Lincoln Cathedral, Lincolnshire UK



Grand Bazaar, Istanbul, Turkey

Why should we be re-using data?


Four types of (proteomics) data re-use


Unbiased proteome-wide (PT)M discovery as example


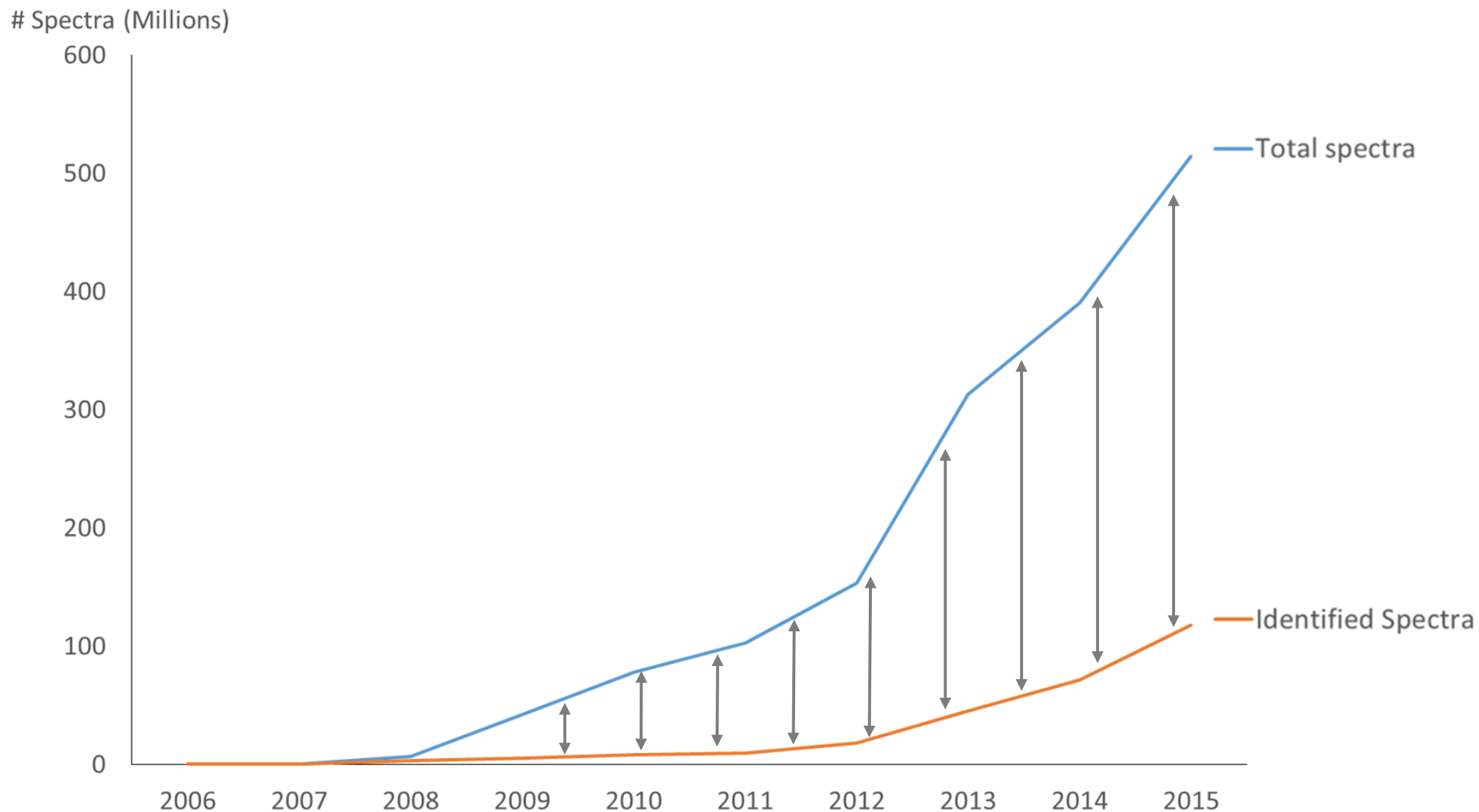A subject of sociological study

# Why should we be re-using data?

Four types of (proteomics) data re-use
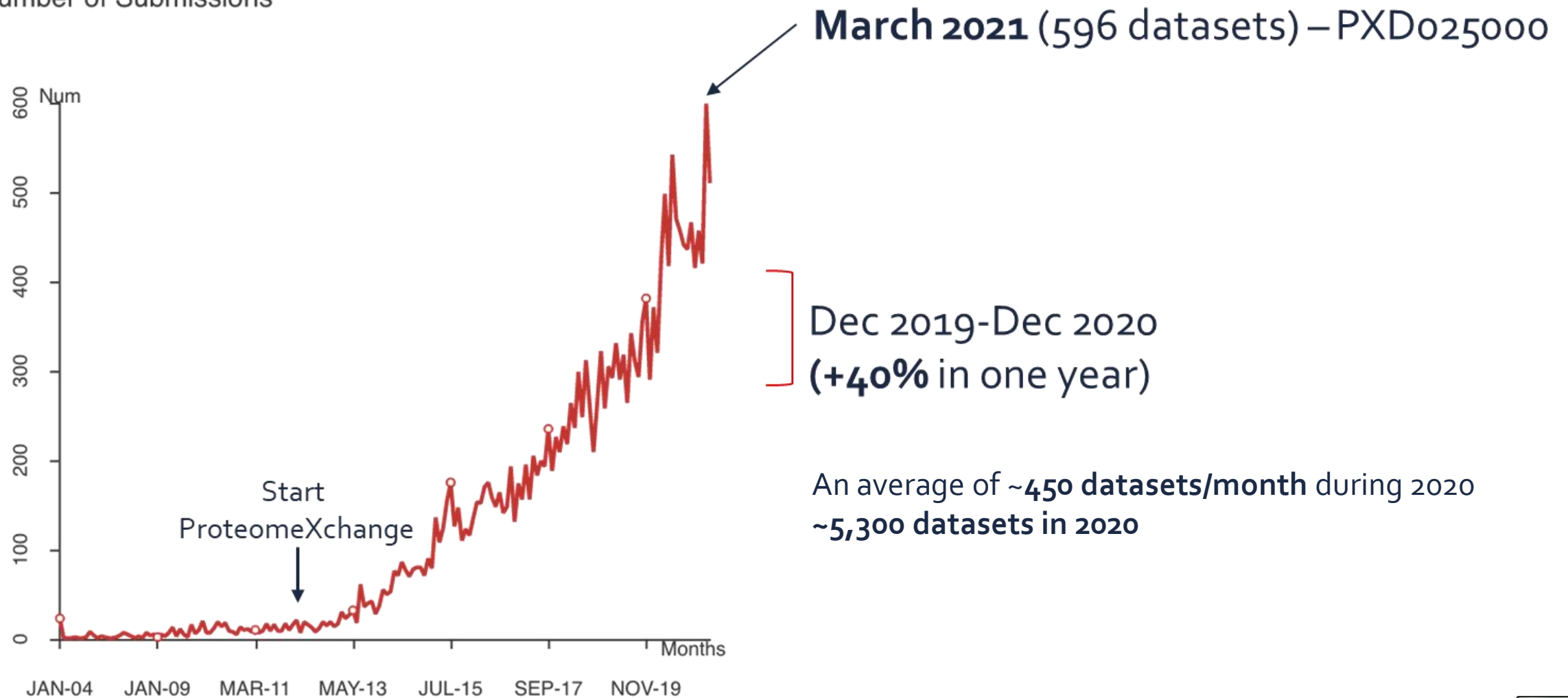
Unbiased proteome-wide (PT)M discovery as example

