# Gene Set Enrichment Analysis (GSEA)

Nicolas Delhomme, Bastian Schiffthaler

# What to do with your list of genes?

- Mostly, you will select a few candidates and do functional characterisation in the lab, right?

- So you went ballistic and did a whole genomic analysis and you dare look at only a couple genes?

# There must be another way

- Well, sort of. GSEA is one.

- GSEA aims at answering the question: is my list of genes (the gene set) associated with experimental condition

  - e.g. are there unusually many de-regulated genes in my gene list

  - e.g. is my DE gene list enriched for some functional processes

# GSEA methods

- Reviewed in Kharti et al., 2012

  - Over-representation analysis (ORA) – are differentially expressed (DE) genes in the set more common than expected?

  - Functional class scoring (FCS) – summarize statistic of DE of genes in a set, and compare to null

  - Pathway topology (PT) – include pathway knowledge in assessing DE of genes in a set

# GSEA you know?

- Gene Ontology Annotation

- KEGG

- reactome

- ...

# GSEA methods (1)

- Competitive methods (Goeman and Bühlmann, 2007): depend on a competitive null hypothesis which assumes the genes in a set do not have a stronger association with the experimental condition compared to randomly chosen genes outside the set.
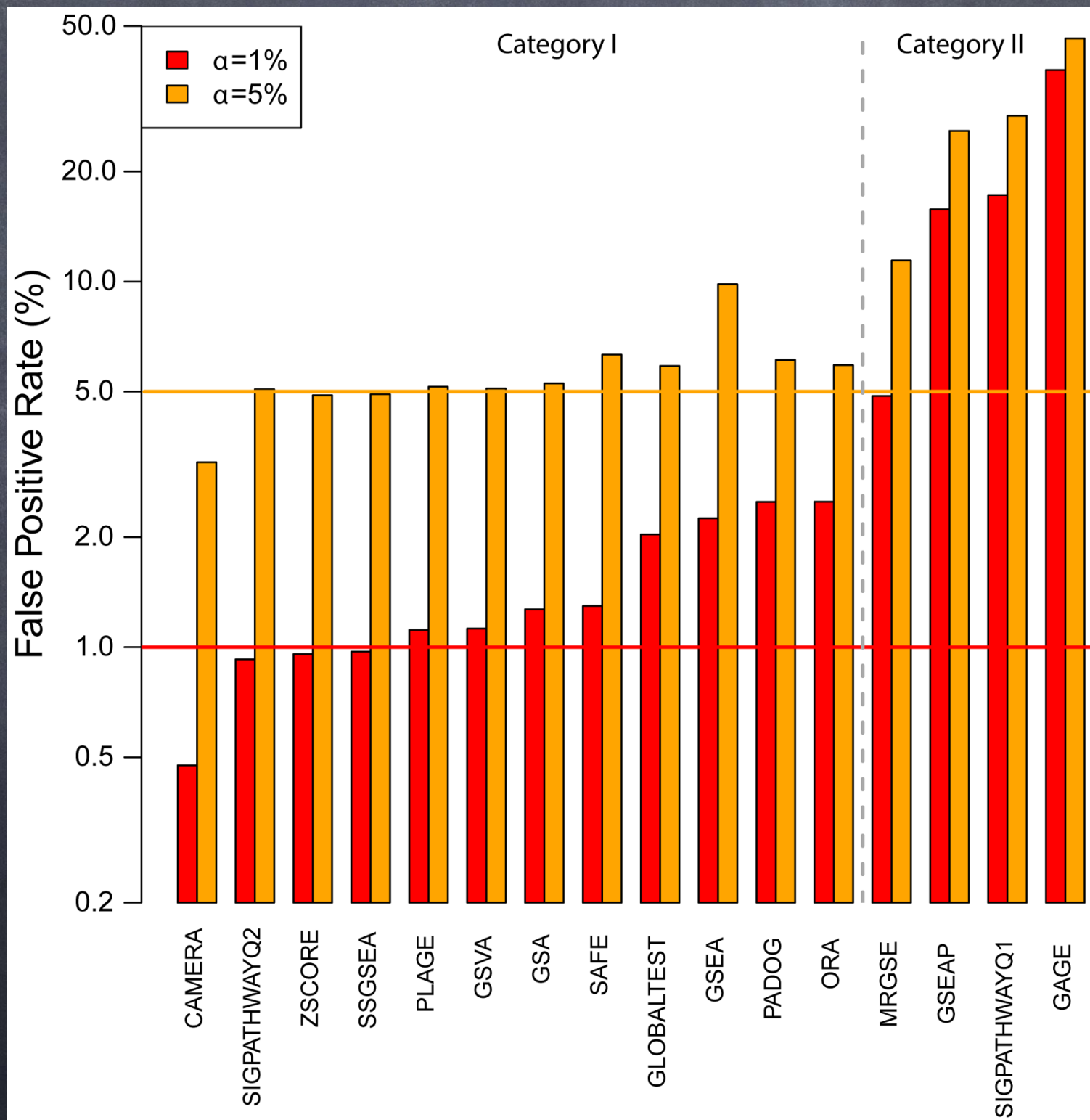
- ORA, GAGE, Camera, GSVA, ...

