# Physicist Meets Biology



http://xkcd.com
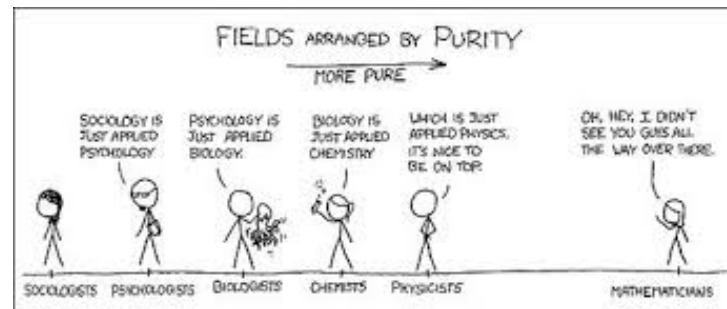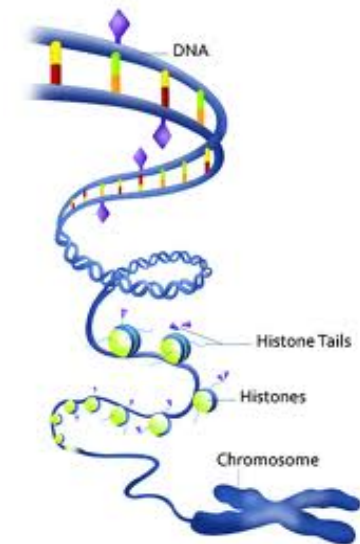
Sara Knaack

Urbana Champaign

4/8/2016

# Outline

- Introduction

- Work on the MuCap Experiment
  - Analysis of complex data in a high-throughput computing environment.

- A brief primer to biology and introduction to the work I do now

- Questions

Thanks to Sushmita Roy and Alireza Siahpirani for contributing slides on gene regulatory networks.

# My own background

- B.S. from the University of Wisconsin at Madison  - Math and Physics
  - Also took an introductory biology curriculum and organic chemistry.

- M.S. from the University of Illinois at Urbana
  - Course work in mathematical methods, quantum mechanics, field theory and statistical mechanics
  - Did beam line simulation work for the g-2 experiment.

- Ph.D. from the University of Illinois at Urbana
  - Work on the MuCap experiment, muon capture on the proton

- **2012 - Present** – Postdoctoral Trainee in Computational Biology at the Wisconsin Institute for Discovery.
  - Research in the regulation of gene expression in the context of evolution and cancer.
  - Funded by the CIBM program – more about that in a moment.
  - Capstone certificate in Bioinformatics – course work in computational biology, statistics and graphical models

# More about the Wisconsin Institute for Discovery.

An inter disciplinary research environment, with many themes focusing on biomedical medical research, but also many other initiatives
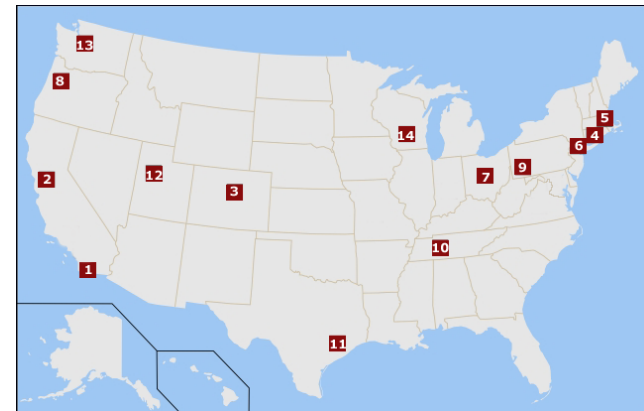




https://discovery.wisc.edu

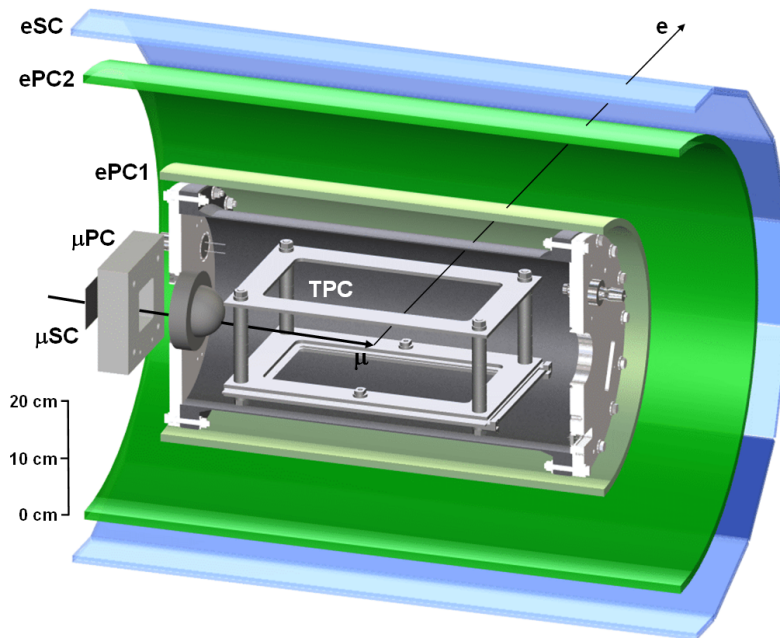# The Computational Informatics in Biology and Medicine program.

- I am a funded by the CIBM program on the campus of the University of Wisconsin, which is a training program through the National Library of Medicine.

- http://www.cibm.wisc.edu

- We are one of multiple training programs run at campuses across the country, including Stanford and Harvard.

- https://www.nlm.nih.gov/ep/GrantTrainInstitute.html

- Provided a core community of fellow trainees and PI's to interact with

- Activities include the annual NLM conferences, and a weekly seminar during the academic year.

Just what kind of training did I come from?

Indulge me three slides on my thesis work…

# The MuCap Experiment



- Muon capture on the proton (MuCap)

- Grew out of the study of hydrogen fusion
    - at the level of fundamental particle interactions.

    - *Physics motivation: quark-gluon substructure of the proton, $g_p$*

- My work was to measure the rate of molecular state formation.

