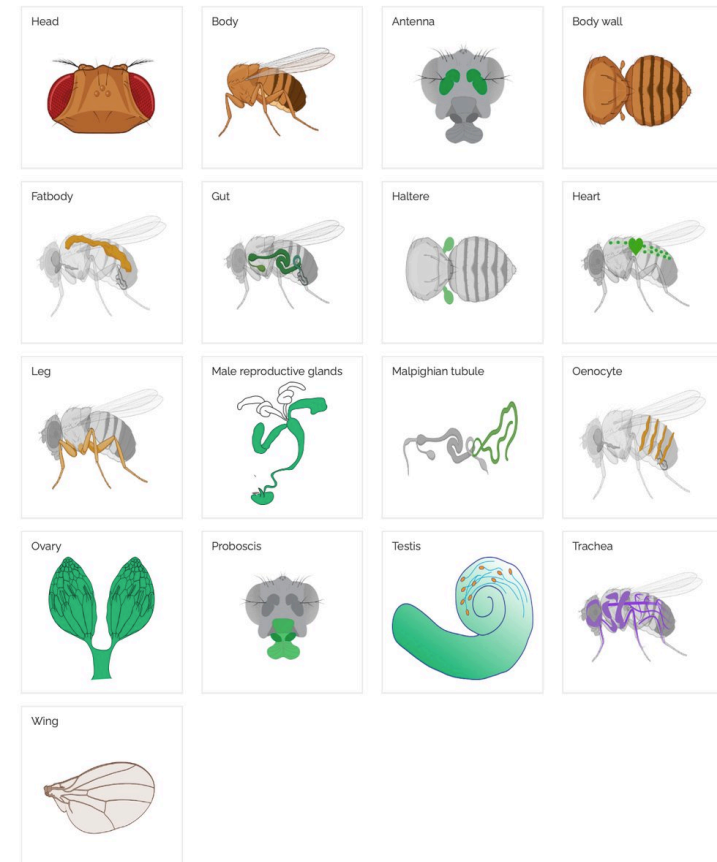
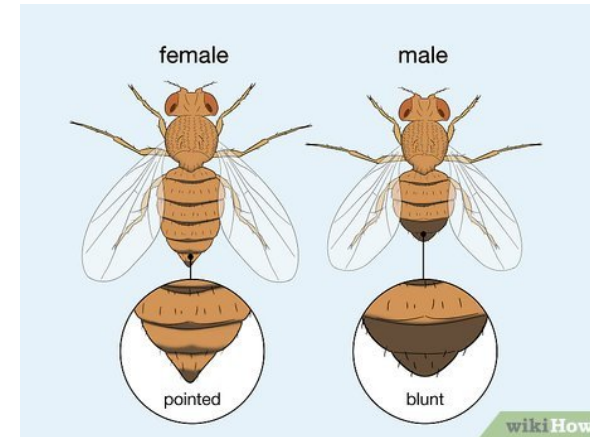


Information theory-based
analyses to find a **minimal non-
redundant** gene set for sex label
classification

Ming
2021

Question:

- Sex dimorphism at **phenotype** level
- Sex dimorphism at transcription, or **gene expression** level
- Is there a molecular readout of sex dimorphism?
- What are the biological processes/gene expression programs associated with gene expression difference between the two sexes?
- In somatic tissues, does every cell have a 'sex identity'?
- Are sex identities of different types of cells mediated/manifested by the same set of genes?



Is the sex differentiation potential be realized?

Turn a biological question into a computational one:

- If there is a gene set most informative about the distinct between female and male cells, does this gene set consistent across diverse cell types in a fly?

How to find a minimal non-redundant
gene set to predict the sex label of every
individual cell

sex labels unknown

	Cell 1	Cell2	Cell3	Cell4	...
Gene 1	0	0	3	10	
Gene 2	24	0	41	12	
Gene 3	175	284	93	162	
Gene 4	0	0	0	0	
Gene 5	36	0	32	21	
...	

