

# **Coexpression and Collaborative Gene Networks**

Hairong Wei

Michigan Technological University

# References

- 1) Persson S. **H. Wei** J. Milne G. Page C. Somerville. 2005. Identification of genes required for cellulose synthesis by regression analysis of public microarray data sets. *Proc Natl Acad Sci USA* 102: 8633-8638. (Faculty 1000 evaluation)
- 2) **Wei H** S. Persson T. Mehta V. Srinivasasainagendra L. Chen G. Page C. Somerville A. Loraine. 2006. Transcriptional coordination of the metabolic network in *Arabidopsis thaliana*. *Plant Physiology*. 142(2):762-74.
- 3) Nie J. R. Stewart F. Ruan J. Thomson H. Zhang X. Cui and **H. Wei**. 2011. TF-Cluster: a pipeline for identifying functionally coordinated transcription factors via network decomposition of the shared coexpression connectivity matrix (SCCM). *BMC Systems Biology* 5:53.
- 4) Ji X. S. Chen J. Li W. Deng Z. Wei and **H. Wei**. 2017. SSGA and MSGA: two seed-growing algorithms for constructing collaborative subnetworks. *Scientific Reports*. (in press)
- 5) Kumari, S., J. Nie,, H.S. Chen, R. Stewart, H. Ma, X. Li, M. Lu, W.M. Taylor, and H. Wei. 2012. Evaluation of gene association methods for coexpression network construction and biological knowledge discovery. *PLoS ONE* 7(11): e50411

# What is coexpression network ?

A **gene co-expression network (GCN)** is an undirected graph where each node corresponds to a gene and a pair of nodes is connected with an edge if there is a significant co-expression relationship between them.

# Methods for associating co-expressed genes

1. Spearman rank correlation
2. Weighted Rank Correlation
3. Kendall rank correlation
4. Hoeffding's D measure
5. Theil-Sen
6. Rank Theil-Sen
7. Distance Covariance
8. Pearson

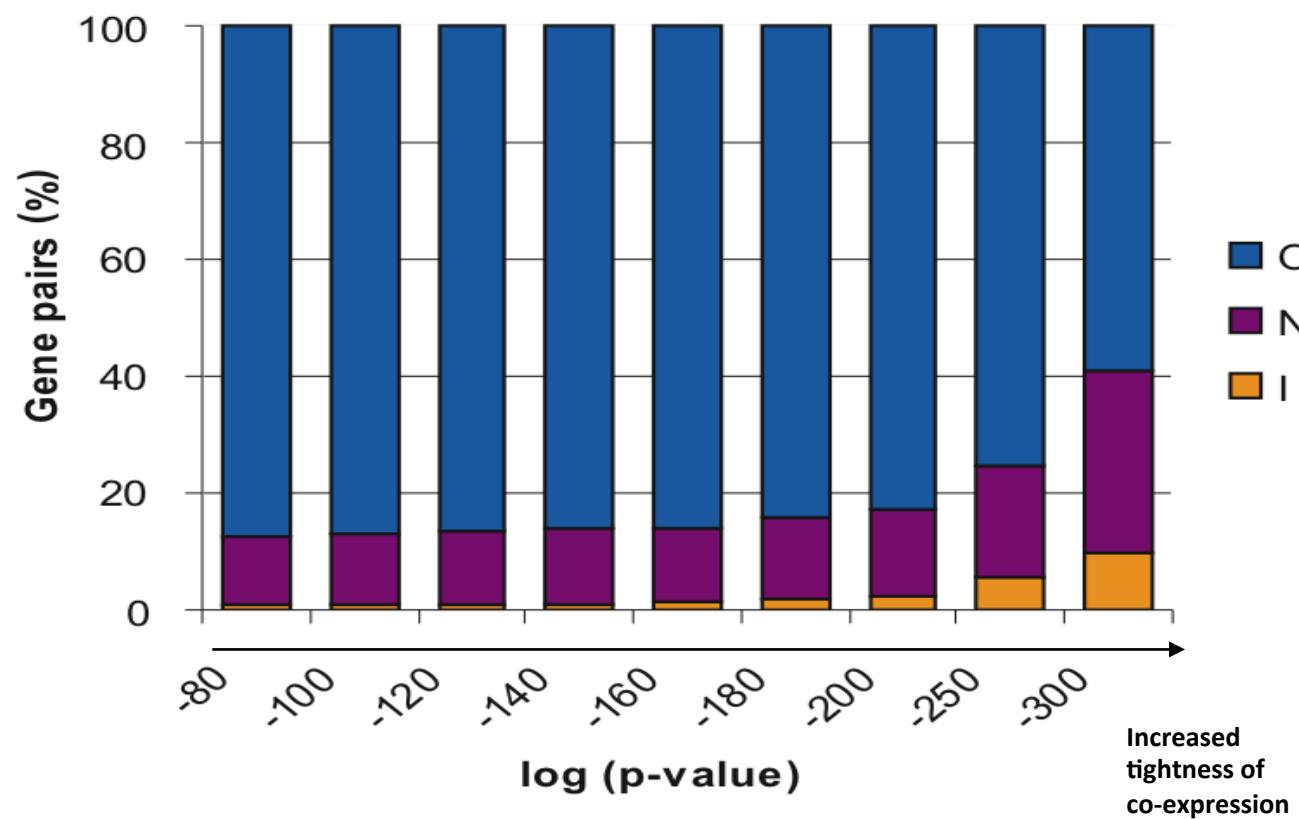
Kumari et al 2012, PLoS ONE. The R codes for above gene associate methods can be found at:

<http://journals.plos.org/plosone/article/file?type=supplementary&id=info:doi/10.1371/journal.pone.0050411.s005>

Supplemental Doc S1

# Are Genes involved in the same biological pathway coexpressed?

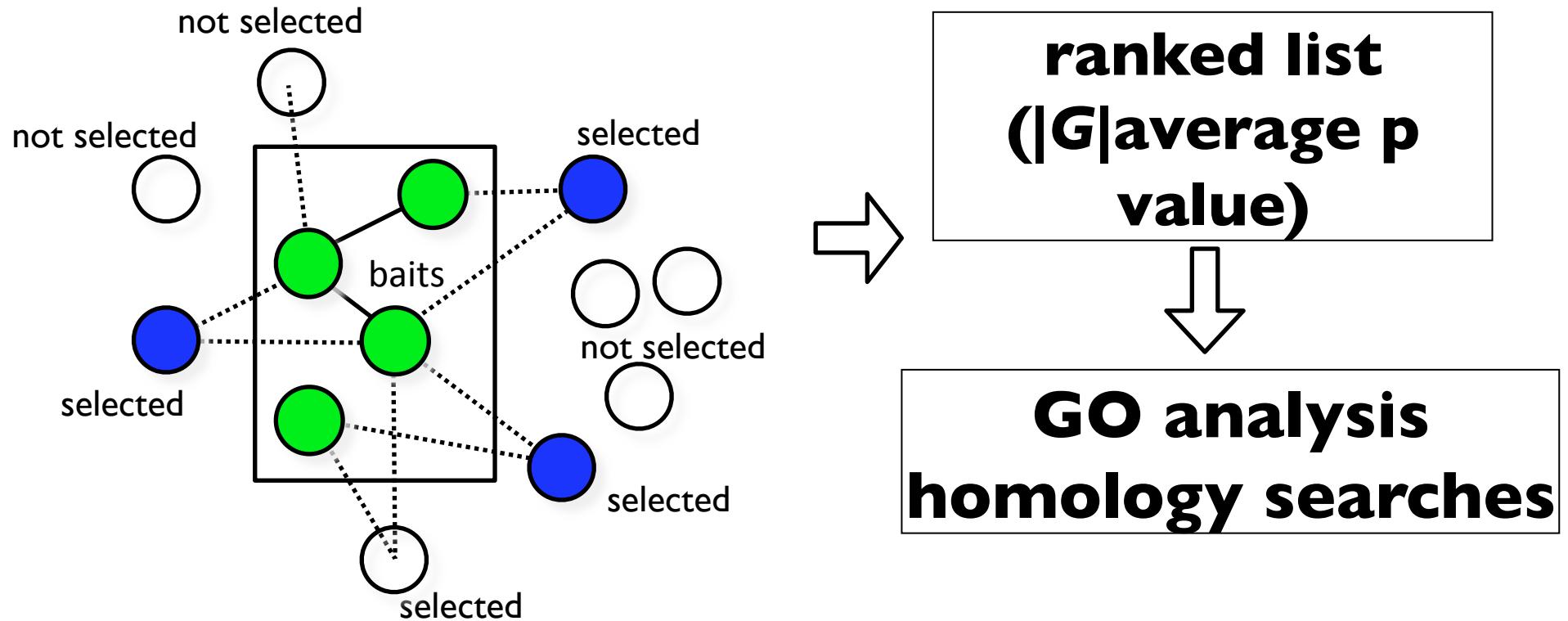
Hypothesis: The genes having functional association e.g. in the same pathway are more tightly co-expressed.



216 biological pathways that contain 1330 genes in Arabidopsis

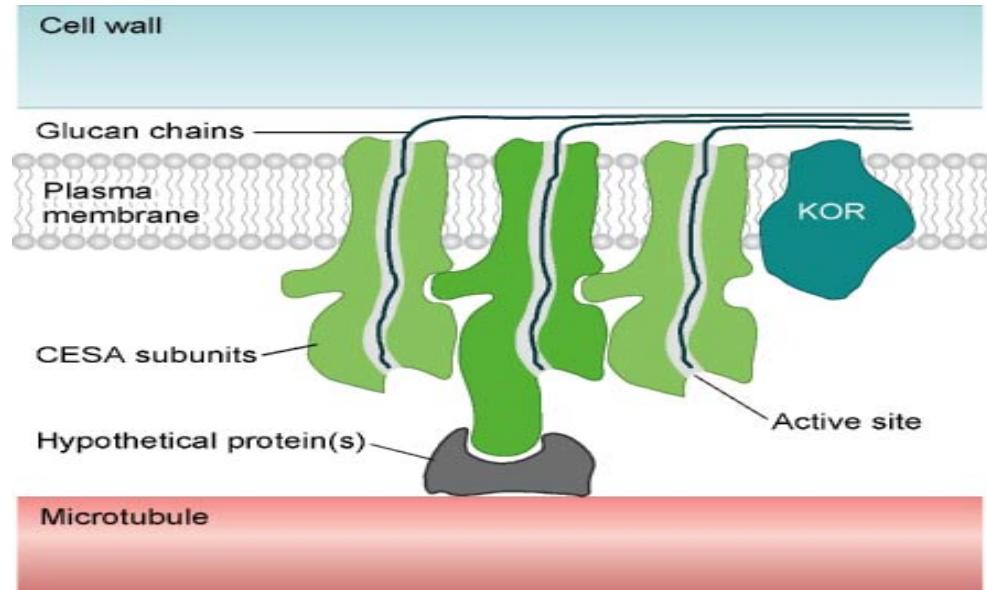
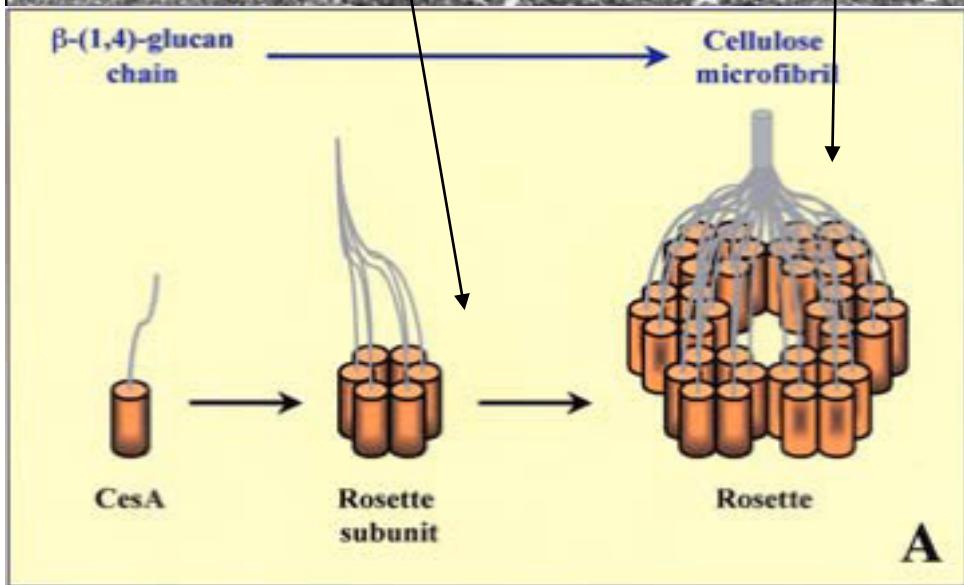
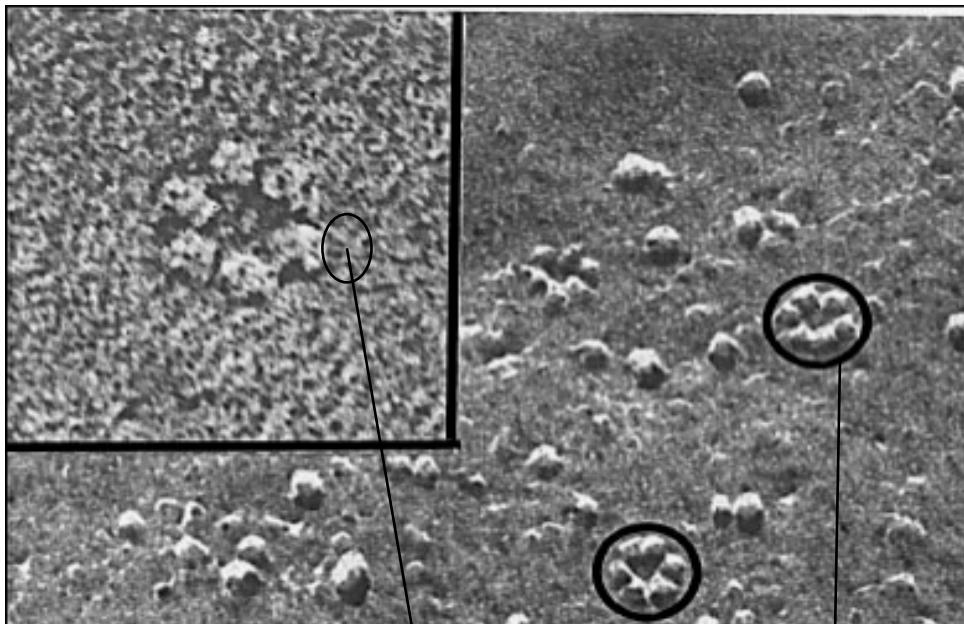
O-outside  
N-different pathway  
I-same pathway

# Intersection of Co-Expression (ICE)



1. Define co-expression network.
2. Select a group of “bait” genes. (e.g. a pathway a complex). Green nodes
3. Perform genome-wide co-expression analysis (correlation or regression).
3. Select genes outside the group of baits that are linked with >1 group members G.
4. Sort them by  $|G|$  then average p value (or  $r^2$ )

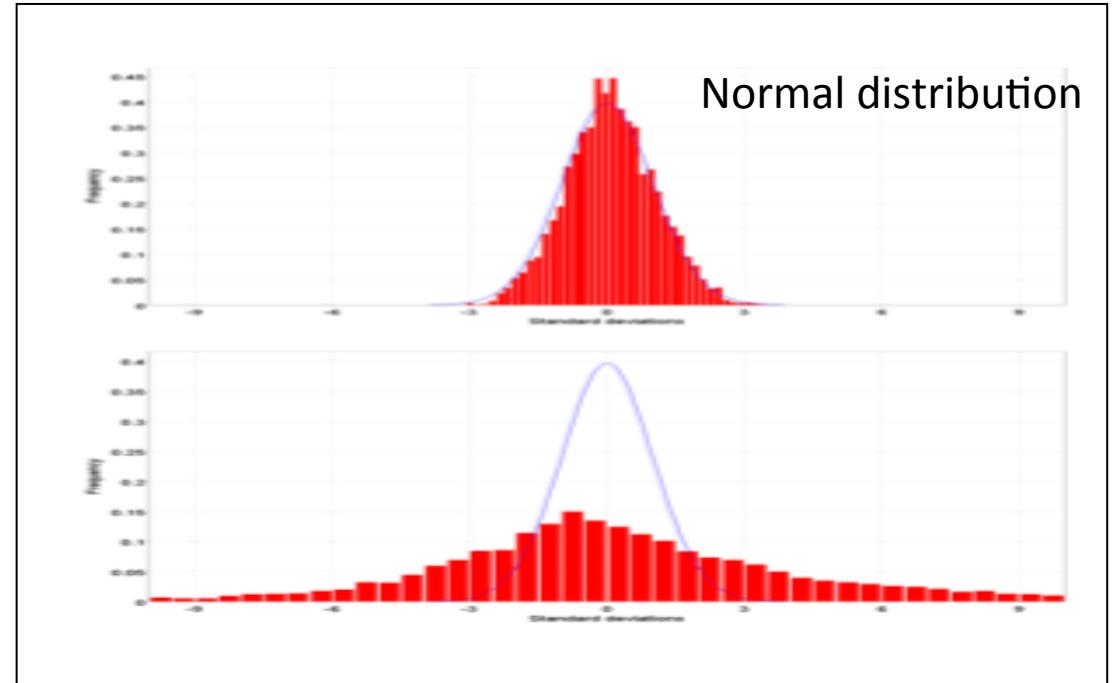
# ICE Implementation 1: finding functionally associated genes



**A very challenging problem:**

What genes in addition to  
CEAS genes ( 4, 7 and 8) are required  
for secondary cell wall synthesis?