# Description of the time distribution

$$n'_{\mu p}(t) = -(\lambda_\mu + \Lambda_{pp\mu} + \Lambda_{pAr} + \Lambda_S + \Lambda_{pf})n_{\mu p}(t),$$

$$n'_{\mu Ar}(t) = \Lambda_{pAr}n_{\mu p}(t) - (h\lambda_\mu + \Lambda_{Ar})n_{\mu Ar}(t),$$

$$n'_{Ortho}(t) = \Lambda_{pp\mu}n_{\mu p}(t) - (\lambda_\mu + \lambda_{op} + \Lambda_O)n_{Ortho}(t),$$

$$n'_{Para}(t) = \Lambda_{pf}n_{\mu p}(t) + \lambda_{op}n_{Ortho}(t) - (\lambda_\mu + \Lambda_P)n_{Para}(t).$$

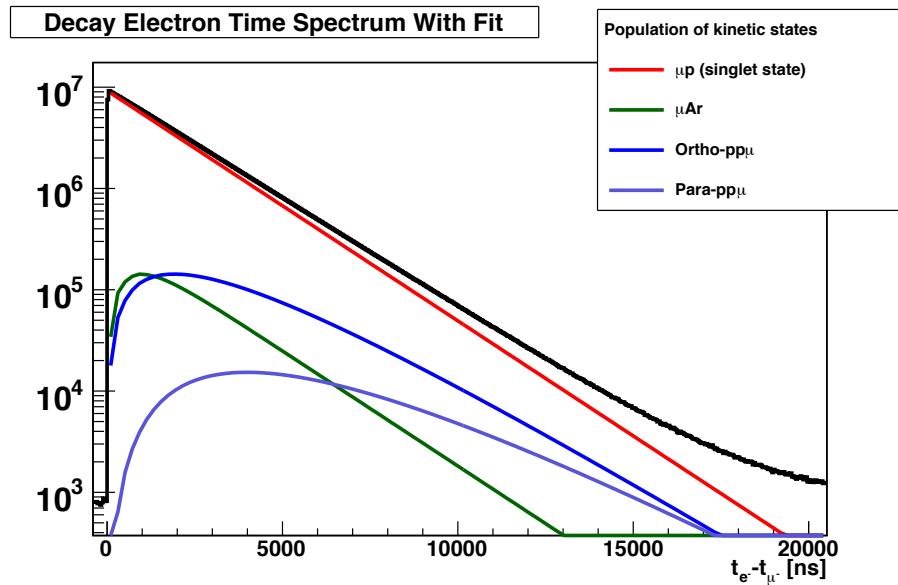$$n_{\mu p}(t=0) = 1 - f$$

$$\text{and } n_{\mu Ar}(t=0) = f,$$

$$\text{where } n_{Ortho}(t=0) = 0$$
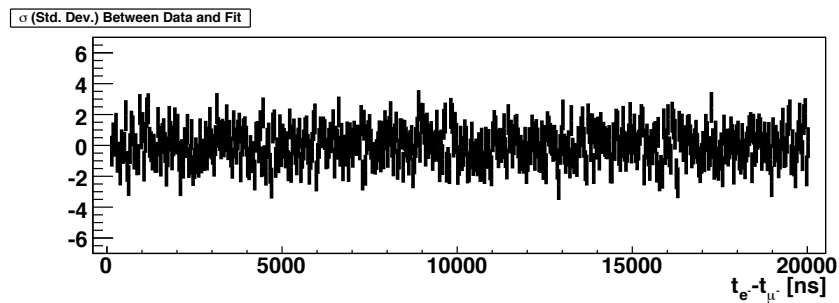
$$\text{and } n_{Para}(t=0) = 0.$$

$$n_e^{Obs.,Ar}(t) = \lambda_\mu \left(n_{\mu p}(t) + n_{Para}(t) + n_{Ortho}(t)\right) + e_{Ar}h\lambda_\mu n_{\mu Ar}(t).$$

- Differential equations, initial conditions, full time distribution.

- Atomic physics parameters f, h, and $e_{Ar}$
  - relative contribution of µAr state decays

- The hydrogen kinetic rates, $\lambda_\mu$, $\Lambda_S$, $\lambda_{op}$, $\Lambda_{pf}$, $\Lambda_O$, and $\Lambda_P$
  - Directly affect the time distribution of events

- The fit function is A $n_e(t)$+B
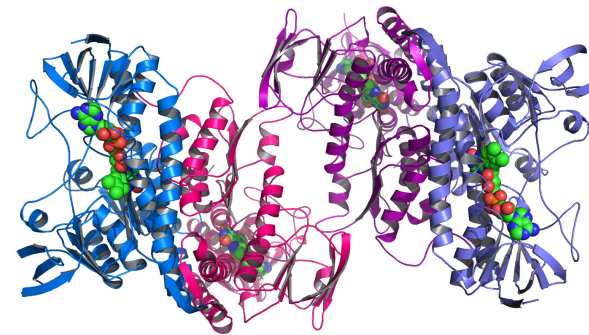
# Fit to the decay electron time distribution



- 4.25 x $10^8$ events

- Basic fit results

  - $\Lambda_{pp\mu}$=2.208(65) x $10^4$ $s^{-1}$

  - $\Lambda_{pAr}$=4.529(15) x $10^4$ $s^{-1}$
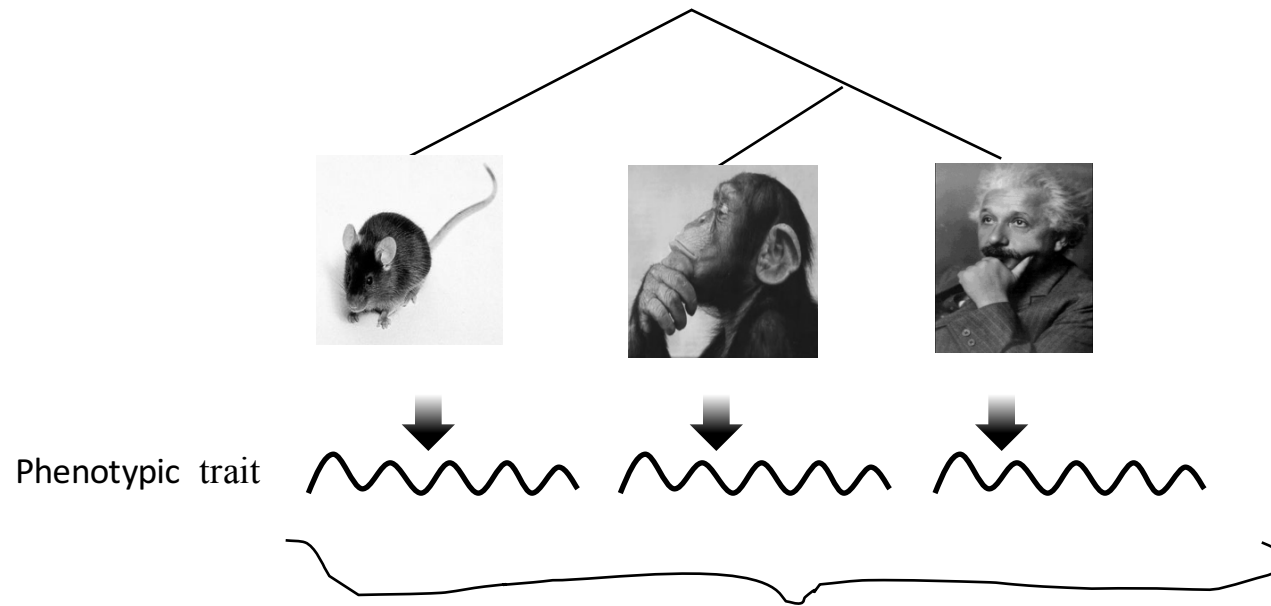
  - $\Lambda_{Ar}$=1.302(14) x $10^6$ $s^{-1}$

- $\chi2/Ndf$=0.983(64)

# Carry-overs in Computational Biology

- Analysis of complex data from a high-throughput computational environment.
    - C++ code development
    - Statistical analysis

- Integrative study of processes on multiple spatial and temporal scales.

