

# AI-Powered Genomics

## Unveiling Hidden Controls of Plant Gene Expression

Jędrzej Jakub Szymański



@NAMlab



[www.szymanskilab.com](http://www.szymanskilab.com)



## **Artificial Intelligence (AI)**

AI refers to computer systems that mimic human intelligence, enabling them to solve problems and understand language.

## **Machine Learning (ML)**

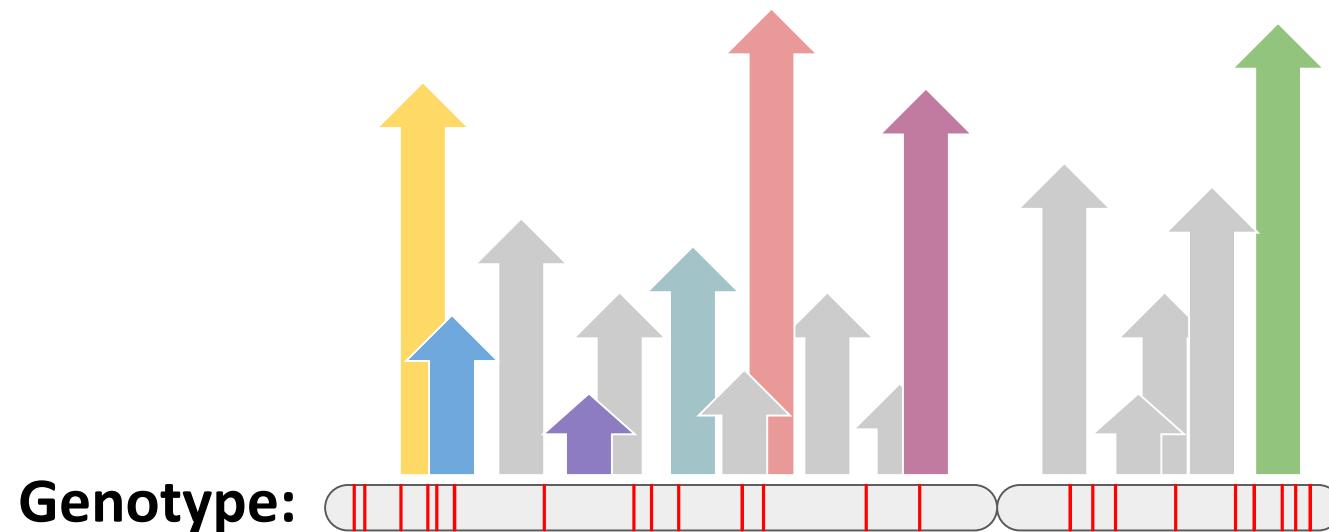
ML is teaching systems to learn from data, enhancing performance without explicit programming.

## **Deep Learning**

Deep Learning employs layered neural networks to find intricate patterns, excelling in tasks like image and speech recognition.

# Genotype to phenotype

**Phenotype:**



# Genotype to phenotype

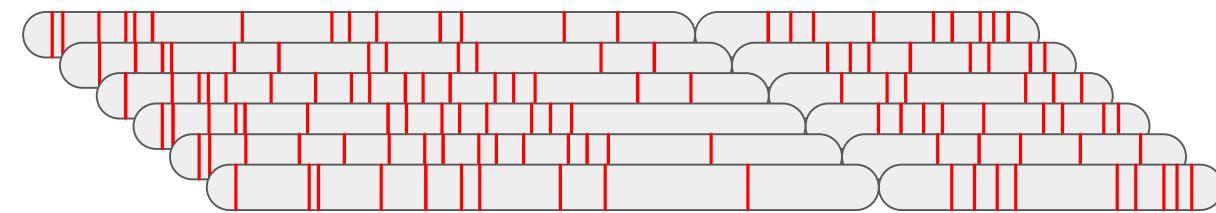
**Phenotypes:**



Photo: Uli Westphal  
(IPK cultivar collections)  
[www.ulwestphal.de](http://www.ulwestphal.de)



**Genotypes:**

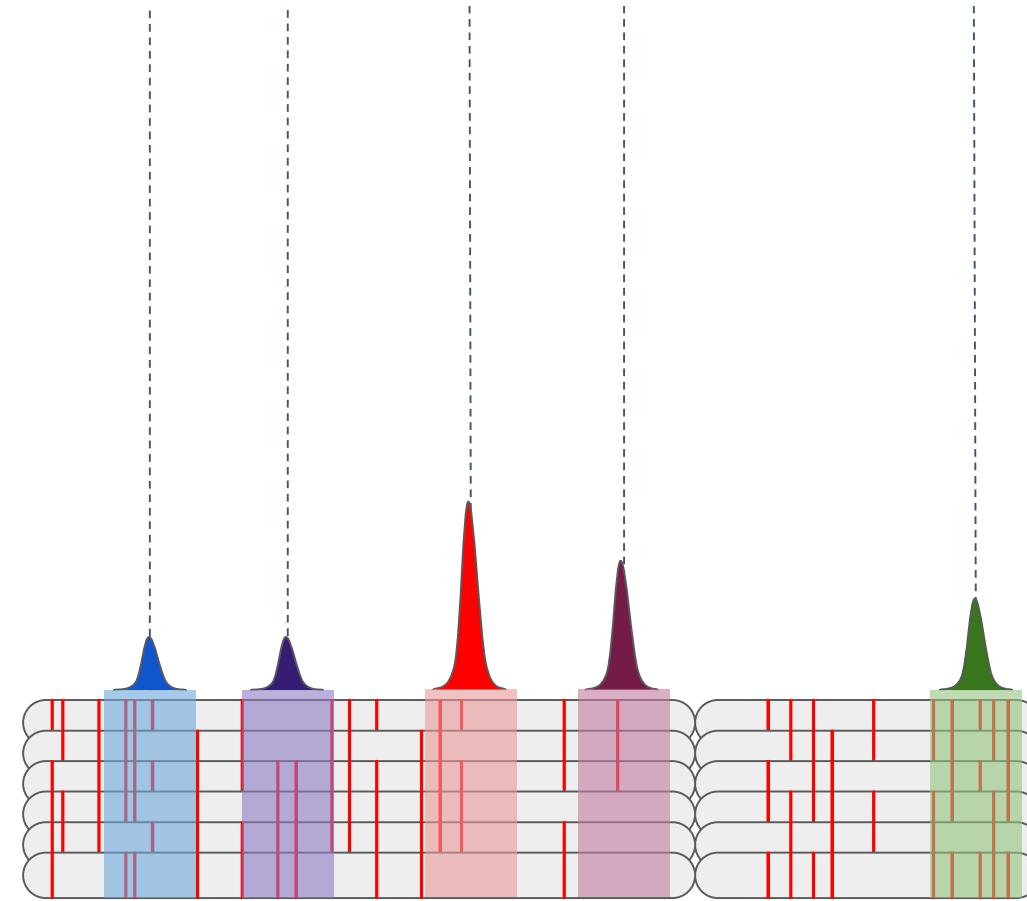


# QTL / GWAS

**Phenotypes:**



**Genotypes:**



# QTL / GWAS / Systems Genetics

Phenotypes:



ENDOPHENOTYPES

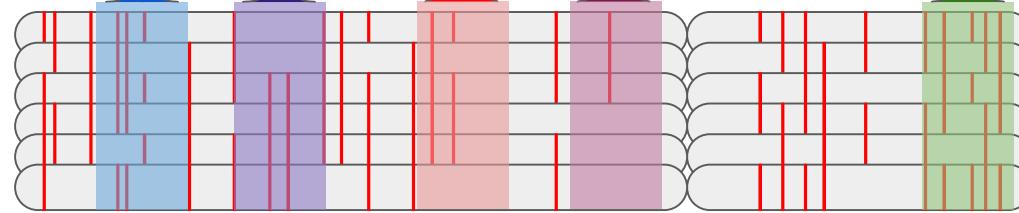
PHENOMICS

METABOLOMICS

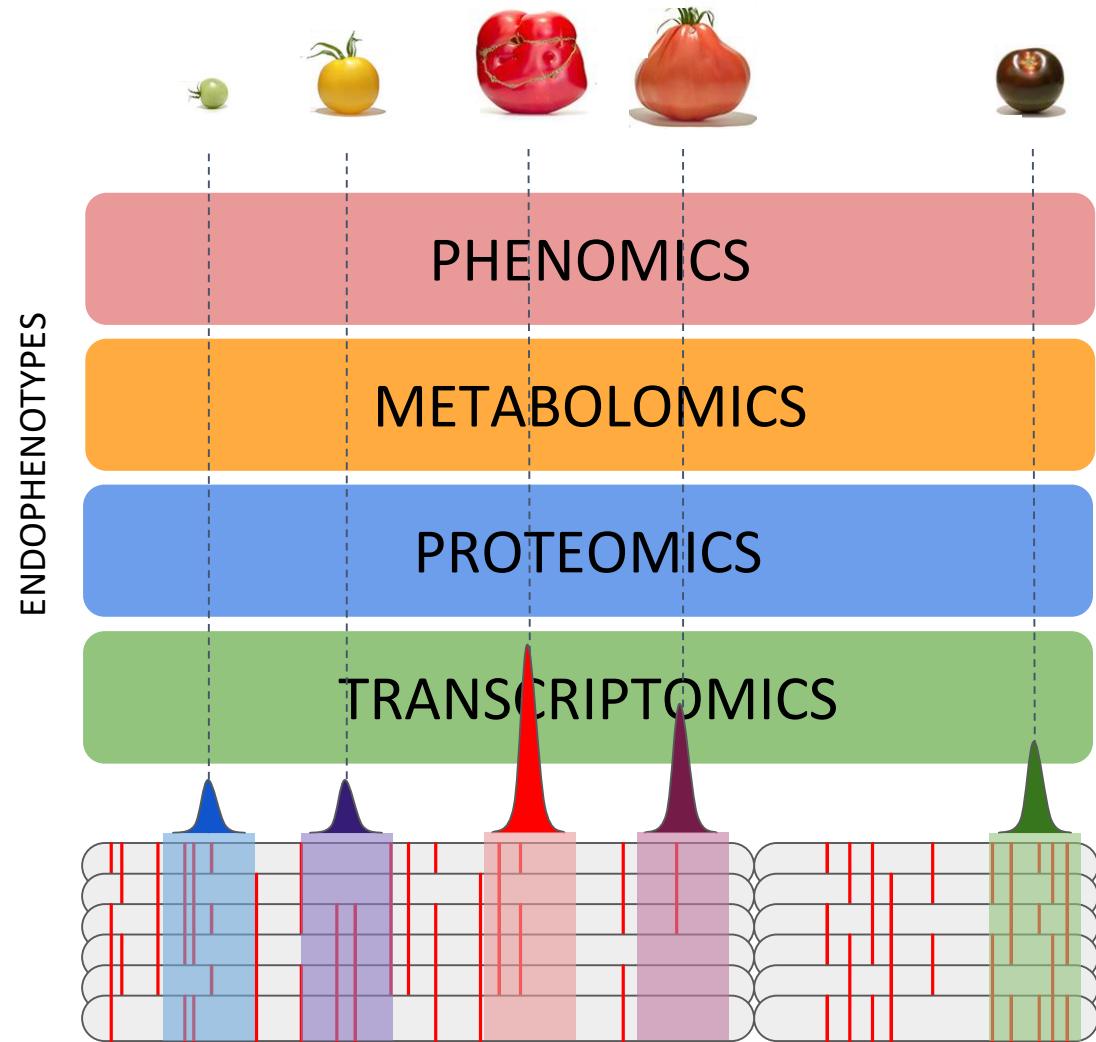
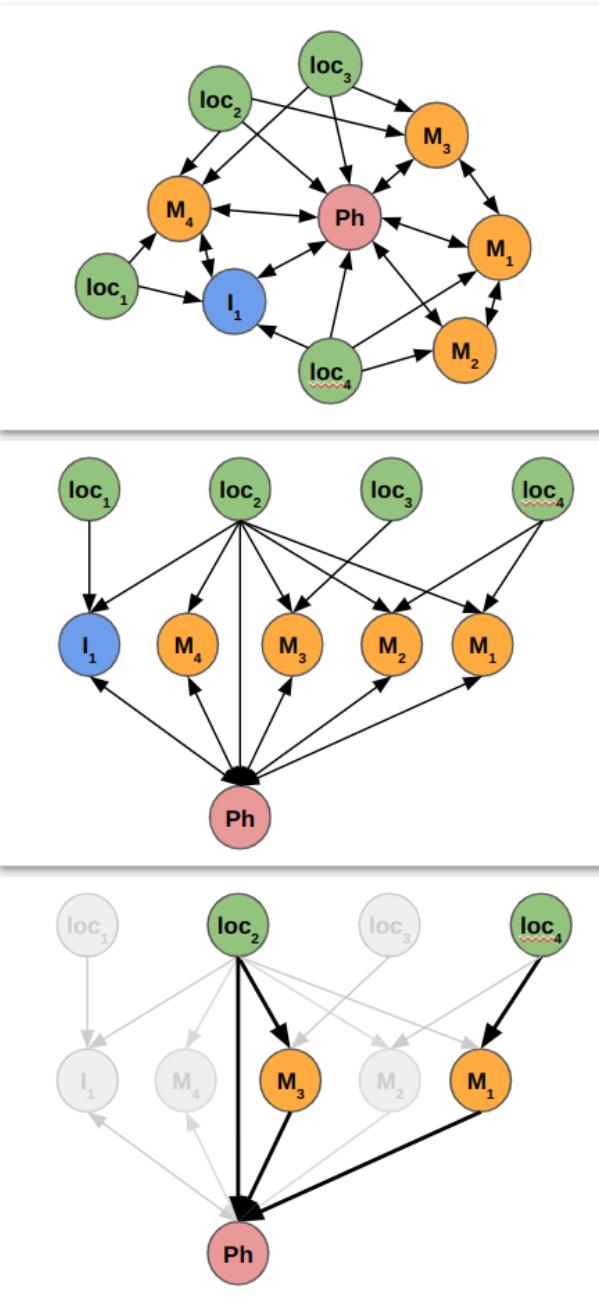
PROTEOMICS

TRANSCRIPTOMICS

Genotypes:



# Genomic networking



# From genomic networking to molecular mechanisms

*S. pennellii*

#WildBeast #Resistant  
#Toxic #StayAway

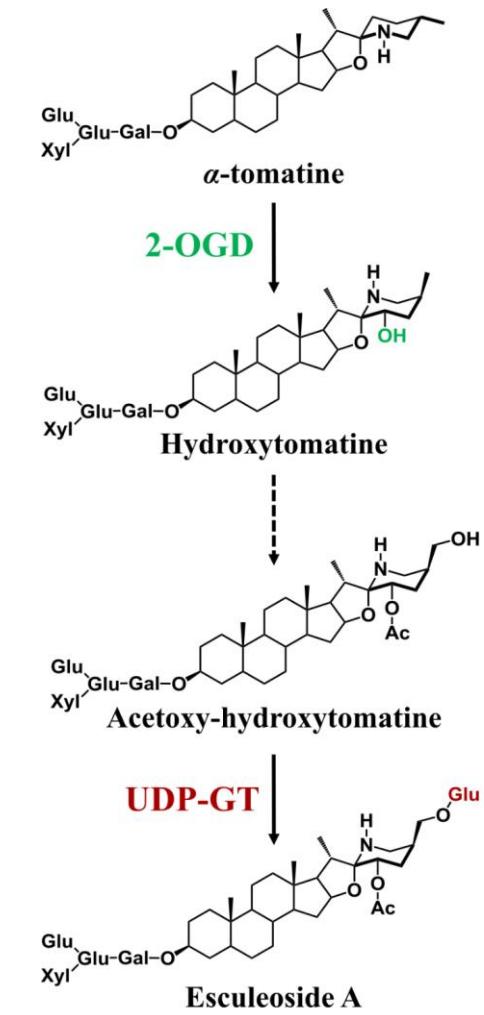
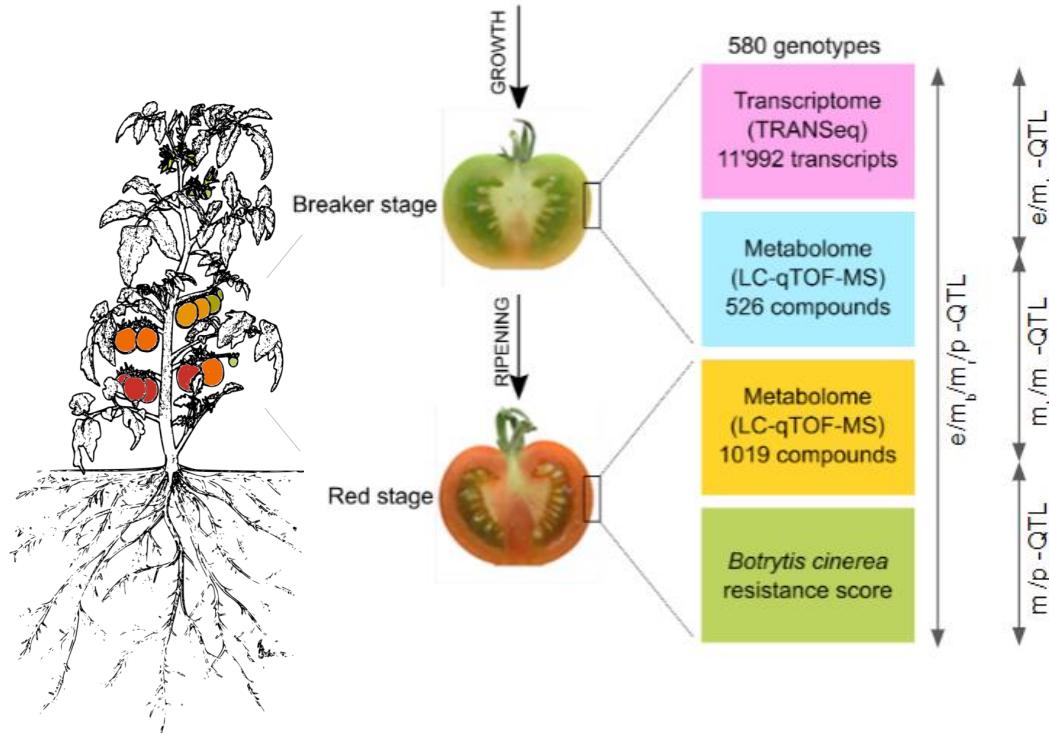
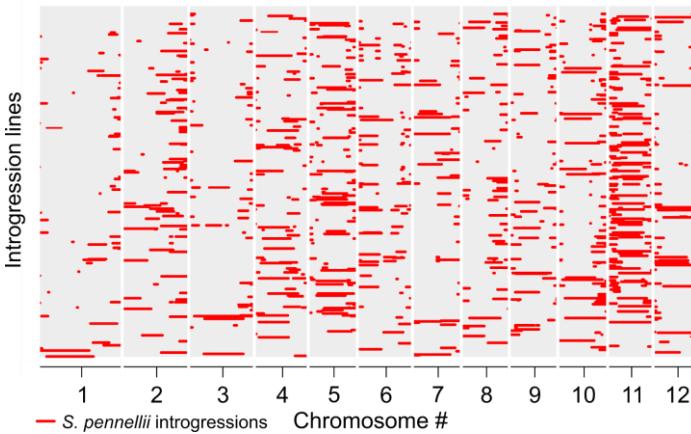


*S. lycopersicum*

#Domesticated #Sensitive  
#Sweet #TakeMeHome

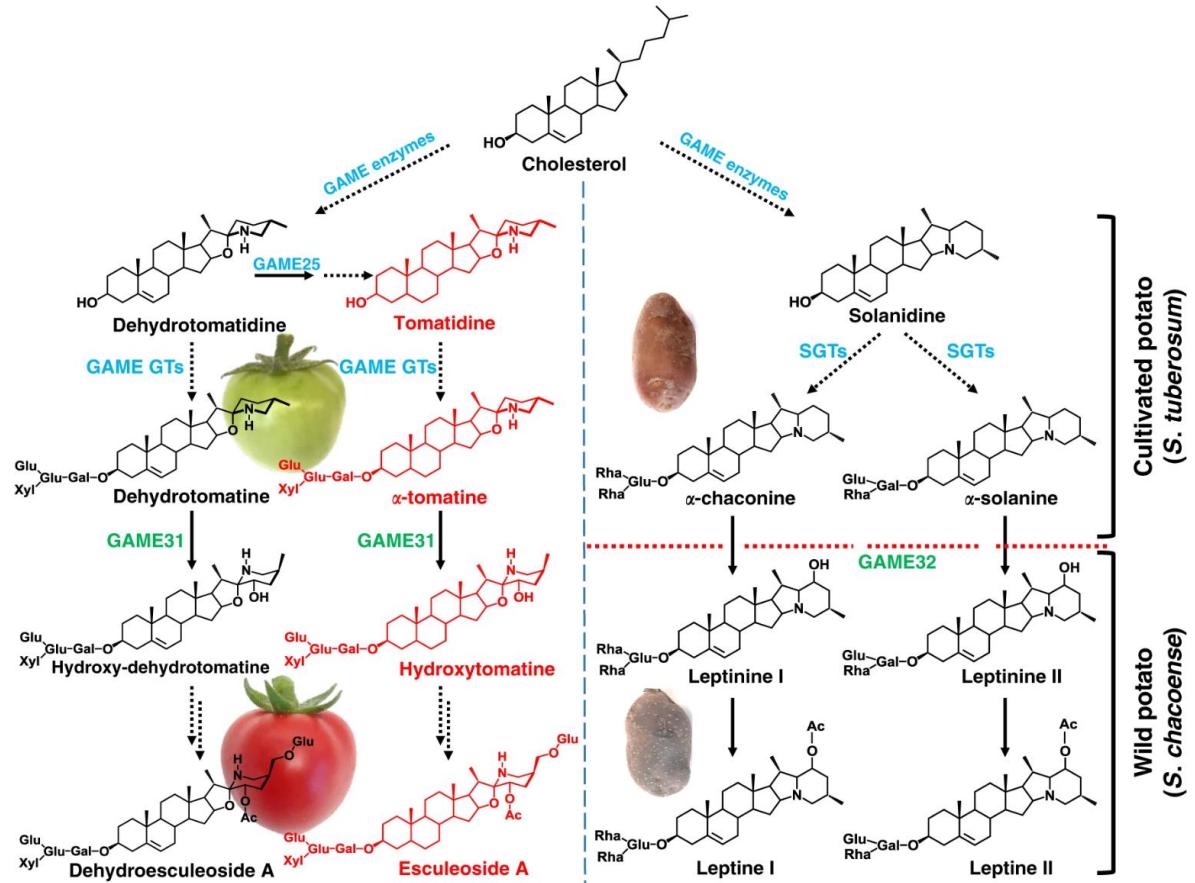


#Domesticated #Resistant  
#Sweet #TakeMeHome



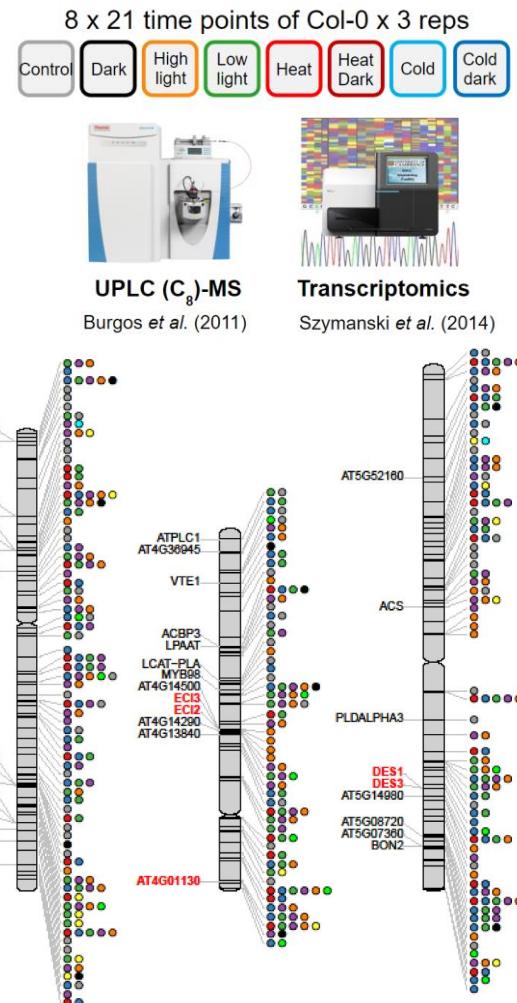
## Metabolic mechanisms of *Botrytis cinerea* resistance in tomato fruits

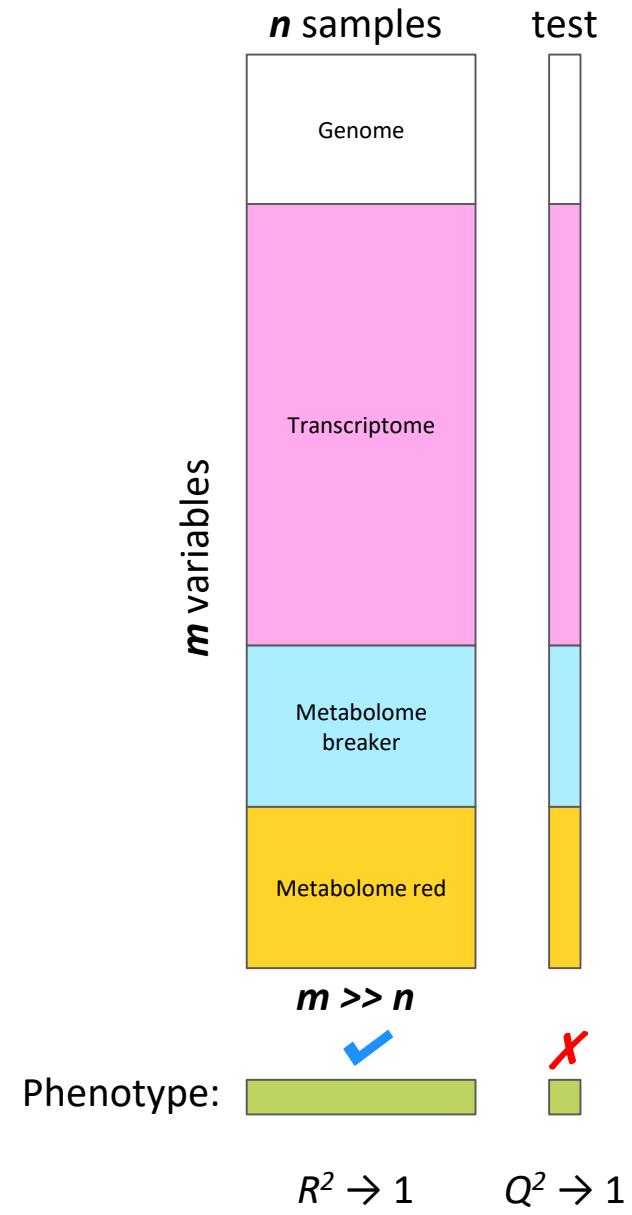
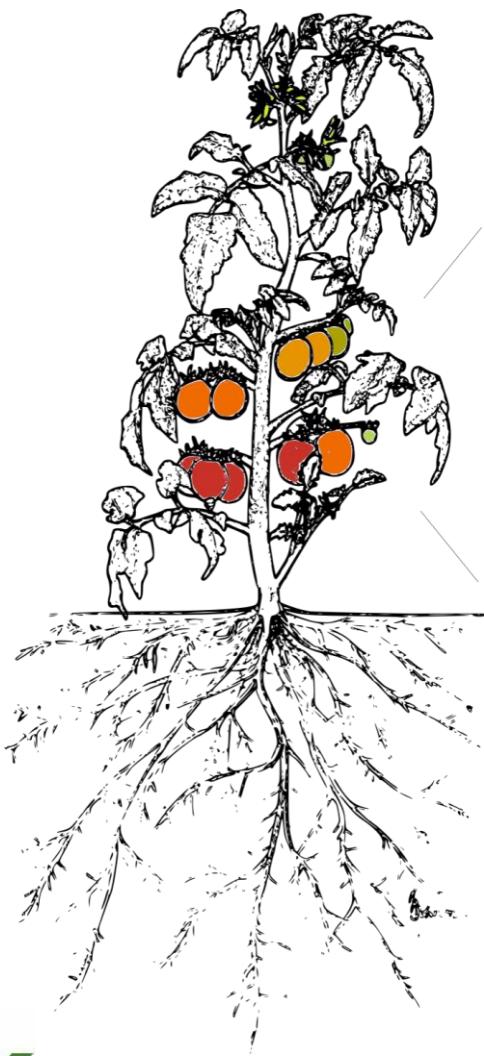
# From genomic networking to molecular mechanisms

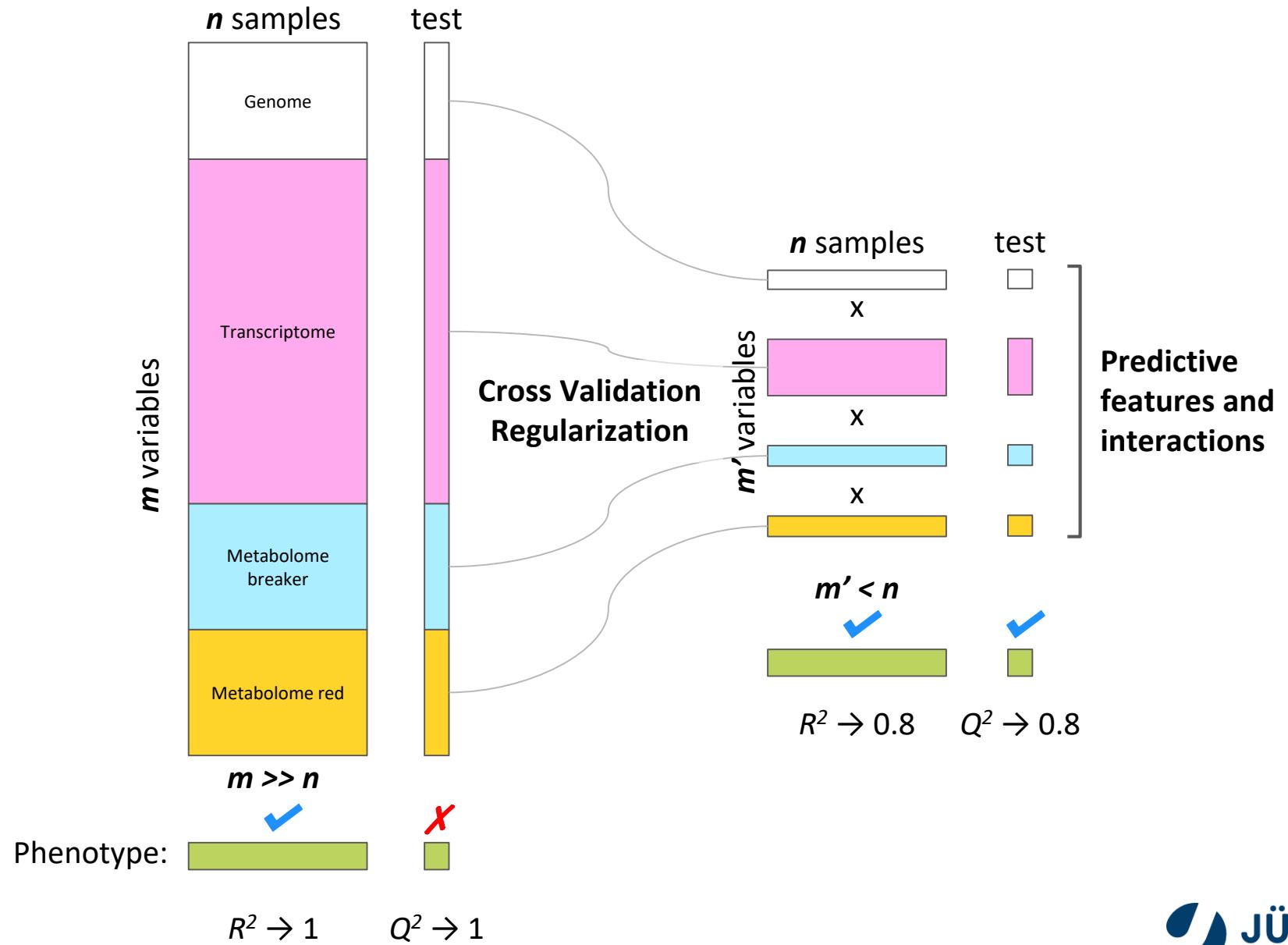
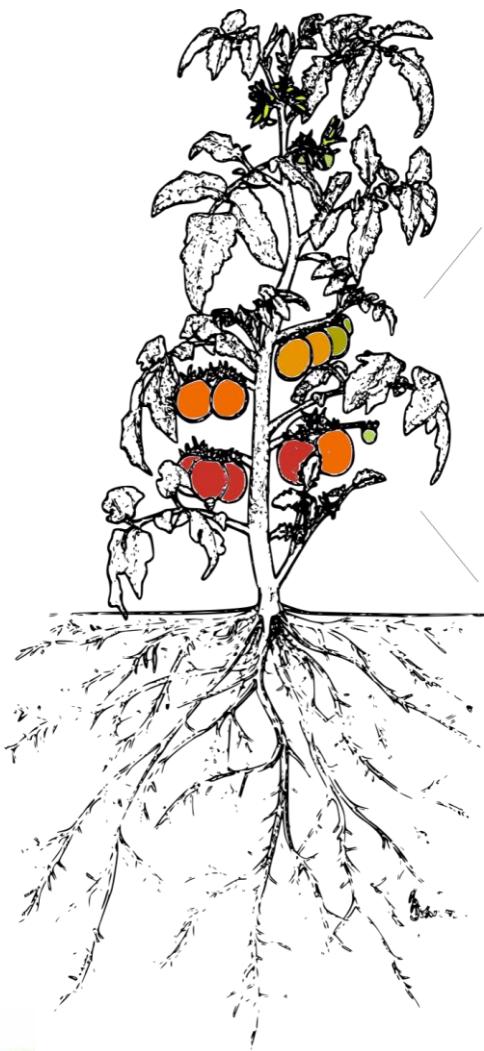


Metabolic mechanisms of herbivore defense in potato

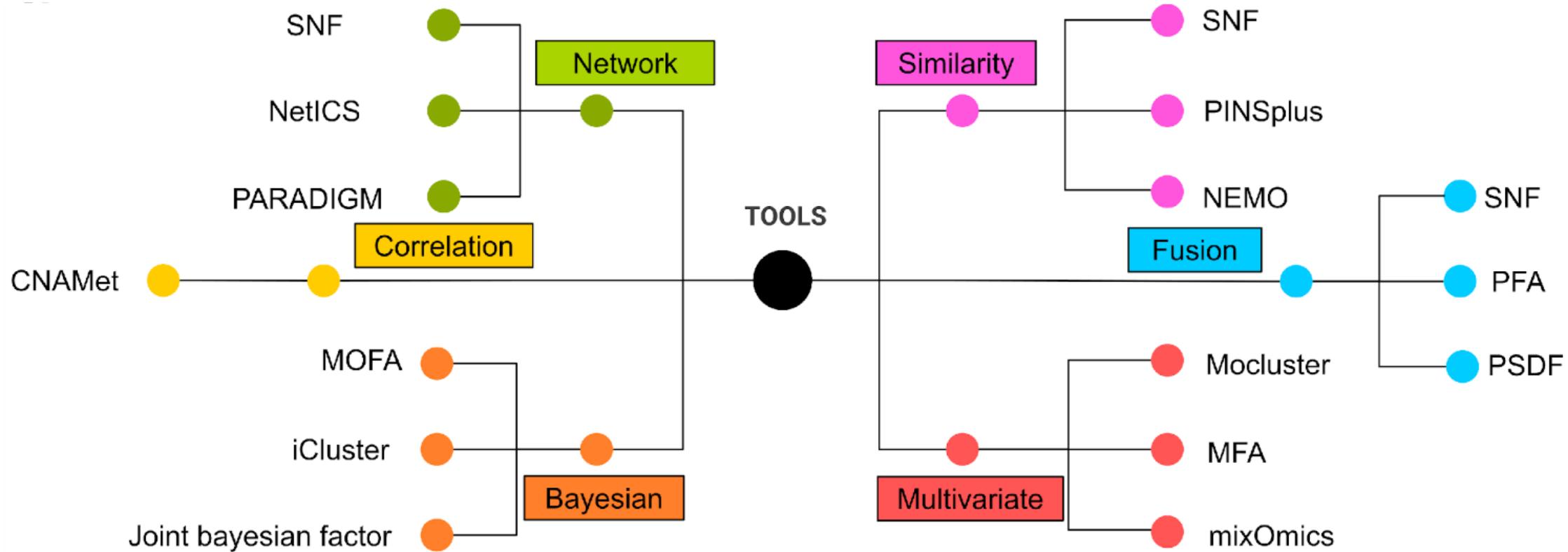
Mechanisms of adaptation to environment in *Arabidopsis*





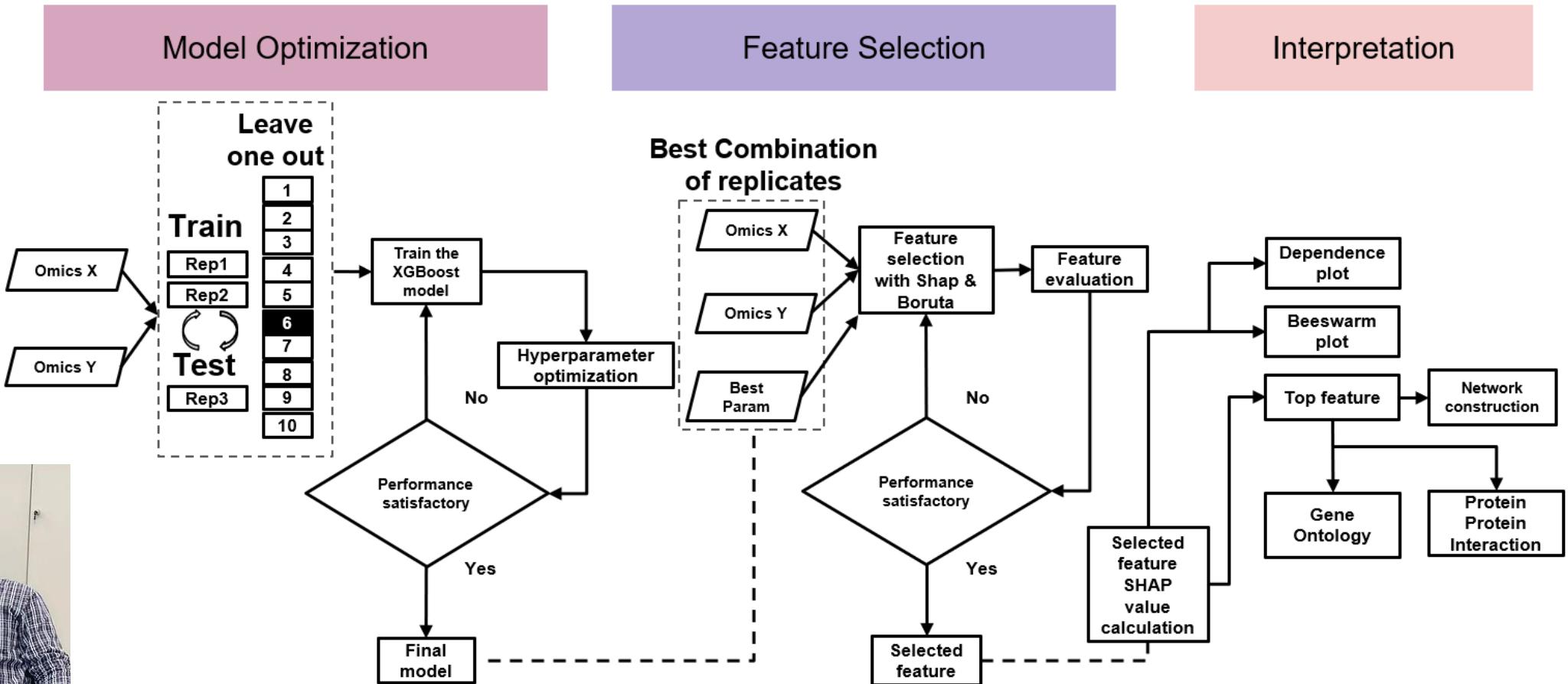


# What kind of machine learning / AI?

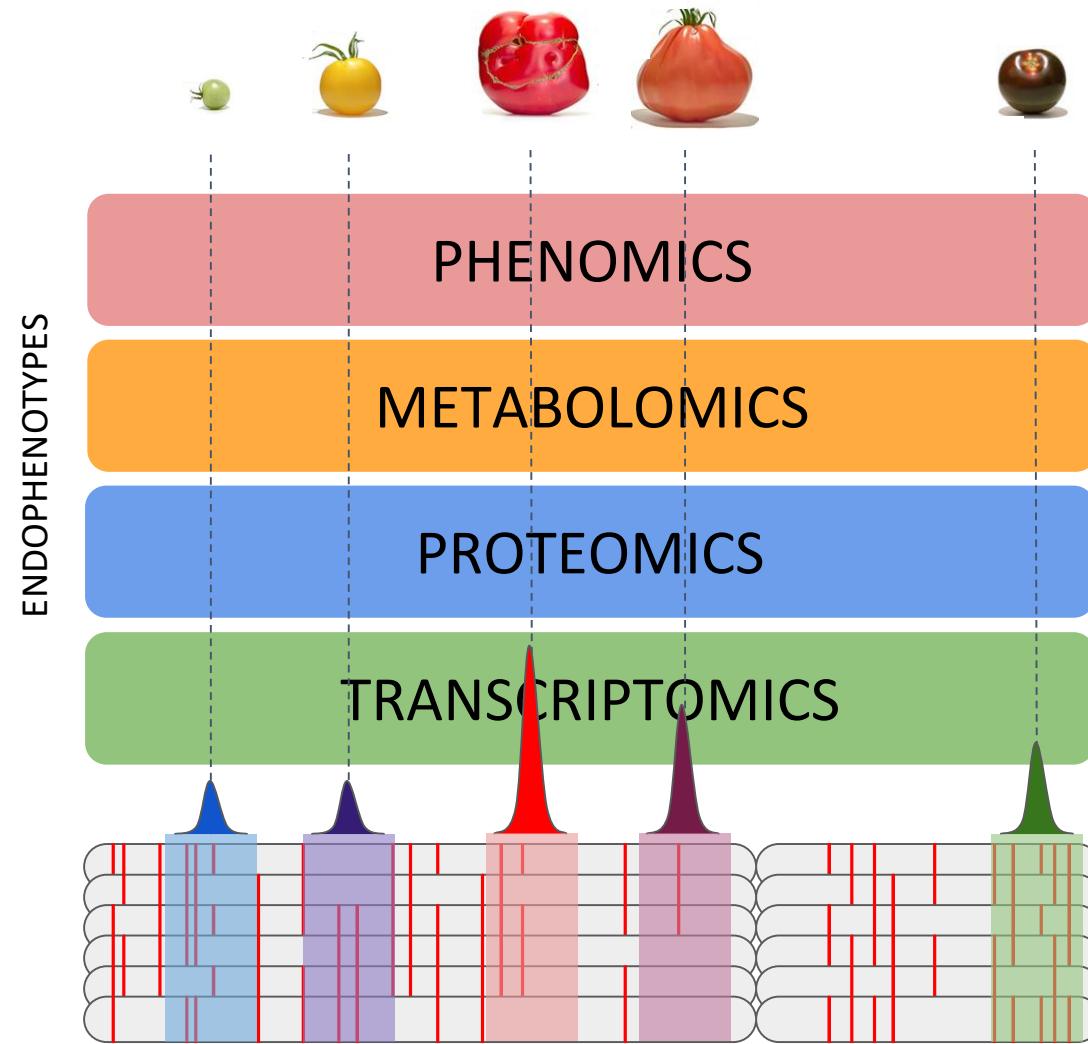




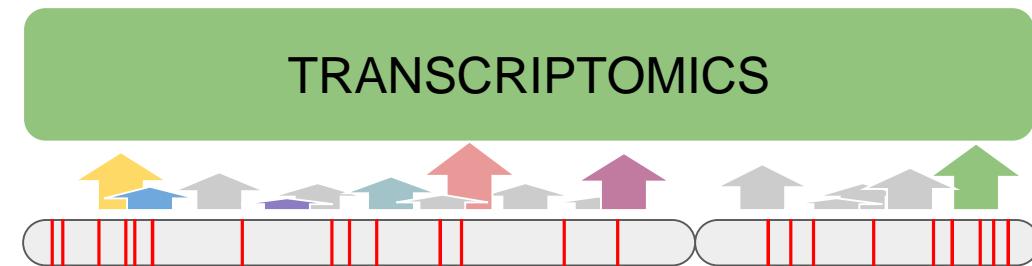
Ankur Sahu



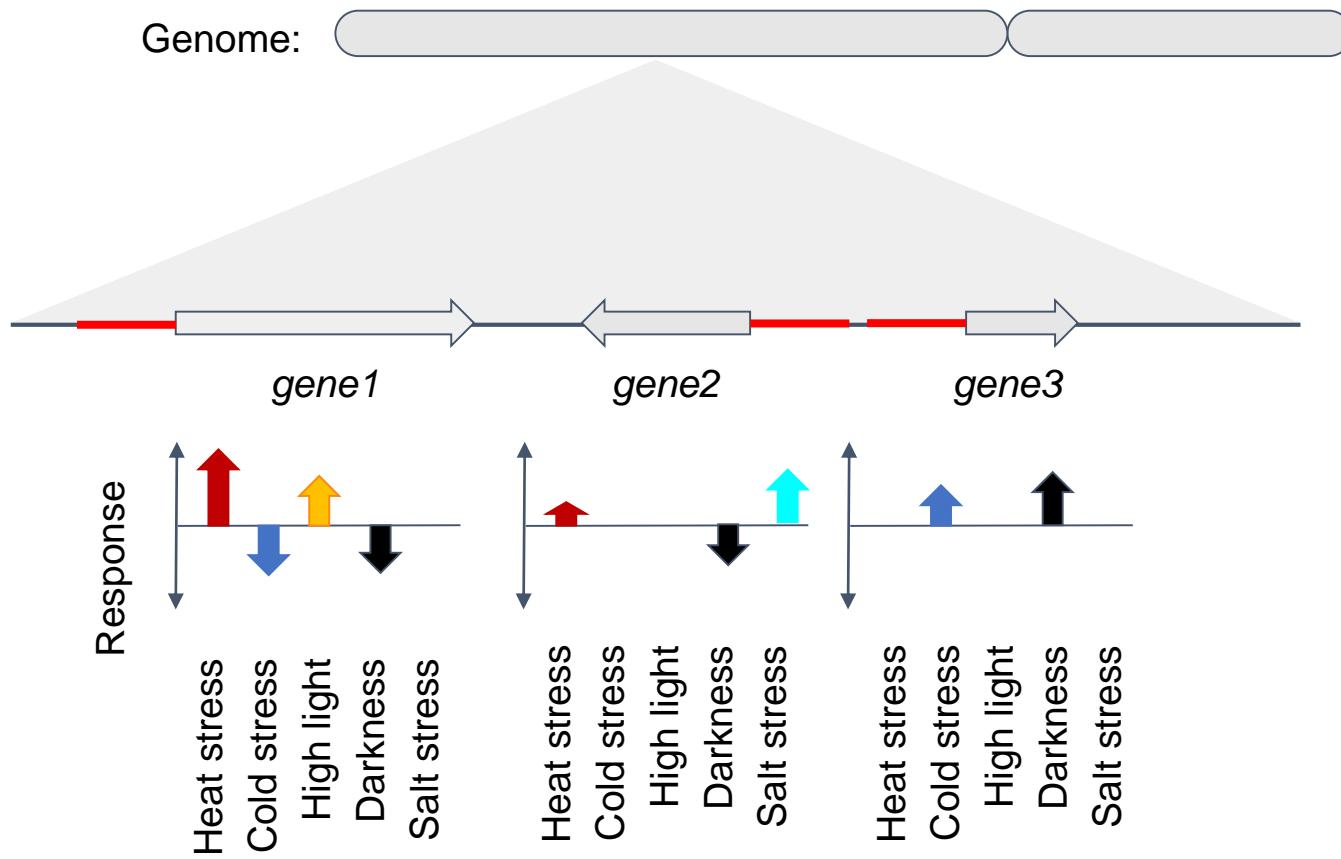
# Systems genetics



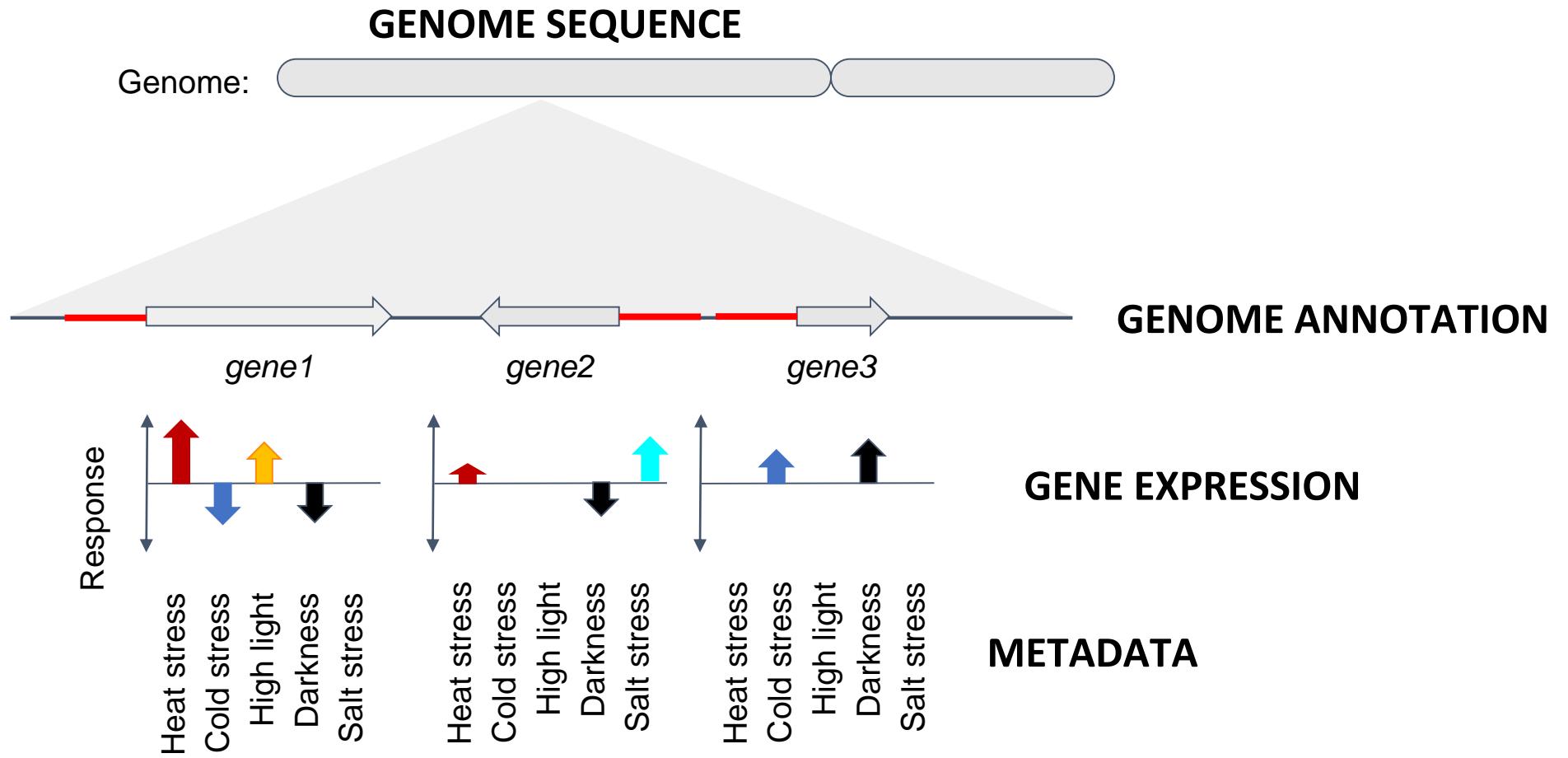
# What about a simpler question?



