

# Automatic identification of experimental conditions relevant to a specific trait

Stavros Makrodimitris

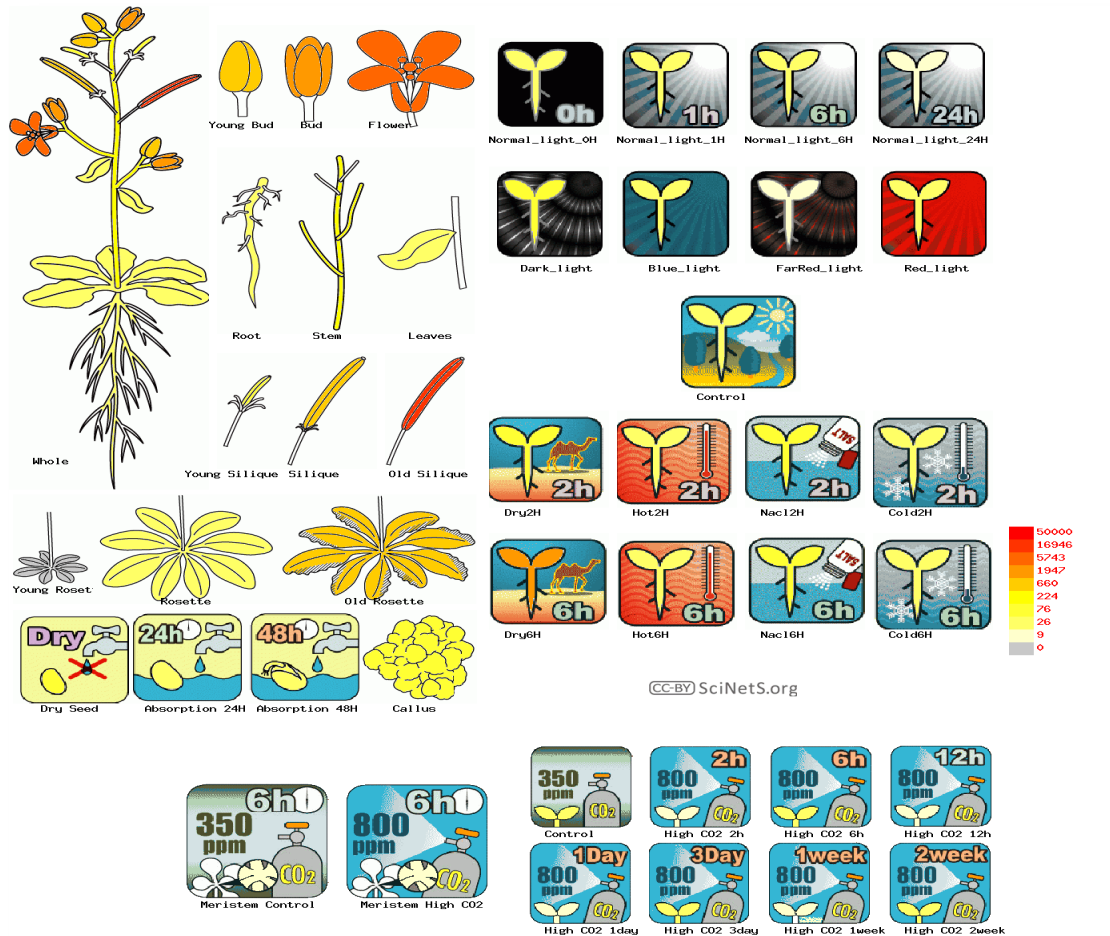
[s.makrodimitris@tudelft.nl](mailto:s.makrodimitris@tudelft.nl)