A subject of sociological study

# Mass spectrometry data is high-content, meaning that much more data is acquired than is used in most papers
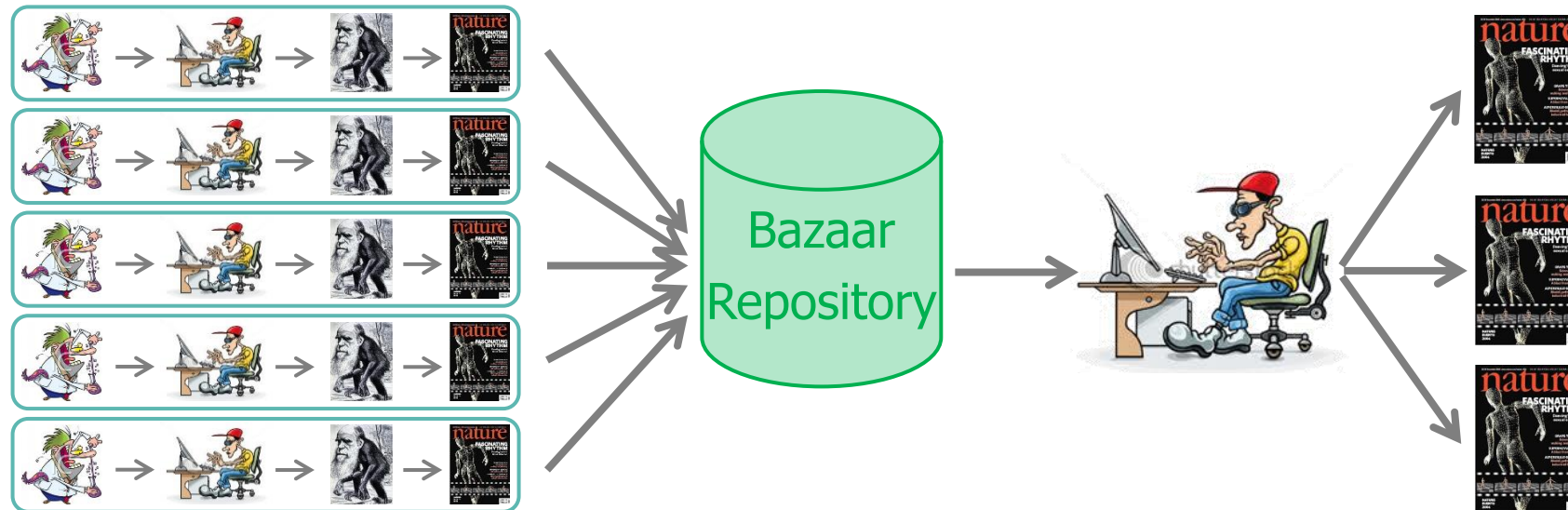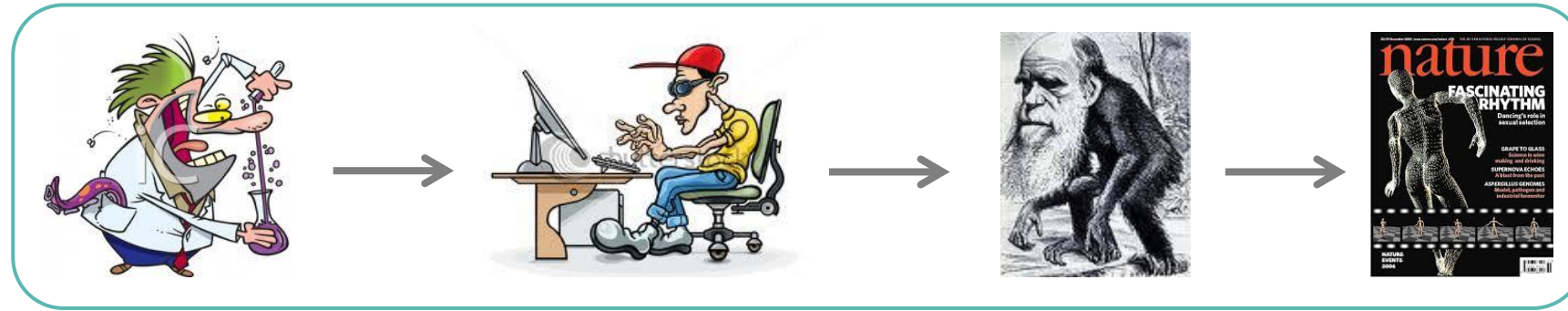
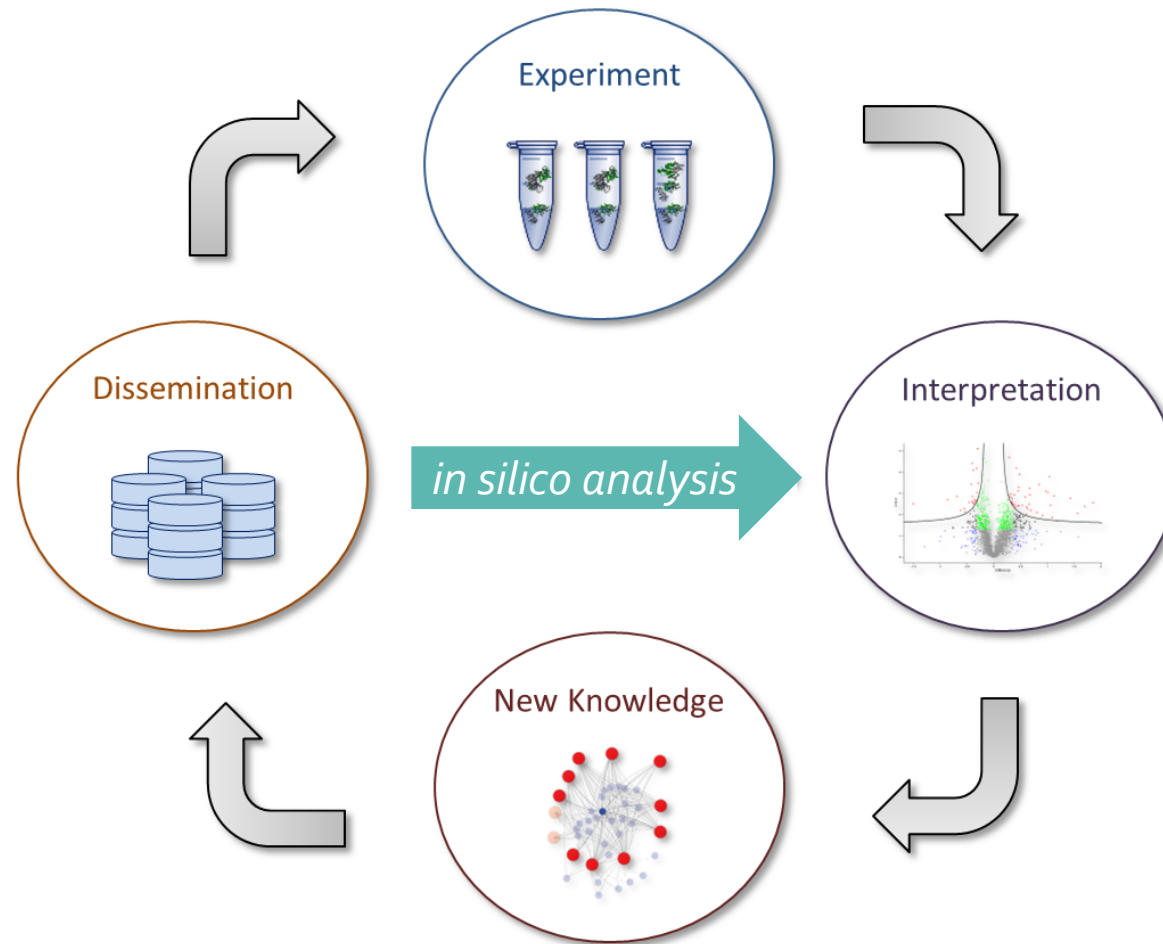# Mass spectrometry is also high throughput, meaning there is lots of data available!



