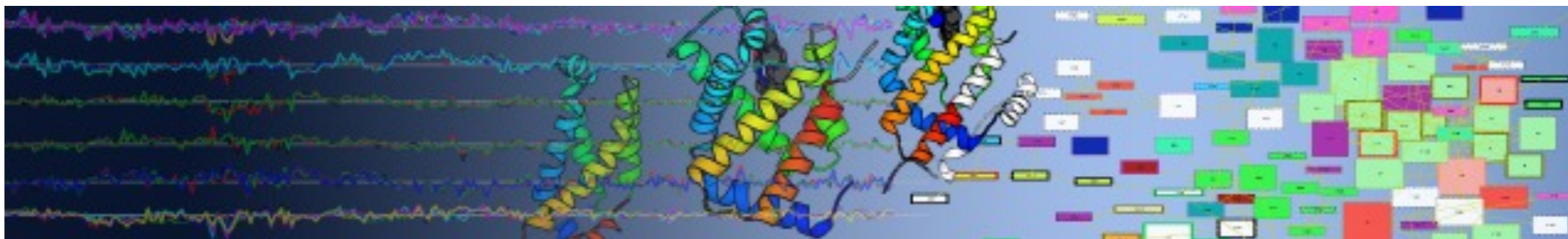


RECONSTRUCTING BIOLOGICAL NETWORKS FROM DATA: PART 1 - cMONKEY



**RICHARD
BONNEAU**

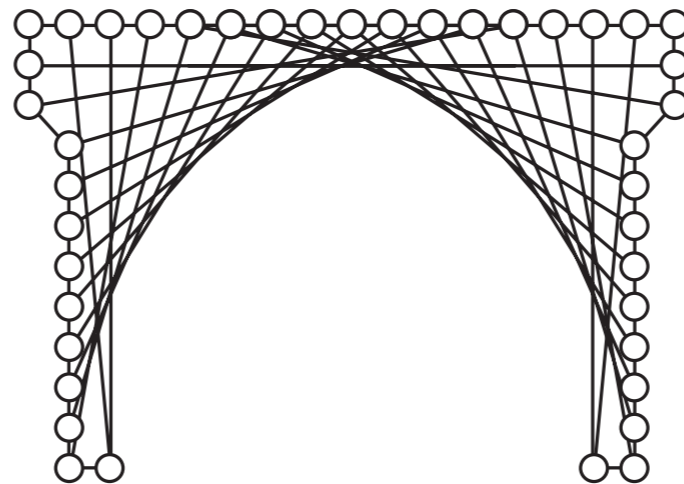
BONNEAU@NYU.EDU

**[HTTP://WWW.CS.NYU.EDU/
~BONNEAU/](http://www.cs.nyu.edu/~bonneau/)**

NEW YORK UNIVERSITY,

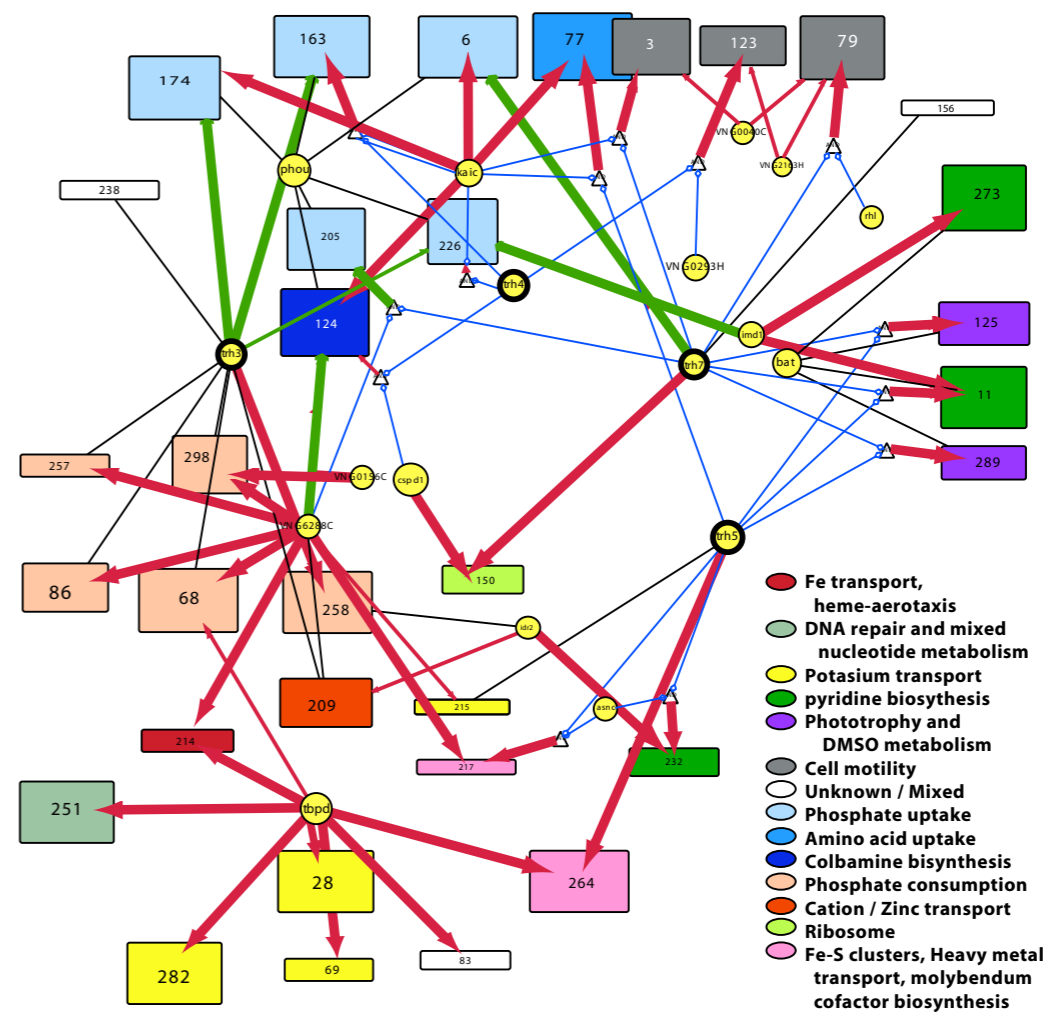
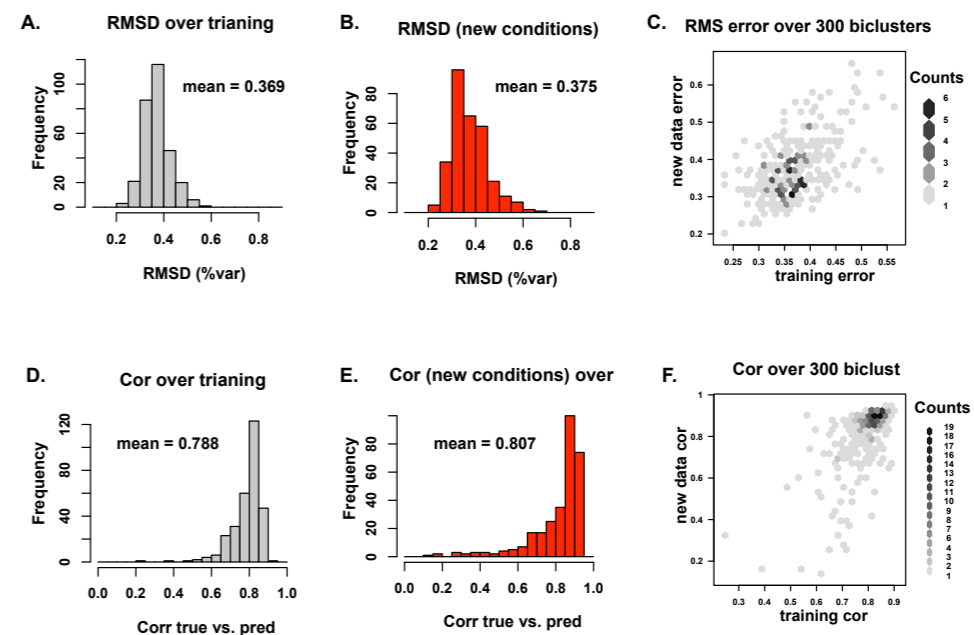
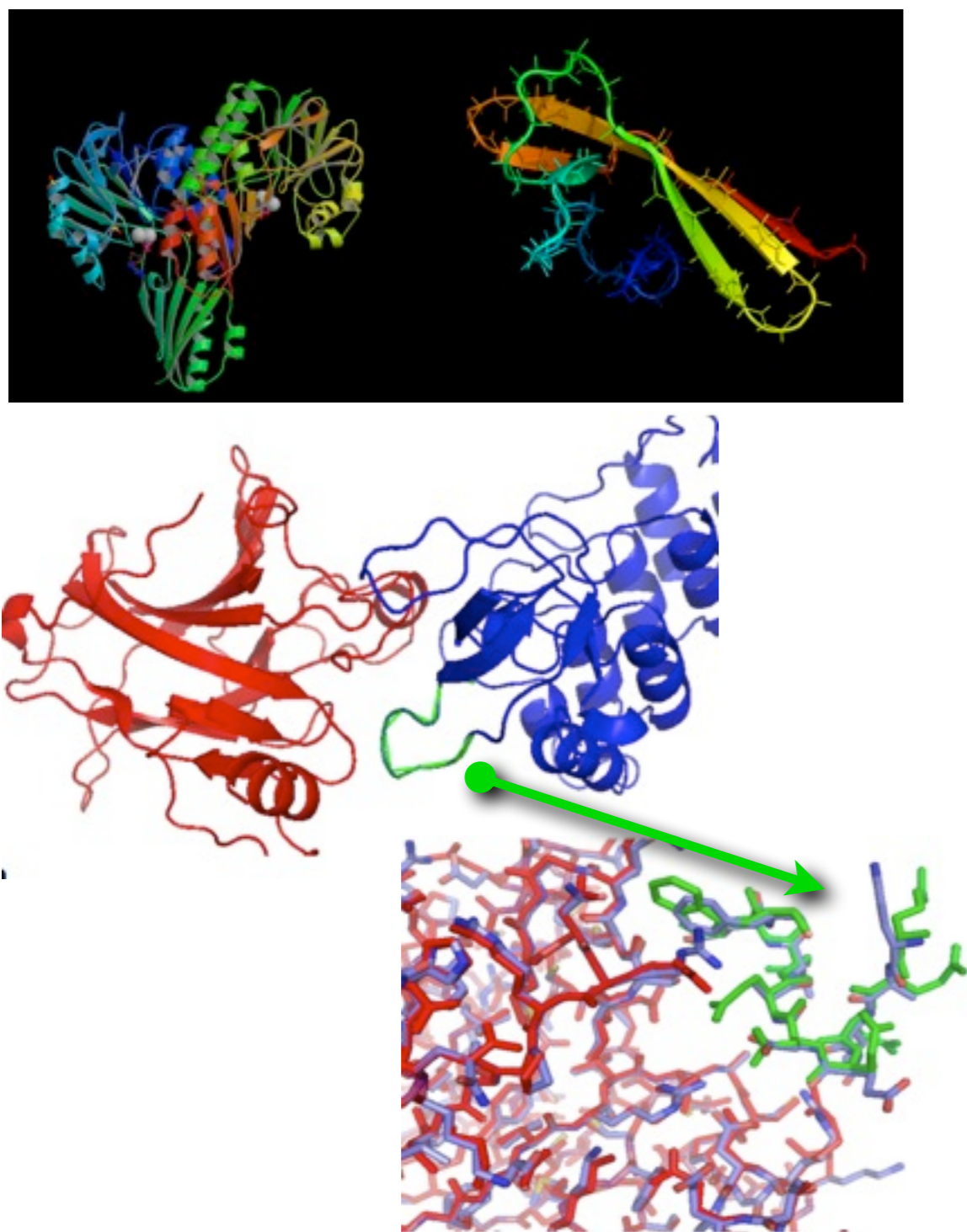
DEPT. OF BIOLOGY &

COMPUTER SCIENCE DEPT.



**CENTER FOR GENOMICS
AND SYSTEMS BIOLOGY
NEW YORK UNIVERSITY**





References

Bonneau, R*, Facciotti, MT, Reiss, DJ, Madar A., et al. , Baliga, NS*. A predictive model for transcriptional control of physiology in a free living cell. (2007) Cell. Dec 131:1354-1365.

cMonkey biclustering and co-regulated modules:

David J Reiss, Nitin S Baliga, Bonneau R. (2006) Integrated biclustering of heterogeneous genome-wide datasets. BMC Bioinformatics. 7(1):280.

Jochen Supper, Claas aufm Kampe, Dierk Wanke, Kenneth W. Berendzen, Klaus Harter, Richard Bonneau, and Andreas Zell. Modeling gene regulation and spatial organization of sequence based motifs. 8th IEEE international conference on Bioinformatics and BioEngineering (BIBE 2008) [In Press].

network inference:

Bonneau R, Reiss DJ, Shannon P, Hood L, Baliga NS, Thorsson V (2006) The Inferelator: a procedure for learning parsimonious regulatory networks from systems-biology data-sets de novo. Genome Biol. 7(5):R36.

Bonneau, R. Learning biological networks: from modules to dynamics (2009). Nature Chemical Biology

Aviv Madar, Alex Greenfield, Harry Ostrer, Eric Vanden-Eijnden and Richard Bonneau, The Inferelator 2.0: a scalable framework for reconstruction of dynamic regulatory network models. IEEE-ECMB09, In Press

visualization:

Iliana Avila-Campillo*, Kevin Drew*, John Lin, David J. Reiss, Richard Bonneau. BioNetBuilder, an automatic network interface. Bioinformatics. (2007) Bioinformatics. Feb 1;23(3):392-3.

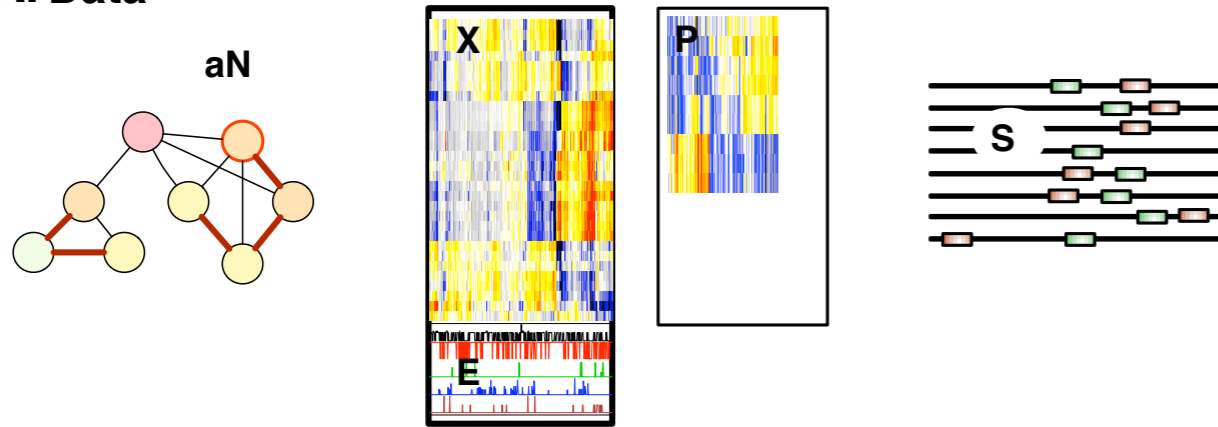
Shannon P, Reiss DJ, Bonneau R, Baliga NS (2006) The Gaggle: A system for integrating bioinformatics and computational biology software and data sources. BMC Bioinformatics. 7:176.

My mentors

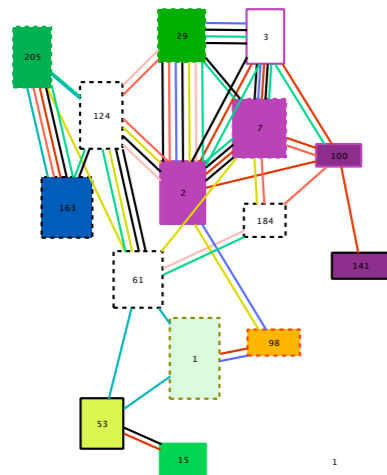
**the PDB,
genomics,
NCBI,
genomes,
etc!**

ME

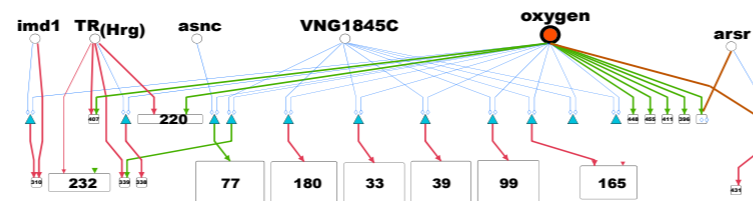
A. Data



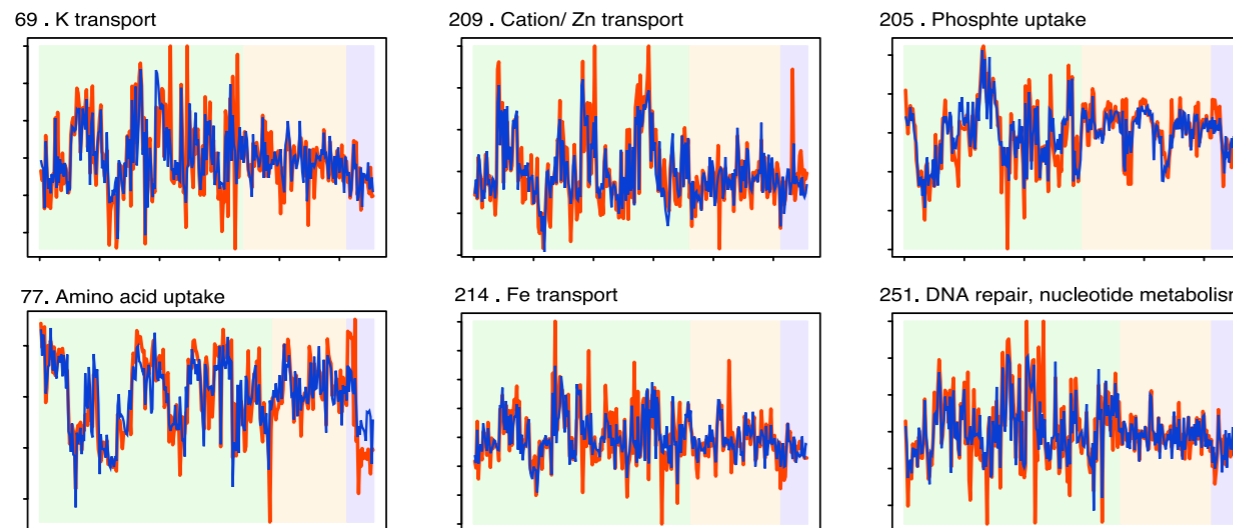
B. (Bi)clustering



C. Dynamical network model



D. Prediction



OVERVIEW

1. CO-REGULATED MODULES
(INTEGRATE DATA TYPES).

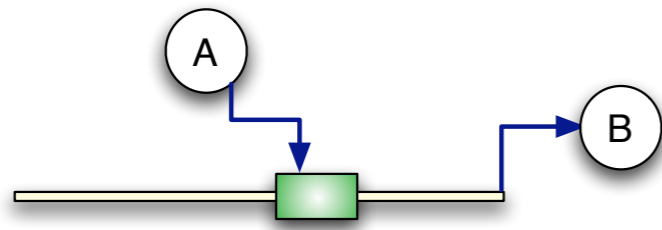
2. LEARN TOPOLOGY AND DYNAMICS WITH GREEDY / LOCAL APROX.
(INFERELATOR 1.0, 1.1)

3. IMPROVING PERFORMANCE OVER MULTIPLE TIME-SCALES
(INFERELATOR 2.X)

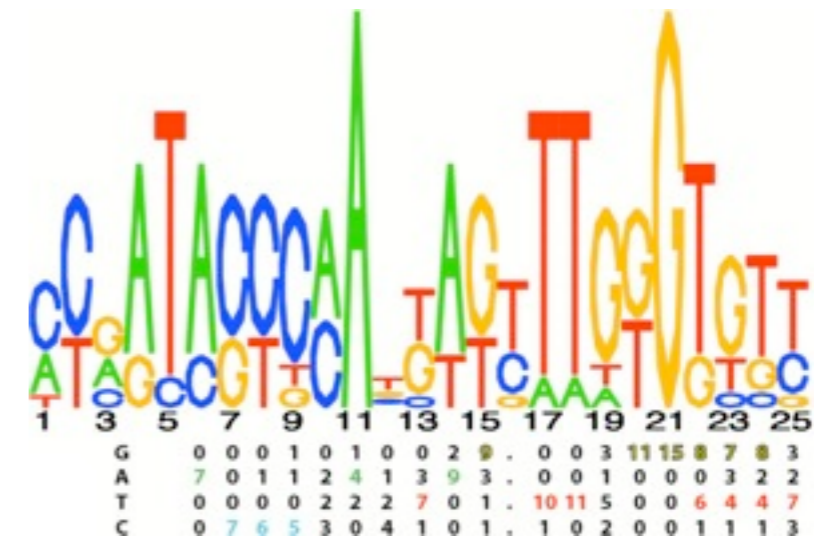
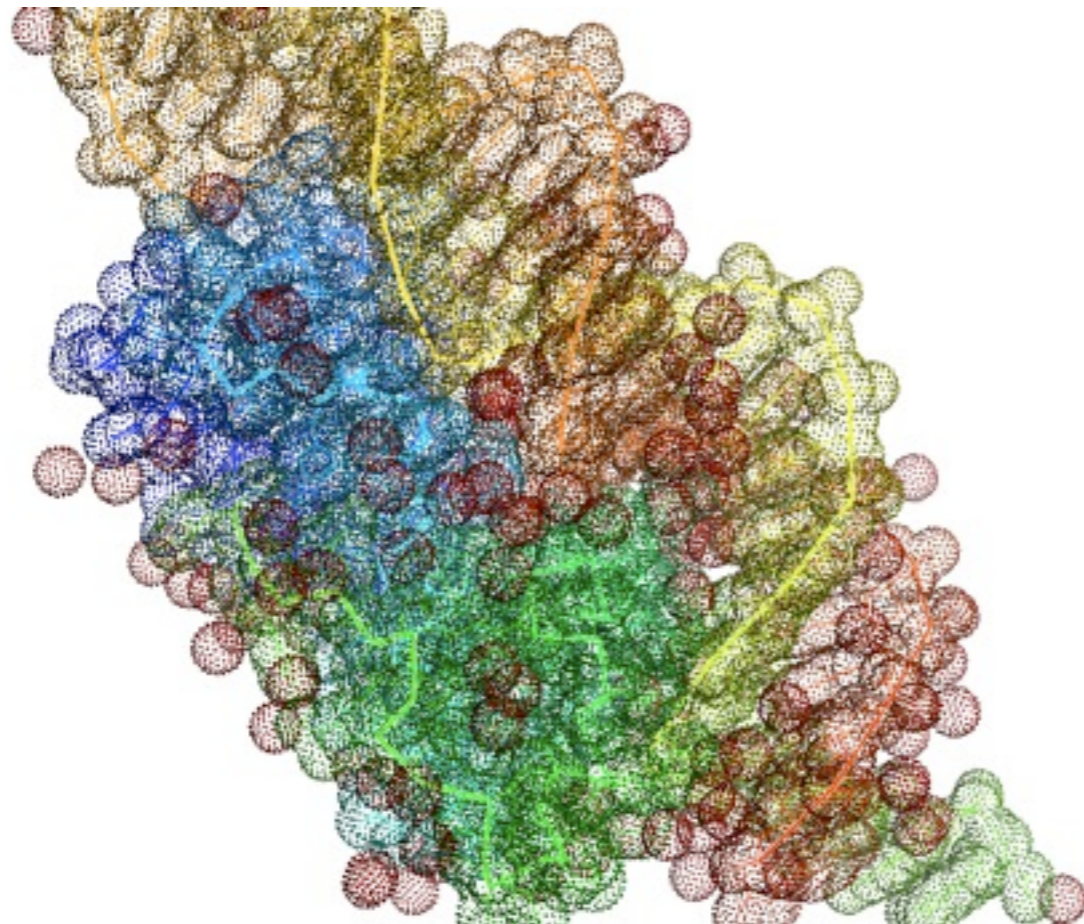
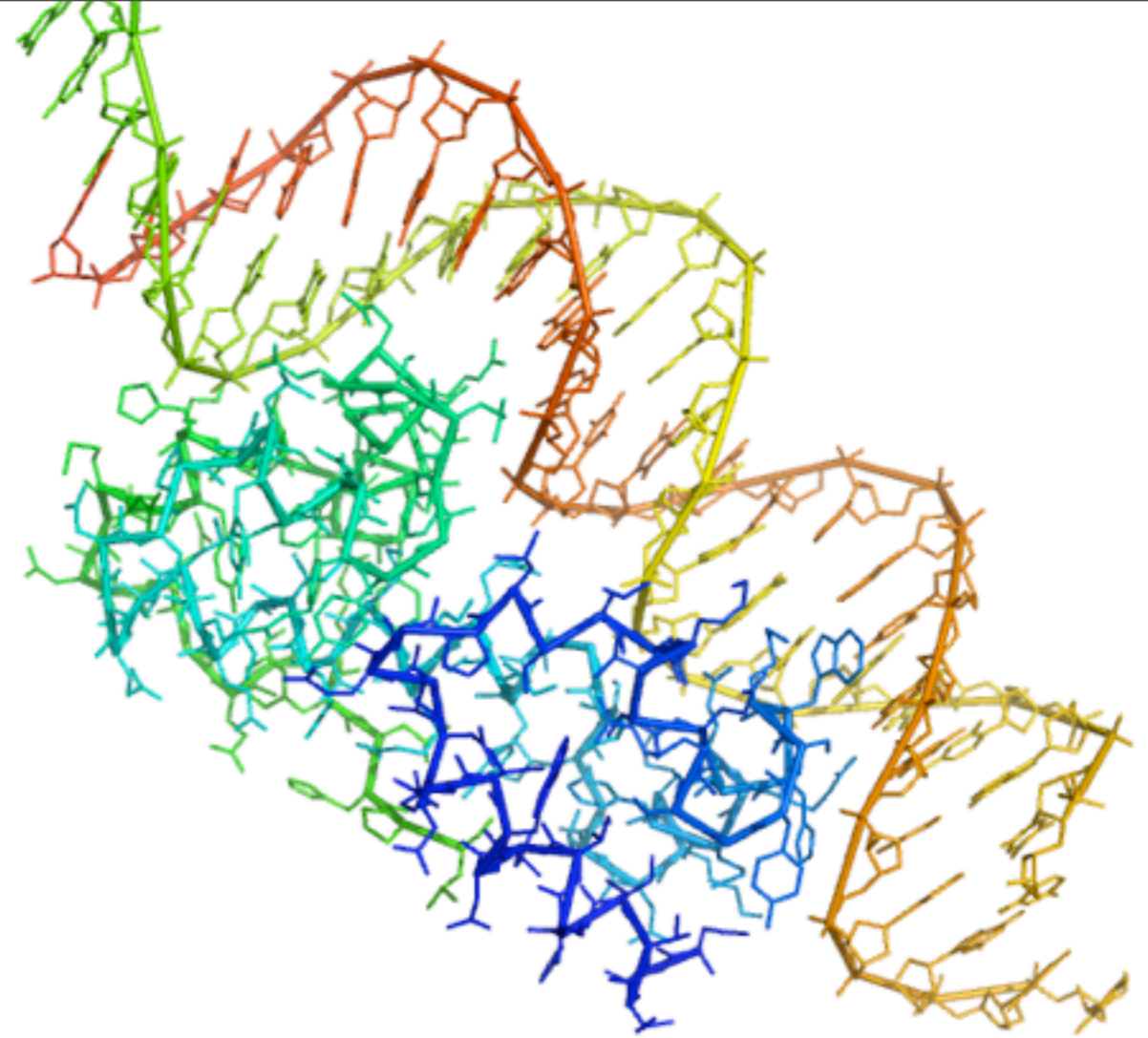
MAIN RESULTS:

- SURPRISING PREDICTIVE PERFORMANCE FOR PROKARYOTIC NETWORKS, T-CELL AND MACROPHAGE DIFFERENTIATION EE NETWORKS
- LONGER TIME SCALE STABILITY
- MODEL FLEXIBILITY

