

Toolchains for integrating heterogeneous algorithms and data: an anecdote

Aaron Brooks



Nitin Baliga



Dave Reiss



Antoine Allard

GENOMES to LIFE

BIOLOGICAL SOLUTIONS FOR ENERGY CHALLENGES

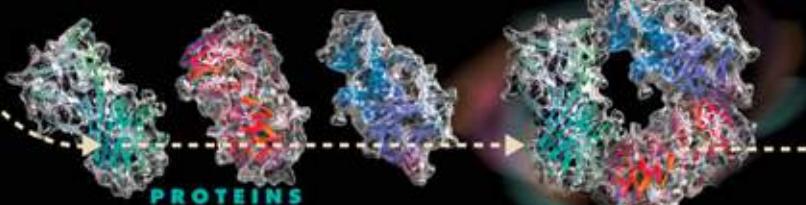
INNOVATIVE APPROACHES ALONG UNCONVENTIONAL PATHS

U.S. DEPARTMENT OF ENERGY

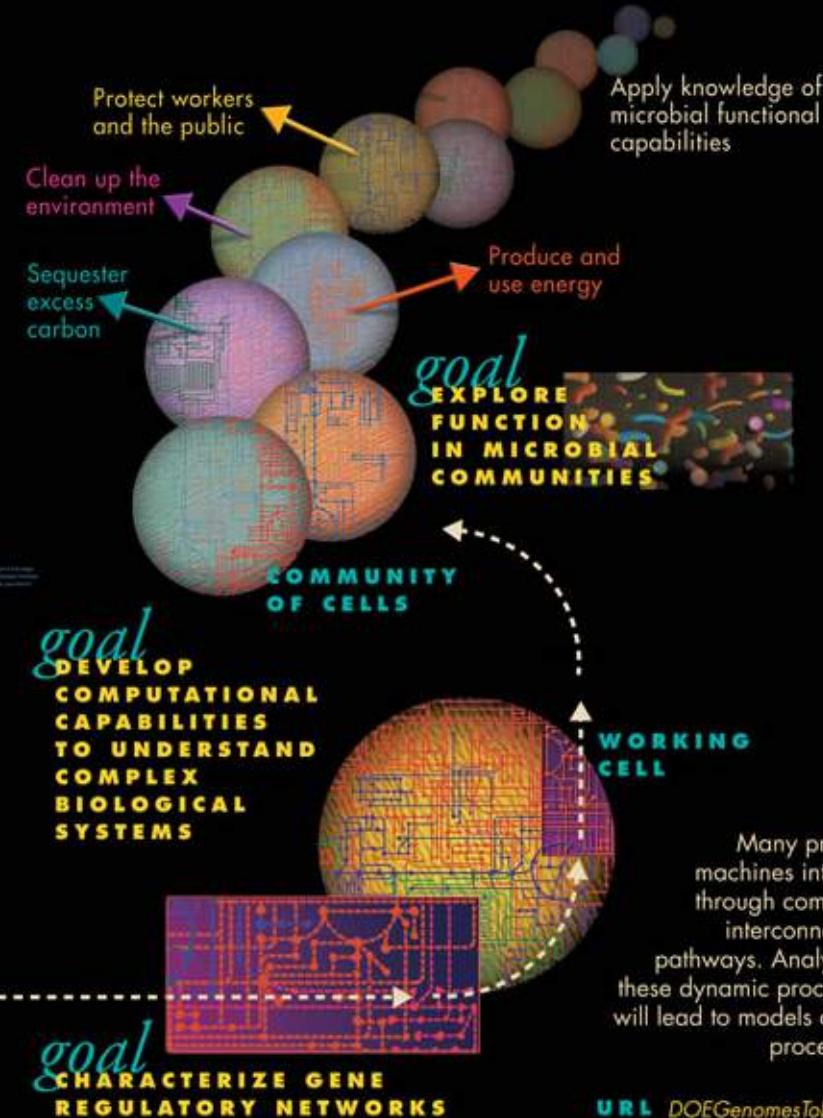


DNA SEQUENCE DATA FROM GENOME PROJECTS

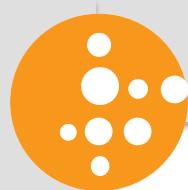
Genes and other DNA sequences contain instructions on how and when to build proteins



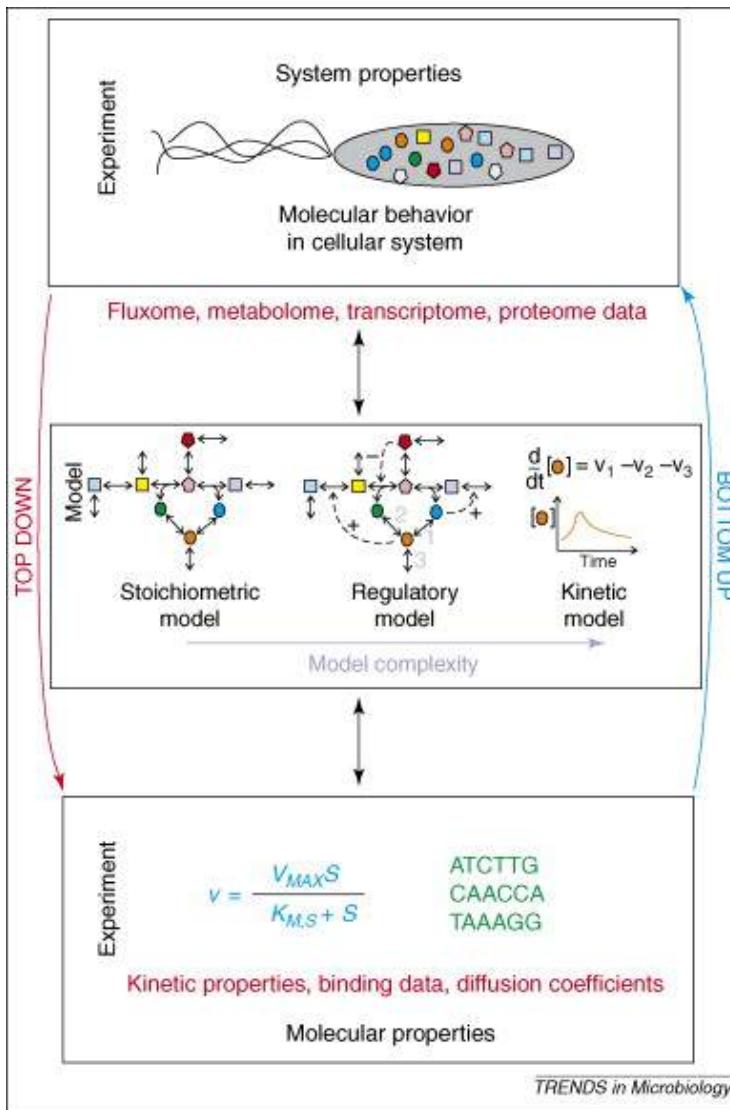
Proteins perform many of life's most essential functions. To carry out their specific roles, they often work together in the cell as protein machines.



URL DOEGenomesToLife.org



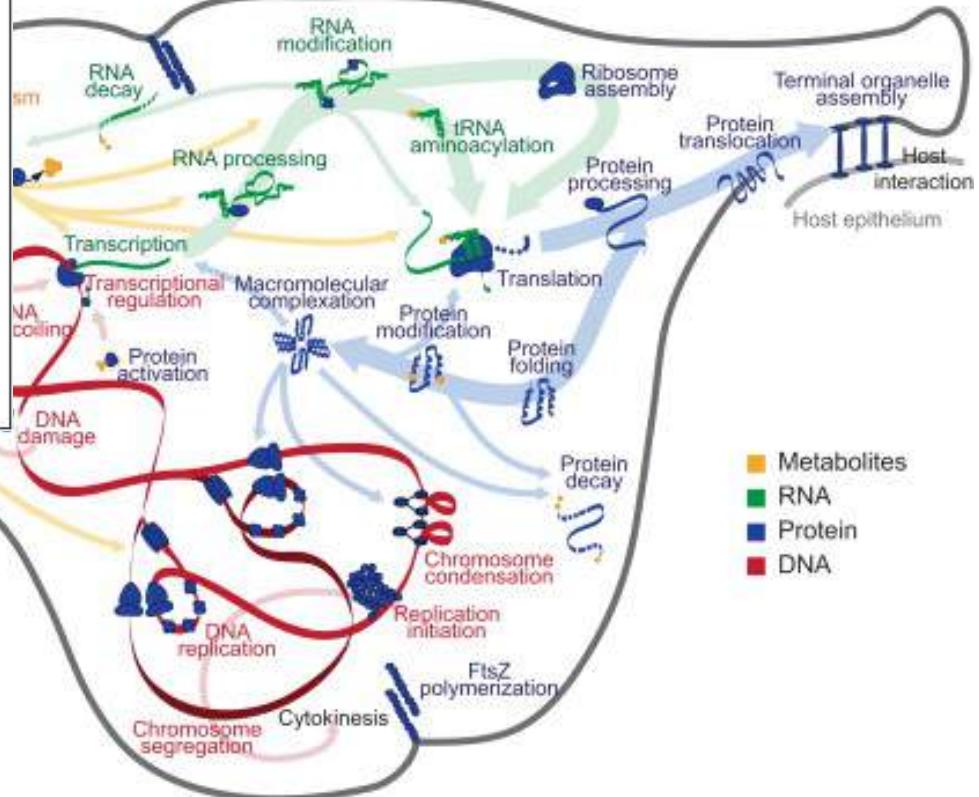
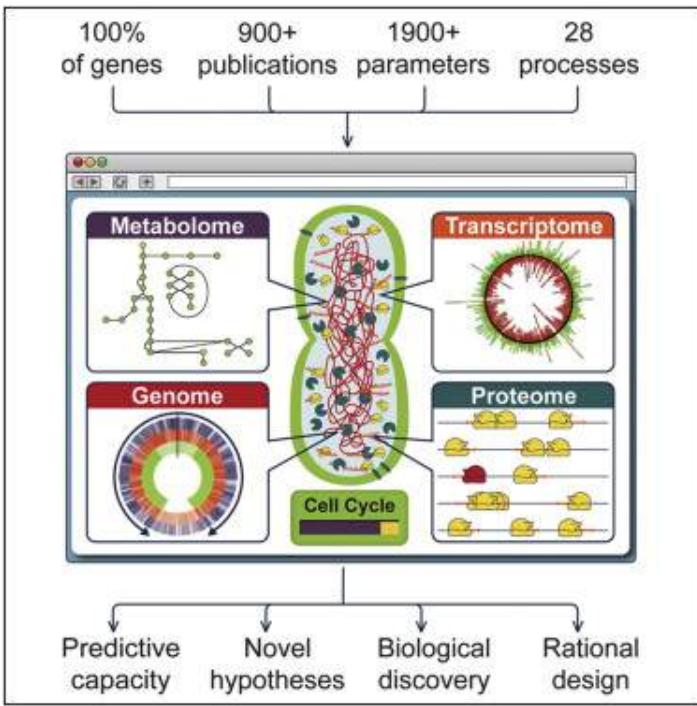
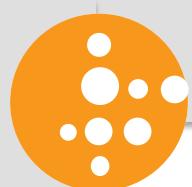
Two conceptual approaches in systems biology



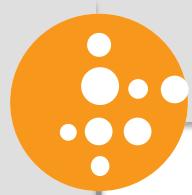
Top down

Bottom up

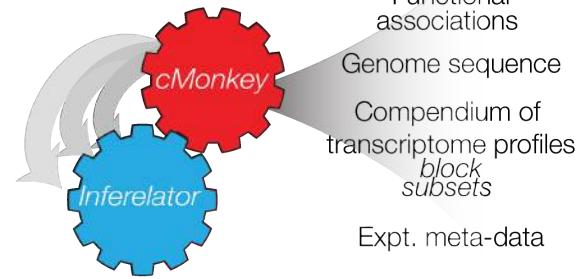
Bottom up



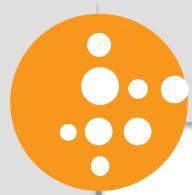
Karr et al 2012



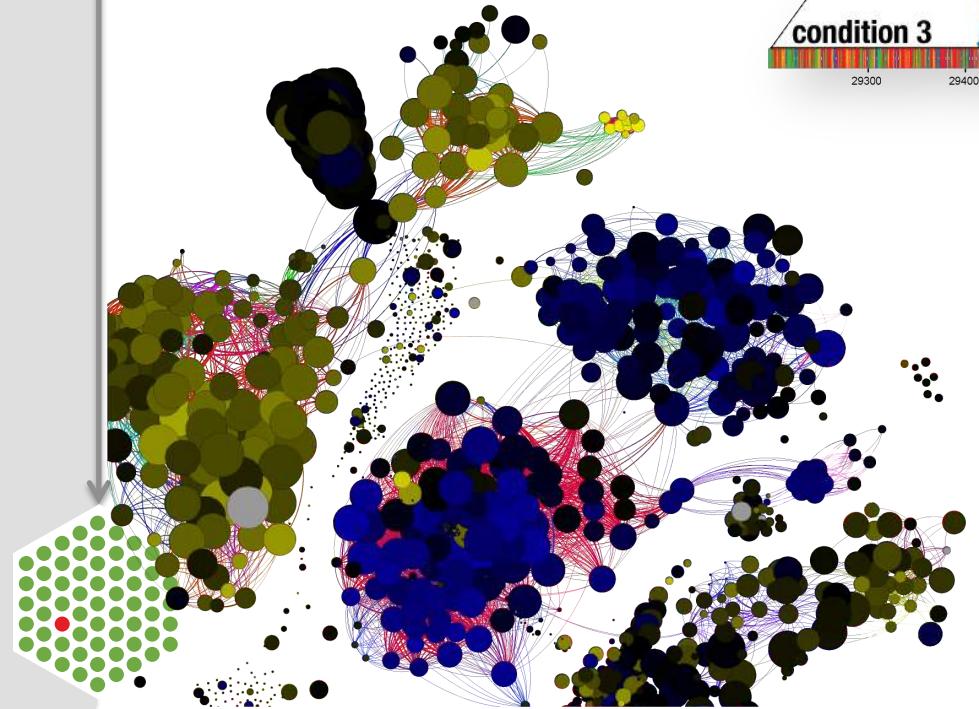
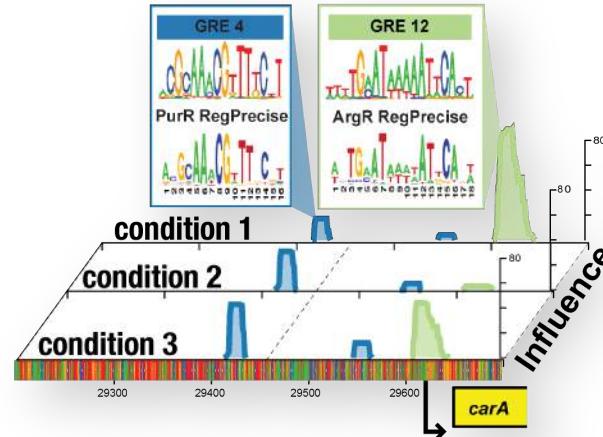
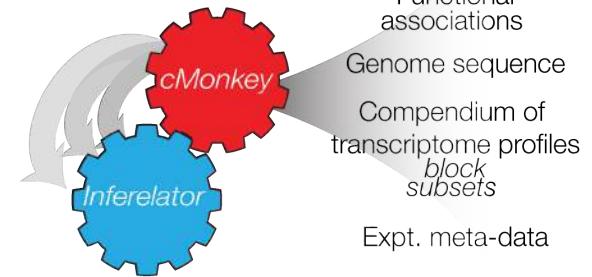
Top down



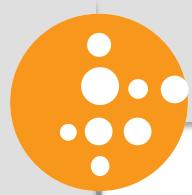
EGRIN2.0
ensemble network inference



Top down

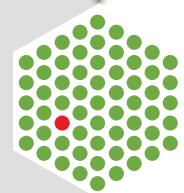


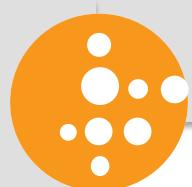
EGRIN2.0
ensemble network inference



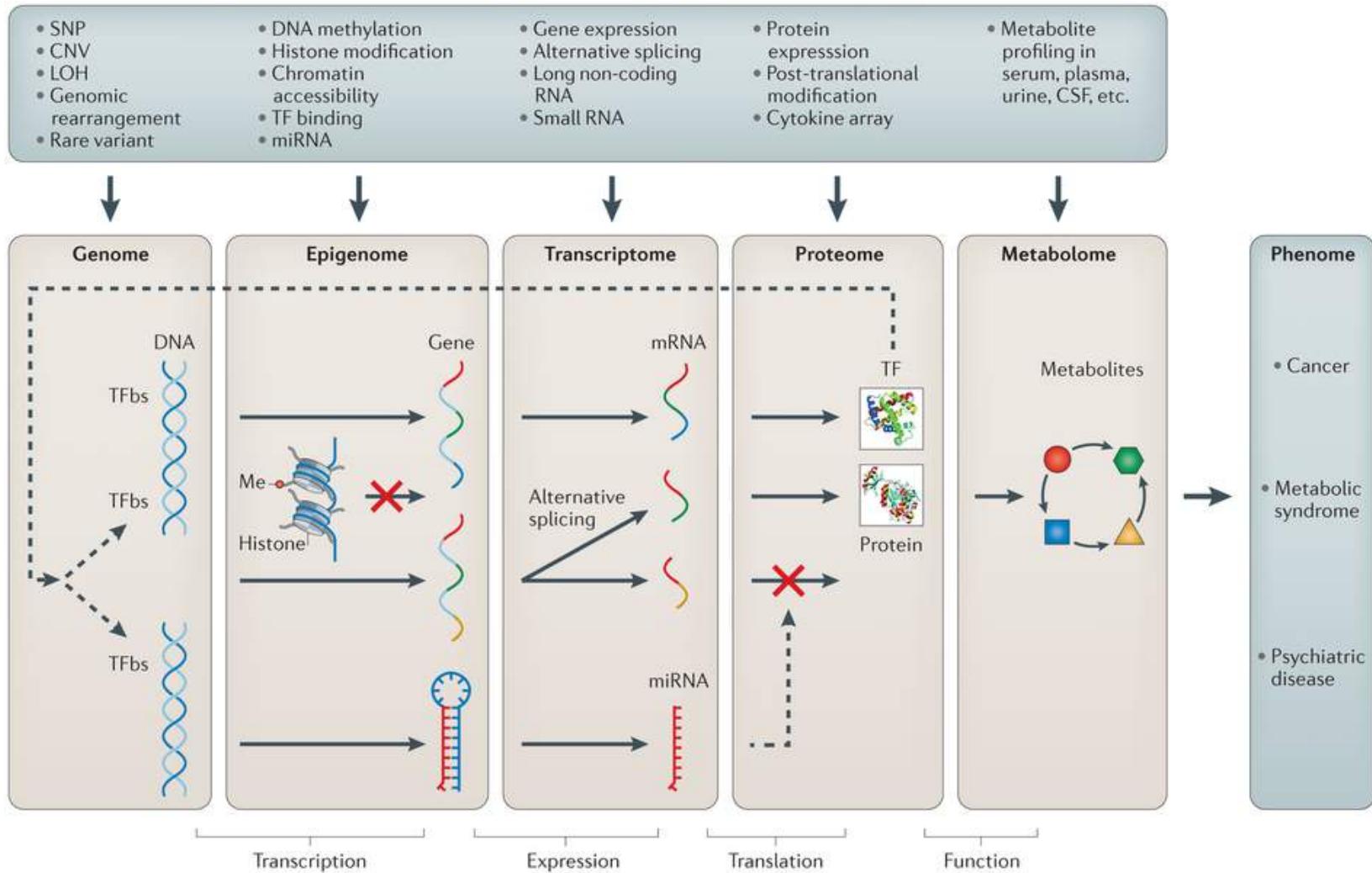
Modeling challenges in systems biology

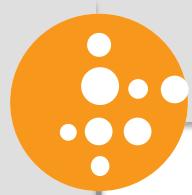
Data is heterogeneous





Challenge: Heterogeneous data

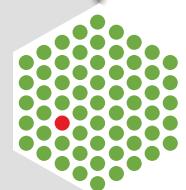


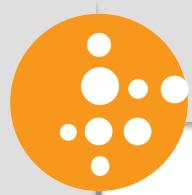


Modeling challenges in systems biology

Data is heterogeneous

Algorithms are numerous and diverse

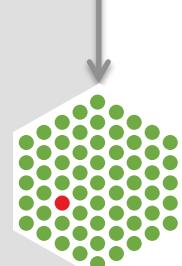




Challenge: Many Algorithms

Table 1 | Network inference methods

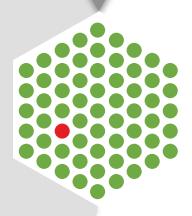
ID Synopsis	Reference
Regression: transcription factors are selected by target gene-specific (i) sparse linear-regression and (ii) data-resampling approaches.	
1 Trustful Inference of Gene REgulation using Stability Selection (TIGRESS): (i) Lasso; (ii) the regularization parameter selects five transcription factors per target gene in each bootstrap sample.	33 ^a
2 (i) Steady-state and time-series data are combined by group Lasso; (ii) bootstrapping.	34 ^a
3 Combination of Lasso and Bayesian linear regression models learned using reversible-jump Markov chain Monte Carlo simulations.	35 ^a
4 (i) Lasso; (i) bootstrapping.	36
5 (i) Lasso; (ii) area under the stability selection curve.	36
6 Application of the Lasso toolbox GENLAB using standard parameters.	37
7 Lasso models are combined by the maximum regularization parameter selecting a given edge for the first time.	36 ^a
8 Linear regression determines the contribution of transcription factors to the expression of target genes.	— ^{a,b}
Mutual information: edges are (i) ranked based on variants of mutual information and (ii) filtered for causal relationships.	
1 Context likelihood of relatedness (CLR): (i) spline estimation of mutual information; (ii) the likelihood of each mutual information score is computed based on its local network context.	11 ^{a,b}
2 (i) Mutual information is computed from discretized expression values.	38 ^{a,b}
3 Algorithm for the reconstruction of accurate cellular networks (ARACNE): (i) kernel estimation of mutual information; (ii) the data processing inequality is used to identify direct interactions.	9 ^{a,b}
4 (i) Fast kernel-based estimation of mutual information; (ii) Bayesian local causal discovery (BLCD) and Markov blanket (HITON-PC) algorithm to identify direct interactions.	39 ^a
5 (i) Mutual information and Pearson's correlation are combined; (ii) BLCD and HITON-PC algorithm.	39 ^a
Correlation: edges are ranked based on variants of correlation.	
1 Absolute value of Pearson's correlation coefficient.	38
2 Signed value of Pearson's correlation coefficient.	38 ^{a,b}
3 Signed value of Spearman's correlation coefficient.	38 ^{a,b}
Bayesian networks: optimize posterior probabilities by different heuristic searches.	
1 Simulated annealing (catnet R package, http://cran.r-project.org/web/packages/catnet/), aggregation of three runs.	—
2 Simulated annealing (catnet R package, hyperlink above).	—
3 Max-min parent and children algorithm (MMPC), bootstrapped data sets.	40
4 Markov blanket algorithm (HITON-PC), bootstrapped data sets.	41
5 Markov boundary induction algorithm (TIE*), bootstrapped data sets.	42
6 Models transcription factor perturbation data and time series using dynamic Bayesian networks (Infer.NET toolbox, http://research.microsoft.com/infernet/).	— ^a
Other approaches: network inference by heterogeneous and novel methods.	
1 GENIE3: a Random Forest is trained to predict target gene expression. Putative transcription factors are selected as tree nodes if they consistently reduce the variance of the target.	19 ^a
2 Codependencies between transcription factors and target genes are detected by the nonlinear correlation coefficient η^2 (two-way ANOVA). Transcription-factor perturbation data are up-weighted.	20 ^a
3 Transcription factors are selected by maximizing the conditional entropy for target genes, which are represented as Boolean vectors with probabilities to avoid discretization.	43 ^a
4 Transcription factors are preselected from transcription-factor perturbation data or by Pearson's correlation and then tested by iterative Bayesian model averaging (BMA).	44
5 A Gaussian noise model is used to estimate whether Marbach et al 2012	45
6 After scaling, target genes are clustered by Pearson's correlation. A neural network is trained (genetic algorithm) and parameterized (back-propagation).	46 ^a