An average of ~**450 datasets/month** during 2020
~**5,300 datasets in 2020**

# As the volume and content of data increases in a field, the role of bioinformatics in that field changes as well
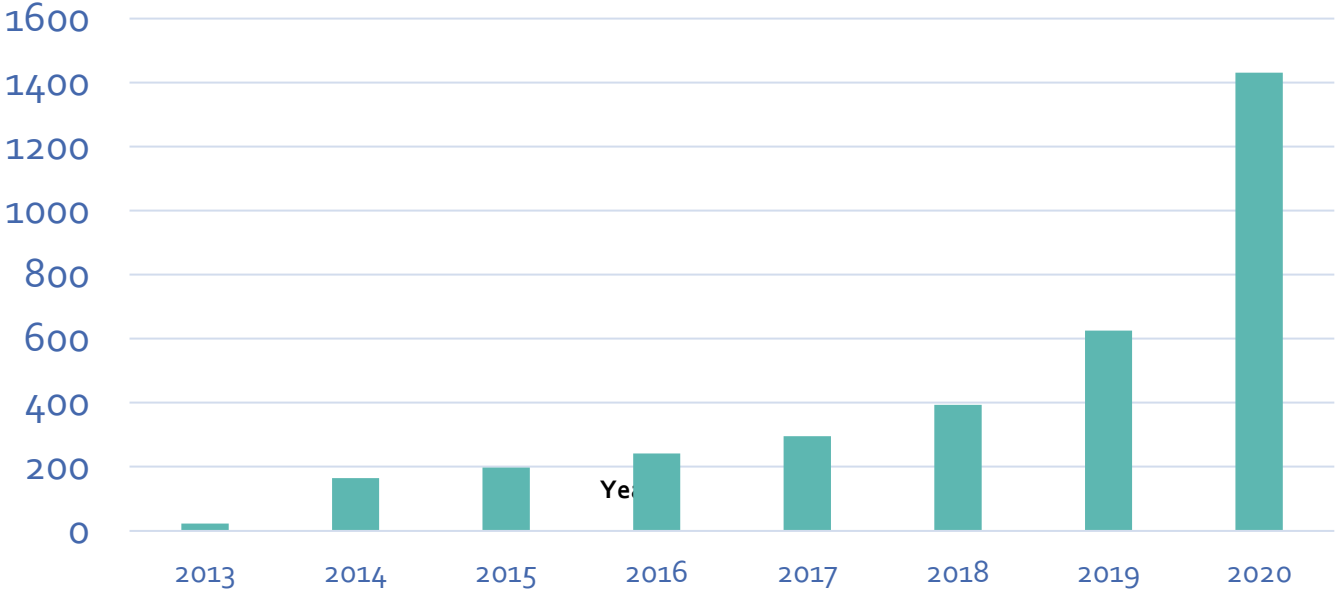
# The data life cycle shows how *in silico* re-use of data fits in with the overall flow of information

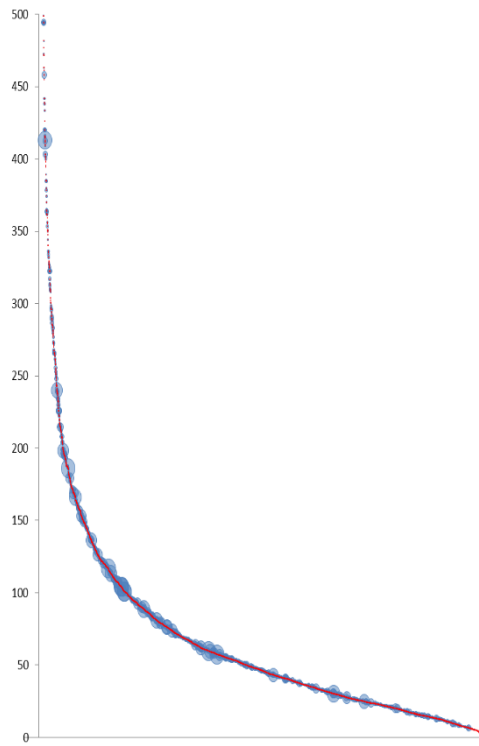# And this is seen in practice in proteomics as data is increasingly re-used
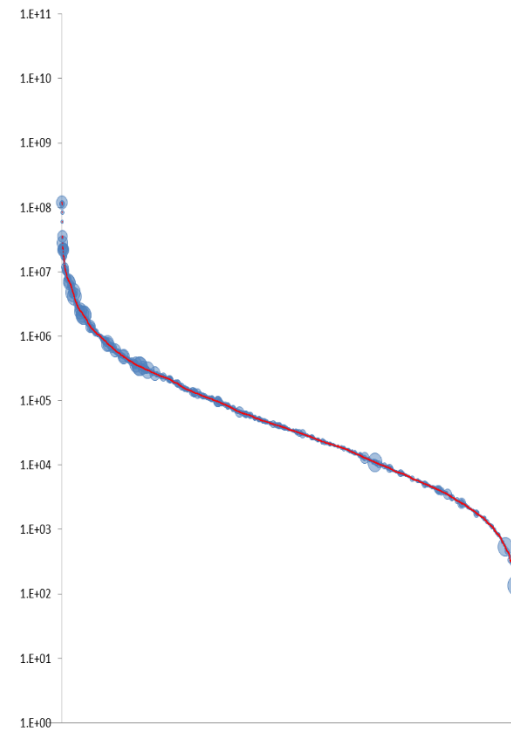


Volume of PRIDE Data Downloads (TB), 2013-2020

# Large-scale data reprocessing can harness heterogeneity to dig very deep into the proteome

Half life (h)

concentration (copies per cell)

concentration (pg/ml)

proteins

proteins

Not Identified
Identified (1 PSM / assay)
Identified (2 PSM / assay)
Identified(3 PSM / assay)

Normal plasma concentration (pg/ml)

Verheggen, Journal of Proteome Research, 2017

Why should we be re-using data?