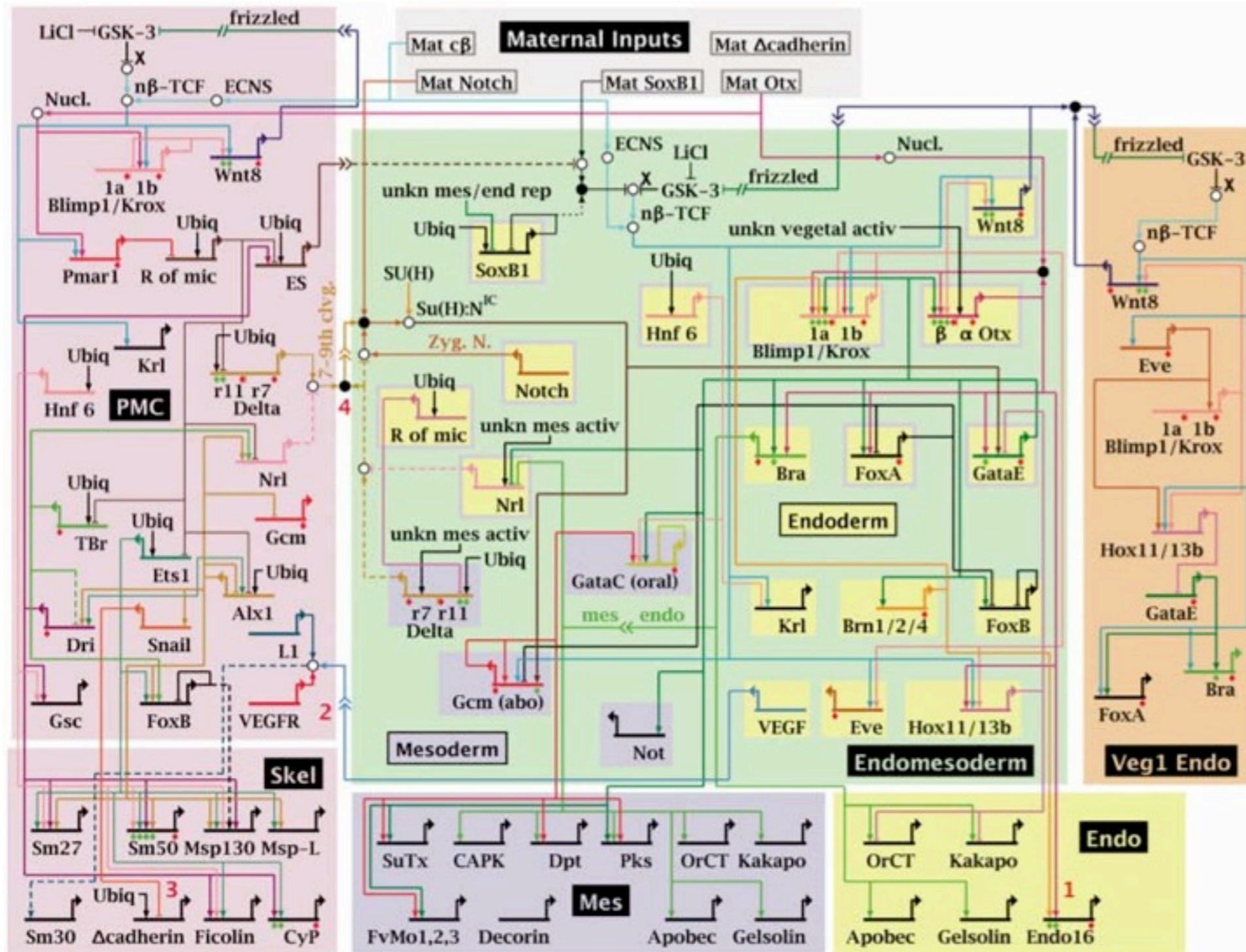
TRANSCRIPTIONAL REGULATION



OR



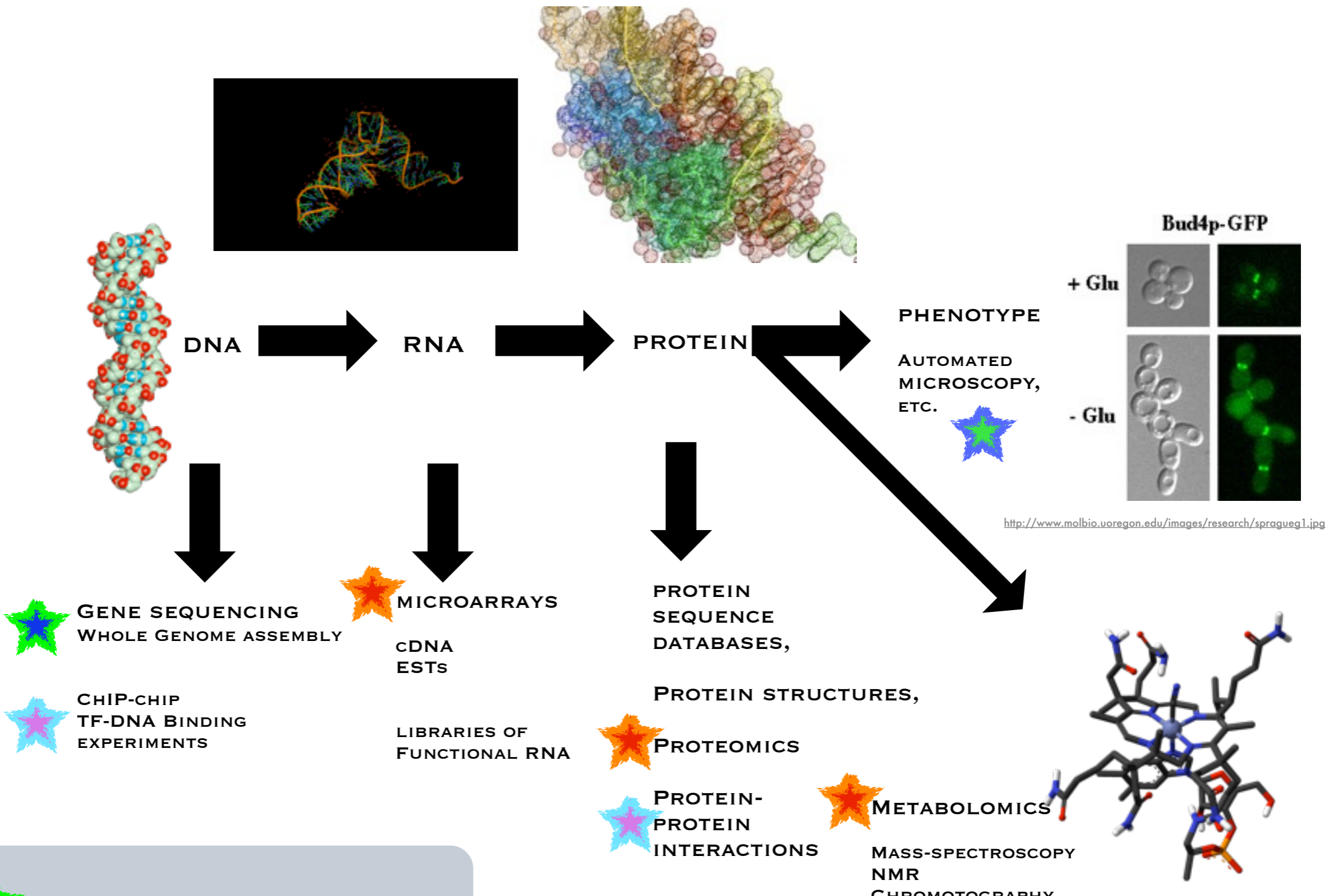
TRANSCRIPTIONAL NETWORKS CONTROLLING DEVELOPMENT

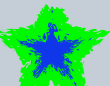


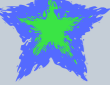


Ubiqu=ubiquitous; Mat = maternal; activ = activator; rep = repressor;
 unkn = unknown; Nucl. = nuclearization; χ = β -catenin source;
 n β -TCF = nuclearized b- β -catenin-Tcf1; ES = early signal;
 ECNS = early cytoplasmic nuclearization system; Zyg. N. = zygotic Notch

Copyright © 2001-2006 Hamid Bolouri and Eric Davidson

BOLOURI, DAVIDSON



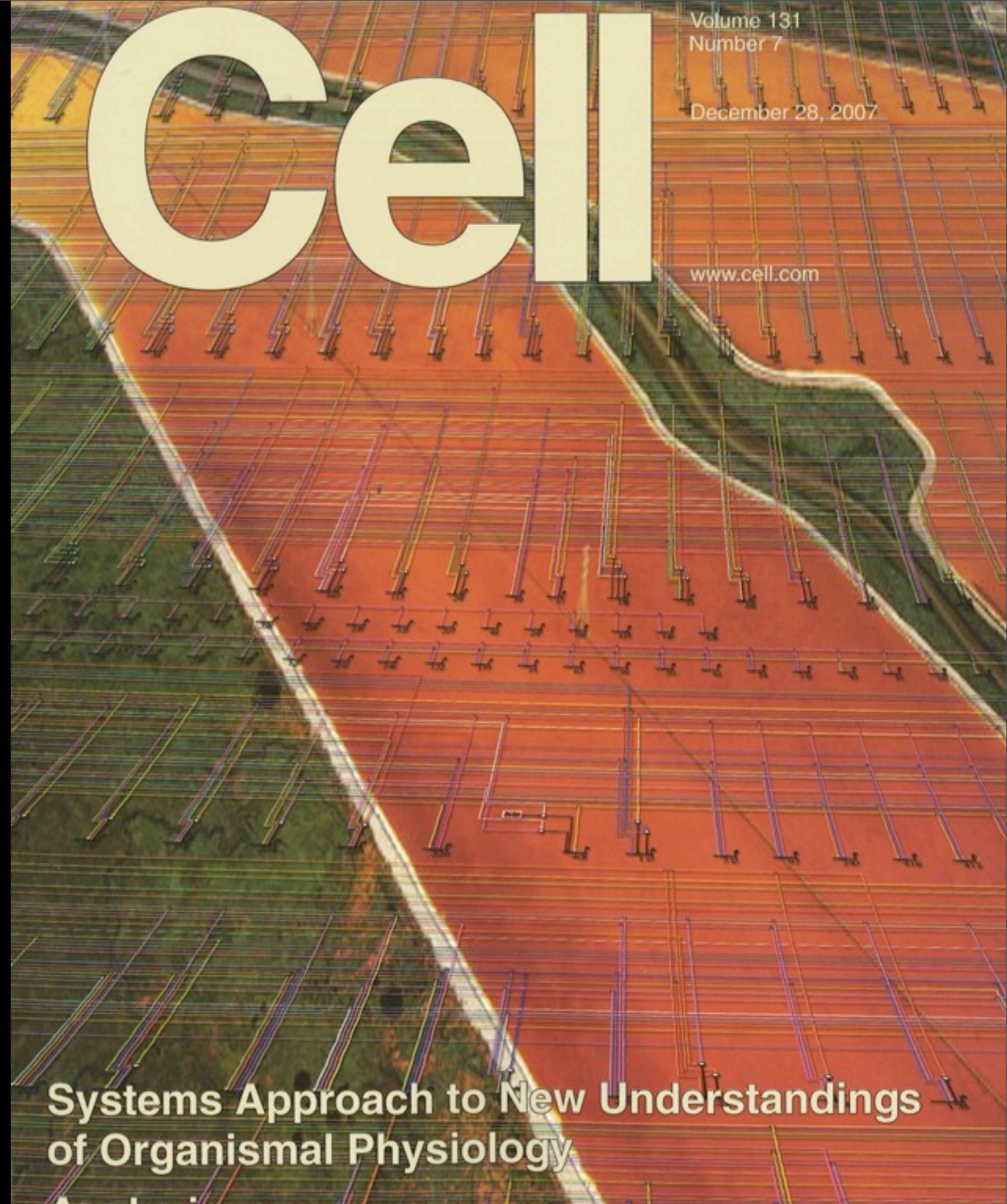
-  GENOTYPE & SEQUENCING
-  MEASURING AFFINITIES / BINDING
-  MEASURING LEVELS
-  ASSAYING FUNCTIONAL OUTCOME

algorithms:

David J. Reiss (cMonkey)
Vesteinn Thorsson (Inferelator)
Richard Bonneau

functional genomics:

Marc T. Facciotti
Amy Schmid,
Kenia Whitehead
Min Pan, Amardeep Kaur,
Leroy Hood
Nitin S. Baliga



Systems Approach to New Understandings
of Organismal Physiology

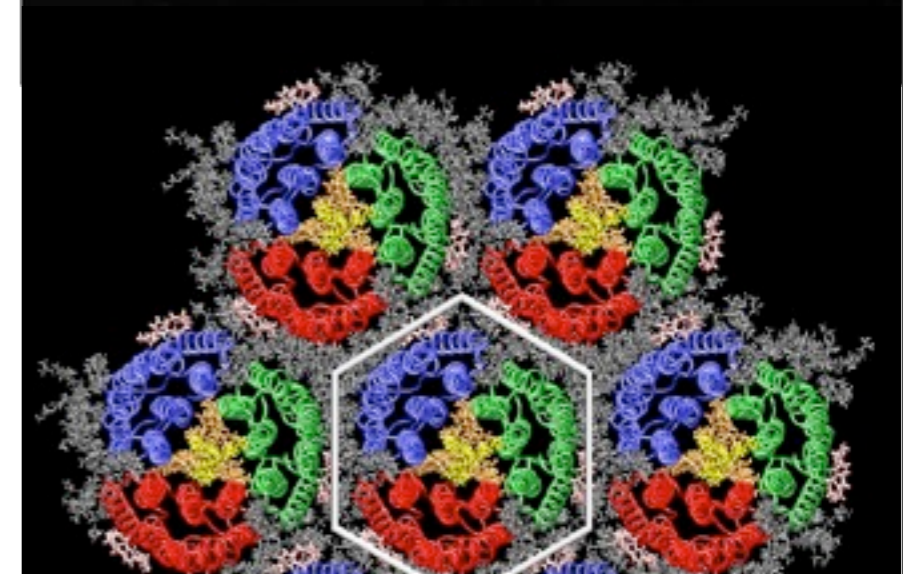
AN EXAMPLE : HALOBACTERIUM

WHY HALOBACTERIUM:

- if your friends are working on halo ... (Hood, Baliga)
- not a “model” system (originally)
- high IQ
- diverse environment
- small genome
- good genetics, cultivable, etc.
- a very tough extremophile, bioengineering

DATA COLLECTION AND MODELING EFFORT

- * genome and genome annotation
- * microarrays
- * genetic and environmental perturbations
- * proteomics
- * ChIP-chip
- * some protein-protein



HALOBACTERIUM DATASET INCLUDING

**>800 MICROARRAYS
TIME SERIES
KNOCK OUTS**

**CHIP-CHIP
EXPERIMENTS**

PROTEOMICS

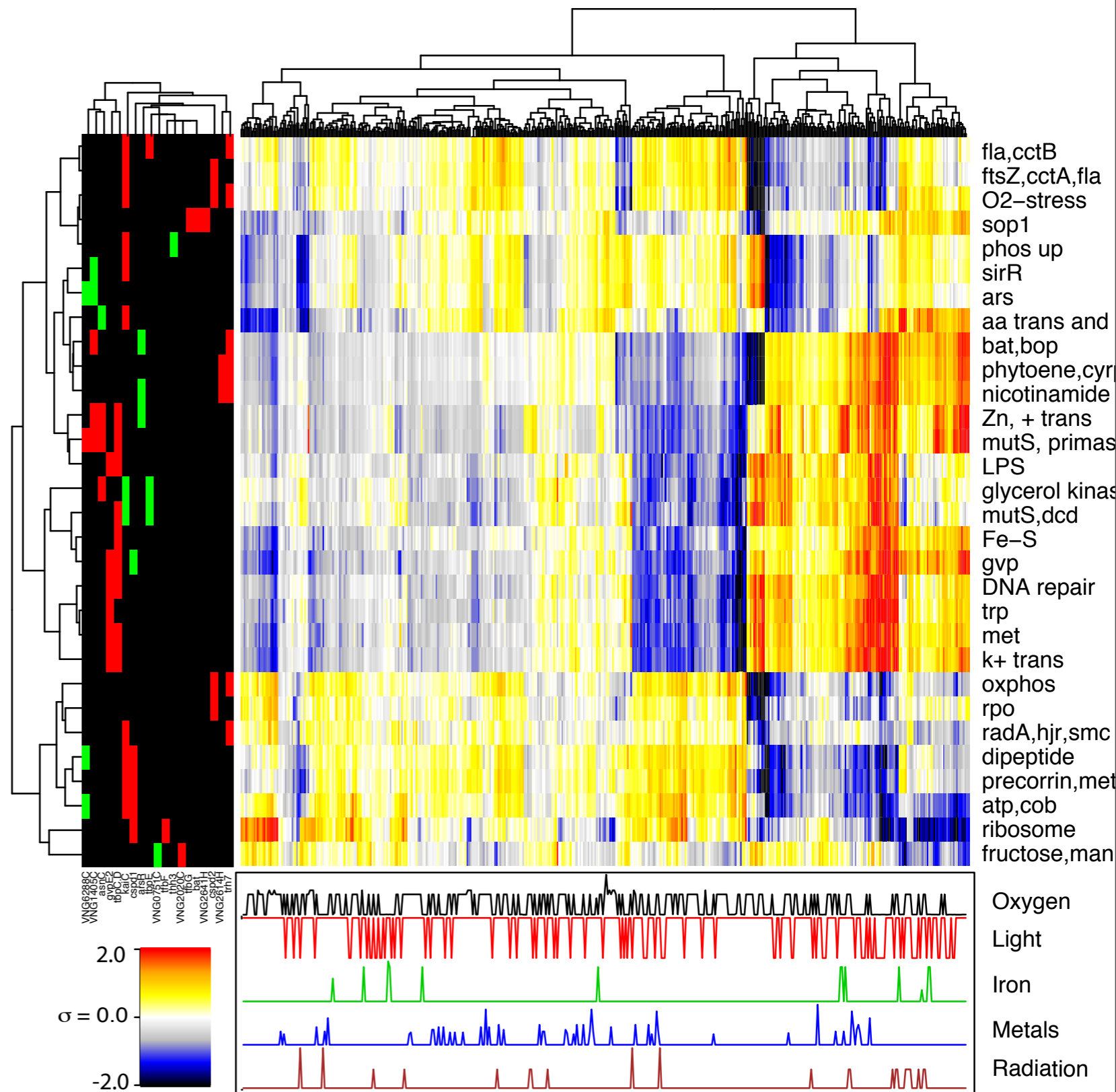
PHENOTYPE

**AMONG THE MOST
COMPLETE
PROKARYOTIC
DATASETS**

M. FACCIOTTI, N. BALIGA



MIN PAN, KENIA WHITEHEAD, AMY SCHMID



OUR APPROACH

BIOLOGICAL MOTIVATION:

Co-regulation dramatically reduces complexity of network inference, and unlike simple co-expression has direct mechanistic relevance to biological control.

Time (explicit learning/modeling of kinetic parameters) helps even in our current state of affairs.

Model must be capable of modeling interactions with bio relevant functional forms.

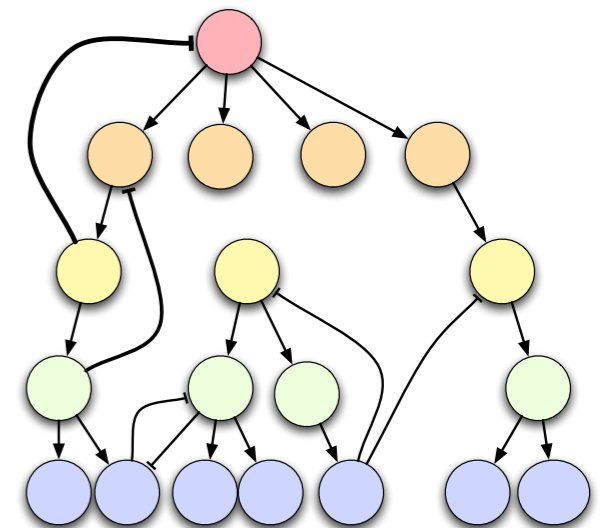
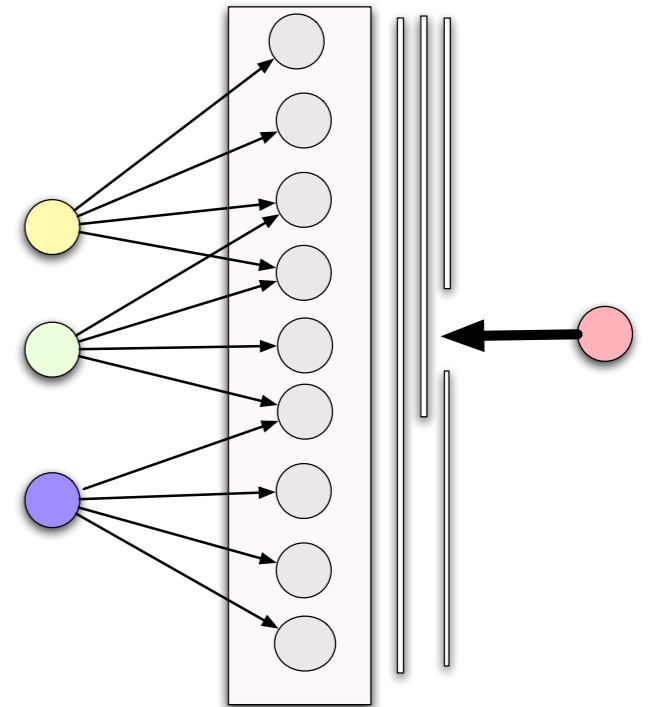
Experimental design is key

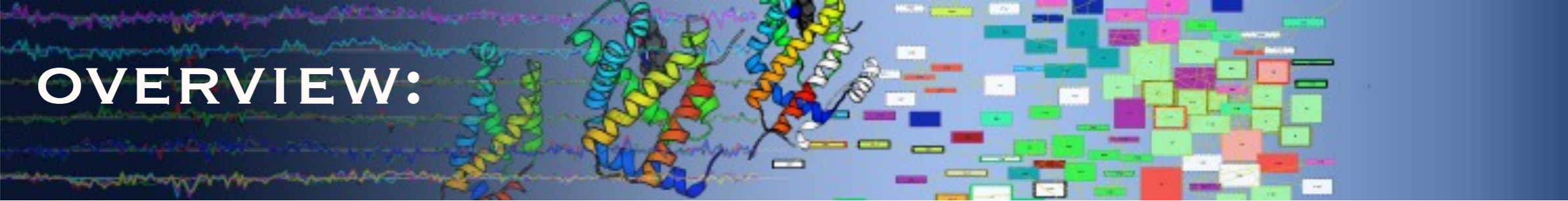
CMONKEY:

- ★ integrate data-types other than expression to constrain search for co-regulated modules
- ★ avoid lossy transformations of the data and derive joint P of gene given bicluster and all datatypes
- ★ derive framework with eye toward flexibility (new datatypes)

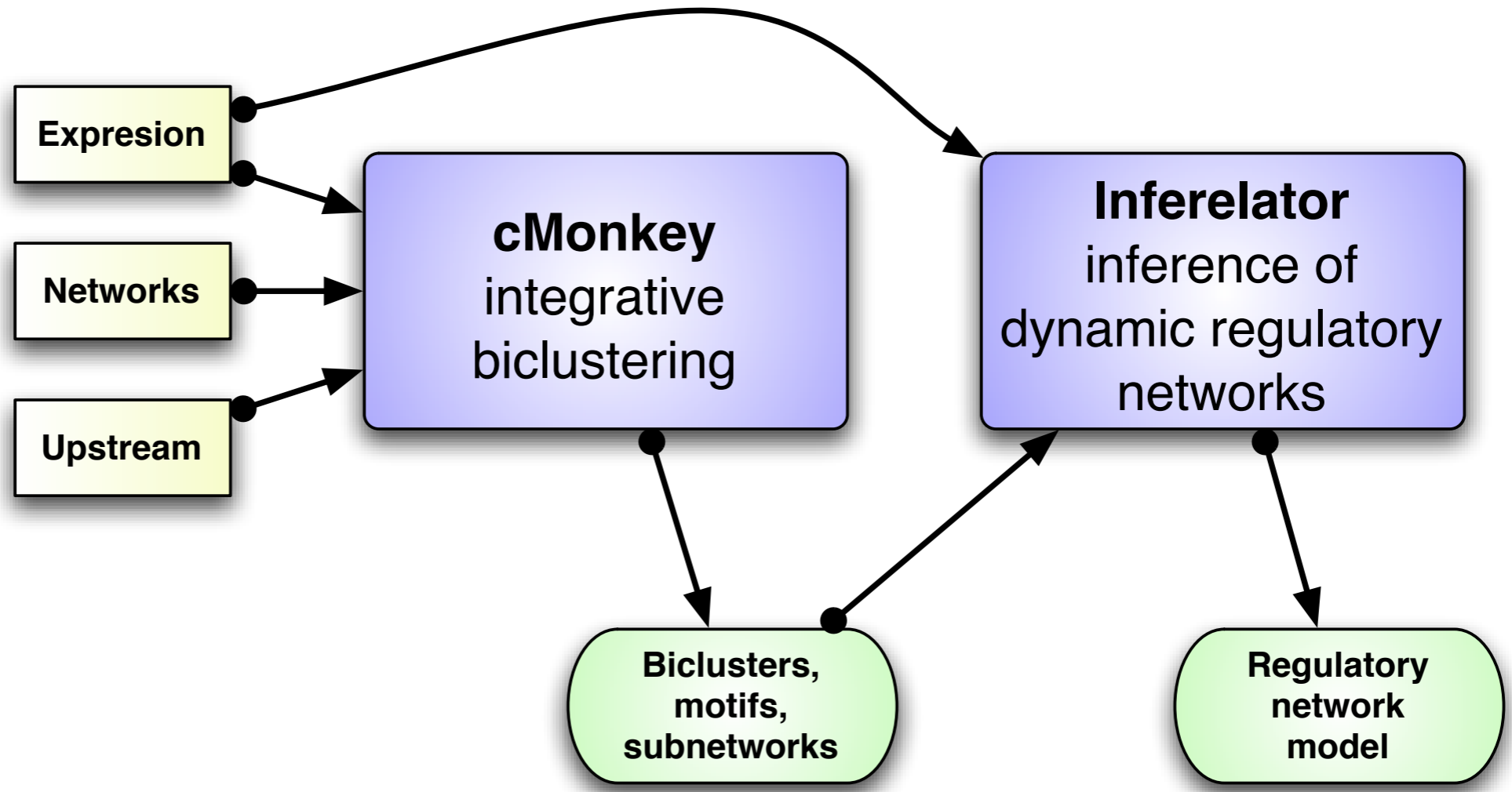
INFRELATOR:

- ★ frame parameterization of global set of ODEs as regression problem
- ★ interactions: map problem onto tropical semi-ring

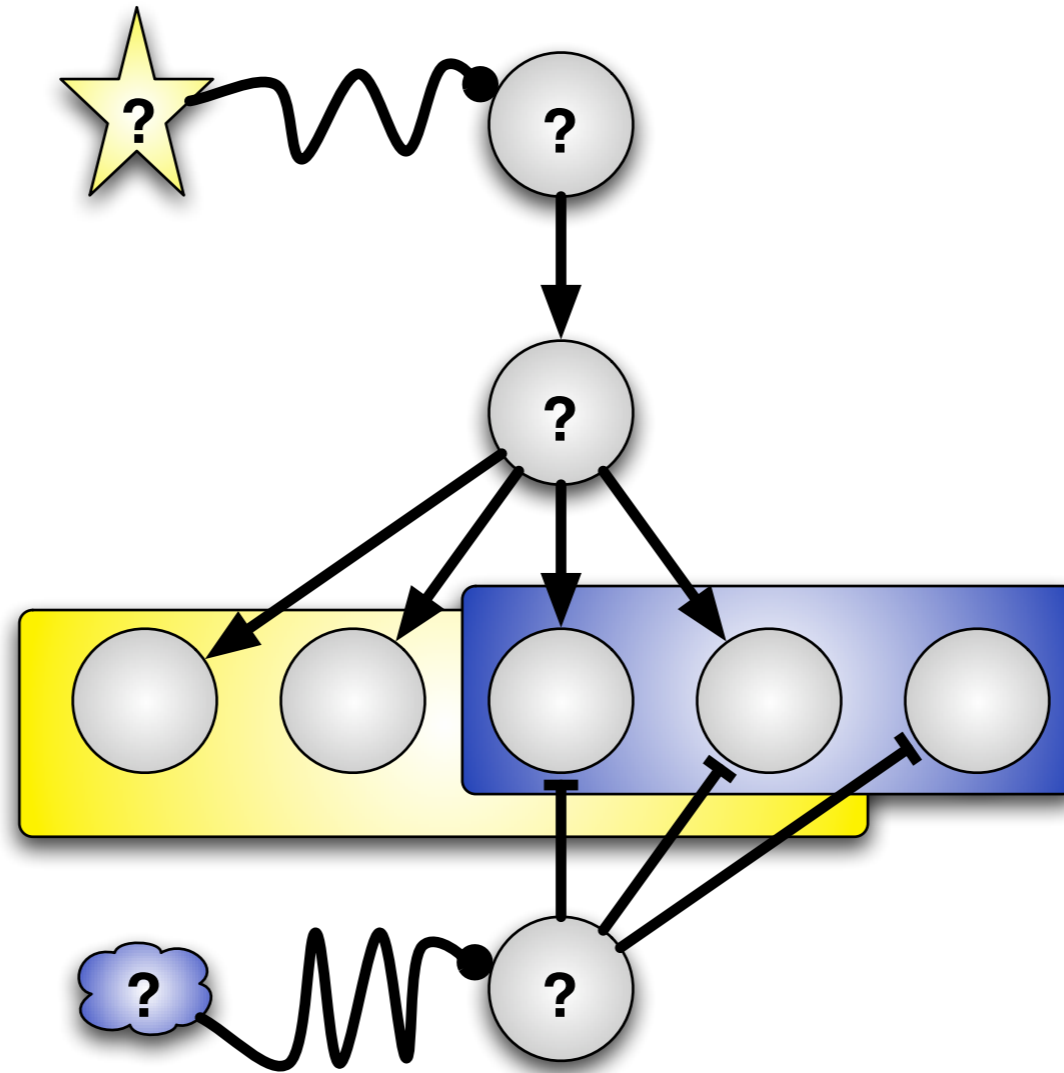




OVERVIEW:



LEARNING CO-REGULATED GROUPS:



What is Biclustering?