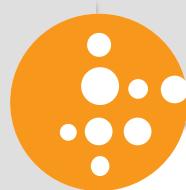


Introduce a generalized strategy for building multistep, data-driven inference toolchains, using:

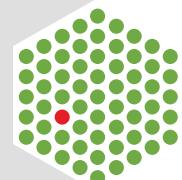
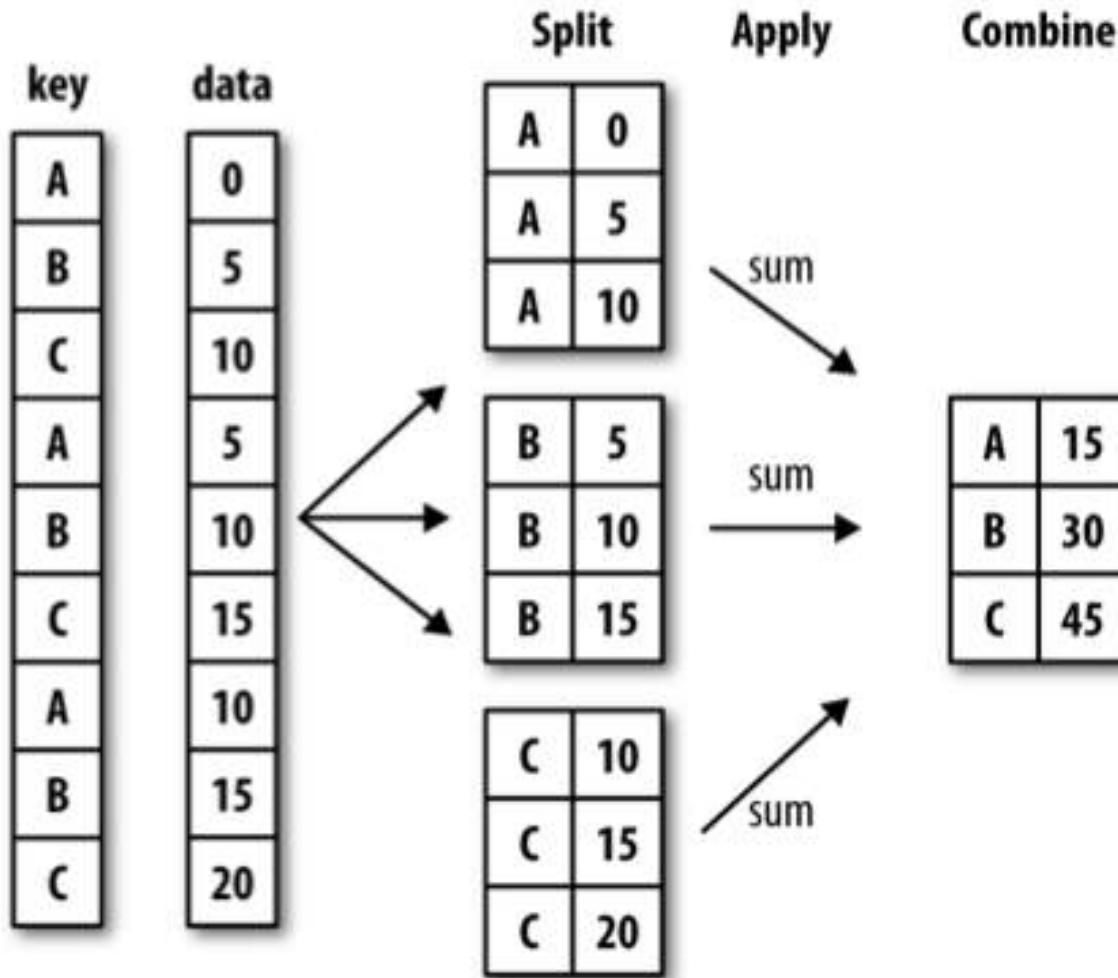
- Large, heterogeneous data sets
- Multiple algorithms combined in an ensemble learning framework
- Graph-based methods to reconcile predictions

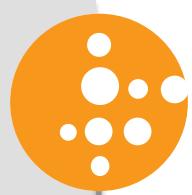
Extension: a suggestion for integrating models across data types and biological processes



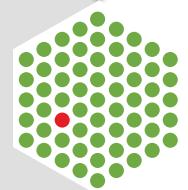
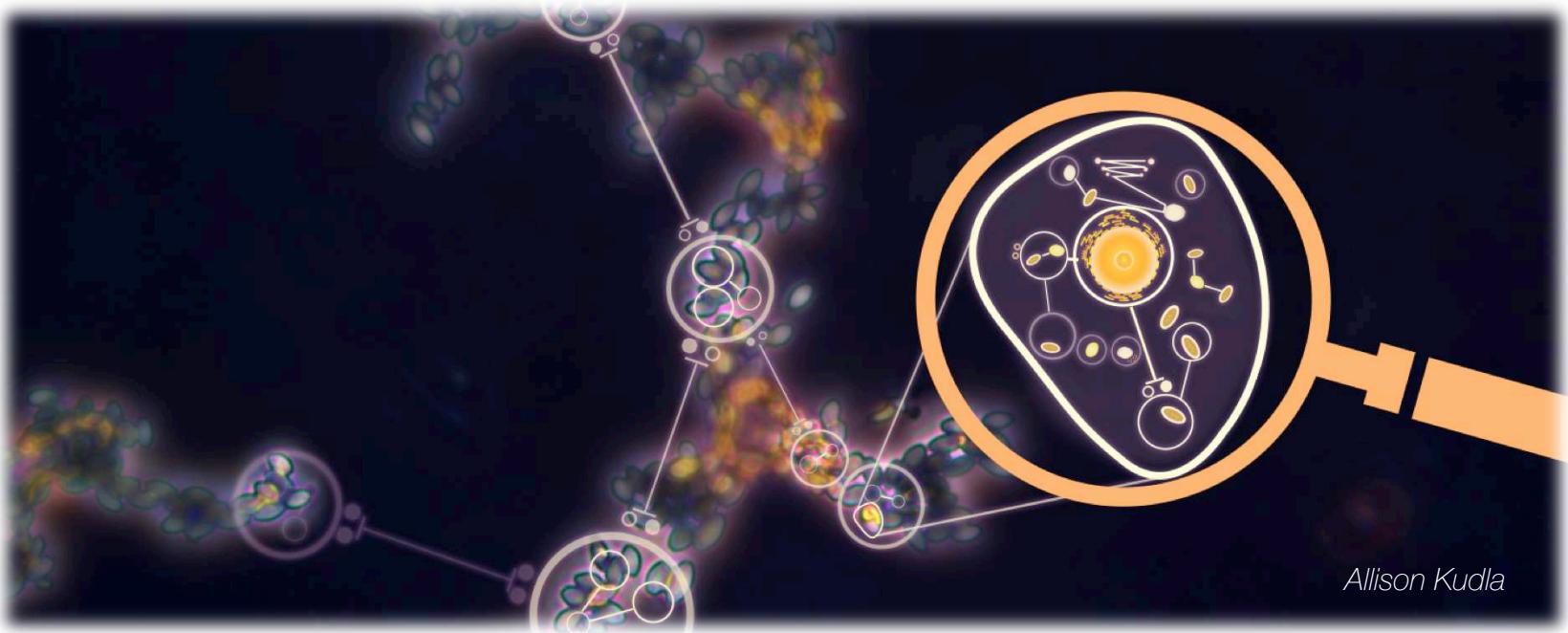


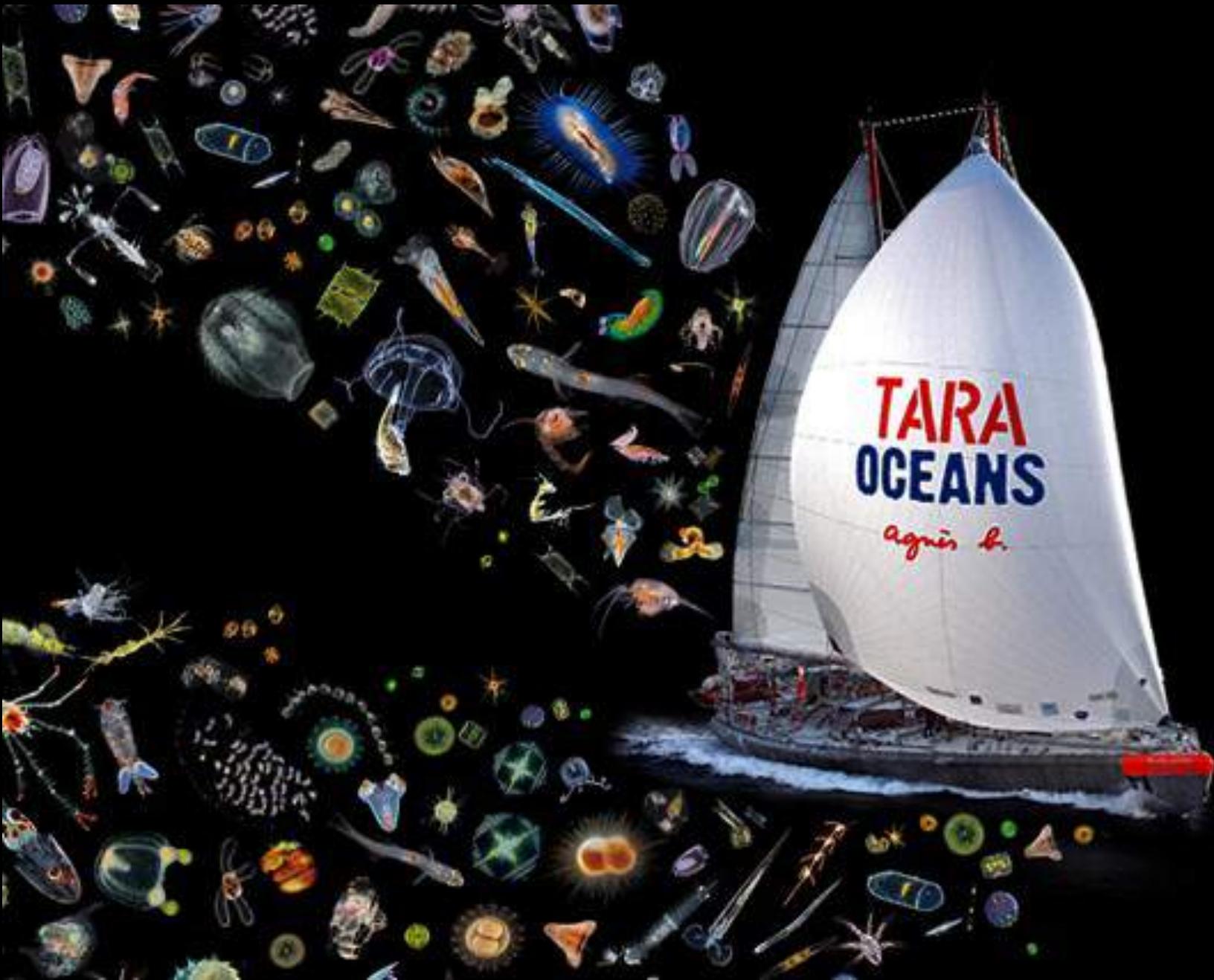
Split-Apply-Combine Approach to Biological Inference

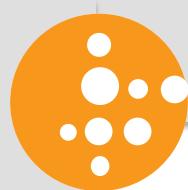




How do microbes control expression of their genomes?

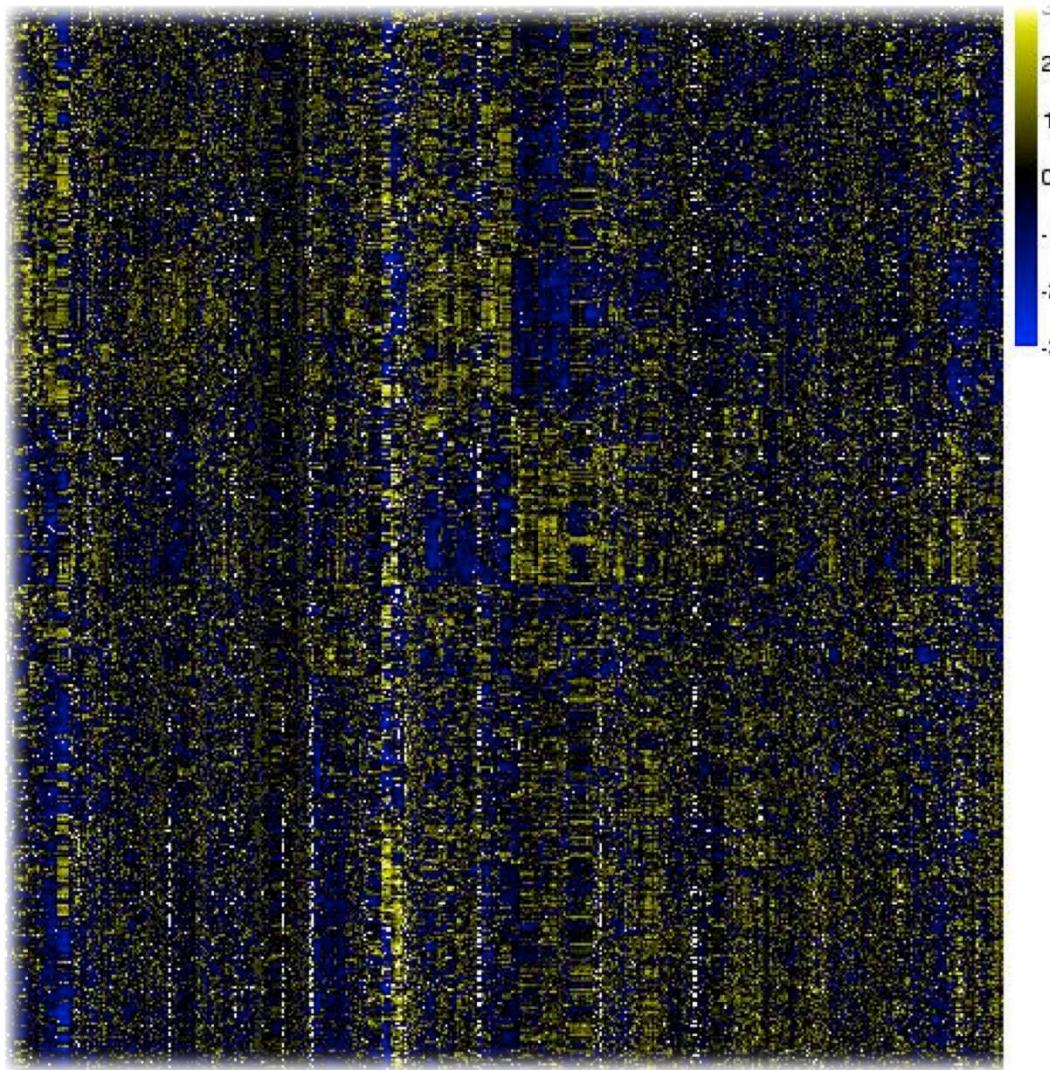




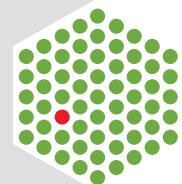


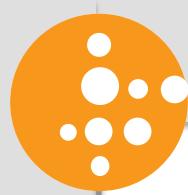
From gene expression to regulatory networks

4,213
genes



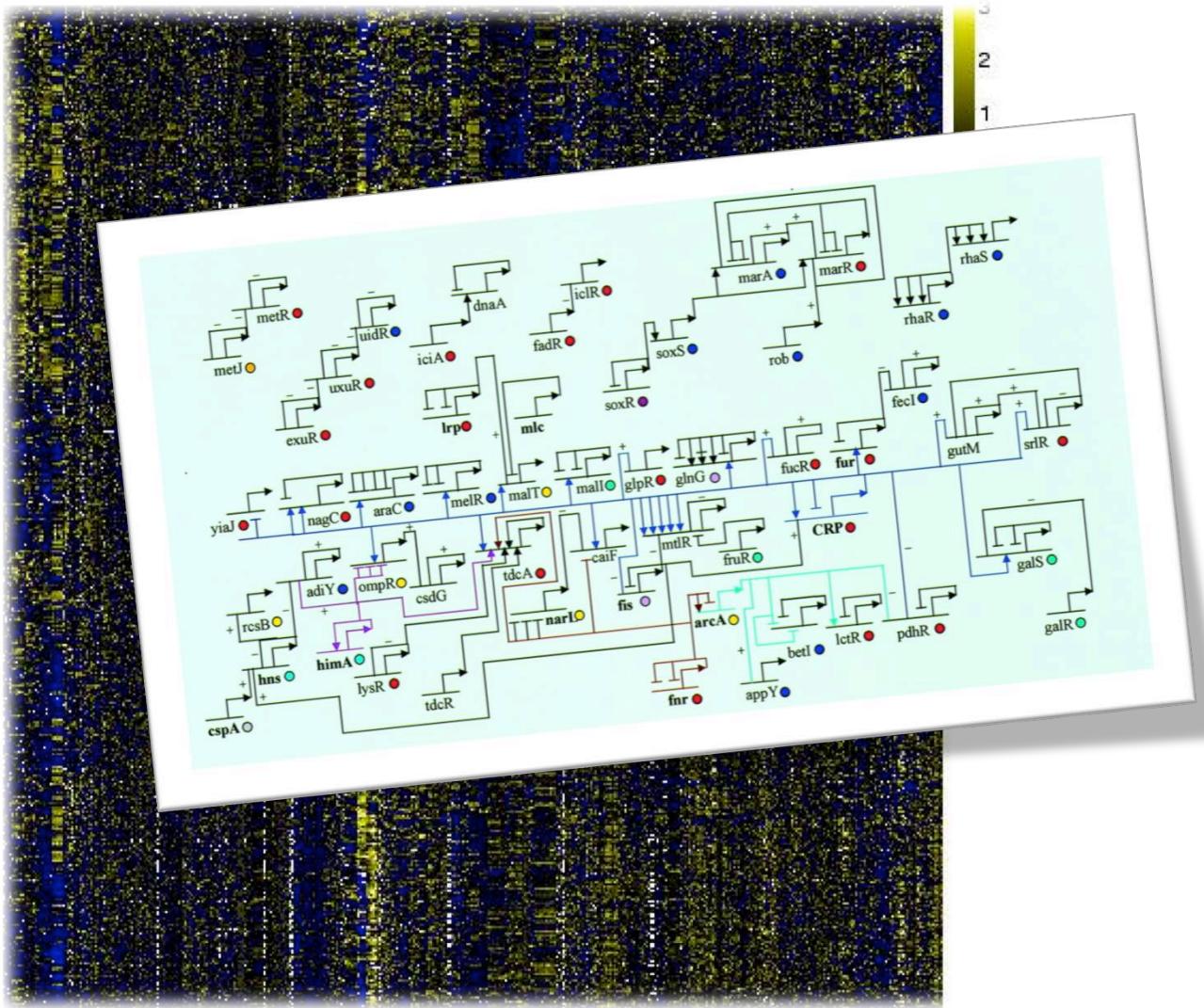
868 *conditions*



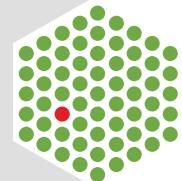


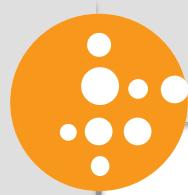
From gene expression to regulatory networks

