

# Probabilistic modelling of protein sub-cellular localisation

Laurent Gatto

[laurent.gatto@uclouvain.be](mailto:laurent.gatto@uclouvain.be)

de Duve Institute – UCLouvain

29 March 2019 – ISBA

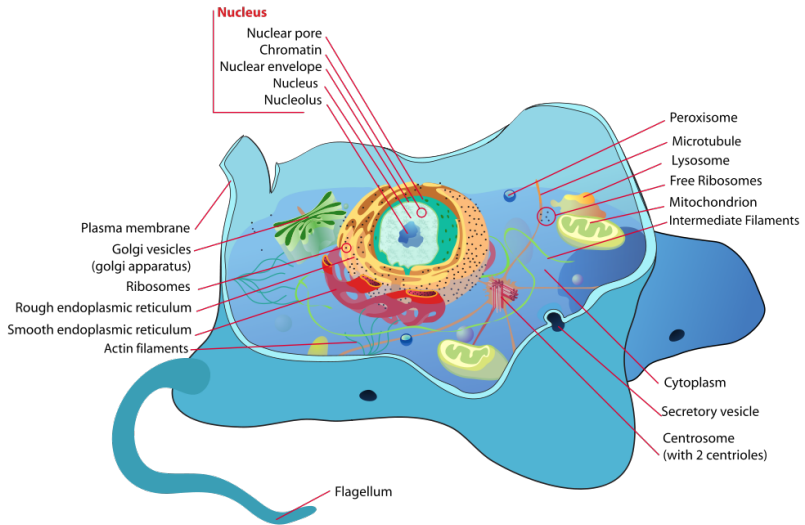
# Abstract

In biology, **localisation is function** - understanding the sub-cellular localisation of proteins is paramount to comprehend the context of their functions. Mass spectrometry-based **spatial proteomics** and contemporary **machine learning** enable to build proteome-wide spatial maps, informing us on the location of thousands of proteins. Nevertheless, while some proteins can be found in a single location within a cell, up to half of proteins may reside in multiple locations, can dynamically re-localise, or reside within an unknown functional compartment, leading to considerable **uncertainty** in associating a protein to their sub-cellular location. Recent advances enable us to **probabilistically** model protein localisation as well as quantify the uncertainty in the location assignments, thus leading to better and more trustworthy biological interpretation of the data.

1. **Use case:** spatial proteomics.
2. Novel **computational biology research and developments** to acquire reliable biological knowledge.
3. **Behind the scenes:** software/data structures and open research practice.

**Use case:** spatial proteomics.

# Cell organisation - regulation of protein localisation



**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** stem cells.
- ▶ **Activation** of biological processes.

# Spatial proteomics - Why?

## Mis-localisation

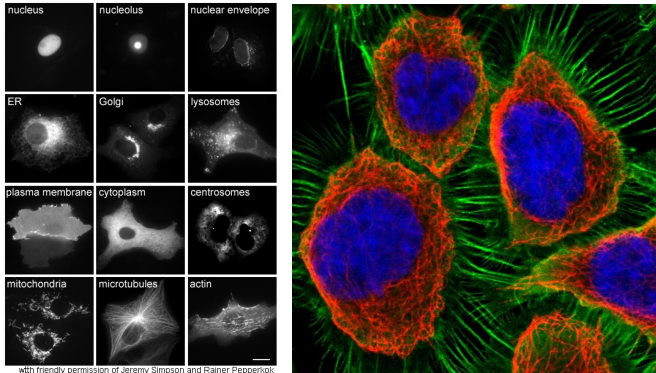
Disruption of the targeting/trafficking process alters proper sub-cellular localisation, which in turn perturb the cellular functions of the proteins.

- ▶ 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](#)).
- ▶ Sub-cellular localisation of MC4R with ADCY3 at neuronal primary cilia underlies a common pathway for genetic predisposition to **obesity** ([Siljee et al., 2018](#)).





