

# From data to biology: using -omics datasets to generate an unbiased hypothesis

#### Katarzyna Kulej

Memorial Sloan Kettering Cancer Center, New York

#### Why do we do it?

- We all produce large amounts of data, and we rarely have time to make sense of them
- The first instinct in proteomics is to calculate the fold change of our proteins between condition A vs condition B, and expect that this unravels the molecular mechanisms of our system
- While there is nothing wrong in doing so, we forget that there are many more perspectives we can observe our data from. These can help us decipher unexpected properties of our sample
- We will initially discuss these perspectives, and then the O'Donovan lab will present a specific application in metabolomics



#### **Quantitative Dimensions in Proteomics**





#### Protein abundance (what it means and how to visualize it)

Protein rank plot from highest to lowest abundant protein, illustrating the dynamic range







- Dynamic range issues
- Abundant proteins might be better biomarkers
- Abundance plots can be used to differentiate sample types (e.g., highly specialized cells have very few very abundant proteins, while stem cells have lots of medium abundant proteins)



#### **Protein absolute intensity**







#### **Protein fold change**

- It is the most intuitive dimension
- To treat carefully: without p-value it might be meaningless, and it does not provide a perspective of how much protein there is







#### **P-value and its meaning**

To consult the statistician after an experiment is finished is often merely to ask him to conduct a postmortem examination. He can perhaps say what the experiment died of.

- R.A. Fisher

The first principle is that you must not fool yourself, and you are the easiest person to fool.

- Richard P. Feynman

Remember, a P-value is not a measure of how right you are or how important a difference is. Instead, think of it as a measure of surprise.

□ If you assume your medication is ineffective and there is no reason other than luck for the two groups to differ, then the smaller the p-value, the more surprising and lucky your results are – or your assumption is wrong, and the medication truly works.



#### **P-value and its meaning**

To consult the statistician after an experiment is finished is often merely to ask him to conduct a postmortem examination. He can perhaps say what the experiment died of.

- R.A. Fisher

The first principle is that you must not fool yourself, and you are the easiest person to fool.

- Richard P. Feynman

□ The P-value reports a probabilistic significance, not a biological one.

- In other words, statistical significance does not mean your result has any practical significance.
- The choice of p < 0.05 as significant is not because of any special logical or statistical reasons but it has become a scientific convention through decades of common use.

## First (obvious) examples of data analysis by combining quantitative dimensions



Volcano plots showing results of comparisons between two conditions (A vs B).

## Fold change AND p-value



fold-change cutoff

 fold-change cutoff is set with regard to the experimental power



## Use the absolute/relative intensity!

- You can find biomarkers
- You can prioritize candidates
- You can interpret protein hindrance

#### Abundance AND fold change AND p-value

Infection of host cells (Mock) with Adenovirus type-5 (Ad5)



#### Less obvious quantitative dimensions



- Protein PTMs have their own abundance, independent from protein abundance
- PTM stoichiometry can be used to assess activity of the enzymes that catalyze them, and thus predict biological targets

#### Protein post-translational modifications (PTMs)

**DNAPK** - DNA-dependent protein kinase catalytic subunit



MSK1 - Ribosomal protein S6 kinase alpha-5





https://www.phosphosite.org/homeAction



Confidence in PTM localization is also quantitative, and it can help to prioritize targets

#### Protein post-translational modifications (PTMs)

**DNAPK** - DNA-dependent protein kinase catalytic subunit



#### Kinase substrate motif DNAPK Ŀ, F og-odds ( 298 input sequences DNA damade MSK1 82 input sequences MAP kinase family – cell stress response

https://www.phosphosite.org/homeAction









#### **Tools for PTM assignment and quantification**

#### **PTM prediction**

NetPhos: Prediction of ph-sites based on the learning model NetPhosK: Prediction of kinase-specific ph-sites Scansite: Motifs likely to be phosphorylated by specific kinases LysAcet: Prediction of acetylation sites PTMfunc: Repository of functional predictions of PTMs iGPS: prediction of site-specific kinase-substrate relations

