

A Bayesian Mixture Modelling Approach For Spatial Proteomics

Oliver M. Crook^{1,2,3}, Claire M. Mulvey², Paul D.W. Kirk³, Kathryn S. Lilley², and Laurent Gatto^{1,2,4,*}

¹ *Computational Proteomics Unit, Department of Biochemistry, University of Cambridge, Cambridge, UK*

² *Cambridge Centre for Proteomics, Department of Biochemistry, University of Cambridge, Cambridge, UK*

³ *MRC Biostatistics Unit, Cambridge Institute for Public Health, Cambridge, UK*

⁴ *Current address: de Duve Institute, UCLouvain, Avenue Hippocrate 75, 1200 Brussels, Belgium*

**laurent.gatto@uclouvain.be*

October 5, 2018

Abstract

Analysis of the spatial sub-cellular distribution of proteins is of vital importance to fully understand context specific protein function. Some proteins can be found with a single location within a cell, but up to half of proteins may reside in multiple locations, can dynamically re-localise, or reside within an unknown functional compartment. These considerations lead to uncertainty in associating a protein to a single location. Currently, mass spectrometry (MS) based spatial proteomics relies on supervised machine learning algorithms to assign proteins to sub-cellular locations based on common gradient profiles. However, such methods fail to quantify uncertainty associated with sub-cellular class assignment. Here we reformulate the framework on which we perform statistical analysis. We propose a Bayesian generative classifier based on Gaussian mixture models to assign proteins probabilistically to sub-cellular niches, thus proteins have a probability distribution over sub-cellular locations, with Bayesian computation performed using the expectation-maximisation (EM) algorithm, as well as Markov-chain Monte-Carlo (MCMC). Our methodology allows proteome-wide uncertainty quantification, thus adding a further layer to the analysis of spatial proteomics. Our framework is flexible, allowing many different systems to be analysed and reveals new modelling opportunities for spatial proteomics. We find our methods perform competitively with current state-of-the art machine learning methods, whilst simultaneously providing more information. We highlight several examples where classification based on the support vector machine is unable to make any conclusions, while uncertainty quantification using our approach provides biologically intriguing results. To our knowledge this is the first Bayesian model of MS-based spatial proteomics data.

Author summary

Sub-cellular localisation of proteins provides insights into sub-cellular biological processes. For a protein to carry out its intended function it must be localised to the correct sub-cellular environment, whether that be organelles, vesicles or any sub-cellular niche. Correct

sub-cellular localisation ensures the biochemical conditions for the protein to carry out its molecular function are met, as well as being near its intended interaction partners. Therefore, mis-localisation of proteins alters cell biochemistry and can disrupt, for example, signalling pathways or inhibit the trafficking of material around the cell. The sub-cellular distribution of proteins is complicated by proteins that can reside in multiple micro-environments, or those that move dynamically within the cell. Methods that predict protein sub-cellular localisation often fail to quantify the uncertainty that arises from the complex and dynamic nature of the sub-cellular environment. Here we present a Bayesian methodology to analyse protein sub-cellular localisation. We explicitly model our data and use Bayesian inference to quantify uncertainty in our predictions. We find our method is competitive with state-of-the-art machine learning methods and additionally provides uncertainty quantification. We show that, with this additional information, we can make deeper insights into the fundamental biochemistry of the cell.

Introduction

Spatial proteomics is an interdisciplinary field studying the localisation of proteins on a large-scale. Where a protein is localised in a cell is a fundamental question, since a protein must be localised to its required sub-cellular compartment to interact with its binding partners (for example, proteins, nucleic acids, metabolic substrates) and carry out its function [1]. Furthermore, mis-localisations of proteins are also critical to our understanding of biology, as aberrant protein localisation have been implicated in many pathologies [2, 3, 4, 5, 6], including cancer [7, 8, 9, 10] and obesity [11].

Sub-cellular localisations of proteins can be studied by high-throughput mass spectrometry (MS) [12]. MS-based spatial proteomics experiments enable us to confidently determine the sub-cellular localisation of thousands of proteins within a cell [13], given the availability of rigorous data analysis and interpretation [12].

In a typical MS-based spatial proteomics experiment, cells first undergo lysis in a fashion which maintains the integrity of their organelles. The cell content is then separated using a variety of methods, such as density separation [14, 13], differential centrifugation [15], free-flow electrophoresis [16], or affinity purification [17]. In LOPIT [18, 14, 19] and *hyper*LOPIT [13, 20], cell lysis is preceded by separation of the content along a density gradient. Organelles and macro-molecular complexes are thus characterised by density-specific profiles along the gradient [21]. Discrete fractions along the continuous density gradient are then collected, and quantitative protein profiles that match the organelle profiles along the gradient, are measured using high accuracy mass spectrometry [20].

The data are first visualised using principal component analysis (PCA) and known sub-cellular compartments are annotated [22]. Supervised machine learning algorithms are then typically employed to create classifiers that associate un-annotated proteins to specific organelles [23], as well as semi-supervised methods that detect novel sub-cellular clusters using both labelled and un-labelled features [24]. More recently, a state-of-the-art transfer learning (TL) algorithm has been shown to improve the quantity and reliability of sub-cellular protein assignments [25]. Applications of such methods have led to organelle-specific localisation information of proteins in plants [14], *Drosophila* [26], chicken [27], human cell lines

[24], mouse pluripotent embryonic stem cells [13] and cancer cell lines [28].

Classification methods which have previously been used include partial least squares discriminant analysis [14], K nearest neighbours [29], random forests [30], naive Bayes [31], neural networks [32] and the support vector machine amongst others (see [23] for an overview). Although these methods have proved successful within the field they have limitations. Typically, such classifiers output an assignment of proteins to discrete pre-annotated sub-cellular locations. However, it is important to note that half the proteome cannot be robustly assigned to a single sub-cellular location, which may be a manifestation of proteins in so far uncharacterised organelles or proteins that are distributed amongst multiple locations. These factors lead to uncertainty in the assignment of proteins to sub-cellular localisations, and thus quantifying this uncertainty is of vital importance [33].

To overcome the task of uncertainty quantification, this article presents a probabilistic generative model for MS-based spatial proteomics data. Our model posits that each annotated sub-cellular niche can be modelled by a multivariate Gaussian distribution. Thus, the full complement of annotated proteins is captured by a mixture of multivariate Gaussian distributions. 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 (TAGM).

Given our model and proteins with known location, we can probabilistically infer the sub-cellular localisation of thousands of proteins. We can perform inference in our model by finding *maximum a posteriori* (MAP) estimates of the parameters. This approach returns the probability of each protein belonging to each annotated sub-cellular niche. These posterior localisation probabilities can then be the basis for classification. In a more sophisticated, fully Bayesian approach to uncertainty quantification, we can additionally infer the entire posterior distribution of localisation probabilities. This allows the uncertainty in the parameters in our model to be reflected in the posterior localisation probabilities. We perform this inference using Markov-chain Monte-Carlo methods; in particular, we provide an efficient collapsed Gibbs sampler to perform inference.

We perform a comprehensive comparison to state-of-the-art classifiers to demonstrate that our method is reliable across 19 different spatial proteomics datasets and find that all classifiers we considered perform competitively. To demonstrate the additional biological advantages our method can provide, we apply our method to a *hyper*LOPIT dataset on mouse pluripotent embryonic stem cells [13]. We consider several examples of proteins that were unable to be assigned using traditional machine-learning classifiers and show that, by considering the full posterior distribution of localisation probabilities, we can draw meaningful biological results and make powerful conclusions. We then turn our hand to a more global perspective, visualising uncertainty quantification for over 5,000 proteins, simultaneously. This approach reveals global patterns of protein organisation and their distribution across sub-cellular compartments.

We make extensive use of the R programming language [34] and existing MS and proteomics packages [35, 36]. We are highly committed to creating open software tools for high quality processing, visualisation, and analysis of spatial proteomics data. We build upon an already extensive set of open software tools [36] as part of the Bioconductor project [37, 38]

and our methods are made available as part of this project.

Results

Application to mouse pluripotent embryonic stem cell data

We model mouse pluripotent embryonic stem cell (E14TG2a) data [13], which contains quantitation data for 5032 proteins. This high-resolution map was produced using the *hy-per*LOPIT workflow [20], which uses a sophisticated sub-cellular fractionation scheme. This fractionation scheme is made possible by the use of Tandem Mass Tag (TMT) 10-plex and high accuracy TMT quantification was facilitated by using synchronous precursor selection MS3 (SPS-MS3) [39], which reduces well documented issues with ratio distortion in isobaric multiplexed quantitative proteomics [40]. The data resolves 14 sub-cellular niches with an additional chromatin preparation resolving the nuclear chromatin and non-chromatin components. Two biological replicates of the data are concatenated, each with 10 fractions along the density gradient. We defined gold standard organelle markers as those with unambiguous single annotation [23]. A protein marker list for the mouse pluripotent embryonic stem cells was manually curated using information from the UniProt database, the Gene Ontology and the literature, as was performed in [13]. The following section applies our statistical methodology to these data and we explore the results.

Maximum a posteriori prediction of protein localisation

This section applies the TAGM model to the mouse pluripotent embryonic stem cell data, by deriving MAP estimates for the model parameters and using these for prediction. Visualisation is important for data analysis and exploration. A simple way to visualise our model is to project probability ellipses onto a PCA plot. 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. Visualising only the first two principal components can be misleading, since proteins can be more (or less) separated in subsequent principal components. We visualise the first two principal components along with the first and fourth principal components as a representative example. For the TAGM model, we derive probability ellipses from the MAP estimates of the parameters.

We now apply the statistical methodology described in section **Model**, to predict the localisation of proteins to organelles and sub-cellular components. In brief, we produce MAP estimates of the parameters by using the expectation-maximisation algorithm, to form the basis of a Bayesian analysis (TAGM-MAP). We run the algorithm for 200 iterations and inspect a plot of the log-posterior to assess convergence of the algorithm (see supporting information, section S3). We confirm that the difference of the log posterior between the final two iterations is less than 10^{-6} and we conclude that our algorithm has converged. The results can be seen in figure 2 (left), where the posterior localisation probability is visualised by scaling the pointer for each protein.

Figure 2 (right) demonstrates a range of probabilistic assignments of proteins to organelles and sub-cellular niches. We additionally consider a full, sampling-based Bayesian analysis

using Markov-chain Monte Carlo (MCMC) to characterise the uncertainty in the localisation probabilities. In our case a collapsed Gibbs sampler is used to sample from the posterior of localisation probabilities. The remainder of this article focus on analysis of spatial proteomics in this fully Bayesian framework.

Uncertainty in the posterior localisation probabilities

This section applies the TAGM model to the mouse pluripotent embryonic stem cell data, by considering the uncertainty in the parameters and exploring how this uncertainty propagates to the uncertainty in protein localisation prediction. In figure 3 we visualise the model as before using the first two principal components along with the first and fourth principal component as a representative example. For the TAGM model, we derive probability ellipses from the expected value of the posterior normal-inverse-Wishart (NIW) distribution.

We apply the statistical methodology detailed in section Model. We perform posterior computation in the Bayesian setting using standard MCMC methods (TAGM-MCMC). We run 6 chains of our Gibbs sampler in parallel for 15,000 iterations, throwing away the first 4,000 iterations for burn-in and retain every 10th sample for thinning. Thus 1,100 sample are retained from each chain. We then visualise the trace plots of our chains; in particular, we monitor the number of proteins allocated to the known components (see supporting information, section S4). We discard 1 chain because we do not consider it to have converged. For the remaining 5 chains we further discard the first 500 samples by visual inspection. We then have 600 retained samples from 5 separate chains. For further analysis, we compute the Gelman-Rubin convergence diagnostic [41, 42], which is computed as $\hat{R} \approx 1.05$. Values of \hat{R} far from 1 indicate non-convergence and since our statistic is less than 1.1, we conclude our chains have converged. The remaining samples are then pooled to produce a single chain containing 3000 samples.

We produce point estimates of the posterior localisation probabilities by summarising samples by their Monte-Carlo average. These summaries are then visualised in figure 2 (right panel), where the pointer is scaled according to the localisation probabilities of the sub-cellular niche with the largest posterior probability. Monte-Carlo based inference also provides us with additional information; in particular, we can interrogate individual proteins and their posterior probability distribution over sub-cellular locations.

Figure 4 illustrates one example of the importance of capturing uncertainty. The E3 ubiquitin-protein ligase TRIP12 (G5E870) is an integral part of ubiquitin fusion degradation pathway and is a protein of great interest in cancer because it regulates DNA repair pathways. The SVM failed to assign this protein to any location, with assignment to the 60S Ribosome falling below a 5% FDR and the MAP estimate assigned the protein to the nucleus non-chromatin with posterior probability < 0.95 . The posterior distribution of localisation probabilities inferred from the TAGM-MCMC model, shown in figure 4, demonstrates that this protein is most probably localised to the nucleus non-chromatin. However, there is some uncertainty about whether it localises to the 40S ribosome. This could suggest a dynamic role for this protein, which could be further explored with a more targeted experiment.

Enrichment analysis of outlier proteins

In previous sections, we demonstrated that we can assign proteins probabilistically to sub-cellular compartment and quantify the uncertainty in these assignments. Some proteins cannot be well described as belonging to any known component and we model this using an additional T-distribution outlier component (see Section [Model](#)).