# Fusion proteins and immunofluorescence



**Figure:** Targeted protein localisation. 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

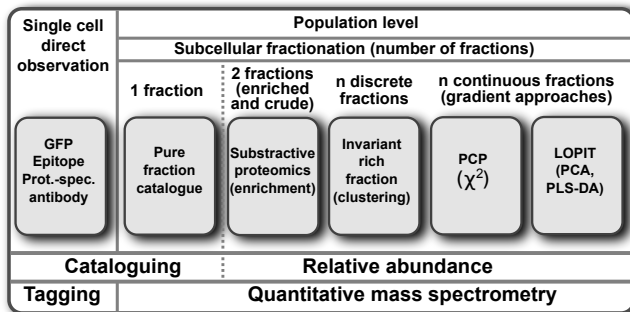
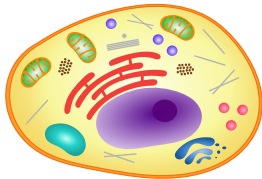


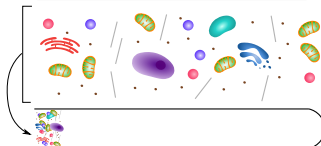
Figure: Organelle proteomics approaches (Gatto et al., 2010).

Gradient approaches: Dunkley et al. (2006), Foster et al. (2006).

Explorative/discovery approaches, steady-state global localisation maps.

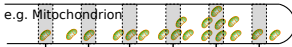


Cell lysis



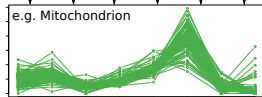
Fractionation/centrifugation

e.g. Mitochondrion



Quantitation/identification  
by mass spectrometry

e.g. Mitochondrion



# Quantitation data

	Fraction <sub>1</sub>	Fraction <sub>2</sub>	...	Fraction <sub>m</sub>
p <sub>1</sub>	q <sub>1,1</sub>	q <sub>1,2</sub>	...	q <sub>1,m</sub>
p <sub>2</sub>	q <sub>2,1</sub>	q <sub>2,2</sub>	...	q <sub>2,m</sub>
p <sub>3</sub>	q <sub>3,1</sub>	q <sub>3,2</sub>	...	q <sub>3,m</sub>
p <sub>4</sub>	q <sub>4,1</sub>	q <sub>4,2</sub>	...	q <sub>4,m</sub>
⋮	⋮	⋮	⋮	⋮
p <sub>j</sub>	q <sub>j,1</sub>	q <sub>j,2</sub>	...	q <sub>j, m</sub>

## Quantitation data and organelle markers

	Fraction <sub>1</sub>	Fraction <sub>2</sub>	...	Fraction <sub>m</sub>	markers
p <sub>1</sub>	q <sub>1,1</sub>	q <sub>1,2</sub>	...	q <sub>1,m</sub>	unknown
p <sub>2</sub>	q <sub>2,1</sub>	q <sub>2,2</sub>	...	q <sub>2,m</sub>	<i>loc<sub>1</sub></i>
p <sub>3</sub>	q <sub>3,1</sub>	q <sub>3,2</sub>	...	q <sub>3,m</sub>	unknown
p <sub>4</sub>	q <sub>4,1</sub>	q <sub>4,2</sub>	...	q <sub>4,m</sub>	<i>loc<sub>i</sub></i>
⋮	⋮	⋮	⋮	⋮	⋮
p <sub>j</sub>	q <sub>j,1</sub>	q <sub>j,2</sub>	...	q <sub>j, m</sub>	unknown

# Data analysis

# Data analysis

- ▶ **Visualisation** (cluster, unsupervised learning)
- ▶ **Classification** (supervised learning)
- ▶ **Novelty detection** (semi-supervised learning)
- ▶ Data integration (transfer learning)
- ▶ **Multi-localisation (Bayesian spatial proteomics)**
- ▶ Spatial dynamics

To uncover and understand biology

# Visualisation

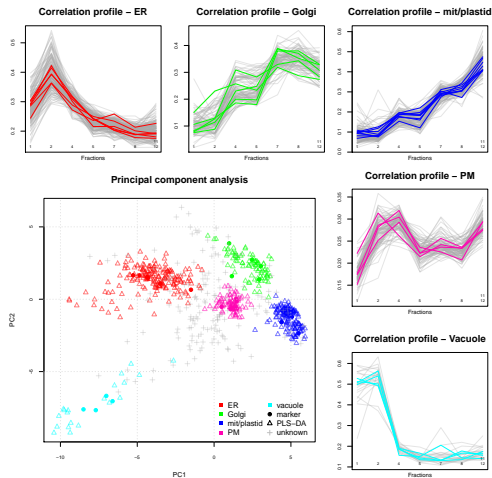
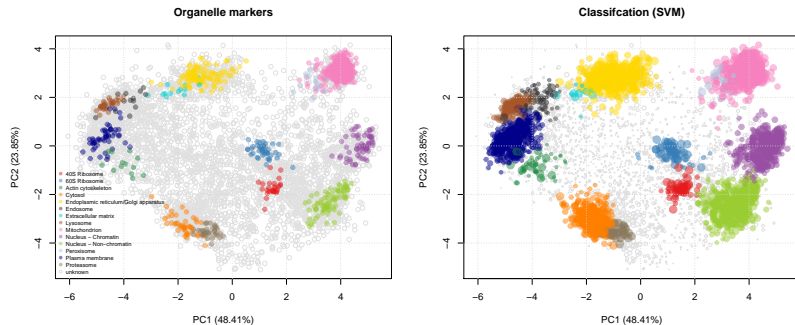