**Primary: CESA136      Secondary: CESA4 7 8**



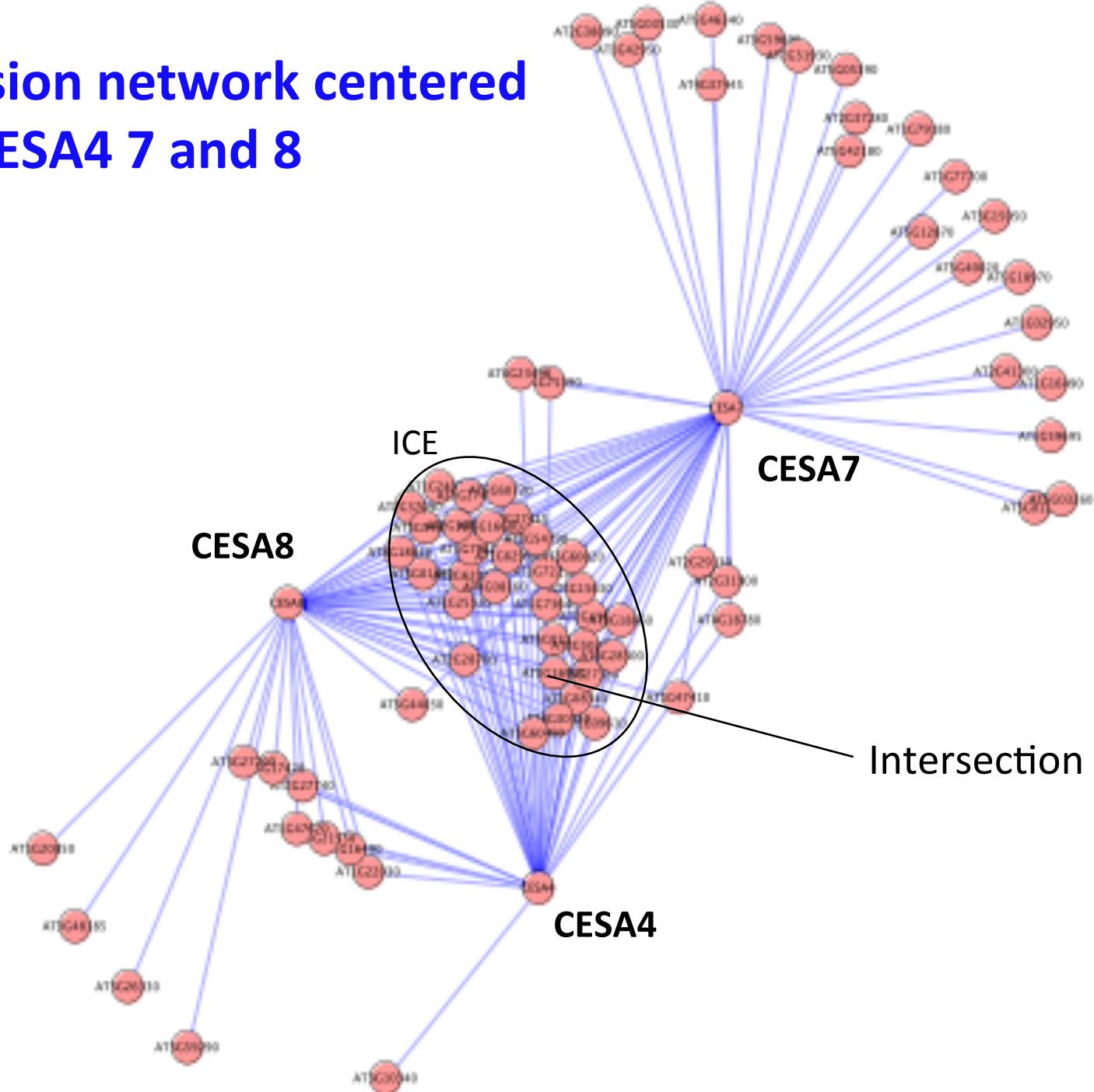
For n genes and m data sets  
deleted residual

$$d_{ij} = X_{ij} - \bar{X}_{i-j},$$

$$d_i^* = d_i / s(d_i)$$

obeys t distribution with m-2 df

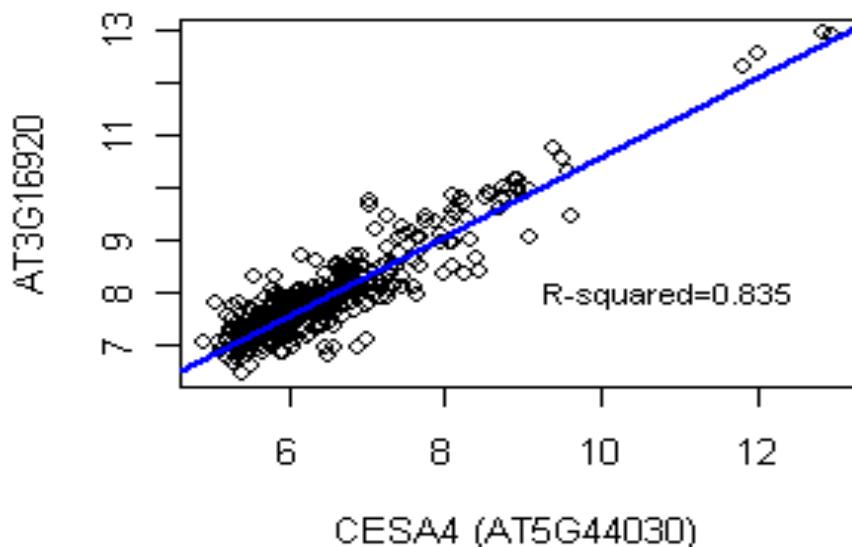
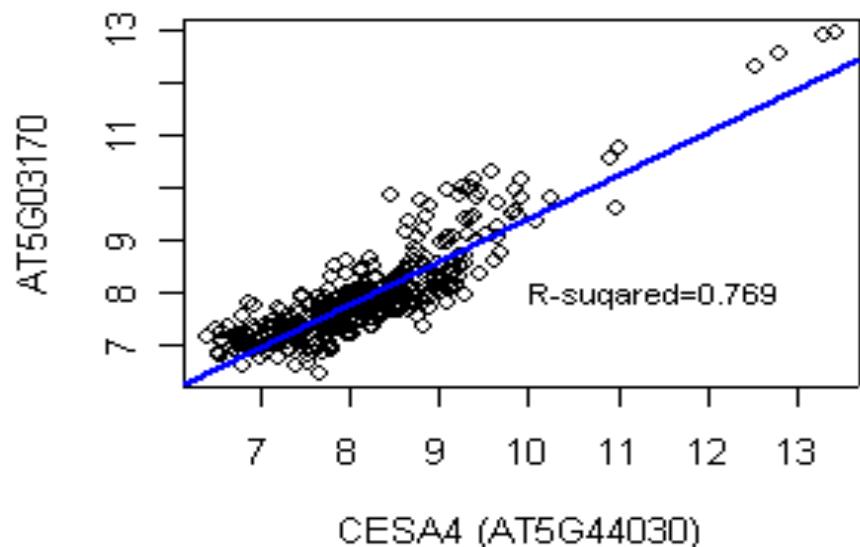
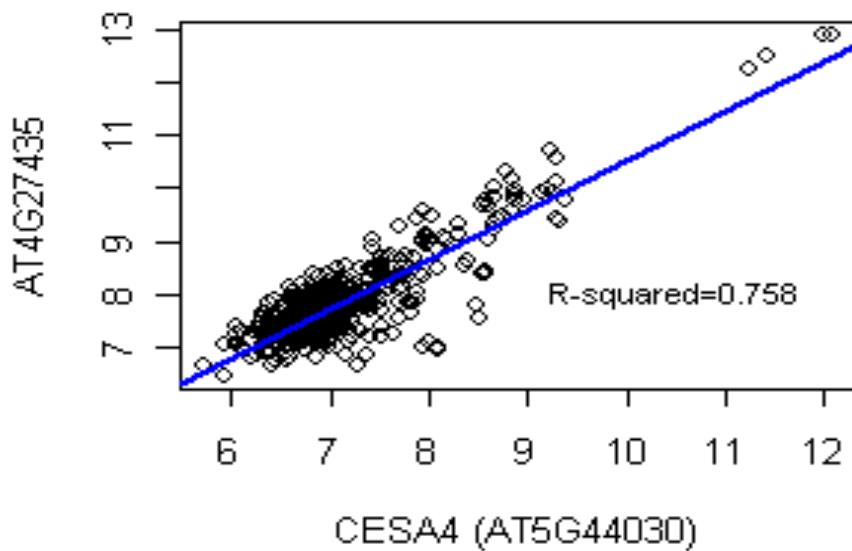
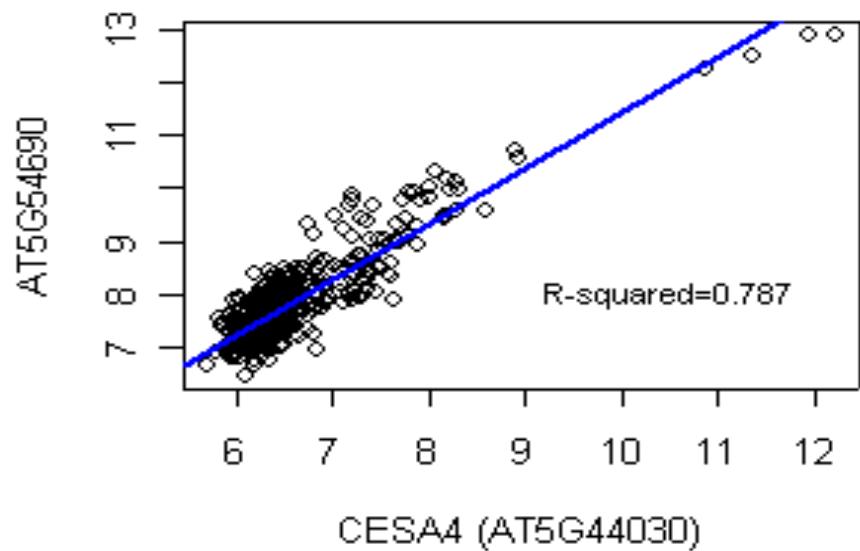
# A coexpression network centered around AtCESA4 7 and 8



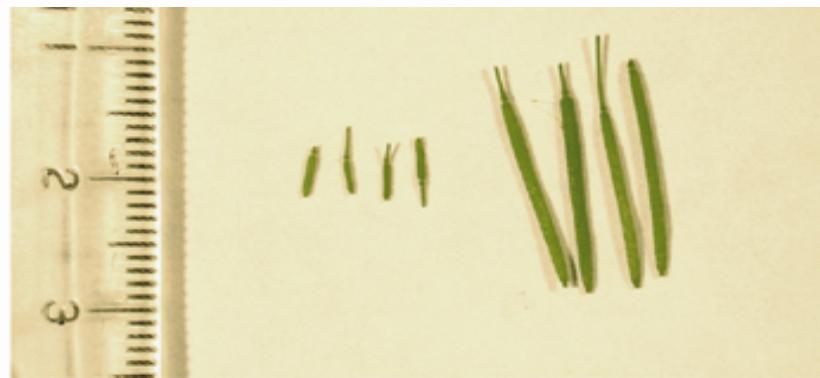
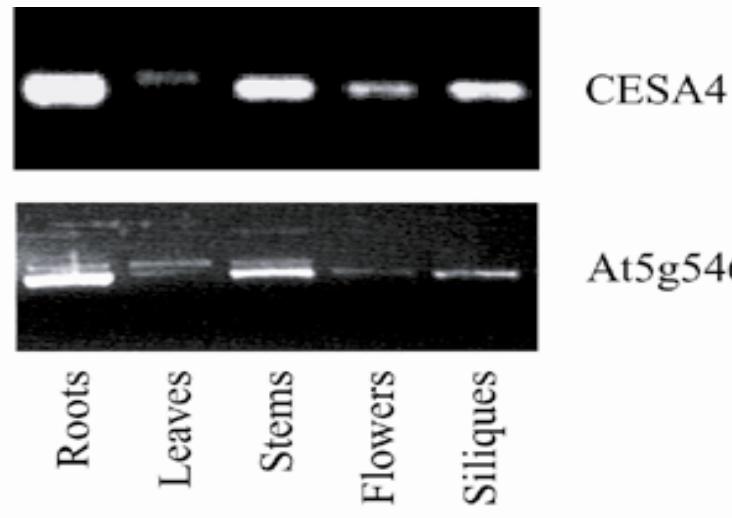
# Genes selected from ICE

AGI	Gene title	Score*	P-value
AT5G44030	cellulose synthase catalytic subunit (IRX5/CESA4)	4	4.1E-35
AT4G18780	cellulose synthase catalytic subunit (IRX1/CESA8)	5	7.1E-33
AT5G54690*	glycosyl transferase family 8 protein	5	4.1E-35
AT5G17420	cellulose synthase catalytic subunit (IRX3/CESA7)	5	3.7E-30
AT3G16920*	CTL1-like	6	9.3E-33
AT2G38080	laccase putative / diphenol oxidase putative	6	1.4E-34
AT5G15630	COBRA cell expansion protein COBL4	7	1.6E-32
AT5G03170*	fasciclin-like arabinogalactan-protein (FLA11)	8	6.7E-32
AT2G37090	glycosyl transferase family 43 protein	10	2.3E-29
AT3G18660	glycogenin glucosyltransferase (glycogenin)-related	29	4.9E-23
AT4G27435*	expressed protein	30	2.6E-24
AT5G60720	expressed protein	37	1.2E-20
AT3G62020	germin-like protein (GLP10)	42	9.2E-21
AT5G60020	laccase putative / diphenol oxidase putative	58	1.1E-17
AT4G28500	no apical meristem (NAM) family protein	62	9.7E-24
AT5G60490	fasciclin-like arabinogalactan-protein (FLA12)	64	7.9E-21