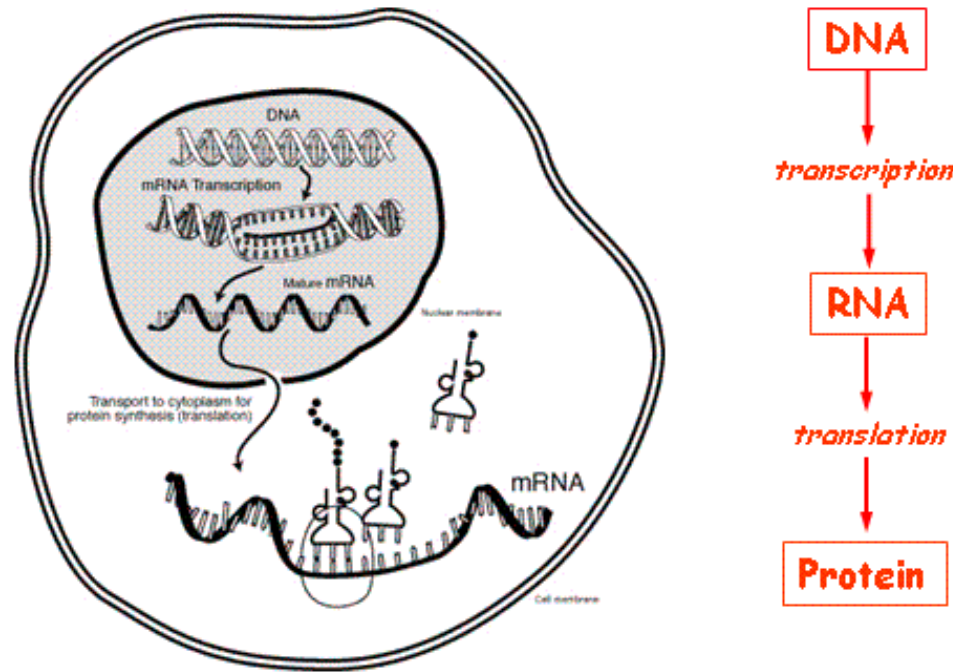- The extraction and interpretation of results from complex systems.

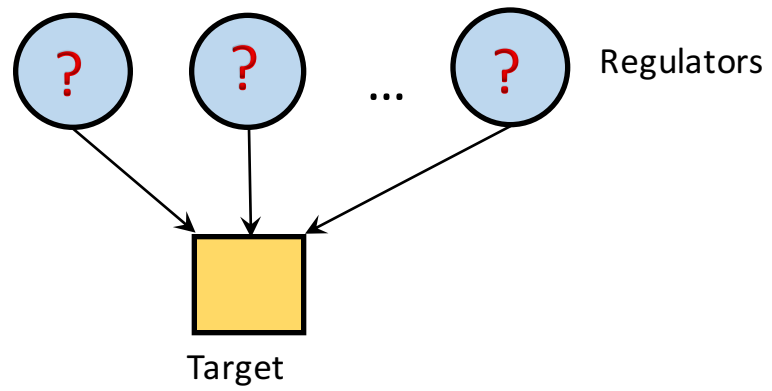# What controls phenotypic diversity?



Phenotypic trait

Changes in gene expression & regulation play a major role in diversifying phenotype.

# How are gene expressed? The central dogma.

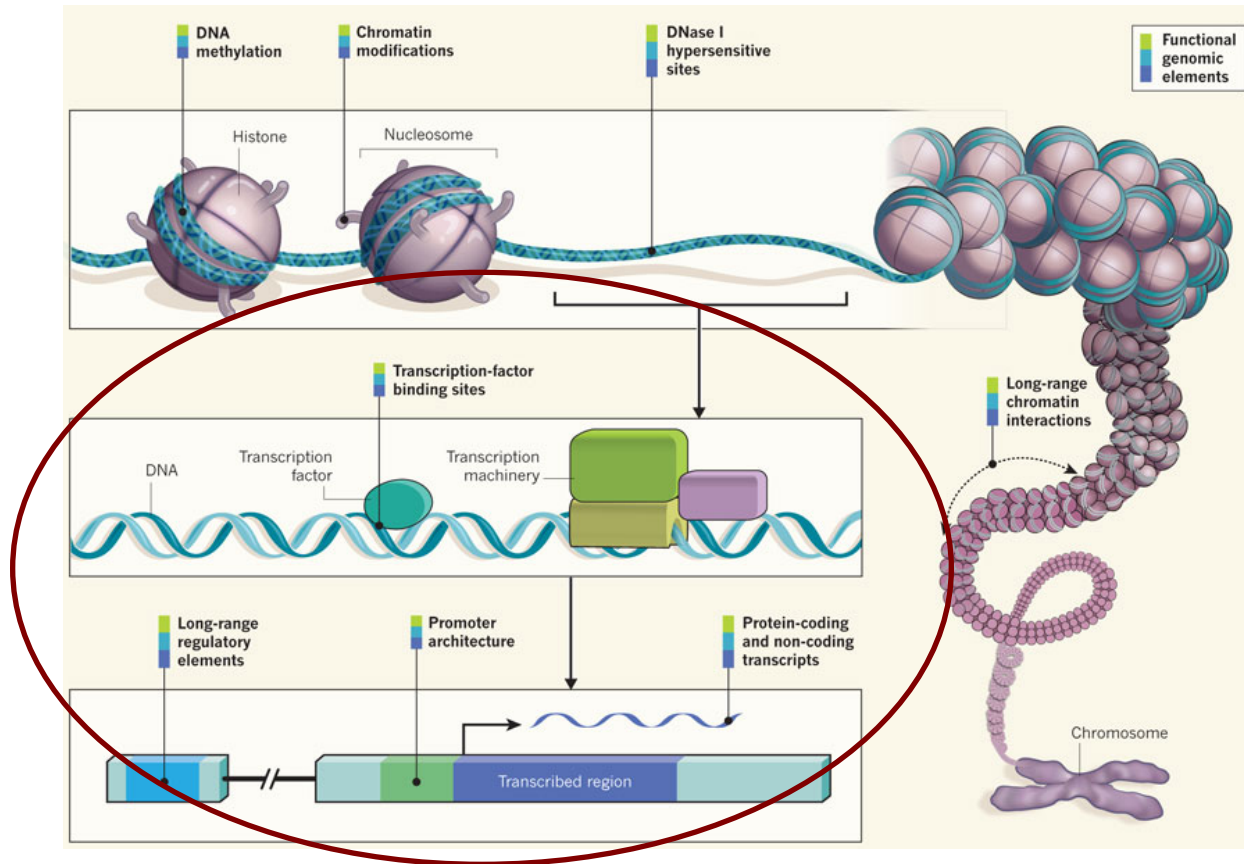# Who regulates whom?