**Four types of (proteomics) data re-use**

Unbiased proteome-wide (PT)M discovery as example

A subject of sociological study

# In general, data re-use can take four distinct forms, all of which are somehow applied in our example
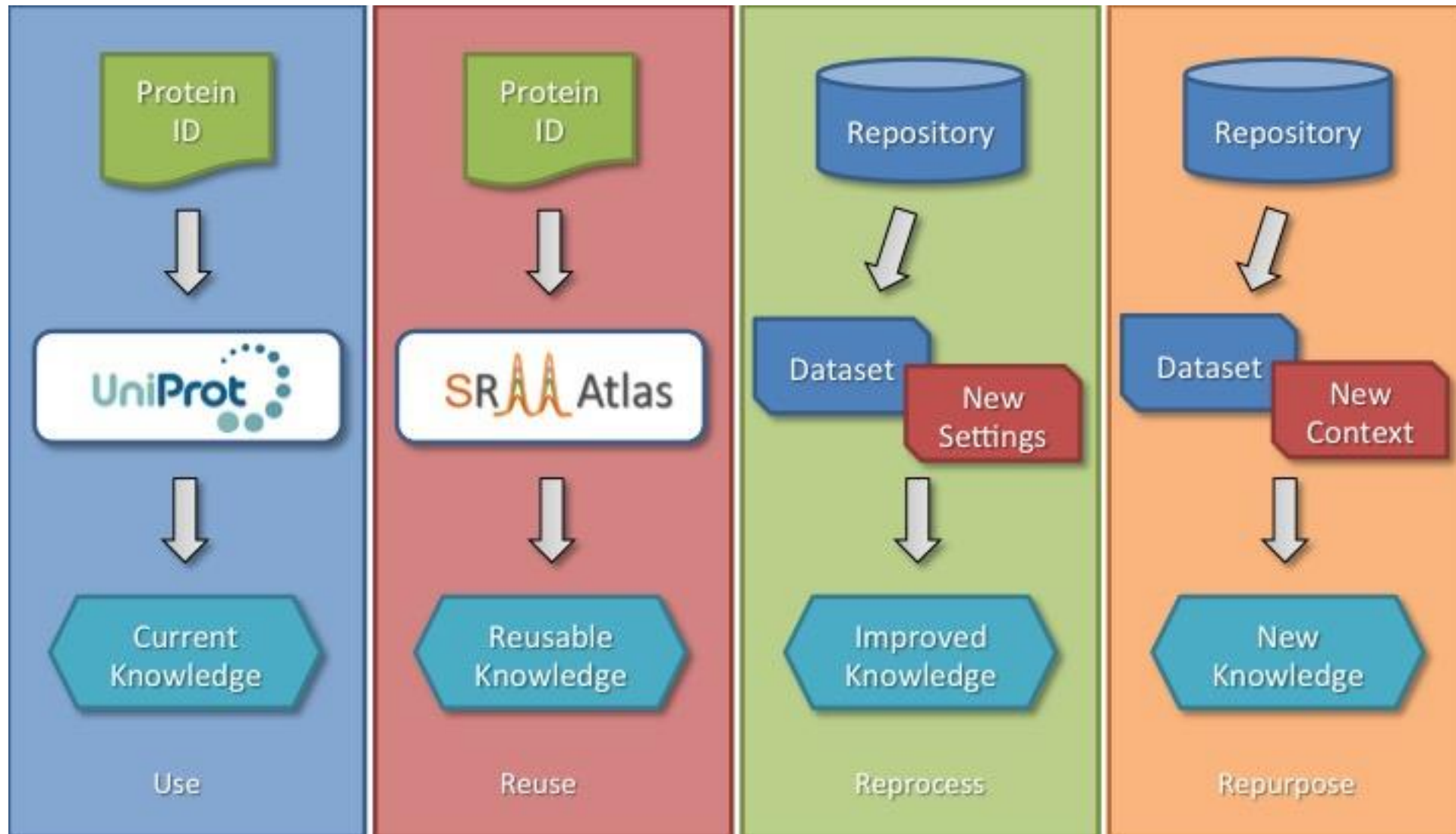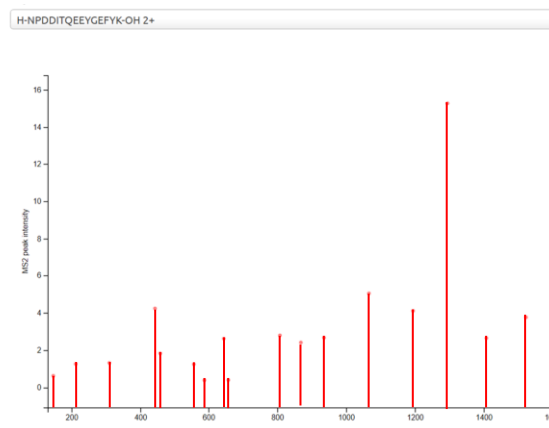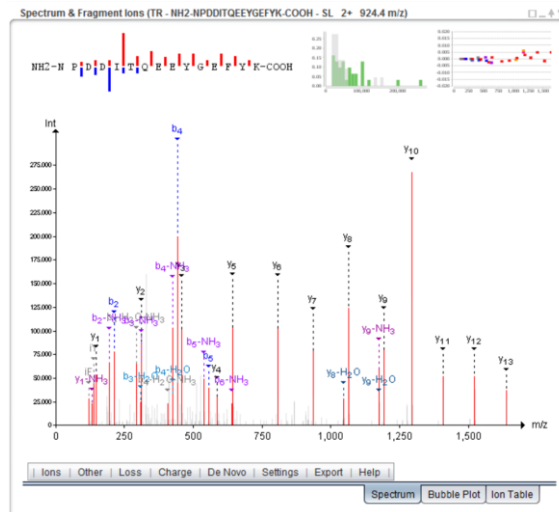
Why should we be re-using data?

Four types of (proteomics) data re-use

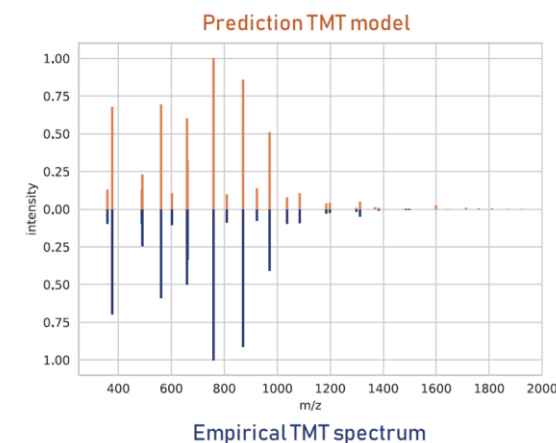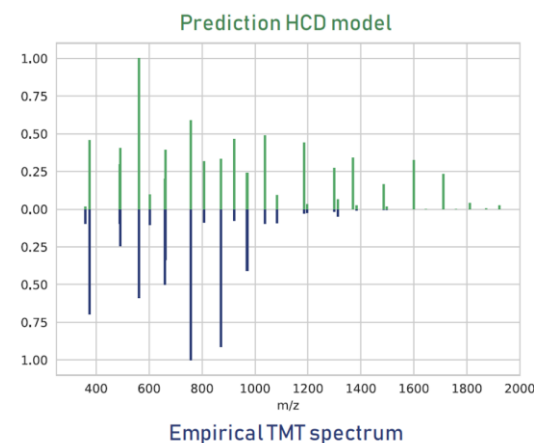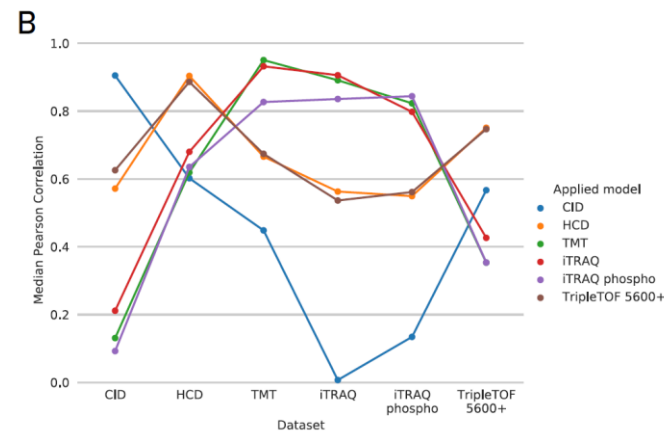**Unbiased proteome-wide (PT)M discovery as example**

A subject of sociological study
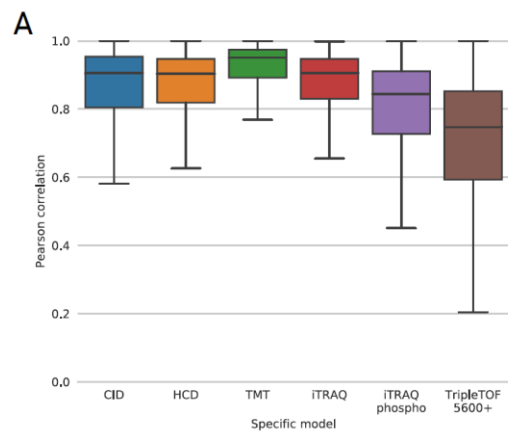
# Our MS2PIP fragmentation model accurately predicts peptide behaviour in varying conditions



