

# **Genome-free mRNA-Seq Data Analysis**

**BayPass:** Outlier testing with SNPs

**WGCNA:** Gene co-expression network analysis

# Genome-free mRNA-Seq Data Analysis

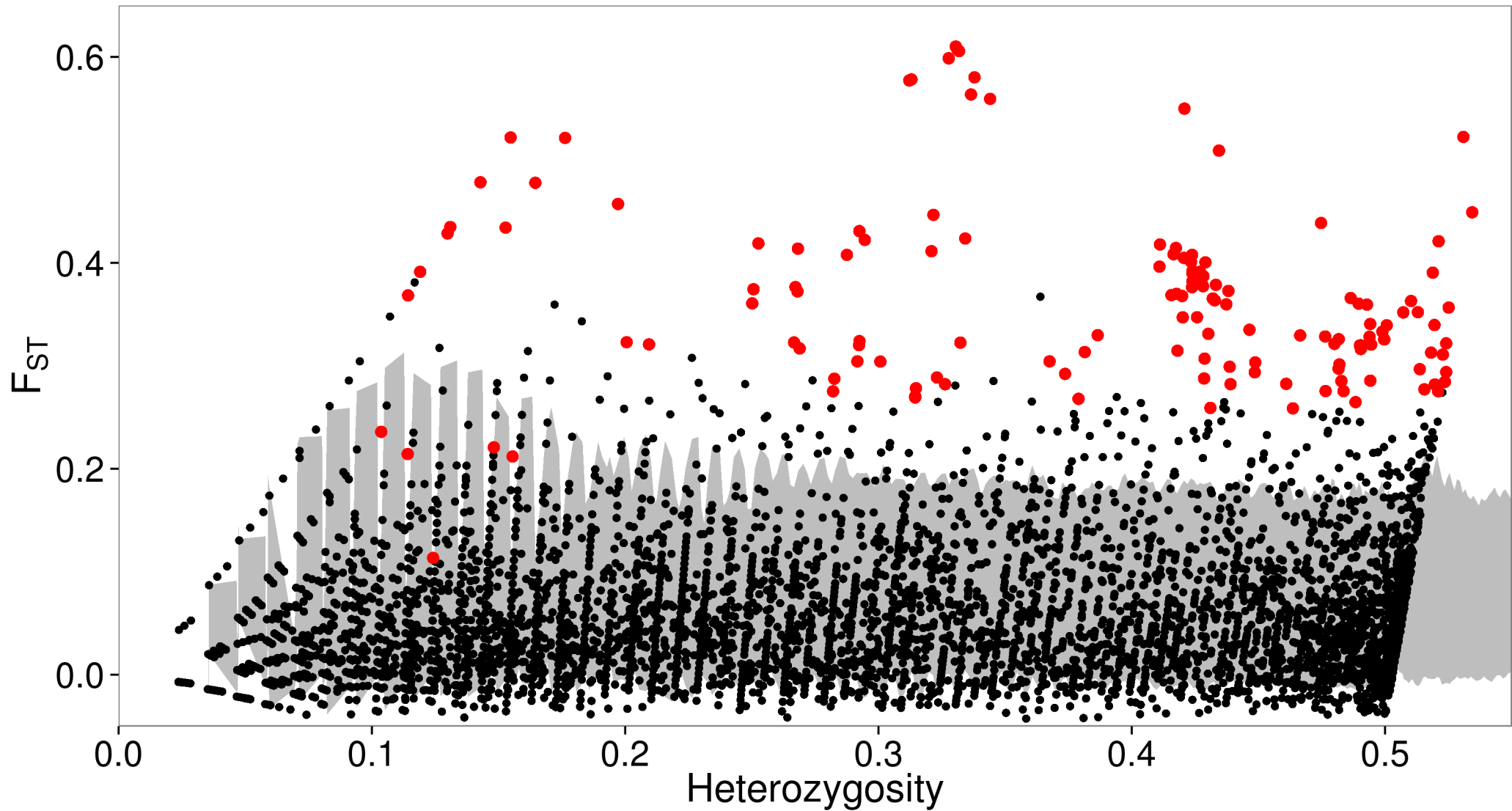
**BayPass:** Outlier testing with SNPs

One species  
Multiple populations  
Biallelic SNPs

Which SNPs might be under selection?

# $F_{ST}$ -based Outlier Testing:

*Test for loci with unusually high  $F_{ST}$*



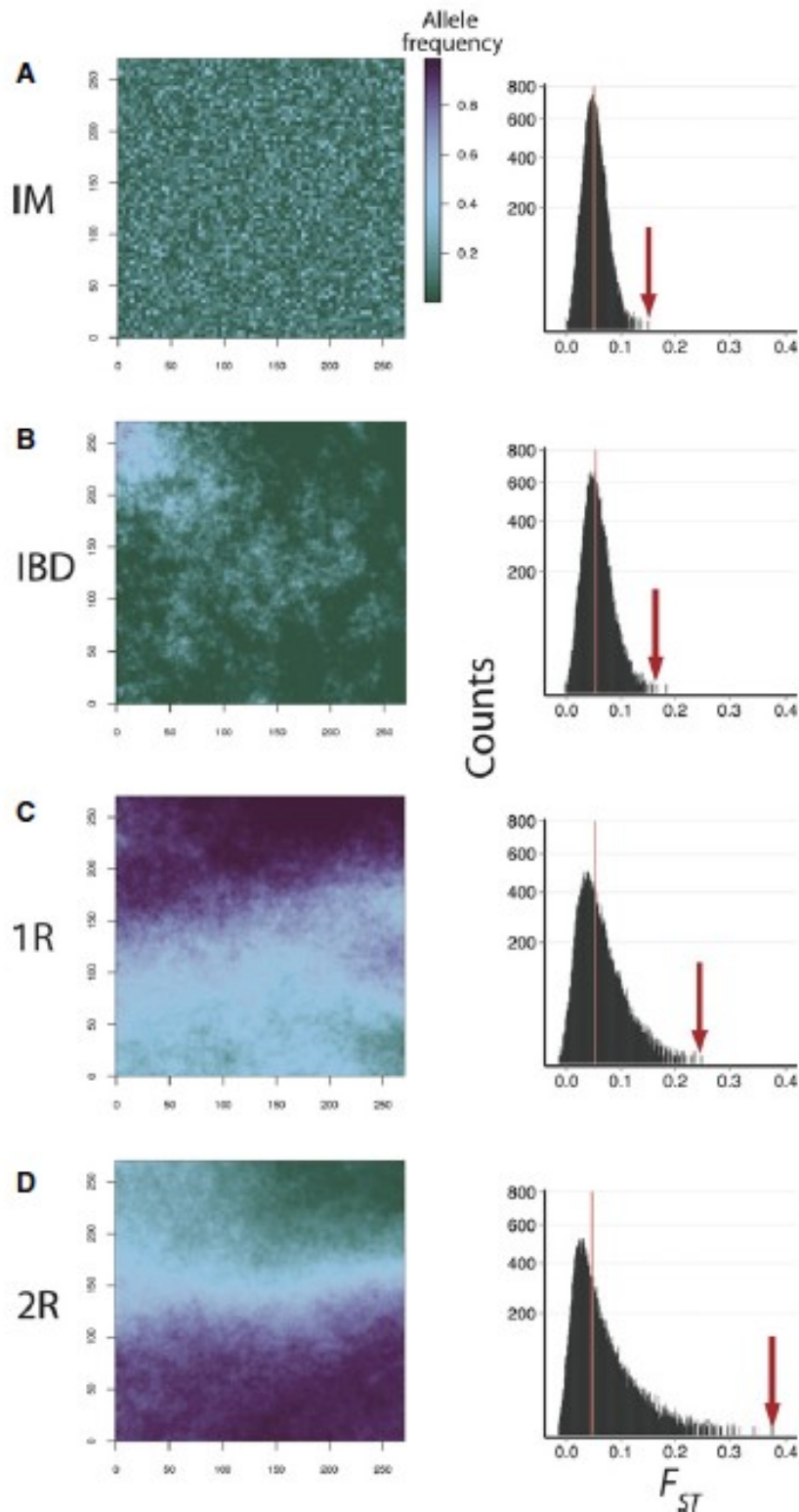
# $F_{ST}$ -based Outlier Testing:

*Test for loci with unusually high  $F_{ST}$*

**Assumes that neutral loci have approximately the same  $F_{ST}$**

“However, in subdivided populations, by chance the measured  $F_{ST}$  can differ substantially from this expectation, causing even neutral genes to vary, sometimes substantially, in their  $F_{ST}$ s.  $F_{ST}$  outlier tests attempt to account for this neutral variation in  $F_{ST}$  and determine which loci have  $F_{ST}$  large enough or small enough to show significant evidence of selection. The challenge with outlier tests is to identify how much variation in  $F_{ST}$  among loci would be expected (i.e. the null distribution of  $F_{ST}$ ) in the absence of selection.”

