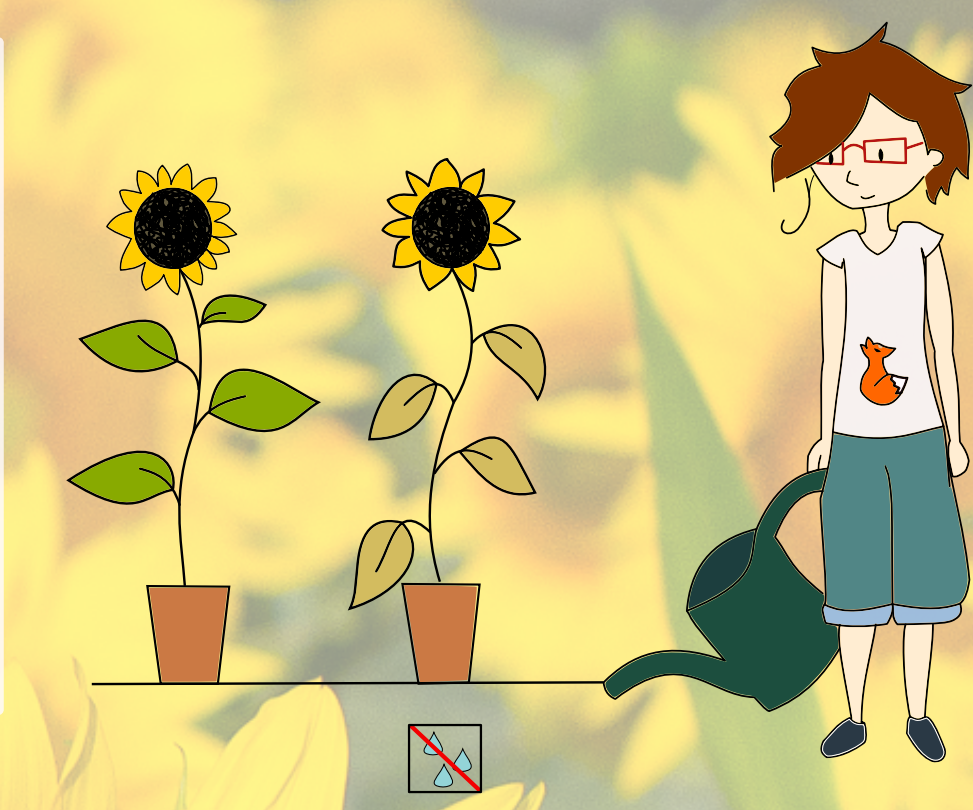


# Building artificial genetical genomic datasets to optimize the choice of gene regulatory inference methods

Lise Pomiès<sup>1</sup>, Louise Gody<sup>2</sup>, Charlotte Penouilh-Suzette<sup>2</sup>, Nicolas Langlade<sup>2</sup>, Brigitte Mangin<sup>2</sup>, Simon de Givry<sup>1</sup>

<sup>1</sup> Unité de Mathématiques et Informatique Appliquées de Toulouse  
MIAT INRA - UR875 - Chemin de Borde Rouge, 31320 Castanet Tolosan, France

<sup>2</sup> Laboratoire des interactions plantes micro-organismes  
LIPM INRA - UMR2594 - Chemin de Borde Rouge, 31320 Castanet Tolosan, France



