# Mapping the sub-cellular proteome

Laurent Gatto

lg390@cam.ac.uk – @lgatt0

http://www.damtp.cam.ac.uk/user/lg390/

Slides @ https://zenodo.org/record/1063508

DOI 10.5281/zenodo.1063508

22 Nov 2017, Cambridge Computational Biology Institute

# Plan

# Regulations

# Cell organisation



**Spatial proteomics** is the systematic study of protein localisations.

# Spatial proteomics - Why?

## Localisation is function

- The cellular sub-division allows cells to establish a range of distinct micro-environments, each favouring different biochemical reactions and interactions and, therefore, allowing each compartment to fulfil a particular functional role.
- Localisation and sequestration of proteins within sub-cellular niches is a fundamental mechanism for the post-translational regulation of protein function.

## Re-localisation in

- Differentiation: Tfe3 in mouse ESC (Betschinger et al., 2013).
- Activation of biological processes.

Examples later.

# Spatial proteomics - Why?

### Mis-localisation
Disruption of the targeting/trafficking process alters proper sub-cellular localisation, which in turn perturb the cellular functions of the proteins.

- Abnormal protein localisation leading to the loss of functional effects in diseases (Laurila and Vihinen, 2009).
- Disruption of the nuclear/cytoplasmic transport (nuclear pores) have been detected in many types of carcinoma cells (Kau et al., 2004).

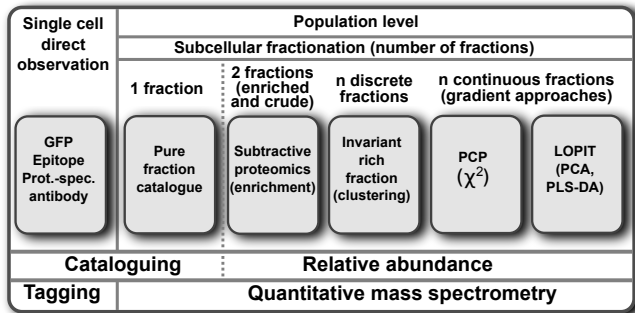# Spatial proteomics - How, experimentally



Figure : Organelle proteomics approaches (Gatto et al., 2010)

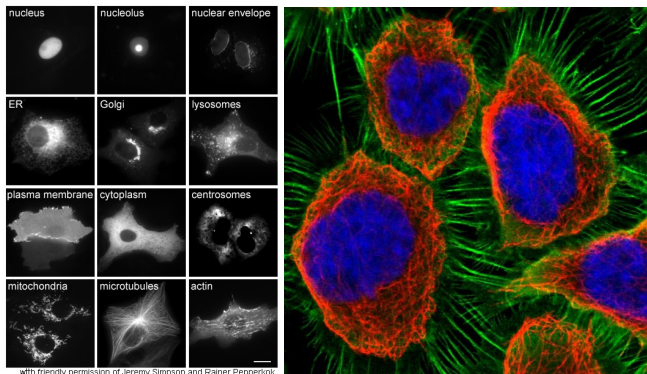# Fusion proteins and immunofluorescence



Figure : Targeted protein localisation.
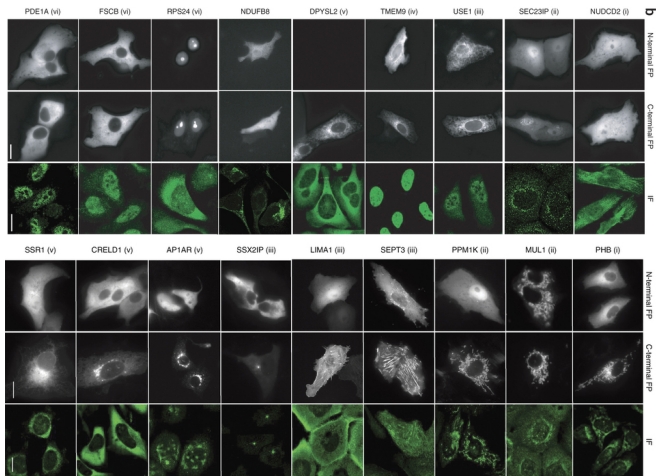
# Fusion proteins and immunofluorescence



Figure : Example of discrepancies between IF and FPs as well as between FP tagging at the N and C termini (Stadler et al., 2013).