# Neutral SNP under different scenarios



# $F_{ST}$ -based Outlier Testing

Older methods: BayeScan, FDIST2

- ♦ Assume evolutionary independence
  - ♦ Null distribution assumes specific demographic history
  - ♦ Independent divergence from a common ancestor
- ♦ Perform poorly w/ IDB, expansion, migration, etc.
- ♦ May give many false positives

# $F_{ST}$ -based Outlier Testing

Newer methods: BayEnv, BayPass

- Account for population structure
  - Estimate coancestry/covariance among populations
  - $X^T X$ :  $F_{ST}$  analog standardized by among-pop covariance
- Perform better in non-equilibrium situations
- Test for covariance with environmental factors

# $F_{ST}$ -based Outlier Testing

Newer methods: BayEnv, **BayPass**

- Account for population structure
  - Estimate coancestry/covariance among populations
  - $X^T X$ :  $F_{ST}$  analog standardized by among-pop covariance
- Perform better in non-equilibrium situations
- Test for covariance with environmental factors



# BayPass Input

*Fortran*

Population-level allele counts (required)

*May be in chromosomal order, if known*

```
--- file begins here ---
```

```
81 19 86 14 2 98 8 92 32 68 23 77
```

```
89 11 81 19 9 91 1 99 27 73 27 73
```

```
89 11 91 9 0 0 15 85 77 23 80 20
```

Covariate data file (optional)

```
--- file begins here ---
```

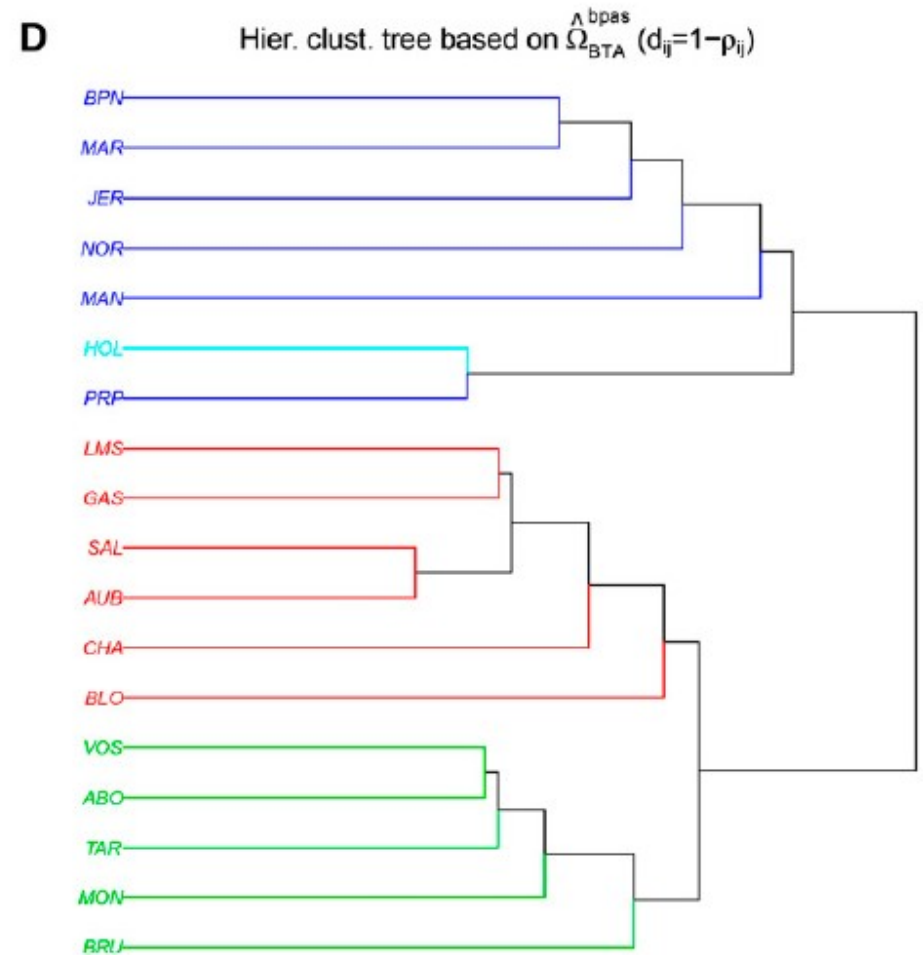
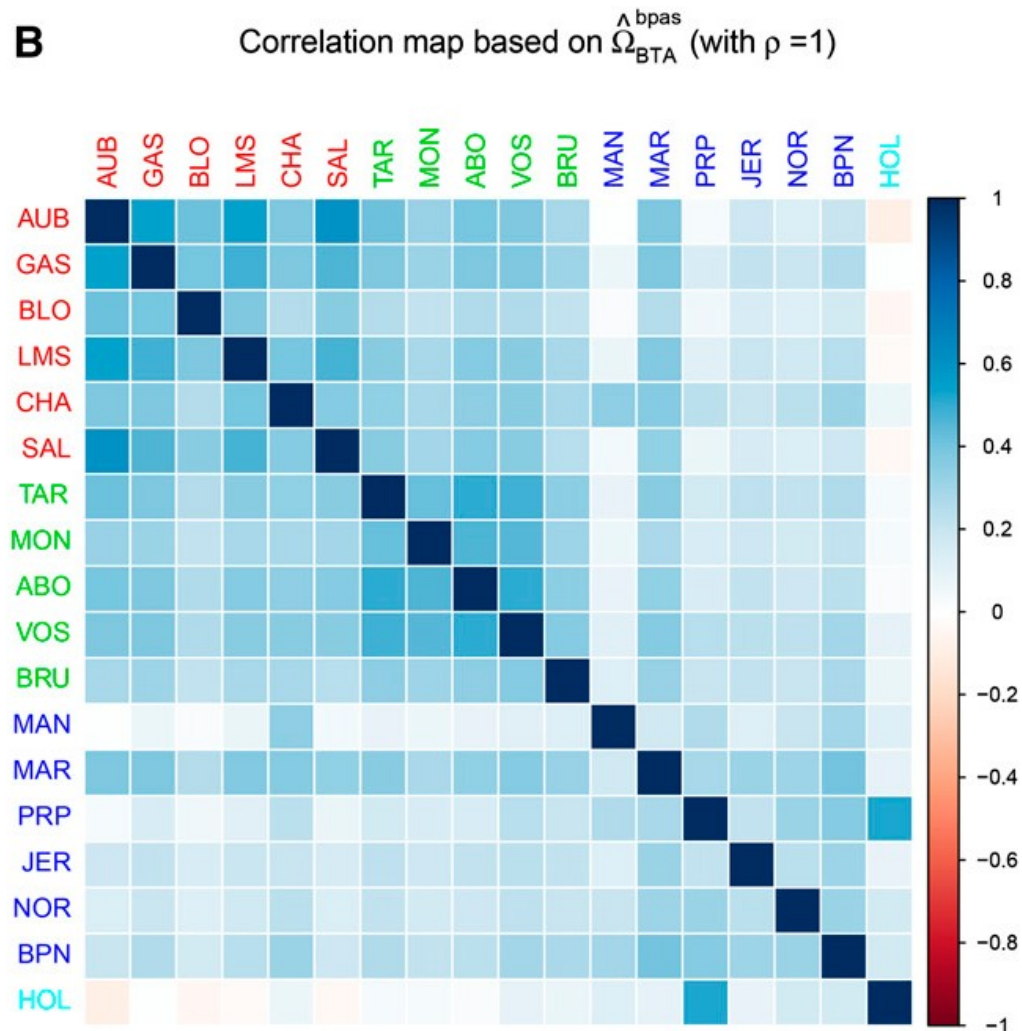
```
150 1500 800 300 200 2500
```

```
181.5 172.6 152.3 191.8 154.2 166.8
```

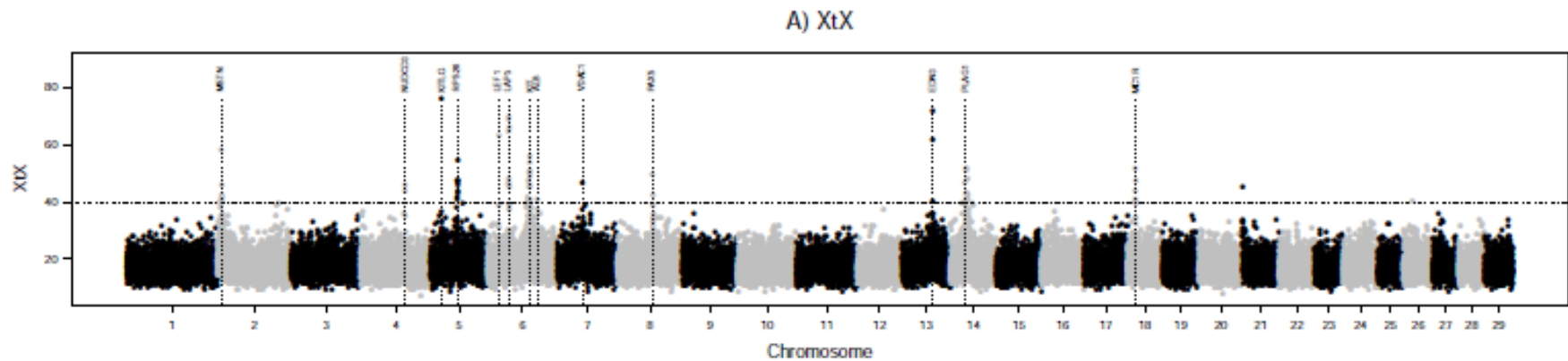
```
1 1 0 0 1 1
```

```
0.1 0.8 -1.15 1.6 0.02 -0.5
```