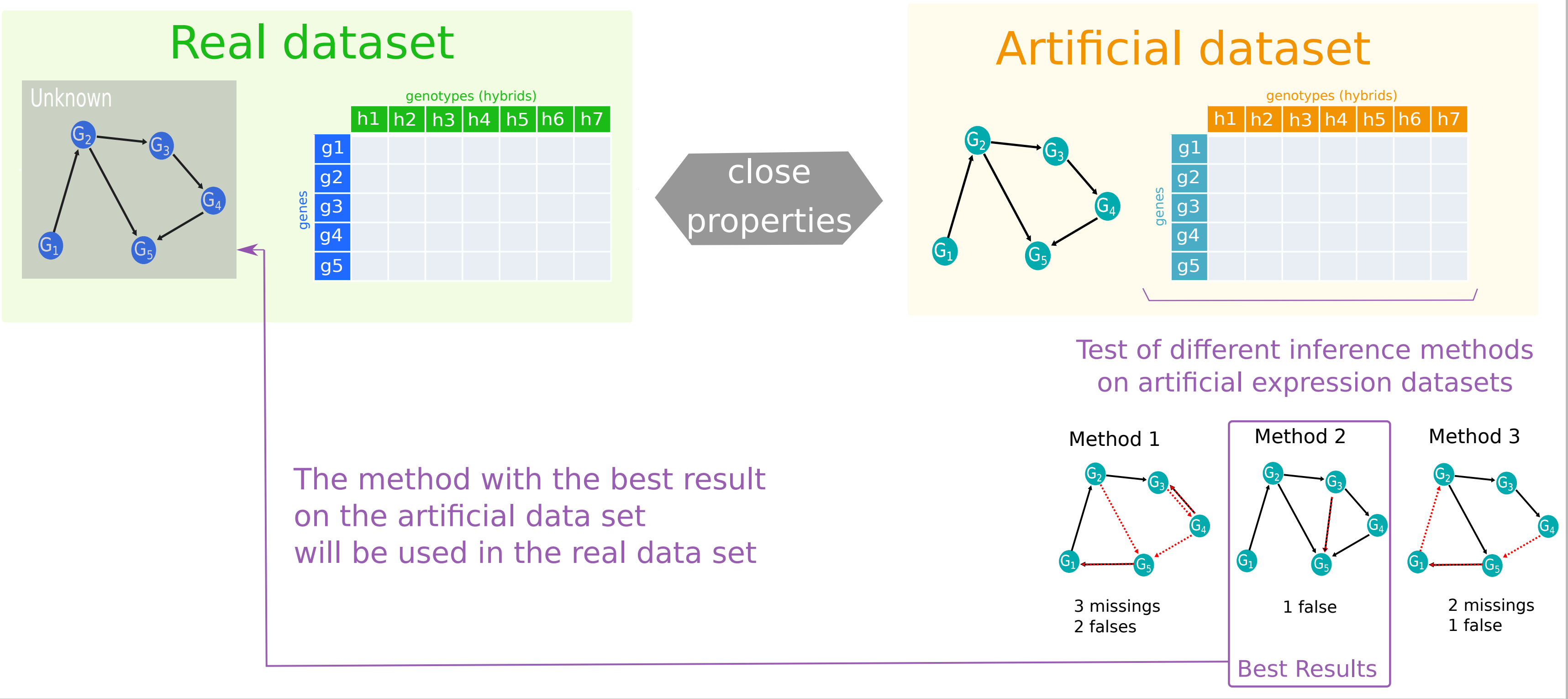
We study the **transcriptomic response of sunflower** to drought combined to the heterosis phenomenon.

**180 gene** expressions were measured in **400 genotypes**. Those genotypes are hybrids coming from a pool of 72 parents. SNP present on the parental genomes were measured.

We want to **infer the gene regulatory network** among those genes. However, because of the **non-independency** of the data, we don't know how inference will work. Therefore, we need to test different methods, to select the best inference method for our biological question.

## Select inference methods based on artificial datasets

**Test different methods** on an **artificial dataset** with known network, expression levels and genotypes (DNA variant). **Biological properties** of this artificial dataset