# Spatial proteomics - How, experimentally
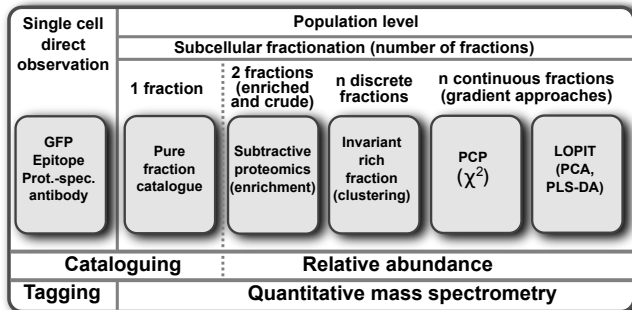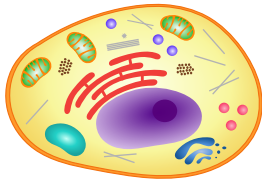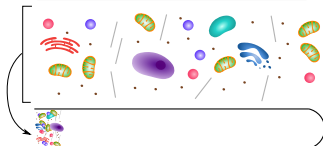


Figure : Organelle proteomics approaches (Gatto et al., 2010). Gradient approaches: Dunkley et al. (2006), Foster et al. (2006).

$\Rightarrow$ **Explorative/discovery approches**, steady-state **global localisation maps**.

# Quantitation data and organelle markers

|  | Fraction$_1$ | Fraction$_2$ | ... | Fraction$_m$ | markers |
|---|---|---|---|---|---|
| $p_1$ | $q_{1,1}$ | $q_{1,2}$ | ... | $q_{1,m}$ | unknown |
| $p_2$ | $q_{2,1}$ | $q_{2,2}$ | ... | $q_{2,m}$ | $loc_1$ |
| $p_3$ | $q_{3,1}$ | $q_{3,2}$ | ... | $q_{3,m}$ | unknown |
| $p_4$ | $q_{4,1}$ | $q_{4,2}$ | ... | $q_{4,m}$ | $loc_i$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $p_j$ | $q_{j,1}$ | $q_{j,2}$ | ... | $q_{j,\,m}$ | unknown |

# Visualisation and classification



Figure : From Gatto et al. (2010), *Arabidopsis thaliana* data from Dunkley et al. (2006)

# Data analysis



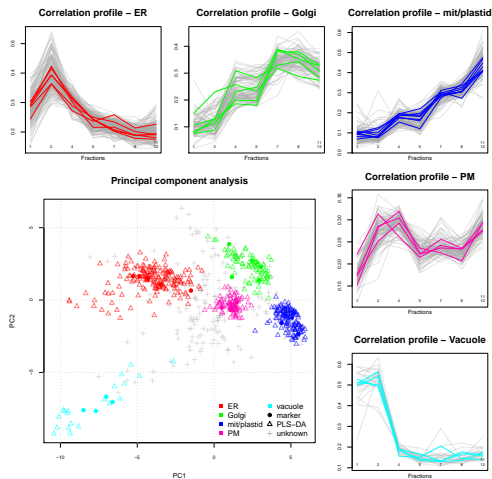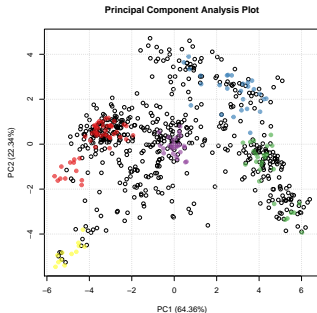|  | Fraction$_1$ | Fraction$_2$ | ... | Fraction$_m$ |  | markers |  |
|---|---|---|---|---|---|---|---|
| prot$_1$ | q$_{1,1}$ | q$_{1,2}$ | ... | q$_{1,\,m}$ | ... | $unknown$ | ... |
| prot$_2$ | q$_{2,1}$ | q$_{2,2}$ | ... | q$_{2,\,m}$ |  | $organelle_1$ |  |
| prot$_3$ | q$_{3,1}$ | q$_{3,2}$ | ... | q$_{3,\,m}$ |  | $unknown$ |  |
| prot$_4$ | q$_{4,1}$ | q$_{4,2}$ | ... | q$_{4,\,m}$ |  | $organelle_2$ |  |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |
| prot$_i$ | q$_{i,1}$ | q$_{i,2}$ | ... | q$_{i,\,m}$ |  | $organelle_k$ |  |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |
| prot$_n$ | q$_{n,1}$ | q$_{n,2}$ | ... | q$_{n,\,m}$ | ... | $unknown$ | ... |
|  | Fraction$_1$ | Fraction$_2$ | ... | Fraction$_m$ |  |  |  |
|  | ... | ... | ... | ... |  |  |  |
|  | ⋮ | ⋮ | ⋮ | ⋮ |  |  |  |
|  | ... | ... | ... | ... |  |  |  |