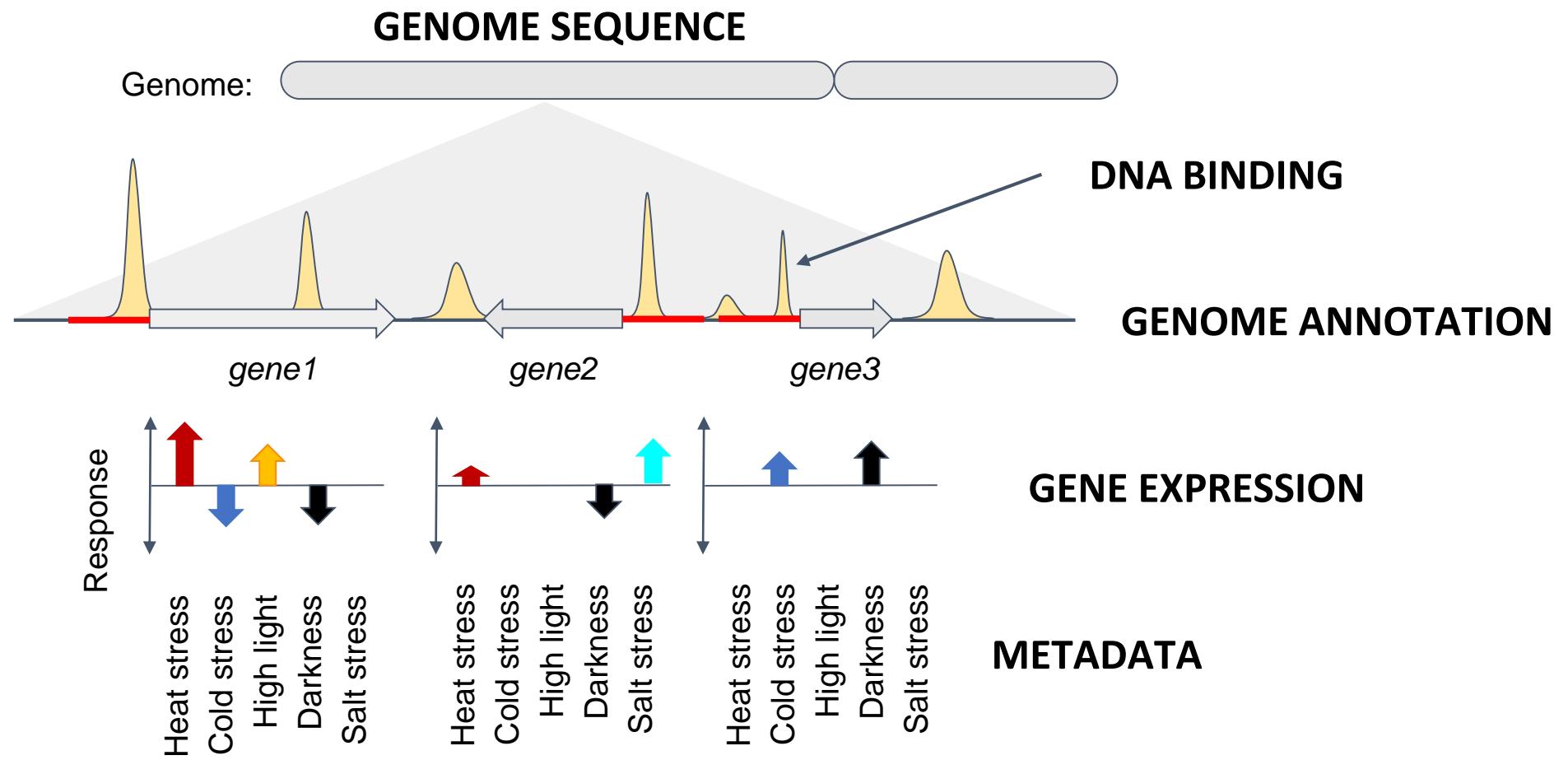
# Gene expression regulation



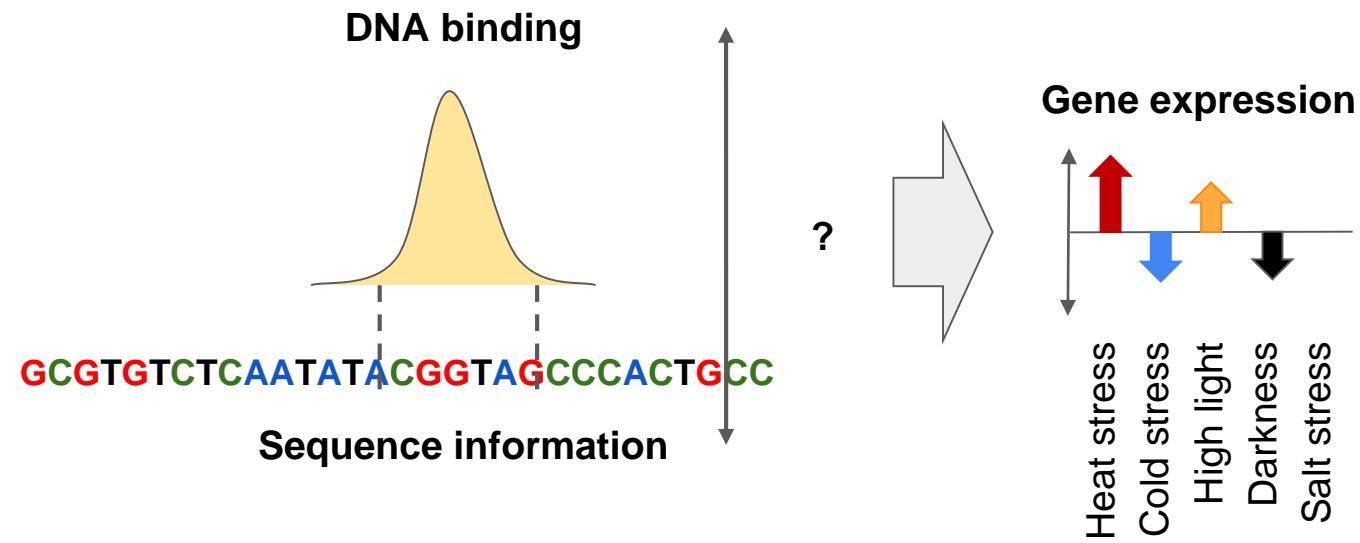
# Gene expression regulation - data



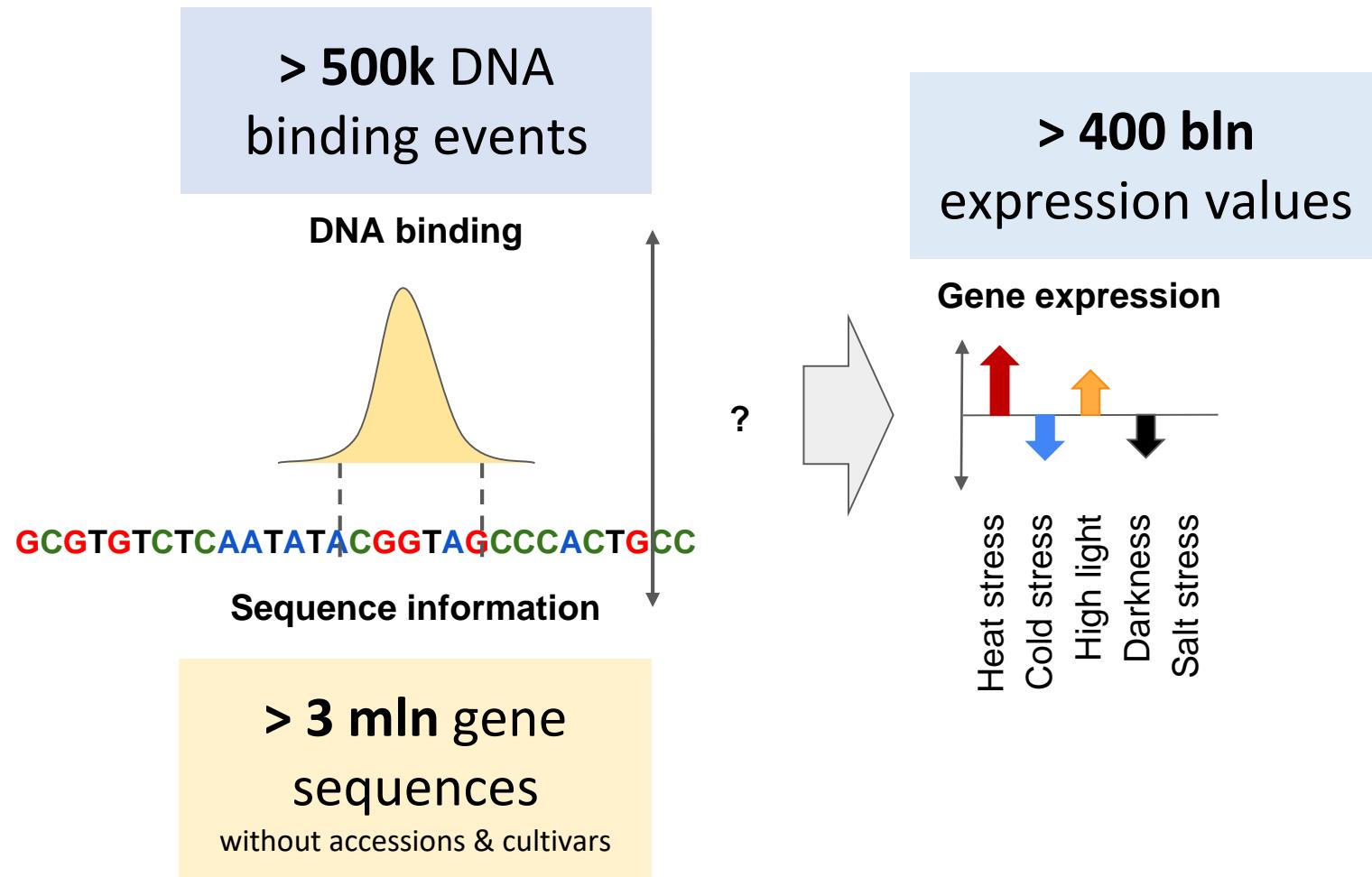
# Gene expression regulation - data



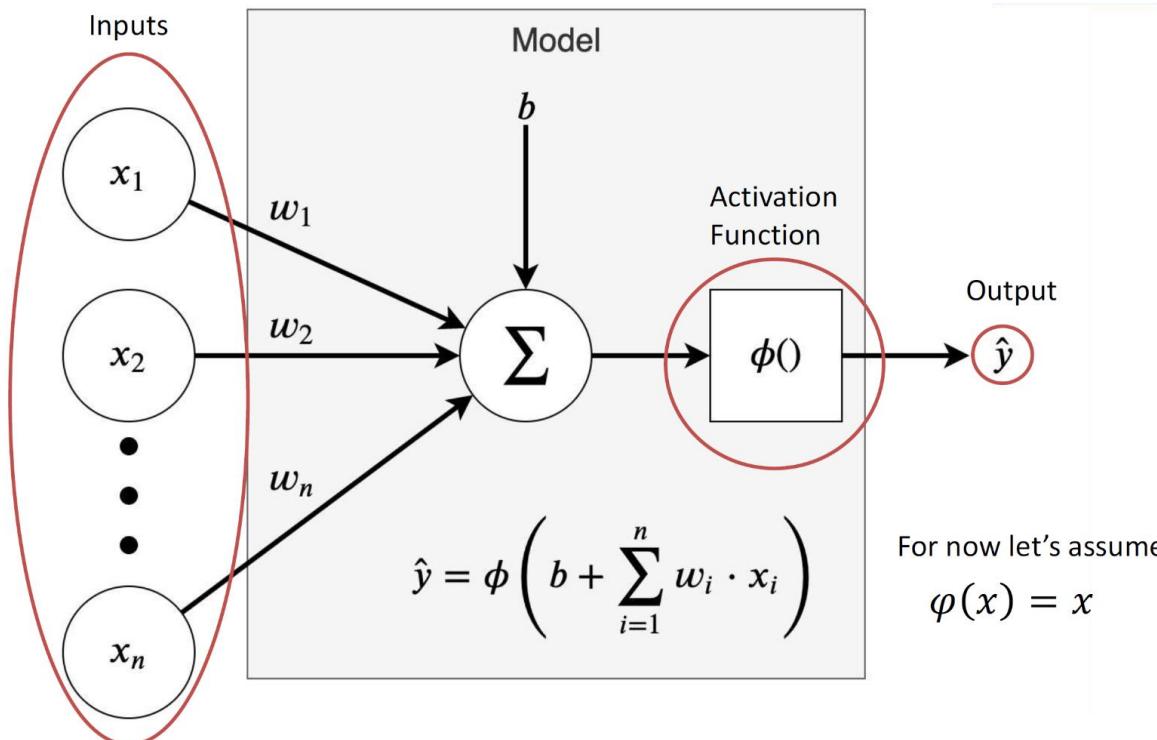
# The task



# The task

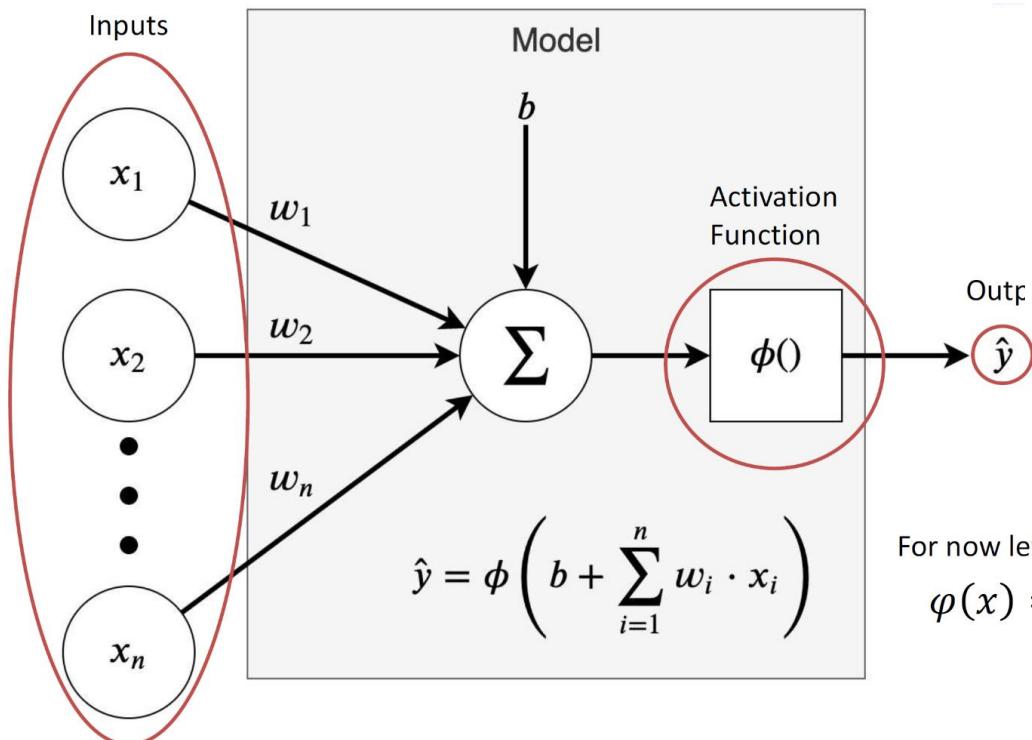


# Artificial Neural Networks



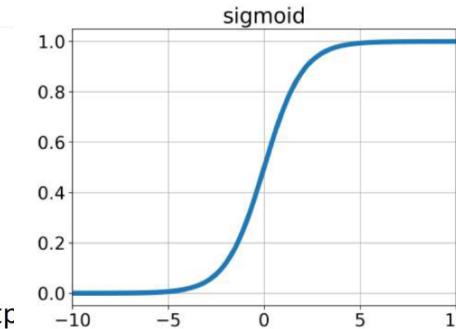
For now let's assume  
 $\varphi(x) = x$

# Artificial Neural Networks

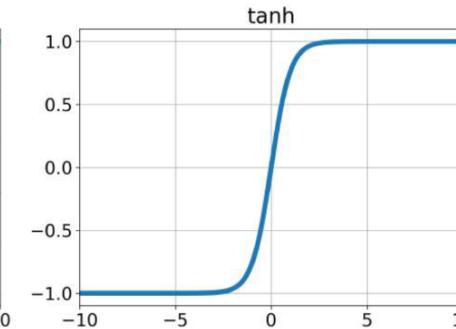


For now let's assume  
 $\varphi(x) = x$

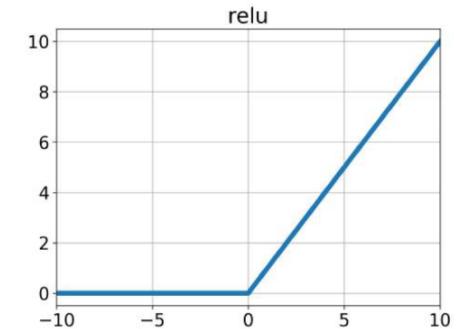
$$\varphi(x) = \frac{1}{1 + e^{-x}}$$



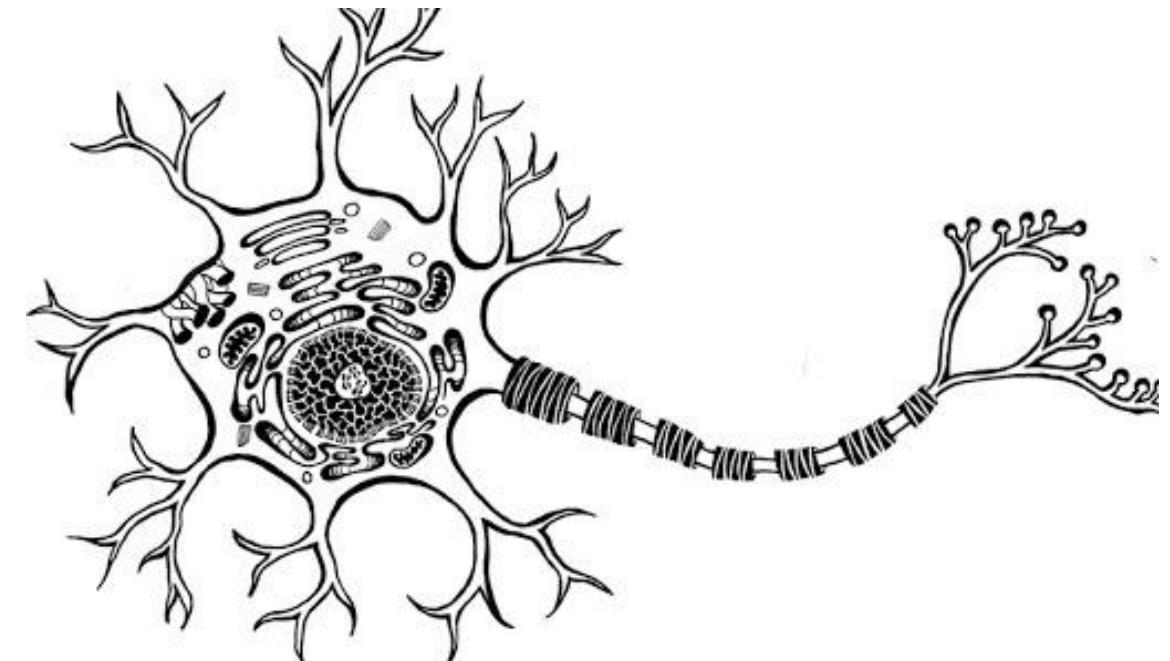
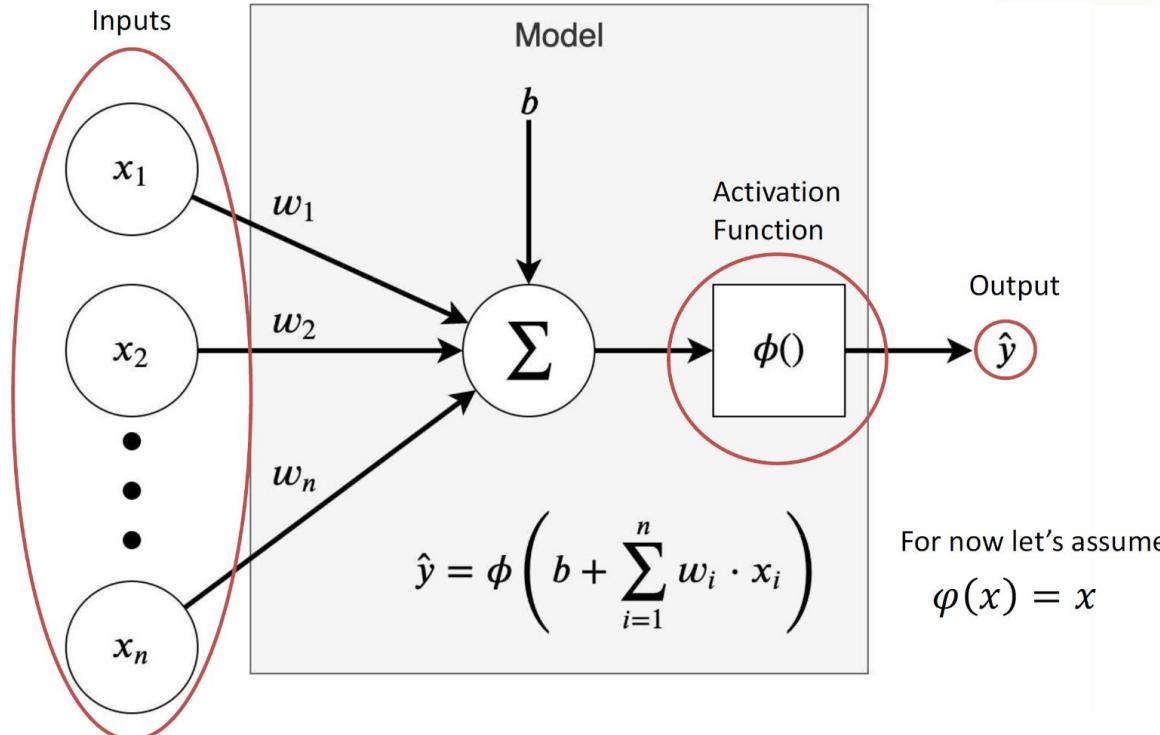
$$\varphi(x) = \frac{e^{2x} - 1}{e^{2x} + 1}$$



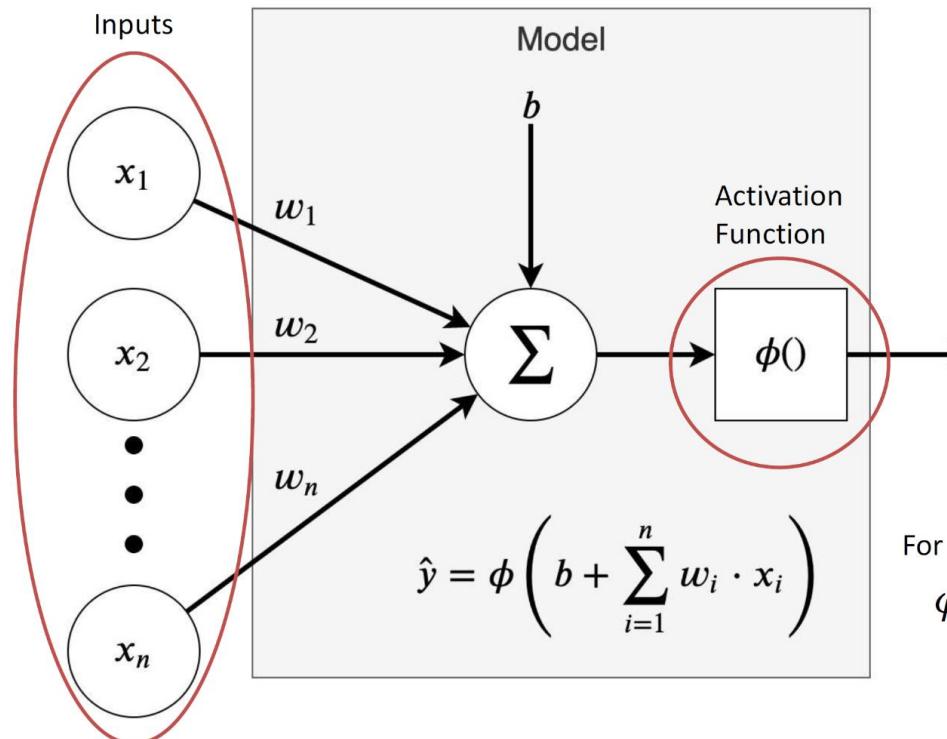
$$\varphi(x) = \max(0, x)$$



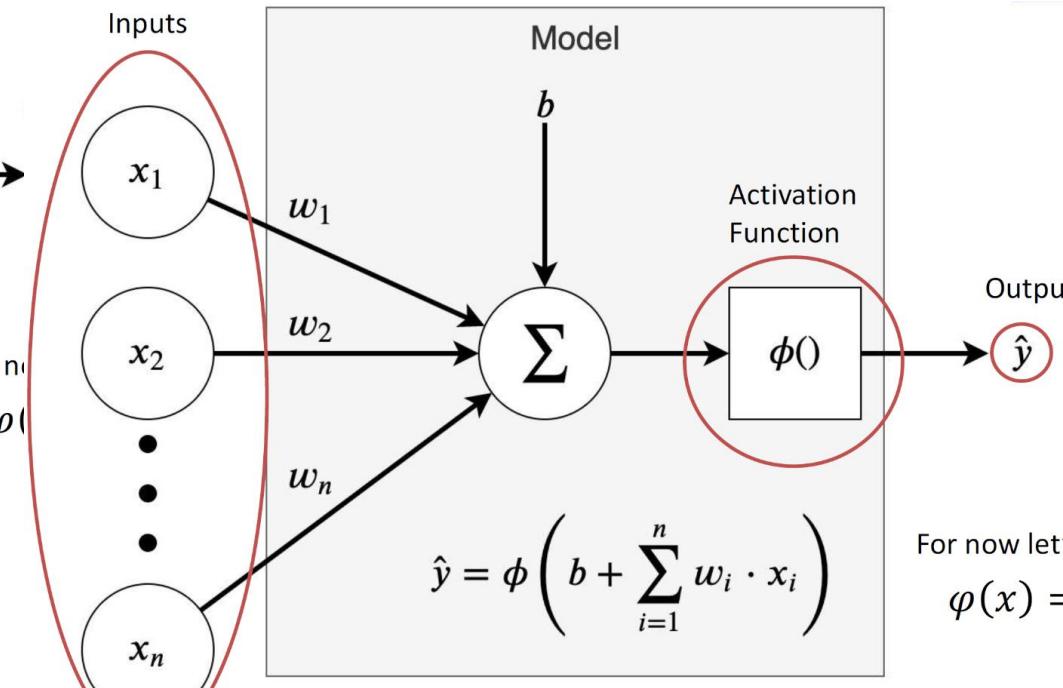
# Artificial Neural Networks



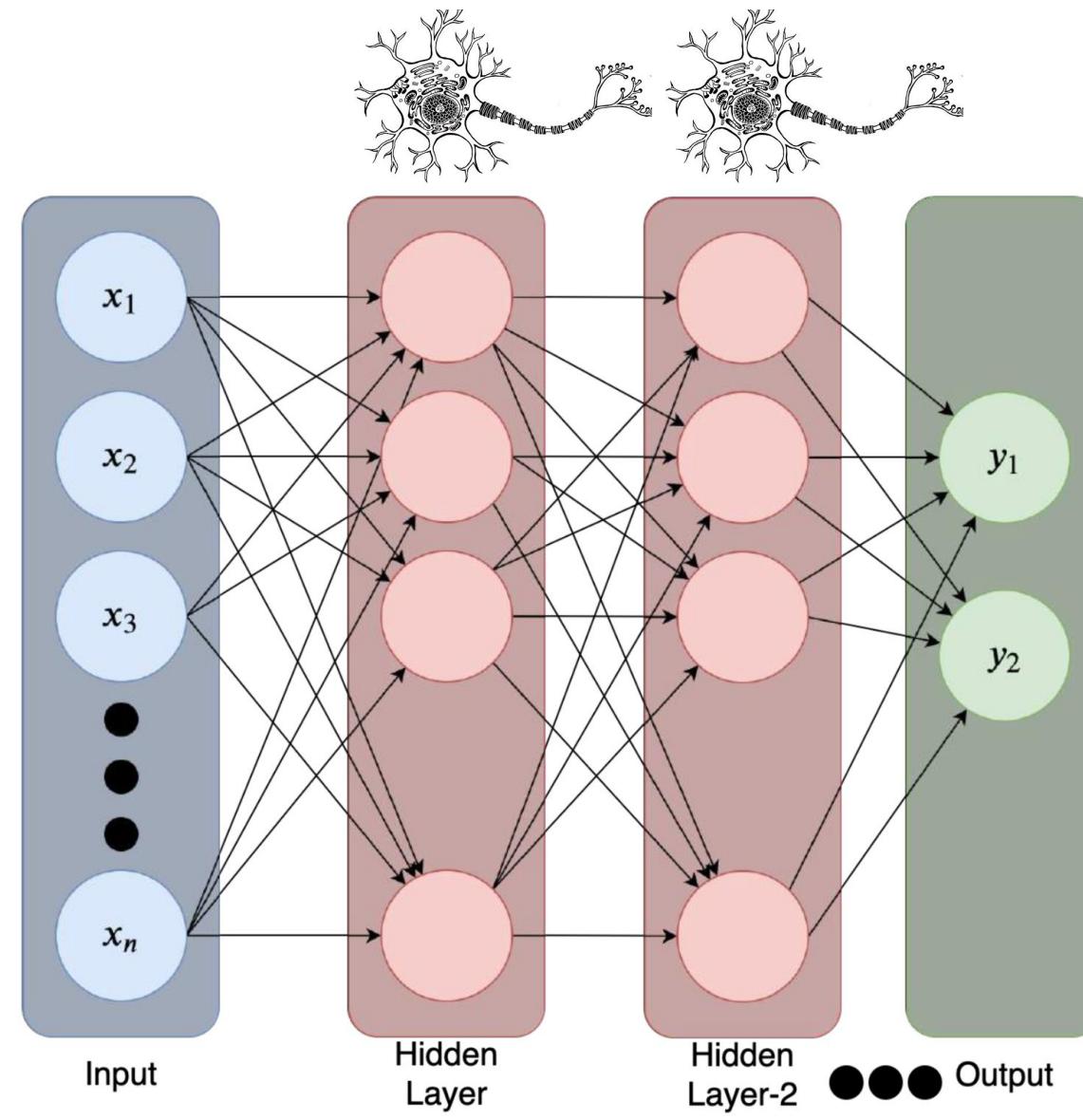
# Artificial Neural Networks



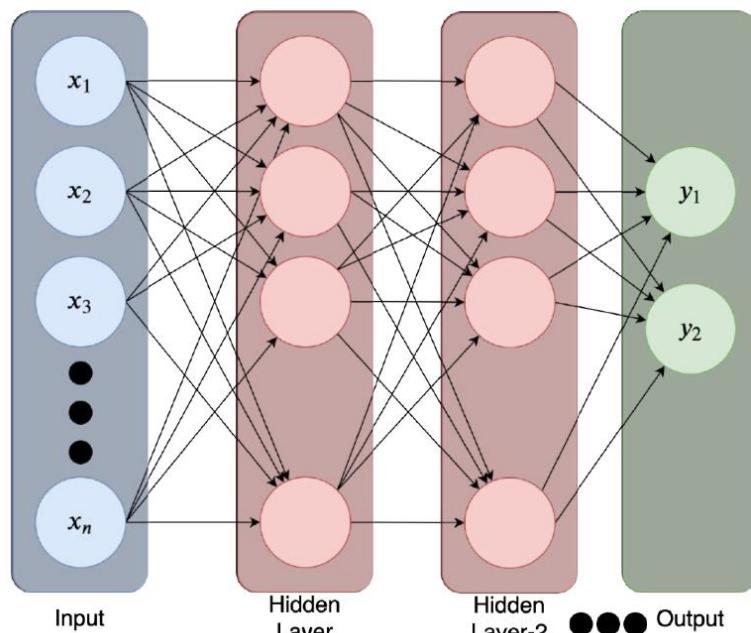
For  
6



For now let's assume



# Training the Artificial Neural Network model



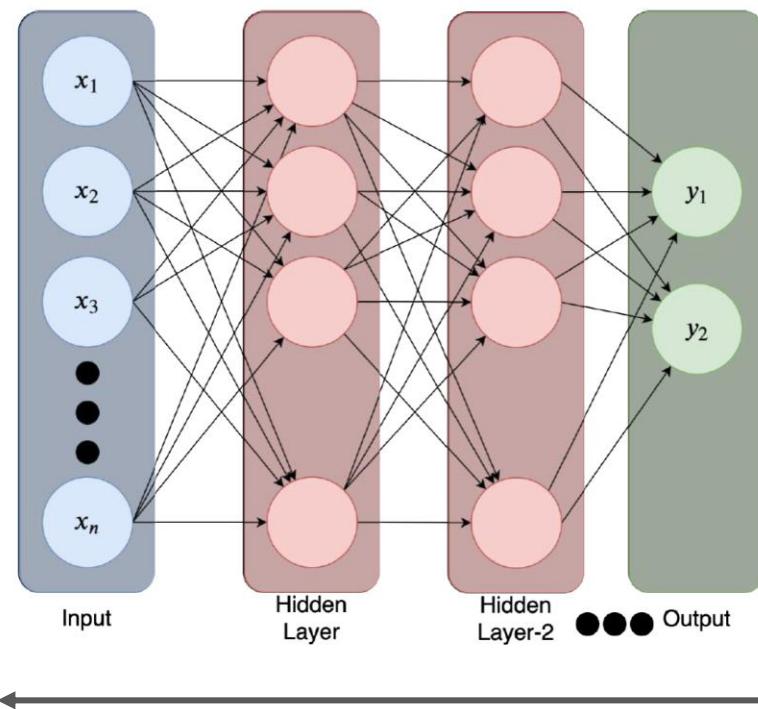
On every layer we multiply the output of the previous layer with the weights and apply activation function  $\varphi(\cdot)$ :

A DOG

$$z^{(l)} = W^{(l)}X^T$$
$$a^{(l)} = \varphi(z^{(l)})$$

Forward pass

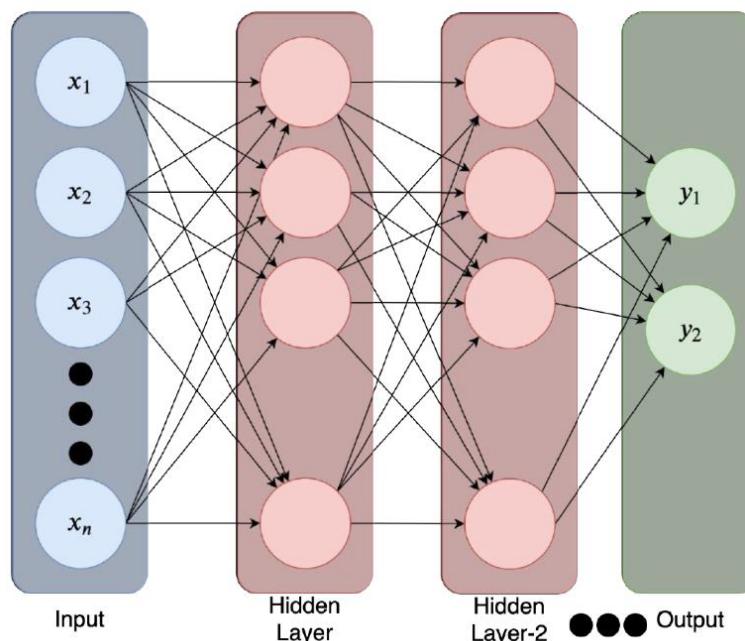
# Training the Artificial Neural Network model



Weights are updated layer by layer to slightly improve the fit between expected and predicted result.

Backpropagation

# Training the Artificial Neural Network model



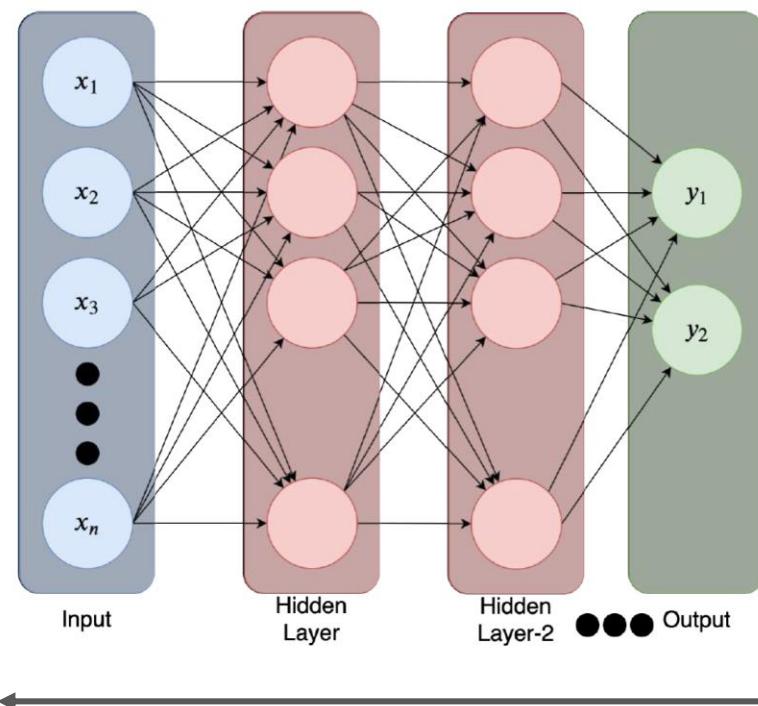
On every layer we multiply the output of the previous layer with the weights and apply activation function  $\varphi(\cdot)$ :

$$z^{(l)} = W^{(l)}X^T$$
$$a^{(l)} = \varphi(z^{(l)})$$

A  
CAT

Forward pass

# Training the Artificial Neural Network model



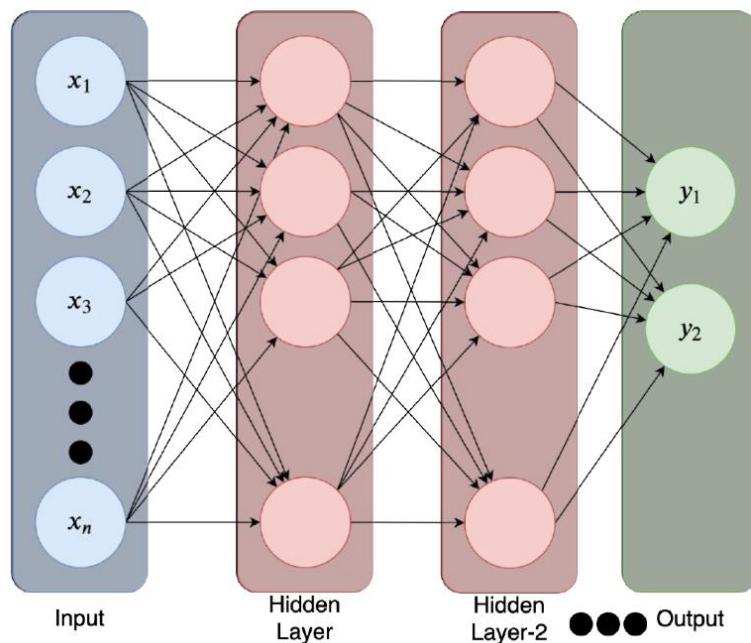
Weights are updated layer by layer to slightly improve the fit between expected and predicted result.

Backpropagation

3 HOURS

LATER . . .

# Training the Artificial Neural Network model



On every layer we multiply the output of the previous layer with the weights and apply activation function  $\varphi(\cdot)$ :

$$z^{(l)} = W^{(l)}X^T$$
$$a^{(l)} = \varphi(z^{(l)})$$

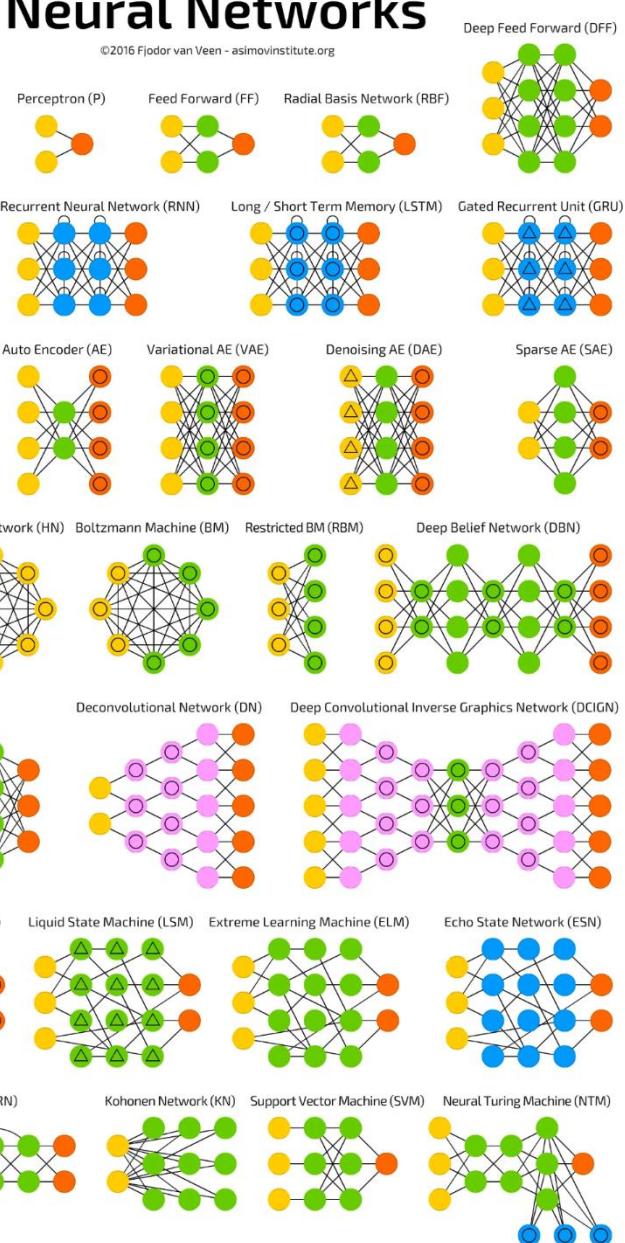
A  
CAT

Forward pass

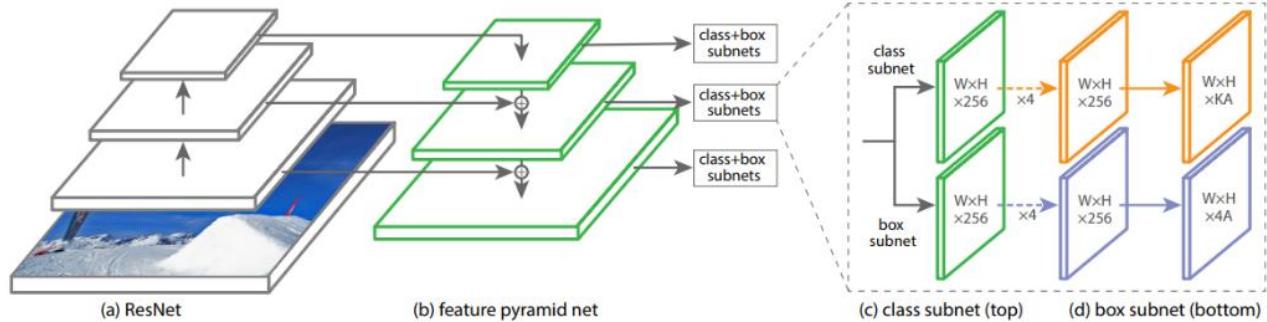
# A mostly complete chart of Neural Networks

- Backfed Input Cell
- Input Cell
- △ Noisy Input Cell
- Hidden Cell
- Probabilistic Hidden Cell
- △ Spiking Hidden Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- △ Different Memory Cell
- Kernel
- Convolution or Pool

©2016 Fjodor van Veen - asimovinstitute.org



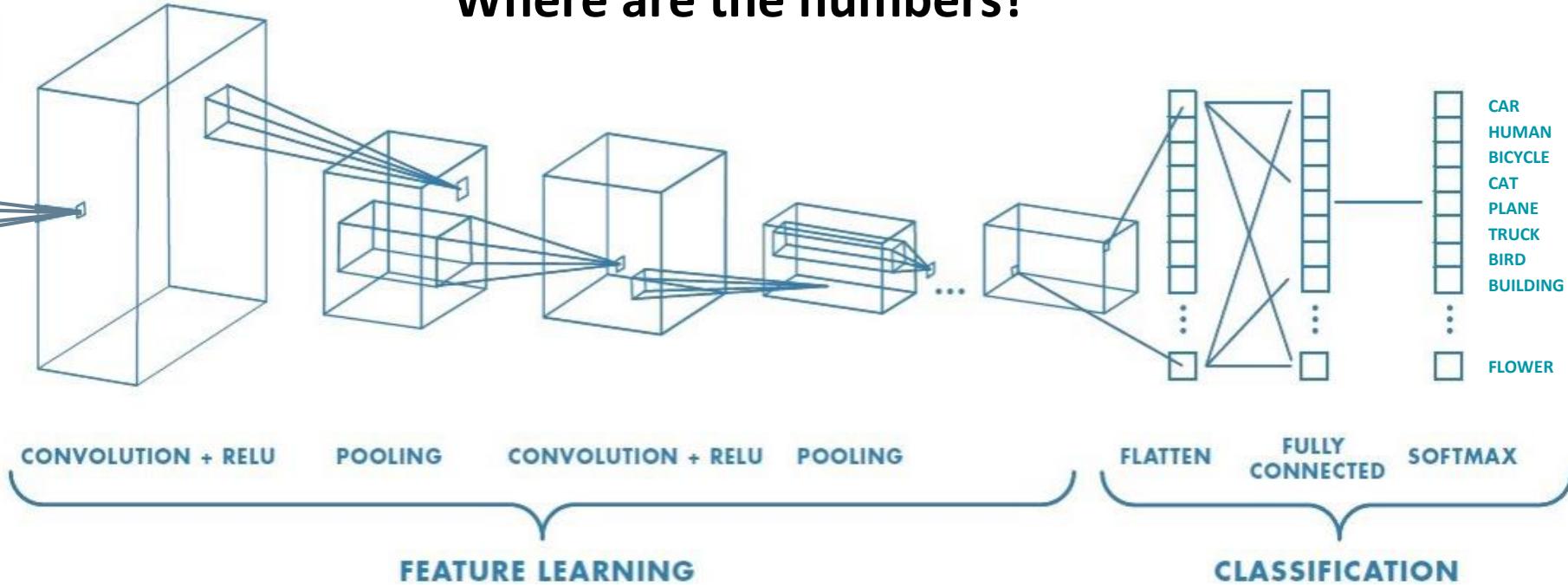
## Custom Architectures - RetinaNet

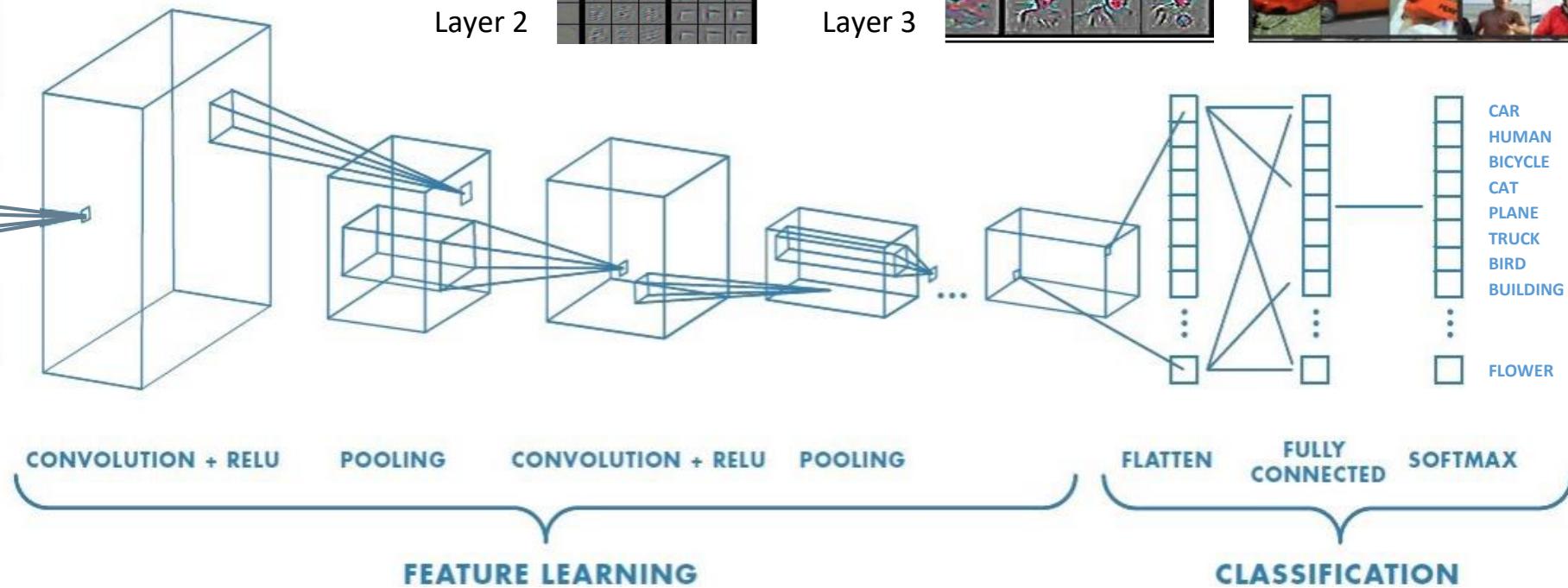
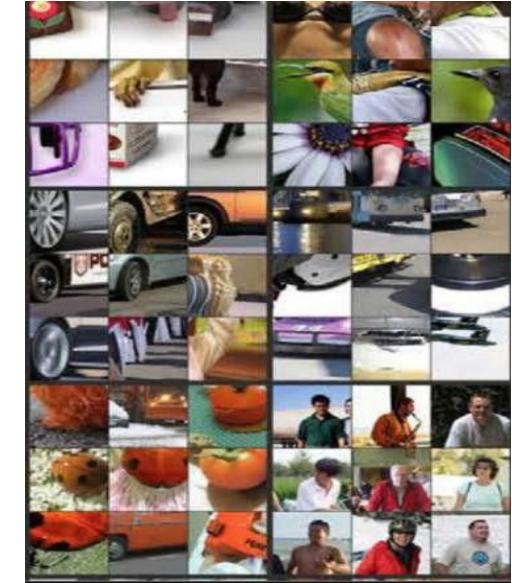
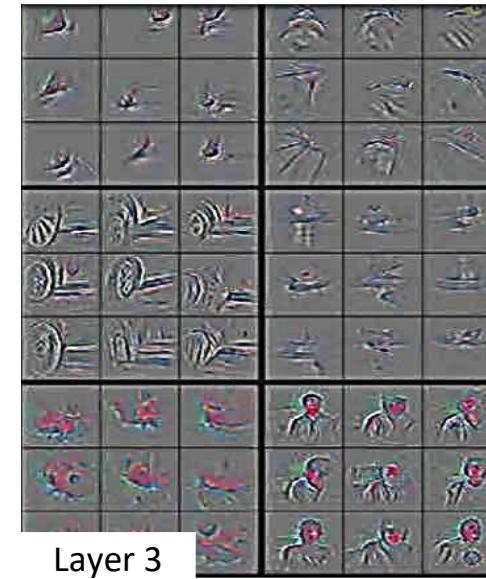
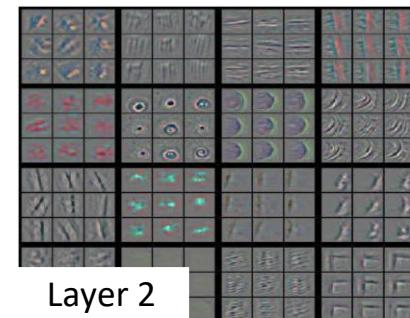




			row	1	2					
			0	.392	.482	.576				
			1	.478	.63	.169	.263	.376		
			2	.580	.79	.263	.44	.306	.376	.451
			0	.373	.60	.376	.443	.478	.561	.674
			1	.443	.569	.443	.443	.569	.674	
			2	.443	.569	.674				

## Where are the numbers?





# Encode a DNA sequence as an image

Input

A	0 0 0 0 0 0 0 0 0 0 1 1 1 0 1 0 1 0 1 0 0 0 0 1
C	0 1 0 0 0 0 1 0 1 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0
G	1 0 1 0 1 0 0 0 0 0 0 0 0 0 0 0 0 1 1 0 0 0
T	0 0 0 1 0 1 0 1 0 0 0 1 0 1 0 0 0 0 1 0

