

HIVE: a new model for Horizontal multi-omics Integration with Variational autoEncoders

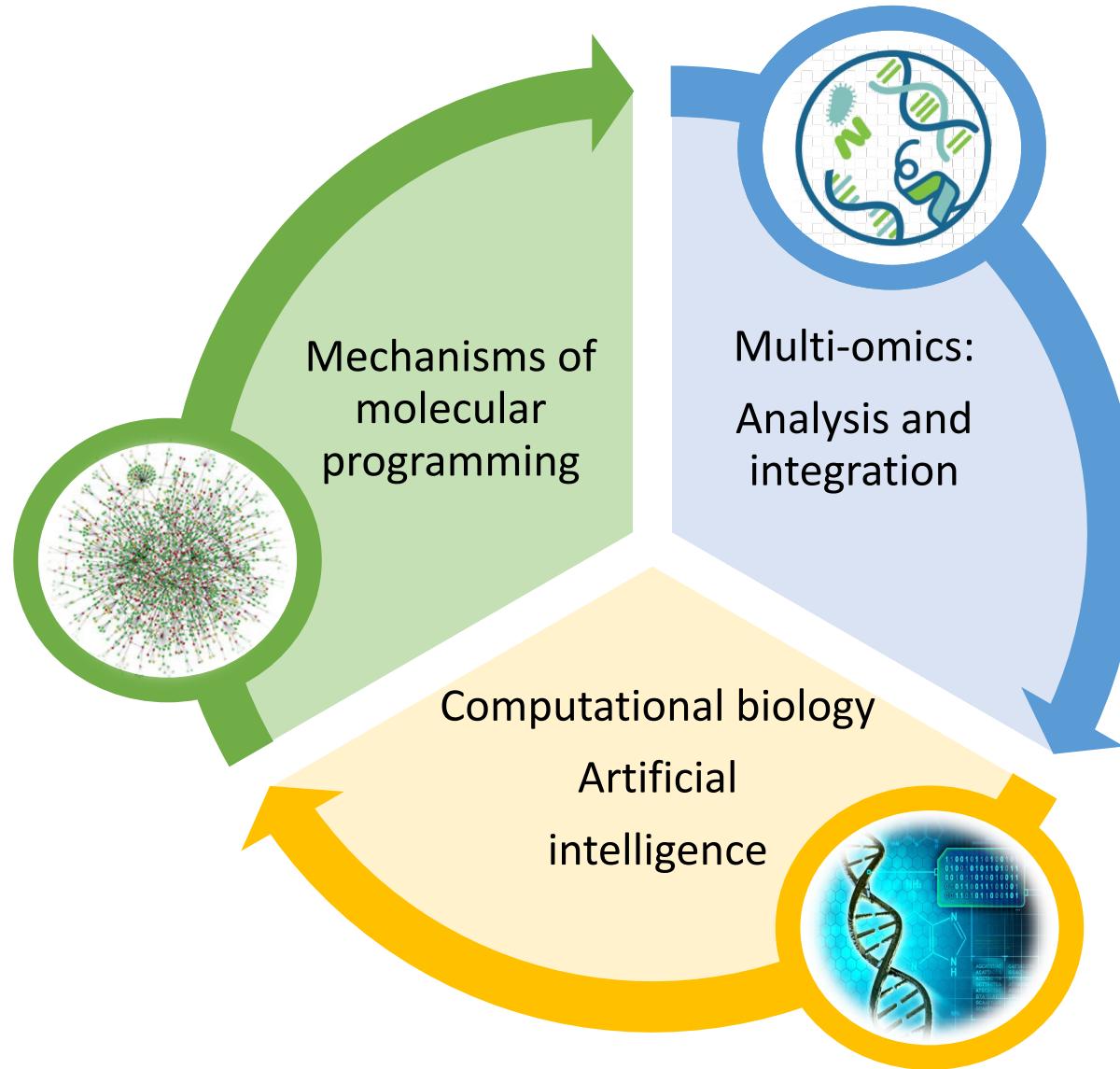


Silvia Bottini, PhD

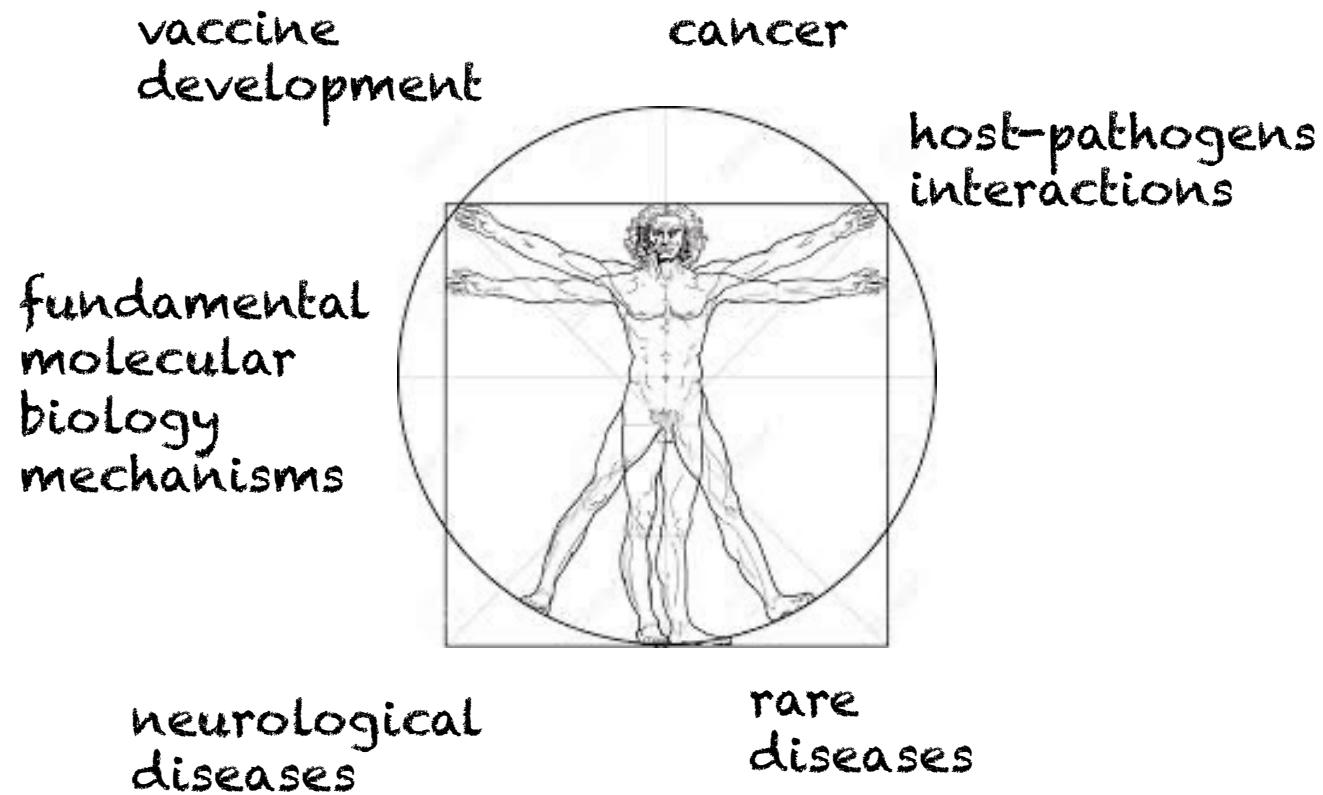
Junior Professor Chair INRAE

Team "M2P2" – Institut Sophia
Agrobiotech – Sophia Antipolis

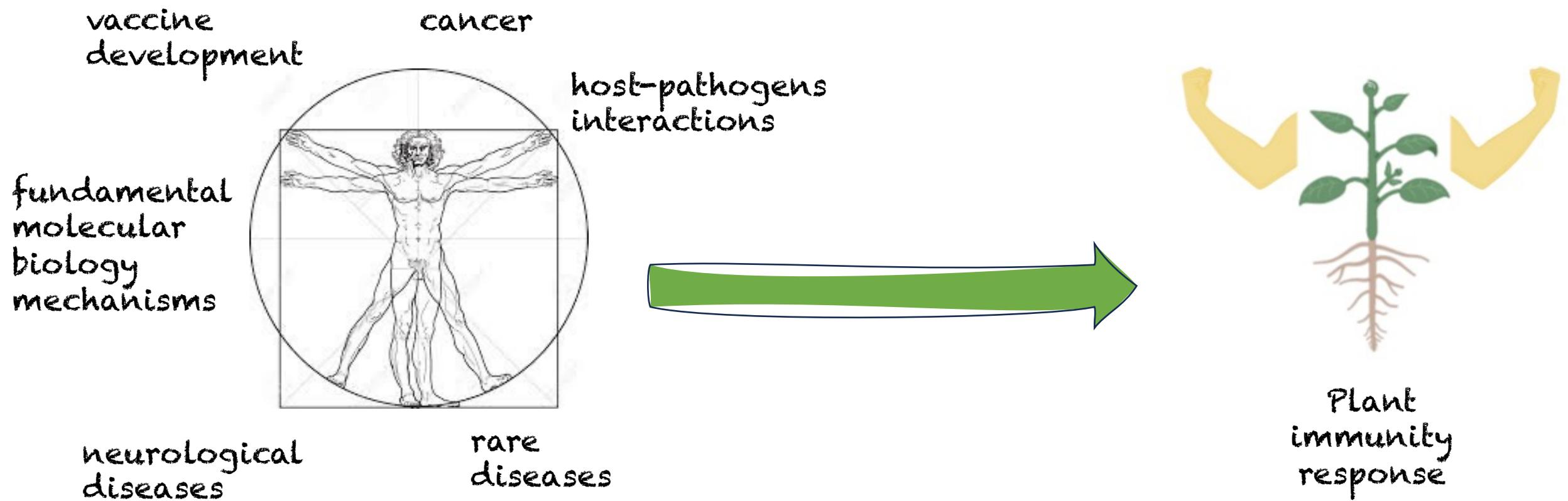
Scientific interests



From human to plant health



From human to plant health



Plants play an important role in human existence



98% of oxygen is produced by plants



82% of the calories in human food supply are provided by terrestrial plants



66% of crop production accounted by **9** species of **plant variety**

Climate change

Environmental degradation
Population growth



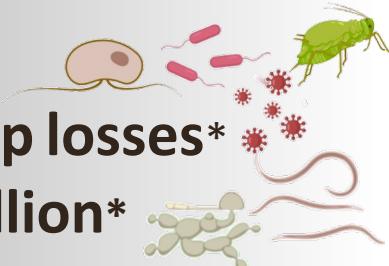
Poor agripractice



Biotic stressors

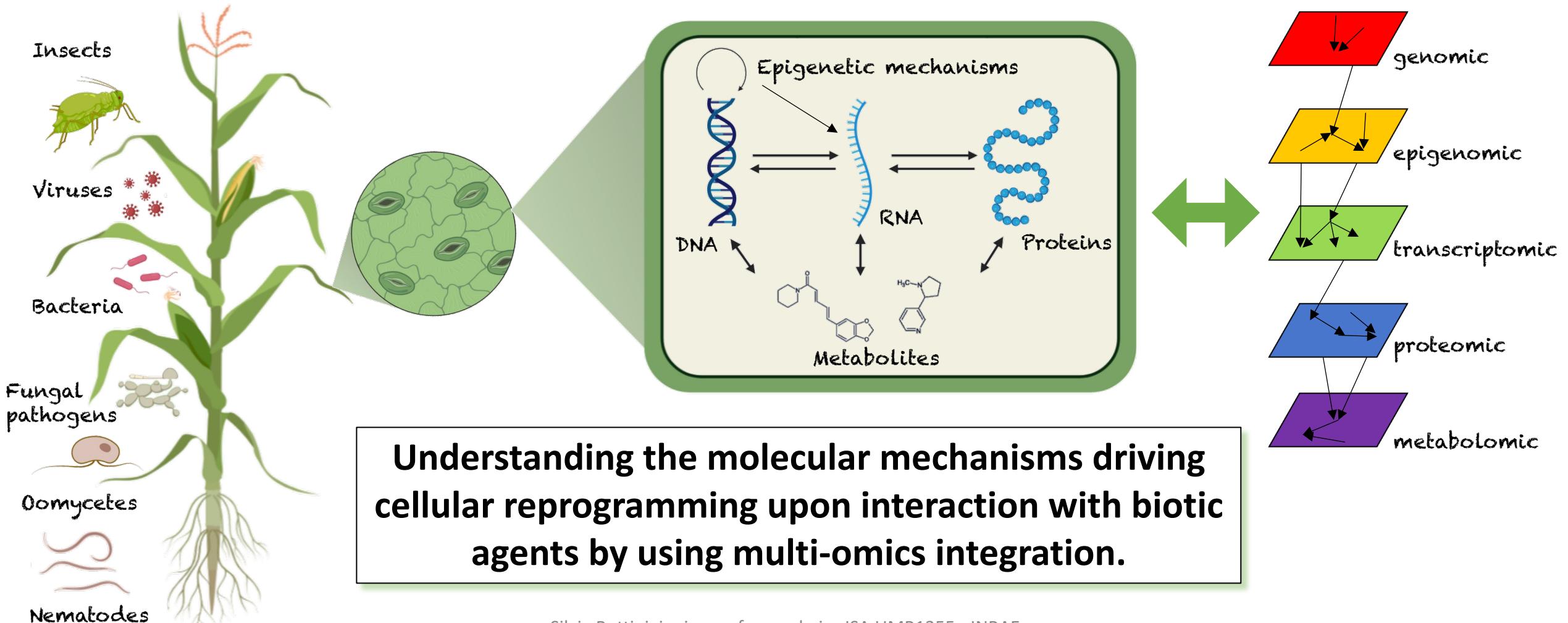
Between **20%** to **40%** of **crop losses***

Economic loss around **\$220 billion***

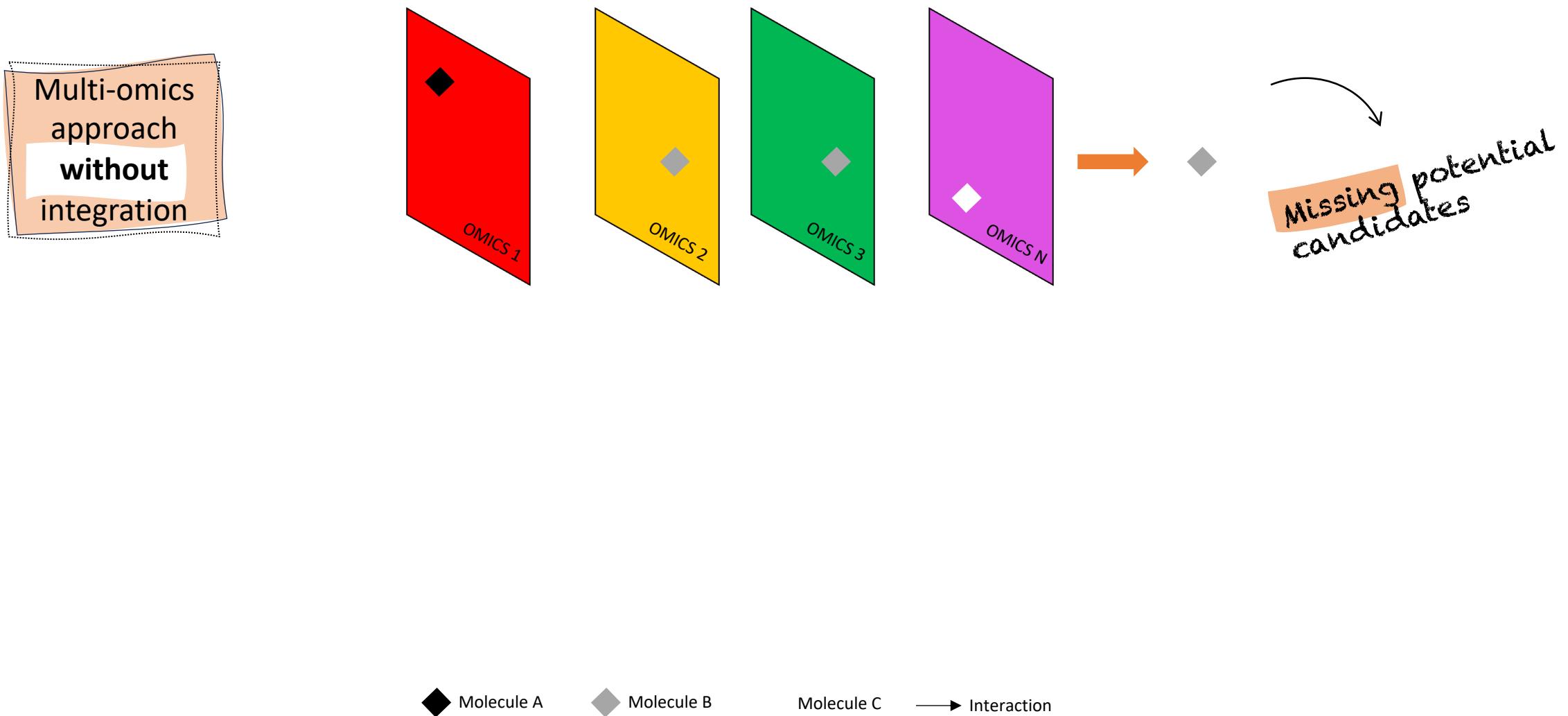


*according to the Food and Agriculture Organization of the United Nations (2023)

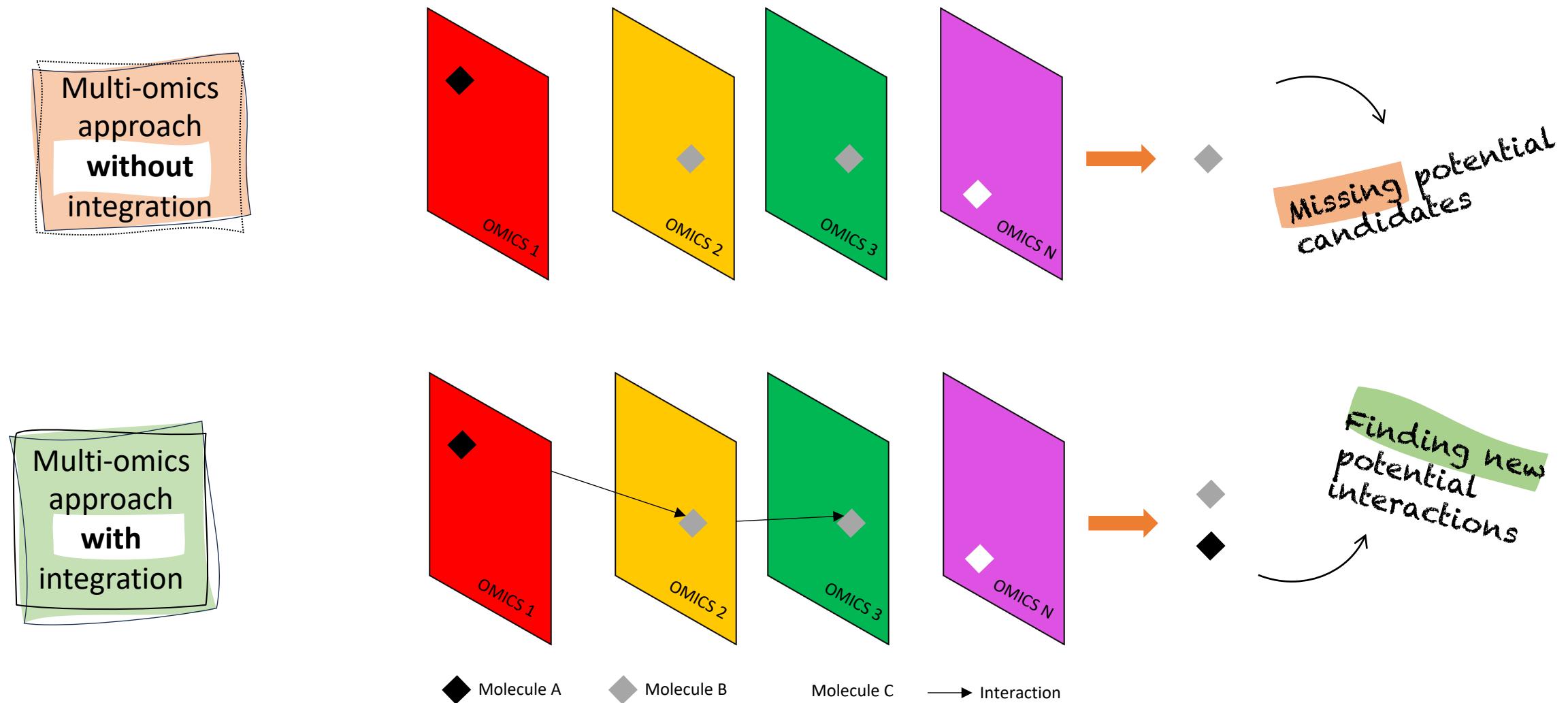
How to improve plant health?



Why multi-omics integration?

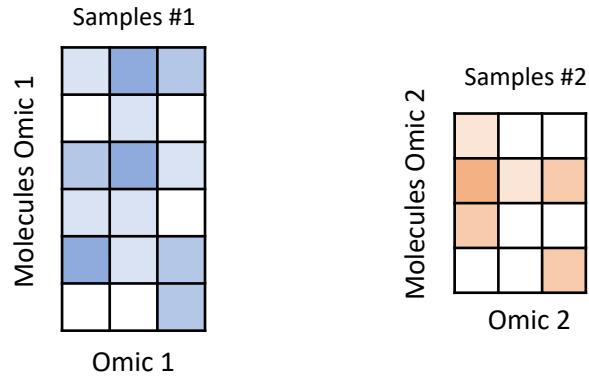


Why multi-omics integration?

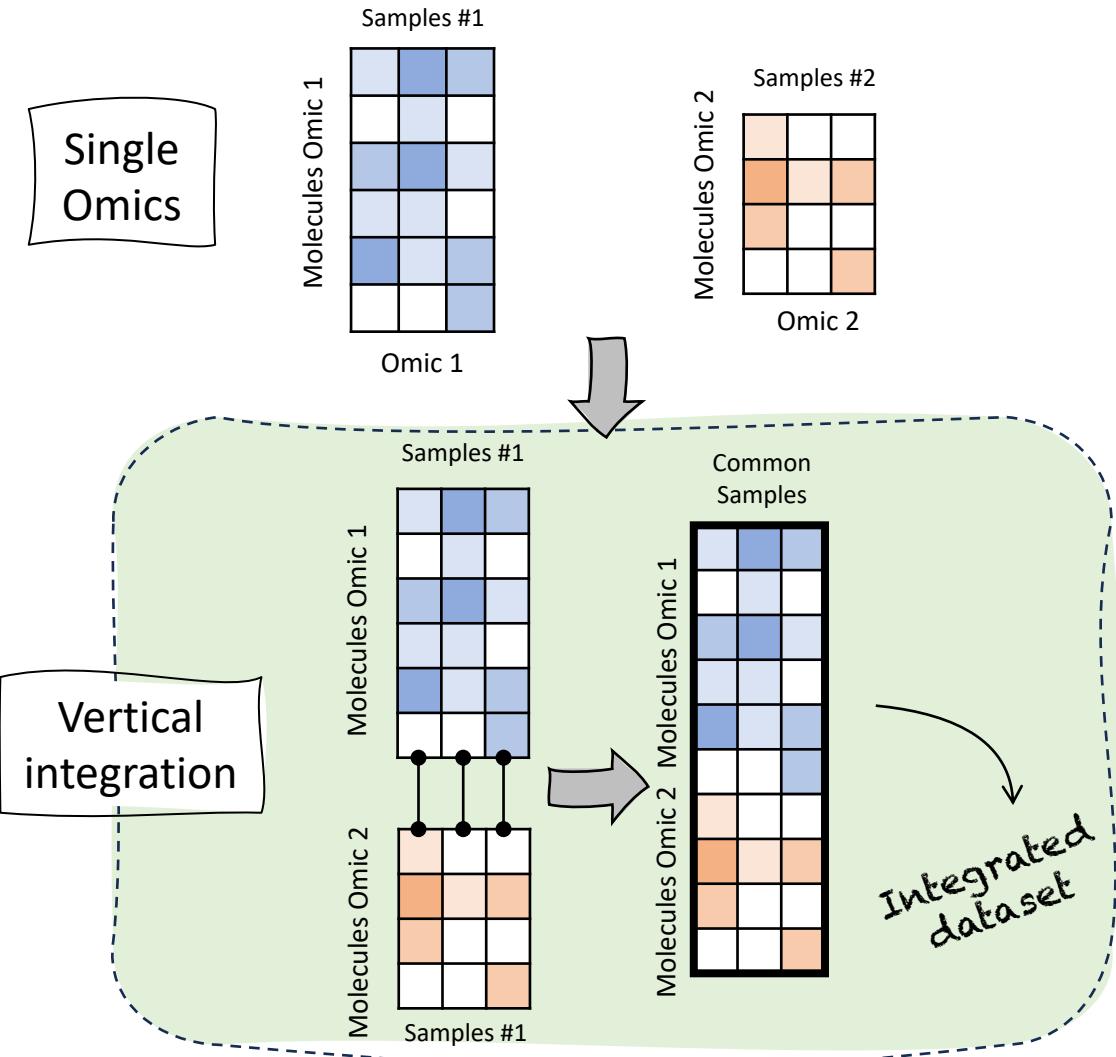


How to do multi-omics integration?