# Core model: Estimate covariance

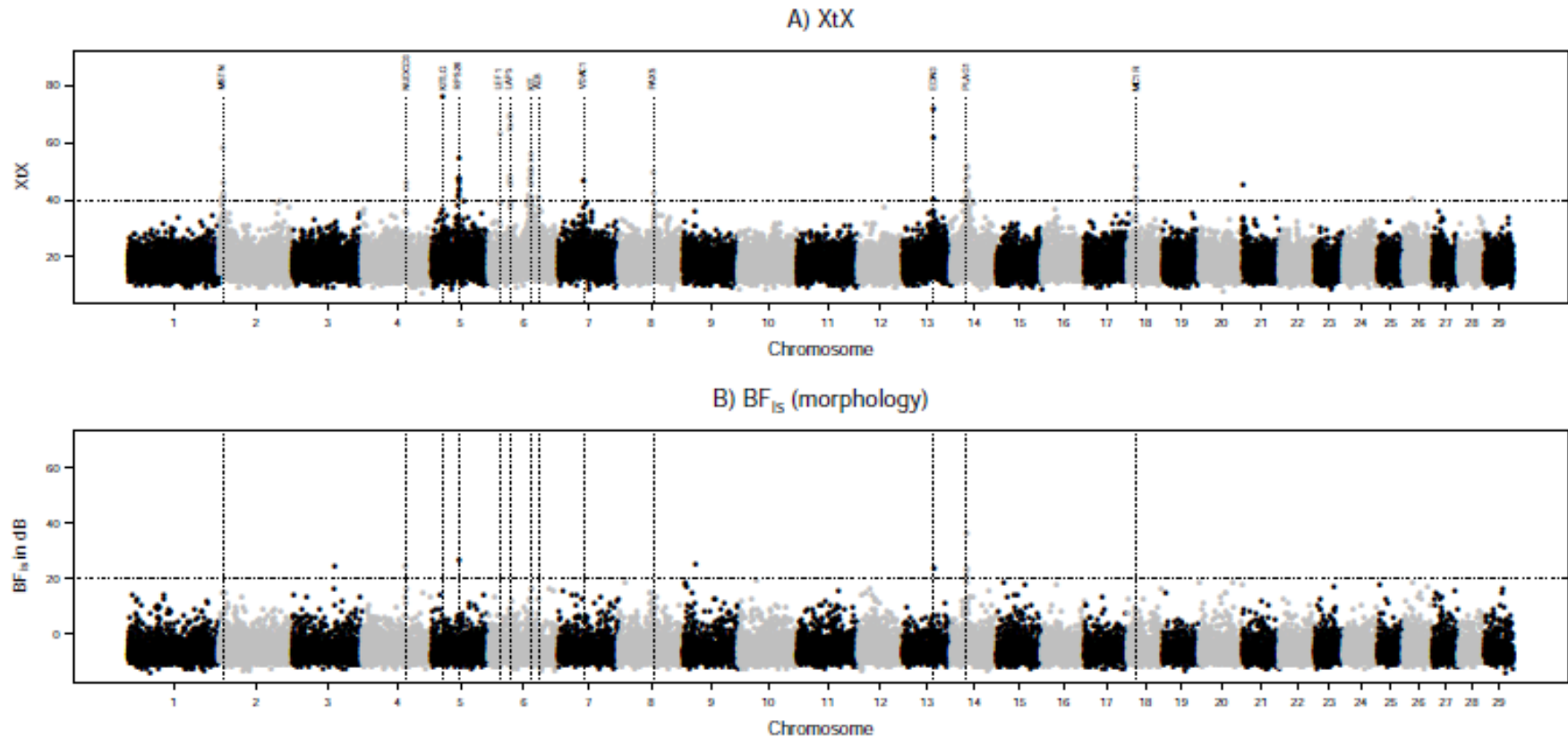


# Core model: Calculate $X^T X$ for all SNPs



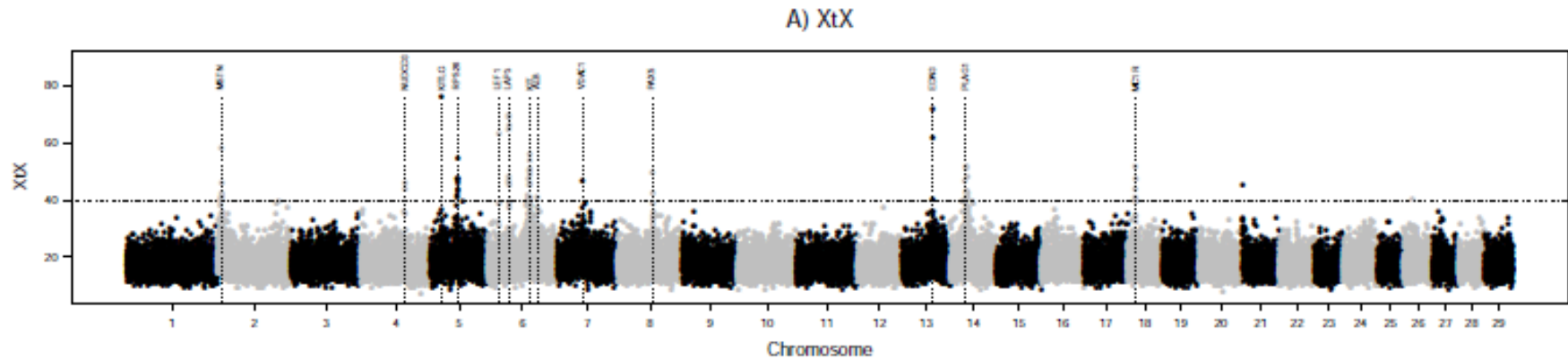
# Calibrate & determine significance via simulated data

# Covariate model: Importance Sampling

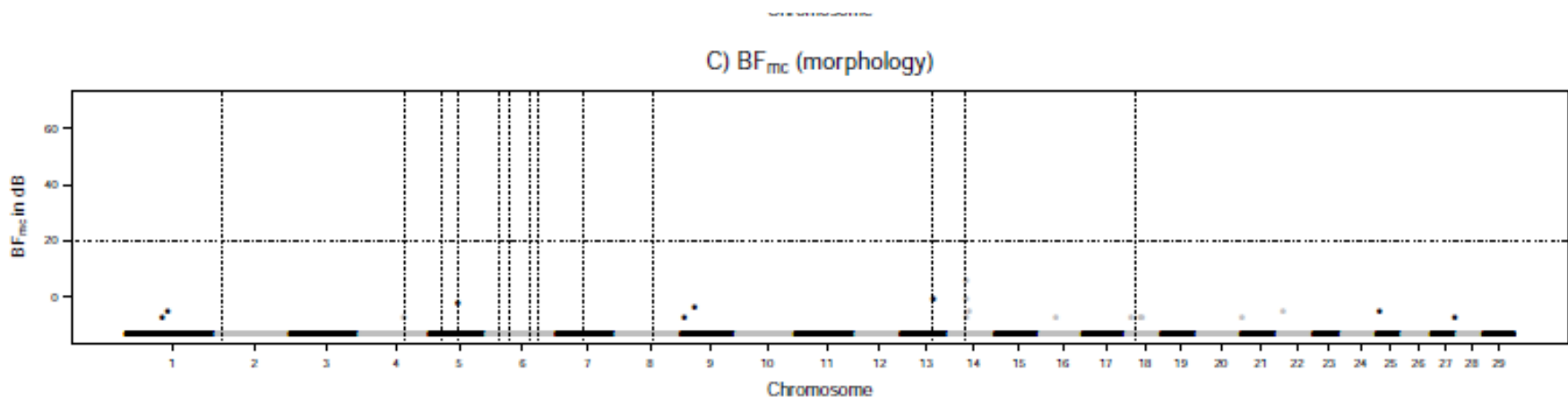


# Estimate Bayes Factor for SNP-covariate association

# Covariate model: Auxiliary Covariate



# Estimate Bayes Factor for SNP-covariate association, allowing incorporation of marker position



# BayPass: What is it good for?

Identifying selection between non-equilibrium populations, with or without marker position information

# Genome-free mRNA-Seq Data Analysis

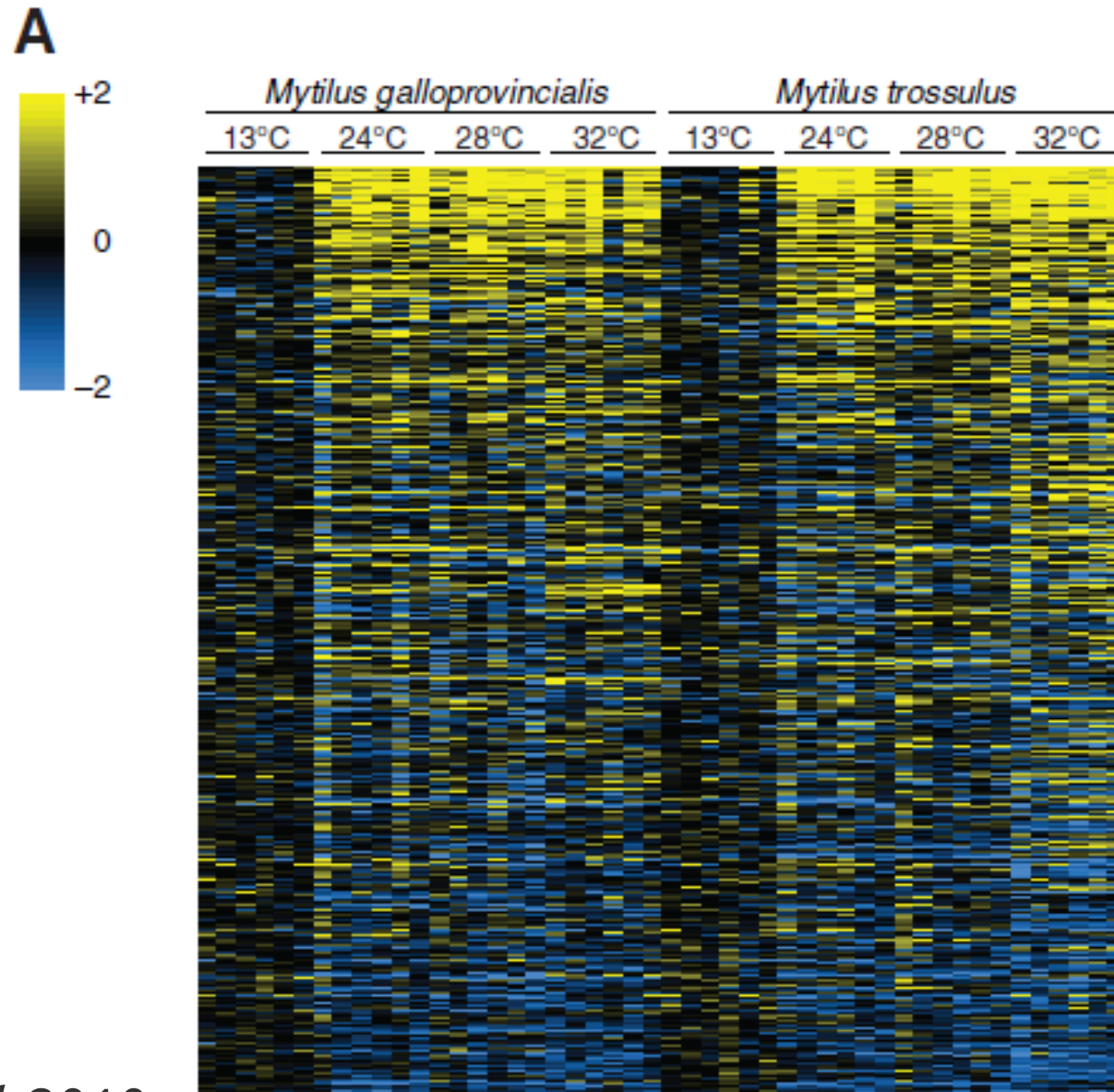
**WGCNA:** Gene co-expression network analysis

One or more species  
Multiple conditions  
Normalized expression data

Which genes are changing expression?