Figure: From Gatto et al. (2010), *Arabidopsis thaliana* data from Dunkley et al. (2006)



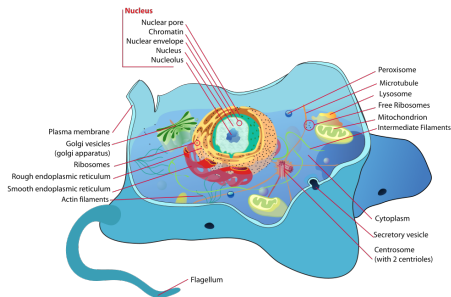
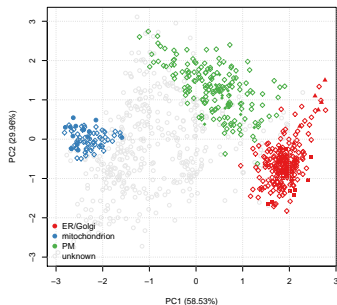
# Supervised Machine Learning



**Figure:** Support vector machines classifier (after 5% FDR classification cutoff) on the embryonic stem cell data from [Christoforou et al. \(2016\)](#).

Novel **computational biology research and developments** to acquire reliable biological knowledge.

# 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

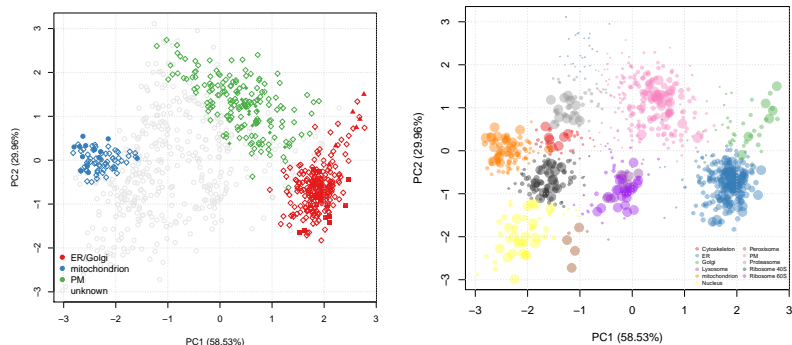
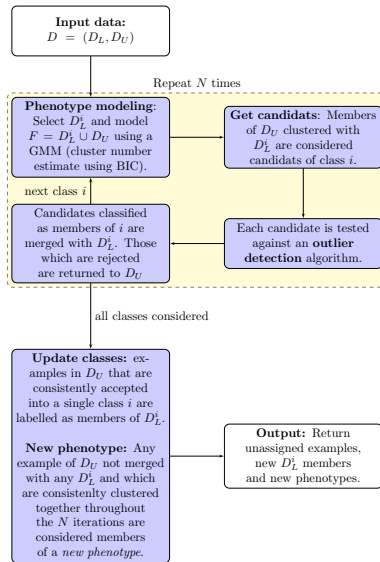


Figure: Left: Original *Drosophila* data from Tan et al. (2009). Right: After semi-supervised learning and classification, Breckels et al. (2013).



**Figure:** The effect of organelle discovery upon sub-cellular protein localisation Breckels et al. (2013).

# Computational advances: Transfer learning

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

Support/complement the **primary** target domain (experimental data) with **auxiliary** data (annotation, imaging, PPI, ...) features without compromising the integrity of our primary data.