# GSEA methods (2)

- Self-contained methods: null hypothesis that only considers genes within a set and again assumes that they have no association with the experimental condition

- SAFE, ZSCORE, SSGSEA, ...
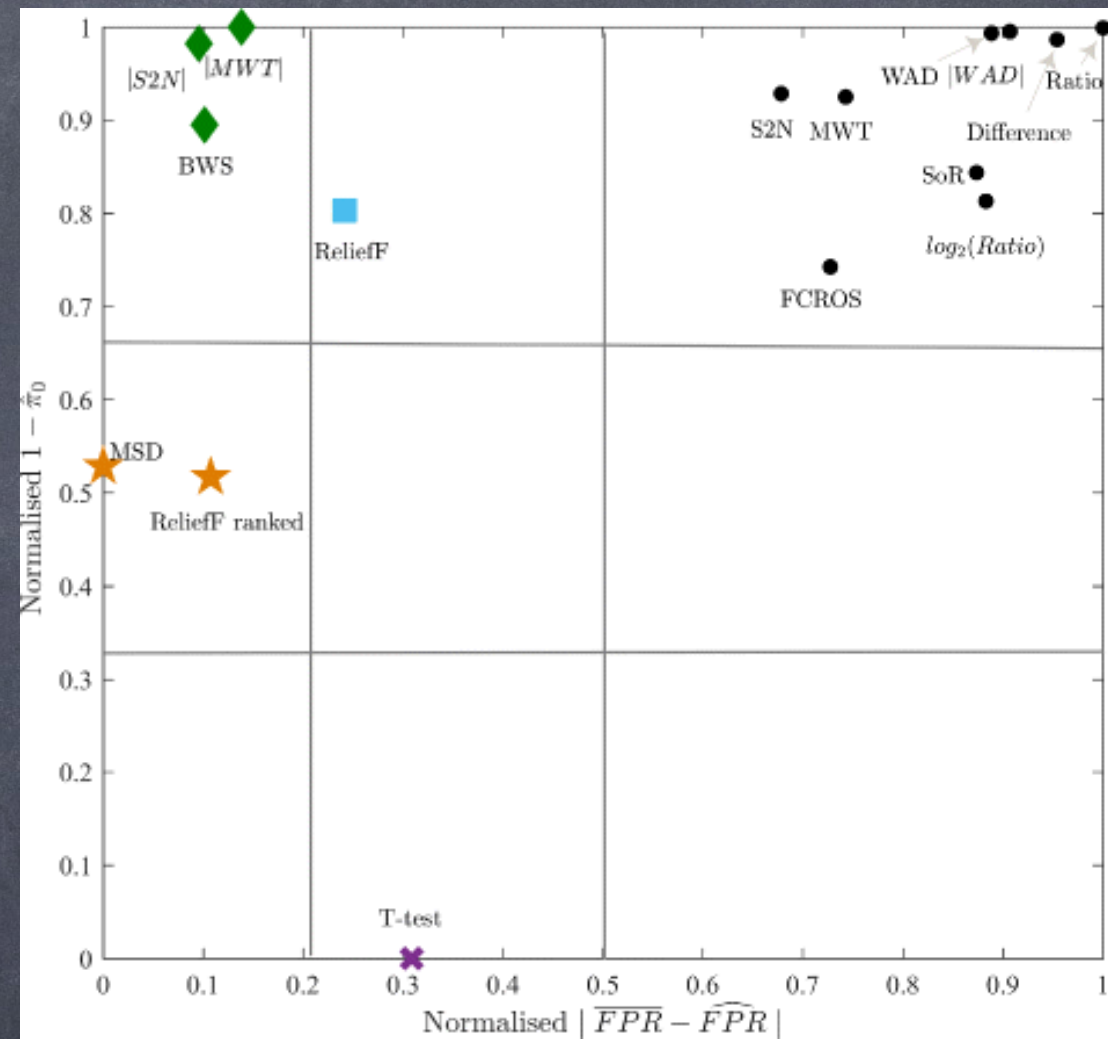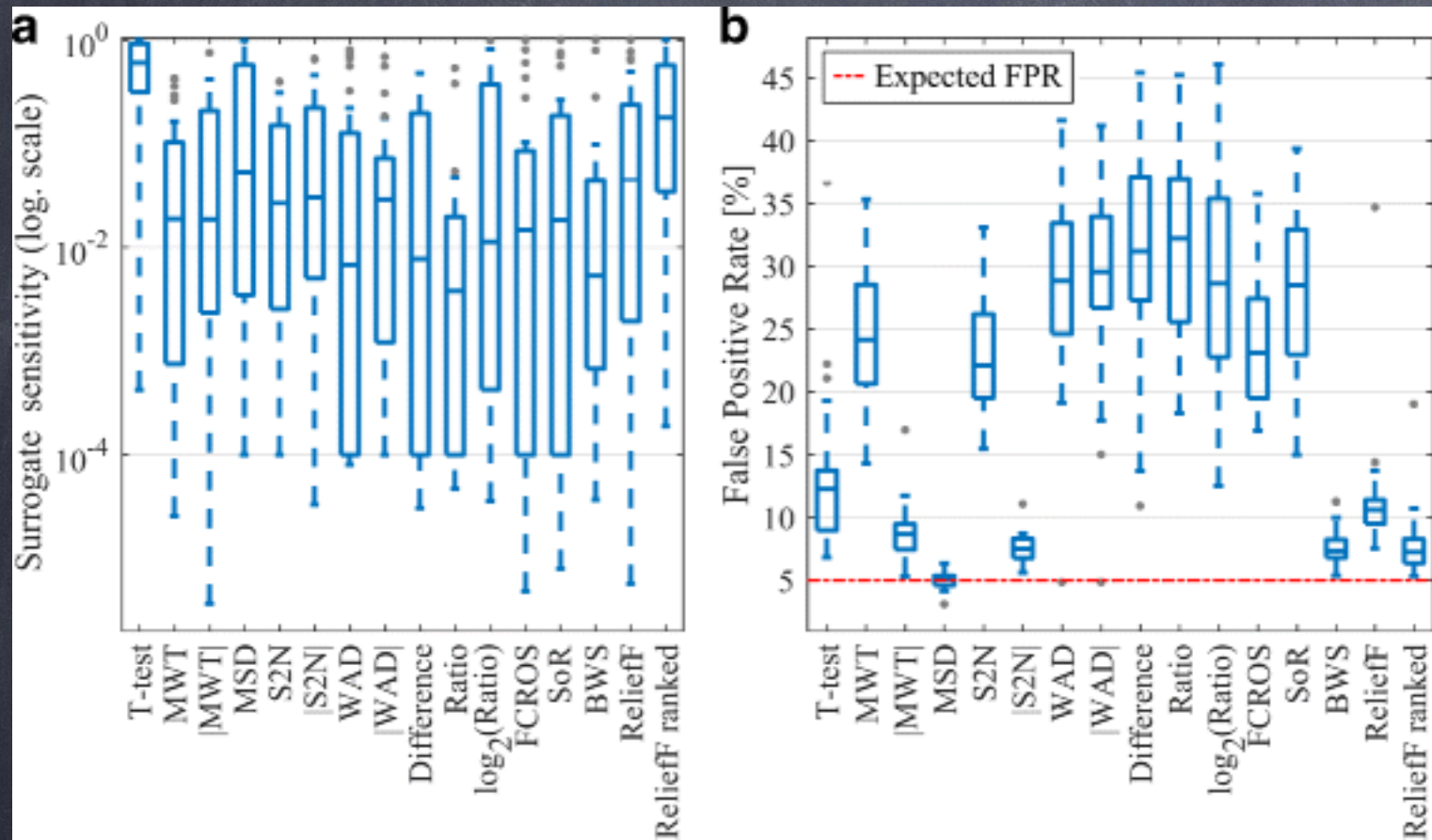
# What to pick?

