*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*
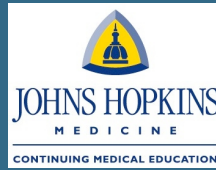
*April 27, 2016*
*John Quackenbush, Ph.D.*

# EXPRESSION ANALYSIS, FUNCTIONAL ENRICHMENT, AND NETWORK INFERENCE

**John Quackenbush**
**Dana-Farber Cancer Institute**
**Harvard T.H. Chan School of Public Health**

## Background and Disclosures

- **Professor of Biostatistics and Computational Biology, Dana-Farber Cancer Institute**
- **Professor of Computational Biology and Bioinformatics, Harvard School of Public Health**
- **Many other academic titles**
- **Numerous advisory boards**
- **Co-Founder of Genospace, a Precision Genomic Medicine Software Company**

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

*Current Topics in Genome Analysis 2016*

*John Quackenbush*

*Genospace, LLC*
*Co-Founder and Board Chair*

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

**Every revolution in science — from Copernican heliocentric model to the rise of statistical and quantum mechanics, from Darwin's theory of evolution and natural selection to the theory of the gene — has been driven by one and only one thing: access to data.**
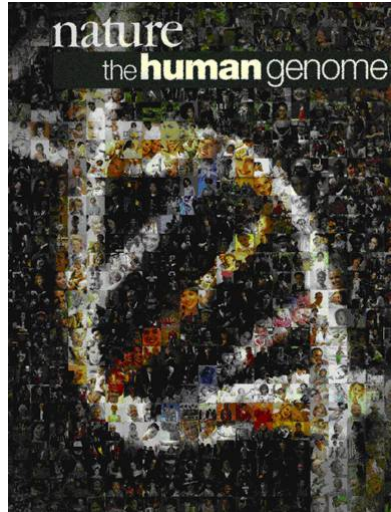
**–John Quackenbush**

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

**@johnquackenbush**-**Every revolution in the history science has been driven by one and only one thing: access to data.**