#### Related kinases/transferases

KinomeXplorer: Predict kinase-substrate interactions (NetworKIN + NetPhorest) <u>PKIS</u>: Computational identification of protein kinases <u>PSEA</u>: Enrichment analysis for prediction of kinases

#### Already reported PTMs

Uniprot: Experimental/putative PTMs

PhosphoELM: Curated database of validated ph-sites

PhosphoSitePlus: Experimentally reported ph-, ac- and ub-sites

OGlycBase: Experimentally reported glycosylation sites

<u>HPRD</u>: Experimentally reported human PTMs

<u>dbPTM</u>: Compendium of experimental and putative PTMs from sources above

PTMcode: Focus on PTM crosstalk

ProteomicsDB: multi-omics and multi-organism resource <u>HPRD:</u> domain architecture, post-translational modifications, interaction networks and disease association for each protein in the human proteome.

<u>PhosphoNetworks</u>: a database for human phosphorylation networks



# HSV-1 infection



Example 6 time points of

#### **Time series**

In most experiments, the proteome (and other biomolecules) is regulated in a dynamic and different manner

Distinguishing early from late responders resolves critical aspects of the biological system

#### Grouping trends with unsupervised clustering (e.g., fuzzy c-means)

#### Time course of HSV-1 infection



Time-resolving protein dynamic changes helps group different stages of regulation

#### Cluster Gene Ontology (GO) Analysis







#### **Combine fold change and turnover**



9hpi of HSV-1 infection, compared to uninfected mock samples



Other quantitative dimensions that assist the interpretation of protein interactions and functions



#### **Co-regulation to identify genetic interactions**

Ad5 – Adenovirus type-5 HSV – Herpes Simplex Virus VACV – Vaccinia virus

#### Concept:

Proteins with the same regulation in a large set of samples have potentially shared functions.



Reves E.D. et al. (2017), Mol Cell Proteomics



 The analysis is based on protein abundance levels from ProteomicsDB. Correlations are shown as clustered heatmaps and as graphs to reveal potential co-regulated protein subgroups.

## **Co-regulation to identify genetic interactions**

Ad5 – Adenovirus type-5 ; HSV – Herpes Simplex Virus; VACV – Vaccinia virus



"-" biotin; "+" biotin

Reyes E.D. et al. (2017), Mol Cell Proteomics



In proteomics, most often, the nodes represent proteins, e.g.

- color of the node represents the fold change enrichment
- size the p-value or abundance
- shape the functional enrichments (e.g., Gene Ontology)

The edges usually represent one of three types of information:

- physical interaction or proximity (interactomics),
- phenotypic similarity (profiling),
- shared annotations (e.g., Gene Ontology)

When creating a network, it is essential to understand the required network type, the key information that should be associated with nodes and edges, and how to represent it in the visualization when creating a network

> Schessner J.P., et al. (2022) Proteomics. Koutrouli, M., et al. (2020), Front Bioeng Biotechnol.



Different tools can be used to create networks depending on the degree of complexity and customization required e.g.,:

Databases for protein-protein interaction networks:

- STRING: edges (connections) given by scores
- BioGRID: curated from publications (nearly complete for yeast)
- **KEGG:** more focused on pathways and metabolism
- **Reactome:** curated annotation of pathways on a diverse set of topics
- NextProt: an integrated collection of interactions between human proteins
- IntAct: curation from literature and submitted information

Databases for protein-PTMs interaction networks:

- **PhosphoSitePlus:** Eexperimentally reported ph-, ac- and ub-sites
- UniProt: experimental/putative PTMs
- PhosphoPath: visualization and analysis of quantitative proteome and phosphoproteome datasets
- MetaCore: data-mining and pathway analysis
- Ingenuity IPA: a tool to integrate and understand complex 'omics data

#### ...endless number of specific databases

Schessner J.P., et al. (2022) Proteomics.





