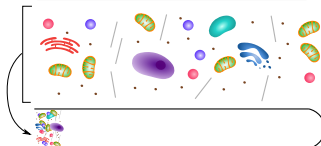
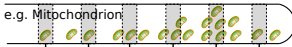


Cell lysis



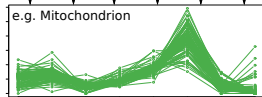
Fractionation/centrifugation

e.g. Mitochondrion



Quantitation/identification
by mass spectrometry

e.g. Mitochondrion



Quantitation data and organelle markers

	Fraction ₁	Fraction ₂	...	Fraction _m	markers
p ₁	q _{1,1}	q _{1,2}	...	q _{1, m}	unknown
p ₂	q _{2,1}	q _{2,2}	...	q _{2, m}	<i>loc₁</i>
p ₃	q _{3,1}	q _{3,2}	...	q _{3, m}	unknown
p ₄	q _{4,1}	q _{4,2}	...	q _{4, m}	<i>loc_i</i>
⋮	⋮	⋮	⋮	⋮	⋮
p _j	q _{j,1}	q _{j,2}	...	q _{j, m}	unknown

Visualisation and classification

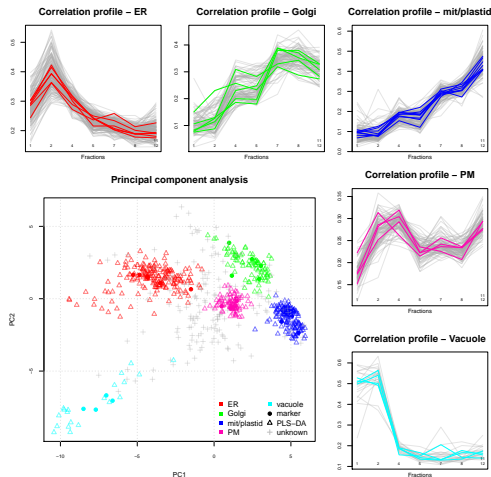


Figure : From Gatto et al. (2010), *Arabidopsis thaliana* data from Dunkley et al. (2006)

What about annotation data from repositories such as GO, sequence features, signal peptide, transmembrane domains, images, protein-protein interactions,

- ▶ From a user perspective: "**free/cheap**" vs. expensive
- ▶ 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**

number GO features \gg experimental fractions
 \Rightarrow dilution of experimental data

Goal

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

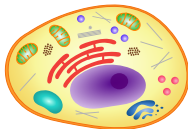
Updated experimental design for

- ▶ primary/experimental data

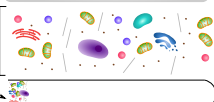
and

- ▶ auxiliary/annotation data

PRIMARY EXPERIMENTAL DATA



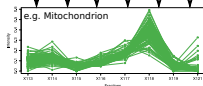
Cell lysis



Fractionation/centrifugation

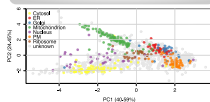


Quantitation/identification by mass spectrometry



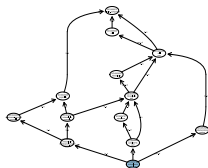
	X110	X114	X115	X116	X117	X118	X119	X121
CD327F7	0.1362	0.1350	0.1062	0.1487	0.2777	0.1429	0.0380	0.00109
PF14486	0.1014	0.1020	0.0946	0.1061	0.1207	0.0996	0.0180	0.00727
CERT3A5	0.1297	0.1201	0.0946	0.1061	0.2962	0.1463	0.0206	0.00962
GRU5C1	0.1008	0.1007	0.0919	0.1061	0.1461	0.1086	0.0002	0.00002