# Candidate Gene Discovery

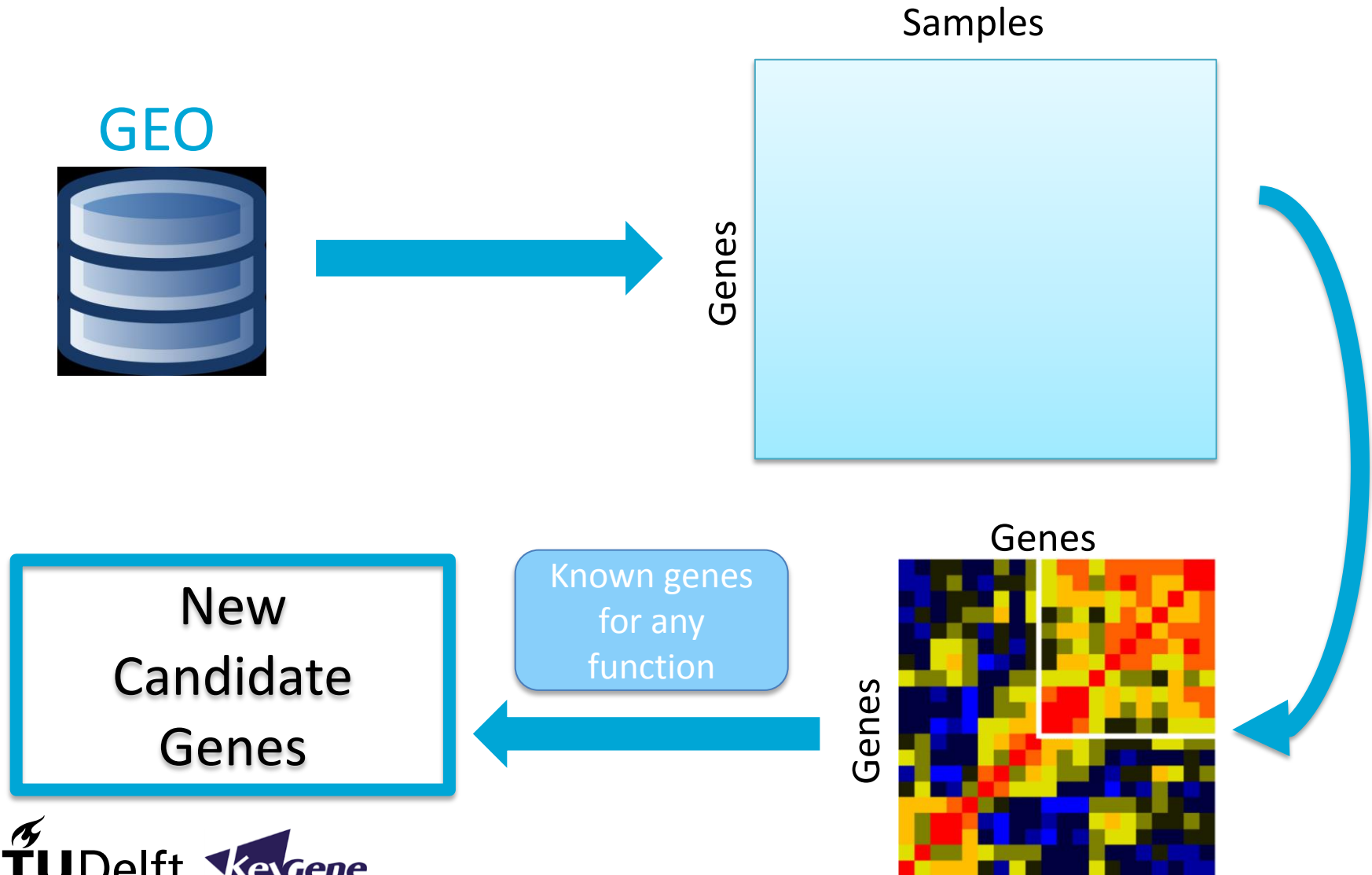
- Which genes affect a given trait?
- For all traits
- Genome-wide gene function prediction
- Experiments are accurate, but difficult

Can we use public data instead?