4,213
genes



868 conditions

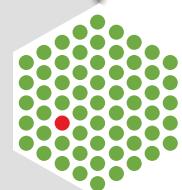


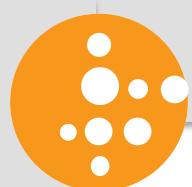


Wisdom of crowds for robust gene network inference

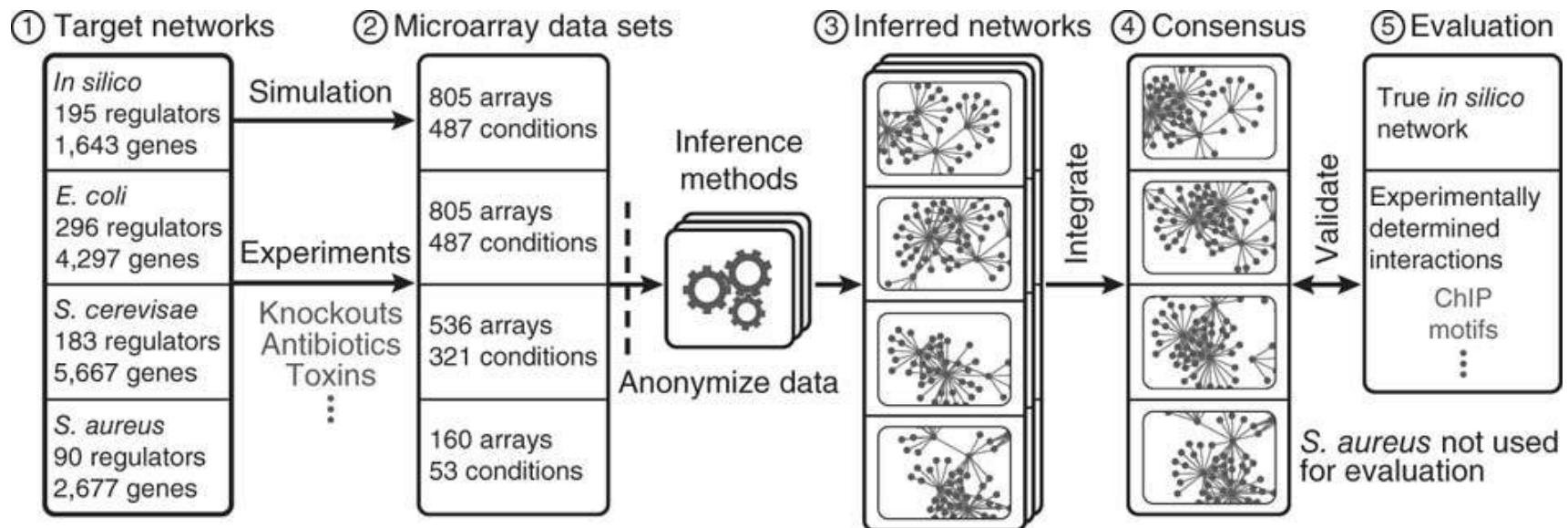


Marbach et al 2012

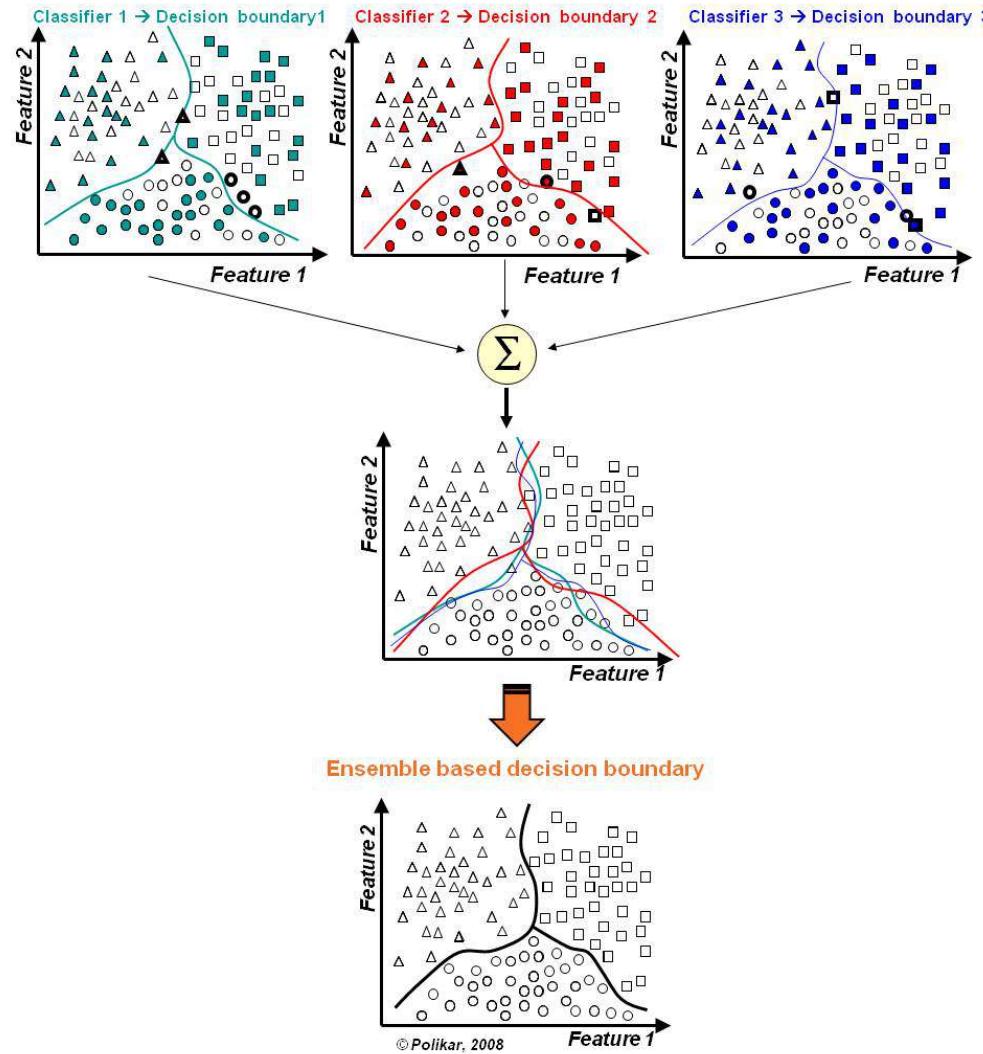


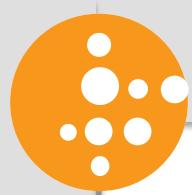


Combine algorithms to increase performance

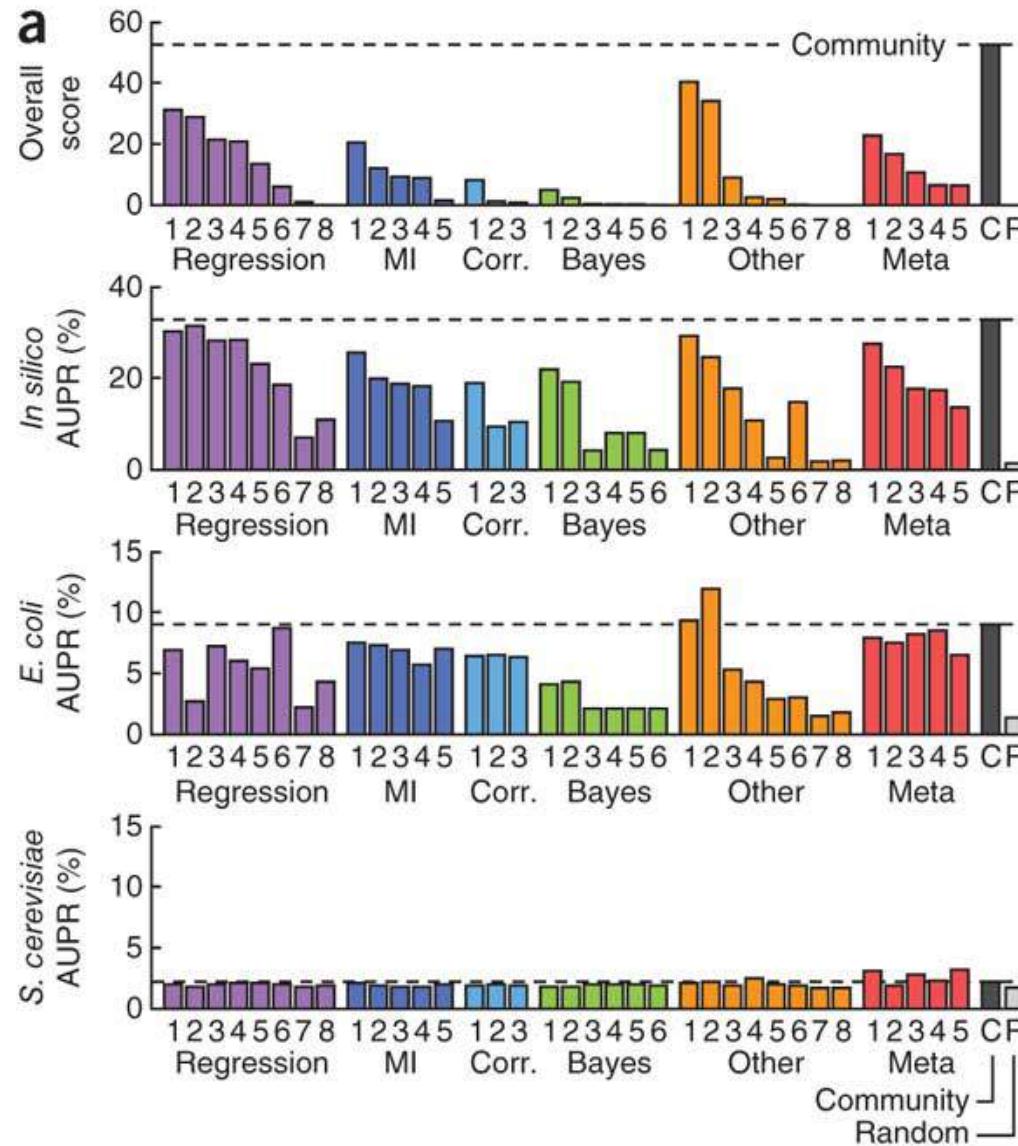


Ensemble learning

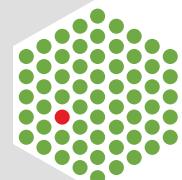


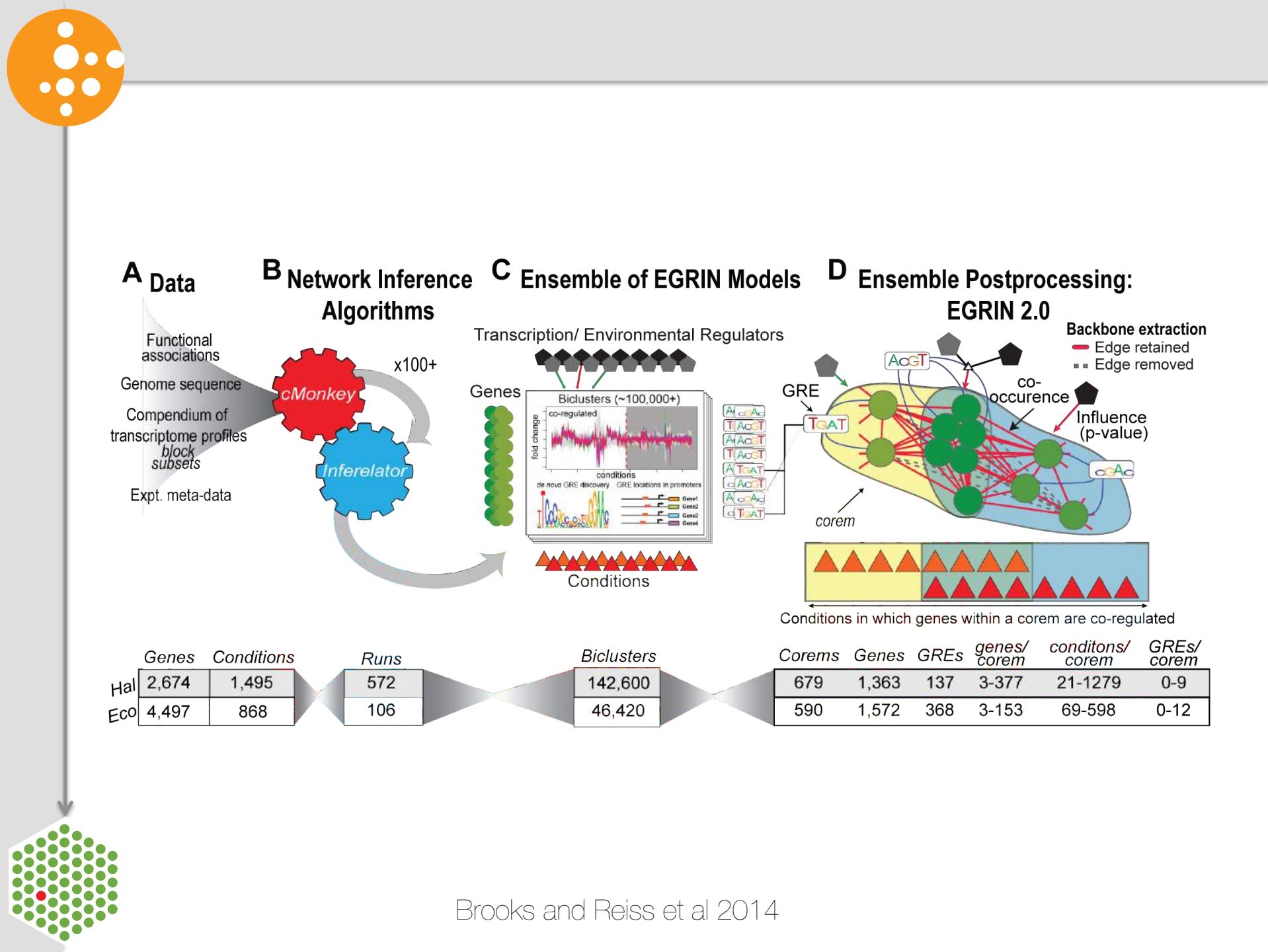


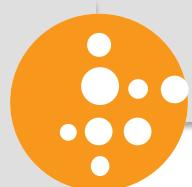
Combining learners increases performance



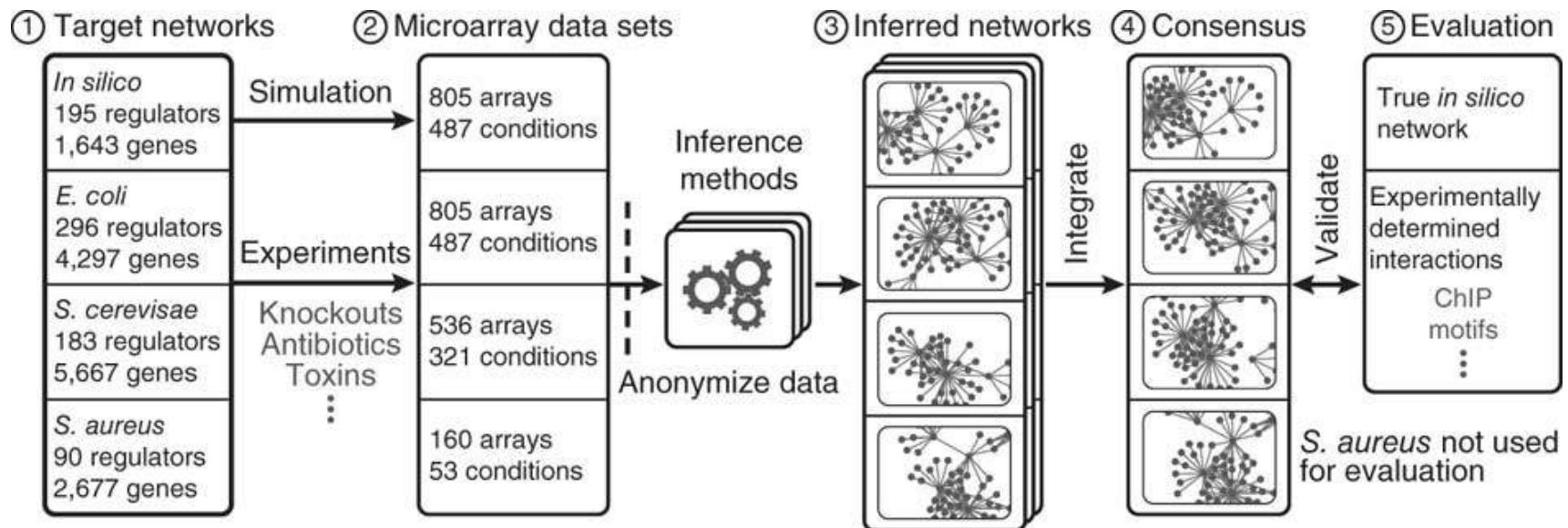
Marbach et al 2012



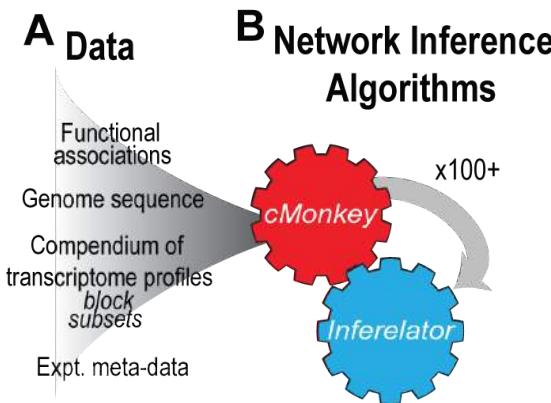




Combine algorithms to increase performance

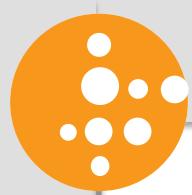


SPLIT: Bagging

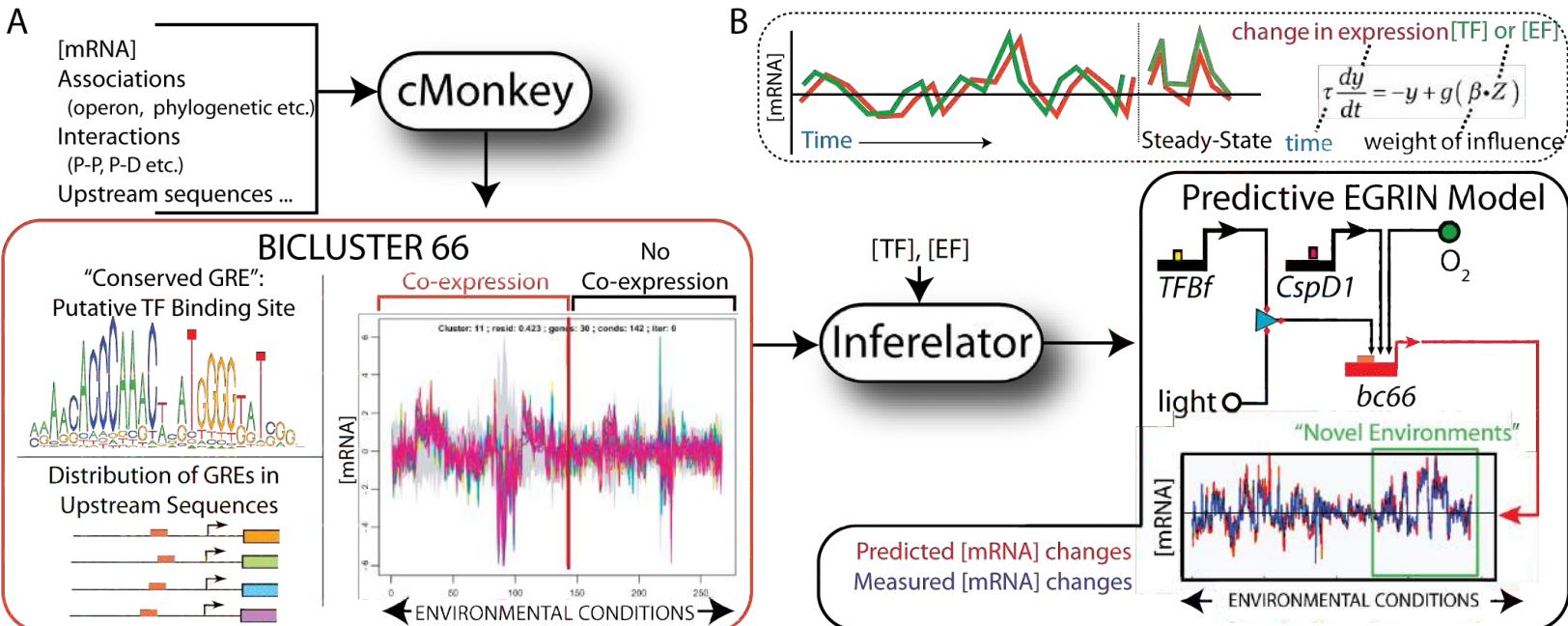


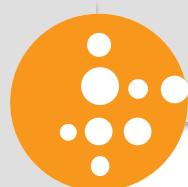
	Genes	Conditions	Runs
Hal	2,674	1,495	572
Eco	4,497	868	106