Vaudel, Nat. Biotech., 2015
PeptideShaker

https://iomics.ugent.be/ms2pip
Degroeve, Bioinformatics, 2013
Degroeve, Nucleic Acids Research, 2015
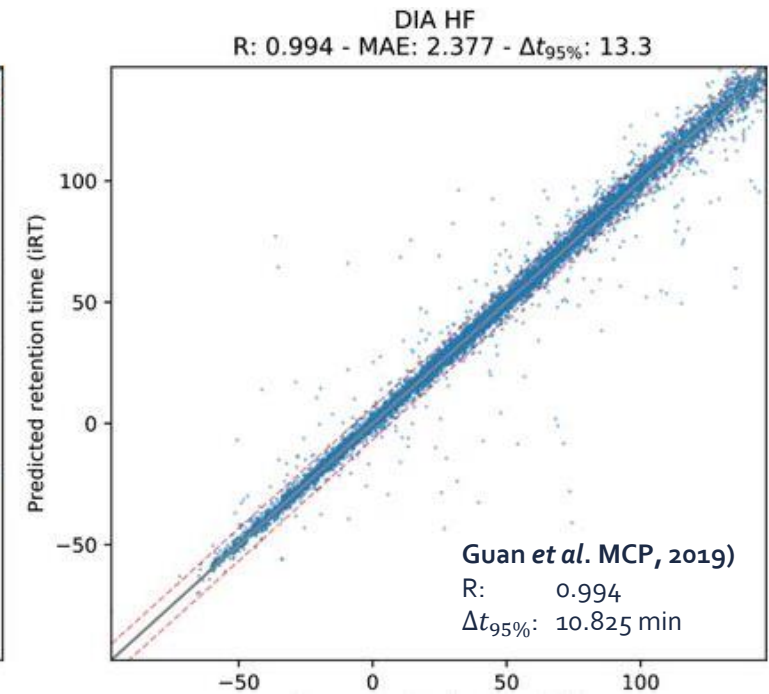
Gabriels, Nucleic Acids Research, 2019

# DeepLC is a retention time predictor that can accurately predict retention times of as-yet unseen modifications

Predicted retention time (min)

SWATH library
R: 0.997 - MAE: 2.543 - $\Delta t_{95\%}$: 14.88

HeLa HF
R: 0.984 - MAE: 0.314 - $\Delta t_{95\%}$: 1.62

DIA HF
R: 0.994 - MAE: 2.377 - $\Delta t_{95\%}$: 13.3

DeepRT+ (Ma, Anal. Chem, 2018)
R: 0.997
$\Delta t_{95\%}$: 13.7 min

Guan *et al.* MCP, 2019)
R: 0.994
$\Delta t_{95\%}$: 10.825 min

Observed retention time (min)

# DeepLC can accurately predict t$_R$ for many modifications, despite never having seen these before



acetylation

succinylation

propionylation

DeepLC prediction
Traditional prediction

phosphorylation

crotonylation

malonylation
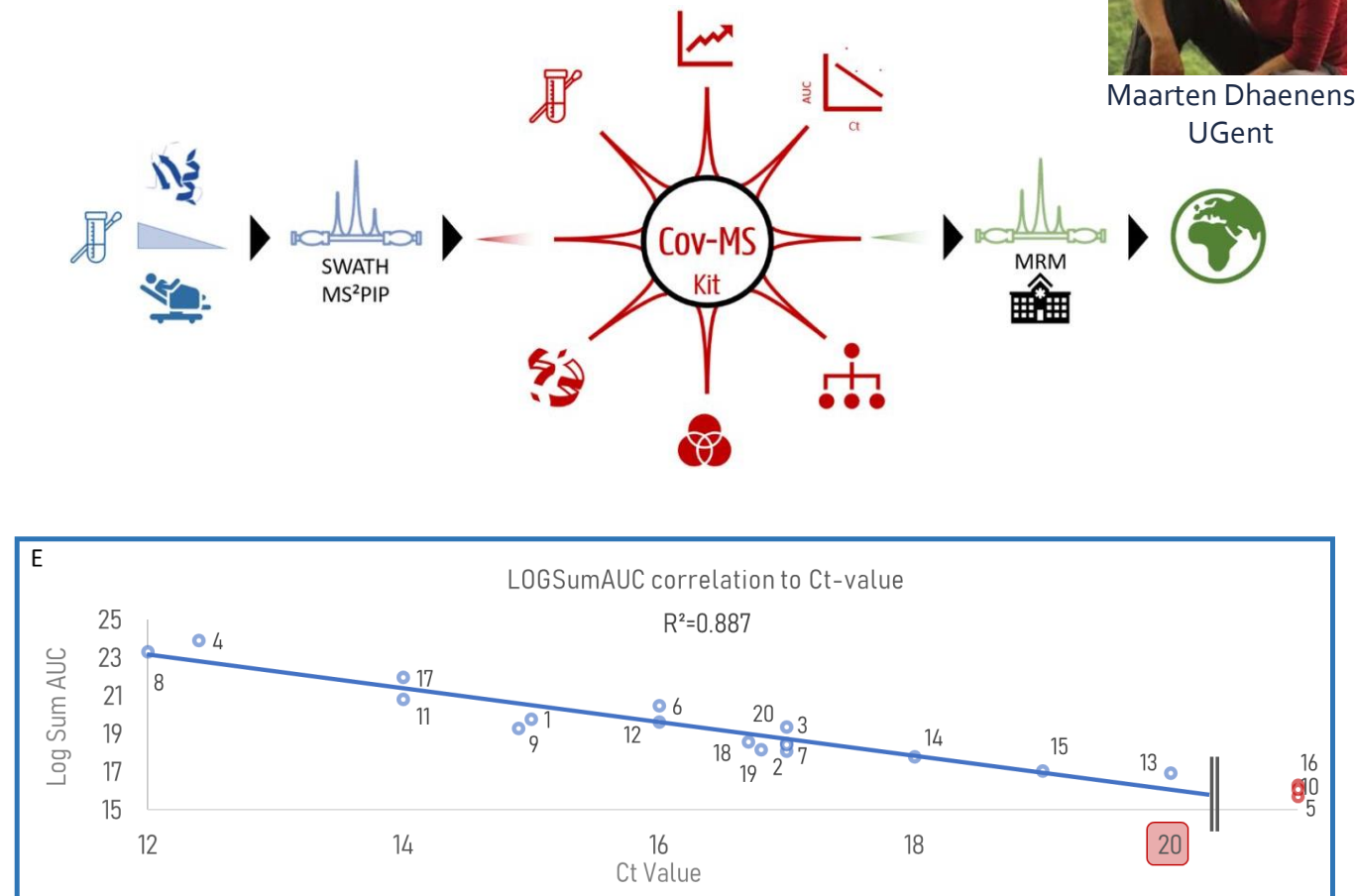
# MS2PIP and DeepLC were crucial in the development of a targeted MS-based COVID-19 test that runs in 38 mins