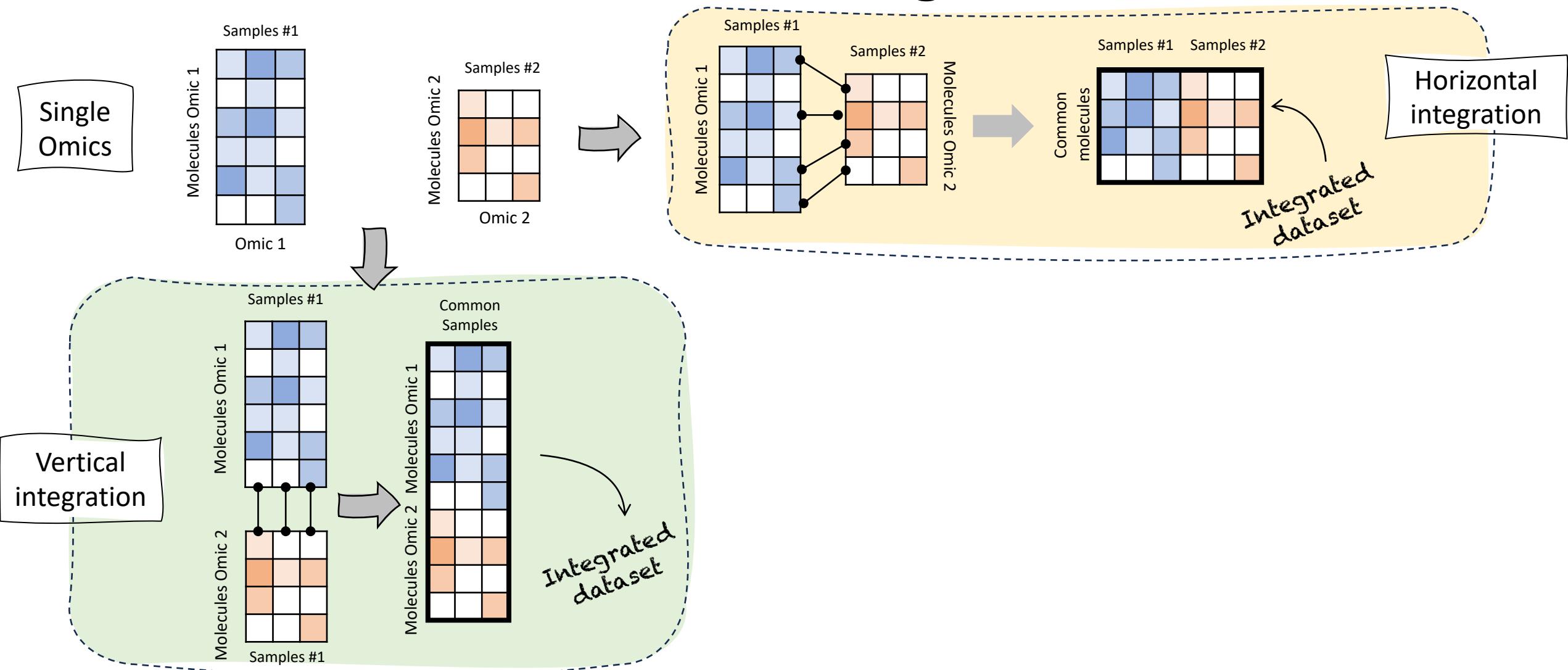
Single
Omics



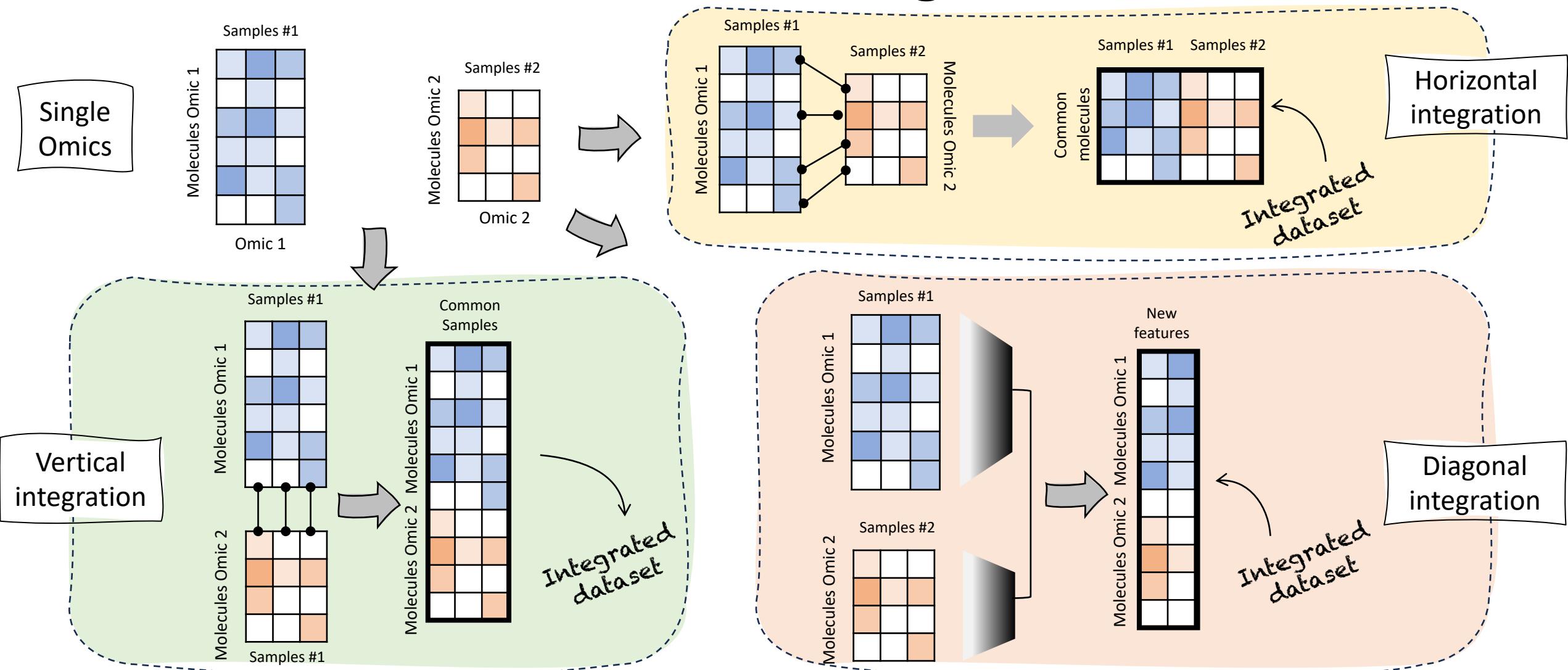
How to do multi-omics integration?



How to do multi-omics integration?

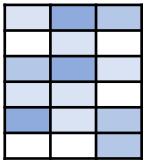
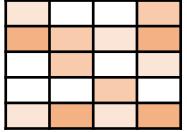


How to do multi-omics integration?



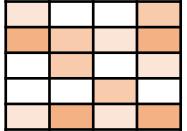
Challenges in multi-omics integration analysis

Challenges in multi-omics integration analysis

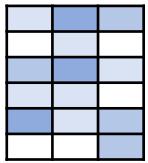


Heterogeneity, sparsity
and uneven datasets.

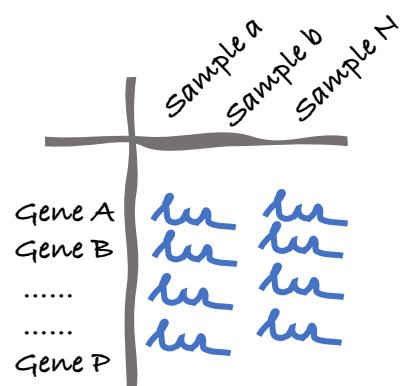
Challenges in multi-omics integration analysis



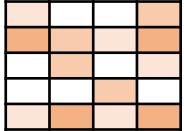
Heterogeneity, sparsity
and uneven datasets.



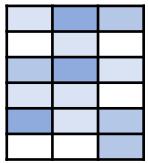
More features than
data ($p \gg n$).



Challenges in multi-omics integration analysis



Heterogeneity, sparsity
and uneven datasets.



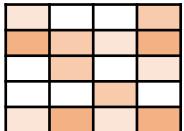
More features than
data ($p \gg n$).

	Sample a	Sample b	Sample c
Gene A	nan	nan	nan
Gene B	nan	nan	nan
.....
.....
Gene P	nan	nan	nan

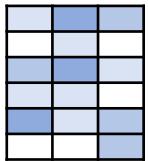
NaN

Missing data.

Challenges in multi-omics integration analysis



Heterogeneity, sparsity
and uneven datasets.

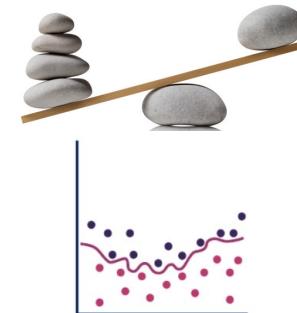


More features than
data ($p \gg n$).

	Sample a	Sample b	Sample c
Gene A	NaN	NaN	NaN
Gene B	NaN	NaN	NaN
.....
.....
Gene P	NaN	NaN	NaN

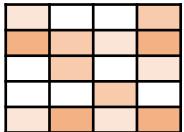
NaN

Missing data.

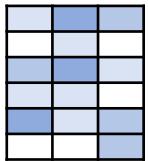


Class imbalance and
overfitting.

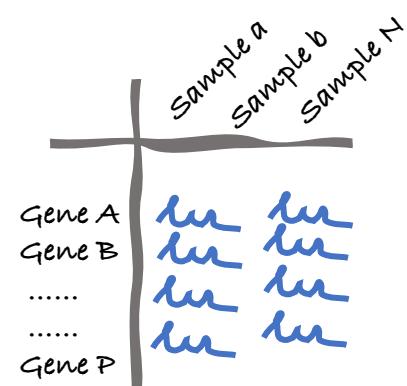
Challenges in multi-omics integration analysis



Heterogeneity, sparsity
and uneven datasets.

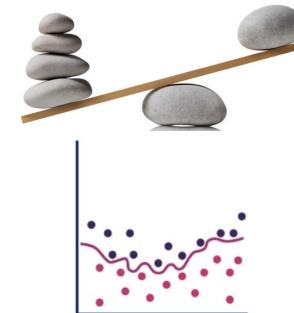


More features than
data ($p \gg n$).



NaN

Missing data.



Class imbalance and
overfitting.



How to choose a good
model to perform the
integration?

Multi-OMICS integration



Data-driven

Concatenation first



- Easy
- Straightforward
- Need proper normalization
- Do not take into account the different distribution of each omic
- Computationally intensive

Feature selection first



- Reduce computational costs
- Inter-omics relationships are lost
- Weak signals could be lost

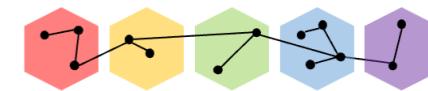
Conversion first



- Allow easy representation of data
- May prevents the identification of indirect mechanisms
- Transformations can be challenging

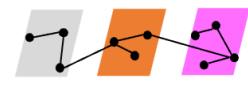
Hypothesis-driven

Interactome first



- Allow inter-omics connections
- Strongly depends on existing knowledge

Pathway as network



- Incorporate pathway topology
- Hard to extend
- Computationally intense

Multi-OMICS integration



Data-driven

Hypothesis-driven

Concatenation
first



- Easy
- Straightforward

- Need proper normalization
- Do not take into account the different distribution of each omic
- Computationally intensive

Model choice
depends on the
available data and
the biological
question.

Pathway as
network



- Incorporate pathway topology

- Hard to extend
- Computationally intense



Peanut – *Arachis spp.*

- Family: fabaceae
- Native: South America
- Important legume in Africa and Asia
- High protein content

Wild species

- More than 80 species
- Originated exclusively from South America
- Important to study as a source of resistance characteristics to biotic and abiotic stresses.



Giulia Calia
PhD student



Ana Zotta Mota
Post-doc



Patricia Messenberg
Guimaraes

Biological model: two wild peanuts species



*Arachis
duranensis*



*Arachis
stenosperma*

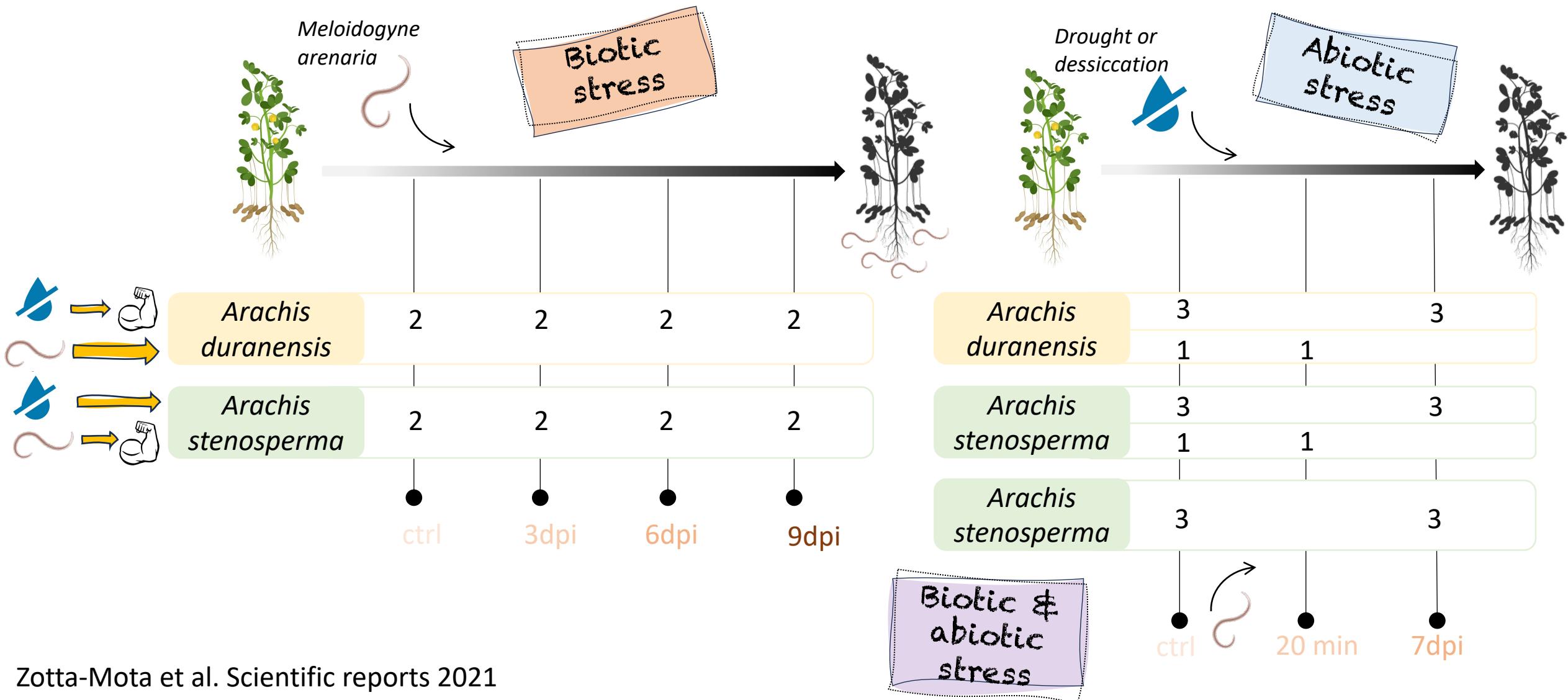
Abiotic stress:
drought



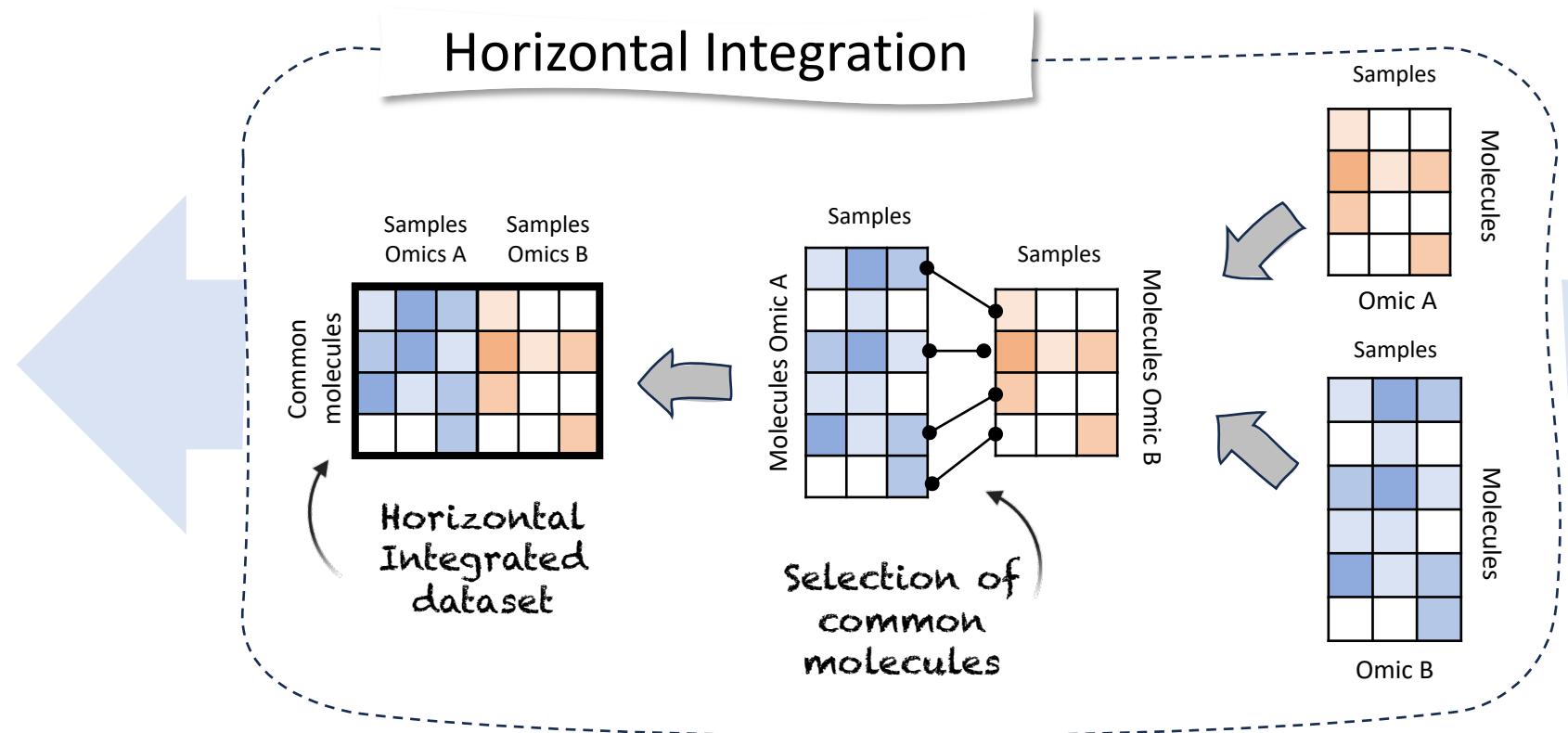
Biotic stress:
root knot nematodes
Meloidogyne arenaria



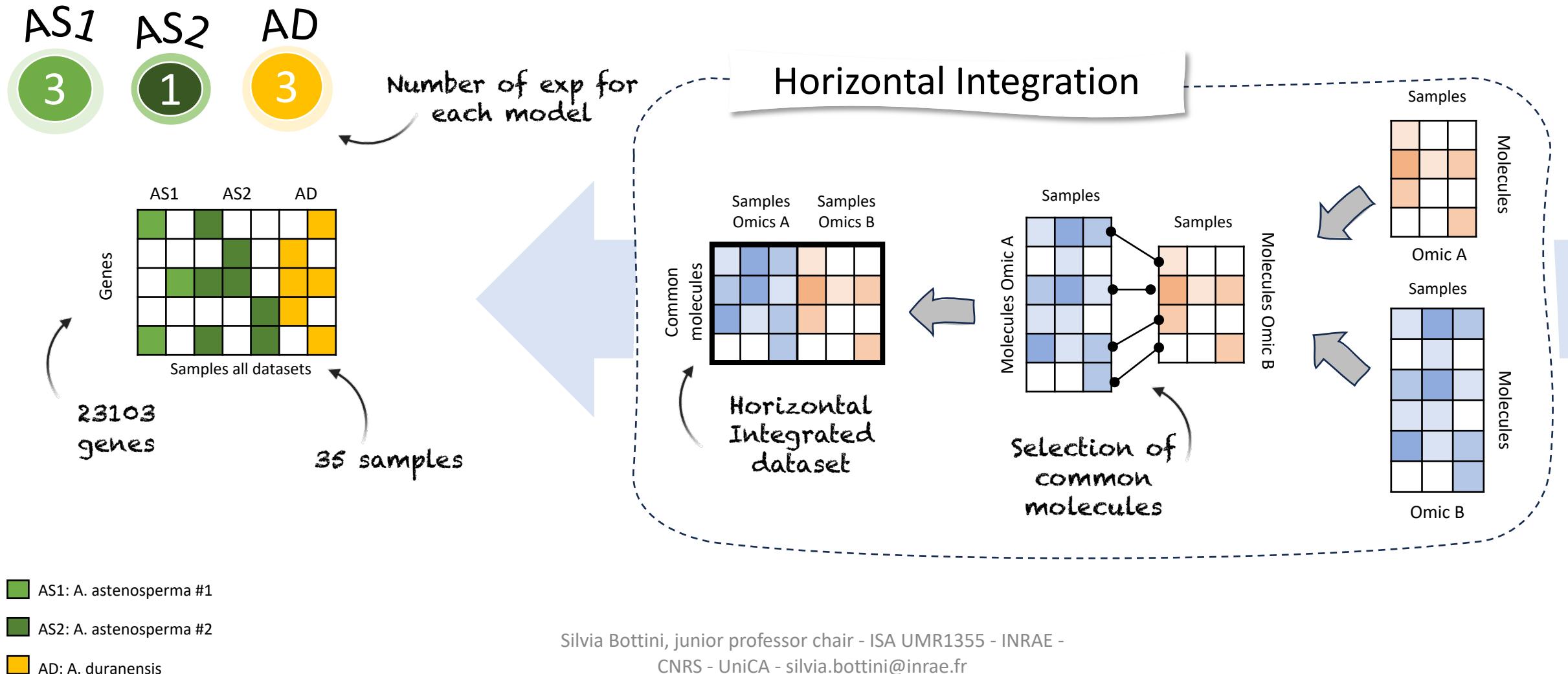
The experimental design



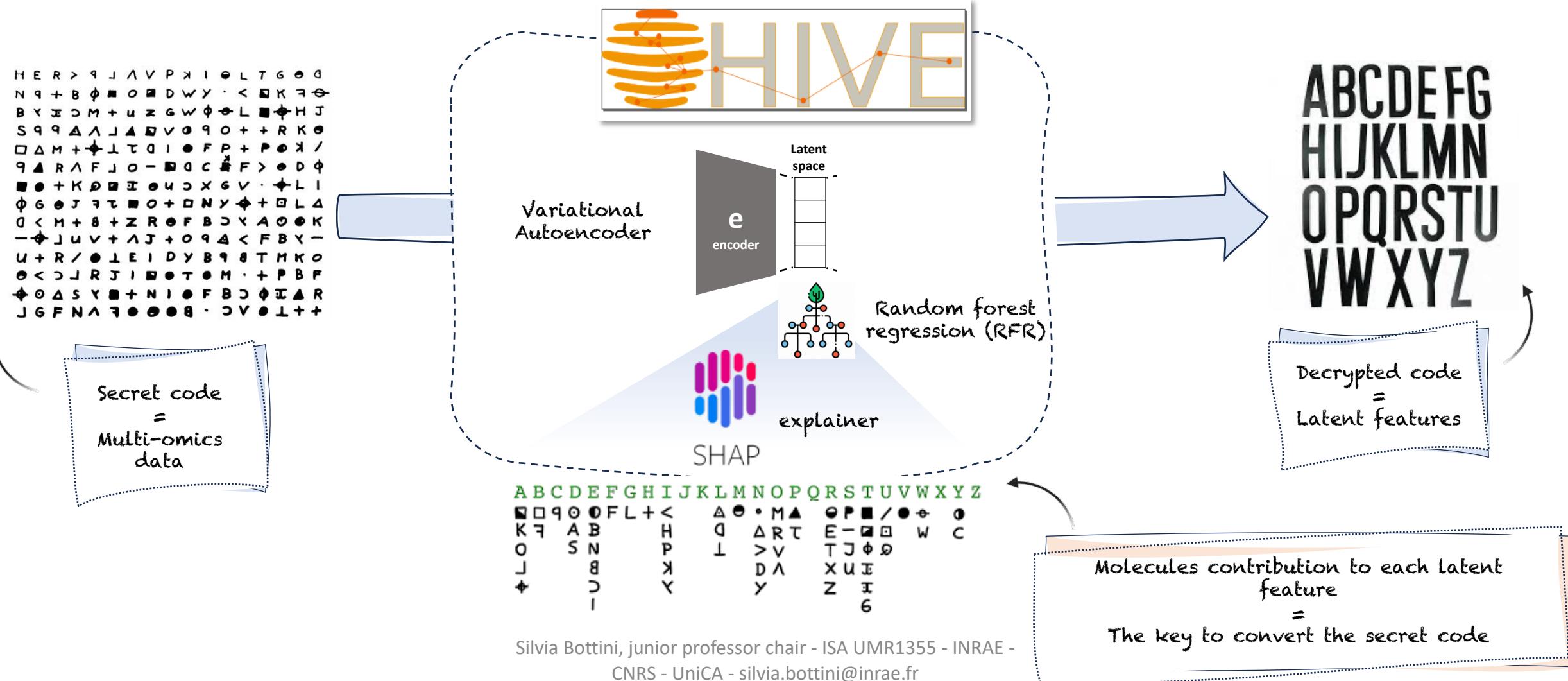
The pilot integrated dataset



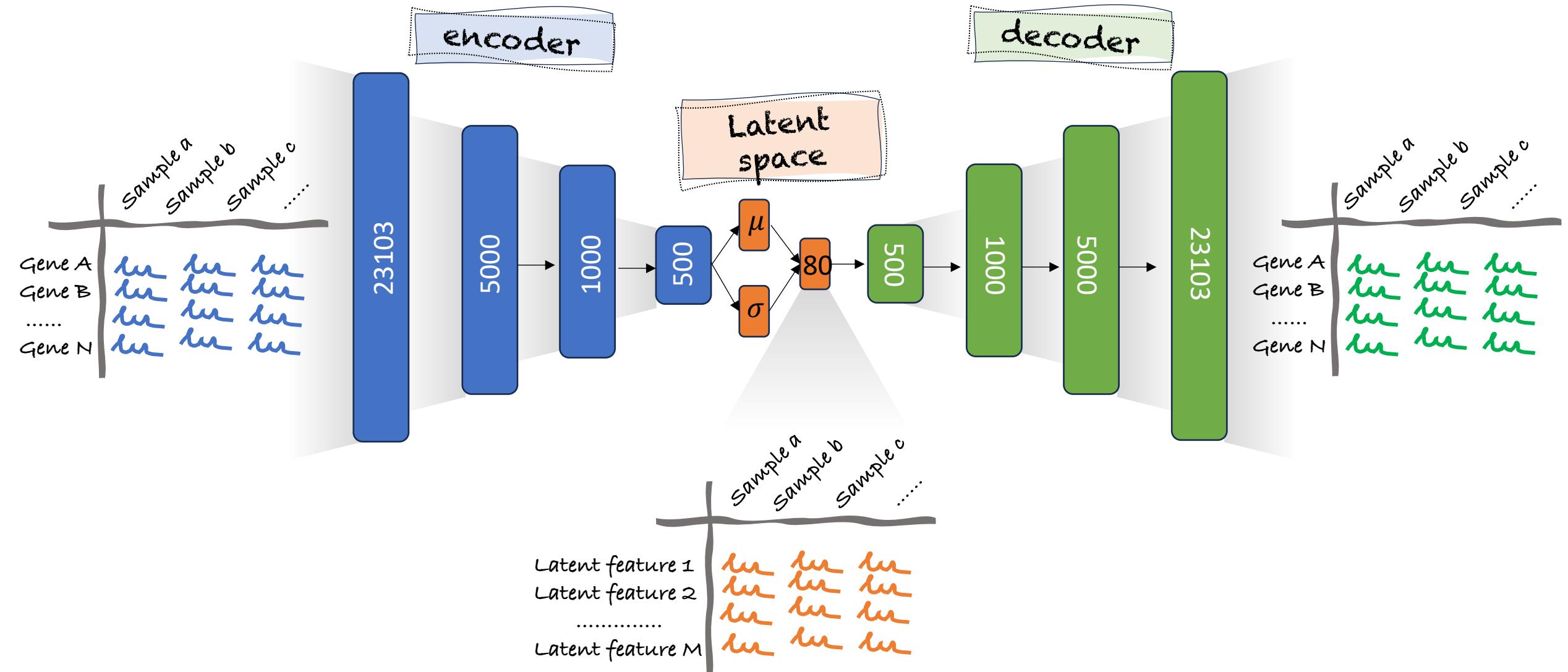
The pilot integrated dataset



HIVE: a general framework to analyse integrated multi-omics data



Variational autoencoder



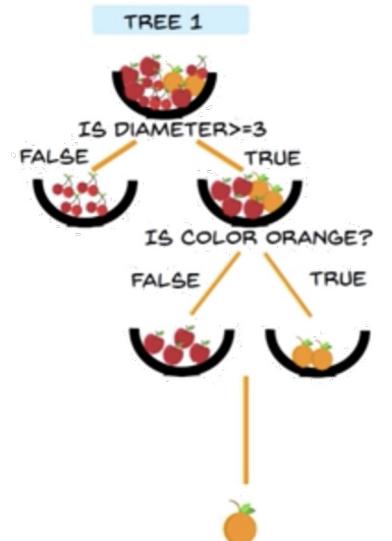
Random forest

The random forest algorithm makes use of **multiple decision trees**.
It can solve both regression and classification problems.