# Gene Expression Analysis

*Which genes are differentially expressed between treatments?*





# Gene Expression Analysis

Gene-by-gene methods: DESeq, EdgeR

- Assess each gene independently
- Enrichment analysis to identify important pathways / processes

**Lots of natural variation, 10K+ tests – what's really important?**

# Gene Expression Analysis

Co-expression network analysis: WGCNA, MEGENA

- Identify modules of co-expressed genes
- Associate modules with traits/conditions of interest
- Enrichment analysis within modules to identify important pathways / processes

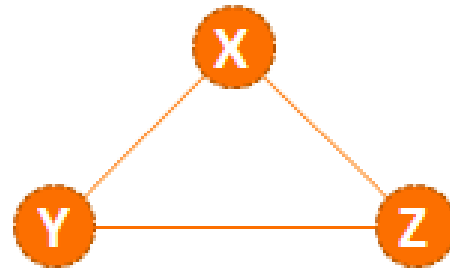
# Gene Expression Analysis

Co-expression network analysis: **WGCNA**, MEGENA

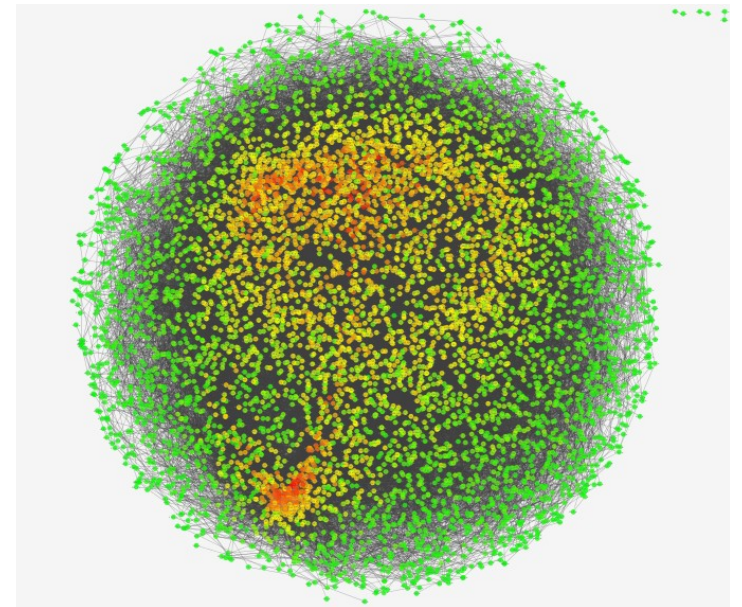
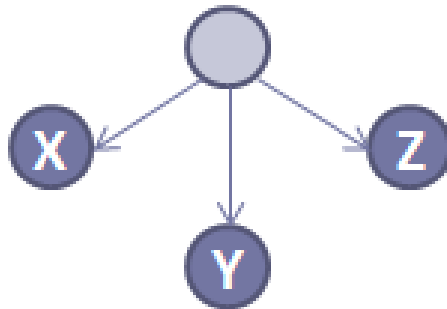
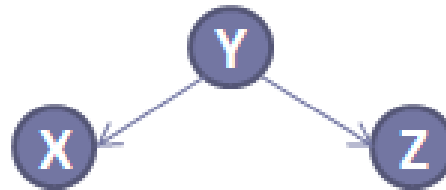
- Identify modules of co-expressed genes
- Associate modules with traits/conditions of interest
- Enrichment analysis within modules to identify important pathways / processes

# Gene Co-expression Analysis

Gene Co-expression



Gene Regulation

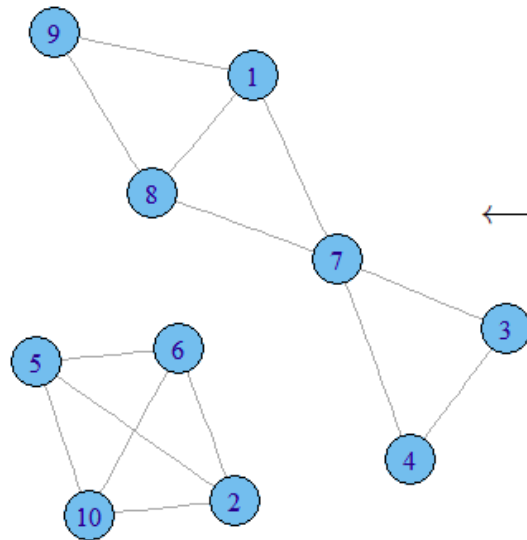


# Gene Co-expression Analysis

	$S_1$	$S_2$	$S_3$		$G_1$	$G_2$	$G_3$	$G_4$	$G_5$	$G_6$	$G_7$	$G_8$	$G_9$	$G_{10}$	
$G_1$	43.26	40.89	5.05	<div><math> r(G_i, G_j) </math></div> <div>Pearson correlation</div>	$G_1$	1.00	0.23	0.61	0.71	0.03	0.35	<b>0.86</b>	<b>1.00</b>	<b>0.97</b>	0.37
$G_2$	166.6	41.87	136.65		$G_2$	0.23	1.00	0.63	0.52	<b>0.98</b>	<b>0.99</b>	0.29	0.30	0.46	<b>0.99</b>
$G_3$	12.53	39.55	42.09		$G_3$	0.61	0.63	1.00	<b>0.99</b>	0.77	0.53	<b>0.93</b>	0.56	0.41	0.51
$G_4$	28.77	191.92	236.56		$G_4$	0.71	0.52	<b>0.99</b>	1.00	0.69	0.41	<b>0.97</b>	0.66	0.52	0.40
$G_5$	114.7	79.7	99.76		$G_5$	0.03	<b>0.98</b>	0.77	0.69	1.00	<b>0.95</b>	0.48	0.09	0.27	<b>0.94</b>
$G_6$	119.1	80.57	114.59		$G_6$	0.35	<b>0.99</b>	0.53	0.41	<b>0.95</b>	1.00	0.17	0.41	0.57	<b>1.00</b>
$G_7$	118.9	156.69	186.95		$G_7$	0.86	0.29	<b>0.93</b>	<b>0.97</b>	0.48	0.17	1.00	<b>0.83</b>	0.72	0.16
$G_8$	3.76	2.48	136.78		$G_8$	<b>1.00</b>	0.30	0.56	0.66	0.09	0.41	0.83	1.00	<b>0.98</b>	0.42
$G_9$	32.73	11.99	118.8		$G_9$	<b>0.97</b>	0.46	0.41	0.52	0.27	0.57	0.72	<b>0.98</b>	1.00	0.58
$G_{10}$	17.46	56.11	21.41		$G_{10}$	0.37	<b>0.99</b>	0.51	0.40	<b>0.94</b>	<b>1.00</b>	0.16	0.42	0.58	1.00

Gene expression values

Similarity (Co-expression) score



	$G_1$	$G_2$	$G_3$	$G_4$	$G_5$	$G_6$	$G_7$	$G_8$	$G_9$	$G_{10}$
$G_1$	0	0	0	0	0	0	1	1	1	0
$G_2$	0	0	0	0	1	1	0	0	0	1
$G_3$	0	0	0	1	0	0	1	0	0	0
$G_4$	0	0	1	0	0	0	1	0	0	0
$G_5$	0	1	0	0	0	1	0	0	0	1
$G_6$	0	1	0	0	1	0	0	0	0	1
$G_7$	1	0	1	1	0	0	0	1	0	0
$G_8$	1	0	0	0	0	0	1	0	1	0
$G_9$	1	0	0	0	0	0	0	1	0	0
$G_{10}$	0	1	0	0	1	1	0	0	0	0

$|r(G_i, G_j)| \geq 0.8$   
 Significance threshold

Network adjacency matrix

# WGCNA Input

*R package*

Normalized gene expression data (required)

*May be transformed, but does not have to be*

	F2_2	F2_3	F2_14	F2_15	F2_19	F2_20	F2_23
MMT00000044	-0.01810	0.0642	6.44e-05	-0.05800	0.04830	-0.15197410	-0.00129
MMT00000046	-0.07730	-0.0297	1.12e-01	-0.05890	0.04430	-0.09380000	0.09340
MMT00000051	-0.02260	0.0617	-1.29e-01	0.08710	-0.11500	-0.06502607	0.00249
MMT00000076	-0.00924	-0.1450	2.87e-02	-0.04390	0.00425	-0.23610000	-0.06900
MMT00000080	-0.04870	0.0582	-4.83e-02	-0.03710	0.02510	0.08504274	0.04450
MMT00000102	0.17600	-0.1890	-6.50e-02	-0.00846	-0.00574	-0.01807182	-0.12500
MMT00000149	0.07680	0.1860	2.14e-01	0.12000	0.02100	0.06222751	0.22600
MMT00000150	0.11000	0.1770	1.22e-01	0.10700	0.11000	0.05107606	0.05600