Visualization and analysis of quantitative proteome and protein PTM datasets, e.g.:







#### **Interaction AND activity**

*Layout apps available in Cytoscape* 





#### **Protein domains: adding structural details to protein networks**



#### Motifs/domains

protein-protein affinity

- <u>Prosite</u>: Biologically significant sites, patterns and profiles <u>Pfam</u>: Collection of protein domain families <u>Blocks</u>: Database of conserved protein "blocks" (short sequences) <u>PRINTS</u>: Protein fingerprints (groups of conserved motifs)
- InterPro: protein sequence analysis & classification searches several databases
- <u>ScanProsite</u>: consists of documentation entries describing protein domains, families and functional sites as well as associated patterns and profiles to identify them
- <u>SMART</u>: (a Simple Modular Architecture Research Tool) allows the identification and annotation of genetically mobile domains and the analysis of domain architectures.

#### Membrane binding domains e.g.,

- C1, PX,
- C2, ENTH,
  - PH, and BAR domains
- FYVE,



#### The transport of molecules between the nucleus and the cytosol e.g.,

- nuclear localization signals (NLSs)
- export peptide signal
- □ import peptide signal



specific peptide tags

Protein secretion e.g.,

#### Protein PTMs binding domains e.g.,

- 14-3-3,
- WW domains,
- LRR domains of F-box proteins,
- FHA domains
- Bromodomains
- methyl-CpG-binding domain (MBD)





#### **Protein PTMs AND translocation**

Upon phosphorylation translocating to the cytosol
Up
translocating to the cytosol
translocation
on stimulation
stimulation
synaptosomes
phosphorylation

Upon dephosphorylation translocating to the cytosol

Low KCI

465

High KCI

## Upon phosphorylation translocating to the membrane



## Upon dephosphorylation translocating to the membrane

stimulation

phosphorylation



#### Conclusions

- We hope we provided an original overview of ideas for experimental design
- Data produced by mass spectrometry provide perspectives beyond fold changes
- Exploiting labeling, time series, modifications and large resources of available data can be integrated in single analyses to resolve multiple properties of the regulated proteome
- As well, co-regulation, rank of abundance, and protein 3D structure are poorly exploited aspects for data interpretation. We hope this works as a light refreshment prior to your evening fun <sup>(C)</sup>

## MetaboLights & Repository Level Workflows

Leveraging computational resource and new annotation logic to draw fresh insight



**Callum Martin** 

**Software Developer** 

ebi.ac.uk/metabolights



## What is MetaboLights?

#### MetaboLights

Open Source Study Repository & Metabolite Knowledgebase

The database is **cross-species**, **cross-technique** and covers **metabolite structures and their reference spectra** as well as their biological roles, locations and concentrations, and **experimental data** from metabolic experiments



https://www.ebi.ac.uk/metabolights

#### Submissions



Geographical distribution of submitted studies (Top 10: China: 1490, USA: 509, UK: 450, Germany: 309, France: 124, Japan: 102, Spain: 87, Italy: 80, Australia: 71, India: 66)





MetaboLights study submission rates per year



## What exactly does our data look like?

Broadly we have two categories of interest - metadata, the information that describes the experiment, and raw and derived data / files, the outputs from the experiment





MTBLS749: Alterations in the tyrosine and phenylalanine pathways revealed by biochemical profiling in cerebrospinal fluid of Huntington's disease subjects

#### Kim Kultima, Stephanie Herman

Huntington's disease (HD) is a severe neurological disease leading to psychiatric symptoms, motor impairment and cognitive decline. The disease is caused by a CAG expansion in the huntingtin (HTT) gene, but how this translates into the clinical phenotype of HD remains elusive. Using liquid chromatography mass spectrometry, we analyzed the metabolome of cerebrospinal fluid (CSF) from premanifest and manifest HD subjects as well as control subjects. Inter-group differences revealed that the tyrosine metabolism, including tyrosine, thyroxine, L-DOPA and dopamine, was significantly altered in manifest compared with premanifest HD. These metabolites demonstrated moderate to strong associations to measures of disease severity and symptoms. Thyroxine and dopamine also correlated with the five year risk of onset in premanifest HD subjects. The phenylalanine and the purine metabolisms were also significantly altered, but associated less to disease severity. Decreased levels of lumichrome were commonly found in mutated HTT carriers and the levels correlated with the five year risk of disease onset in premanifest carriers. These biochemical findings demonstrates that the CSF metabolome can be used to characterize molecular pathogenesis occurring in HD, which may be essential for future development of novel HD therapies.

