



HAL
open science

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

Lise Pomies, Louise Gody, Charlotte Penouilh-Suzette, Nicolas Langlade,
Brigitte Mangin, Simon de Givry

► To cite this version:

Lise Pomies, Louise Gody, Charlotte Penouilh-Suzette, Nicolas Langlade, Brigitte Mangin, et al.. Building artificial genetical genomic datasets to optimize the choice of gene regulatory network inference methods. 17th European Conference on Computational Biology (ECCB 2018), Sep 2018, Athènes, Greece. 2018. hal-02734486

HAL Id: hal-02734486

<https://hal.inrae.fr/hal-02734486>

Submitted on 2 Jun 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

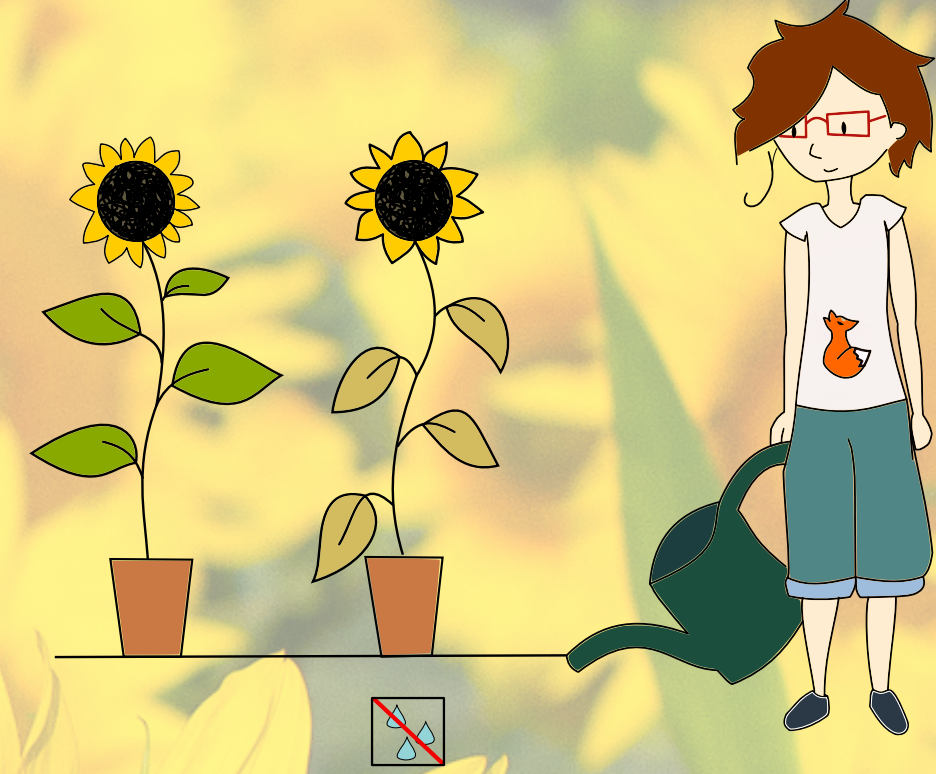
L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

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

Lise Pomiès¹, Louise Gody², Charlotte Penouilh-Suzette²,
Nicolas Langlade², Brigitte Mangin², Simon de Givry¹