## How to build an artificial dataset with the same biological properties as our real one?

### 1. Build artificial network

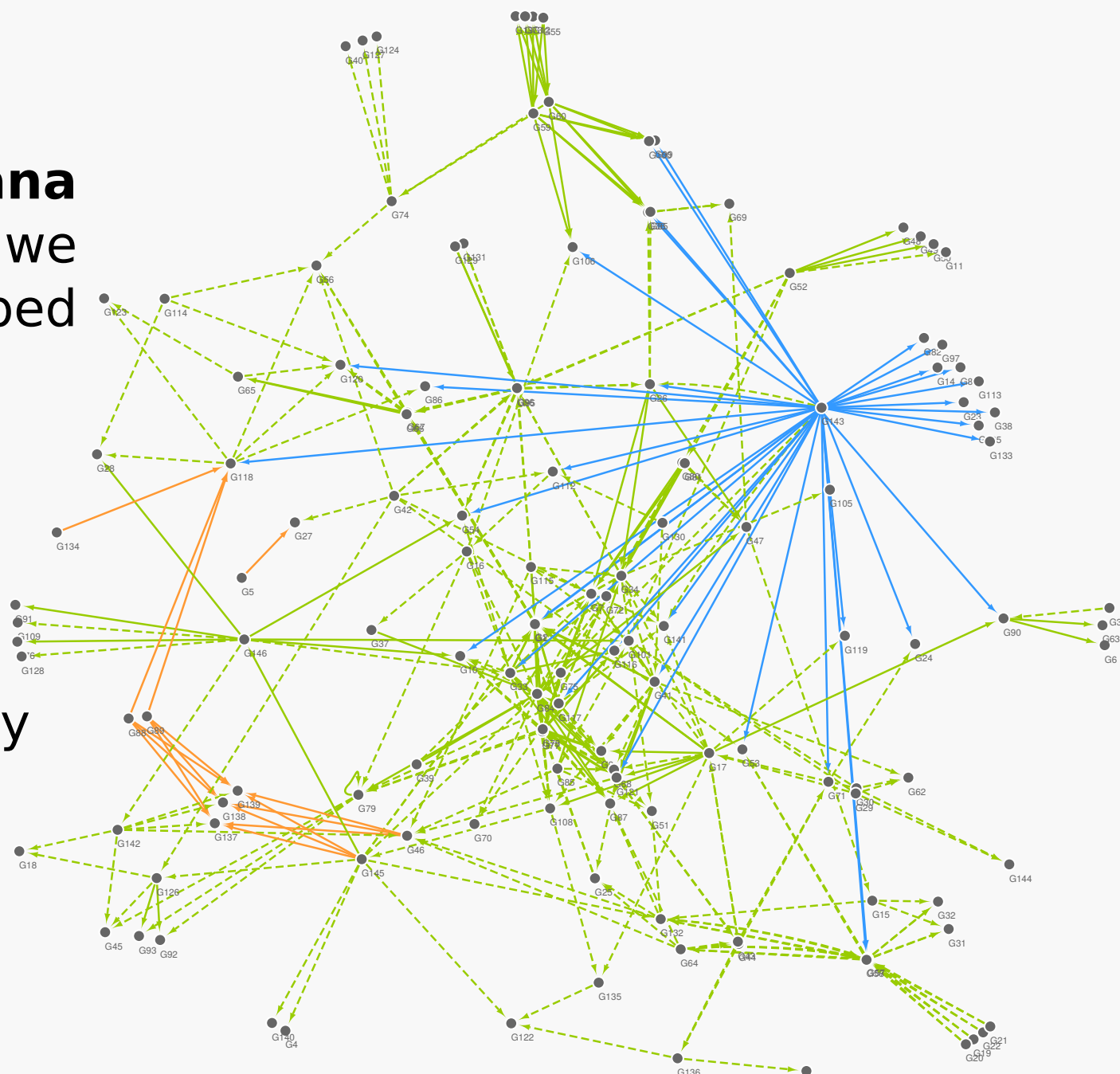
based on real biological information available for the same biological process on a close organism

For the **homologues** on **A. thaliana** of our 180 genes of interest, we selected **regulations** described between them in 3 databases.

**AtRegNet** [1]  
**AtPID** [2]  
**PlantRegMap** [3]

Regulations can be supported by experiments or predicted (----)