Brooks and Reiss et al 2014

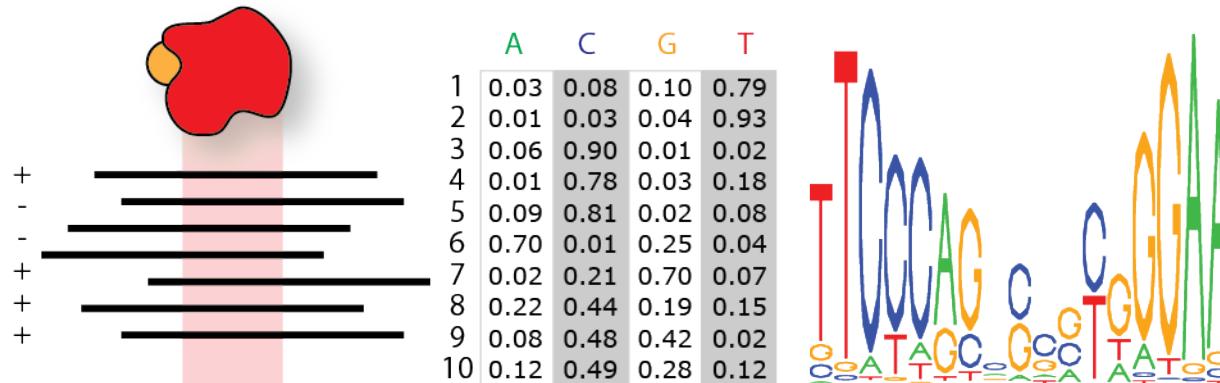
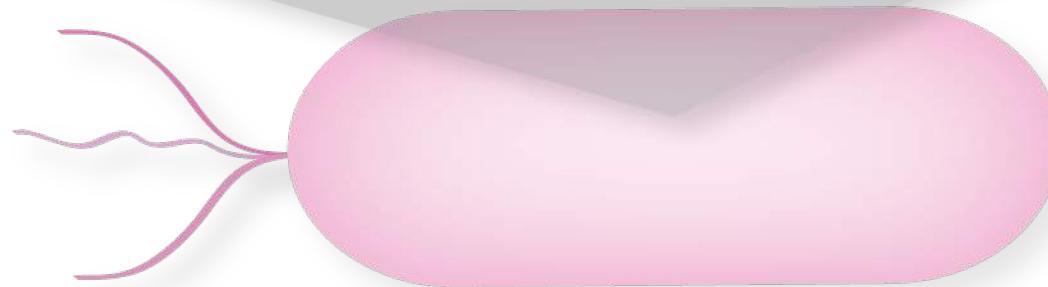
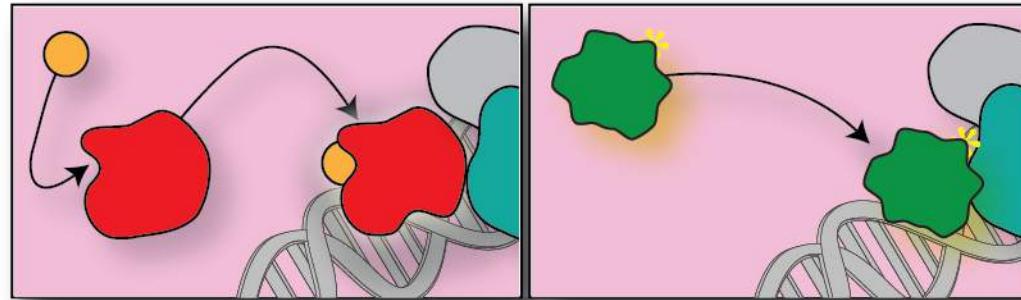


APPLY: cMonkey and Inferelator

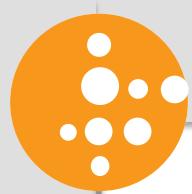




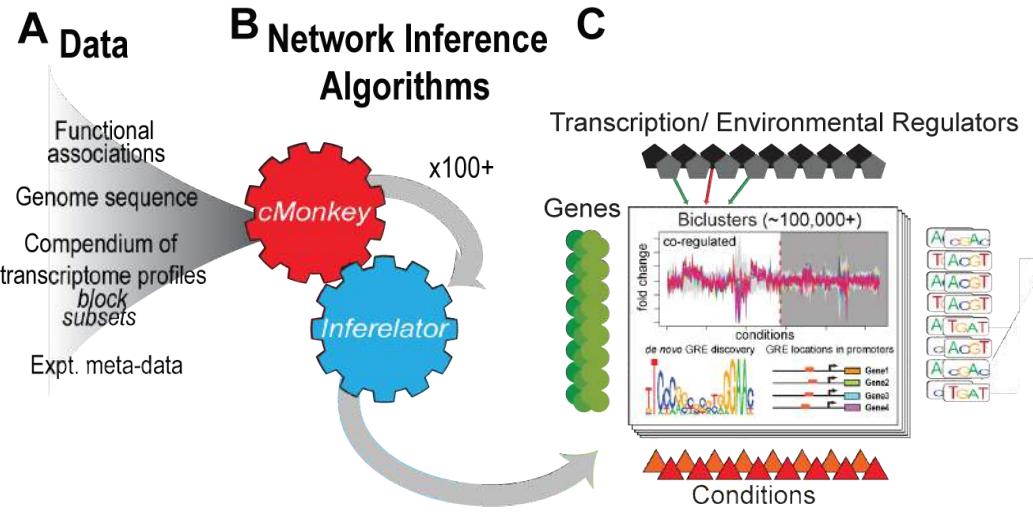
cMonkey clusters and detects motifs de novo



Brooks 2014



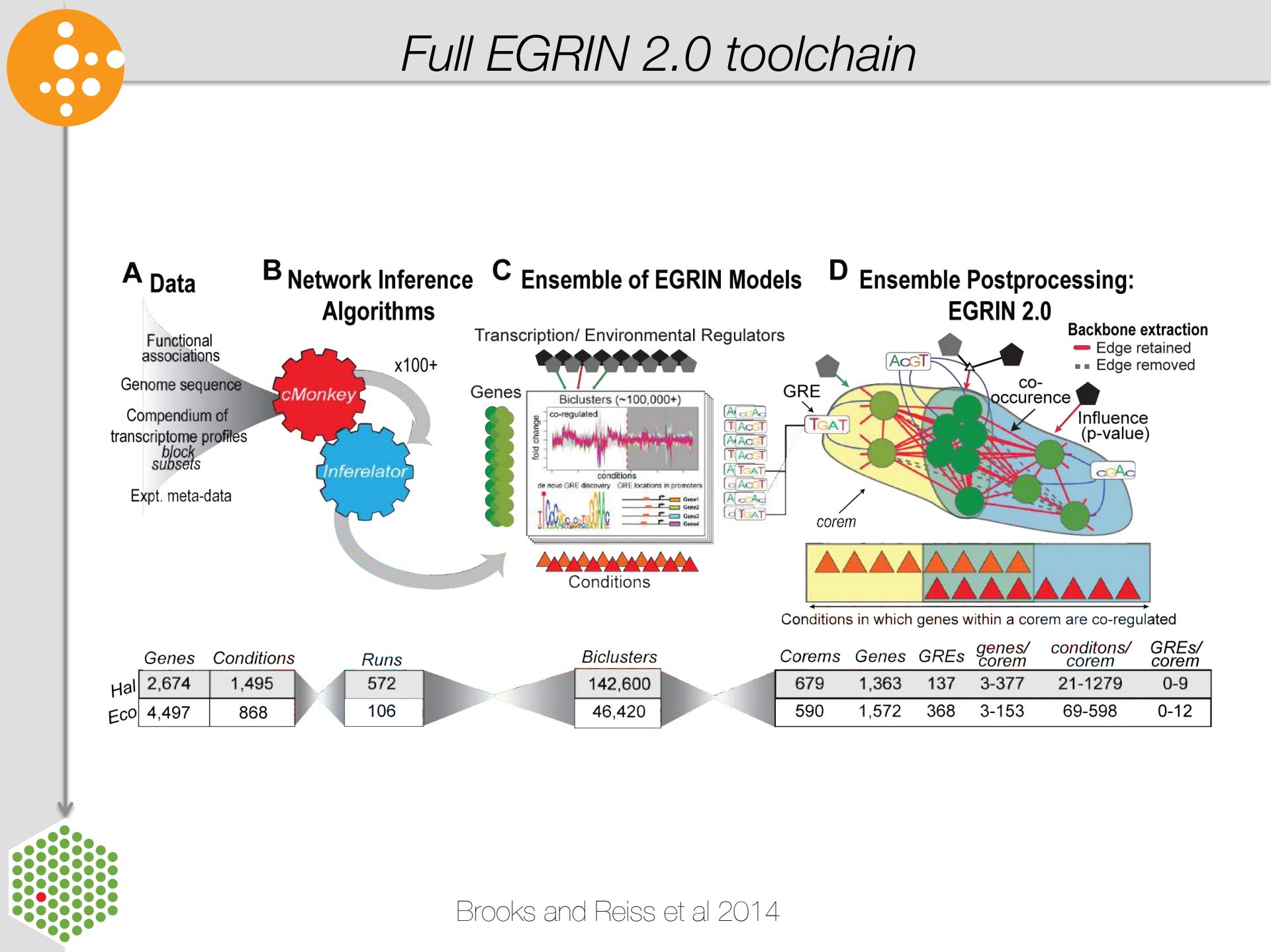
Full EGRIN 2.0 toolchain



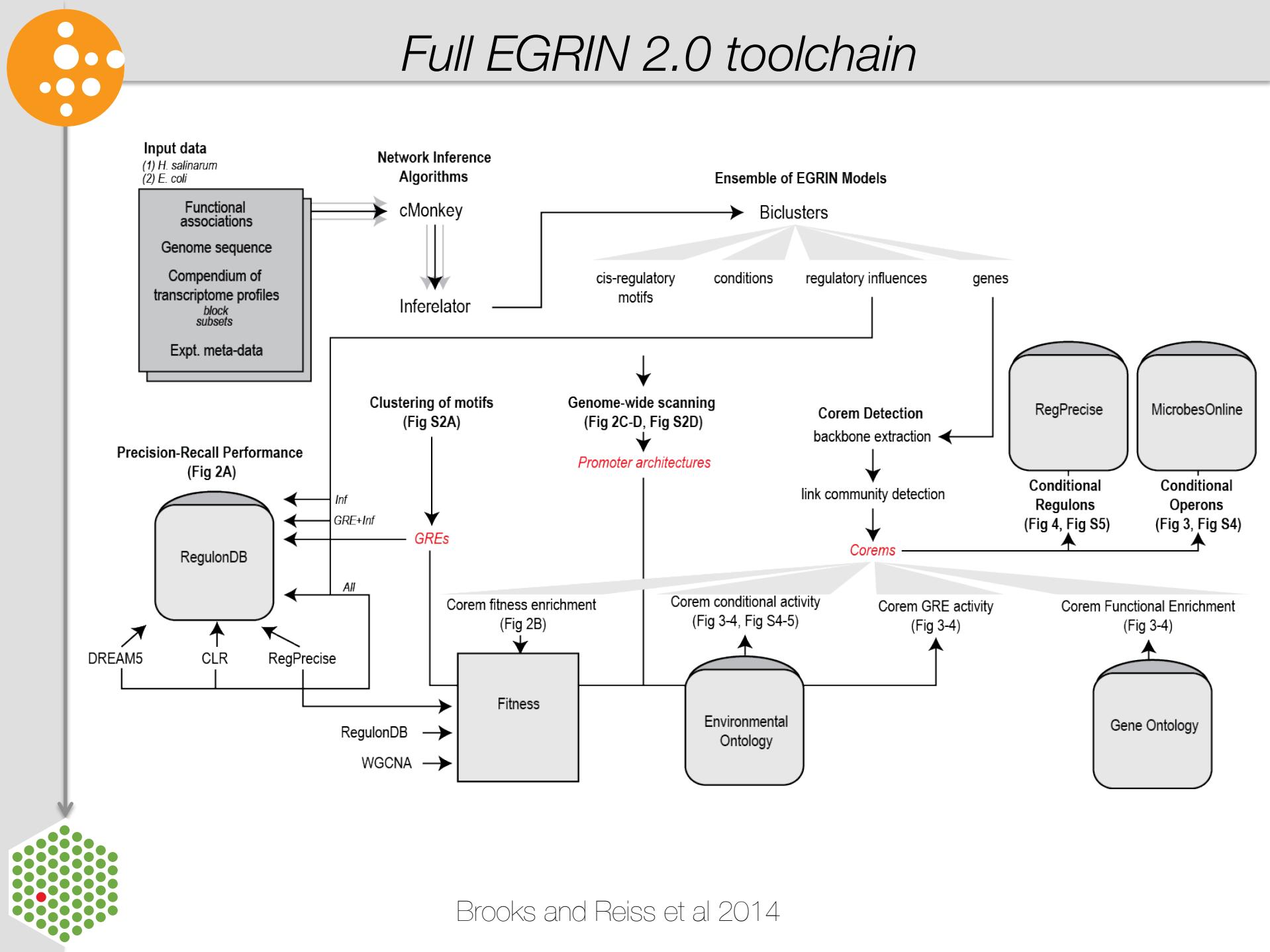
	Genes	Conditions	Runs	Biclusters
Hal	2,674	1,495	572	142,600
Eco	4,497	868	106	46,420

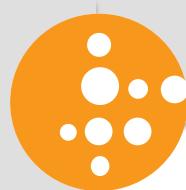
Brooks and Reiss et al 2014

Full EGRIN 2.0 toolchain

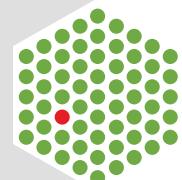
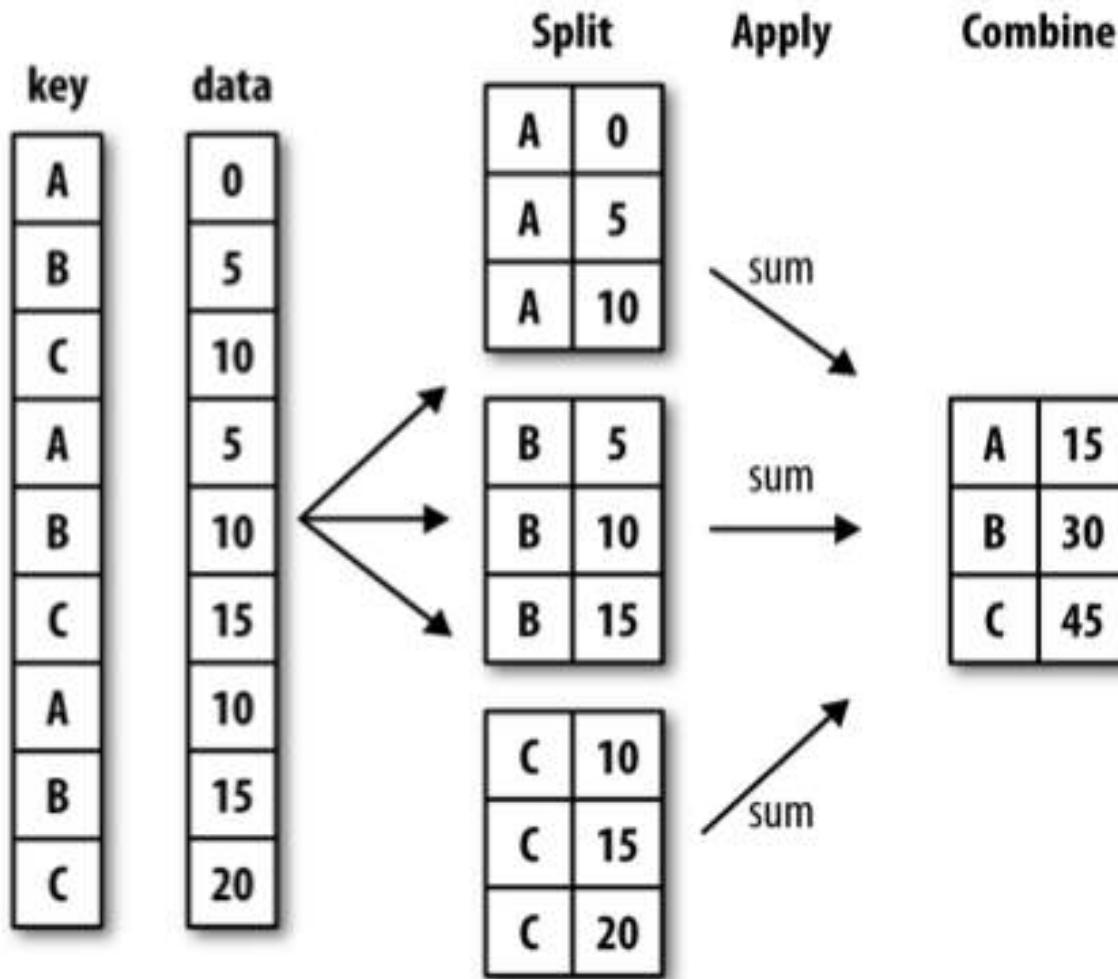


Full EGRIN 2.0 toolchain

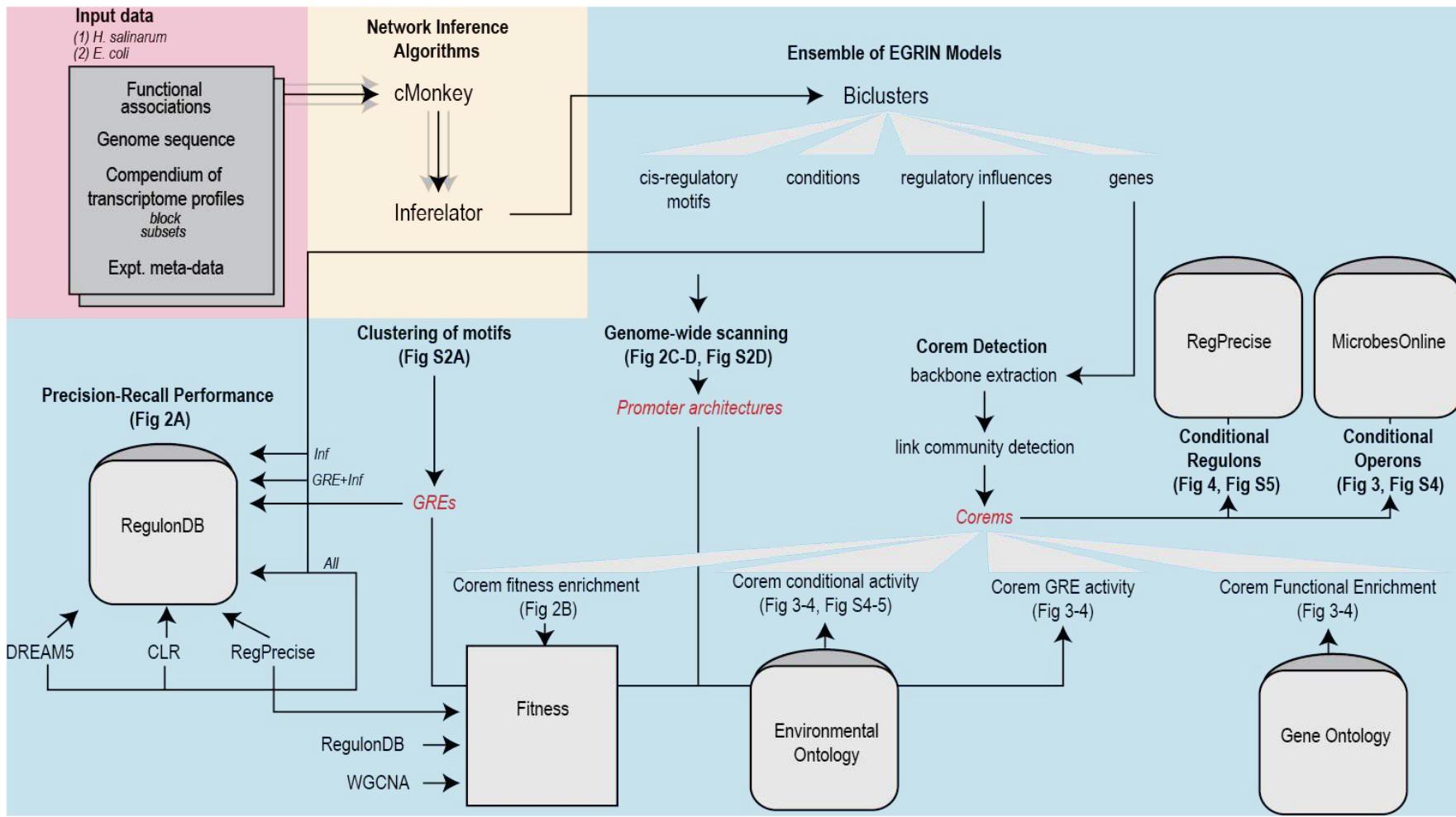




Split-Apply-Combine Approach to Biological Inference

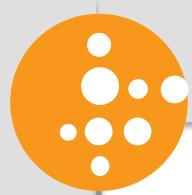


Full EGRIN 2.0 toolchain



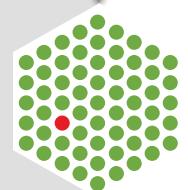
SPLIT → **APPLY** → **COMBINE**

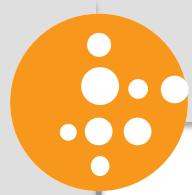
Brooks and Reiss et al 2014



Generalized pattern for reconciling predictions

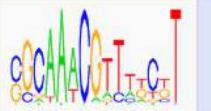
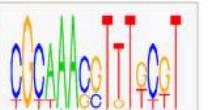
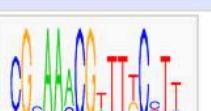
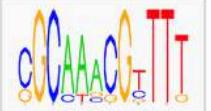
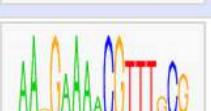
- Combine predictions in a graph
(clean as needed)
- Interpret with graph-based methods
(clustering)