These genes highlighted in blue are four selected for experiment validation

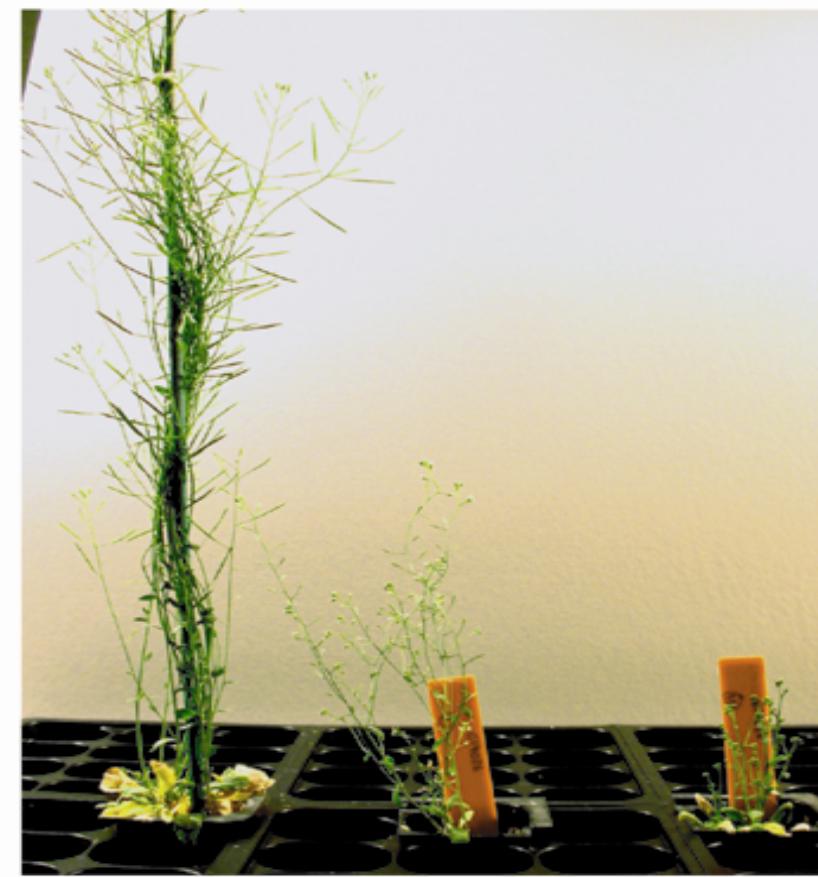


# Phenotypic changes in mutant lines



*At5g54690*

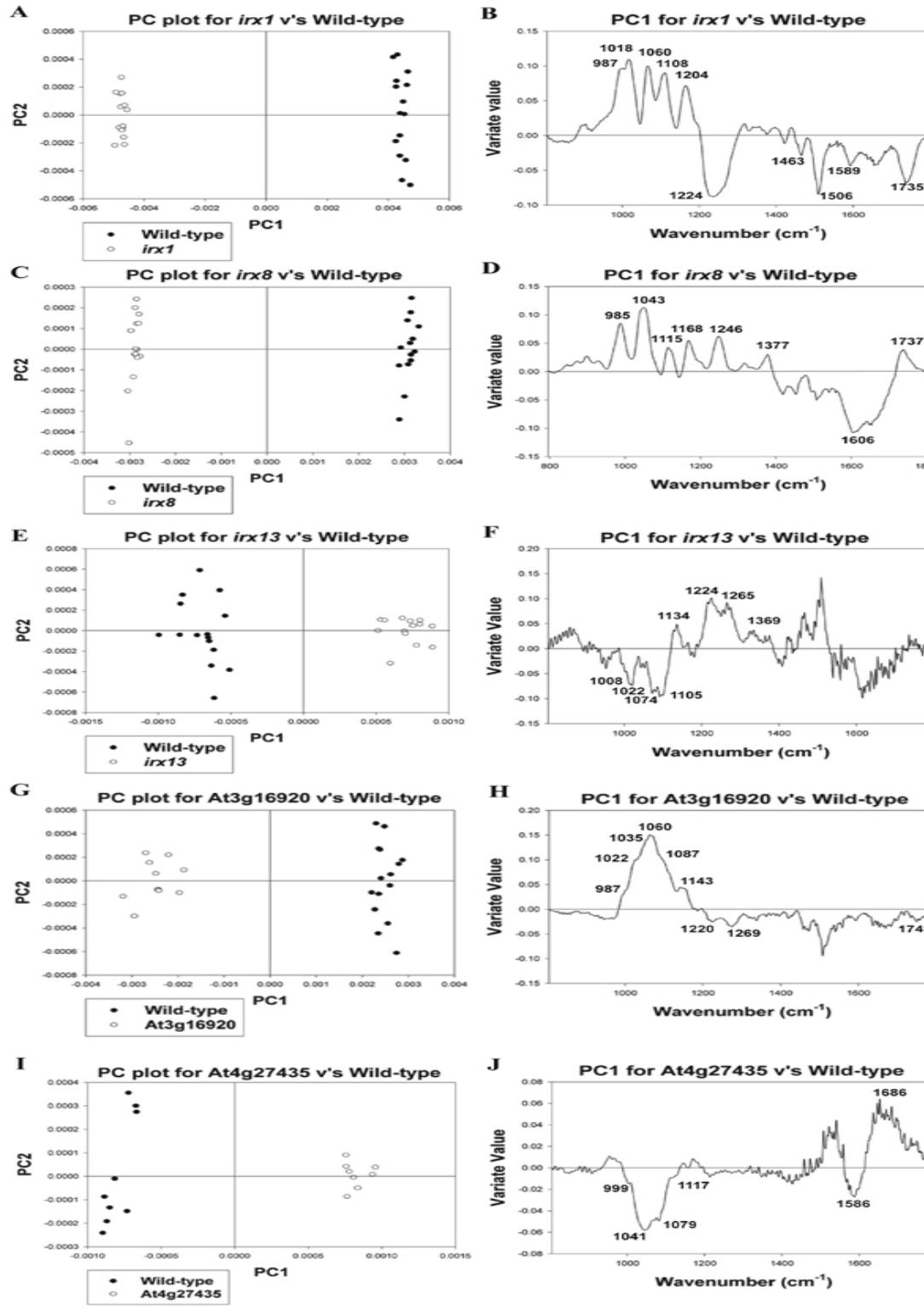
Wild-type



Wild-type

*At5g54690\_1*

*At5g54690\_2*

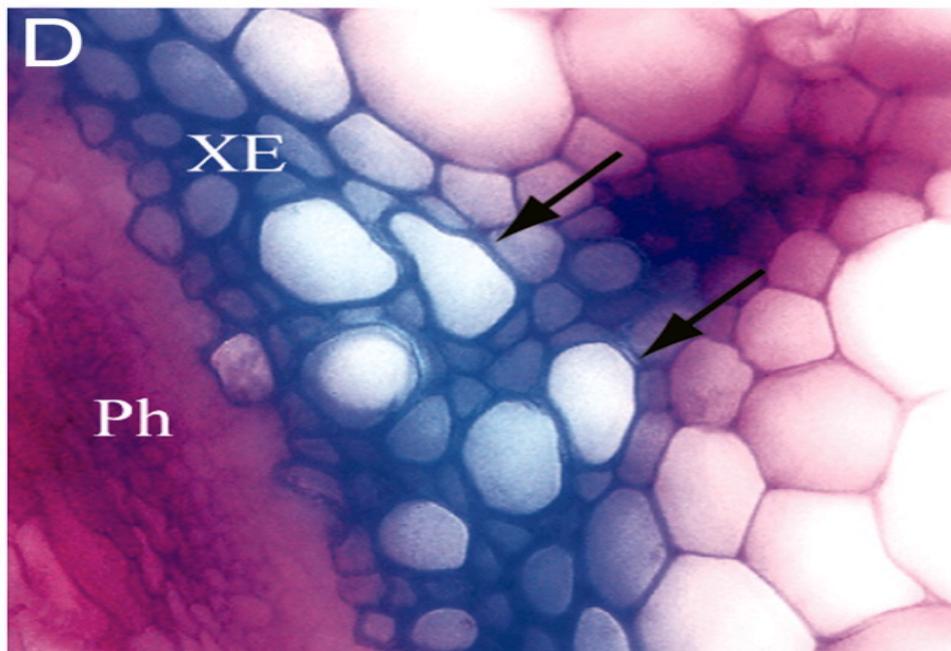
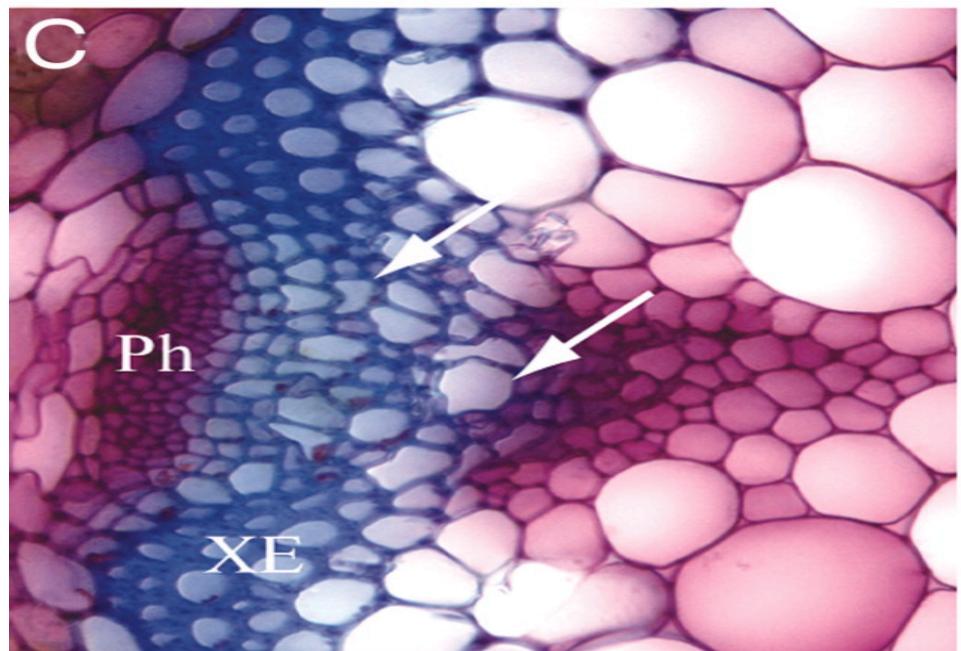
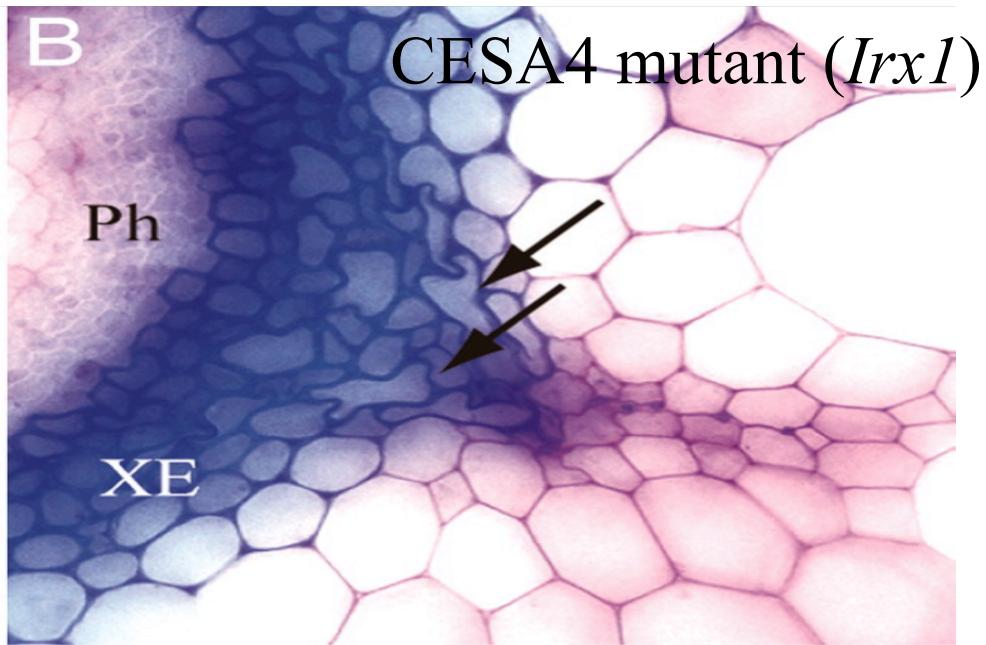
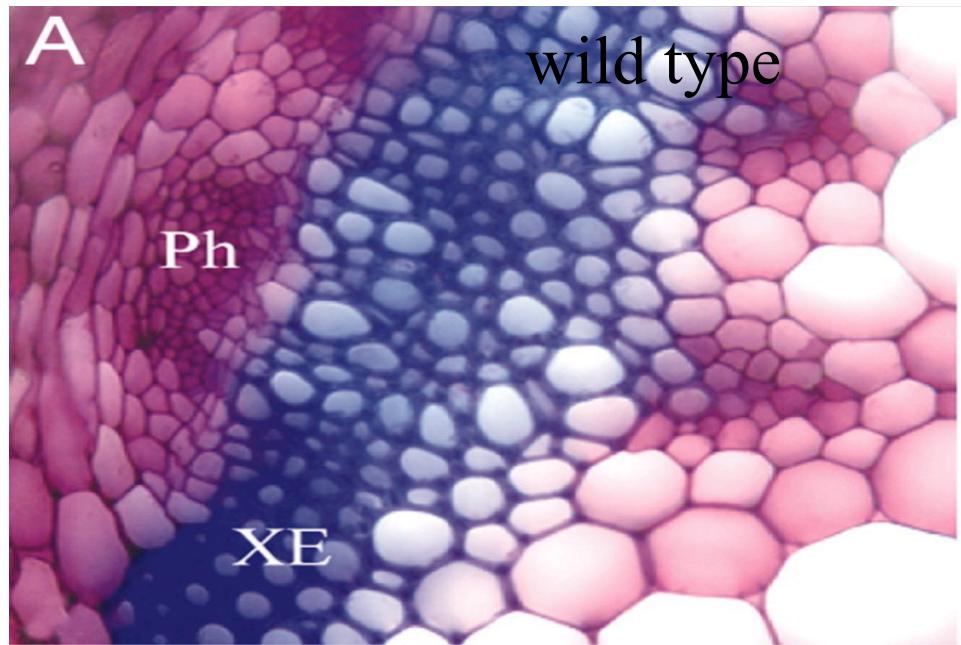


## FTIR Analysis

Fourier Transform Infrared Spectroscopy is a technique which is used to obtain an infrared spectrum of absorption or emission of a solid liquid or gas.

Cell wall analyses by FTIR. (A-F) Principal components plots for mutant [*irx1* (*CESA4* mutant) *irx8* (SALK\_014026 At5g54690) *irx13* (SALK\_046976 At5g03170 ) At3g16920 (SALK\_055713) and At4g27435 (SALK\_137109)] vs. wild-type spectra.

Spectra for all mutants show a clear separation from wild-type spectra based on principal component (PC) 1. In all cases the mutant spectra show differences from wild-type spectra in the carbohydrate fingerprint regions that correspond to the deformations of cellulosic and noncellulosic polymers.