## Supervised machine learning

Using labelled marker proteins to match unlabelled proteins (of
unknown localisation) with similar profiles and classify them as
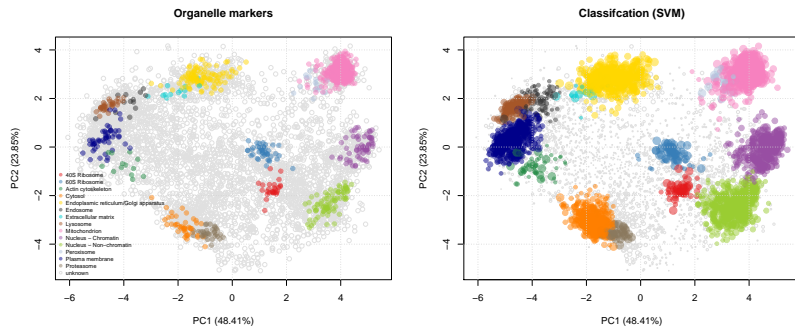residents to the markers organelle class.
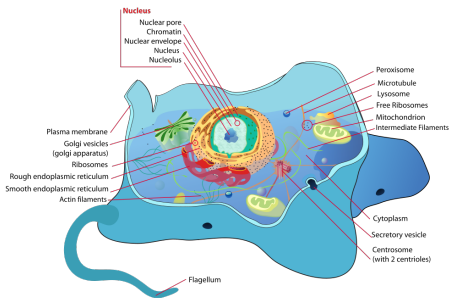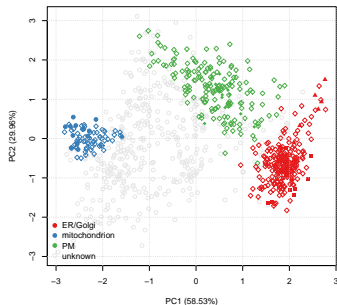
# Supervised ML



Figure : Support vector machines classifier (after classification cutoff) on the embryonic stem cell data from Christoforou et al. (2016).

# Importance of annotation



Incomplete annotation, and therefore lack of training data, for many/most organelles. *Drosophila* data from Tan et al. (2009).
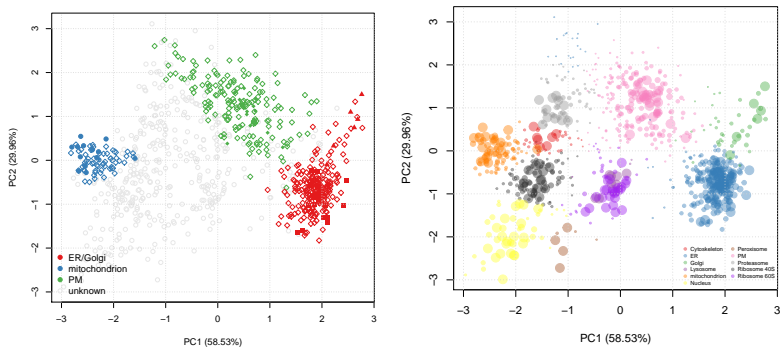
# Semi-supervised learning: novelty detection



Figure : Left: Original *Drosophila* data from Tan et al. (2009). Right: After semi-supervised learning and classification, Breckels et al. (2013).

# Improving on LOPIT

Improving is obtaining better sub-cellular resolution to increase the number of protein that can be **confidently** assigned to a sub-cellular niche.



Figure : E14TG2a embryonic stem cells: old (left) *vs.* new, better resolved (right) experiments (Christoforou et al. (2016)).