- Concurrent clustering of both rows & conditions
- Given an $n \times m$ matrix, A , find a set of submatrices, B_k , such that the contents of each B_k follow a desired pattern, i.e. gene co-expression.



Based on lecture notes from Kai Li:

<http://www.cs.princeton.edu/courses/archive/spr05/cos598E/Biclustering.pdf>

Reasons to Biclust in Biology & Bioinformatics

- Genes not regulated under all conditions
 - ⇒ patterns of correlation may exist only under subsets of conditions
- Genes can participate in multiple modules or processes
 - ⇒ exclusive clustering algorithms (HAC, K-means) will miss valid clusterings

I. cMONKEY: INTEGRATIVE BICLUSTERING

BIOLOGICAL MOTIVATION:

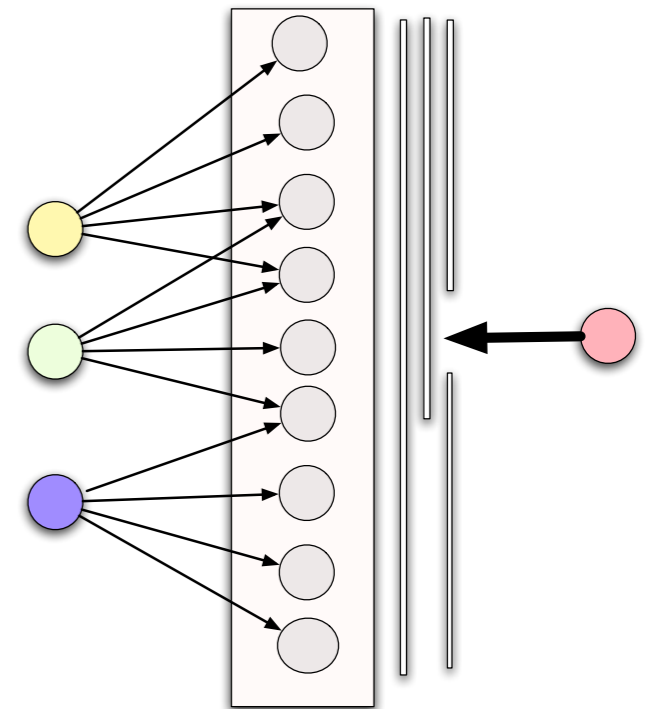
Co-regulation dramatically reduces complexity of network inference, and unlike simple co-expression has direct mechanistic relevance to biological control.

STRATEGY:

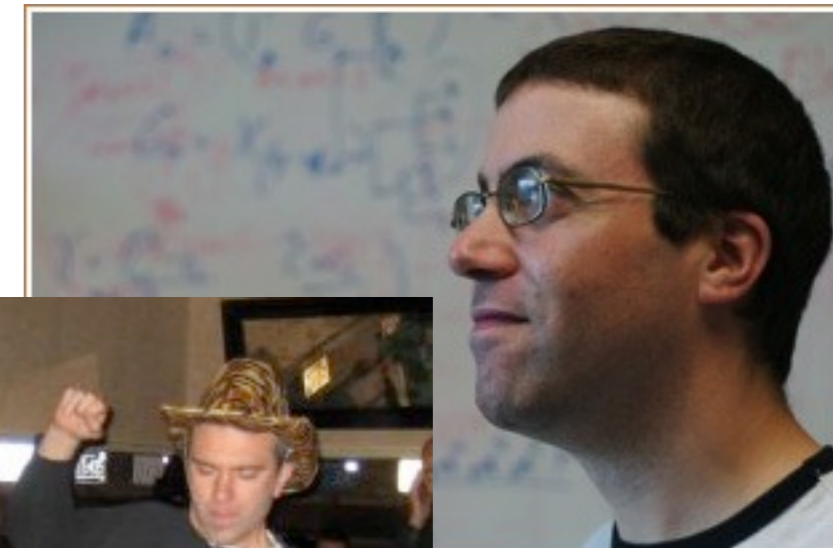
- ★ integrate data-types other than expression to constrain search for co-regulated modules
- ★ avoid lossy transformations of the data and derive joint P of gene given bicluster and all datatypes
- ★ derive framework with eye toward flexibility (new datatypes)

CHALLENGES:

- 💣 overlapping (genes participate in multiple functions)
- 💣 diverse data types
- 💣 mix of well studied and completely unknown genes
- 💣 **many think of this as a solved problem...why?**
- 💣 Resultant models are a complex low-level abstraction of the systems behavior (functional modules, complexes, annotations, etc. are linked to clusters).



DAVE REISS



PETER
WALTMAN



CMONKEY: MCMC OPTIMIZATION OF A MULTI-DATA LIKELIHOOD

OTHER DATA:

EXP-LIKE:

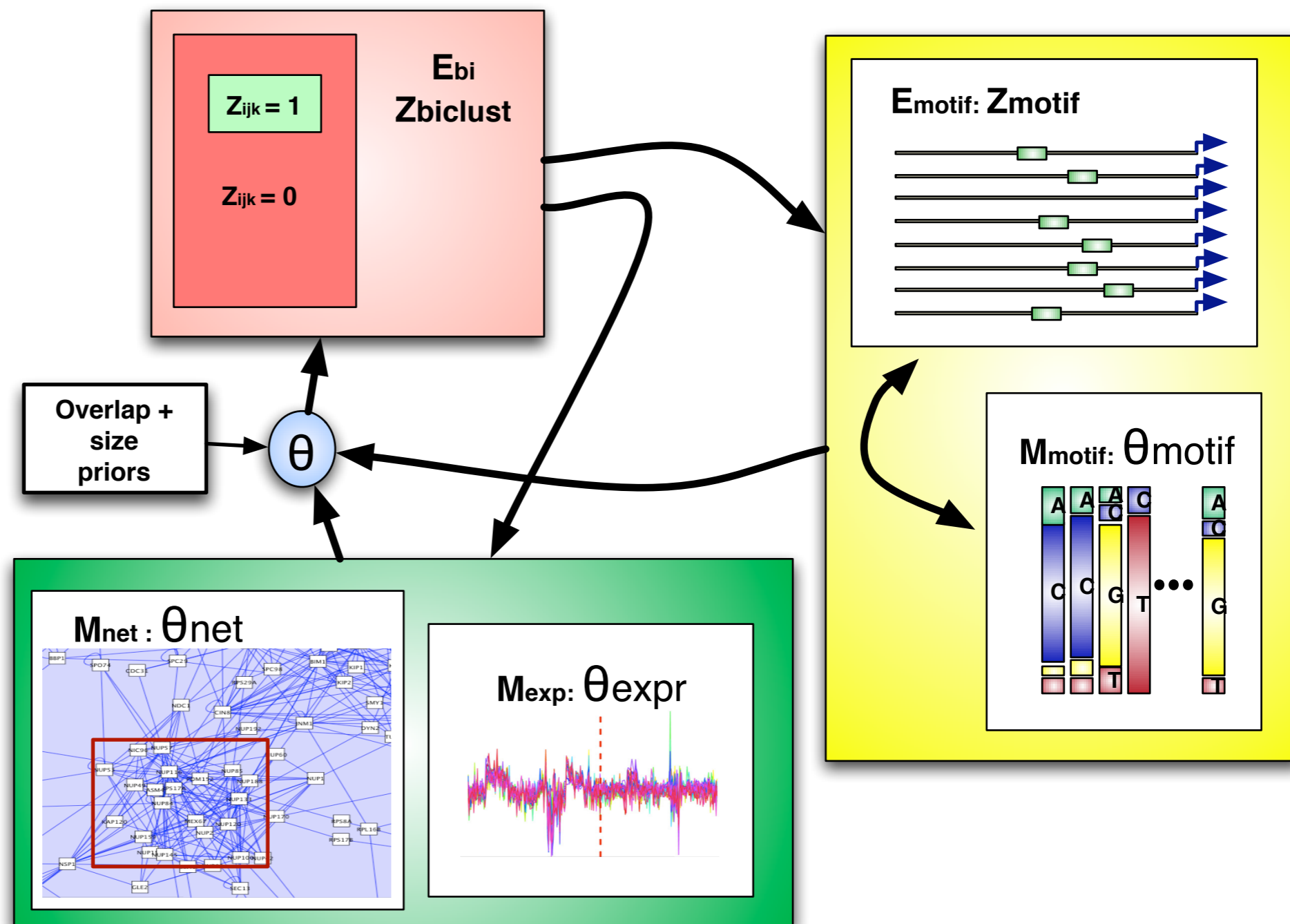
[GWA,
COPY NUMBER,
PHENOTYPE]

NETS:

[CHIP-SEQ,
ETC.]

SEQ:

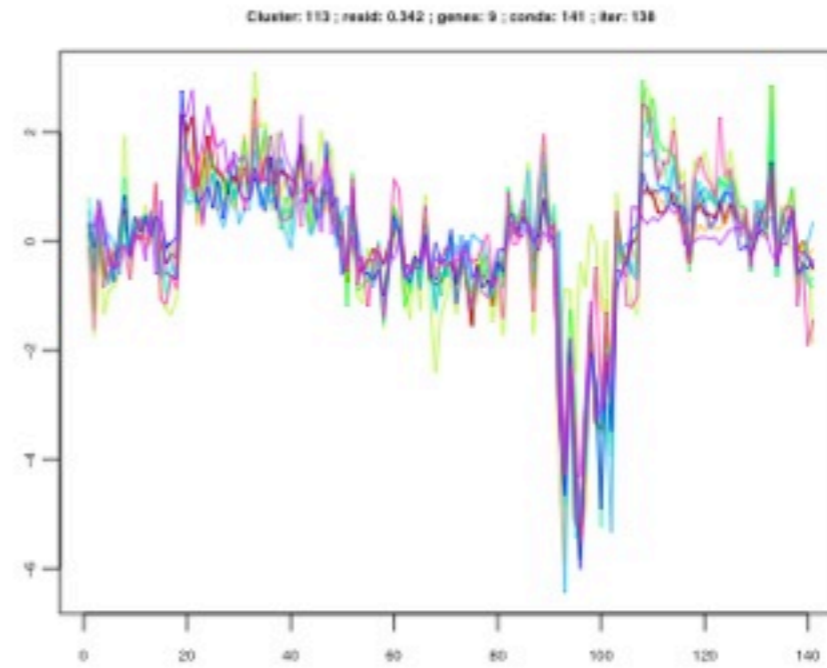
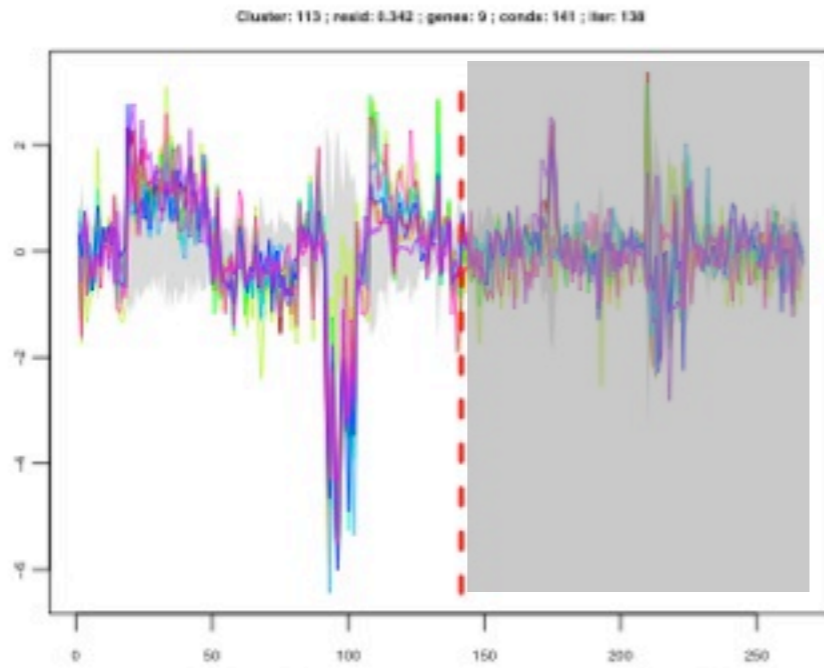
[UTR, KNOWN SITES]



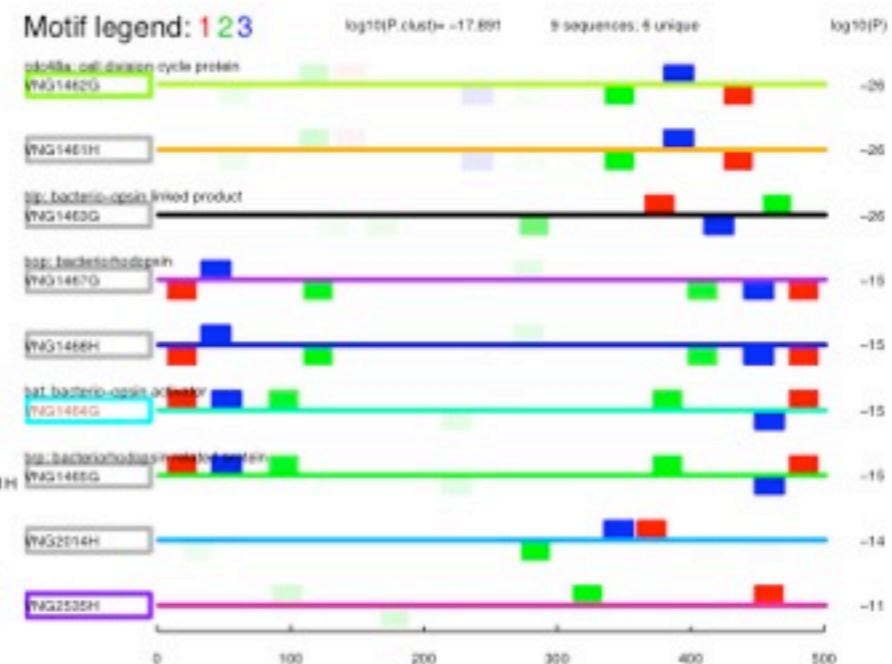
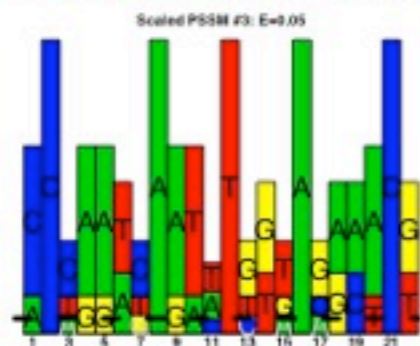
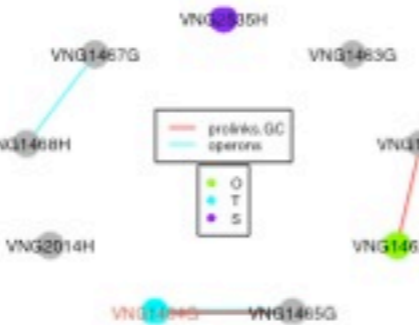
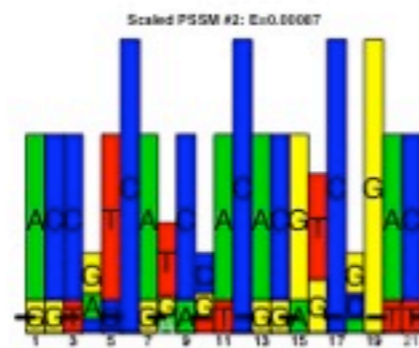
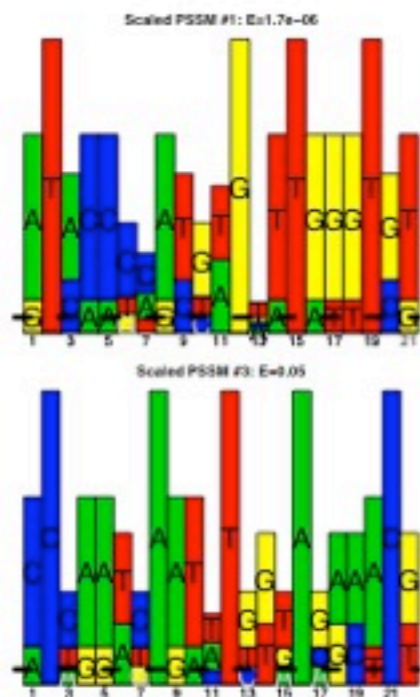
THE REAL
ADVANTAGE OF
CMONKEY IS ITS
LACK OF LOSSY
TRANSFORMATIONS

ARCHAEA: BOP/BAT-ASSOCIATED REGULON [HALOBACTERIUM NRC-1]

EXPRESSION

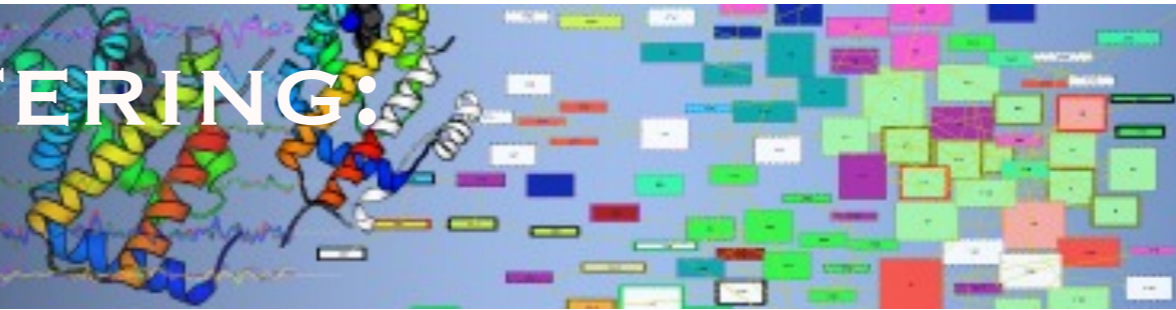


MOTIFS



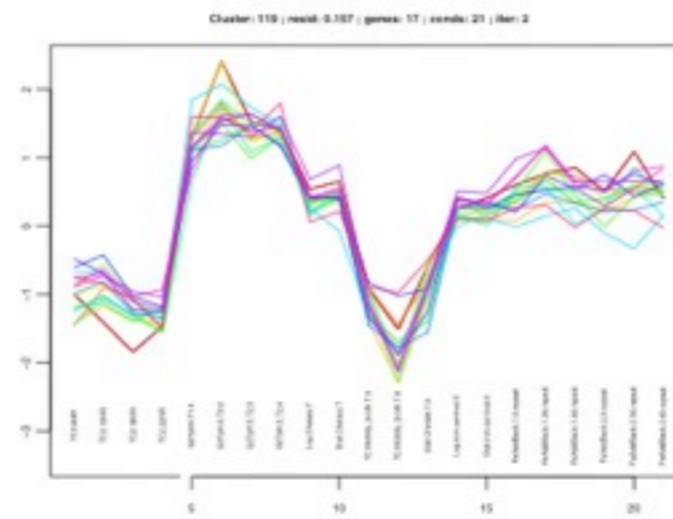
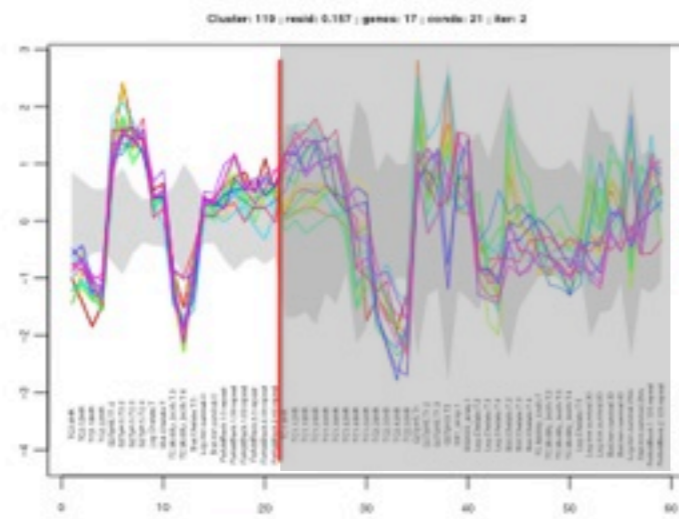
BALIGA, ET AL.
(1999,2000)

MULTI-BICLUSTERING: MULTISPECIES



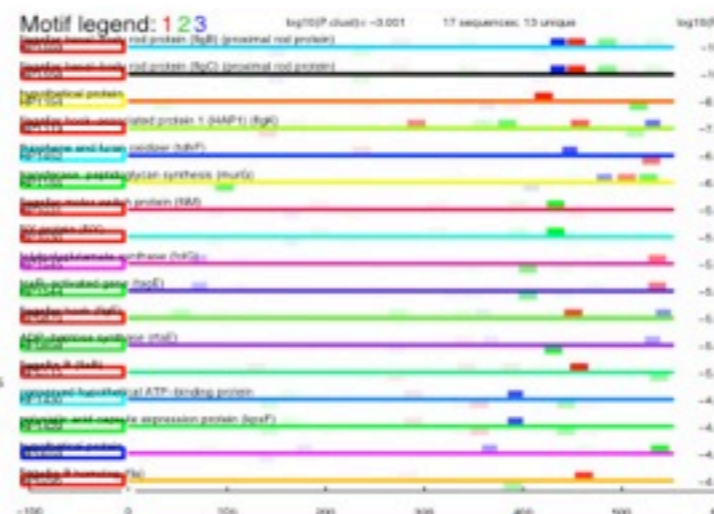
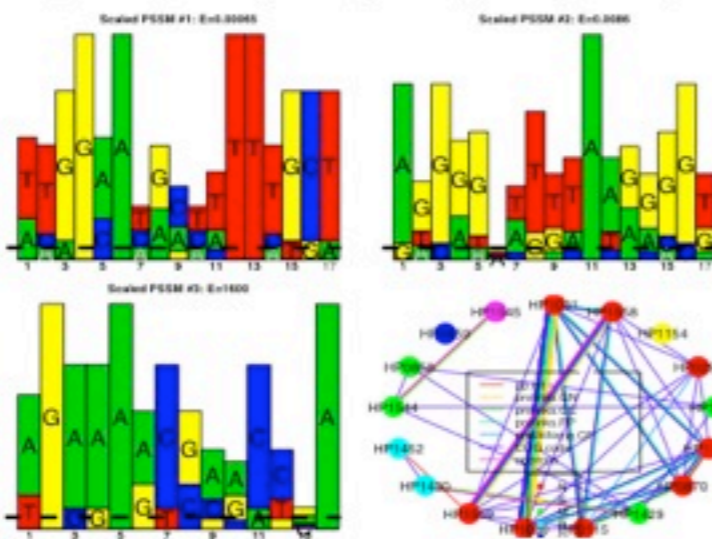
**BACTERIA: RPON-ASSOCIATED FLAGELLAR REGULON [H. PYLORI]
---> [ALSO IN E. COLI]**

EXPRESSION



**NIEHUS, ET AL.
(2004)**

MOTIFS

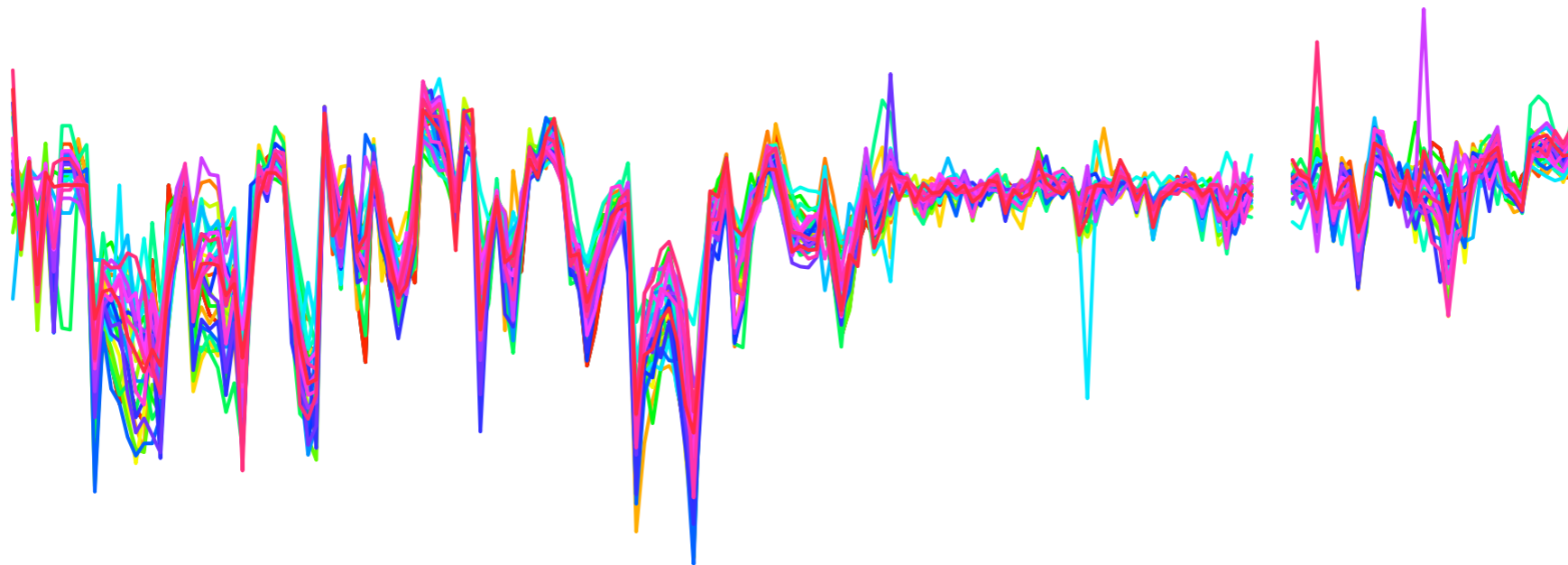


UPSTREAM

**W/ PATRICK EICHENBERGER, NYU; HARRY OSTRER, NYU-MED
ERIC ALM, MIT, BROAD**

SCORE COMPONENT I: R, EXPRESSION [LEVELS]

$$p(x_{ij}) = \frac{1}{\sqrt{2\pi(\sigma_j + \epsilon)^2}} \exp \left[-\frac{1}{2} \left(\frac{x_{ij} - \bar{x}_{jk} + \epsilon}{\sigma_j + \epsilon} \right)^2 \right]$$



Reiss, Shannon, Baliga, Bonneau, 2006

SCORE COMPONENT II:

P, MOTIF DETECTION AND CO-OCCURRANCE

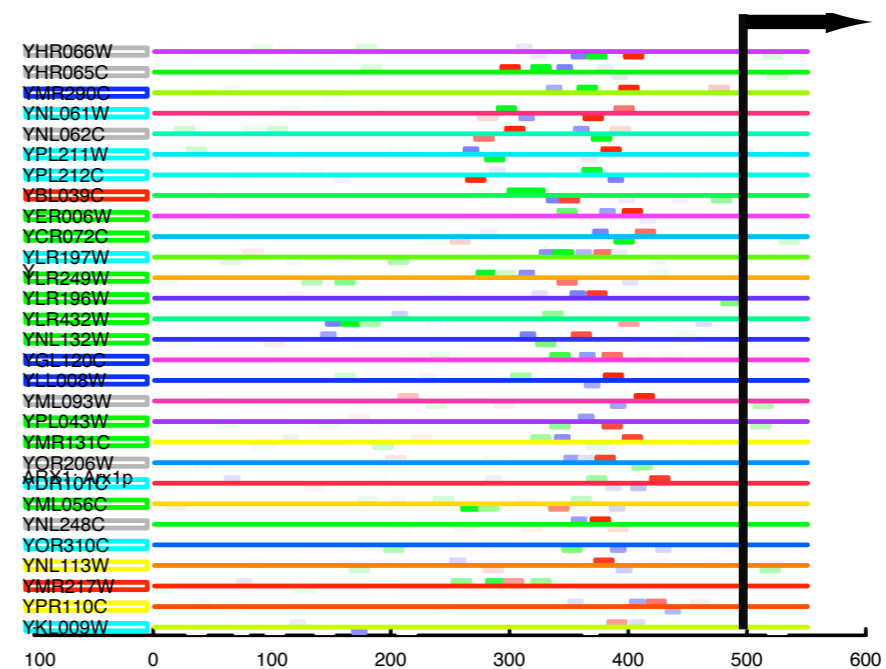
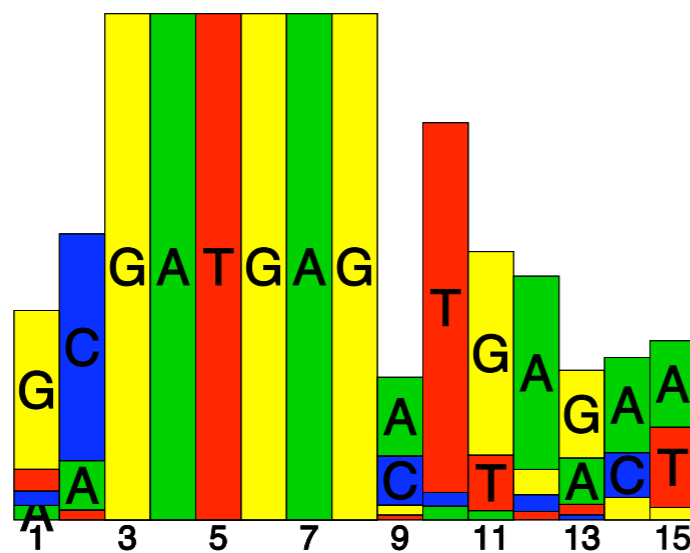
[SHORT SEQUENCES]

motif models:

MEME, Weeder, known

->

cis, trans, UTR

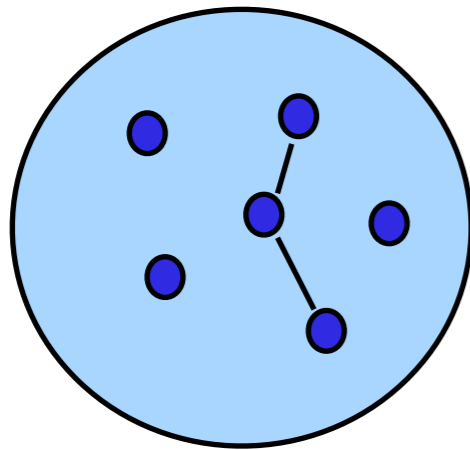


Reiss, Shannon, Baliga, Bonneau, 2006

SCORE COMPONENT III:

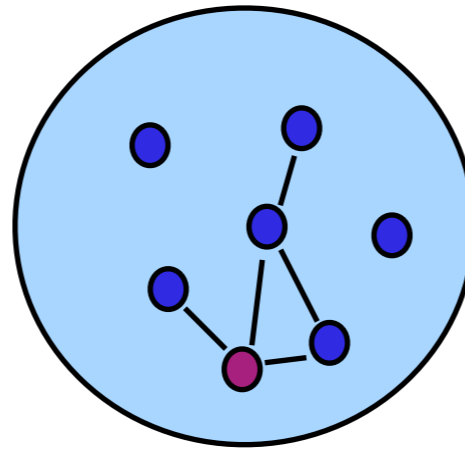
Q, NETWORKS [ASSOCIATIONS]

Before adding gene



$$p = 0.16$$

After adding gene



$$p = 0.012$$

Reward addition of genes to bicluster that share edges with other genes in bicluster.