Gene set

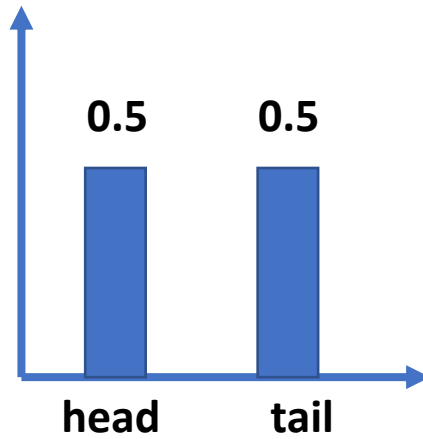
Female male female male

	Cell 1	Cell2	Cell3	Cell4	...
Gene 1	0	0	3	10	
Gene 2	24	0	41	12	
Gene 3	175	284	93	162	
Gene 4	0	0	0	0	
Gene 5	36	0	32	21	
...	

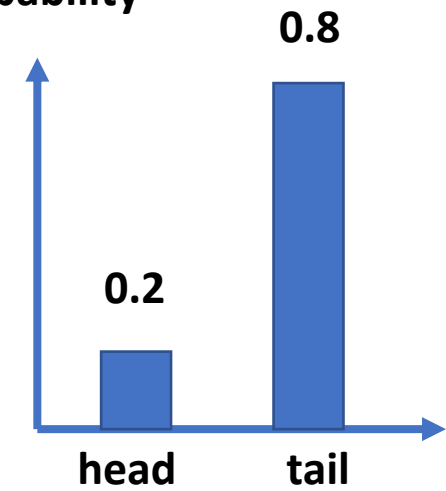
Background

Flip a coin, guess head or tail
which one is more uncertain?

Probability



Probability

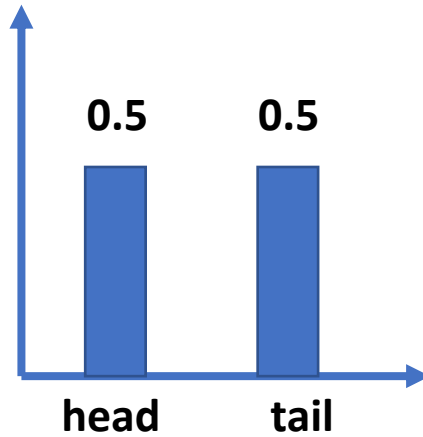


Background

- Information entropy, a measure of uncertainty

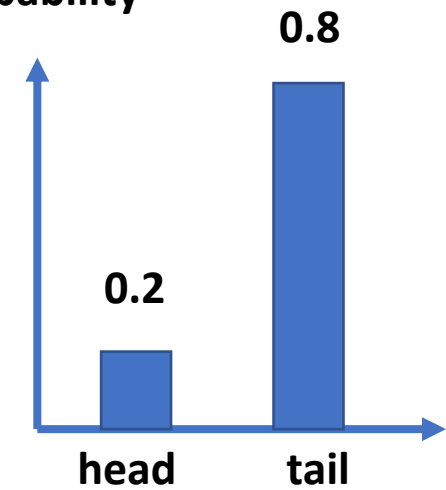
$$H(X) = - \sum_{i=1}^n P(x_i) \log_b P(x_i)$$

Probability



Flip a coin, guess head or tail
which one is more uncertain?

Probability



$$-1 * \{0.5 * \log(0.5) + 0.5 * \log(0.5)\} = \mathbf{0.6931}$$

$$-1 * \{0.2 * \log(0.2) + 0.8 * \log(0.8)\} = \mathbf{0.5004}$$

Background

- **Information entropy**, the uncertainty of a variable's possible outcomes.
Larger value ~ higher uncertainty

$$H(X) = - \sum_{i=1}^n P(x_i) \log_b P(x_i)$$

- **Conditional entropy**, given a second variable (X), what's the amount of the uncertainty still existing in the first variable (Y).

$$H(Y|X = x) = - \sum_{y \in \mathcal{Y}} \Pr(Y = y|X = x) \log_2 \Pr(Y = y|X = x).$$

Y	1	1	1	1	1	0	0	0	0	0
---	---	---	---	---	---	---	---	---	---	---

$$H(Y) = -1 * \{ 0.5 * \log(0.5) + 0.5 * \log(0.5) \} = 0.6931$$

Y	1	1	1	1	1	0	0	0	0	0
X	1	1	1	0	0	0	0	0	0	1

$$H(Y|X) = \Pr(X=0) * H(Y|X=0) + \Pr(X=1) * H(Y|X=1) = 0.6068$$

Mutual information

How much uncertainty about Y decreases when we observe X.