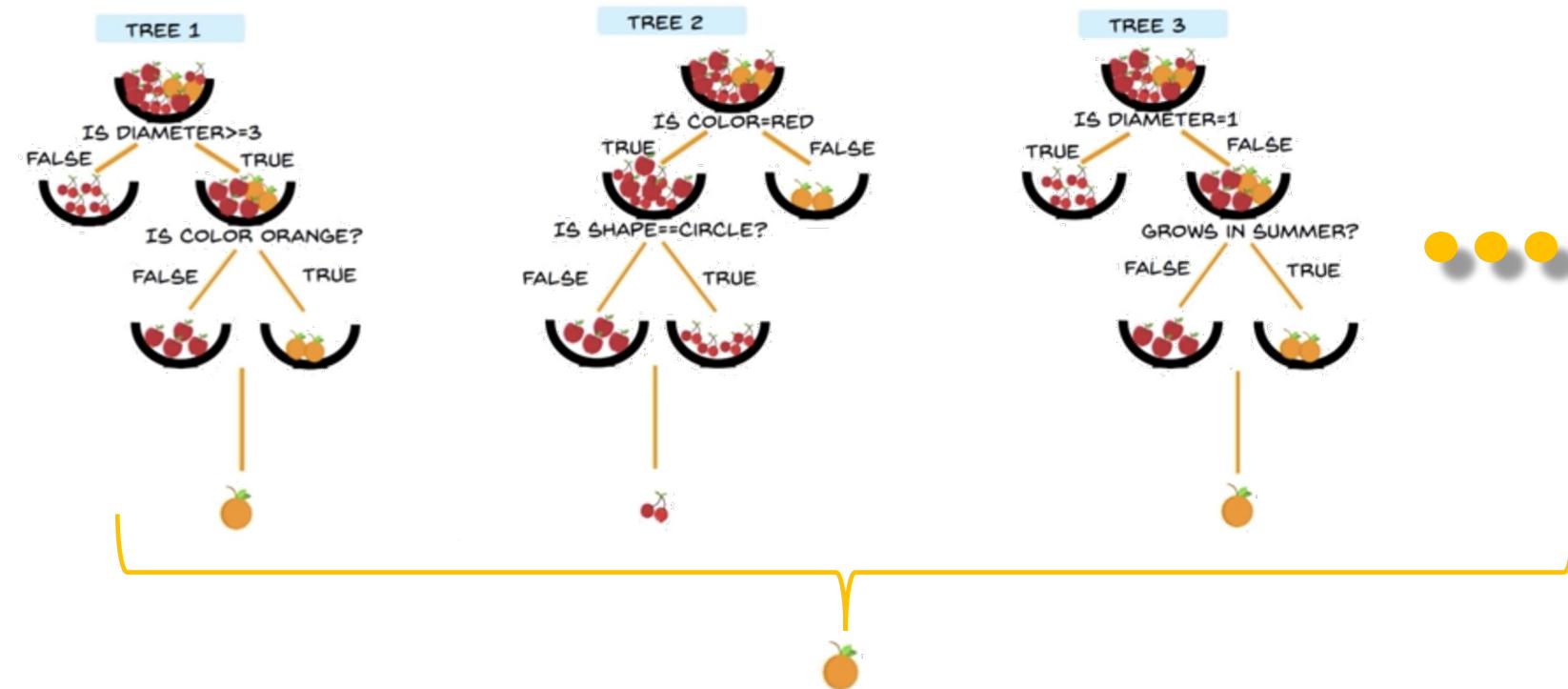
Random forest

The random forest algorithm makes use of **multiple decision trees**.
It can solve both regression and classification problems.



Random forest

The random forest algorithm makes use of **multiple decision trees**.
It can solve both regression and classification problems.



Robustness

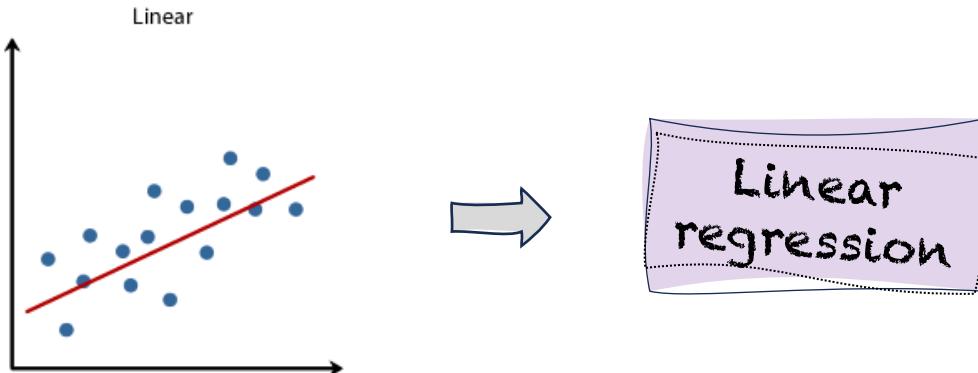


Computationally efficient



Immune to the curse of dimensionality

Random forest regression



$$Y = \alpha + \beta x + \varepsilon$$

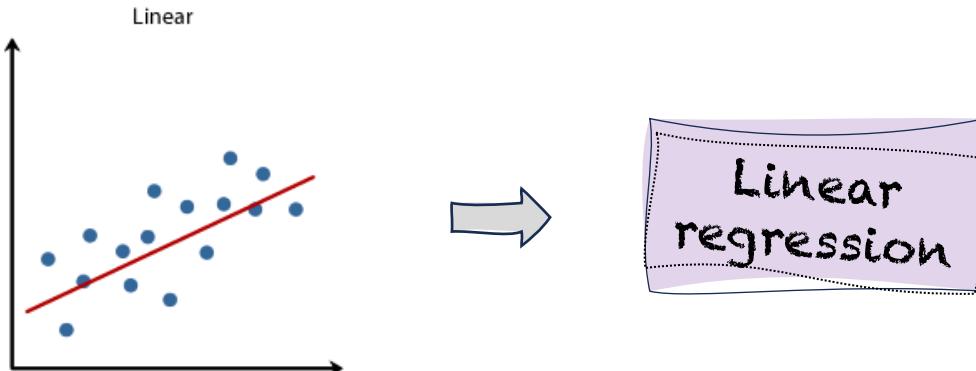
Diagram illustrating the linear regression equation:

- Y : Response, dependent variable (circled in red)
- α : coefficients (circled in green)
- βx : explanatory, independent variable (circled in green)
- ε : noise (circled in blue)

✗ Large sample size and low variance in the data set.

✗ Assumptions about the relations between predictors and target variables, such as constant variance, normality of errors, etc.

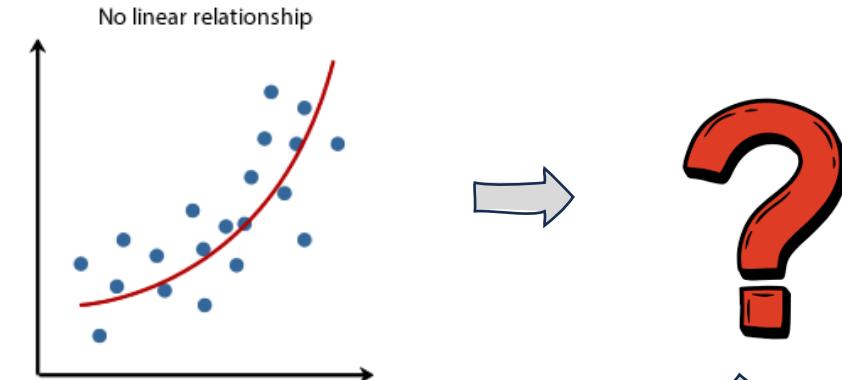
Random forest regression



$$Y = \alpha + \beta x + \varepsilon$$

coefficients
explanatory, independent variable
Response, dependent variable
noise

Annotations explain the components of the linear regression equation: Y is the response variable, x is the explanatory variable, α and β are coefficients, and ε represents noise.



A yellow box containing the text "Random forest regression". An arrow points from the "No linear relationship" plot to this box.

- ✖ Large sample size and low variance in the data set.
- ✖ Assumptions about the relations between predictors and target variables, such as constant variance, normality of errors, etc.

- ✓ Does not make any assumption about the underlying data distribution.
- ✓ Handle nonlinear relationships better.
- ✓ Less prone to overfitting.

The Shapley value

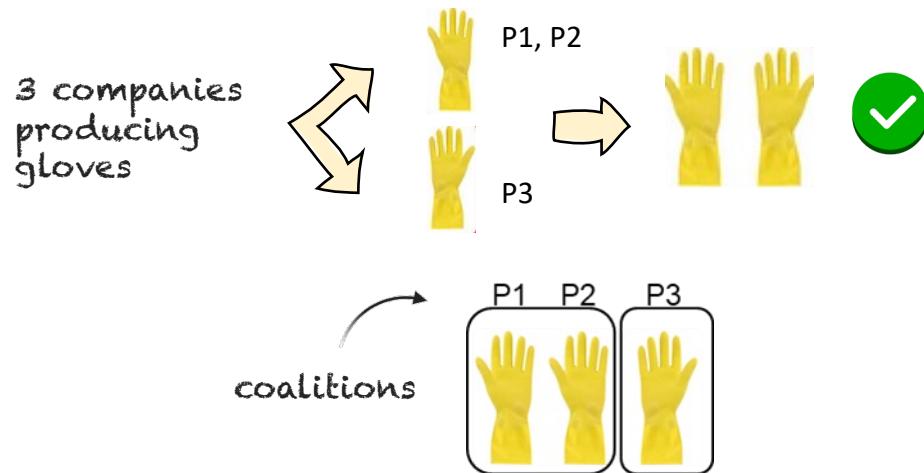
Def. In a cooperative game, “players” have the possibility to forge coalitions to achieve a common goal. One difficulty in the theory of cooperative games is the distribution of benefits among the players.

Lloyd Shapley, introduced the concept in 1953 and received the Nobel Prize in Economics in 2012.



The Shapley value

Def. In a cooperative game, “players” have the possibility to forge coalitions to achieve a common goal. One difficulty in the theory of cooperative games is the distribution of benefits among the players.



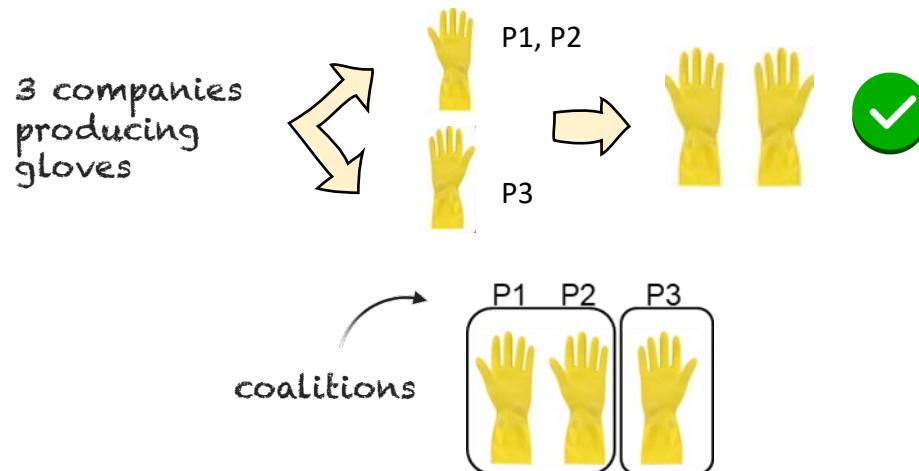
Lloyd Shapley, introduced the concept in 1953 and received the Nobel Prize in Economics in 2012.



The Shapley value

Def. In a cooperative game, “players” have the possibility to forge coalitions to achieve a common goal. One difficulty in the theory of cooperative games is the distribution of benefits among the players.

Lloyd Shapley, introduced the concept in 1953 and received the Nobel Prize in Economics in 2012.



More formally:

$$v(P1) = v(P2) = v(P3) = 0$$

$$v(P1, P2) = 0$$

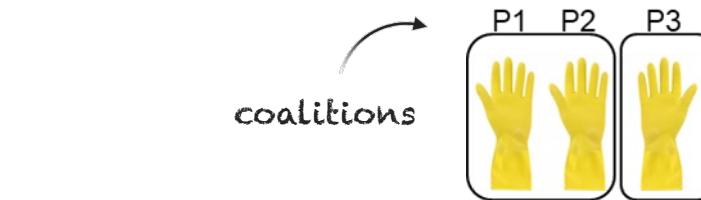
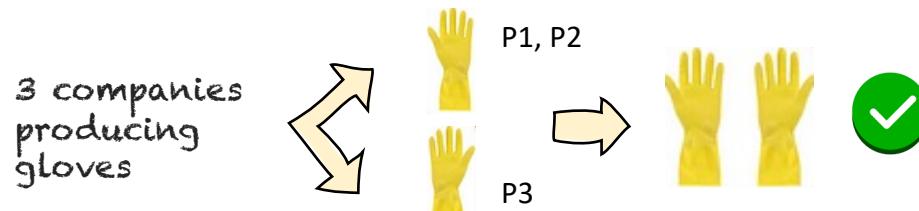
$$v(P1, P3) = v(P2, P3) = v(P1, P2, P3) = 1$$



Maximum gain

The Shapley value

Def. In a cooperative game, “players” have the possibility to forge coalitions to achieve a common goal. One difficulty in the theory of cooperative games is the distribution of benefits among the players.



More formally:

$$v(P1) = v(P2) = v(P3) = 0$$



$$v(P1, P2) = 0$$



$$v(P1, P3) = v(P2, P3) = v(P1, P2, P3) = 1$$



Maximum gain

Lloyd Shapley, introduced the concept in 1953 and received the Nobel Prize in Economics in 2012.

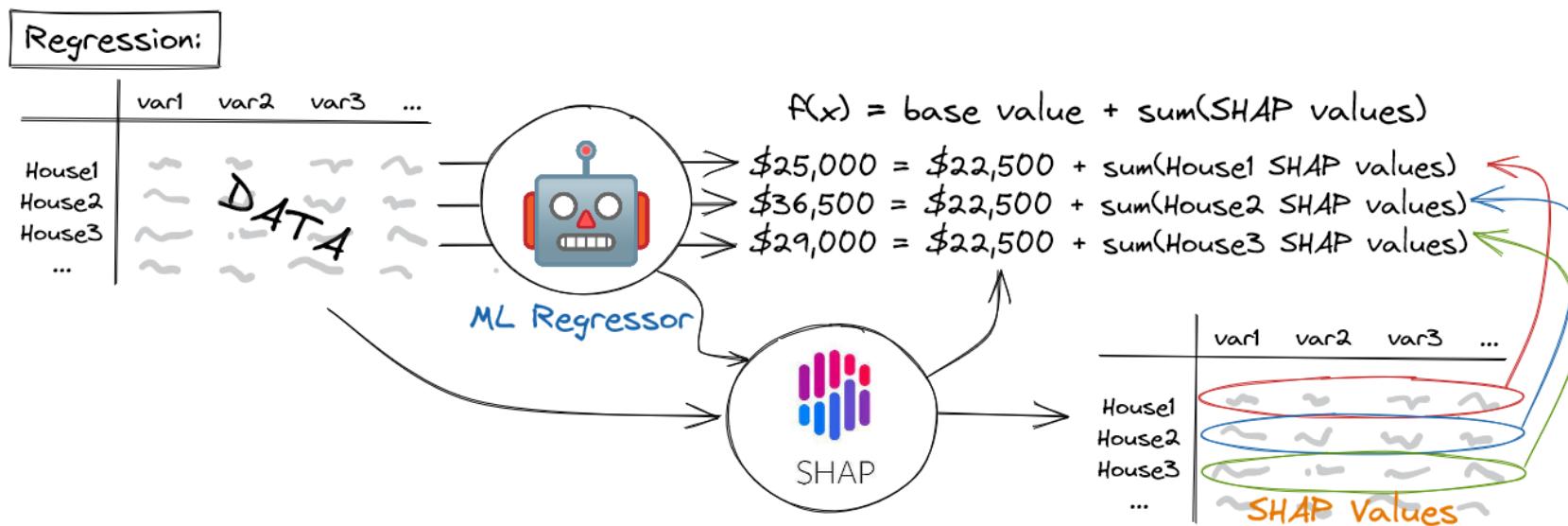


How can the revenue be redistributed fairly?

σ	Left-handed glove	Right-handed glove
P1, P2, P3	0	0
P1, P3, P2	0	0
P2, P1, P3	0	0
P2, P3, P1	0	0
P3, P1, P2	1	0
P3, P2, P1	0	1
$\emptyset_i(N, v)$	1/6	1/6
		4/6

SHAP explainer

SHAP is an explainable artificial intelligence technique based on Shapley value and used in machine learning to determine how input variables contribute to output predictions.



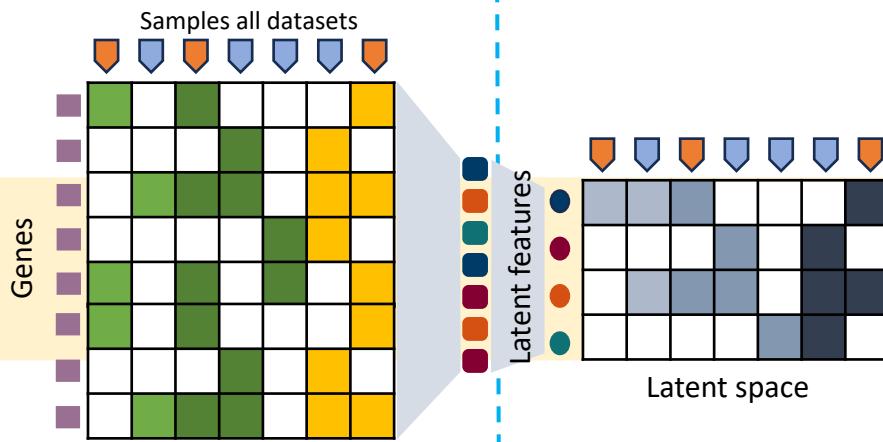
✓ SHAP can be applied to any machine learning model as a *post hoc* interpretation technique

✓ it is particularly efficient to compute SHAP for tree-based models, such as random forests and gradient boosted trees.

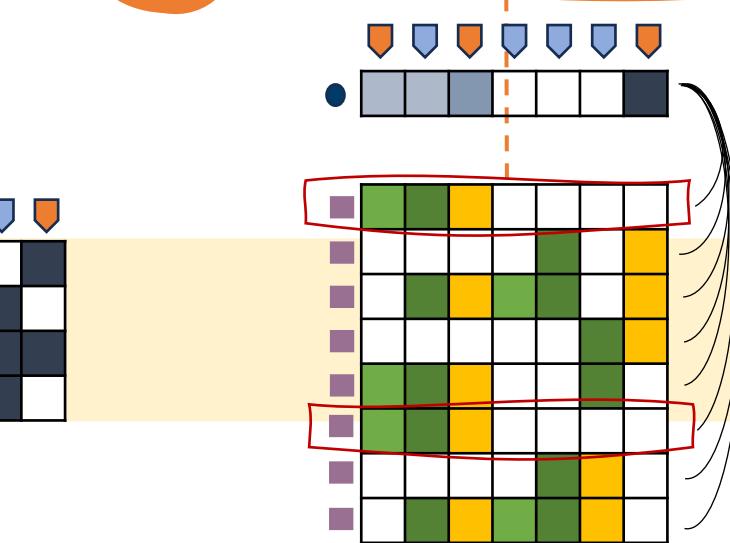


HIVE pipeline

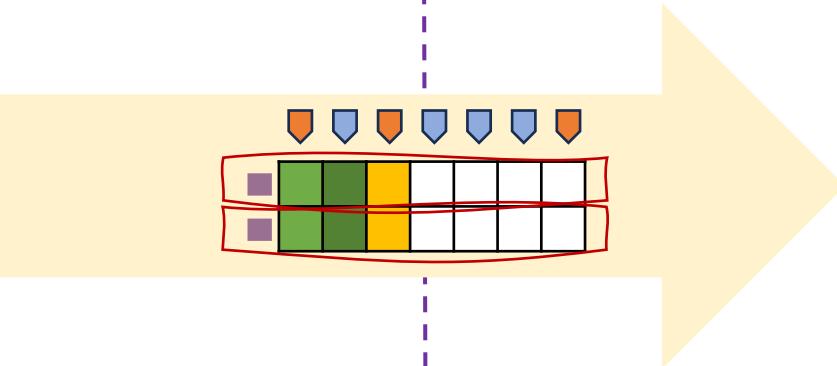
1 Removing batch effects



2 Opening the black box



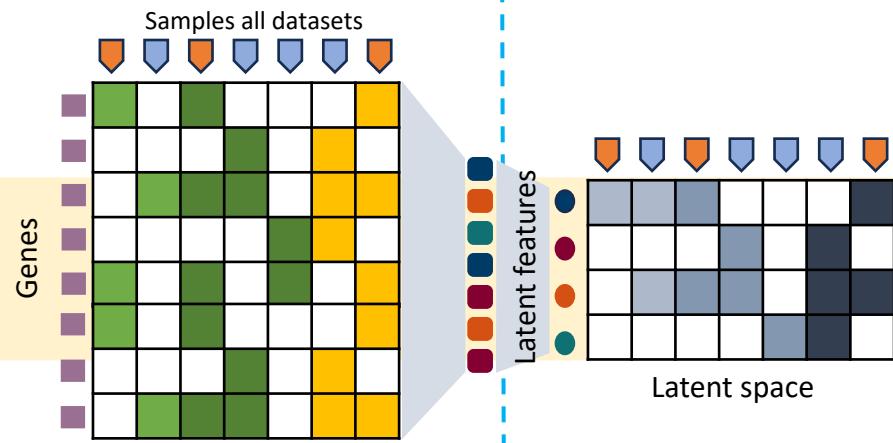
3 Finding important genes



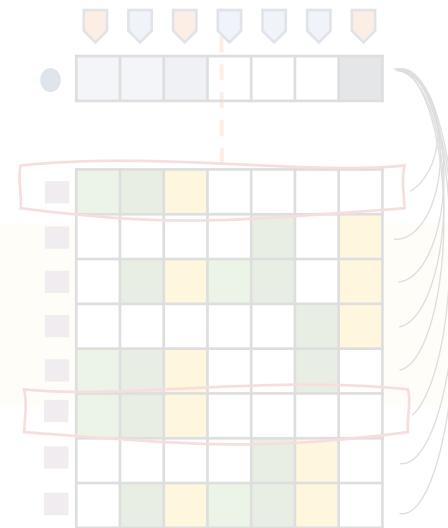


HIVE pipeline

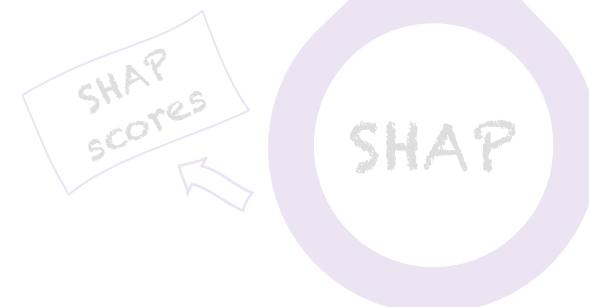
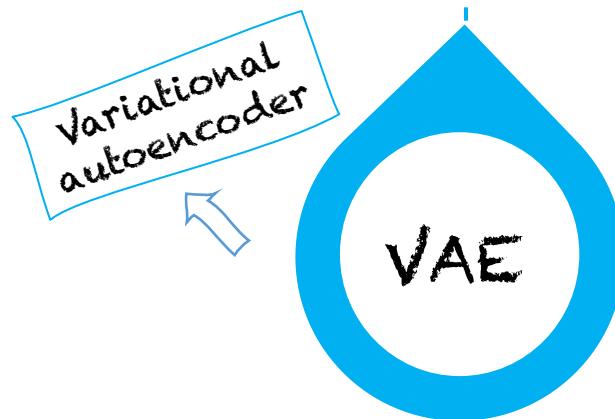
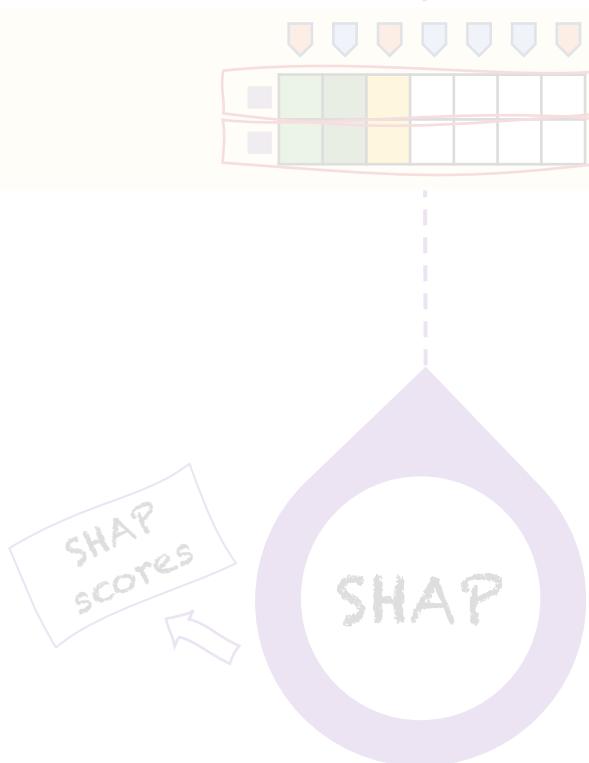
1 Removing batch effects