# Improving on LOPIT

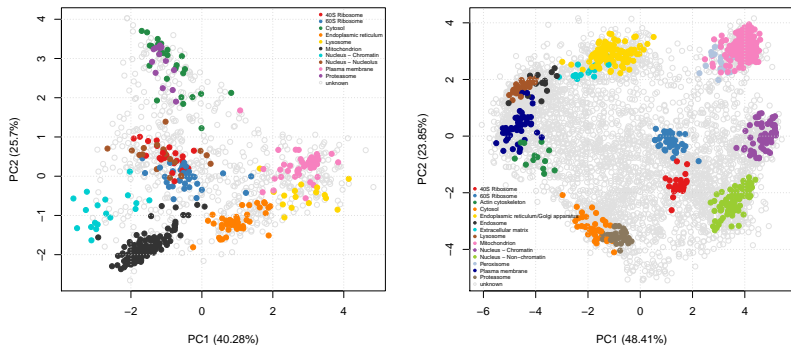Improving is obtaining better sub-cellular resolution to increase the number of protein that can be **confidently** assigned to a sub-cellular niche $\Rightarrow$ **biological discoveries**.

| | |
|---|---|
| LOPIT<br>Dunkley et al. (2006)<br>Gatto et al. (2014a) | **Computational**:<br>*transfer learning*<br>Breckels et al. (2016a) |
| **Experimental**:<br>*hyperLOPIT*<br>Christoforou et al. (2016)<br>Mulvey et al. (2017)<br>Breckels et al. (2016b) | Biological<br>discoveries |

# Experimental advances: hyperLOPIT



Figure : From Mulvey et al. (2017) *Using hyperLOPIT to perform high-resolution mapping of the spatial proteome.*
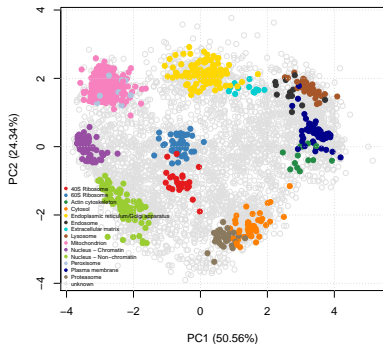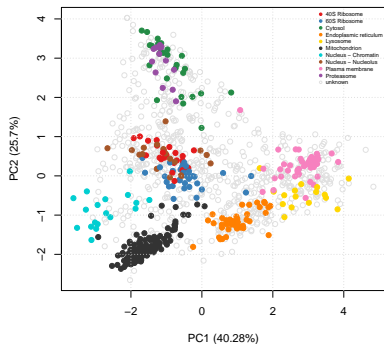
Figure : E14TG2a LOPIT on 8 fractions (using iTRAQ 8-plex) and 1109 proteins *vs.* hyperLOPIT on 10 fractions (using TMT 10-plex) and SPS-MS$^3$ for 5032 proteins.

# Computational advances: Transfer learning

What about using **addition data**, such as annotations from the Gene Ontogy (GO), sequence features (pseudo aminoacid composition), signal peptide, trans-membrane domains (length, number, ...), images (IF, FP), prediction software, . . .

- From a user perspective: **"free/cheap"** vs. expensive and time-consuming experiments.
- Abundant (all proteins, 100s of features) vs. (experimentally) limited/**targeted** (1000s of proteins, 6 − 20 of features)
- For localisation in system at hand: *low* vs. high **quality**
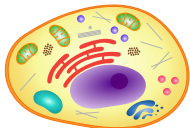- **Static** vs. **dynamic**

# Transfer learning

What about annotation data from repositories such as the
Gene Ontology (GO), sequence features, signal peptide,
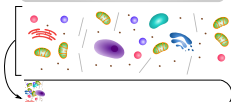transmembrane domains, images, prediction software, . . .