It is biologically interesting to decipher what functional role proteins that are far away from known components play. We perform an over-representation analysis of gene ontology (GO) terms to assess the biological relevance of the outlier component [43, 44]. We take 1111 proteins that were allocated to known components with probability less than 0.95. Note that these 1111 proteins exclude proteins that are likely to belong to a known location, but we are uncertain about which localisation. We then perform enrichment analysis against the set of all proteins quantified in the *hyper*LOPIT experiment. We search against the cellular compartment, biological process and molecular function ontologies.

Supplementary figure 3 shows this outlier component is enriched for cytoskeletal part ($p < 10^{-7}$) and microtubule cytoskeleton ($p < 10^{-7}$). Cytoskeleton proteins are found throughout the cell and therefore we would expect them to be found in every fraction along the density gradient. We also observe enrichment for highly dynamic sub-cellular processes such as cell division ($p < 10^{-6}$) and cell cycle processes ($p < 10^{-6}$), again these proteins are unlikely to have steady-state locations within a single component. We also see enrichment for molecular functions such as transferase activity ($p < 0.005$), another highly dynamic process. These observations justify including an additional outlier component in our mixture model.

Comparison with other classifiers

In this section, we assess the generalisation performance of our methods on several datasets, by computing performance metrics associated with each classifier as detailed in section [Classifier assessment](#). We compare the SVM and KNN classifiers alongside the MAP and MCMC approaches detailed in the methods section. We compute the F1 score and quadratic loss over 100 rounds of stratified 5-fold cross-validation. The hyperparameter for the KNN algorithm, the number of nearest neighbours, is optimised via an additional internal 5-fold cross-validation and the hyperparameters for the SVM, sigma and cost, are also optimised via internal 5-fold cross validation [45].

We test our methods on the following datasets *Drosophila* [26], chicken [27], mouse pluripotent embryonic stem cells from [13] and [25], the human bone osteosarcoma epithelial (U2-OS) cell line [28], the HeLa cell line of [15], the 3 HeLa cell lines from [46] and 10 primary fibroblast datasets from [47]. These datasets represent a great variety of spatial proteomics experiments across many different workflows.

The two *hyper*LOPIT datasets on mouse pluripotent embryonic stem cells and the U2-OS cell line use TMT 10-plex labelling and contain the greatest number of proteins. Earlier LOPIT experiments on the *Drosophila* and chicken use iTRAQ 4-plex labelling, whilst another LOPIT mouse pluripotent embryonic stem cell dataset uses iTRAQ 8-plex. The datasets of [15] and [46] employ a different methodology completely - separating cellular content using differential centrifugation (as opposed to along a density-gradient). Furthermore, the methods use SILAC rather than iTRAQ or TMT for labelling. The experiments of [46] were

designed to explore the functional role of AP-5 by coupling CRISPR-CAS9 knockouts with spatial proteomics methods. We analysed all three datasets from [46], which includes a wild type HeLa cell line as a control, as well as two CRISPR-CAS9 knockouts: AP5Z1-KO1 and AP5Z1-KO2 respectively.

In addition, we analyse the spatio-temporal proteomics experiments of [47], which uses TMT-based MS quantification. This experiment explored infecting primary fibroblasts with Human cytomegalovirus (HCMV) and the goal of these experiments was to explore the dynamic perturbation of host proteins during infection, as well as the sub-cellular localisation of viral proteins through the HCMV life-cycle. They produced spatial maps at different time points: 24, 48, 72, 96, 120 hours post infection (hpi), as well as mock maps at these same time points to serve as a control - this results in 10 different spatial proteomics maps.

In each case, a dataset specific marker list was used, which is curated specifically for the each cell line. We removed "high-curvature ER" annotations from the HeLa dataset [15], as well as the "ER Tubular", "Nuclear pore complex" and "Peroxisome" annotations from the HeLa CRISPR-CAS9 knockout experiments [46] as there are too few proteins to correctly perform cross-validation. Table 1 summarises these datasets, including information about number of quantified proteins, the workflow used and the number of fractions.

MS-based Spatial Proteomics datasets				
Cell line or organism	Workflow	Labelling	Fractions (including combined replicates)	Proteins
<i>Drosophila</i>	LOPIT	iTRAQ	4	888
Chicken DT40	LOPIT	iTRAQ	16	1090
Mouse pluripotent E14TG2a stem cell	HyperLOPIT	TMT	20	5032
HeLa (Itzhak et al.)	Organeller Maps	SILAC	30	3766
HeLa (Hirst et al.)	Organeller Maps	SILAC	15	2046
U2-OS cell line	HyperLOPIT	TMT	37	5020
Primary Fibroblast	Spatio-Temporal Methods	TMT	6	2196
E14TG2a (Breckels et al.)	LOPIT	iTRAQ	8	2031

Table 1: Summary of spatial proteomics datasets used for comparisons

Figure 5 compares the Macro-F1 scores across the datasets for all classifiers and demonstrates that no single classifier consistently outperforms any other across all datasets, with results being highly consistent across all methods, as well as across datasets. We perform

a pairwise unpaired t-test with multiple testing correction applied using the Benjamini-Höcherberg procedure [48] to detect differences between classifier performance.

In the *Drosophila* dataset only the KNN algorithm outperforms the SVM at significance level of 0.01, whilst no other significant differences exist between the classifiers. In the chicken DT40 dataset only the MCMC method outperforms the KNN classifier at significance level of 0.01, no other significant conclusion can be drawn. In the mouse dataset the MAP based method outperforms the MCMC method at significance level of 0.01, no other significant conclusions can be drawn. In the HeLa dataset all classifiers are significantly different at a 0.01 level. These differences may exist because the dataset does not fit well with our modelling assumptions; in particular, this dataset set has been curated to have a class called "Large Protein Complex", which likely describes several sub-cellular structures. These might include nuclear compartments and ribosomes, as well as any cytosolic complex and large protein complex which pellets during the centrifugation conditions used to capture this mixed sub-cellular fraction. Moreover, the cytosolic and nuclear fraction were processed separately leading to possible imbalance with comparisons with other datasets. Thus, the large protein complex component might be better described as itself a mixture model or more detailed curation of these data may be required. We do not consider further modelling of this dataset in this manuscript. For the U2-OS all classifiers are significantly different at a significance level of 0.01 except for the SVM classifier and the MCMC method, with the MAP method performing the best. Figure 5 shows that for this dataset all classifiers are performing extremely well. In the three Hirst datasets the MAP method significantly outperforms all other methods ($p < 0.01$), whilst in the wild type HeLa and in the CRISPR-CAS9 KO1 there is no significant difference between the KNN and MCMC method. In the CRISPR-CAS9 KO2 the MCMC method outperforms the SVM and KNN methods ($p < 0.01$). In the interest of brevity, the remaining results for the t-tests can be found in tables in supporting information, section S5.

The Macro-F1 scores do not take into account that whilst the TAGM model may misclassify, it may do so with low confidence. We therefore additionally compute the quadratic loss, which allows us to make use of the probabilistic information provided by the classifiers. The lower the quadratic loss the closer the probabilistic prediction is to the true value. We plot the distributions of quadratic losses for each classifier in figure 6. We observe highly consistent performance across all classifiers across all datasets. Again, we perform a pairwise unpaired t-test with multiple testing correction.

We find that in 16 out of 19 datasets (all of those except HeLa Wild type, HeLa KO1 and HeLa KO2) the MCMC methods achieves the lowest quadratic loss at a significance level < 0.0001 over the SVM and KNN classifiers. In 6 out of these 16 datasets there is no significant difference between the MCMC and the MAP methods. In the three Hirst datasets in which the MCMC did not achieve the lowest quadratic loss, the SVM outperformed. However, in two of these datasets (HeLa Wild type and KO1) the MAP method and SVM classifier were not significantly different. In the Hirst KO2 dataset there were no significant differences between the MAP and MCMC methods.

In the vast majority of cases, we observe that if the TAGM model, using the MCMC methodology, makes an incorrect classification it does so with lower confidence than the SVM classifier, the KNN classifier and the MAP based classifier, whilst if it is correct in its assertion it does so with greater confidence. Additionally, a fully Bayesian methodology provides us

with not only point estimates of classification probabilities but uncertainty quantification in these allocations, and we show in the following section that this provides deeper insights into protein localisation.

Computing distributions of F1 scores and quadratic losses, which can only be done on the marker proteins, can help us understand whether a classifier might have greater generalised performance accuracy. However, we are interested in whether there is a large disagreement between classifiers when prediction is performed on proteins for which we have no withheld localisation information. This informs us about a systematic bias for a particular classifier or whether a classifier ensemble could increase performance. To maintain a common set of proteins we set thresholds for each classifier in turn and compare to the other classifier without thresholding. Firstly, we set a global threshold of 0.95 for the TAGM-MCMC and then for these proteins plot a contingency table against the classification results from the SVM. Secondly, we set a 5% FDR for the SVM and then for these proteins plot a contingency table against the classification results from the TAGM-MCMC. We visualise the contingency tables as heat plots in figure 7.

In general, we see an extremely high level of coherence between the TAGM and the SVM, with almost all proteins predicted to concordant sub-cellular compartments. Figure 7 shows there is some disagreement between assigning proteins to the lysosome and plasma membrane, to the cytosol and proteasome, and between the large and small ribosomal subunits. However, we have not used the uncertainty in the probabilistic assignments to produce the contingency tables above. In the next sections, we explore examples of proteins with uncertainty in their posterior localisation probabilities. Selecting biologically relevant thresholds is important for any classifier and exploring uncertainty is of vital importance when drawing biological conclusions.

Interpreting and exploring uncertainty

Protein sub-cellular localisation can be uncertain for a number of reasons. Technical variations and unknown biological novelty, such as yet uncharacterised functional compartments, can be some of the reasons why a protein might have an unknown or uncertain localisation. Furthermore many proteins are known to reside in multiple locations with possibly different functional duties in each location [49]. With these considerations in mind, it is pertinent to quantify the uncertainty in our allocation of proteins to organelles. This section explores several situations where proteins display uncertain localisation and considers the biological factors that influence uncertainty. We later explore and visualise whole proteome uncertainty quantification.

Exportin 5 (Q924C1) forms part of the micro-RNA export machinery of the nucleus, 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. Exportin 5 can then continue to mediate further transport between nucleus and cytoplasm. The SVM was unable to assign a localisation of Exportin 5, with its assignment falling below a 5% FDR to wrongly assign this protein to the proteasome. This incorrect assertion by the SVM was confounded by the similarity between the cytosol and proteasome profiles. Figure 8 demonstrates, according to the TAGM-MCMC model, that Exportin 5 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 this uncertainty is a manifestation of the fact that the function of this protein is to shuttle between the cytosol and nucleus.

The Phenylalanine-tRNA ligase beta subunit protein (Q9WUA2) has an uncertain localisation between the 40S ribosome and the nucleus non-chromatin demonstrated in figure 9. This protein was left unclassified by the SVM because its score fell below a 5% FDR threshold to assign it to the 40S ribosome. Considering that this protein is involved in the acylation of transfer RNA (tRNA) with the amino acid phenylalanine to form tRNA-Phe to be used in translation of proteins, it is therefore unsurprising that this protein’s steady state location is ribosomal. Whilst the SVM is unable to make an assignment, TAGM-MCMC is able to suggest an assignment and quantify our uncertainty.

Relatively little is known about the Dedicator of cytokinesis (DOCK) protein 6 (Q8VDR9), a guanine nucleotide exchange factor for CDC42 and RAC1 small GTPases. The SVM could not assign localisation to the ER/Golgi, since its score fell below a 5% FDR. Furthermore, the TAGM-MCMC model assigned this DOCK 6 to the outlier component with posterior probability > 0.95 . Figure 10 shows possible localisation to several components along the secretory pathway. As an activator for CDC42 and RAC1 we may expect to see them with similar localisation. CDC42, a plasma membrane associated protein, regulates cell cycle and division and is found with many localisations. Furthermore RAC1, a small GTPase, also regulates many cellular processes and is found in many locations. Thus the steady-state distribution of DOCK6 is unlikely to be in a single location, since its interaction partners are found in many locations. This justifies including an outlier component in our model, else we may erroneously assign such proteins to a single location.

Visualising whole sub-cellular proteome uncertainty

The advantage of the TAGM-MCMC model is its ability to provide proteome wide uncertainty quantification. Regions where organelle assignments overlap are areas where uncertainty is expected to be the greatest, as well as areas with no dominant component. We take an information theoretic approach to summarising uncertainty in protein localisation by computing the Shannon entropy [50] for each Monte-Carlo sample $t = 1, \dots, T$ of the posterior localisation probabilities of each protein

$$\left\{ H^{(t)} = - \sum_{k=1}^K p_{ik}^{(t)} \log \left(p_{ik}^{(t)} \right) \right\}_{t=1}^T, \quad (1)$$

where $p_{ik}^{(t)}$ denotes the posterior localisation probability of protein i to component k at iteration t . We then summarise this as a Monte-Carlo averaged Shannon entropy. The greater the Shannon entropy the more uncertainty associated with the assignment of this protein. The lower the Shannon entropy the lower the uncertainty associated with the assignment of this protein. In figure 11 panel (a), we visualise the Shannon entropy of each protein in a PCA plot, by scaling the pointers in accordance to this metric. We also note that while localisation probability (of a protein to its most probable location) and the Shannon entropy

are correlated, figure 11 panel (c), it is not perfect. Thus it is important to use both the localisation probabilities and the uncertainty in these assignments to make conclusions.