Or the amount of information **gained** knowing X to predict Y

$$I(Y;X) = H(Y) - H(Y|X) = 0.6931 - 0.6068 = 0.0863$$

Sex label

1	1	1	1	1	0	0	0	0	0
---	---	---	---	---	---	---	---	---	---

$$H(Y) = -1 * \{ 0.5 * \log(0.5) + 0.5 * \log(0.5) \} = 0.6931$$

Sex label
Binary gene
expression

1	1	1	1	1	0	0	0	0	0
1	1	1	0	0	0	0	0	0	1

$$H(Y|X) = \Pr(X=0) * H(Y|X=0) + \Pr(X=1) * H(Y|X=1) \\ = 0.6068$$

Mutual information

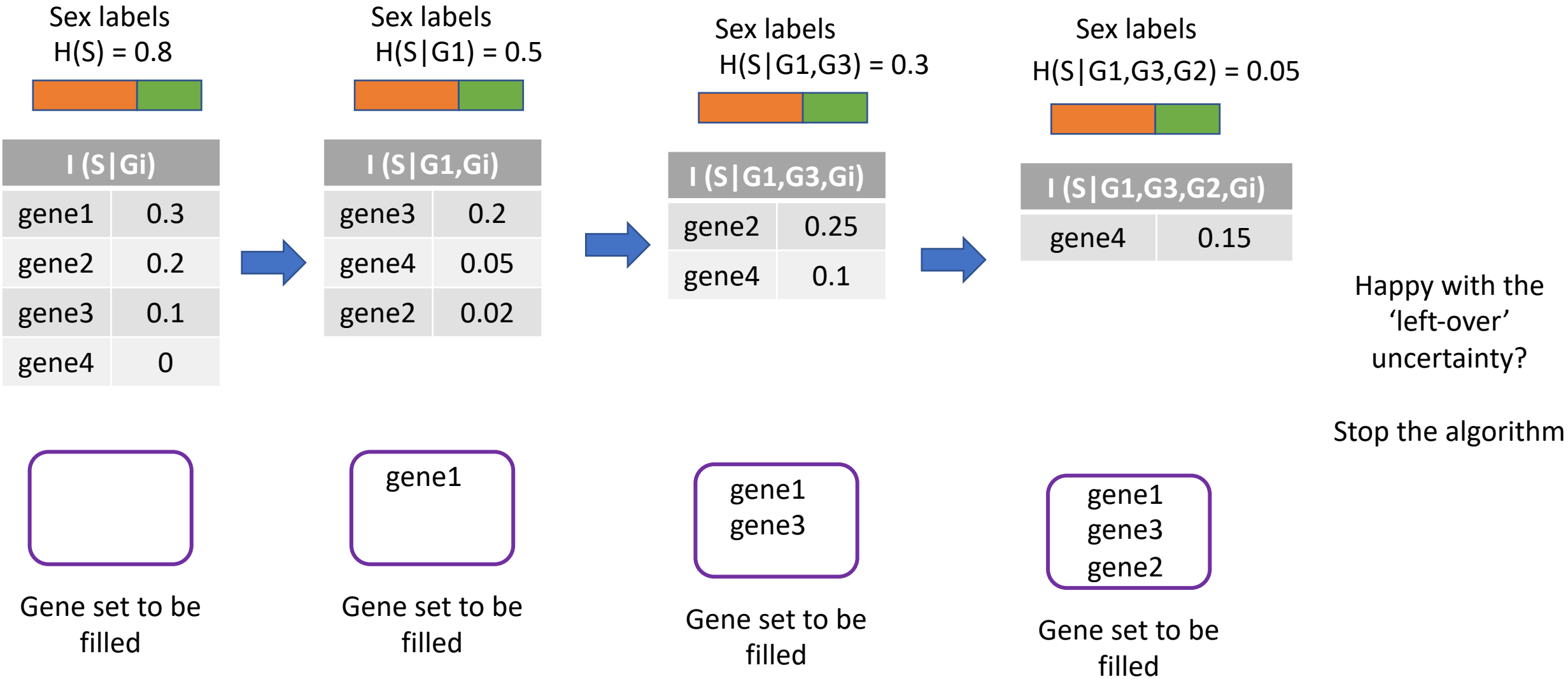
How much uncertainty about Y decreases when we observe X.