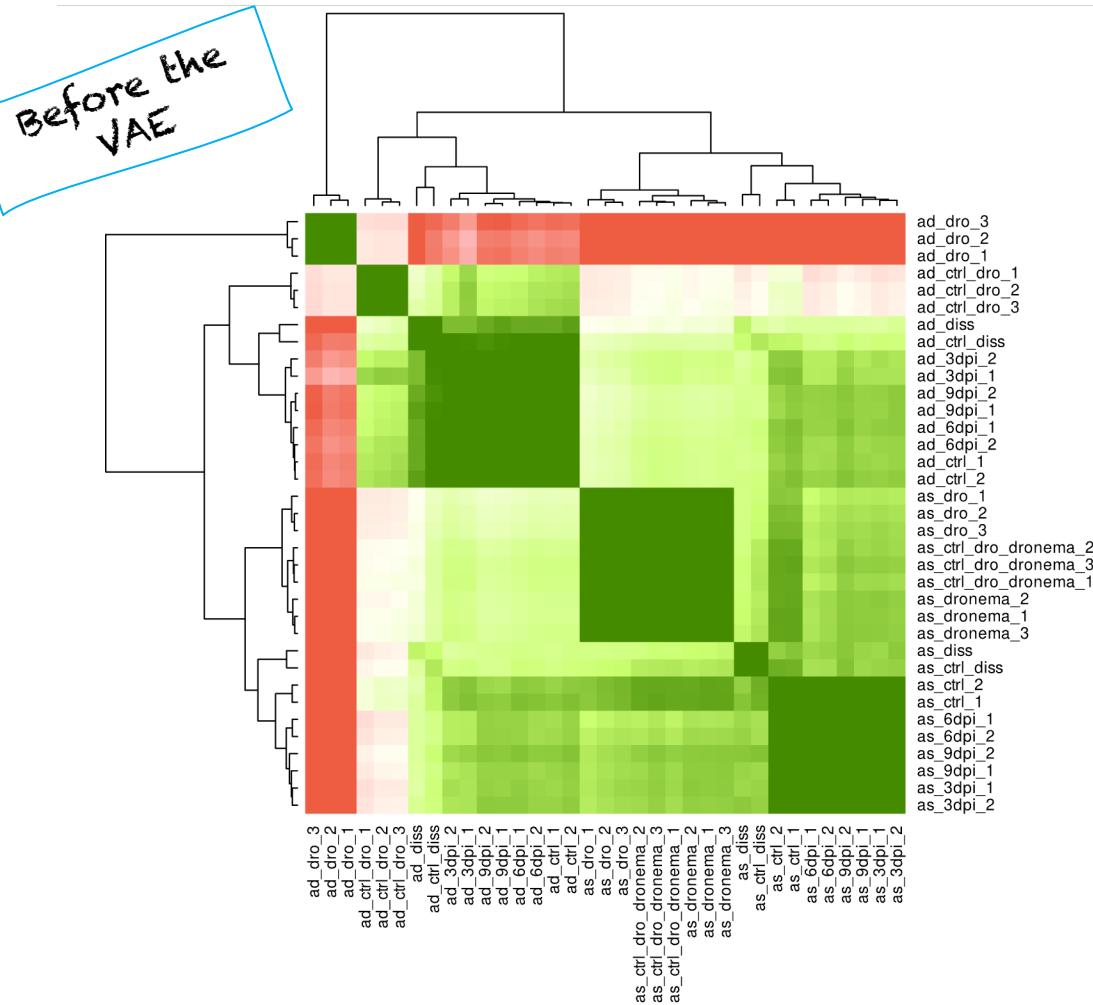
2 Opening the black box



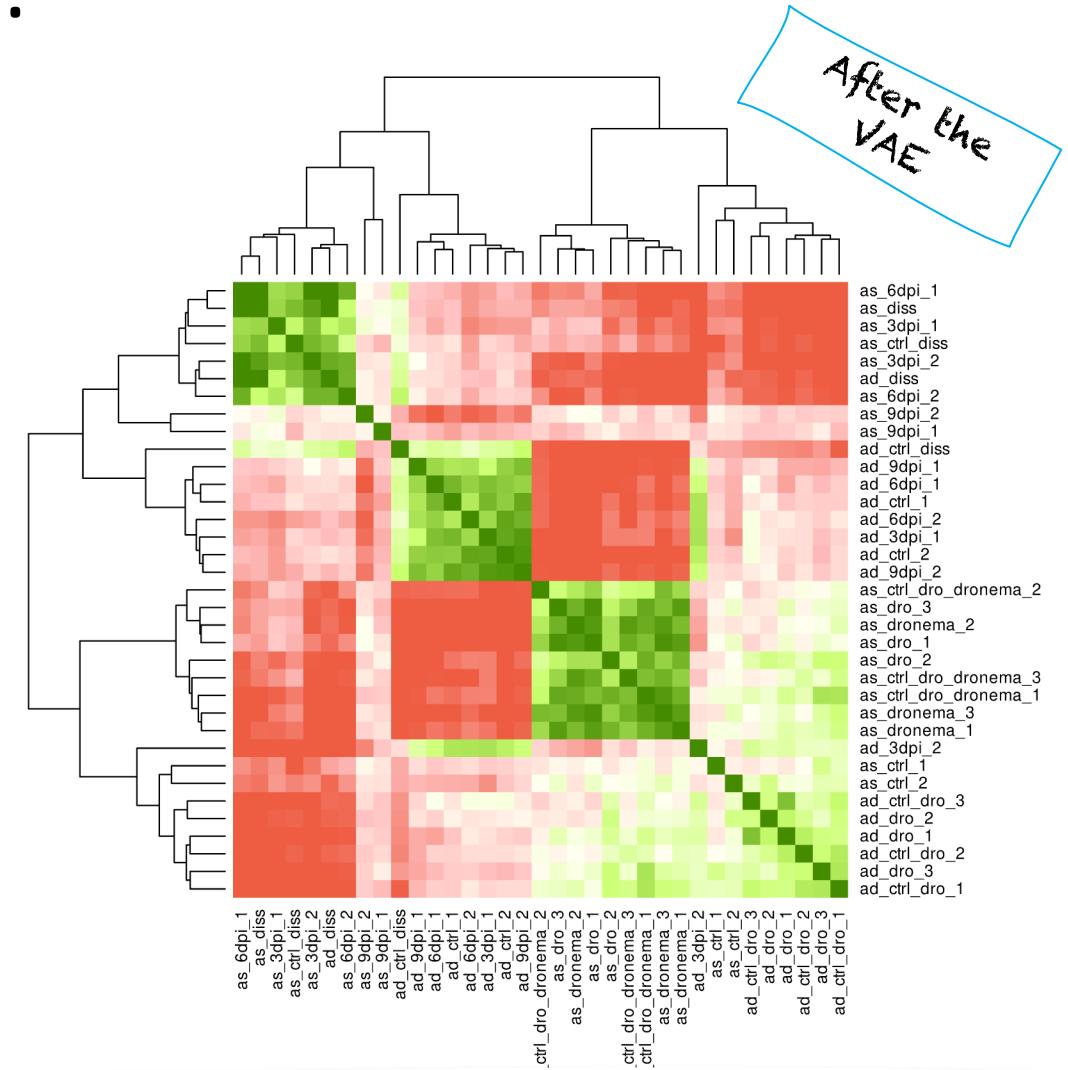
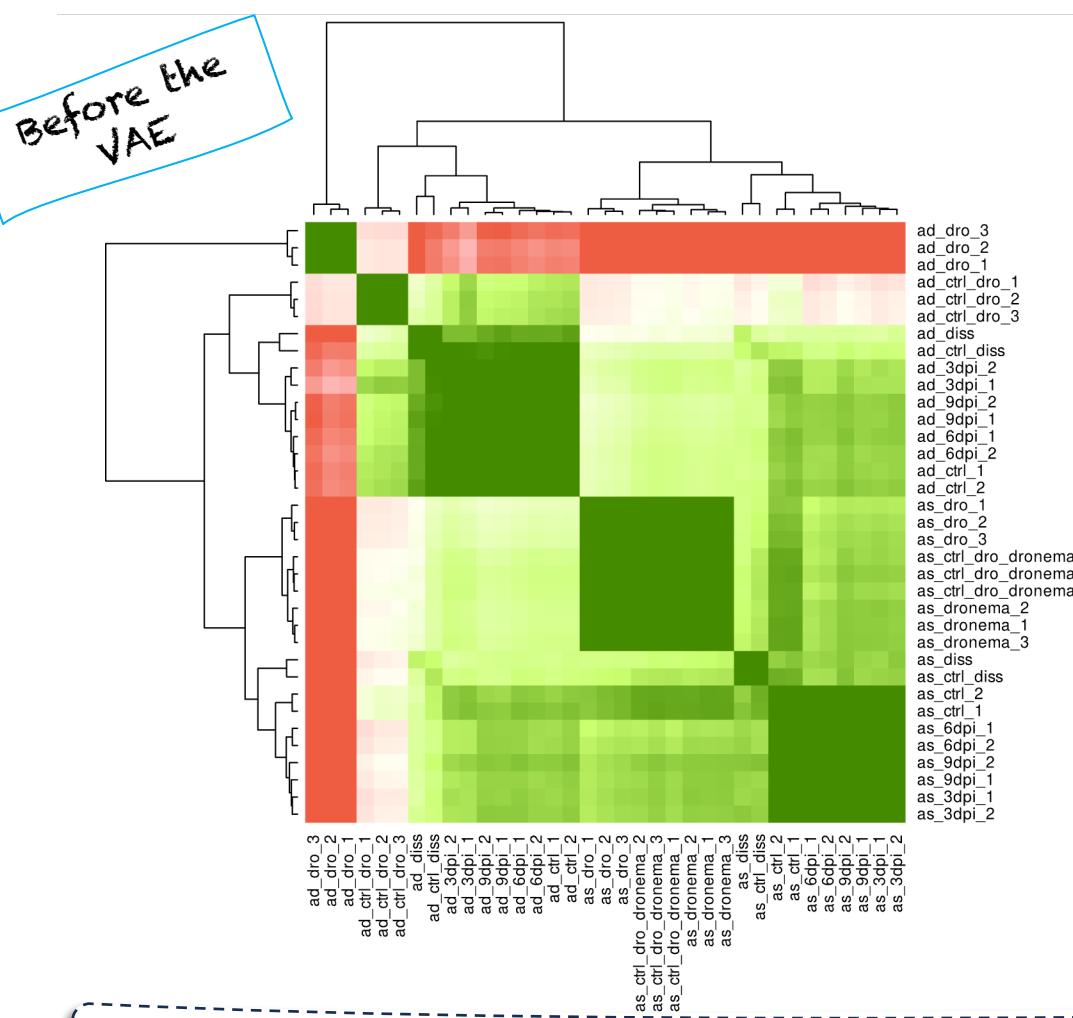
3 Finding important genes



Why do we need the VAE?



Why do we need the VAE?

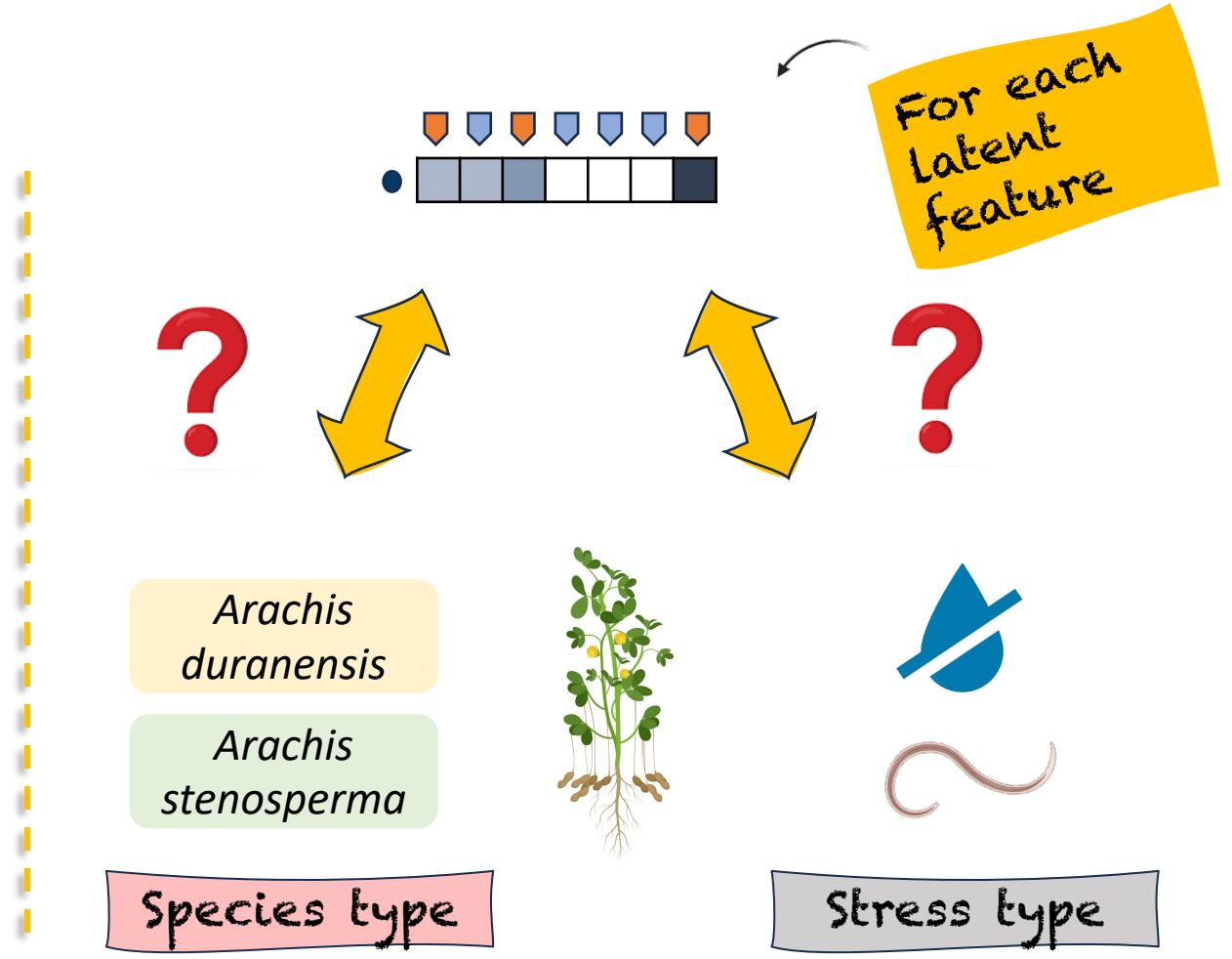
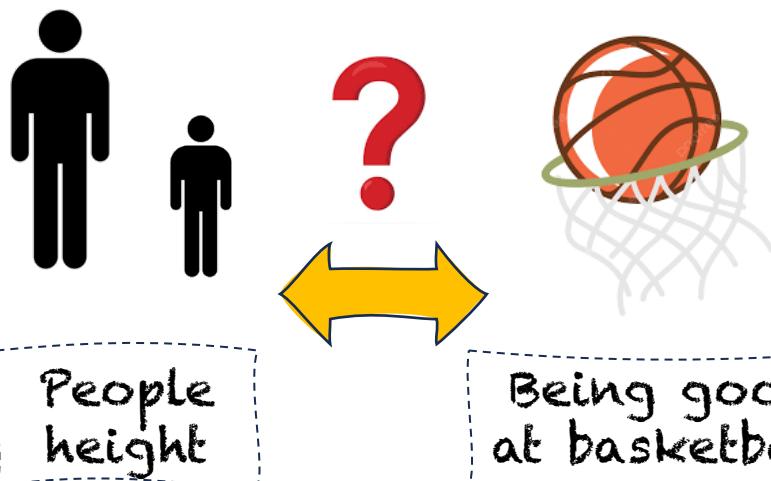


The VAE reduces the batch effects allowing unpaired experiments integrated analysis.

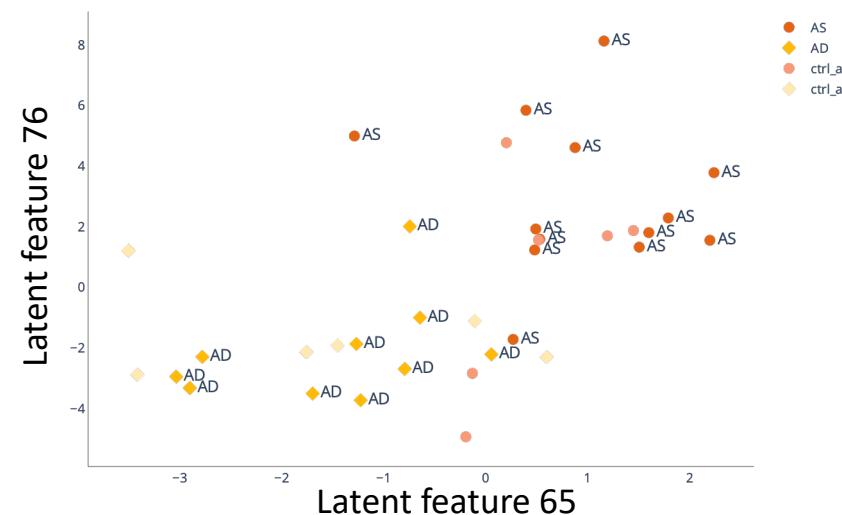
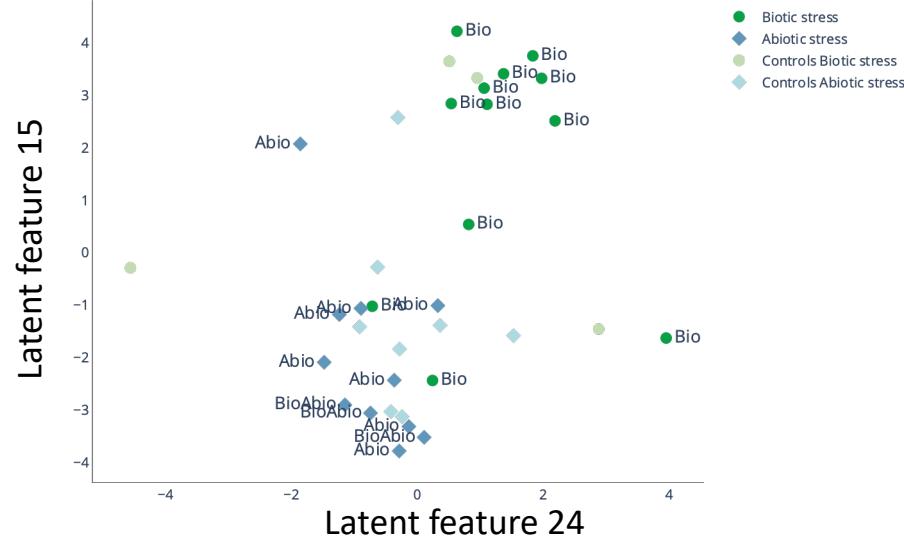
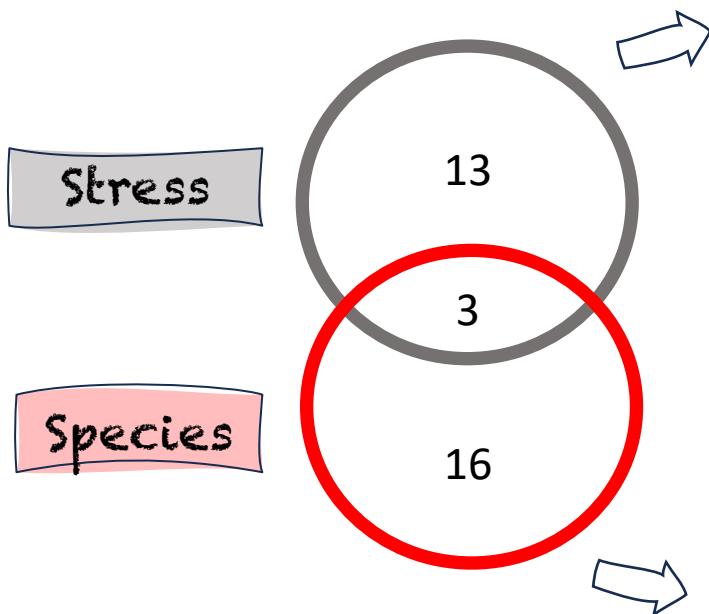
Which is the meaning of the latent features?

Point bi-serial correlation

is a special case of the Pearson Correlation and is used when you want to measure the relationship between a continuous variable and a dichotomous (or binary) variable.