Figure 11 demonstrates that the regions of highest uncertainty are those in regions where organelles assignments overlap. The conclusions from this plot are manifold. Firstly, many proteins are assigned unambiguously to sub-cellular localisations; that is, not only are some proteins assigned to organelles with high probability but also with low uncertainty. Secondly, there are well defined regions with high uncertainty, for example proteins in the secretory pathway or proteins on the boundary between cytosol and proteasome. Finally, some organelles, such as the mitochondria, are extremely well resolved. This observed uncertainty in the secretory pathway and cytosol could be attributed to the dynamic nature of these parts of the cell with numerous examples of proteins that traffic in and out of these sub-cellular compartments as part of their biological role. Moreover, the organelles of the secretory pathway share similar and overlapping physical properties making their separation from one another using biochemical fractionation more challenging. Furthermore, there is a region located in the centre of the plot where proteins simultaneously have low probability of belonging to any organelle and high uncertainty in their localisation probability. This suggests that these proteins are poorly described by any single location. These proteins could belong to multiple locations or belong to undescribed sub-cellular compartments. The information displayed in these plots and the conclusion therein would be extremely challenging to obtain without the use of Bayesian methodology.

Discussion

We have demonstrated that a Bayesian framework, based on Gaussian mixture models, for spatial proteomics can provide whole sub-cellular proteome uncertainty quantification on the assignment of proteins to organelles and such information is invaluable. Performing MAP inference using our generative model provides fast and straightforward approach, which is vital for quality control and early data exploration. Full posterior inference using MCMC provides not only point estimates of the posterior probability that a protein belongs to a particular sub-cellular niche, but uncertainty in this assignment. Then, this uncertainty can be summarised in several ways, including, but not limited to, equi-tailed credible intervals of the Monte-Carlo samples of posterior localisation probabilities. Posterior distributions for individual proteins can then be rigorously interrogated to shed light on their biological mechanisms; such as, transport, signalling and interactions.

As well as the local uncertainty seen by exploring individual proteins, we further explored using a Monte-Carlo averaged Shannon entropy to visualise global uncertainty. Regions of high uncertainty, as measured using this Shannon entropy, reflect highly dynamics regions of the sub-cellular environment. Hence, biologists can now explore uncertainty at different levels and then are able to make quantifiable conclusions and insights about their data. Furthermore, our Bayesian model is interpretable and our inferences are fully conditional on our data, allowing them to be easily modified with changing experimental design.

In addition, we produced competitive classifier performance to the state-of-the-art classifiers. We considered two traditional machine-learning methods: the SVM and KNN classifiers; as well as two classifiers based on our model: a MAP classifier and classification based

on MCMC. We compared all methods on 19 different spatial proteomics datasets, across four different organisms. When considering the macro-F1 score as a performance metric, no single classifier outperformed another across all datasets. However, using MCMC based inference our method significantly outperforms the SVM and KNN classifiers with respect to the quadratic loss in 16 out of 19 datasets. This allows us to have greater confidence in our conclusions when they are drawn from our Bayesian inferences. Furthermore, using MCMC provides a wealth of additional information, and so becomes the method of choice for analysing spatial proteomics data.

Analysis of a *hyper*LOPIT experiment applied to mouse pluripotent embryonic stem cells demonstrated that the additional layer of information that our model provides is biologically relevant and provides further avenues for additional exploration. Moreover, applying our method to a biologically significant dataset now provides the scientific community with localisation information on up to 4000 proteins for the mouse pluripotent stem cell proteome. Figure 12 demonstrates that from an initial input of roughly 1000 marker proteins with *a priori* known location and 4000 unknown proteins with unknown location, SVM and TAGM-MCMC can provide rigorous localisation information on roughly 2000 proteins. However, our methodology, by also considering uncertainty, allows us to obtain information on another 1000 proteins. Thus, we have augmented this dataset by providing uncertainty quantification on the localisation of proteins to their sub-cellular niches, which had been previously unavailable. We note that our method is general enough to be applied to many MS-based spatial proteomics protocols including: LOPIT, *hyper*LOPIT, protein correlation profiling (PCP) [69], differential centrifugation approaches and spatio-temporal proteomics methods. In our flexible software implementation, all hyperparameters for the priors can be changed if users have precise priors they wish to specify.

We have also provided a new set of visualisation methods to accompany our model, which allow us to easily interrogate our data. High quality visualisation tools are essential for rigorous quality control and sound biological conclusions. Our methods have been developed in the R statistical programming language and we continue to contribute to the Bioconductor project [37, 38] with inclusion of our methods within the pRoloc package ($\geq 1.21.1$) [36]. The underlying source code used to generate this document is available at <https://github.com/lgatto/2018-TAGM-paper>.

Currently, our model does not integrate localisation information from different data sources, nor does it explicitly model proteins with multiple localisation. However, one (of many) biological explanations for the uncertainty that we model in the allocation probabilities is provided by multiple localisation. Thus a protein for which it is uncertain to which two sub-cellular niches it is resident within it is perhaps resident of both niches. In further work, we plan to explicitly model such cases to deconvolute different sources of uncertainty. In addition, extensions to semi-supervised non-parametric methods are under consideration to detect novel sub-cellular niches. These are the subjects of further work.

Model

We describe in this section the probabilistic model that uses the labelled data to associate un-annotated proteins to specific organelles or sub-cellular compartments.

Mixture models for spatial proteomic data

We observe N protein profiles each of length L , corresponding to the number of quantified fractions along the gradient density, including combining replicates. For $i = 1, \dots, N$, we denote the profile of the i -th protein by $\mathbf{x}_i = [x_{1i}, \dots, x_{Li}]$. We suppose that there are K known sub-cellular compartments to which each protein could localise (e.g. cytoplasm, endoplasmic reticulum, mitochondria, ...). Henceforth, we refer to these K sub-cellular compartments as *components*, and introduce component labels z_i , so that $z_i = k$ if the i -th protein localises to the k -th component. We denote by X_L the set of proteins whose component labels are known, and by X_U the set of unlabelled proteins. If protein i is in X_U , we desire the probability that $z_i = k$ for each $k = 1, \dots, K$. That is, for each unlabelled protein, we want the probability of belonging to each component (given a model and the observed data).

We initially model the distribution of profiles associated with proteins that localise to the k -th component as multivariate normal with mean vector $\boldsymbol{\mu}_k$ and covariance matrix Σ_k , so that:

$$\mathbf{x}_i | z_i = k \sim \mathcal{N}(\boldsymbol{\mu}_k, \Sigma_k). \quad (2)$$

For any i , we define the prior probability of the i -th protein localising to the k -th component to be $p(z_i = k) = \pi_k$. Letting $\boldsymbol{\theta} = \{\boldsymbol{\mu}_k, \Sigma_k\}_{k=1}^K$ denote the set of all component mean and covariance parameters, and $\boldsymbol{\pi} = \{\pi_k\}_{k=1}^K$ denote the set of all mixture weights, it follows (from the law of total probability) that:

$$p(\mathbf{x}_i | \boldsymbol{\theta}, \boldsymbol{\pi}) = \sum_{k=1}^K \pi_k f(\mathbf{x}_i | \boldsymbol{\mu}_k, \Sigma_k), \quad (3)$$

where $f(\mathbf{x} | \boldsymbol{\mu}, \Sigma)$ denotes the density of the multivariate normal with mean vector $\boldsymbol{\mu}$ and covariance matrix Σ evaluated at \mathbf{x} .

Equation (3) defines a generative probabilistic model known as a *mixture model*. Such models are useful for describing populations that are composed of a number of distinct homogeneous subpopulations. In our case, we model the full complement of measured proteins as being composed of K subpopulations, each corresponding to a different organelle or sub-cellular compartment. The literature of mixture model applications to biology is rich and some recent example include applications to retroviral integration sites [51], genome-wide associations studies [52], single-cell transcriptomics [53] and affinity purification MS proteomics [54].

Though some proteins are well described as belonging to a single component, many proteins multi-localise or might belong to uncharacterised organelles. In order to allow the

model to better account for these "outliers" that cannot be straightforwardly allocated to any single known component, we extend it by introducing an additional "outlier component". To do this, we augment our model by introducing a further indicator latent variable ϕ . Each protein \mathbf{x}_i is now described by an additional variable ϕ_i , with $\phi_i = 1$ indicating that protein \mathbf{x}_i belongs to a organelle derived component and $\phi_i = 0$ indicating that protein \mathbf{x}_i is not well described by these known components. This outlier component is modelled as a multivariate T distribution with degrees of freedom κ , mean vector \mathbf{M} , and scale matrix V . Thus equation (2) becomes

$$\mathbf{x}_i|z_i = k, \phi_i \sim \mathcal{N}(\boldsymbol{\mu}_k, \Sigma_k)^{\phi_i} \mathcal{T}(\kappa, \mathbf{M}, V)^{1-\phi_i}. \quad (4)$$

Further let $g(\mathbf{x}|\kappa, \mathbf{M}, V)$ denote the density of the multivariate T-distribution so that Equation (3) becomes:

$$p(\mathbf{x}_i|\boldsymbol{\theta}, \boldsymbol{\pi}, \phi_i, \kappa, \mathbf{M}, V) = \sum_{k=1}^K \pi_k (f(\mathbf{x}_i|\boldsymbol{\mu}_k, \Sigma_k)^{\phi_i} g(\mathbf{x}_i|\kappa, \mathbf{M}, V)^{1-\phi_i}). \quad (5)$$

For any i , we define the prior probability of the i -th protein belonging to the outlier component as $p(\phi_i = 0) = \epsilon$.

We can then rewrite equation (5) in the following way:

$$p(\mathbf{x}_i|\boldsymbol{\theta}, \boldsymbol{\pi}, \kappa, \epsilon, \mathbf{M}, V) = \sum_{k=1}^K \pi_k ((1 - \epsilon)(f(\mathbf{x}_i|\boldsymbol{\mu}_k, \Sigma_k) + \epsilon g(\mathbf{x}_i|\kappa, \mathbf{M}, V))), \quad (6)$$

Throughout we take $\kappa = 4$, \mathbf{M} as the global mean, and V as half the global variance of the data, including labelled and unlabelled proteins. The reason for formulating the model as in equation (5) is because it leads to a flexible modelling framework. Furthermore, ϕ has an elegant model selection interpretation, since it decides whether \mathbf{x}_i is better modelled by the known components or the outlier component. It is important to note that f and g could be replaced by many combinations of distributions and thus could be valuable in modelling other datasets. The choice of parameters for the multivariate T-distribution was decided so that it mimicked a multivariate normal component with the same mean and variance but with heavier tails to better capture dispersed proteins, which we refer to as outlier proteins throughout the text. The variance of the multivariate T-distribution is designed to be large such that is relatively flat when compared with multivariate Gaussian distributions which describe annotated components. Similar approaches for modelling outliers have been explored in the literature and often the outlier term is considered constant or as a Poisson process, independent of the observation [55, 56, 57, 58].

Model fitting

We adopt a Bayesian approach toward inferring the unknown parameters, $\boldsymbol{\theta} = \{\boldsymbol{\mu}_k, \Sigma_k\}_{k=1}^K$, $\boldsymbol{\pi} = \{\pi_k\}_{k=1}^K$, and ϵ of the mixture model presented in Equation (5). For $\boldsymbol{\pi}$, we take a

conjugate symmetric Dirichlet prior with parameter β , so that $\pi_1, \dots, \pi_K \sim \text{Dirichlet}(\beta)$; and for the component-specific parameters $\boldsymbol{\mu}_k$ and Σ_k we take conjugate normal-inverse-Wishart (NIW) priors with parameters $\{\boldsymbol{\mu}_0, \lambda_0, \nu_0, S_0\}$, so that:

$$\mu_k, \Sigma_k \sim \mathcal{N}\left(\boldsymbol{\mu}_k | \boldsymbol{\mu}_0, \frac{\Sigma_k}{\lambda_0}\right) IW(\Sigma_k | \nu_0, S_0). \quad (7)$$

We also place a conjugate Beta prior on ϵ with parameters u and v , so that $\epsilon \sim \mathcal{B}(u, v)$. Allowing ϵ to be random allows us to infer the number of proteins that are better described by an outlier component rather than any known component.

The full model, which we henceforth refer to as a T-augmented Gaussian Mixture model (TAGM), can then be summarised by the plate diagram shown in Figure 13.

To perform inference for the parameters, we make use of both the labelled and unlabelled data. For the labelled data X_L , since z_i and ϕ_i are known for these proteins, we can update the parameters with their data analytically by exploiting conjugacy of the priors [see, for example, 59]. For the unlabelled data we do not have such information and so in the next sections we explain how to make inferences of the latent variables.