### Transfer learning

Support/complement the **primary** target domain (experimental
data) with **auxiliary** data (annotation, imaging, PPI, ...) features
without compromising the integrity of our primary data (Breckels
et al., 2016a).
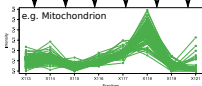
PRIMARY EXPERIMENTAL DATA
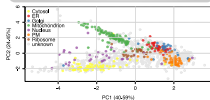
Cell lysis

Fractionation/centrifugation

e.g. Mitochondrion
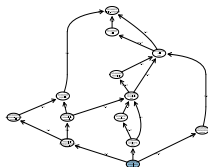
Quantitation/identification by mass spectrometry

e.g. Mitochondrion

Visualisation

AUXILIARY DRY DATA

Database query

Extract GO CC terms

Convert terms to binary

Visualisation

**Transfer learnig**, based on Wu and Dietterich (2004):

Class-weighted kNN

$$V(c_i)_j = \theta^* n_{ij}^P + (1 - \theta^*) n_{ij}^A$$



Linear programming SVM

$$f(\mathbf{x}, \mathbf{v}; \boldsymbol{\alpha}_P, \boldsymbol{\alpha}_A, b) = \sum_{l=1}^{m} y_l \left[ \alpha_l^P K^P(\mathbf{x}_l, \mathbf{x}) + \alpha_l^A K^A(\mathbf{v}_l, \mathbf{v}) \right] + b$$

Data from mouse stem cells (E14TG2a).

Figure : From Breckels et al. (2016a) *Learning from heterogeneous data sources: an application in spatial proteomics*.

# Biological discoveries

- Multi-localisation
- Trans-localisation

**Dependent on good sub-cellular resolution.**

**Dual-localisation** Proteins may be present simultaneously in several organelles (e.g. trafficking). Simulation on *A. thaliana* data from Dunkley et al. (2006) (Gatto et al., 2014b) (left). Example from embryonic stem cells (Christoforou et al., 2016) (right).

**Dual-localisation** Proteins may be present simultaneously in several organelles (e.g. trafficking). Simulation on *A. thaliana* data from Dunkley et al. (2006) (Gatto et al., 2014b) (left). Example from embryonic stem cells (Christoforou et al., 2016) (right).



From Betschinger et al. (2013)

# Spatial dynamics

### Trans-localisation event during monocyte to macrophage differenciation

Investigate the effect of LPS-mediated inflammatory response in human monocytic cells (THP-1)

### Data

- Triplicate **temporal** profiling (0, 2, 4, 6, 12, 24 hours).
- Triplicate **spatial** profiling (0 vs 12 hours) - early trafficking, before actual morphological differentiation at 24h.

Work lead by **Dr Claire Mulvey**, Cambridge Centre for Proteomics.

Figure : Spatial maps: unstimulated and LPS-treated.

Figure : Relocation of Protein Kinase C alpha and beta from the cytosol to the plasma membrane, **driving maturation into a differentiated macrophage phenotype**.

Figure : Relocation of Signal transducer and activator of transcription 6 (STAT6) from the cytosol to the Nucleus, **activating anti-bacterial and anti-viral-like response**. Validated by microscopy and see also Chen et al. (2011).

# Beyond organelles: application to PPI/Protein complexes



Figure : Data on proteasome complexes from Fabre *et al.* Mol Syst Biol
(2015), DOI: 10.15252/msb.20145497

# Plan

Open development: R/Bioconductor software

But none of this would matter if it wasn't **reproducible**!

Try it out yourselves:

```
> source("http://www.bioconductor.org/biocLite.R")
Bioconductor version 3.6 (BiocInstaller 1.28.0), ?biocLite for help
> BiocInstaller::biocLite(c("pRoloc", "pRolocdata"))
BioC_mirror: https://bioconductor.org
Using Bioconductor 3.6 (BiocInstaller 1.28.0), R 3.4.2 Patched (2017-10-12
  r73548).
*** output flushed ***
> library("pRoloc")
> library("pRolocdata")
> data(hyperLOPIT2015)
> plot2D(hyperLOPIT2015)
```

R/Bioconductor:

- ► Software for spatial proteomics.
- ► Ecosystem for high throughput biology data analysis and comprehension.

# Software for mass spectrometry and (spatial) proteomics

**Bioconductor** Open source, enable **reproducible research**, enables understanding of the data (not a black box) and **drive scientific innovation**.

- ▶ `mzR` – low level access to raw and identification mass spectrometry data (Chambers and et al., 2012)

- ▶ `MSnbase` – infrastructure to handle quantitative data and meta-data (Gatto and Lilley, 2012) (~500 unique IP download/month in 2016).