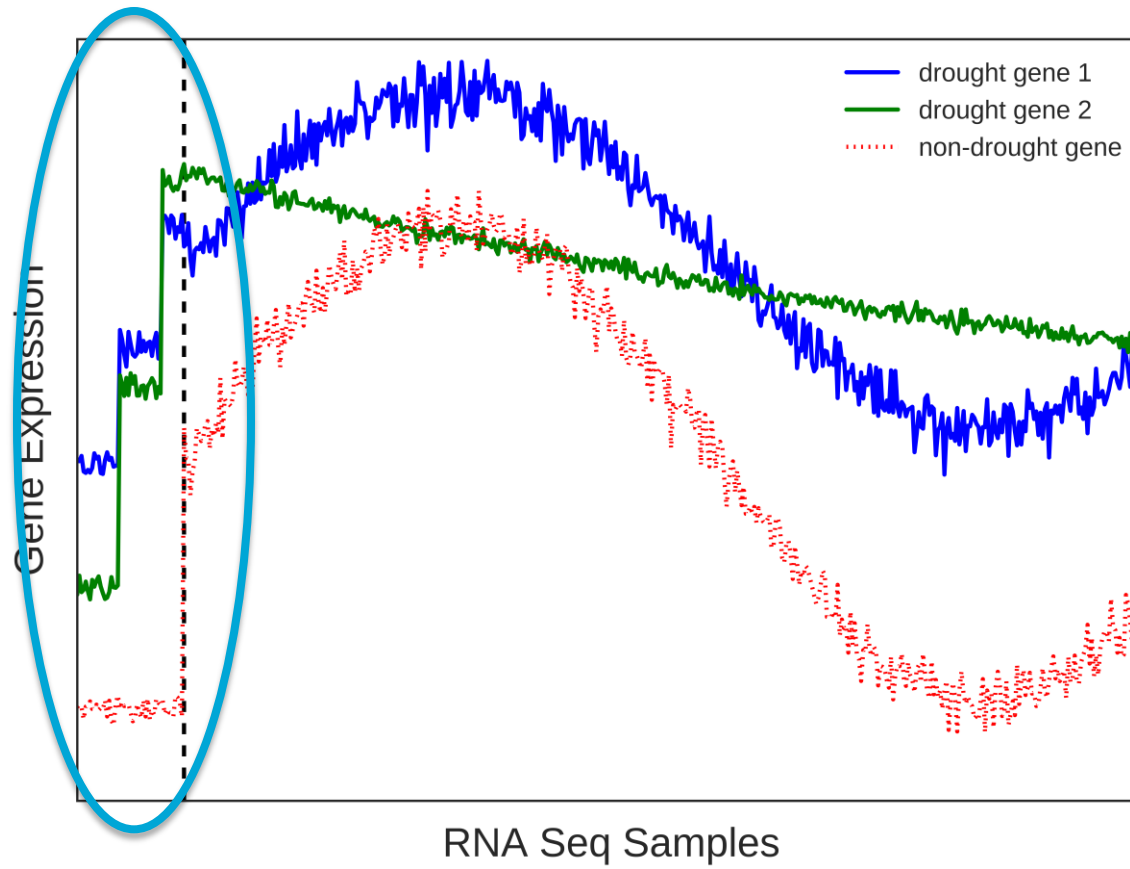
# Many functions, Many conditions



# Co-expression $\rightarrow$ Similar Function



# Most conditions are irrelevant!



# Manual Sample Selection is Hard

- Poor meta-data
- Not obvious what might be relevant

Let the computer decide!