Prediction of localisation of unlabelled proteins

Having obtained the posterior distribution of the model parameters analytically using, at first, the labelled data only, we wish to predict the component to which each of the unlabelled proteins belongs. The probability that a protein belongs to any of the K known components, that is $z_i = k$ and $\phi_i = 1$, is given by (see supporting information section S1 for derivations):

$$p(\phi_i = 1, z_i = k | \mathbf{x}_i, \boldsymbol{\theta}, \boldsymbol{\pi}, \epsilon, \kappa, \mathbf{M}, V) = \frac{\pi_k (1 - \epsilon) f(\mathbf{x}_i | \boldsymbol{\mu}_k, \Sigma_k)}{\sum_{k=1}^K \pi_k ((1 - \epsilon) f(\mathbf{x}_i | \boldsymbol{\mu}_k, \Sigma_k) + \epsilon g(\mathbf{x}_i | \kappa, \mathbf{M}, V))}, \quad (8)$$

whilst on the other hand,

$$p(\phi_i = 0, z_i = k | \mathbf{x}_i, \boldsymbol{\theta}, \boldsymbol{\pi}, \kappa, \epsilon, \mathbf{M}, V) = \frac{\pi_k \epsilon g(\mathbf{x}_i | \kappa, \mathbf{M}, V)}{\sum_{k=1}^K \pi_k ((1 - \epsilon) f(\mathbf{x}_i | \boldsymbol{\mu}_k, \Sigma_k) + \epsilon g(\mathbf{x}_i | \kappa, \mathbf{M}, V))}. \quad (9)$$

Processing of the unlabelled data can be done by inferring *maximum a posteriori* (MAP) estimates for the parameters. However, this approach fails to account for the uncertainty in the parameters, thus we additionally explore inferring the distribution over these parameters.

Maximum a posteriori prediction

We use the Expectation-Maximisation (EM) algorithm [60] to find *maximum a posteriori* (MAP) estimates for the parameters [see, for example, 61]. To specify the parameters of the prior distributions, we use a simple set of heuristics provided by [62]. By defining the following quantities

$$\begin{aligned}
a_{ik} &= p(z_i = k, \phi_i = 1 | \mathbf{x}_i), b_{ik} = p(z_i = k, \phi_i = 0 | \mathbf{x}_i) \\
w_{ik} &= p(z_i = k | x_i) = a_{ik} + b_{ik} \\
a_k &= \sum_{i=1}^n a_{ik}, a = \sum_{k=1}^K a_k \\
b_k &= \sum_{i=1}^n b_{ik}, b = \sum_{k=1}^K b_k \\
r_k &= \sum_{i=1}^n w_{ik},
\end{aligned} \tag{10}$$

we can compute

$$\begin{aligned}
\lambda_k &= \lambda_0 + a_k, \\
\nu_k &= \nu_0 + a_k, \\
m_k &= \frac{a_k \bar{\mathbf{x}}_k + \lambda_0 \mu_0}{\lambda_k}, \\
S_k^{-1} &= S_0^{-1} + \frac{\lambda_0 a_k}{\lambda_k} (\bar{\mathbf{x}}_k - \mu_0)^T (\bar{\mathbf{x}}_k - \mu_0) + \sum_{i=1}^n a_{ik} (x_i - \bar{\mathbf{x}}_k)^T (x_i - \bar{\mathbf{x}}_k).
\end{aligned} \tag{11}$$

Then the parameters of the posterior mode are:

$$\begin{aligned}
\hat{\mu}_k &= m_k \\
\hat{\Sigma}_k &= \frac{1}{\nu_k + D + 2} S_k^{-1}.
\end{aligned} \tag{12}$$

We note if \mathbf{x}_i is a labelled protein then $a_{ik} = 1$ and these parameters can be updated without difficulty. The above equation constitutes a backbone of the E-step of the EM algorithm, with the entire algorithm specified by the following summary:

E-Step: Given the current parameters compute the values given by equations (10), with formulae provided in equations (8) and (9).

M-Step: Compute

$$\epsilon = \frac{u + b - 1}{(a + b) + (u + v) - 2},$$

and

$$\pi_k = \frac{r_k + \beta_k - 1}{N + \sum \beta_k - K},$$

as well as

$$\bar{\mathbf{x}}_k = \frac{1}{a_k} \left(\sum_{i=1}^n a_{ik} \mathbf{x}_i \right).$$

Finally, compute the MAP estimates given by equations (12). These estimates are then used in the following iteration of the E-step.

Denoting by Q the expected value of the log-posterior and letting t denote the current iteration of the EM algorithm, we iterate until $|Q(\boldsymbol{\theta}|\boldsymbol{\theta}_t) - Q(\boldsymbol{\theta}|\boldsymbol{\theta}_{t-1})| < \delta$ for some pre-specified $\delta > 0$. Once we have found MAP estimates for the parameters $\boldsymbol{\theta}_{MAP}$, $\boldsymbol{\pi}_{MAP}$ and ϵ_{MAP} we proceed to perform prediction. We plug the MAP parameter estimates into Equation (8) in order to obtain the posterior probability of protein i localising to component k , $p(z_i = k, \phi = 1|\mathbf{x}_i, \boldsymbol{\theta}_{MAP}, \boldsymbol{\pi}_{MAP}, \epsilon_{MAP}, \kappa, \mathbf{M}, V)$. To make a final assignment, we may allocate each protein according to the component that has maximal probability. A full technical derivation of the EM algorithm can be found in the supporting information (section S1).

Uncertainty in the posterior localisation probabilities

The MAP approach described above provides us with a probabilistic assignment, $p(z_i = k, \phi = 1|\mathbf{x}_i, \boldsymbol{\theta}_{MAP}, \boldsymbol{\pi}_{MAP}, \epsilon_{MAP}, \kappa, \mathbf{M}, V)$, of each unlabelled protein to each component. However, it fails to account for the uncertainty in the parameters $\boldsymbol{\theta}$, $\boldsymbol{\pi}$ and ϵ . To address this, we can sample parameters from the posterior distribution.

Let $\{\boldsymbol{\theta}^{(t)}, \boldsymbol{\pi}^{(t)}, \epsilon^{(t)}\}_{t=1}^T$ be a set of T sampled values for the parameters $\boldsymbol{\theta}$, $\boldsymbol{\pi}$, ϵ , drawn from the posterior.

The assignment probabilities can then be summarised by the Monte-Carlo average:

$$p(z_i = k, \phi = 1|\mathbf{x}_i, \epsilon, \mathbf{M}, V) \approx T^{-1} \sum_{t=1}^T p(z_i = k, \phi = 1|\mathbf{x}_i, \boldsymbol{\theta}^{(t)}, \boldsymbol{\pi}^{(t)}, \epsilon^{(t)}, \kappa, \mathbf{M}, V).$$

Other summaries of the assignment probabilities can be determined in the usual ways to obtain, for example, interval-estimates. We summarise interval-estimates using the 95% equi-tailed interval, which is defined by the 0.025 and 0.975 quantiles of the distribution of assignment probabilities, $\{p(z_i = k, \phi = 1|\mathbf{x}_i, \boldsymbol{\theta}^{(t)}, \boldsymbol{\pi}^{(t)}, \epsilon^{(t)}, \mathbf{M}, V)\}_{t=1}^T$.

Sampling parameter values in our model requires us to compute the required conditional probabilities and then a straightforward Gibbs sampler can be used to sample in turn from these conditionals. In addition, we can bypass sampling the parameters by exploiting the conjugacy of our priors. By marginalising parameters in our model we can obtain an efficient collapsed Gibbs sampler and therefore only sample the component allocation probabilities and the outlier allocation probabilities. The derivations and required conditionals can be found in the supporting information (section S2).

Classifier assessment

We compared the classification performance of the two above learning schemes to the K-nearest neighbours (KNN) and the weighted support vector machine (SVM) classifiers.

The following schema was used to assess the classifier performance of all methods. We split the marker sets for each experiment into a class-stratified training (80%) and test (20%) partitions, with the separation formed at random. The true classes of the test profiles are

withheld from the classifier, whilst the algorithm is trained. The algorithm is then assessed on its ability to predict the classes of the proteins in the test partition for generalisation accuracy. How each classifier is trained is specific to that classifier. The KNN and SVM have hyperparameters optimised using 5-fold cross-validation. This 80/20 data stratification is performed 100 times in order to produce 100 sets of macro-F1 [63] scores and class specific F1 scores [25]. The F1 score is the harmonic mean of the precision and recall, more precisely:

$$\text{precision} = \frac{tp}{tp + fp}, \text{recall} = \frac{tp}{tp + fn}.$$

tp denotes the number of true positives; fp the number of false positives and fn the number of false negatives. Thus

$$F1 = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}.$$

High Macro F1 scores indicates that marker proteins in the test dataset are consistently correctly assigned by the classifier. We note that accuracy alone is an inadequate measure of performance, since it fails to quantify false positives.

However, a Bayesian Generative classifier produces probabilistic assignment of observations to classes. Thus while the classifier may make an incorrect assignment it may do so with low probability. The F1 score is unforgiving in this situation and will not use this information. To measure this uncertainty, we introduce the quadratic loss which allows us to compare probabilistic assignments [64]. For the SVM, a logistic distribution is fitted using maximum likelihood estimation to the decision values of all binary classifiers. Then, the membership probabilities for the multi-class classification is calculated using quadratic optimisation. The logistic regression model assumes errors which are distributed according to a centred Laplace distribution for the predictions, where maximum likelihood estimation is used to estimate the scale parameter [65]. For the KNN classifier, we interpret the proportion of neighbours belonging to each class as a non-parametric posterior probability. To avoid non-zero probabilities for classes we perform Laplace smoothing; that is, the posterior allocation probability is given by

$$p(z_i = k|x_i) = \frac{N_{ik} + \alpha d_k C}{K + \alpha C}, \quad (13)$$

where N_{ik} is the number of neighbours belonging to class k in the neighbourhood of x_i , C is the number of classes, K is the number of nearest neighbours (optimised through 5-fold cross validation) and d_k is the incidence rate of each class in the training set. Finally, $\alpha > 0$ is the pseudo-count smoothing parameter. Motivated by a Bayesian interpretation of placing a Jeffrey's type Dirichlet prior over multinomial counts, we choose $\alpha = 0.5$ [66, 67, 68]. The quadratic loss is given by the following formula:

$$Q_2 = \sum_{i=1}^N \|q_i - p_i\|_2^2, \quad (14)$$

where $\|\cdot\|_2$ is the l_2 norm and q_i is the true classification vector and p_i is a vector of predicted assignments to each class. It is useful to note that the corresponding risk function is the mean square error (MSE), which is the expected value of the quadratic loss.

Acknowledgments

The authors would also like to thank Dr Sean B. Holden, University of Cambridge, for helpful discussions.

References

- [1] Gibson T. Cell regulation: determined to signal discrete cooperation. *Trends in biochemical sciences*. 2009;34(10):471–482.
- [2] Olkkonen V, Ikonen E. When intracellular logistics fails-genetic defects in membrane trafficking. *Journal of cell science*. 2006;119(24):5031–5045.
- [3] Luheshi L, Crowther D, Dobson C. Protein misfolding and disease: from the test tube to the organism. *Current opinion in chemical biology*. 2008;12(1):25–31.
- [4] Laurila K, Vihinen M. Prediction of disease-related mutations affecting protein localization. *BMC genomics*. 2009;10(1):122.
- [5] De Matteis M, Luini A. Mendelian disorders of membrane trafficking. *New England Journal of Medicine*. 2011;365(10):927–938.
- [6] Cody N, Iampietro C, Lécuyer E. The many functions of mRNA localization during normal development and disease: from pillar to post. *Wiley Interdisciplinary Reviews: Developmental Biology*. 2013;2(6):781–796.
- [7] Kau T, Way J, Silver P. Nuclear transport and cancer: from mechanism to intervention. *Nature Reviews Cancer*. 2004;4(2):106–117.
- [8] Rodriguez J, Au W, Henderson B. Cytoplasmic mislocalization of BRCA1 caused by cancer-associated mutations in the BRCT domain. *Experimental cell research*. 2004;293(1):14–21.
- [9] Latorre I, Roh M, Frese K, Weiss R, Margolis B, Javier R. Viral oncoprotein-induced mislocalization of select PDZ proteins disrupts tight junctions and causes polarity defects in epithelial cells. *Journal of cell science*. 2005;118(18):4283–4293.
- [10] Shin S, Smith J, Rezniczek G, Pan S, Chen R, Brentnall T, et al. Unexpected gain of function for the scaffolding protein plectin due to mislocalization in pancreatic cancer. *Proceedings of the National Academy of Sciences*. 2013;110(48):19414–19419.
- [11] Siljee J, Wang Y, Bernard A, Ersoy B, Zhang S, Marley A, et al. Subcellular localization of MC4R with ADCY3 at neuronal primary cilia underlies a common pathway for genetic predisposition to obesity. *Nat Genet*. 2018 Jan;.
- [12] Gatto L, Vizcaíno J, Hermjakob K, Huber W, Lilley K. Organelle proteomics experimental designs and analysis. *Proteomics*. 2010;10(22):3957–3969.