Twitter version, 115 characters with spaces

# A Brief History of Expression Analysis

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

## February 2001: Completion of the Draft Human Genome



**Public HGP**

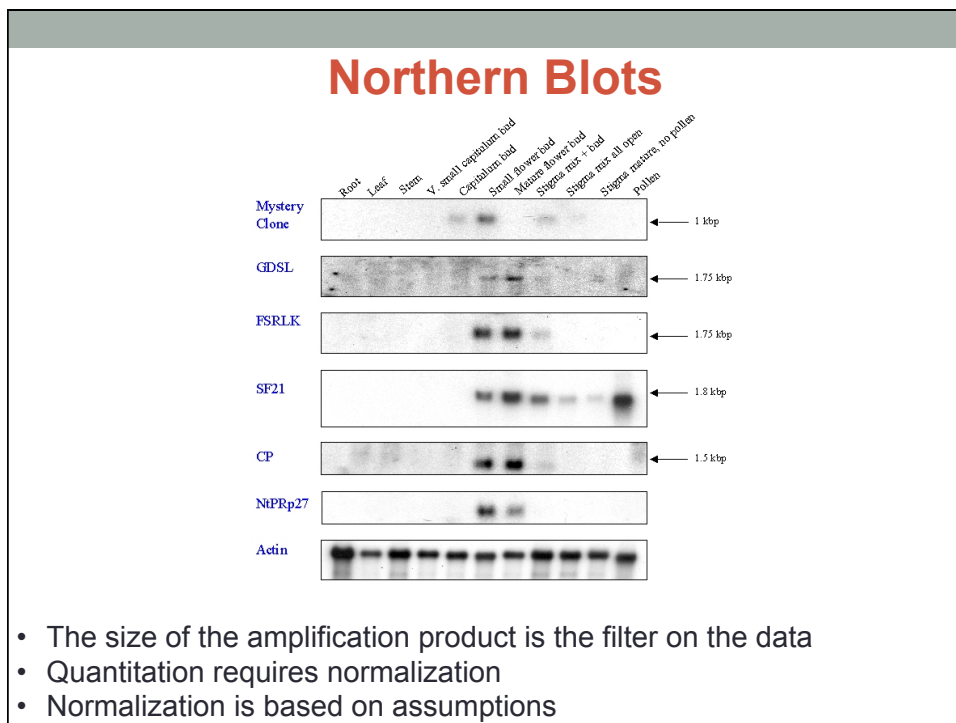
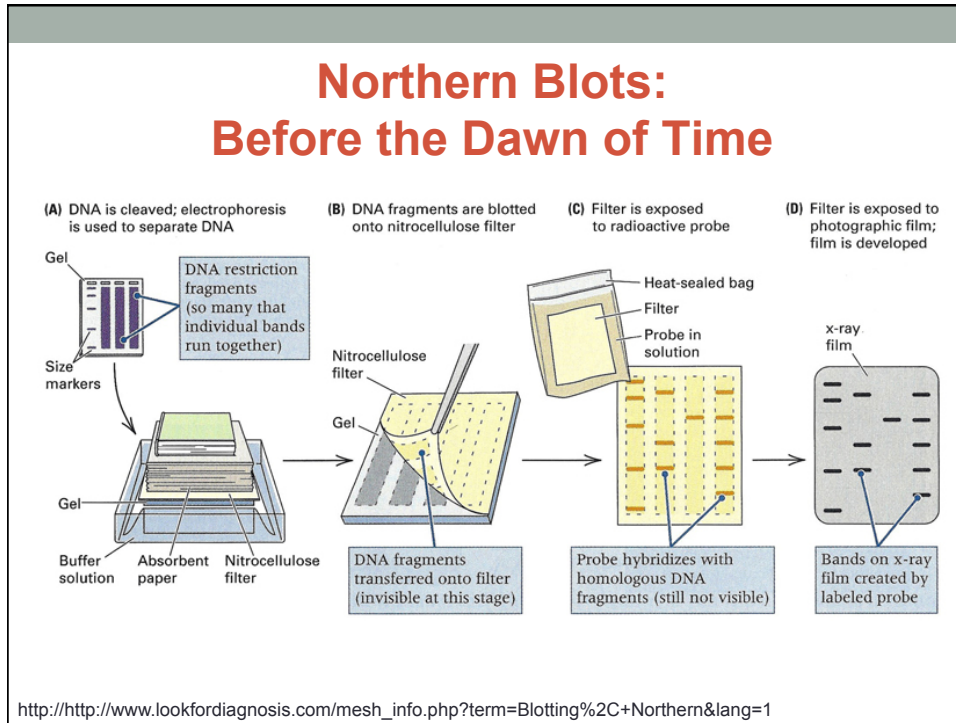
**Celera Genomics**

## Molecular Biology in 7 Words

Gene ——————→ Protein

RNA

Regulation

Folding

Function ←—————— Structure

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

# The Genome Project has provided a "parts list" for a human cell



# Different cell types express different sets of genes



- Neuron
- Thyroid Cell
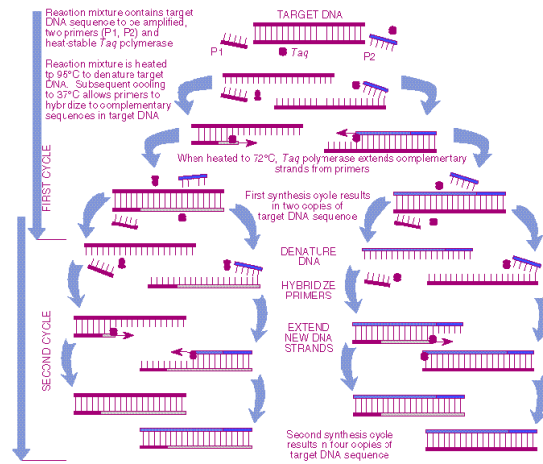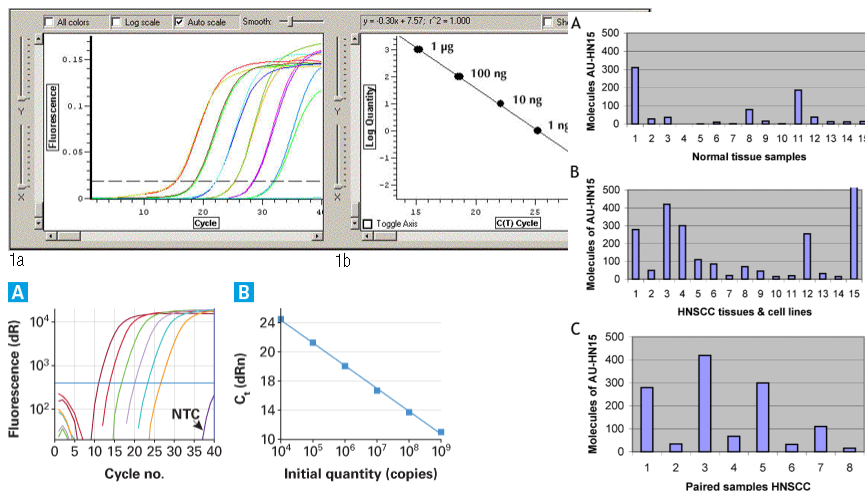- Lung Cell
- Cardiac Muscle
- Pancreatic Cell
- Kidney Cell
- Skeletal Muscle
- Skin Cell