**G C G T G T C T C A A T A T A C G G T A**

# Add more data?

TF binding?	0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 0 0 0 0 0
Input	A 0 0 0 0 0 0 0 0 0 0 1 1 1 0 1 0 1 0 1 0 0 0 0 1
	C 0 1 0 0 0 0 0 1 0 1 0 0 0 0 0 0 0 1 0 0 0 0 0 0
	G 1 0 1 0 1 0 0 0 0 0 0 0 0 0 0 0 0 1 1 1 0 0
	T 0 0 0 1 0 1 0 1 0 0 0 1 0 1 0 0 0 0 1 0

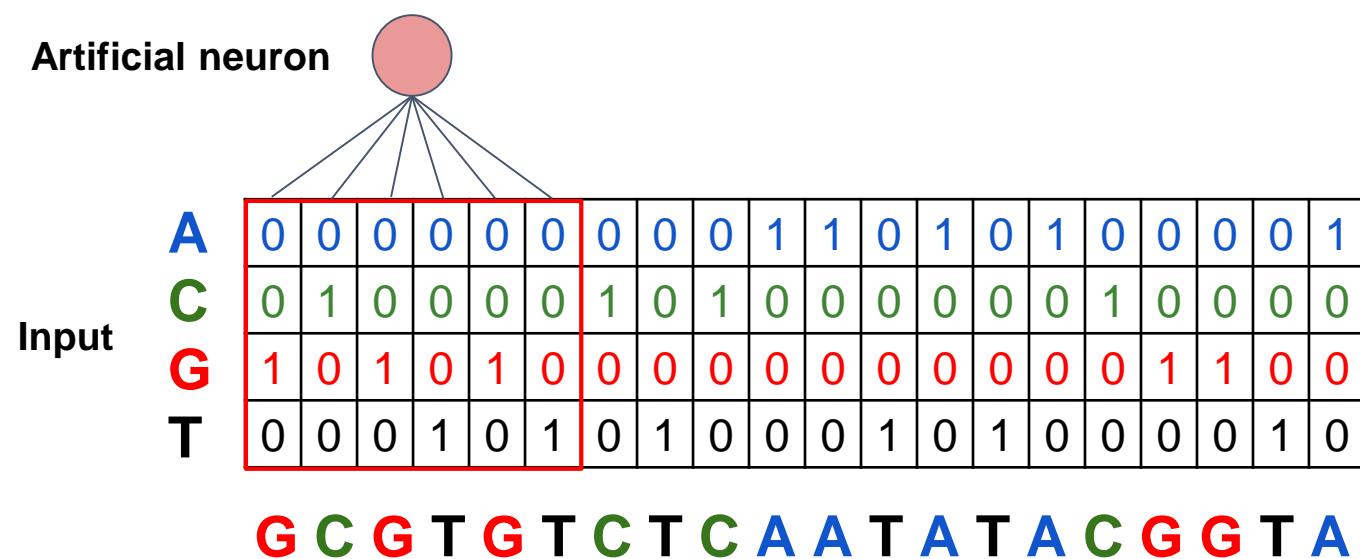
**G C G T G T C T C A A T A T A C G G T A**

# Add more data?

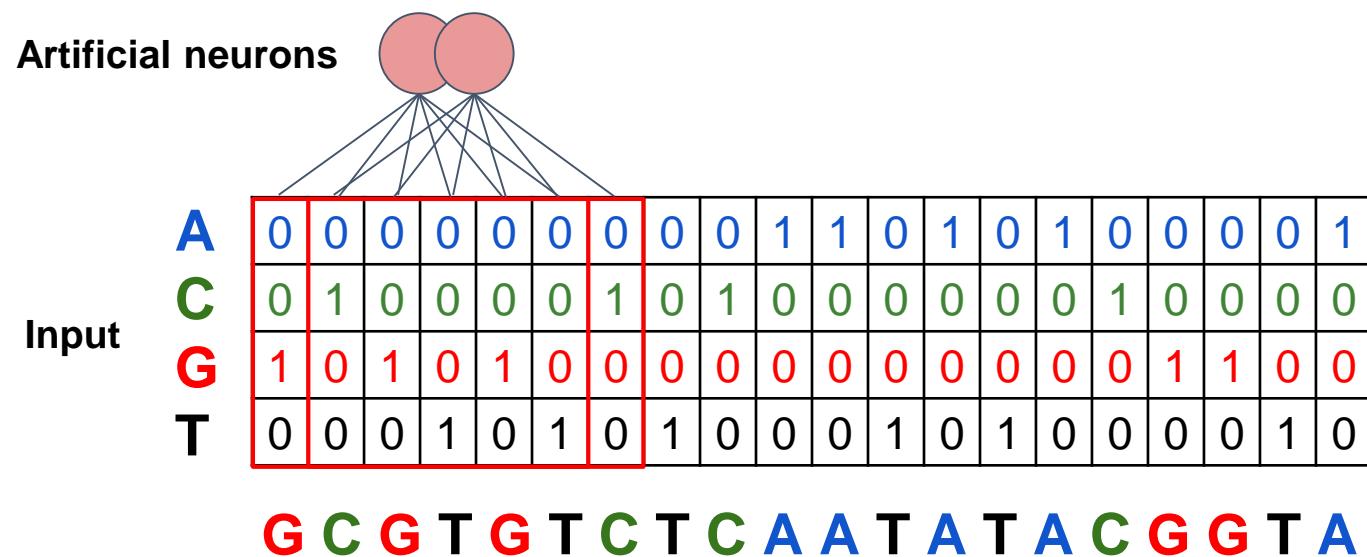
Methylation?	0   0   0   0   0   0   0   0   0   1   0
TF binding?	0   0   0   1   0
Input	A C G T
	0   0   0   0   0   0   0   0   0   1   1   0   1   0   1   0   1   0   0   0   0   0   1   1   0   0   0   0   1
	0   1   0   0   0   0   1   0   1   0   0   0   0   0   0   0   1   0   0   0   0   0   1   0   0   0   0   0   0
	1   0   1   0   1   0   0   0   0   0   0   0   0   0   0   0   0   0   0   0   0   0   1   1   0   0   0   0   0
	0   0   0   1   0   1   0   1   0   0   0   0   1   0   1   0   1   0   0   0   0   0   1   0   0   0   0   1   0

**G C G T G T C T C A A T A T A C G G T A**

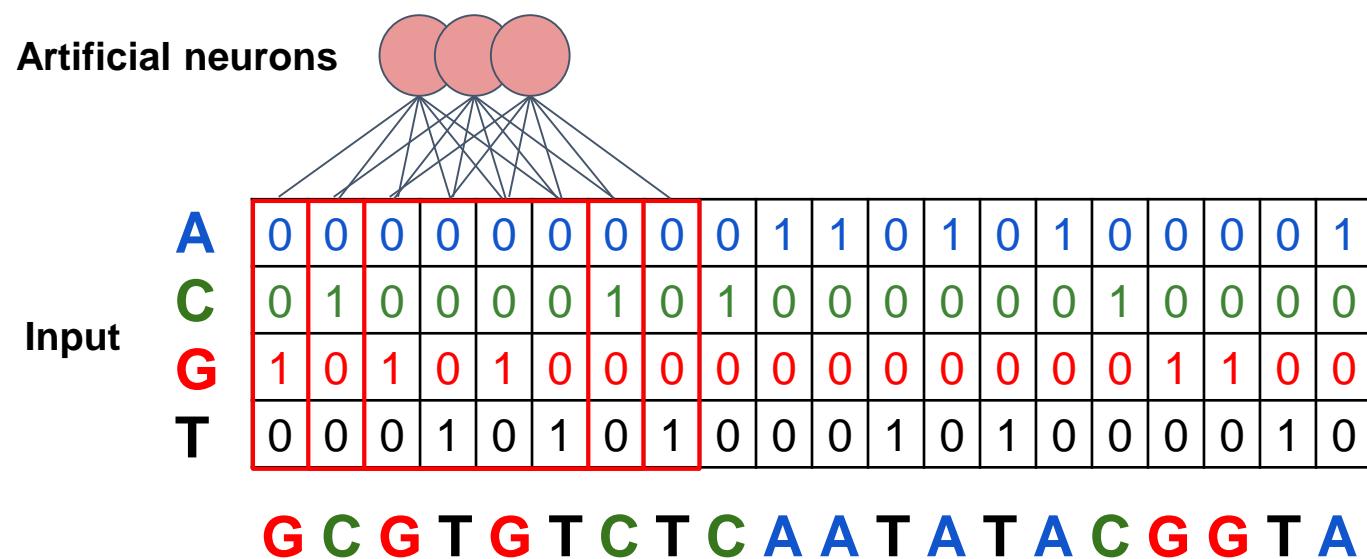
# Build the artificial neural network



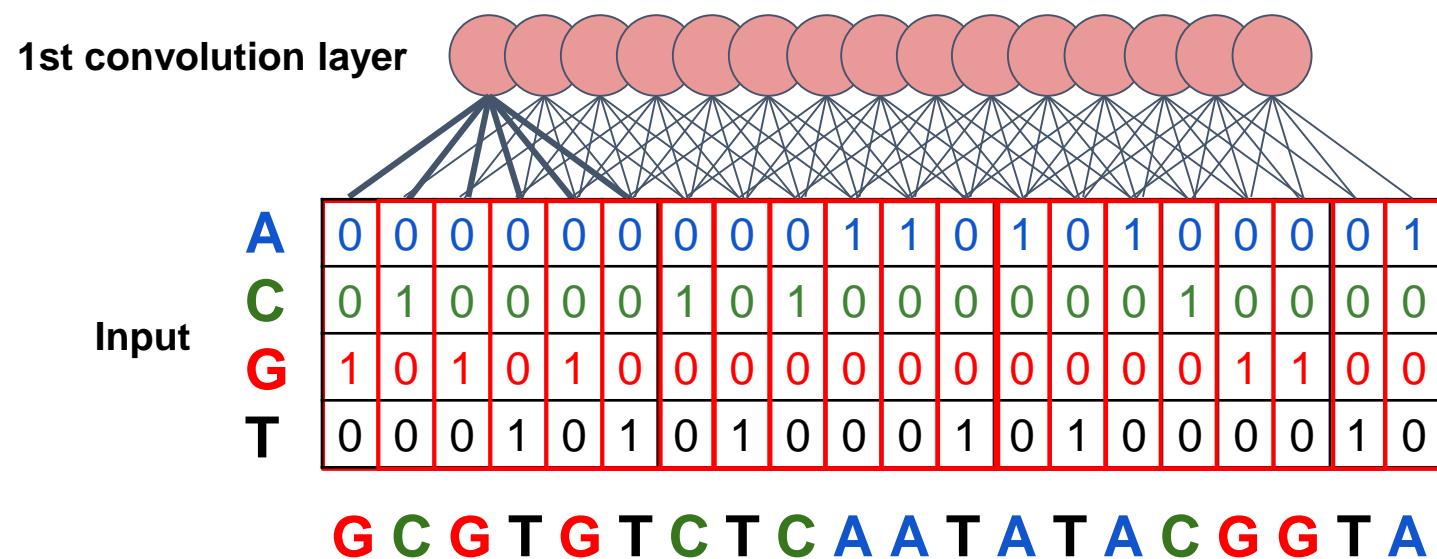
# Build the artificial neural network



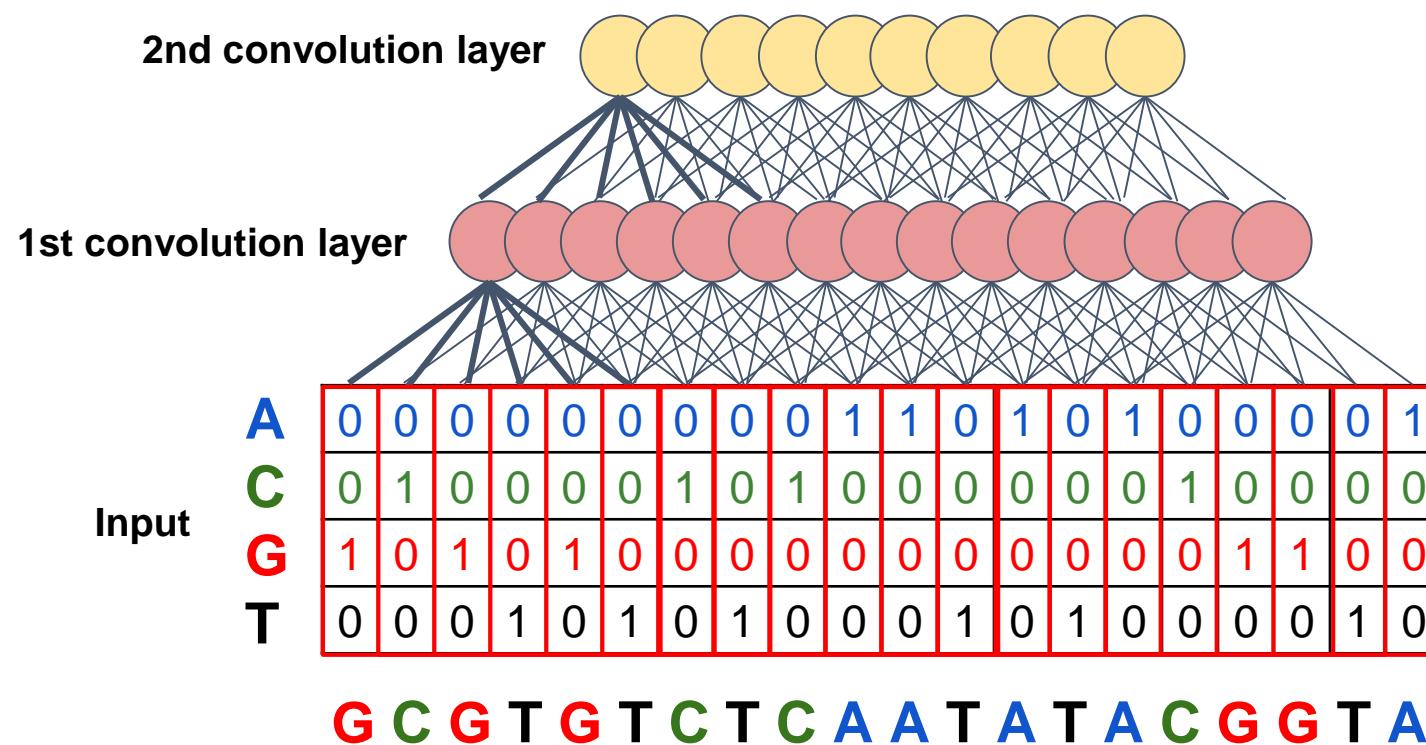
# Build the artificial neural network

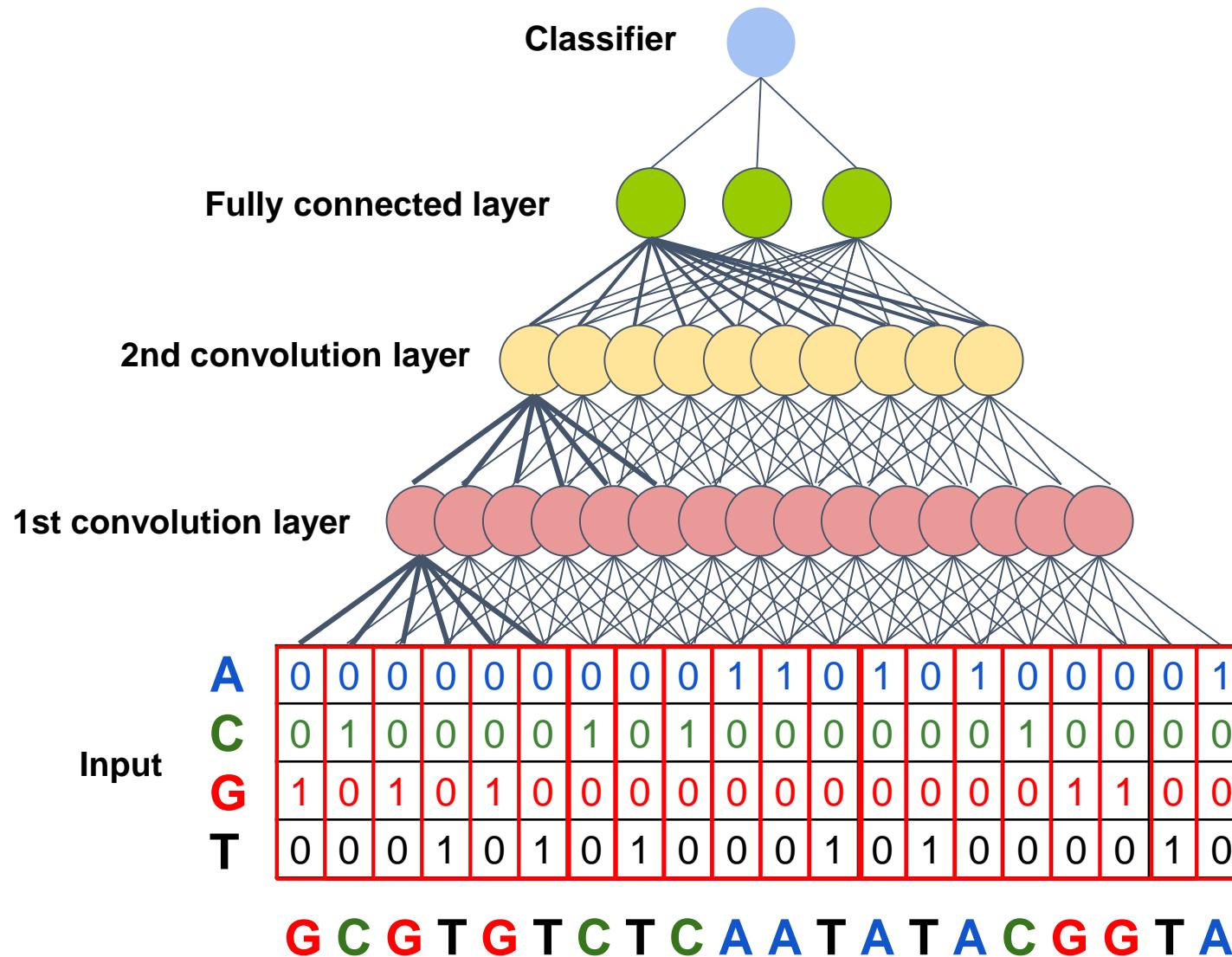


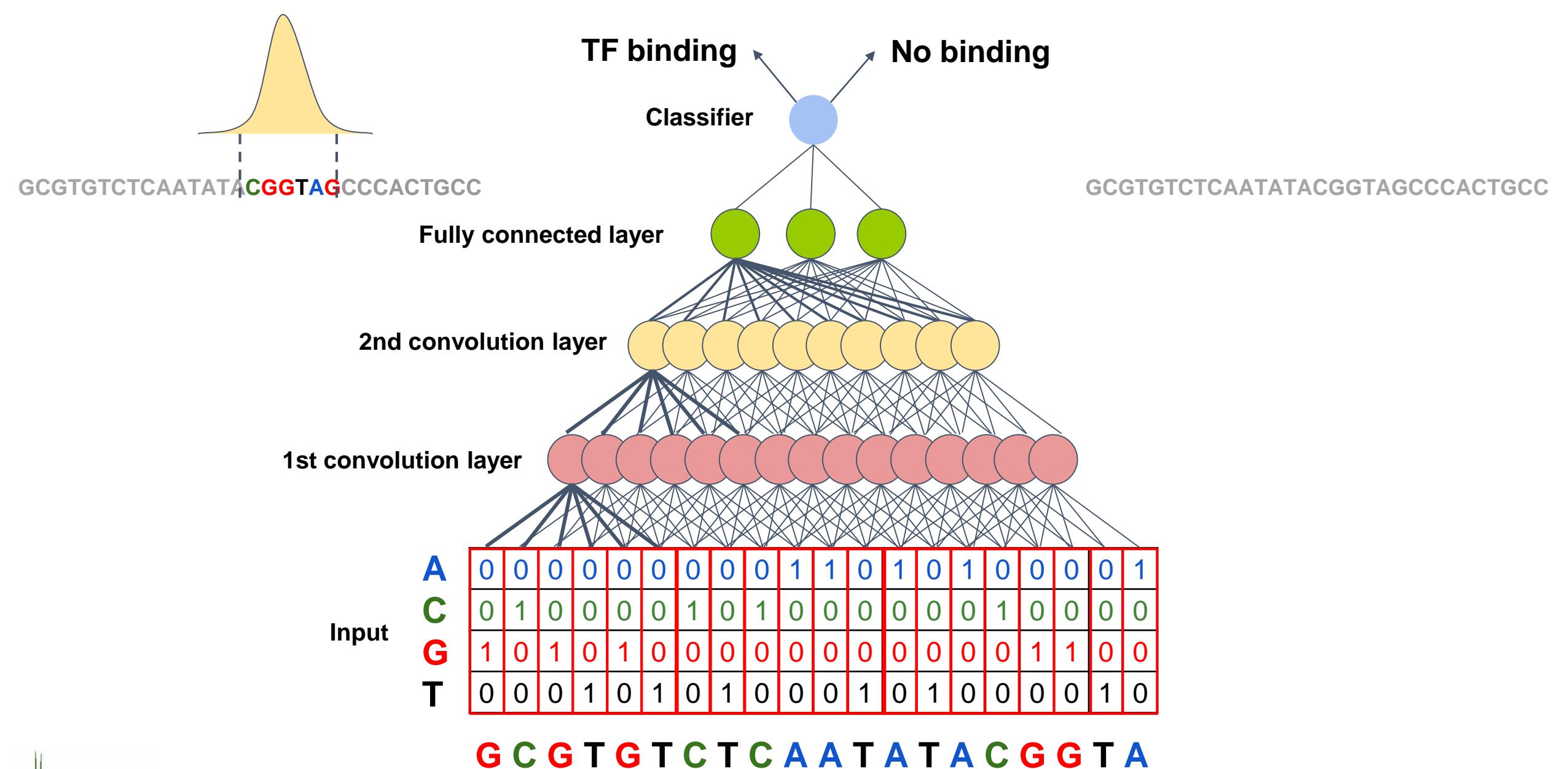
# The first convolution layer

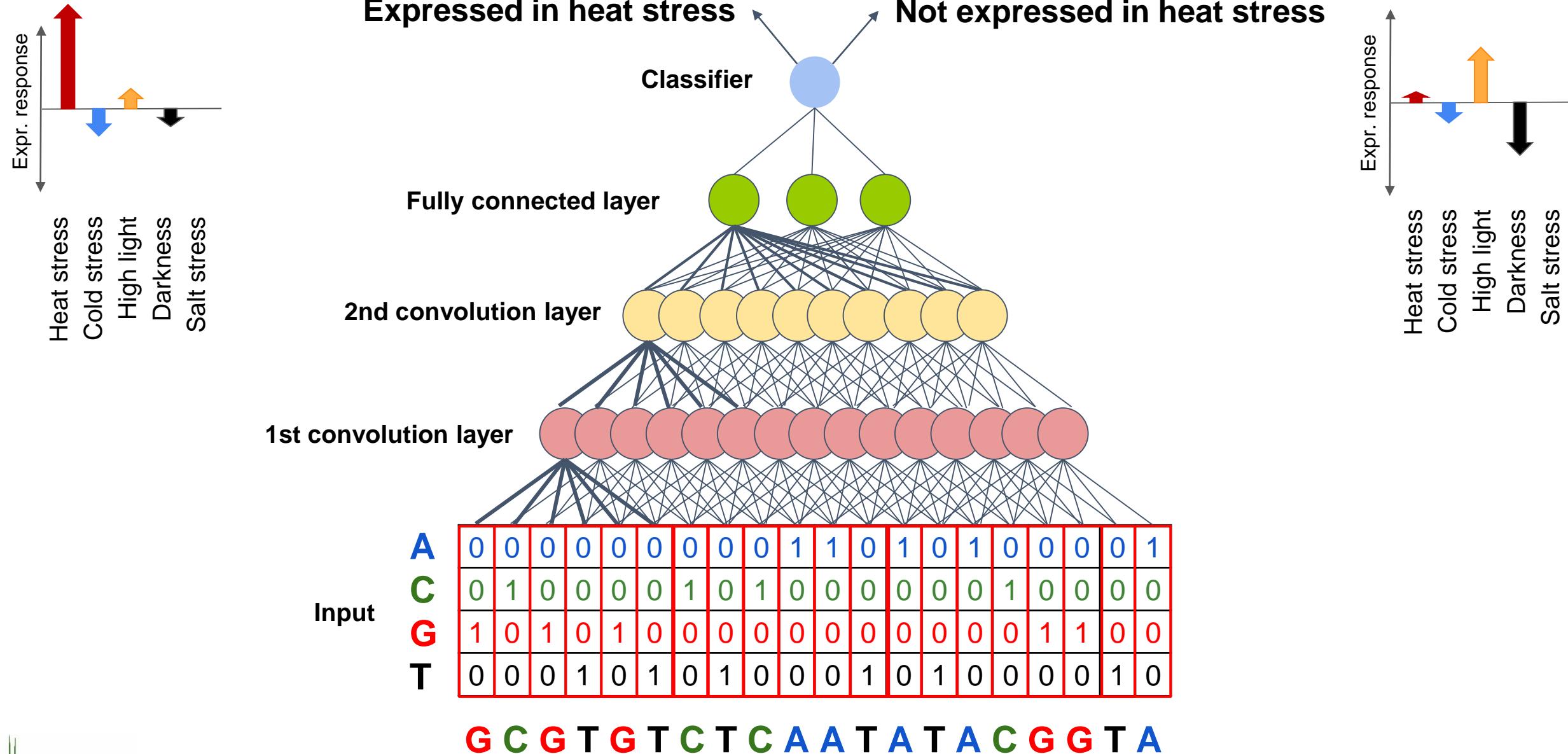


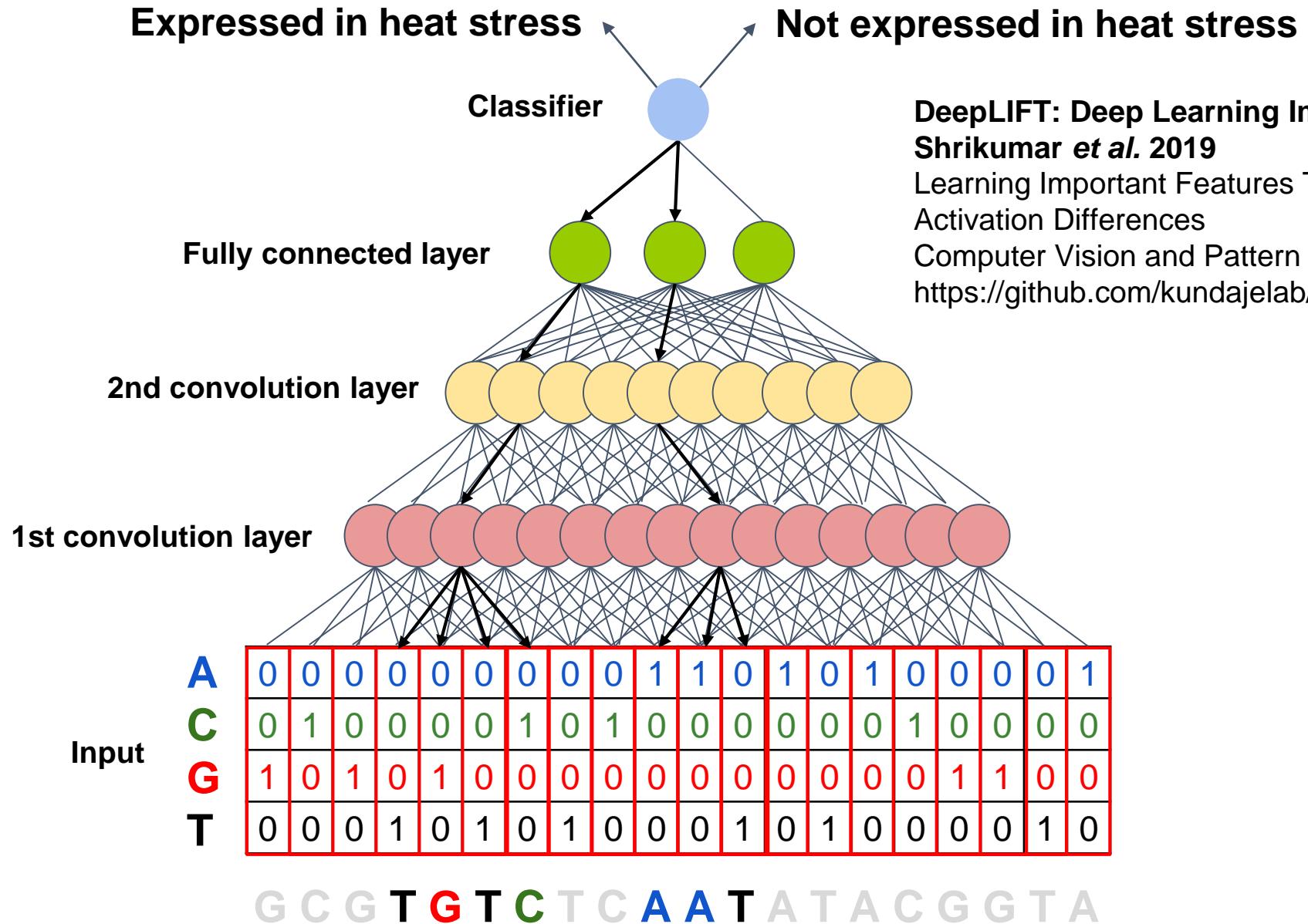
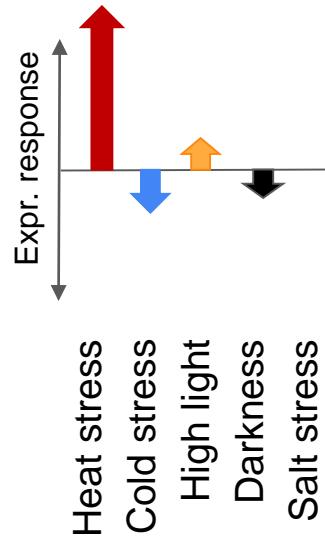
# Add more convolution layers





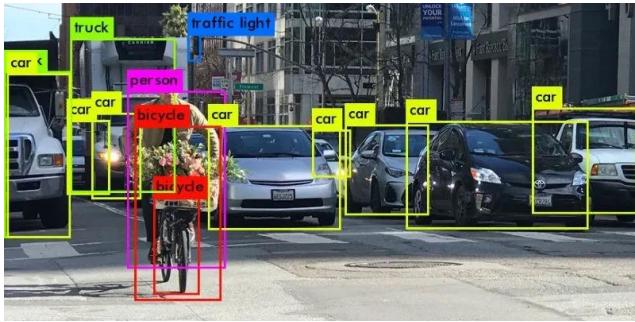






**DeepLIFT: Deep Learning Important FeaTures**  
**Shrikumar et al. 2019**  
Learning Important Features Through Propagating Activation Differences  
Computer Vision and Pattern Recognition  
<https://github.com/kundajelab/deeplift>

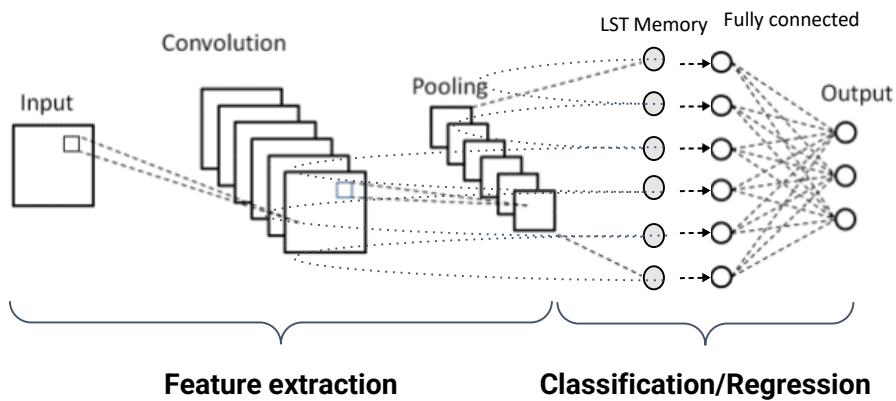
# Implementation



gene  
sequences

## CNN-LSTM

Convolutional Neural Networks - Long Short Term Memory



*Arabidopsis thaliana*



*Solanum lycopersicum*



*Sorghum bicolor*



*Zea mays*



expression values

# Implementation

Can we detect **cis**-regulatory sequence features predictive for gene expression?



Simon Zumkeller



Fritz Peleke



*Arabidopsis thaliana*



*Solanum lycopersicum*



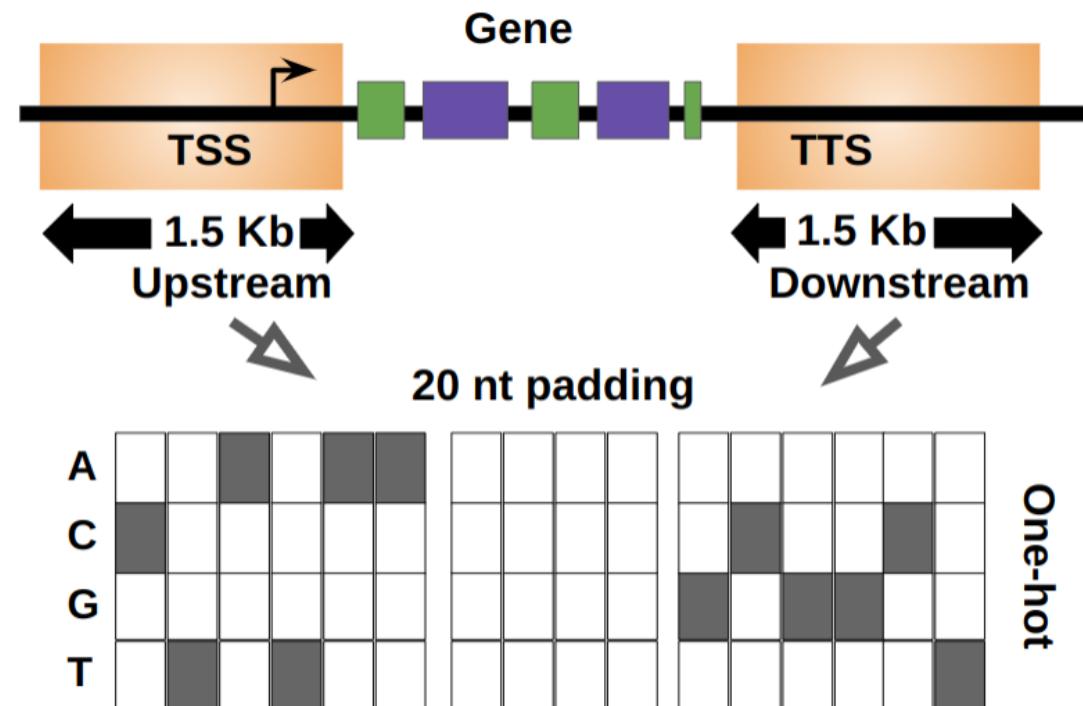
*Sorghum bicolor*



*Zea mays*

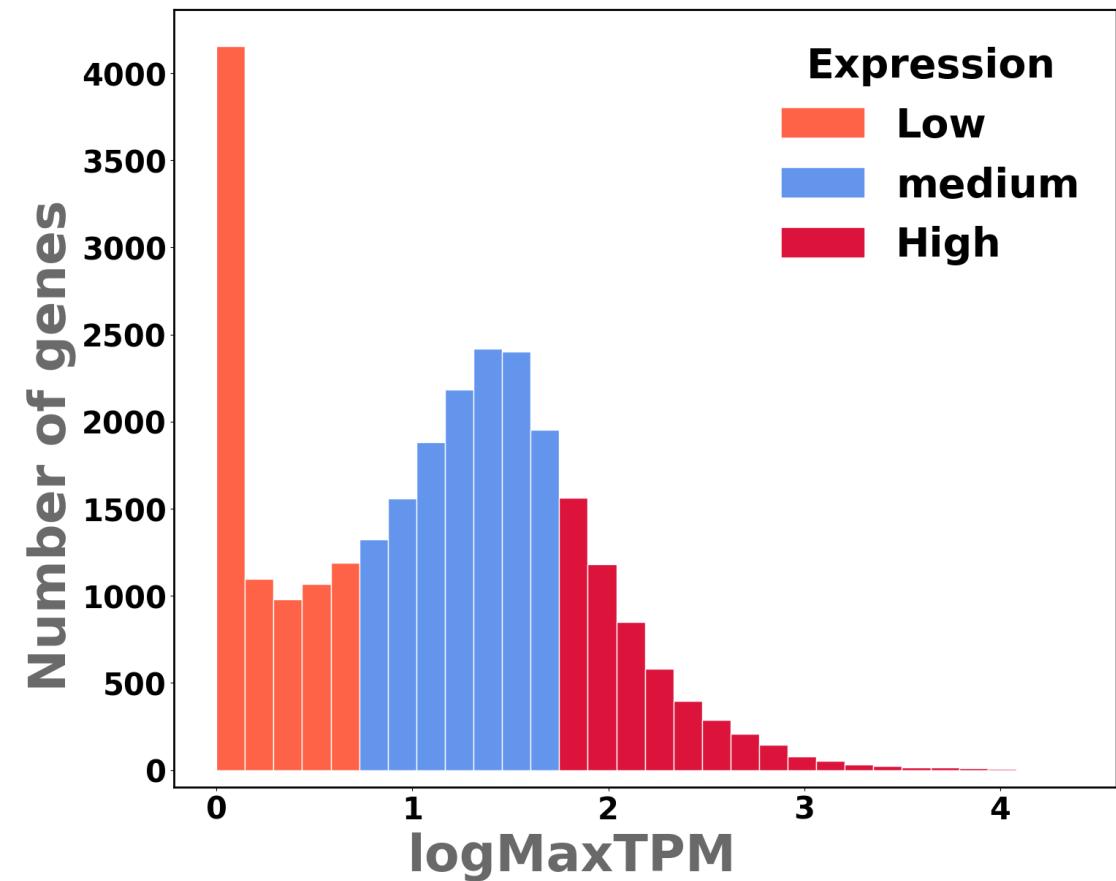
# Input data

Can we detect cis-regulatory sequence features predictive for gene expression?

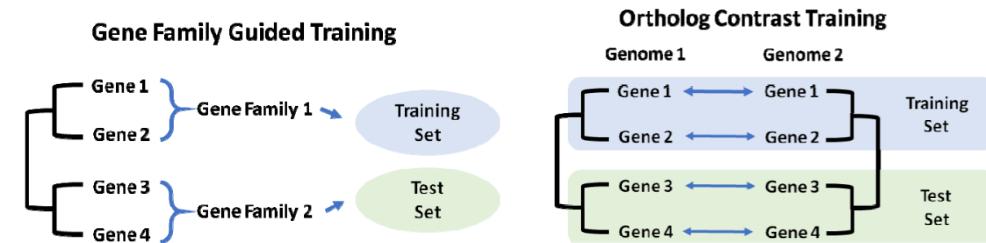
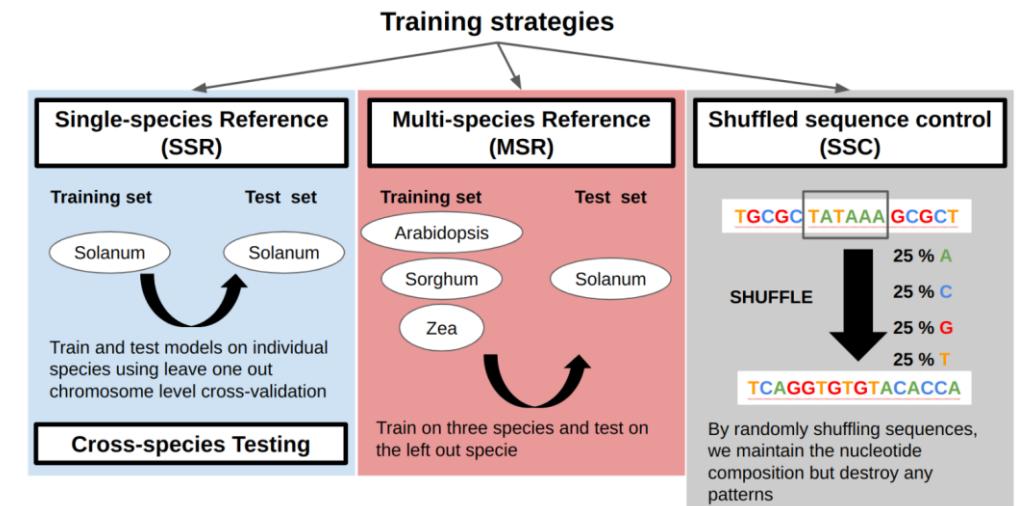
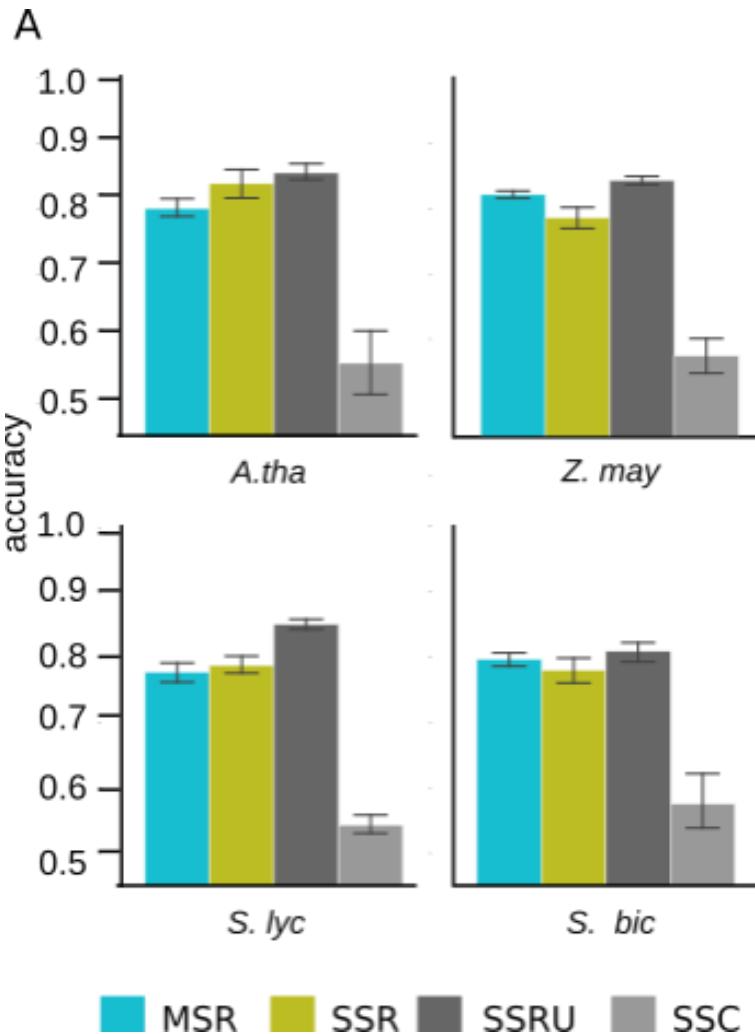


# Input data

Can we detect cis-regulatory sequence features predictive for gene expression?

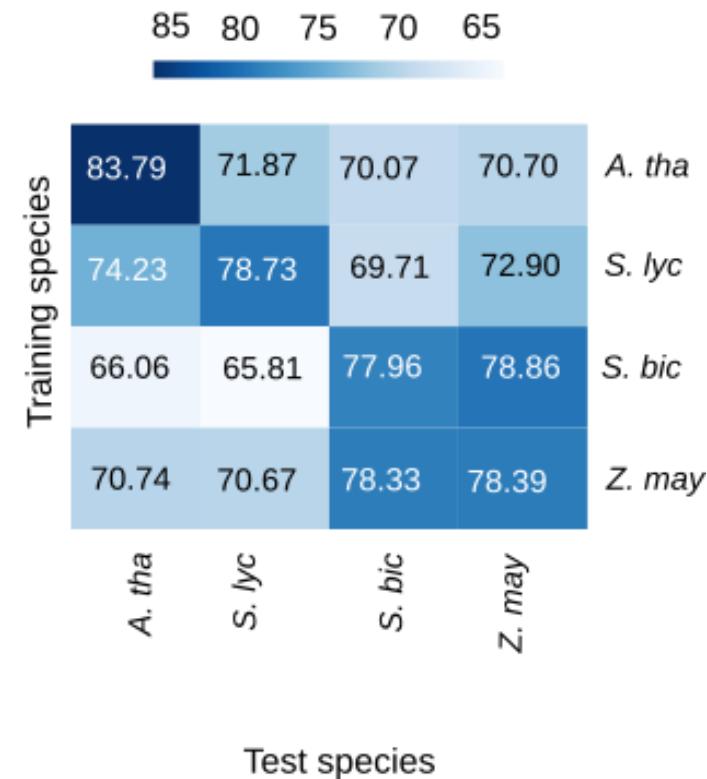
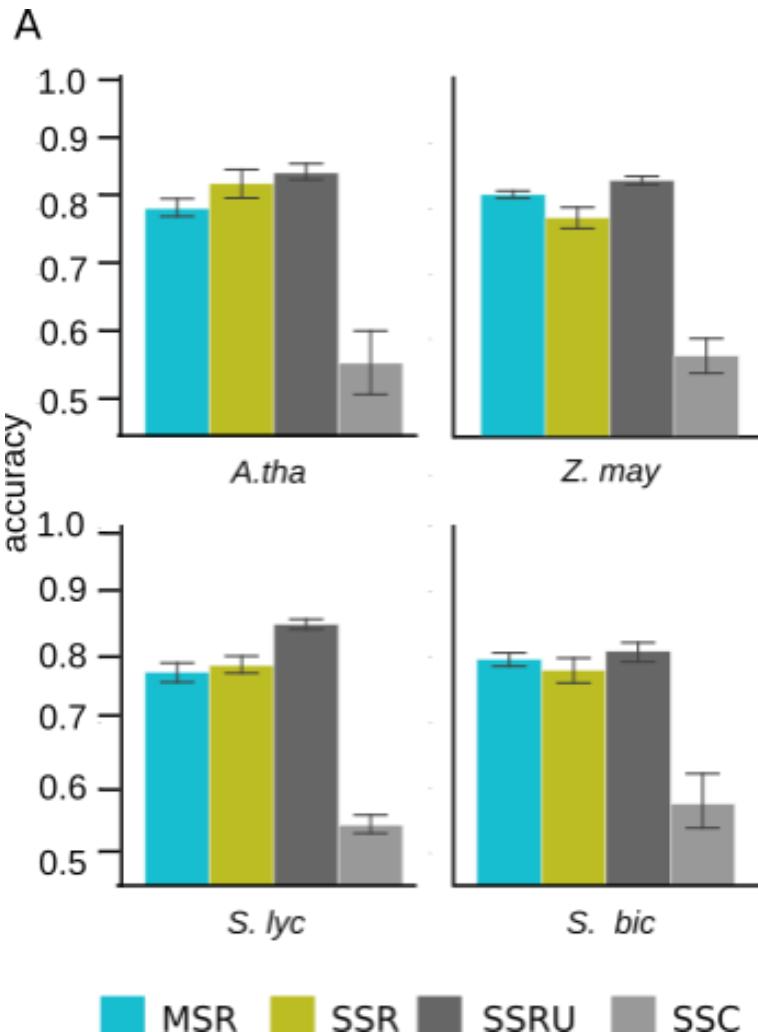


# Prediction

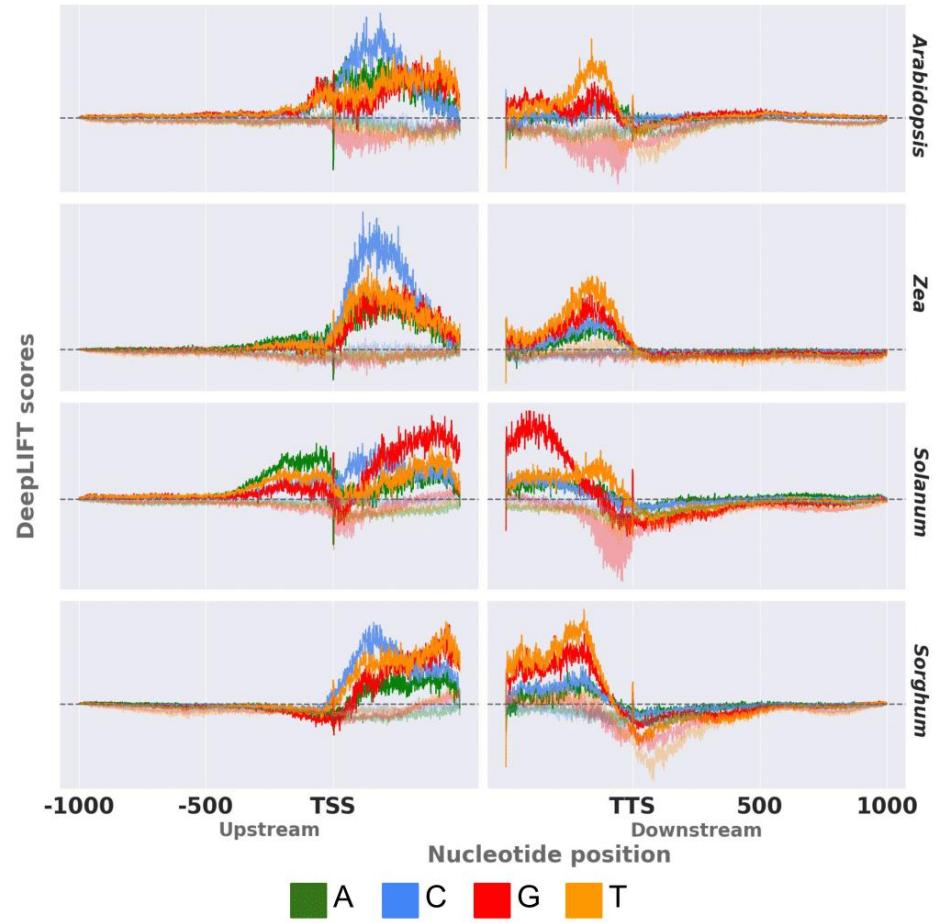
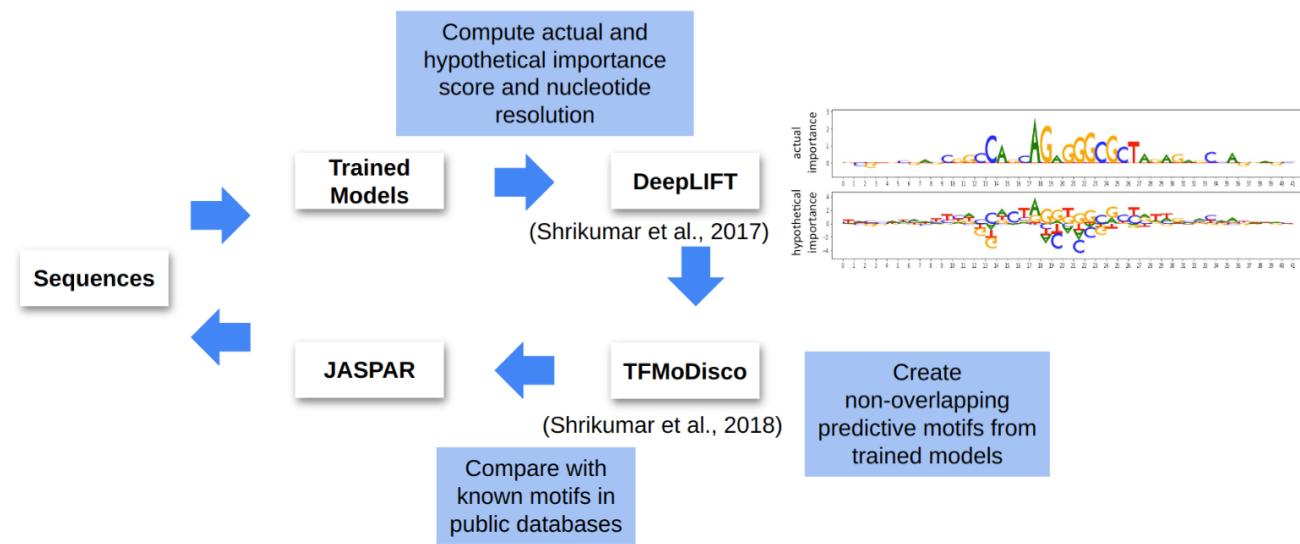