Latent features are associated with phenotypic characteristics

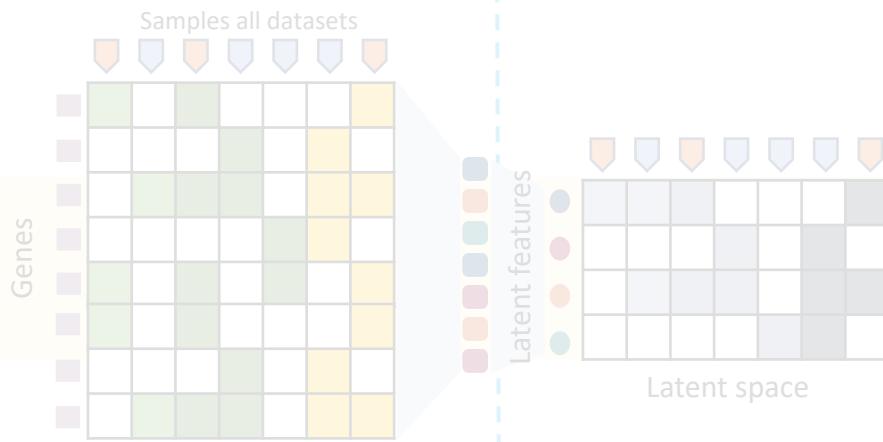


Top 2 associated latent features

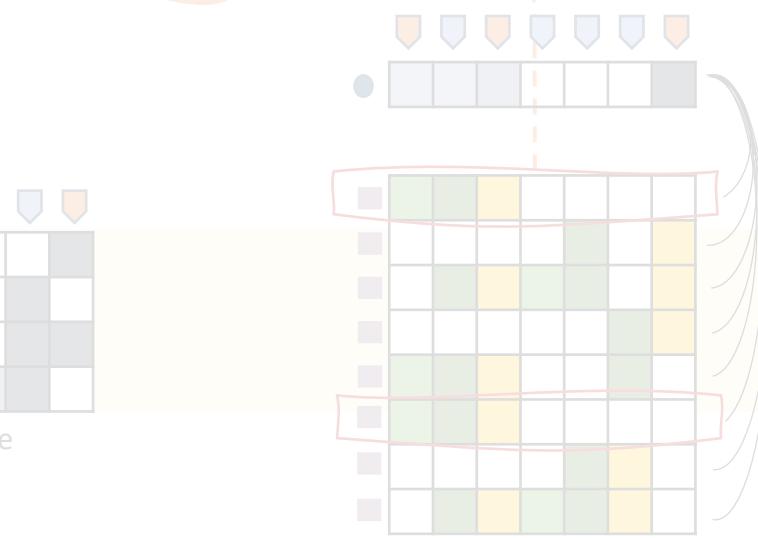


HIVE pipeline

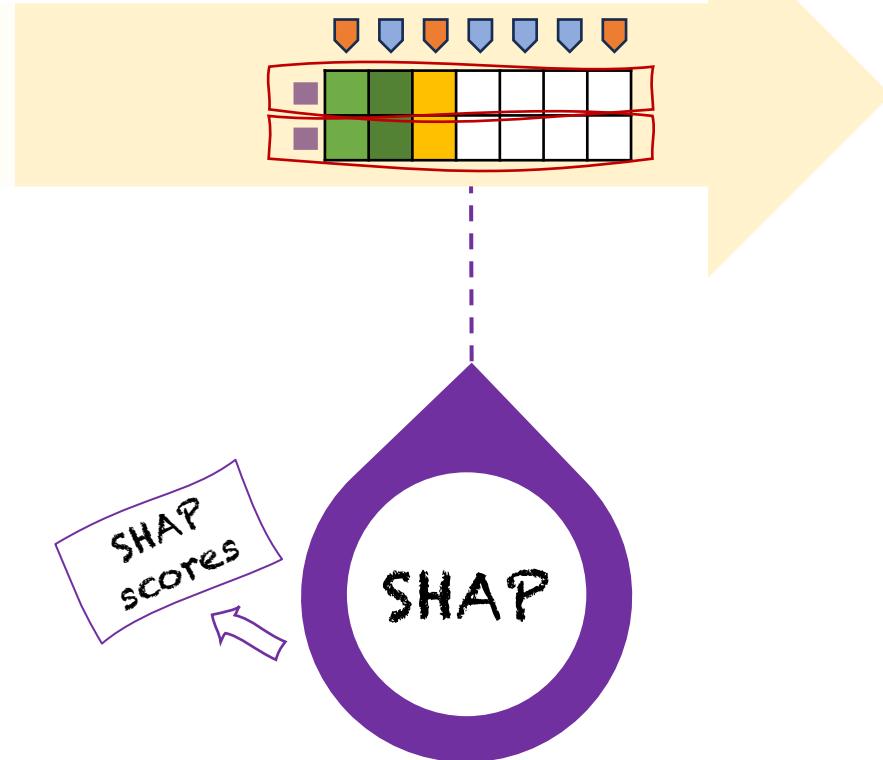
1 Removing batch effects



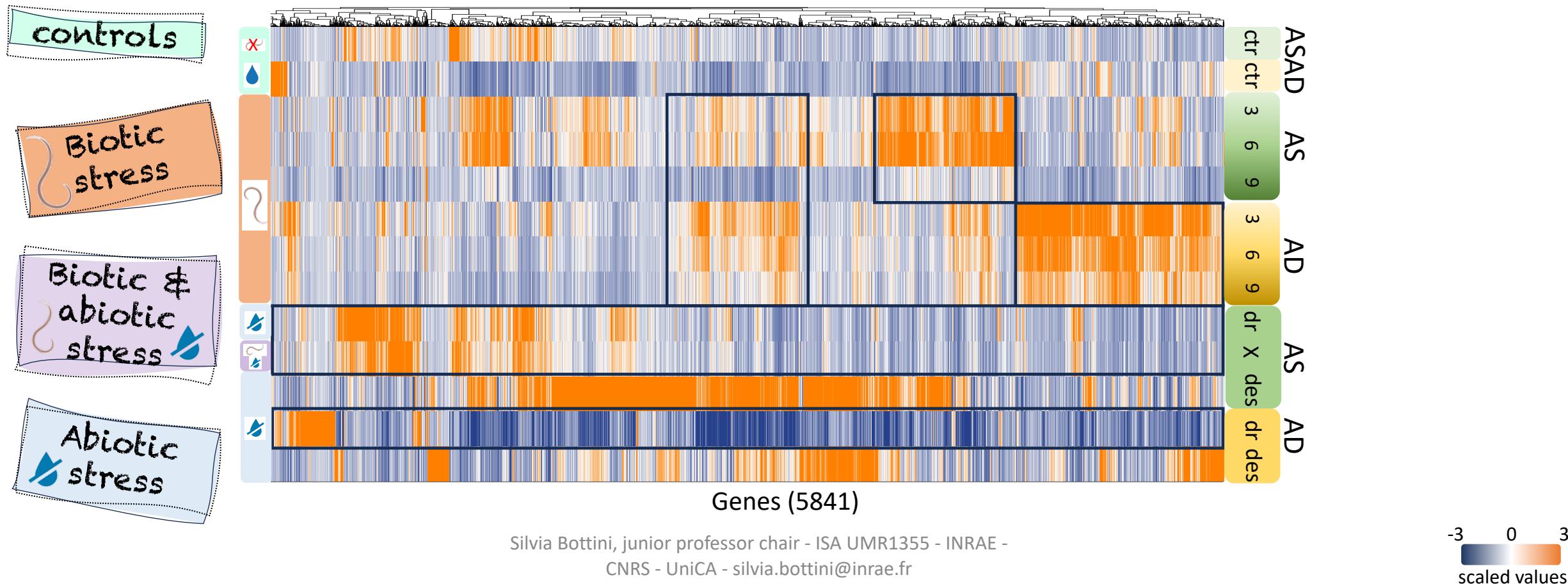
2 Opening the black box



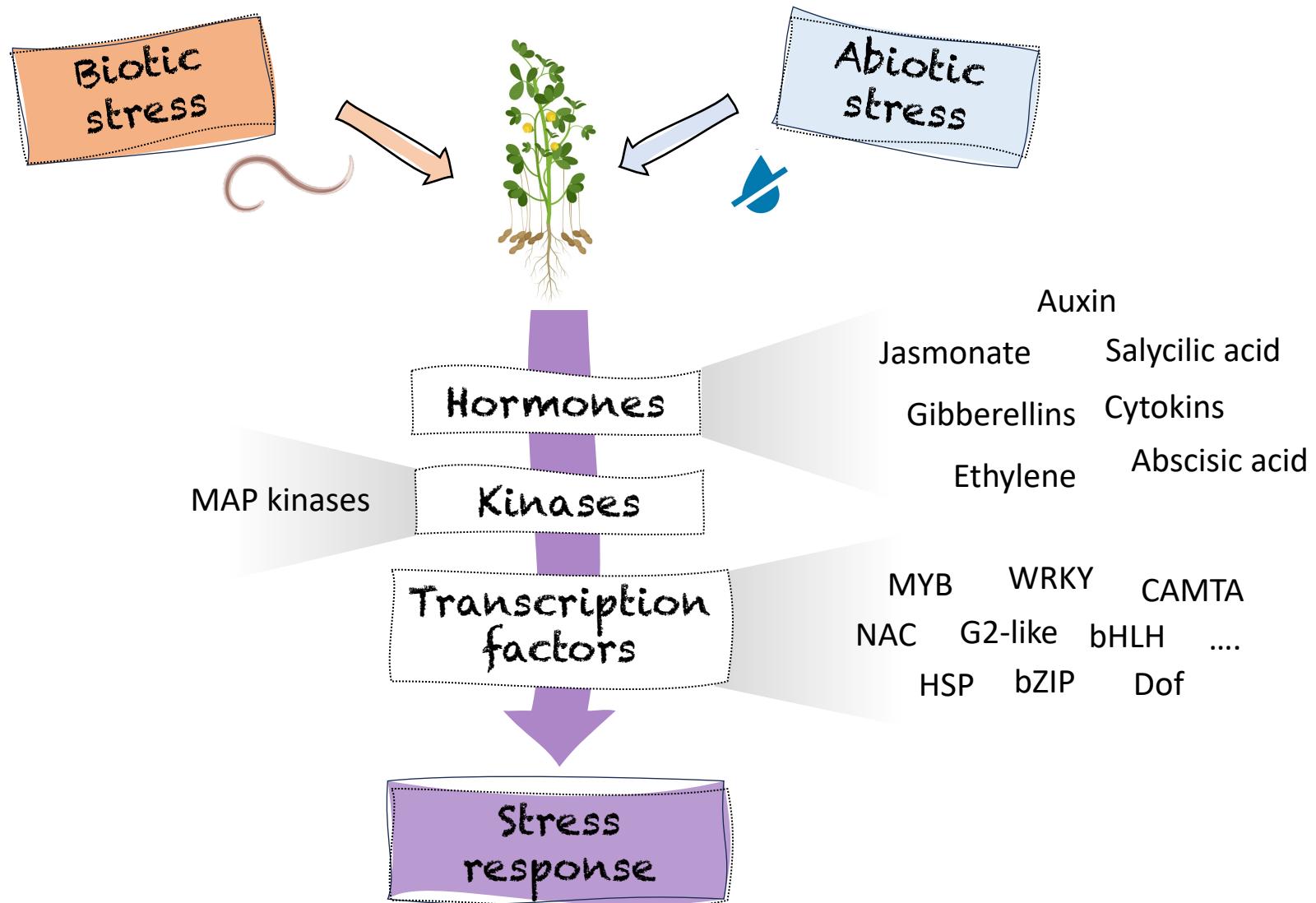
3 Finding important genes



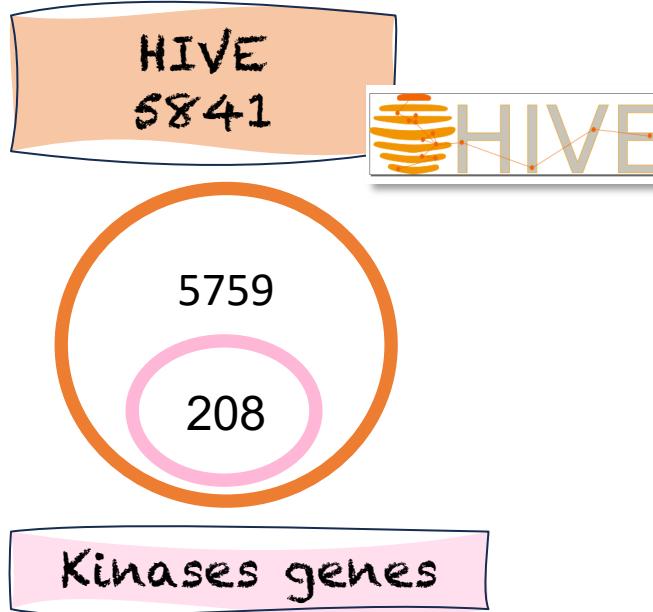
Global expression profiles of genes selected by HIVE



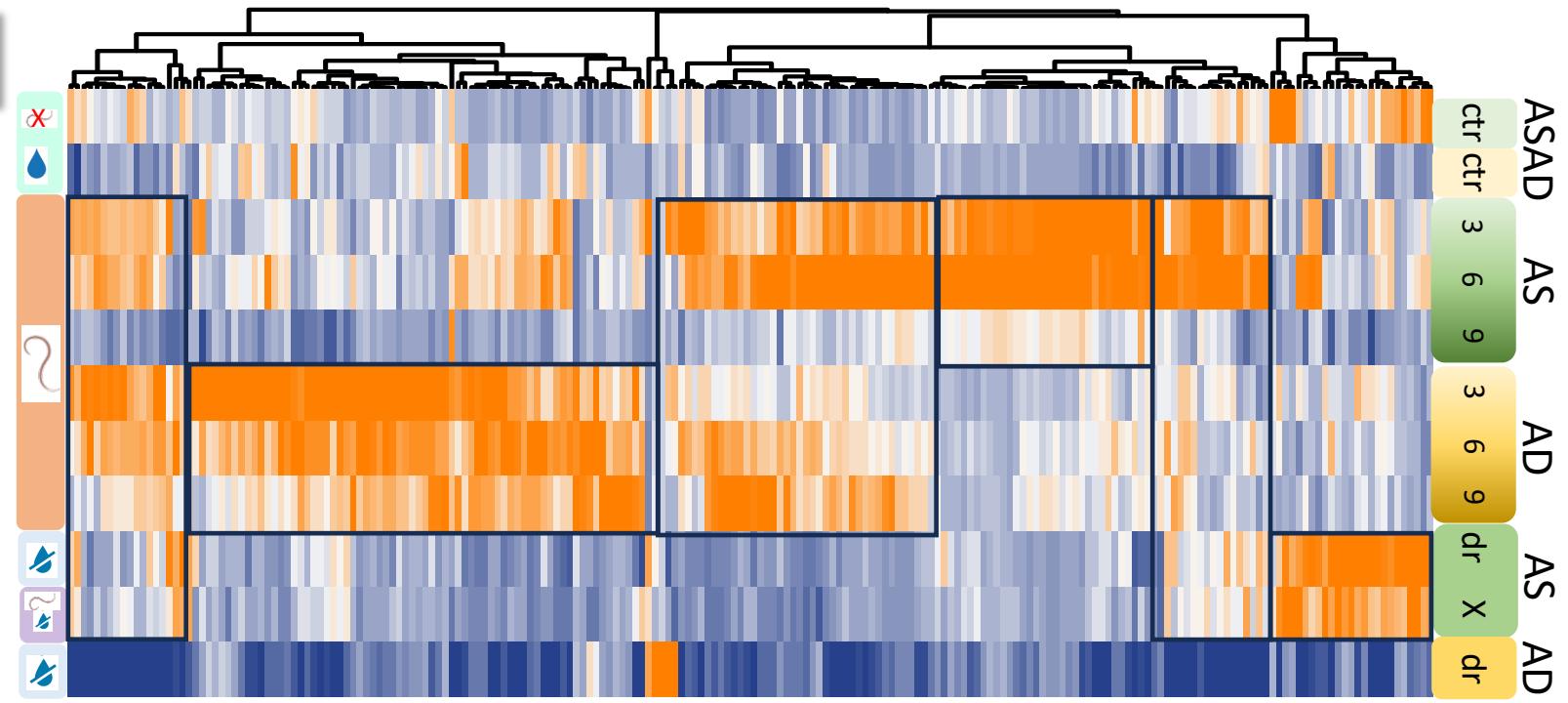
Plant molecular response to stress



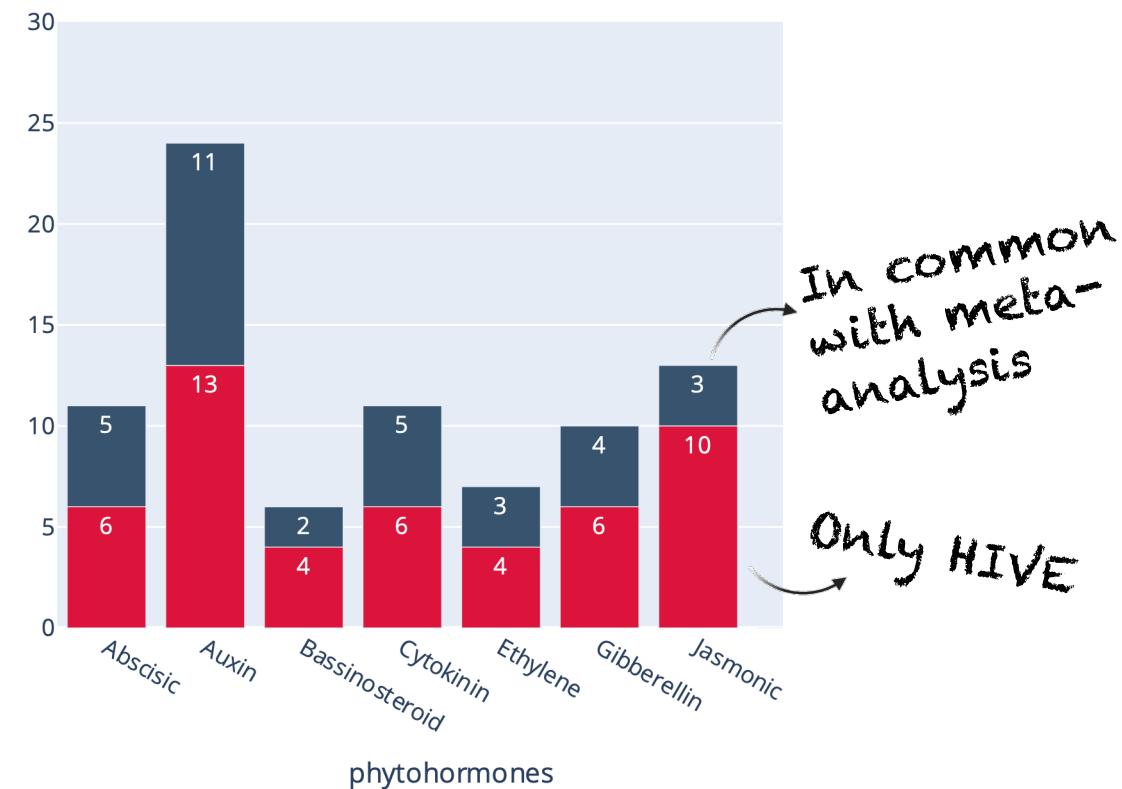
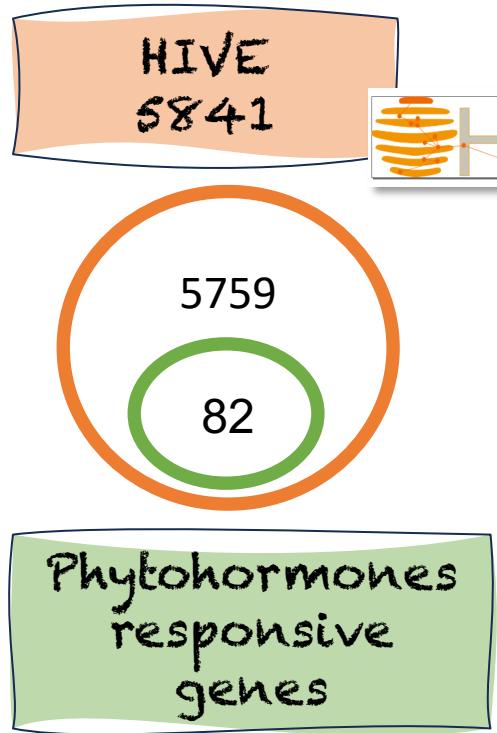
HIVE identifies genes coding for kinases



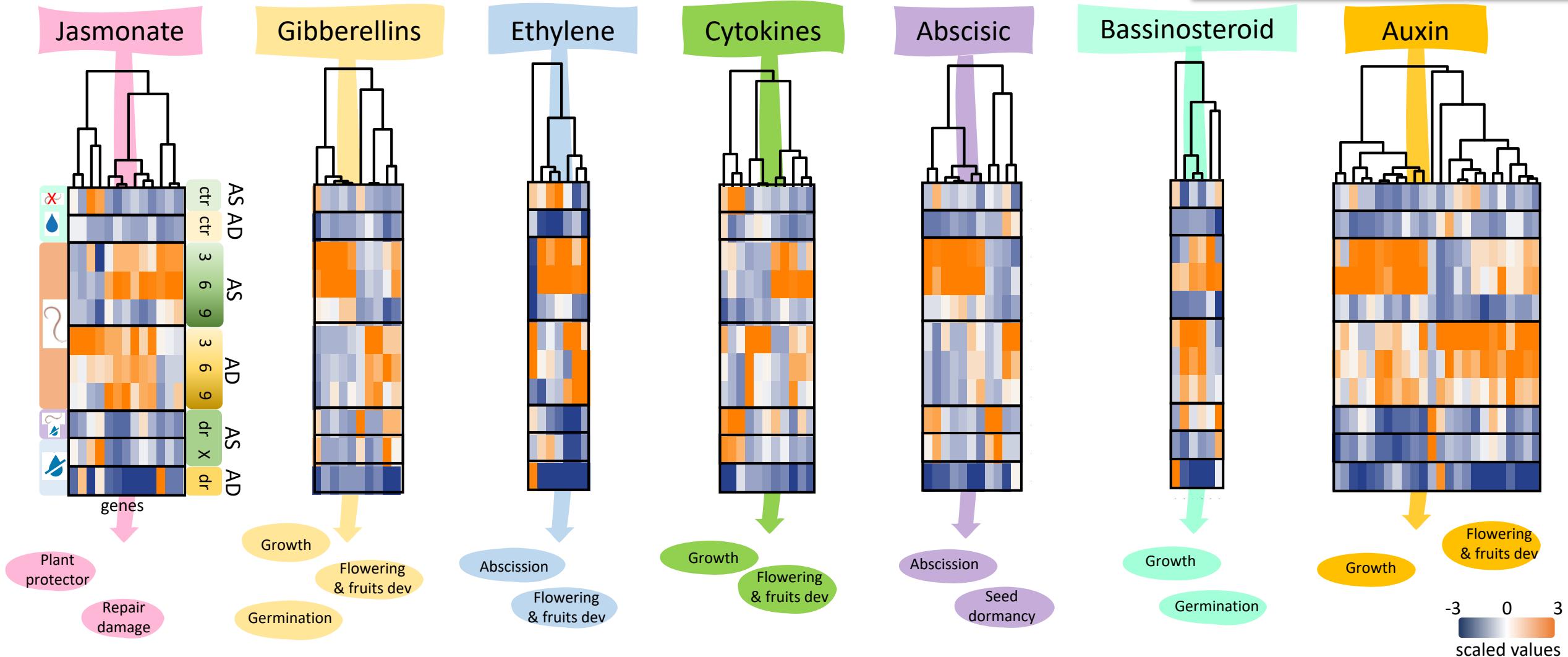
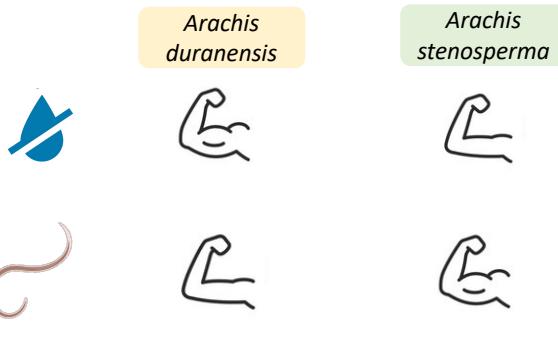
Stress sensing and signal transduction by protein kinases play crucial roles in plant responses to different stress conditions



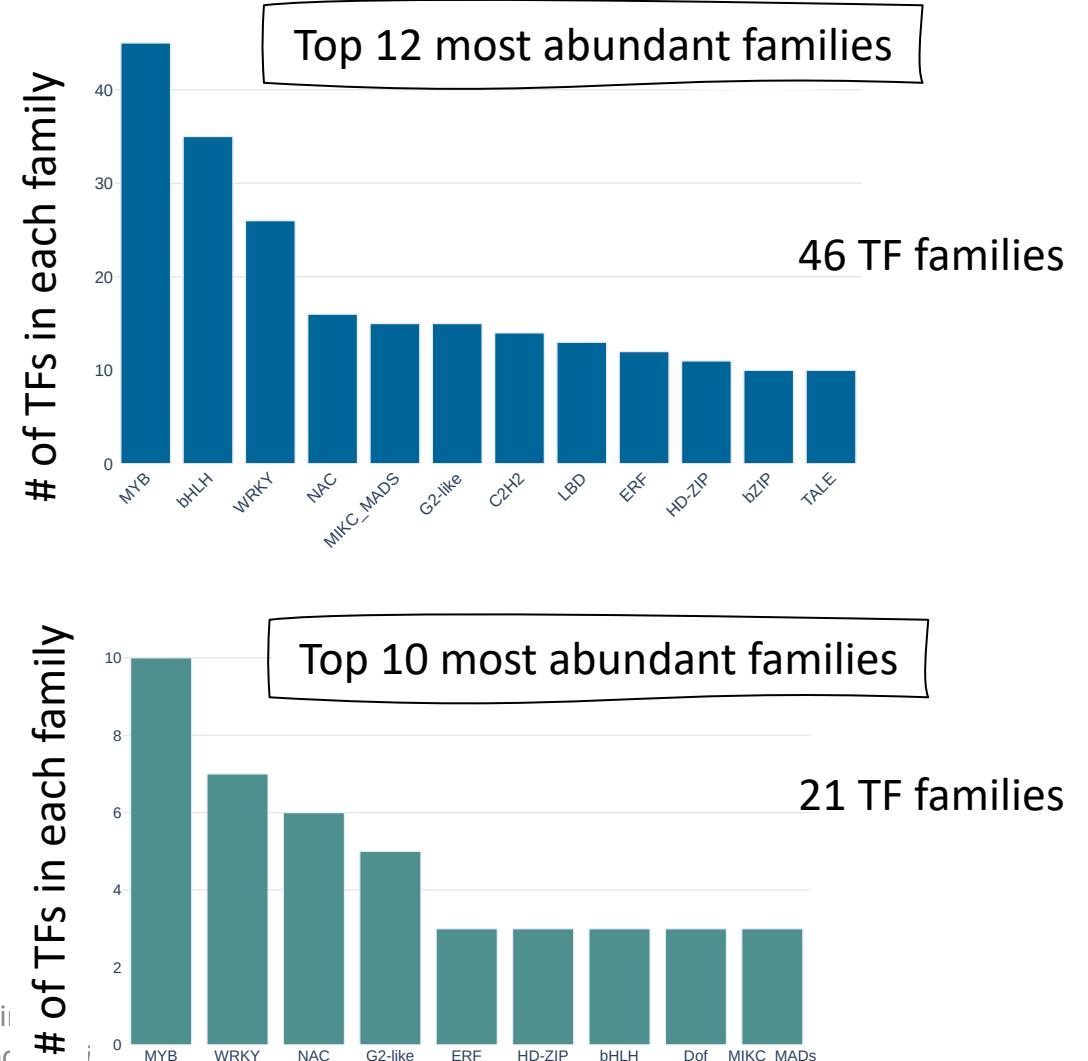
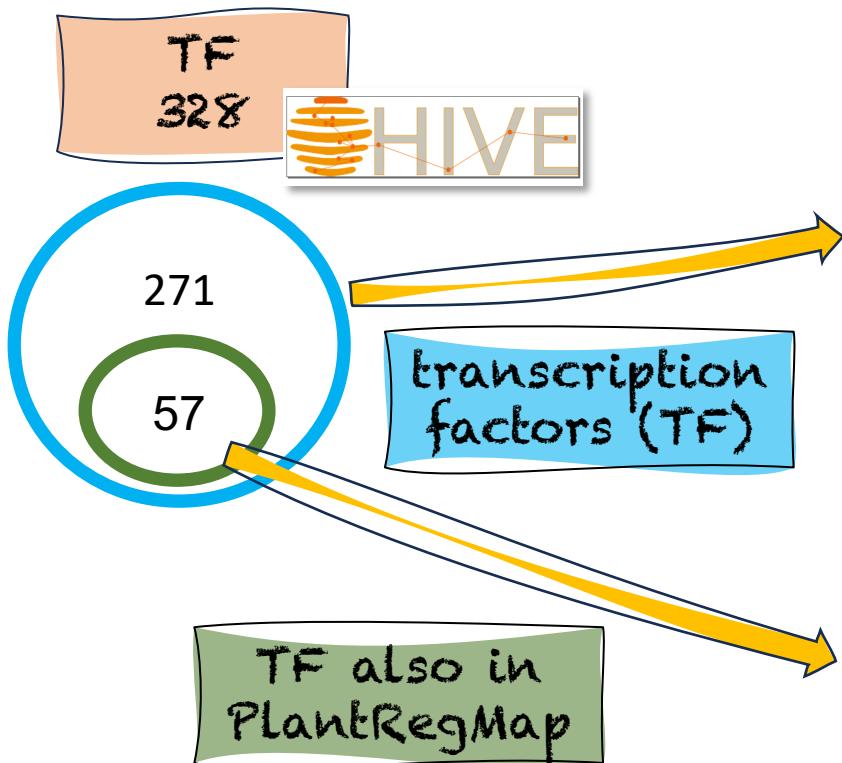
HIVE identifies genes involved in the main hormone signaling pathways in response to stress



Phytohormones responsive related genes expression profiles



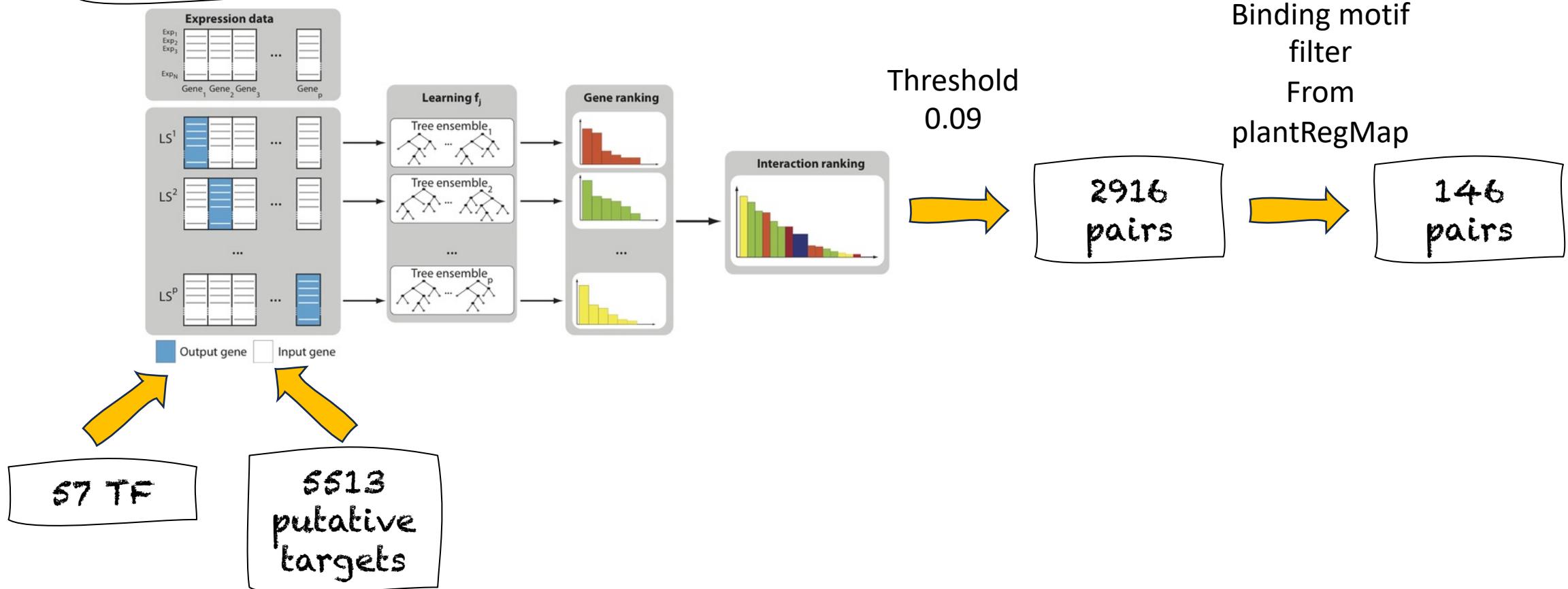
The integrated analysis allowed to retrieve transcription factors induced upon biotic and/or abiotic stress