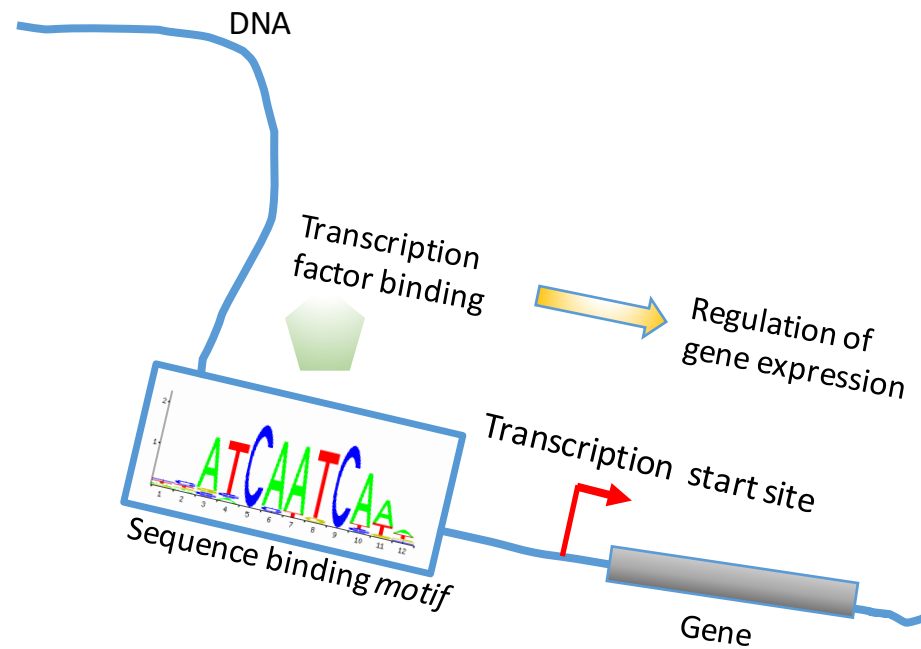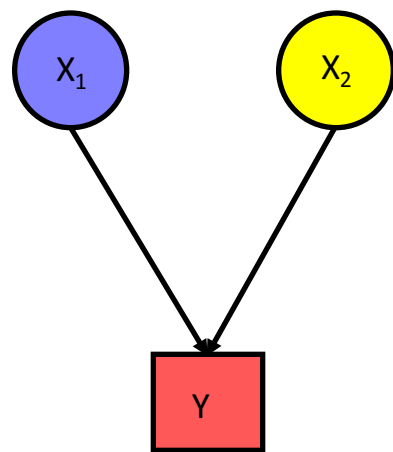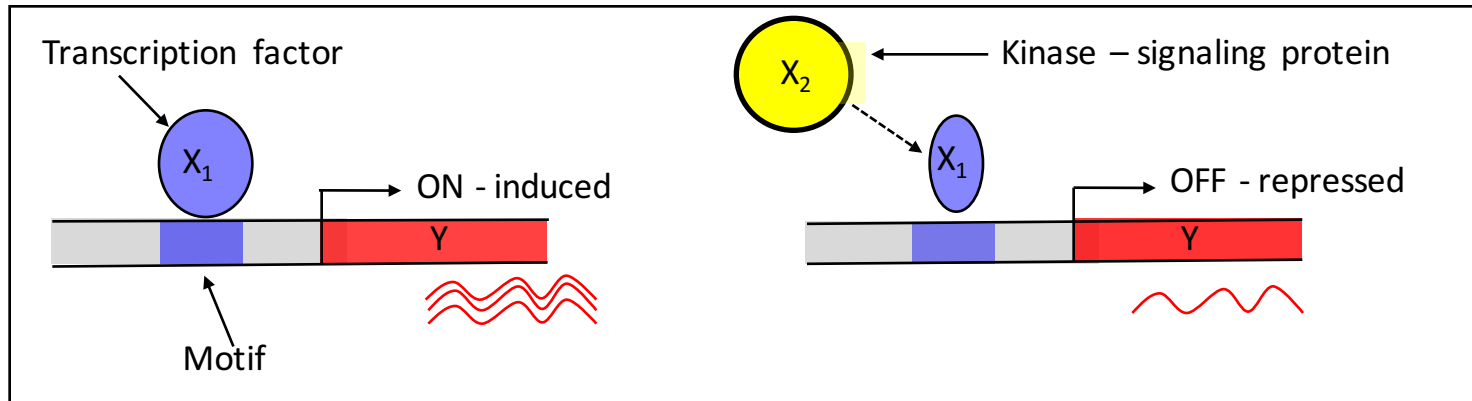# Regulation of Gene Expression is Multilayered



Image: ENCODE Consortium

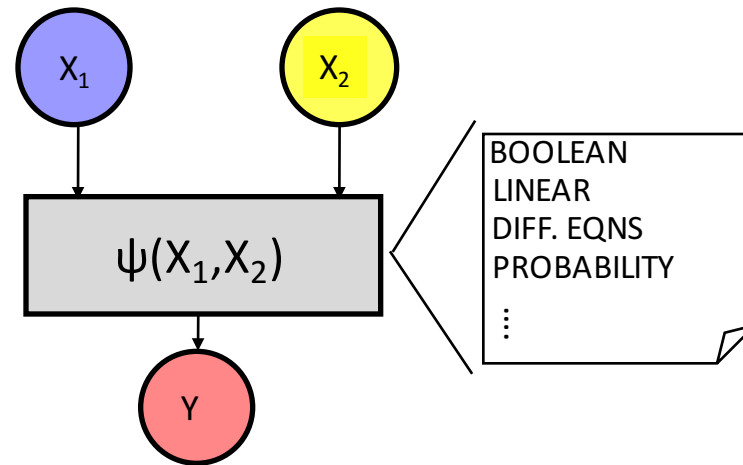# One mode of gene regulation: transcription factor binding at cis-regulatory elements in the genome.

DNA

Transcription factor binding

Regulation of gene expression

Transcription start site

ATCAATCAA

Sequence binding motif

Gene

# Modeling a regulatory network



Transcription factor

$X_1$

ON - induced

Y

Motif

Kinase – signaling protein

$X_2$

$X_1$

OFF - repressed

Y

$X_1$  $X_2$

Y

$X_1$  $X_2$

$\psi(X_1,X_2)$

BOOLEAN
LINEAR
DIFF. EQNS
PROBABILITY
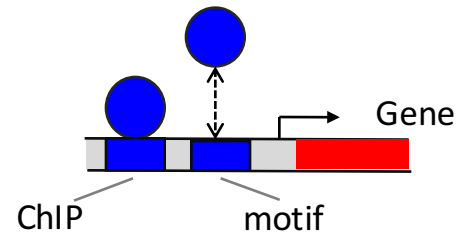⋮

Y

Structure

Who are the regulators?
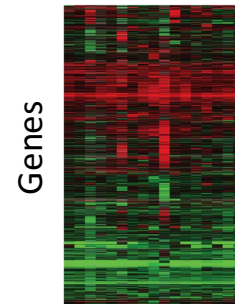
Function

How do they determine expression levels?