Trait data file (optional)

	weight_g	length_cm	ab_fat	other_fat	total_fat	X100xfat_weight	Trigly
F2_2	38.0	10.5	3.81	2.78	6.59	17.342105	14
F2_3	33.5	10.8	1.70	2.05	3.75	11.194030	109
F2_14	33.9	10.0	1.29	1.67	2.96	8.731563	2
F2_15	44.3	10.3	3.62	3.34	6.96	15.711061	71
F2_19	32.9	9.7	2.08	1.85	3.93	11.945289	55
F2_20	44.8	10.3	3.72	3.20	6.92	15.446429	34

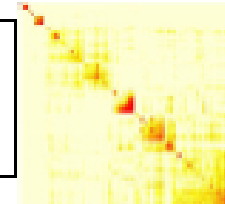
Gene annotation information (optional)

# WGCNA Overview

## Construct a gene co-expression network

**Rationale:** make use of interaction patterns among genes

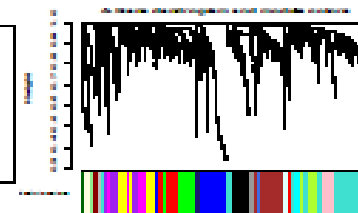
**Tools:** correlation as a measure of co-expression



## Identify modules

**Rationale:** module (pathway) based analysis

**Tools:** hierarchical clustering, Dynamic Tree Cut

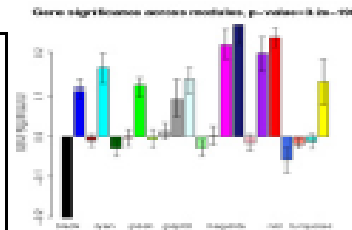


## Relate modules to external information

Array Information: clinical data, SNPs, proteomics

Gene Information: ontology, functional enrichment

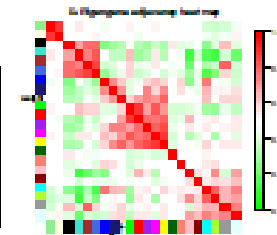
**Rationale:** find biologically interesting modules



## Study module relationships

**Rationale:** biological data reduction, systems-level view

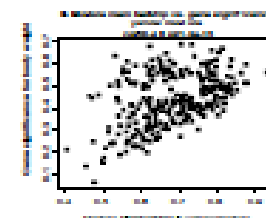
**Tools:** Eigengene Networks



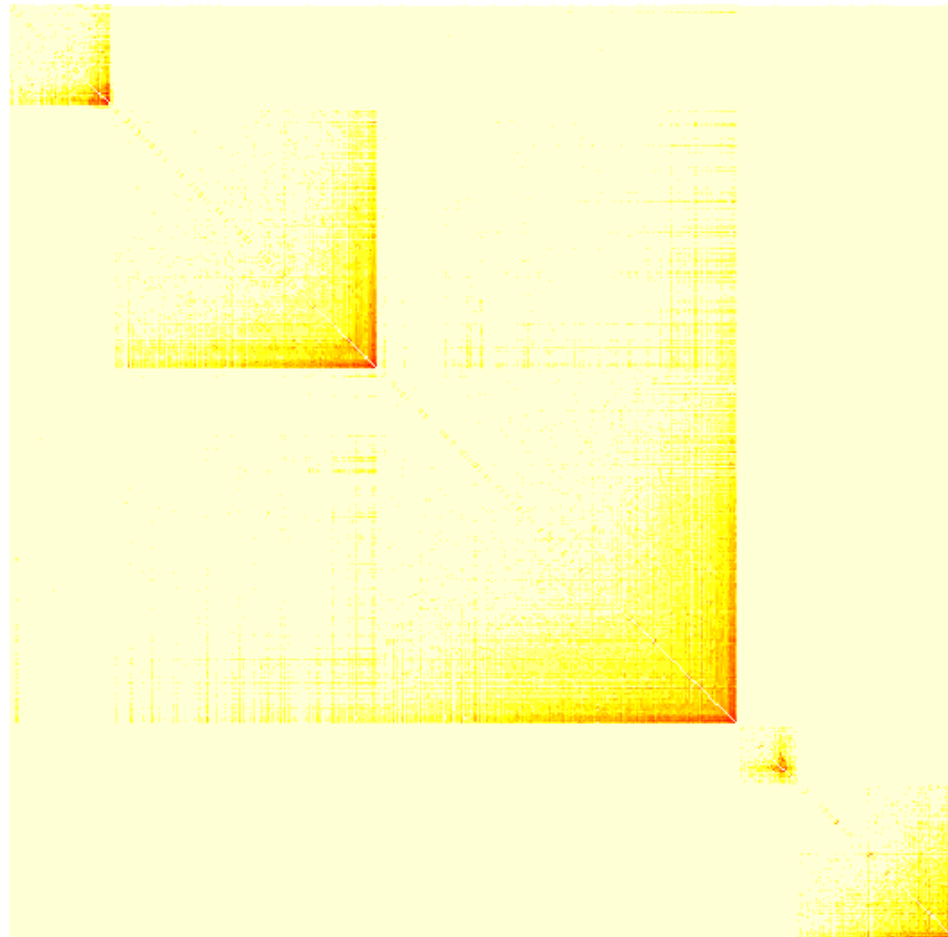
## Find the key drivers in *interesting* modules

**Rationale:** experimental validation, biomarkers

**Tools:** intramodular connectivity, causality testing



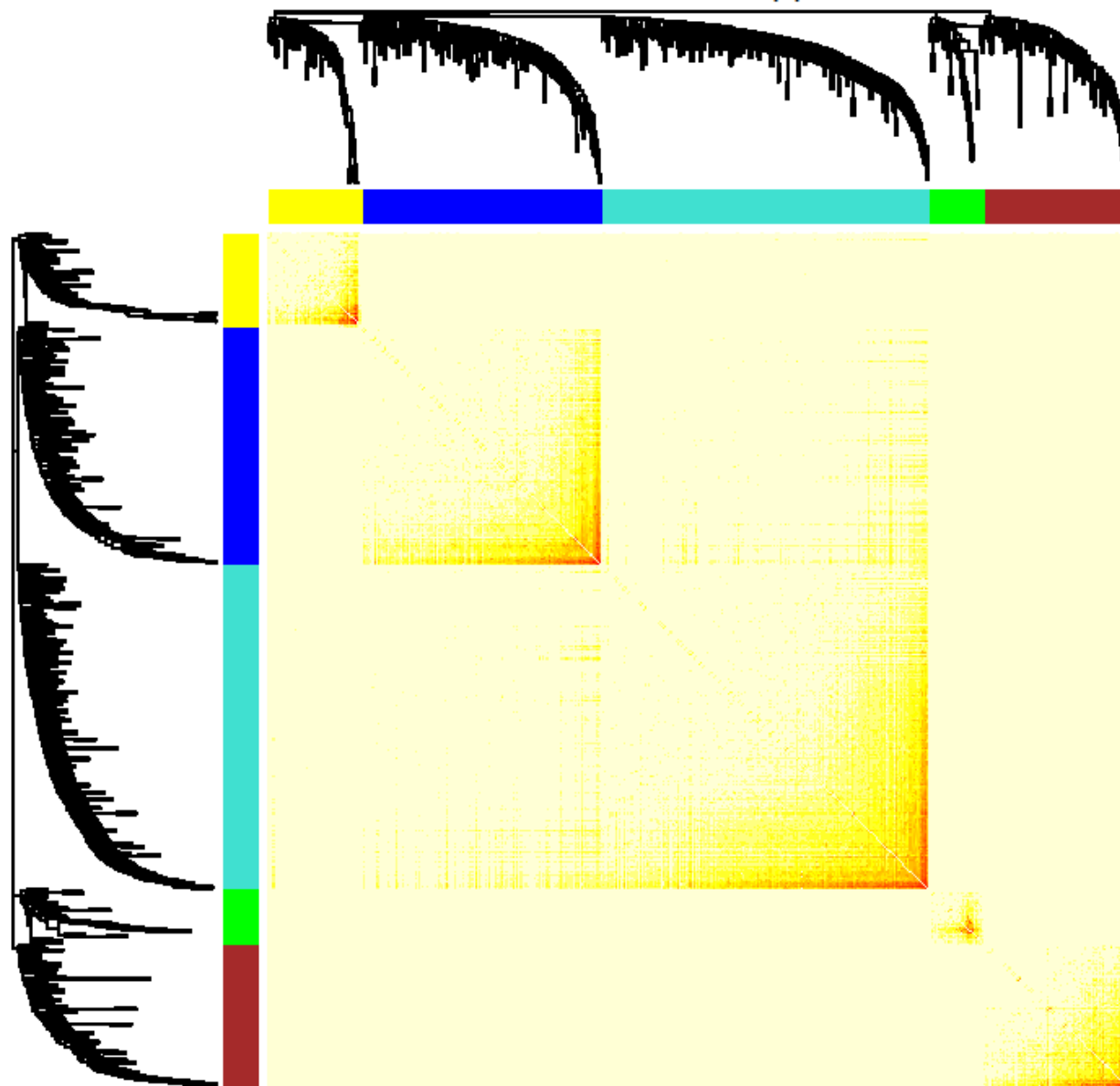
# Step 1: Construct co-expression network