- ▶ `pRoloc` and `pRolocGUI` – dedicated visualisation and ML infrastructure for spatial proteomics (Gatto et al., 2014a) (~200 unique IP download/month in 2016). Try it out at `https://lgatto.shinyapps.io/christoforou2015/`

- ▶ `pRolocdata` – structured and annotated spatial proteomics data (Gatto et al., 2014a).

- ▶ And more generally `RforProteomics` (Gatto and Christoforou, 2014) (~160 unique IP download/month in 2016).

http://www.bioconductor.org

**Bioconductor** Open source, and **coordinated open development**, enabling **reproducible research**, enables understanding of the data (not a black box) and **drive scientific innovation**.

- ▶ Bioconductor core team (lead by Dr. Martin Morgan)
- ▶ Common infrastructure
- ▶ Common documentation standards
- ▶ Common testing infrastructure
- ▶ Open package technical peer review

Quick getting started guide: [https://lgatto.github.io/2017_11_09_Rcourse_Jena/navigating-the-bioconductor-project.html](https://lgatto.github.io/2017_11_09_Rcourse_Jena/navigating-the-bioconductor-project.html)
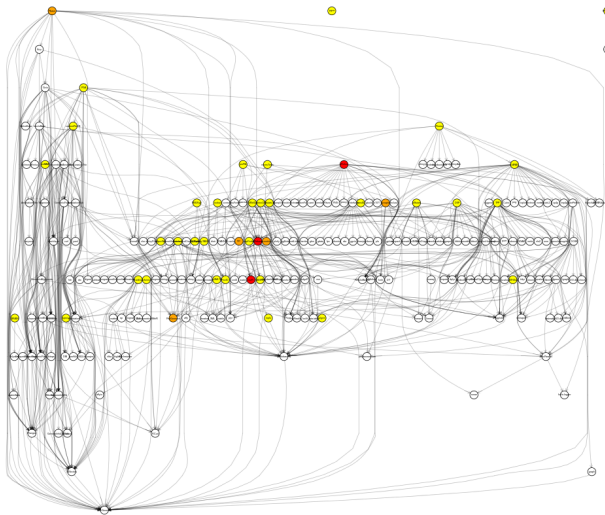
Figure : Dependency graph containing 41 MS and proteomics-tagged packages (out of 100+) and their dependencies. Showing all packages and deps would produce a big hairball.
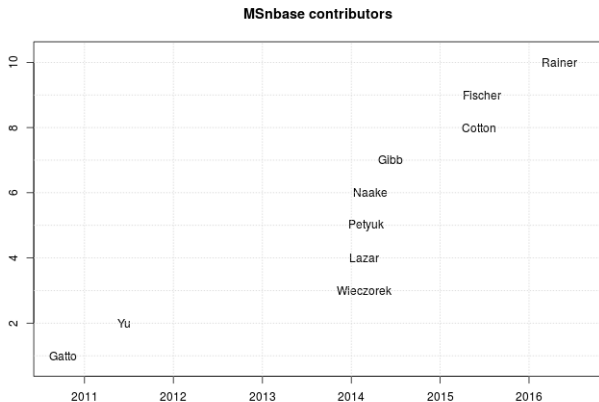
# MSnbase example



Figure : Contributions to the `MSnbase` package since its creation, the last one leading to **common proteomics/metabolomics infrastructure**. More details: https://lgatto.github.io/msnbase-contribs/

# References I

J Betschinger, J Nichols, S Dietmann, P D Corrin, P J Paddison, and A Smith. Exit from pluripotency is gated by intracellular redistribution of the bhlh transcription factor tfe3. *Cell*, 153(2):335–47, Apr 2013. doi: 10.1016/j.cell.2013.03.012.

L M Breckels, S B Holden, D Wojnar, C M Mulvey, A Christoforou, A Groen, M W Trotter, O Kohlbacher, K S Lilley, and L Gatto. Learning from heterogeneous data sources: An application in spatial proteomics. *PLoS Comput Biol*, 12(5):e1004920, May 2016a. doi: 10.1371/journal.pcbi.1004920.

LM Breckels, L Gatto, A Christoforou, AJ Groen, KS Lilley, and MW Trotter. The effect of organelle discovery upon sub-cellular protein localisation. *J Proteomics*, 88:129–40, Aug 2013.