Washburn et al. 2019 PNAS

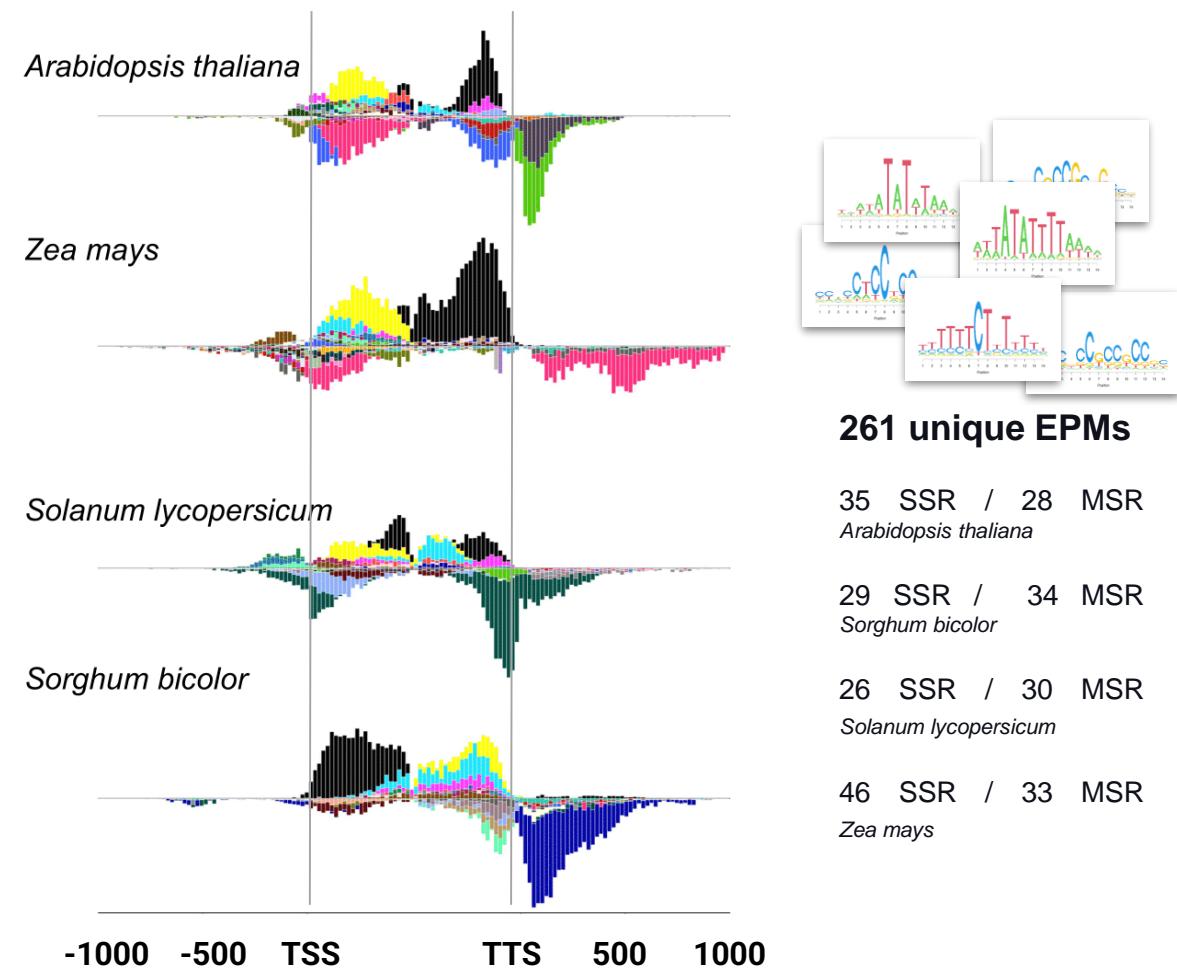
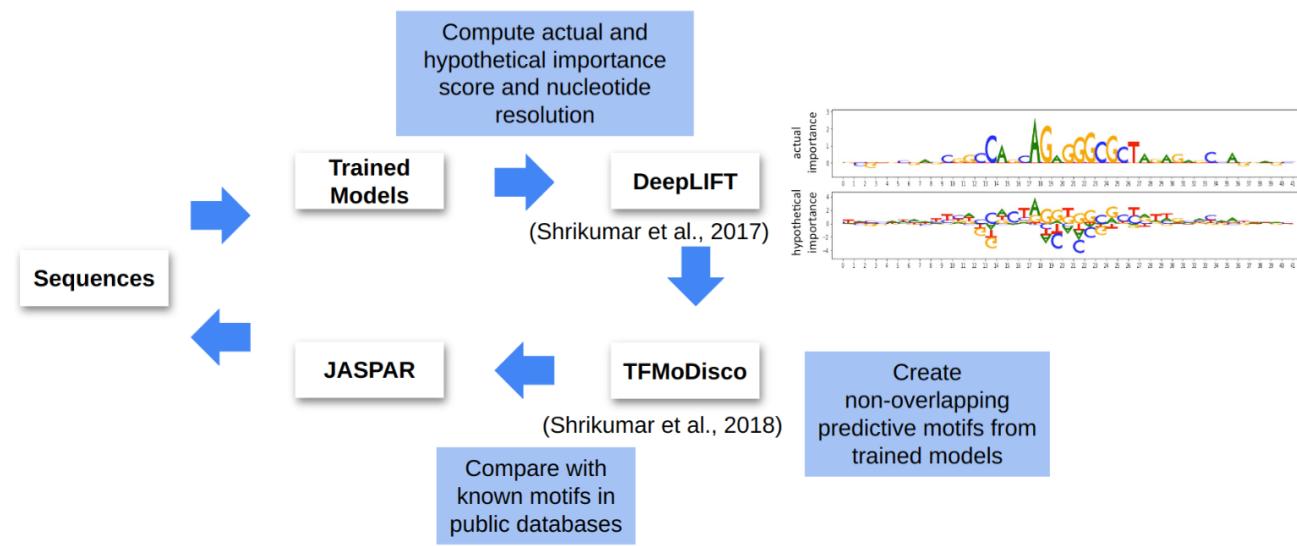
# Prediction



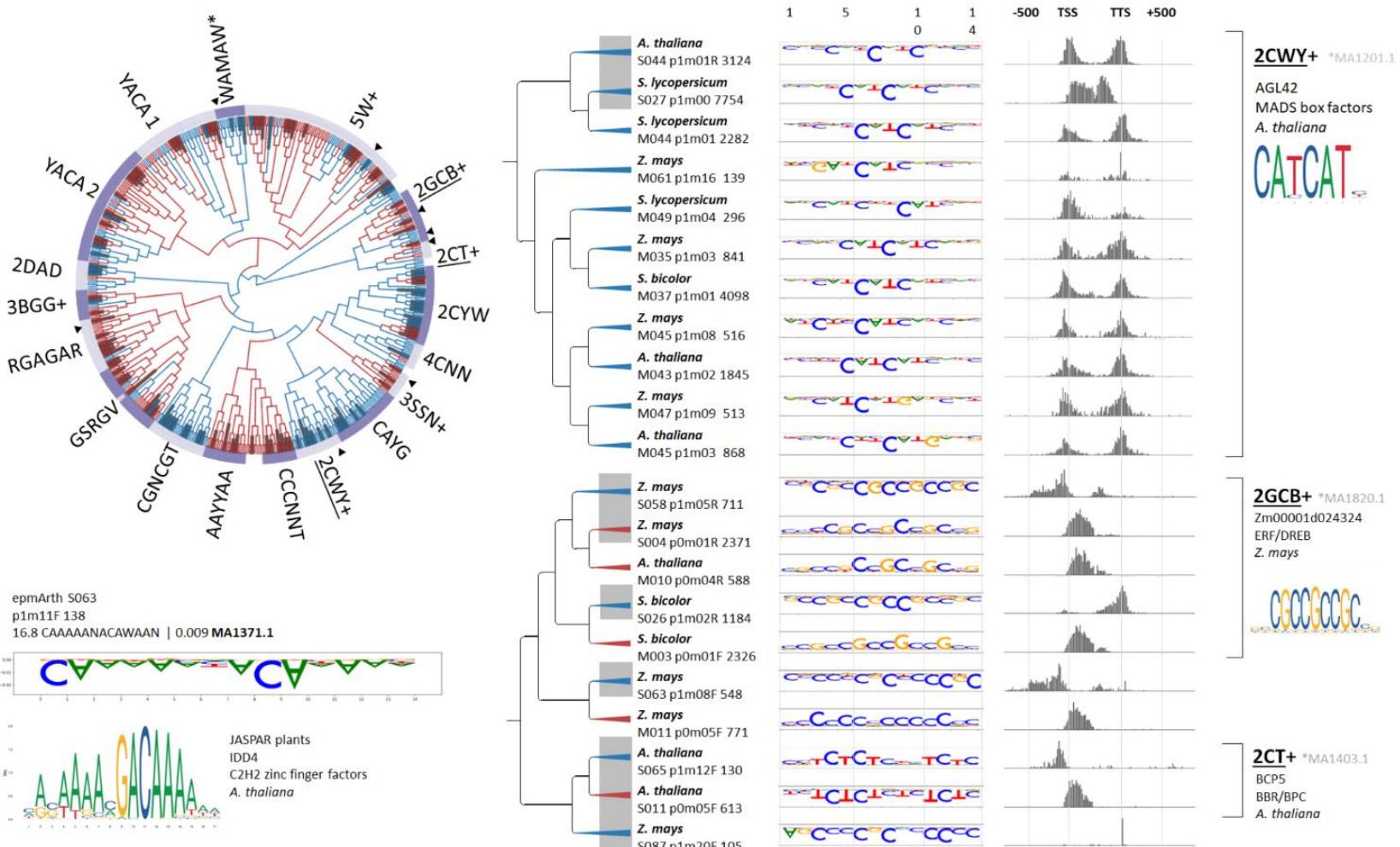
# Feature selection



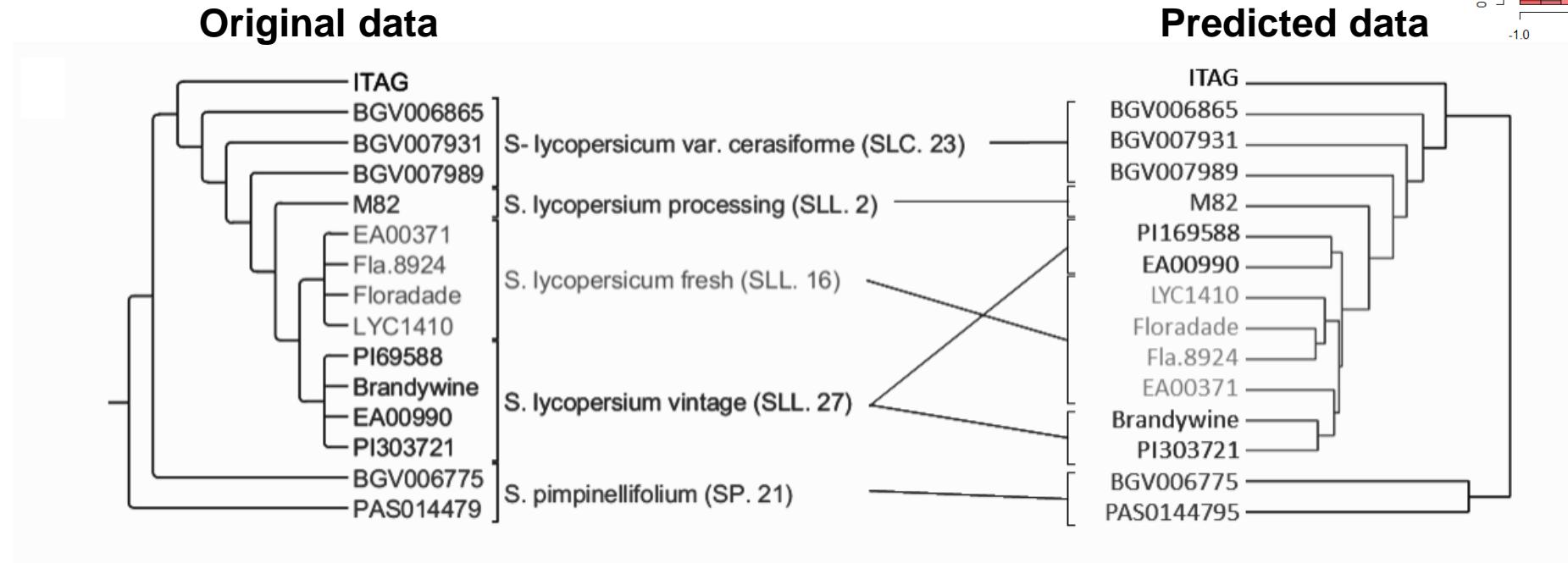
# Feature selection



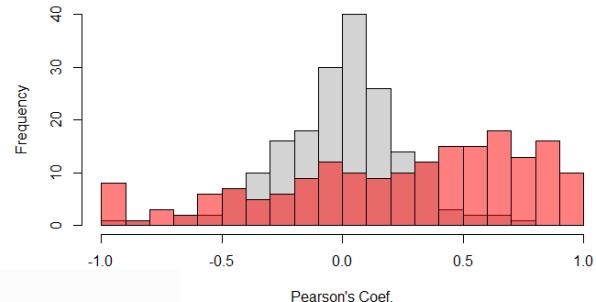
# Expression Predictive Motifs (EPM)



# Annotate effects of genetic variation



Cor MSR, 15 genotypes

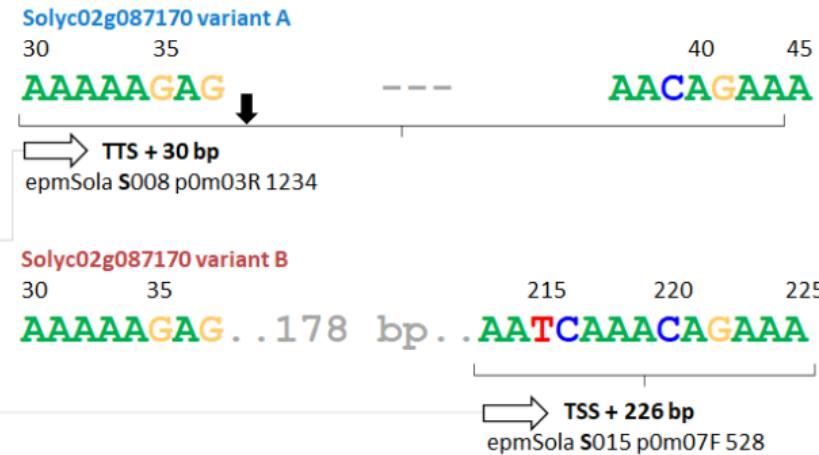
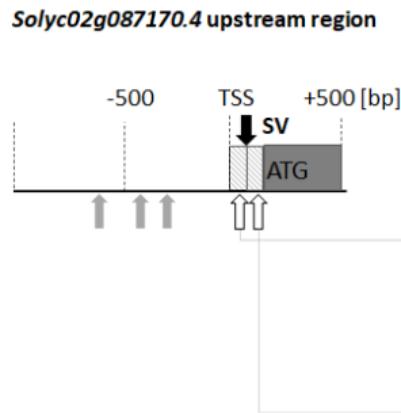
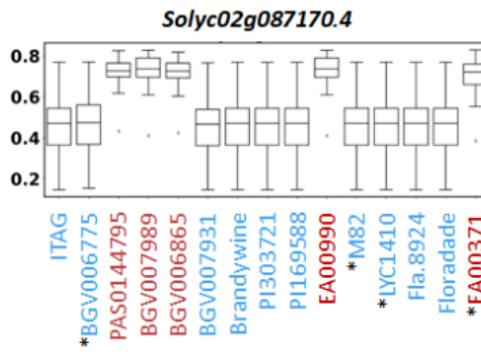


Prediction accuracies  
>75%

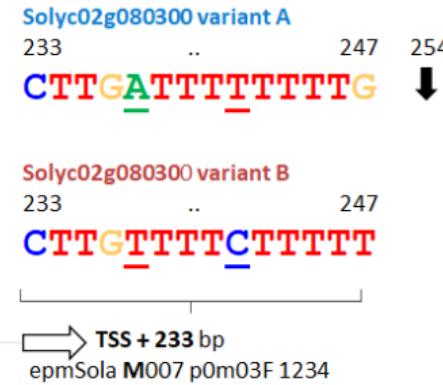
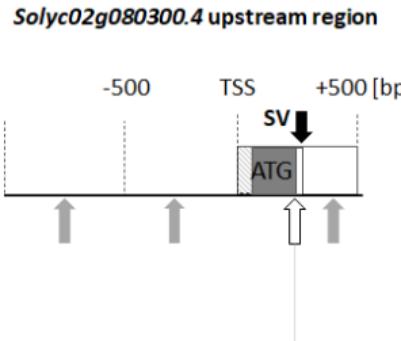
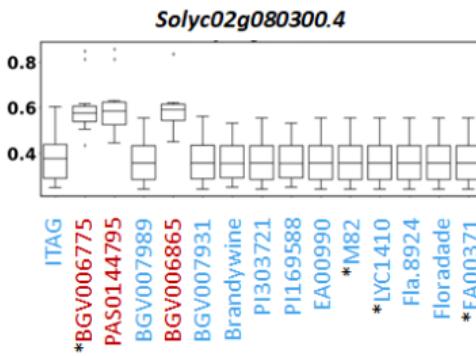
MSR model that  
never saw tomato

Alonge *et al* 2020, Major Impacts of Widespread Structural Variation on Gene Expression and Crop Improvement in Tomato. *Cell.* 182, 145–161.e23 (2020).

# Annotate effects of genetic variation



Introduction of functional EPM



Edit from neutral to functional EPM

# DeepCRE - framework and toolbox: <https://deepcre.ipk-gatersleben.de/>

The screenshot displays the DeepCRE framework and toolbox interface. On the left, a sidebar shows genome files (Arabidopsis\_thaliana.TA... 36.5MB and 11.2MB), genes (genes.csv 50.0KB), and a deepCRE model (arabidopsis\_thaliana\_leaf). The main area has an 'analysis' dropdown set to 'manual'. A gene 'AT1G79700' is selected. The 'Mutate Promoter' section shows a sequence from position 0 to 1500 with a 'submit' button. The 'Mutate Terminator' section shows a sequence from position 1520 to 3020 with a 'submit' button. Below these are target regions and 'Mutate' / 'Reset' buttons. At the bottom is a bar chart comparing 'Probability of high expression' between 'AT1G79700' and 'AT1G79700: Mutated' (both ~0.5) and a line graph of 'Saliency Score' across the sequence, with peaks at TSS and TTS.



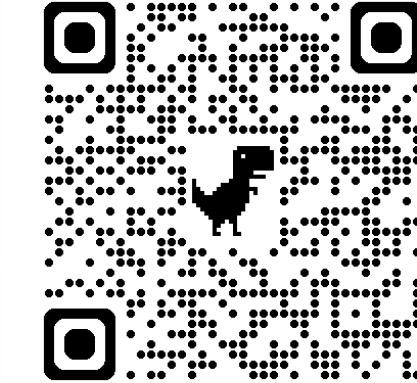
Merle Stein



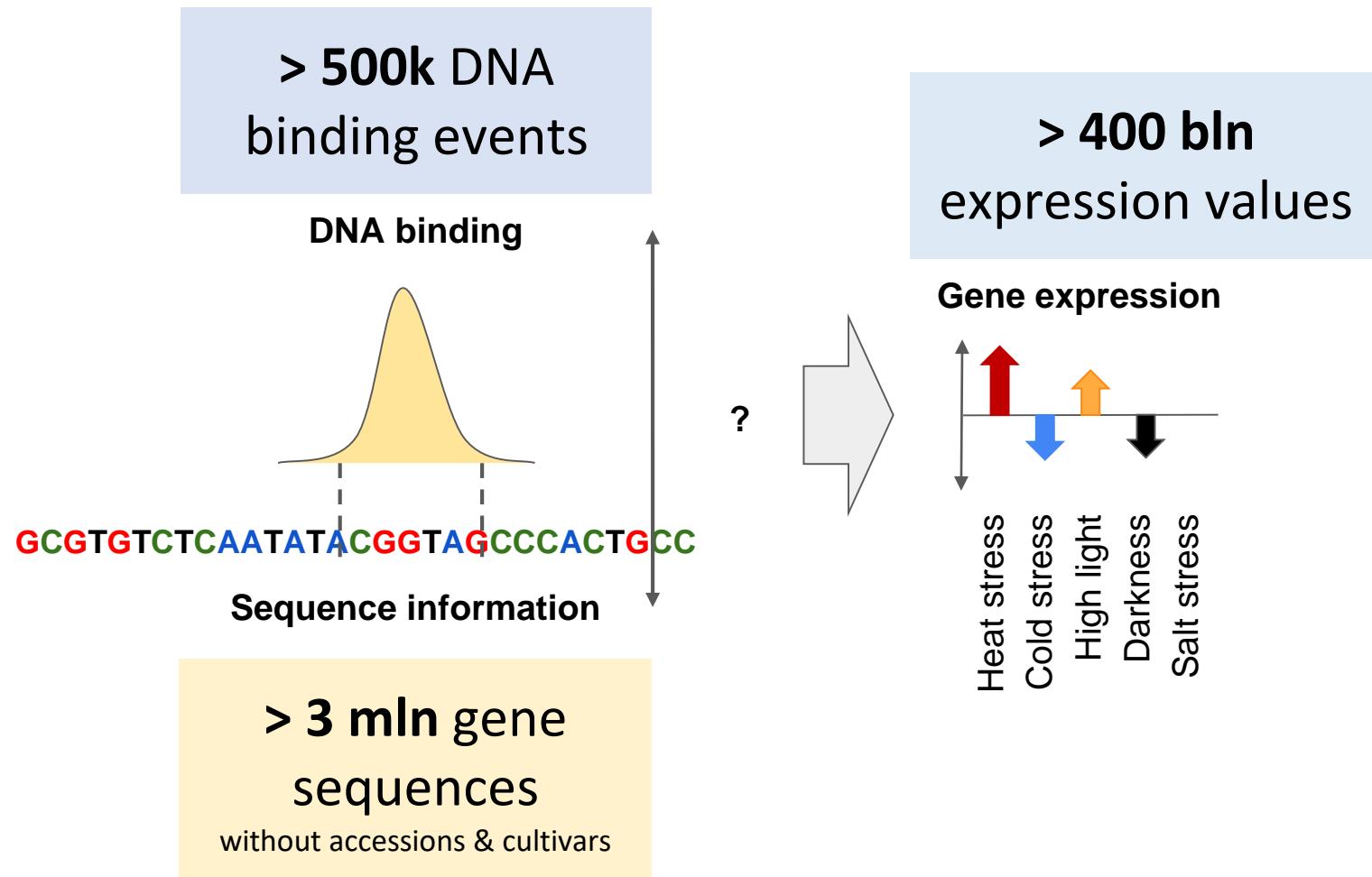
Gernot Schmitz



Fritz Peleke



# The task

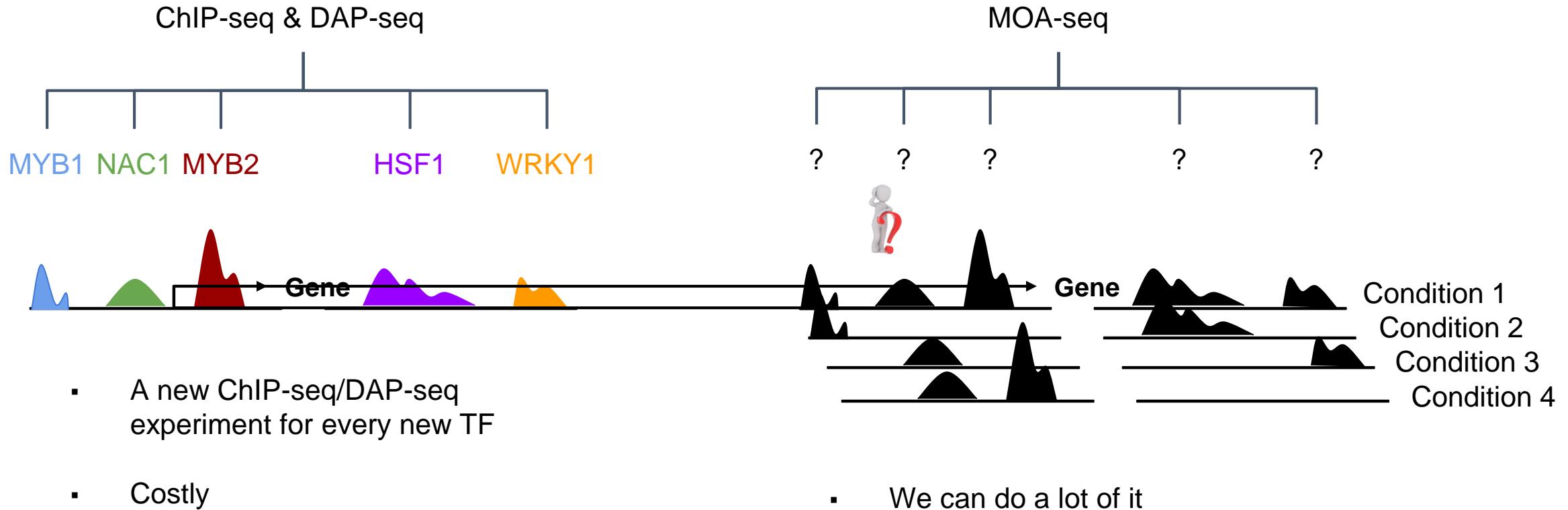


# Beyond gene expression

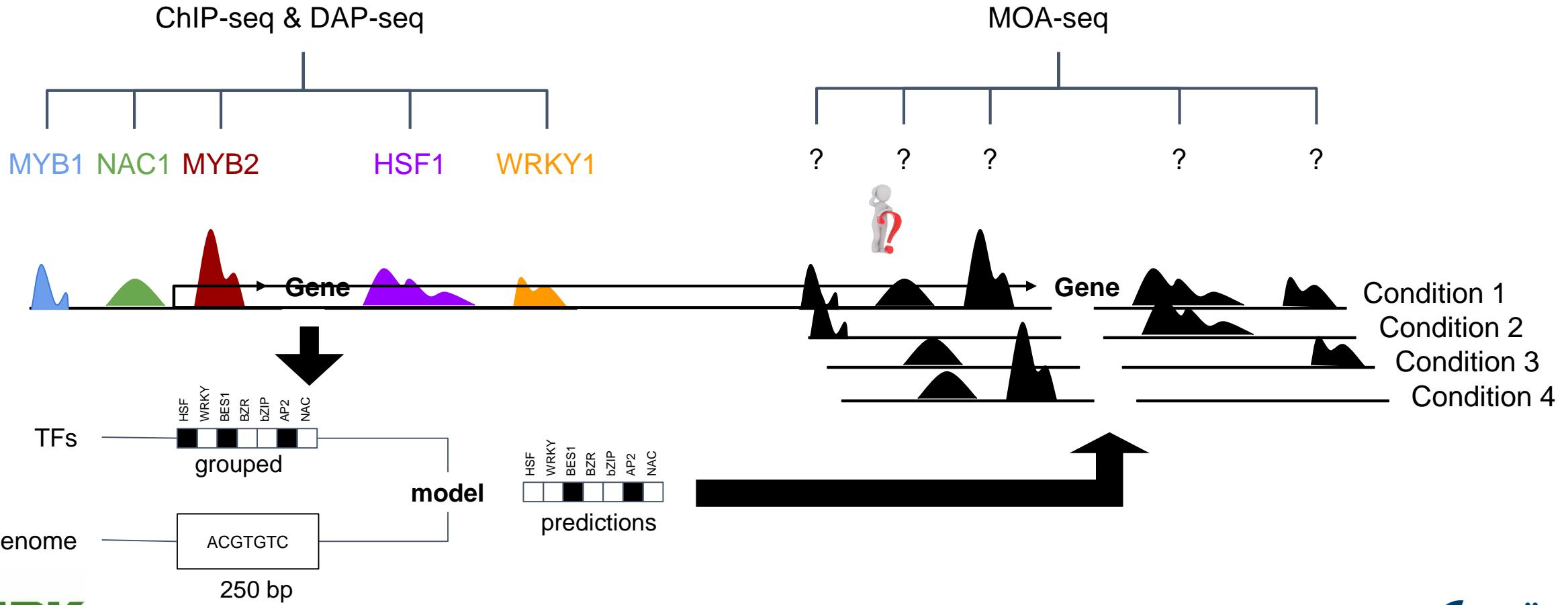


- A new ChIP-seq/DAP-seq experiment for every new TF
- Costly
- An experiment for all TFs
- Not sure what the peaks are

# Beyond gene expression



# Beyond gene expression

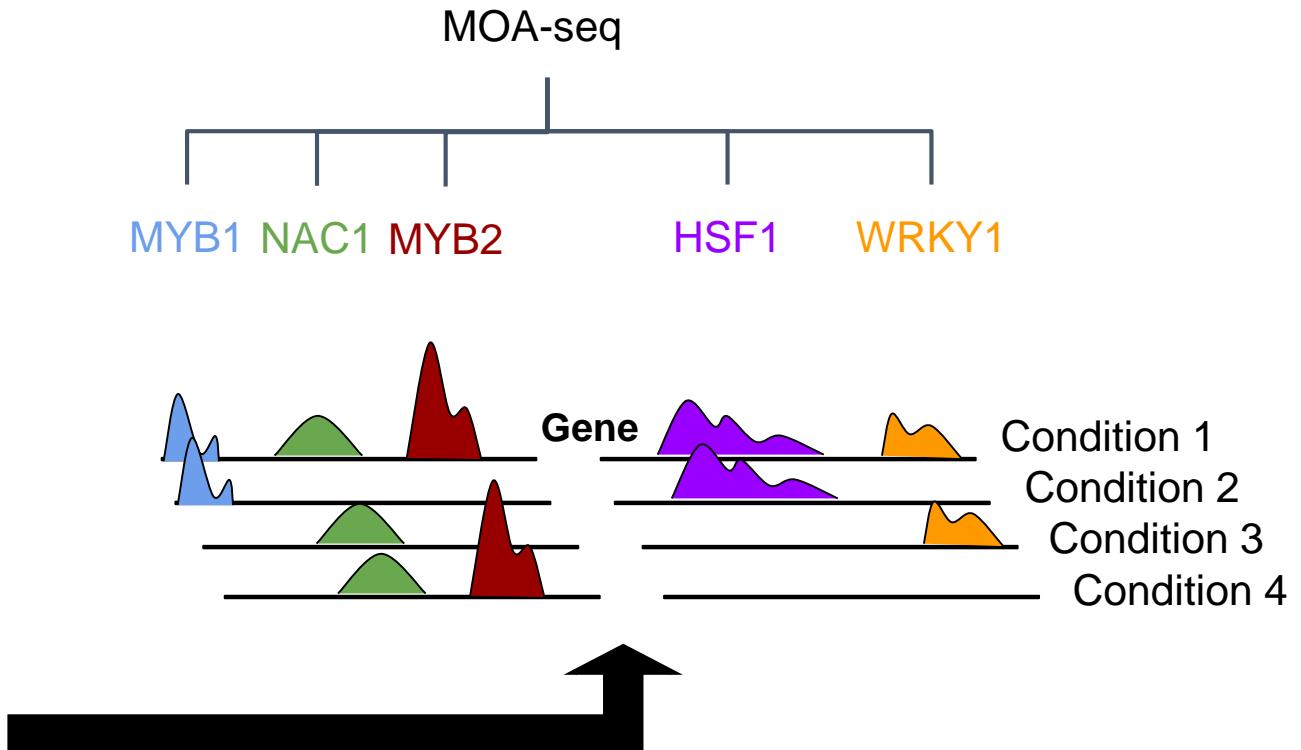
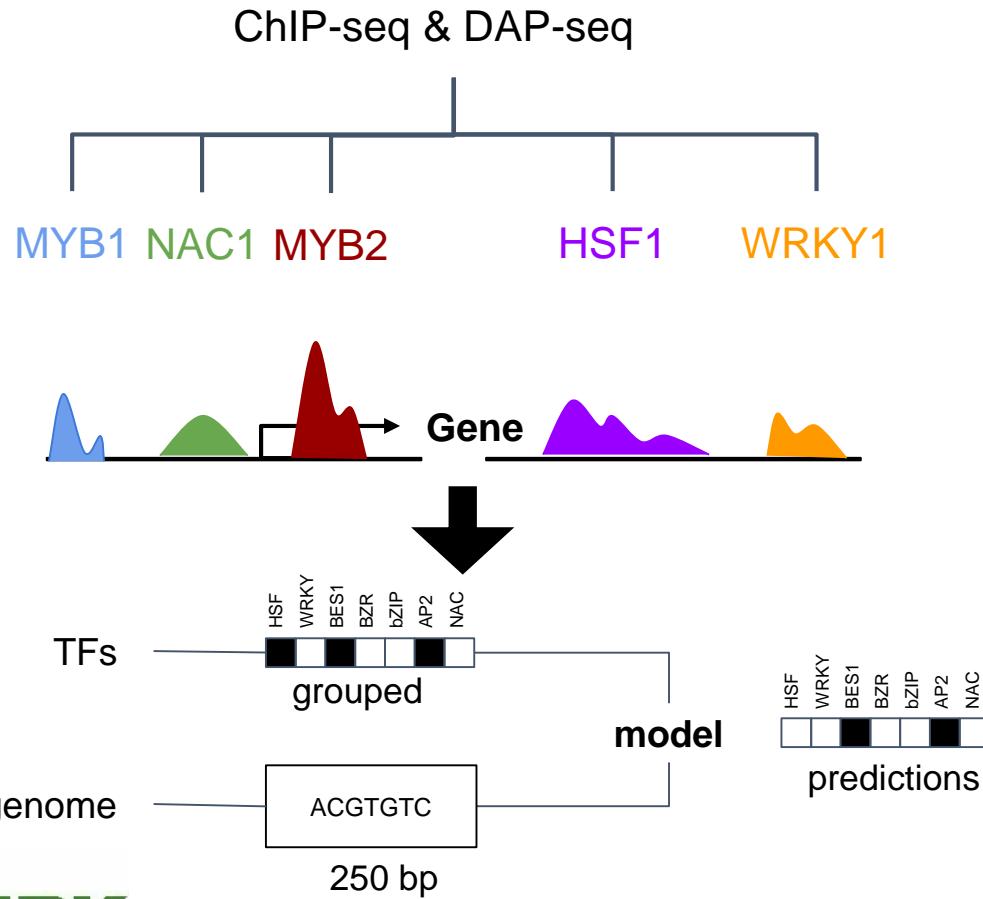


# Beyond gene expression

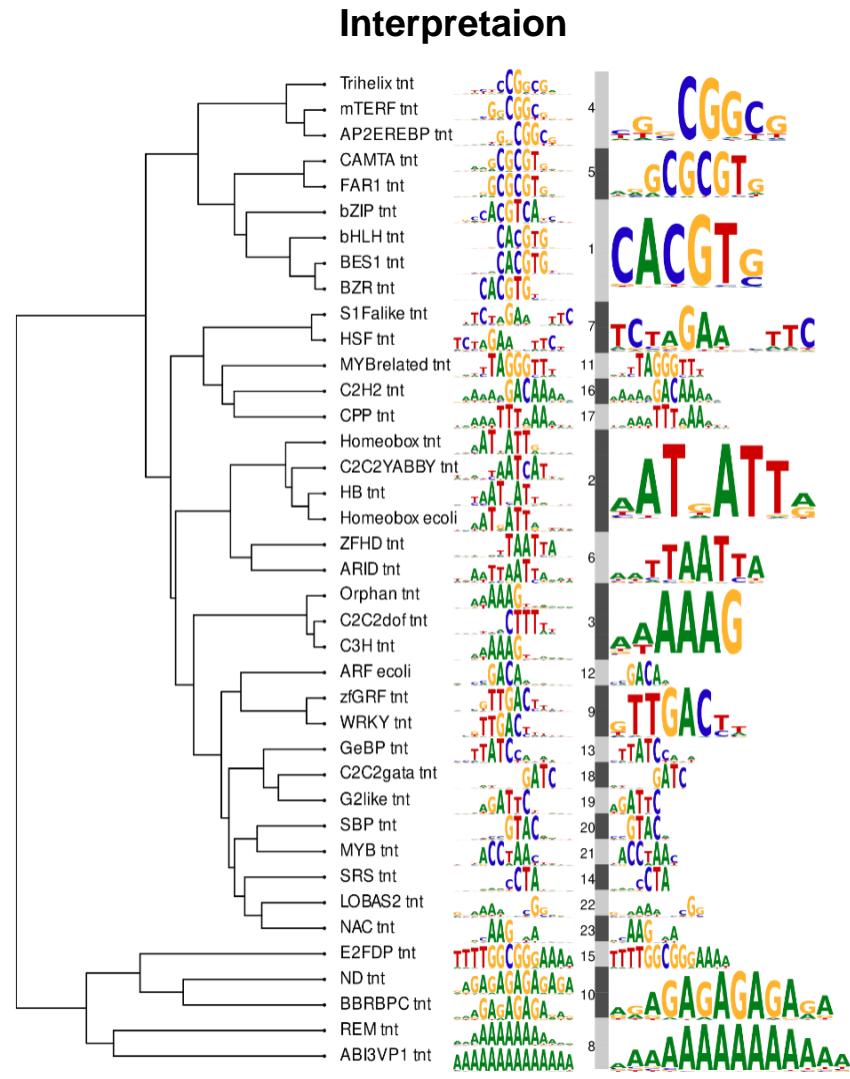
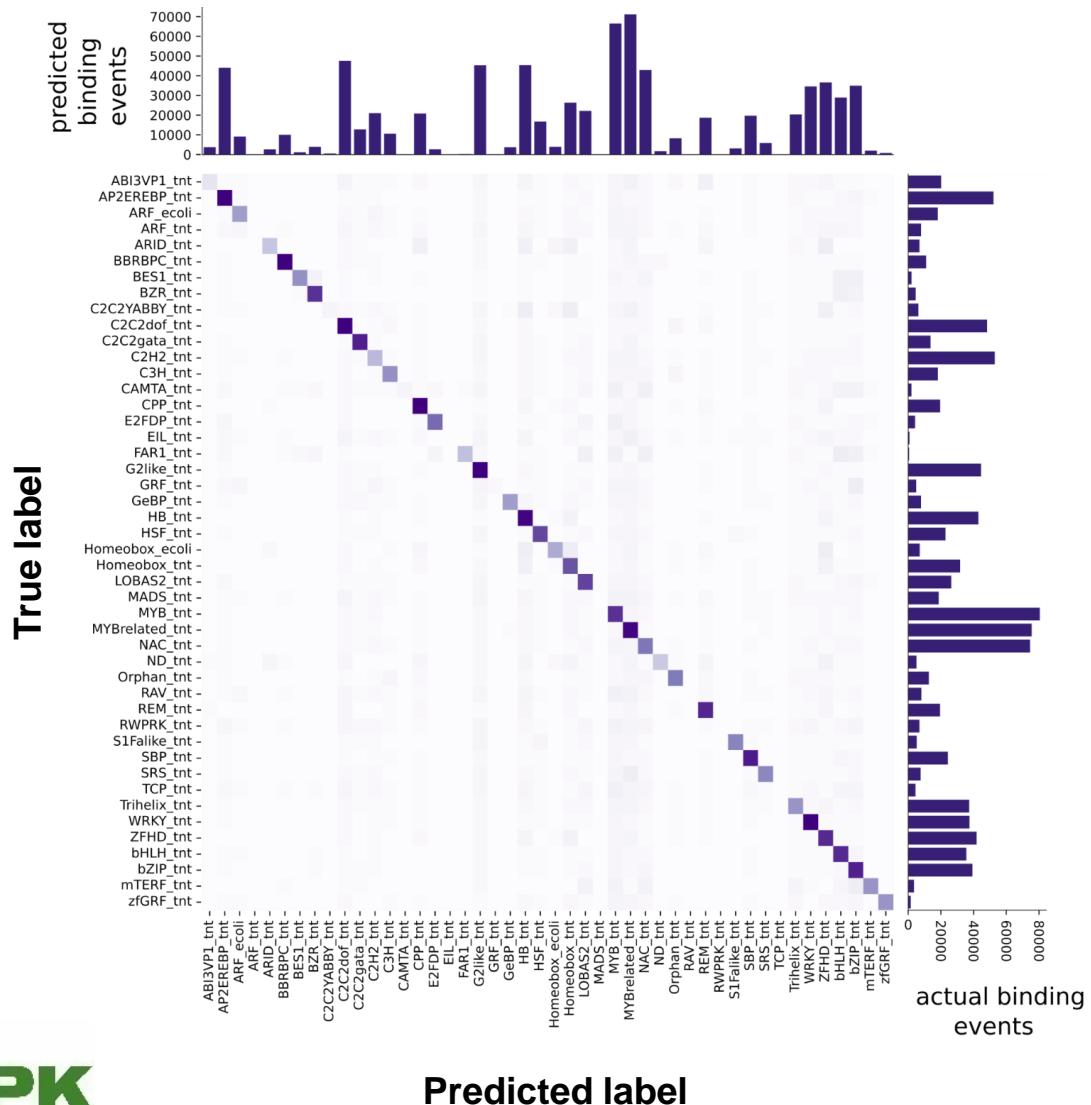


Simon Zumkeller

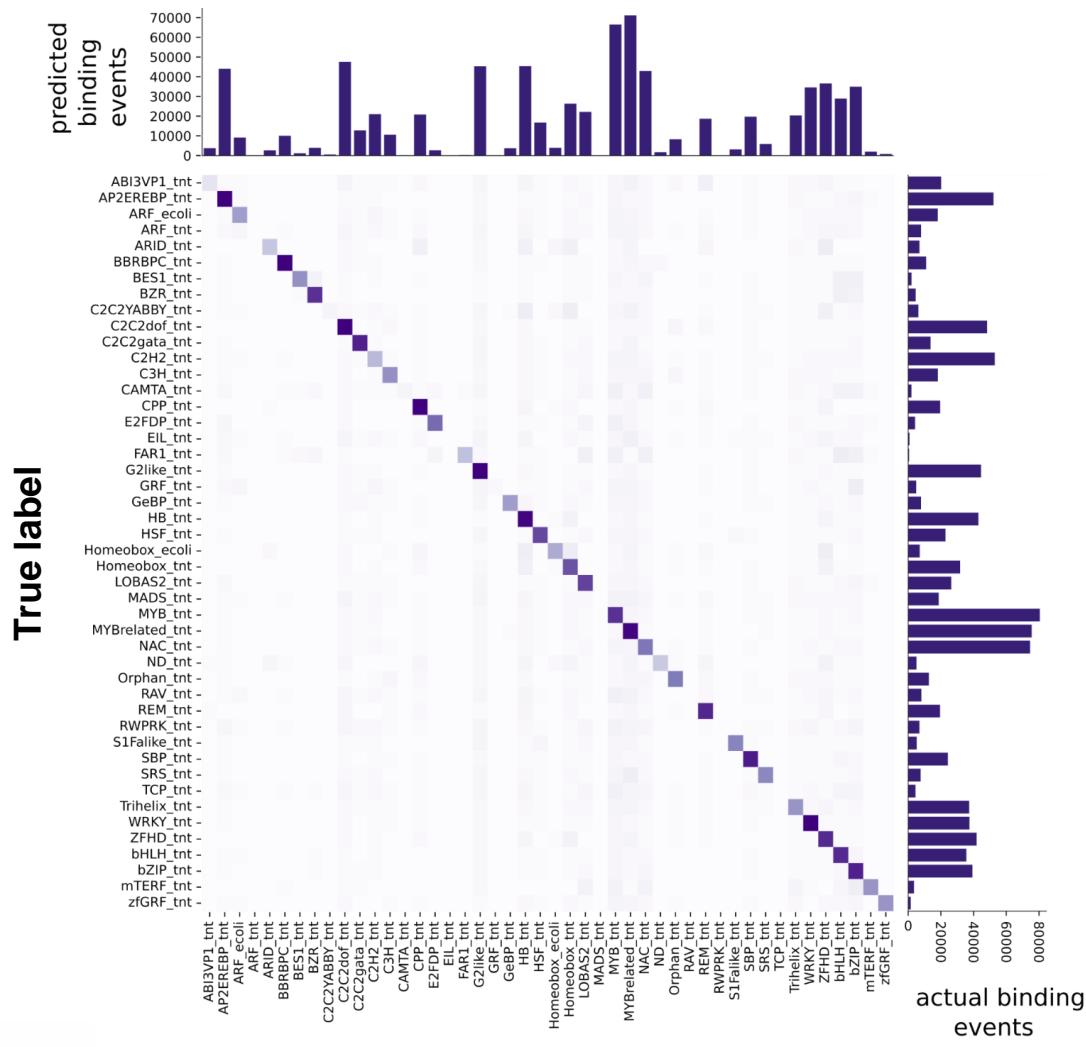
Fritz Peleke



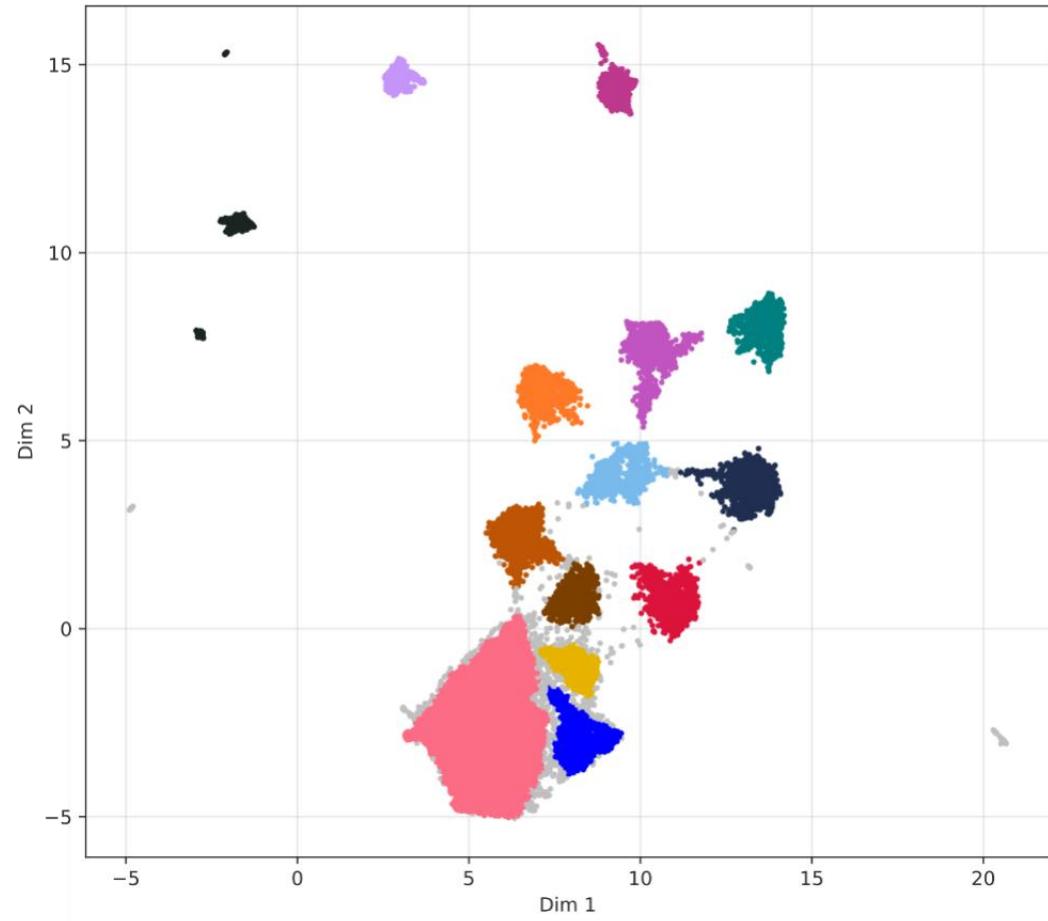
# DNA binding prediction



# DNA binding prediction



Gene regulatory profiles



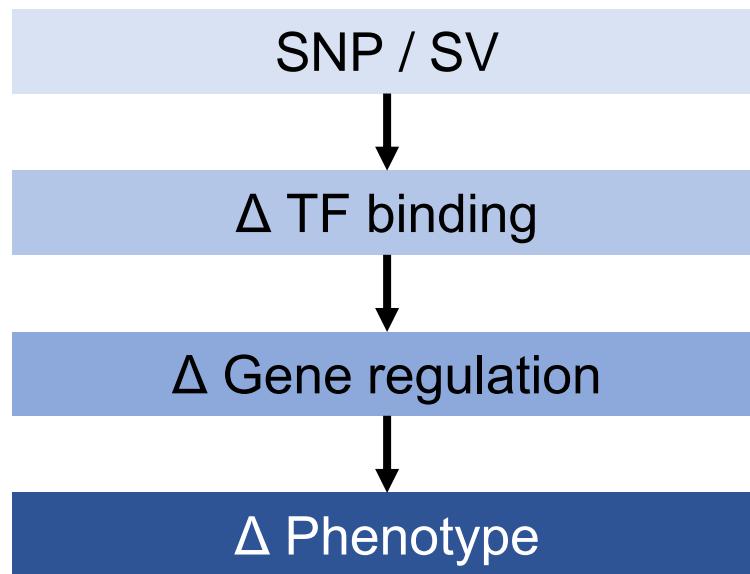
# From genetic variation to phenotypic traits



1001 Genomes

AraGWAS Catalog

36878 phenotype associations



Simon Zumkeller



Gernot Schmitz

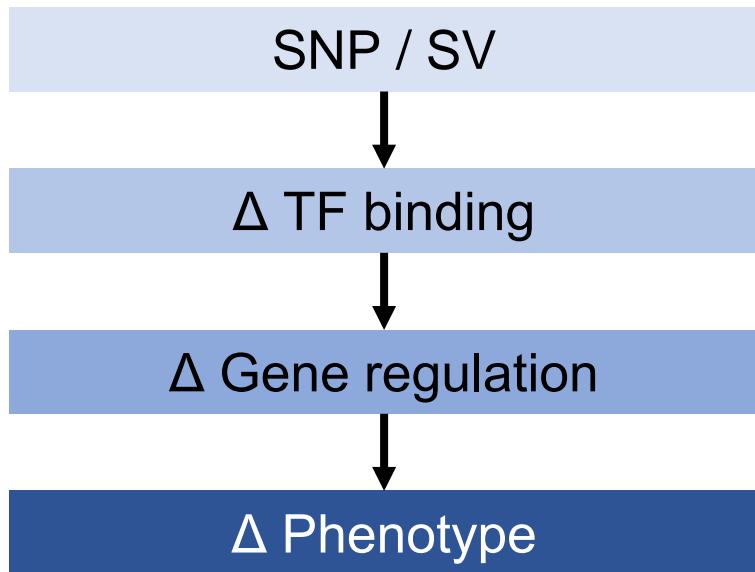
# From genetic variation to phenotypic traits