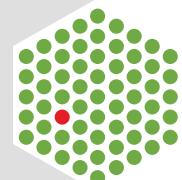


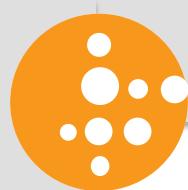
Determine consensus GRE

*GRE = gene regulatory element
(putative TF binding site)*

CRM ID	Bicluster	PSSM	eval
eco_2716_1	eco_2716		0
eco_1863_1	eco_1863		0
eco_2223_1	eco_2223		0
eco_2666_1	eco_2666		0
eco_859_1	eco_859		0
eco_74_1	eco_74		0
eco_357_1	eco_357		0
eco_2076_1	eco_2076		0
eco_3077_1	eco_3077		0
eco_2741_1	eco_2741		0

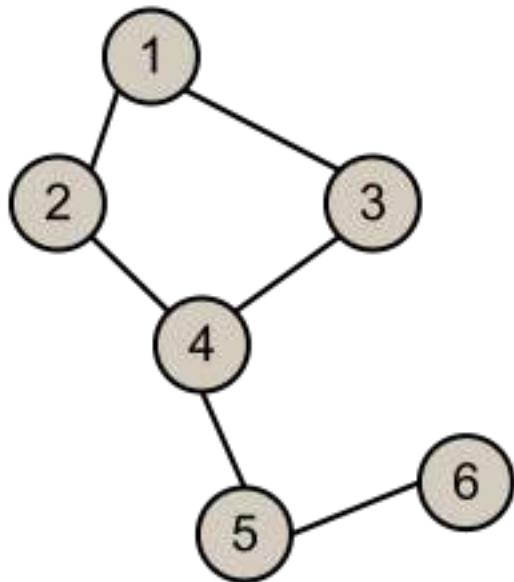
Showing 1 to 10 of 331 entries





Combine: network representation for GReS

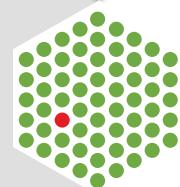
Undirected Graph & Adjacency Matrix

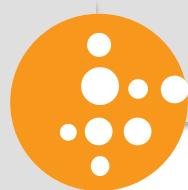


Undirected Graph

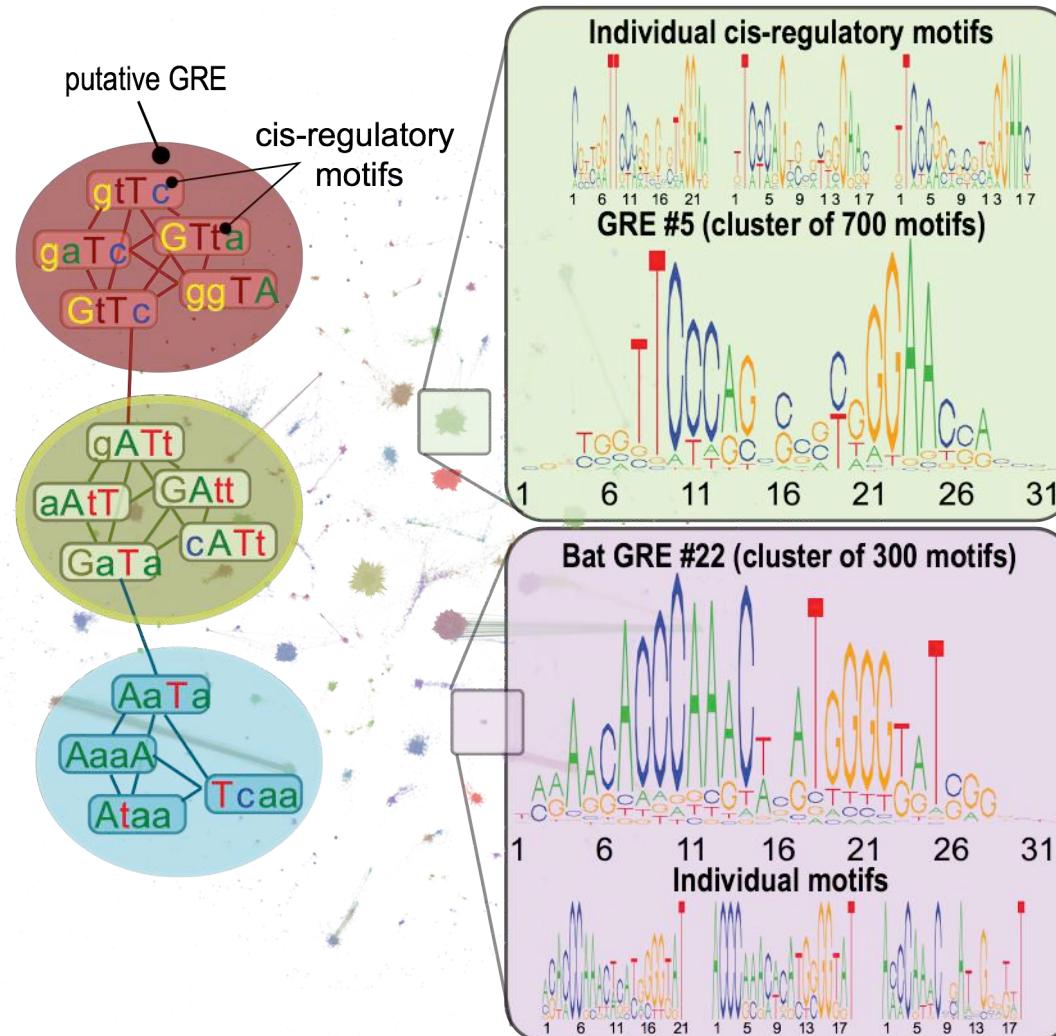
1	2	3	4	5	6
1	0	1	1	0	0
2	1	0	0	1	0
3	1	0	0	1	0
4	0	1	1	0	1
5	0	0	0	1	0
6	0	0	0	0	1

Adjacency Matrix

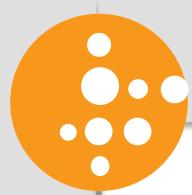




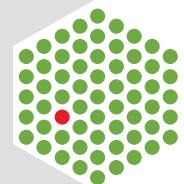
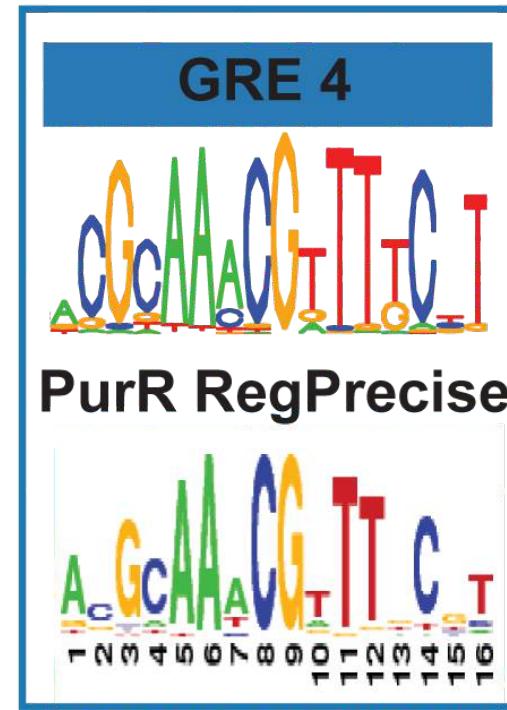
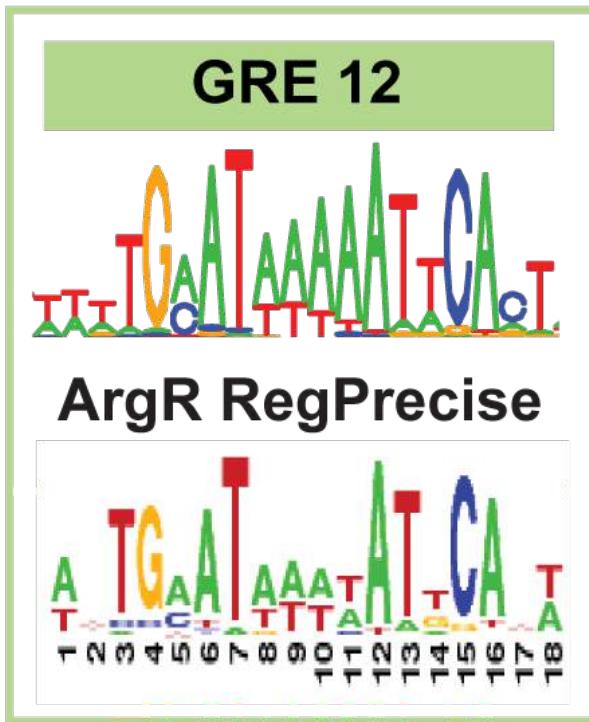
Combine: GRE clustering

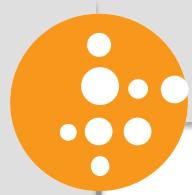


Brooks and Reiss et al 2014

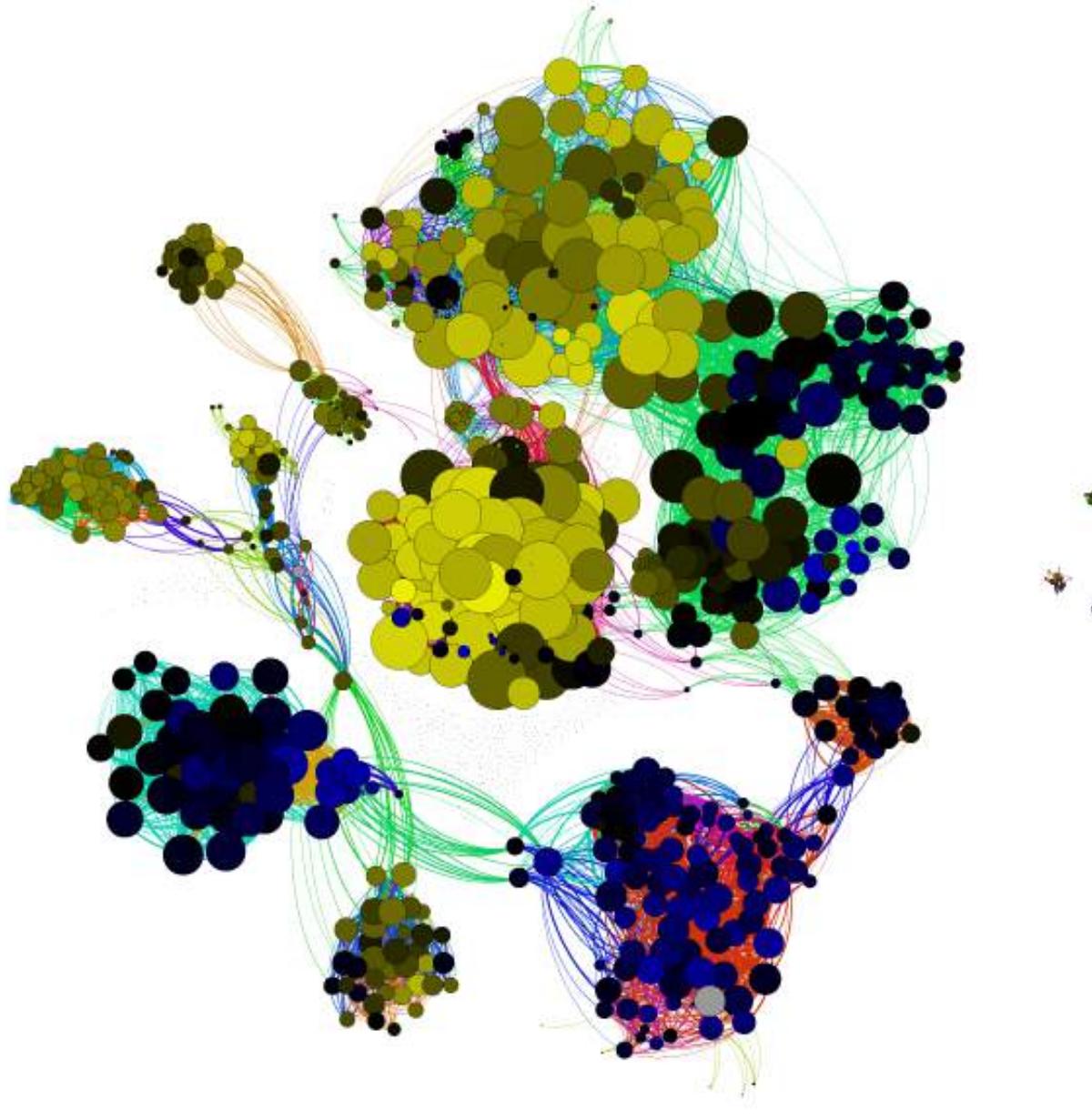


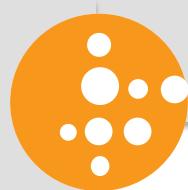
Combine: Inferred GREs resemble known TFBs



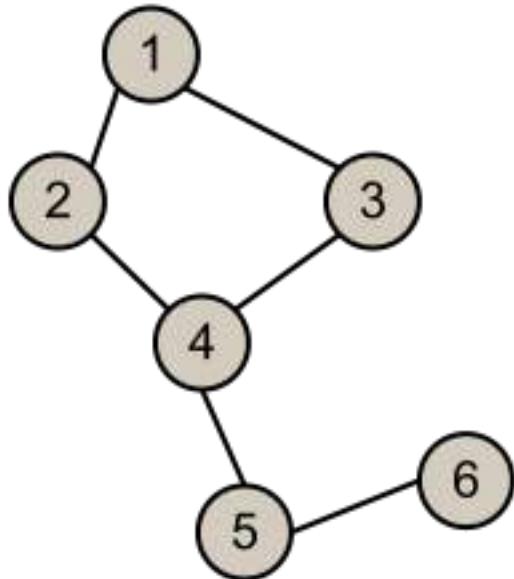


Discover co-regulatory modules





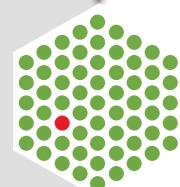
Undirected Graph & Adjacency Matrix

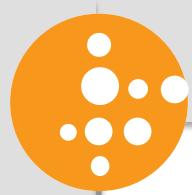


Undirected Graph

1	2	3	4	5	6	
1	0	1	1	0	0	0
2	1	0	0	1	0	0
3	1	0	0	1	0	0
4	0	1	1	0	1	0
5	0	0	0	1	0	1
6	0	0	0	0	1	0

Adjacency Matrix



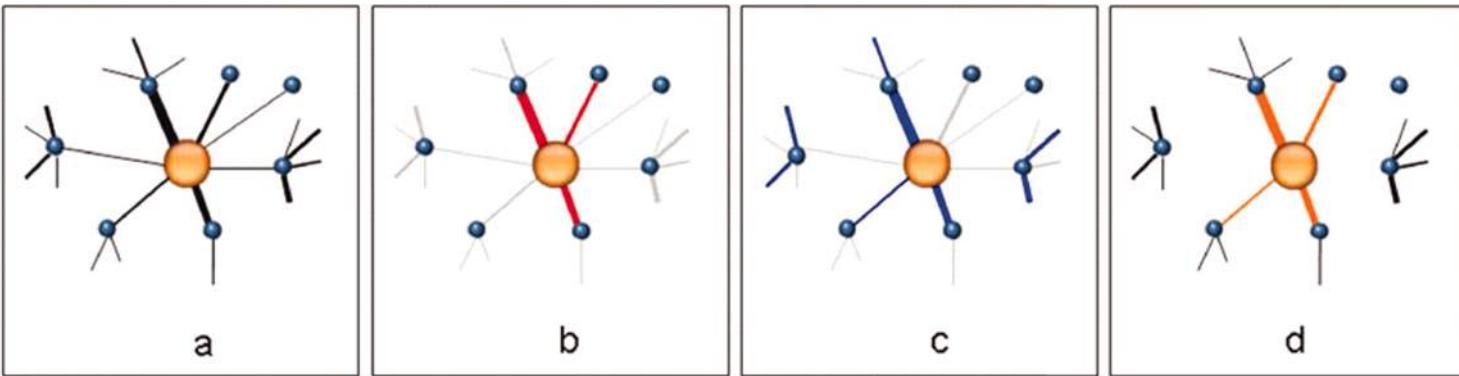


Combine: network backbone extraction

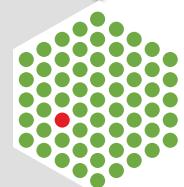
Original



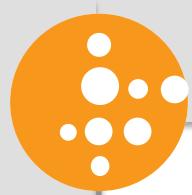
Final



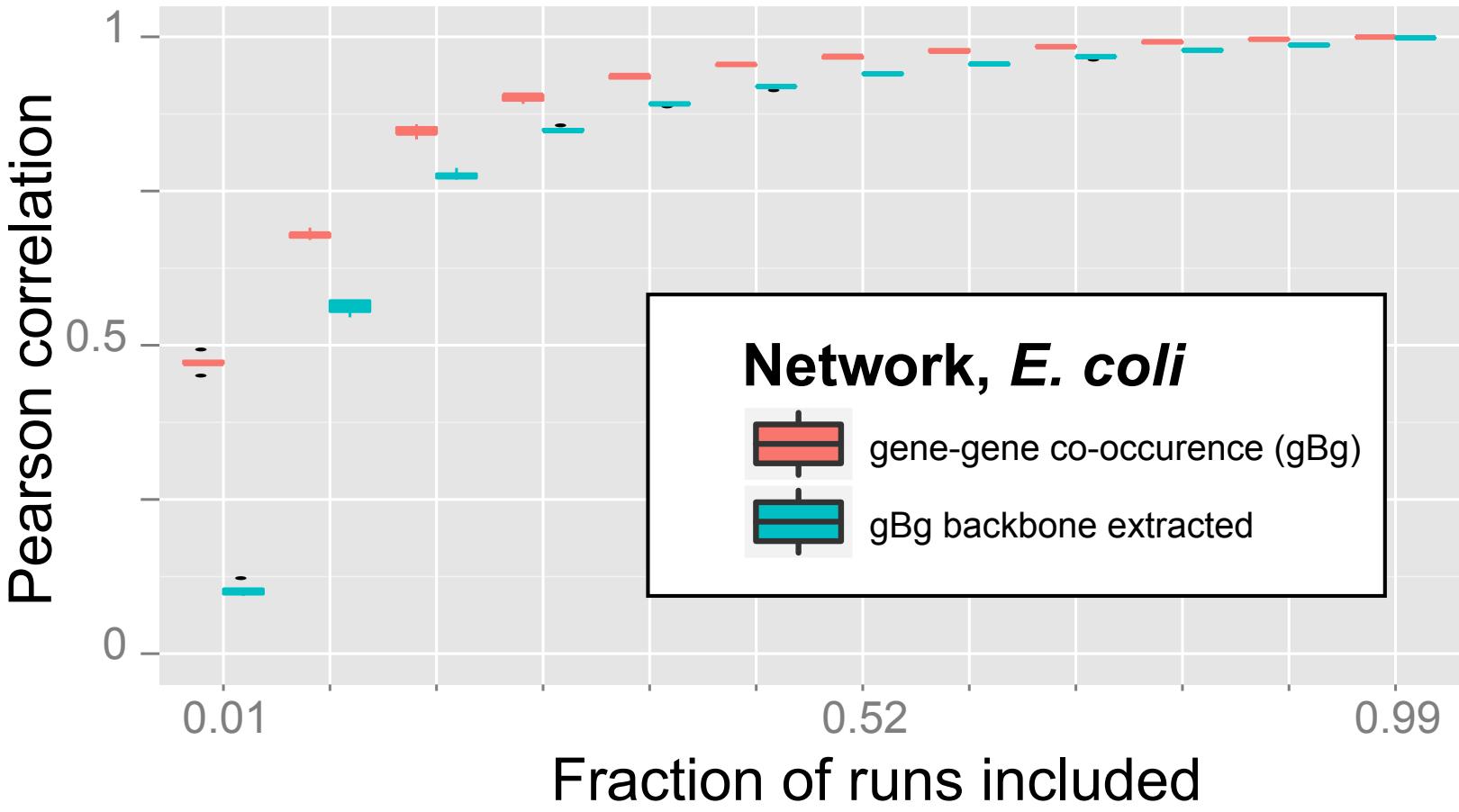
$$\alpha_{ij} = 1 - (k-1) \int_0^{w_{ij}} (1-x)^{k-2} dx \leq 0.05,$$