At5g54690 knockout (*Irx8*)

At5g03170 knockout (*Irx13*)

# Related studies

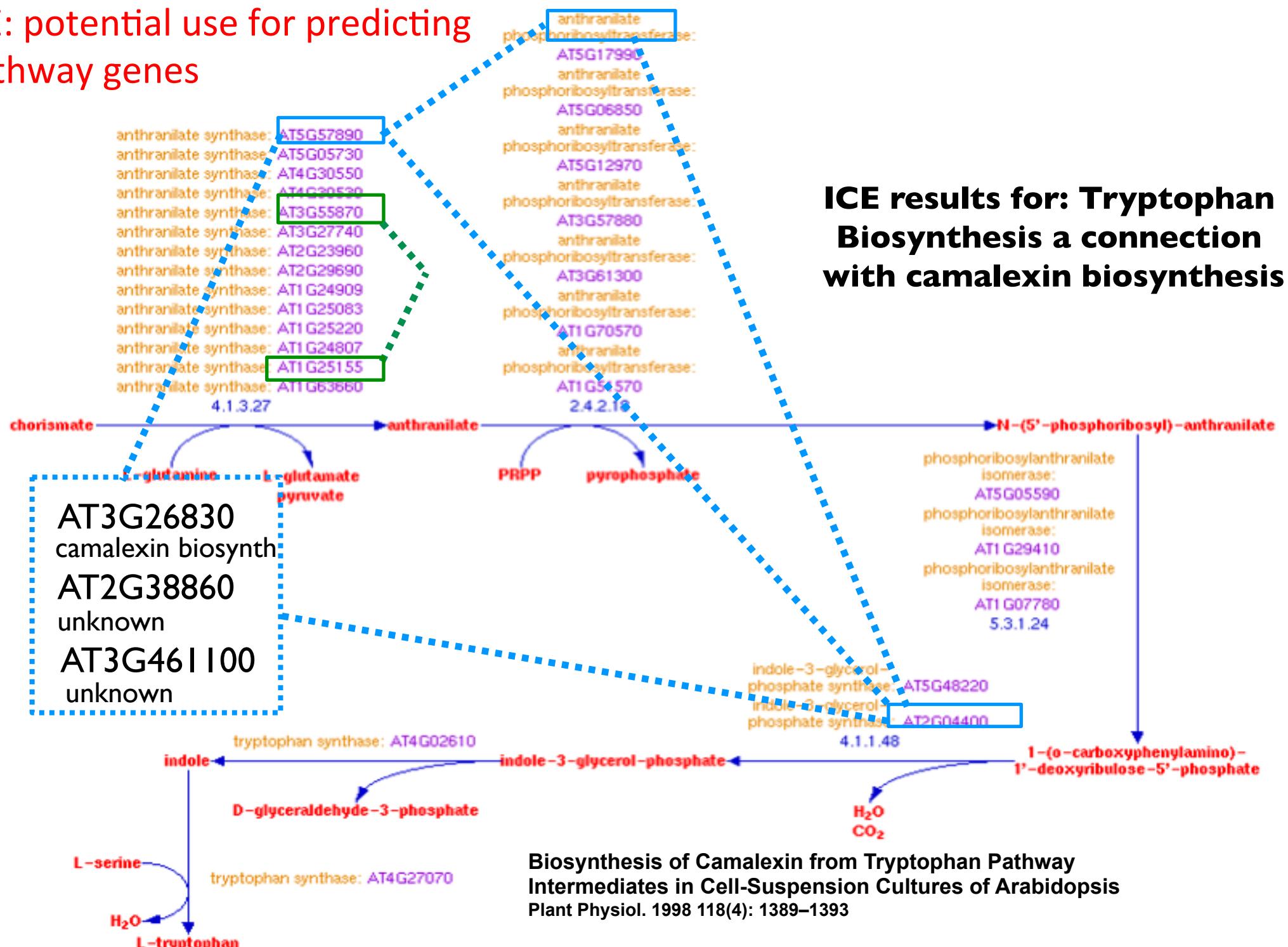
- 1. Among the 25 genes 16 are the same as those genes identified by Dr. Simon Turner via forward genetics**

Brown DM Zeef LA Ellis J Goodacre R Turner SR (2005) Identification of novel genes in *Arabidopsis* involved in secondary cell wall formation using expression profiling and reverse genetics. *The Plant cell* 17(8): 2281-95.

- 2. *Irx8* (At5g54690) protein has been shown to target to Golgi where glucuronoxylans (GXs) are synthesized. *Irx8* is implicated in affecting GX synthesis. GX cellulose and lignin are three major components of secondary cell walls in woody plants.**

Pena et al *Arabidopsis irregular xylem8 and irregular xylem9: Implications for the Complexity of Glucuronoxylan Biosynthesis* The Plant Cell 19:549-563 (2007)

# ICE: potential use for predicting pathway genes



## Summary for coexpression network:

1. Co-expression network can be used to predict genes involved in various pathways and biological processes. It is especially suitable for predicting functionally associated genes metabolic pathway genes structural genes that are usually located at the low hierarchical terminals of gene regulatory networks. These genes are often co-regulated which is why they are often highly correlated.
2. Due to the large number of terminal genes the identification of true genes is often interfered by noise. The prediction of pathway genes in many cases is not accurate. However it indeed works in some cases.

# What is gene collaborative network ?

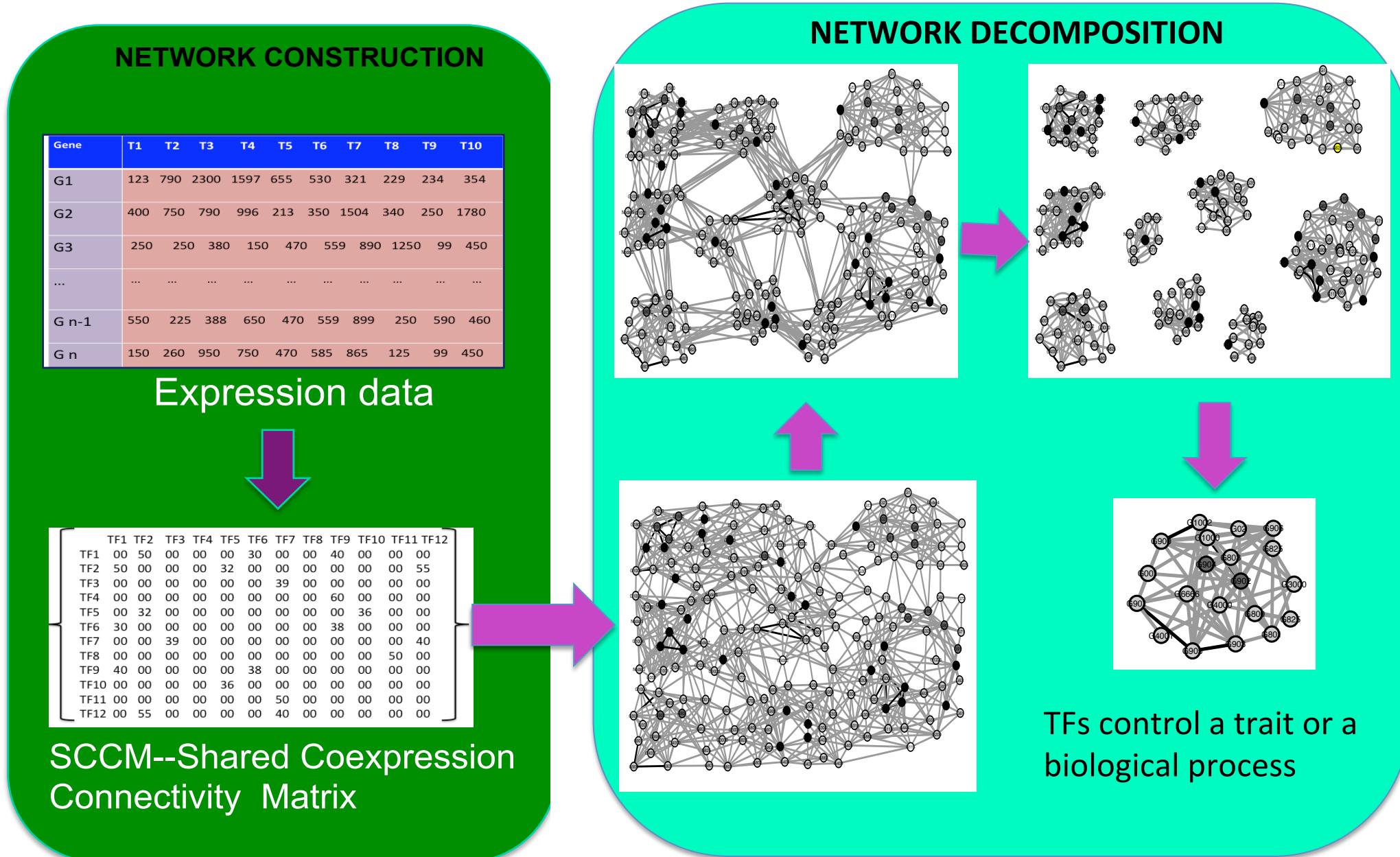
A **gene collaboration network (GCN)** is an undirected graph where each node corresponds to a regulatory gene (e.g. TF) and a pair of nodes is connected with an edge if there is a significant **collaboration** relationship between them.

# Microarray or RNA-seq data

Gene	T1	T2	T3	T4	T5	T6	T7	T8	T9	T10
G1	123	790	2300	1597	655	530	321	229	234	354
G2	400	750	790	996	213	350	1504	340	250	1780
G3	250	250	380	150	470	559	890	1250	99	450
...	...	...	...	...	...	...	...	...	...	...
G n-1	550	225	388	650	470	559	899	250	590	460
G n	150	260	950	750	470	585	865	125	99	450

Hypothesis: Regulatory genes (e.g. TFs) controlling a biological process or complex trait are collaborative.

# Identification of Regulatory Genes That Collectively Govern a Biological Process or Complex Trait?

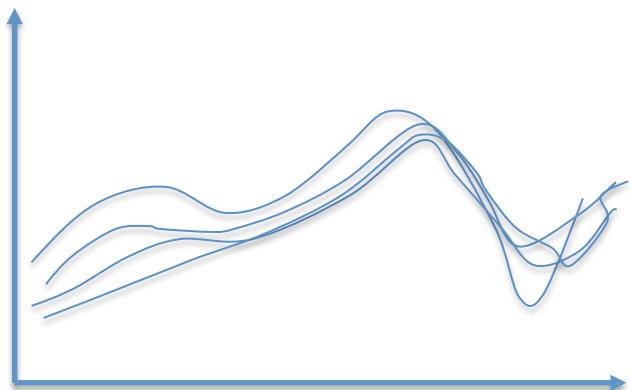


# How can we identify regulatory genes that are collaborative?

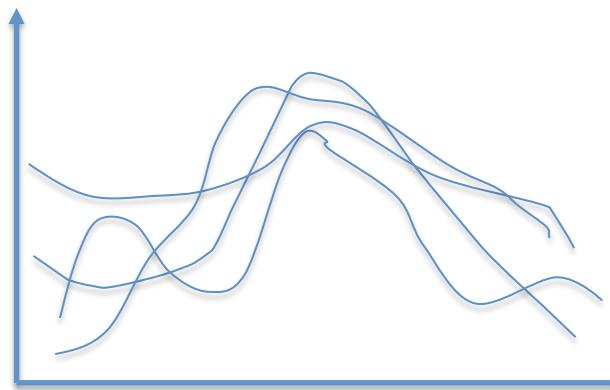
## Hypothesis

Regulatory genes that collaboratively control the same complex trait are **coordinated in gene expression profiles**.

Coexpressed genes refer to those that have higher concordance in expression profiles whereas coordinated genes refer to those that have relatively looser concordance in expression profiles.