LM Breckels, CM Mulvey, KS Lilley, and L Gatto. A bioconductor workflow for processing and analysing spatial proteomics data [version 1; referees: awaiting peer review]. *F1000Research*, 5(2926), 2016b. doi: 10.12688/f1000research.10411.1.

MC Chambers and et al. A cross-platform toolkit for mass spectrometry and proteomics. *Nat Biotechnol*, 30(10): 918–20, Oct 2012.

H Chen, H Sun, F You, W Sun, X Zhou, L Chen, J Yang, Y Wang, H Tang, Y Guan, W Xia, J Gu, H Ishikawa, D Gutman, G Barber, Z Qin, and Z Jiang. Activation of stat6 by sting is critical for antiviral innate immunity. *Cell*, 147(2):436–46, Oct 2011. doi: 10.1016/j.cell.2011.09.022.

A Christoforou, C M Mulvey, L M Breckels, A Geladaki, T Hurrell, P C Hayward, T Naake, L Gatto, R Viner, A Martinez Arias, and K S Lilley. A draft map of the mouse pluripotent stem cell spatial proteome. *Nat Commun*, 7:8992, Jan 2016. doi: 10.1038/ncomms9992.

TPJ Dunkley, S Hester, IP Shadforth, J Runions, T Weimar, SL Hanton, JL Griffin, C Bessant, F Brandizzi, C Hawes, RB Watson, P Dupree, and KS Lilley. Mapping the Arabidopsis organelle proteome. *PNAS*, 103(17): 6518–6523, Apr 2006.

LJ Foster, CL de Hoog, Y Zhang, Y Zhang, X Xie, VK Mootha, and M Mann. A mammalian organelle map by protein correlation profiling. *Cell*, 125(1):187–199, Apr 2006.

L Gatto and A Christoforou. Using R and Bioconductor for proteomics data analysis. *Biochim Biophys Acta*, 1844 (1 Pt A):42–51, Jan 2014.

L Gatto and KS Lilley. MSnbase - an R/Bioconductor package for isobaric tagged mass spectrometry data visualization, processing and quantitation. *Bioinformatics*, 28(2):288–9, Jan 2012.

L Gatto, JA Vizcaino, H Hermjakob, W Huber, and KS Lilley. Organelle proteomics experimental designs and analysis. *Proteomics*, 2010.

L Gatto, L M Breckels, S Wieczorek, T Burger, and K S Lilley. Mass-spectrometry based spatial proteomics data analysis using pRoloc and pRolocdata. *Bioinformatics*, Jan 2014a.

L Gatto, LM Breckels, T Burger, DJ Nightingale, AJ Groen, C Campbell, N Nikolovski, CM Mulvey, A Christoforou, M Ferro, and KS Lilley. A foundation for reliable spatial proteomics data analysis. *MCP*, 13(8): 1937–52, Aug 2014b.

TR Kau, JC Way, and PA Silver. Nuclear transport and cancer: from mechanism to intervention. *Nat Rev Cancer*, 4(2):106–17, Feb 2004.

K Laurila and M Vihinen. Prediction of disease-related mutations affecting protein localization. *BMC Genomics*, 10:122, 2009.

C M Mulvey, L M Breckels, A Geladaki, N K Britov?ek, DJH Nightingale, A Christoforou, M Elzek, M J Deery, L Gatto, and K S Lilley. Using hyperlopit to perform high-resolution mapping of the spatial proteome. *Nat Protoc*, 12(6):1110–1135, Jun 2017. doi: $10.1038/\text{nprot}.2017.026$.

DJL Tan, H Dvinge, A Christoforou, P Bertone, A Arias Martinez, and KS Lilley. Mapping organelle proteins and protein complexes in Drosophila melanogaster. *J Proteome Res*, 8(6):2667–2678, Jun 2009.

P Wu and TG Dietterich. Improving svm accuracy by training on auxiliary data sources. In *Proceedings of the Twenty-first International Conference on Machine Learning*, ICML '04, New York, NY, USA, 2004. ACM.

## Acknowledgements

- **Lisa Breckels**, Computational Proteomics Unit, Cambridge (ML, algo)
- **Kathryn Lilley**, Cambridge Centre of Proteomics (Proteomics)
- **Funding**: BBSRC, Wellcome Trust

- Slides: https://zenodo.org/record/1063508
- License:

**Thank you for your attention**