Maarten Dhaenens
UGent

# MS2PIP and DeepLC power ionbot, a novel and extensible open modification search engine with high reliability

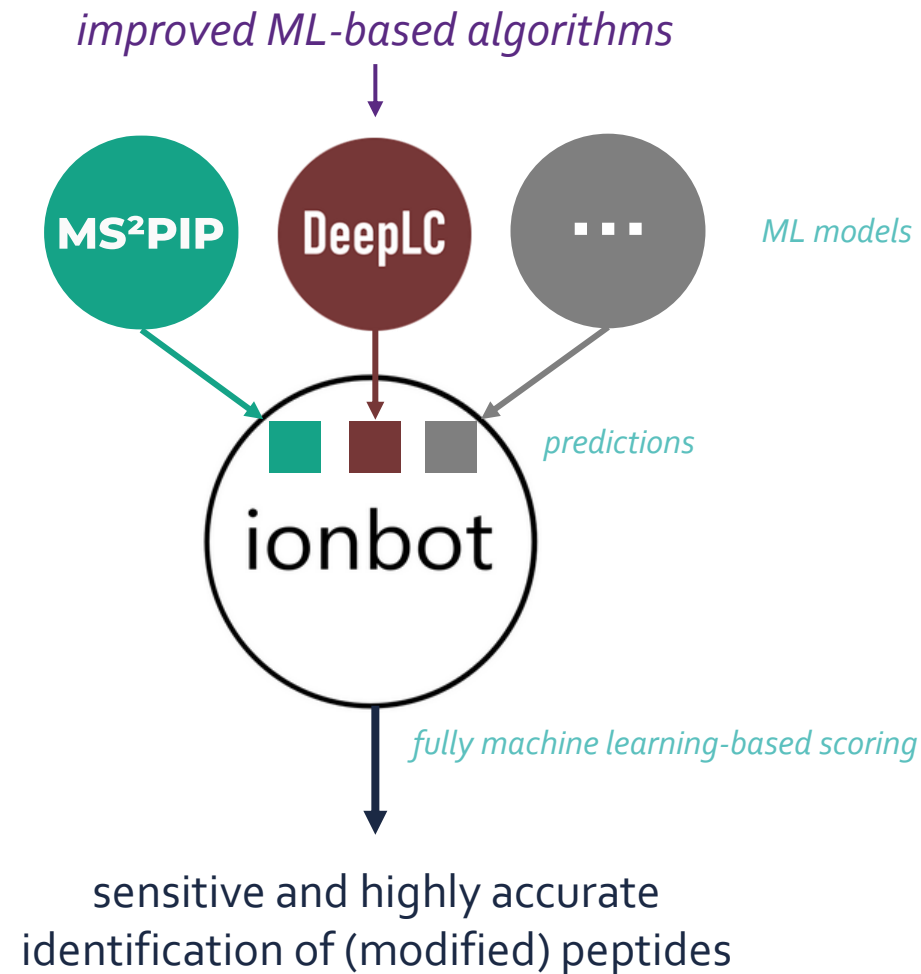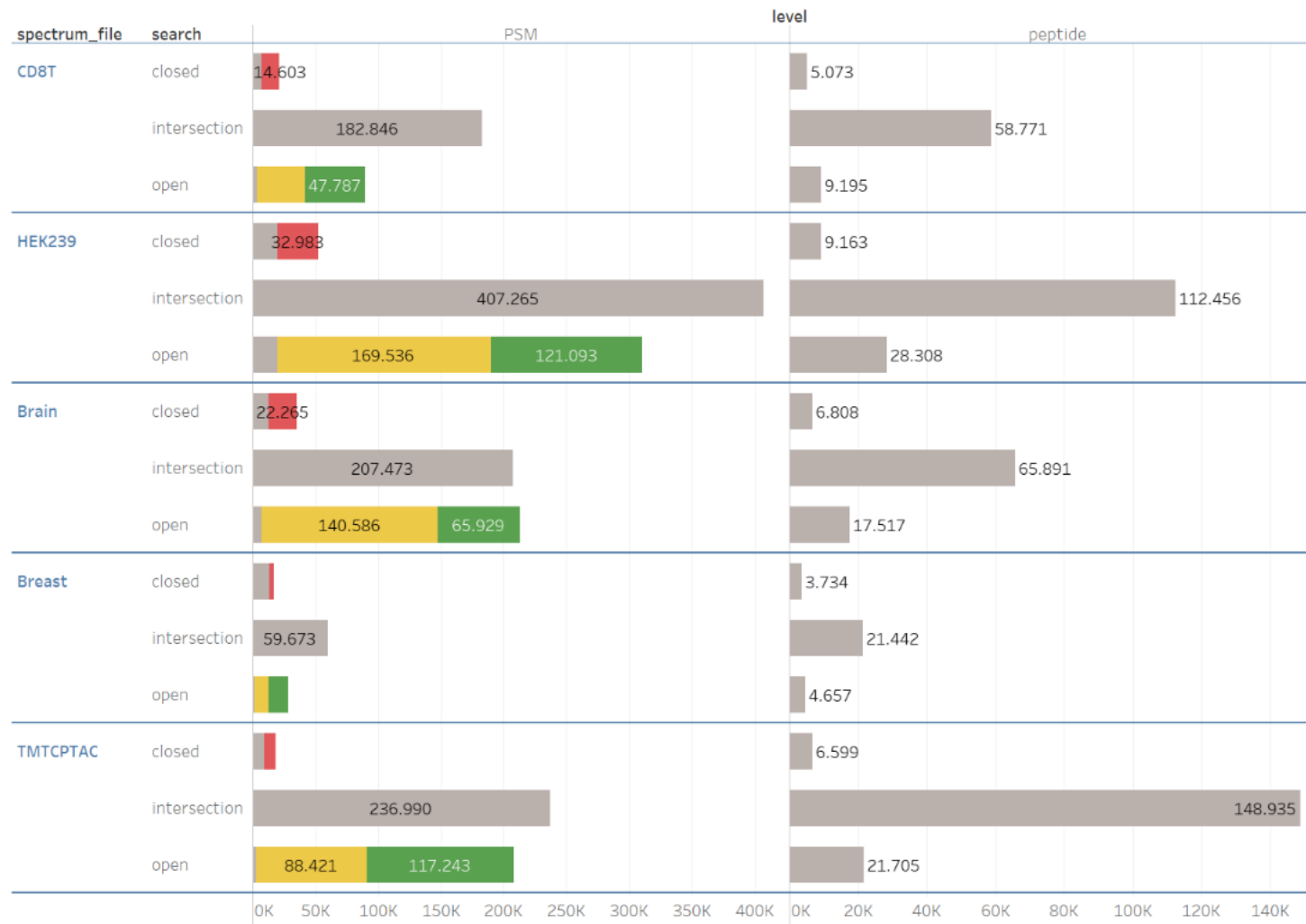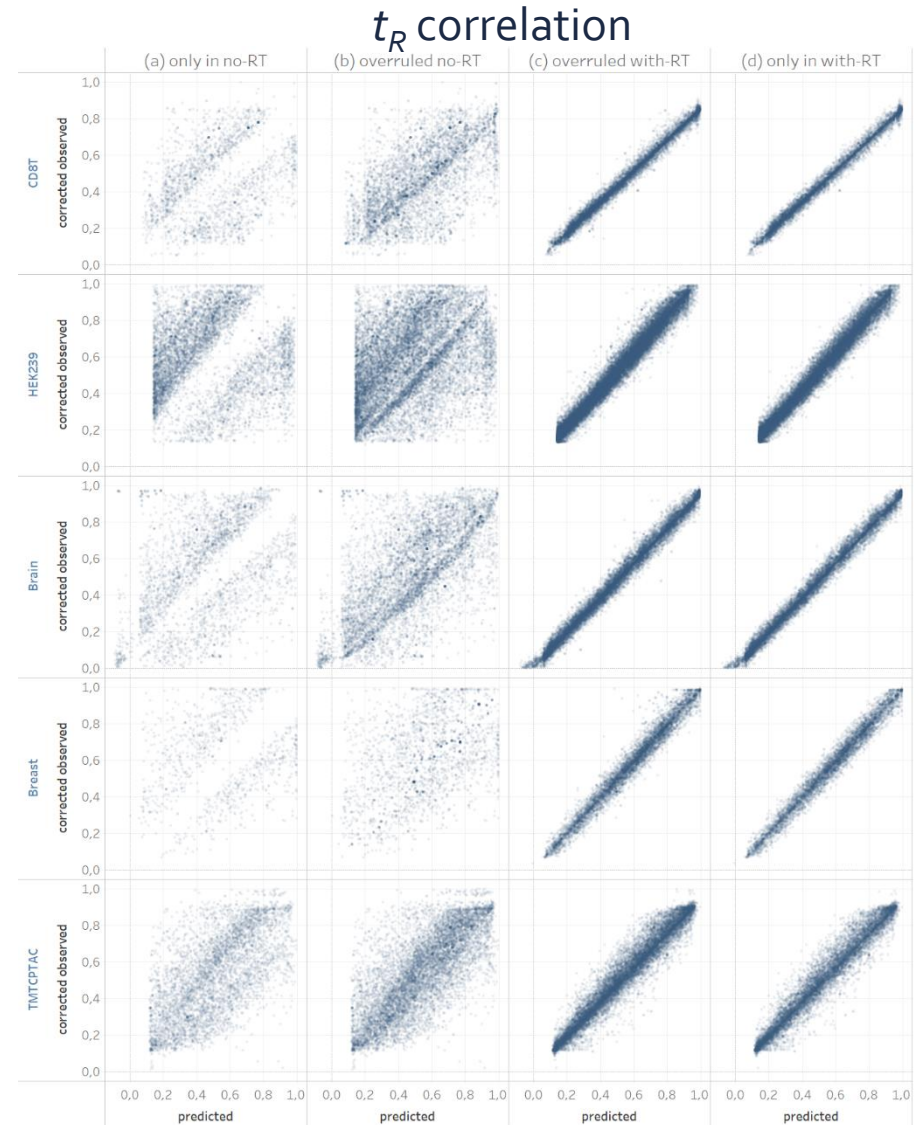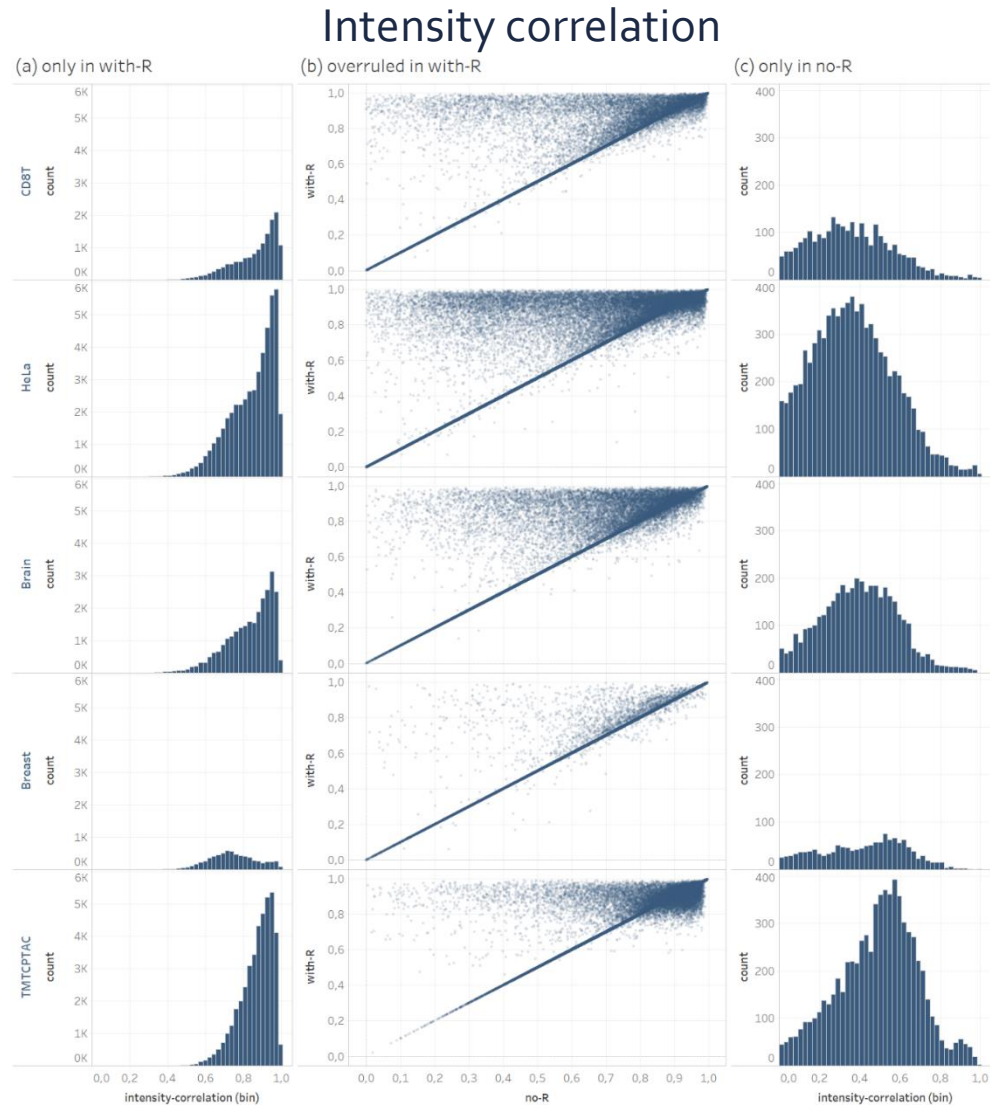# Ionbot shows the value of open modification searches, as well as the value of accurate prediction models