Visualisation



Database query

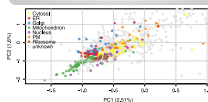
Extract GO CC terms



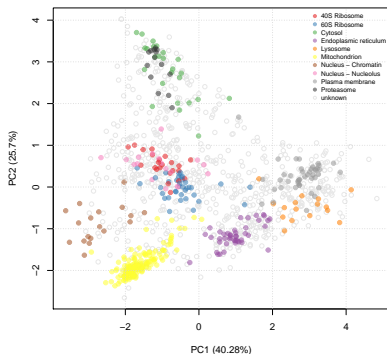
Convert terms to binary

	GO:	GO:0005822	GO:0005789	GO:0005783	GO:
CD327F7	0	1	1	1	...
PF14486	1	1	1	1	...
CERT3A5	0	0	0	0	...
GRU5C1	0	0	0	0	...

Visualisation



AUXILIARY DRY DATA



Data from mouse stem cells (E14TG2a)

We use a **class-weighted** kNN transfer learning algorithm to combine primary and auxiliary data, based on Wu and Dietterich (2004):

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

Classes and weights

$$\mathbb{C} = \{c_{i=1}, \dots, c_{i=l}\}; \Theta = \{0, 0.5, 1\}$$

Primary data

$$L_P = \begin{bmatrix} q_{1,1} & q_{1,2} & \dots & q_{1,m} \\ q_{2,1} & q_{2,2} & \dots & q_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ q_{j,1} & q_{j,2} & \dots & q_{j,m} \end{bmatrix}; \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_j \end{bmatrix}; k_P$$

Auxiliary data

$$L_A = \begin{bmatrix} b_{1,1} & b_{1,2} & \dots & \dots & b_{1,n} \\ b_{2,1} & b_{2,2} & \dots & \dots & b_{2,n} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ b_{j,1} & b_{j,2} & \dots & \dots & b_{j,n} \end{bmatrix}; \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_j \end{bmatrix}; k_A$$

Neighbour matrices

$$N_P = \begin{bmatrix} c_{i=1} & \dots & c_{i=l} \\ n_{1,1}^P & \dots & n_{1,l}^P \\ n_{2,1}^P & \dots & n_{2,l}^P \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \end{bmatrix}; N_A = \begin{bmatrix} c_{i=1} & \dots & c_{i=l} \\ n_{1,1}^A & \dots & n_{1,l}^A \\ n_{2,1}^A & \dots & n_{2,l}^A \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \end{bmatrix}$$

Classes and weights

$$\mathbb{C} = \{c_{i=1}, \dots, c_{i=l}\}; \Theta = \{0, 0.5, 1\}$$

Primary data

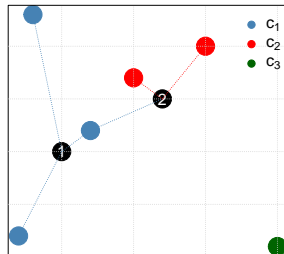
$$L_P = \begin{bmatrix} q_{1,1} & q_{1,2} & \dots & q_{1,m} \\ q_{2,1} & q_{2,2} & \dots & q_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ q_{j,1} & q_{j,2} & \dots & q_{j,m} \end{bmatrix}; \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_j \end{bmatrix}; k_P$$

Auxiliary data

$$L_A = \begin{bmatrix} b_{1,1} & b_{1,2} & \dots & \dots & b_{1,n} \\ b_{2,1} & b_{2,2} & \dots & \dots & b_{2,n} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ b_{j,1} & b_{j,2} & \dots & \dots & b_{j,n} \end{bmatrix}; \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_j \end{bmatrix}; k_A$$