Or the amount of information **gained** knowing X to predict Y

$$I(Y;X) = H(Y) - H(Y|X) = 0.6931 - 0.6068 = 0.0863$$

Algorithm

H(.) information entropy, or uncertainty
I(.) mutual information, or information gained given knowing something



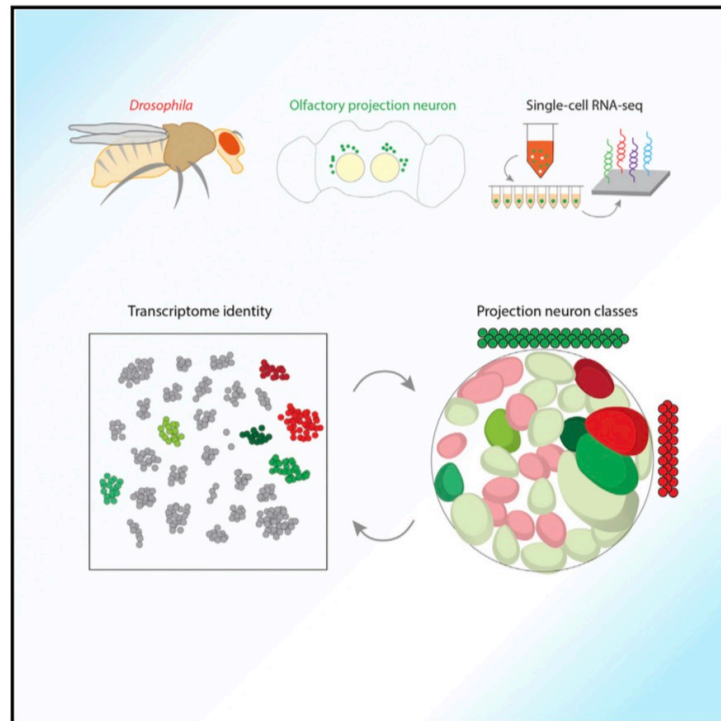
- This algorithm was proposed in this following paper to classify neuron types
- I implemented it in R to derive a gene signature of sex labels of single cells

Resource

Cell

Classifying *Drosophila* Olfactory Projection Neuron Subtypes by Single-Cell RNA Sequencing

Graphical Abstract



Authors

Hongjie Li, Felix Horns, Bing Wu, ...,
David J. Luginbuhl, Stephen R. Quake,
Liqun Luo

Correspondence

quake@stanford.edu (S.R.Q.),
lluo@stanford.edu (L.L.)

In Brief

Single-cell RNA sequencing establishes that transcriptomic identity of *Drosophila* olfactory projection neurons corresponds with connectivity and function and identifies transcription factors and cell-surface molecules as highly informative in encoding neuronal identity.

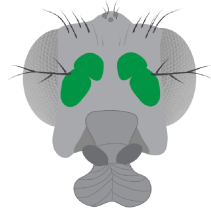
Dataset

- For each cell type in a tissue, I took the top **100** most informative genes as the input gene list, to initiate the 'minimal non-redundant gene set' searching journey.
- Cell type considered: ≥ 20 cells in female and ≥ 20 cells in male
- For each cell type in a tissue, if I want a gene set that could explain 0.99 uncertainty in the classification of sex labels, the resulting gene set would have a certain size, for example, 2 genes, 3 genes, etc.