- [13] Christoforou A, Mulvey C, Breckels L, Geladaki A, Hurrell T, Hayward P, et al. A draft map of the mouse pluripotent stem cell spatial proteome. *Nature communications*. 2016;7:9992.
- [14] Dunkley T, Hester S, Shadforth I, Runions J, Weimar T, Hanton S, et al. Mapping the Arabidopsis organelle proteome. *Proceedings of the National Academy of Sciences*. 2006;103(17):6518–6523.
- [15] Itzhak D, Tyanova S, Cox J, Borner G. Global, quantitative and dynamic mapping of protein subcellular localization. *Elife*. 2016;5:e16950.
- [16] Parsons H, Fernández-Niño S, Heazlewood J. Separation of the plant Golgi apparatus and endoplasmic reticulum by free-flow electrophoresis. *Methods in molecular biology* (Clifton, NJ). 2014;1072:527.
- [17] Heard W, Sklenář J, Tome D, Robatzek S, Jones A. Identification of regulatory and cargo proteins of endosomal and secretory pathways in Arabidopsis thaliana by proteomic dissection. *Molecular & Cellular Proteomics*. 2015;14(7):1796–1813.
- [18] Dunkley T, Watson R, Griffin J, Dupree P, Lilley K. Localization of organelle proteins by isotope tagging (LOPIT). *Molecular & Cellular Proteomics*. 2004;3(11):1128–1134.
- [19] Sadowski P, Dunkley T, Shadforth I, Dupree P, Bessant C, Griffin J, et al. Quantitative proteomic approach to study subcellular localization of membrane proteins. *Nature protocols*. 2006;1(4):1778–1789.
- [20] Mulvey C, Breckels L, Geladaki A, Britovšek N, Nightingale D, Christoforou A, et al. Using hyperLOPIT to perform high-resolution mapping of the spatial proteome. *Nature Protocols*. 2017;12(6):1110–1135.
- [21] De Duve C, Beaufay H. A short history of tissue fractionation. *The Journal of cell biology*. 1981;91(3):293.
- [22] Breckels L, Mulvey C, Lilley K, Gatto L. A Bioconductor workflow for processing and analysing spatial proteomics data. *F1000Research*. 2016;5.
- [23] Gatto L, Breckels L, Burger T, Nightingale D, Groen A, Campbell C, et al. A foundation for reliable spatial proteomics data analysis. *Molecular & Cellular Proteomics*. 2014;p. mcp-M113.
- [24] Breckels L, Gatto L, Christoforou A, Groen A, Lilley K, Trotter M. The effect of organelle discovery upon sub-cellular protein localisation. *Journal of proteomics*. 2013;88:129–140.
- [25] Breckels L, Holden S, Wojnar D, Mulvey C, Christoforou A, Groen A, et al. Learning from heterogeneous data sources: an application in spatial proteomics. *PLoS computational biology*. 2016;12(5):e1004920.

- [26] Tan D, Dvinge H, Christoforou A, Bertone P, Martinez Arias A, Lilley K. Mapping organelle proteins and protein complexes in *Drosophila melanogaster*. *Journal of proteome research*. 2009;8(6):2667–2678.
- [27] Hall S, Hester S, Griffin J, Lilley K, Jackson A. The organelle proteome of the DT40 lymphocyte cell line. *Molecular & Cellular Proteomics*. 2009;8(6):1295–1305.
- [28] Thul P, Akesson L, Wiking M, Mahdessian D, Geladaki A, Ait Blal H, et al. A subcellular map of the human proteome. *Science*. 2017;.
- [29] Groen A, Sancho-Andres G, Breckels L, Gatto L, Aniento F, Lilley K. Identification of trans-Golgi network proteins in *Arabidopsis thaliana* root tissue. *Journal of proteome research*. 2014;13(2):763–776.
- [30] Ohta S, Bukowski-Wills J, Sanchez-Pulido L, de Lima Alves F, Wood L, Chen Z, et al. The protein composition of mitotic chromosomes determined using multiclassifier combinatorial proteomics. *Cell*. 2010;142(5):810–821.
- [31] Nikolovski N, Rubtsov D, Segura M, Miles G, Stevens T, Dunkley T, et al. Putative glycosyltransferases and other plant Golgi apparatus proteins are revealed by LOPIT proteomics. *Plant physiology*. 2012;160(2):1037–1051.
- [32] Tardif M, Atteia A, Specht M, Cogne G, Rolland N, Brugière S, et al. PredAlgo: a new subcellular localization prediction tool dedicated to green algae. *Molecular biology and evolution*. 2012;29(12):3625–3639.
- [33] Kirk P, Babbie A, Stumpf M. Systems biology (un) certainties. *Science*. 2015;350(6259):386–388.
- [34] R Core Team. R: A Language and Environment for Statistical Computing. Vienna, Austria; 2017. Available from: <https://www.R-project.org/>.
- [35] Gatto L, Lilley K. MSnbase - an R/Bioconductor package for isobaric tagged mass spectrometry data visualization, processing and quantitation. *Bioinformatics*. 2012;28:288–289.
- [36] Gatto L, Breckels L, Wieczorek S, Burger Y, Lilley K. Mass-spectrometry based spatial proteomics data analysis using pRoloc and pRolocdata. *Bioinformatics*. 2014;.
- [37] Gentleman R, Carey V, Bates D, Bolstad B, Dettling M, Dudoit S, et al. Bioconductor: open software development for computational biology and bioinformatics. *Genome biology*. 2004;5(10):R80.
- [38] Huber W, Carey V, Gentleman R, Anders S, Carlson M, Carvalho B, et al. Orchestrating high-throughput genomic analysis with Bioconductor. *Nature methods*. 2015;12(2):115.
- [39] McAlister G, Nusinow D, Jedrychowski M, Wu L, Huttlin E, Erickson B, et al. MultiNotch MS3 enables accurate, sensitive, and multiplexed detection of differential expression across cancer cell line proteomes. *Analytical chemistry*. 2014;86(14):7150–7158.

- [40] Ting L, Rad R, Gygi S, Haas W. MS3 eliminates ratio distortion in isobaric multiplexed quantitative proteomics. *Nature methods*. 2011;8(11):937.
- [41] Gelman A, Rubin D. Inference from iterative simulation using multiple sequences. *Statistical science*. 1992;p. 457–472.
- [42] Brooks S, Gelman A. General methods for monitoring convergence of iterative simulations. *Journal of computational and graphical statistics*. 1998;7(4):434–455.
- [43] Boyle E, Weng S, Gollub J, Jin H, Botstein D, Cherry M, et al. GO::TermFinder – open source software for accessing Gene Ontology information and finding significantly enriched Gene Ontology terms associated with a list of genes. *Bioinformatics*. 2004;20(18):3710–3715.
- [44] Yu G, Wang L, Han Y, He Q. clusterProfiler: an R package for comparing biological themes among gene clusters. *Omics: a journal of integrative biology*. 2012;16(5):284–287.
- [45] Hsu C, Chang C, Lin C. A practical guide to support vector classification; 2010.
- [46] Hirst J, Itzhak D, Antrobus R, Borner G, Robinson M. Role of the AP-5 adaptor protein complex in late endosome-to-Golgi retrieval. *PLoS biology*. 2018;16(1):e2004411.
- [47] Beltran P, Mathias RA, Cristea I. A portrait of the human organelle proteome in space and time during cytomegalovirus infection. *Cell systems*. 2016;3(4):361–373.
- [48] Benjamini Y, Hochberg Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the royal statistical society Series B (Methodological)*. 1995;p. 289–300.
- [49] Jeffery C. Moonlighting proteins - an update. *Molecular BioSystems*. 2009;5(4):345–350.
- [50] Shannon C. A mathematical theory of communication. *The Bell System Technical Journal*. 1948 July;27(3):379–423.
- [51] Kirk P, Huvet M, Melamed A, Maertens G, Bangham C. Retroviruses integrate into a shared, non-palindromic DNA motif. *Nature microbiology*. 2016;2:16212.
- [52] Liley J, Todd J, Wallace C. A method for identifying genetic heterogeneity within phenotypically defined disease subgroups. *Nature genetics*. 2017;49(2):310.
- [53] LÄnnberg T, Svensson V, James K, Fernandez-Ruiz D, Sebina I, Montandon R, et al. Single-cell RNA-seq and computational analysis using temporal mixture modeling resolves TH1/TFH fate bifurcation in malaria. *Science Immunology*. 2017;2(9).
- [54] Choi H, Kim S, Gingras A, Nesvizhskii A. Analysis of protein complexes through model-based biclustering of label-free quantitative AP-MS data. *Molecular Systems Biology*. 2010;6(1):385. Available from: <https://onlinelibrary.wiley.com/doi/abs/10.1038/msb.2010.41>.

- [55] Banfield J, Raftery A. Model-based Gaussian and non-Gaussian clustering. *Biometrics*. 1993;p. 803–821.
- [56] Cooke E, Savage R, Kirk P, Darkins R, Wild D. Bayesian hierarchical clustering for microarray time series data with replicates and outlier measurements. *BMC bioinformatics*. 2011;12(1):399.
- [57] Coretto P, Hennig C. Robust improper maximum likelihood: tuning, computation, and a comparison with other methods for robust Gaussian clustering. *Journal of the American Statistical Association*. 2016;111(516):1648–1659.
- [58] Hennig C. Breakdown points for maximum likelihood estimators of location-scale mixtures. *Annals of Statistics*. 2004;p. 1313–1340.
- [59] Gelman A, Carlin J, Stern H, Rubin D. *Bayesian Data Analysis*. London: Chapman & Hall; 1995.
- [60] Dempster A, Laird N, Rubin D. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the royal statistical society Series B (methodological)*. 1977;p. 1–38.
- [61] Murphy K. *Machine learning: a probabilistic perspective*; 2012.
- [62] Fraley C, Raftery A. Bayesian regularization for normal mixture estimation and model-based clustering. *Journal of Classification*. 2007;24(2):155–181.
- [63] He H, Garcia E. Learning from imbalanced data. *IEEE Transactions on knowledge and data engineering*. 2009;21(9):1263–1284.
- [64] Gneiting T, Raftery A. Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*. 2007;102(477):359–378.
- [65] Meyer S, Dimitriadou E, Hornik K, Weingessel A, Leisch F. e1071: Misc Functions of the Department of Statistics, Probability Theory Group (Formerly: E1071), TU Wien; 2018. R package version 1.7-0. Available from: <https://CRAN.R-project.org/package=e1071>.
- [66] Hazimeh H, Zhai C; ACM. Axiomatic analysis of smoothing methods in language models for pseudo-relevance feedback. *Proceedings of the 2015 International Conference on The Theory of Information Retrieval*. 2015;p. 141–150.
- [67] Valcarce D, Parapar J, Barreiro A. Additive Smoothing for Relevance-Based Language Modelling of Recommender Systems. *Proceedings of the 4th Spanish Conference on Information Retrieval*. 2016;p. 1–8. Available from: <http://doi.acm.org/10.1145/2934732.2934737>.
- [68] Manning C, Raghavan P, Schütze H. *Introduction to Information Retrieval*. New York, NY, USA: Cambridge University Press; 2008.
- [69] Foster J, de Hoog C, Zhang Y, Zhang Y, Xie X, Mootha V, et al. A mammalian organelle map by protein correlation profiling. *Cell*. 2006;125(1):187–199.

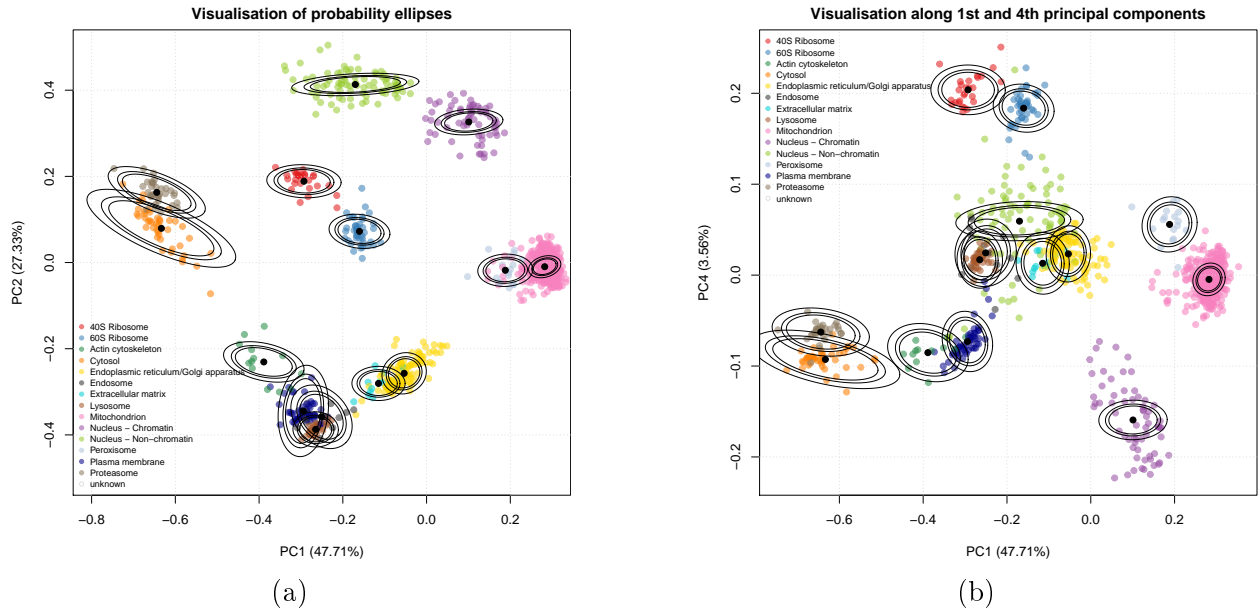


Figure 1: (a) PCA plot of the 1st and 2nd principal components for the curated marker proteins of the mouse stem cell data. The organelles are, in general, well separated. Though some organelles overlap, they are separated along different principal components. The densities used to produce the ellipses are derived from the MAP estimates. (b) Marker resolution along the 1st and 4th principal components show that the mitochondrion and peroxisome markers are well resolved, despite overlapping in the 1st and 2nd component. We also see that the ER/Golgi apparatus markers are better separated from the extracellular matrix markers.