Hypergeometric distribution to derive p -values:

$$p(n_{i \rightarrow I_k} | n_{i \rightarrow I'_k}, n_{I_k \rightarrow I_k}, n_{I_k \rightarrow I'_k}) = \frac{\binom{n_{i \rightarrow I_k} + n_{I_k \rightarrow I_k}}{n_{i \rightarrow I_k}} \binom{n_{i \rightarrow I'_k} + n_{I_k \rightarrow I'_k}}{n_{i \rightarrow I'_k}}}{\binom{n_{i \rightarrow I_k} + n_{I_k \rightarrow I_k} + n_{i \rightarrow I'_k} + n_{I_k \rightarrow I'_k}}{n_{i \rightarrow I_k} + n_{i \rightarrow I'_k}}}$$

cMonkey continued

- Combine 3 likelihoods into a joint log-likelihood:

$$g_{ik} = r_0 \log(\tilde{r}_{ik}) + s_0 \log(\tilde{s}_{ik}) + \sum_{n \in N} q_0^n \log(\tilde{q}_{ik}^n)$$

where r_0 , s_0 and q_0 are “mixing parameters” –

Pre-selected and set according to an annealing schedule

- Logistic regression to discriminate between genes in/out of bicluster:

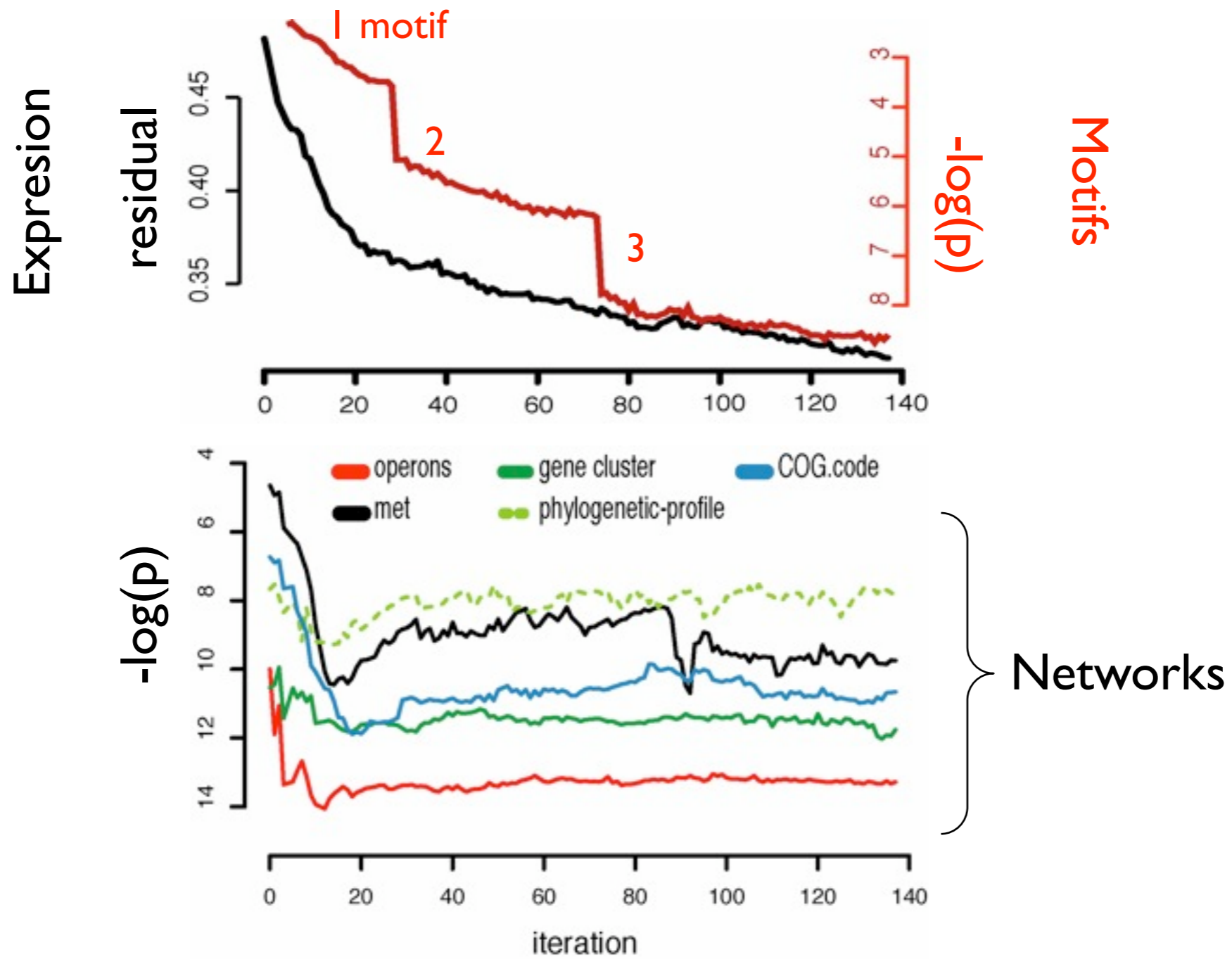
$$\pi_{ik} \equiv p(y_{ik} | X_k, S_i, M_k, N) = \frac{e^{(\beta_0 + \beta_1 g_{ik})}}{1 + e^{(\beta_0 + \beta_1 g_{ik})}}$$

where $p(y_{ik}=1)$ indicates likelihood of membership of gene i to cluster k

- Monte Carlo, annealing of the biclusters:

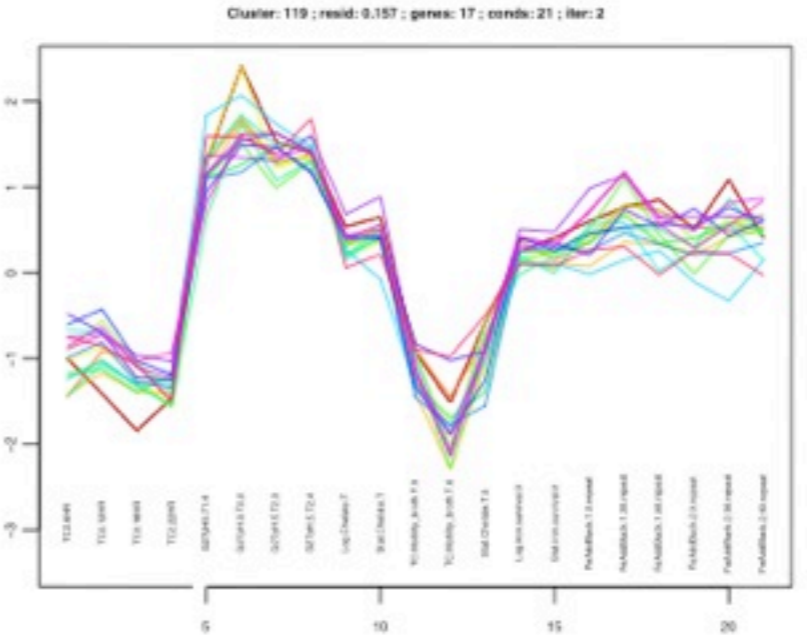
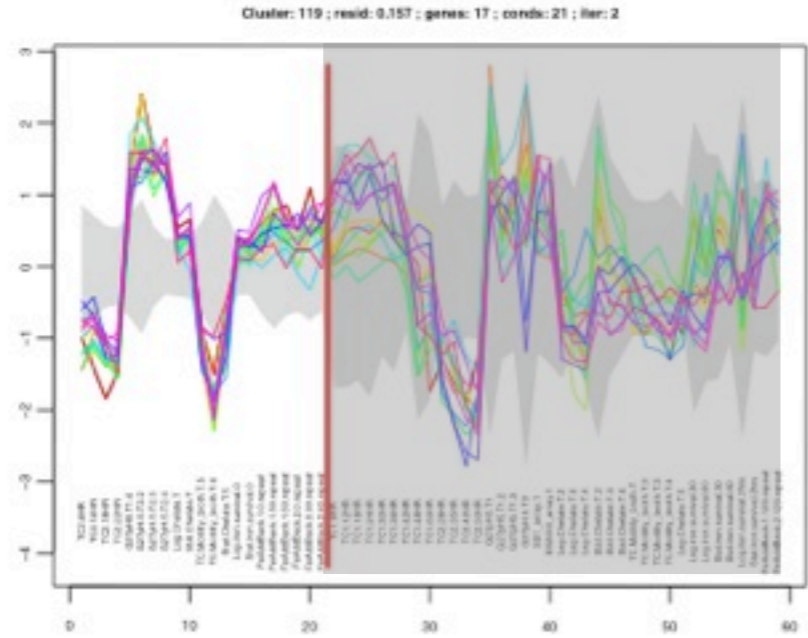
$$p(\text{add}_{ik} | \pi_{ik}) = e^{\frac{-(1-\pi_{ik})}{T}}; \quad p(\text{drop}_{ik} | \pi_{ik}) = e^{\frac{-\pi_{ik}}{T}}$$

OPTIMIZATION OF SCORE ELEMENTS

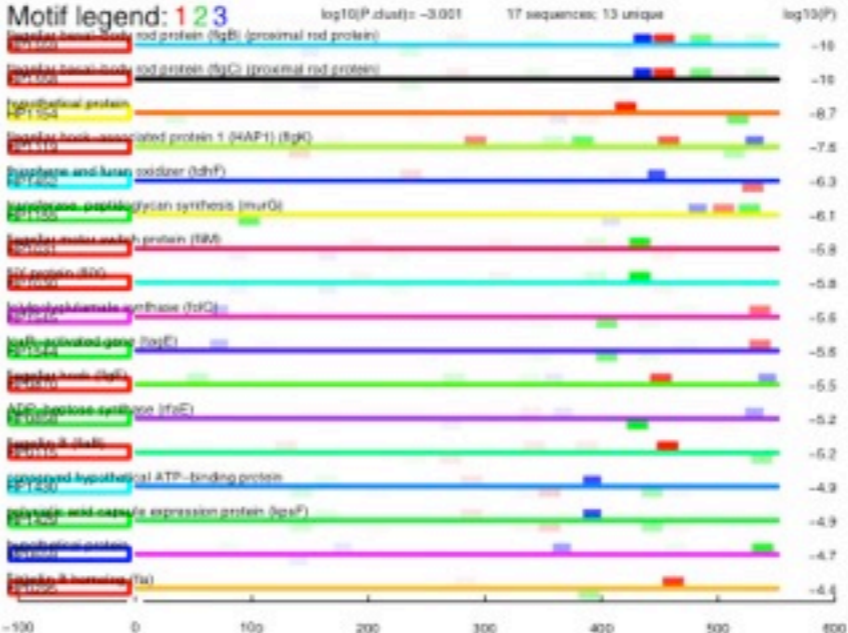
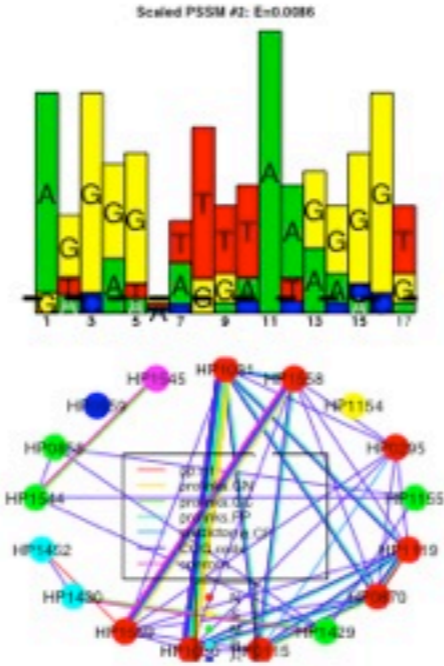
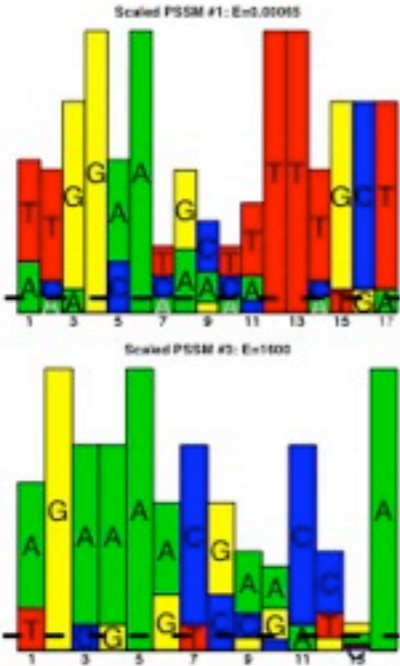


BACTERIA: RPON-ASSOCIATED FLAGELLAR REGULON [H. PYLORI] ---> [ALSO IN E. COLI]

EXPRESSION



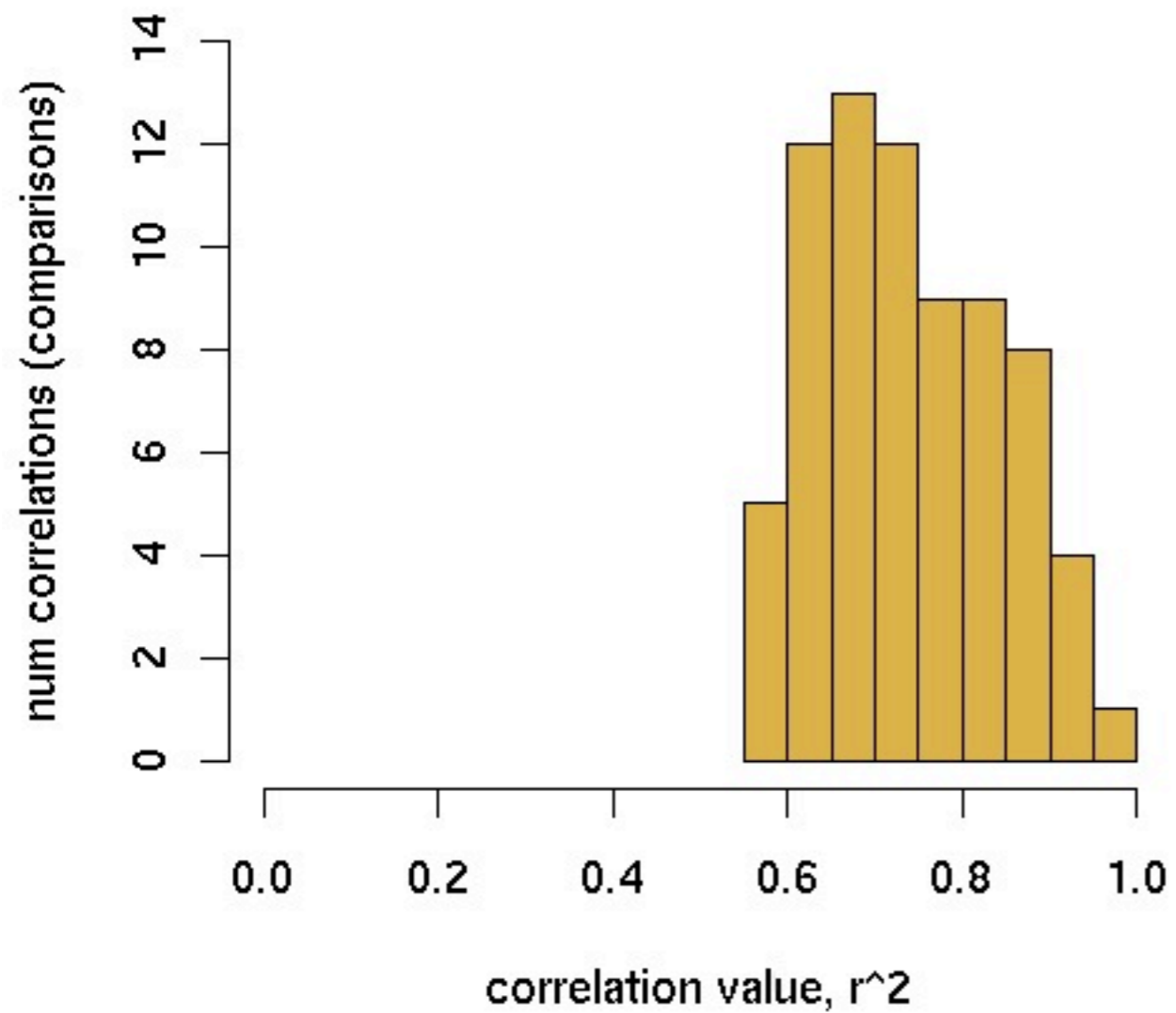
MOTIFS



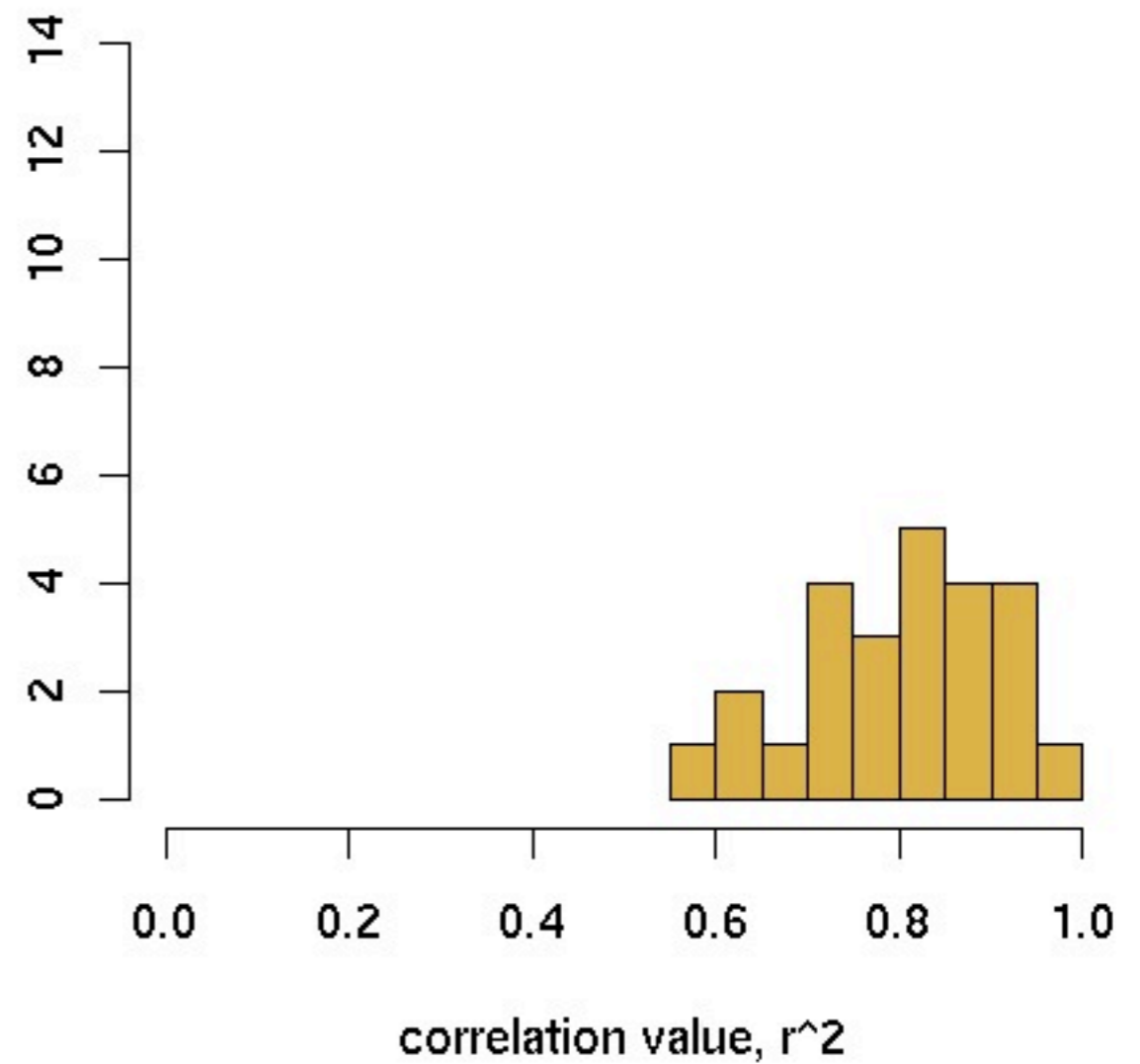
UPSTREAM

NIEHUS, et al. (2004)

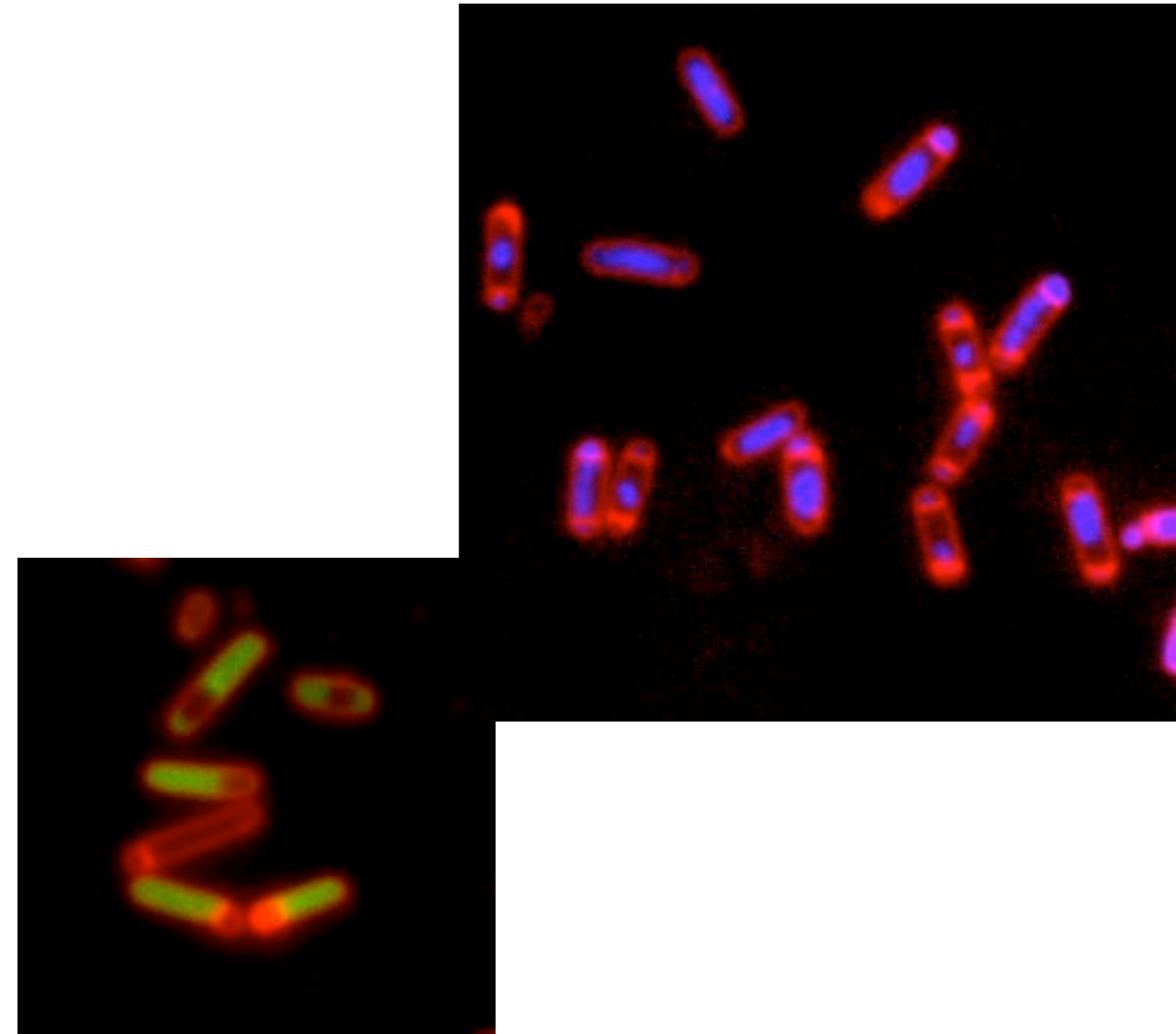
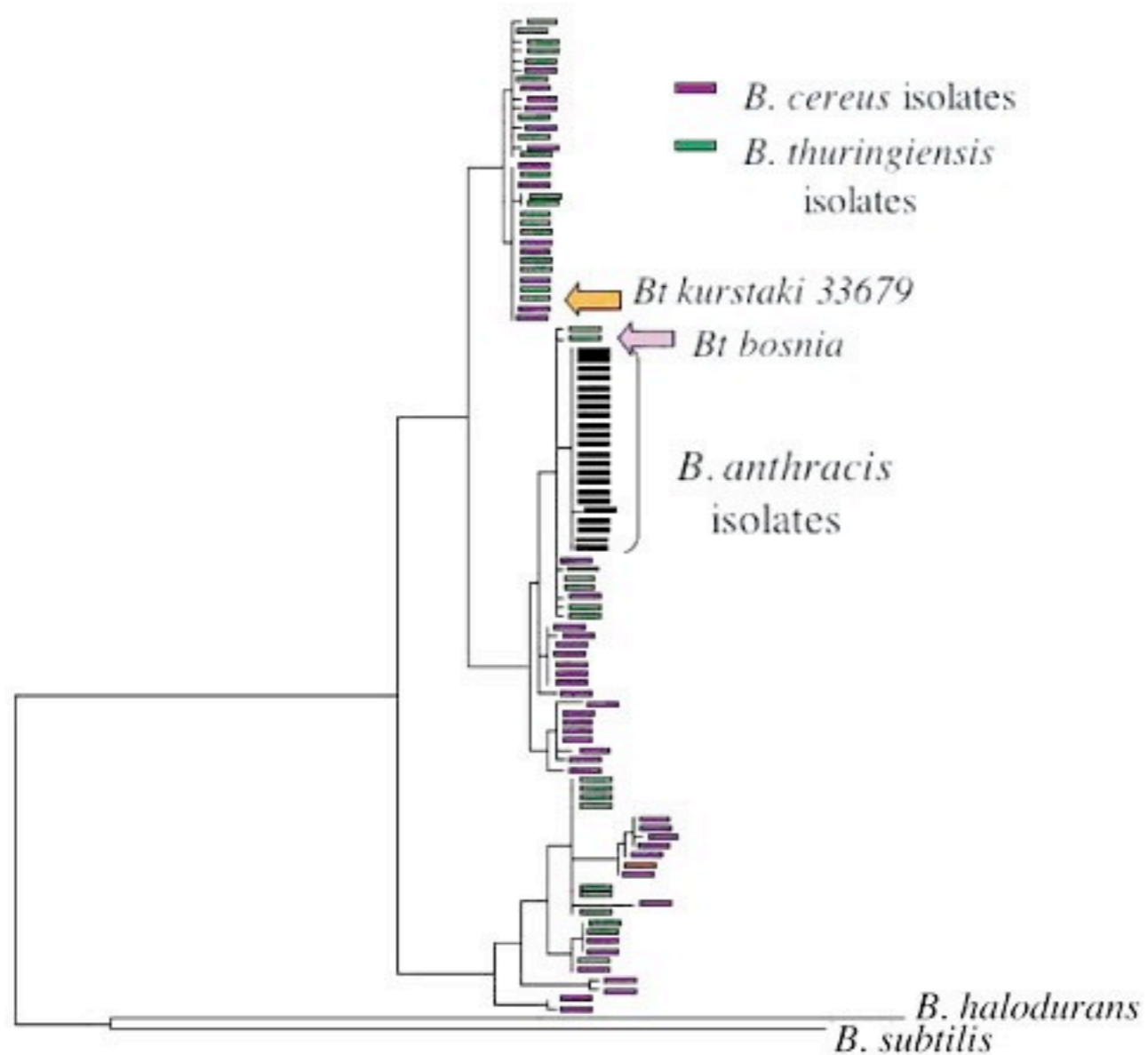
E. coli motif comparisons