Langfelder & Horvath 2008



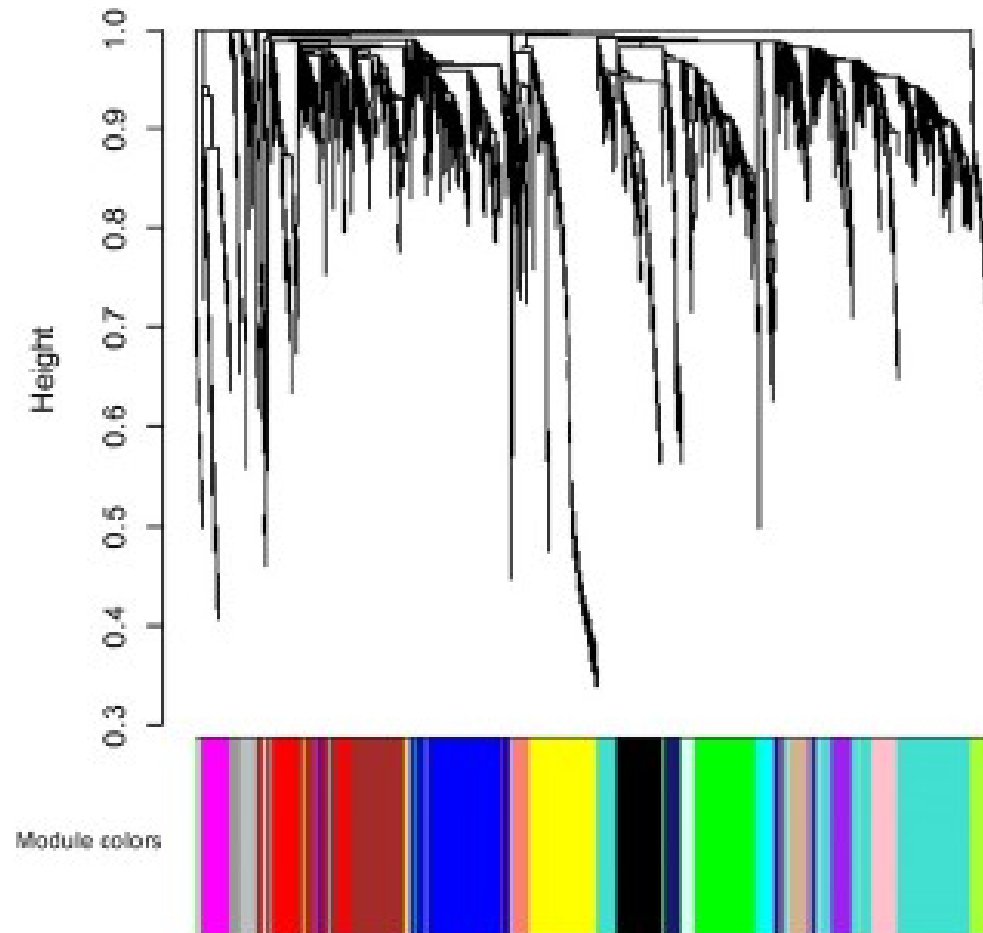
## Step 1: Construct co-expression network



## Langfelder & Horvath 2008

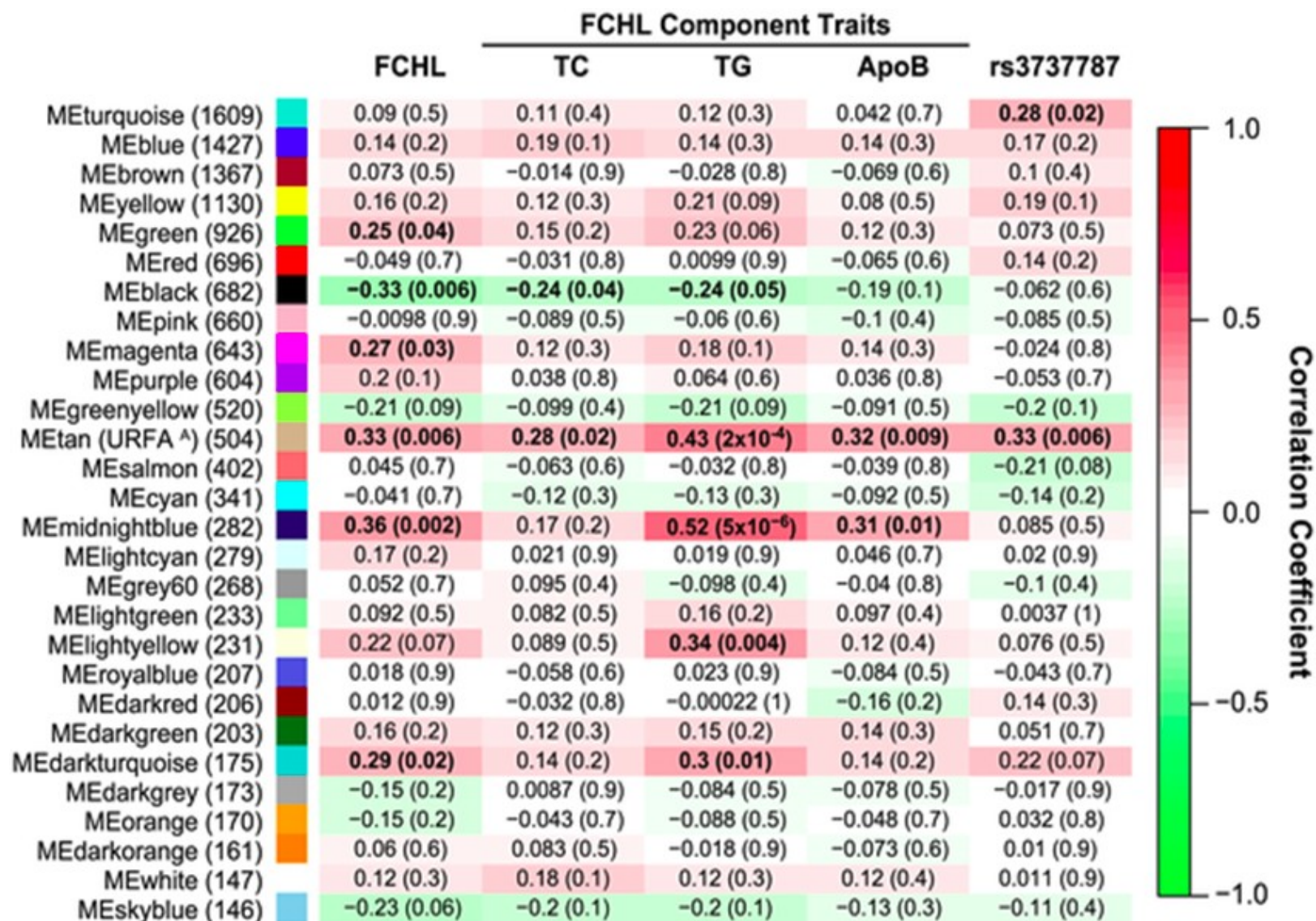
# Step 1: Construct co-expression network

Computationally intensive – can perform all at once (ideal),  
or blockwise with gene subsets