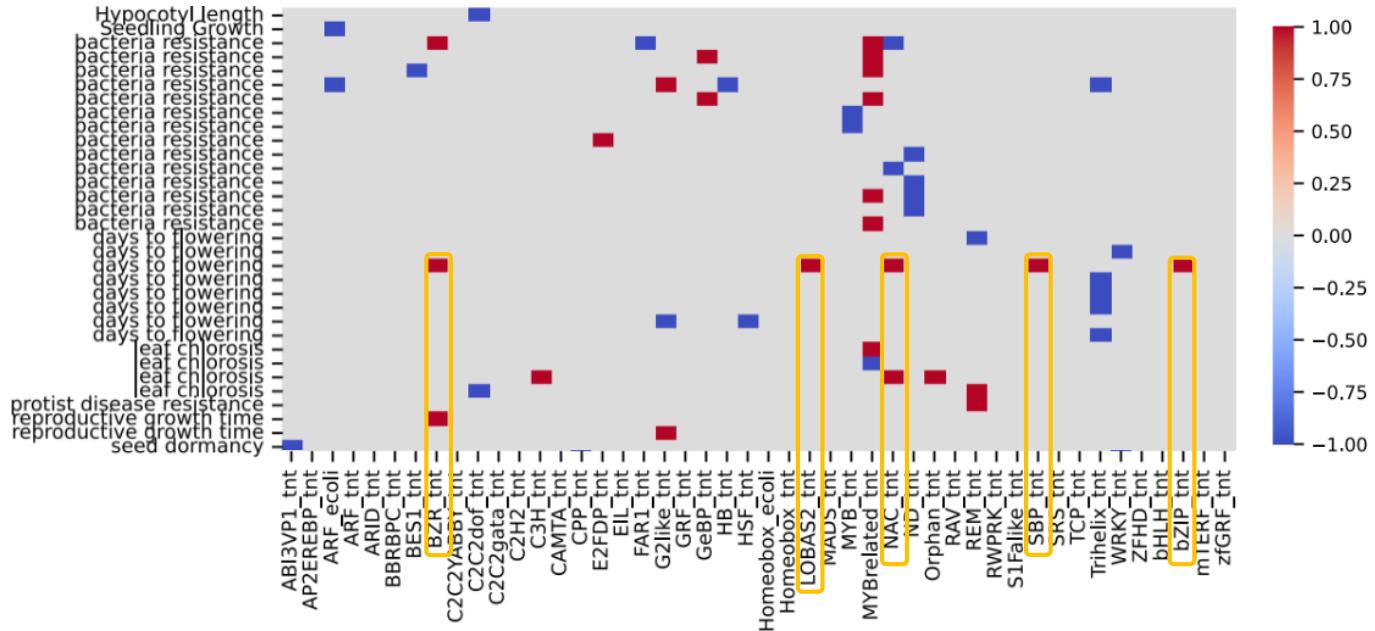
1001 Genomes

AraGWAS Catalog

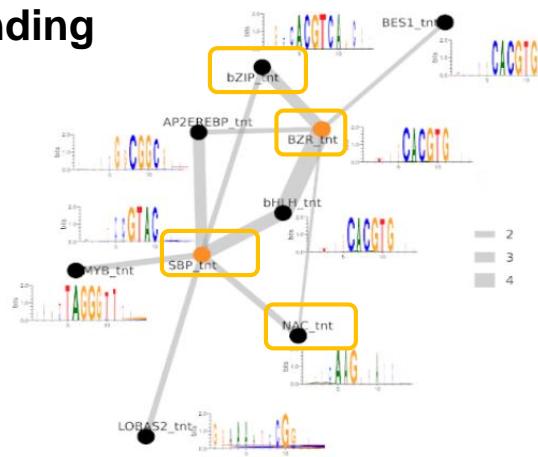
36878 phenotype associations



TF binding – phenotype associations



Combinatorial TF binding



## AI for detective work

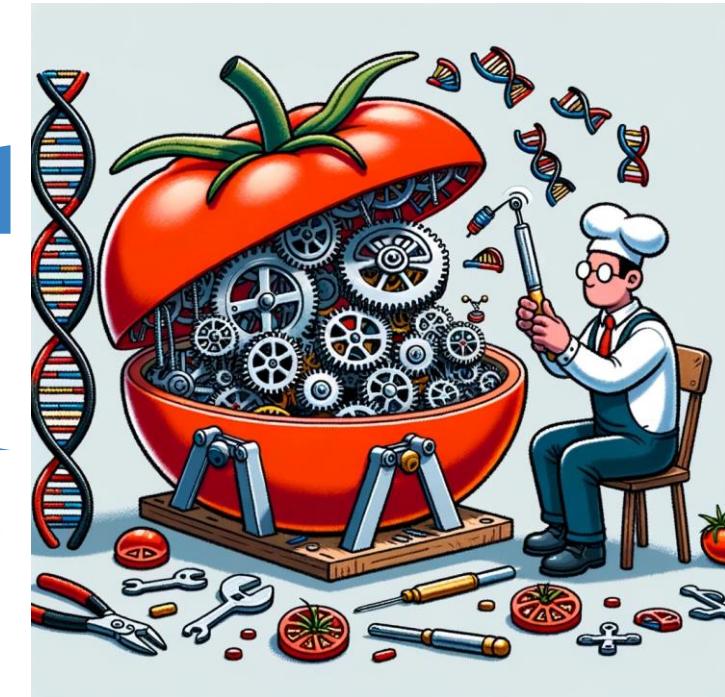


Observations vs variables

System scale



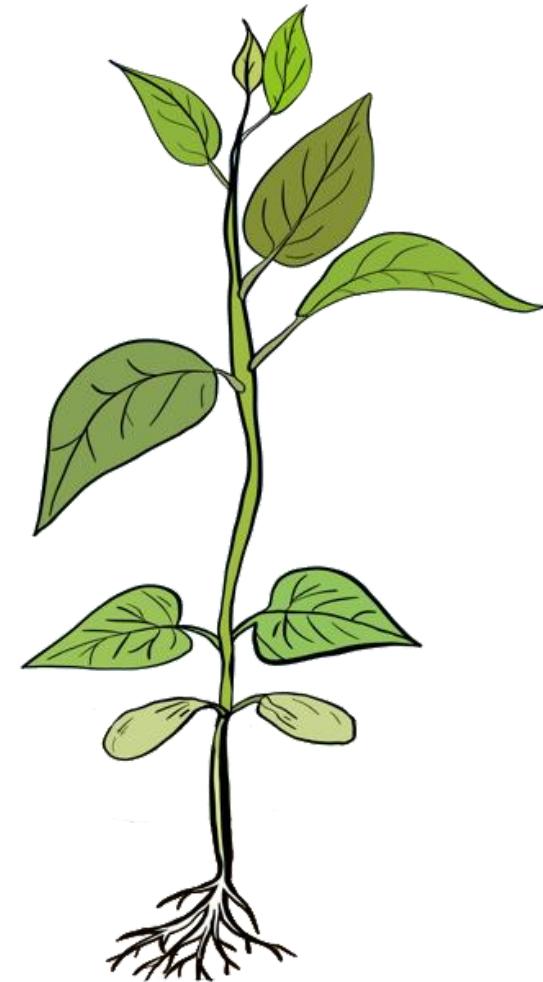
## AI for reverse engineering



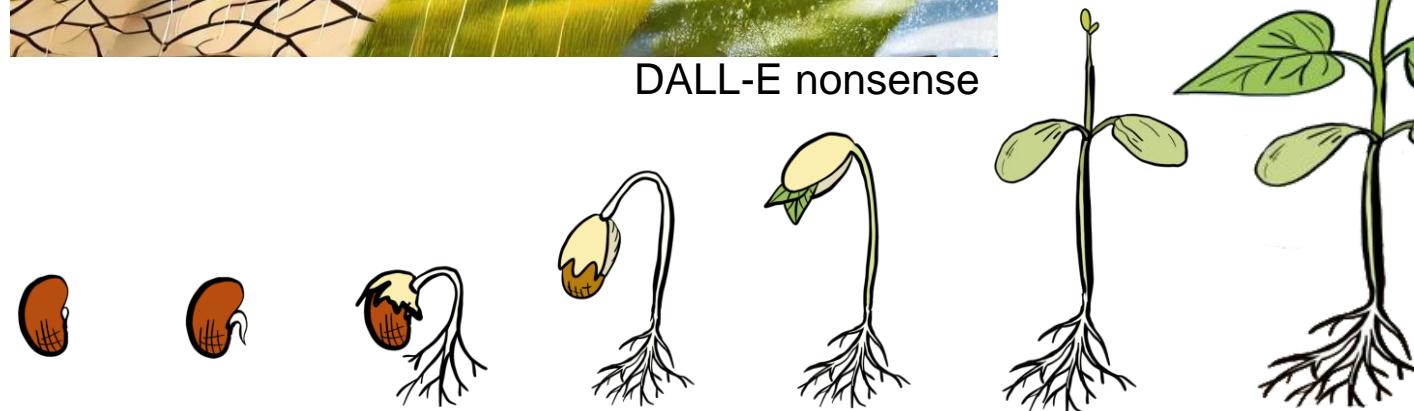
## AI for computer games



# The green game of survival



# The green game of survival



# The green game of survival

Balance resources

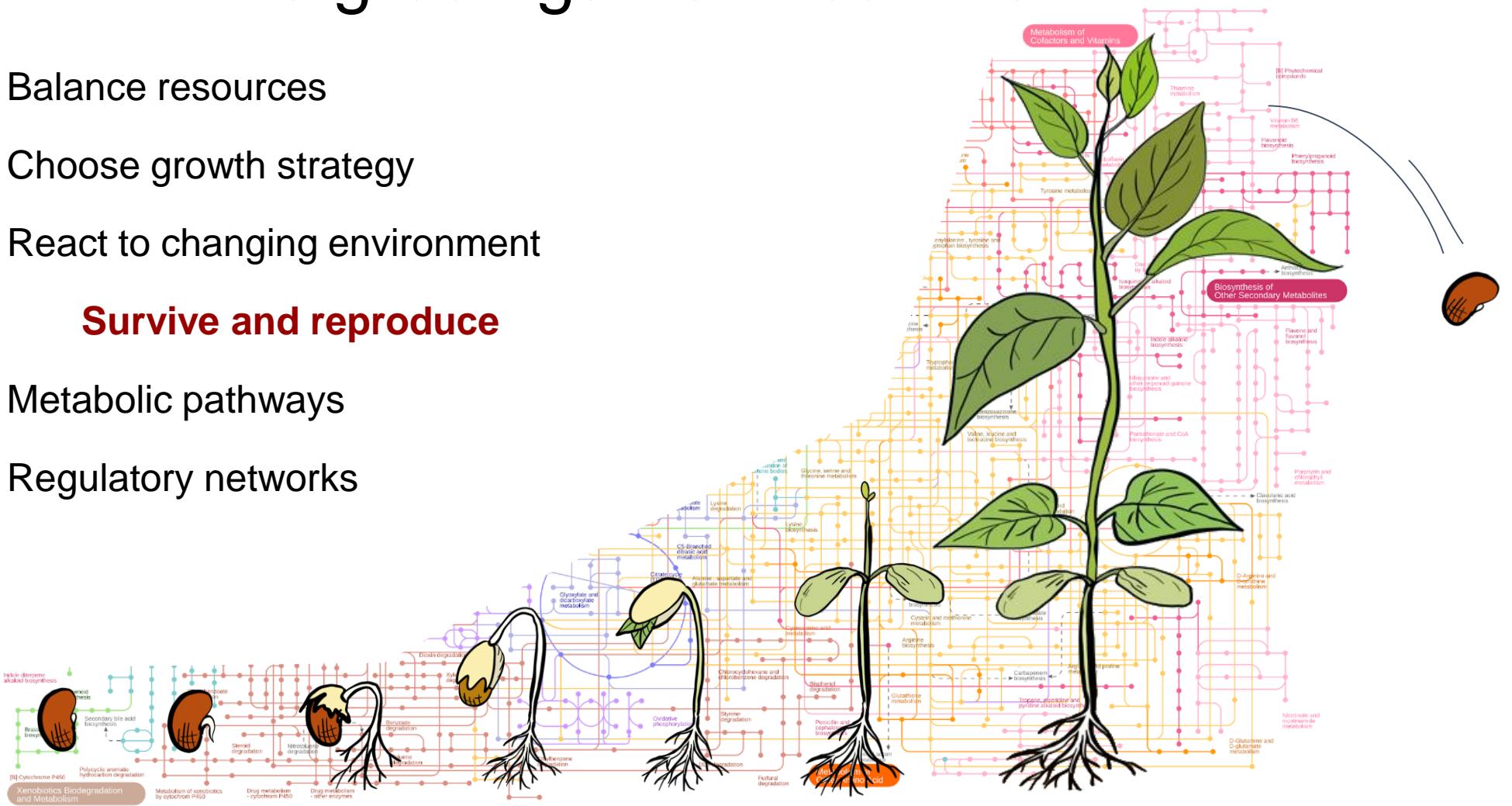
Choose growth strategy

React to changing environment

## Survive and reproduce

Metabolic pathways

Regulatory networks



# The green game of survival

## So what if...

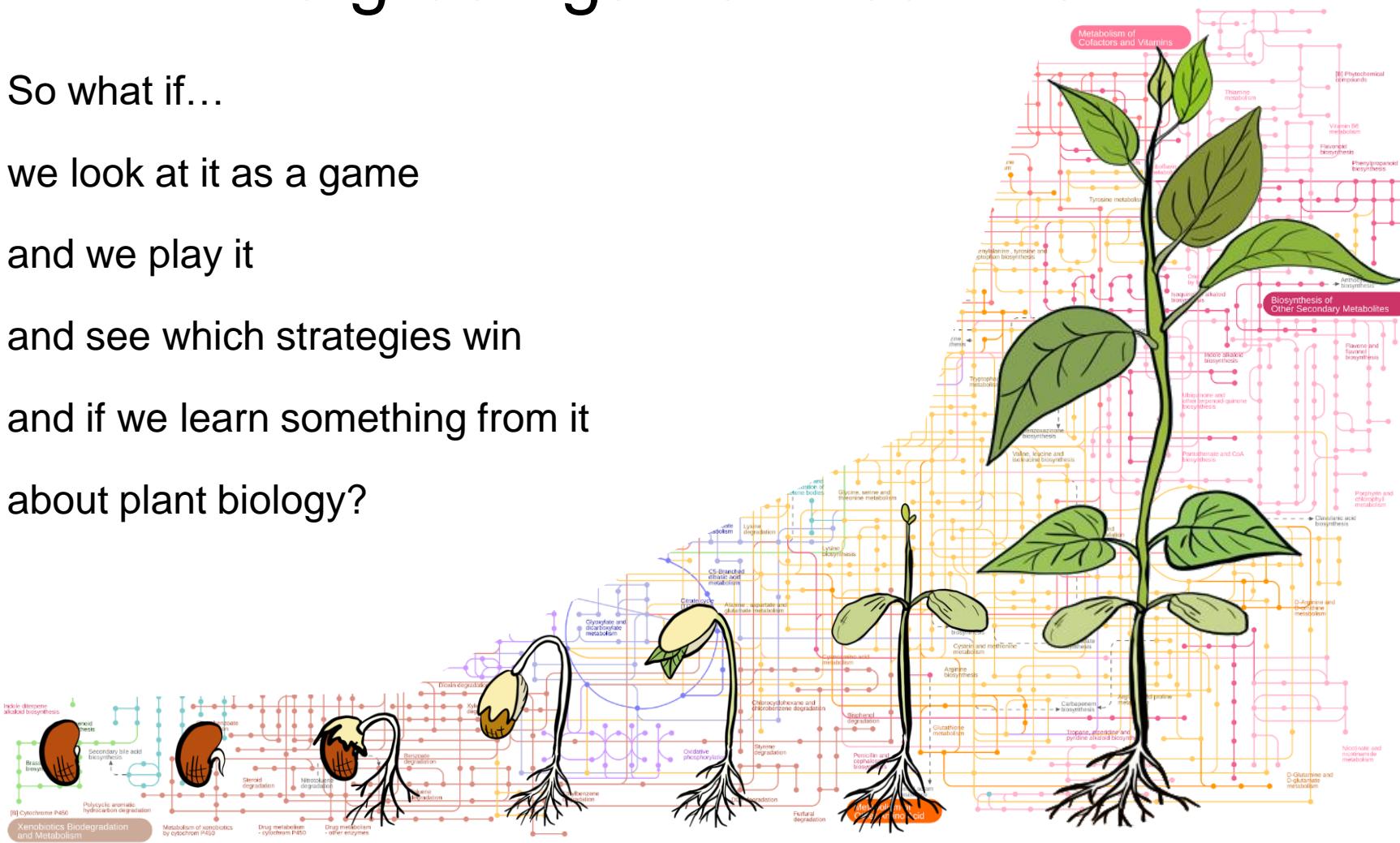
we look at it as a game

and we play it

and see which strategies win

and if we learn something from it

# about plant biology?



# The green game of survival

## So what if...

we look at it as a game

and we play it

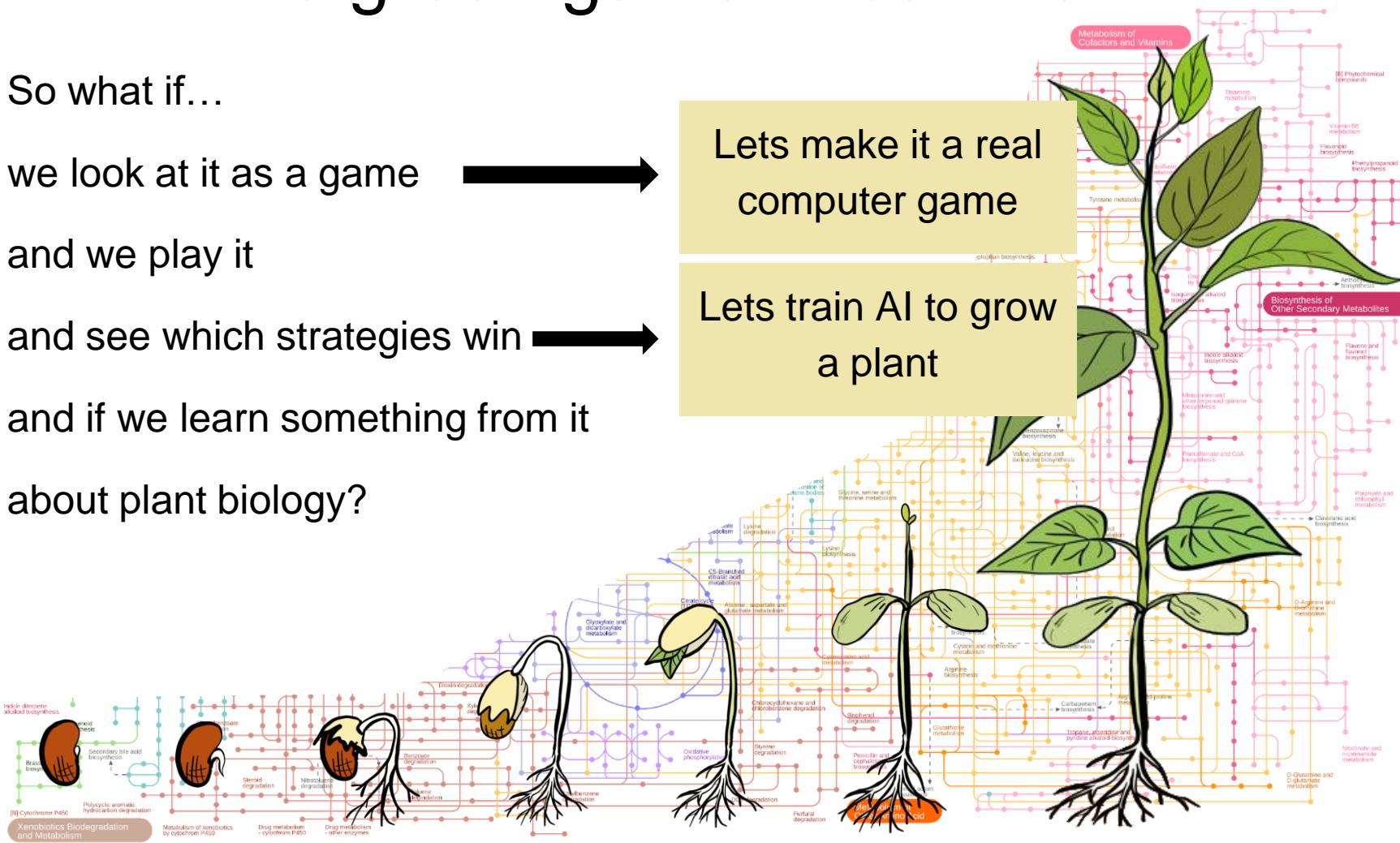
and see which strategies win.

and if we learn something from it

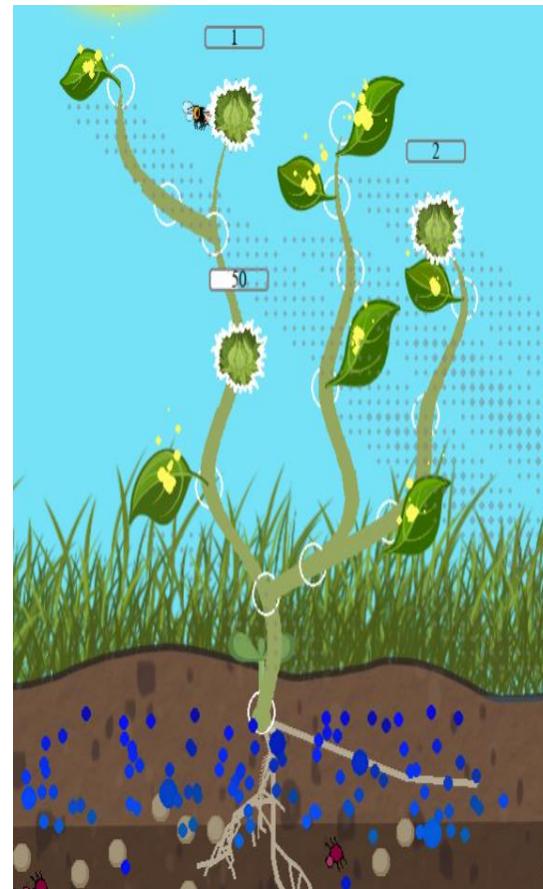
# about plant biology?

# Lets make it a real computer game

# Lets train AI to grow a plant



# A digital plant in a digital environment



- Light intensity
- Light spectrum
- Temperature
- Humidity
- Water availability
- Nutrients
- .
- .
- $\text{CO}_2$
- $\text{O}_2$
- Time of the day
- Day/night cycle
- Season
- etc.
- .
- .
- .

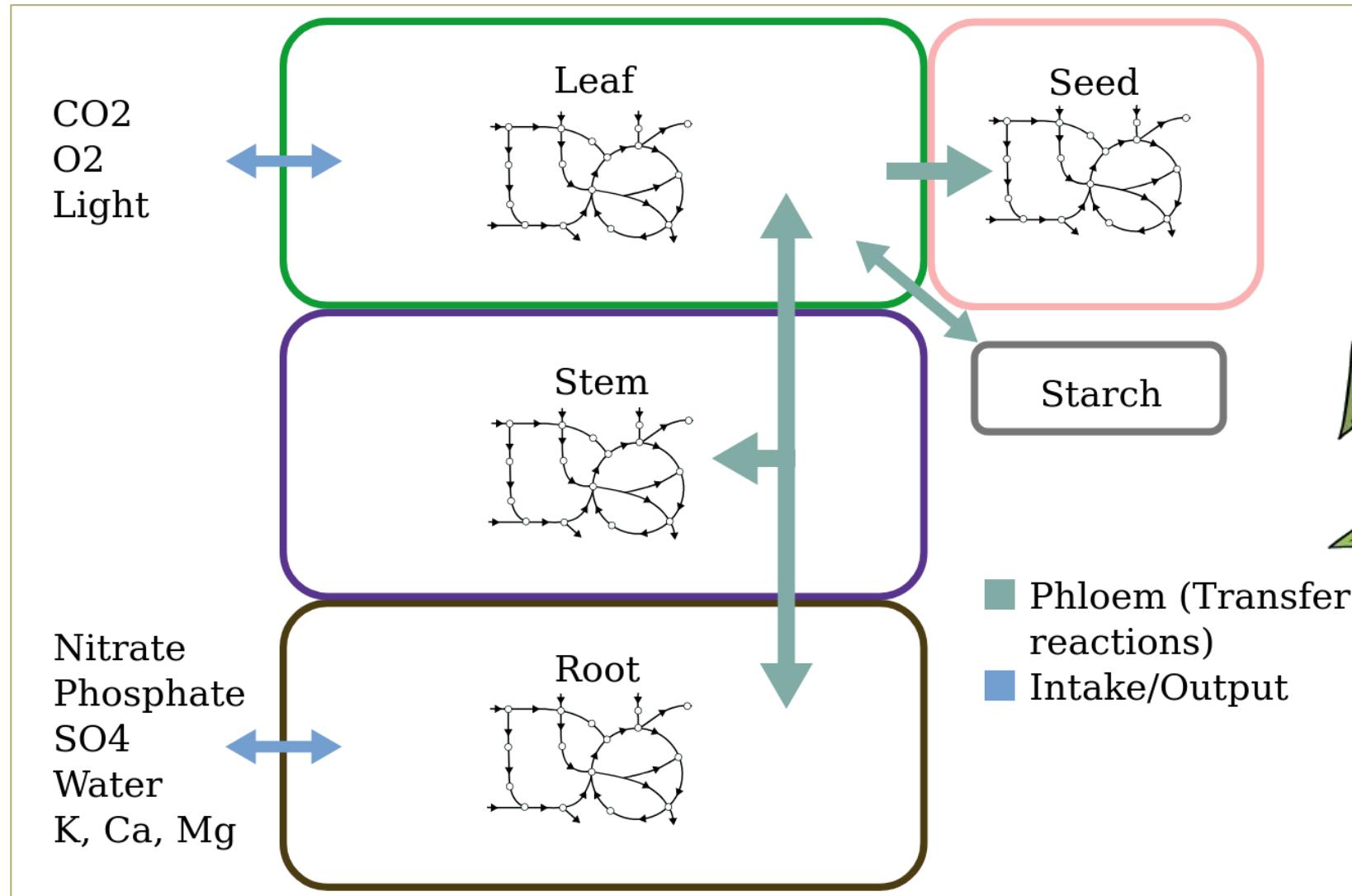
# Whole-plant metabolic model



Stefano La Cruz  
Uni Köln



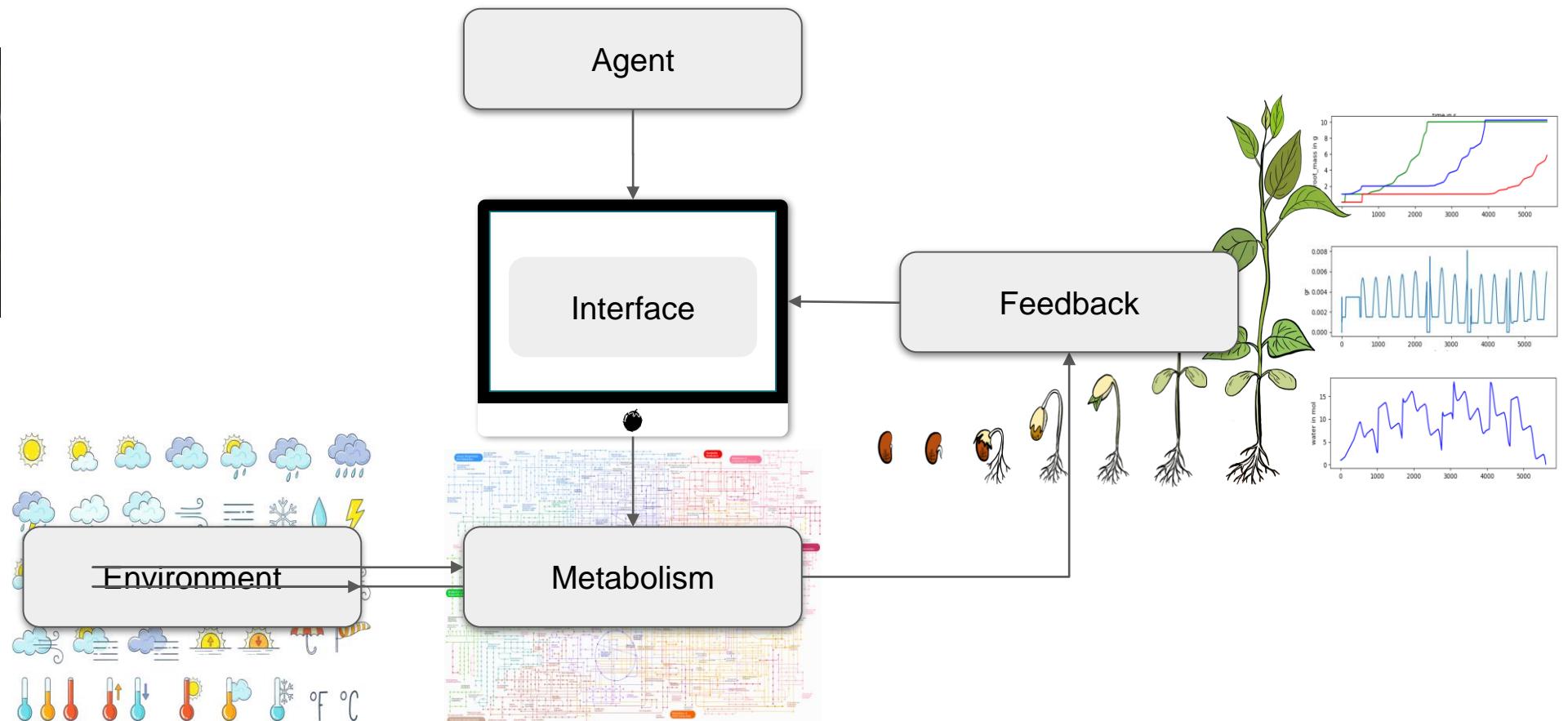
Nadine Töpfer  
Uni Köln



# The game loop



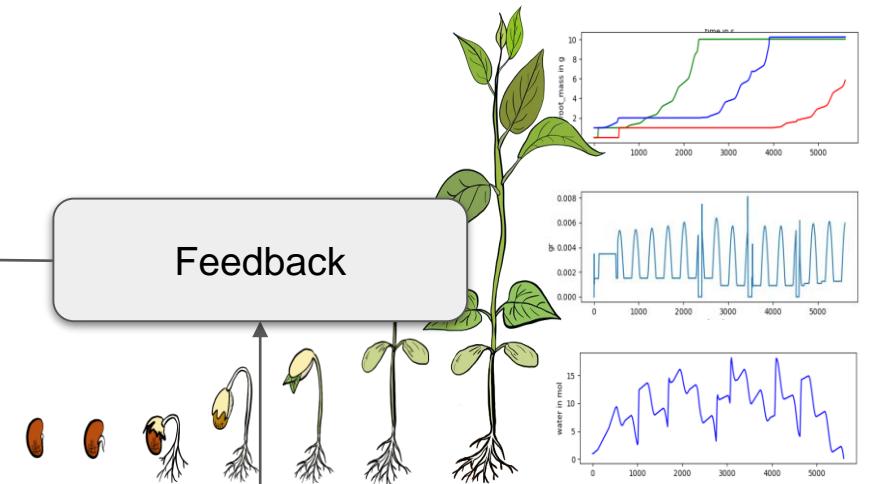
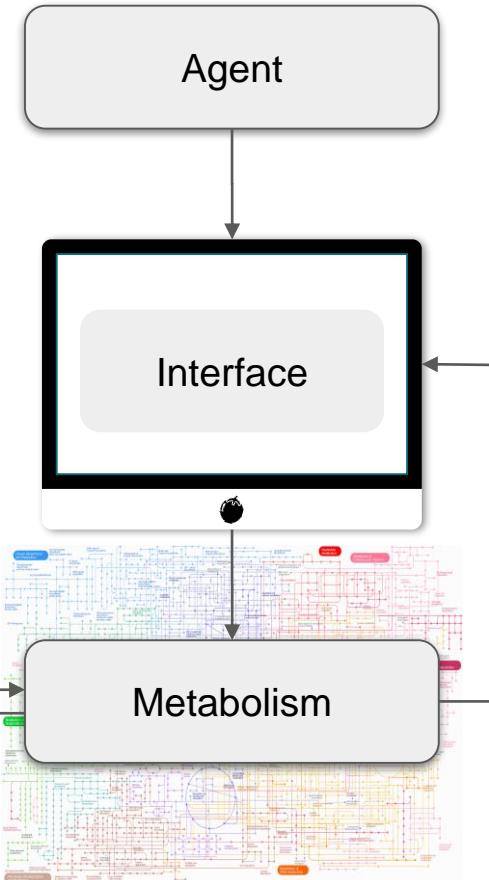
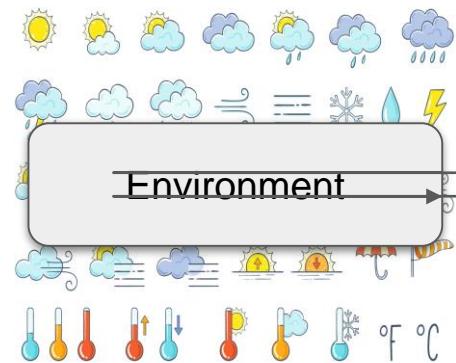
Daniel Koch