B. subtilis motif comparisons



MULTI-BICLUSTERING: MULTISPECIES



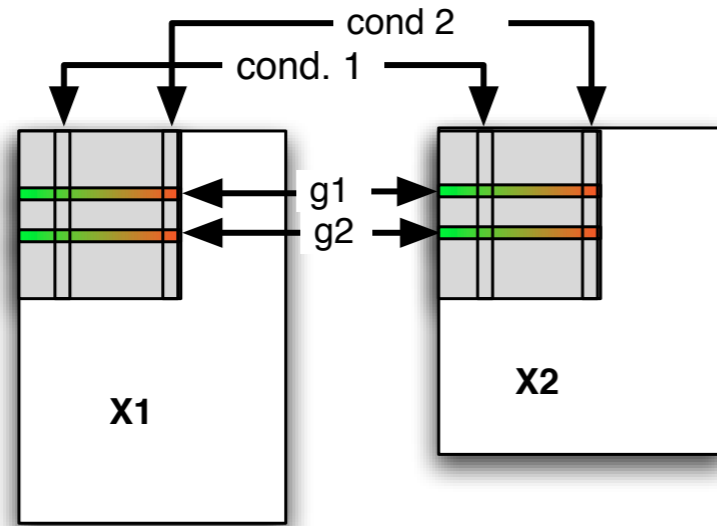
W/ PATRICK EICHENBERGER, NYU
W/ ERIC ALM, MIT, BROAD
W/ HARRY OSTRER

Previous Multi-Species Comparisons

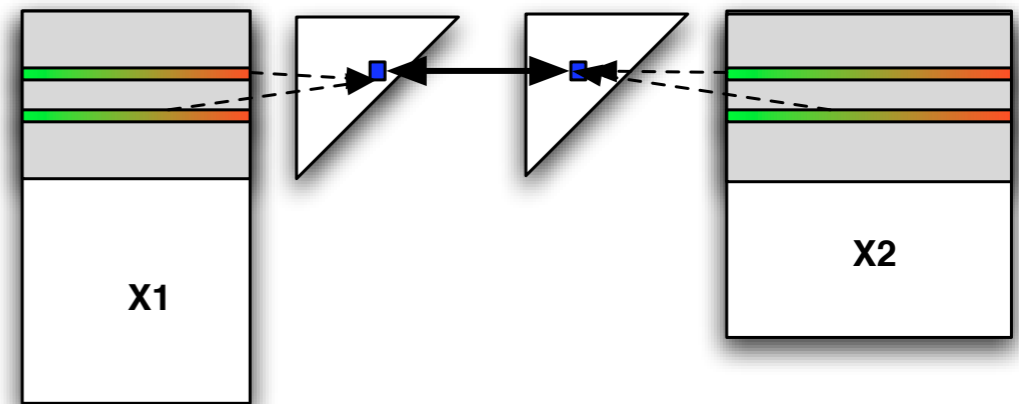
- McCarroll, Murphy, Zou, et al (2004, Nature Genetics)
- Ihmels, Bergmann, Berman, Barkai (2005, PLoS Genetics)
- **Tirosh, Barkai (2007, Genome Biology)**

3 classes of multi-species comparisons

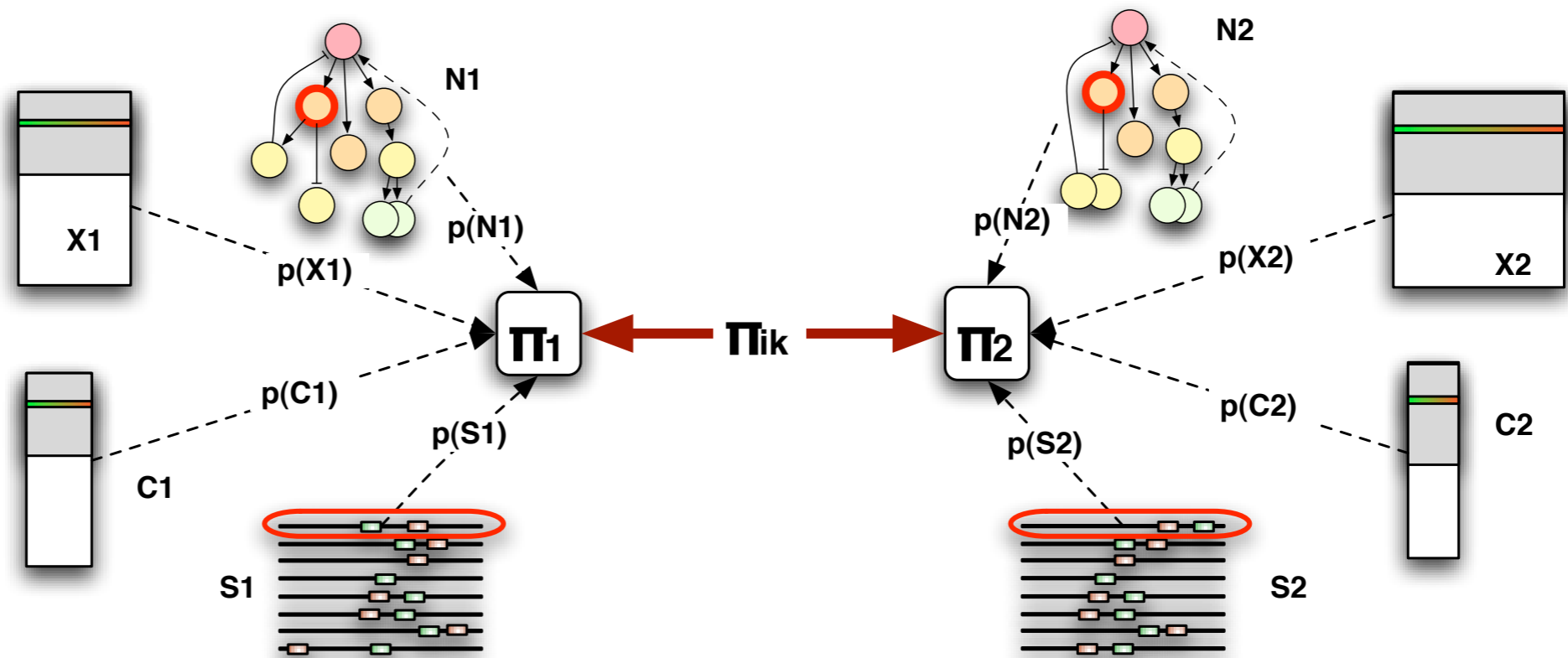
A. Class I. Matched conditions



B. Class II. Co-expression

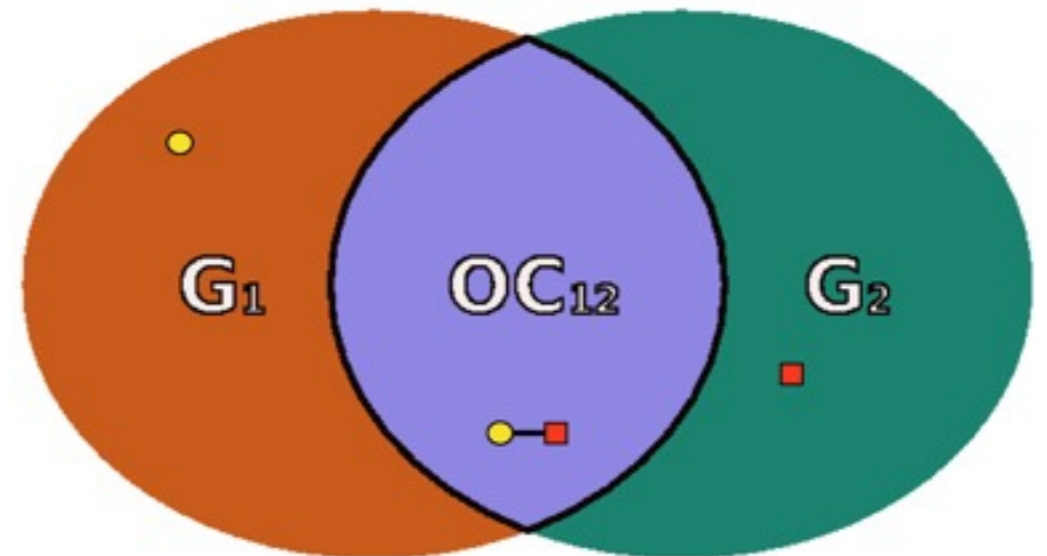


C. Multi-data+multi-species cMonkey:



Proposed Multi-species cMonkey model

- Given 2 genomes, G_1 & G_2 :
 - OC_1 & OC_2 as the set of genes in G_1 & G_2 with orthologs in the other
 - Define OC_{12} as the set of all putative orthologous pairs, including:
 - 1-to-1
 - 1-to-many
 - many-to-many

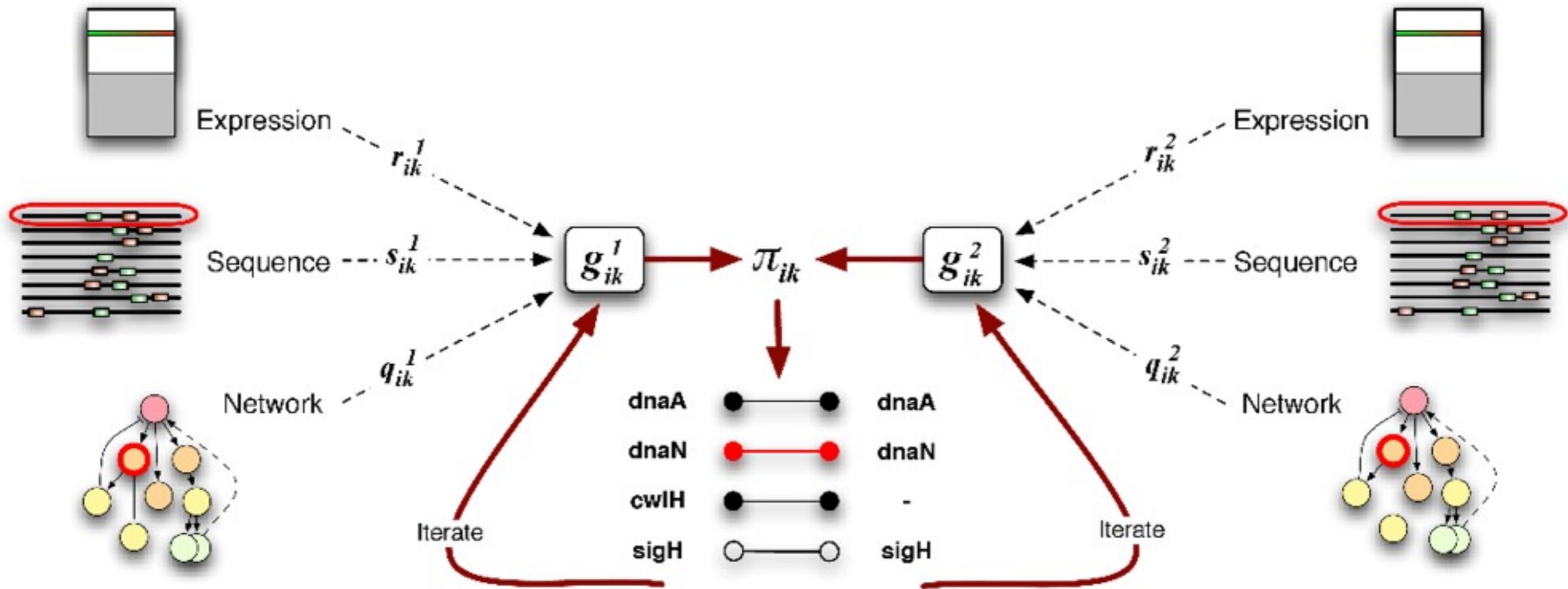


Proposed Multi-species cMonkey model

Algorithm outline:

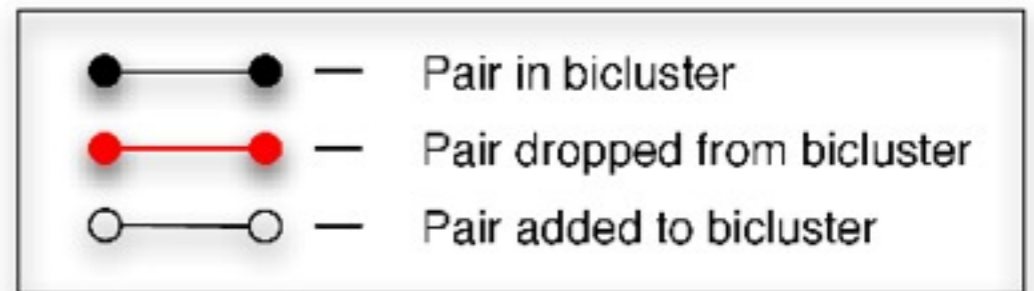
- Shared-space search: optimize biclusters in OC_{12} space
 - Optimize each OC_{12} bicluster within “species data space”
don't merge data
 - Add/drop a gene-pair from OC_{12} based on evolving single species models
 - *What to do if a gene exhibits correlation to bicluster in one species, but its ortholog in other does not? (answer coming)*

Shared Optimization



$$\pi_{ik} \equiv p(y_{ik}^1, y_{ik}^2 | g_{ik}^1, g_{ik}^2) \propto e^{(\beta_0 + \beta_1 (g_{ik}^1 + g_{ik}^2))}$$

$$g_{ik}^j \equiv r_0 \log(\tilde{r}_{ik}^j) + s_0 \log(\tilde{s}_{ik}^j) + \sum_{n \in N} q_0^n \log(\tilde{q}_{nik}^j)$$

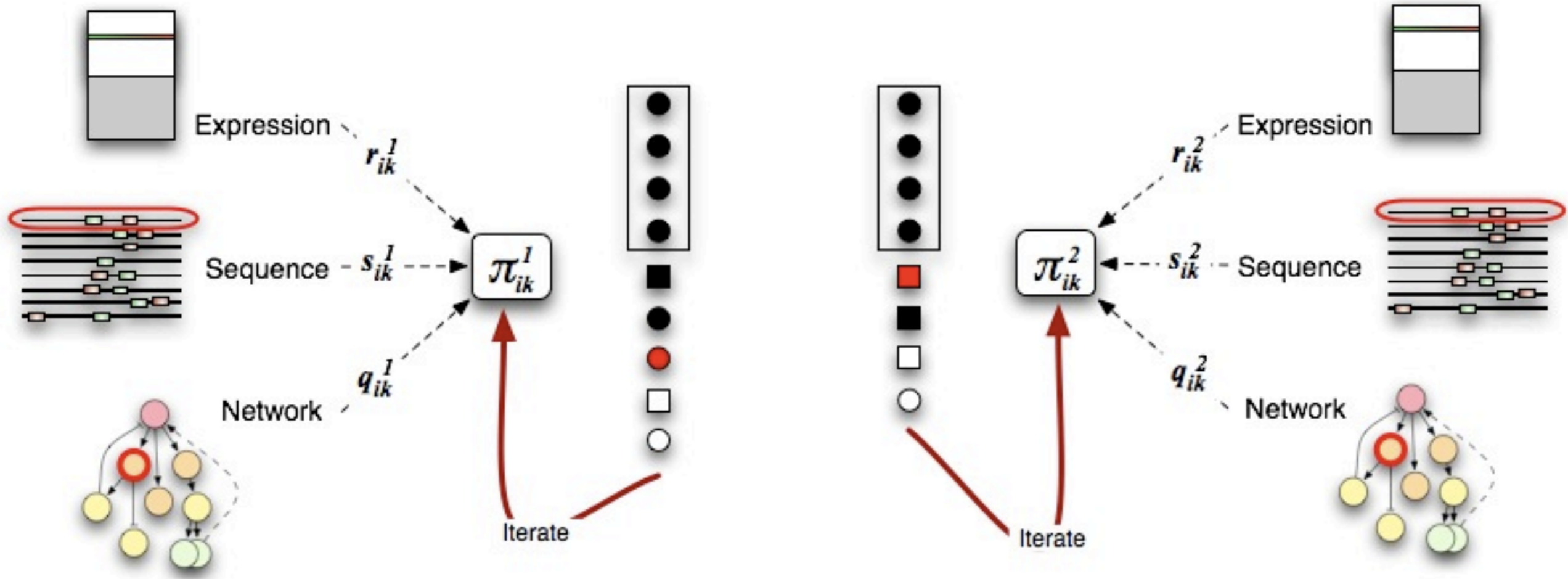


Proposed Multi-species cMonkey model

Algorithm outline:

- Elaborate: optimize OC_{12} biclusters in each organism's "species space"
 - Seed with genes from pairs in the OC_{12} biclusters
 - Use original single-species *cMonkey* to optimize the OC_{12} seeds:
 - Cannot drop genes from original OC_{12} gene-pairs
 - Allow genes from entire genome to be added, i.e. species-specific, "orthologous core" and paralogs.

Elaboration



$$\pi_{ik}^j \equiv p(y_{ik} | X_k^j, S_i^j, N_{ik}^j) \propto e^{(\beta_0 + \beta_o g_{ik}^j)}$$

- — Orthologous core gene, from original shared bicluster
- — Orthologous core gene, in elaborated bicluster
- — Orthologous core gene, added to elaborated bicluster
- — Orthologous core gene, added to elaborated bicluster
- — Species-specific gene in bicluster
- — Species-specific gene dropped from bicluster
- — Species-specific gene added to bicluster

Proposed Multi-species cMonkey model

Algorithm outline:

- Extend: find new biclusters for G_j in its own “species space”
 - Seed & optimize new clusters following original *cMonkey* single-species model
 - Allow extend step to consider genes from orthologous core (*OC*)?:
 - Yes (*currently, we allow overlap potential to reduce over-sampling of explored modules*)
 - No (*possible future direction to force identification of species-specific modules*)

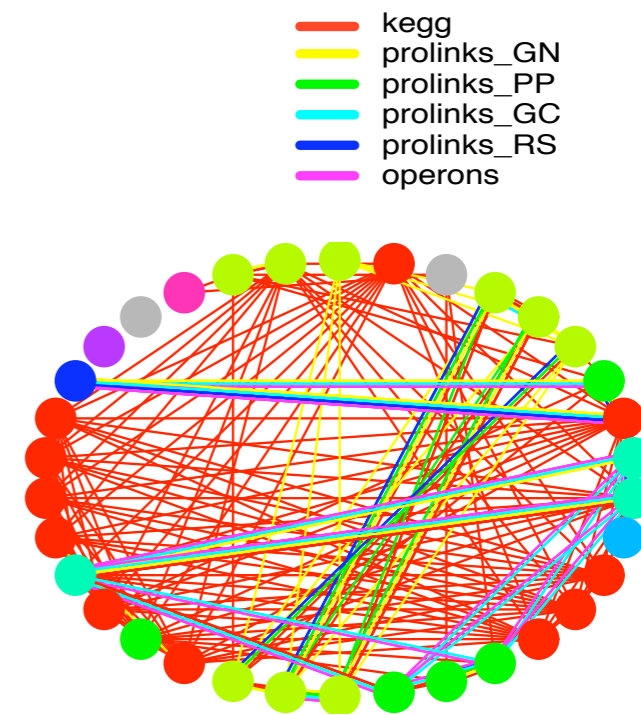
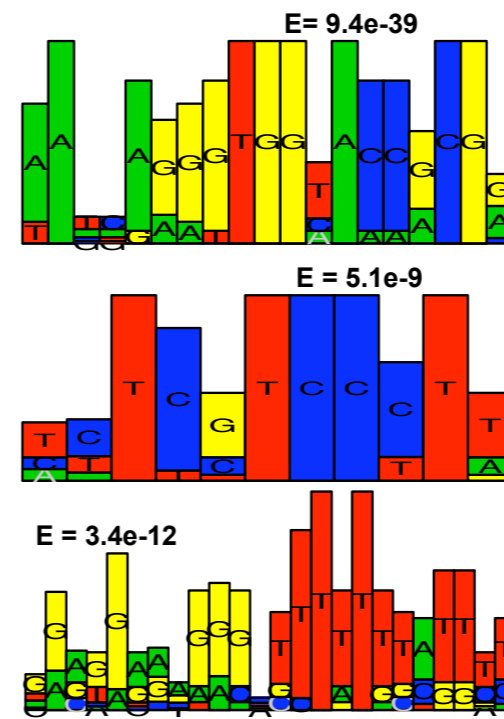
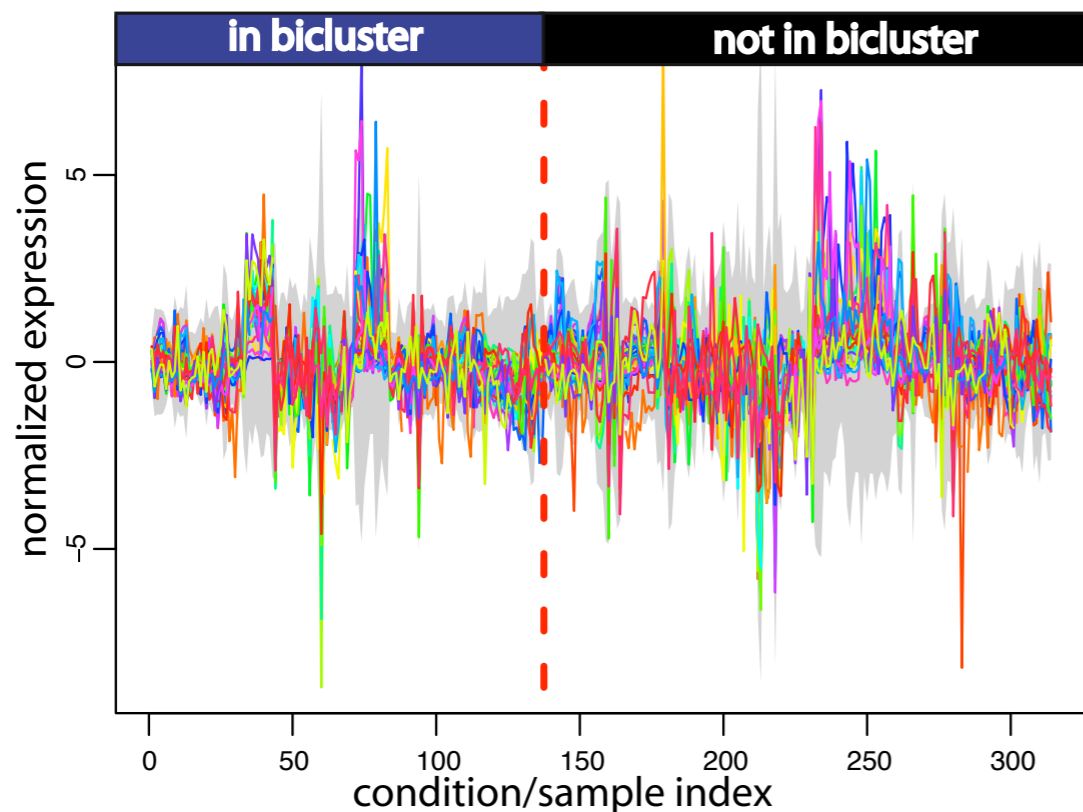
Species Analyzed

- Compared 3 bacterial species:
 - *Bacillus subtilis*
 - *Bacillus anthracis* (Anthrax)
 - *Listeria monocytogenes* (Listeriosis)
- 3 organisms → 3 pairings
 - Inparanoid to identify orthologs and orthologous families
 - 150 biclusters generated per pairing