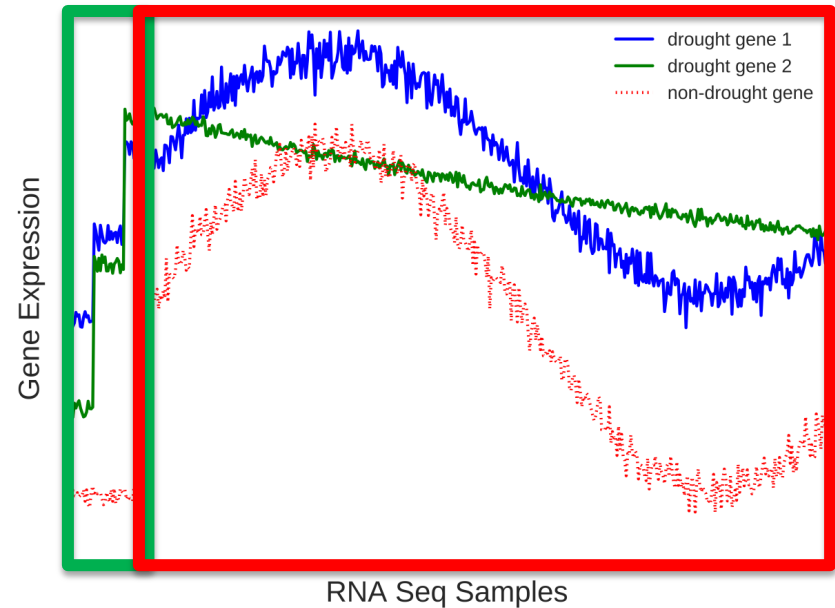
# Score samples based on relevance

$$\text{Coex}(x_i x_j) = \sum_m x_{im} x_{jm}$$

$$\text{Coex}(x_i x_j, W) = \sum_m w_m x_{im} x_{jm}$$

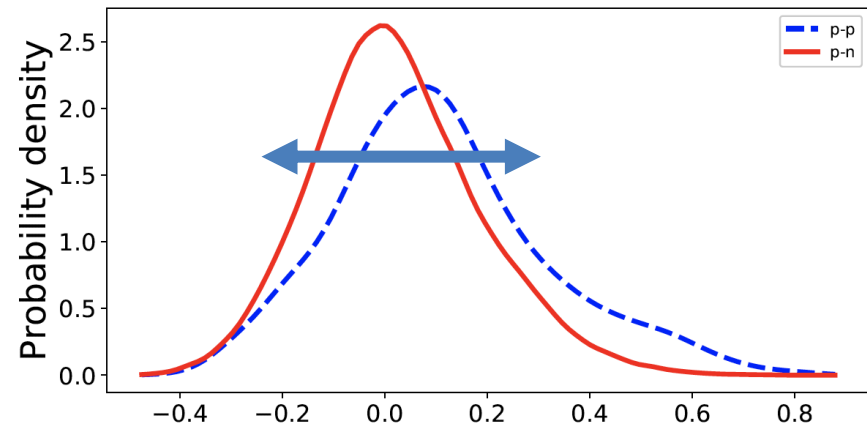
$w =$   
large

$w = 0$

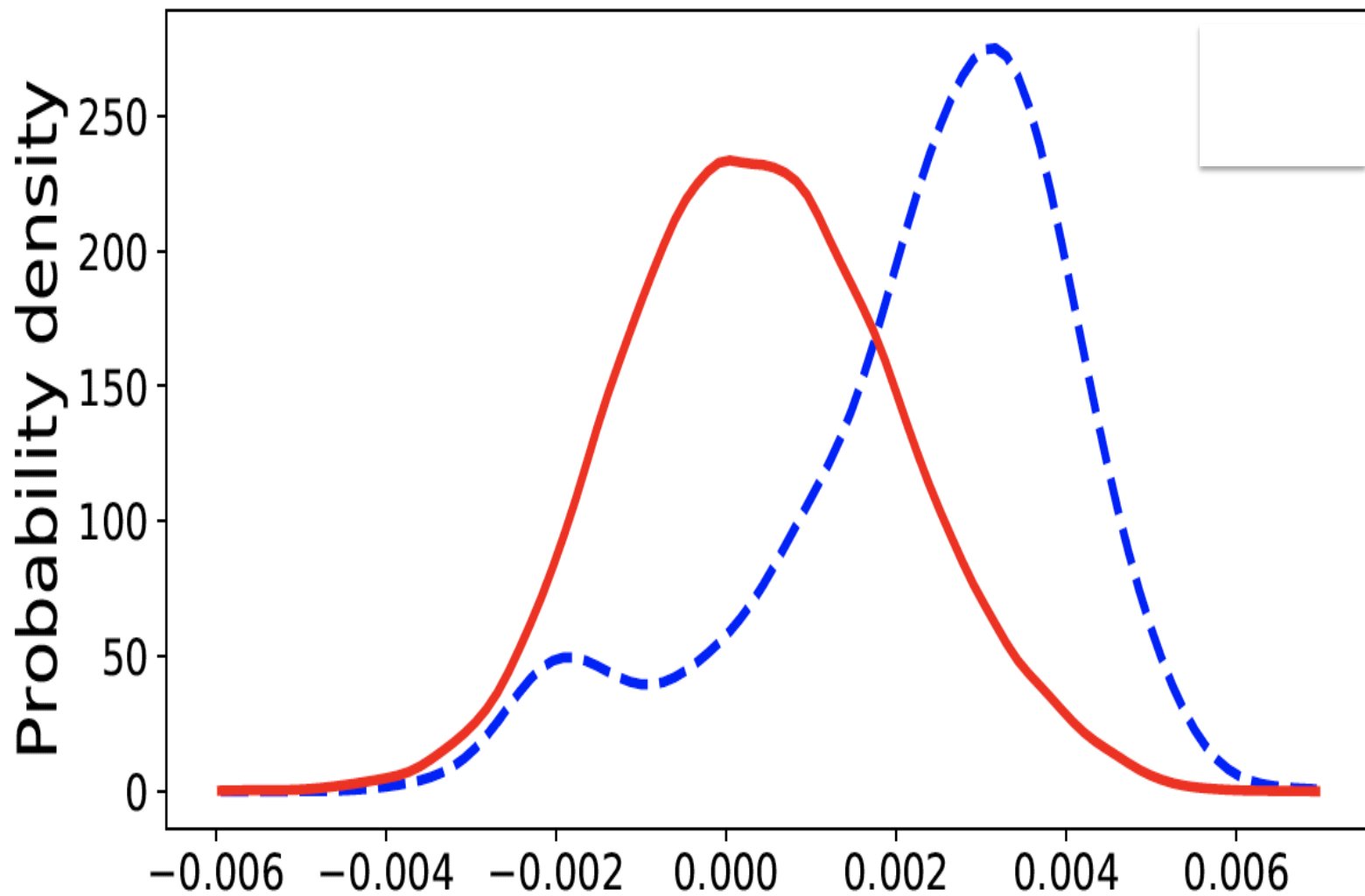


# Automatic Sample Selection

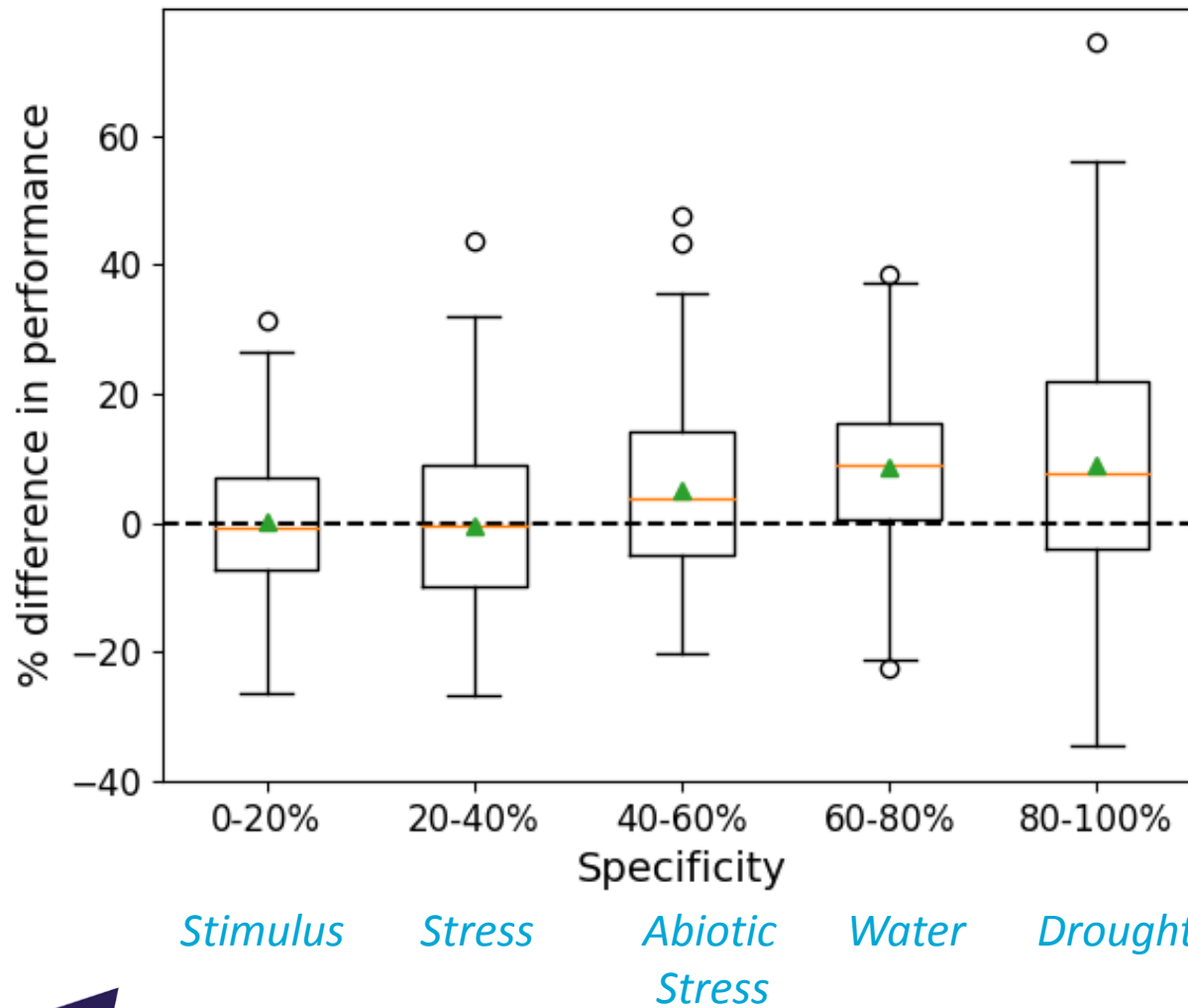
- ML model learns the  $w$ 's
- Genes that have the function should:
  - Co-express with each other
  - Not co-express with other genes







# Results - Arabidopsis thaliana



# Summary - MLC

- Trait-specific candidate gene discovery
- From a large collection of samples
- Learns from the data
- Recovers relevant samples
- Best at discovering specific knowledge

# Questions?

## Acknowledgements

Global Engage

Marcel Reinders

Roeland van Ham



## Contact

[s.makrodimitris@tudelft.nl](mailto:s.makrodimitris@tudelft.nl)

 @sta\_makro