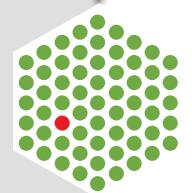
Serrano et al 2009

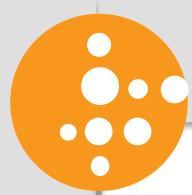


Network structure converges

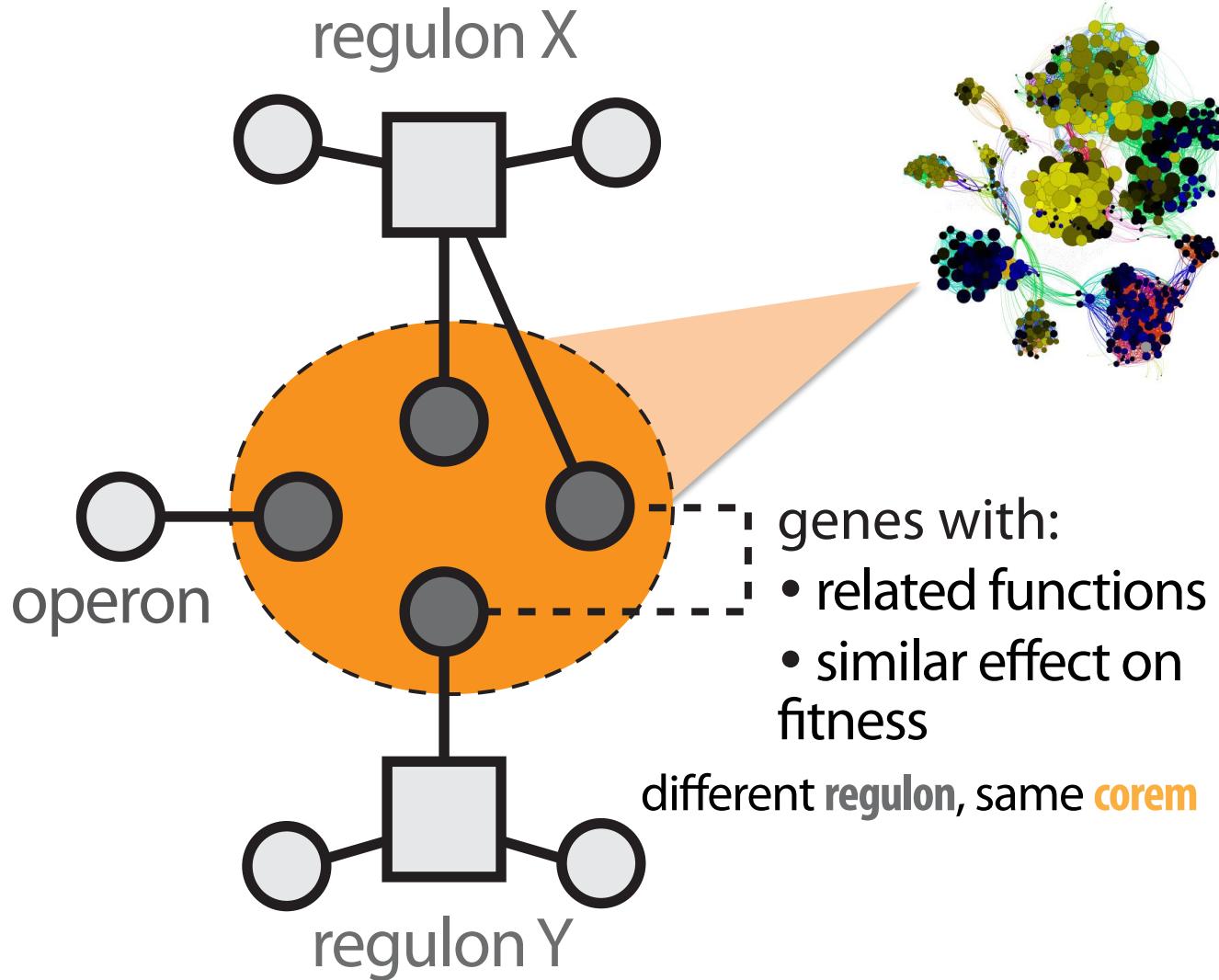


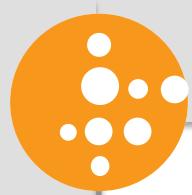
Brooks and Reiss et al 2014



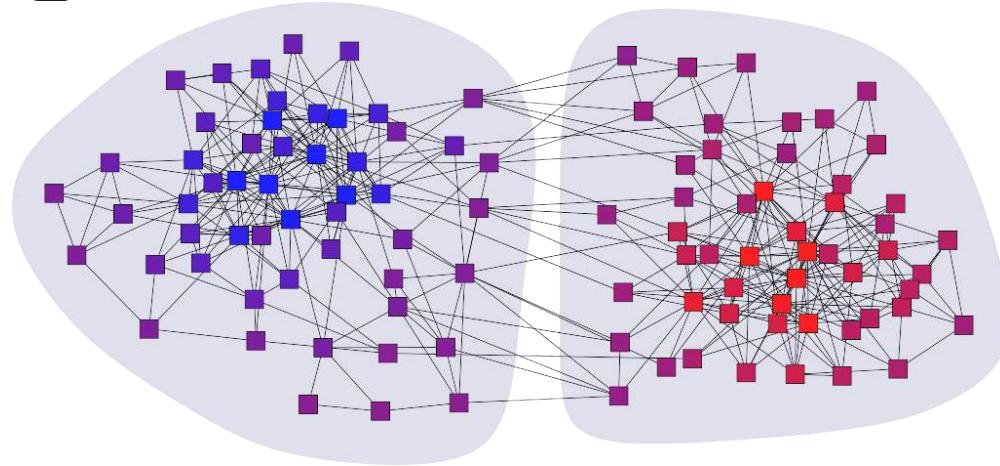
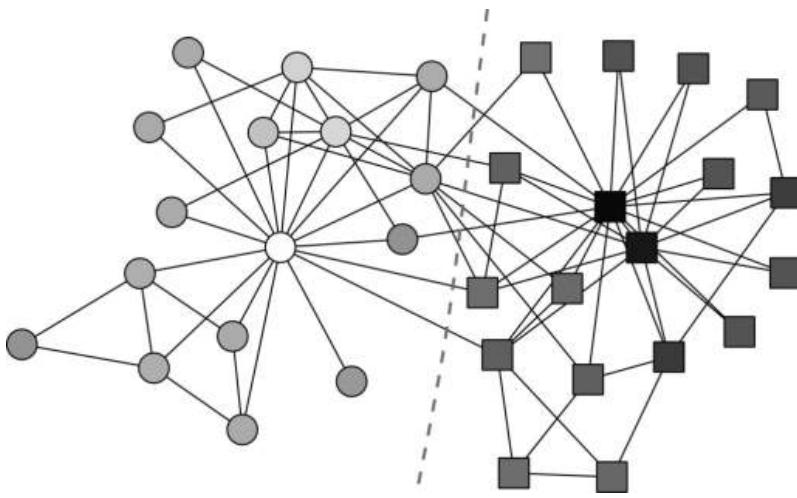


Corems: conditionally co-regulated modules

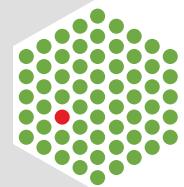


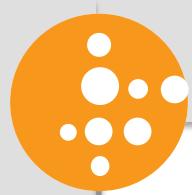


Community Detection

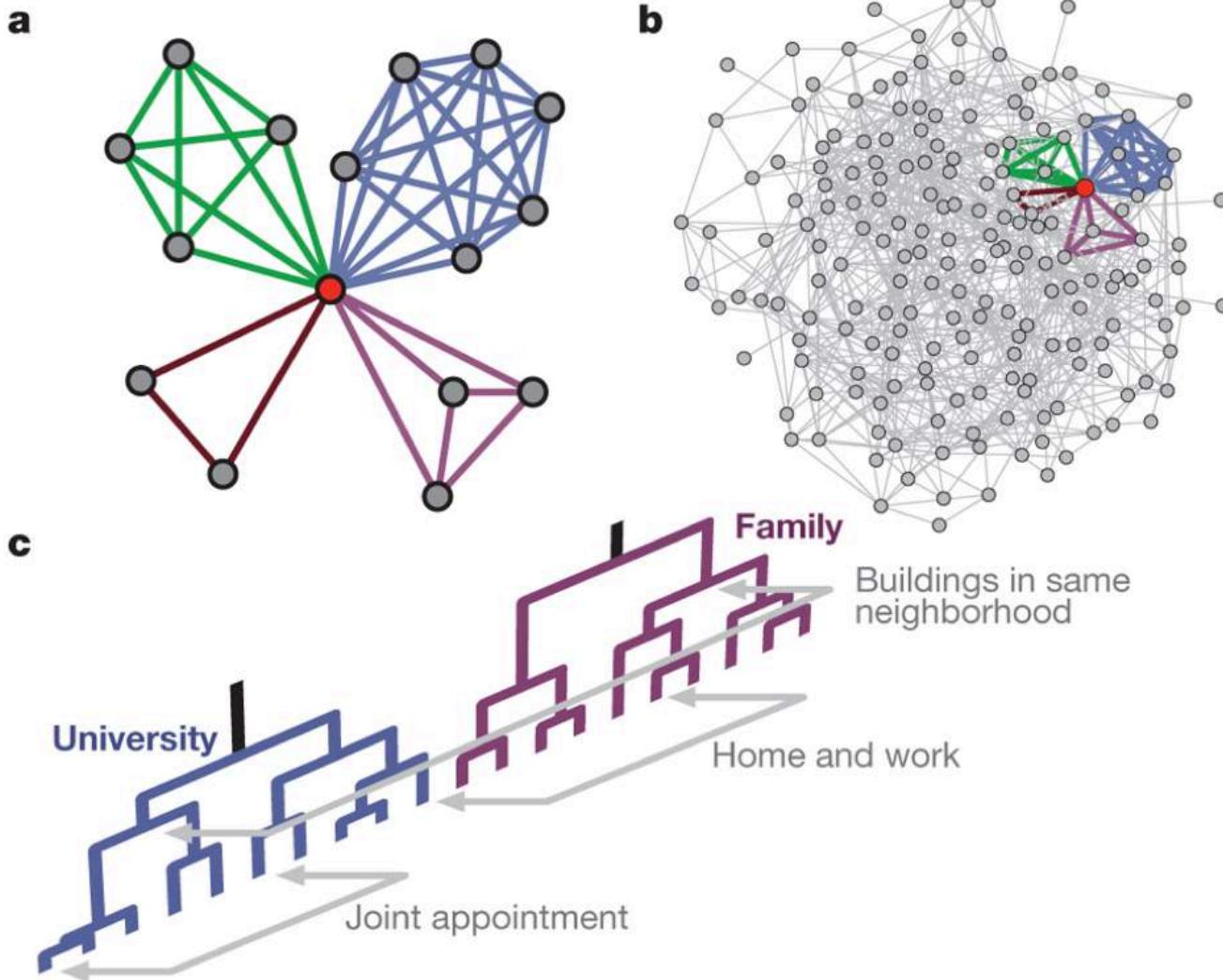


Newman 2006

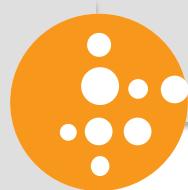




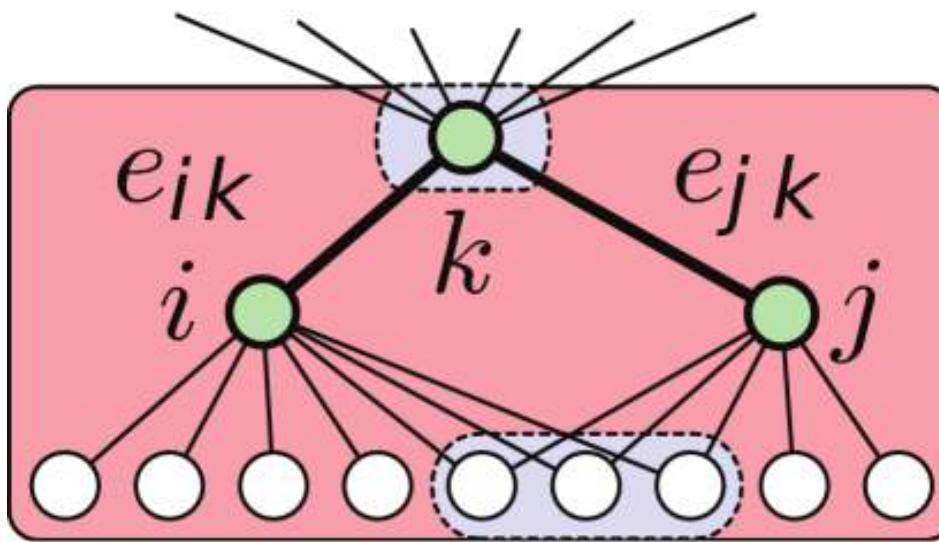
Link Community Detection



Ahn et al 2010

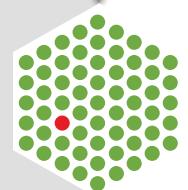


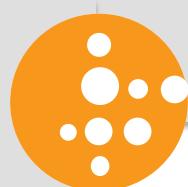
Compute weighted similarity for every pair of edges



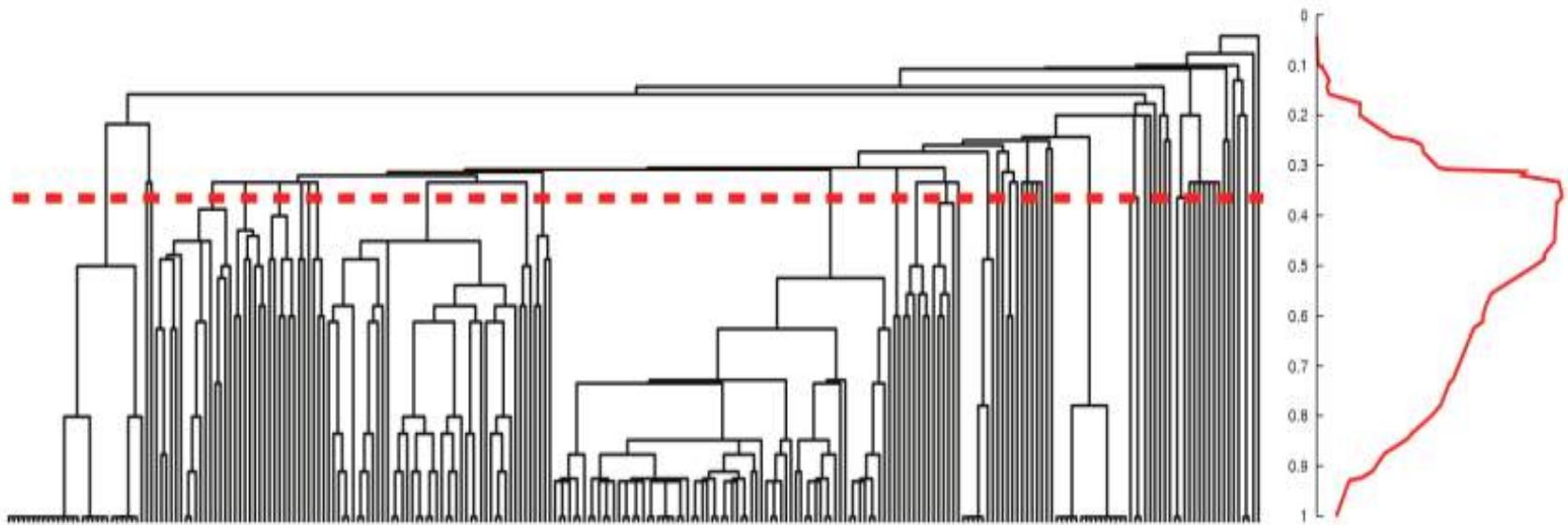
$$\tilde{A}_{ij} = \frac{1}{k_i} \sum_{i' \in n(i)} w_{ii'} \delta_{ij} + w_{ij}$$

$$S(e_{ik}, e_{jk}) = \frac{\mathbf{a}_i \cdot \mathbf{a}_j}{|\mathbf{a}_i|^2 + |\mathbf{a}_j|^2 - \mathbf{a}_i \cdot \mathbf{a}_j}$$



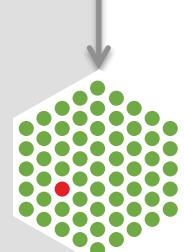


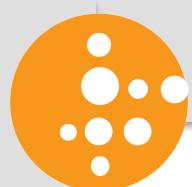
Cluster edges to maximize intra-cluster density



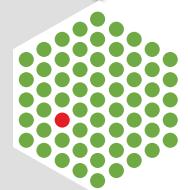
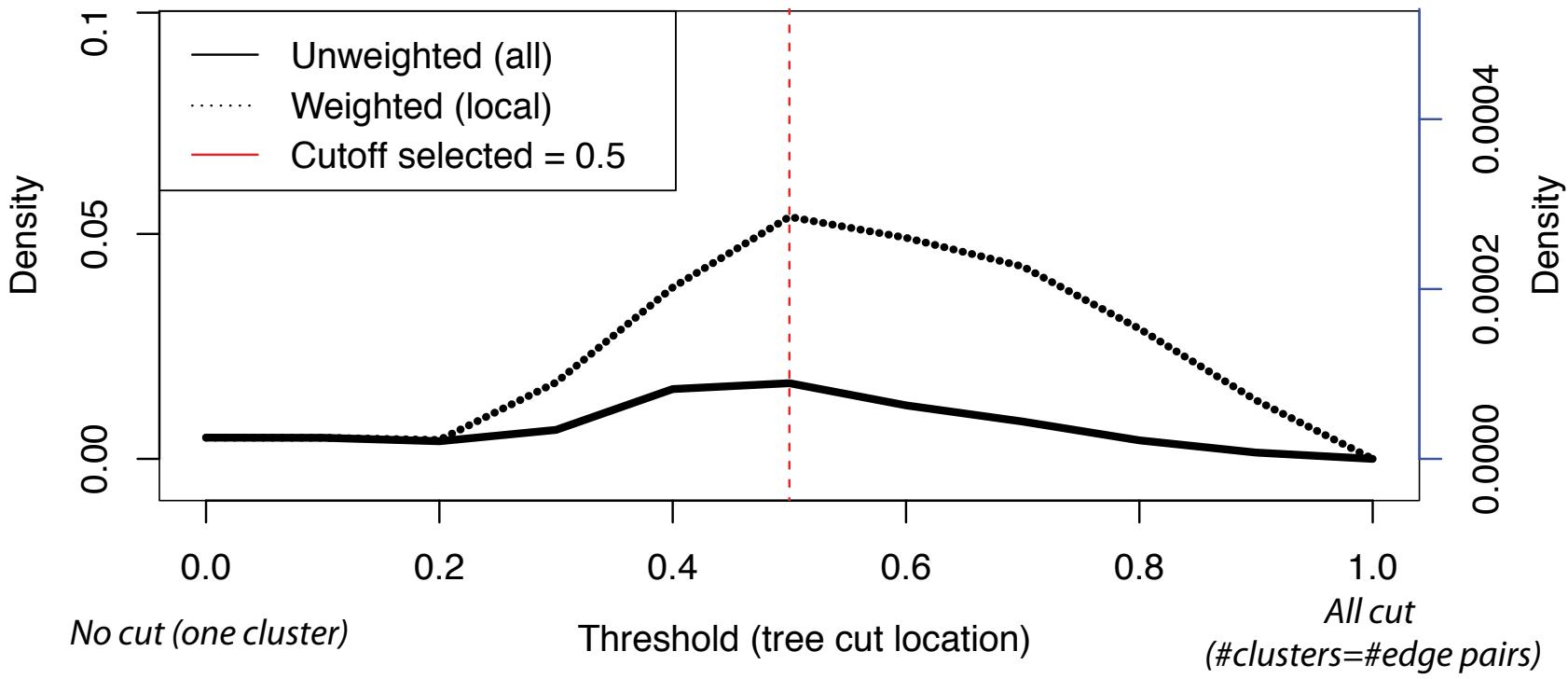
N	Total number of nodes in the network
M	Total number of edges in the network
$\langle w \rangle$	Average weight of edges of the network
\mathcal{C}	Partition of the network into edge communities
\mathcal{E}	Set of all edges in the network ($ \mathcal{E} = M$)
n_c	Number of nodes in community c
m_c	Number of edges in community c
$\langle w \rangle_c$	Average weight of edges in community c

$$D^{(5)} = \frac{1}{M\langle w \rangle} \sum_{c \in \mathcal{C}} m_c \langle w \rangle_c \left(\langle w \rangle_c \left[\frac{m_c - (n_c - 1)}{n_c(n_c - 1)/2 - (n_c - 1)} \right] \right)$$