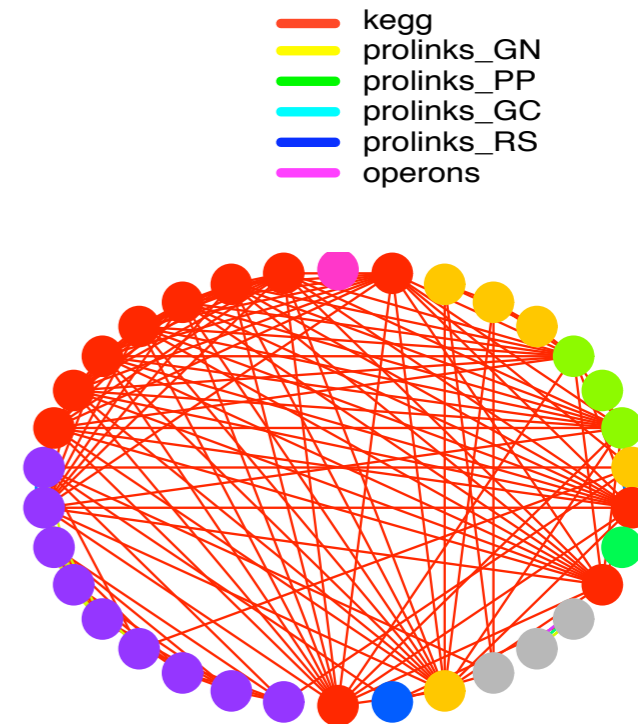
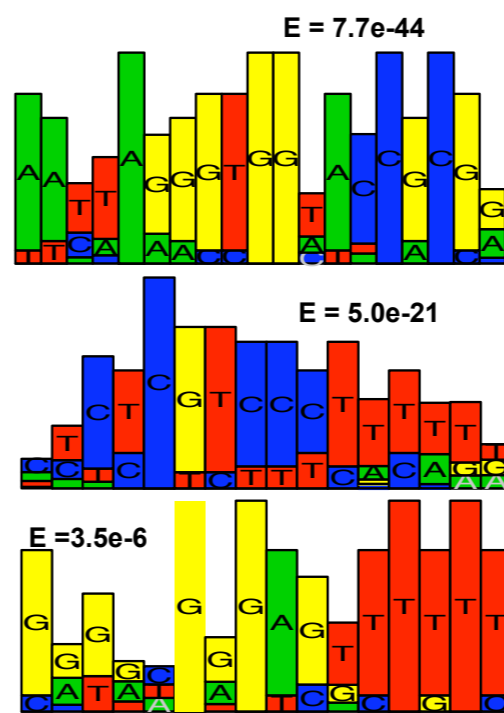
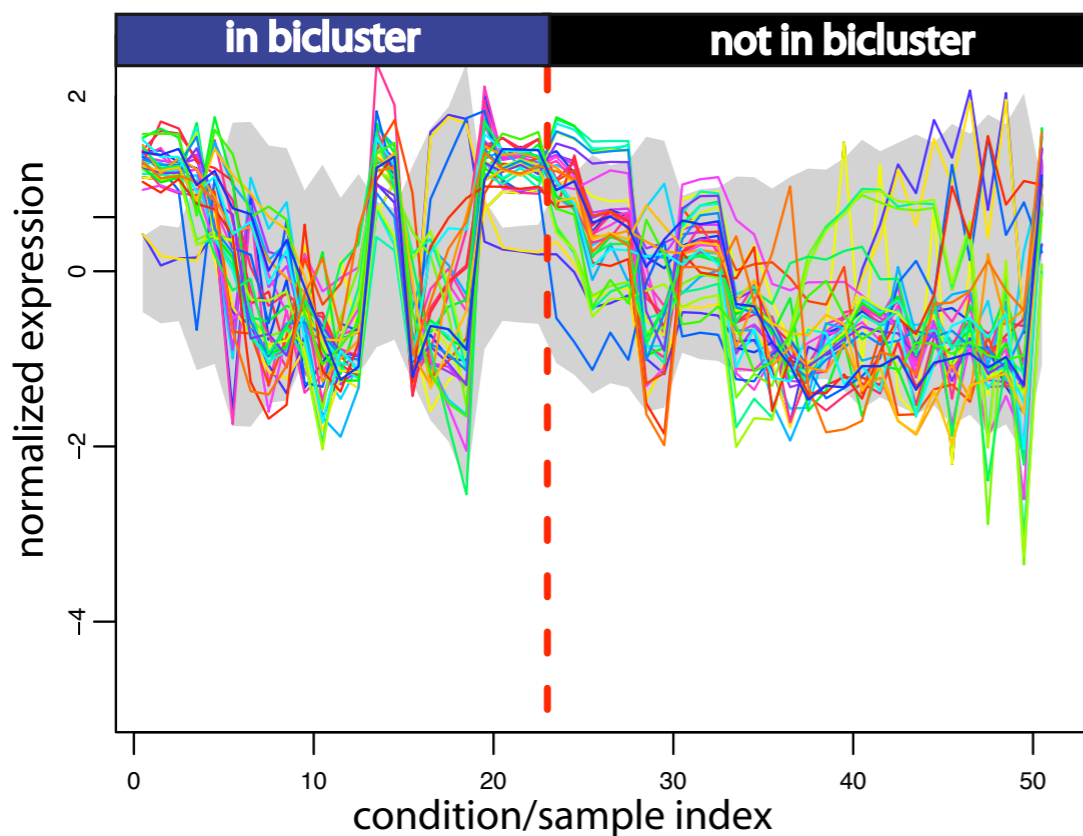
Number of:

	<i>B. subtilis</i> – <i>B. anthracis</i>	<i>B. subtilis</i> – <i>L. monocytogenes</i>	<i>B. anthracis</i> – <i>L. monocytogenes</i>
orthologous groups	2225	1439	1494
orthologous pairs	2443	1564	1690
unique genes (per organism)	<i>B. subtilis</i> : 2279/3928 <i>B. anthracis</i> : 2339/5865	<i>B. subtilis</i> : 1519/3928 <i>L. mono</i> : 1478/2795	<i>B. anthracis</i> : 1634/5865 <i>L. mono</i> : 1537/2795

A. *B. subtilis*:



B. *B. anthracis*:



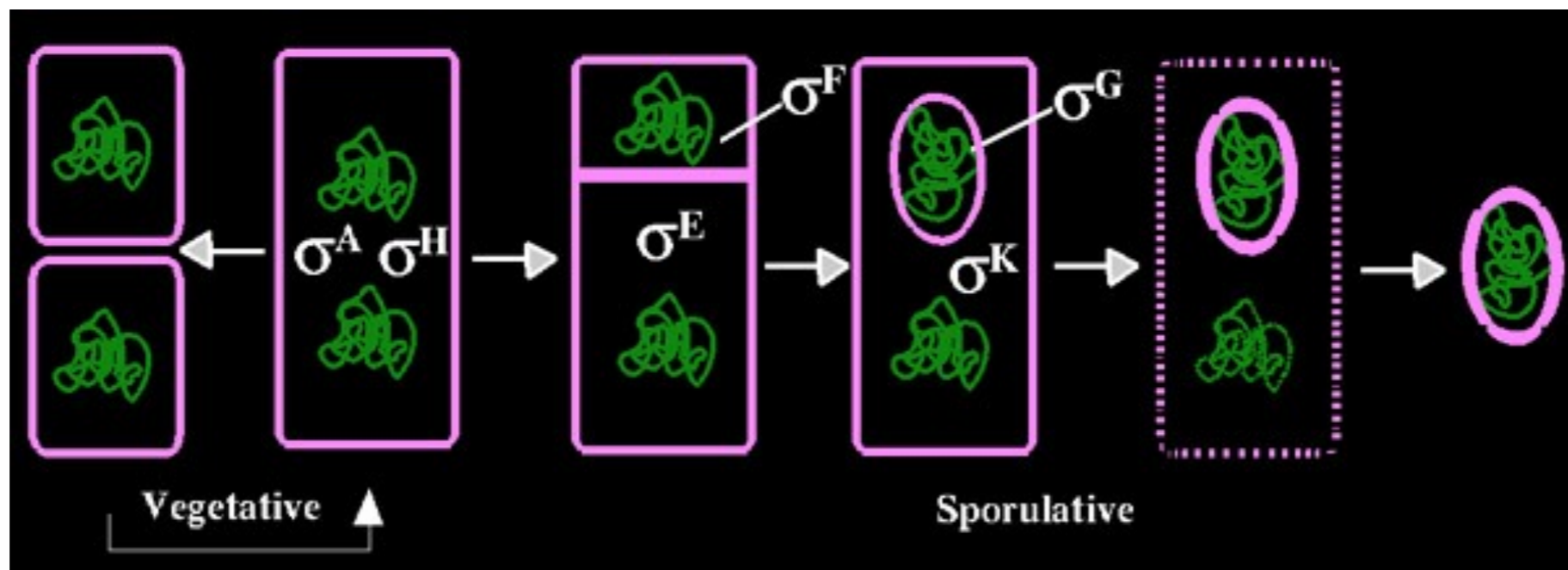
PETER WALTMAN

σ^E biclusters

(*B. subtilis* – *B. anthracis*)

Significantly enriched for sporulation genes (σ^E regulated):

- Bicluster 17 (includes Metabolism, Glutamine Transport, Transporters genes)
- Bicluster 35 (includes Metabolism, Glutamine Transport, Detoxification, Transporters genes)
- Bicluster 84 (includes Metabolism, Glutamine Transport genes)



http://biology.kenyon.edu/courses/biol114/Chap11/spore_cycle.jpg

• *B. anthracis* Waves of Gene Expression

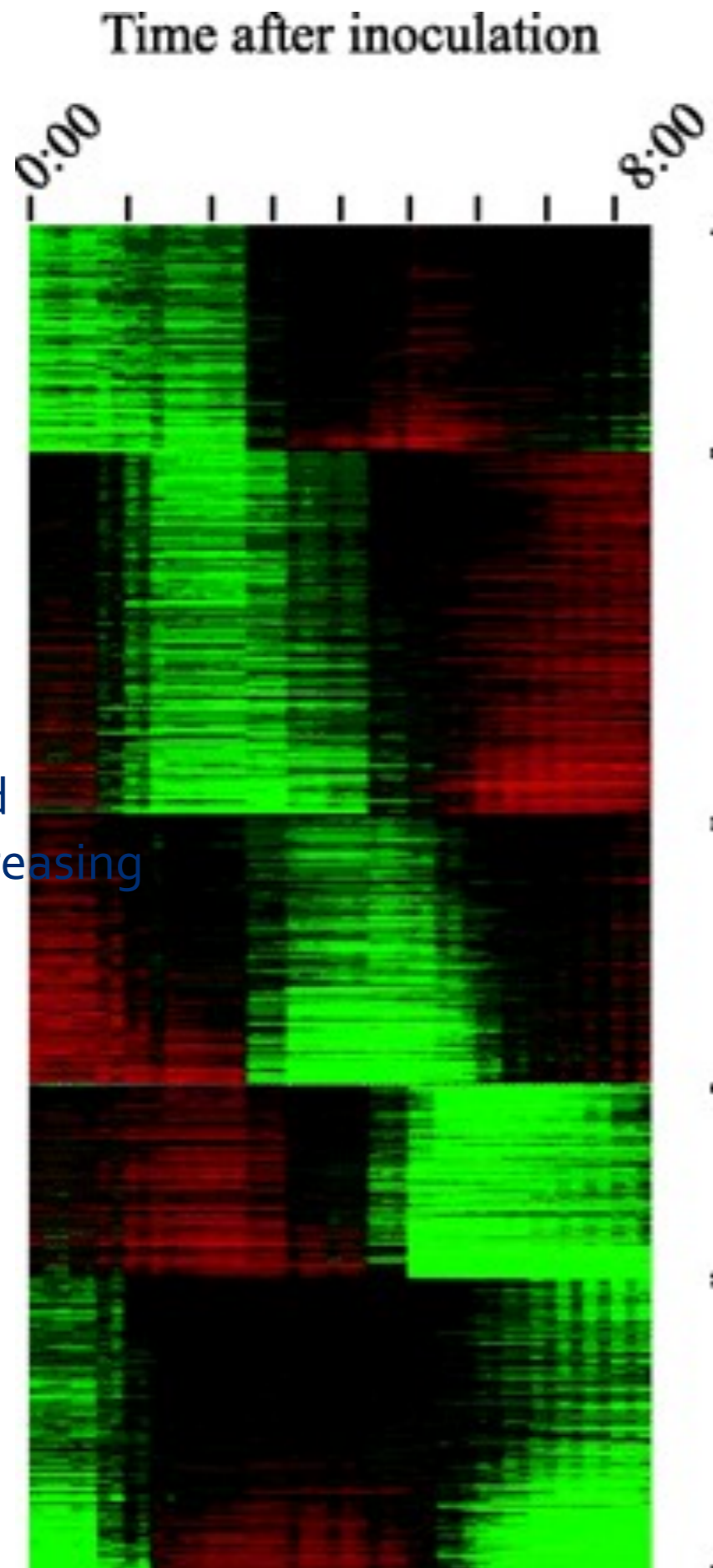
• Germination and early outgrowth

• Rapid Growth

• Rapid Growth and Responding to increasing toxic environment

• Sporulation and Oxidative Stress

• Sporulation and early germination and outgrowth



Wave 1

Wave 2

Wave 3

Wave 4

Wave 5

• Biclusters 8₄ : 3 into 4

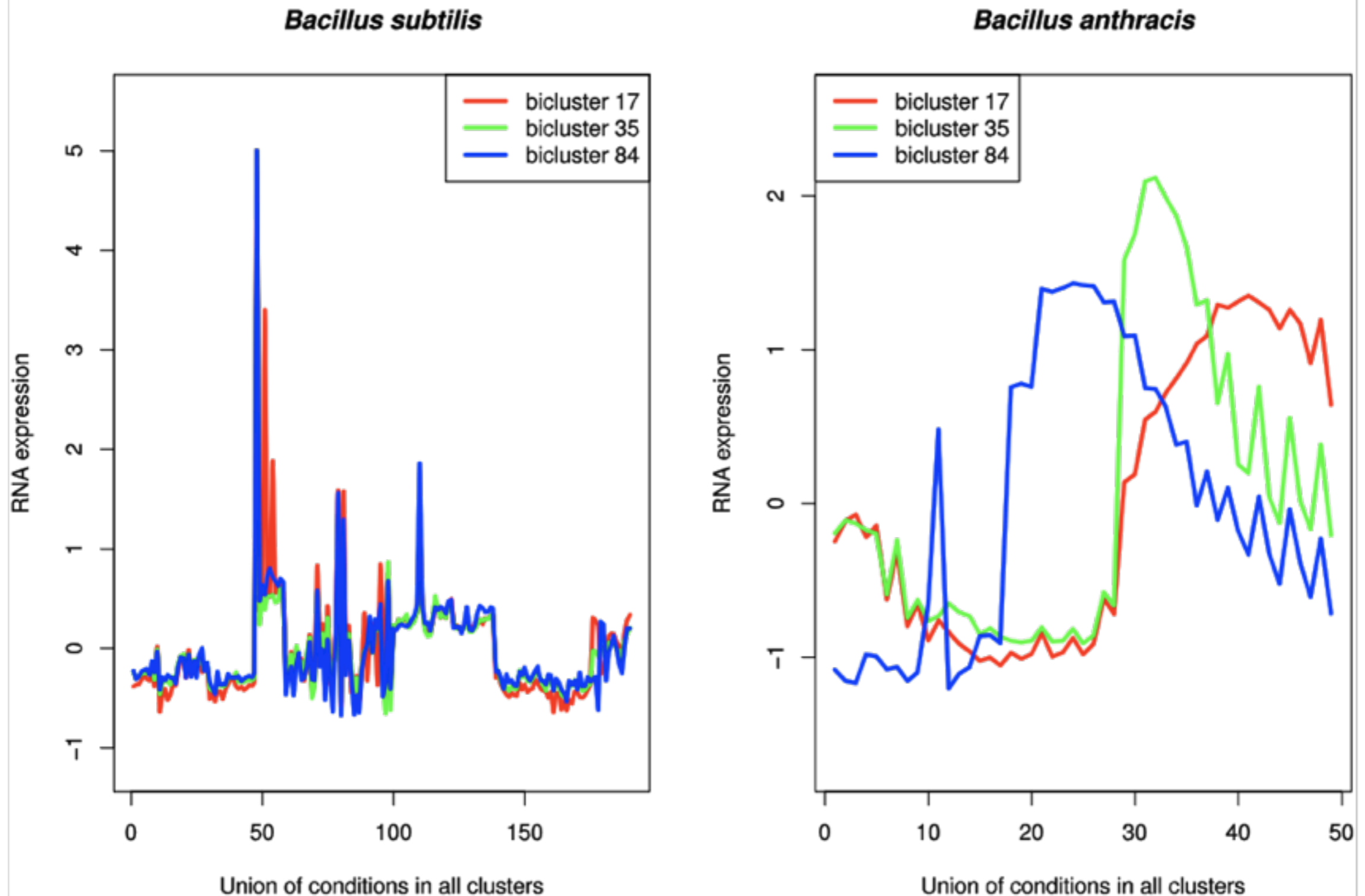
• Biclusters 3₅ : 4

• Biclusters 1₇ : 4 into 5

• (Bergmen *et al.*, 2006)

Sporulation Biclusters

Comparing Mean Expression of clusters 17,35,84

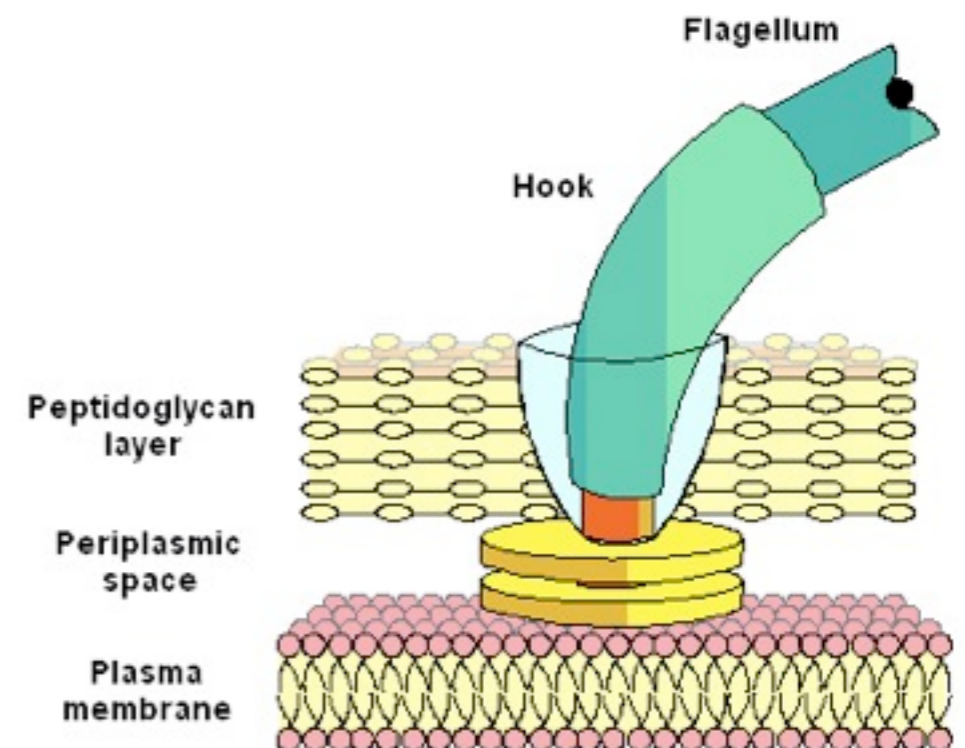


Flagellar Assembly Biclusters

- Flagellar Assembly biclusters for all 3 organisms
- *B. anthracis* thought to be non-motile:
 - Missing σ^D (flagellar TF in *B. subtilis*)
 - frameshifts to 4 critical flagellar genes:
 - *cheA*
 - *flgL*
 - *fliF* (MS ring)
 - *fliM* (C ring component)
- *B. anthracis* biclusters also enriched for:
 - Chemotaxis
 - Type III secretion system*

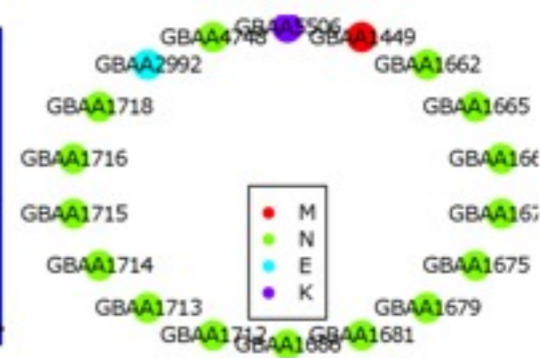
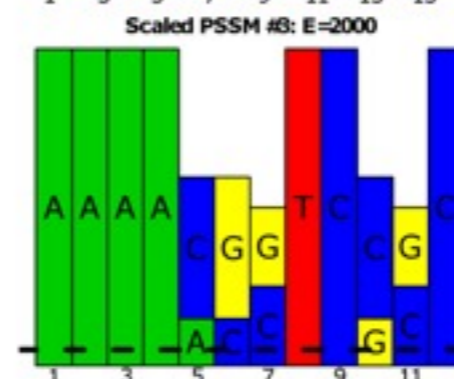
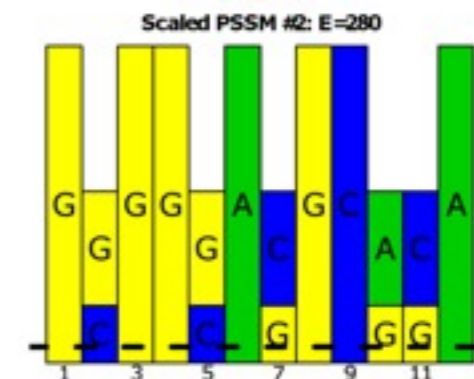
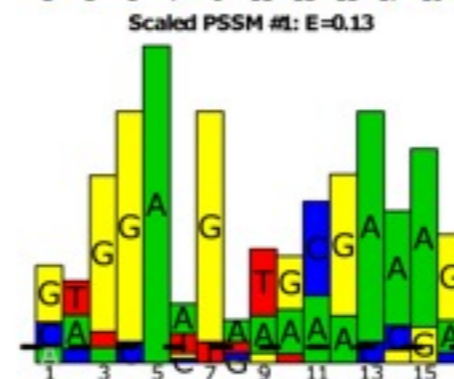
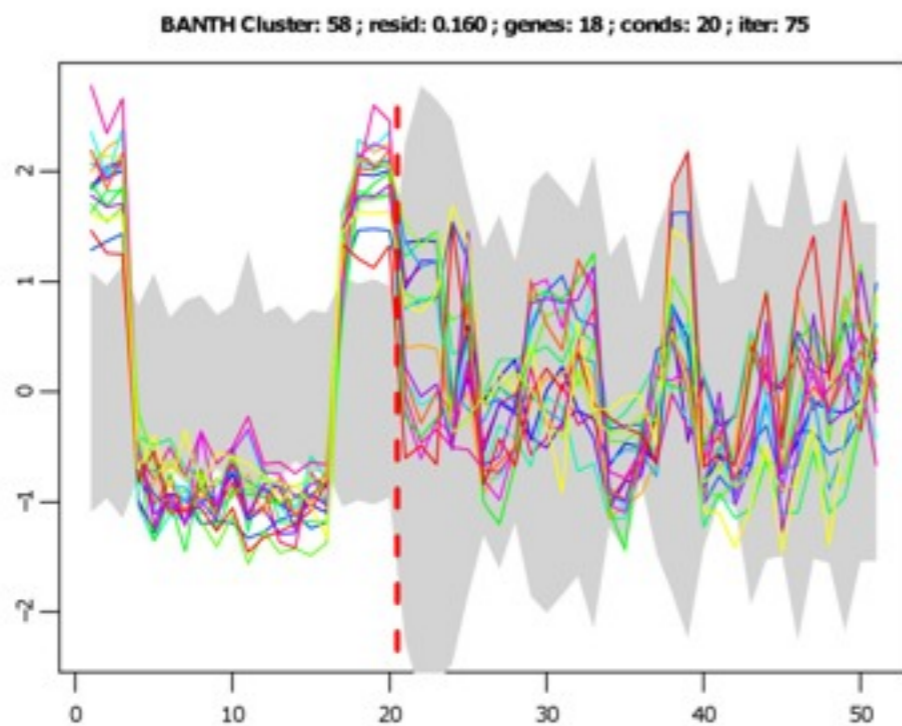
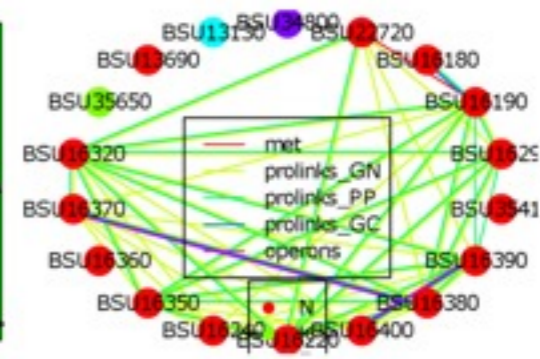
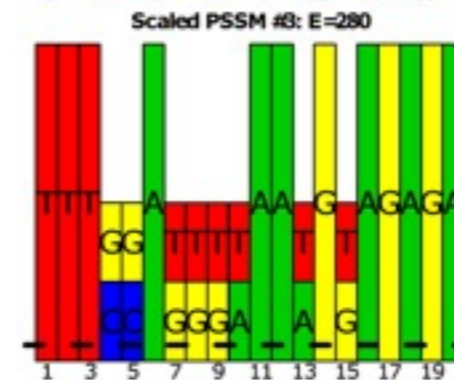
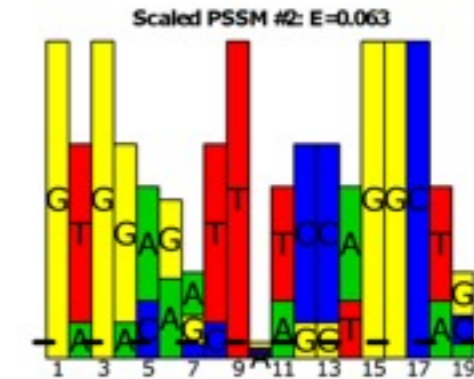
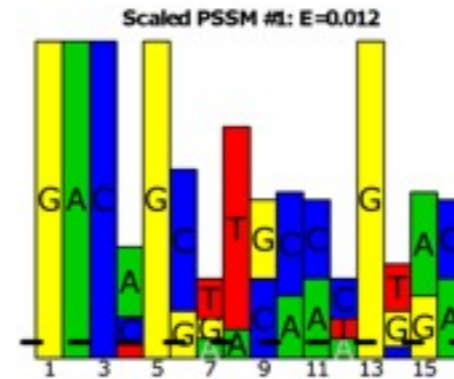
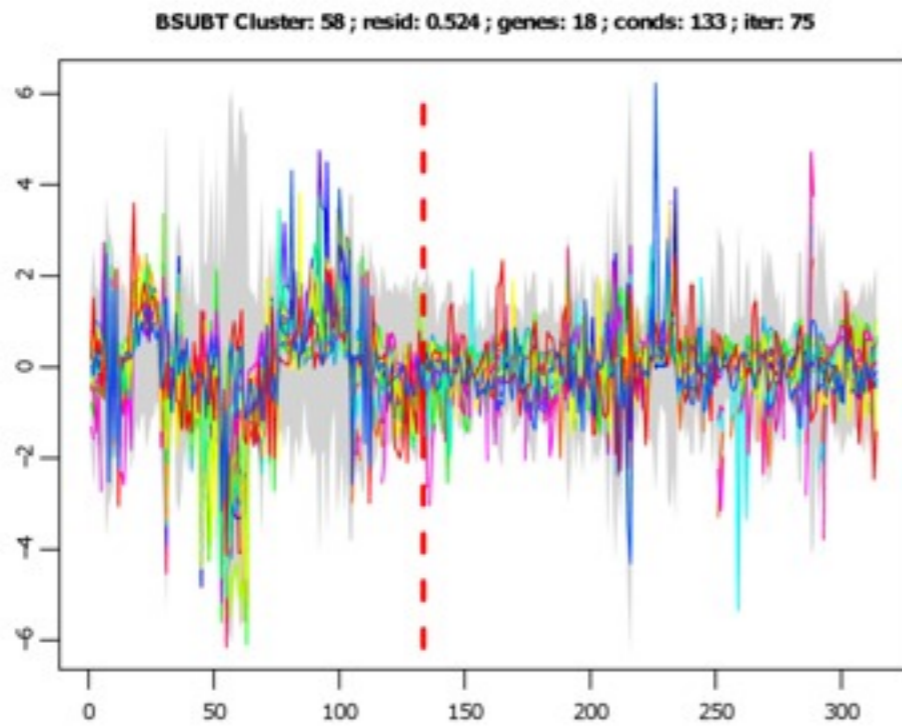
* *B. subtilis* – *B. anthracis* only

Ultrastructure of Gram-Positive Bacterial Flagella



<http://www.conceptdraw.org/sampletour/medical/GPositiveBFlagella.gif>

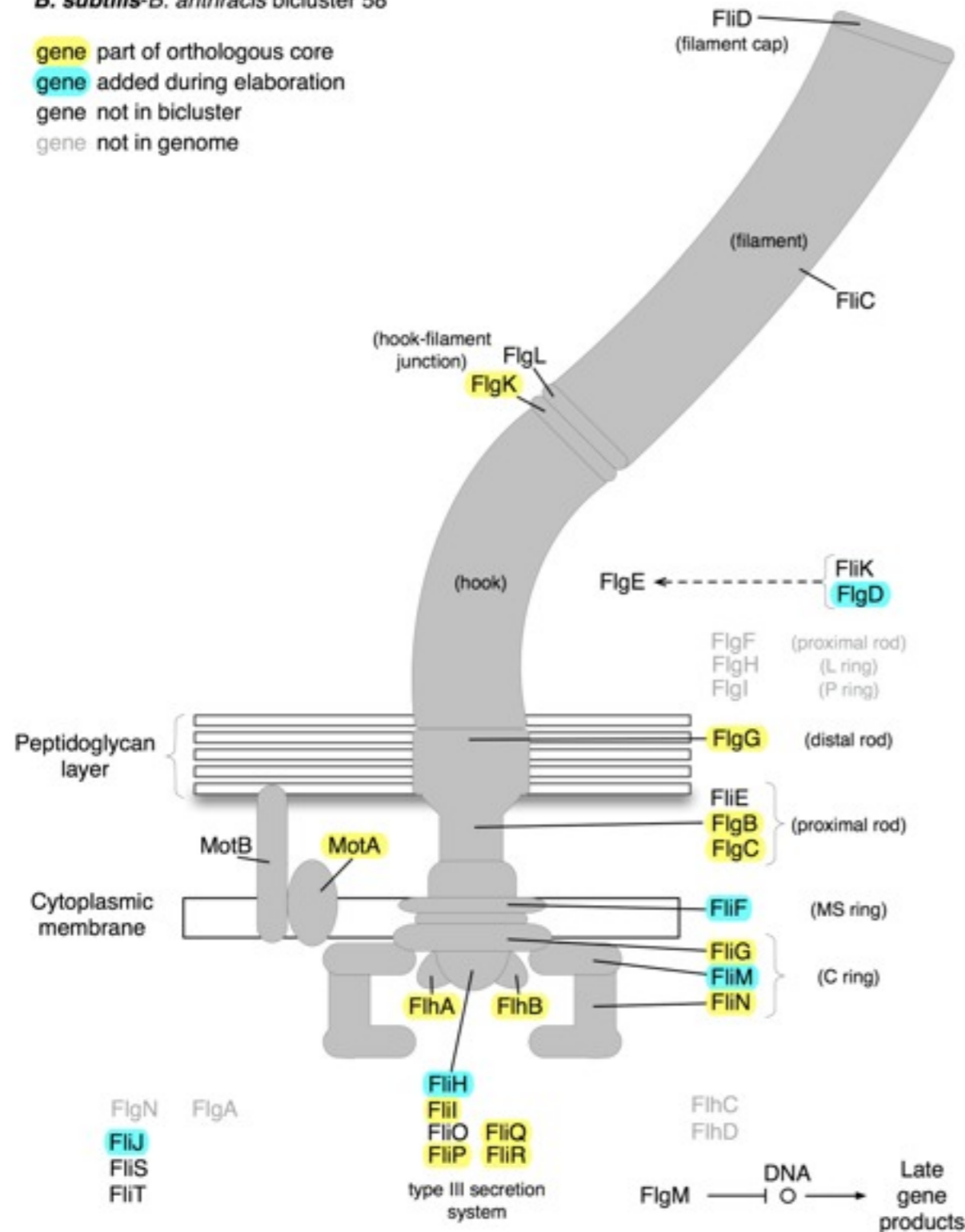
Shared *B. subtilis* - *B. anthracis* Flagellar Assembly Bicluster