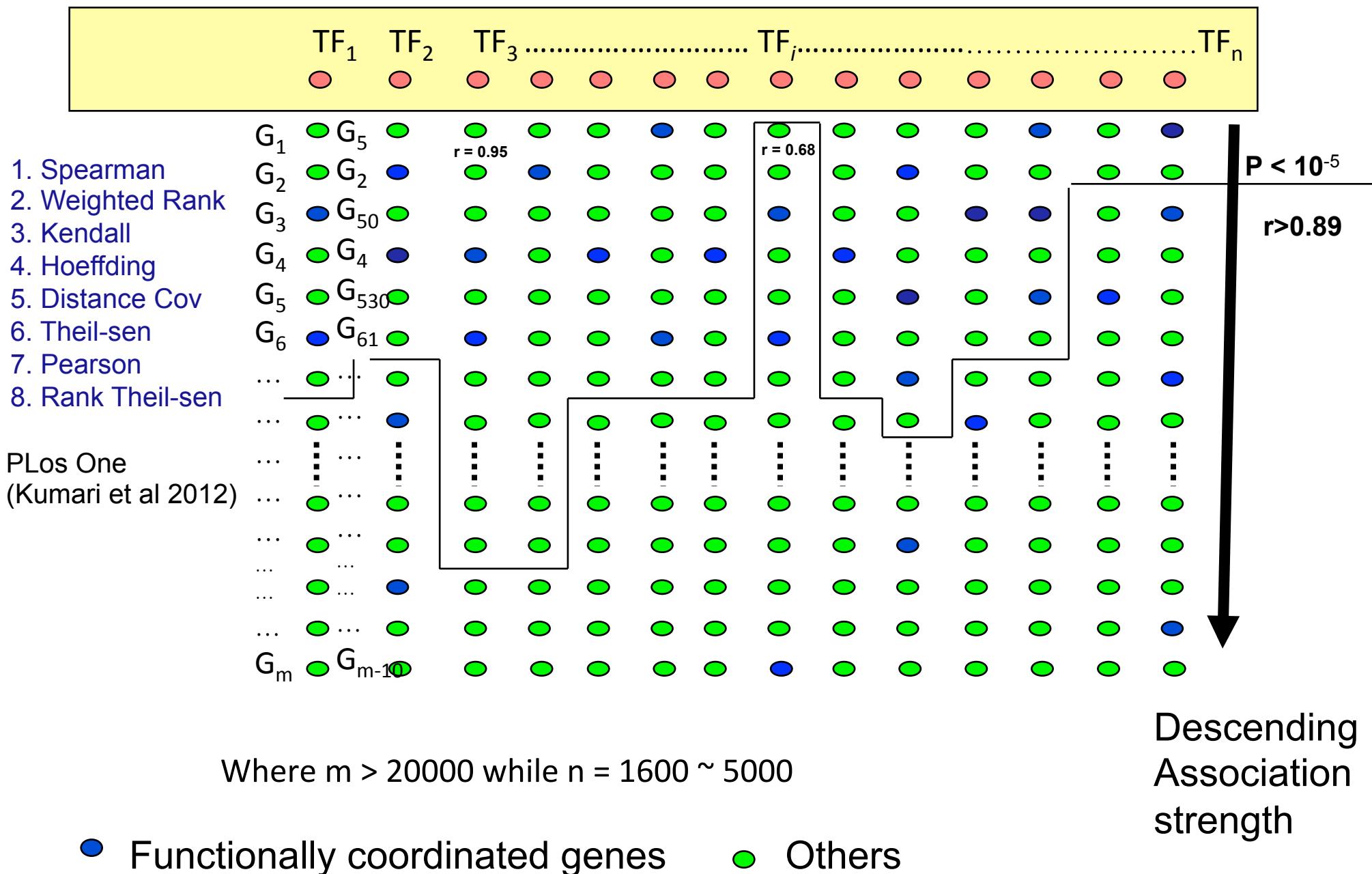


Coexpressed genes



Collaborative genes

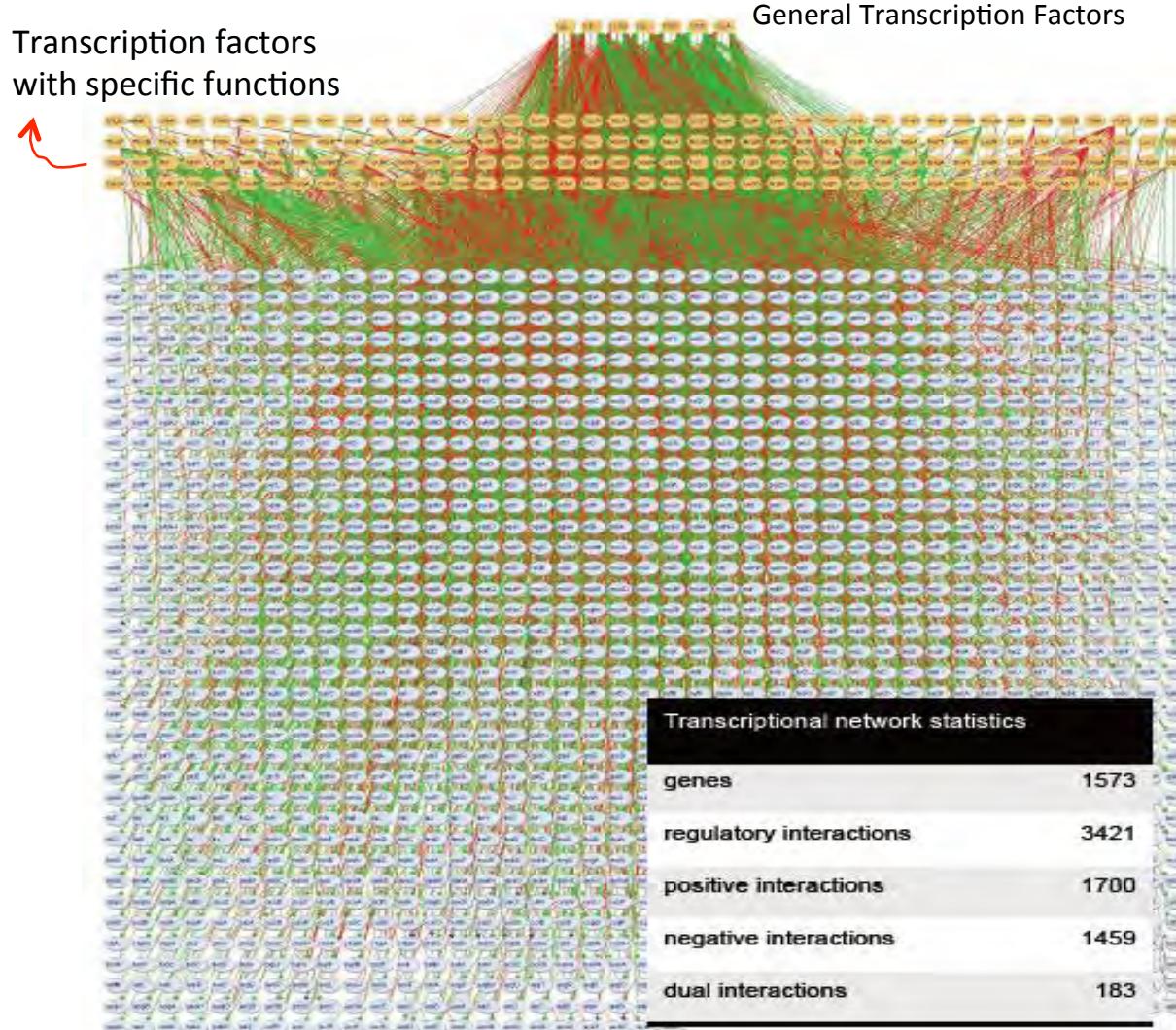
# Genome-wide coexpression analysis



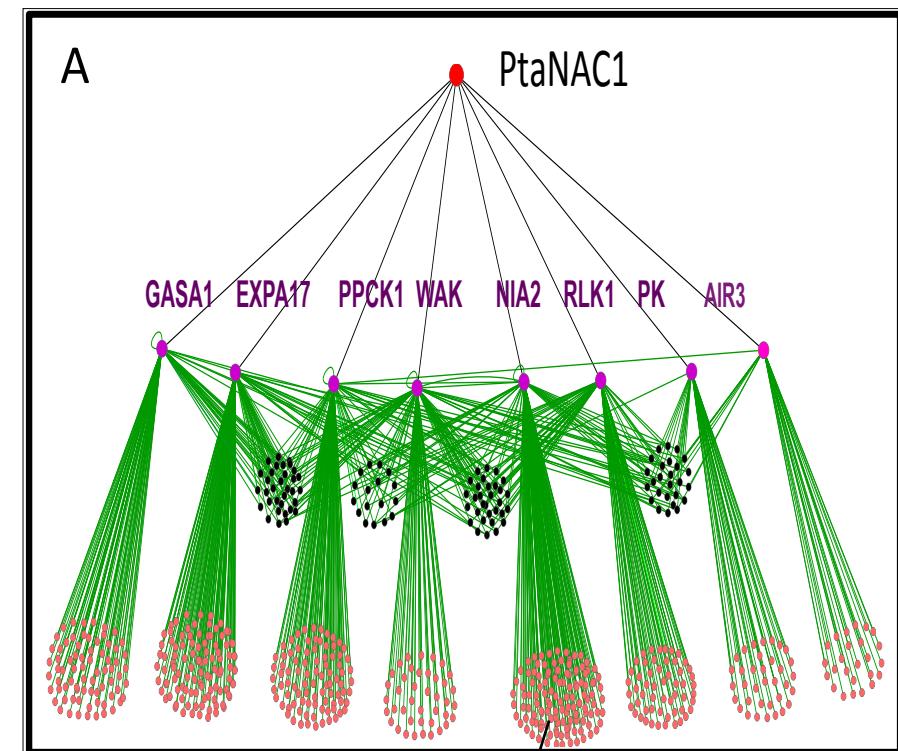
# Examples for illustration of TFs' hierarchies and their possible impact on coexpression

## *Escherichia coli* transcriptional regulatory network

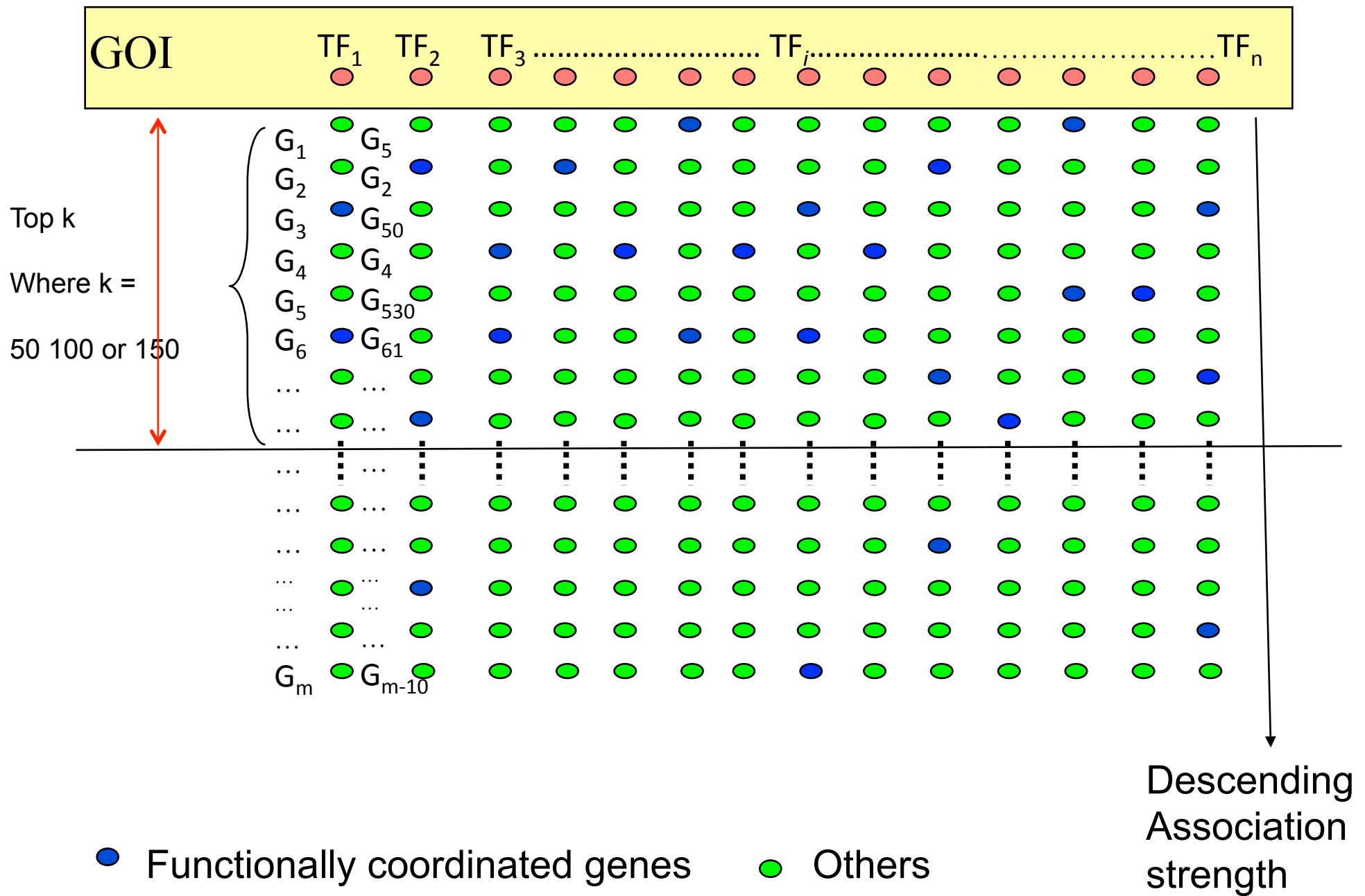
Network Biology 2011, 1(1): 21-33



## A Hierarchical Gene Regulatory Network in poplar

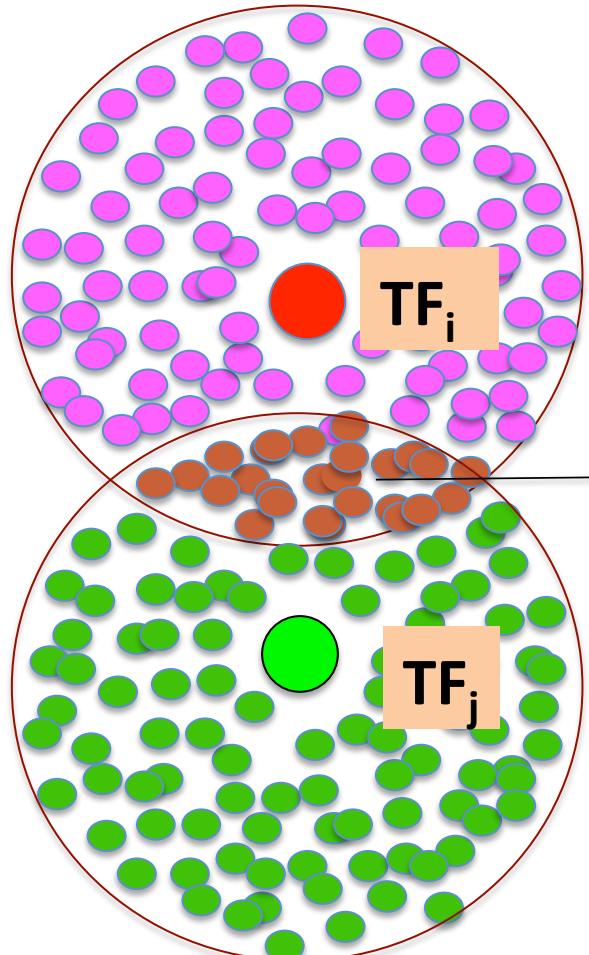


# Genome-wide coexpression analysis



# How to measure the loosely coordination between two TFs?

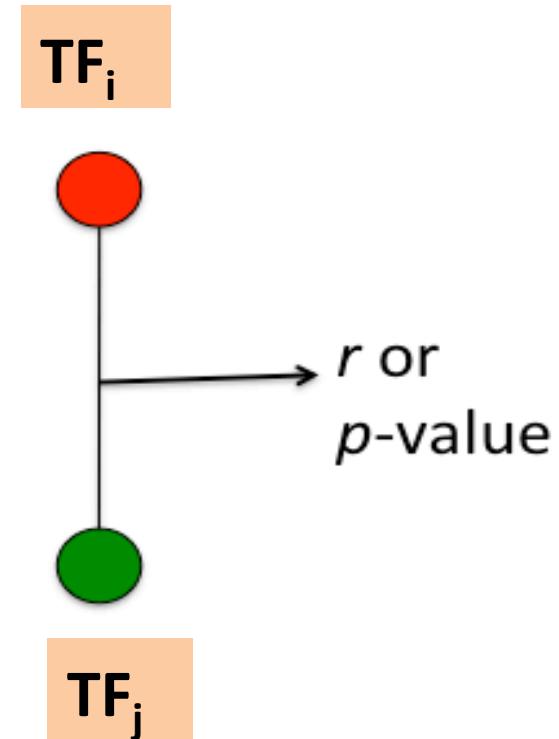
100 genes that are  
most tightly coexpressed to  $\text{TF}_1$



100 genes that are  
most tightly coexpressed to  $\text{TF}_j$

$$n_c$$

Number of  
shared genes



Kumari S Nie J Chen H-S Ma H et al. (2012) Evaluation of Gene Association Methods for Coexpression Network Construction and Biological Knowledge Discovery. PLoS ONE 7(11): e50411. doi:10.1371/journal.pone.0050411  
<http://www.plosone.org/article/info:doi/10.1371/journal.pone.0050411>

# Shared Coexpression Connectivity Matrix (SCCM)

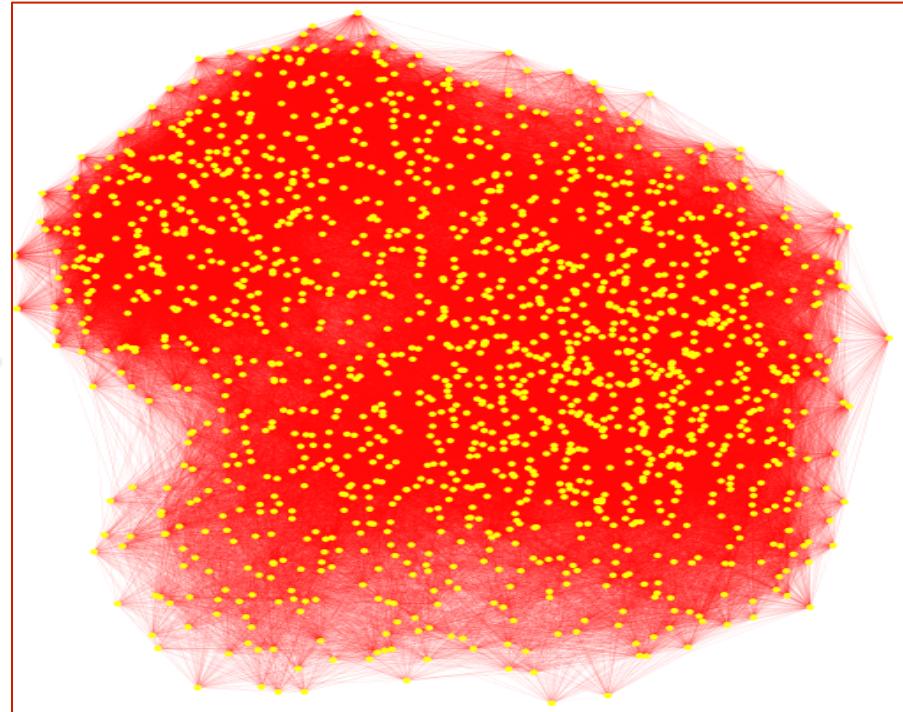
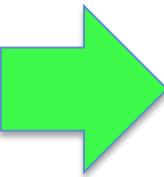
TFs