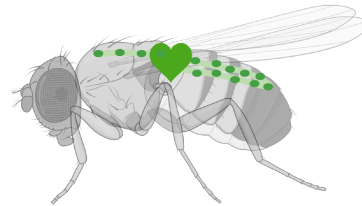
head



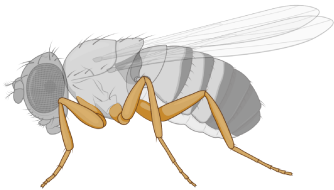
antenna



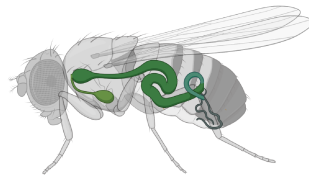
heart



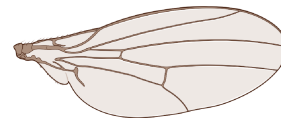
leg



gut



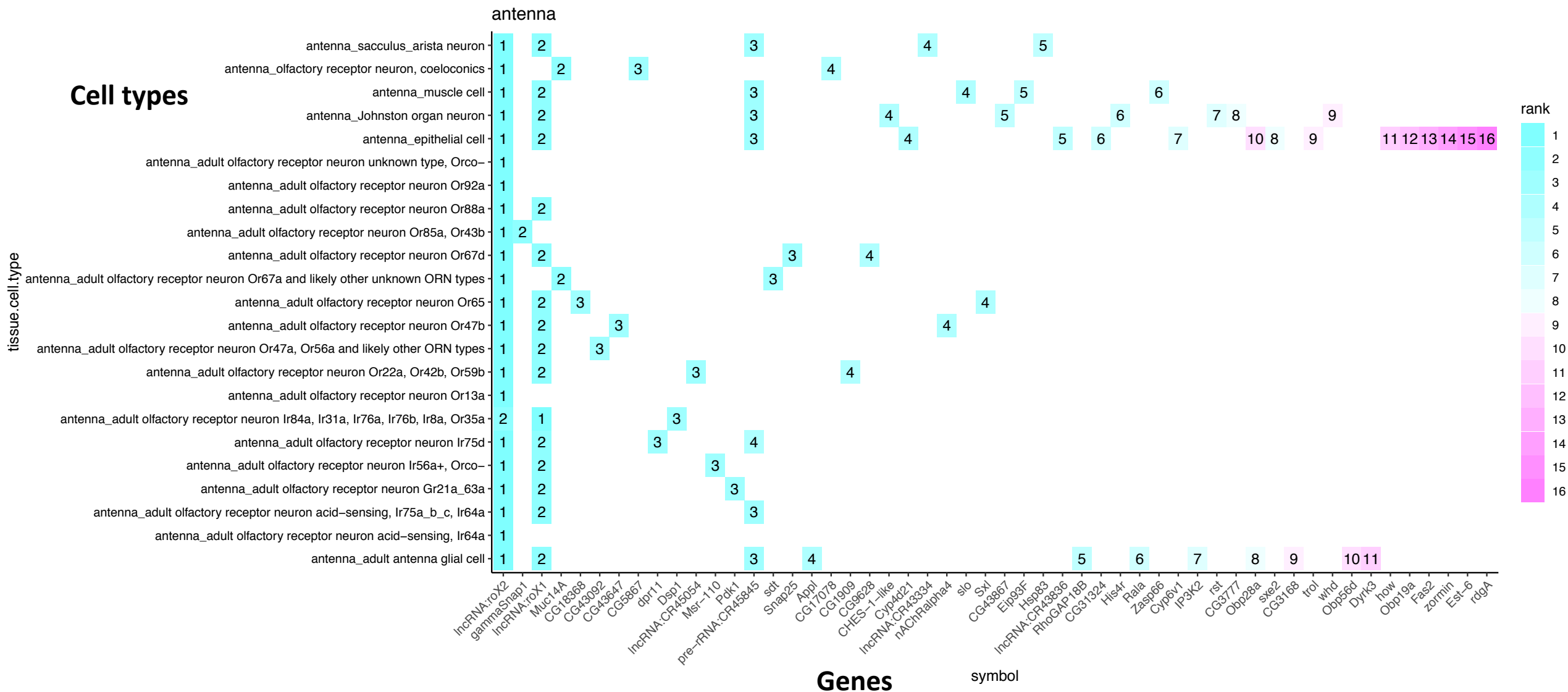
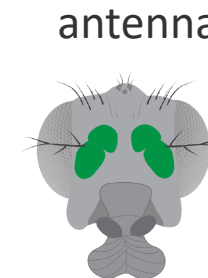
wing