Elaborated *B. subtilis* - *B. anthracis* Flagellar Assembly Bicluster

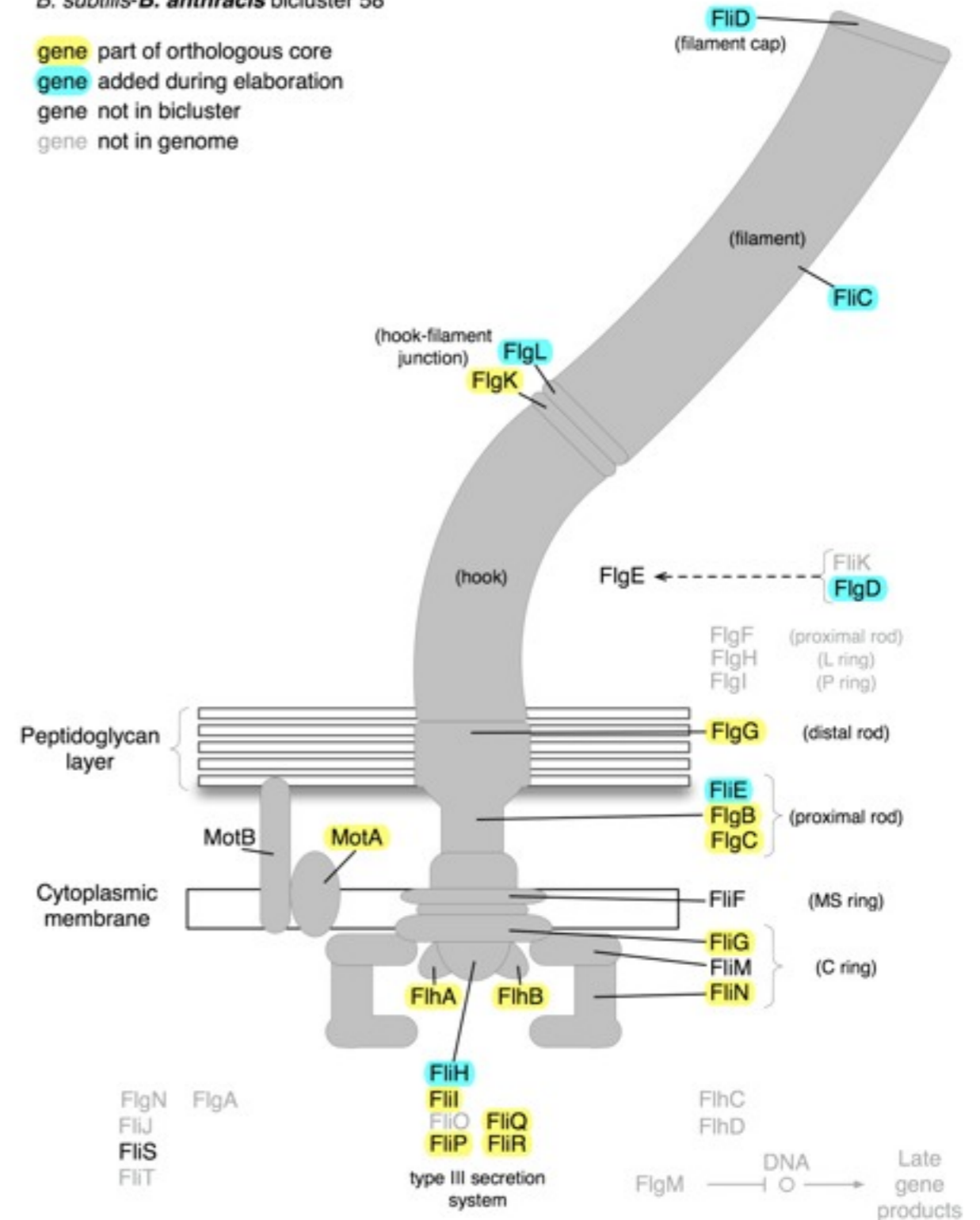
B. subtilis-*B. anthracis* bicluster 58

gene part of orthologous core
 gene added during elaboration
 gene not in bicluster
 gene not in genome



B. subtilis-*B. anthracis* bicluster 58

gene part of orthologous core
 gene added during elaboration
 gene not in bicluster
 gene not in genome



Globally Validating Multi-Species method

- **Issues**

- No solved organism as validation - only partial solutions available
- Large number of results (12) to validate:
 - 3 organism-pairs → 6 results (2 for each pair)
 - 2 steps (shared & elaboration opt's) → 12 total
- No existing metric for measuring quality & conservation
 - Could we use either DCA or ECC for a metric?
 - DCA not a genuine clustering method & no metric provided
 - ECC gave inconsistent results in our own tests

How to measure or compare conservation & quality?

- Conservation metric
- Compare Shared & Elaboration optimizations with biclusters from *ideal* single-species *cMonkey*
 - Expression (*residuals*)
 - Networks (*association p-values*)
 - Sequence:
 - *Motif E-values*
 - *Sequence p-values*
 - Enrichments of:
 - GO terms
 - KEGG pathways

New Conservation Metric

$$Cons(bicl_i^1, bicl_j^2) = \frac{|matches(oc(bicl_i^1), oc(bicl_j^2))|}{|oc(bicl_i^1)|}$$

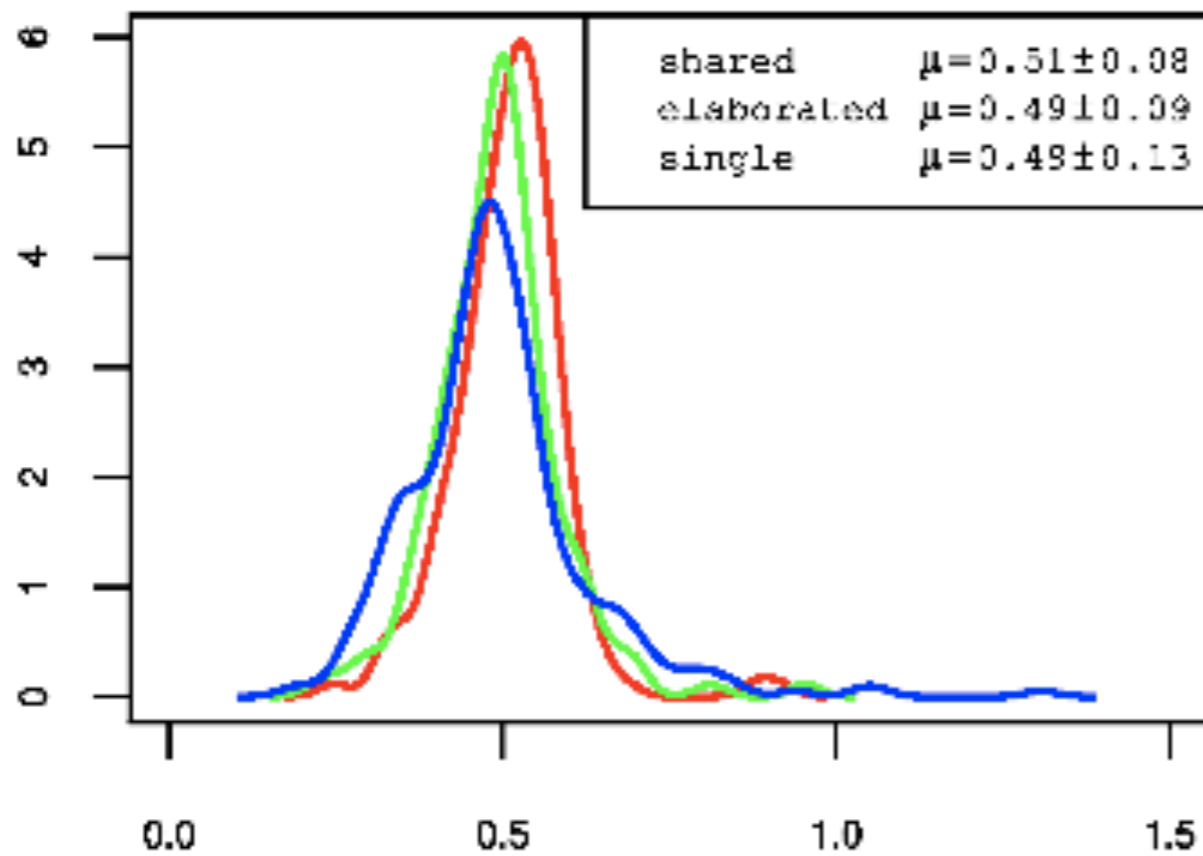
Find $\max(Cons(bicl_i^1, bicl_j^2))$ for each bicluster and average over all biclusters

Conservation Metric

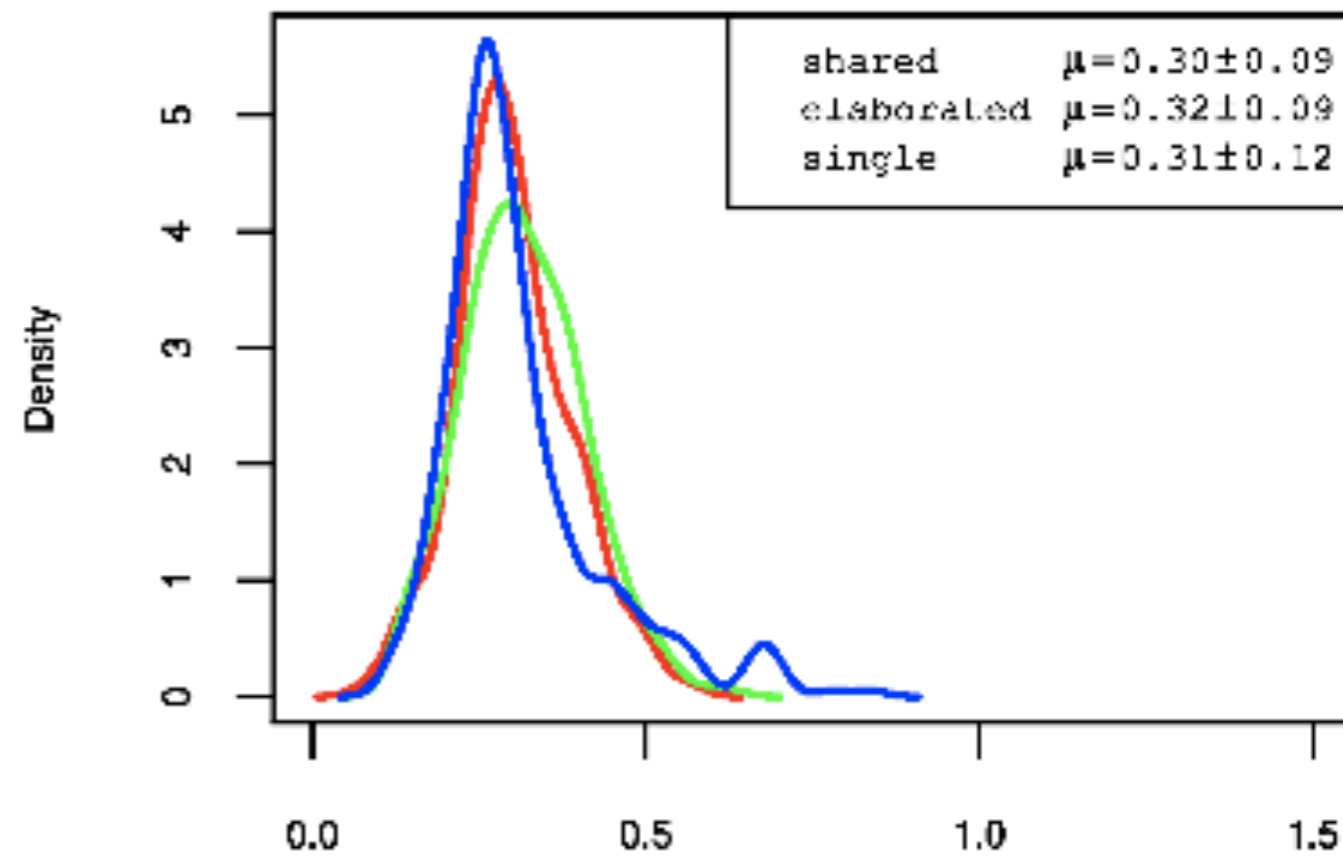
	<i>B. subtilis</i> - <i>B. anthracis</i>	<i>B. subtilis</i> - <i>L. monocytogenes</i>	<i>B. anthracis</i> - <i>L. monocytogenes</i>
Single	0.218	0.235	0.177
Elaborated	0.825	0.883	0.922
Shared	1	1	1

Residuals

Bacillus subtilis
(bsubt-barth)



Bacillus anthracis
(bsubt-barth)



Shared

Elaborated

Single

P-values from two-sided
Wilcoxon's rank test ($\alpha=0.01$):

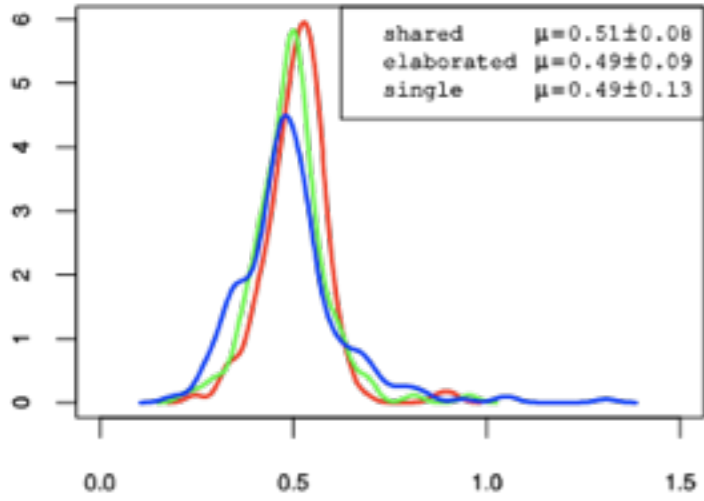
B. subtilis - *B. anthracis*

	<i>B. subt.</i>	<i>B. anth.</i>
Shared-Single	1.03E-03	0.46
Elaborated-Single	0.27	0.04

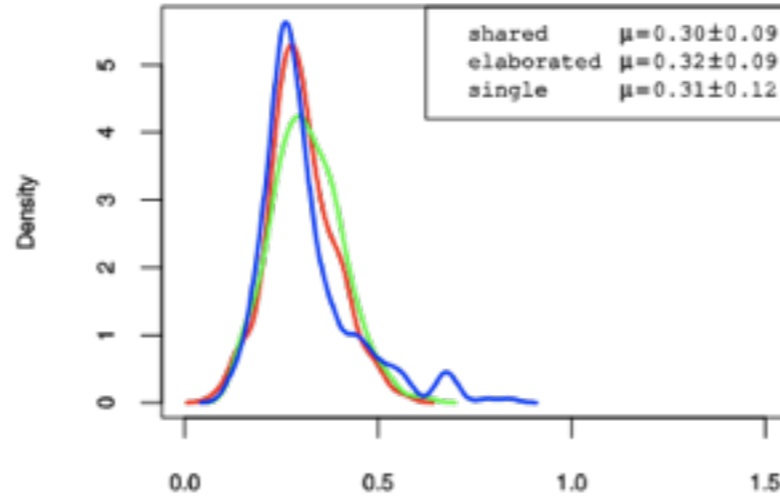
Residuals



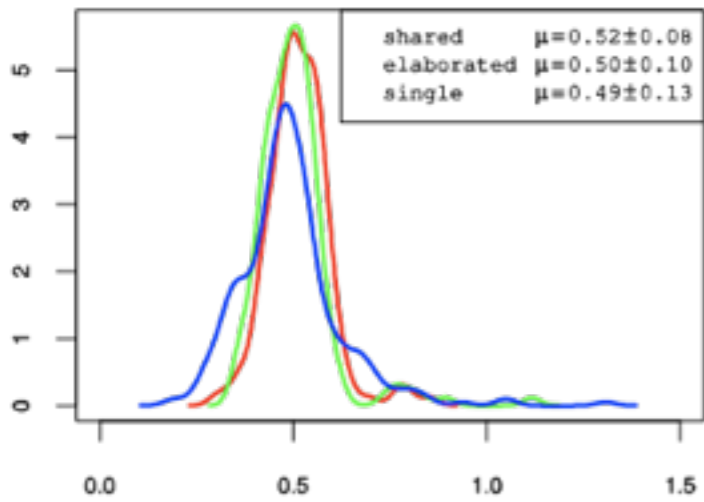
Bacillus subtilis
(bsubt-banth)



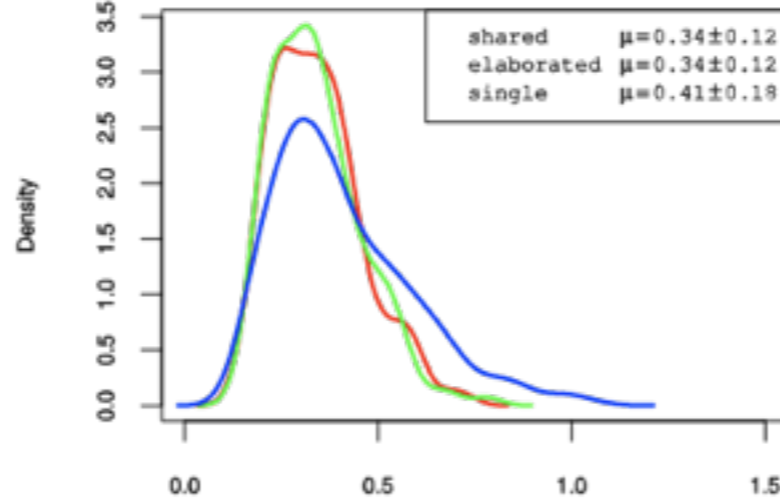
Bacillus anthracis
(bsubt-banth)



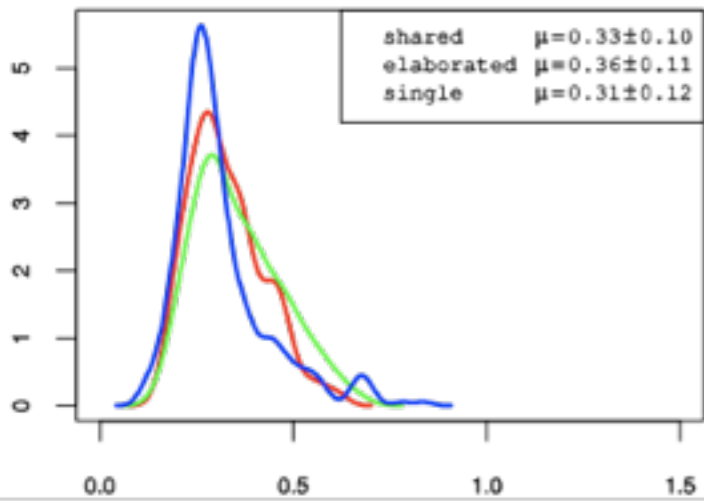
Bacillus subtilis
(bsubt-lmo)



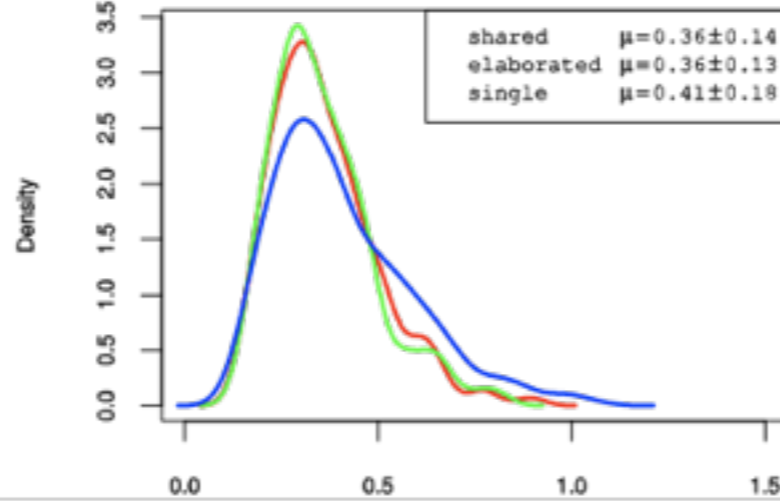
Listeria monocytogenes
(bsubt-lmo)



Bacillus anthracis
(banth-lmo)



Listeria monocytogenes
(banth-lmo)



$$resid(I, J) = \frac{\frac{1}{|I||J|} \sum_{i \in I, j \in J} abs(x_{ij} - x_{i\cdot} - x_{\cdot j} + x_{\cdot\cdot})}{\frac{1}{|I|} \sum_{i \in I} var_j(x_i)}$$

P-values from two-sided Wilcoxon's rank test ($\alpha=0.01$):

B. subtilis - B. anthracis

	<i>B. subt.</i>	<i>B. anth.</i>
Shared-Single	1.03E-03	0.46
Elaborated-Single	0.27	0.04

B. subtilis - L. monocytogenes

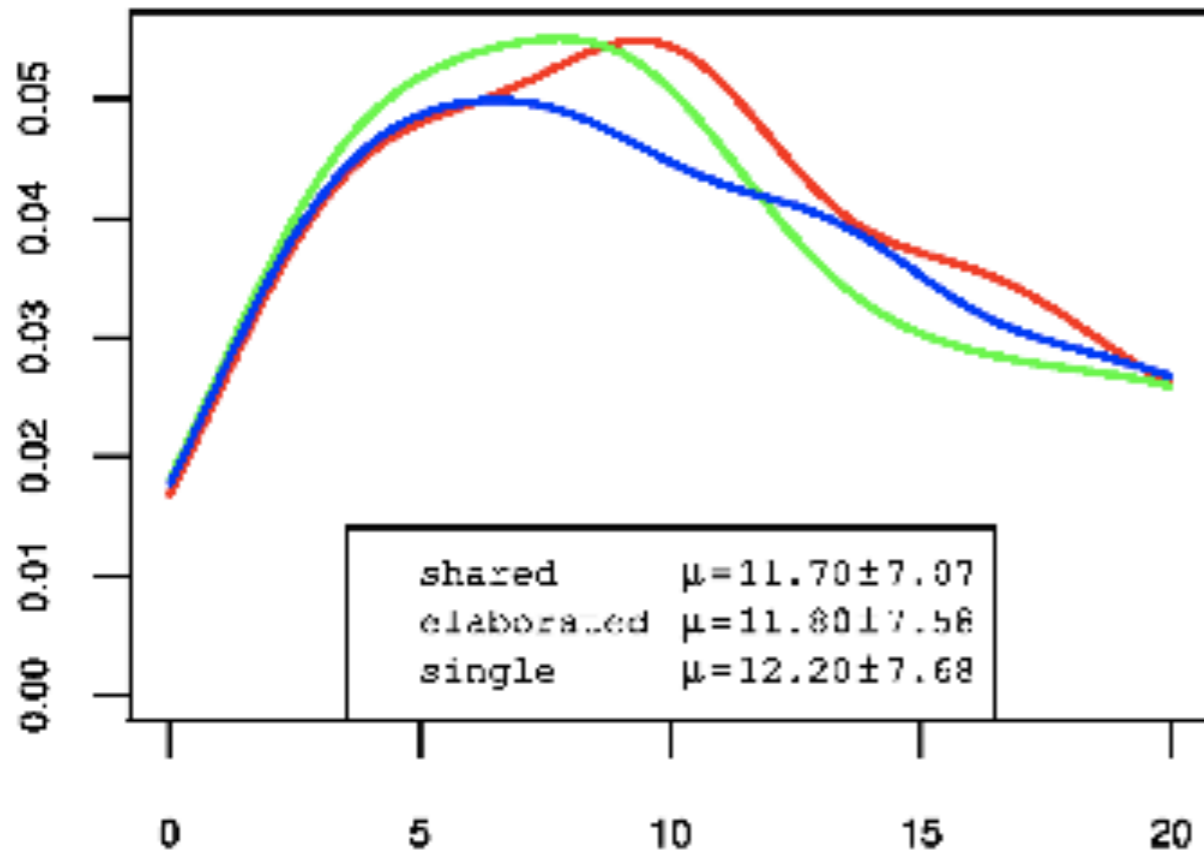
	<i>B. subt.</i>	<i>L. mono.</i>
Shared-Single	2.82E-04	7.81E-04
Elaborated-Single	0.09	3.98E-04

B. anthracis - L. monocytogenes

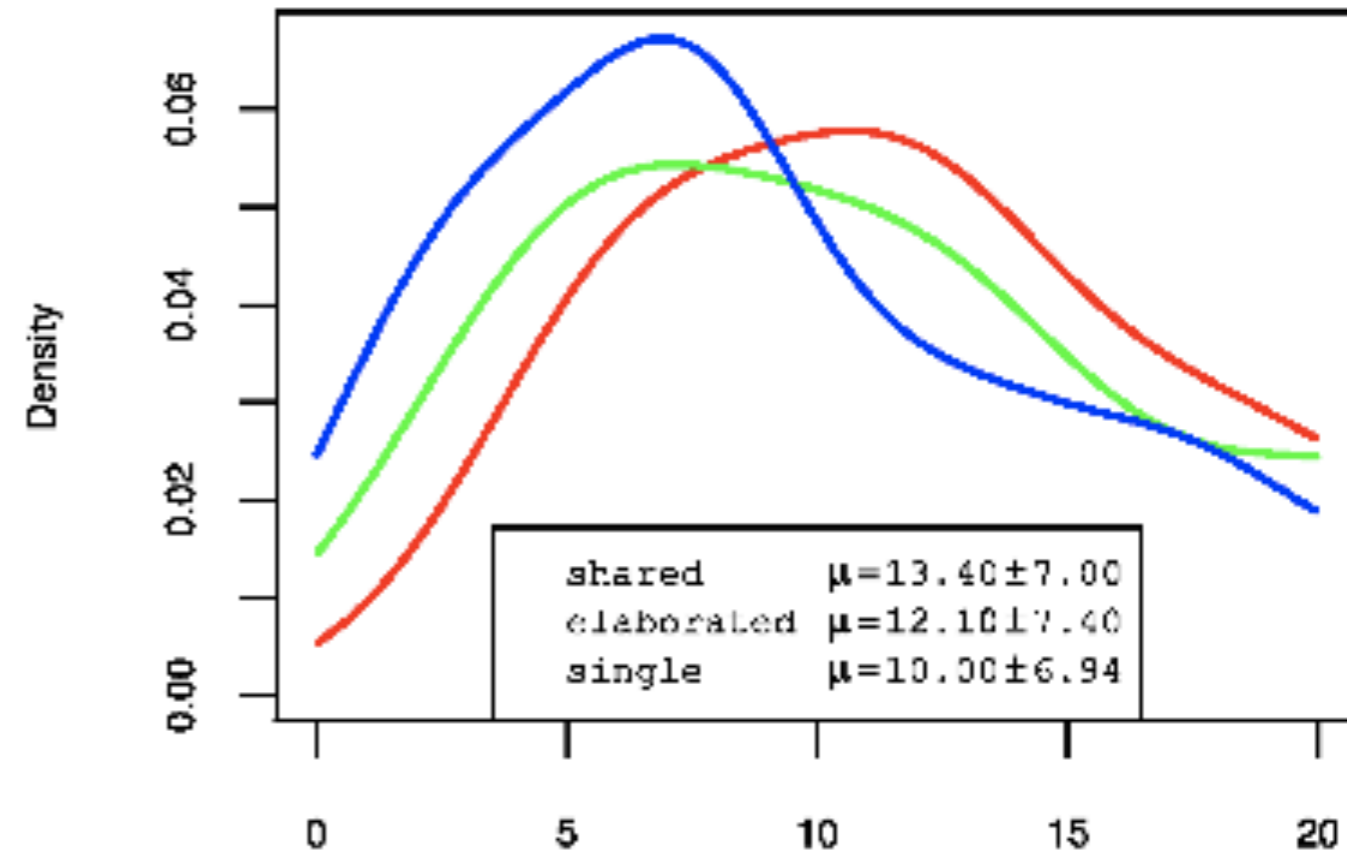
	<i>B. anth.</i>	<i>L. mono</i>
Shared-Single	0.01	0.02
Elaborated-Single	1.68E-06	0.01