Degroeve, https://www.biorxiv.org/content/10.1101/2021.07.02.450686v1

# Interestingly, many identifications from the closed search are overruled by ionbot in the open search



Legend:
- **overruled by ionbot** (red)
- **unexpected modification** (green)
- **large mass window** (yellow)

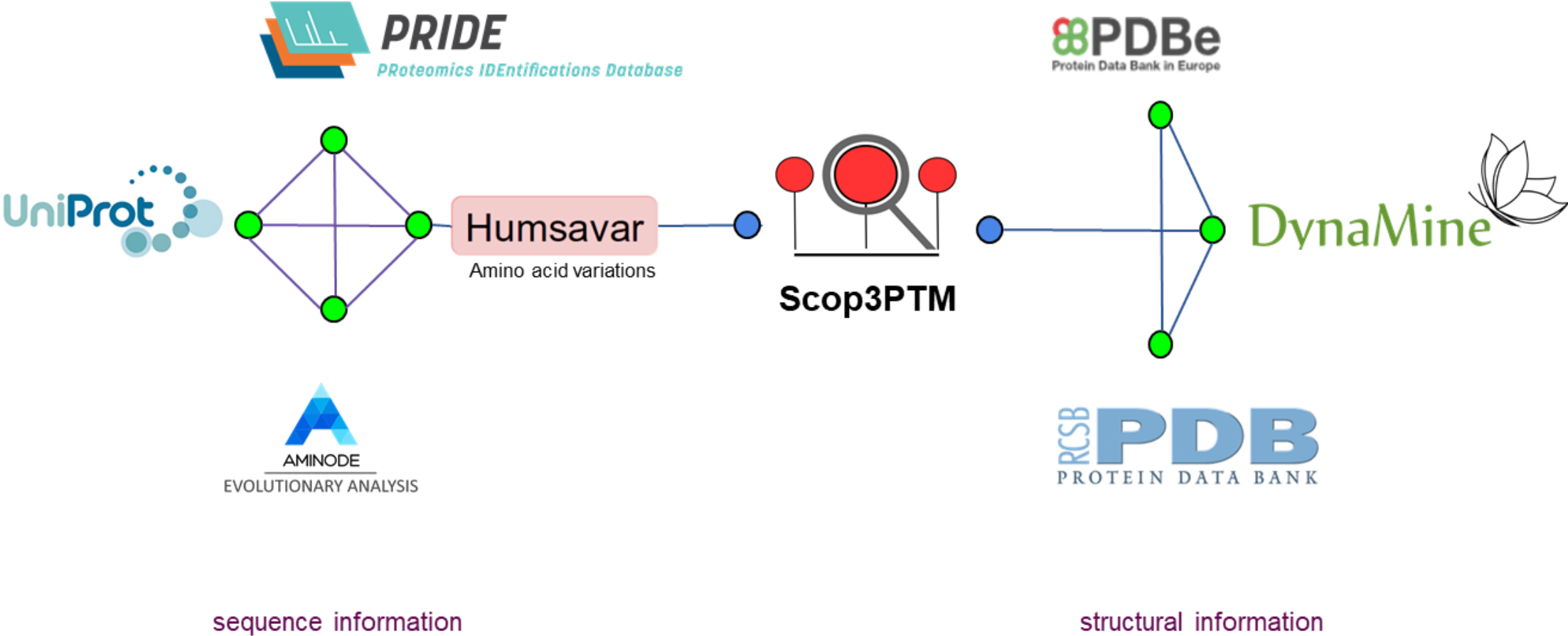| spectrum_file | search | PSM | peptide |
|---|---|---|---|
| CD8T | closed | 14.603 | 5.073 |
| | intersection | 182.846 | 58.771 |
| | open | 47.787 | 9.195 |
| HEK239 | closed | 32.983 | 9.163 |
| | intersection | 407.265 | 112.456 |
| | open | 169.536 / 121.093 | 28.308 |
| Brain | closed | 22.265 | 6.808 |
| | intersection | 207.473 | 65.891 |
| | open | 140.586 / 65.929 | 17.517 |
| Breast | closed | | 3.734 |
| | intersection | 59.673 | 21.442 |
| | open | | 4.657 |
| TMTCPTAC | closed | | 6.599 |
| | intersection | 236.990 | 148.935 |
| | open | 88.421 / 117.243 | 21.705 |