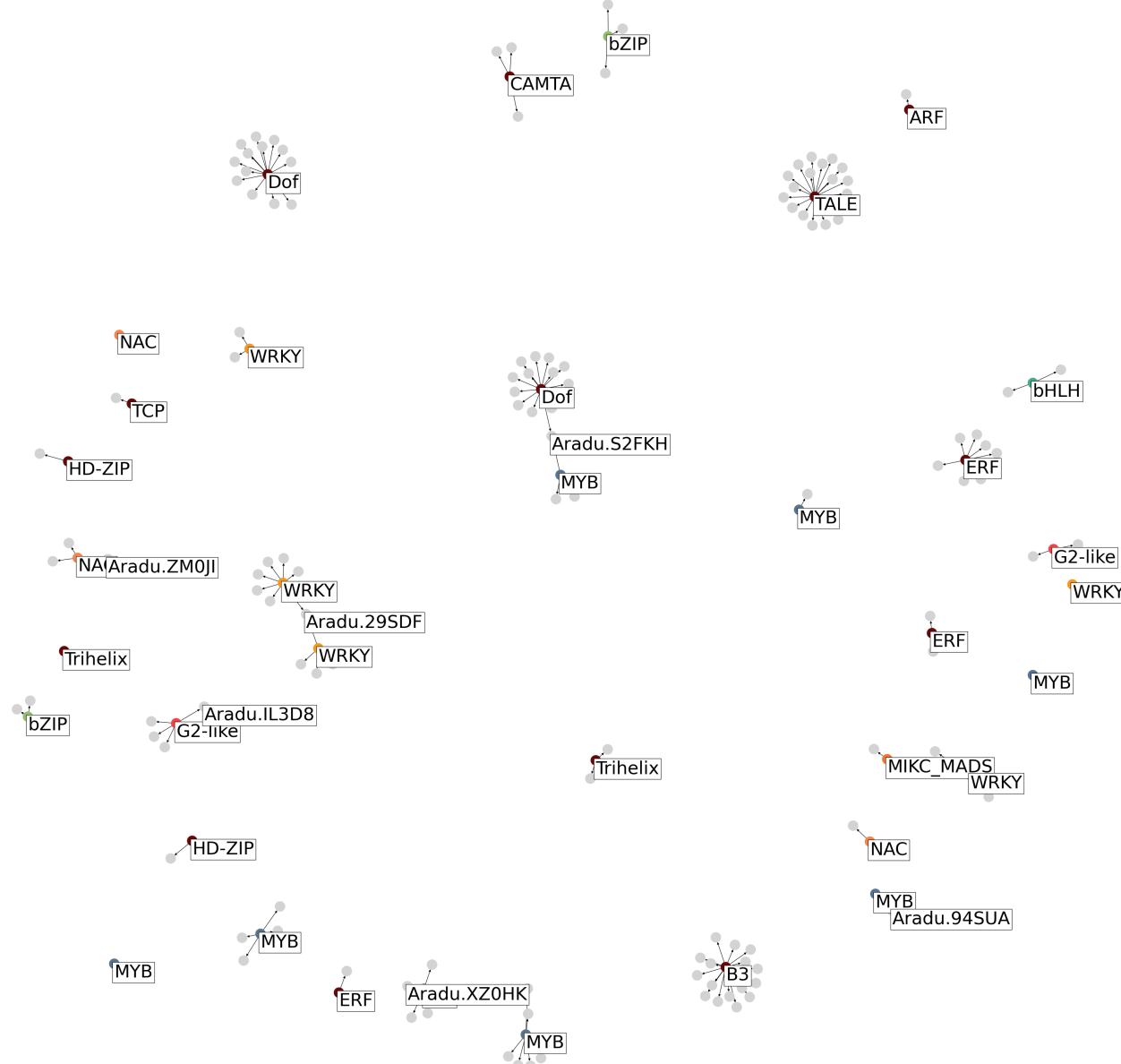


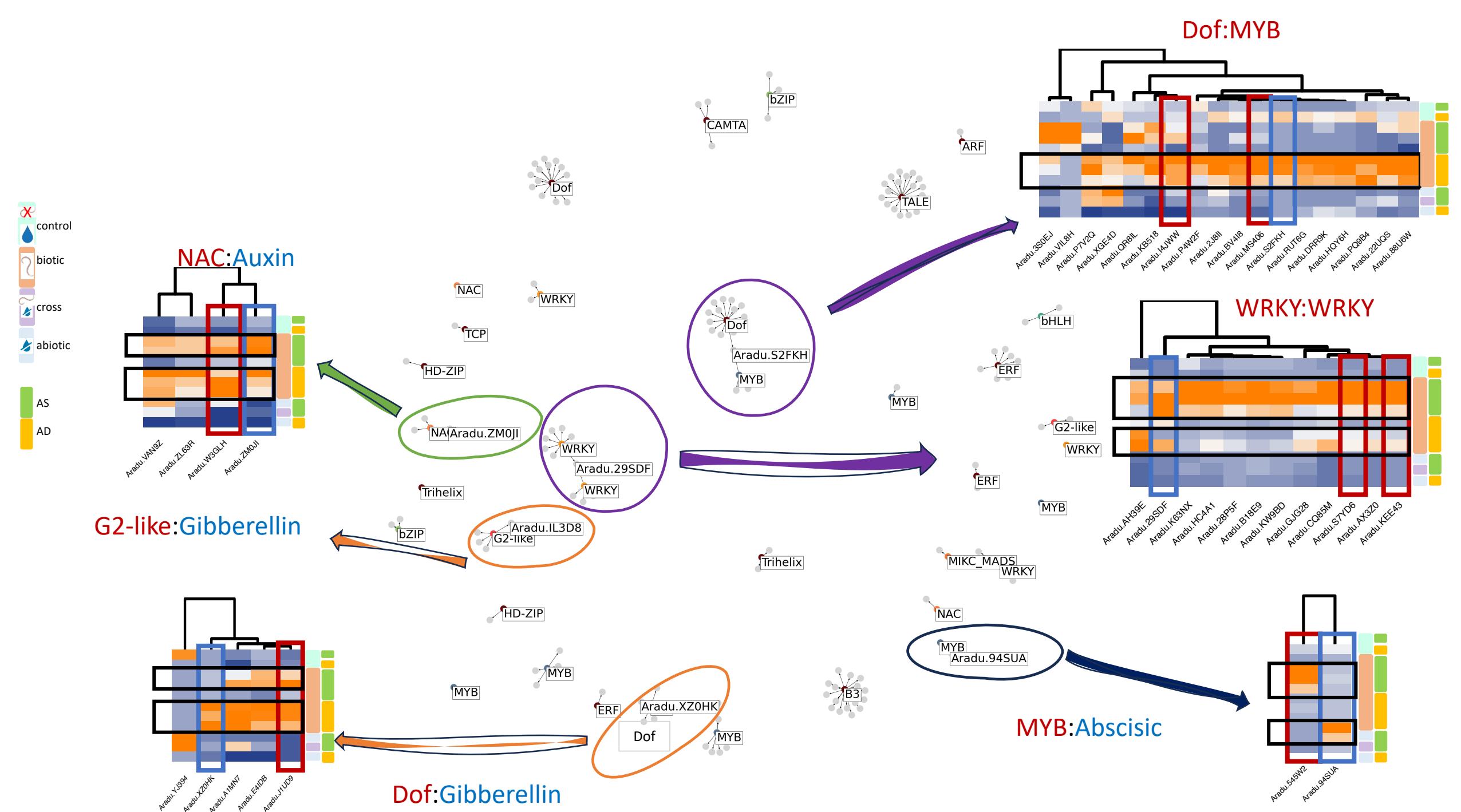
Finding THE targets

GENIE3

Huynh-Thu et al. Plos One 2010



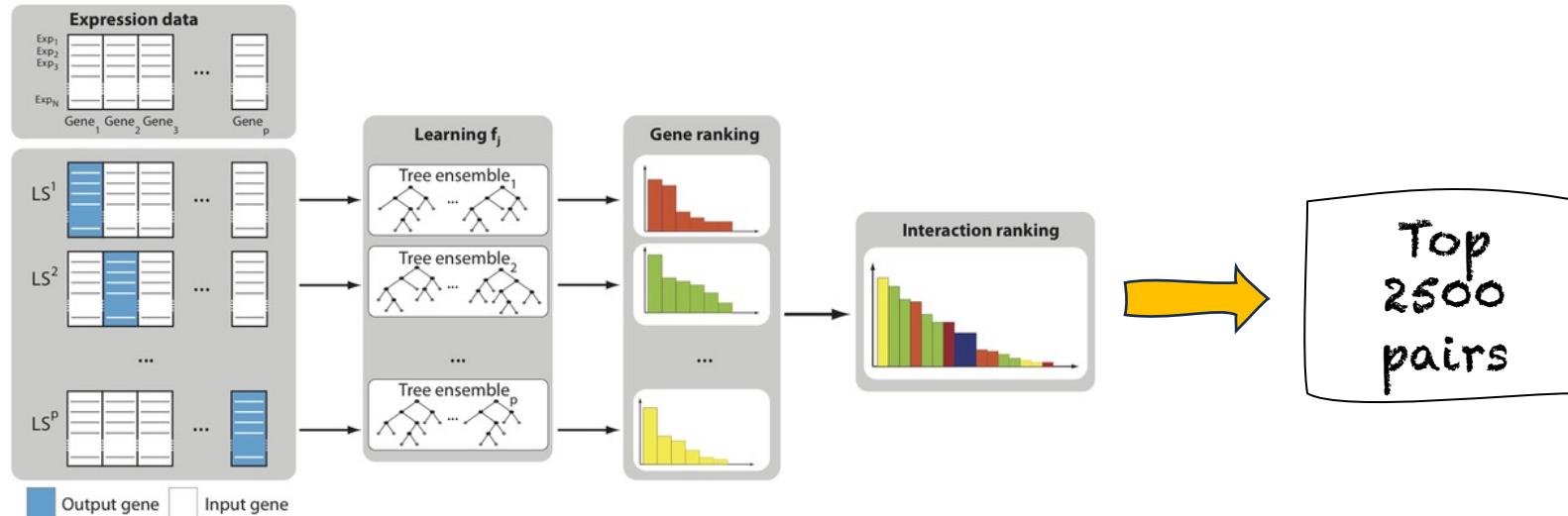




Finding THE targets

GENIE3

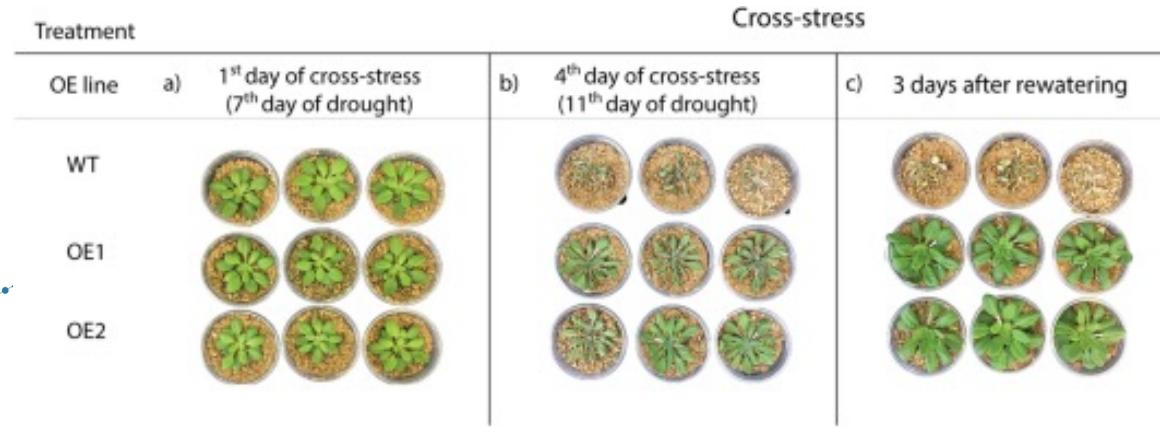
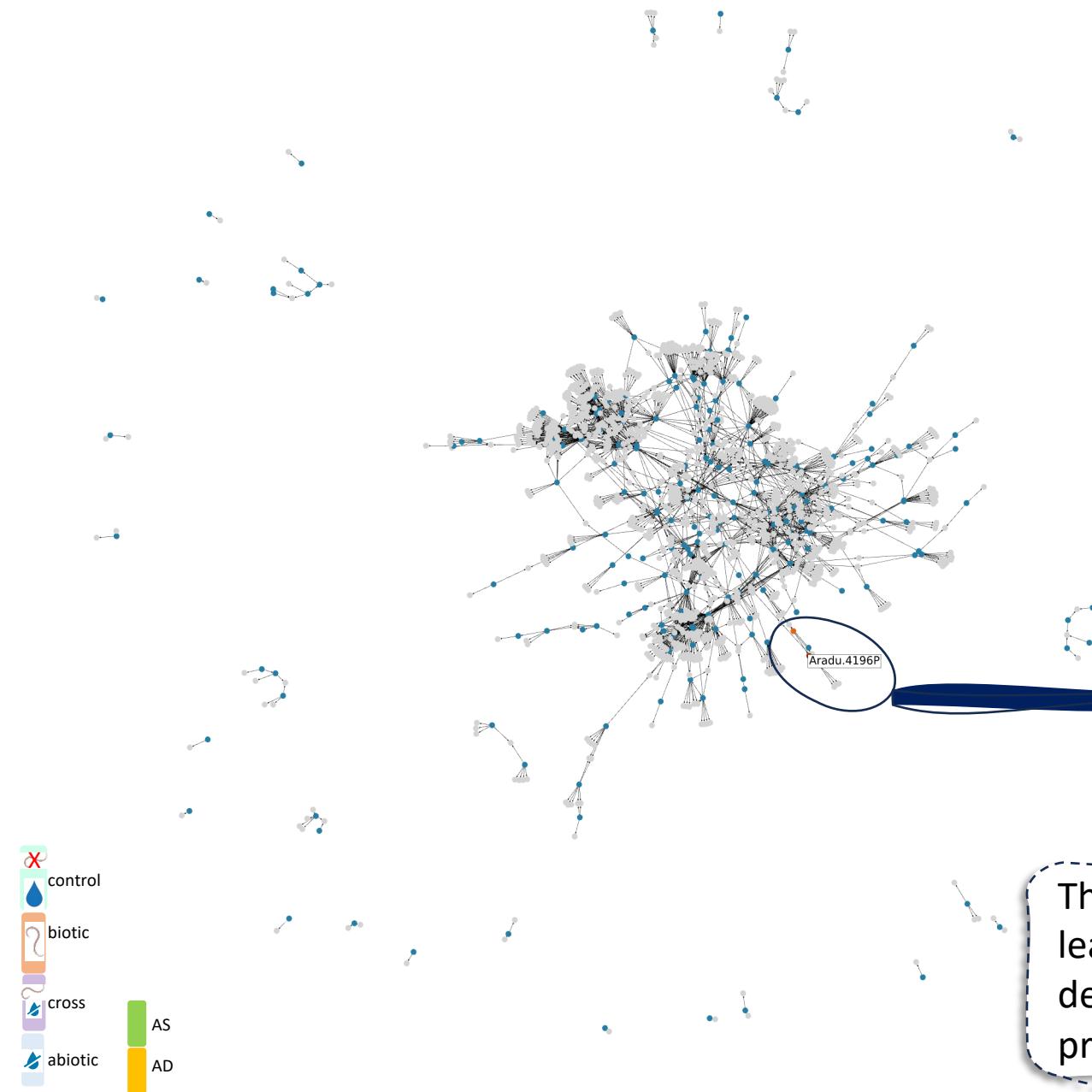
Huynh-Thu et al. Plos One 2010



271 TF

5513
putative
targets

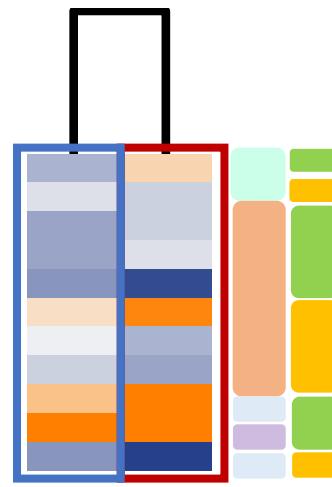




Overexpression *in planta* showed high resistance to both stresses, separately or combined.

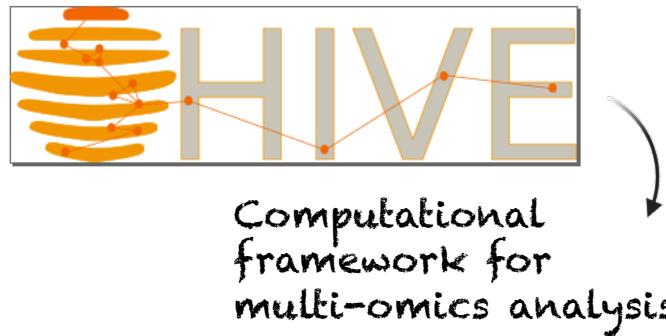
G2-like:Endochitinase

The induction of chitinases by specific agents leads to simultaneous enhancement of other defense reactions and pathogenesis related proteins in the same plant tissues.

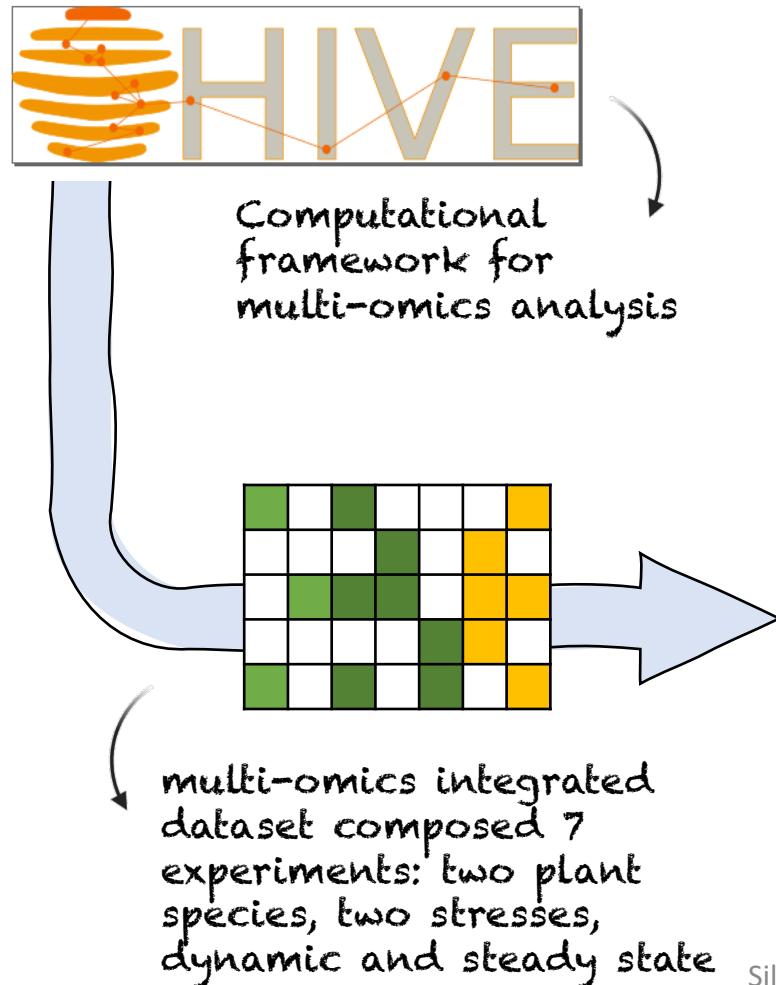


Aradu.4196P
Aradu.Y7NUP

Conclusions



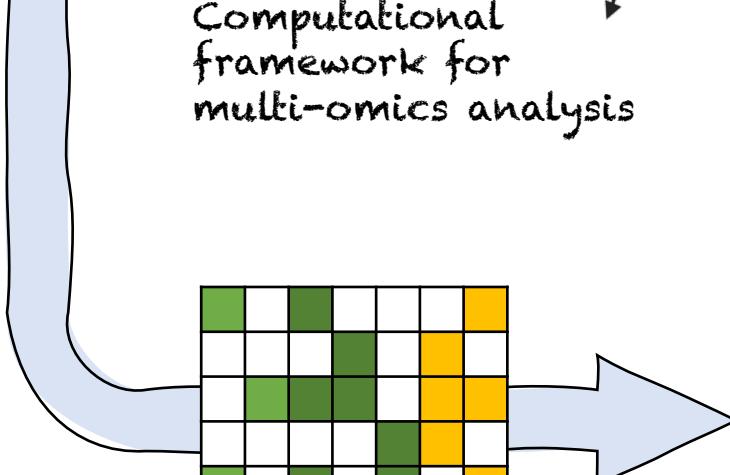
Conclusions



Conclusions



Computational
framework for
multi-omics analysis



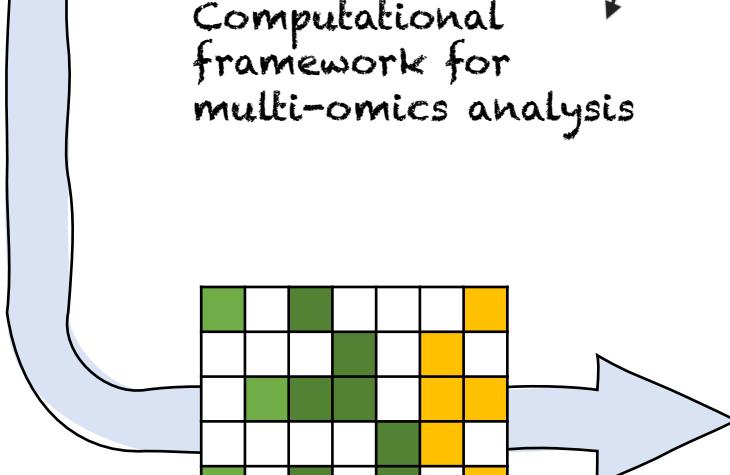
multi-omics integrated
dataset composed 7
experiments: two plant
species, two stresses,
dynamic and steady state

Identification of **5841 genes**
deregulated in at least one
condition, including **several**
experimentally validated genes.

Conclusions



Computational
framework for
multi-omics analysis



multi-omics integrated
dataset composed 7
experiments: two plant
species, two stresses,
dynamic and steady state

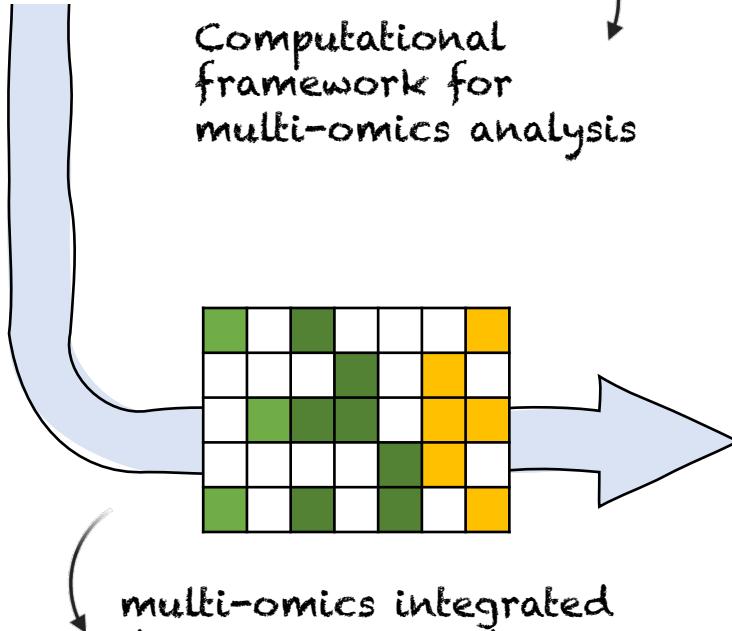
Identification of **5841 genes**
deregulated in at least one
condition, including **several
experimentally validated genes.**

Identification of several **genes**
connected with **plant-stress
response** (phytohormones,
kinases, transcription factors,
R-genes).

Conclusions



Computational framework for multi-omics analysis



multi-omics integrated dataset composed of 7 experiments: two plant species, two stresses, dynamic and steady state

Identification of **5841 genes** deregulated in at least one condition, including **several experimentally validated genes**.

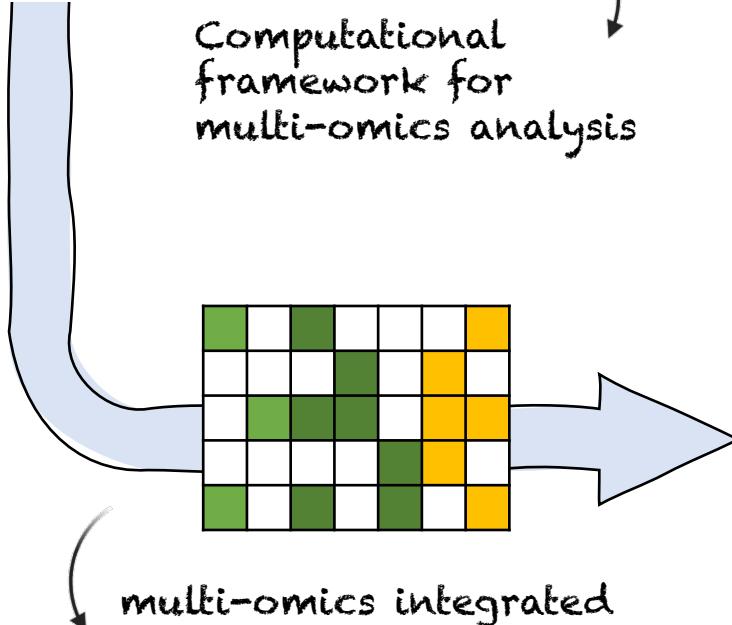
Identification of several genes connected with **plant-stress response** (phytohormones, kinases, transcription factors, R-genes).

Identification of **potential candidates** to better understand the responses to stresses.