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

# Northern Blots: Before the Dawn of Time



http://http://www.lookfordiagnosis.com/mesh_info.php?term=Blotting%2C+Northern&lang=1

# Northern Blots



- The size of the amplification product is the filter on the data
- Quantitation requires normalization
- Normalization is based on assumptions

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

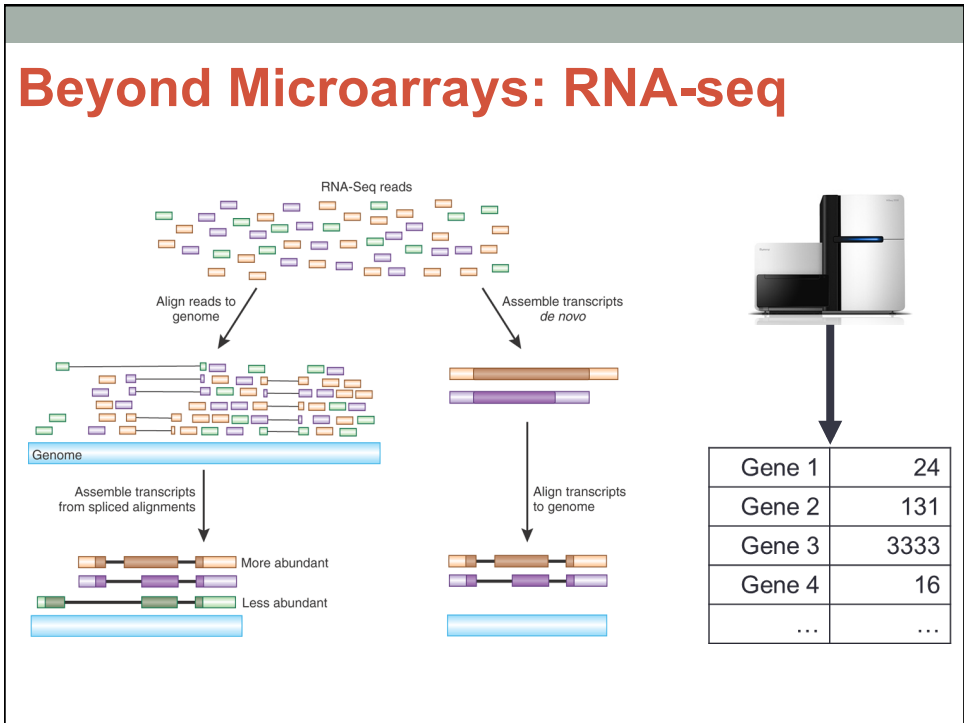# Quantitative RT-PCR: The Ancient World



# Quantitative PCR



- Paired hybridization of two primers is the filter on the data
- Quantitation requires normalization (comparison to standard curves)
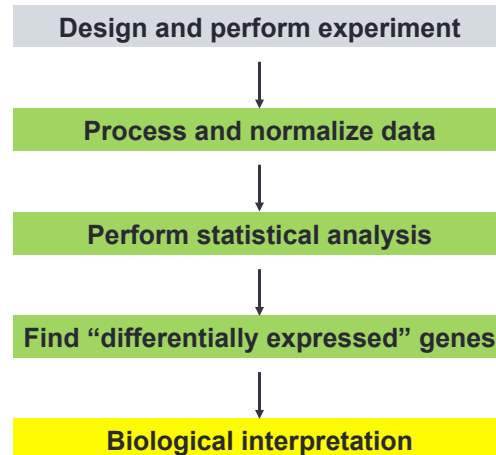- Normalization is based on assumptions