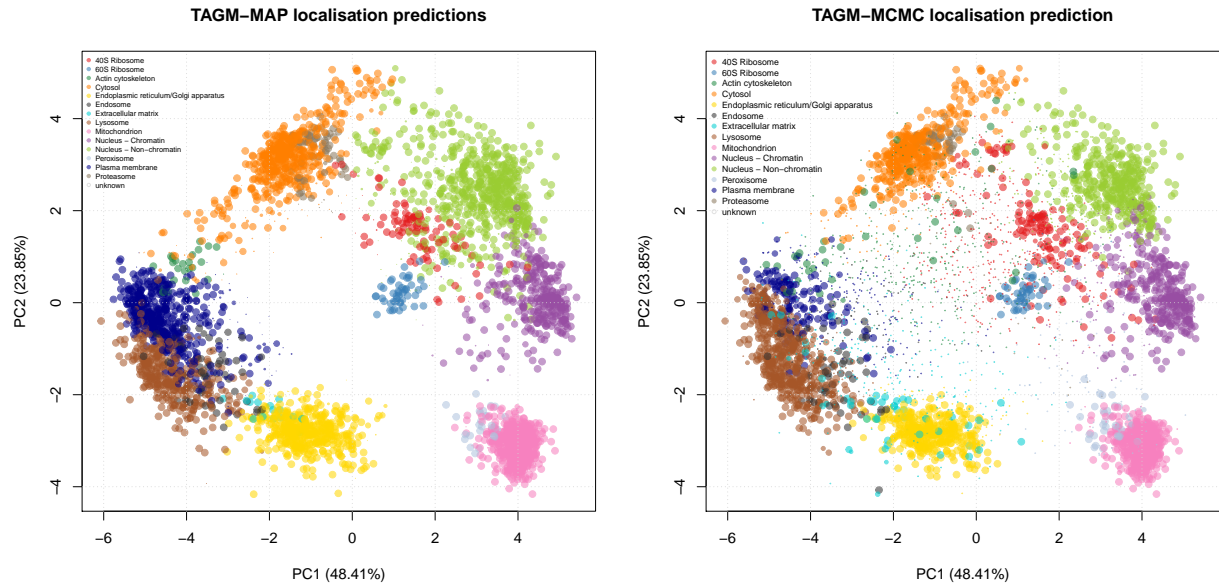


Figure 2: PCA plot of the protein quantitation data with colours representing the predicted class (5032 proteins) illustrating protein localisation predictions using TAGM-MAP (left) and TAGM-MCMC (right) respectively. The pointer size of a protein is scaled to the probability that particular protein was assigned to that organelle. Markers, proteins whose localisations are already known, are automatically assigned a probability of 1 and the size of the pointer reflects this.

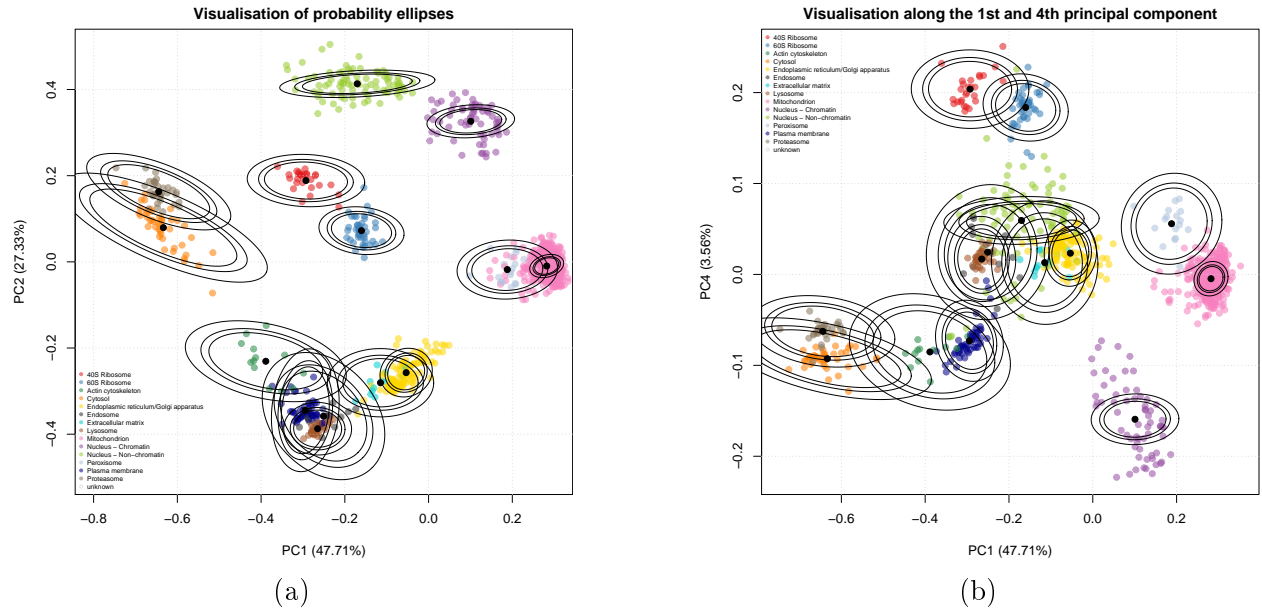


Figure 3: (a) Probability ellipses produced from using the MCMC method. The density is the expected value from the NIW distribution. (b) Probability ellipses visualised along the 1st and 4th principal component also from the MCMC method.

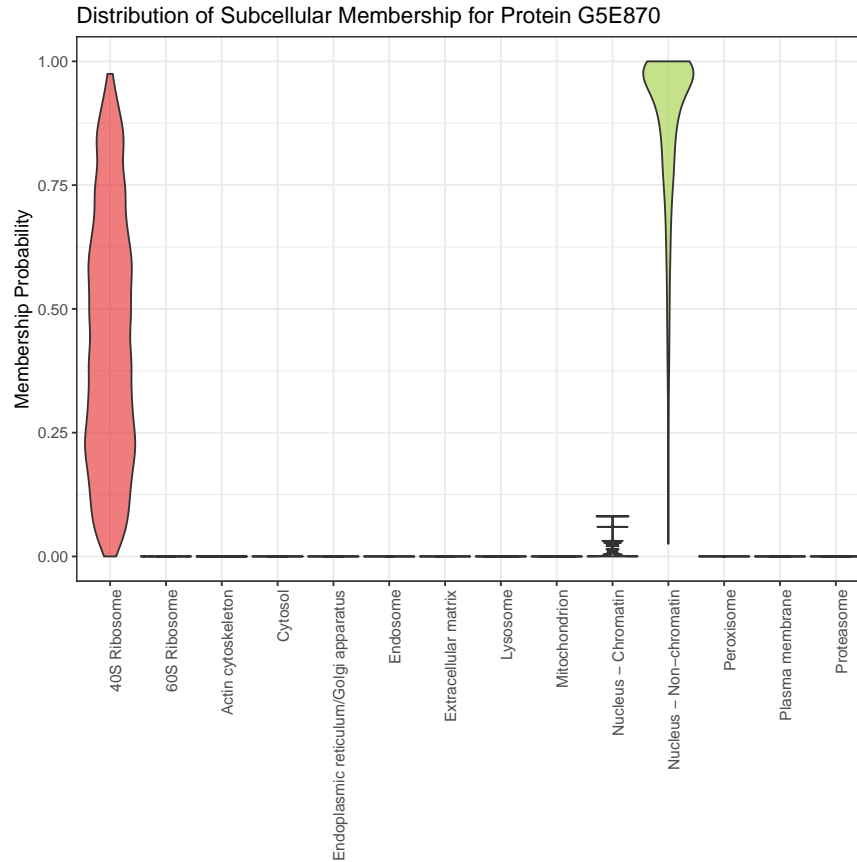


Figure 4: Violin plot revealing the posterior distribution of localisation probabilities of protein E3 ubiquitin-protein ligase (G5E870) to organelles and sub-cellular niches. The most probable localisation is nucleus non-chromatin, however there is uncertainty associated with this assignment.

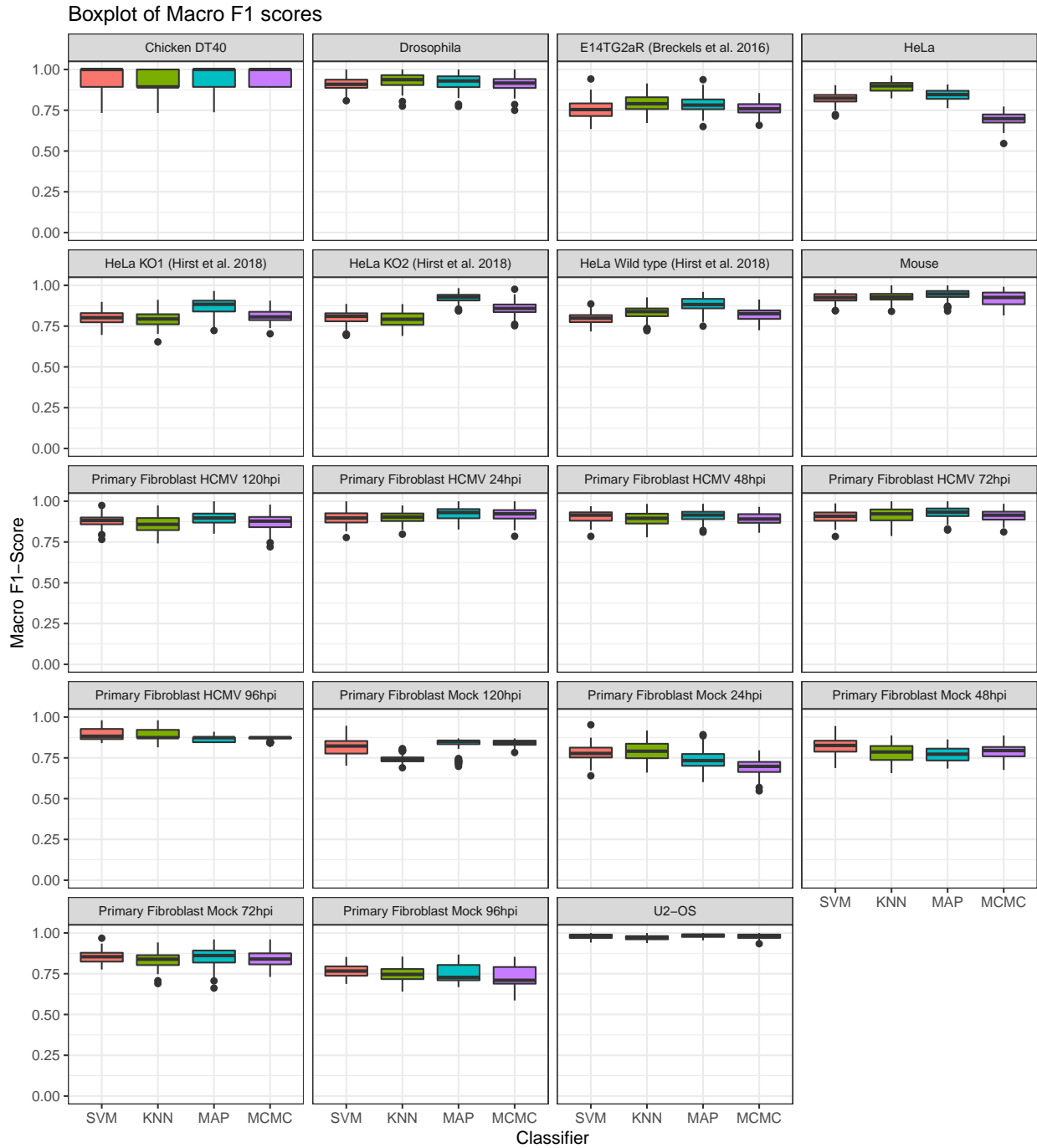


Figure 5: Boxplots of the distributions of Macro F1 scores for all spatial proteomics datasets.