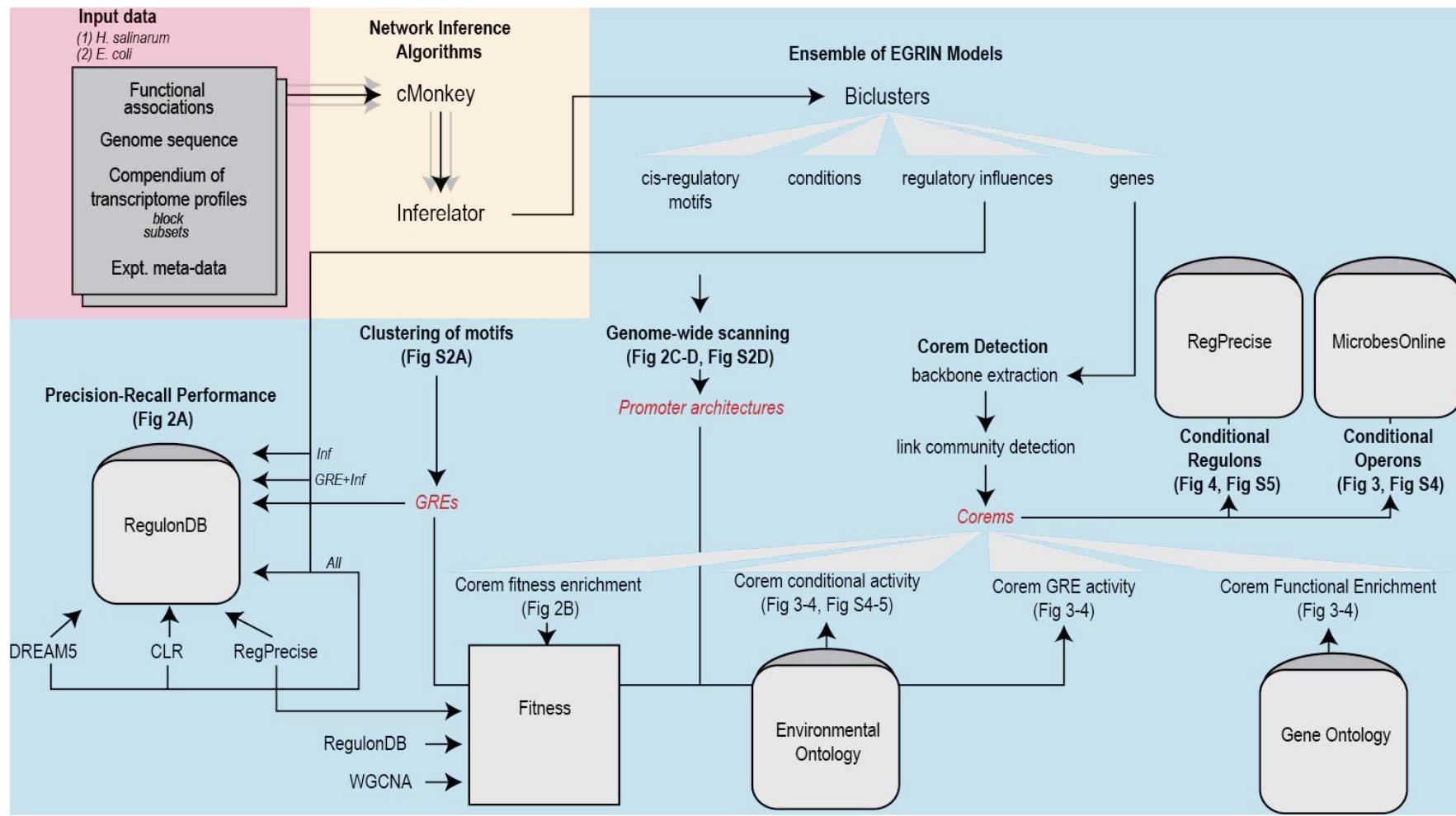




Network density given threshold selection for corem selection, *E. coli*

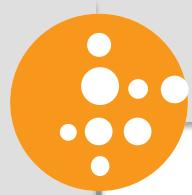


Full EGRIN 2.0 toolchain

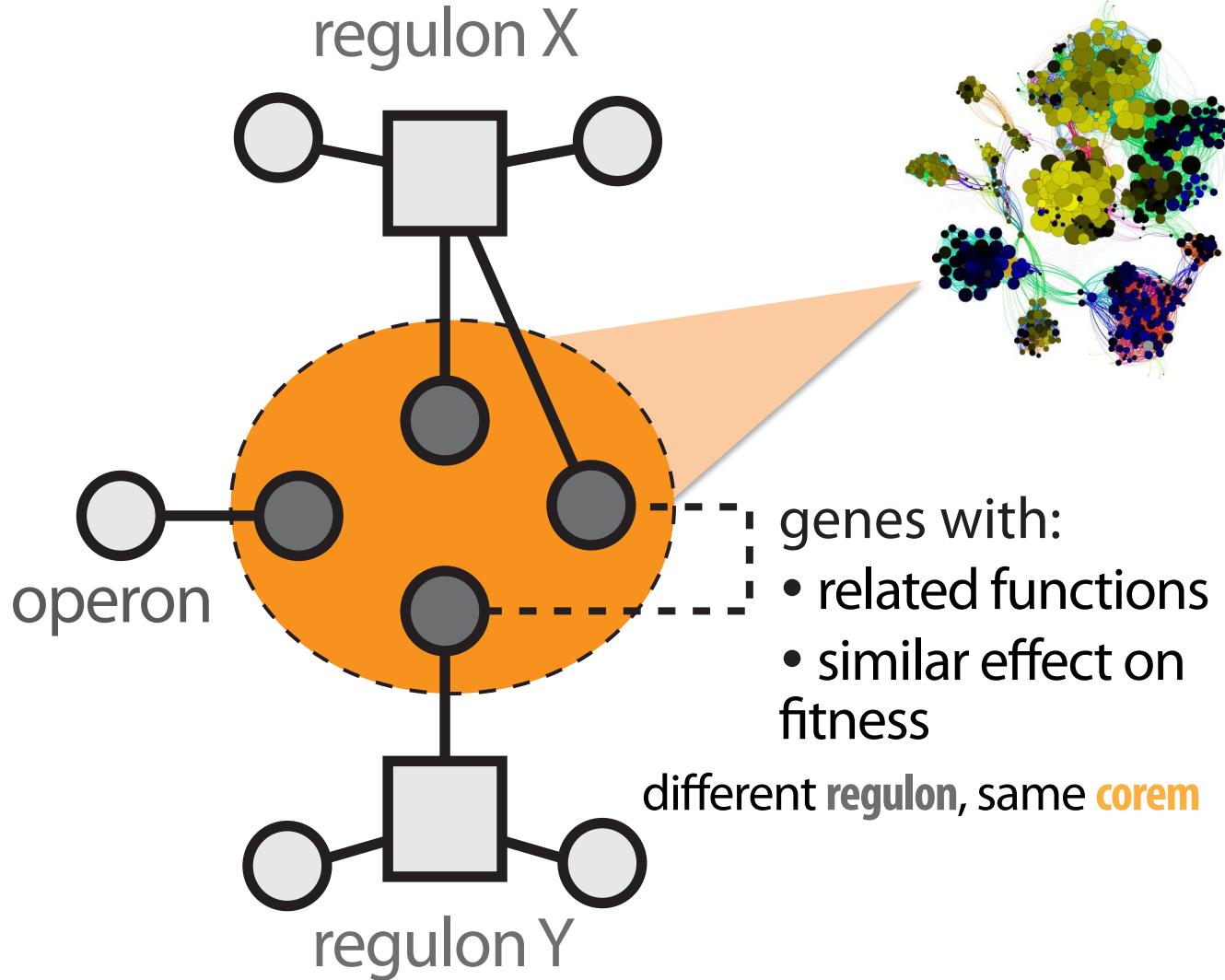


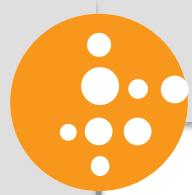
SPLIT → **APPLY** → **COMBINE**

Brooks and Reiss et al 2014

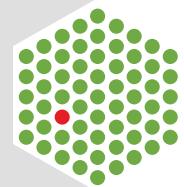
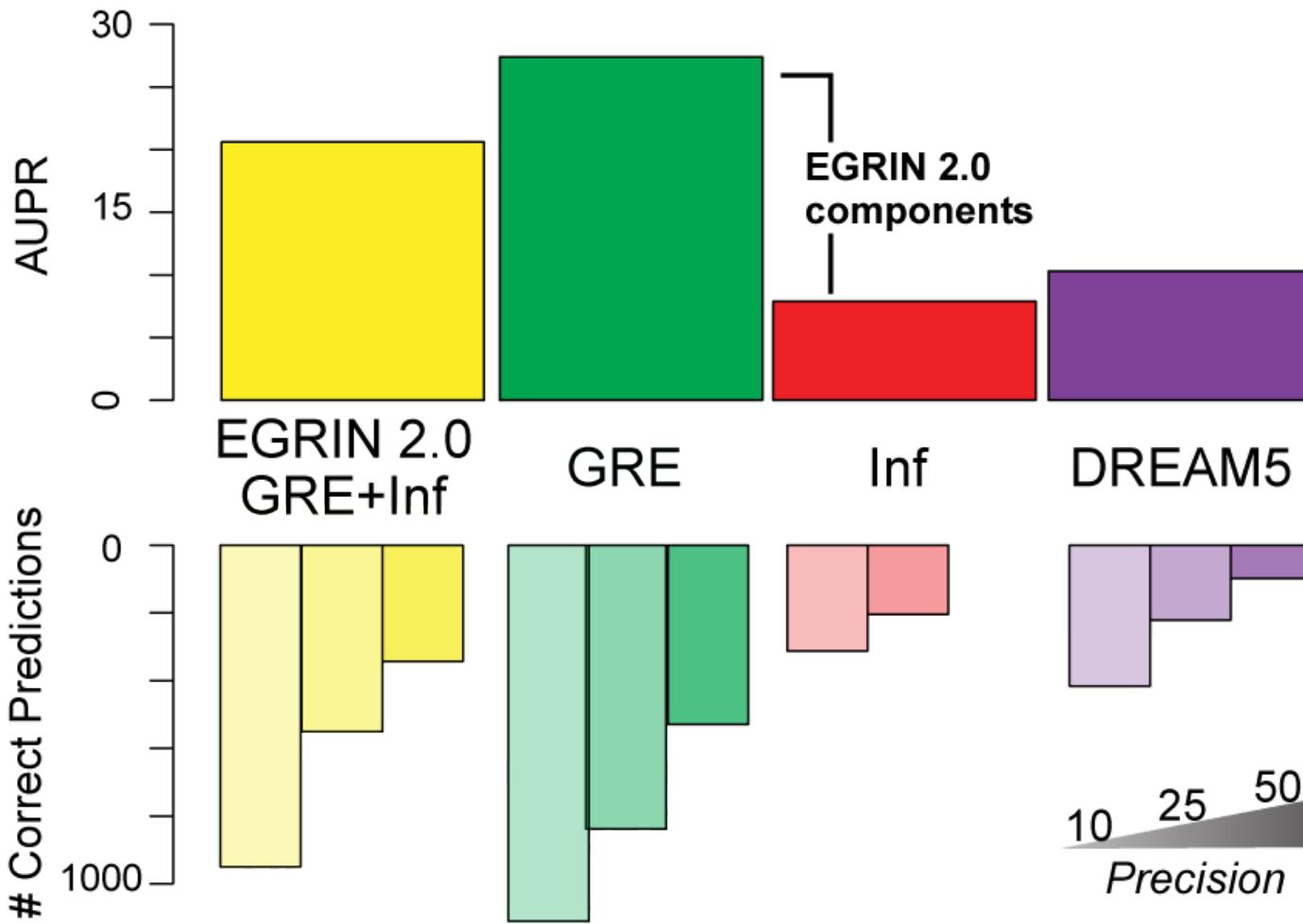


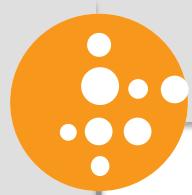
Corems: conditionally co-regulated modules





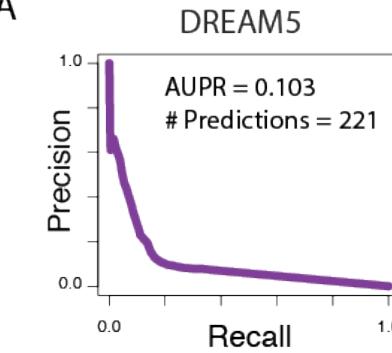
EGRIN 2.0 has ~3-fold increased performance



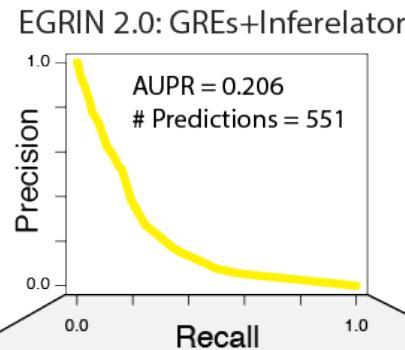


Performance increase due to motif detection

A

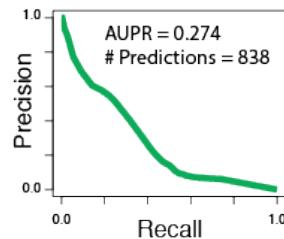


B



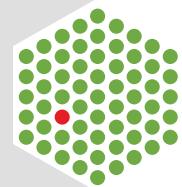
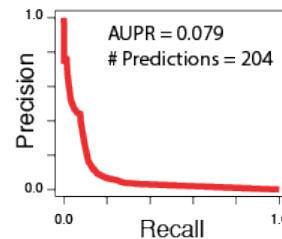
EGRIN 2.0 Components

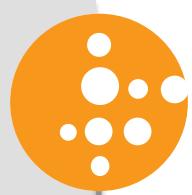
GREs



+

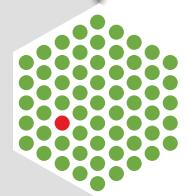
Inferelator



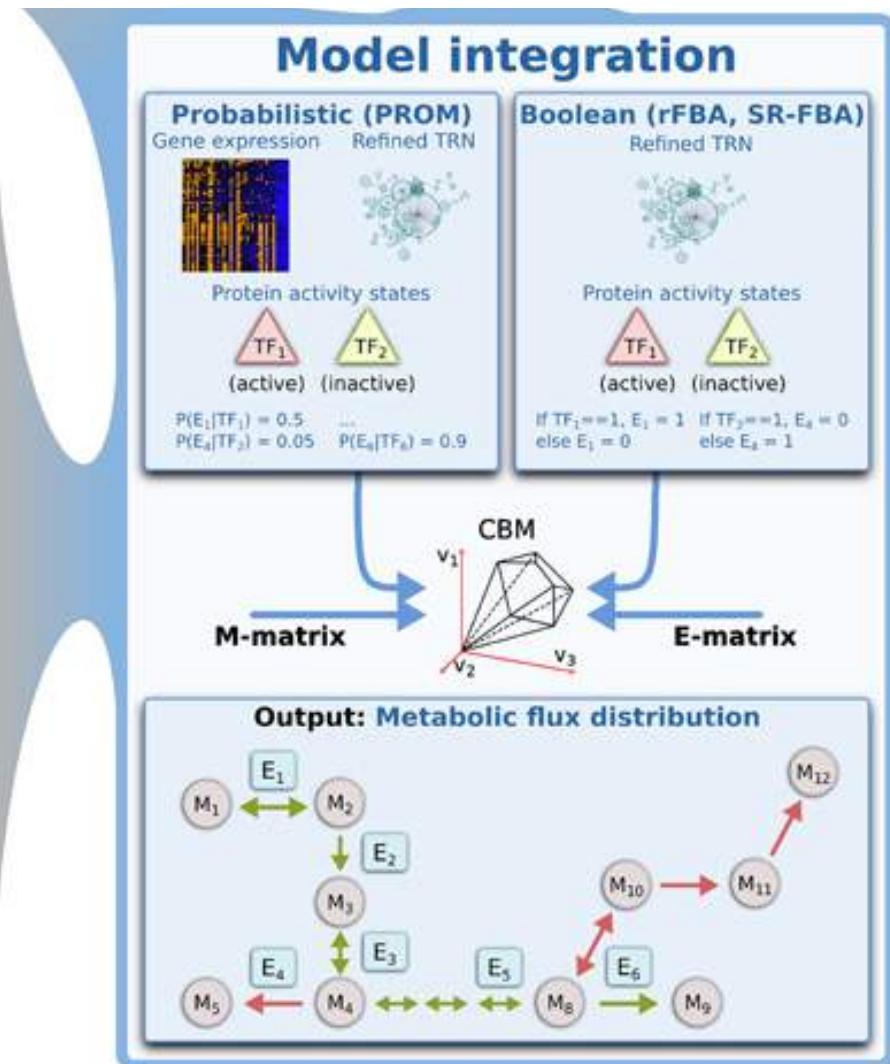
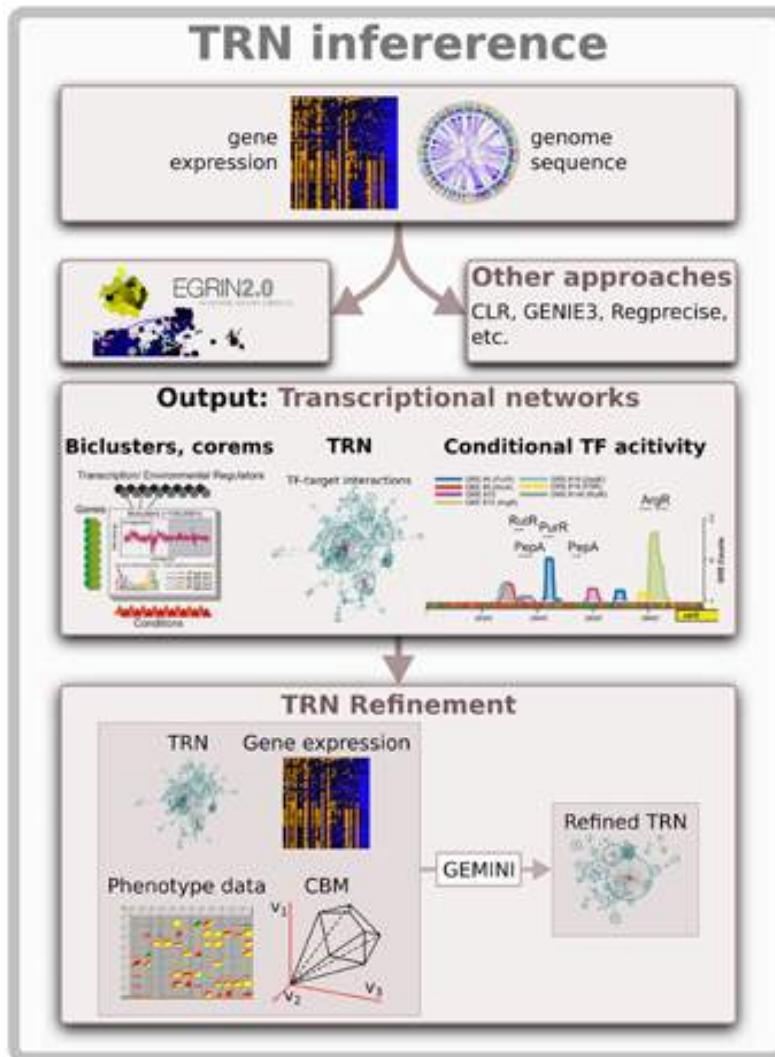


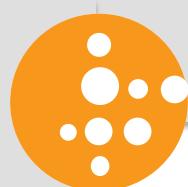
Extensions

Integration with metabolic models

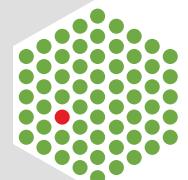
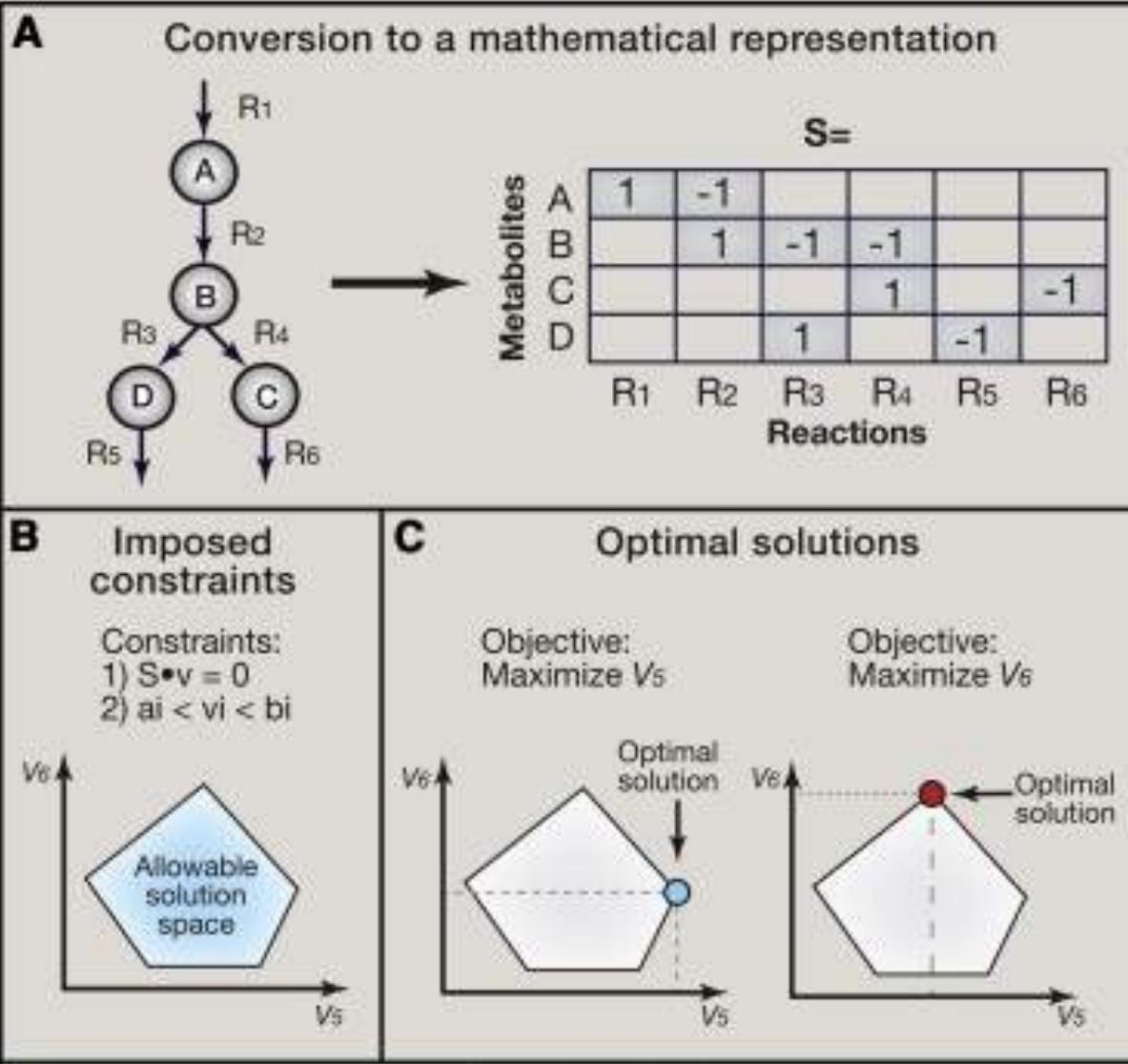