# Transfer learning results

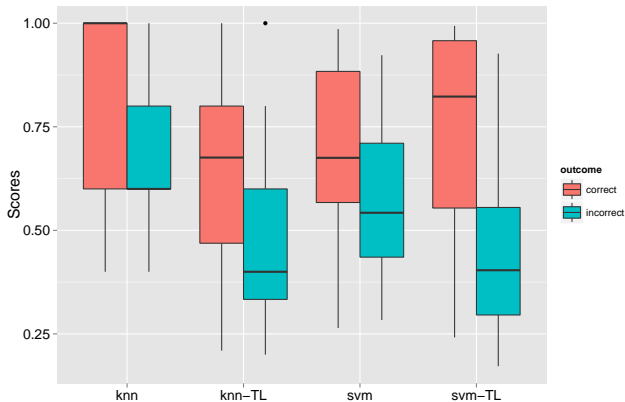
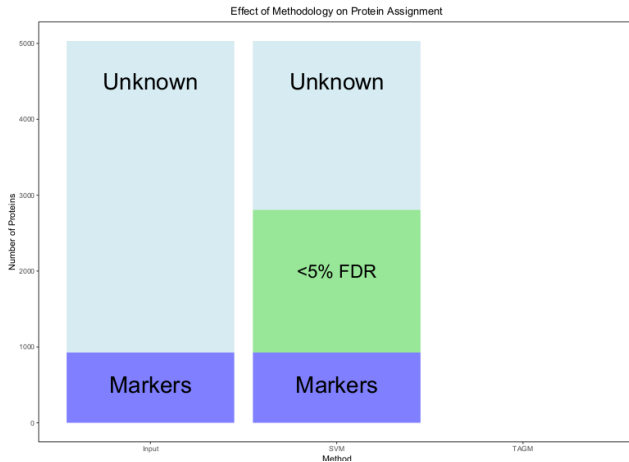


Figure: From Breckels et al. (2016) *Learning from heterogeneous data sources: an application in spatial proteomics*.

# How much do we learn? How much do we miss?



# A Bayesian Mixture Modelling Approach For Spatial Proteomics

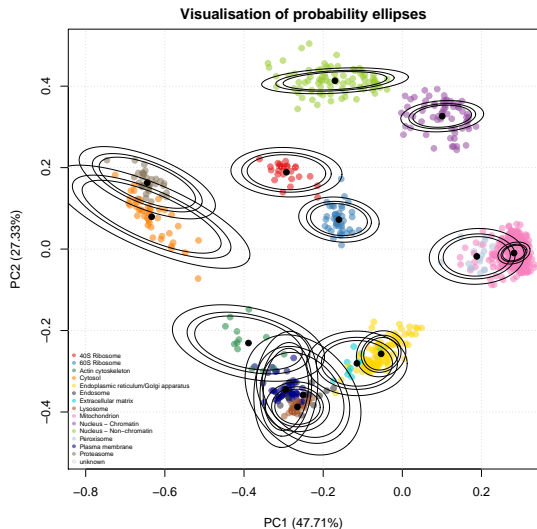
- ▶ *T Augmented Gaussian Mixture model (TAGM)* is a **multivariate Gaussian generative model** for MS-based spatial proteomics data. It posits that each annotated sub-cellular niche can be modelled by a multivariate Gaussian distribution.

# A Bayesian Mixture Modelling Approach For Spatial Proteomics

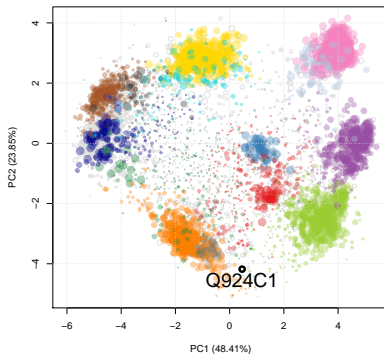
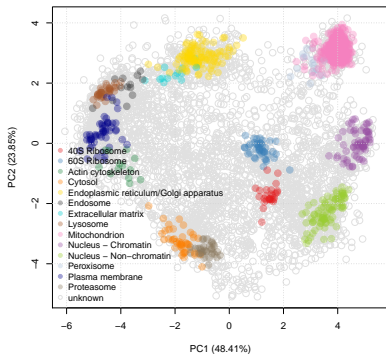
- ▶ *T Augmented Gaussian Mixture model (TAGM)* is a **multivariate Gaussian generative model** for MS-based spatial proteomics data. It posits that each annotated sub-cellular niche can be modelled by a multivariate Gaussian distribution.
- ▶ With the prior knowledge that many proteins are not captured by known sub-cellular niches, we augment our model with an **outlier component**. Outliers are often dispersed and thus this additional component is described by a heavy-tailed distribution: the multivariate Student's t-distribution, leading us to a *T Augmented Gaussian Mixture model*.