Result

A gene set that could explain **0.99** uncertainty in the classification of sex labels for each cell type

23 cell types in the fly antenna

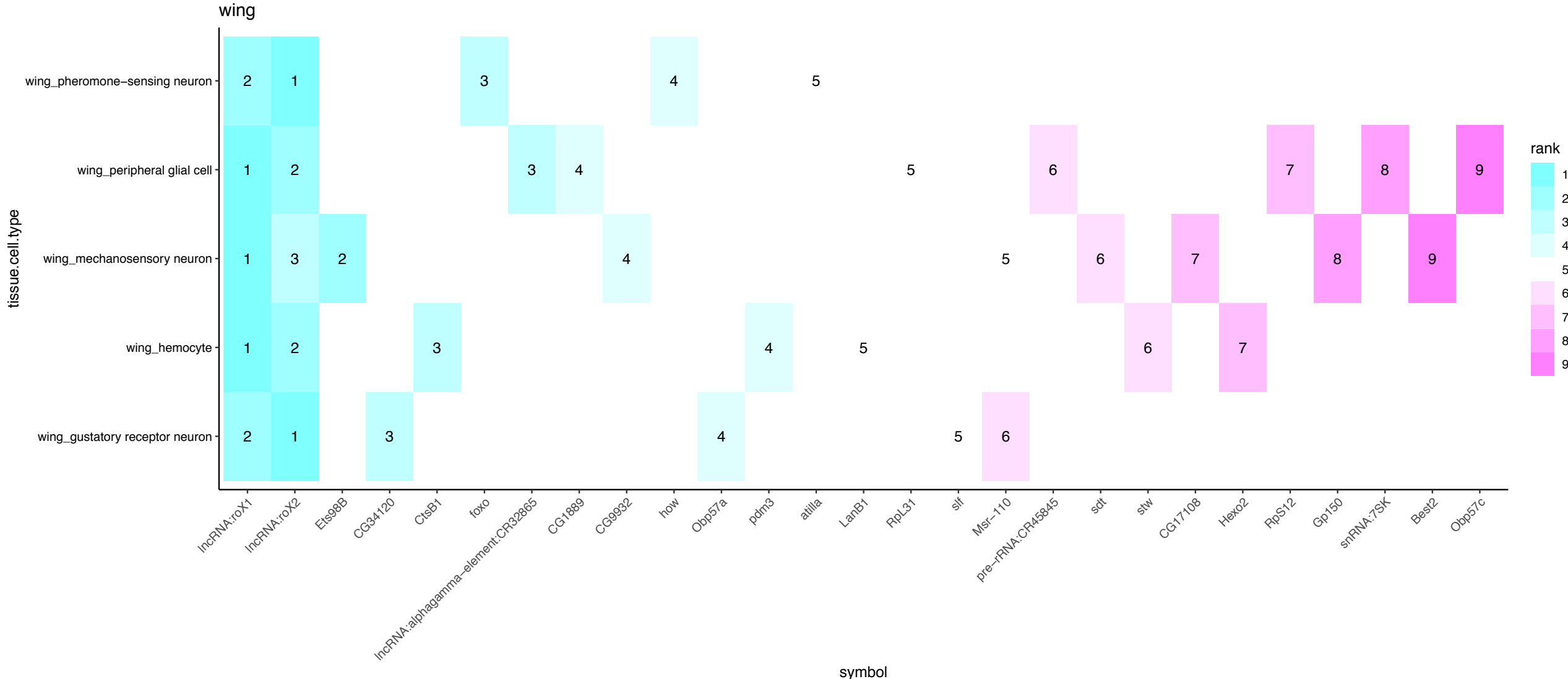
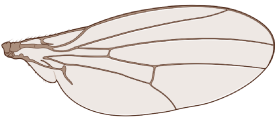