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

# Beyond qRT-PCR: Microarrays



**Spatial position is the filter on the data.**

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

# Beyond Microarrays: RNA-seq



# Disease Progression and Precision Care



Adapted from a slide by Peter van der Spek, Erasmus University

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

# Experimental Overview

## Expression Analysis Pipeline: Microarrays

Design and perform experiment

↓

Process and normalize data

↓

Perform statistical analysis

↓

Find "differentially expressed" genes

↓

Biological interpretation

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

# Design the Experiment

## Why Design an Experiment?

- The goal of an experiment dictates everything from how the samples are collected to how the data are generated

- The design of the analytical protocol should be reflected in the design
    - Do we have enough replicates?
    - Do we have sufficient controls?
    - Do we collect samples and data to avoid confounding and batch effects?

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

## Basis of Experimental Design

- In biology, "traditional" approaches to inquiry involved hypothesis testing.
  - We identify a problem and postulate a mechanism
  - We design an experiment in which we perturb the system and then look for changes
  - The response of the system either validates or invalidates our hypothesis

- In these types of experiments, we attempt to tightly control the variables so as to carefully measure the influence of these, perturbing a single parameter at a time

- Good experimental design requires sufficient replication to estimate the effects we wish to measure

## Basis of Experimental Design

- Functional genomics technologies have dramatically changed the way in which we approach biological questions
  - We can now survey the responses of thousands of genes, proteins, or metabolites in a particular system and look for patterns of expression
  - These "hypothesis generating" experiments do not (necessarily) require a mechanistic hypothesis ahead of time
  - However, this does not mean we do not have to carefully design our experiment and analyze the data

- Here, we attempt to control the variables so as to carefully measure the influence of these, perturbing a single parameter at a time

- Good experimental design requires sufficient replication to estimate the effects we wish to measure

13

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

## Types of Experiments

- Class Comparison
  - Can I find genes that distinguish between two classes, such as tumor and normal?

- Class Discovery
  - Given what I think is a uniform group of samples, can I find subsets that are biologically meaningful?

- Classification
  - Given a set of samples in different classes, can I assign a new, unknown sample to one of the classes?

- Large-scale Functional Studies
  - Can I discover a causative mechanism associated with the distinction between classes?

These are often not completely distinct and a single dataset can often be used for multiple purposes

# Normalization

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

## Why Normalize Data?

- The goal of normalization is to remove systematic variation from the data and scale it so that comparisons can be made across studies

## RMA Background correction

- Expression= Background ($N(0,\sigma^2)$) + Signal ($Exp(\alpha)$)

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

# RMA Normalization

- Force the empirical distribution of probe intensities to be the same for every chip in an experiment
- The common distribution is obtained by averaging each *quantile* across chips:

  *Quantile Normalization*

# One distribution for all arrays: the black curve



Density of PM probe intensities for Spike-In chips

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

# RMA: Probe set summary

- Robustly fit a two-way model yielding an estimate of $\log_2$(signal) for each probe set
- Fit may be by
  - median polish (quick) or by
  - Mestimation (slower but yields standard errors and good quality
- RMA reduces variability without loosing the ability to detect differential expression

# RMA: Before and After



http://www.slideshare.net/wjjessen/covance-talk

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

# Ratio-Intensity: Before



# Ratio-Intensity: After

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

# Normalization

- There are many, many methods
- All attempt to do the same thing, but all have their own assumptions that may or may not be violated
- RMA is widely accepted as the standard for microarrays
- There is less consensus on what works best for RNA-seq
- We constantly have to test our assumptions, even with normalization

# GTEx: Complex data requires complex methods



**Raw expression**

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

# Clustering: Finding Patterns

## Hierarchical Clustering

1. Calculate the distance between all genes. Find the smallest distance. If several pairs share the same similarity, use a predetermined rule to decide between alternatives.