#### PUBLICATIONS

Alterations in the tyrosine and phenylalanine pathways revealed by biochemical profiling i...

🚓 Herman Stephanie, Niemelä Valter, Emami Khoonsari Payam, Sun...

Ð

## Sample Information

#### Information pertaining to each individual sample processed in a study

| Descriptors Prot   | tocols Sampl | es Assays   | Metabolites | Files        |                     |           |           |            |
|--|--------------|-------------|-------------|--------------|---------------------|-----------|-----------|------------|
| File: s_MTBLS749.txt     Items per page: 494      1 - 494 of 494     1 < < > > |              |             |             |              |                     |           |           |            |
| Filter   |              |             |             |              |                     |           |           |            |
| Protocol   | REF          | Sample I    | Name        | Organism     | Organism part       | Phenotype | Replicate | MS format  |
| Sample colle   | ection 1_Re  | p1_NEG      |             | Homo sapiens | cerebrospinal fluid | control   | 1         | MS1 format |
| Sample colle   | ection 1_Re  | p1_POS      |             | Homo sapiens | cerebrospinal fluid | control   | 1         | MS1 format |
| Sample colle   | ection 1_Re  | p2_NEG      |             | Homo sapiens | cerebrospinal fluid | control   | 2         | MS1 format |
| Sample colle   | ection 1_Re  | 1_Rep2_POS  |             | Homo sapiens | cerebrospinal fluid | control   | 2         | MS1 format |
| Sample colle   | ection 10_F  | 10_Rep1_NEG |             | Homo sapiens | cerebrospinal fluid | control   | 1         | MS1 format |
| Sample colle   | ection 10_F  | 10_Rep1_POS |             | Homo sapiens | cerebrospinal fluid | control   | 1         | MS1 format |
| Sample colle   | ection 10_F  | ep2_NEG     |             | Homo sapiens | cerebrospinal fluid | control   | 2         | MS1 format |
| Sample colle   | ection 10_F  | ep2_POS     |             | Homo sapiens | cerebrospinal fluid | control   | 2         | MS1 format |
| Sample colle   | ection 11_F  | ep1_NEG     |             | Homo sapiens | cerebrospinal fluid | control   | 1         | MS1 format |
| Sample colle   | ection 11_F  | ep1_POS     |             | Homo sapiens | cerebrospinal fluid | control   | 1         | MS1 format |
| Sample colle   | ection 11_F  | ep2_NEG     |             | Homo sapiens | cerebrospinal fluid | control   | 2         | MS1 format |
| Sample colle   | ection 11_F  | ep2_POS     |             | Homo sapiens | cerebrospinal fluid | control   | 2         | MS1 format |