	AATF	ABCA11	ABCA11	ABCA11	ABT1	ACVR2A	ADNP	ADPGK	AEBP1	AEBP2	AFF1
AATF	100	0	0	0	3	0	0	0	0	0	0
ABCA11	0	100	0	13	0	0	0	0	0	0	0
ABCA11	0	0	100	11	0	0	0	0	0	0	0
ABCA11	0	13	11	100	0	0	0	0	0	0	0
ABT1	3	0	0	0	100	0	0	0	0	0	11
ACVR2A	0	0	0	0	0	100	8	8	0	0	0
ADNP	0	0	0	0	0	8	100	0	0	0	0
ADPGK	0	0	0	0	0	8	0	100	0	0	0
AEBP1	0	0	0	0	0	0	0	0	100	0	0
AEBP2	0	0	0	0	0	0	0	0	0	100	0
AFF1	0	0	0	0	11	0	0	0	0	0	100
AFF2	0	0	0	0	3	0	0	0	0	0	11
AFF3	0	0	0	0	0	0	0	0	0	0	0
AHR	0	0	0	0	0	0	0	0	0	0	0
AIP	3	0	0	0	13	0	0	0	0	0	15
AIP	0	0	0	0	0	0	0	0	0	4	0
AIRE	0	0	0	0	0	0	0	0	0	0	0
AK096318	0	0	0	0	0	0	0	0	0	0	0
AK126463	0	0	0	0	0	0	0	0	5	0	0
AL834146	0	0	0	0	0	4	11	3	0	0	0
ALF	0	0	0	0	0	0	0	0	0	0	0
ALX3	0	0	0	0	0	0	0	0	6	0	0
ALX4	0	0	0	0	0	0	0	0	0	0	0
AMOT	0	0	0	0	0	0	0	0	0	0	0
ANKRD30A	0	0	0	0	0	0	0	0	0	0	0
ANP32A	0	0	0	0	0	0	0	0	0	0	0
APBA2BP	0	0	0	0	0	0	0	0	0	0	0
APEX1	0	0	0	0	0	0	0	0	0	0	0
AR	0	0	0	0	0	0	0	0	0	0	0
ARC	0	0	0	0	0	0	0	0	0	0	0
ARID1A	0	0	0	0	0	0	17	0	0	0	0
ARID1B	0	0	0	0	8	0	0	0	0	0	19

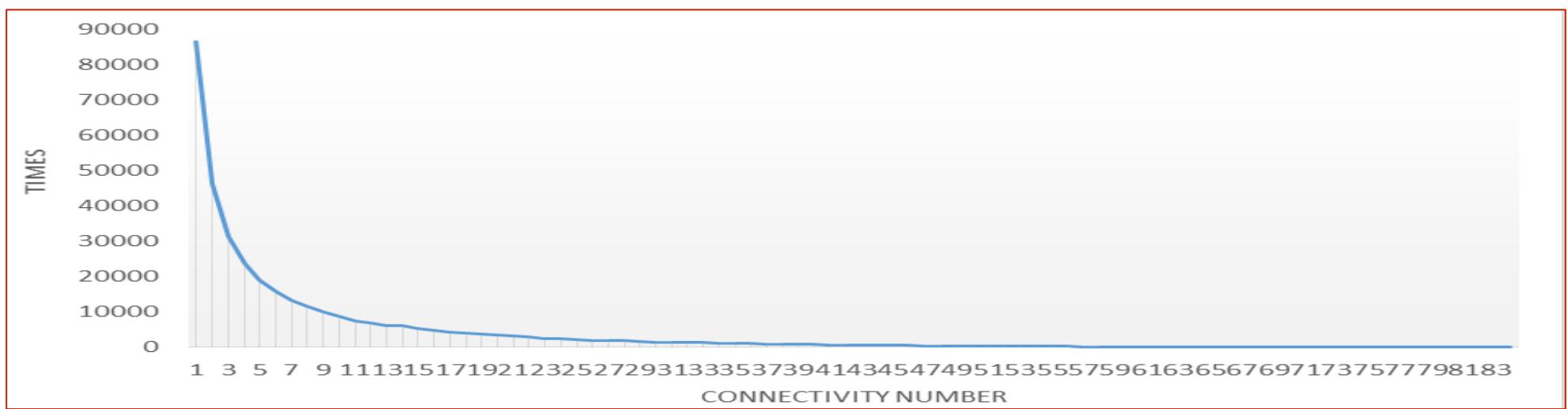
# Collaborative Network of TF

	AATF	ABCA11	ABCA11	ABCA11	ABT1	ACVR2A	ADNP	ADPGK	AEBP1	AEBP2	AFF1
AATF	100	0	0	0	3	0	0	0	0	0	0
ABCA11	0	100	0	13	0	0	0	0	0	0	0
ABCA11	0	0	100	11	0	0	0	0	0	0	0
ABCA11	0	13	11	100	0	0	0	0	0	0	0
ABT1	3	0	0	0	100	0	0	0	0	0	11
ACVR2A	0	0	0	0	0	100	8	8	0	0	0
ADNP	0	0	0	0	0	8	100	0	0	0	0
ADPGK	0	0	0	0	0	8	0	100	0	0	0
AEBP1	0	0	0	0	0	0	0	0	100	0	0
AEBP2	0	0	0	0	0	0	0	0	0	100	0
AFF1	0	0	0	0	11	0	0	0	0	0	100
AFF2	0	0	0	0	3	0	0	0	0	0	11
AFF3	0	0	0	0	0	0	0	0	0	0	0
AHR	0	0	0	0	0	0	0	0	0	0	0
AIP	3	0	0	0	13	0	0	0	0	0	15
AIP	0	0	0	0	0	0	0	0	0	4	0
AIRE	0	0	0	0	0	0	0	0	0	0	0
AK096318	0	0	0	0	0	0	0	0	0	0	0
AK126463	0	0	0	0	0	0	0	0	5	0	0
AL834146	0	0	0	0	0	4	11	3	0	0	0
ALF	0	0	0	0	0	0	0	0	0	0	0
ALX3	0	0	0	0	0	0	0	0	6	0	0
ALX4	0	0	0	0	0	0	0	0	0	0	0
AMOT	0	0	0	0	0	0	0	0	0	0	0
ANKRD30A	0	0	0	0	0	0	0	0	0	0	0
ANP32A	0	0	0	0	0	0	0	0	0	0	0
APBA2BP	0	0	0	0	0	0	0	0	0	0	0
APEX1	0	0	0	0	0	0	0	0	0	0	0
AR	0	0	0	0	0	0	0	0	0	0	0
ARC	0	0	0	0	0	0	0	0	0	0	0
ARID1A	0	0	0	0	0	0	17	0	0	0	0
ARID1B	0	0	0	0	8	0	0	0	0	0	19

shared coexpression connectivity matrix (SCCM)



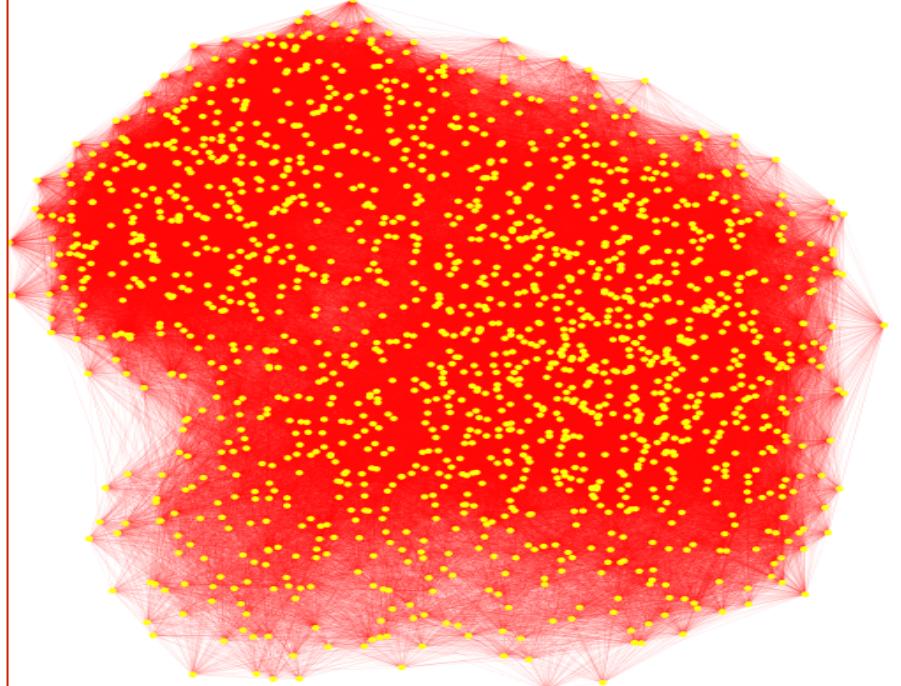
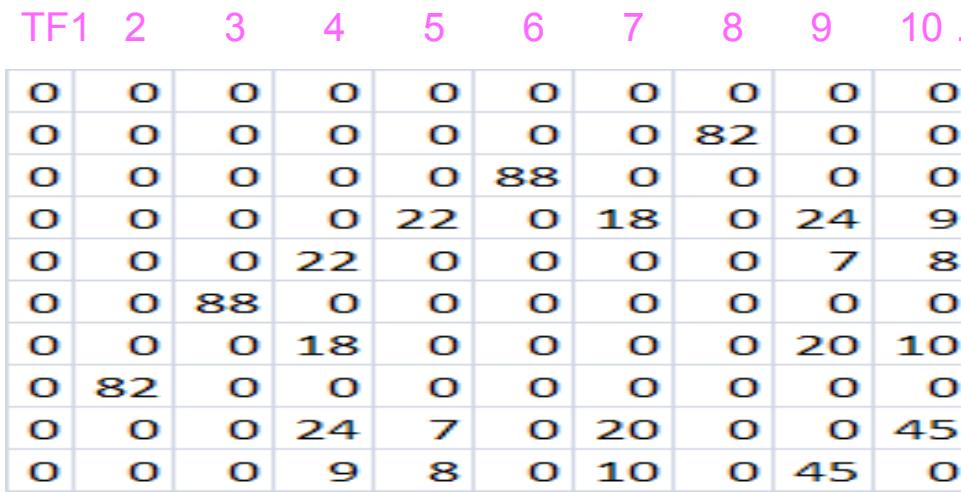
Collaborative network of all TFs



Gene	N6_R1	N6_R2	N12_R1	N12_R2	H24_R1	H24_R2	N48_R1	N48_R2	N96_R1	N96_R2	N94_R1	N94_R2	CTRL_R1	CTRL_R2
PgkHfr_1001_1.41_at	7.555895	7.388326	7.292131	4.469743	7.145161	7.288017	7.300384	7.449564	7.462835	7.240070	7.474262	7.474262	7.031193	7.031193
PgkHfr_1002_1.41_at	3.555152	3.892755	3.893563	1.893653	2.700424	6.909127	2.474975	3.071873	1.748641	3.178731	1.748641	3.178731	2.850748	2.850748
PgkHfr_1003_1.41_at	9.736727	9.389712	9.507334	6.638864	9.813783	9.273278	8.718356	9.474073	9.192535	9.703574	9.038595	9.038595	9.265185	9.265185
PgkHfr_1004_1.41_at	8.479371	8.058951	8.058951	8.058951	8.058951	8.173563	8.173563	8.173563	8.109951	8.109951	9.703574	9.703574	9.703574	9.703574
PgkHfr_1008_1.41_at	11.043953	11.040865	11.227214	10.850549	10.267324	10.360719	9.750629	10.028971	9.54526729	9.397994	9.54526729	9.54526729	9.6507046	9.4086651
PgkHfr_1009_1.41_at	7.014704	7.884052	7.5098454	7.212688	7.014704	7.381651	7.403056	7.482754	7.514725	7.610581	7.884260	7.884260	7.740513	7.740513
PgkHfr_1010_1.41_at	1.467576	1.1174528	1.2021072	1.790493	1.459346	1.767876	1.317876	1.300734	1.300734	1.309759	1.309759	1.309759	1.317296	1.317296
PgkHfr_1011_1.41_at	1.5014096	1.081026	1.3901555	1.408281	1.408281	2.171262	1.269045	1.699416	1.776051	1.477315	1.527459	1.527459	1.267042	1.267042
PgkHfr_1012_1.41_at	1.576794	1.5993246	1.5993246	1.5993246	1.5993246	1.700075	1.700075	1.700075	1.645215	1.645215	1.700075	1.700075	1.700075	1.700075
PgkHfr_1013_1.41_at	1.3194411	1.0287298	1.3101685	1.3843655	1.7657661	1.4032052	1.5646050	1.15181941	1.7802874	1.4586354	1.7802874	1.4586354	1.9961515	1.30574988
PgkHfr_1014_1.41_at	1.8938869	4.9624626	7.1783408	4.7033269	4.41191717	9.734966	1.16191928	1.703685	1.75352154	6.199352	1.7419459	1.7419459	1.920448	1.920448
PgkHfr_1015_1.41_at	1.881018	1.51270502	1.51270502	1.51270502	1.51270502	1.594033	1.594033	1.594033	1.594033	1.594033	1.594033	1.594033	1.588145	1.588145
PgkHfr_1016_1.41_at	1.021727	1.8546368	1.8033761	4.8619137	4.498136	8.742783	1.3463293	1.408033	1.408033	2.850048	1.408033	1.408033	1.408033	1.408033
PgkHfr_1017_1.41_at	5.640293	5.7631207	5.7631207	5.7631207	5.7631207	5.4903581	5.4903581	5.4903581	5.4903581	5.4903581	5.4903581	5.4903581	5.715817	5.715817
PgkHfr_1018_1.41_at	6.717046	4.084254	2.079123	4.7033269	5.583393	4.342324	1.323808	4.051675	4.110758	6.004085	4.110758	4.110758	5.491582	5.491582
PgkHfr_1019_1.41_at	1.001404	1.488293	1.488293	1.488293	1.488293	1.333948	1.333948	1.333948	1.333948	1.333948	1.333948	1.333948	1.4027375	1.4027375
PgkHfr_1020_1.41_at	1.001071	1.487208	1.487208	1.487208	1.487208	1.333948	1.333948	1.333948	1.333948	1.333948	1.333948	1.333948	1.4027375	1.4027375
PgkHfr_1021_1.41_at	1.117398	1.156902	1.3794067	4.383864	7.827109	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1022_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1023_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1024_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1025_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1026_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1027_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1028_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1029_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1030_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1031_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1032_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1033_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1034_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1035_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1036_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1037_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1038_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1039_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1040_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1041_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1042_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1043_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1044_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1045_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1046_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1047_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1048_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1049_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1050_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1051_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1052_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1053_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1054_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1055_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1056_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1057_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1058_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1059_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.428704	1.388278	1.388278
PgkHfr_1060_1.41_at	1.117398	1.2717857	1.2717857	1.2717857	1.2717857	1.428704	1.428704	1						

How to decompose  
this collaborative  
network into  
subnetworks each  
of which controls a  
biological process  
or complex trait?