# Overruled identifications are better, as shown for results obtained with, and without predictions provided to ionbot



Intensity correlation

$t_R$ correlation

Degroeve, https://www.biorxiv.org/content/10.1101/2021.07.02.450686v1

# We reprocessed a large amount of phosphoproteomics data using ionbot, and made it available through Scop3P



www.ebi.ac.uk/pride

*2016 RAW files*
*60.2 million spectra*

*1490 UniMod modifications*
*all possible AA mutations*

iomics.ugent.be/scop3p

*19.2 million PSMs*
*139 048 P-sites*
*94 111 sites, PhosphoRS > 0.5*
*14 261 phosphoproteins*

# Scop3PTM integrates protein information at the residue level from a variety of resources

# Scop3P shows these results interactively on the web, and presents REST APIs for 3rd party re-use
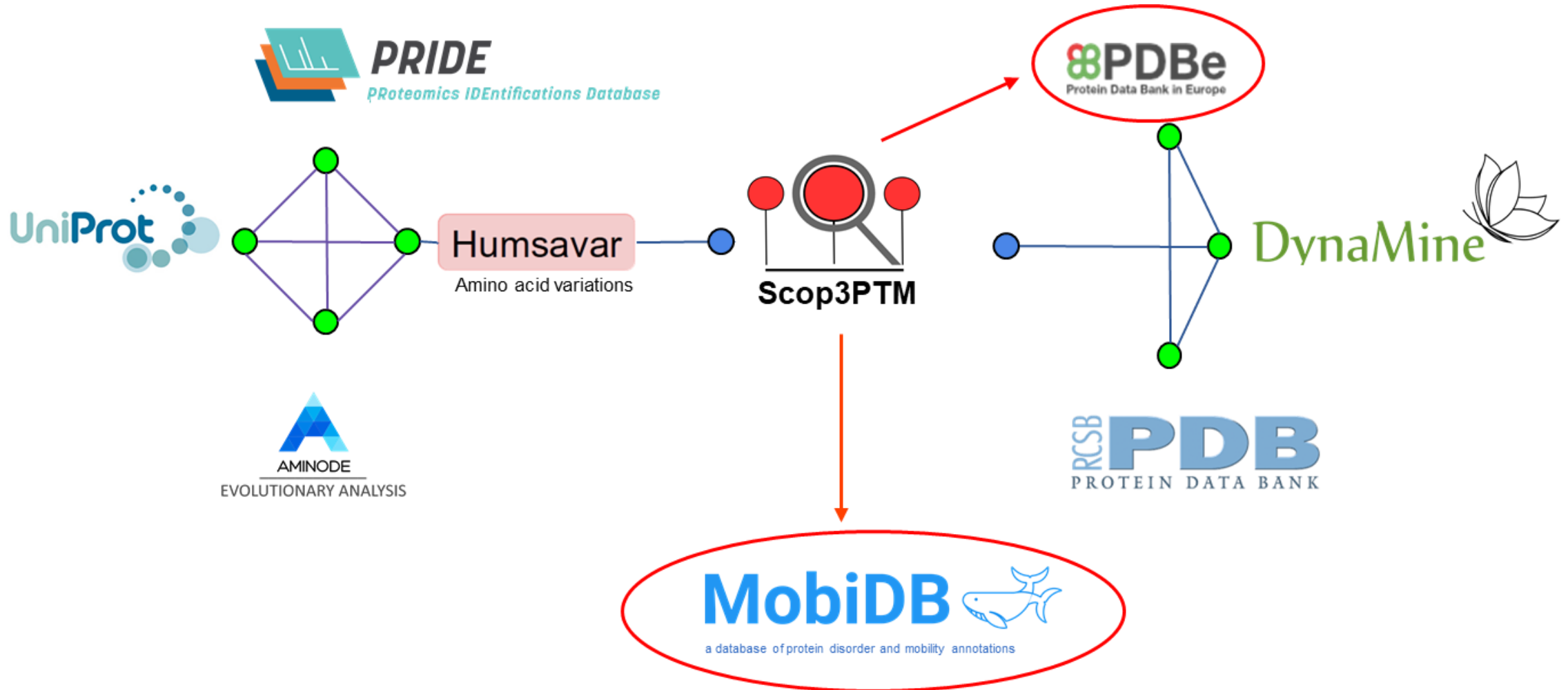
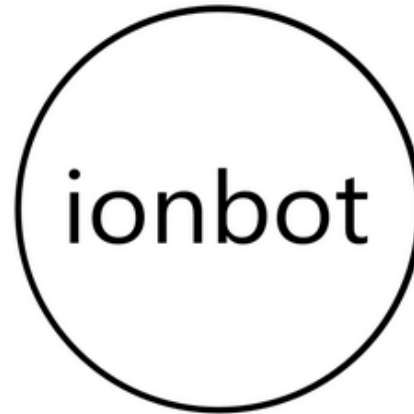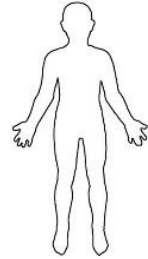https://iomics.ugent.be/scop3p
Ramasamy, Journal of Proteome Research, 2020

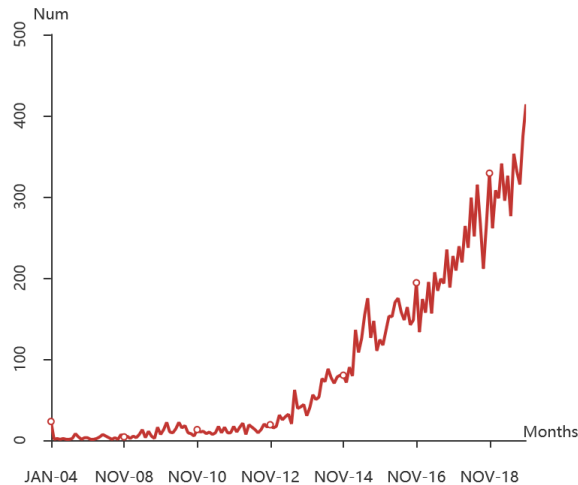# Closing the circle, Scop3P itself becomes a resource for use and re-use by others!

# And we are now running ionbot on all human spectra contained in the PRIDE database



**PRIDE**
PRoteomics IDEntifications Database

Number of Submissions

**25 314 RAW files**
**1 *billion* spectra**

**ionbot**

*1490 UniMod modifications*
*all possible AA mutations*

215,046,444  PSMs
742,422  unique peptides
20,246  proteins
444  modifications
4,869,660  modified residues

# Scop3PTM will become a proteome-wide PTM detection knowledgebase

Why should we be re-using data?

Four types of (proteomics) data re-use

Unbiased proteome-wide (PT)M discovery as example

**A subject of sociological study**

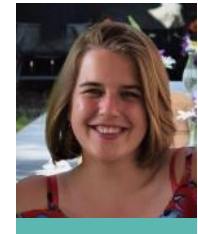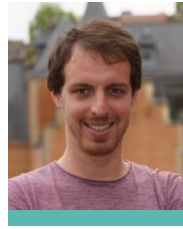# A sociological take on work of my group highlights the key benefits of open science

"This desire to reactivate data is widespread, and Klie et al. are not alone in wanting to show that 'far from being places where data goes to die' (Klie et al., 2007: 190), **such data collections can be mined for valuable information that could not be obtained in any other way**."

"In attempting to **reactivate sedimented data** in order to enable its re-use, their first step was …"

"… they are experiments in seeing, in furnishing ways of seeing how data on proteins could become re-usable, could be reactivated as **collective property rather than the by-product of publication**."

Mackenzie and McNally, *Theory, Culture and Society*, 2013

Gift shop at Chester Cathedral, UK

# www.compomics.com
## @compomics