## Assay Information

| escriptors Protocols Samples Assays Metabolites Files  |  |  |   |  |   |                            |       |                      |  |   |                                    |
|--|--|--|---|--|---|----------------------------|-------|----------------------|--|---|------------------------------------|
| Assay Sheet 1       Assay Sheet 2       Assay Sheet 3       Assay Sheet 4         File: a_MTBLS1987_LC-MS_positive_hilic_metabolite_profiling.txt       Items per page: 69        1 - 69 of 69        < < >> |  |  |   |  |   |                            |       |                      |  |   |                                    |
| Filter<br>Protocol REF   | Parameter Value<br>-<br>Chromatography<br>Instrument | Parameter<br>Value -<br>Autosampler<br>model | Parameter<br>Value -<br>Column<br>model                                       | Parameter<br>Value -<br>Column<br>type | Parameter<br>Value -<br>Guard<br>column | Labeled<br>Extract<br>Name | Label | Protocol REF         | Parameter<br>Value -<br>Scan<br>polarity | Parameter<br>Value -<br>Scan m/z<br>range | Parameter<br>Value -<br>instrument |
| Chromatography   | Waters<br>ACQUITY<br>UPLC H-<br>Class System         |  | ACQUITY<br>UPLC<br>BEH<br>Amide<br>(1.7 µm,<br>2.1 mm x<br>100 mm;<br>Waters) | HILIC                                  |   |                            |       | Mass<br>spectrometry | Positive                                 | 50-1000                                   | Waters<br>Xevo<br>G2-S<br>QTof     |
| Chromatography   | Waters<br>ACQUITY<br>UPLC H-<br>Class System         |  | ACQUITY<br>UPLC<br>BEH<br>Amide<br>(1.7 µm,<br>2.1 mm x<br>100 mm;<br>Waters) | HILIC                                  |   |                            |       | Mass<br>spectrometry | Positive                                 | 50-1000                                   | Waters<br>Xevo<br>G2-S<br>QTof     |
| Chromatography   | Waters<br>ACQUITY<br>UPLC H-<br>Class System         |  | ACQUITY<br>UPLC<br>BEH<br>Amide<br>(1.7 µm,<br>2.1 mm x<br>100 mm;<br>Waters) | HILIC                                  |   |                            |       | Mass<br>spectrometry | Positive                                 | 50-1000                                   | Waters<br>Xevo<br>G2-S<br>QTof     |

## Metabolite Annotation Information

| escriptors | Protocols Sa                     | amples Assays             | Metabolites         | Files                 |   |                           |                   |
|------------|----------------------------------|---------------------------|---------------------|-----------------------|---|---------------------------|-------------------|
| MAF She    | eet 1 MAF Sheet                  | t 2 MAF Sheet 3           | MAF Sheet 4         | 4                     |   |                           |                   |
| File: m_MT | BLS1987_LC-MS_posit              | tive_hilic_metabolite_pro | filing_v2_maf.tsv   |                       | Items per page: 10  | 1 – 10 of 13   <          | < > >I            |
| Filter     |                                  |                           |                     |                       |   |                           |                   |
|            | Structure                        | Database<br>identifier    | Chemical<br>formula | SMILES                | InChl   | Metabolite identification | Mass to<br>charge |
|            | OH<br>NH                         | CHEBI:26271               |                     | OC(=O)C1CCCN1         | InChI=1S/C5H9NO2/c7-<br>5(8)4-2-1-3-6-4/h4,6H,1-<br>3H2,(H,7,8)               | Proline                   | 116.07029         |
|            |                                  | CHEBI:17750               |                     | C[N+](C)(C)CC([O-])=O | InChI=1S/C5H11NO2/c1-<br>6(2,3)4-5(7)8/h4H2,1-3H3                             | Betaine                   | 118.08747         |
| ,          | H <sub>L</sub> N NH <sub>2</sub> | он СНЕВІ:28300            |                     | NC(CCC(N)=0)C(0)=0    | InChI=1S/C5H10N2O3/c6-<br>3(5(9)10)1-2-4(7)8/h3H,1-<br>2,6H2,(H2,7,8)(H,9,10) | Glutamine                 | 147.07607         |
|            | Hgc S NH2                        | он СНЕВІ:16811            |                     | CSCCC(N)C(O)=O        | InChI=1S/C5H11NO2S/c1-9-<br>3-2-4(6)5(7)8/h4H,2-<br>3,6H2,1H3,(H,7,8)         | Methionine                | 150.05818         |



## MetaboLights Data & Galaxy Workflows

#### Galaxy is an open source web based computational platform that facilitates dataintensive biomedical research and analysis

- It provides a user-friendly interface for data-intensive biomedical research and analysis
- It supports reproducible research by capturing and recording the entire analysis process
- Galaxy integrates a comprehensive collection of bioinformatics and data analysis tools
- Researchers can create and execute complex computational workflows using a visual interface
   Galaxy Metabolights Labs
- We have our own instance at metabolights-labs.org
- (but you can deploy it anywhere including locally)