5 cell types

female male

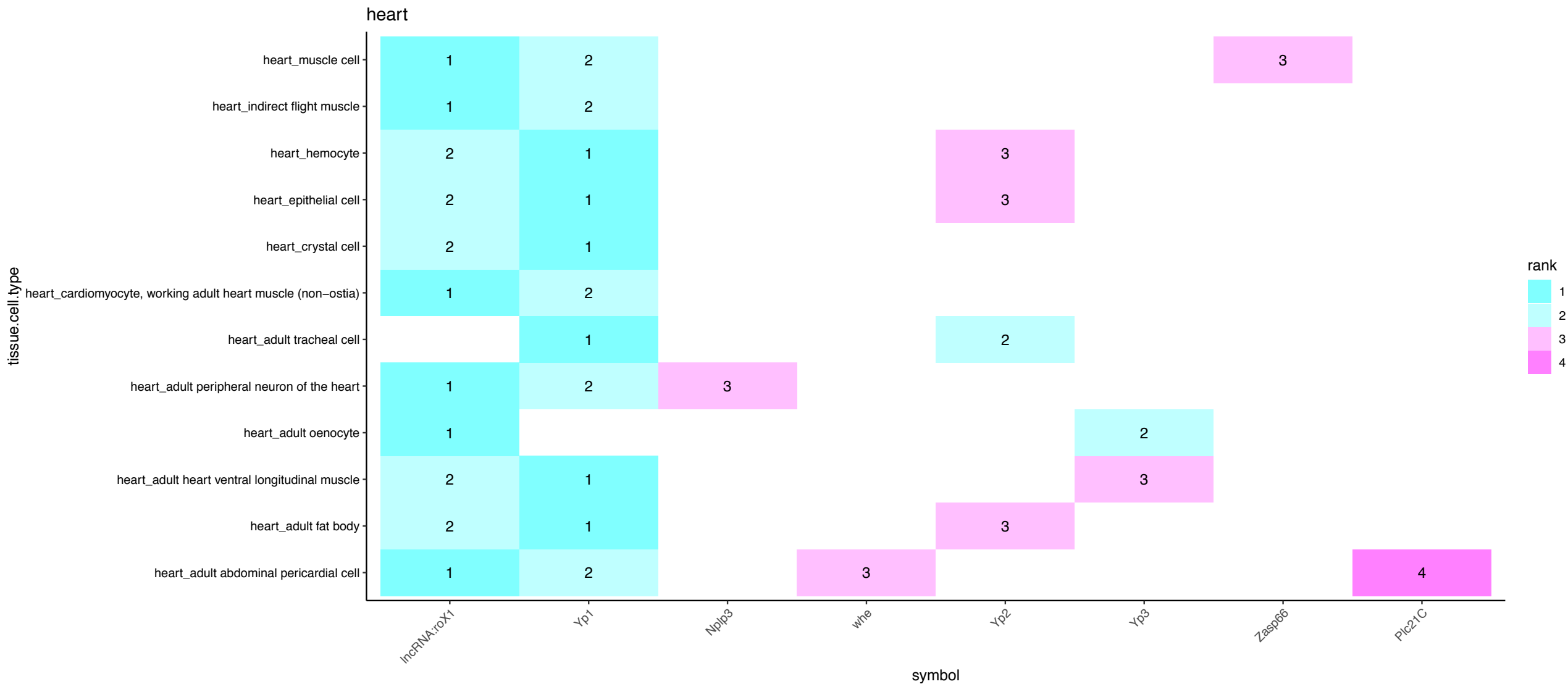
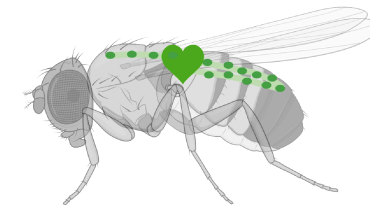
Muscle cells 19 125
nociceptive neuron 9 30

wing



12 cell types

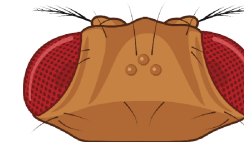
heart



A detailed illustration of a housefly (Musca domestica) from a side profile. The fly has a grey, segmented body with a darker, almost black, thorax and abdomen. Its head is large and rounded with prominent eyes. It has two large, transparent wings with visible veins. The legs are thick and yellowish-brown. The fly is shown in a slightly angled position, facing left.

[illegible]

head



72 cell types

