

and all C-HPP collaborators