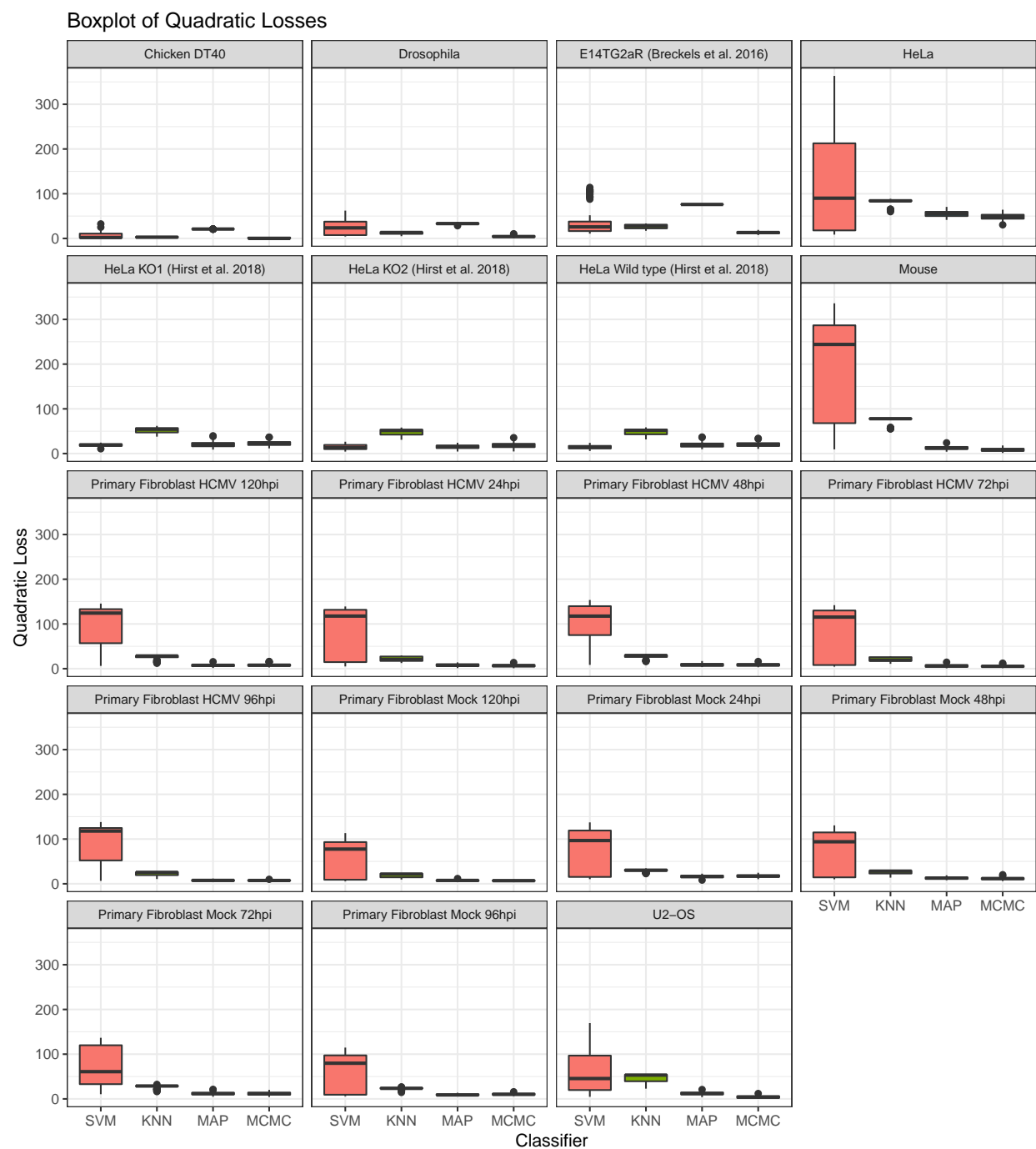


Figure 6: Boxplots of the distributions of Quadratic losses for all spatial proteomics datasets.

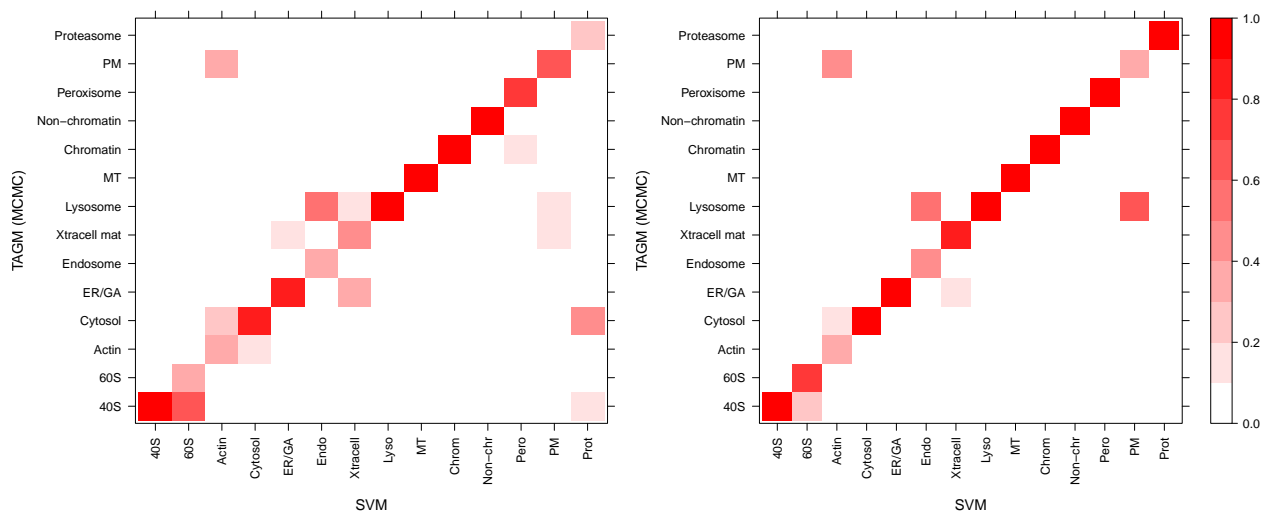


Figure 7: A heatmap representation of a contingency table, where we compare assignment results for proteins with unknown protein localisation using the TAGM-MCMC and SVM. The scale ranges from 0 to 1 with values indicating the proportion of assigned proteins to that sub-cellular location. Values along the diagonal represent agreement between classifiers whilst other values represent disagreement. The coherence between the classifiers is very high. (a) In this case we set a probability threshold of 0.95 for the TAGM assignments with no threshold for the SVM. (b) In this case we set a 5% FDR threshold for the SVM and no threshold for the TAGM-MCMC.

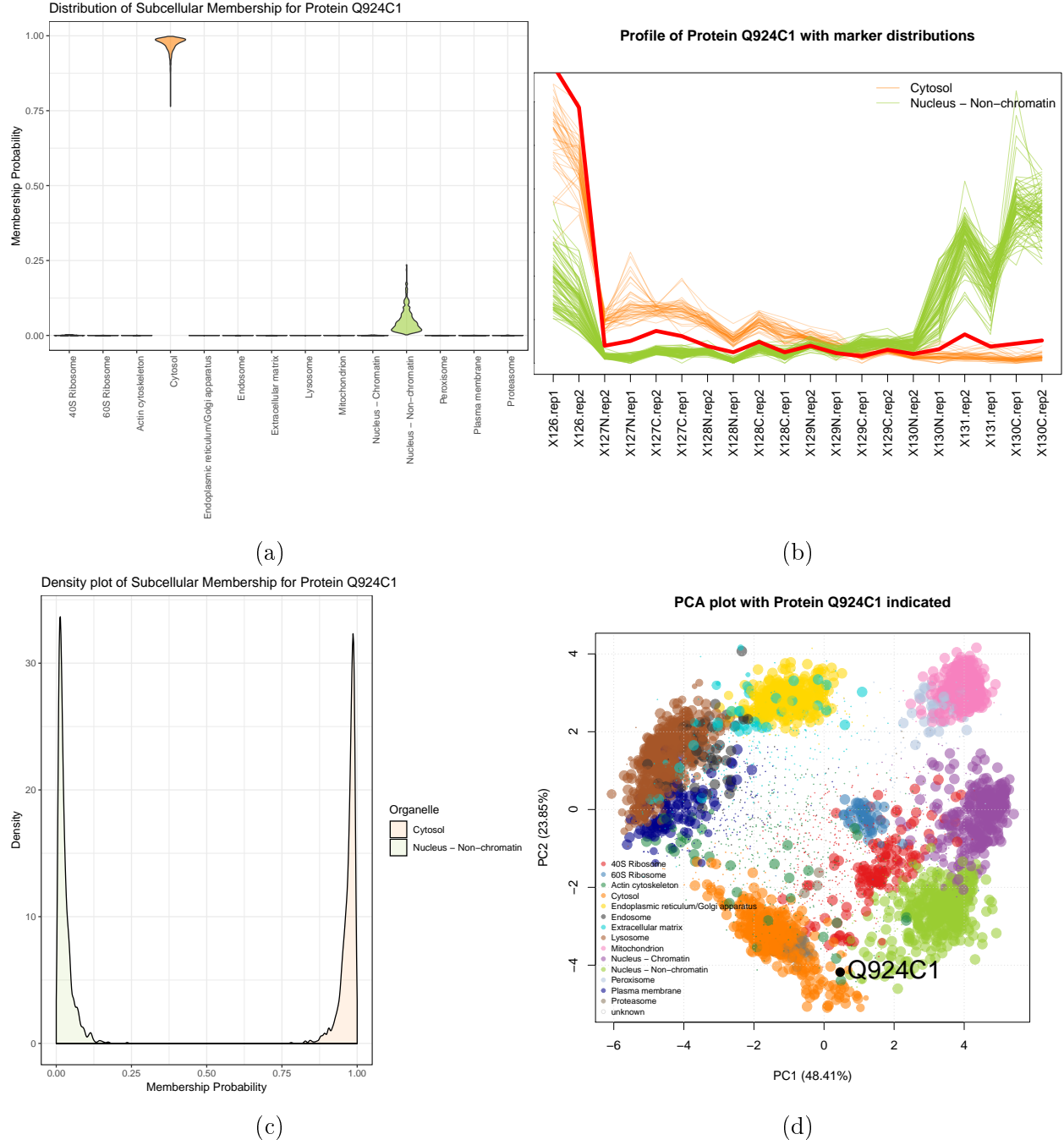


Figure 8: Exportin 5 (Q924C1) showing localisation to the cytosol with some uncertainty about association to the nucleus non-chromatin. (a) The violin plot shows uncertain localisation between these two sub-cellular localisations. (b) The quantitative profile of this protein shows mixed profile between the profiles of the organelle markers. (c) The density plot shows a complex distribution over localisations for this protein. (d) The protein Q924C1 has steady state distribution between the cytosol and nucleus non-chromatin.

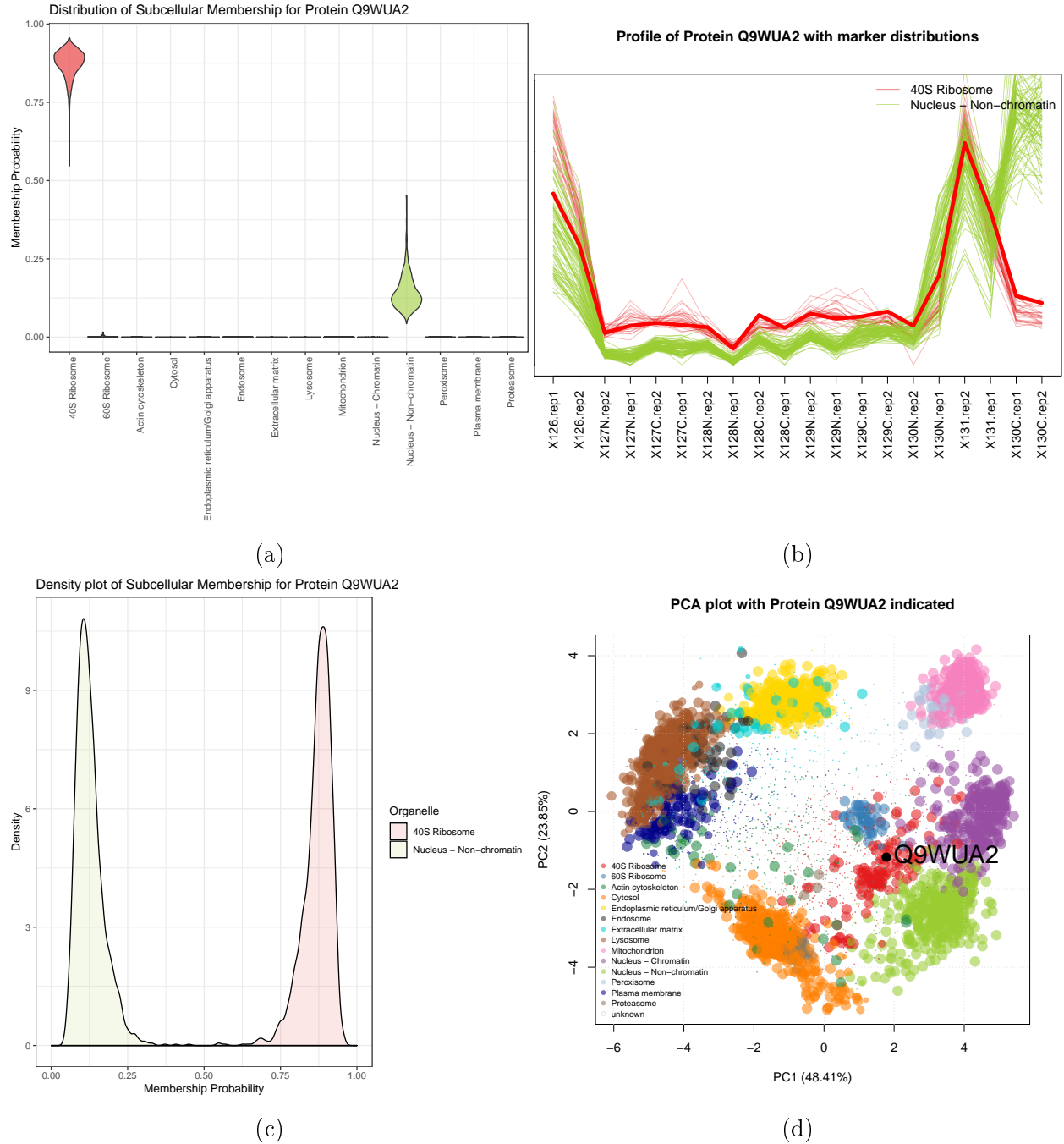


Figure 9: Phenylalanine-tRNA ligase beta subunit protein TRIP12 (Q9WUA2) showing localisation to the 40S Ribosome with some uncertainty about association to the nucleus non-chromatin. (a) The violin plot shows uncertain localisation between these two subcellular localisations. (b) The quantitative profile of this protein shows mixed profile between the profiles of the organelle markers. (c) The density plot shows a complex distribution over localisations for this protein. (d) The protein Q9WUA2 has steady state distribution skewed towards the 40S Ribosome and close to the nucleus non-chromatin.