# Types of data for reconstructing networks

- Physical
  - ChIP-chip and ChIP-seq
  - Sequence-specific binding - motifs
  - Regulator centric
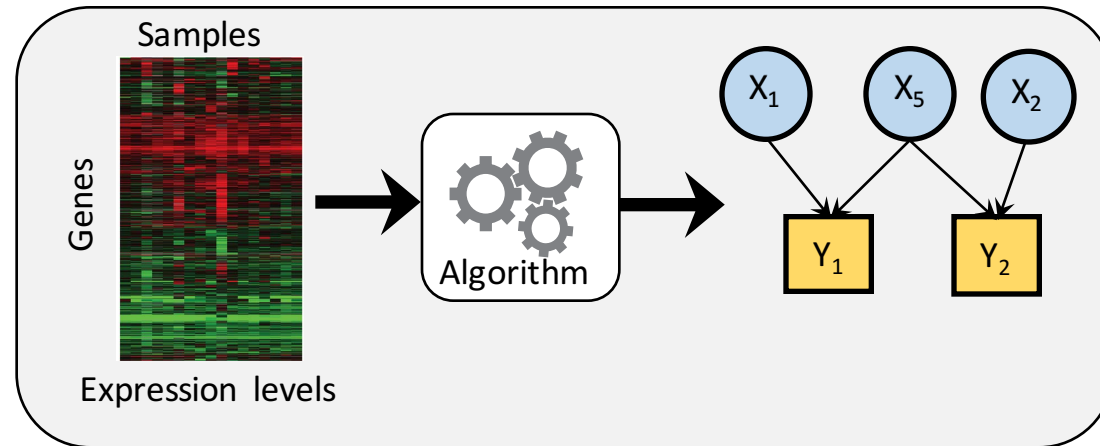


Gene

ChIP          motif

- Functional
  - Gene expression
  - Measure dynamic information
  - Can potentially recover genome-wide regulatory networks

Samples/Conditions
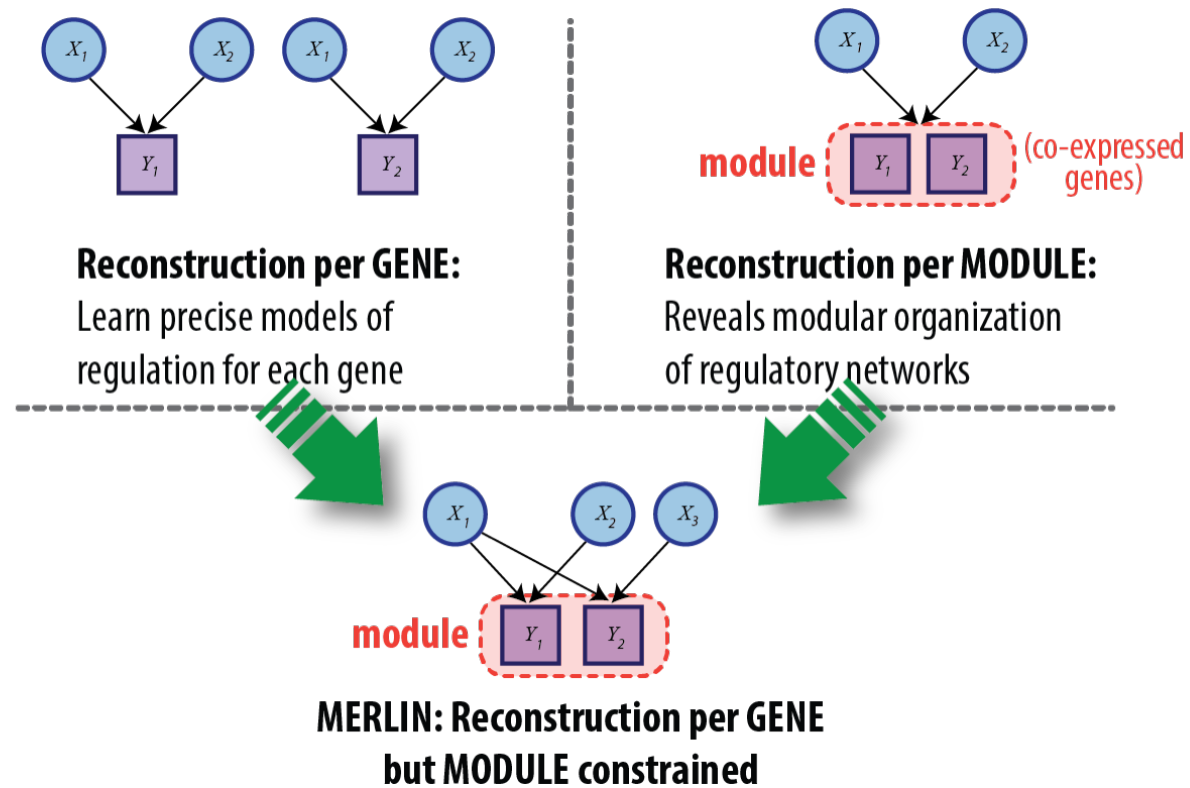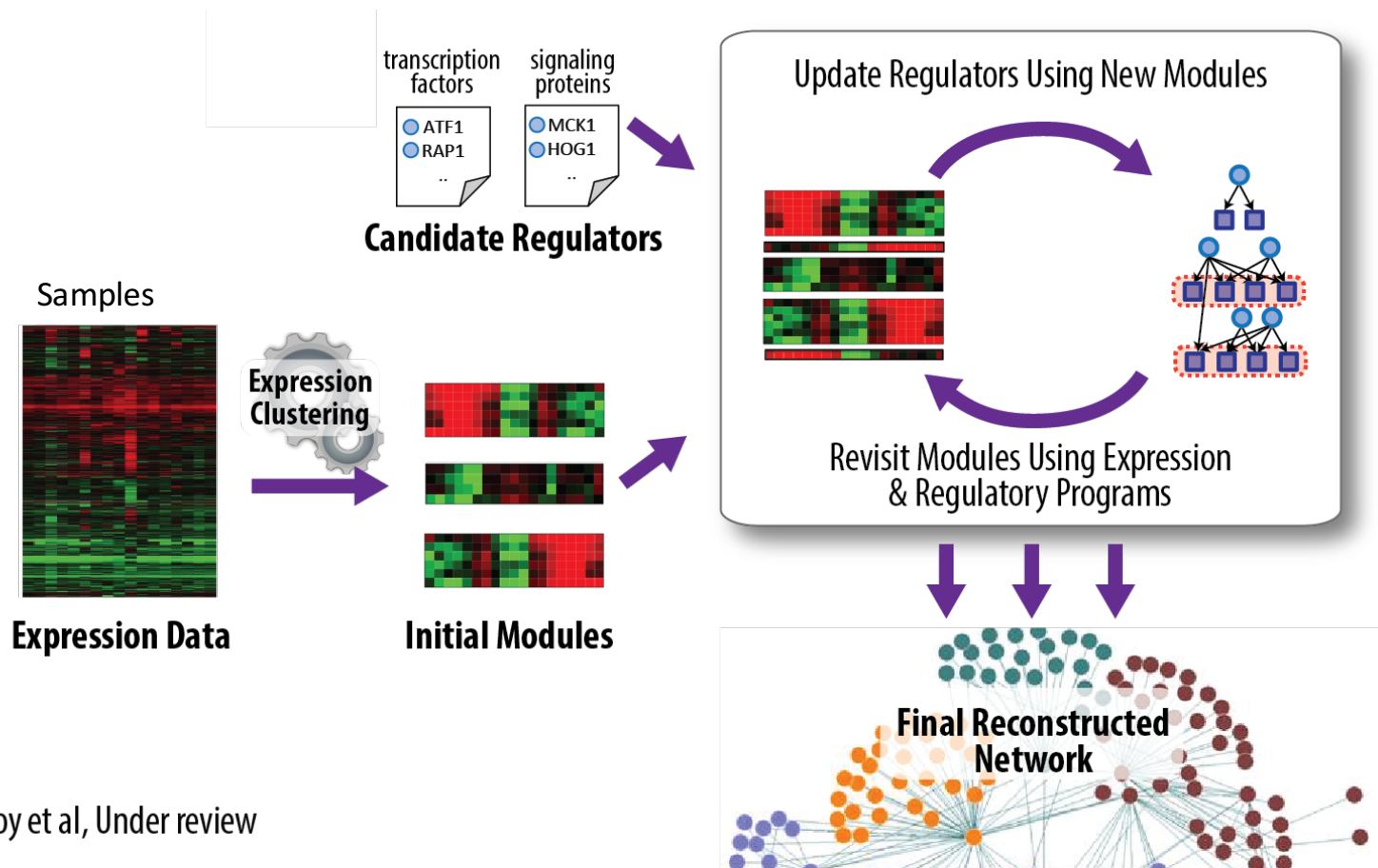