| Galaxy MetaboLig                               | jhts Labs | ☆ Workflow Visualize Shared Data ▼ Help ▼ Login   |  |  |  |  |
|--|-----------|---|--|--|--|--|
| Tools  | •         |   |  |  |  |  |
| search tools                                   | ×         | MetaboLights  |  |  |  |  |
| 1 Upload Data                                  |           | Welcome to the MetaboLights Labs Galaxy Instance!   |  |  |  |  |
| Get Data                                       |           |   |  |  |  |  |
| Send Data                                      |           | Take an interactive tour: Galaxy UI History Scratchbook   |  |  |  |  |
| Get Data Metabolomics                          |           |   |  |  |  |  |
| Data Handling Metabolomics                     |           | MetaboLights Labs aims at building, testing and delivering cloud-based Galaxy workfl  |  |  |  |  |
| Prepocessing LCMS                              |           | datasets with 1000's of samples and simultaneously capture all metadata associated<br>MetaboLights Labs infrastructure is formulated based on several community agreed proces |  |  |  |  |
| Annotation LCMS FIAMS<br>Quality Processing MS |           | tools. The necessary functionality is incorporated so that it can be readily used in conjunct   |  |  |  |  |
|  |           | MetaboLights.   |  |  |  |  |

## MetaboLights Workflows

C

metabolights-labs.org

#### "Standardised metabolite annotation workflows for enhancing biological interpretation in metabolomic data repositories"

| Galaxy MetaboLights   | Labs 🛪 Workflow Visualize Shared Data - Admin Help - User - 🞓 🌲 🏢   | Us   | ing 30%   |  |
|---|---|--|---|--|
| Fools 🗠 🔹   |   | ← → C  ⓐ metabolights-labs.org   |   | 🖞 🖈 🗯 🗊 🚺 🕝 (Update 🔅                                      |
| search tools ×  | MetaboLights  | Galaxy MetaboLights  | Labs 🗥 Workflow Visualize Shared Data - Admin Help - User - 🞓 🌲 🏢   | Using 30%  |
| 1 Upload Data   | Welcome to the MetaboLights Labs Galaxy Instance!   | Tools     ☆       search tools     ×   |   | History + = ·  |
| end Data<br>iet Data Metabolomics<br>lata Handling Metabolomics | Take an interactive tour: Galaxy UI History Scratchbook   | L Upload Data  | Image: Constraint of the state of the s | Unnamed history  |
| repocessing LCMS  | MetaboLights Labs aims at building, testing and delivering cloud-based Galaxy workflow(s), with the<br>computational capacity to process datasets with 1000's of samples and simultaneously capture all<br>metadata associated to allow ricorous reproducibility and reporting.                                     | xcms process history Create a<br>summary of XCMS analysis  | D     D     Nothing selected  |  |
| unnotation LCMS FIAMS<br>Quality Processing MS                  | MetaboLights Labs infrastructure is formulated based on several community agreed processing, feature<br>extraction and compound identification tools. The necessary functionality is incorporated so that it can be<br>readily used in conjunction with popular online metabolomics databases such as MetaboLights. | Annotation LCMS FIAMS<br>BEAMSpy - Birmingham mEtabolite   | Ion mode Positive   | 161 : ebi_codon_params.ya 🐵 🖋 👕<br>ml                      |
| Quality Processing All  |   | Annotation for Mass SpectroMetry<br><b>bank_inhouse</b> search by accurate<br>mass (and by Retention time) on a          | Library adducts       D     D     Nothing selected <ul> <li>E</li> <lie< li=""> <li>E</li> <lie< li=""></lie<></lie<></ul>  | 160 : MTBLS876_1r_noesy<br>gppr1d_spectralMatrix.RD<br>S   |
| ext Manipulation<br>VORKFLOWS                                   | tanipulation<br>=Lows 4   | IDCAI DARK<br>Lipidmaps : search on LIPID MAPS<br>Structure Database (LMSD) online<br>with masses and its Text/Ontology- | Group Features 🗞  | 159 : MTBLS863_1r_noesy                                    |
|   |   | based search engine.<br>MassBank spectrum searches :<br>Search by pseudo-spectra on a High                               | Maximum retention time difference (sec) 5.0 Grouping method   | 158 : MTBLS862_1r_noesy<br>gppr1d_spectralMatrix.RD<br>S   |
|   |   | Quality Mass Spectral Database.<br>Bih4MaConDa : Utility to detect<br>potential contaminants in your peak                | Pearson correlation         Coefficient threshold   | 157 : MTBLS841_1r_cpmgp ④ 🖋 🔋<br>r1d_spectralMatrix.RDS    |
|   |   | HR2 formula find a chemical  | 0.7<br>D-value threshold  | 156 : MTBLS806_1r_zgpr_s  ? 156 : MTBLS806_1r_zgpr_s       |
|   |   | LCMS matching Annotation of<br>LCMS peaks using matching on a in-  | 0.0001  | 155 : MTBLS781_1r_cpmgp 	④ ✔ 	■<br>r1d_spectralMatrix.RDS  |
|   |   | house spectra database or on<br>PeakForest spectra database.   | Annotate Peak Patterns 🗞  | 154 : MTBLS723_1r_noesyg                                   |
|   |   | Quality Processing All   | Ppm error tolerance   | 153 : MTBLS675_1r_noesyg 💿 🖍 🧃<br>ppr1d.comp_spectralMatri |
| htt   | tps://metabolights-labs.org   | Statistics All   |   | x.RDS  |