# A Bayesian Mixture Modelling Approach For Spatial Proteomics

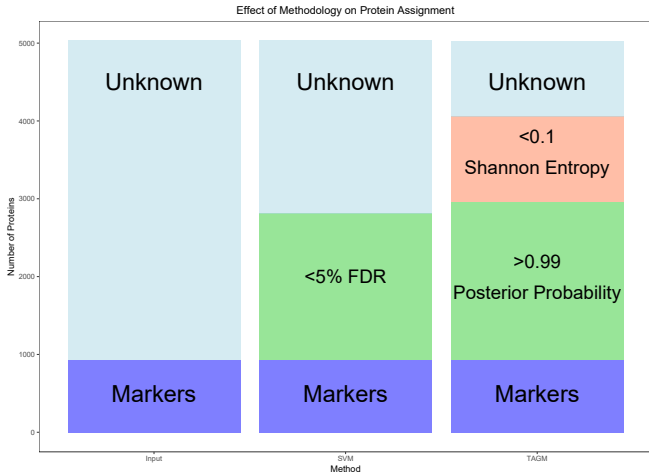
- ▶ *T Augmented Gaussian Mixture model (TAGM)* is a **multivariate Gaussian generative model** for MS-based spatial proteomics data. It posits that each annotated sub-cellular niche can be modelled by a multivariate Gaussian distribution.
- ▶ With the prior knowledge that many proteins are not captured by known sub-cellular niches, we augment our model with an **outlier component**. Outliers are often dispersed and thus this additional component is described by a heavy-tailed distribution: the multivariate Student's t-distribution, leading us to a *T Augmented Gaussian Mixture model*.
- ▶ This methodology allows proteome-wide **uncertainty quantification** (Shannon entropy), thus adding a further layer to the analysis of spatial proteomics.



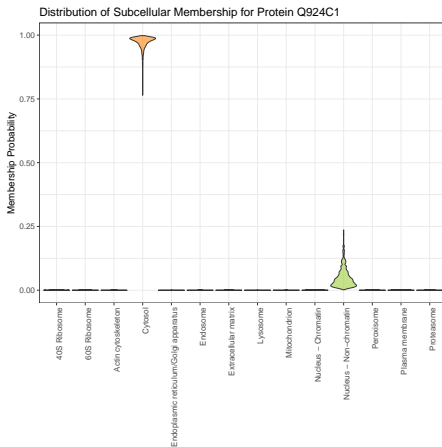
**Figure:** Illustration of how the TAGM model describes the pluripotent mouse embryonic stem cell data. Each ellipse contains a proportion of total probability of a particular multivariate Gaussian density. The outer ellipse contains 99% of the total probability whilst the middle and inner ellipses contain 95% and 90% of the probability respectively.



**Figure:** Assignment of proteins of *unknown* location to one of the annotated classes. The dots are scaled according to the protein assignment probabilities.

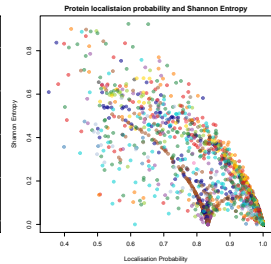
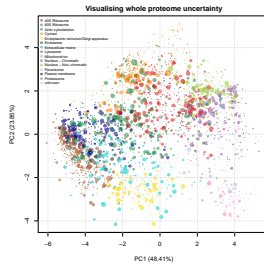
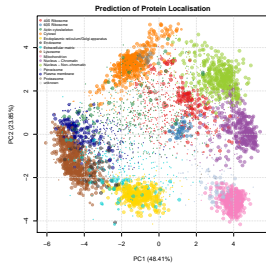






**Figure:** Exportin 5 (Q924C1) forms part of the micro-RNA export machinery, transporting miRNA from the nucleus to the cytoplasm for further processing. It then translocates back through the nuclear pore complex to return to the nucleus to mediate further transport between nucleus and cytoplasm. The model correctly infers that it most likely localises to the cytosol but there is some uncertainty with this assignment. This uncertainty is reflected in possible assignment of Exportin 5 to the nucleus non-chromatin and reflects the multi-location of the protein.