Genes

Expression levels

# Expression-based network inference

# MERLIN: A network reconstruction method to predict regulators of genes and modules



**Reconstruction per GENE:**
Learn precise models of
regulation for each gene

**Reconstruction per MODULE:**
Reveals modular organization
of regulatory networks

module (co-expressed genes)

module

**MERLIN: Reconstruction per GENE
but MODULE constrained**
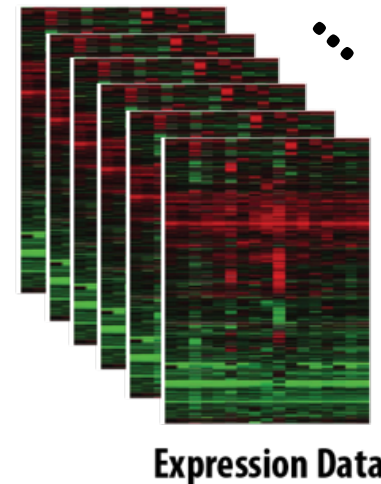
Roy et al, 2013 Plos comp bio
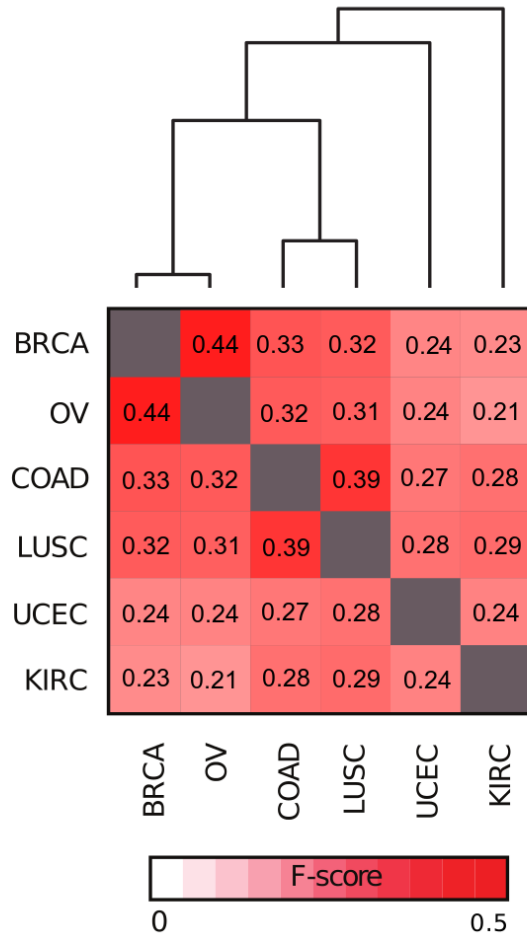
# MERLIN learning algorithm



Roy et al, Under review

# Data from The Cancer Genome Atlas

- Cancer Genome Atlas Research Network:
    - Weinstein et al. *Nat Genet*. 2013

- Microarray gene expression data for 6 cancers:
    - (1) Breast (BRCA)
    - (2) Colon (COAD)
    - (3) Kidney ma (KIRC)
    - (4) Lung (LUSC)
    - (5) Ovarian (OV)
    - (6) Uterine (UCEC).

- 54 (UCEC) to 598 (OV) patient samples

- 8499 genes were selected
    - Variation in expression across patient samples in each data set
    - Any gene annotated in curated NCI cancer pathways

- 1050 were known transcription factors TFs and kinases − regulators
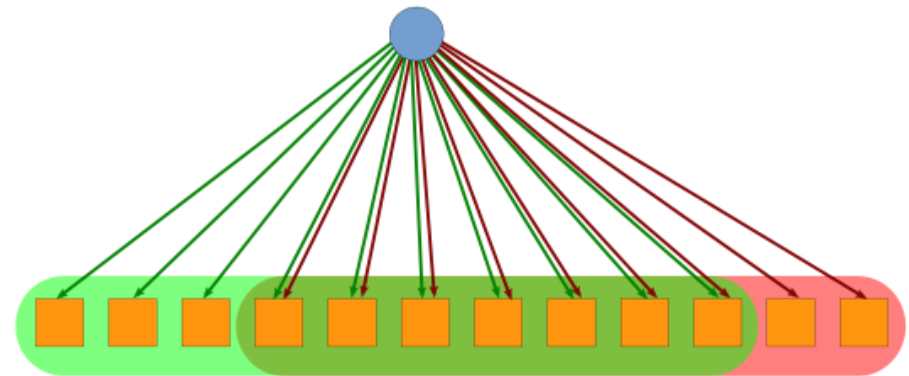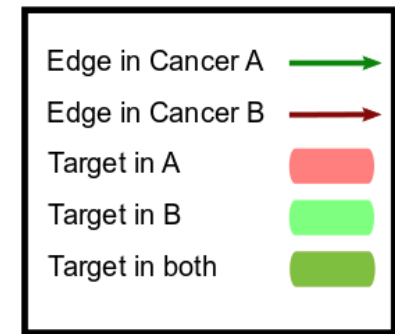