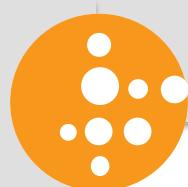
Integration with models of metabolism



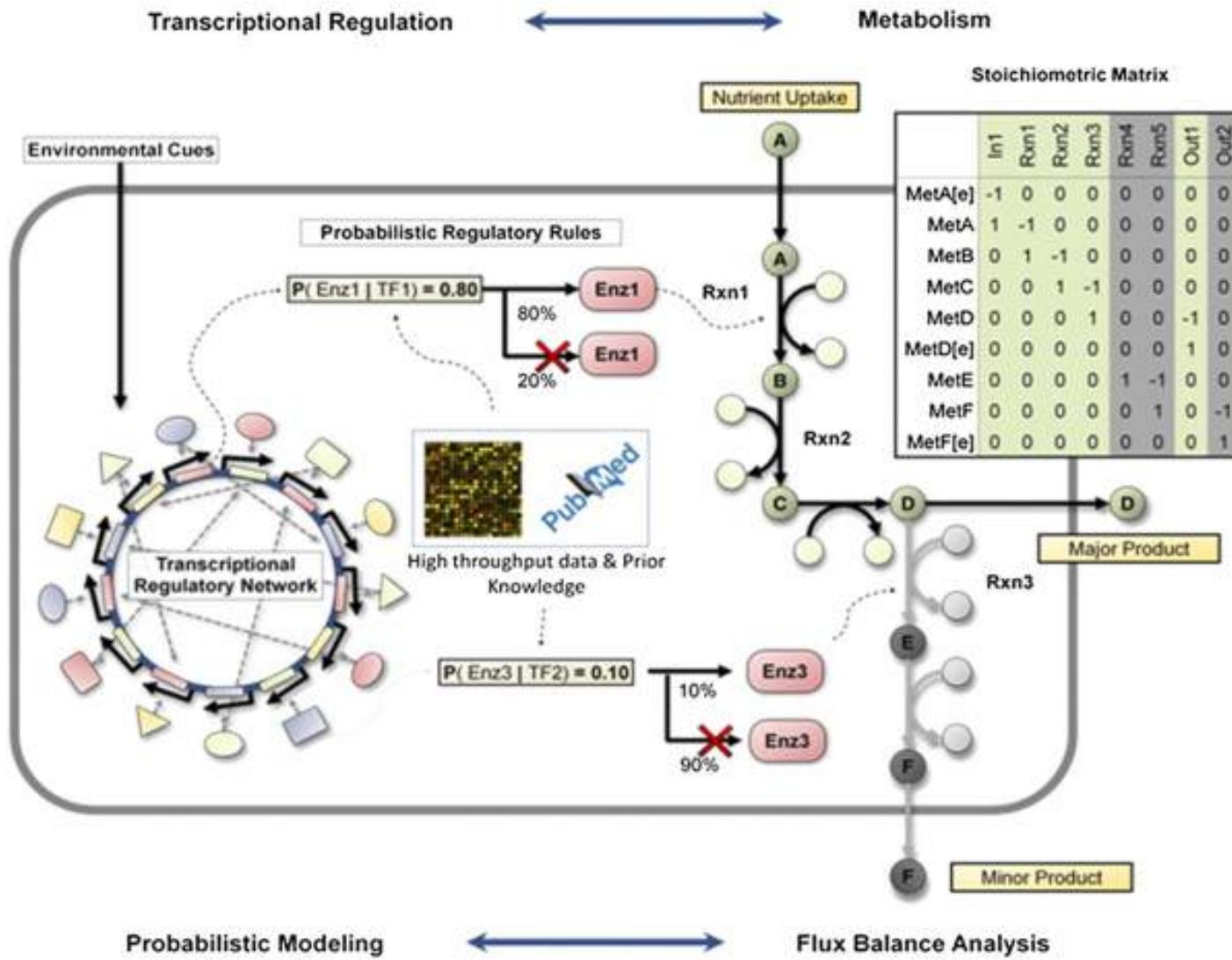


Genome-scale models of metabolism

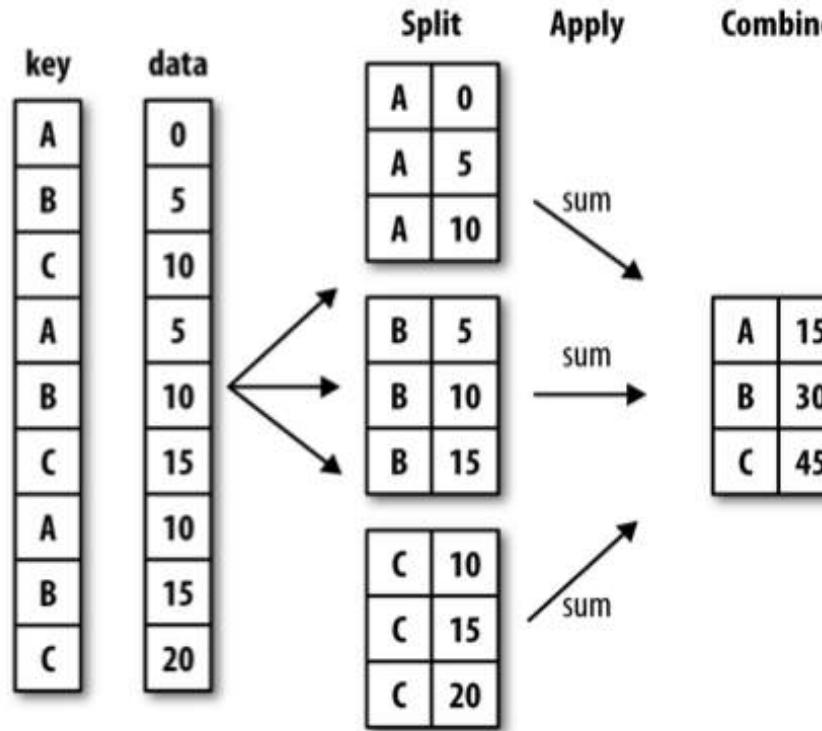




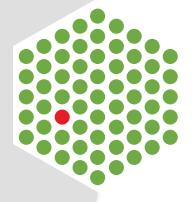
Probabilistic Regulation of Metabolism (PROM)



Split-Apply-Combine-like strategy for building multistep, data-driven inference toolchains



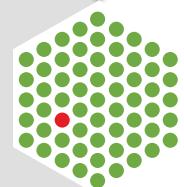
SPLIT → **APPLY** → **COMBINE**

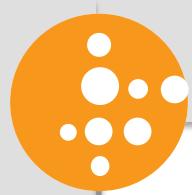


Split-Apply-Combine-like strategy for building multistep, data-driven inference toolchains

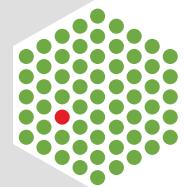
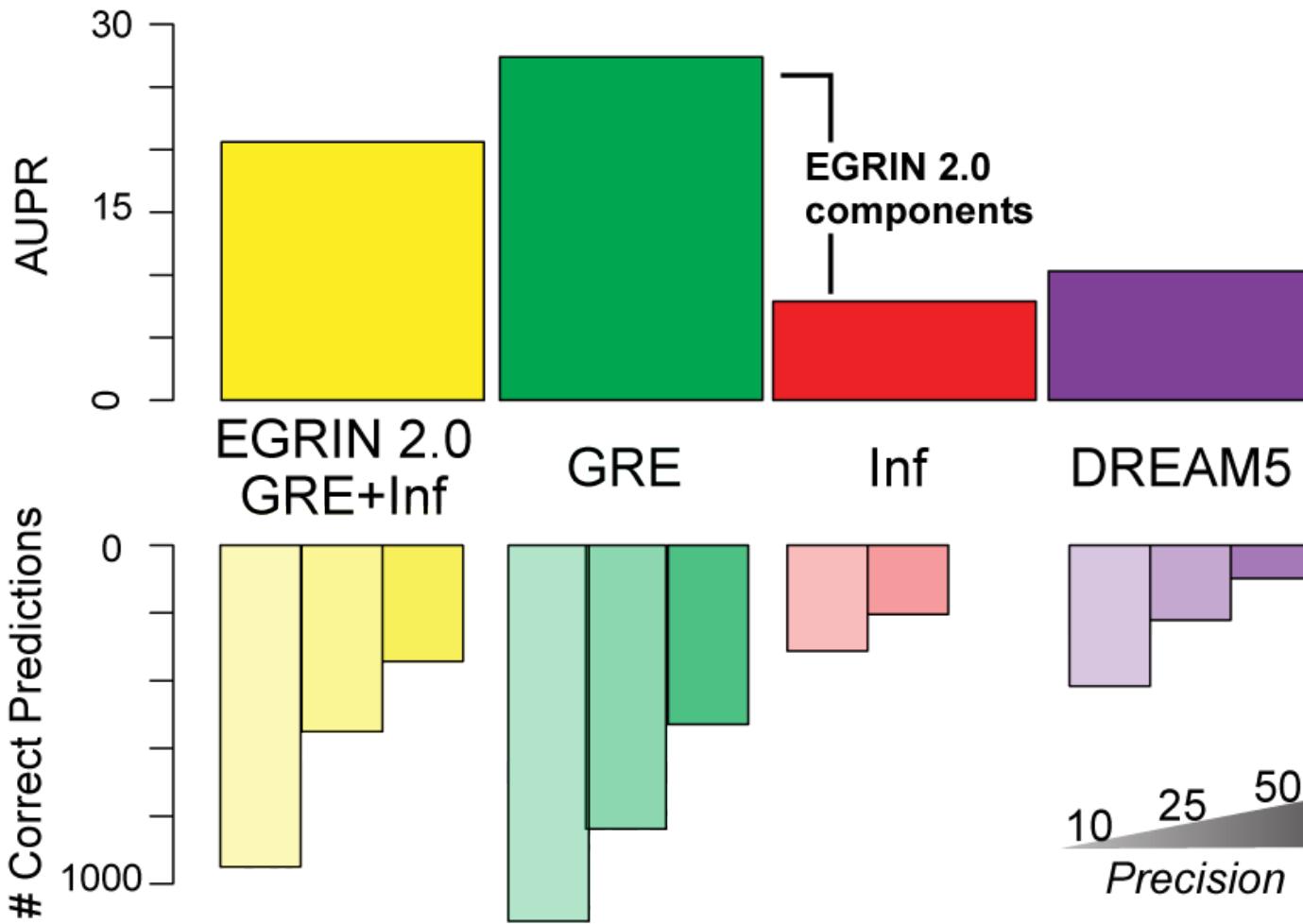
- Heterogeneous data can (and should) be integrated

Increased predictive performance by simultaneous motif detection



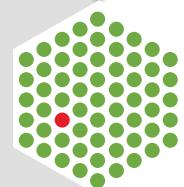


EGRIN 2.0 has ~3-fold increased performance

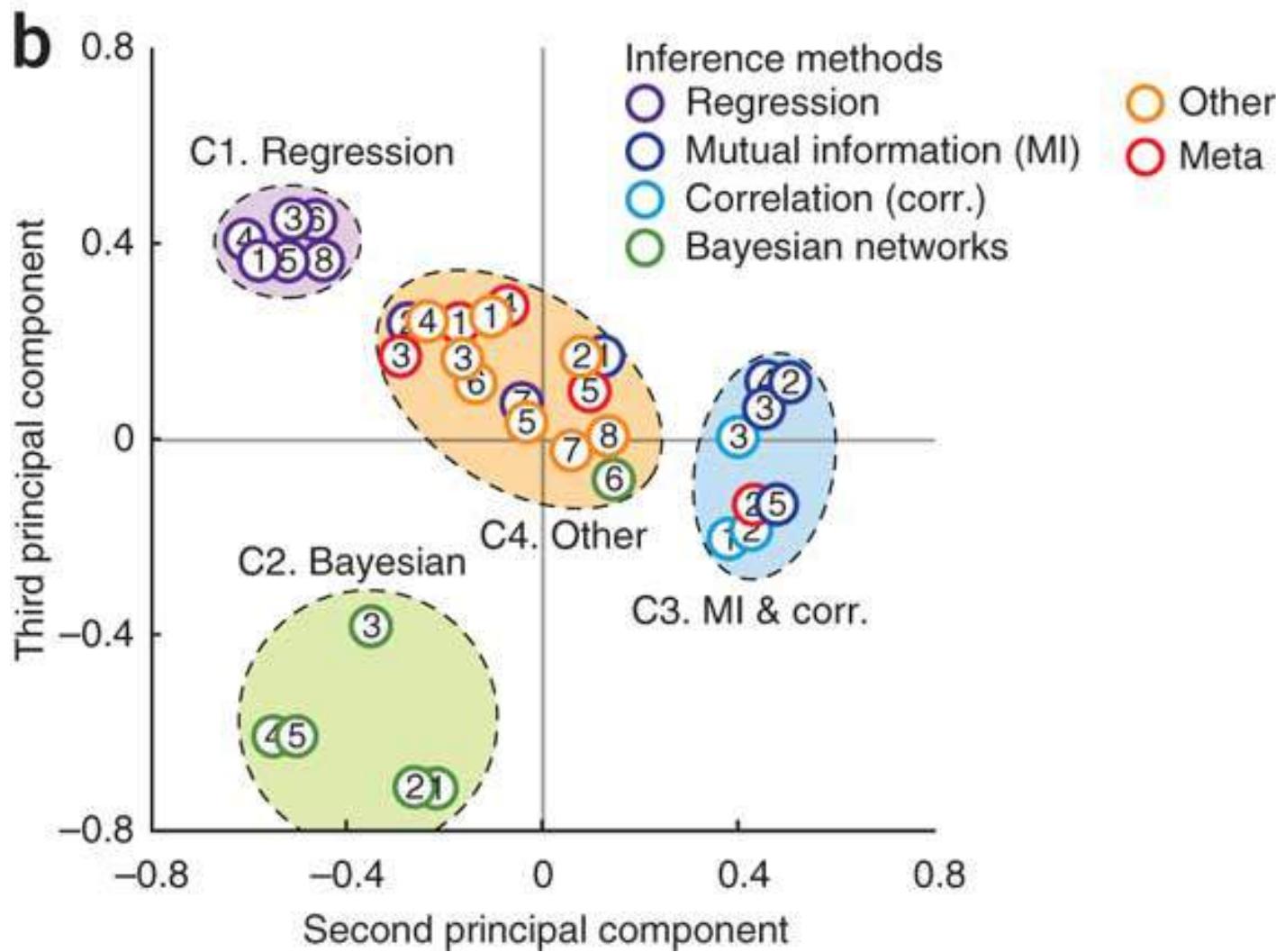


Split-Apply-Combine-like strategy for building multistep, data-driven inference toolchains

- Heterogeneous data can (and should) be integrated
Increased predictive performance by simultaneous motif detection
- Ensemble learning frameworks can harness model diversity
Combining models increases overall model performance

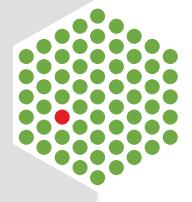


Algorithms learn different network features

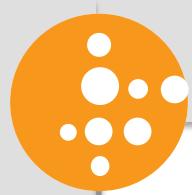


Split-Apply-Combine-like strategy for building multistep, data-driven inference toolchains

- Heterogeneous data can (and should) be integrated
Increased predictive performance by simultaneous motif detection
- Ensemble learning frameworks can harness model diversity
Combining models increases overall model performance
- Graph-based methods can flexibly encode and reconcile predictions



Examples of model reconciliation with graphs: GRE and corem detection



New tools for SAC-biology

The screenshot shows a web browser window with the URL nbviewer.ipython.org/github/baliga-lab/egrin2-tools/blob/master/doc/index.ipynb. The browser interface includes standard navigation buttons, a search bar, and a menu bar with items like Apps, Bookmarks, Logon, Cooking, Tools, Spritzlet, Other Books, nbviewer, FAQ, IPython, and Jupyter. Below the menu is a breadcrumb trail: egrin2-tools / doc. The main content area features a large, abstract visualization of overlapping colored circles (yellow, purple, blue) on a light background. The text "Welcome to EGRIN 2.0" is centered above the visualization. Below the visualization, there is descriptive text and three icons: a puzzle piece labeled "BUILD", another puzzle piece labeled "ASSEMBLE", and a flask labeled "QUERY". A small arrow points from the bottom left towards the "BUILD" icon.

Welcome to EGRIN 2.0

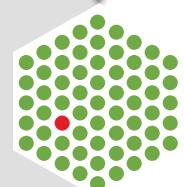
These documents will introduce you to the EGRIN 2.0 ensemble and all of its associated tools.

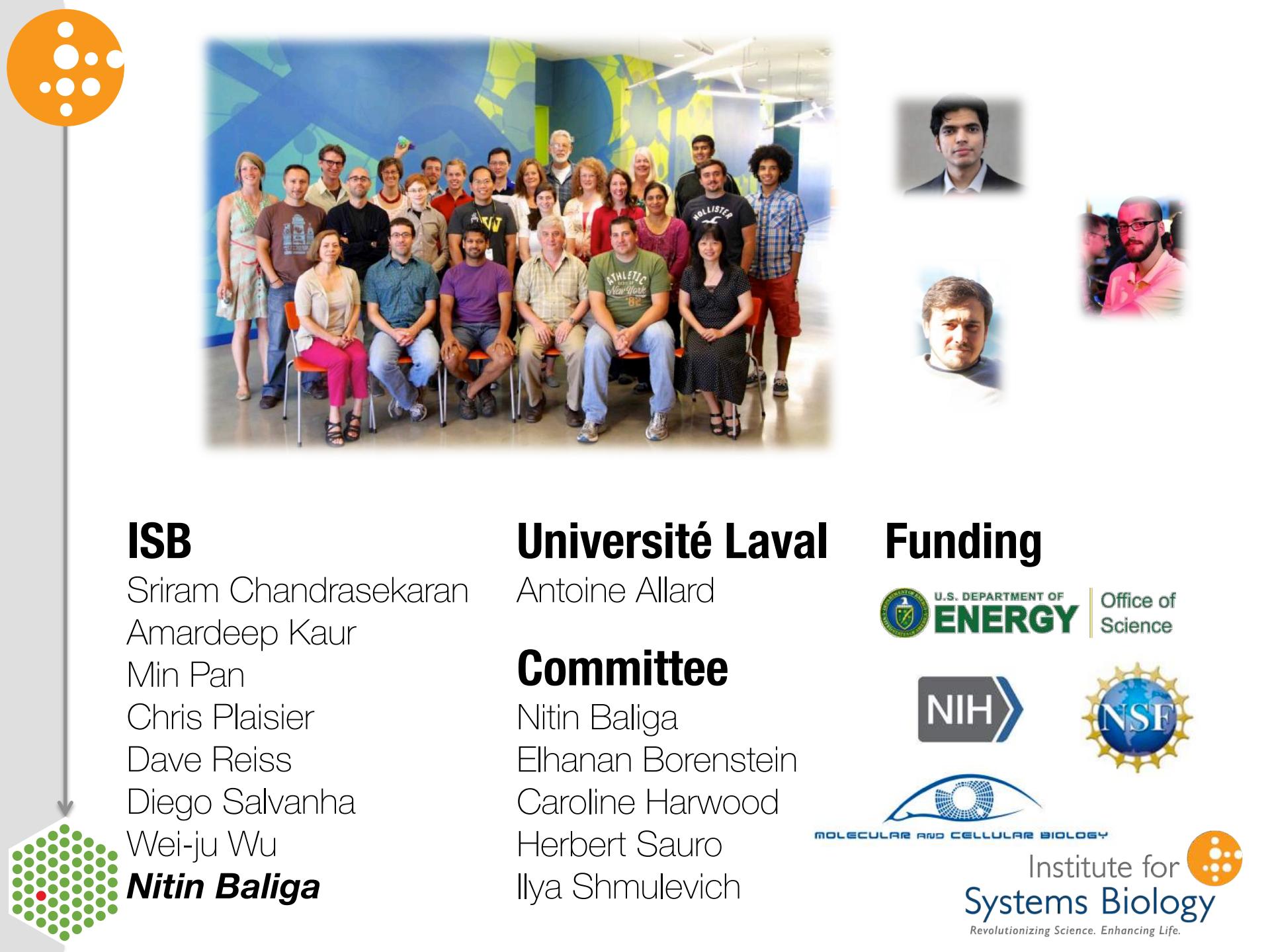
Given a little bit of your time and patience, you will learn how to:

an EGRIN 2.0 ensemble.

You can click on one of the images above to begin your exploration of the tools or read on for more information.

github.com/baliga-lab/egrin2-tools





ISB

Sriram Chandrasekaran

Amardeep Kaur

Min Pan

Chris Plaisier

Dave Reiss

Diego Salvanha

Wei-ju Wu

Nitin Baliga

Université Laval

Antoine Allard

Committee

Nitin Baliga

Elhanan Borenstein

Caroline Harwood

Herbert Sauro

Ilya Shmulevich

Funding



Office of
Science



MOLECULAR AND CELLULAR BIOLOGY

Institute for
Systems Biology
Revolutionizing Science. Enhancing Life.