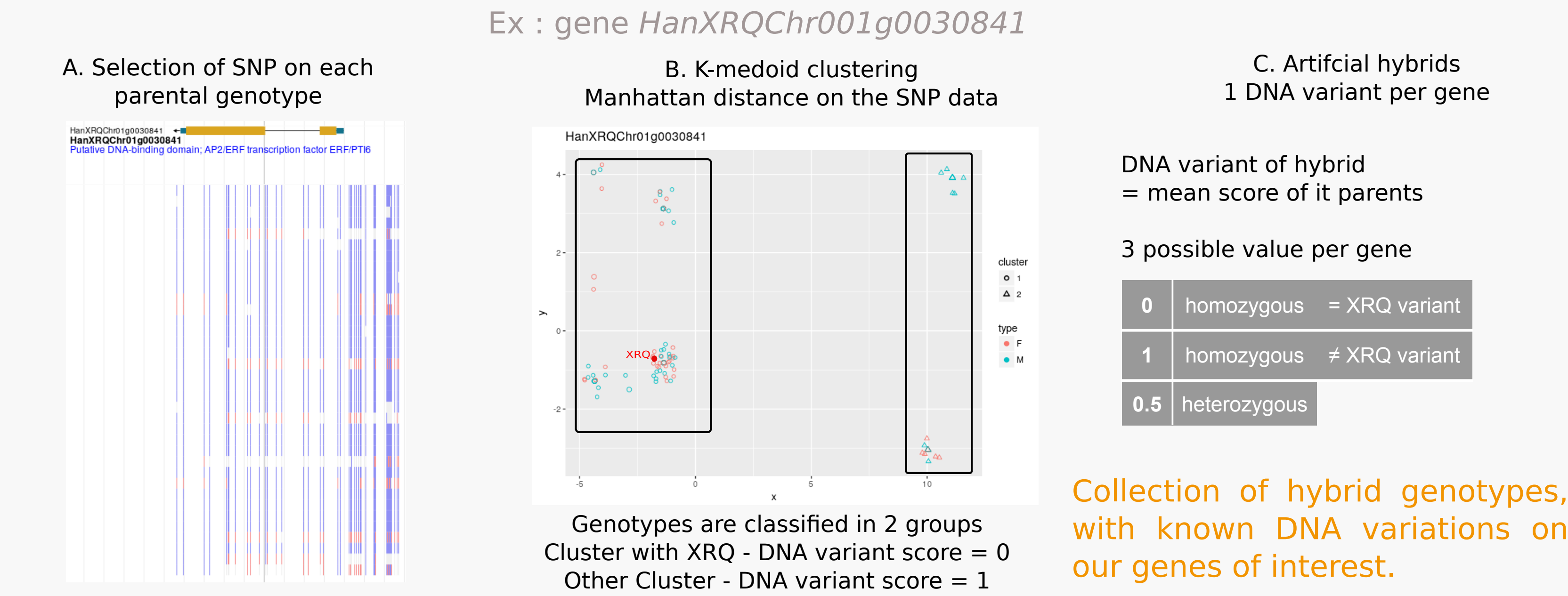
Our artificial network based on biological information is composed of 137 genes linked by 364 edges



### 2. Create artificial hybrid genotypes

based on genomic information available for the real hybrids used in the experiment

**SNP** for each **parental genotype** are associated to a score 0 if it is like in XRQ-line or 1 if different. The hybrids SNP are obtained by combining locus-per-locus SNP of their parents. We created **artificial hybrids** associated to **DNA variant** on each measured gene. We considered **one variant per gene**, those DNA variants are based on SNP of the real data.



### 3. Select and adapt an existing gene expression simulator

emulating the same type of experiment that the one we performed, with steady state measurements on different genotypes

**SysGenSIM** simulates steady state gene expressions using ordinary differential equations. Simulation is based on a gene network topology and DNA variant for each gene. Work only on RIL (both allele of a gene are identical) [4]

$$\frac{dG_g}{dt} = Z_g^c \cdot V_g \cdot \theta_g^{syn} \cdot \prod_k \left( 1 - A_{kg} \frac{G_k^{h_{kg}}}{G_k + (K_{kg}/Z_k^c)^{h_{kg}}} \right) - \lambda_g \cdot \theta_g^{deg} \cdot G_g$$

basal transcription rate of gene g  
SNP(c) noise  
effect of gene g SNP(c) on gene g expression  
role of gene k on gene g  $\in \{-1, 0, 1\}$   
[mRNA] of gene k  
degradation rate constant of gene g  
degradation noise  
min([mRNA] of gene k for k to have an effect on gene g)  
effect of SNP(t) of k on its activity (more or less efficient regulator)  
[mRNA] of gene g

We **modified the simulator** to use our **heterogeneous hybrids**, and mimetized the allelic dominance caused by the **heterosis phenomenon**.