2. Fuse the two selected clusters to produce a new cluster that now contains at least two objects. Calculate the distance between the new cluster and all other clusters.

3. Repeat steps 1 and 2 until only a single cluster remains.

4. Draw a tree representing the results.

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

**Hierarchical Clustering**

g1 is most like g8

g4 is most like {g1, g8}

(HCL2)



**Hierarchical Clustering**

g5 is most like g7

{g5,g7} is most like {g1, g4, g8}

(HCL3)

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

**Hierarchical Clustering**

(HCL4)

**Agglomerative Linkage Methods**

Linkage methods are rules or metrics that return a value that can be used to determine which elements (clusters) should be linked.

Three linkage methods that are commonly used are:

- Single Linkage
- Average Linkage
- Complete Linkage

(HCL6)

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

# Single Linkage

Cluster-to-cluster distance is defined as the *minimum distance* between members of one cluster and members of the another cluster. Single linkage tends to create 'elongated' clusters with individual genes chained onto clusters.

$$D_{AB} = \min ( d(u_i, v_j) )$$

where $u \in A$ and $v \in B$
for all $i = 1$ to $N_A$ and $j = 1$ to $N_B$



**(HCL7)**

# Average Linkage

Cluster-to-cluster distance is defined as the *average distance* between all members of one cluster and all members of another cluster. Average linkage has a slight tendency to produce clusters of similar variance.

$$D_{AB} = 1/(N_A N_B) \, \Sigma \, \Sigma \, ( d(u_i, v_j) )$$

where $u \in A$ and $v \in B$
for all $i = 1$ to $N_A$ and $j = 1$ to $N_B$



**(HCL8)**

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

# Complete Linkage

Cluster-to-cluster distance is defined as the *maximum distance* between members of one cluster and members of the another cluster. Complete linkage tends to create clusters of similar size and variability.

$$D_{AB} = \max ( d(u_i, v_j) )$$

where $u \in A$ and $v \in B$
for all $i = 1$ to $N_A$ and $j = 1$ to $N_B$



$D_{AB}$

(HCL9)

# Comparison of Linkage Methods



Single          Average          Complete

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

49

# *K*-means/*K*-medians Clustering (KMC)

1. Specify number of clusters, e.g., 5.

2. Randomly assign genes to clusters.

50

# KMC, continued

3. Calculate mean / median expression profile of each cluster.

4. Select a gene and move it to the cluster having the closest mean profile.

5. If the gene is shifted to a new cluster, recalculate means for the winning and losing clusters.

6. Repeat steps 4 and 5 until genes cannot be shuffled around any more, OR a userspecified number of iterations has been reached.

*k*means is most useful when the user has an *a priori* hypothesis about the number of clusters the genes should belong to.

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

# Finding Differentially Expressed Genes

# Lies, Damn Lies, and Statistics

## Finding Significant Genes

*t*-test for each gene

- Tests whether the difference between the mean of the query and reference groups are the same
- Essentially measures signal-to-noise
- Calculate *p*-value (permutations or distributions)
- May suffer from intensity-dependent effects

*t* = signal = difference between means = <Xq> – <Xc>
   noise        variability of groups      SE(XqXc)

$$t = \frac{\langle Xq \rangle - \langle Xc \rangle}{\sqrt{\dfrac{\sigma_q^2}{n_q} + \dfrac{\sigma_c^2}{n_c}}}$$

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

## *t*-tests

A significant difference

Probably not

## *limma*: The standard for microarray analysis

**Schematic of the major components that are central to any limma analysis.**

Matrix of expression values
(from RNA-seq / microarray)

| Gene ID | LSK_1 | LSK_2 | CMP_1 | CMP_2 |
|---------|-------|-------|-------|-------|
| 11303 | 478 | 619 | 4830 | 7165 |
| 11305 | 27 | 20 | 48 | 55 |
| 11306 | 132 | 200 | 560 | 408 |
| 11307 | 42 | 60 | 131 | 99 |
| .... | | ... tens of thousands more | ... | |