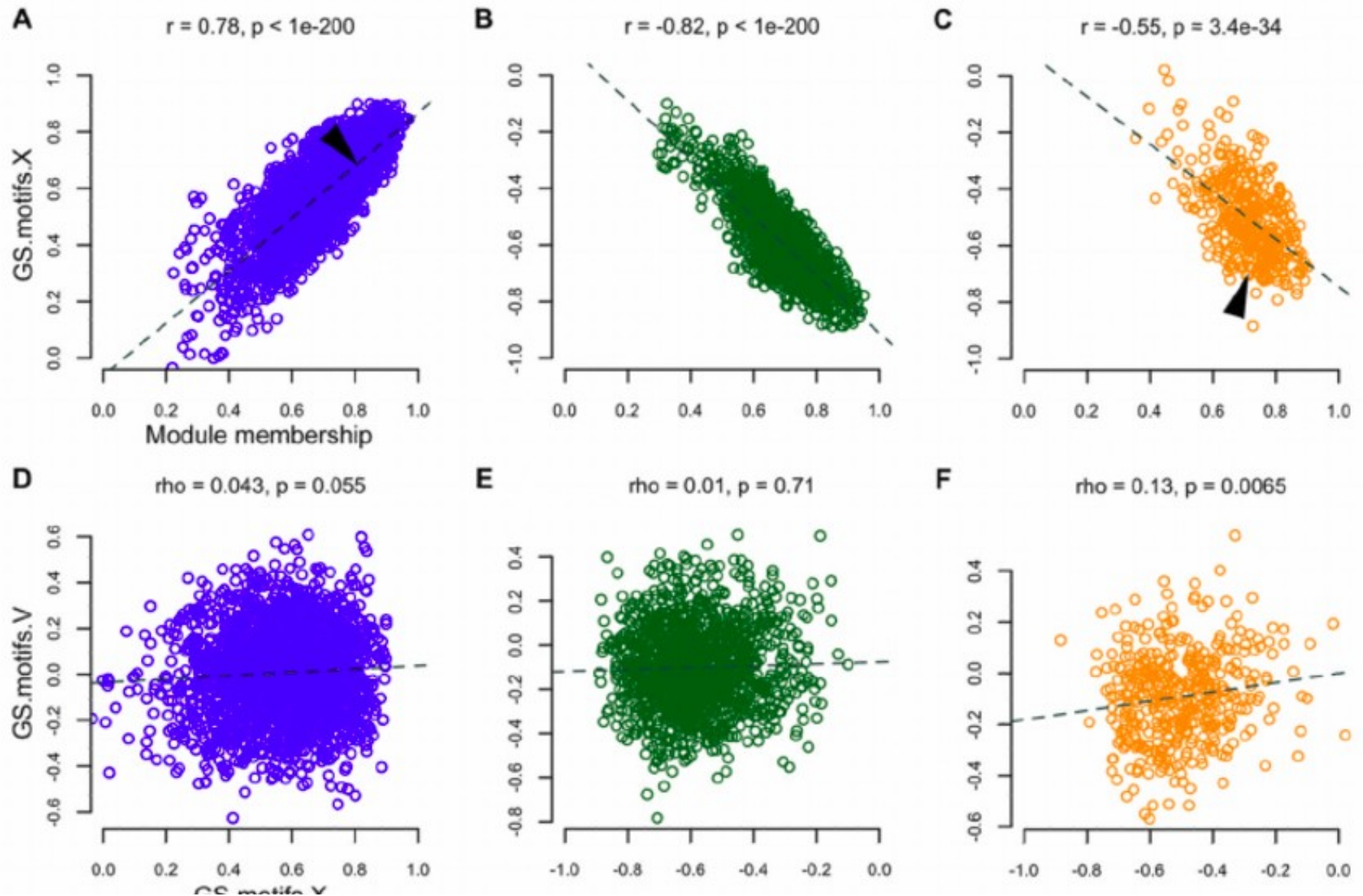


Langfelder & Horvath 2008

# Step 2: Relate Modules with Traits



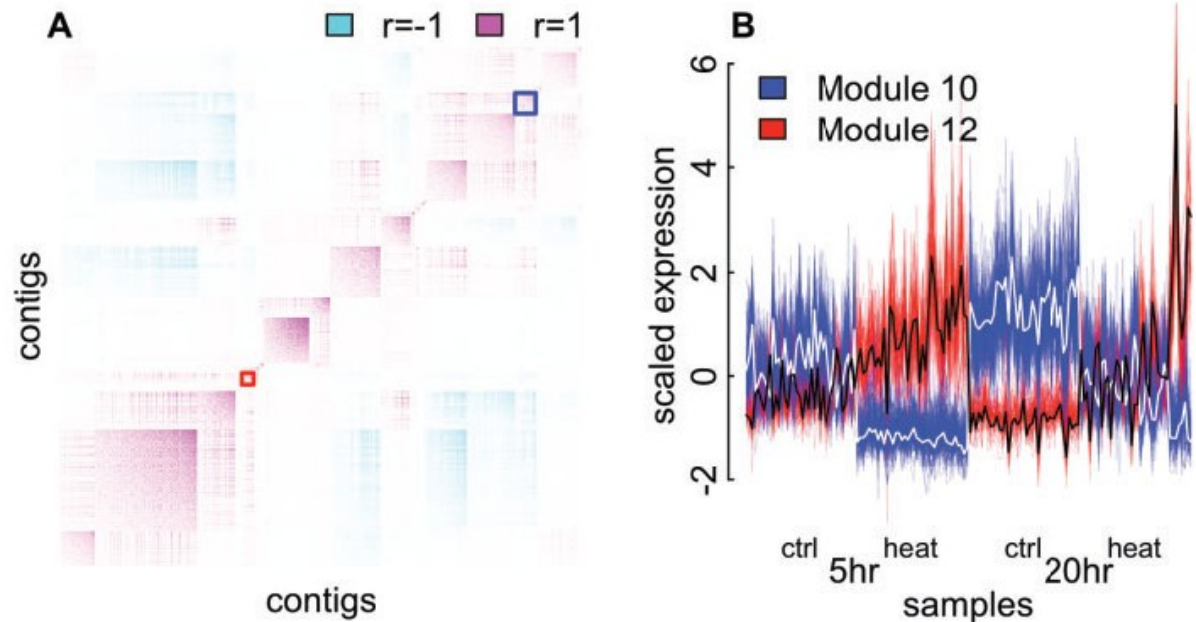
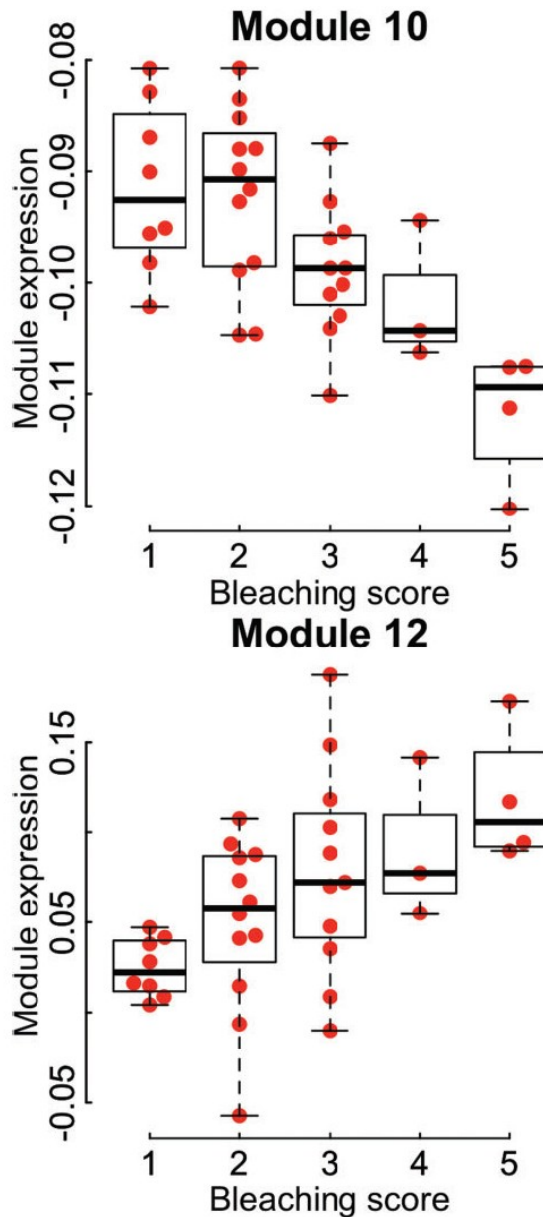
# Step 2: Relate Modules with Traits





# Step 3: Explore Gene Networks

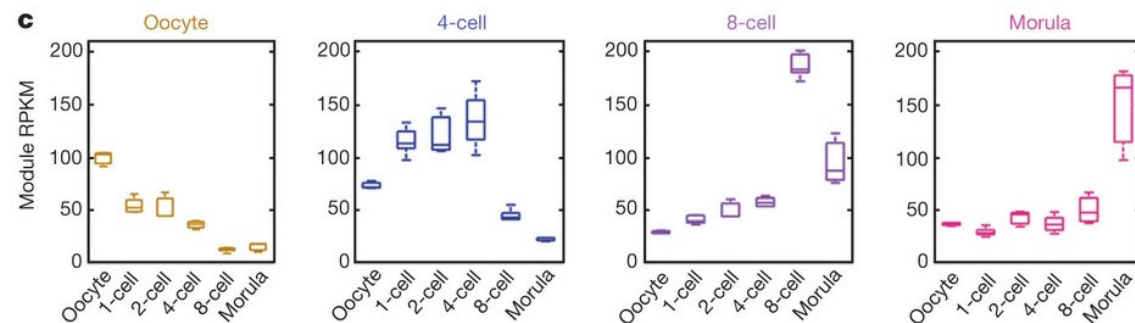
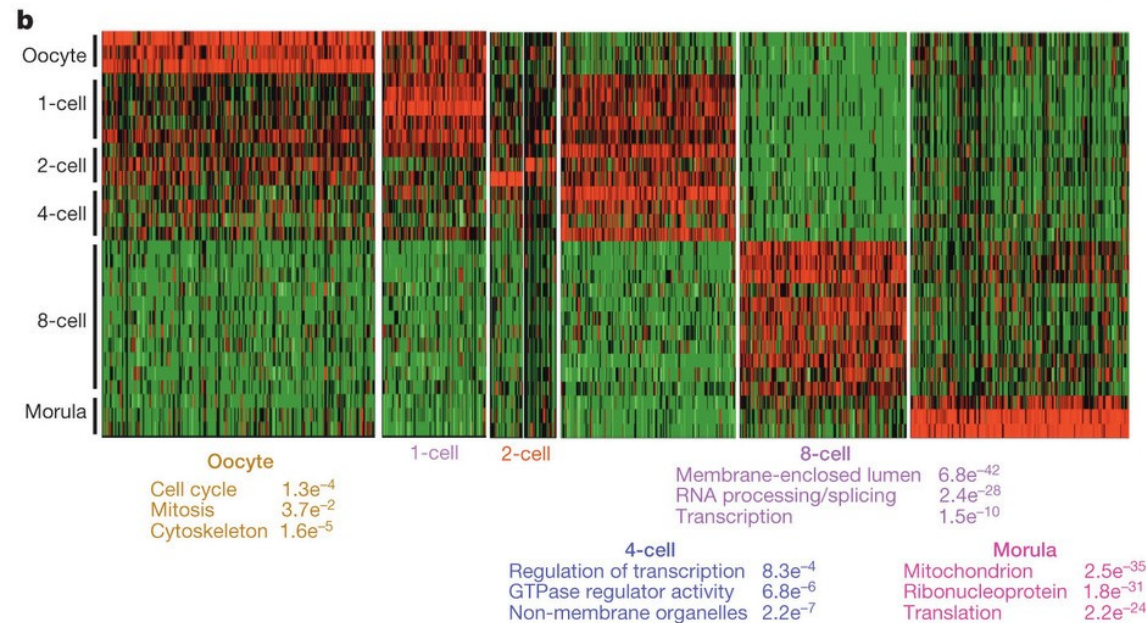
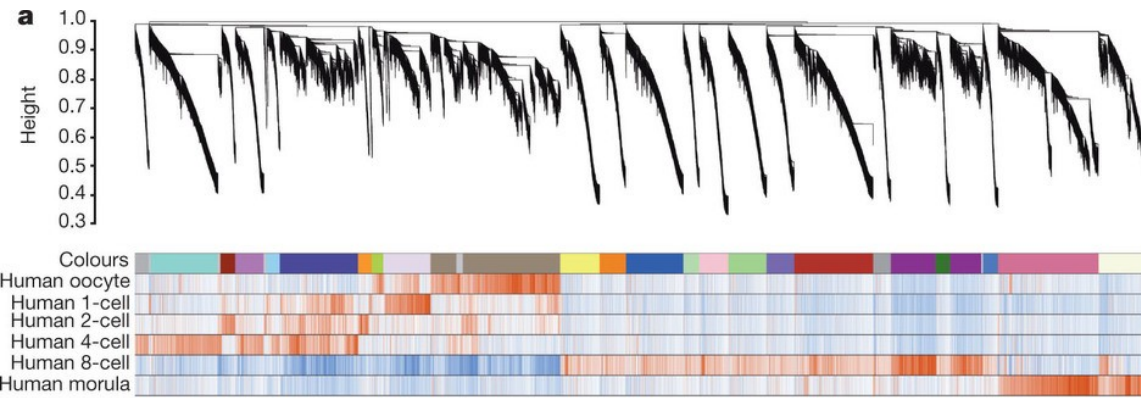
Identify expression patterns associated with a specific biological response.