Tout Maninulation

Maximum tolerated m/z deviation in parts per million (ppm)

#### MetaboLights Labs - Workflows - MDP + XCMS + BEAMSpy



Rick Dunn (PI), Tim Ebbels (PI),

Claira O'Danayan (DI)







## MetaboLights Data Provider (MDP) Tool

- Researchers can select files from MetaboLights and apply custom msConvert filters
  - Studies: 596, Assays: 1271, Assay Files: 174565
- Creates an open-source MS standard / file for preprocessing





## MetaboLights Data Provider (MDP) Tool

#### Organism & Organism Parts & File Formats

| Homo sapiens                             | 57987<br>23914 |
|--|----------------|
|  | 23914          |
| Lolium perenne                           |                |
| Mus musculus                             | 7355           |
| Saccharomyces cerevisiae                 | 5984           |
| reference compound                       | 4623           |
| blank                                    | 2717           |
| sea water                                | 2654           |
| Vitis vinifera                           | 2193           |
| Ovis aries                               | 1948           |
| Rattus norvegicus                        | 1943           |
| groundwater                              | 1859           |
| Camellia sinensis                        | 1788           |
| Solanum lycopersicum x Solanum pennellii | 1694           |
| Gallus gallus                            | 1562           |
| Arabidopsis thaliana                     | 1465           |
| Bos taurus                               | 1427           |
| Chenopodium quinoa                       | 1407           |
| Escherichia coli                         | 1326           |
| Drosophila melanogaster                  | 1047           |
|  |                |

| Organism Parts              | Count of files |
|-----------------------------|----------------|
| leaf                        | 28844          |
| blood serum                 | 16528          |
| blood plasma                | 15909          |
| whole organism              | 10813          |
| urine                       | 6445           |
| feces                       | 3754           |
| exometabolome               | 3471           |
| endometabolome              | 2898           |
| pure substance              | 2443           |
| Colon                       | 2361           |
| mixture                     | 2293           |
| cervical mucus              | 2164           |
| blank                       | 2088           |
| Plasma                      | 1992           |
| Serum                       | 1950           |
| liver                       | 1941           |
| fruit                       | 1896           |
| Sweat                       | 1774           |
| Seed                        | 1727           |
| rosette leaf                | 1197           |
| supragingival dental plaque | 1156           |
| Wine                        | 1087           |
|                             |                |

| File Format | Count of files |
|-------------|----------------|
| .raw        | 72826          |
| .mzML       | 52556          |
| .mzXML      | 10590          |
| .d          | 3907           |
| raw.zip     | 1262           |

## Galaxy Workflow: XCMS (& MSnbase)

- Individual XCMS steps are structured into a workflow to extract and match peaks across files and generate data matrix (m/z vs RT vs intensity)
- 4 parameter sets based on 'assay type' & 'column model'
  - 50 cm length, UPLC
  - 100 cm length, UPLC
  - 150 cm length, UHPLC
  - 150 cm length, HPLC
- Optimisation in collaboration with Rick Dunn (UoL)





## Galaxy Workflow: XCMS (& MSnbase)

#### XCMS Workflow for mzML Dataset Collection - Default



## BEAMSpy

BEAMSpy is a Python package that includes **several automated and seamless computational modules** that are applied to putatively annotate metabolites detected in **untargeted ultra (high) performance** liquid chromatography-mass spectrometry or **untargeted direct infusion** mass spectrometry metabolomic assays.