Conclusions



Computational
framework for
multi-omics analysis



multi-omics integrated
dataset composed of
7 experiments: two plant
species, two stresses,
dynamic and steady state

Identification of **5841 genes**
deregulated in at least one
condition, including **several
experimentally validated genes**.

Identification of several **genes**
connected with **plant-stress
response** (phytohormones,
kinases, transcription factors,
R-genes).

Identification of **potential
candidates** to better
understand the responses
to stresses.

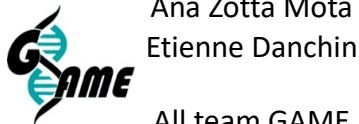
HIVE can be applied on
any biological model!

Acknowledgments



Giulia Calia
Justine Labory

All team M2P2



Ana Zotta Mota
Etienne Danchin

All team GAME

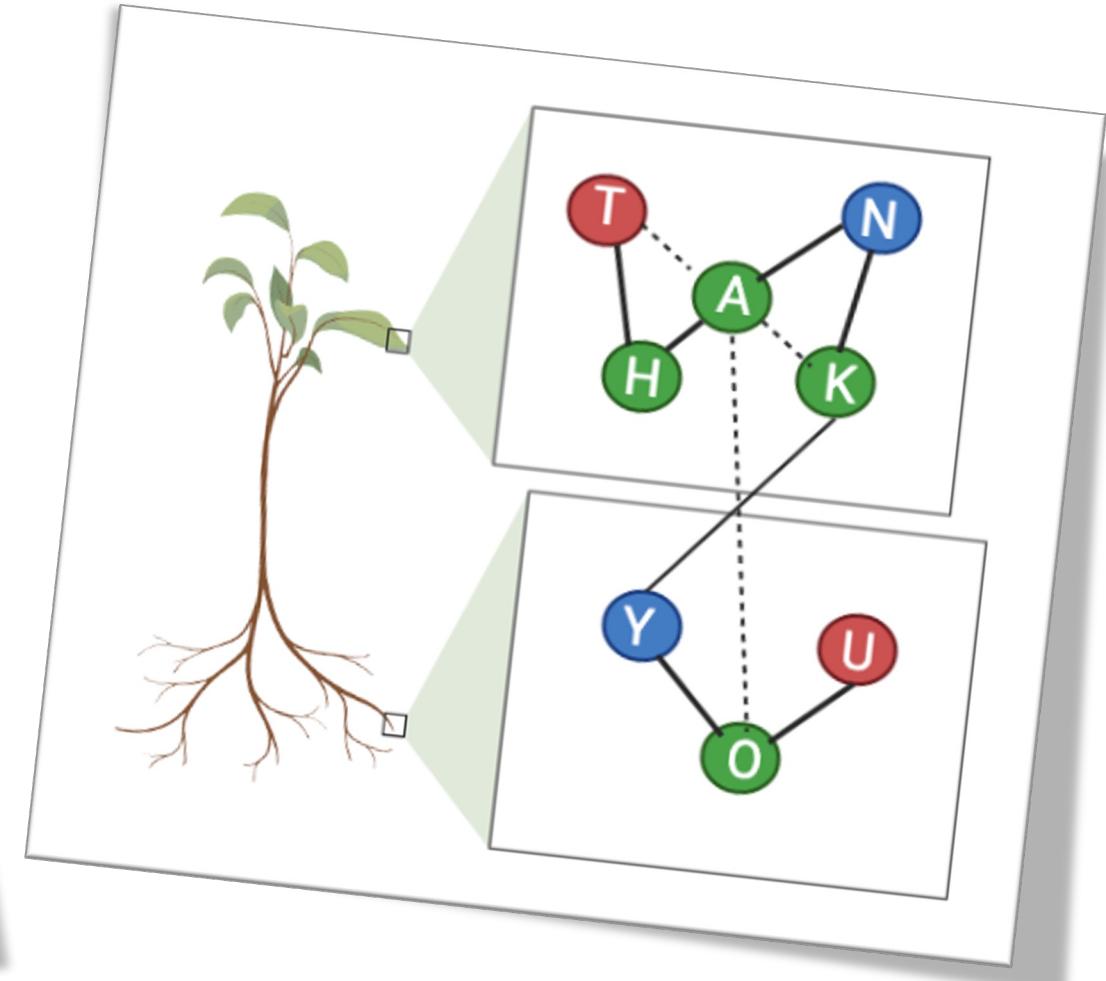


Patricia Messenberg
Guimaraes
& all the team



Le comité de pilotage du
métaprogramme DIGIT-BIO





Email: silvia.bottini@inrae.fr

Benchmarking tools for genes selection from integrated datasets



PLS-DA Partial Least Squares
Discriminant Analysis

Supervised method

Iterative method that constructs H successive artificial (latent) components, where each component is a linear combination of the variables.

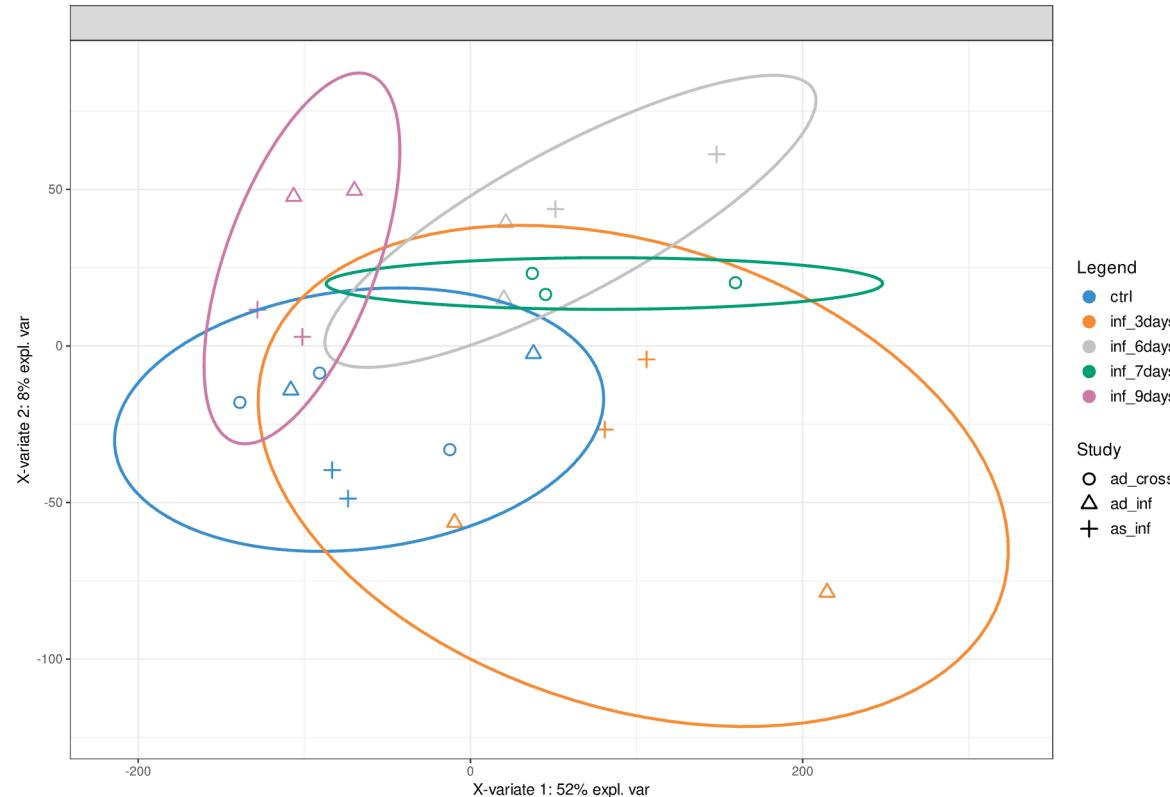
Benchmarking tools for genes selection from integrated datasets



PLS-DA Partial Least Squares
Discriminant Analysis

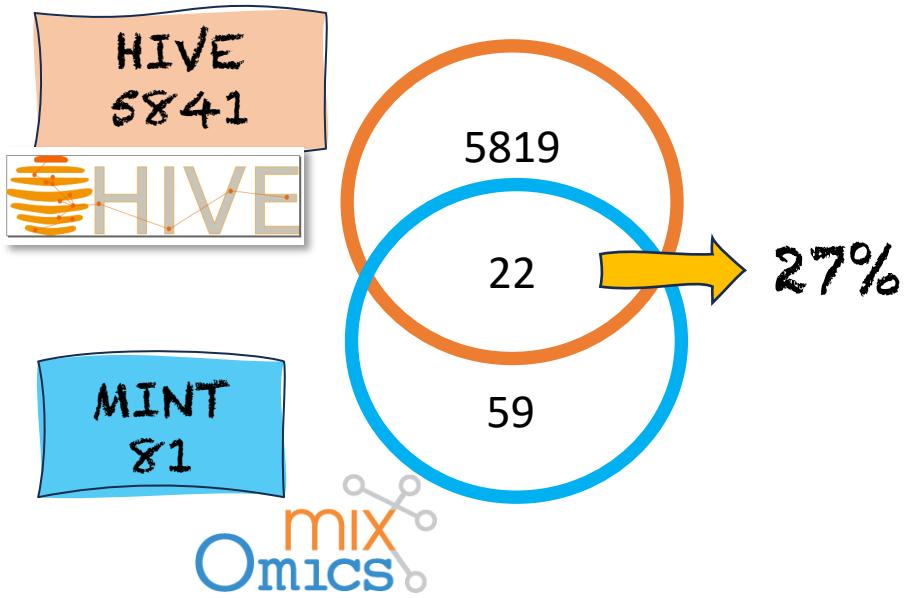
Supervised method

Iterative method that constructs H successive artificial (latent) components, where each component is a linear combination of the variables.



Component 1 = 22 genes
Component 2 = 59 genes
Total = 81 genes

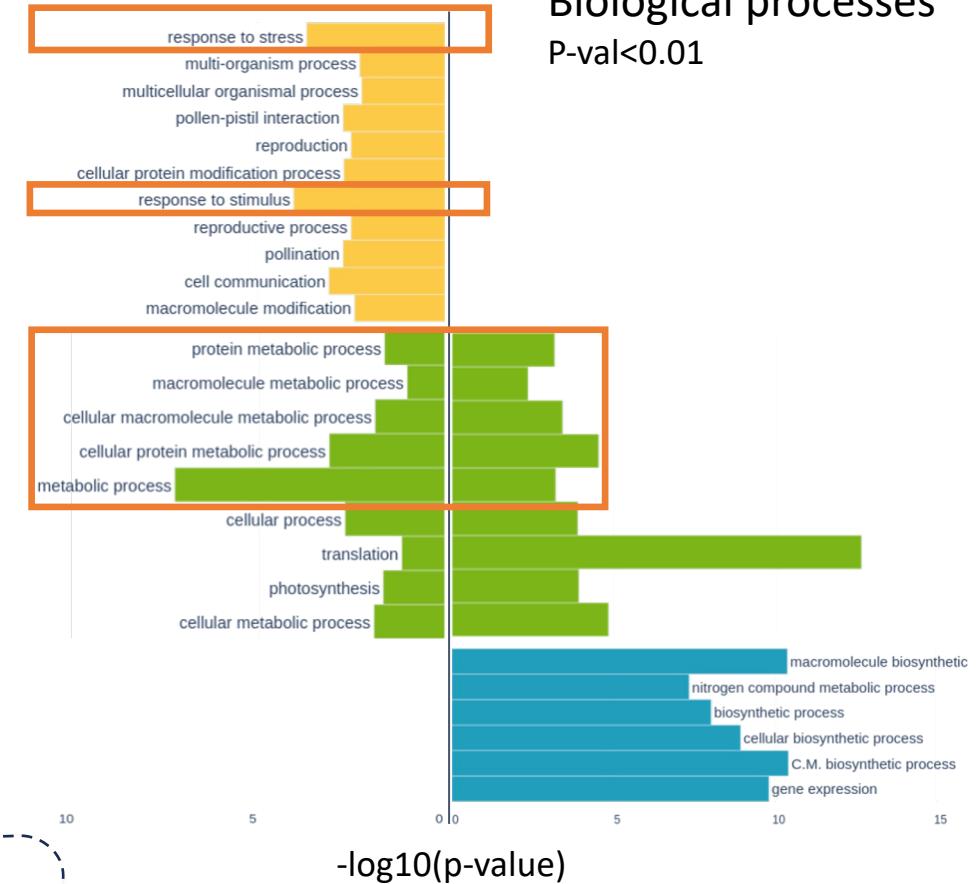
HIVE vs MINT



Only
HIVE

both

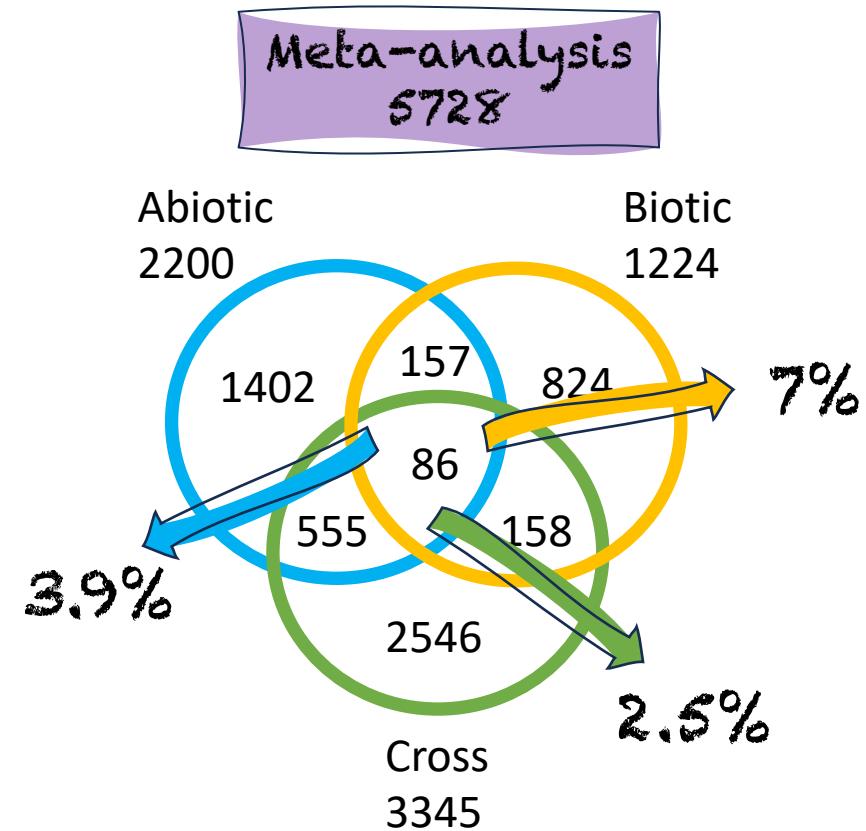
Only
MINT



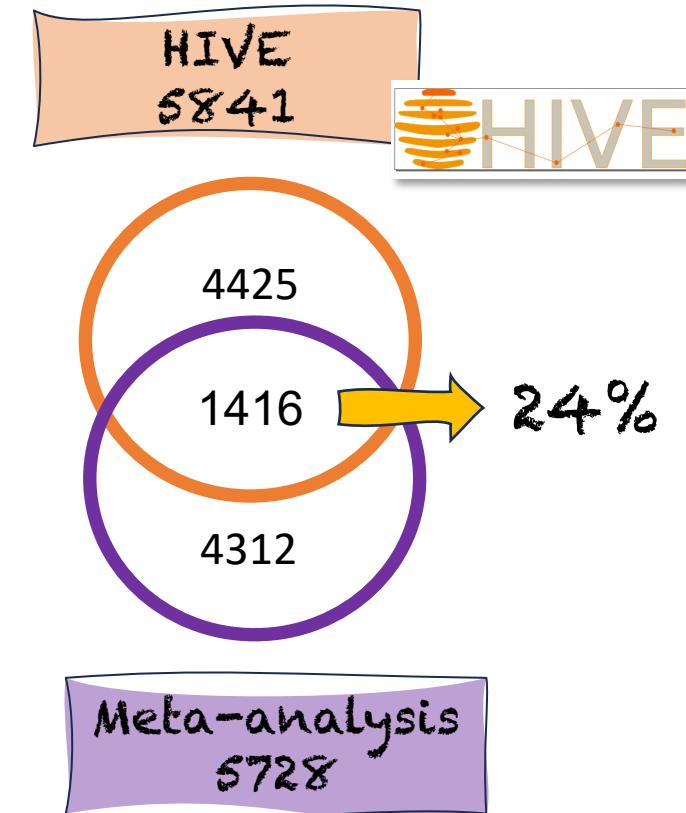
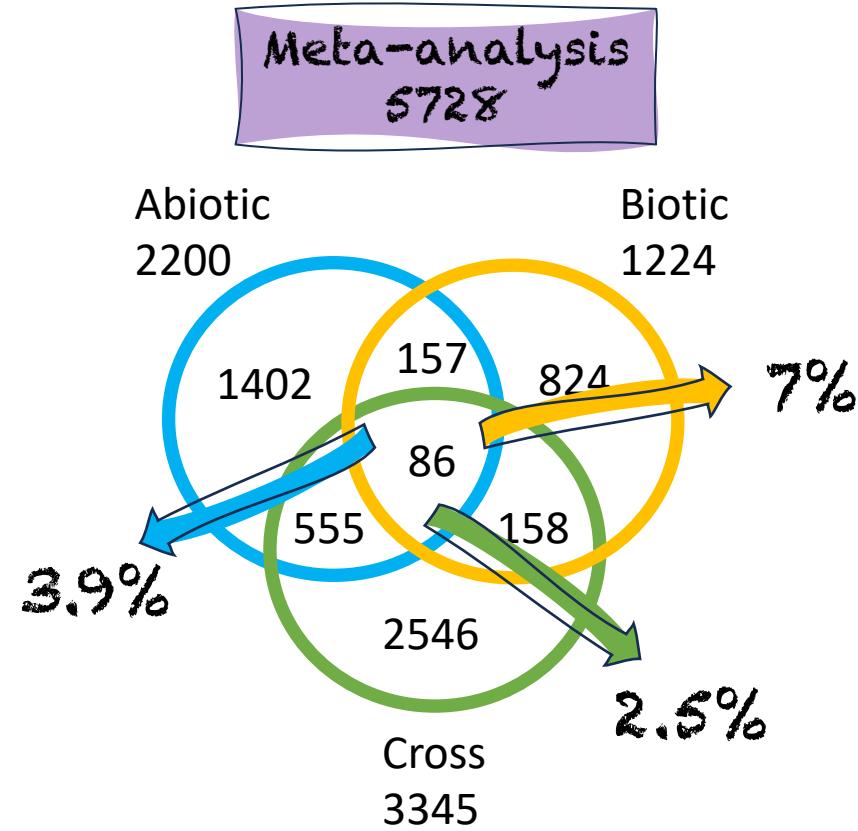
The possible limitations in the performances of MINT on this case study:

- the complex experimental design making difficult to choose the classes to perform the classification
- the limited number of replicates of some classes

HIVE vs meta-analysis



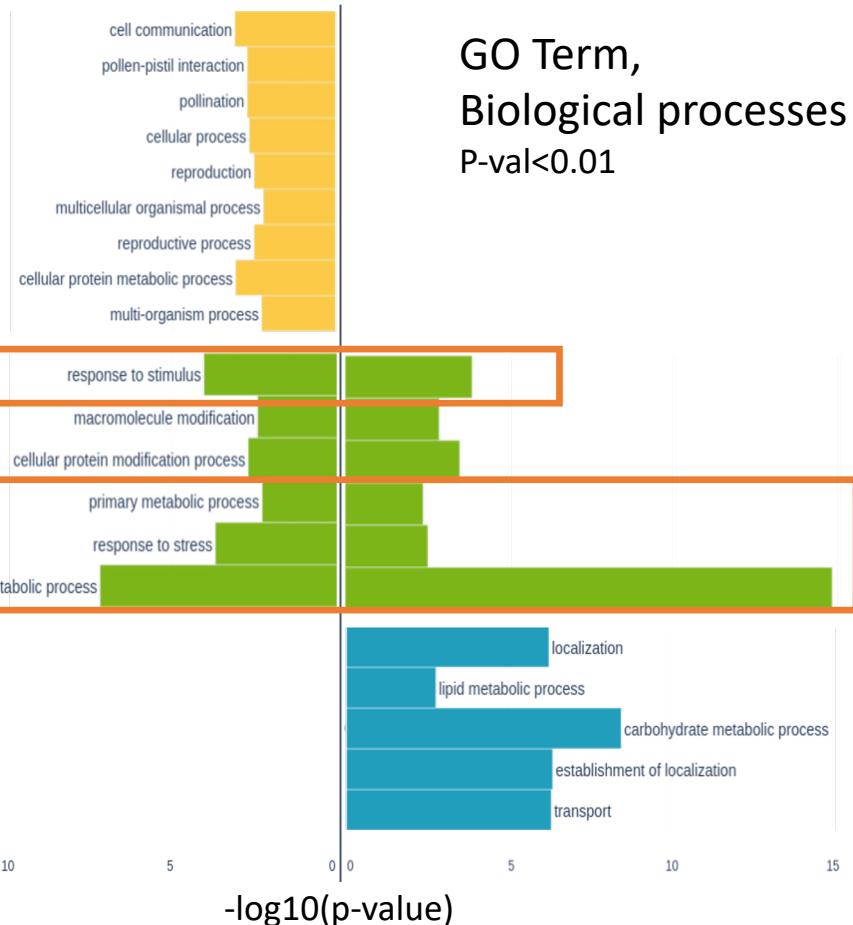
HIVE vs meta-analysis



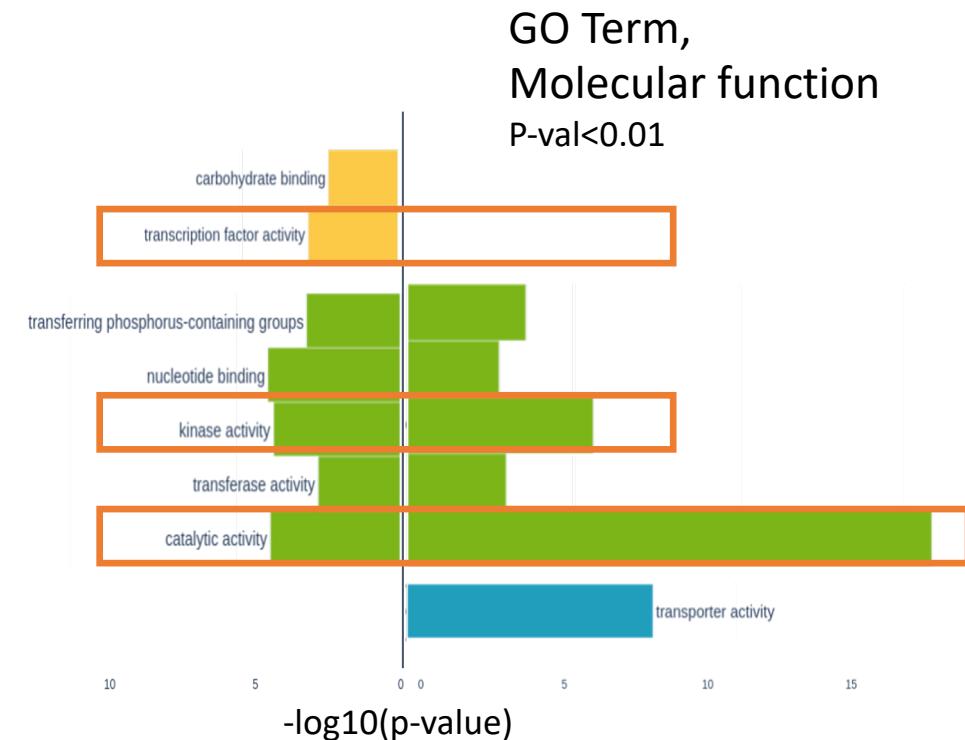
HIVE and meta-analysis have a higher overlap than
the three separated meta-analysis

HIVE vs meta-analysis

Only
HIVE



both



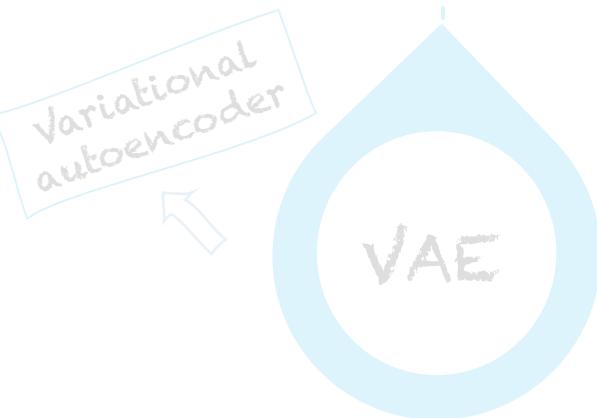
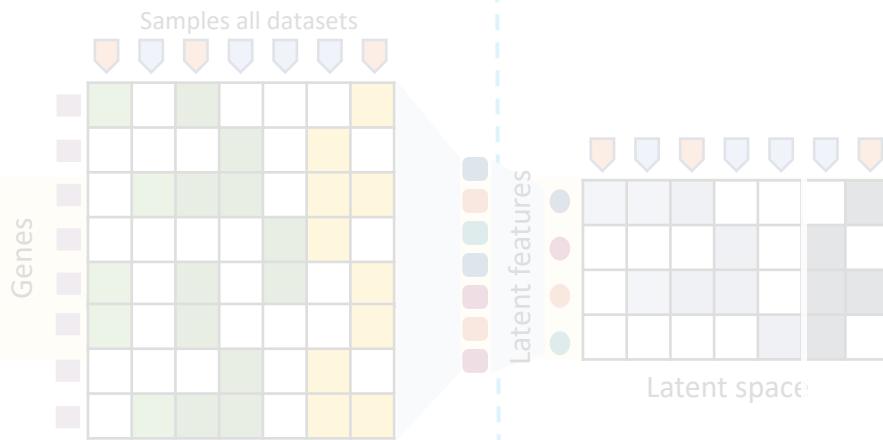
Only meta-
analysis

HIVE selected the genes in common with the meta-analysis that bring most of the biological information and found additional genes with novel biological information

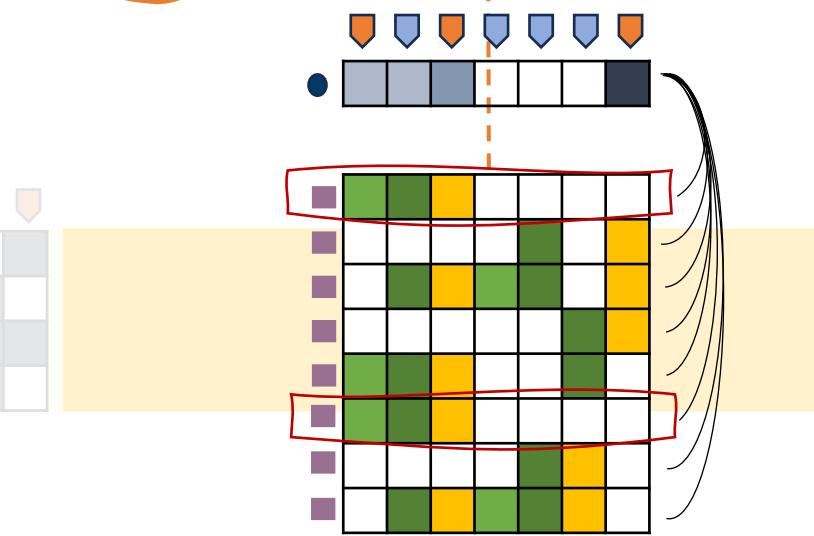


HIVE pipeline

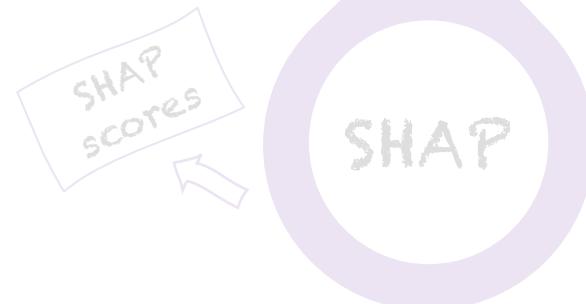
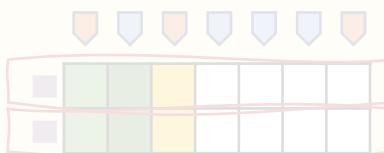
1 Removing batch effects



2 Opening the black box



3 Finding important genes

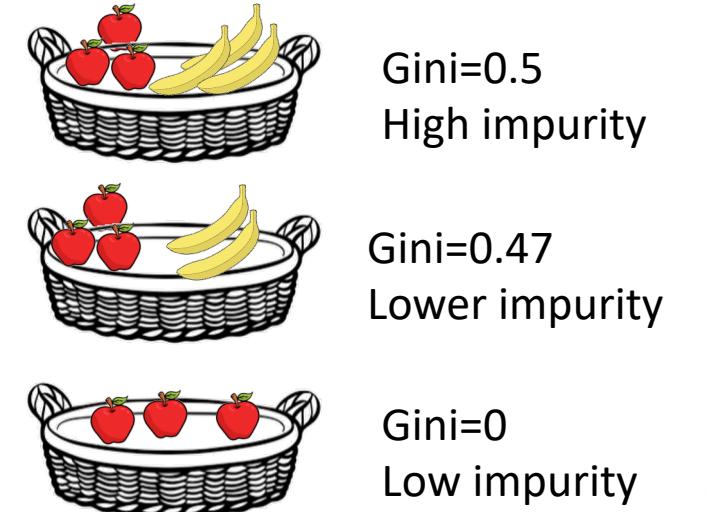


Feature importance

Random forest uses Mean Decrease in Impurity (MDI) to calculate Feature importance.