## Whole sub-cellular proteome uncertainty



# Spatial dynamics

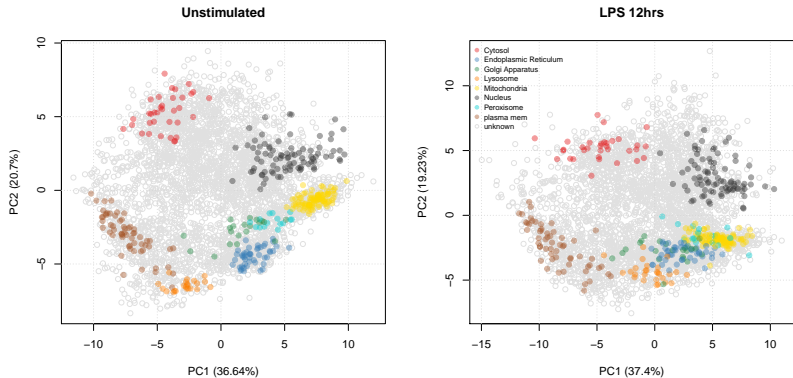
## Trans-localisation event during monocyte to macrophage differentiation

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

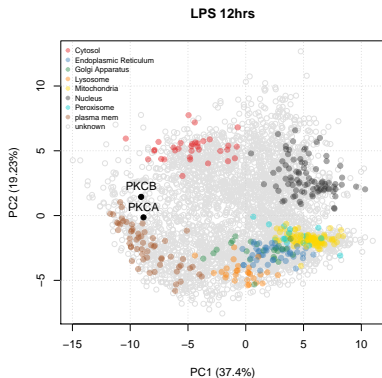
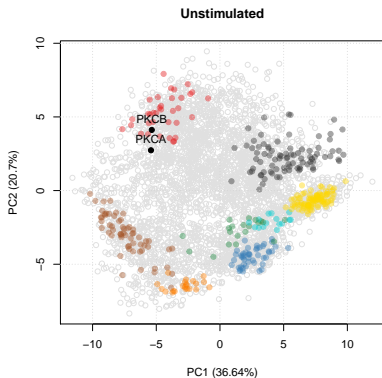
### Data

- ▶ 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.

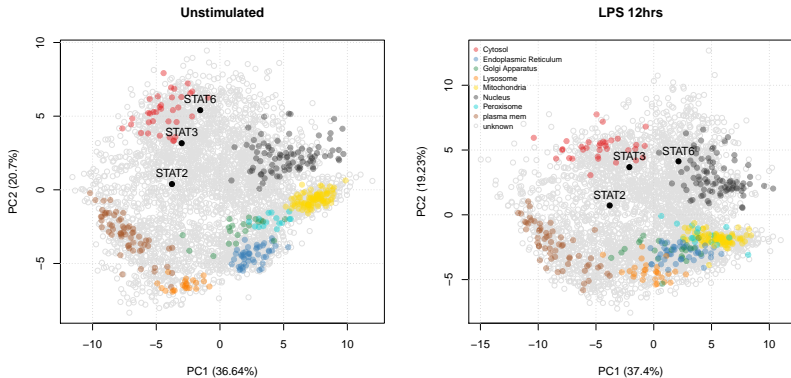
Work lead by **Dr Claire Mulvey** at the Cambridge Centre for Proteomics.



**Figure:** Spatial maps of unstimulated and LPS-treated cells (combined triplicates).



**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. \(2011\)](#).

**Behind the scenes:** software/data structures and open research practice.

## Beyond the figures<sup>1</sup>

- ▶ Software: **infrastructure** (MSnbase, Gatto and Lilley (2012)), **dedicated machine learning** (pRoloc, Gatto et al. (2014b)), **interactive visualisation**<sup>2</sup> (pRolocGUI, Breckels et al. (2017)) and **data** (pRolocdata, Gatto et al. (2014b)) for spatial proteomics.

---

<sup>1</sup>... which are all reproducible, by the way.

<sup>2</sup><https://lgatto.shinyapps.io/christoforou2015/>

<sup>3</sup>between and within domains/software



## Beyond the figures<sup>1</sup>

- ▶ Software: **infrastructure** (MSnbase, Gatto and Lilley (2012)), **dedicated machine learning** (pRoloc, Gatto et al. (2014b)), **interactive visualisation**<sup>2</sup> (pRolocGUI, Breckels et al. (2017)) and **data** (pRolocdata, Gatto et al. (2014b)) for spatial proteomics.
- ▶ The **Bioconductor** (Huber et al., 2015) ecosystem for high throughput biology data analysis and comprehension: **open source**, and **coordinated and collaborative**<sup>3</sup> **open development**, enabling **reproducible research**, enables understanding of the data (not a black box) and **drive scientific innovation**.