# Collaborative Network



TF	TF1	2	3	4	5	6	7	8	9	10	...
1	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	82	0	0	0
3	0	0	0	0	0	88	0	0	0	0	0
4	0	0	0	0	22	0	18	0	24	9	0
5	0	0	0	22	0	0	0	0	7	8	0
6	0	0	88	0	0	0	0	0	0	0	0
7	0	0	0	18	0	0	0	0	20	10	0
8	0	82	0	0	0	0	0	0	0	0	0
9	0	0	0	24	7	0	20	0	0	45	0
10	0	0	0	9	8	0	10	0	45	0	0

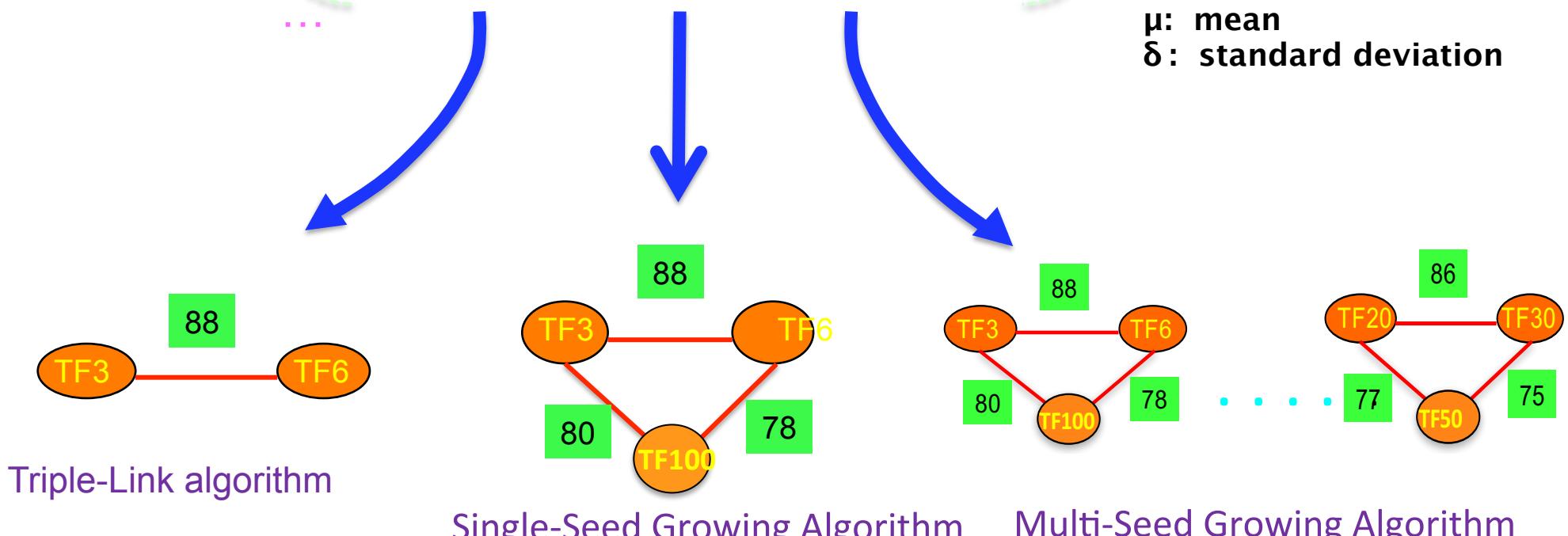
$$\alpha = \mu + \delta$$

$$\beta = \mu + 2\delta$$

$$\gamma = \mu + 3\delta$$

$\mu$ : mean

$\delta$ : standard deviation



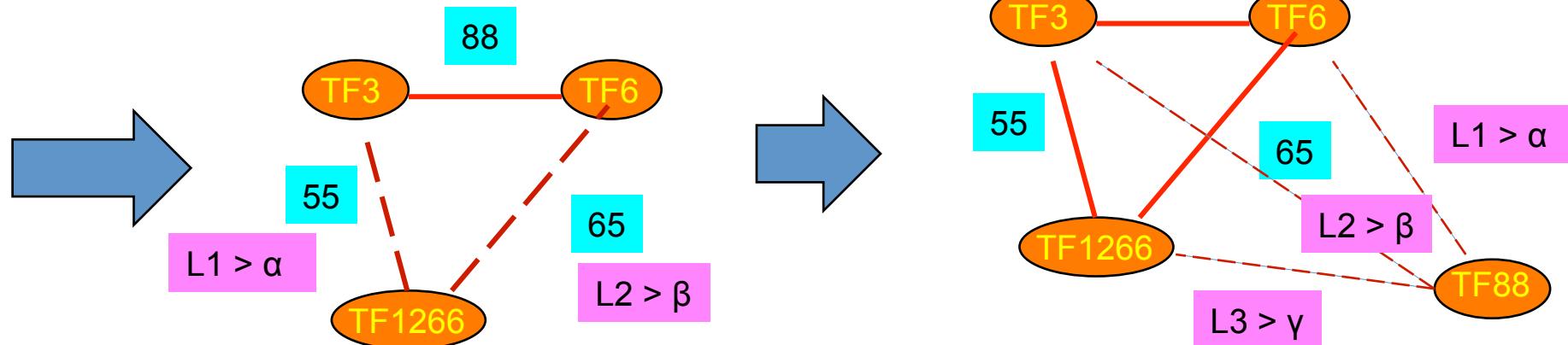
# Triple-link algorithm

TF	TF1	2	3	4	5	6	7	8	9	10	...
1	0	0	0	0	0	0	0	0	0	0	
2	0	0	0	0	0	0	0	82	0	0	
3	0	0	0	0	0	88	0	0	0	0	
4	0	0	0	0	22	0	18	0	24	9	
5	0	0	0	22	0	0	0	0	7	8	
6	0	0	88	0	0	0	0	0	0	0	
7	0	0	0	18	0	0	0	0	20	10	
8	0	82	0	0	0	0	0	0	0	0	
9	0	0	0	24	7	0	20	0	0	45	
10	0	0	0	9	8	0	10	0	45	0	
...											

$$\alpha = \mu + \delta$$

$$\beta = \mu + 2\delta$$

$$\gamma = \mu + 3\delta$$



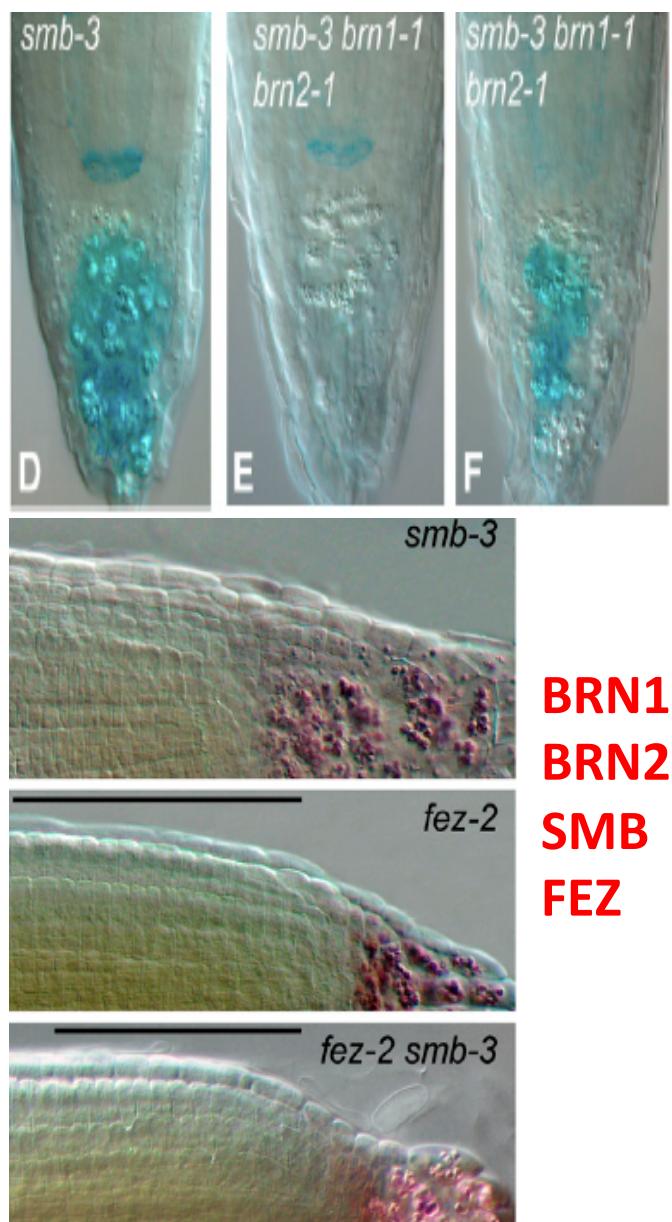
# Test 1: Arabidopsis Data

## **Arabidopsis salt stress experiments**

- 6 Microarray experiments of Arabidopsis
- 108 chips (downloaded from GEO)

Question: What TFs control root growth under salt stress condition in Arabidopsis?

Cluster 2: A cluster of TFs control root cap development (stem cells of roots)			
AT1G33280	BRN1	BRN1, SMB control root cap maturation	[66]
AT4G10350	BRN2	BRN2, SMB control root cap maturation	[66]
AT1G79580	SMB	FEZ and SMB control root stem cells	[67]
AT5G39820	ANAC094	Apical meristem protein, function unknown	[59]
AT1G26870	FEZ	FEZ and SMB control root stem cells in cap	[67]
AT1G74500	TOM7	Embryonic root initiation	[68]
AT3G27010	TCP20	Postembryonic cell division in root	[70]
AT2G30340	LBD13	Expressed in cells at the adaxial base of lateral roots	[104]
AT2G40470	LBD15	Expressed in cells at the adaxial base of lateral roots	[104]
AT1G51190	PLT2	Control root stem cell activity near cap	[69]
AT1G66350	RGL1	Root epidermal differentiation	[105]
AT2G37260	TTG2	Differentiation of trichomes and root hairless cells	[106]
AT5G57420	IAA33	IAA is involved in root development	[107, 108]
AT2G29060		scarecrow transcription factor family protein	



# Test 2: Data Set from humans

What TFs are necessary for pluripotency maintenance in human stem cells?

Present knowledge about regulation of human stem cell pluripotency

NANOG + POU5F1 + SOX2

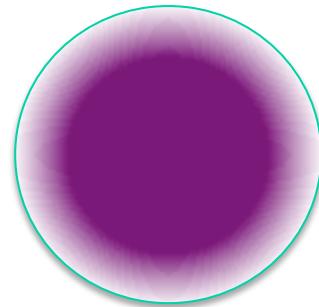
SALL2  
PHC1  
PRDM14  
DNMT3A

CSRP1  
NANOGP8  
LOC653441  
SOX3

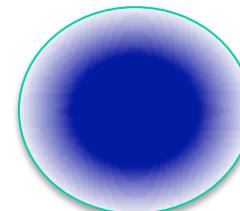
LIN28  
MYCL1  
DNMT3B  
RORB

Identified by many biologists from 1998 to 2008

The Power of TFs: Reprogramming in somatic cells



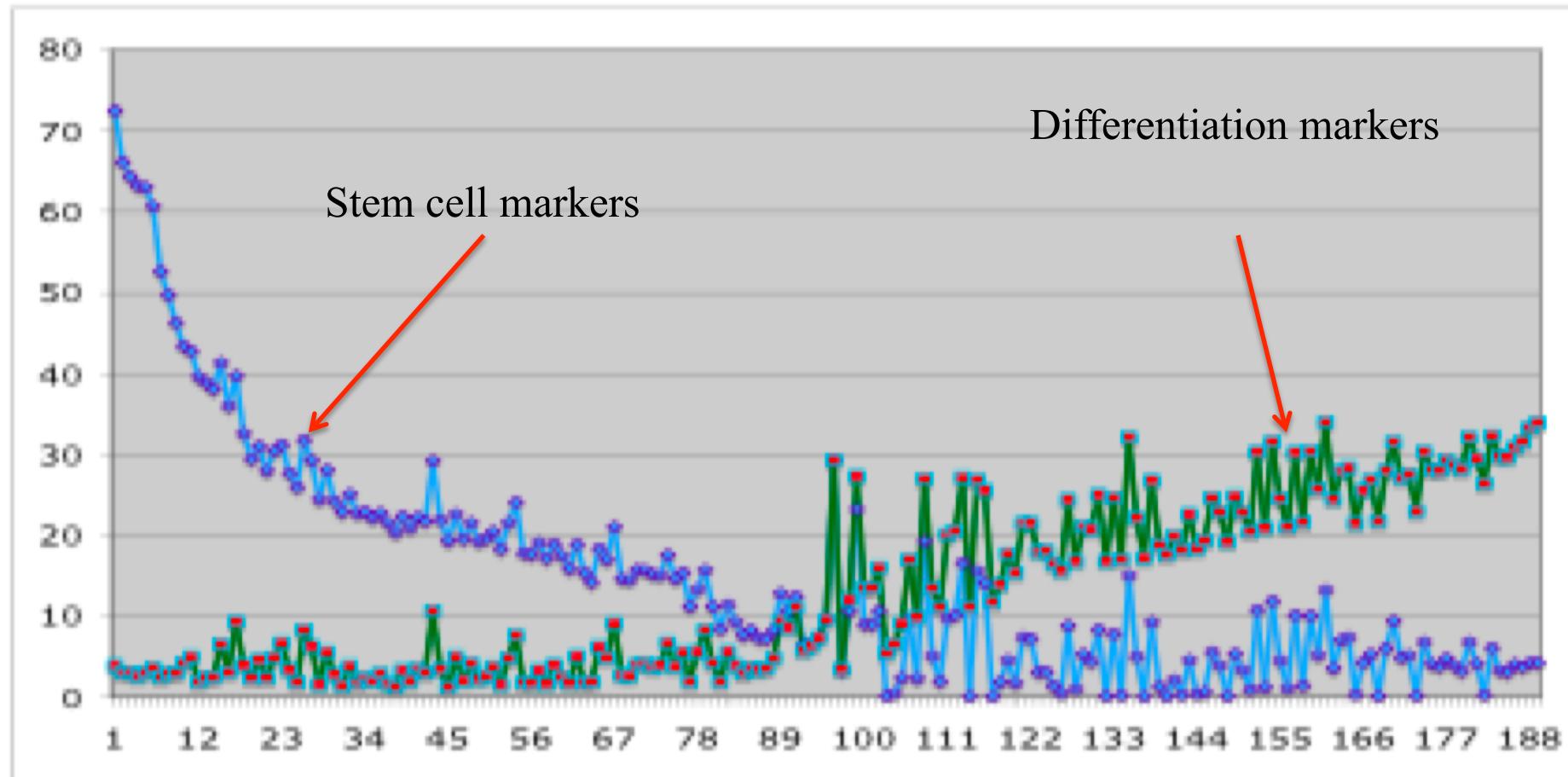
Skin Cell



Induced Pluripotency Cell (IPC)

NANOG  
POU5F1  
SOX2  
+ MYC2

Input data: 189 microarray chips from human stem cells that underwent differentiation from 0 ~ 96 hours.