Rose et al. 2015

# Step 3: Explore Gene Networks

Annotate genes  
add GO terms, etc.  
Within-module  
enrichment analyses.

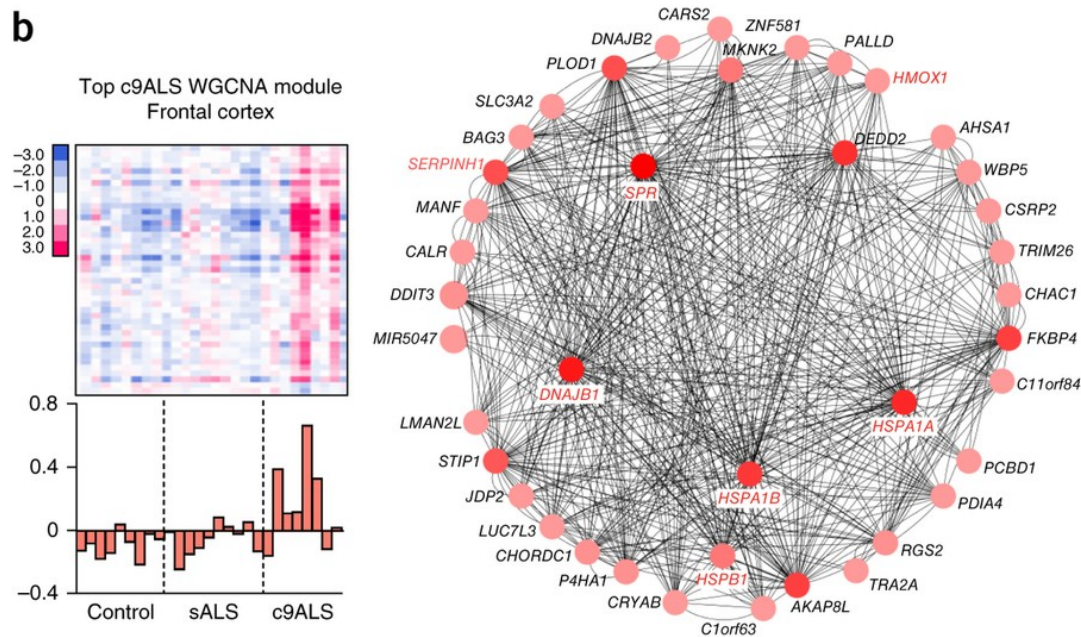
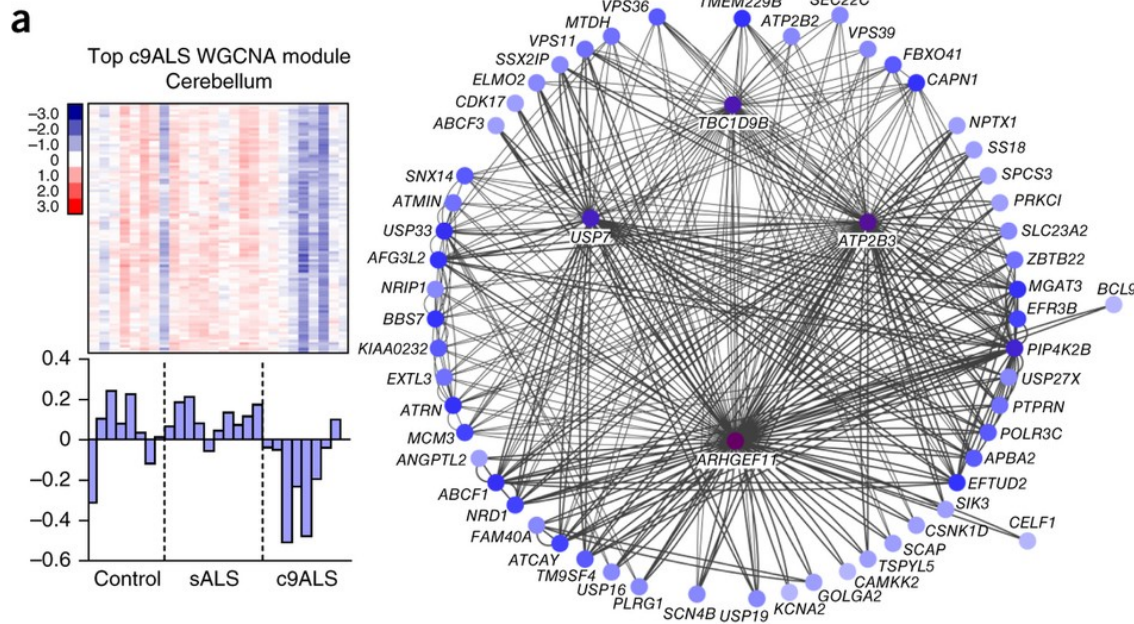


Xue et al. 2013



## Step 3: Explore Gene Networks

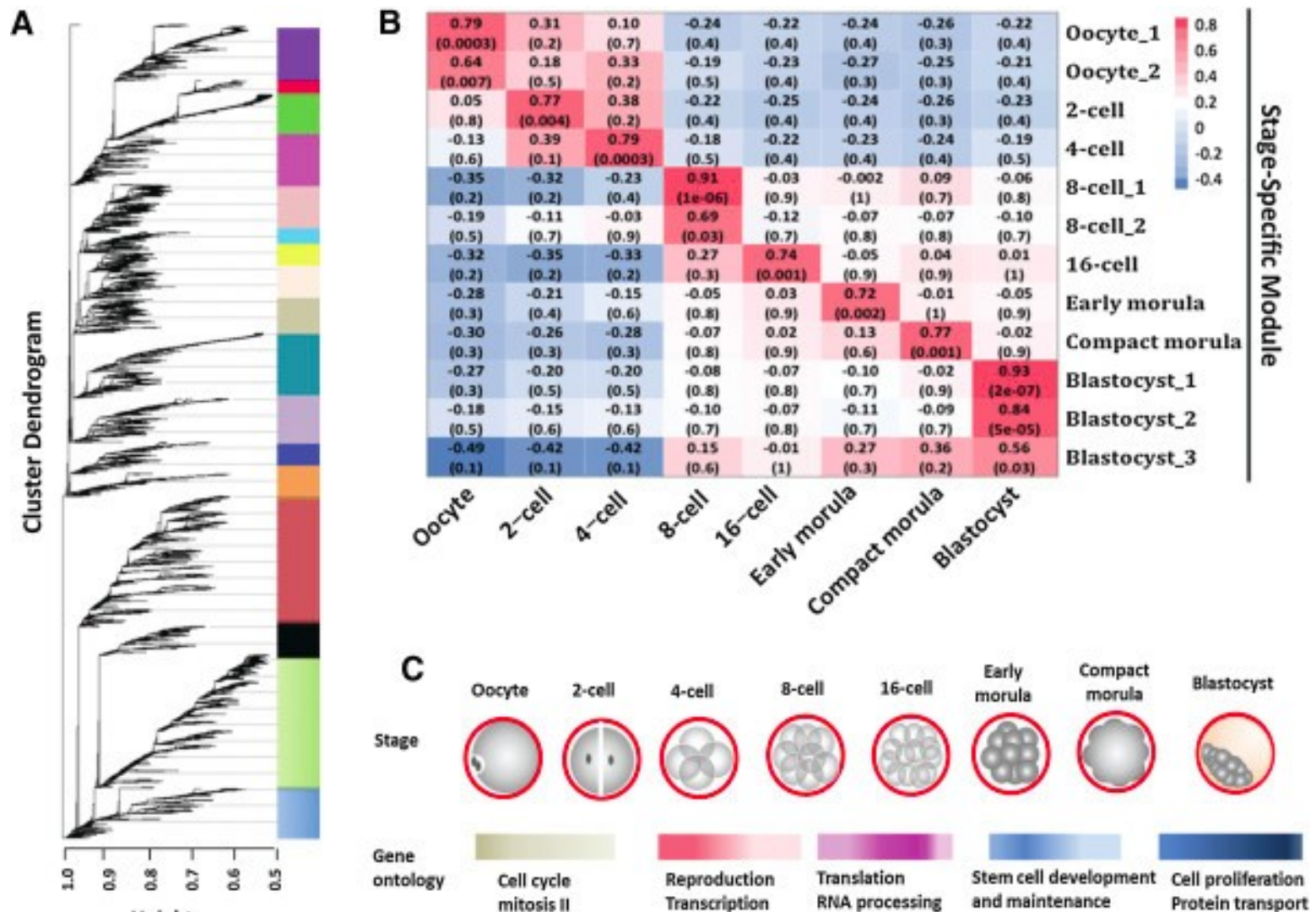
## Plot networks & identify key “hub” genes.



Prudencio et al. 2015

# Step 3: Explore Gene Networks

Compare expression similarity among different conditions.





# **WGCNA: What is it good for?**

Identifying & exploring the gene networks  
underlying regulatory variation.