Neighbour matrices

$$N_P = \begin{bmatrix} c_{i=1} & \dots & c_{i=l} \\ n_{1,1}^P & \dots & n_{1,l}^P \\ n_{2,1}^P & \dots & n_{2,l}^P \\ \vdots & \vdots & \vdots \end{bmatrix}; N_A = \begin{bmatrix} c_{i=1} & \dots & c_{i=l} \\ n_{1,1}^A & \dots & n_{1,l}^A \\ n_{2,1}^A & \dots & n_{2,l}^A \\ \vdots & \vdots & \vdots \end{bmatrix}$$



$$N_P = \begin{matrix} & c_1 & c_2 & c_3 \\ \begin{matrix} p_1 \\ p_2 \end{matrix} & \begin{bmatrix} \frac{3}{3} & 0 & 0 \\ \frac{1}{3} & \frac{2}{3} & 0 \\ \vdots & \vdots & \vdots \end{bmatrix} \end{matrix}$$

Classes and weights

$$\mathbb{C} = \{c_{i=1}, \dots, c_{i=I}\}; \Theta = \{0, 0.5, 1\}$$

Primary data

$$L_P = \begin{bmatrix} q_{1,1} & q_{1,2} & \dots & q_{1,m} \\ q_{2,1} & q_{2,2} & \dots & q_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ q_{j,1} & q_{j,2} & \dots & q_{j,m} \end{bmatrix}; \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_j \end{bmatrix}; k_P$$

Auxiliary data

$$L_A = \begin{bmatrix} b_{1,1} & b_{1,2} & \dots & \dots & b_{1,n} \\ b_{2,1} & b_{2,2} & \dots & \dots & b_{2,n} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ b_{j,1} & b_{j,2} & \dots & \dots & b_{j,n} \end{bmatrix}; \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_j \end{bmatrix}; k_A$$

Neighbour matrices

$$N_P = \begin{bmatrix} c_{i=1} & \dots & c_{i=I} \\ n_{1,1}^P & \dots & n_{1,I}^P \\ n_{2,1}^P & \dots & n_{2,I}^P \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \end{bmatrix}; N_A = \begin{bmatrix} c_{i=1} & \dots & c_{i=I} \\ n_{1,1}^A & \dots & n_{1,I}^A \\ n_{2,1}^A & \dots & n_{2,I}^A \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \end{bmatrix}$$

Weights matrix (labelled)

$$\begin{matrix} & c_1 & c_2 & c_3 \\ \theta_1 & \begin{bmatrix} 0 & 0 & 0 \end{bmatrix} \\ \theta_2 & \begin{bmatrix} 0 & 0 & 1 \end{bmatrix} \\ \theta_i & \begin{bmatrix} \vdots & & \vdots \end{bmatrix} \\ \vdots & \begin{bmatrix} 1 & 1 & 0 \end{bmatrix} \\ \theta_{\Theta^I} & \begin{bmatrix} 1 & 1 & 1 \end{bmatrix} \end{matrix} \begin{bmatrix} F_{1_1} \\ F_{1_2} \\ F_{1_i} \\ \vdots \\ F_{1_{\Theta^I}} \end{bmatrix}$$

$$\theta^* = \{1, 0, 1\}$$

(♥ BiocParallel)

$$\mathbb{C} = \{c_{j=1}, \dots, c_{j=J}\}; \Theta = \{0, 0.5, 1\}$$
$$L_P = \begin{bmatrix} q_{1,1} & q_{1,2} & \dots & q_{1,m} \\ q_{2,1} & q_{2,2} & \dots & q_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ q_{i,1} & q_{i,2} & \dots & q_{i,m} \end{bmatrix}; \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_i \end{bmatrix}; k_P$$
$$L_A = \begin{bmatrix} b_{1,1} & b_{1,2} & \dots & \dots & b_{1,n} \\ b_{2,1} & b_{2,2} & \dots & \dots & b_{2,n} \\ \vdots & \vdots & & & \vdots \\ b_{i,1} & b_{i,2} & \dots & \dots & b_{i,n} \end{bmatrix}; \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_j \end{bmatrix}; k_A$$
$$N_P = \begin{bmatrix} c_{i=1} & \dots & c_{i=l} \\ n_{1,1}^P & \dots & n_{1,l}^P \\ n_{2,1}^P & \dots & n_{2,l}^P \\ \vdots & & \vdots \\ \vdots & & \vdots \end{bmatrix}; N_A = \begin{bmatrix} c_{i=1} & \dots & c_{i=l} \\ n_{1,1}^A & \dots & n_{1,l}^A \\ n_{2,1}^A & \dots & n_{2,l}^A \\ \vdots & & \vdots \\ \vdots & & \vdots \end{bmatrix}$$
$$V(c_i)_j = \theta^* n_{ij}^P + (1 - \theta^*) n_{ij}^A$$