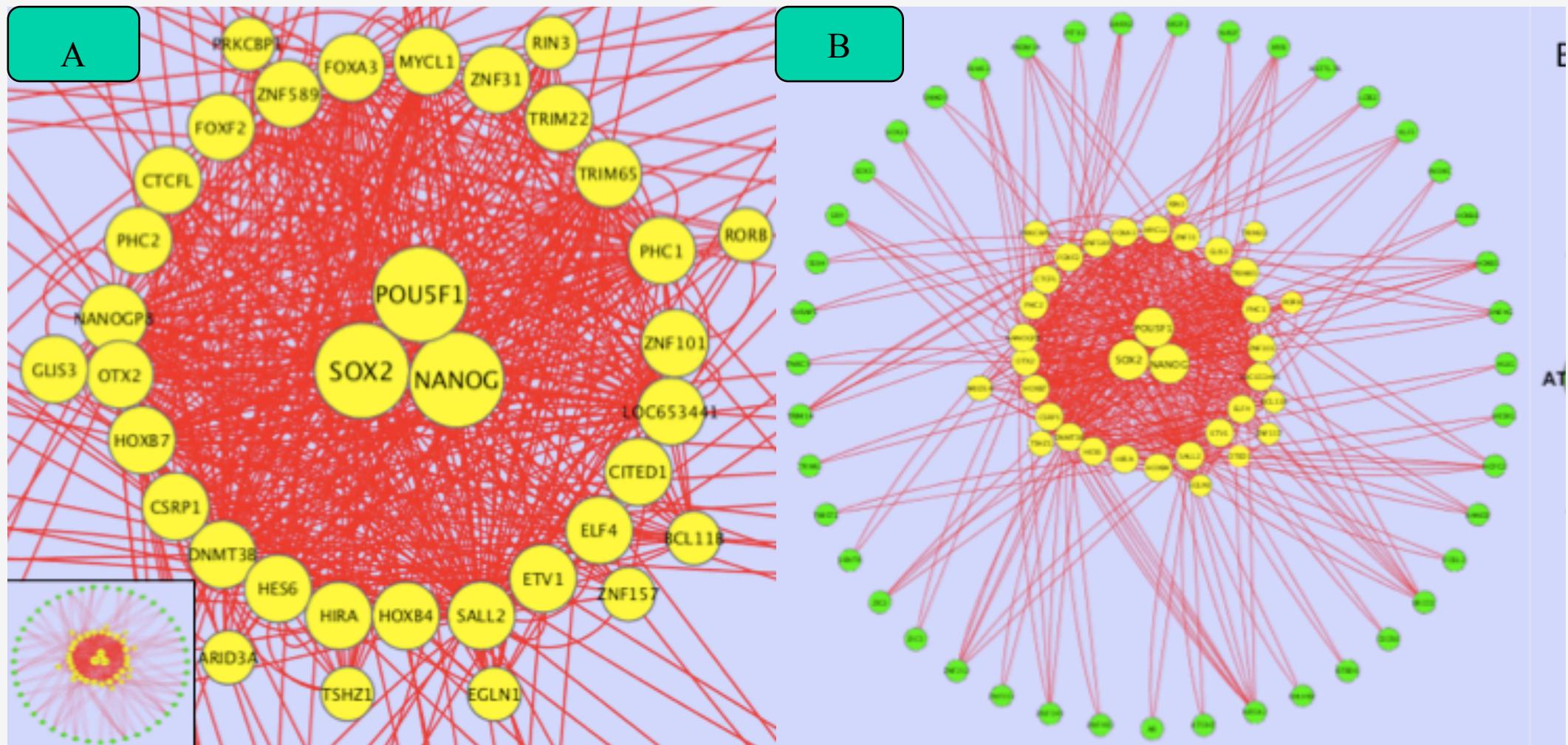
Microarray chips (189) provided by James Thomson @ U of Wisconsin

# Subnetwork 1 : Output from TF-cluster

Genes	Symbol	Description	Evidence
NM_024865	NANOG	Nanog homeobox	[14]
BC099704	NANOGP8	Nanog homeobox pseudogene 8	Pseudogene with similarity to Nanog.
NM_003106	SOX2	SRY box 2	[14]
NM_002701	POU5F1	POU class 5 homeobox 1	[14]
NM_006892	DNMT3B	DNS methyltransferase 3 beta	[15, 16]
NM_004078	CSRP1	cysteine-rich protein	
NM_080618	CTCFL	CCCTC-binding factor (zinc finger protein)-like	[17]
NM_016089	ZNF589	Zinc finger 589	[18, 19]
NM_004426	PHC1	Polyhomeotic homolog 1	[20]
NM_005407	SALL2	SAL2 like	[18]
NM_018645	HES6	hairy and enhancer of split 6	
NM_173547	TRIM65	Tripartite motif containing 65	
NM_004427	PHC2	Polyhomeotic homolog 2	[21]
NM_032805	ZFP206	Zinc finger protein 206 (ZSCAN10)	[22, 23]
NM_001421	ELF4	ETS domain TF	
NM_003325	HIRA	HIR Histone Cell Cycle regulator	[24]
NM_033204	ZNF101	Zinc finger protein 101	
BC098403	ETV1	ets variant 1	[25, 26]
NM_006079	CITED2	Cbp/p300-interacting transactivator	[27, 28]
NM_021728	OTX2	orthodenticle homeobox 2	
NM_024015	HOXB4	homeobox B4	
NM_006074	TRIM22	tripartite motif-containing 22	[19]
XM_929986	LOC653441	similar to polyhomeotic 1-like	Gene with sequence similarity to PHC1
NM_004497	FOXA3	Forkhead box 3	

Among 24 predicted TFs 17 have literature evidence (~70%)

# The coordination between TFs within Subnetwork 1



- A. Intra-connections among the genes within Collaborative Subnetwork 1
- B. Comparison of intra-subnetwork connections with outward-connections

# Test 3: Axolotl RNA-Seq

Total 74 samples

## Axolotl transcriptome

S, Z, A, Digits,

30 day blastema|

RNA-Seq (on GA IIx)

Eval for patterns in adult limb  
and Blastema-specific genes

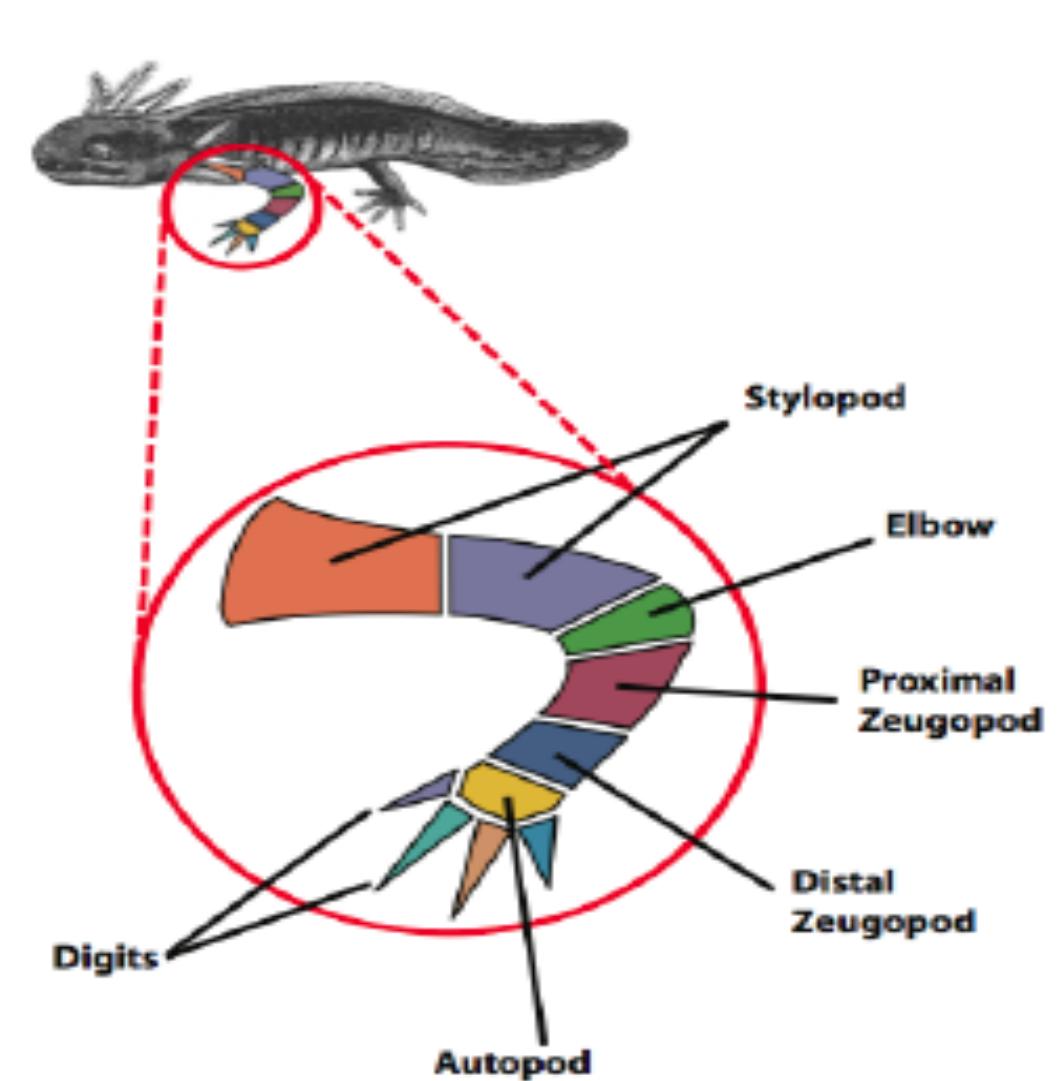
### Challenges

No Genome,

Poor gene annotation

Genome 30GB,

Evol distance to Human/Mouse

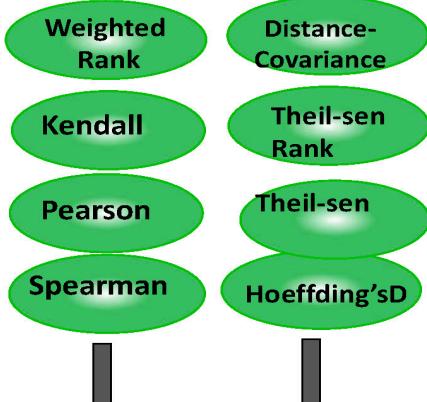
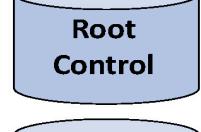
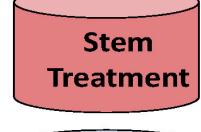
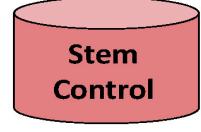


# Cluster 2: Limb development

- 1. HES4: multicellular organismal development
- 2. HOXB1: anterior/posterior pattern formation
- 3. BARHL2: neuron differentiation
- 4. CDX2: caudal type homeobox 2
- 5. CIC
- 6. FOXH1 nodal expression; Lower-limb Deep Vein Thrombosis
- 7. FOXJ1: motile cilia lupus erythematosus
- 8. FOXN1: skin nail development
- 9. GLIS3: limb mesenchymal patterns
- 10. HES6: limb buds
- 11. HESX1: multicellular organismal development
- 12. HNF4A: limb mesoderm
- 13. HOXA5: skeletal and organ morphogenesis
- 14. HOXB6: Development
- 15. OLIG1: oligodendroglial differentiation
- 16. PAX2:
- 17. PLAGL1
- 18. PPARGC1B
- RIN3
- 20. SCML4
- 21. T: anterior/posterior pattern formation  
tail morphogenesis notochord development
- 22. TAF6L
- 23. TAL2
- 24. TGIF2
- 25. TRAF4: development of the axial skeleton  
the closure of the neural tube
- 26. ZBTB16
- 27. ZBTB4
- 28. ZFPM2: hematopoiesis and cardiogenesis
- 29. ZNF154

# Collaborative Network Pipeline

High-throughput  
Gene Expression  
Data

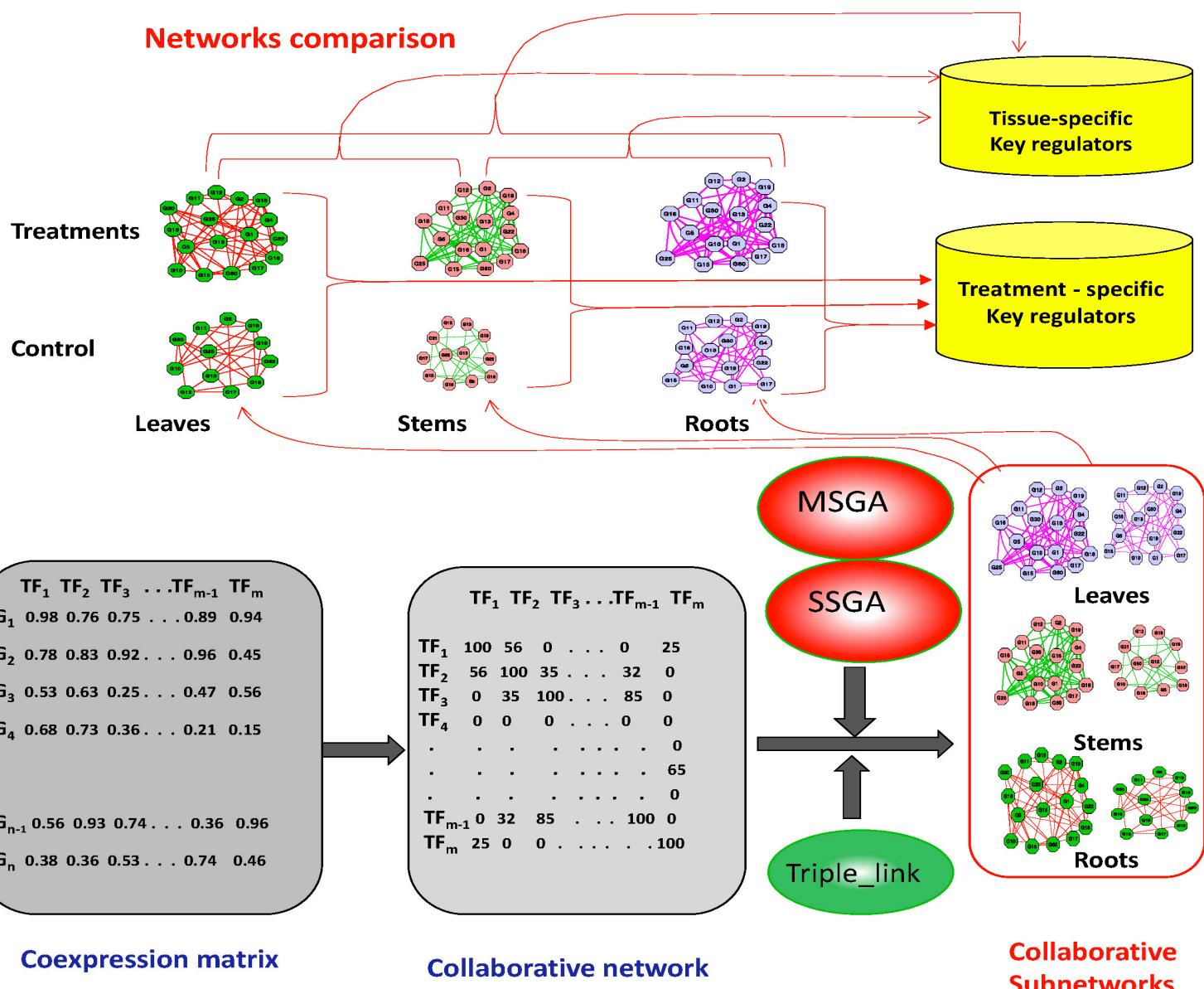


Genome-Wide  
Coexpression Analysis

	$TF_1$	$TF_2$	$TF_3$	$\dots$	$TF_{m-1}$	$TF_m$
$G_1$	0.98	0.76	0.75	$\dots$	0.89	0.94
$G_2$	0.78	0.83	0.92	$\dots$	0.96	0.45
$G_3$	0.53	0.63	0.25	$\dots$	0.47	0.56
$G_4$	0.68	0.73	0.36	$\dots$	0.21	0.15
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$
$G_{n-1}$	0.56	0.93	0.74	$\dots$	0.36	0.96
$G_n$	0.38	0.36	0.53	$\dots$	0.74	0.46

Coexpression matrix

Ji et al. *Scientific Reports* (2017) and Nie et al 2011



## Summary for Collaborative Gene Network:

1. Construction of collaborative subnetworks of TFs is novel approach for identifying collaborative TFs that govern a biological process or complex trait using gene expression data. Compared to markers-assistant approaches for identifying genes that control a trait, collaborative subnetwork construction is a more direct one. It appears to be highly efficient because it is based on the footprints of regulatory events in gene expression profiles rather than the information encoded in DNA sequences, for instance, genome-wide gene association by SNPs.
2. As gene expression data explode we believe construction of collaborative network will find its way into a wide range of applications for novel biological knowledge discovery.