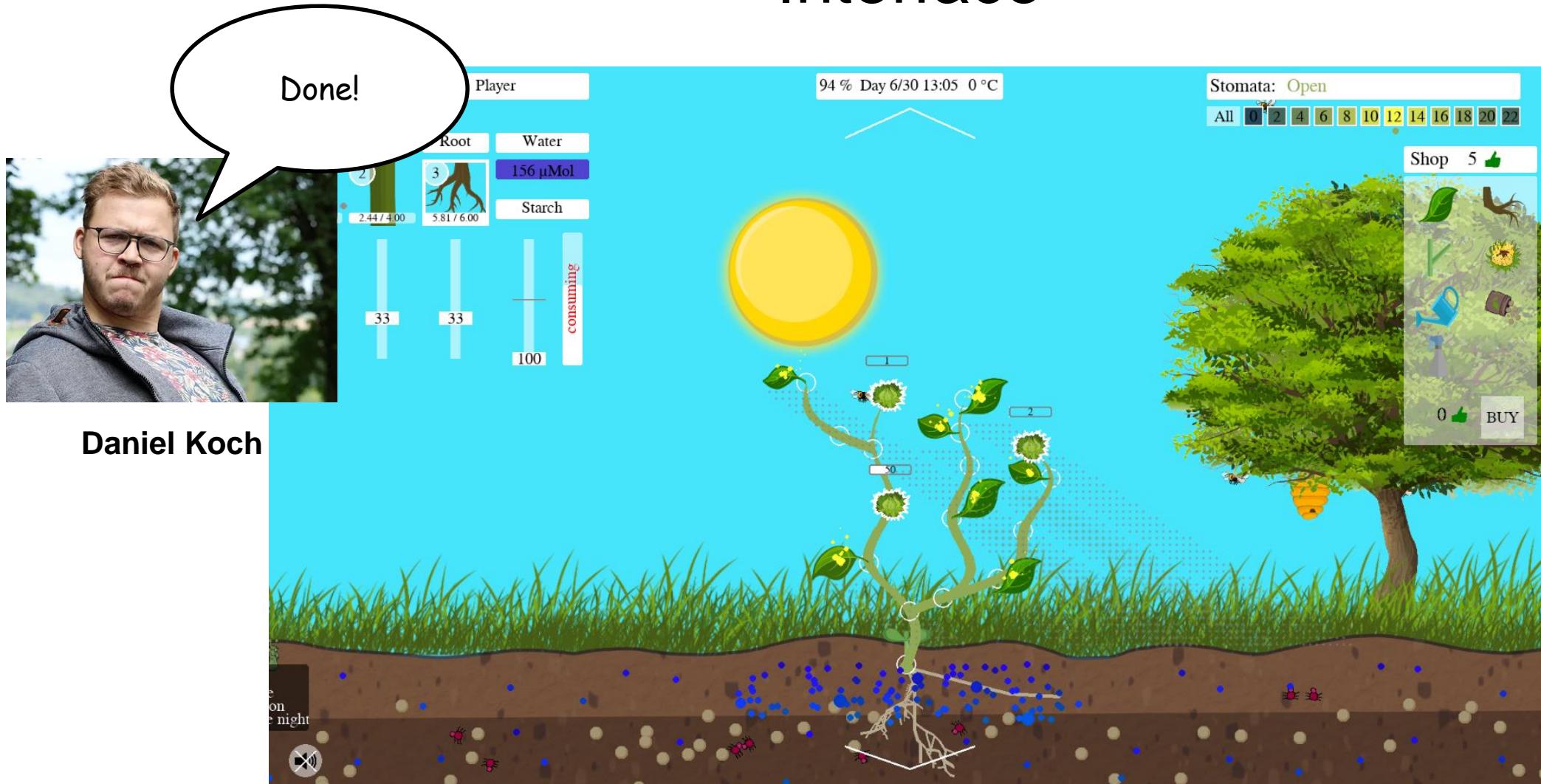
# The game loop



Daniel Koch



# Interface



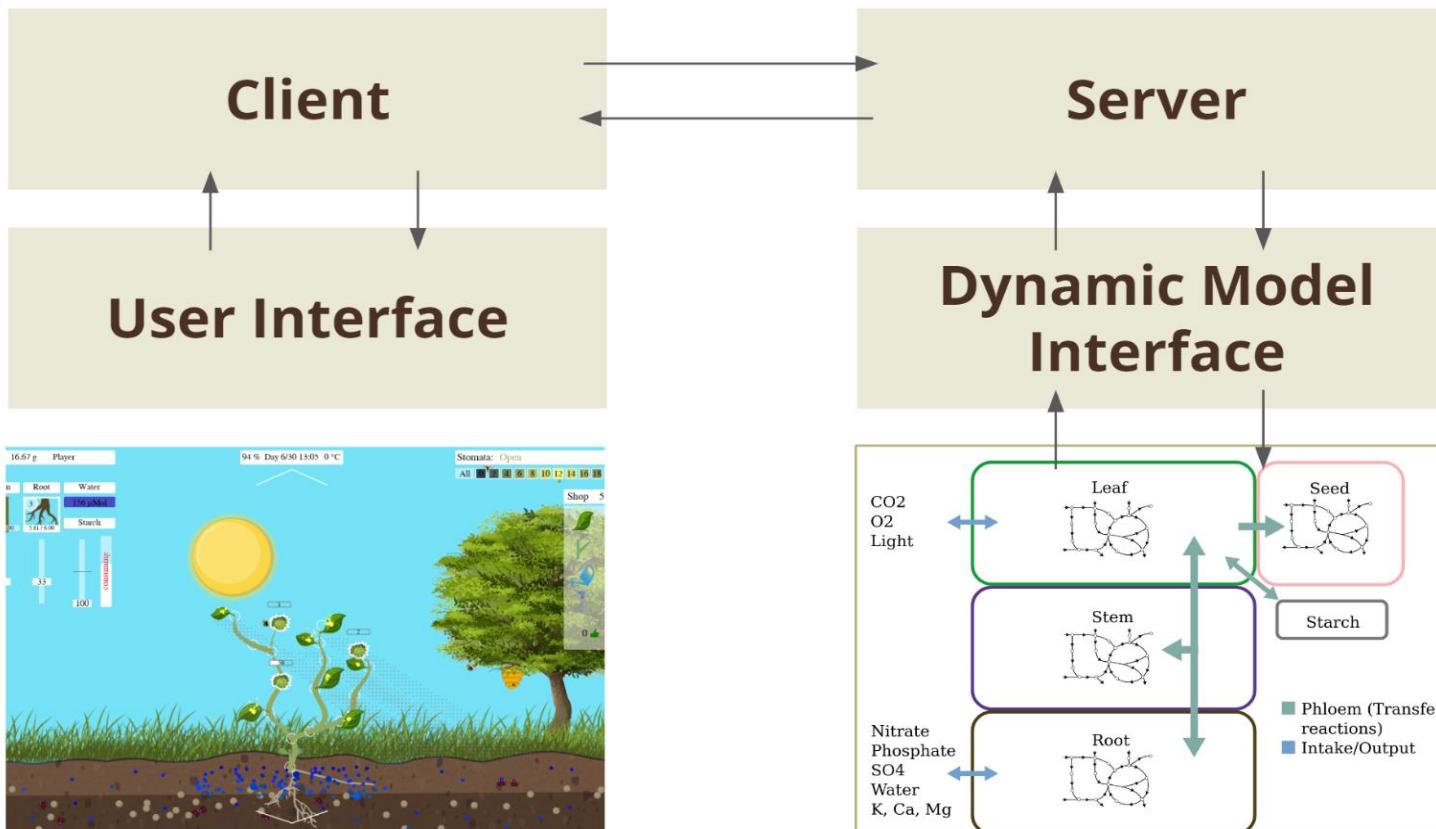
Daniel Koch

<https://danielkoch.itch.io/planted>

# Implementation



Daniel Koch

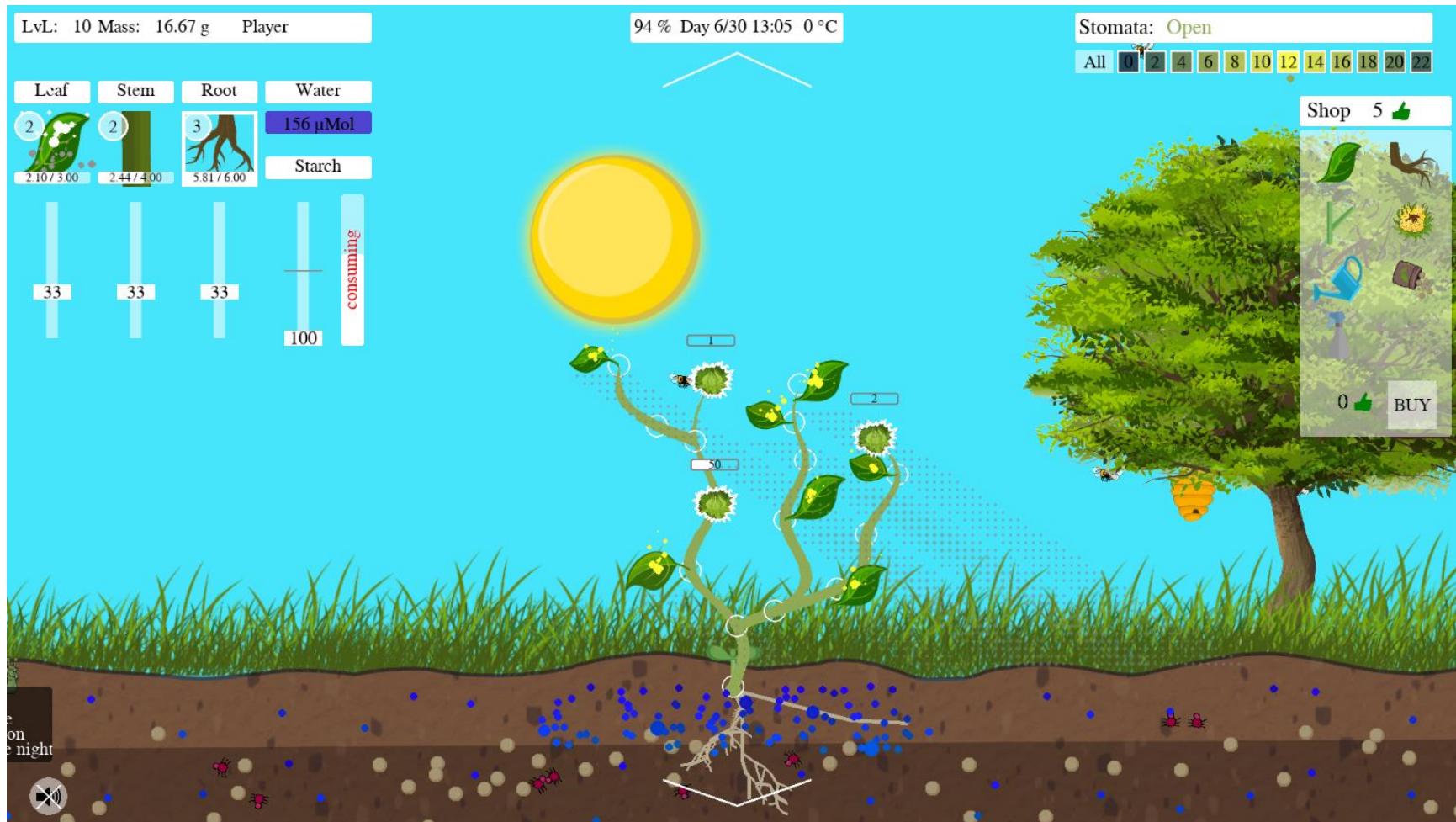


Jan-Niklas Weder  
Uni Köln



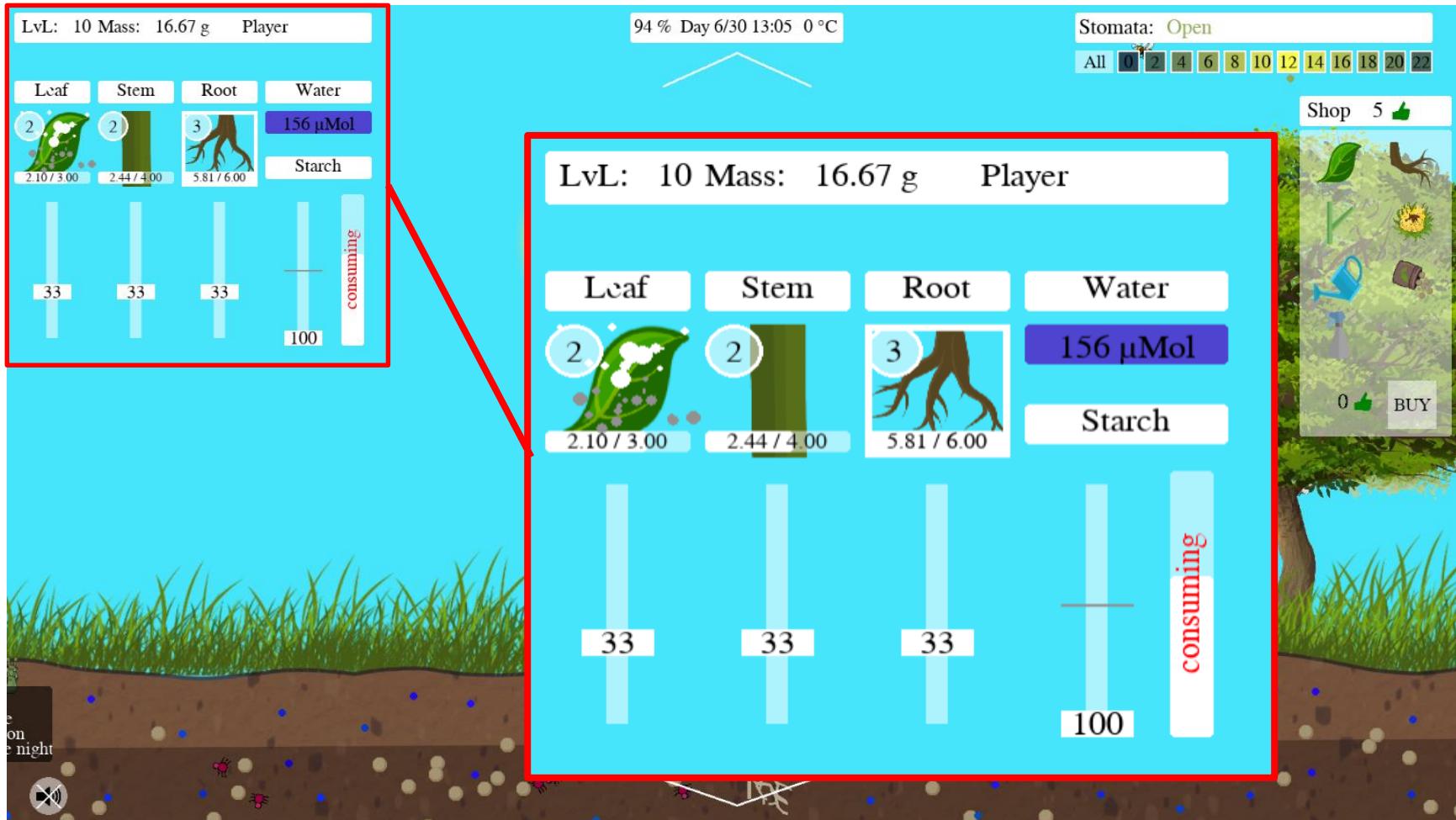
Dennis Psaroudakis  
IPK Gatersleben

# Interface



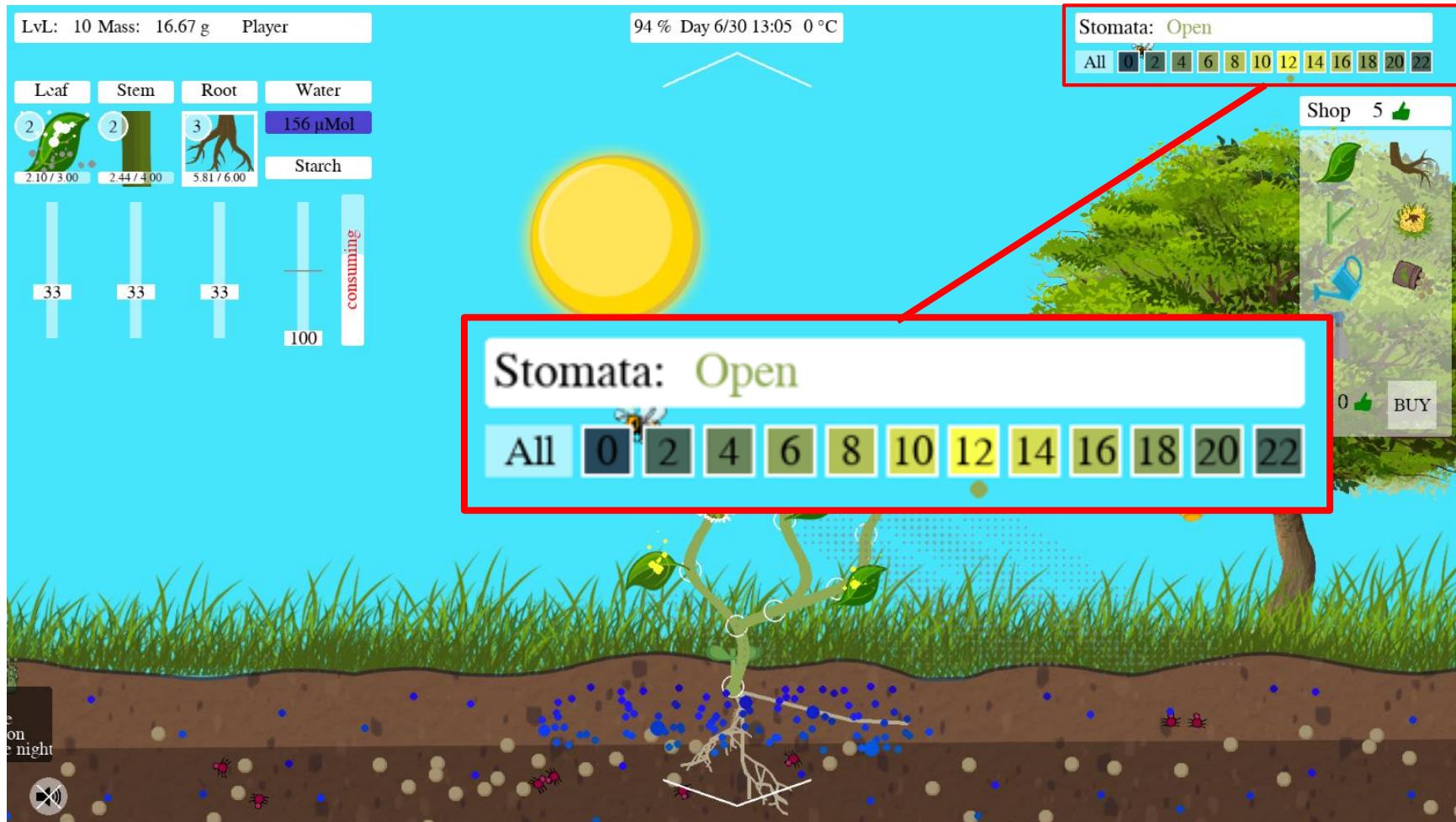
<https://danielkoch.itch.io/planted>

# Interface



<https://danielkoch.itch.io/planted>

# Interface

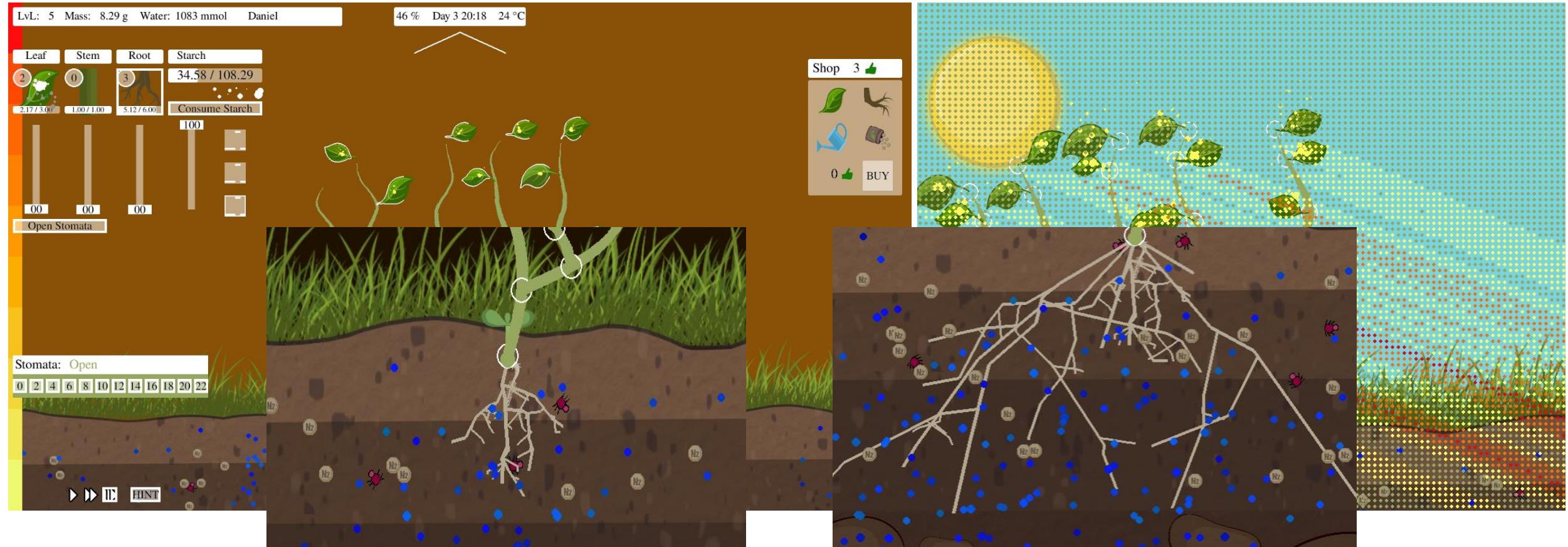


<https://danielkoch.itch.io/planted>

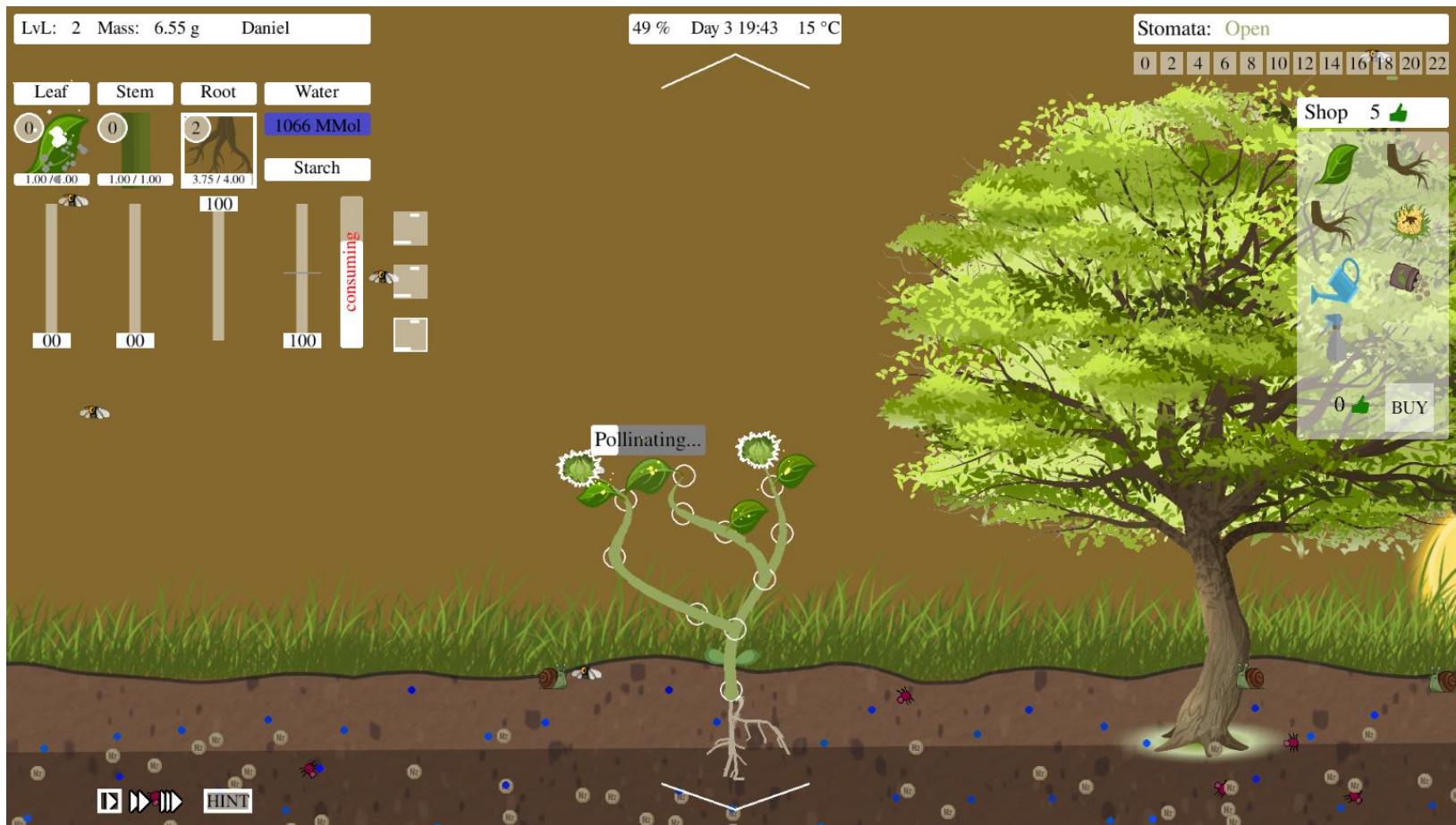
# Interface



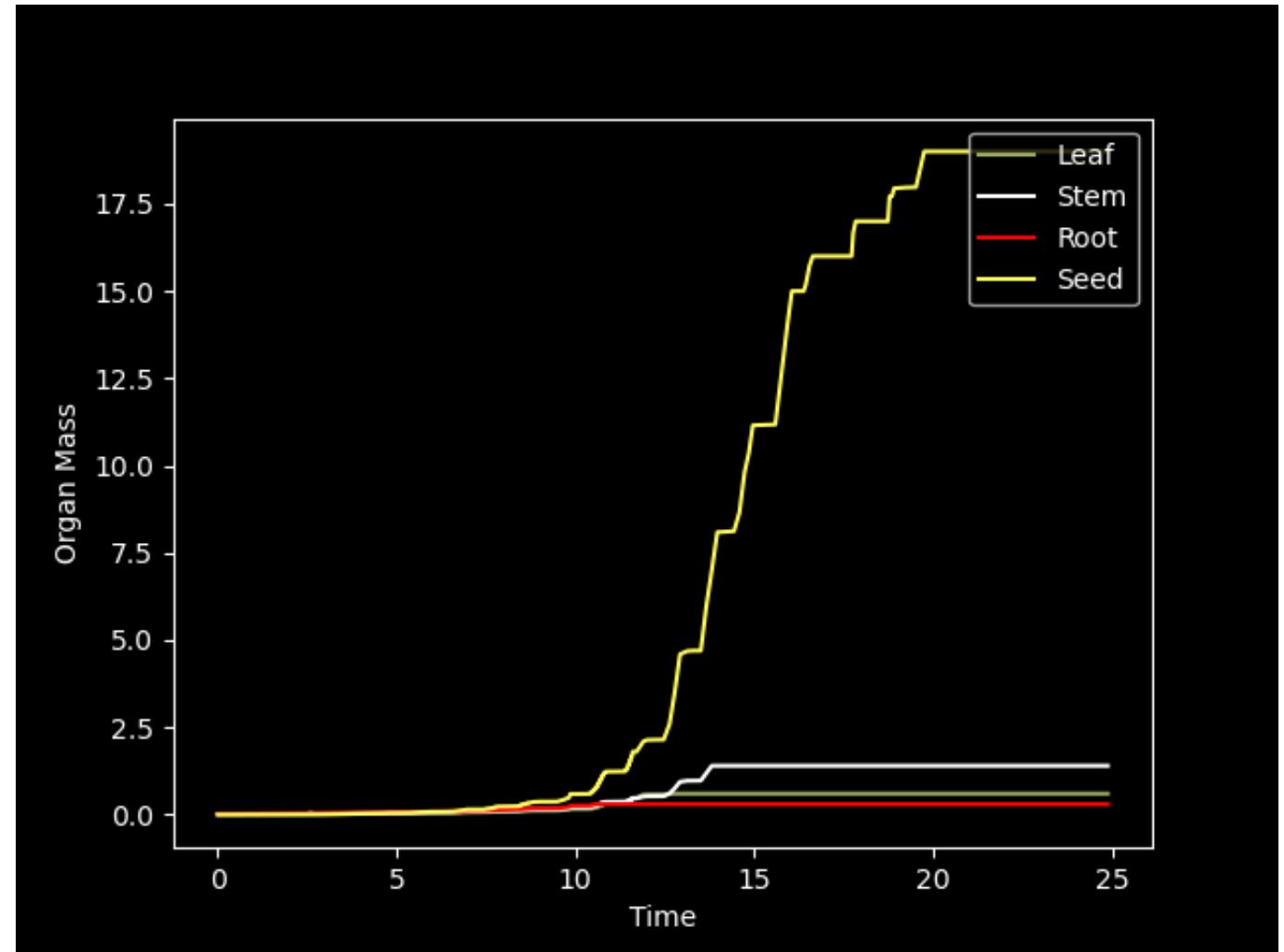
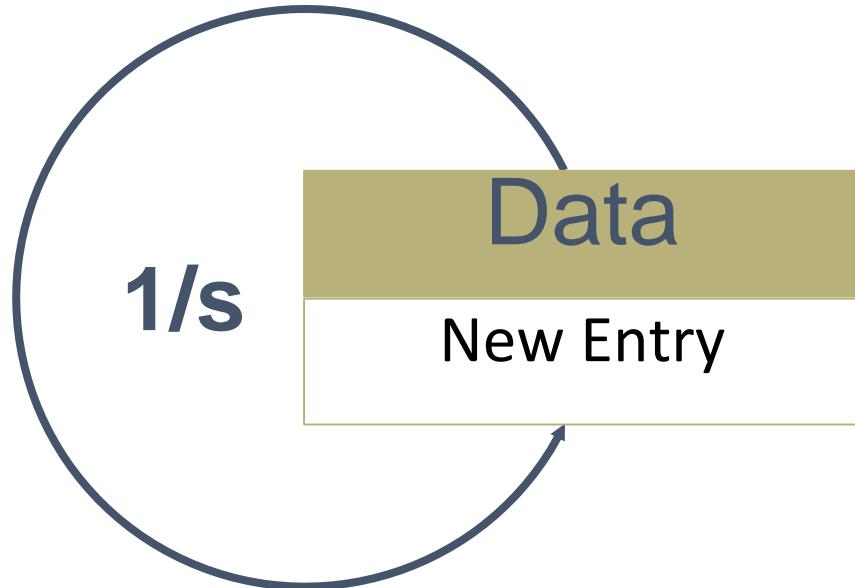
# Interface



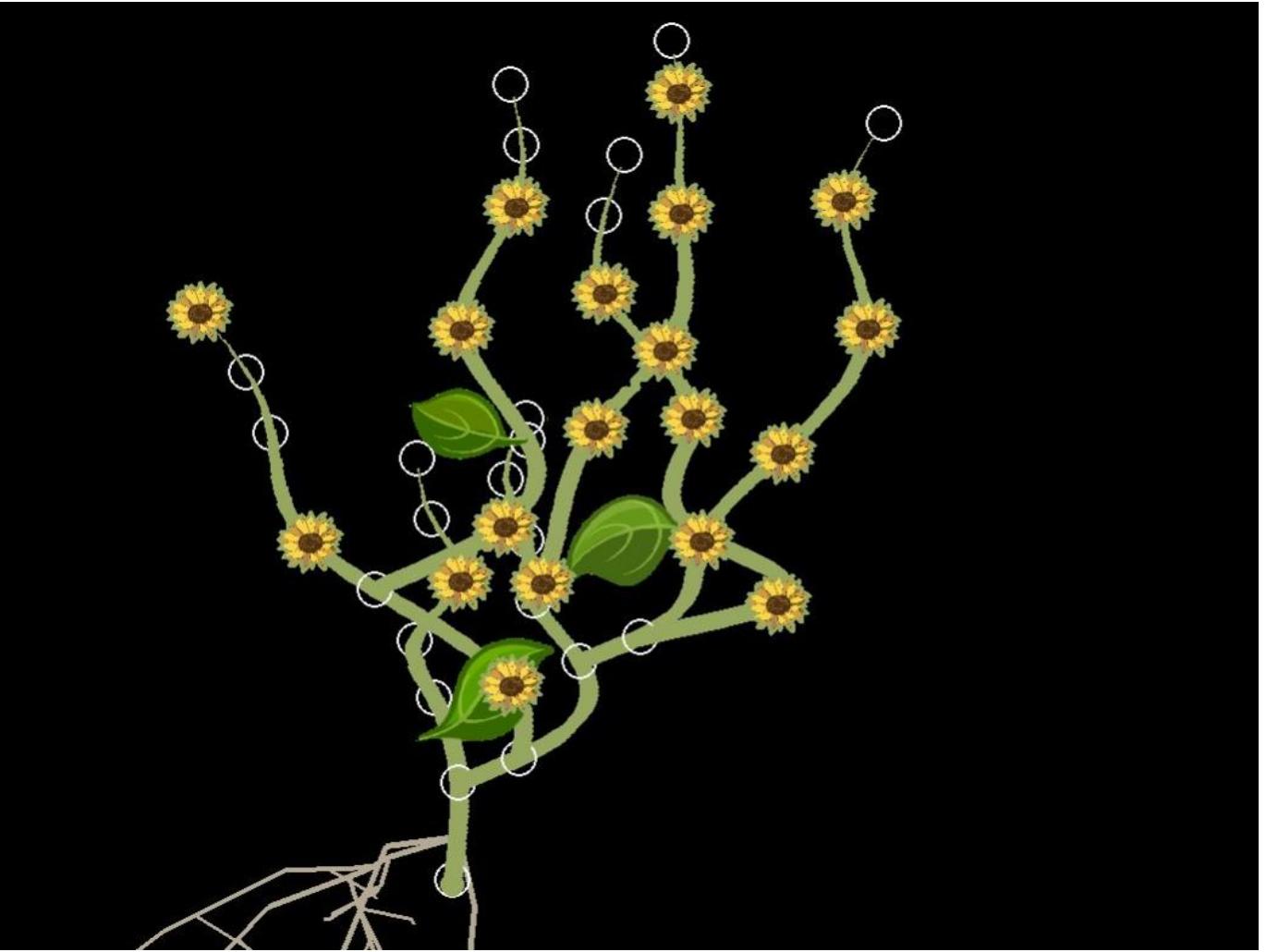
# Interface

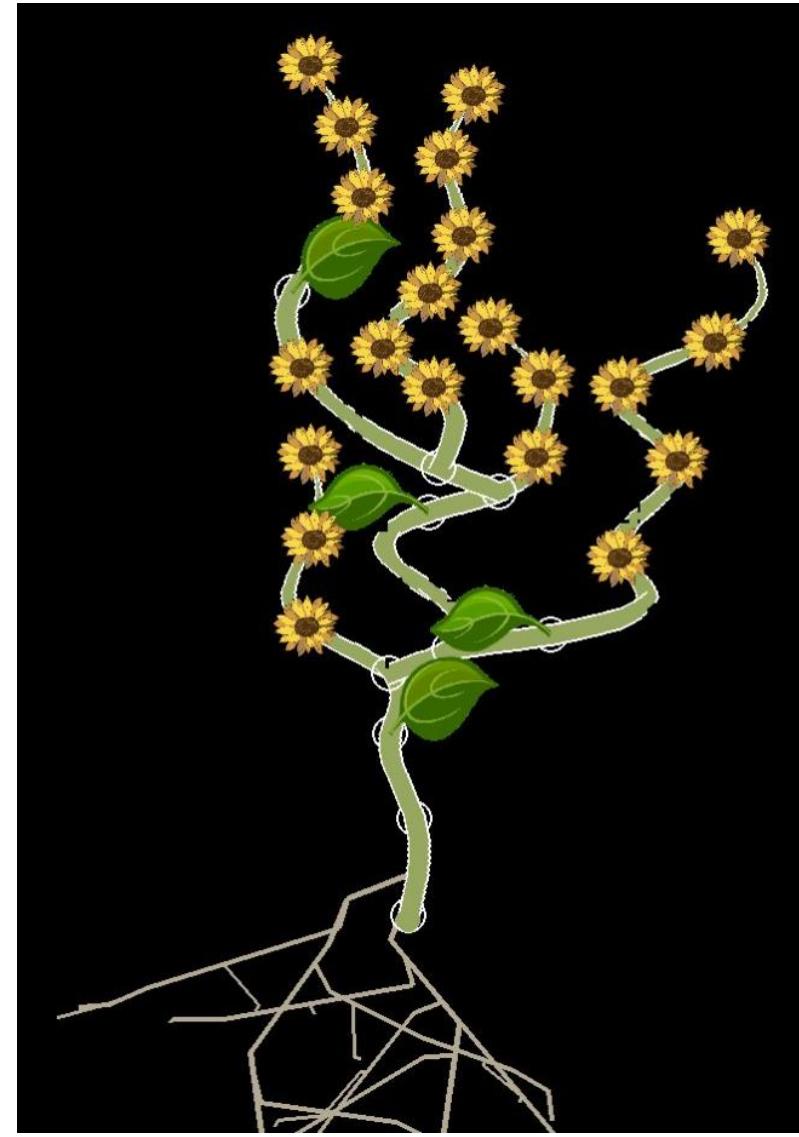
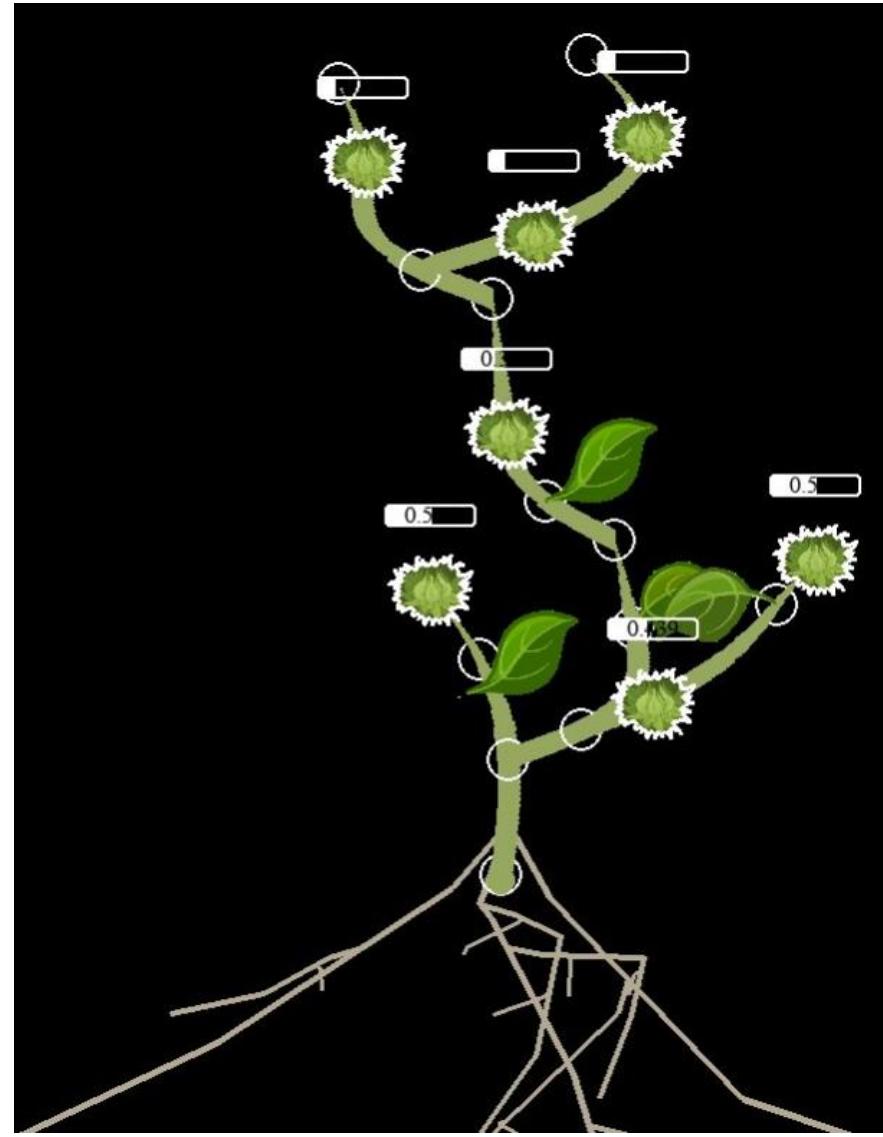
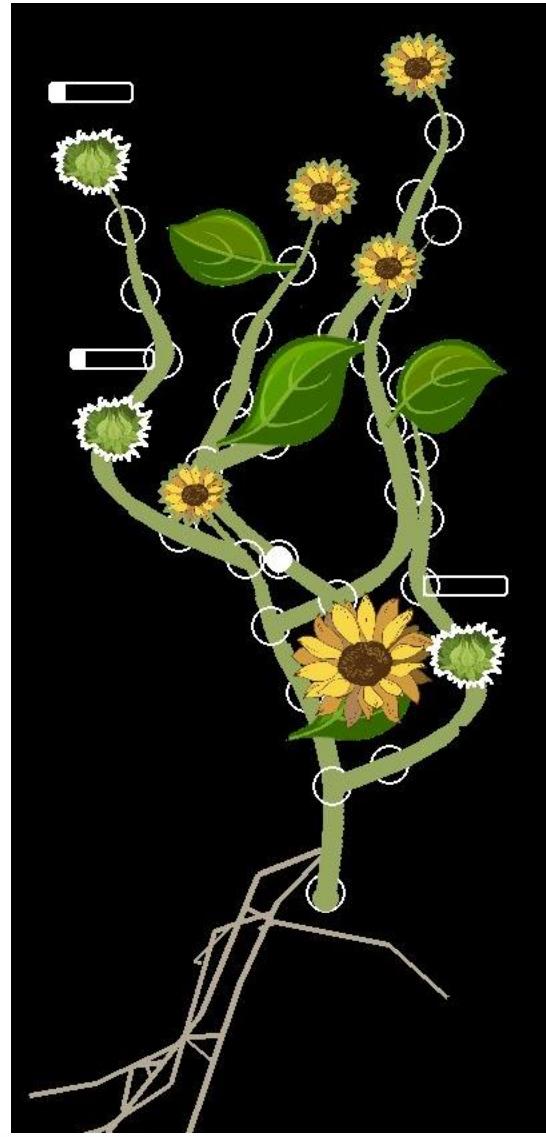


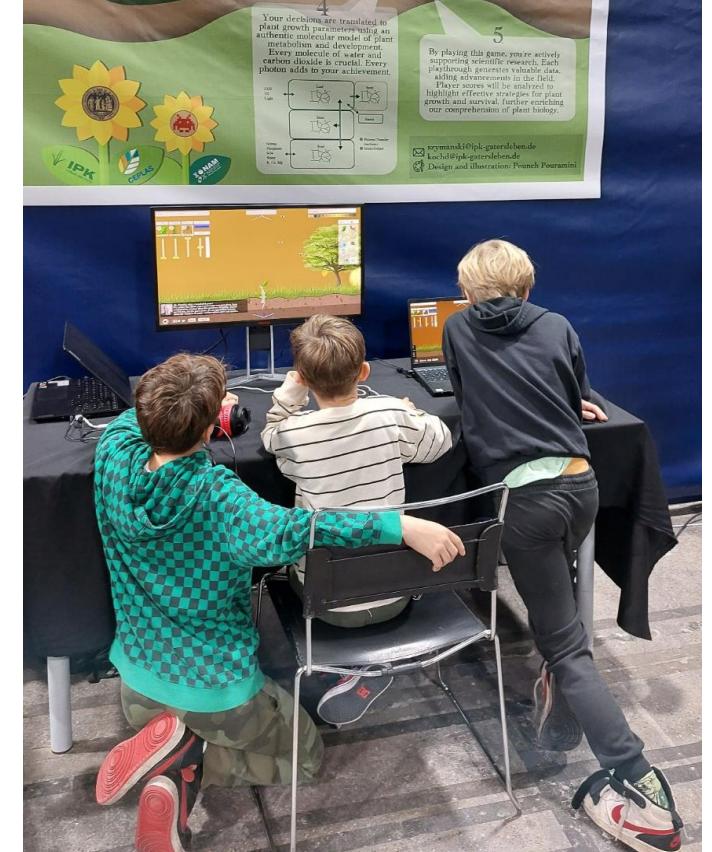
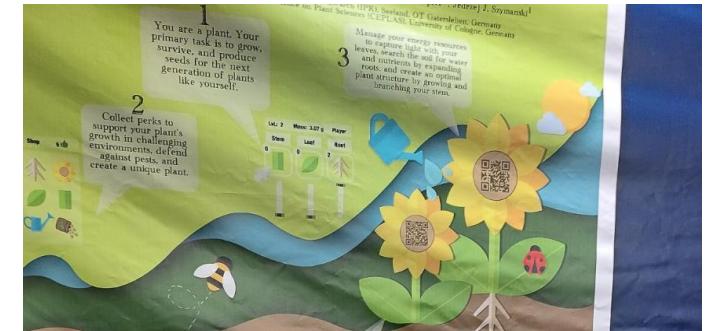
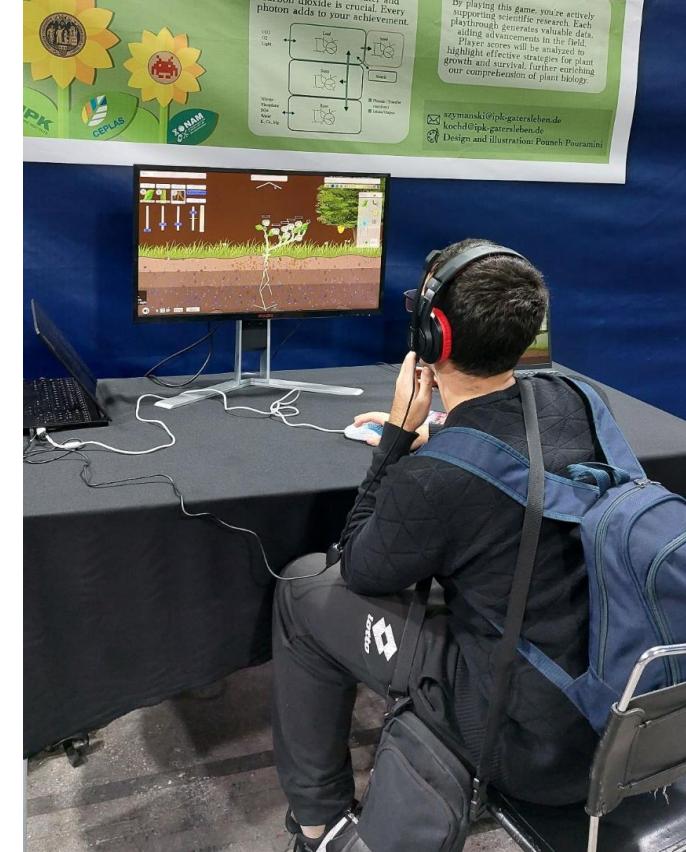
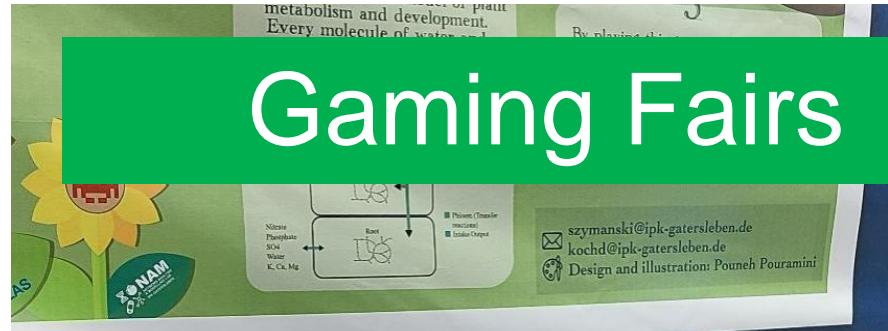
# Data collection



# Data collection



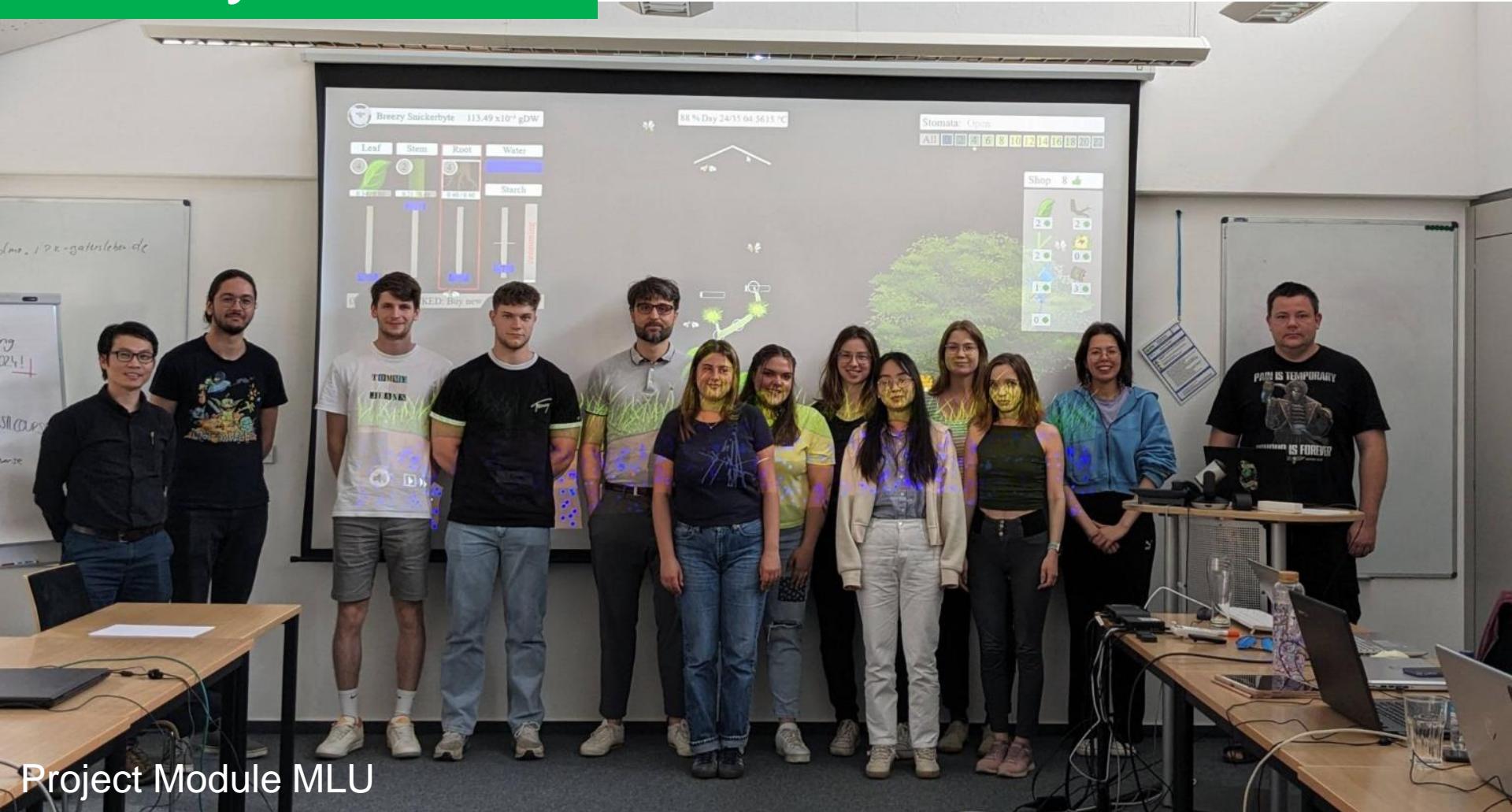




# Schools



# University courses



Project Module MLU

<https://danielkoch.itch.io/planted>



## PlantEd

### PlantEd beta

PlantEd is an engaging plant growth simulation game where players nurture a seed or seedling into maturity. Strategically allocate biomass to develop leaves, branches, roots, and flowers, aiming to maximize seed mass within a 35-day timeframe. Navigate dynamic weather conditions, manage resources efficiently, and interact with fun elements like snails, bees, and mites. Challenge yourself to optimize growth strategies and achieve the highest scores. Immerse into the fascinating world of plant biology with PlantEd!

[More information](#) ▾

### Download

[Download](#) PlantEd\_V1.0.0.zip 221 MB

[Download](#) PlantEd\_V1.0.0.tar.xz 183 MB

### Install instructions

Tipps and tricks:

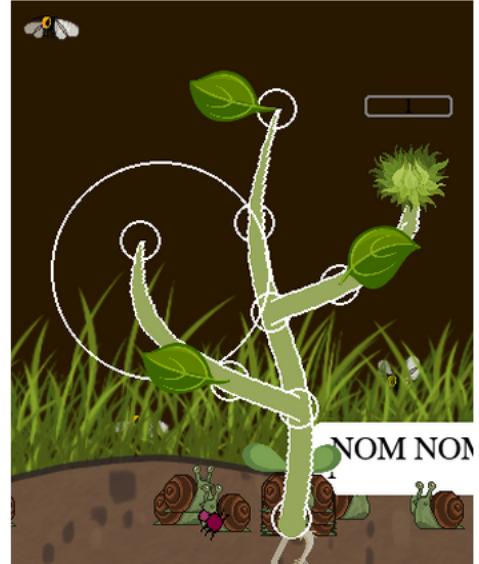
**1. Early Focus on Leaves:**

- Begin by investing in leaves to maximize photon intake, as they are crucial for the plant's energy production through photosynthesis.
- Higher photon intake will contribute to increased energy availability for other growth processes.

**2. Strategic Branch Development:**

- Once a solid leaf foundation is established, invest in branches strategically. Branches enable more leaves and branches to grow, amplifying the overall photon intake.
- Consider balancing branch and leaf growth to optimize energy production.

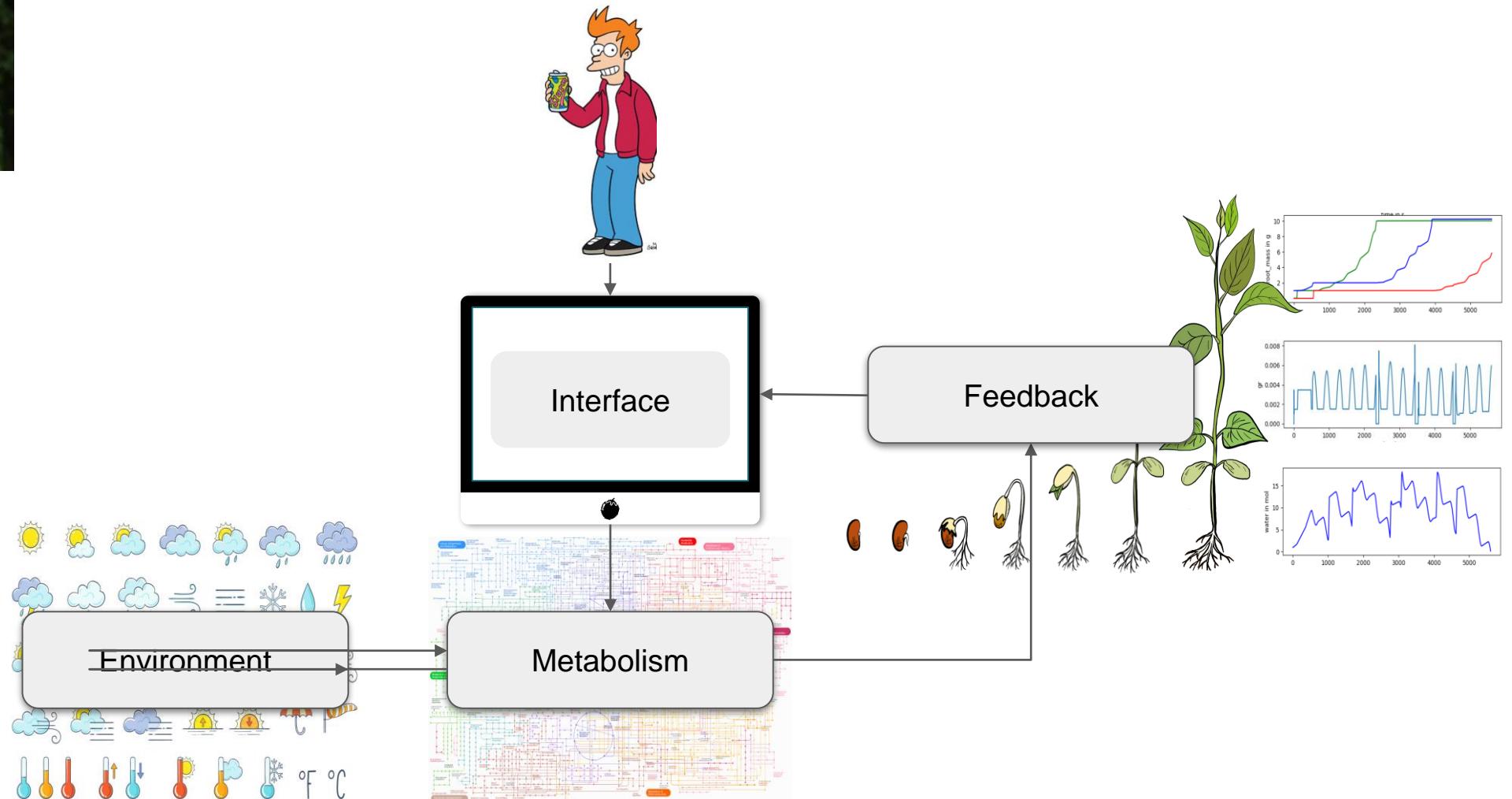
**3. Root Expansion for Resources:**





Daniel Koch

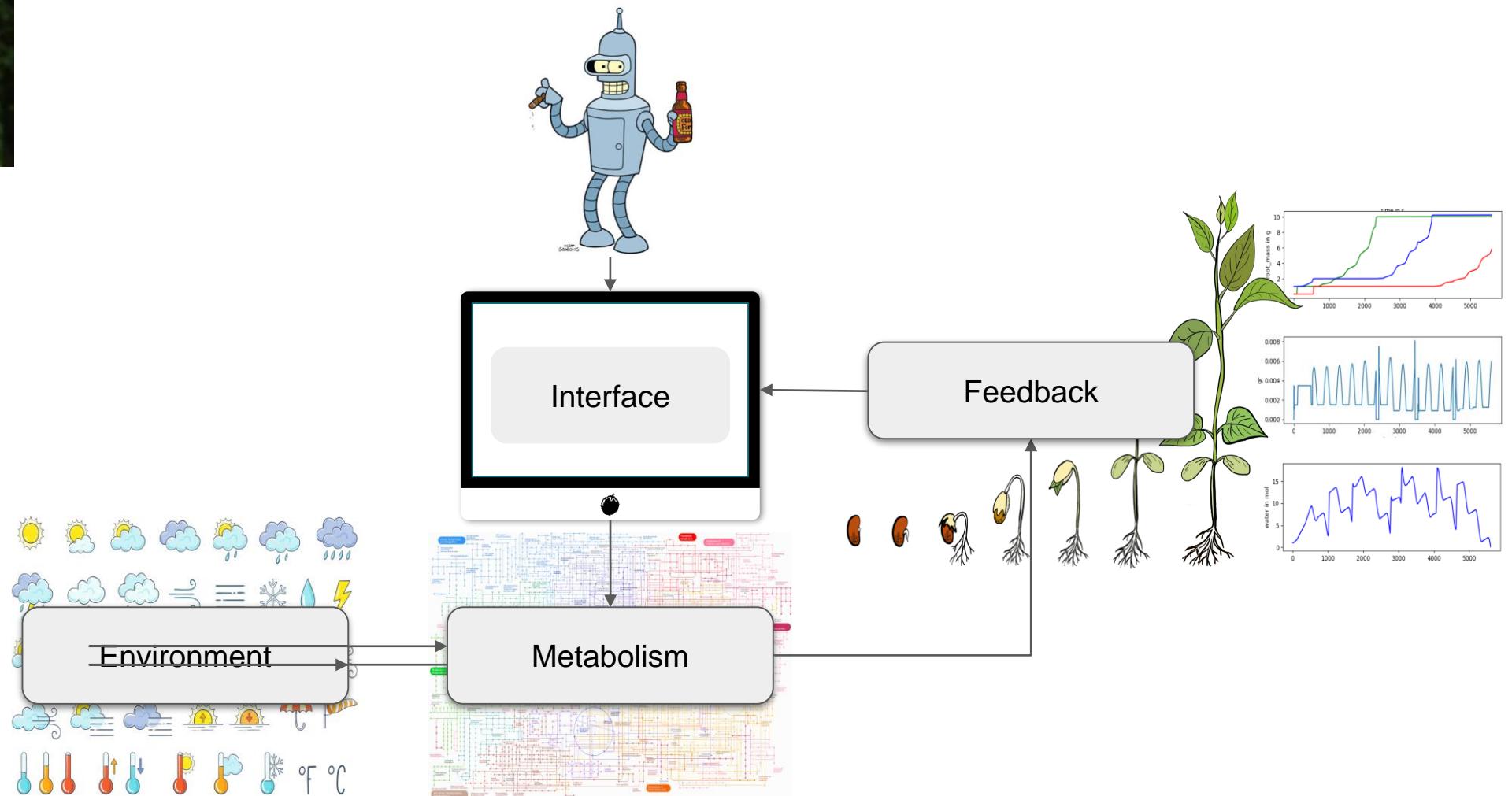
# Human agent





Daniel Koch

# AI agent

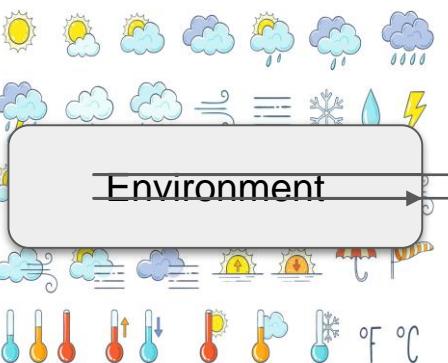




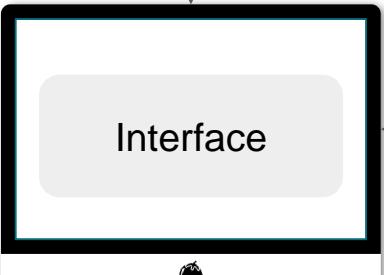
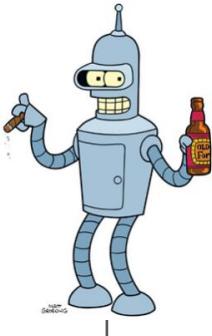
Daniel Koch



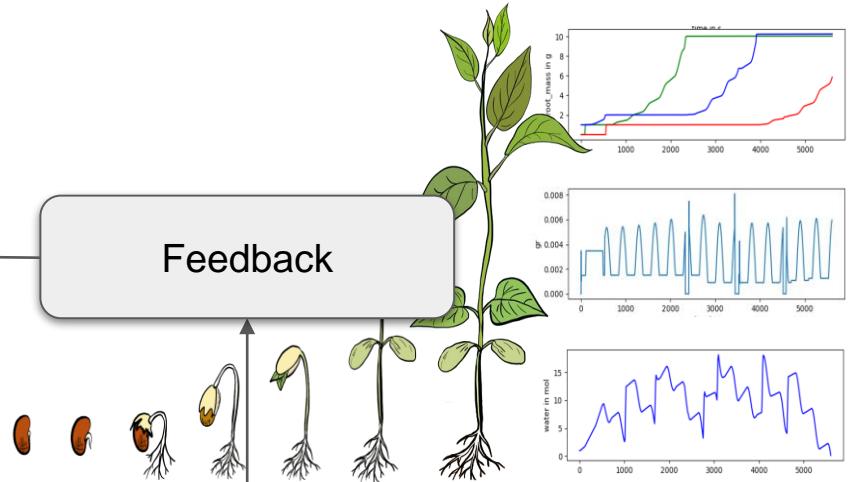
Dennis Psaroudakis



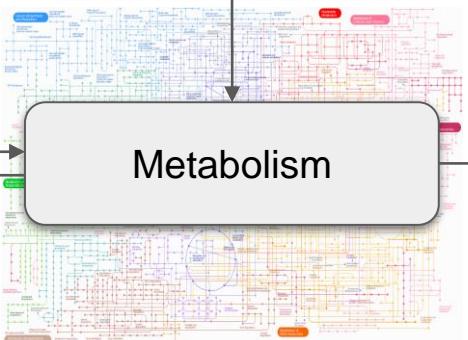
# AI agent



Feedback



Metabolism

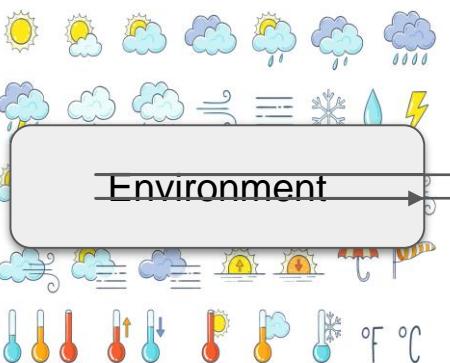




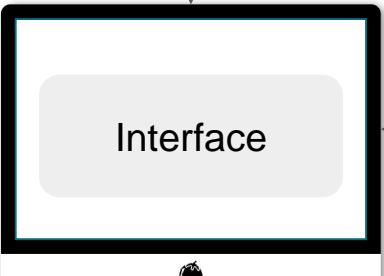
Daniel Koch



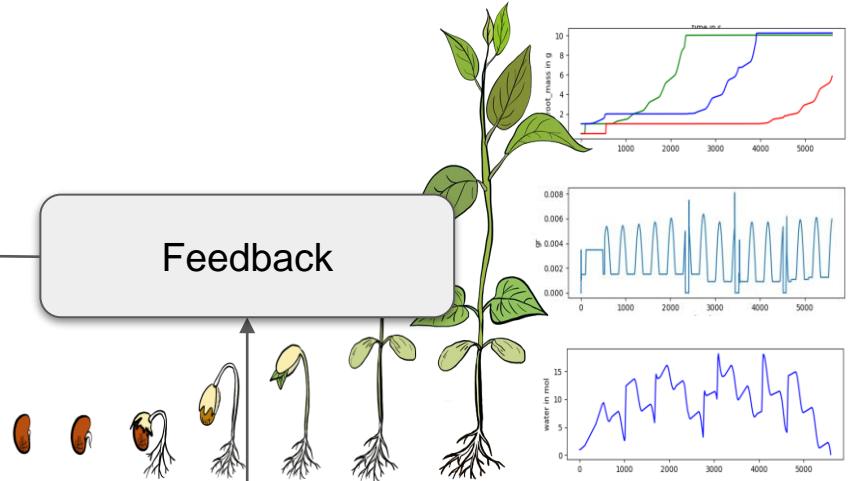
Dennis Psaroudakis



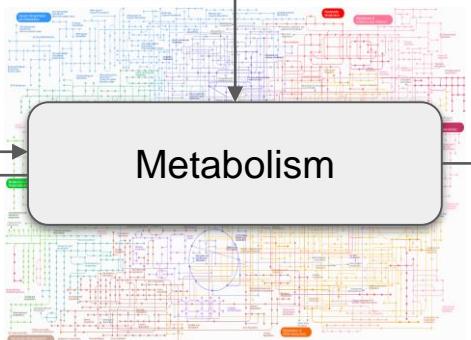
# AI agent



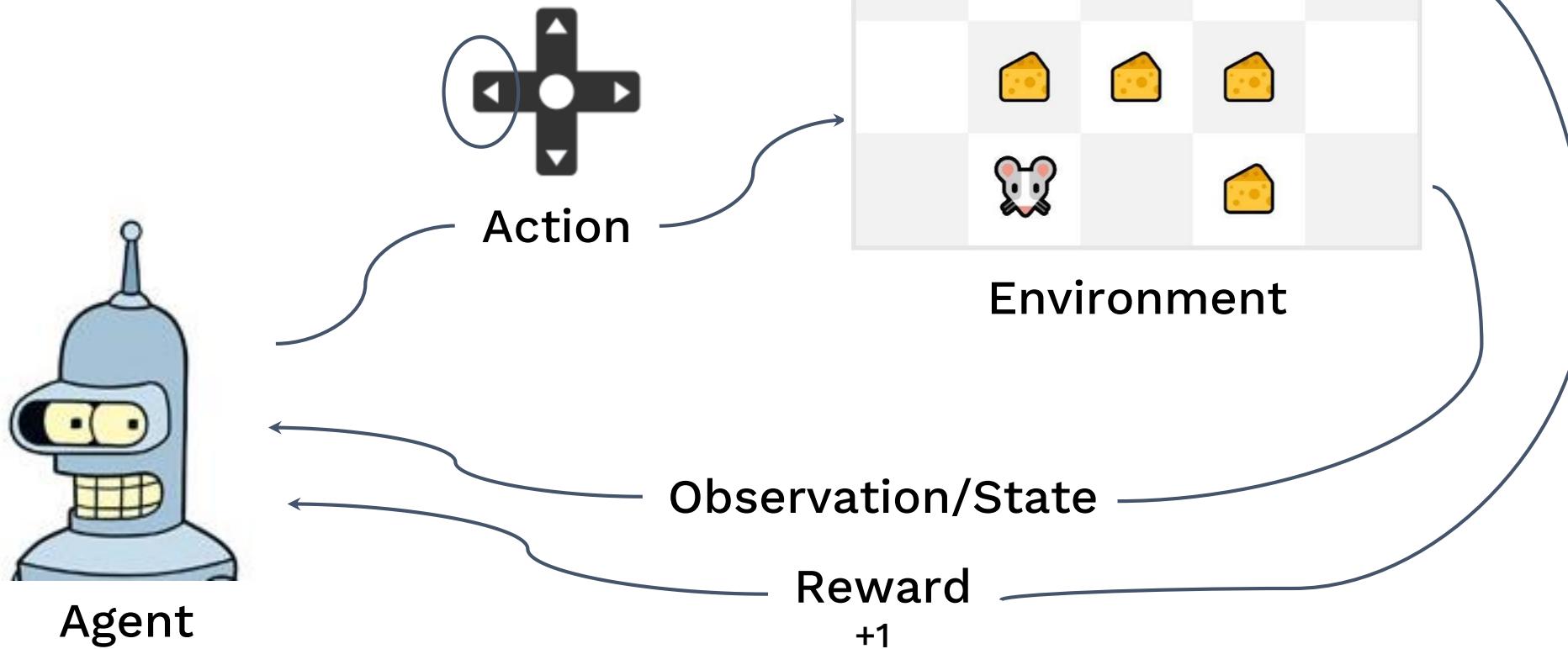
Feedback



Metabolism



# Reinforcement Learning

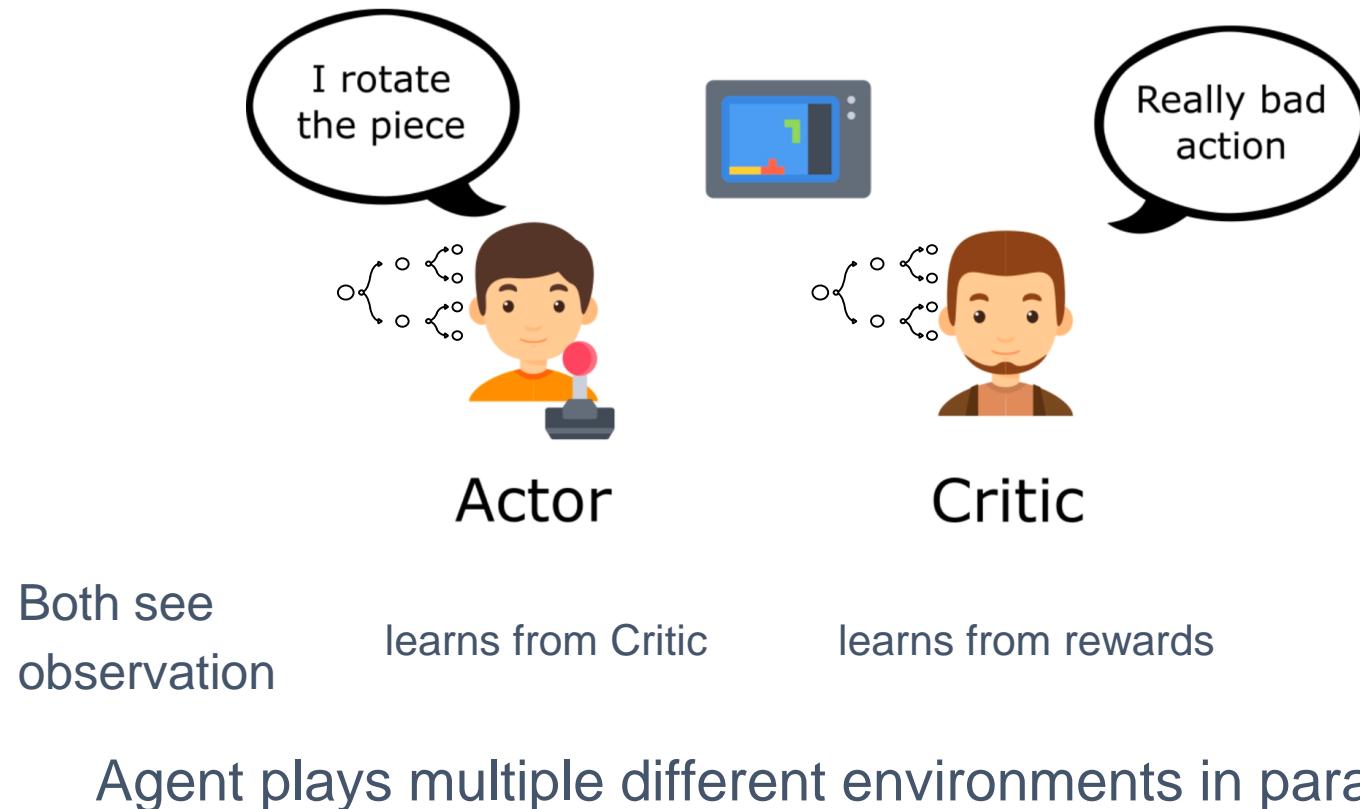


# Our Agent

## Action Space

- Grow Leaves
- Grow Stem
- Grow Roots
- Produce Starch
- Open Stomata
- Close Stomata
- Buy new Leaf
- Buy new Branch
- Buy new Root
- Buy new Flower