- Maciejewski, 2013

  - sample randomisation based methods are better

  - popular methods (GSEA, SAFE) do not strictly test the competitive null hypothesis: A significant result from these methods does not necessarily mean that the gene set of interest contains more genes associated with the phenotype than its complement, but it rather means that either the gene set or its complement are associated with the phenotype.
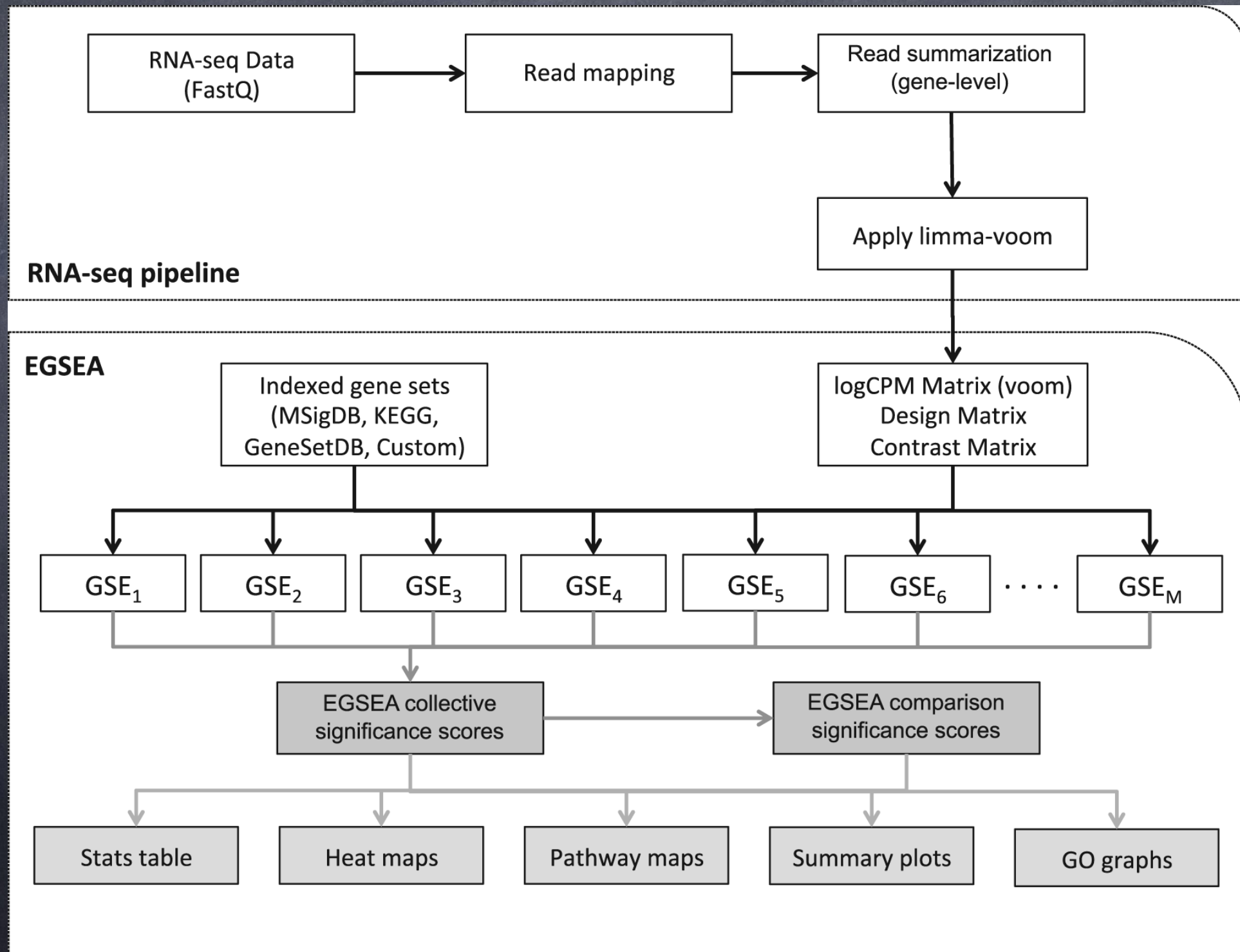
Tarca et al., 2013

# Zyla et al., 2017

# Alhamdoosh et al., 2017
# Ensemble GSEA

**A**

IL-13 vs. Control

IL-13R Antagonist vs. IL13

Comparison Summary Plot

**B**

hsa05310 Asthma

hsa05134 Viral Myocarditis

hsa04060 Cyt.-Cyt. Recept. Inter.

Color key and histogram of gene fold changes

hsa05310

hsa05134

hsa04060

Significance of DE

FDR <= 0.05 for at least one

FDR > 0.05 for all contrasts

# Let's give it a shot

- http://bioconductor.org/help/workflows/EGSEA123/

- This is still an unrefined workflow (to say the least, but it does look promising)

- Open the R script EGSEAWorkflow.Rmd in RStudio (you copied it earlier)

# References

- Some slides inspired from Martin Morgan (Bioconductor) – https://www.bioconductor.org/help/course-materials/2015/CSAMA2015/lect/L16a-gene-set-enrichment-theory-morgan.pdf

- Alhamdoosh et al, 2017, Bioinformatics, 10.1093/bioinformatics/btw623

- Glass and Girvan, 2014, Scientific Reports, 10.1038/srep04191

- Goeman & Bühlmann, 2007, Bioinformatics 23.8: 980-987.

- Grossmann et al., 2007, Bioinformatics, 10.1093/bioinformatics/btm440

- Khatri et al., 2012, PLoS Comp Biol 8.2: e1002375.

- Maciejewski, 2013, Brief. in Bioinf., 10.1093/bib/bbt002

- Tarca et al, 2013. PLOS ONE, 10.1371/journal.pone.0079217

- Zyle et al., 2017, BMC Bioinformatics 10.1186/s12859-017-1674-0