**Expression Data**
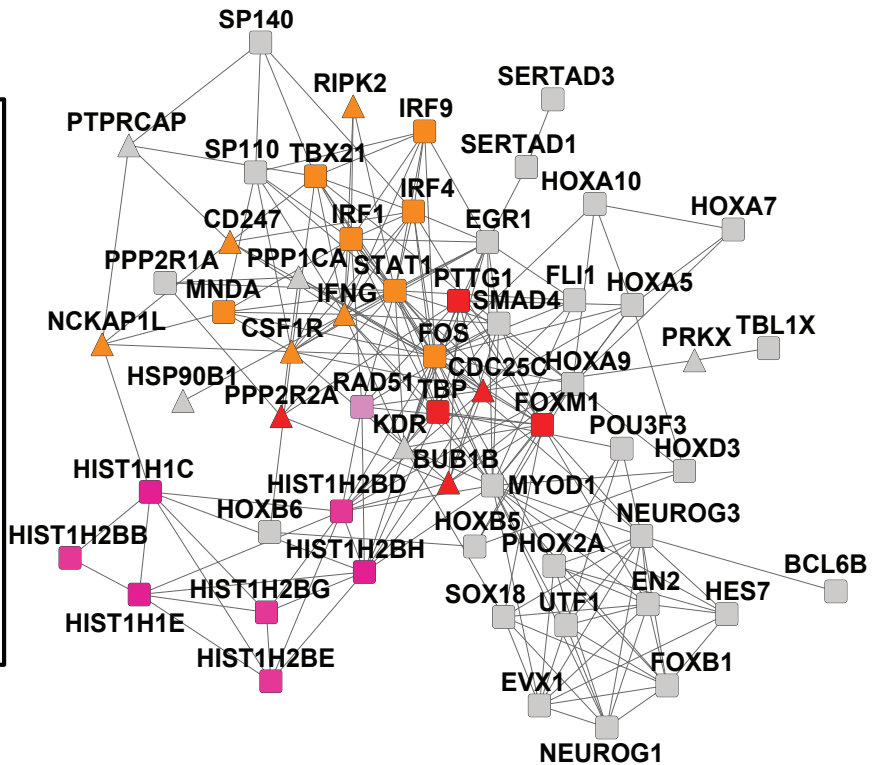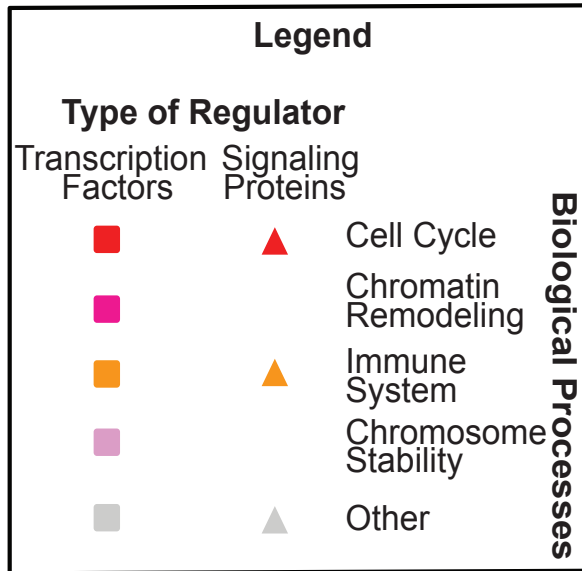
# The inferred networks are distinctly different



F-score or harmonic mean for the set of edges in the network

**How great is the overlap of edges in the network from cancer A and the network of cancer B?**

# Common regulators are associated with chromatin, cell cycle and immune response

Consists of edges between 75 regulatory proteins and 156 target genes

# Expression in consensus modules



**BRCA**
(590 Samples)
(52 Modules)

Genes (1909)

**COAD**
(181 Samples)
(55 Modules)

Genes (1863)

**KIRC**
(74 Samples)
(24 Modules)

Genes (1223)

**OV**
(598 Samples)
(44 Modules)

**LUSC**
(161 Samples)
(35 Modules)

Genes (1116)

Genes (1481)

**UCEC**
(54 Samples)
(9 Modules)

Genes (597)

Row-zero-mean expresssion values

-3    0    3

# Validating that our modules are biologically coherent



Here we have several sets of annotated genes, and each set provides us lists of gens with a certain biological significance.

We look to find if our modules are significantly enriched in genes from any of these annotated sets using a Hypergeometric test

Above we count the fraction of modules that have an enrichment with 0.05 significance in the results from each data set.

# An example of a module
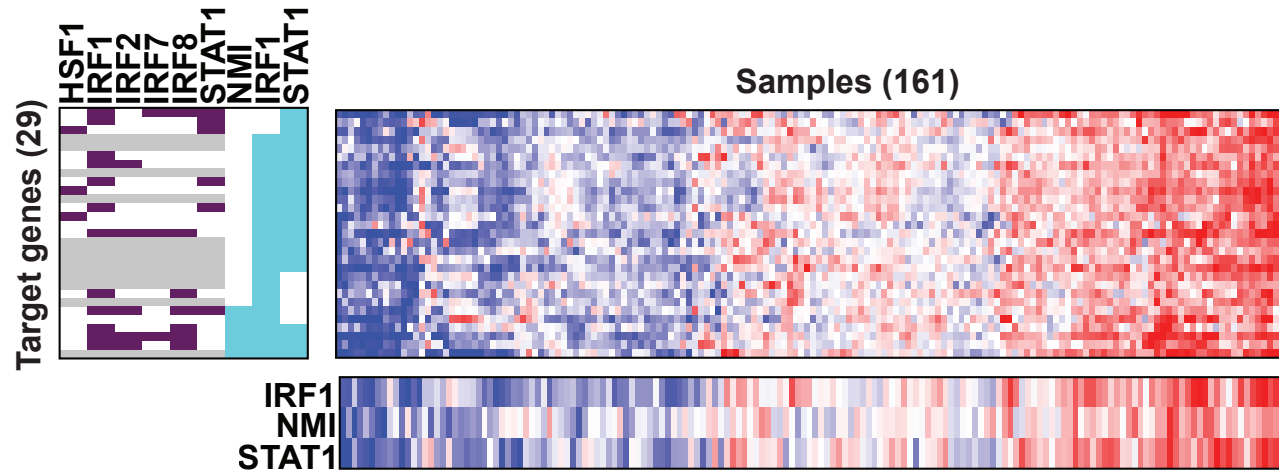


LUSC (module 14)