## Algorithm: Proximal Policy Optimization (PPO)



PPO Algorithm: Schulman et al. 2017. Implementation: Stable Baselines 3, Raffin et al. 2021

# Just grow!

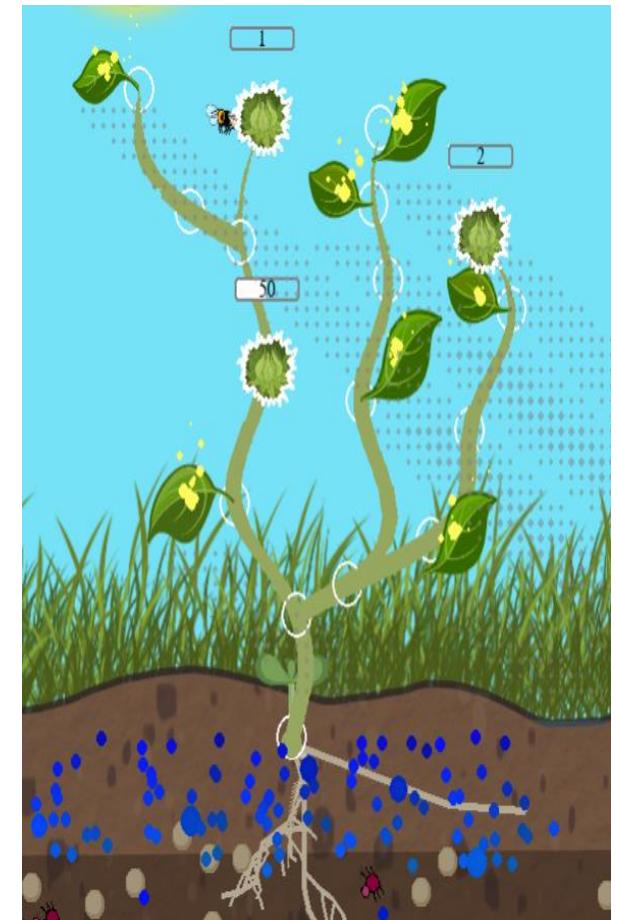
Level 1

## Action Space

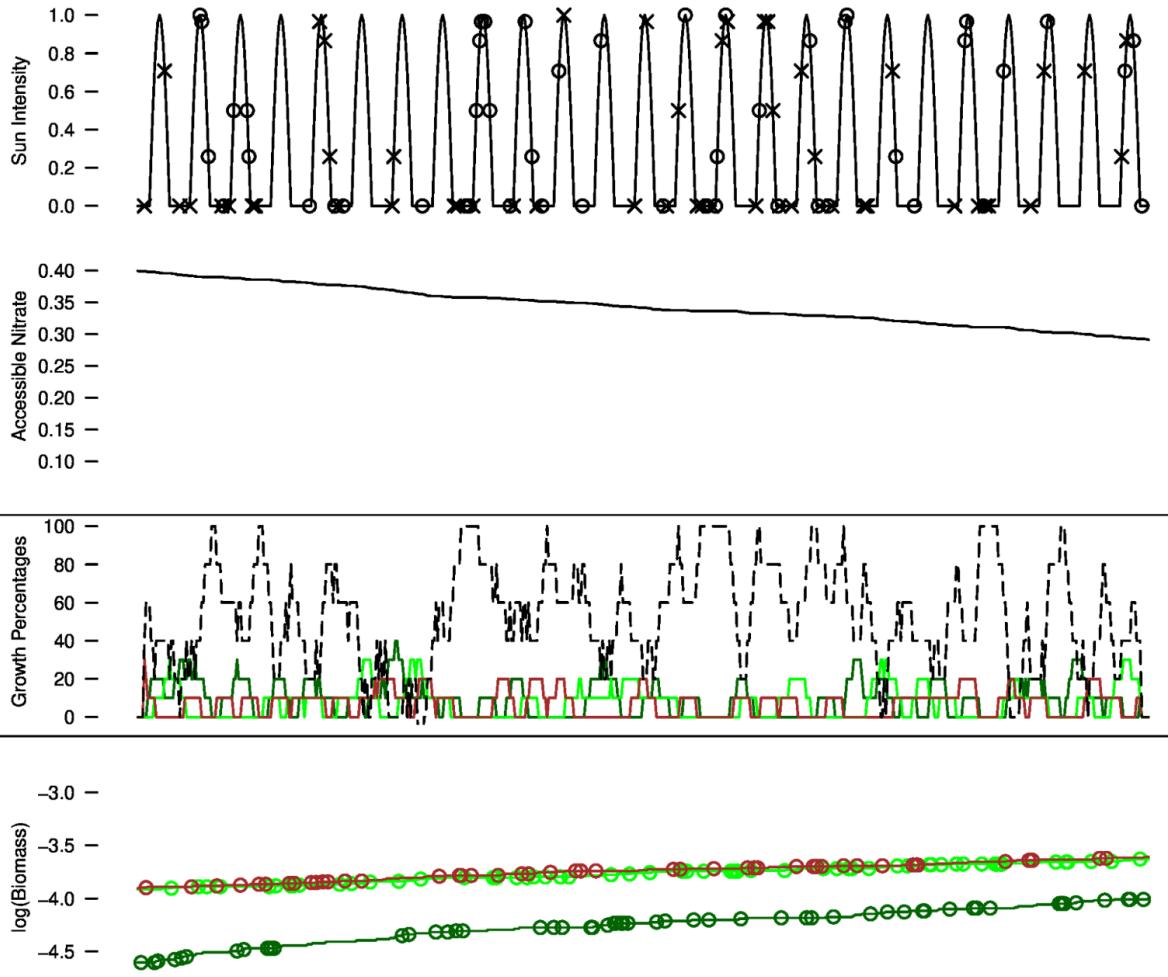
### Reward Function

Difference in total biomass  
(g) since last step

- Grow Leaves
- Grow Stem
- Grow Roots
- Produce Starch
- Open Stomata
- Close Stomata
- Buy new Leaf
- Buy new Branch
- Buy new Root
- Buy new Flower



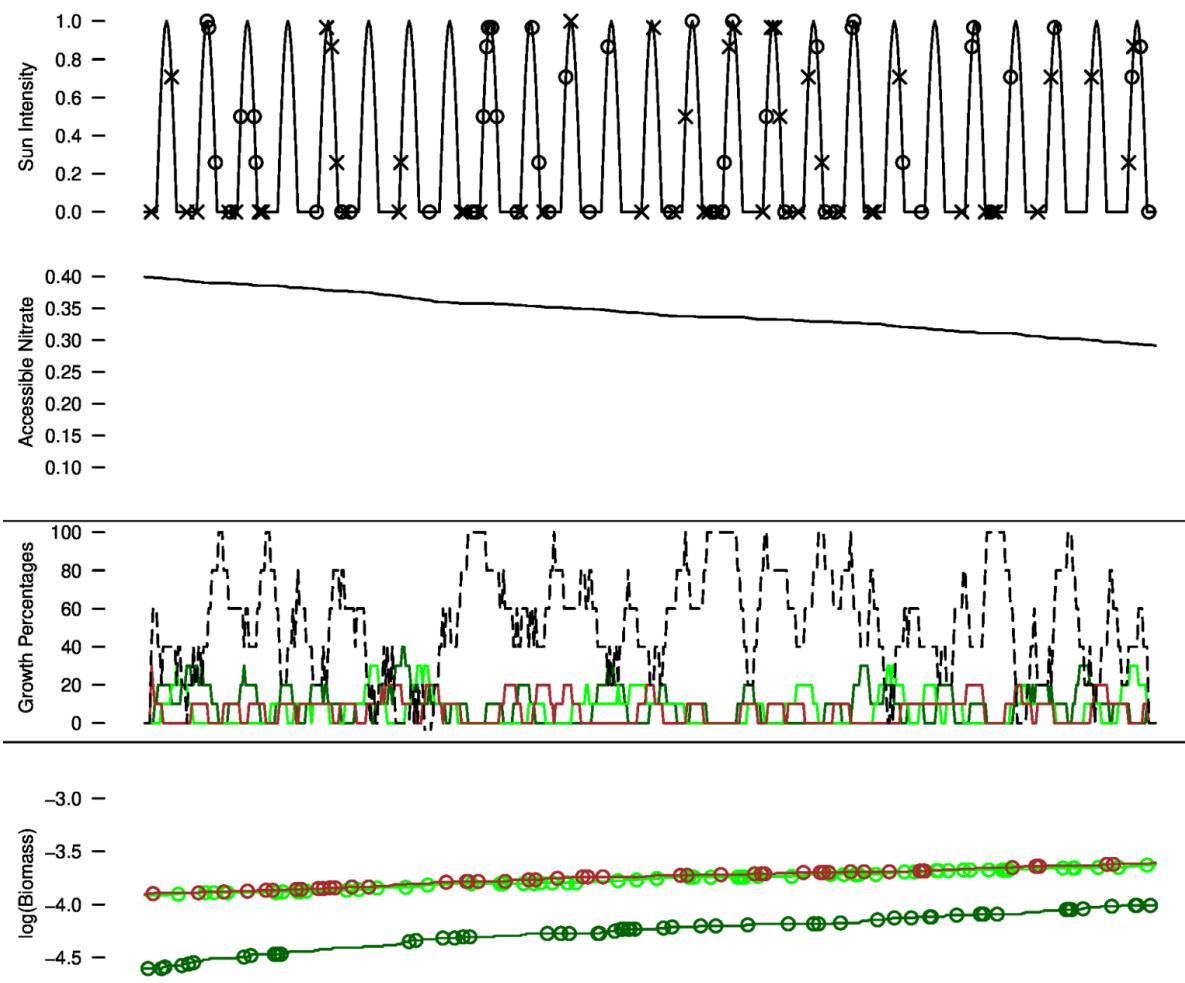
# Low Nitrate



Episode 1

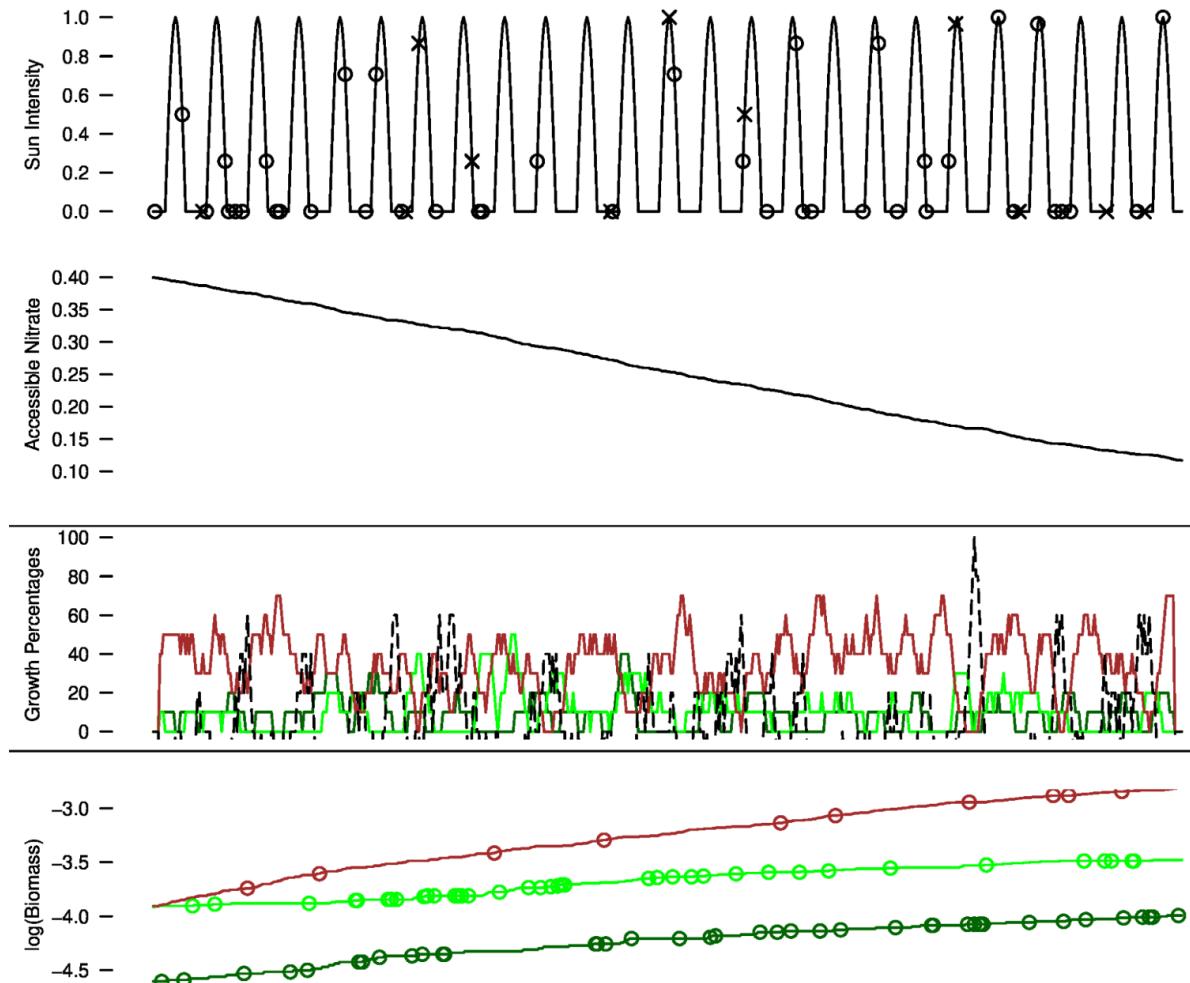
root  
stem  
leaf

# Low Nitrate



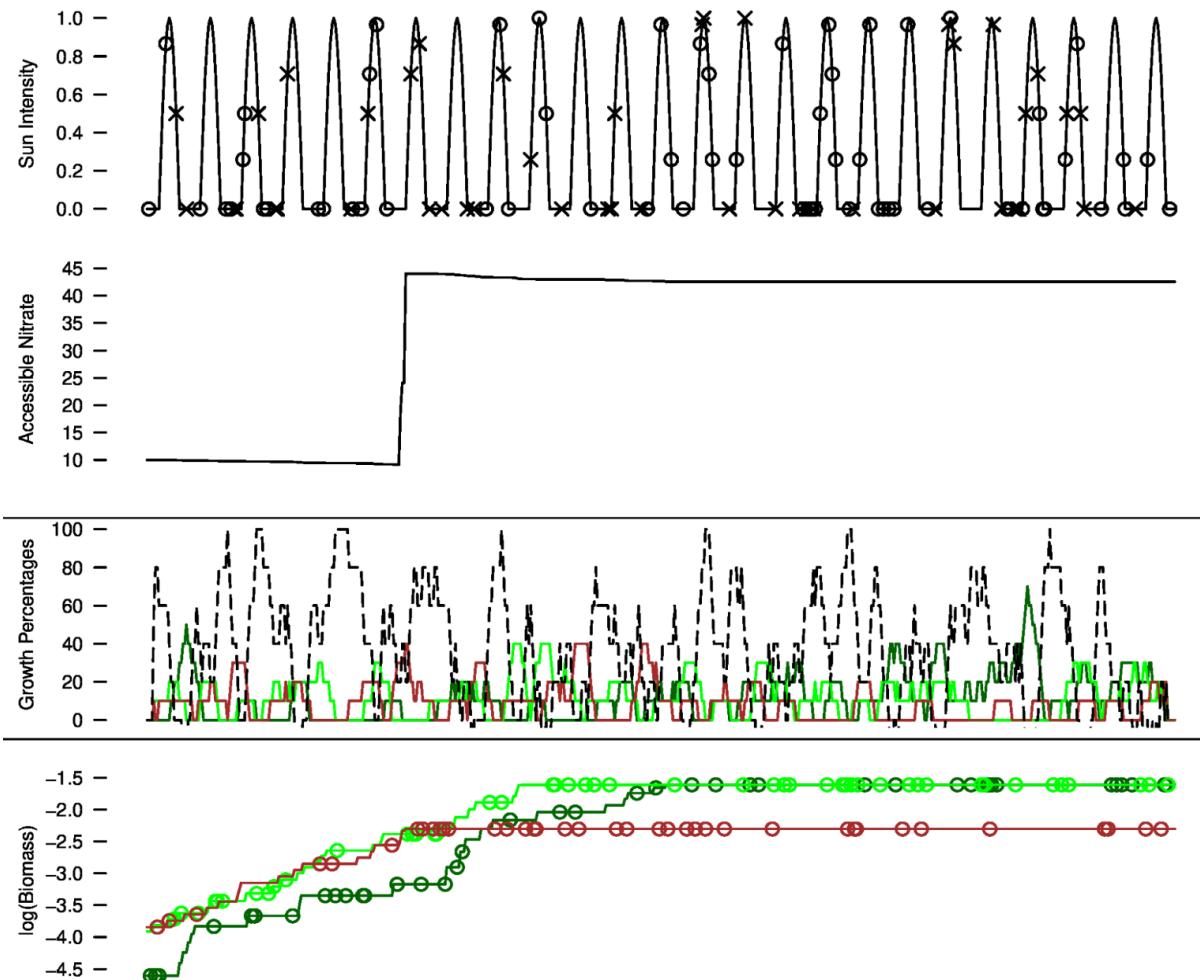
Episode 1

root  
stem  
leaf



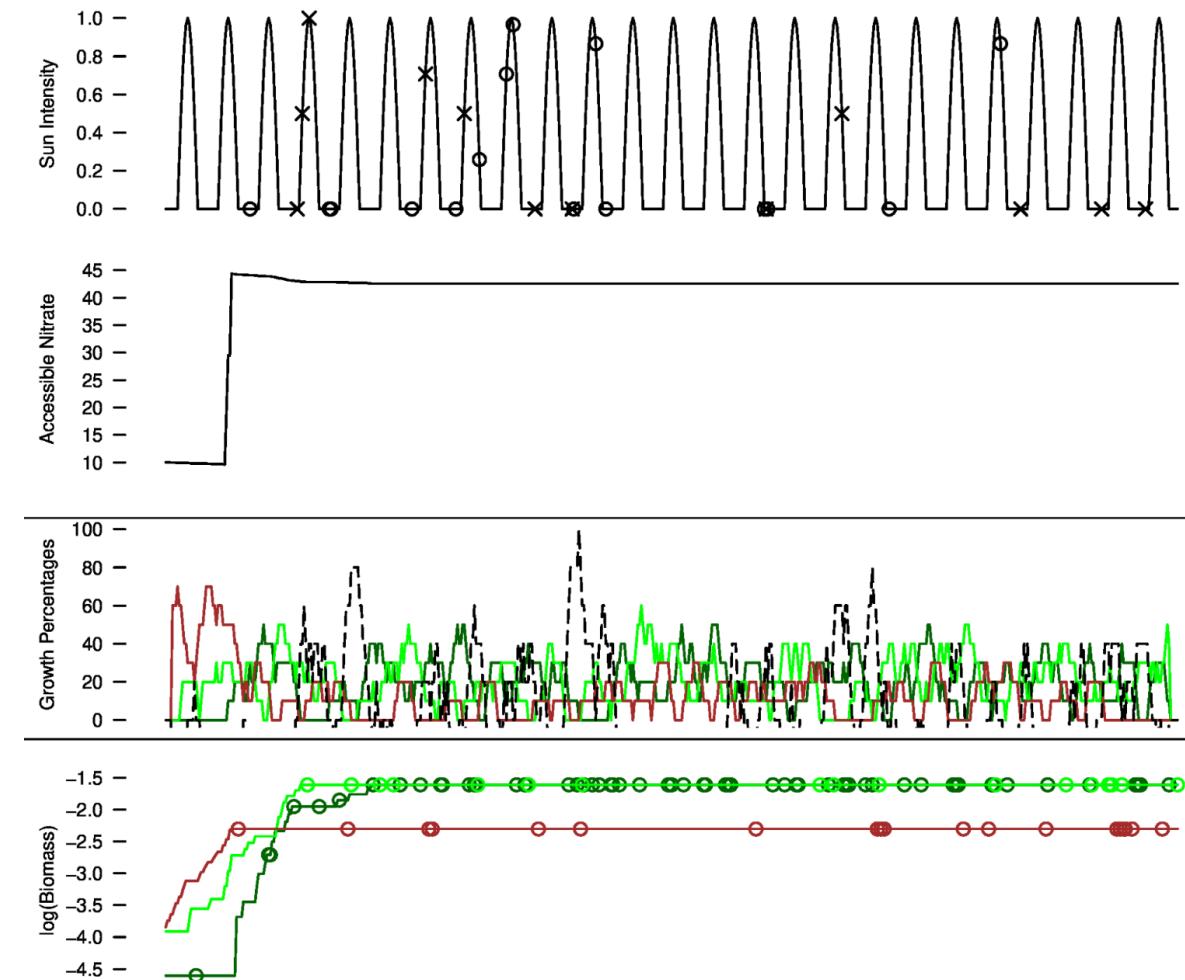
Episode 9

# High Nitrate



Episode 1

root  
stem  
leaf



Episode 15

# Let's go to Seed!

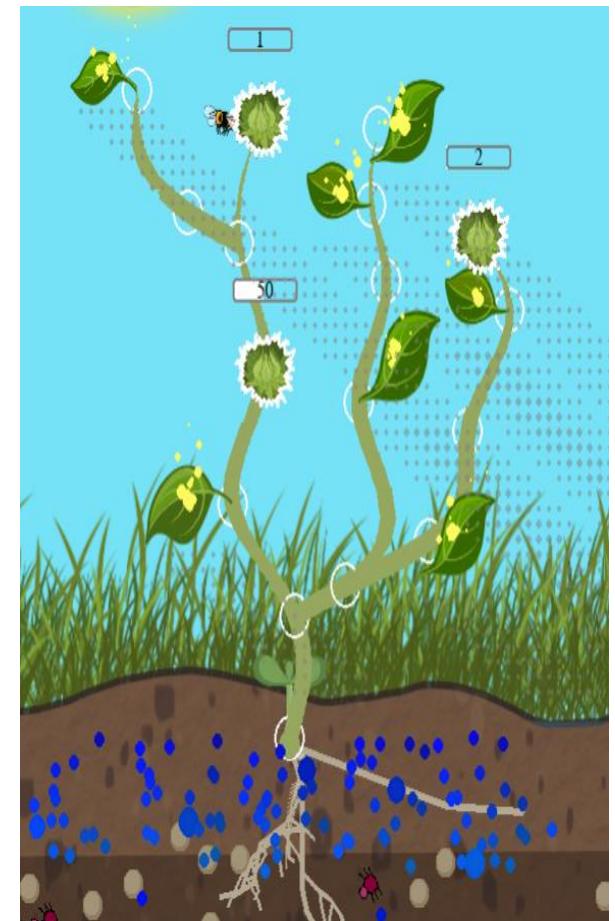
## Level 2

Double Reward for seed mass!

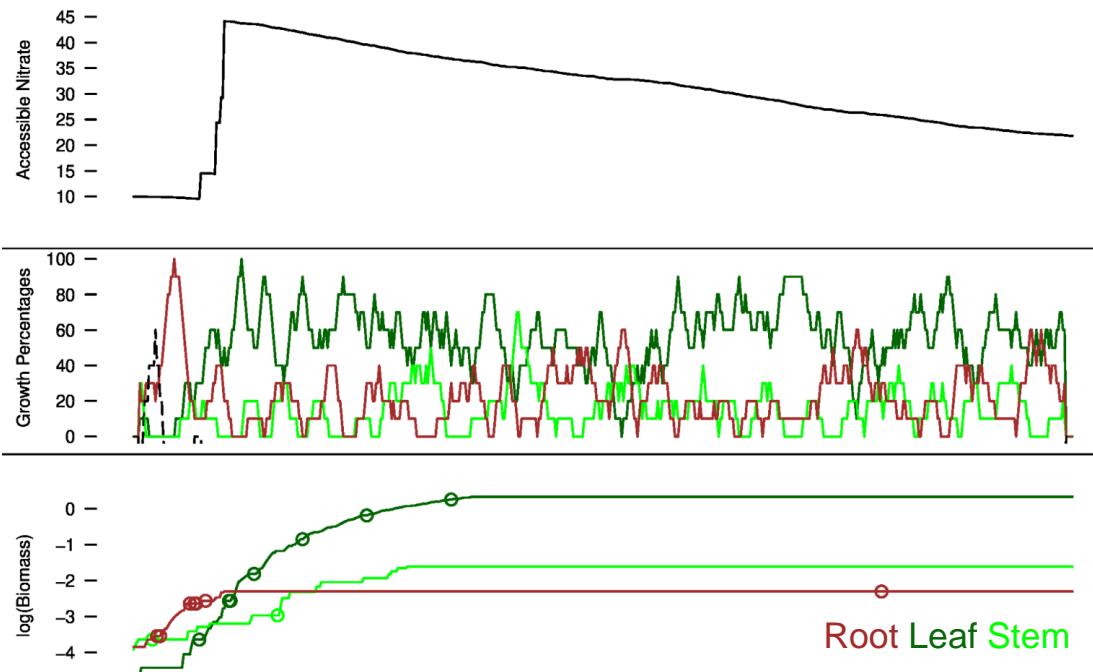
Penalties (in relation to current biomass)

### Action Space

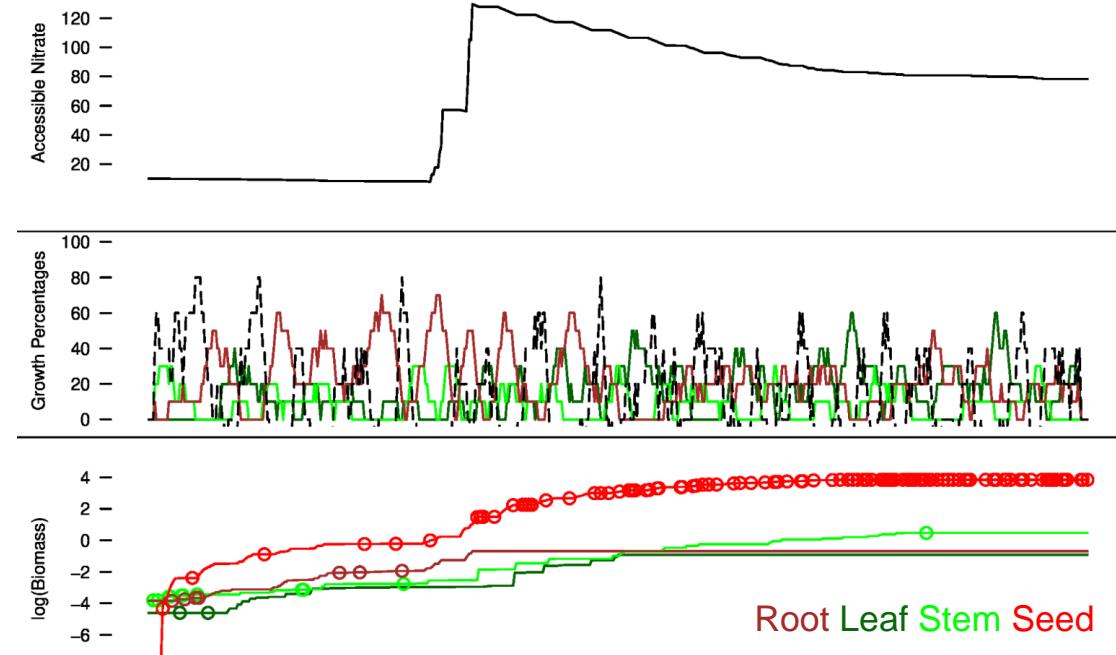
- Grow Leaves
- Grow Stem
- Grow Roots
- Produce Starch
- Open Stomata
- Close Stomata
- Buy new Leaf
- Buy new Branch
- Buy new Root
- Buy new Flower



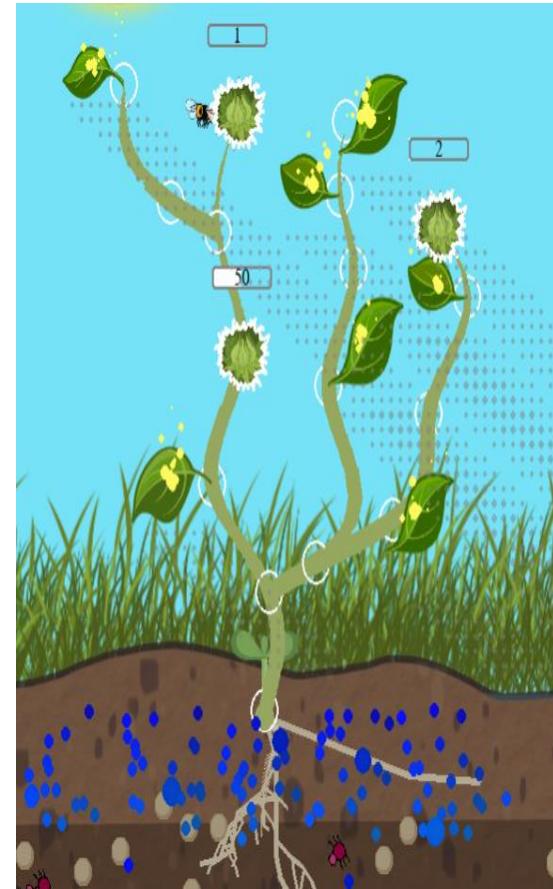
# Spring (High Nitrate) max biomass



# Spring (High Nitrate) max seed production



# A digital twin



- Light intensity
- Light spectrum
- Temperature
- Humidity
- Water availability
- Nutrients
- .
- .
- $\text{CO}_2$
- $\text{O}_2$
- Time of the day
- Day/night cycle
- Season
- etc.
- .
- .
- .

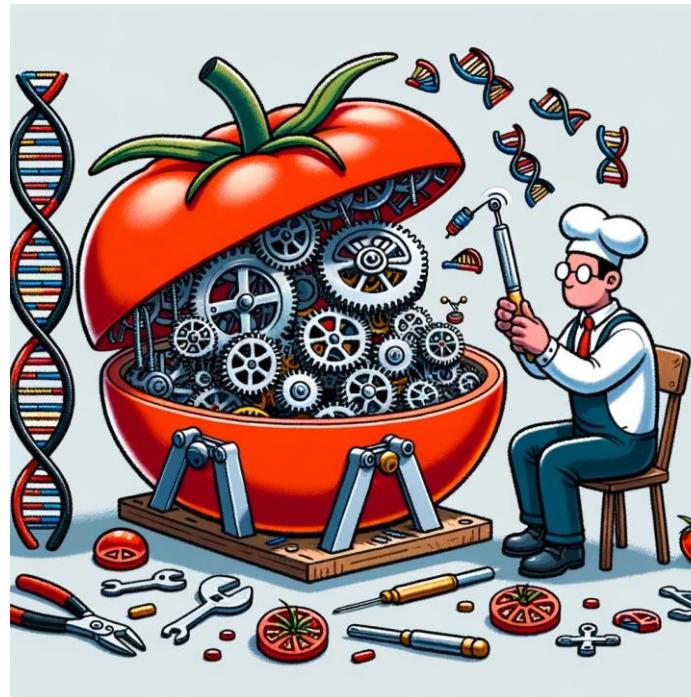


# Summary

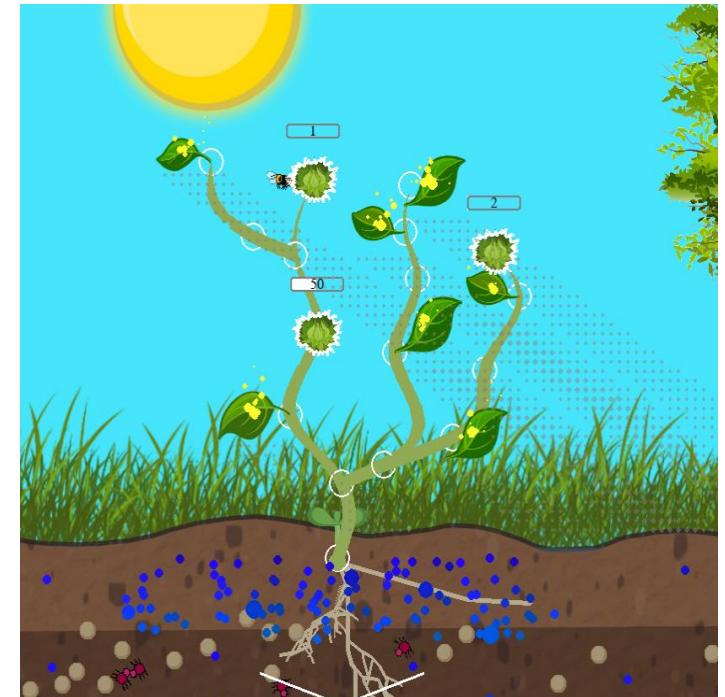
AI for detective work



AI for reverse engineering



AI for computer games



# Thanks to:



Simon Zumkeller



Fritz Peleke



Dennis Psaroudakis



Ankur Sahu



Uni Cologne  
Nadine Töpfer



MPIPZ @Cologne  
Thomas Hartwig and Julia  
Engelhorn



Merle Stein



Gernot Schmitz



Daniel Koch



Kalyan Pininti



Uni Cologne  
Julietter de Meaux



IPK  
Gatersleben  
Jozefus  
Schippers



Sana Sangeen



Adeboye Adejoro



Max Opitz



FZ Jülich  
HHU  
Björn  
Usadel

..and many many others!



# Thank you for your attention!



@NAMlab



[www.szymanskilab.com](http://www.szymanskilab.com)

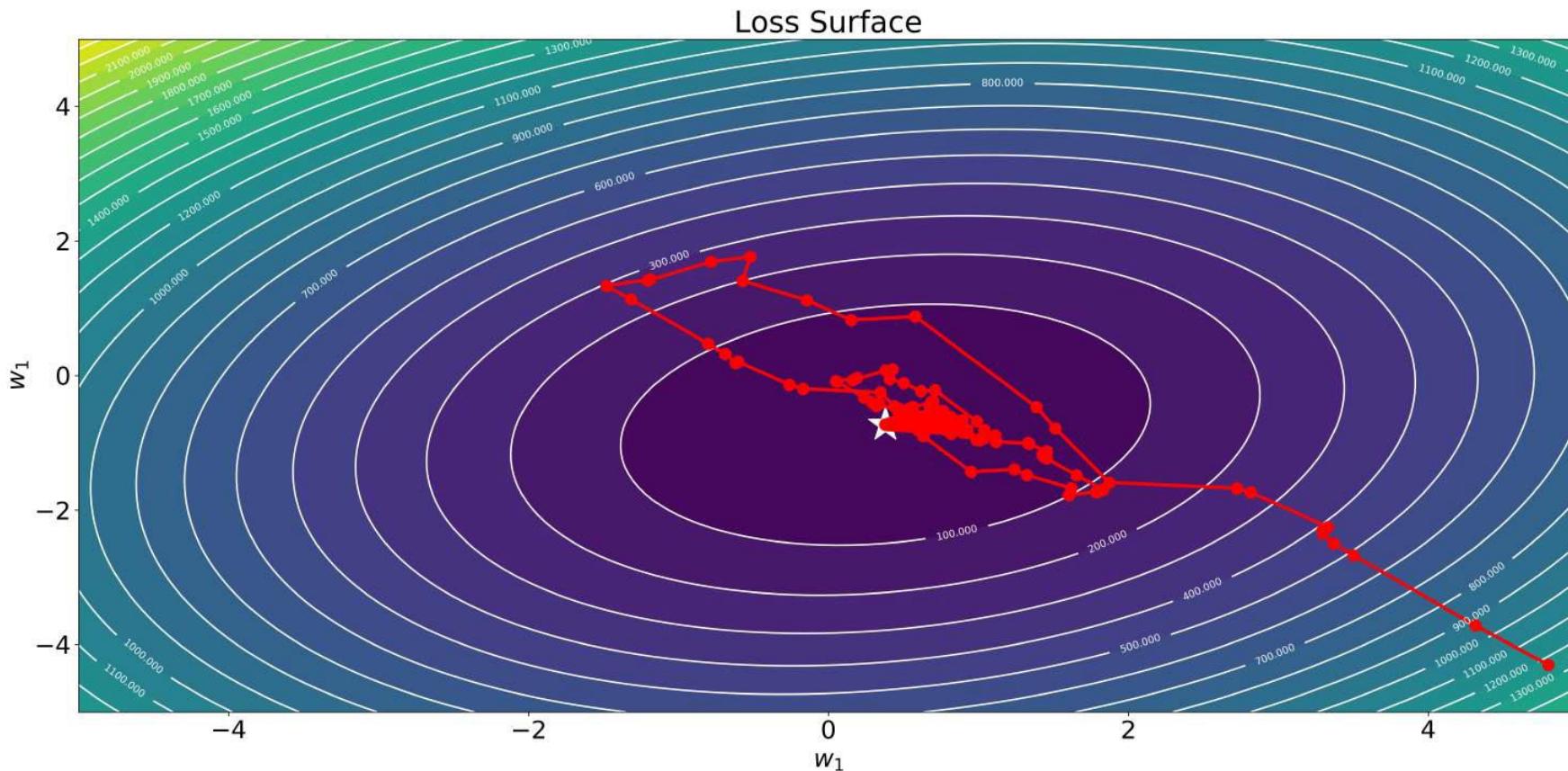
[szymanski@ipk-gatersleben.de](mailto:szymanski@ipk-gatersleben.de)



## WE ARE HIRING!

postdoc & PhD positions in  
machine learning and gene  
regulation

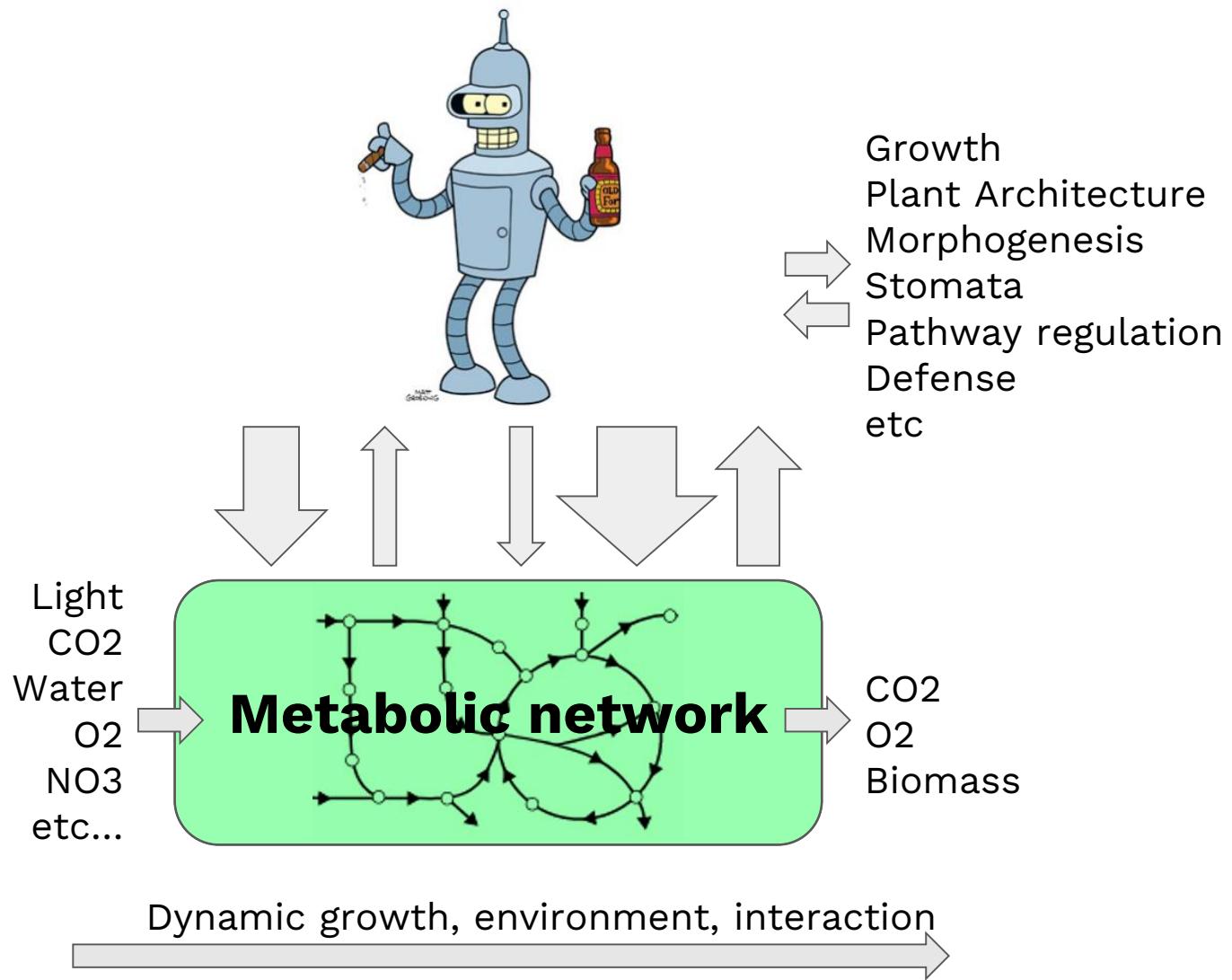
# Gradient descent – tracking the learning progress



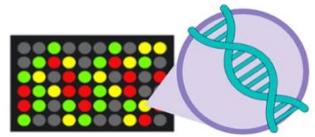


Dennis Psaroudakis  
IPK Gatersleben  
FZ Jülich

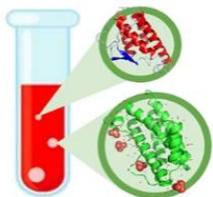
AI



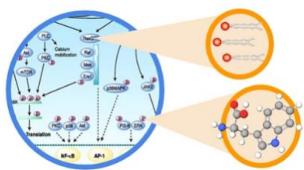
# What kind of machine learning / AI?



Transcriptomics  
~ 100k transcripts



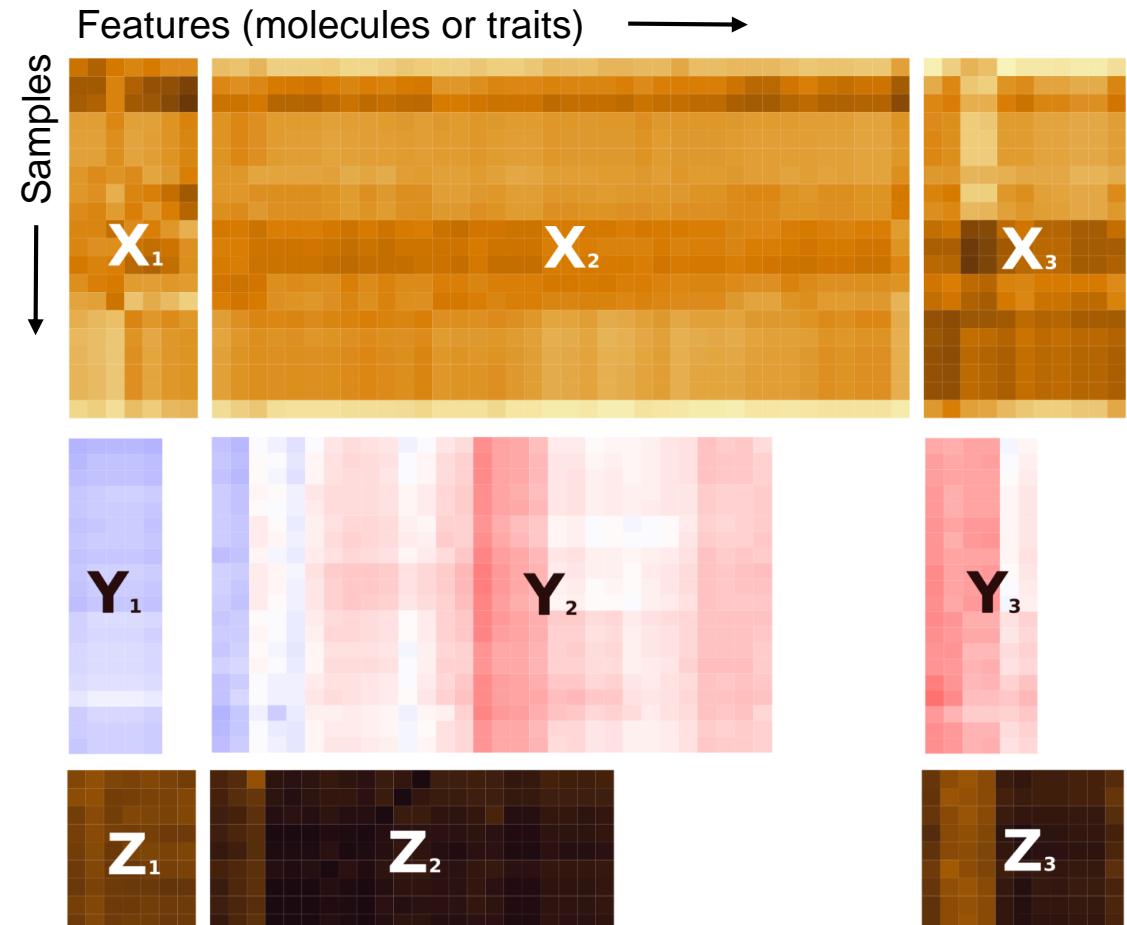
Proteomics  
~ 10k proteins



Metabolomics  
~ 3k metabolites



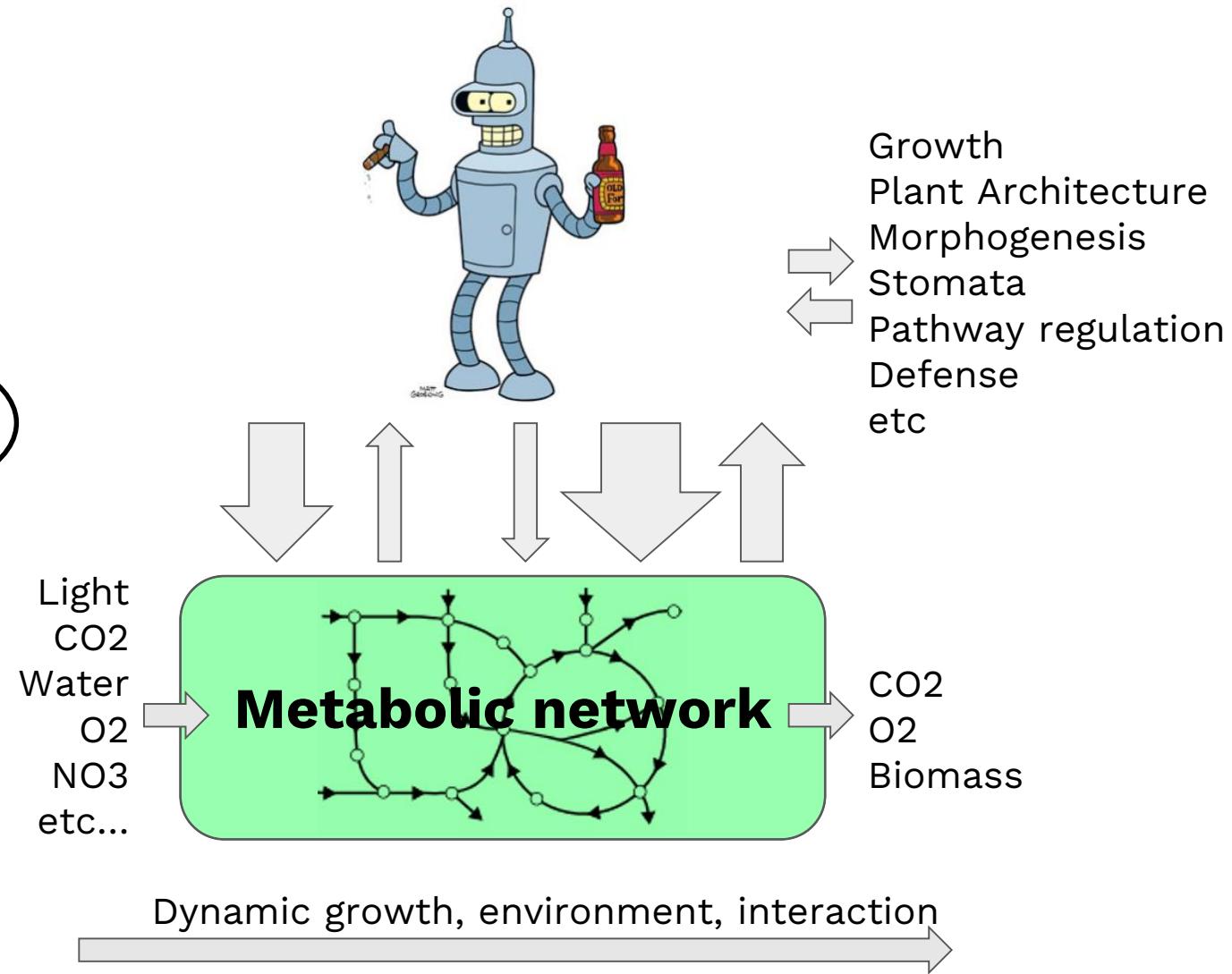
Phenomics  
~ 300 phenotypes



AI



Dennis Psaroudakis  
IPK Gatersleben  
FZ Jülich



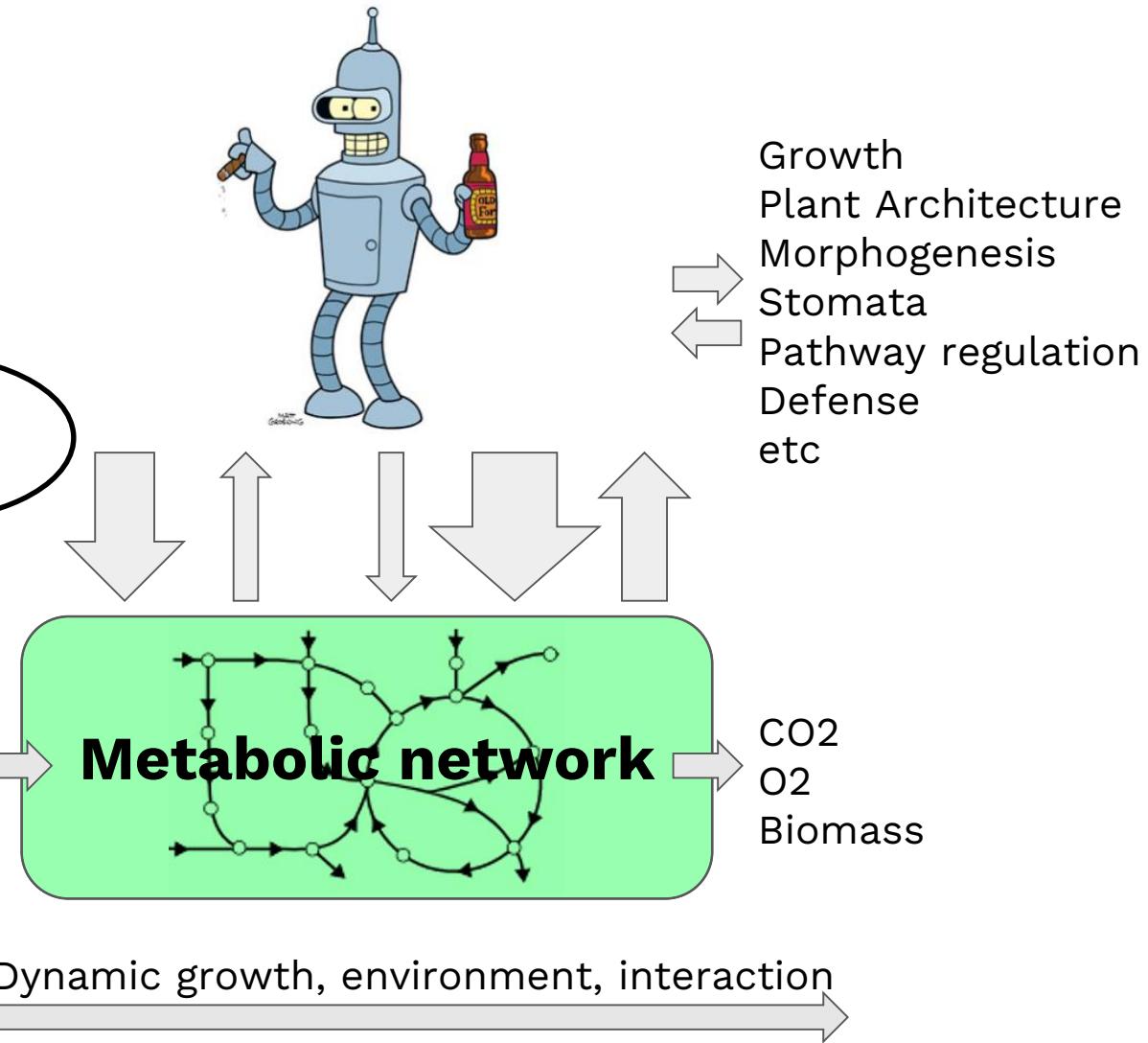


Dennis Psaroudakis  
IPK Gatersleben  
FZ Jülich

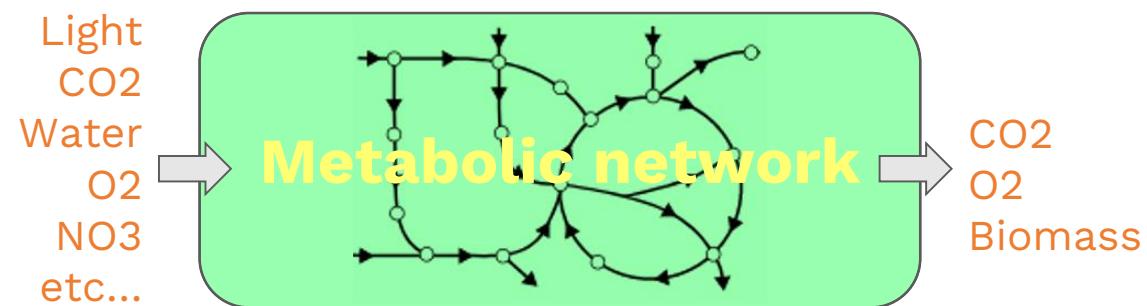
AI

Done

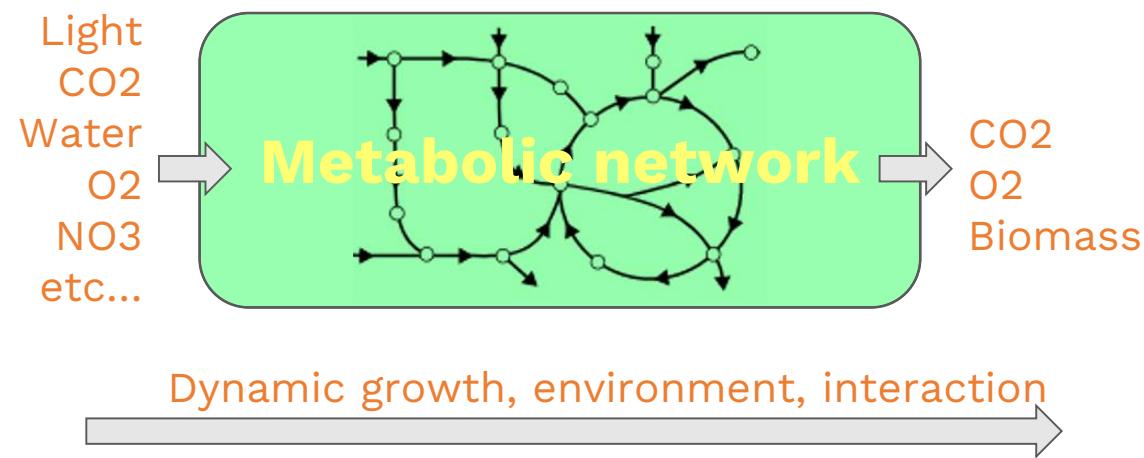
Light  
CO<sub>2</sub>  
Water  
O<sub>2</sub>  
NO<sub>3</sub>  
etc...