---

<sup>1</sup>... which are all reproducible, by the way.

<sup>2</sup><https://lgatto.shinyapps.io/christoforou2015/>

<sup>3</sup>between and within domains/software

# Open research: open source software

The screenshot shows the GitHub repository for **lgatto / pRoloC**. The repository is described as "A unifying bioinformatics framework for organelle proteomics" with a link to <http://lgatto.github.io/pRoloC/>. It has 2,051 commits, 10 branches, 25 releases, 1 environment, and 14 contributors. The repository is in the **master** branch. A table of files is shown below:

File	Description	Latest commit	Time
R	fix to make work with devel	6 day	
data	Merge branch 'master' into devel	2 year	
inst	updated documentation	6 month	
man	new logPosterior accessor	2 month	
src	more C exporting fuss	10 month	
tests	add test back	6 day	
vignettes	automatic fix of indentation	6 day	
_Rbuildignore	ignore .editorconfig when building	3 month	
.editorconfig	fix tab with spaces and add editorconfig	4 month	
.gitignore	ignore docs, update news	5 month	
.travis.yml	Fix indentation	5 month	
CONDUCT.md	Merge branch 'master' into devel	3 year	
DESCRIPTION	bump devel version on gh	6 day	
NAMESPACE	fix notes and warnings	7 day	
NEWS	bump devel version on gh	6 day	
NEWS.md	bump devel version on gh	6 day	

The screenshot shows the Bioconductor website for the **pRoloC** package. The Bioconductor logo is at the top left, and the navigation bar includes **Home**, **Install**, and **Help**. The breadcrumb trail is **Home > Bioconductor 3.8 > Software Packages > pRoloC**. The package name **pRoloC** is displayed in green. Below it, there are statistics: **platforms** (all), **Rank** (254 / 1649), **posts** (1 / 2 / 2 / 0), and **in Bioc** (6 years). There are also buttons for **build**, **warnings**, and **updated < 1 week**. The DOI is [10.18129/BIOC.P01603](https://doi.org/10.18129/BIOC.P01603). The description states: "A unifying bioinformatics framework for spatial proteomics". The Bioconductor version is Release (3.8). The package description is: "The pRoloC package implements machine learning and visualisation methods for the analysis and interpretation of quantitative mass spectrometry data to reliably infer protein sub-cellular localisation." The authors are Laurent Gatto, Oliver Crook and Lisa M. Breckels, with contributions from Thomas Burger and Samuel Wleczorek. The maintainer is Laurent Gatto <laurent.gatto@ulb.ac.be>. The citation is: "Gatto L, Breckels LM, Wleczorek S, Burger T, Lilley KS (2014). 'Mass-spectrometry based spatial proteomics data analysis using pRoloC and pRoloCdata.' Bioinformatics." The package is also cited in several other publications, including those by Breckels LM, Gatto L, Christoforou A, Groen AJ, Lilley KS, Trotter MW (2013), Gatto L, Breckels LM, Burger T, Nightingale DJ, Groen AJ, Campbell C, Mulvey CM, Christoforou A, Ferro M, Lilley KS (2014), Breckels LM, Holden S, Worger D, Mulvey CM, Christoforou A, Groen A, Trotter MW, Kohlbacker O, Lilley KS, Gatto L (2016), and Breckels LM, Mulvey CM, Lilley KS, Gatto L (2016).

Figure: Gatto et al. (2014a) Left: Public repository for the pRoloC software (<https://github.com/lgatto/pRoloC>). Right: official Bioconductor page.

# Open and reproducible research

The image displays three web pages related to the research paper "Assessing sub-cellular resolution in spatial proteomics experiments" by Gatto et al. (2018).

**Left Panel (GitHub):** Shows the repository `lgatto / QSep-manuscript`. It lists files and folders including `data`, `figure`, `github`, `travis.yml`, `MultiFile`, `README.md`, `cover.pdf`, `cover.tex`, `cover2.pdf`, `etichetta.xls`, `HiRes.xls`, `risk.R`, `reknetch.ppt.pdf`, `reknetch-qap.pdf`, `reknetch-qap.R`, `qsep.R`, `qsep.Rnw`, `qsep.bb`, `qsep.pdf`, `qsep.tex`, `sims.pdf`, and `sims.R`. Each file has a brief description of its content.