Samples (161)

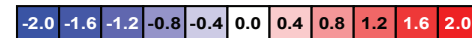Target genes (29)

HSF1 IRF1 IRF2 IRF7 IRF8 STAT1 NMI IRF1 STAT1

IRF1
NMI
STAT1

■ Targets of regulators predicted by MERLIN

■ Targets of regulators predicted by MERLIN
■ Targets of regulators f MSigDB motifs
■ Not annotated

**Row-zero-meaned expression values**

| -2.0 | -1.6 | -1.2 | -0.8 | -0.4 | 0.0 | 0.4 | 0.8 | 1.2 | 1.6 | 2.0 |

# Enrichments for motifs of immune system regulators

- interferon regulatory factor (IRF) family – all cancers
- regulatory factor X (RFX) – five cancers
- signal transducer and activator of transcription (STAT1) – five cancers

- All regulators of the immune system

# Immune system function is over-represented

# Is the immune system induced or repressed?
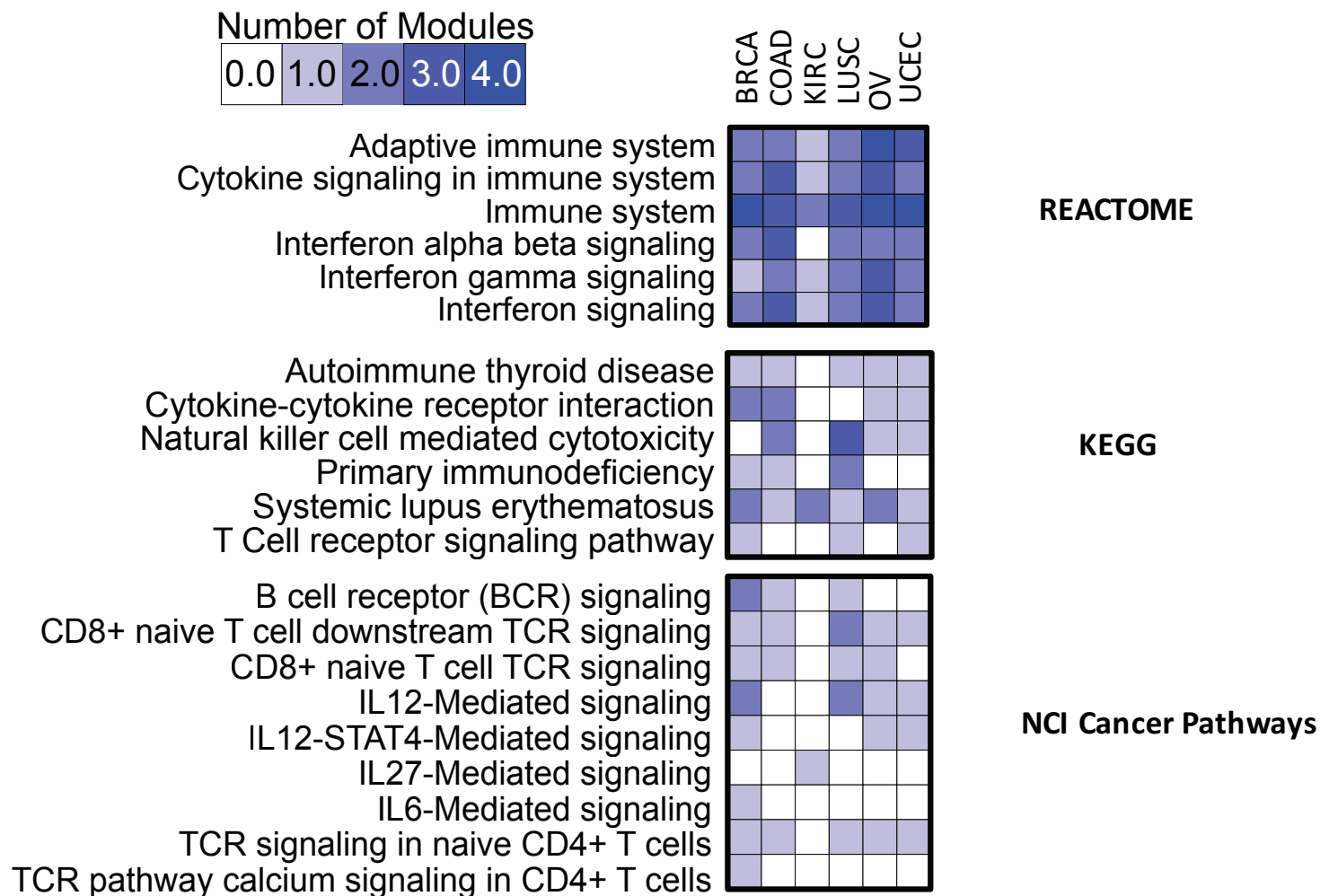
Modules associated with the immune system:

Samples

Genes

Consistent with observations of immune system activation in anti-cancer therapy
- Apetoh et al. Nat. Med. 2007

Per-gene, per-sample expression values

- Known cross-talk mechanism of activation between STAT3 and interleukin 6 signaling Lee et al. Nat. Med. 2010

Modules not associated with the immune system:

Samples

Genes



| Average | 1.16 | 0.05 | 1.02 | 0.06 | 1.63 | 0.06 | 1.6 | 0.52 | 1.08 | -0.18 | 0.46 | 0.12 |

BRCA In, BRCA Out, COAD In, COAD Out, KIRC In, KIRC Out, LUSC In, LUSC Out, OV In, OV Out, UCEC In, UCEC Out

Per-sample Expression Values

# Summary

- We have introduced stability-selection into our MERLIN-based approach to infer regulatory networks across different conditions.

- Our approach builds on the idea that both module- and network-based characterization of transcriptional programs are important.

- Our methods can be extended with additional data types.

Knaack SA, Siahpirani AF, Roy S. A pan-cancer modular regulatory network analysis to identify common and cancer-specific network components. *Cancer Inform*. 2014 Oct 28;13(Suppl. 5):69-84.

doi: 10.4137/CIN.S14058
PMID: 25374456 [PubMed]

Cancer Informatics

## Conclusions

Work with beautiful complex systems, rich for exploration and discovery.

Computational biology is a fast paced field, with many emerging technologies and methods.

The mindset towards measurements is different than what you are used to from physics.

It's a field that will have an increasing impact on medicine and human health as we learn more.

# Acknowledgements

- My mentor,  Sushmita Roy

- Members of the Roy Group

- Funding support from

  - National Science Foundation, CAREER grant (SR)

  - National Library of Medicine training grant NLM5T15LM007359 (SK)

- You!

Thanks to Sushmita Roy and Alireza Siahpirani for contributing slides on gene regulatory networks.