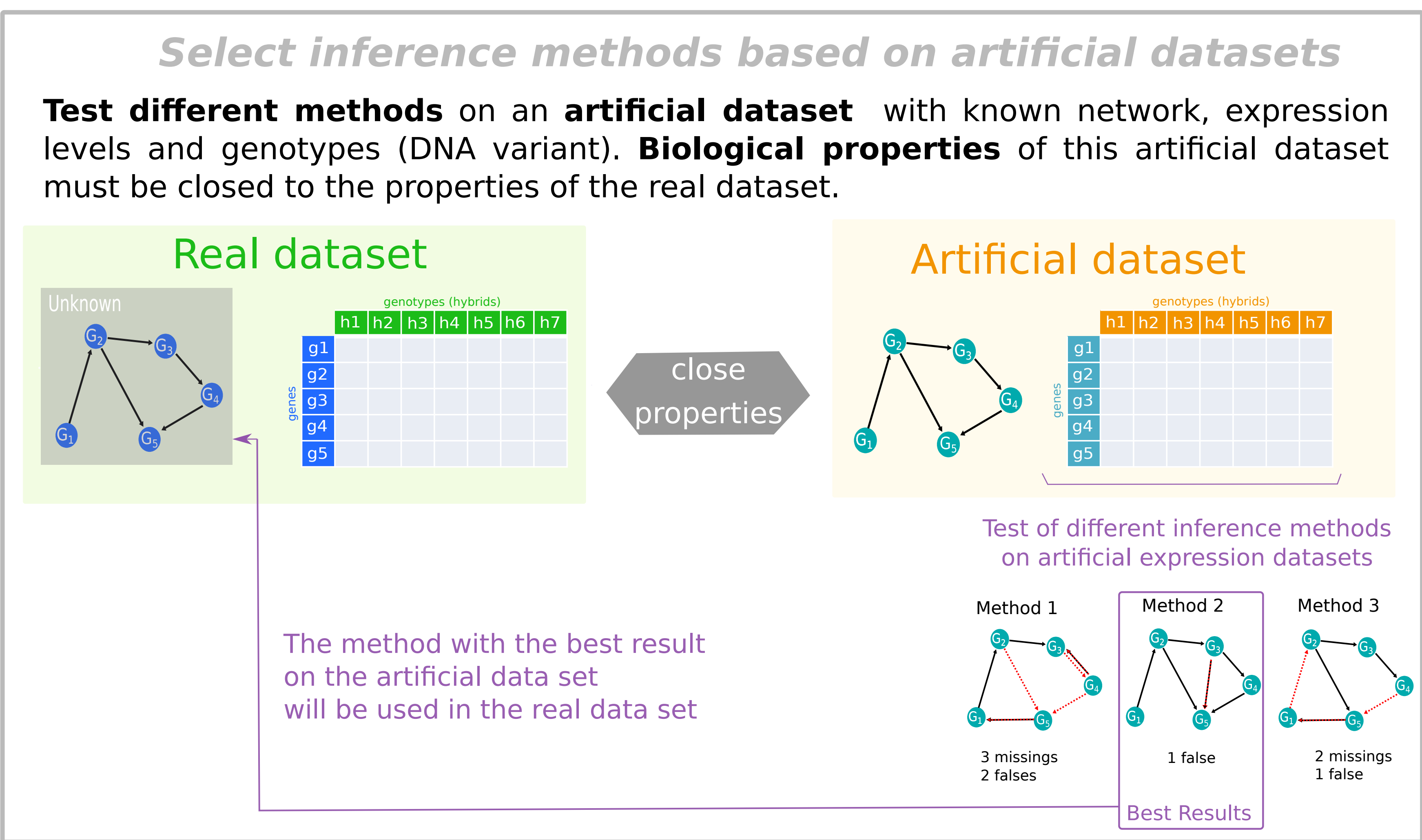
¹ Unité de Mathématiques et Informatique Appliquées de Toulouse
MIAT INRA - UR875 - Chemin de Borde Rouge, 31320 Castanet Tolosan, France

² 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, across **180 gene** expressions on **400 hybrids genotypes**, coming from a pool of 72 parents. SNP present on the parental genomes were measured.

Our goal is to **infer the gene regulatory network** among those genes. However, because of the **non-independency** of the data, accuracy of inference results is unpredictable. Therefore, we need to test different methods, to select the best inference method for our biological question.



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.

Ex : gene *HanXRQChr001g0030841*

DNA variant of hybrid = mean score of its parents

3 possible values per gene

0	homozygous = XRQ variant
1	homozygous ≠ XRQ variant
0.5	heterozygous

Collection of hybrid genotypes, with known DNA variations on our genes of interest.

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]

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

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

SysGenSIM		Modified SysGenSIM	
Z parameter 2 possible values		Z parameter 3 possible values	
wt-wt	1	wt - wt	1
m-m	0.75	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

expression variance (bar chart for hybrid)

heritability (% of expression variance of a gene explained by its parent genotypes)

We use a mixed model to estimate the heritability [5]

Artificial dataset
SysGenSIM parameters adjusted to obtain the same heritability distribution

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.