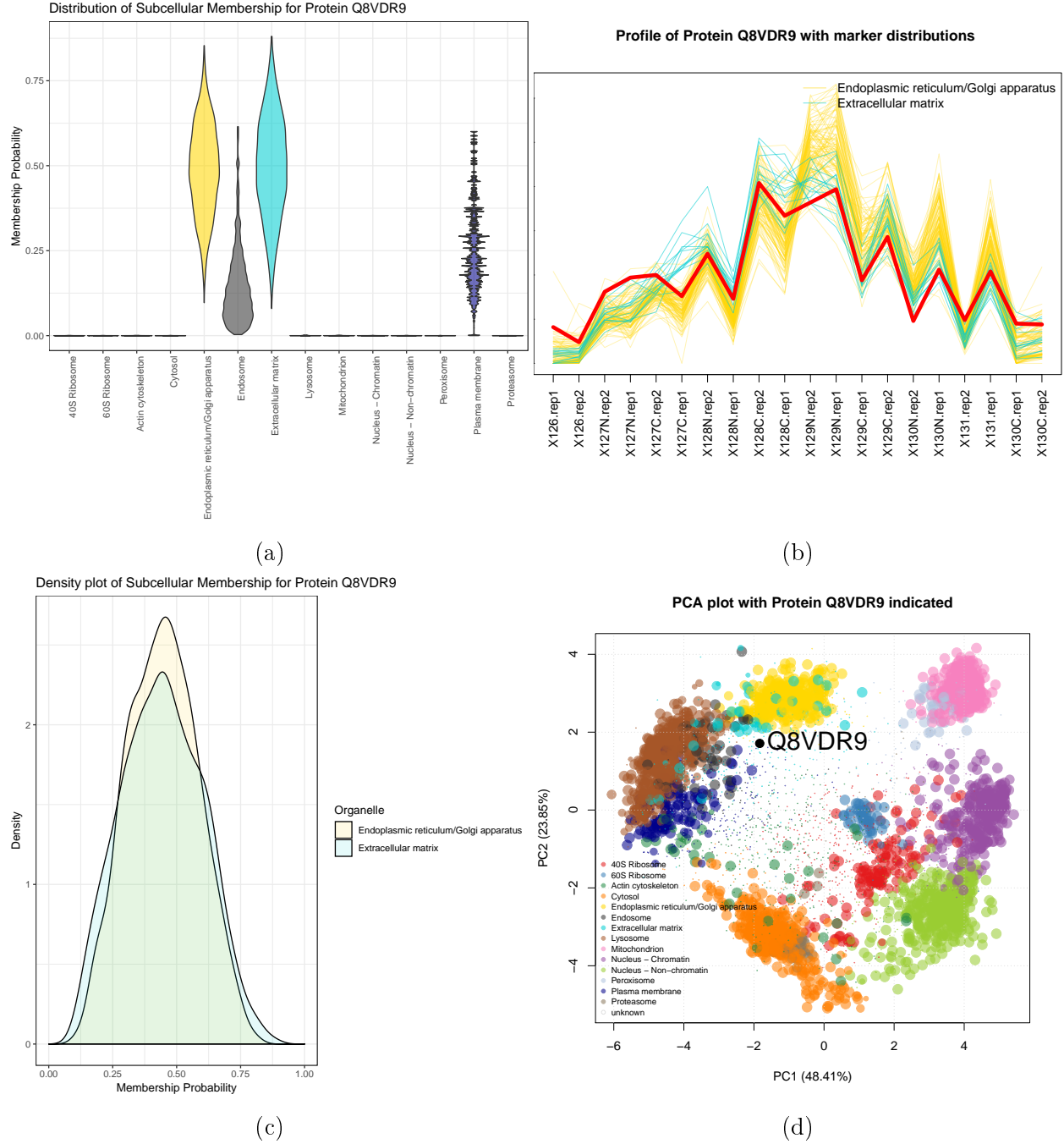


Figure 10: Q8VDR9 showing localisation to the outlier component. (a) The violin plot shows uncertain localisation between several sub-cellular niches. (b) The quantitative profile of this protein shows mixed profile between the profiles of the organelle markers. (c) The density plot shows a similar localisation probabilities for both the ER/Golgi and Extracellular matrix. (d) The protein Q8VDR9 has steady state distribution in the centre of the plot skewed toward the secretory pathway; in particular, the ER/Golgi and Extracellular matrix components.

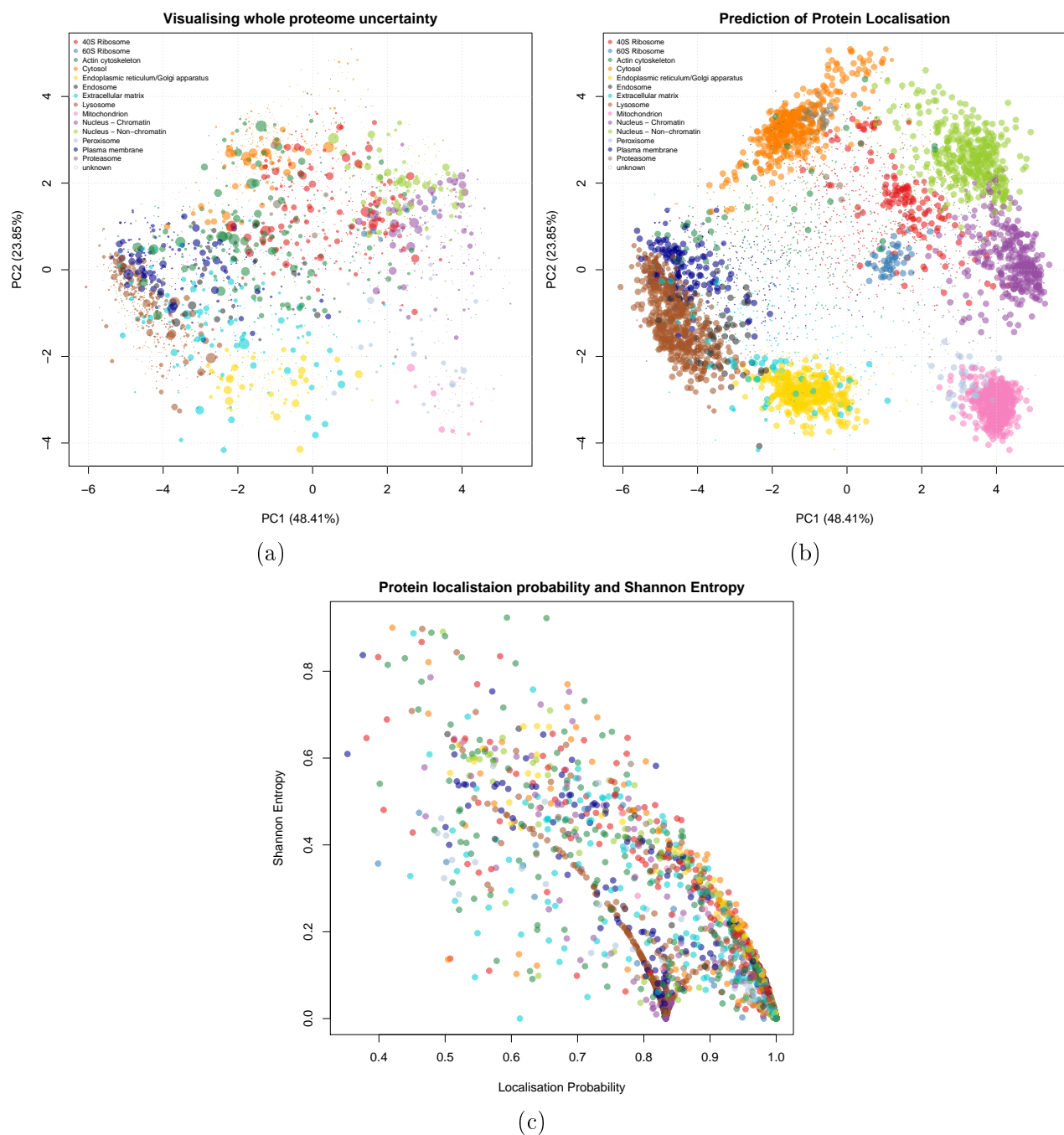


Figure 11: PCA plots of the mouse pluripotent embryonic stem cell data, where each point represents a protein and is coloured to its (probabilistically-)assigned organelle. (a) In this plot, the pointer is scaled to the Shannon entropy of this protein, with larger pointers indicating greater uncertainty. (b) In this plot, the pointer is scaled to the probability of that protein belonging to its assigned organelle. (c) We plot the localisation probabilities against the Shannon entropy with each protein.

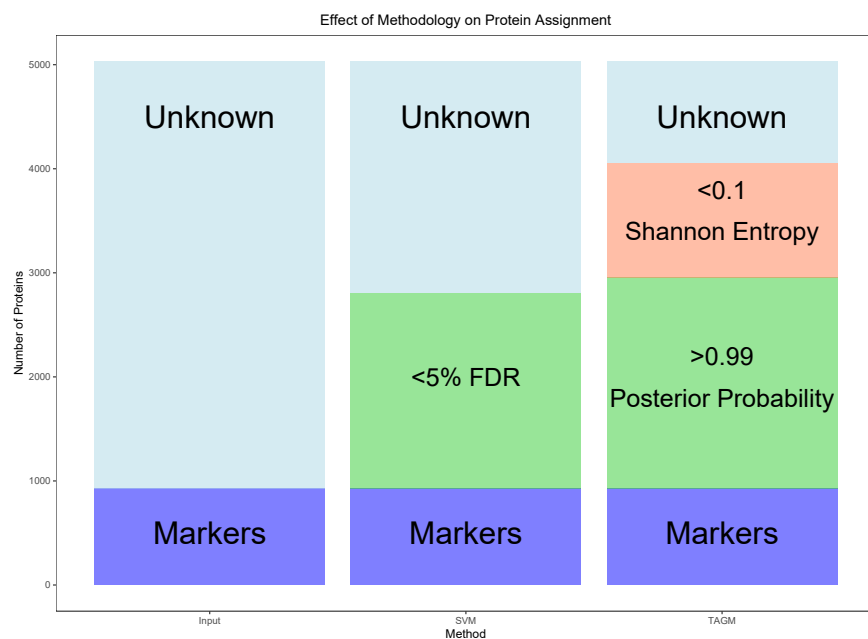


Figure 12: The barplot demonstrates the effect of applying different methodologies on protein assignment when applied the mouse pluripotent embryonic stem cell data. Roughly 2000 proteins are classified using either SVM and TAGM-MCMC; however, TAGM-MCMC can draw additional conclusions about an extra 1000 proteins by quantifying uncertainty.

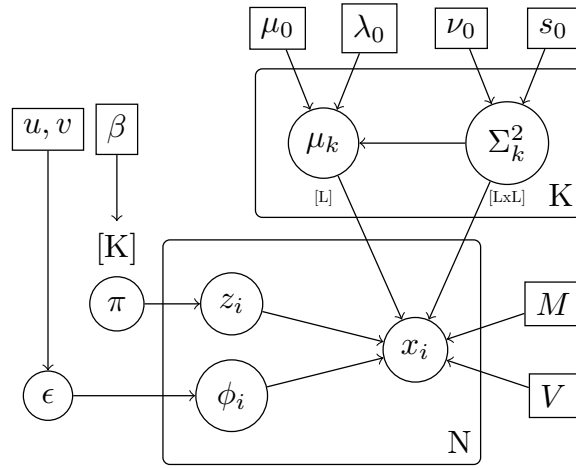


Figure 13: Plate diagram for TAGM model. This diagram specifies the conditional independencies and parameters in our model.