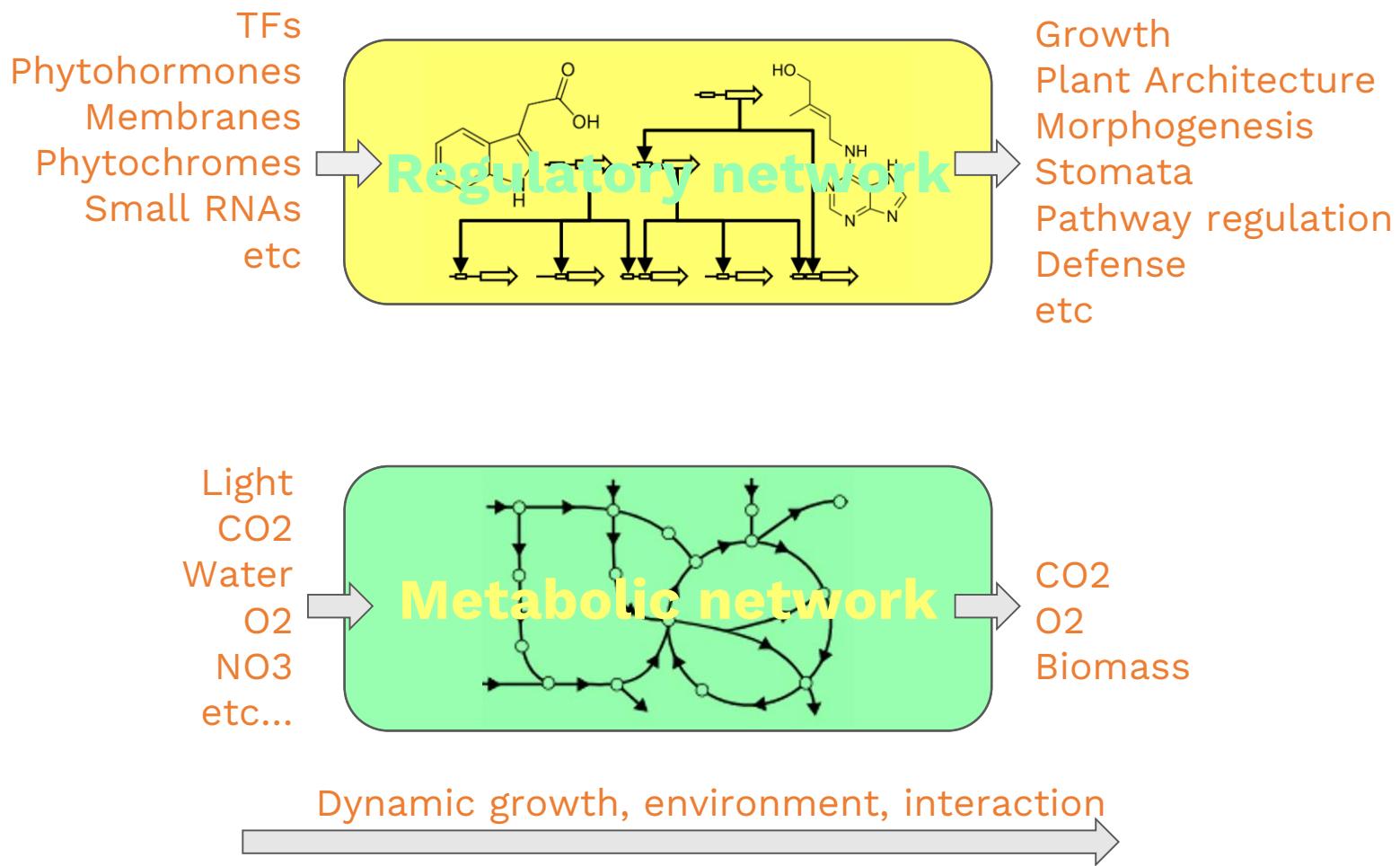
# The challenge



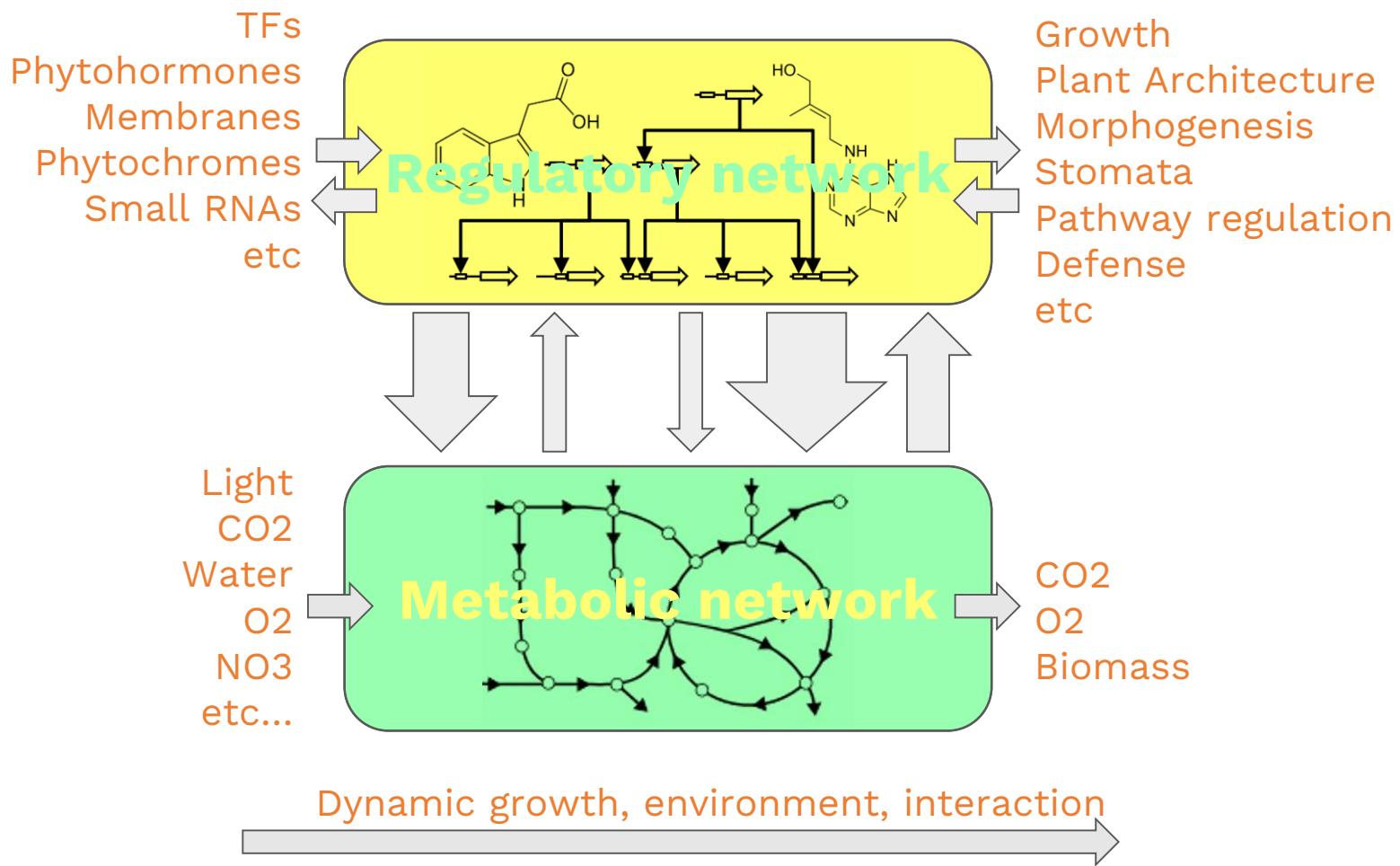
# The challenge



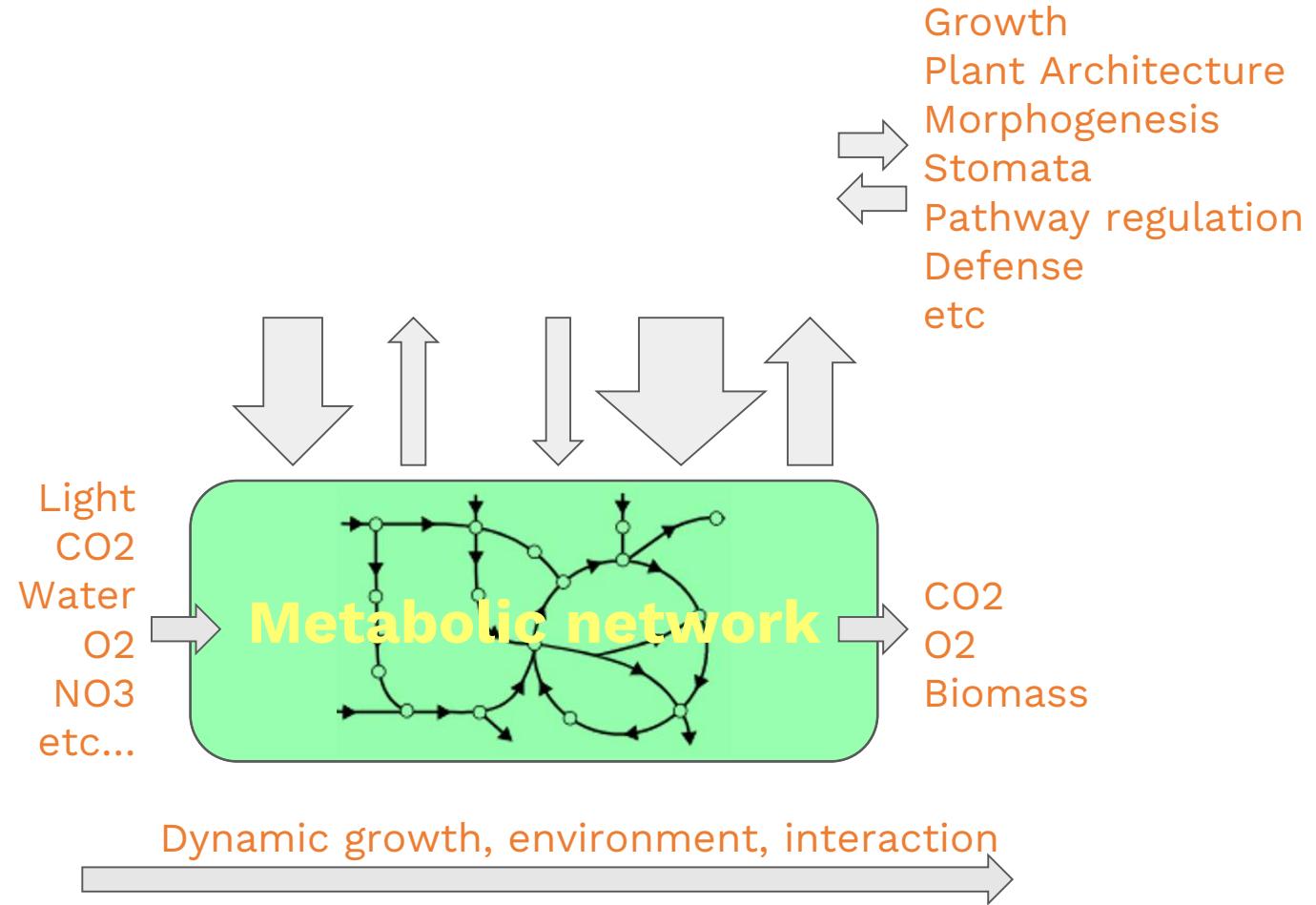
# The challenge



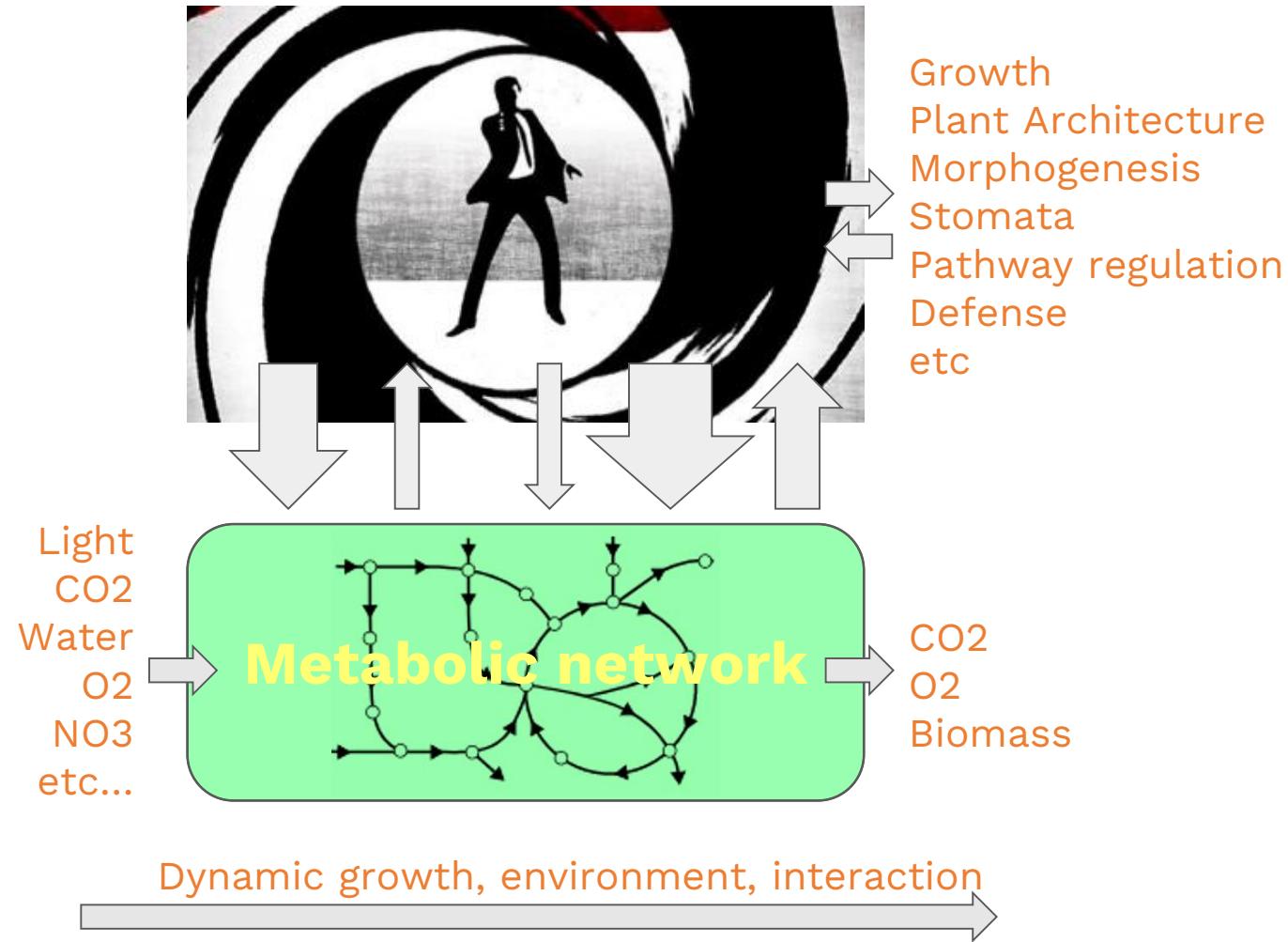
# The challenge



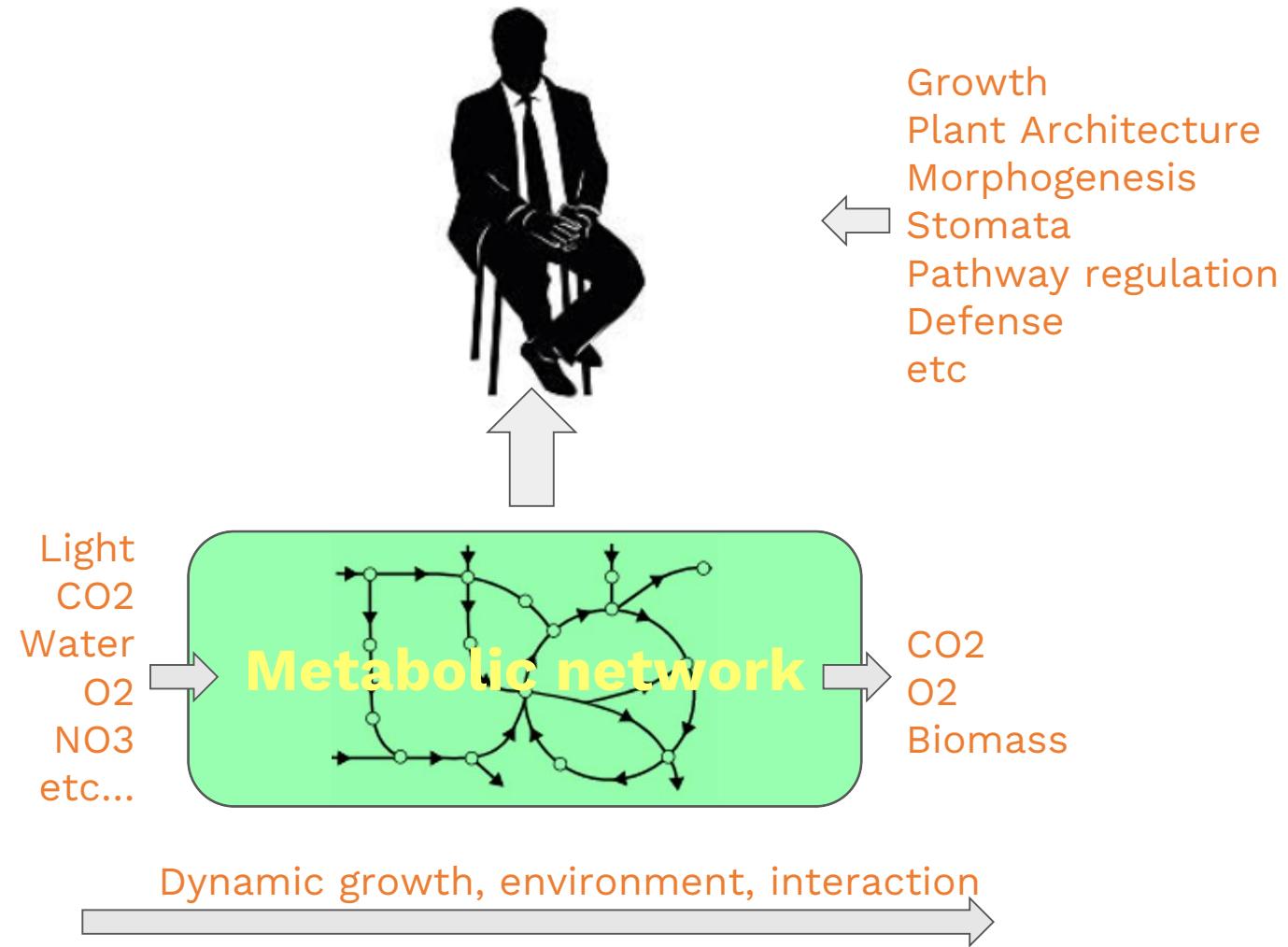
# The challenge



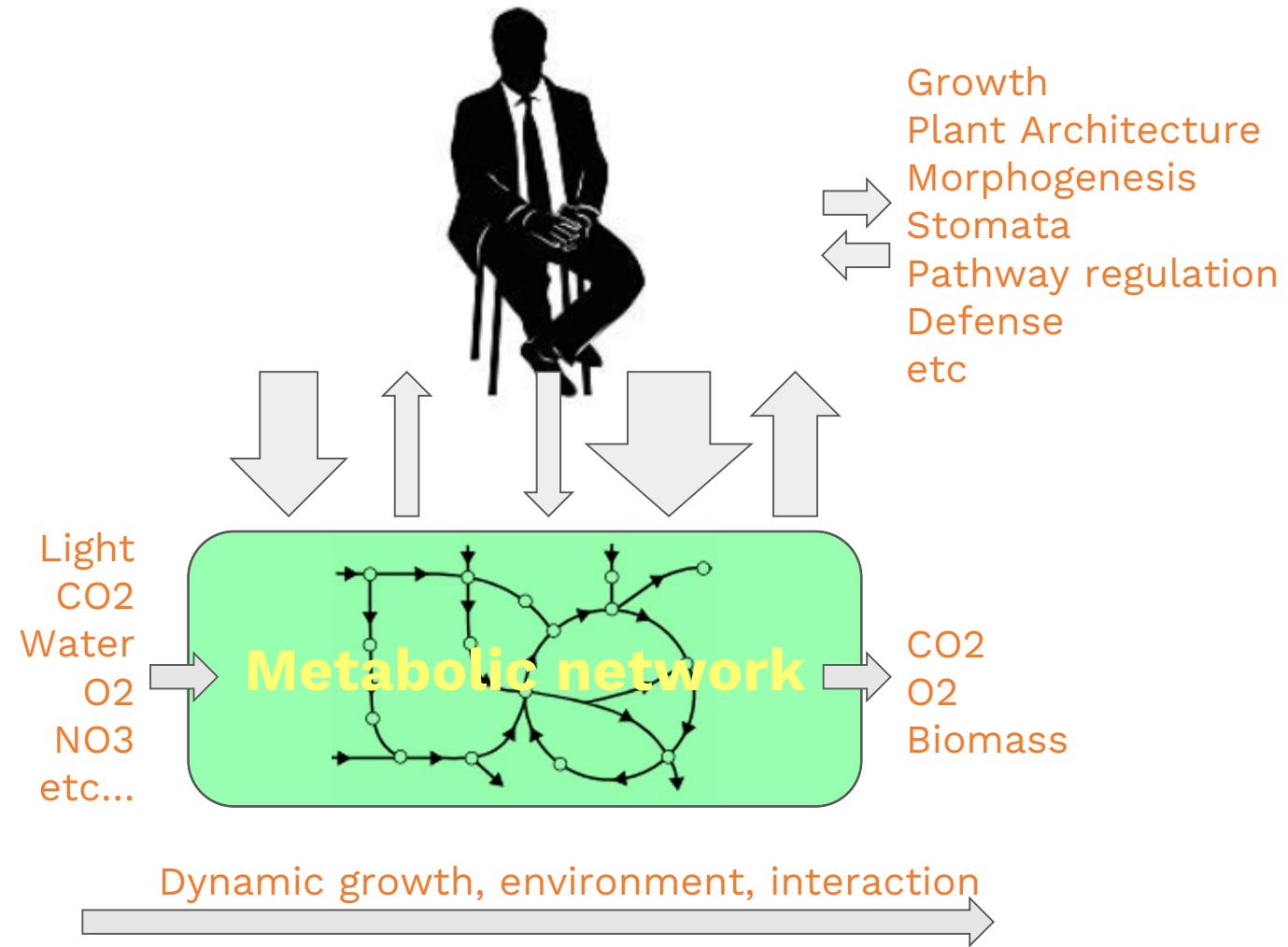
# Agent



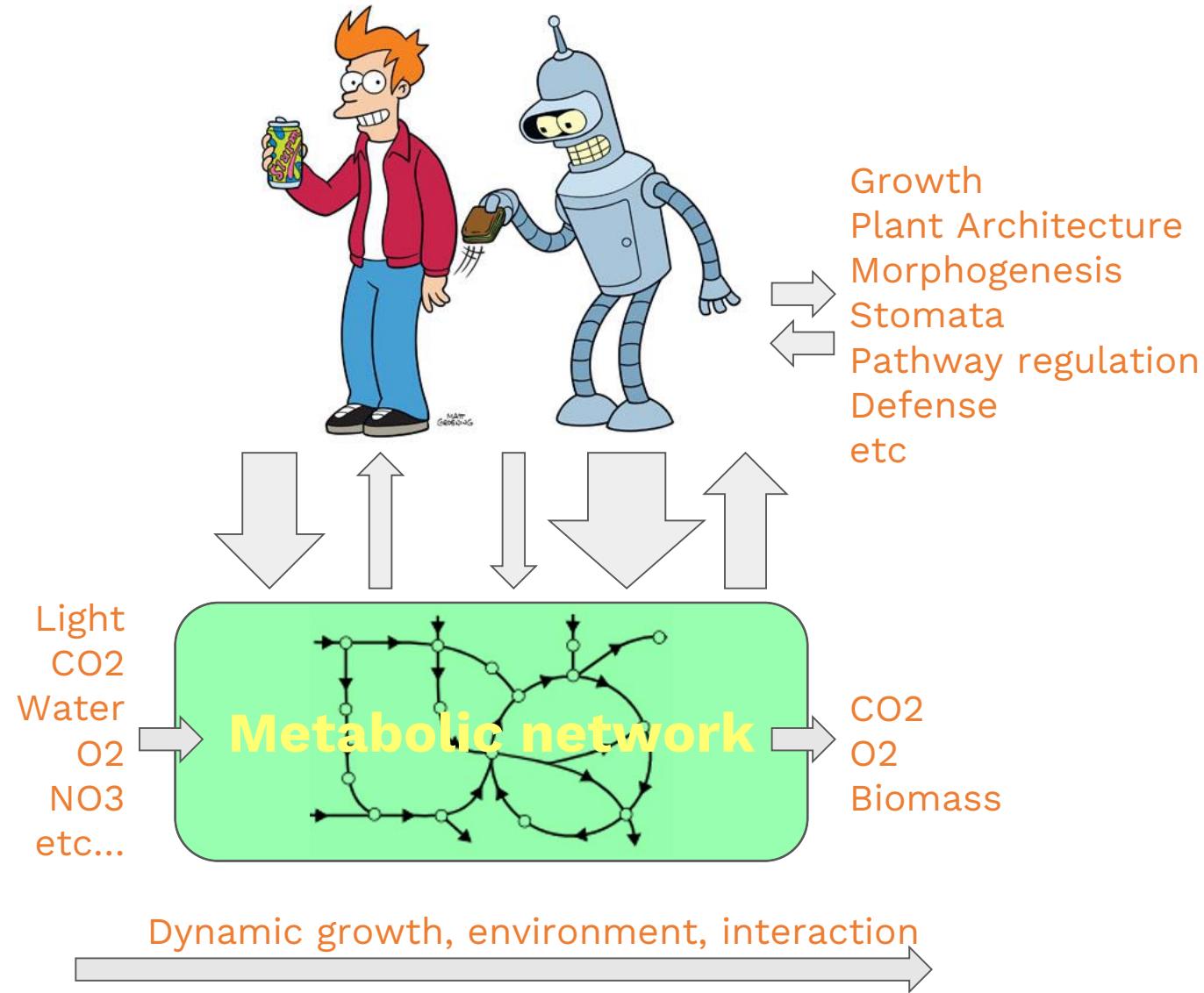
# Agent



# Agent

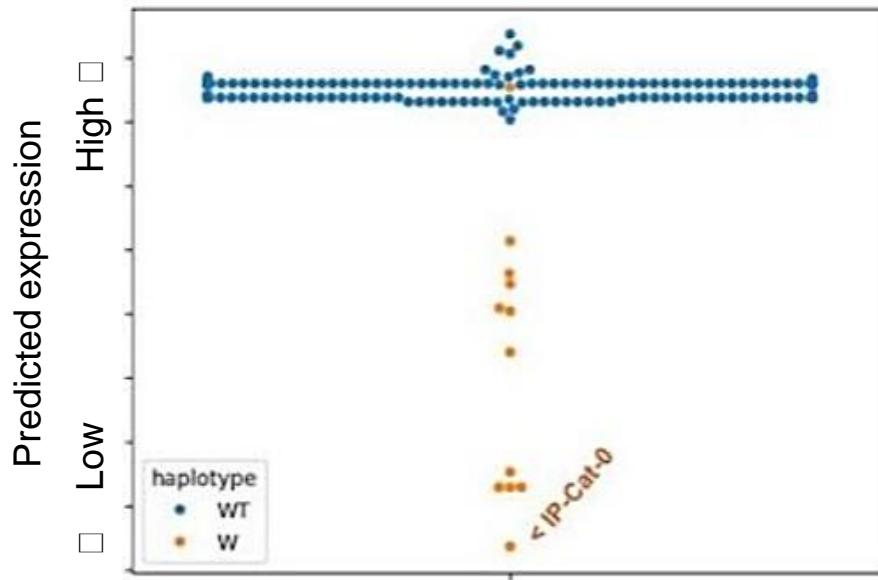


# Human AI



# Regulatory gene variants

SNP effects on expression of key regulatory genes



Example gene RAP2.2 in *A. thaliana*

Related to SUB1A in rice

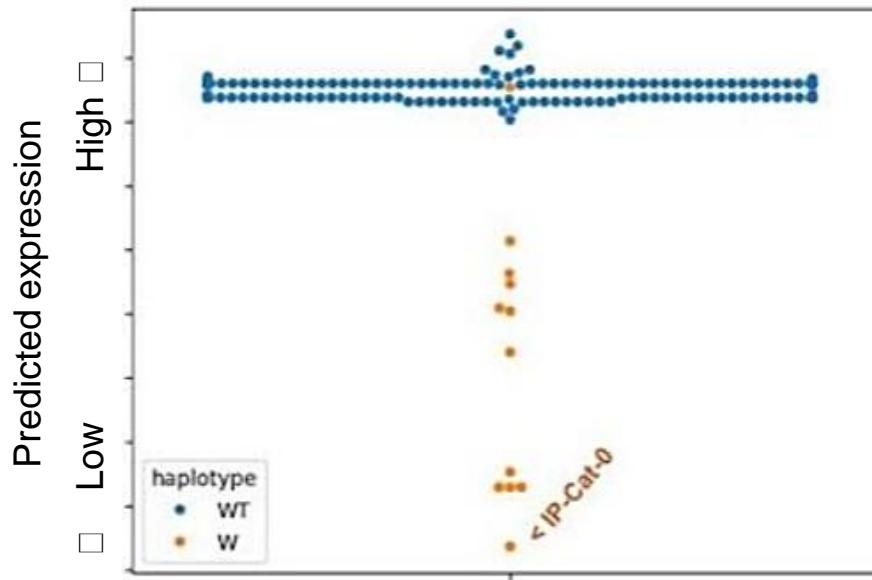
Responsible for survival during hypoxia / flooding



Gernot Schmitz  
(PhD Student)

## Regulatory gene variants

SNP effects on expression of key regulatory genes



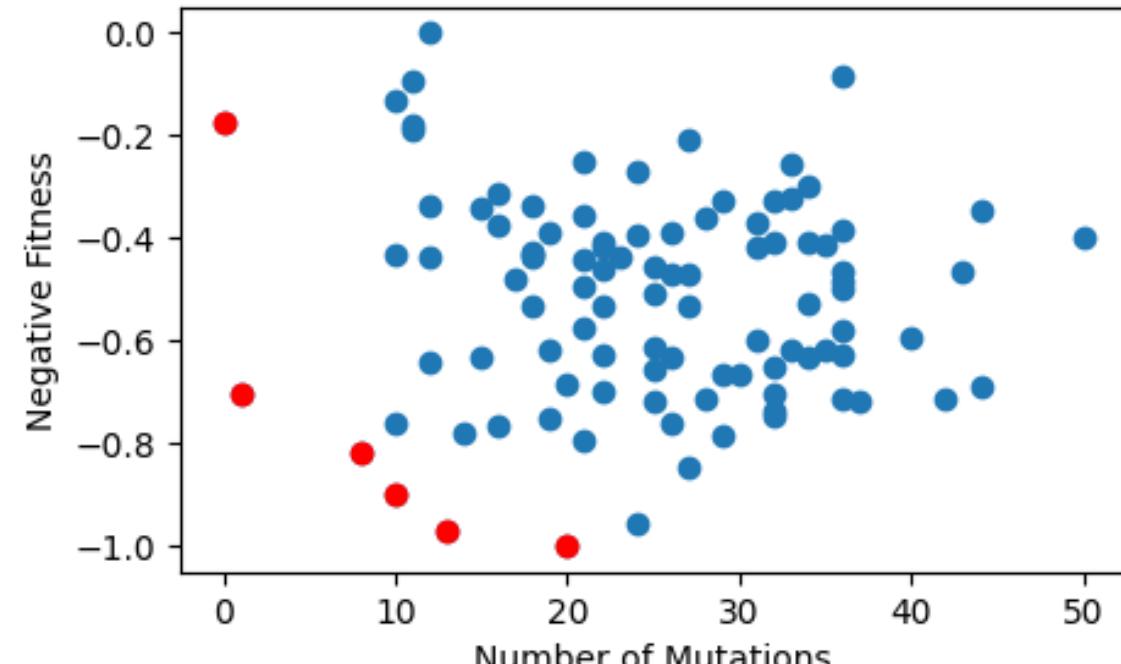
Example gene RAP2.2 in *A. thaliana*

Related to SUB1A in rice

Responsible for survival during hypoxia / flooding

## Promoter optimization

Smart gene editing for max effect with minimum mutations

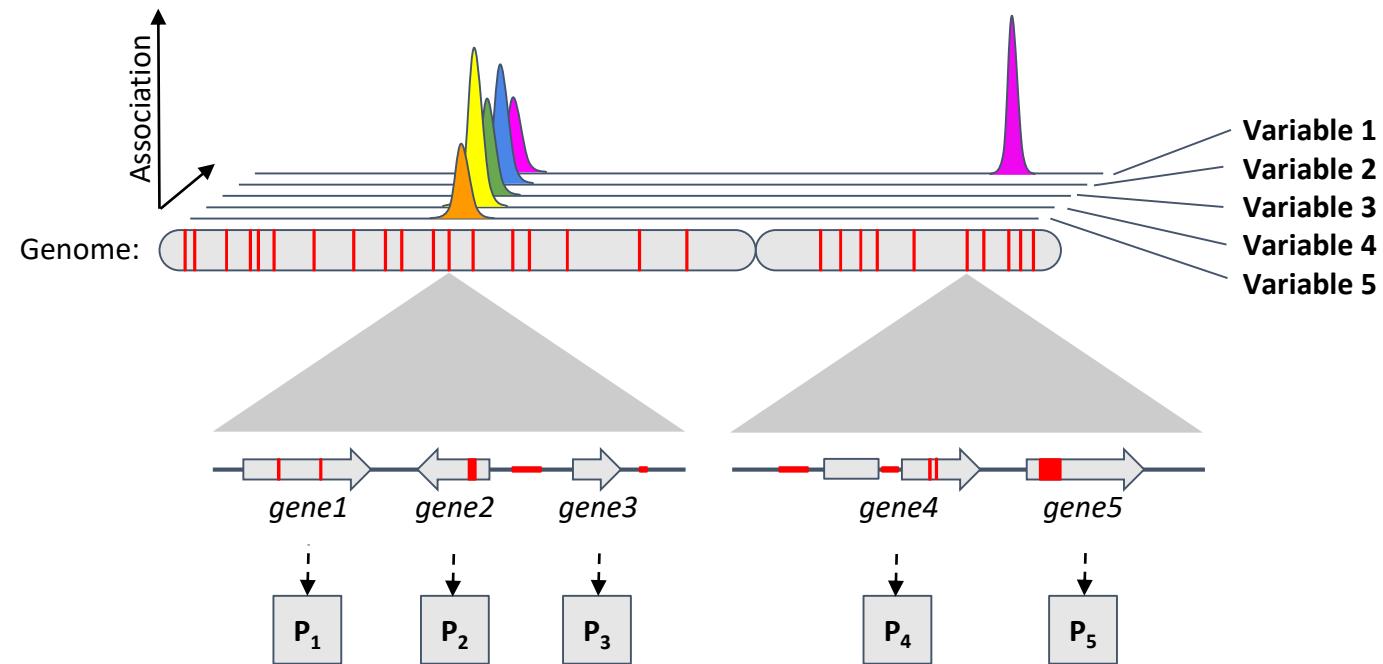
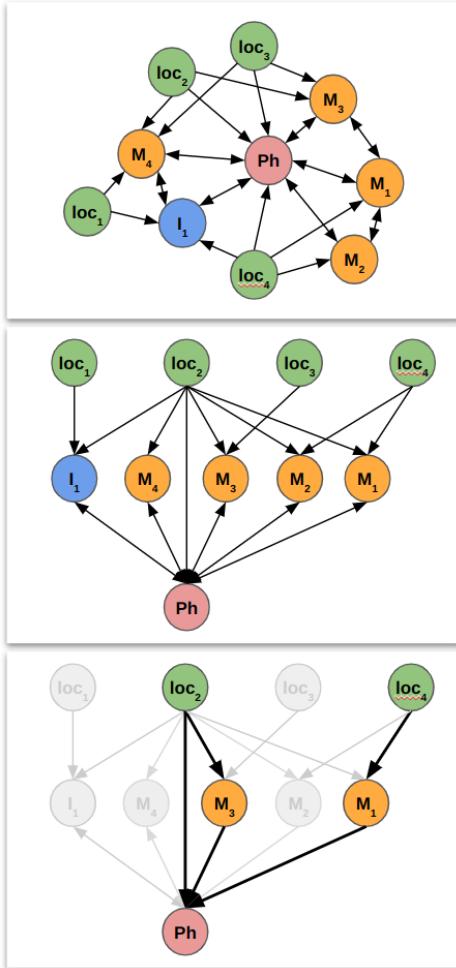


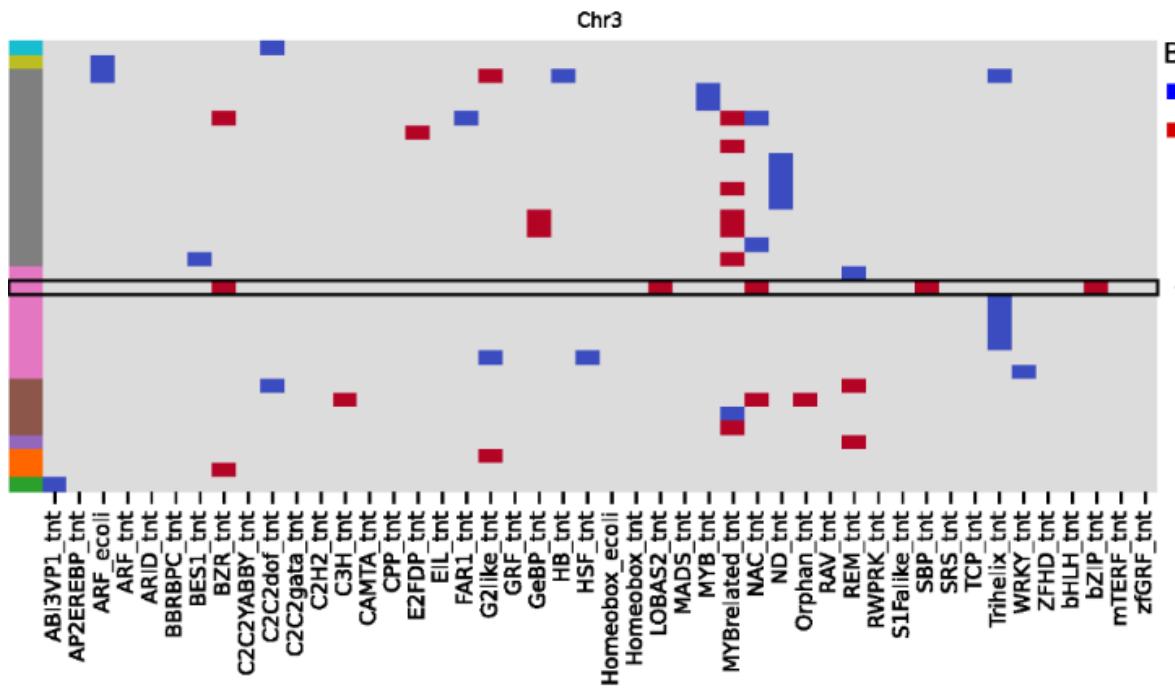
Gernot Schmitz  
(PhD Student)

# Genomic networking

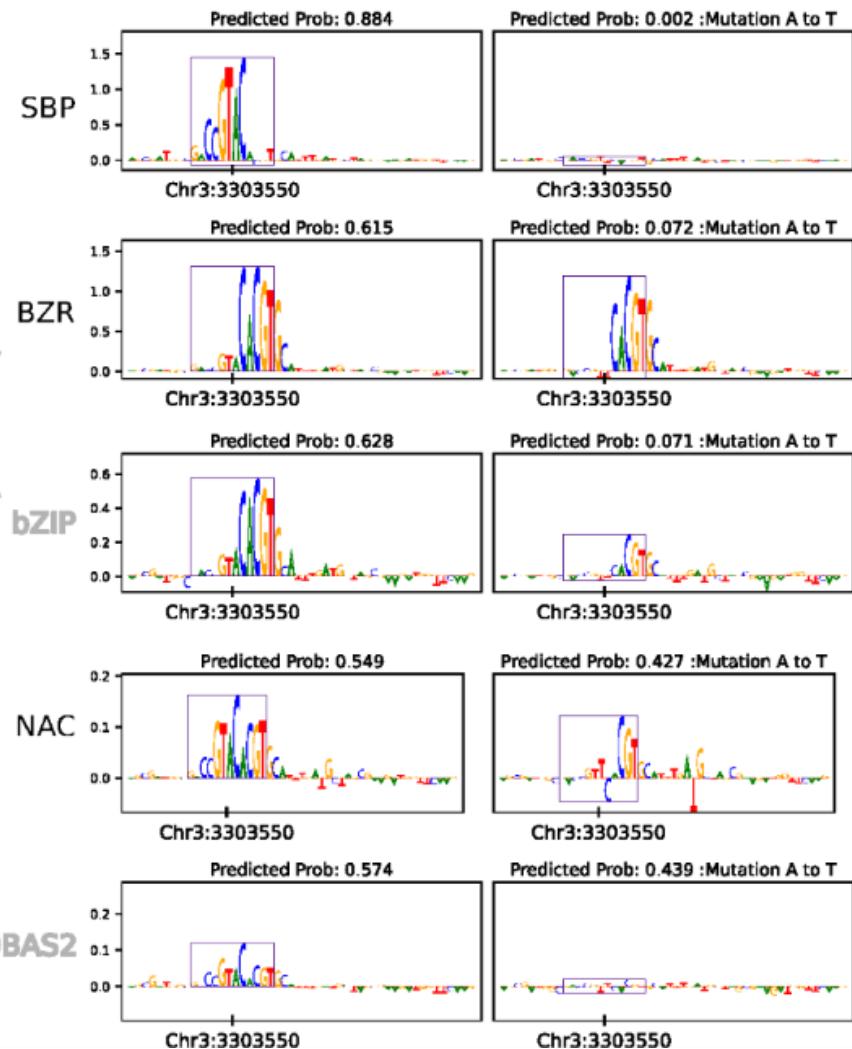
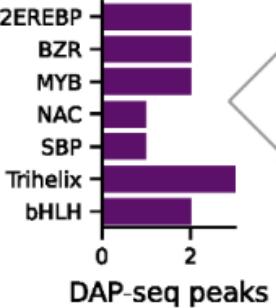
## Genomic networking

multivariate GWAS, path analysis, mixed graphical models





SBP  
BZR  
MYB  
NAC  
SBP  
Trihelix  
bHLH

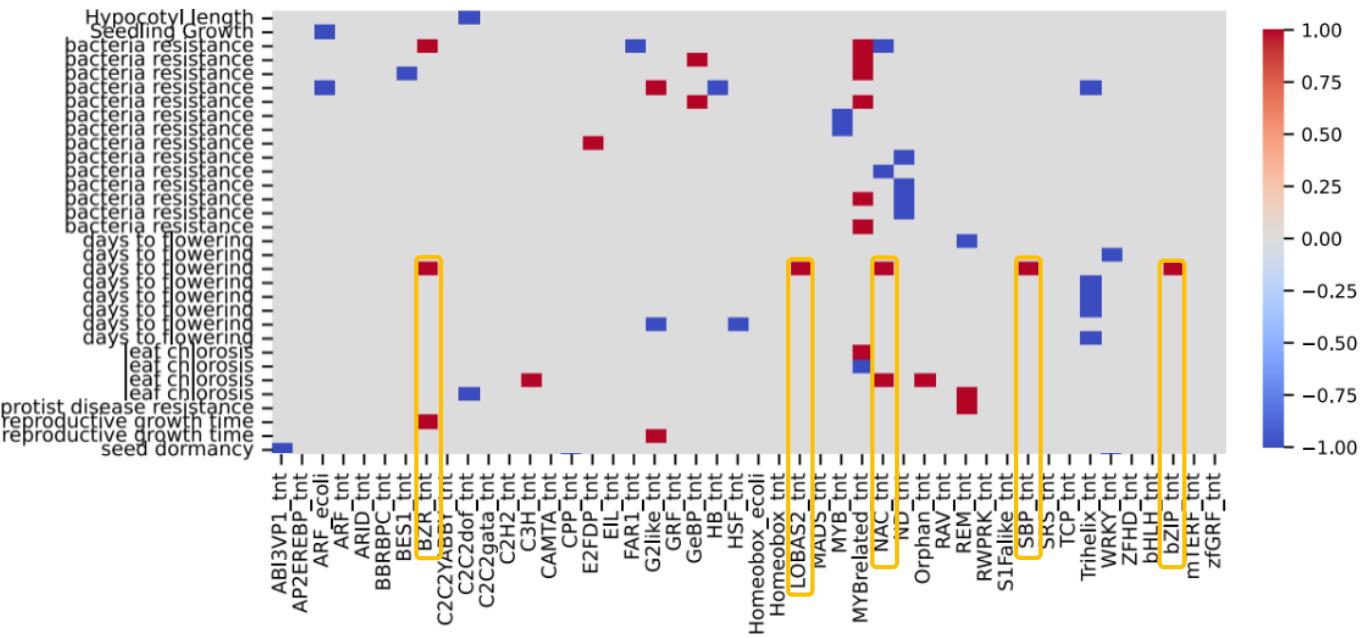
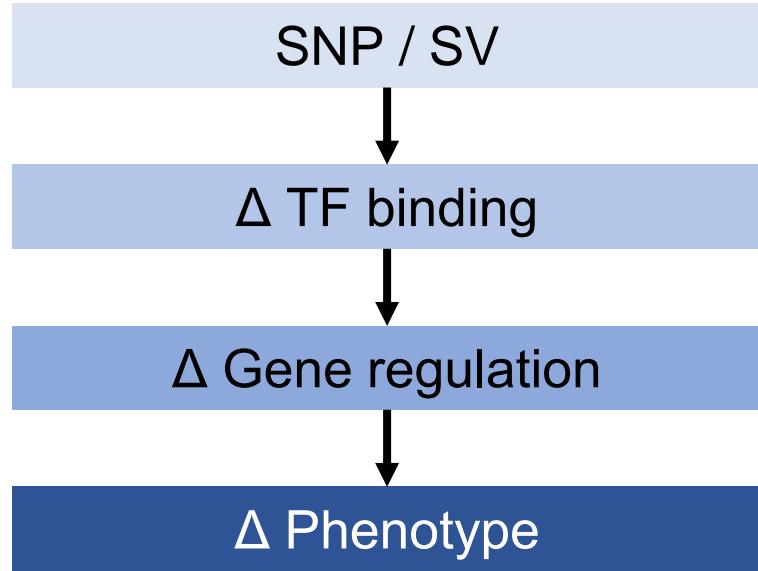




1001 Genomes

AraGWAS Catalog

36878 phenotype associations



### Combinatorial TF binding