Association p-values (-log10)

Bacillus subtilis
(bsubt-banth)



Bacillus anthracis
(bsubt-banth)



Shared

Elaborated

Single

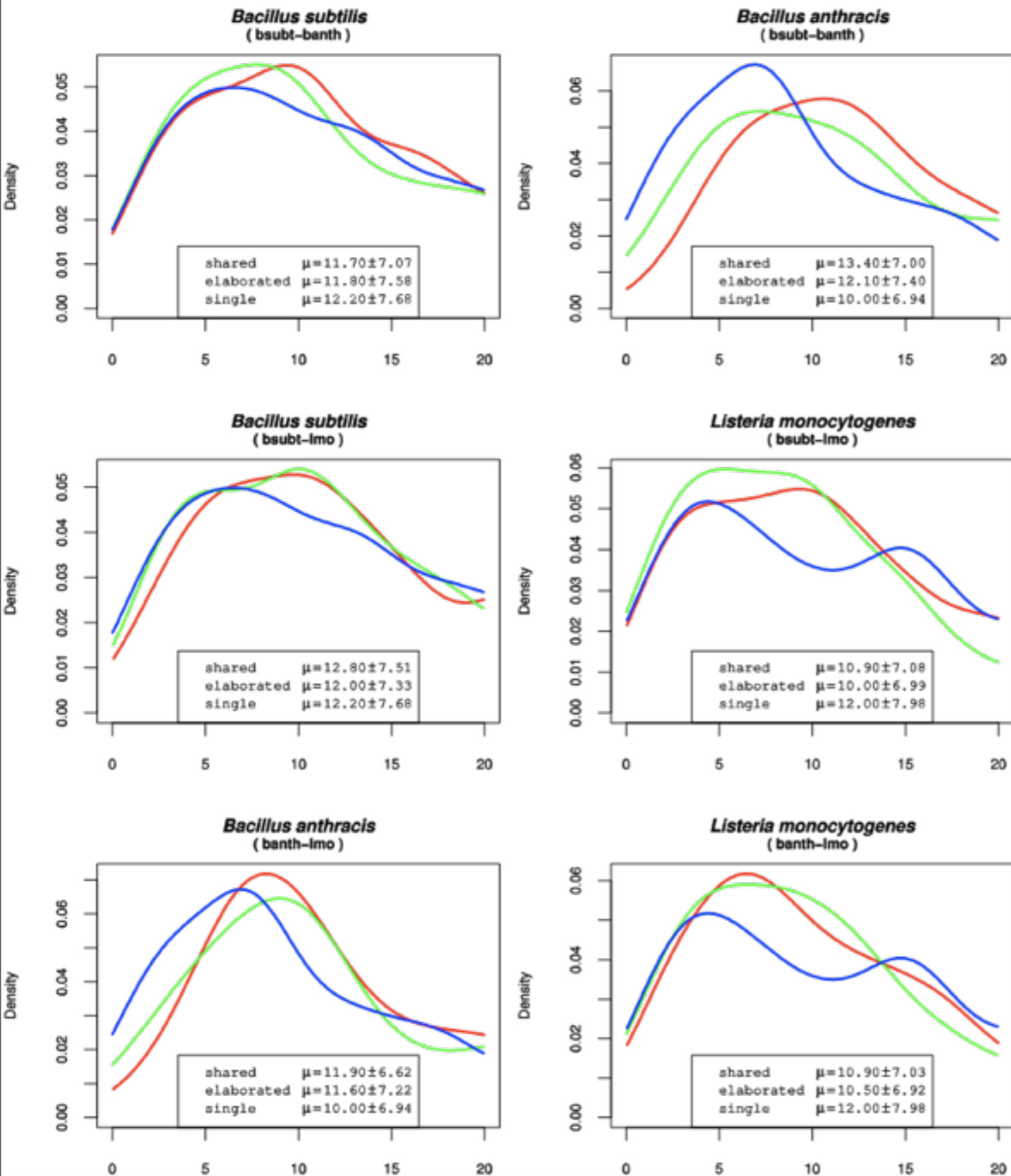
P-values from two-sided
Wilcoxon's rank test ($\alpha=0.01$):

B. subtilis - *B. anthracis*

	<i>B. subt.</i>	<i>B. anth.</i>
Shared-Single	0.40	1.52E-05
Elaborated-Single	0.34	0.01

Association P-values (-log10)

— SHARED — ELABORATED — SINGLE



$$pvalue(b_k, N) = \frac{\binom{|N|}{n_{b_k \rightarrow b_k}} (\text{poss}(G) - |N|)}{\binom{\text{poss}(G)}{\text{poss}(b_k)}}$$

P-values from two-sided Wilcoxon's rank test (α=0.01):

B. subtilis - *B. anthracis*

	<i>B. subt.</i>	<i>B. anth.</i>
Shared-Single	0.40	1.52E-05
Elaborated-Single	0.34	0.01

B. subtilis - *L. monocytogenes*

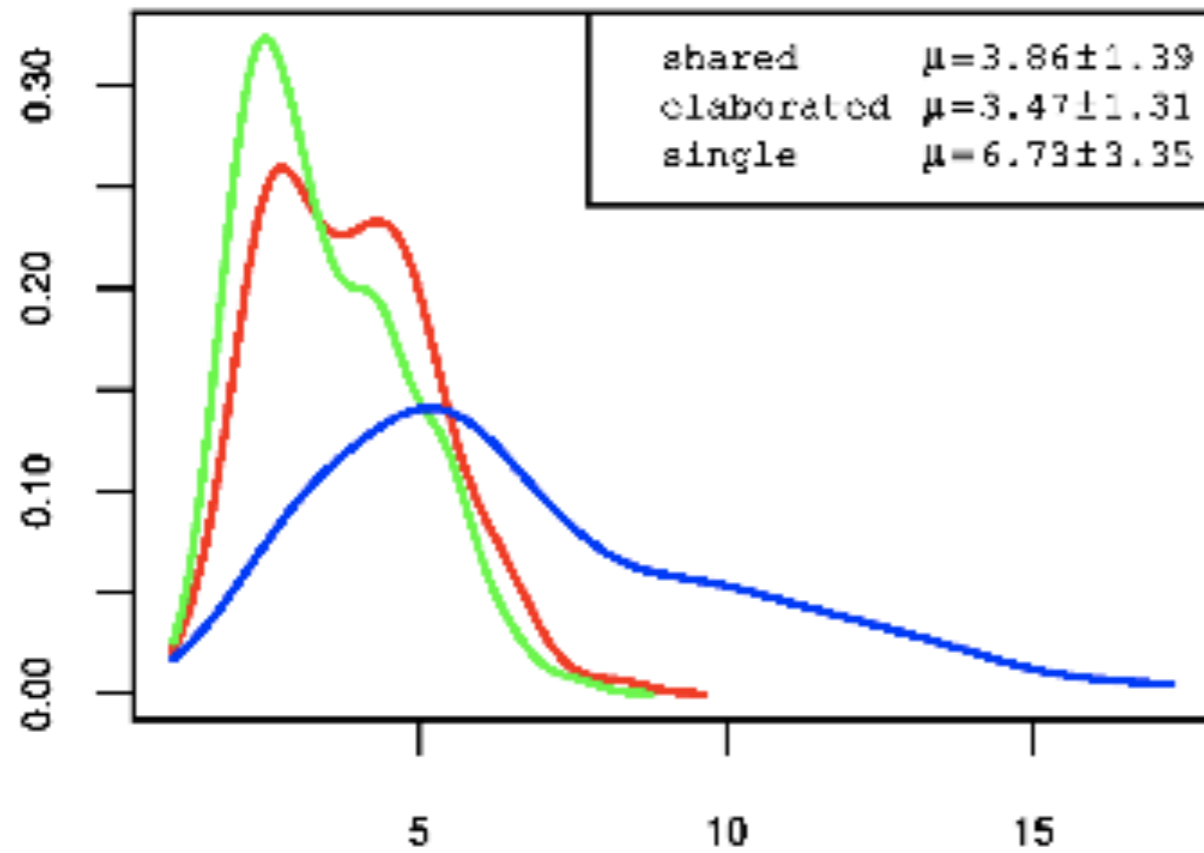
	<i>B. subt.</i>	<i>L. mono.</i>
Shared-Single	0.28	0.14
Elaborated-Single	0.73	1.64E-03

B. anthracis - *L. monocytogenes*

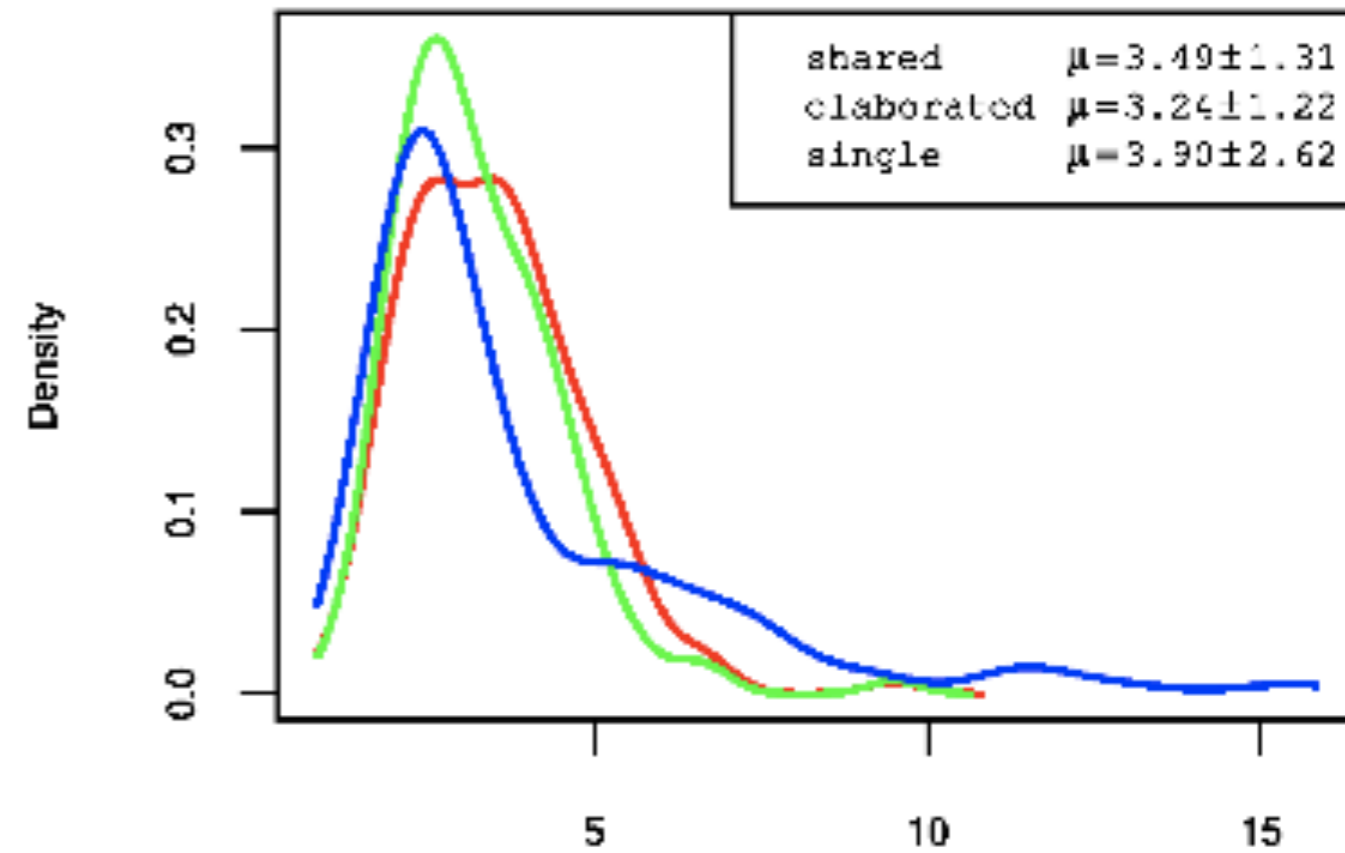
	<i>B. anth</i>	<i>L. mono</i>
Shared-Single	0.01	0.18
Elaborated-Single	0.03	0.03

Sequence p-values (-log10)

Bacillus subtilis
(bsubt-banth)



Bacillus anthracis
(bsubt-banth)



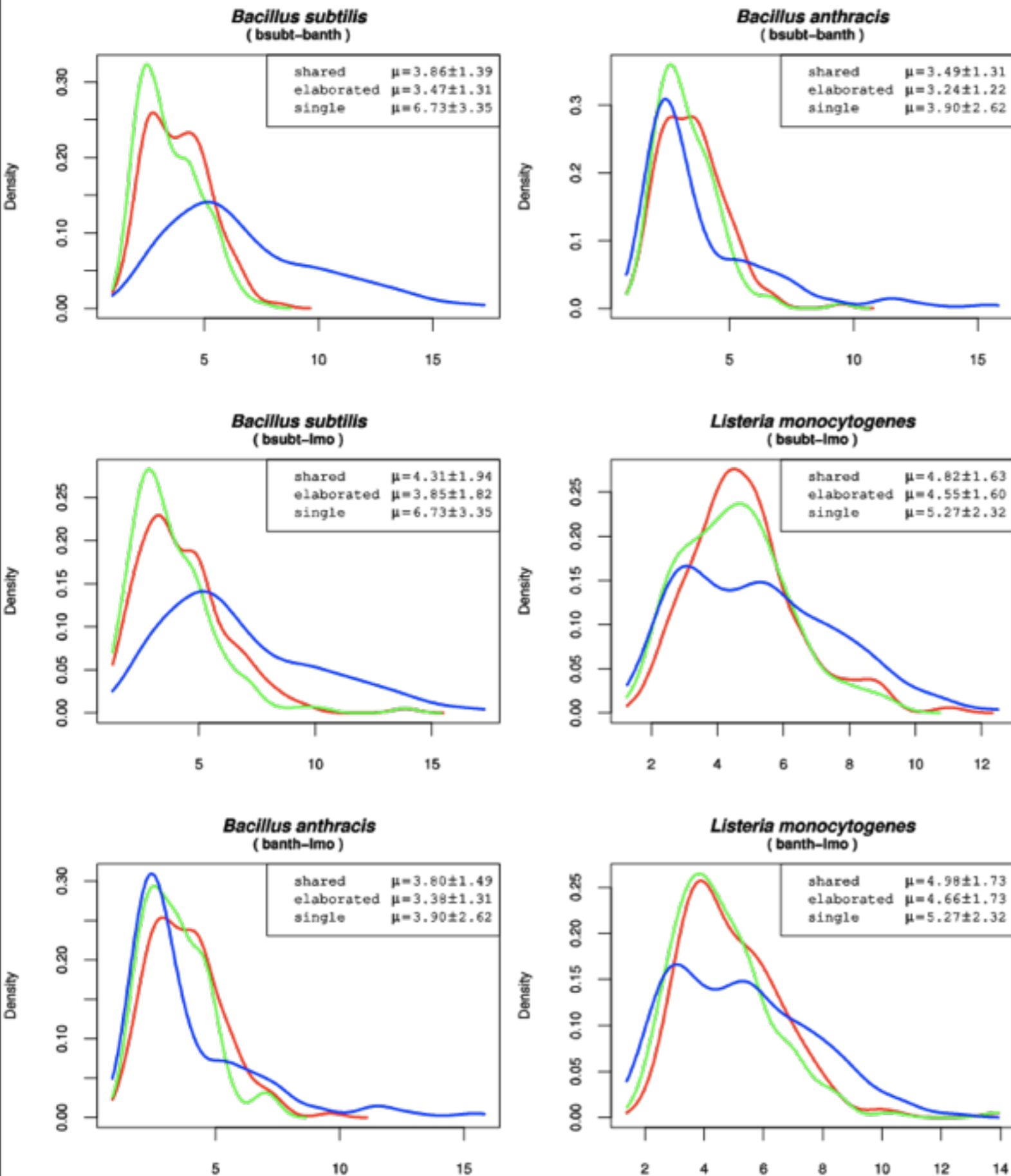
Shared
Elaborated
Single

P-values from two-sided
Wilcoxon's rank test ($\alpha=0.01$):

B. subtilis - *B. anthracis*

	<i>B. subt.</i>	<i>B. anth.</i>
Shared-Elaborated	0.01	0.04
Shared-Single	1.42E-22	0.85
Elaborated-Single	1.70E-29	0.36

Motif P-values (-log10)



$$p(seq) = \prod_{m \in motifs} p(seq | mod el_m)$$

$$p(seq | mod el) = \prod_{i=1}^L q_i$$

P-values from two-sided Wilcoxon's rank test ($\alpha=0.01$)

B. subtilis - *B. anthracis*

	<i>B. sub</i>	<i>B. anth</i>
Shared-elab	0.01	0.04
Shared-single	1.42E-22	0.85
Elaborated-single	1.70E-29	0.36

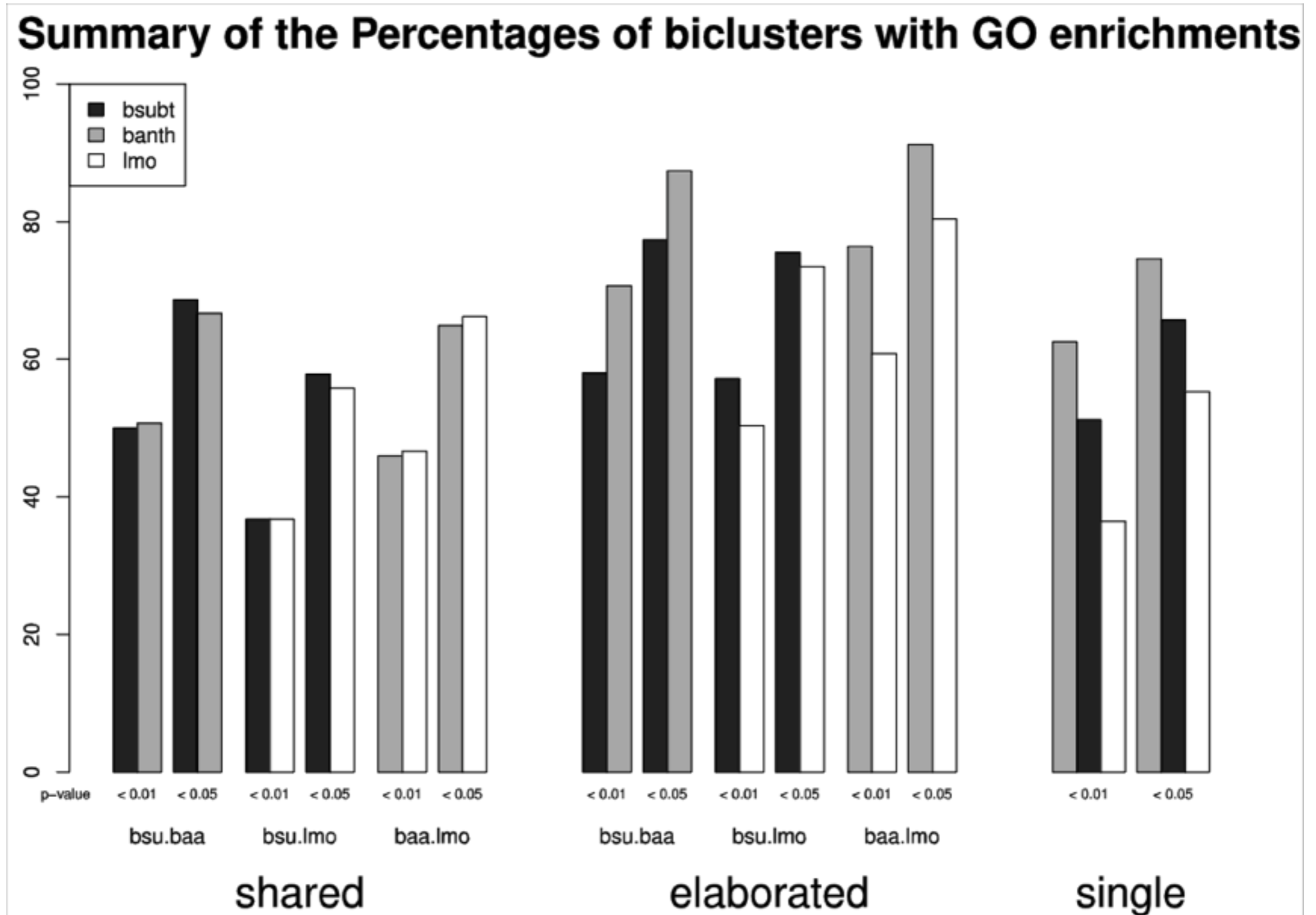
B. subtilis - *L. monocytogenes*

	<i>B. sub</i>	<i>L. mono</i>
Shared-elab	0.01	0.1
Shared-single	3.37E-15	0.07
Elaborated-single	7.36E-23	9.23E-03

B. anthracis - *L. monocytogenes*

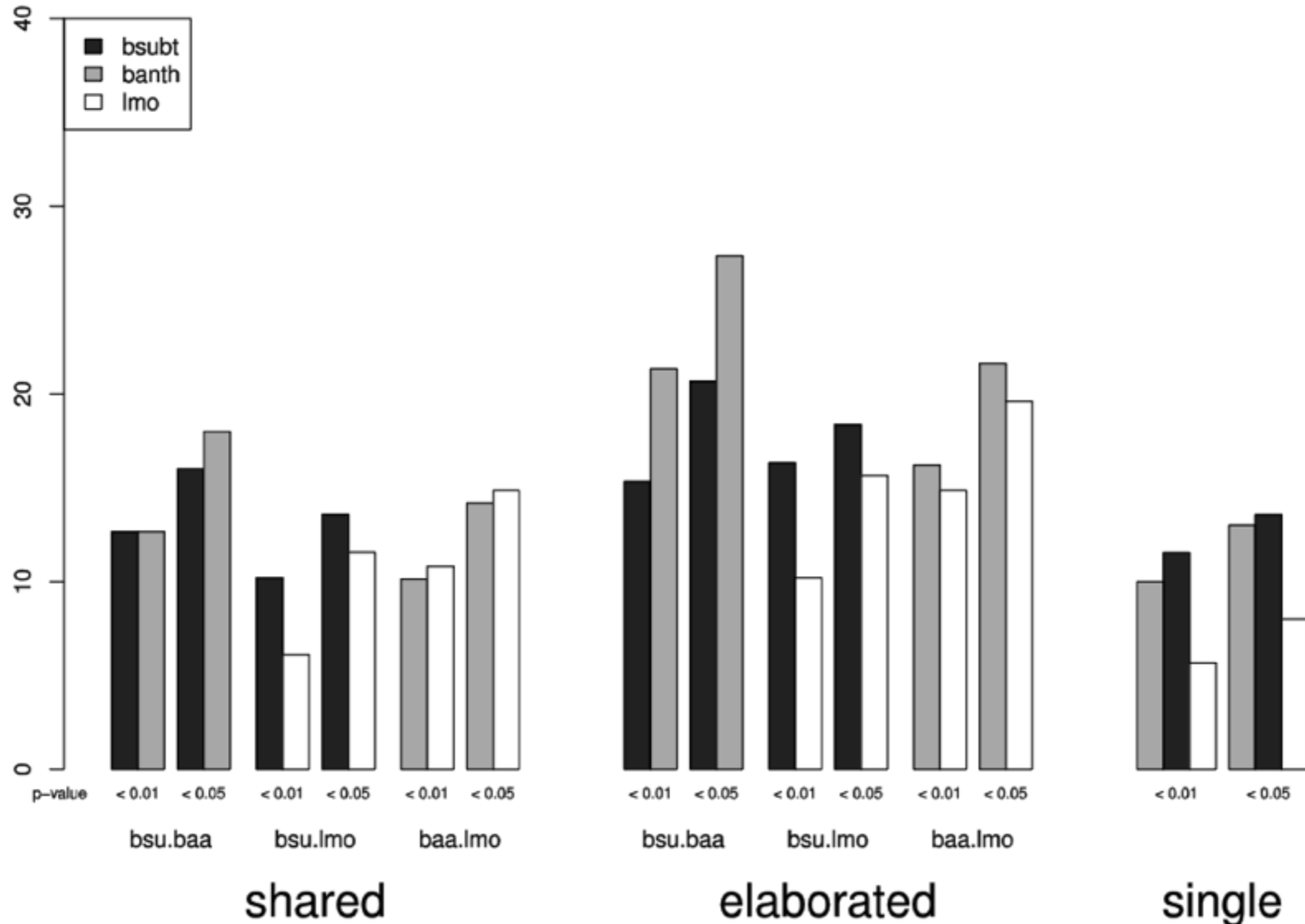
	<i>B. anth</i>	<i>L. mono</i>
Shared-elab	0.01	0.03
Shared-single	0.99	0.24
Elaborated-single	0.55	0.01

Multi-species retrieves more biologically meaningful results



Multi-species retrieves more biologically meaningful results

Summary of the Percentages of biclusters with KEGG enrichment



Conclusions

- Multi-species *cMonkey* improves bicluster quality over conserved modules
 - Expression (residuals)
 - Networks (association p-values)
 - Motifs are area of potential improvement
- Retrieves more
 - biologically significant results (GO/KEGG)
 - conserved modules

Proposed Multi-species cMonkey model

- Explore the optimization further:
 - Alternate objective functions for OC_{12} optimization:

- Bi-variate model:
$$\pi_{ik} \equiv p(y_{ik}^1, y_{ik}^2 | g_{ik}^1, g_{ik}^2) \propto e^{(\beta_0 + \beta_1 g_{ik}^1 + \beta_2 g_{ik}^2)}$$

- Co-reference model:
$$\pi_{ik} = \pi_{ik}^1 \pi_{ik}^2$$

If, we let:

$$\pi_{ik}^1 \equiv p(y_{ik}^1 | X_k^1, S_i^1, M_k^1, N_{i \otimes k}^1) \propto e^{(\beta_0^1 + \beta_1^1 g_{ik}^1)}$$

$$\pi_{ik}^2 \equiv p(y_{ik}^2 | X_k^2, S_i^2, M_k^2, N_{i \otimes k}^2) \propto e^{(\beta_0^2 + \beta_1^2 g_{ik}^2)}$$

- Application to different species/data sets
 - Cancer: human-mouse, cancer-normal
 - Additional triplets, i.e. (*E. coli*, *Salmonella*, *Vibrio* (already have preliminary results))

Acknowledgments

Bonneau lab:

Glenn Butterfoss
Kevin Drew
Aviv Madar
Peter Waltman
Thadeous Kacmarczyk
Shailla Musharof
Devorah Kengmana
Chris Poultny (Shasha)
Irina Nudelman
Alex Pearlman (Ostrer)
Alex Pine

NYU:

Eric Vanden-Eijnden
Harry Ostrer
Mike Purugganan
Patrick Eichenberger
Dennis Shasha

Tacitus-

Howard Coale

• **IBM**

- Robin Wilner
- Bill Boverman
- Viktors Berstis
- Rick Alther
- ETH Zurich
 - Reudi Aebersold
 - Lars Malmstroem

Mike Boxem

Marc Vidal

Dave Goodlett

Jochen Supper (Zell Lab)

- ISB

- Nitin Baliga (&lab)
- Leroy Hood
- Marc Facciotti
- David Reiss
- Vesteynn Thorsson
- Paul Shannon
- Iliana Avila-Campillo (MERC)
- Alan Aderem

Rosetta Commons

Charlie Strauss (Los Alamos)
David Baker (UW Seattle)

**DOD-computing
and society,**

NSF ABI,

NSF Plant genome

NSF DBI,

DOE GTL