SysGenSIM  
Z parameter 2 possible values

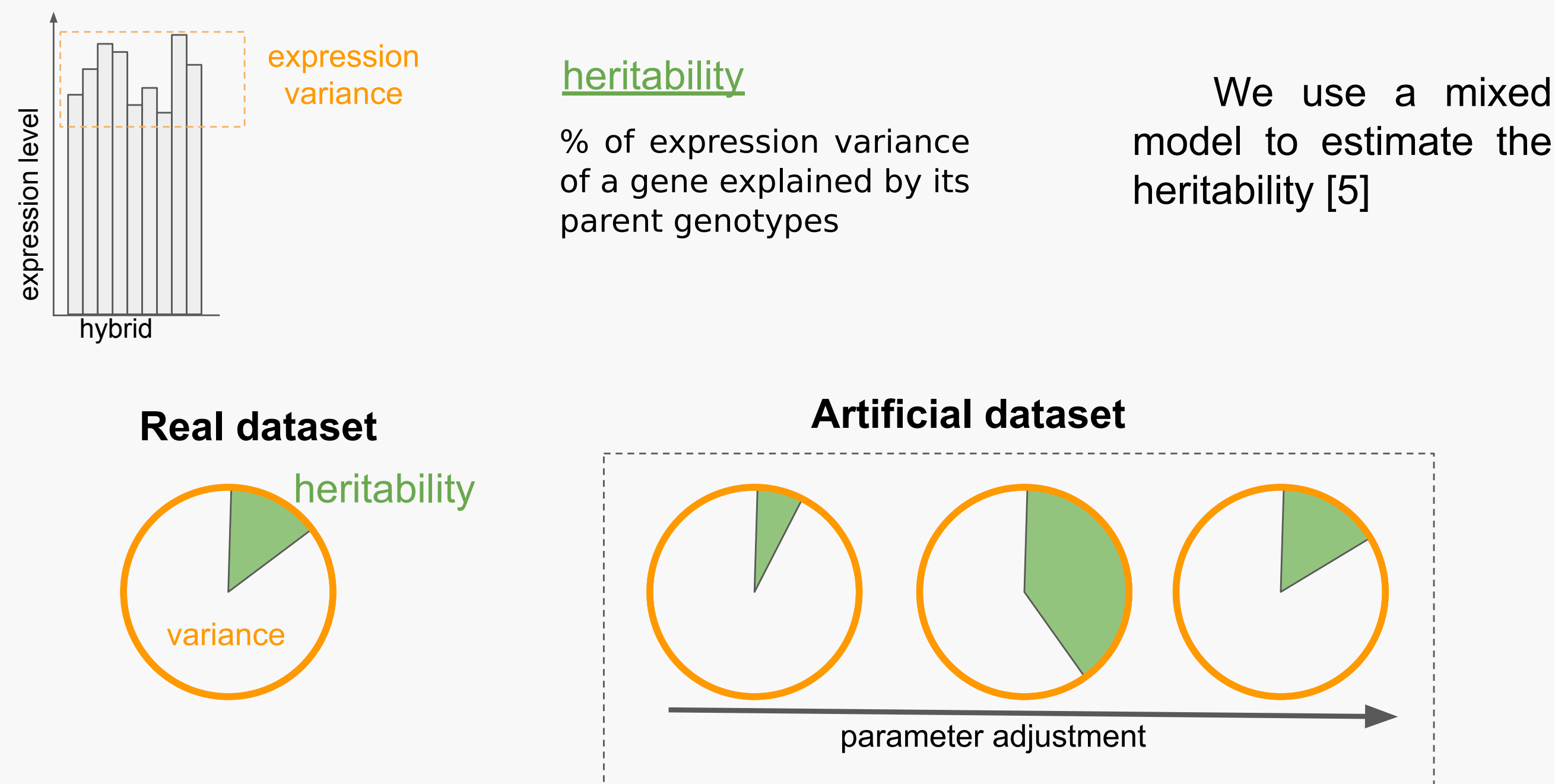
wt-wt	1
m-m	0.75

Modified SysGenSIM  
Z parameter 3 possible values

wt - wt	1
m - m	0.75
wt - m	0.75 mutated dominance (10%) 0.87 additif effect (80%) 1 wt dominance (10%)

### 4. Comparison of biological score

obtained on real and simulated datasets (in our case the heritability score) to adjust parameters of the simulator



**The artificial dataset produced have the same biological properties as our real dataset.** We can now test different methods of network inference and test the accuracy of these methods by comparing networks inferred by the algorithms to the artificial network. Network inference methods with the best results will be used on the experimental dataset to answer our biological question.