> ∃ README.rst 0 BEAMSpy - Birmingham mEtabolite Annotation for Mass Spectrometry (Python package) pypi v1.2.0 python 3.8 | 3.9 | 3.10 repository GitHub install with bioconda C beamspy passing License GPL v3 🤮 launch binder documentation RTD codecov 84% BEAMSpy (Birmingham mEtabolite Annotation for Mass Spectrometry) is a Python package that includes several automated and seamless computational modules that are applied to putatively annotate metabolites detected in untargeted ultra (high) performance liquid chromatography-mass spectrometry or untargeted direct infusion mass spectrometry metabolomic assays. All reported metabolites are annotated to level 2 or 3 of the Metabolomics Standards Initiative (MSI) reporting standards (Metabolomics. 2007 Sep; 3(3): 211-221. doi: 10.1007/s11306-007-0082-2). The package is highly flexible to suit the diversity of sample types studied and mass spectrometers applied in untargeted metabolomics studies. The user can use the standard reference files included in the package or can develop their own reference files.





MTBLS749: Alterations in the tyrosine and phenylalanine pathways revealed by biochemical profiling in cerebrospinal fluid of Huntington's disease subjects

#### Kim Kultima, Stephanie Herman

Huntington's disease (HD) is a severe neurological disease leading to psychiatric symptoms, motor impairment and cognitive decline. The disease is caused by a CAG expansion in the huntingtin (HTT) gene, but how this translates into the clinical phenotype of HD remains elusive. Using liquid chromatography mass spectrometry, we analyzed the metabolome of cerebrospinal fluid (CSF) from premanifest and manifest HD subjects as well as control subjects. Inter-group differences revealed that the tyrosine metabolism, including tyrosine, thyroxine, L-DOPA and dopamine, was significantly altered in manifest compared with premanifest HD. These metabolites demonstrated moderate to strong associations to measures of disease severity and symptoms. Thyroxine and dopamine also correlated with the five year risk of onset in premanifest HD subjects. The phenylalanine and the purine metabolisms were also significantly altered, but associated less to disease severity. Decreased levels of lumichrome were commonly found in mutated HTT carriers and the levels correlated with the five year risk of disease onset in premanifest carriers. These biochemical findings demonstrates that the CSF metabolome can be used to characterize molecular pathogenesis occurring in HD, which may be essential for future development of novel HD therapies.

#### PUBLICATIONS

Alterations in the tyrosine and phenylalanine pathways revealed by biochemical profiling i...

🚓 Herman Stephanie, Niemelä Valter, Emami Khoonsari Payam, Sun...

Ð

#### MetaboLights Labs - Workflows - MDP + XCMS + BEAMSpy



Rick Dunn (PI), Tim Ebbels (PI),

Claira O'Danayan (DI)







## In the Future

In the future we plan to:

- Integrate our workflows directly and seamlessly into submission process
- Develop workflows for the processing and annotation of MS/MS data, and facilitate additional analyses such as molecular networking
- Integrate MetaboLights and GNPS (as well as other repositories such as Metabolomics Workbench) to enable cross repository analysis
- Work with other communities to map standards (data and metadata) for example proteomics and SDRF (PRIDE)





## Acknowledgements

**EMBL-EBI** Metabolomics team

- Claire O'Donovan (PI)
- Mark Williams
- Thomas Payne
- Noemi Tejera Hernandez
- Felix Amaladoss
- Callum Martin
- Ozgur Yurekten



The Metabolomics team's activities were supported in the past year by the above funding bodies