Gini
index



It calculates for each feature the mean decrease in impurity it introduced across all the decision trees while constructing them.

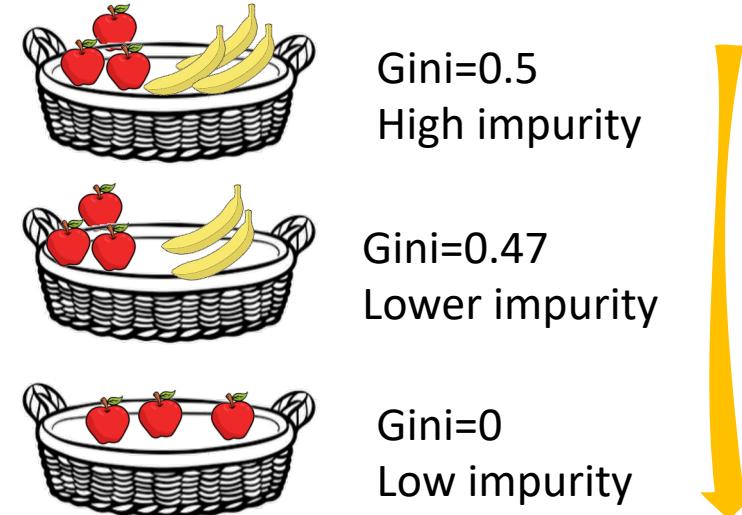
Feature importance

Random forest uses Mean Decrease in Impurity (MDI) to calculate Feature importance.



It calculates for each feature the mean decrease in impurity it introduced across all the decision trees while constructing them.

For every feature, it tries every split and calculates the Gini impurity, compares it to the current impurity, and picks the best one that reduces the impurity.



Feature importance

Random forest uses Mean Decrease in Impurity (MDI) to calculate Feature importance.



Gini index

It calculates for each feature the mean decrease in impurity it introduced across all the decision trees while constructing them.

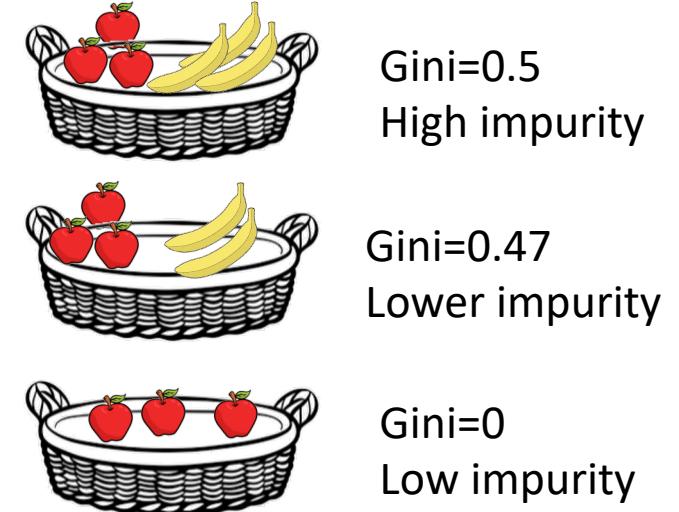
For every feature, it tries every split and calculates the Gini impurity, compares it to the current impurity, and picks the best one that reduces the impurity.



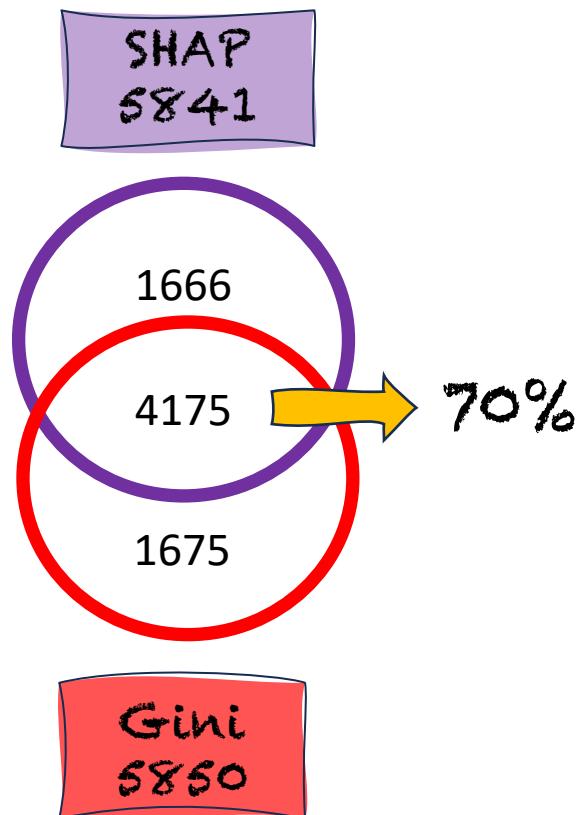
Two features, one has more splitting points than the other, but they both decrease the impurity by the same amount, they will get the same importance!



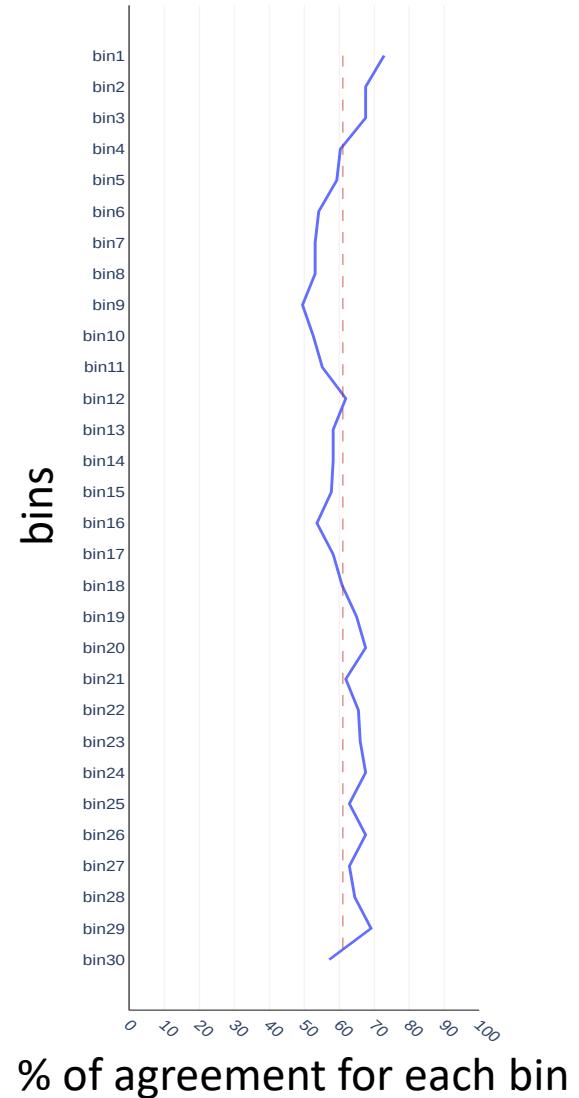
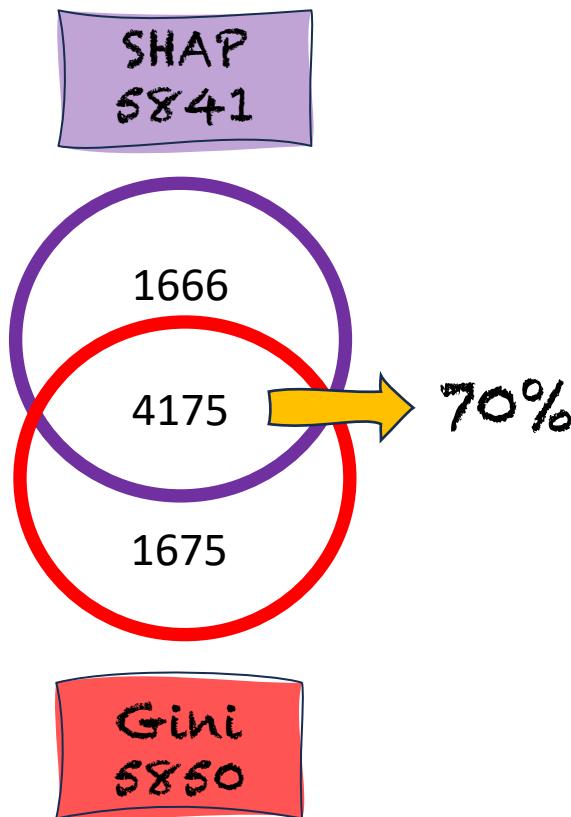
biased towards high cardinal features



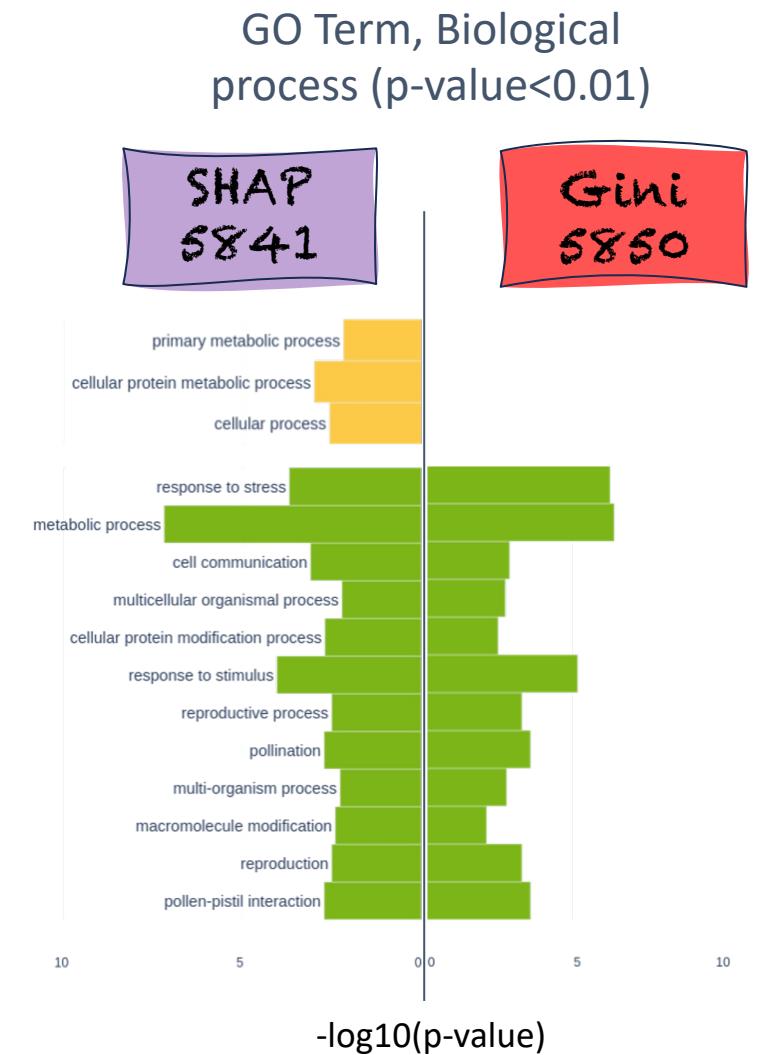
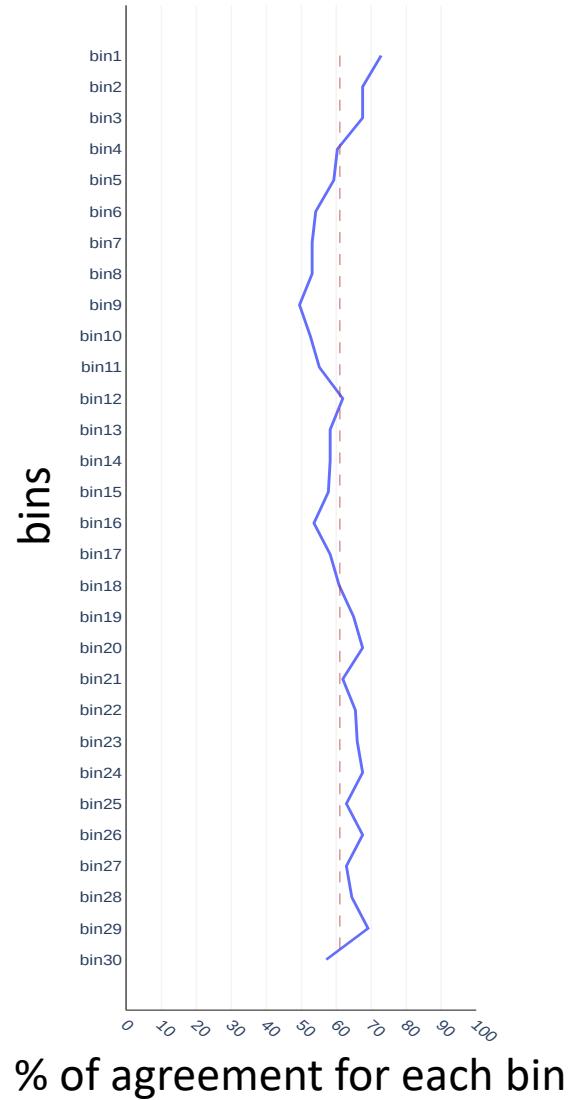
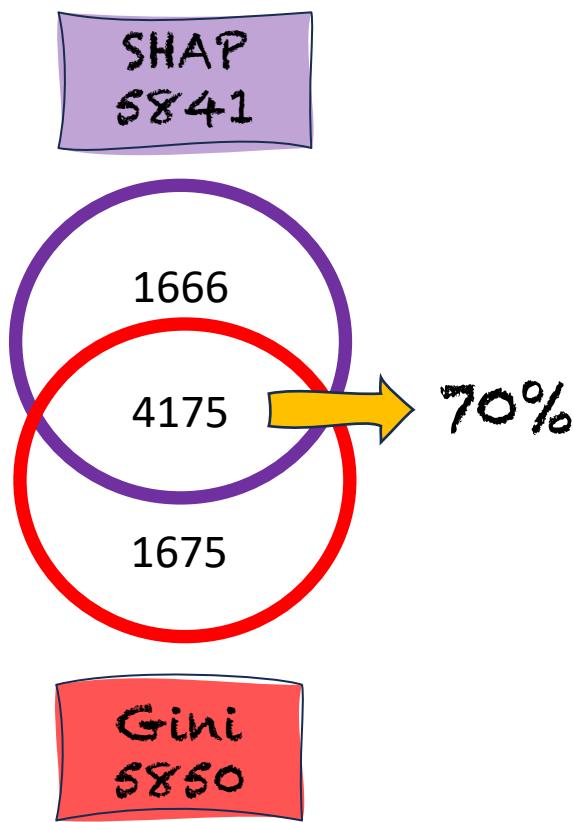
Comparing importance scores



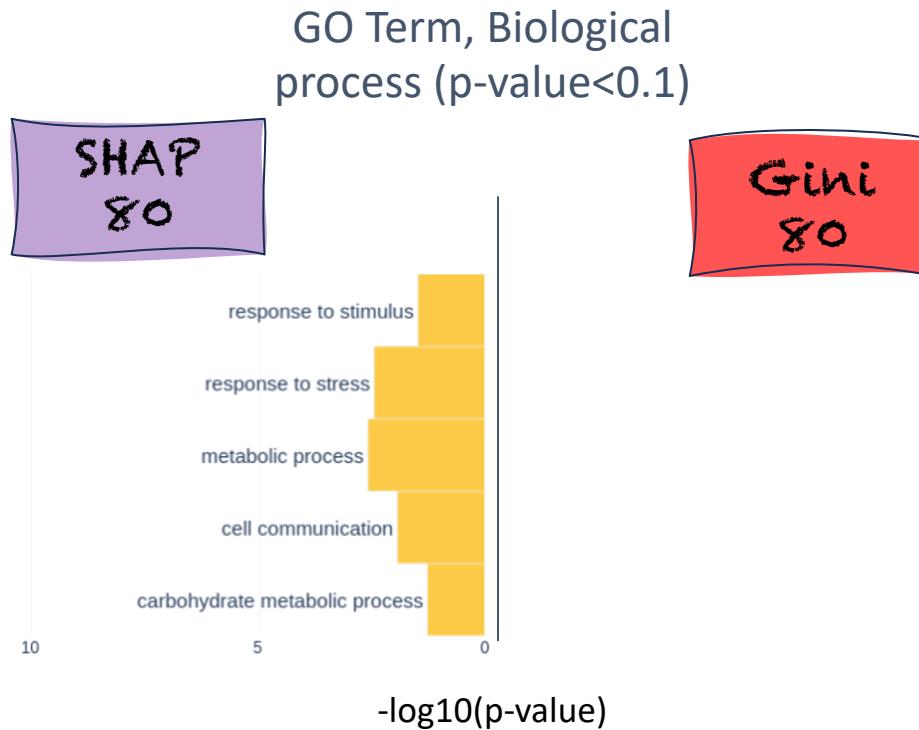
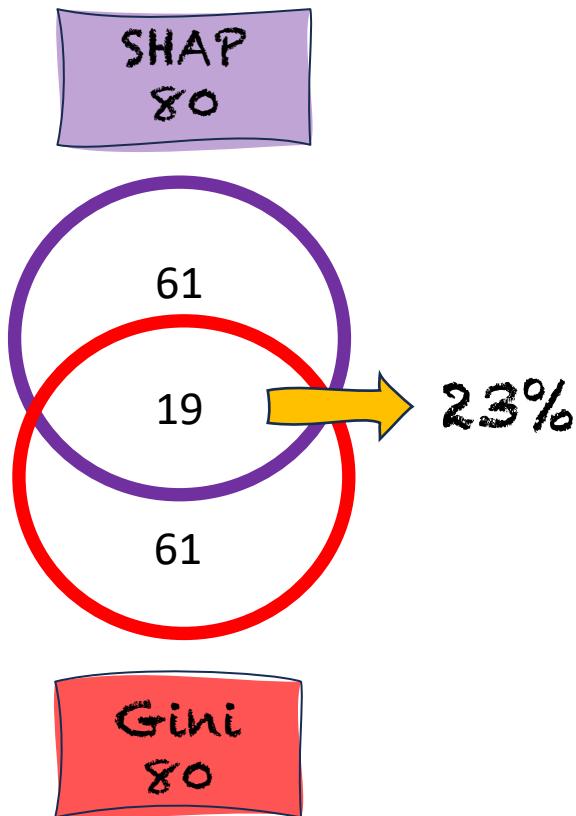
Comparing importance scores



Comparing importance scores

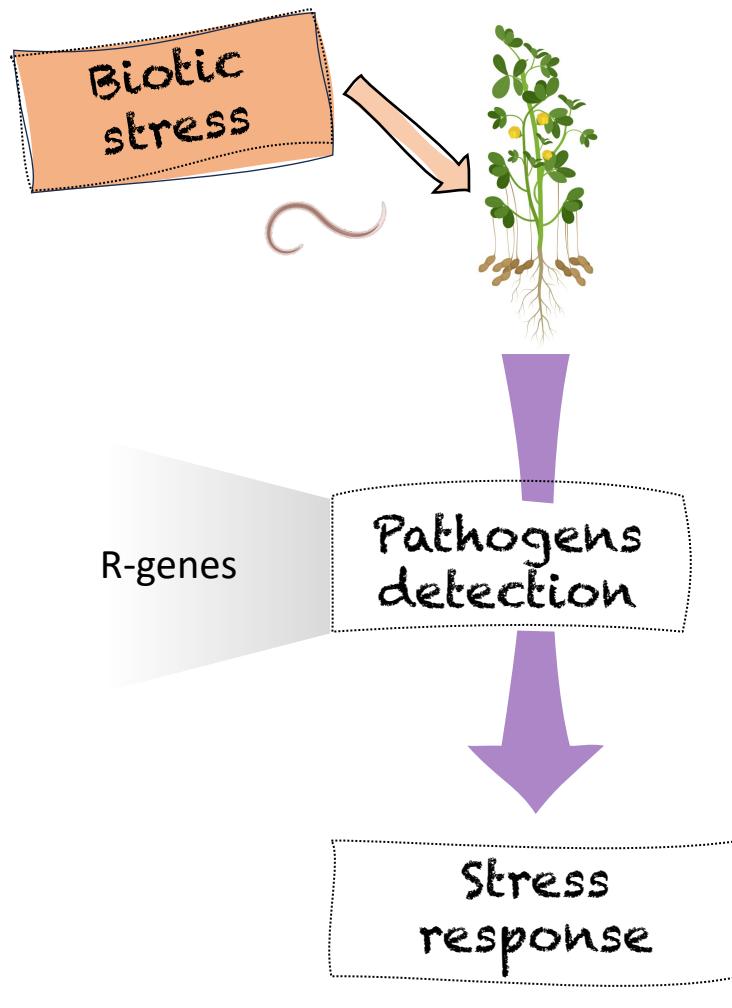


SHAP score performs better in selecting genes with biological meaning

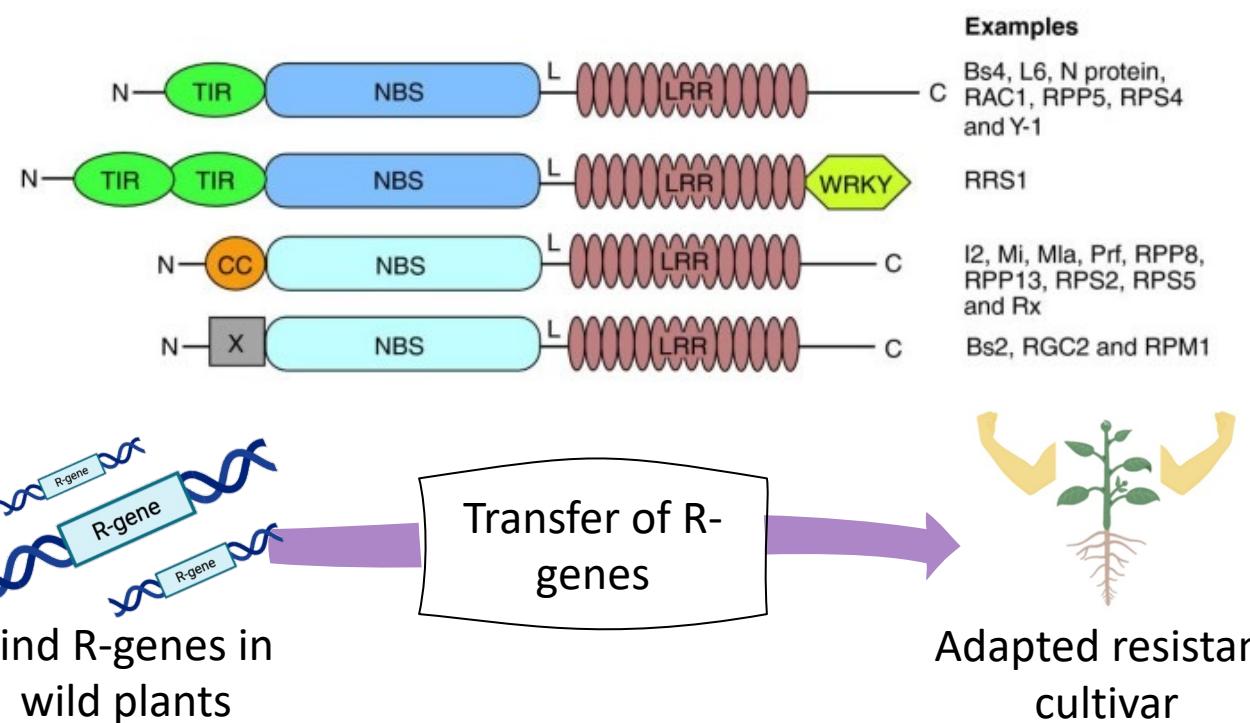


... soon:
LIME & SAGE

Plant NBS-LRR proteins: adaptable guards



Most of disease resistance genes (R-genes) in plants encode nucleotide-binding site leucine-rich repeat (NBS-LRR) proteins. This large family is encoded by hundreds of diverse genes per genome and can be subdivided into the functionally distinct TIR-domain-containing (TNL) and CC-domain-containing (CNL) subfamilies.



NBS-LRR genes