**Middle Panel (bioRxiv):** Displays the preprint version of the paper. The title is "Assessing sub-cellular resolution in spatial proteomics experiments". The authors are Laurent Gatto, Lisa M. Breckels, and Kathryn S. Lilley. The paper is marked as "New Results". The bioRxiv logo and navigation links are visible at the top.

**Right Panel (ScienceDirect):** Shows the published version of the paper in the journal "Current Opinion in Chemical Biology". The title is "Assessing sub-cellular resolution in spatial proteomics experiments". The authors are Laurent Gatto<sup>1,2,3,4</sup>, Lisa M. Breckels<sup>1,2</sup>, and Kathryn S. Lilley<sup>2</sup>. The paper is marked as "New Results". The ScienceDirect logo and navigation links are visible at the top.

Figure: Gatto et al. (2018) reproducible document  
(<https://github.com/lgatto/QSep-manuscript>), preprint  
(<https://doi.org/10.1101/377630>) and paper  
(<https://doi.org/10.1016/j.cbpa.2018.11.015>).

Working with open and reproducible research in mind doesn't mean releasing everything prematurely, it means

- ▶ managing research in a way one can find data and results at every stage
- ▶ one can reproduce results, re-run/compare them with new data or different methods/parameters, and
- ▶ one can release data (or parts thereof) when/if appropriate.

# Conclusions

- ▶ Protein sub-cellular localisation: technologies (hyperLOPIT) and opportunities.
- ▶ Reliance on computational biology and dedicated software (pRoLoc *et al.*) to interpret data and acquire biological knowledge.
- ▶ Rigorous computational infrastructure and sound data analysis and interpretation is a **long term investment**.

# References I

- 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 2016. doi: 10.1371/journal.pcbi.1004920.
- Lisa Breckels, Thomas Naake, and Laurent Gatto. *pRolocGUI: Interactive visualisation of spatial proteomics data*, 2017. URL <http://ComputationalProteomicsUnit.github.io/pRolocGUI/>. R package version 1.11.2.
- 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.
- 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 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 Wiczorek, 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.

# References II

- Laurent Gatto, Lisa M Breckels, and Kathryn S Lilley. Assessing sub-cellular resolution in spatial proteomics experiments. *bioRxiv*, 2018. doi: 10.1101/377630.
- W Huber, V J Carey, R Gentleman, S Anders, M Carlson, B S Carvalho, H C Bravo, S Davis, L Gatto, T Girke, R Gottardo, F Hahne, K D Hansen, R A Irizarry, M Lawrence, M I Love, J MacDonald, V Obenchain, A K Oleś, H Pagès, A Reyes, P Shannon, G K Smyth, D Tenenbaum, L Waldron, and M Morgan. Orchestrating high-throughput genomic analysis with Bioconductor. *Nat Methods*, 12(2):115–21, Jan 2015. doi: 10.1038/nmeth.3252.
- 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.
- J E Siljee, Y Wang, A A Bernard, B A Ersoy, S Zhang, A Marley, M Von Zastrow, J F Reiter, and C Vaisse. Subcellular localization of mc4r with adcy3 at neuronal primary cilia underlies a common pathway for genetic predisposition to obesity. *Nat Genet*, Jan 2018. doi: 10.1038/s41588-017-0020-9.
- C Stadler, E Rexhepaj, V R Singan, R F Murphy, R Pepperkok, M Uhlén, J C Simpson, and E Lundberg. Immunofluorescence and fluorescent-protein tagging show high correlation for protein localization in mammalian cells. *Nat Methods*, 10(4):315–23, Apr 2013.
- 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.

## Acknowledgements

- ▶ **Mr Oliver Crook and Dr Lisa Breckels**, (U of Cambridge): spatial proteomics, machine learning, software.
- ▶ **Dr Sebastian Gibb and Dr Johannes Rainer**: MS and proteomics software.
- ▶ Prof Kathryn Lilley (U of Cambridge), Dr Claire Mulvey, (CRUK Cambridge Institute): data.
- ▶ Funding: BBSRC, Wellcome Trust

**Thank you for your attention**