$$V(c_i) = \begin{bmatrix} 1 \\ 2 \\ 3 \\ \vdots \\ j \end{bmatrix}$$

$$y_i = \operatorname{argmax}(V(c_i)_i)$$

Class-weighted classifier (unlabelled)

$$\theta^* = \{1, 0, 1\} \quad N_P = \begin{matrix} & c_1 & c_2 & c_3 \\ p_1 & \frac{3}{3} & 0 & 0 \\ p_2 & \frac{1}{3} & \frac{2}{3} & 0 \\ & \vdots & \vdots & \vdots \end{matrix}$$

$$V(c_1)_1 = 1 \times \frac{3}{3} + (1 - 1) \times n_{1,1}^A$$

$$V(c_2)_1 = 0 \times 0 + (1 - 0) \times n_{1,2}^A$$

$$V(c_3)_1 = 1 \times 0 + (1 - 1) \times n_{1,3}^A$$

$$V(c_1)_2 = 1 \times \frac{1}{3} + (1 - 1) \times n_{1,1}^A$$

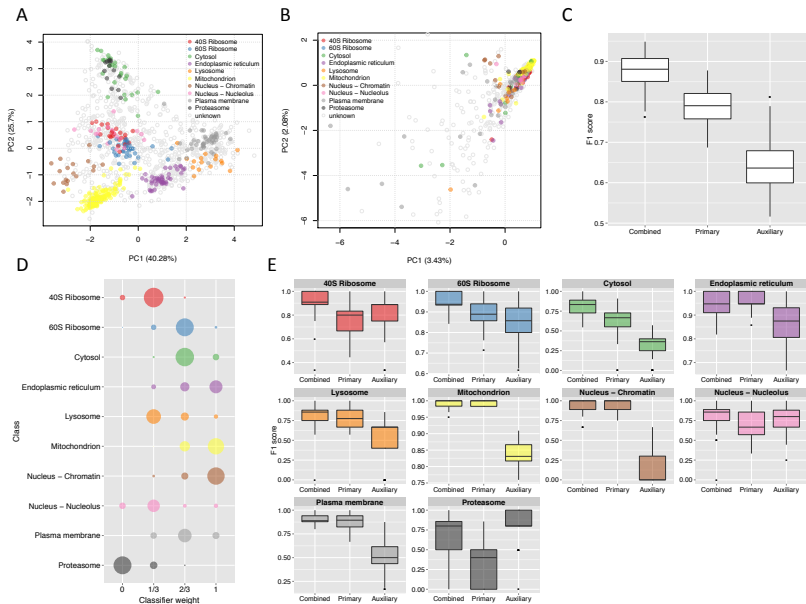
$$V(c_2)_2 = 0 \times \frac{2}{3} + (1 - 0) \times n_{1,2}^A$$

$$V(c_3)_2 = 1 \times 0 + (1 - 1) \times n_{1,3}^A$$

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

$$\begin{matrix} & c_1 & c_2 & c_3 \\ 1 & V(c_1)_1 & V(c_2)_1 & V(c_3)_1 \\ 2 & V(c_1)_2 & V(c_2)_2 & V(c_3)_2 \\ \vdots & & \vdots & \\ j & & & \end{matrix}$$

$$y_j = \operatorname{argmax}(V(c_i)_j)$$



Data from mouse stem cells (E14TG2a).