Gene-wise linear models
$$E(y_g) = X\beta_g$$
$$\mathrm{var}(y_{gj}) = \sigma_g^2 / w_{gj}$$

**Advanced statistical algorithms in *limma* that allow...**

- *Information Borrowing*
- *Quantitative Weighting*
- *Variance Modelling*
- *Data Pre-processing*

*limma* delivers powerful inference for differential expression analysis

$$\hat{\beta}_g, s_g^2 *$$

Estimated gene-specific parameters used for gene prioritization and gene set testing

Matthew E. Ritchie et al. Nucl. Acids Res. 2015;nar.gkv007

**Nucleic Acids Research**

27

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

# Biological Interpretation

# What do the genes in this list do?

# Tell me a story, Grampa

## Biological Interpretation

- An obvious way to gain biological insight is to assess the differentially expressed genes in terms of their known function(s)

- Requires an automated and objective (statistical) approach

- Functional profiling or pathway analysis

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

# Early functional analyses

- Manually annotate list of differentially expressed (DE) genes

- Extremely time-consuming, not systematic, user-dependent

- Group together genes with similar function

- Conclude functional categories with most DE genes important in disease/condition under study

- BUT… it may not be the right conclusion

- This is what we call "Biopoetry."

# GO and functional analysis



| Functional category | Number of sig genes |
|---------------------|---------------------|
| Immune response | 40 |
| Metabolism | 20 |
| Transcription | 20 |
| Energy production | 10 |
| Neurotransmission | 5 |
| Protein transport | 5 |
| **TOTAL** | **100** |

Immune response category contains 40% of all significant genes - by far the largest category.

Reasonable to conclude that immune response may be important in the condition being studied?

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

# However …

- What if 40% of the genes on the array were involved in immune response?

- Only detected as many significant immune response genes as expected by chance

- Need to consider not only the number of significant genes for each category, but also total number on the array

# Same example, relative to background

| Functional category | Number of genes on array | Observed number of significant genes | Expected number of significant genes |
|---|---|---|---|
| Immune response | 8000 | 40 | 40 |
| Metabolism | 4000 | 20 | 20 |
| Transcription | 2000 | 10 | 10 |
| Energy production | 4000 | 30 | 20 |
| Neurotransmission | 200 | 5 | 1 |
| Protein transport | 1800 | 5 | 9 |
| | | | |
| ALL | 20000 | 100 | |

Expected number of significant genes for category X is
(num sig genes ÷ total genes on array)*(num genes in category X on array)

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

# Same example, relative to background

| Functional category | Number of genes on array | Observed number of significant genes | Expected number of significant genes |
|---|---|---|---|
| Immune response | 8000 | 40 | 40 |
| Metabolism | 4000 | 20 | 20 |
| Transcription | 2000 | 10 | 10 |
| Energy production | 4000 | 30 | 20 |
| Neurotransmission | 200 | 5 | 1 |
| Protein transport | 1800 | 5 | 9 |
|  |  |  |  |
| ALL | 20000 | 100 |  |

- Now, energy production and neurotransmission categories appear more interesting as many more significant genes were observed than expected by chance

- Largest categories are not necessarily the most interesting!



https://david.ncifcrf.gov/

31

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

http://software.broadinstitute.org/gsea/index.jsp

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

# KEGG pathway database



# WikiPathways

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

# Biological Networks

# Can we make this more complicated?

# How *NOT* to do Network Analysis



**Conditions**

Genes

Expression data
(**Phenotype I**)

**Statistical Comparison**

**Conditions**

Genes

Expression data
(**Phenotype II**)

**Differentially Expressed Genes**

**Add Protein-Protein Interaction Network**

**Color and Start Bio-Poetry**

- Should things that are differentially expressed be connected?
- Is the PPI network even "relevant"?

35

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

# How *NOT* to do Network Analysis (2)



**Conditions**

**Genes**

**Expression data (Phenotype I)**

**Statistical Analysis**

"r²"

"t-test"

**Conditions**

**Genes**

**Expression data (Phenotype II)**

**Correlation Network**

**Differentially Expressed Genes**

**Color and Start Bio-Poetry**

- Are things that are correlated functionally connected?
- Are correlations the same in different phenotypes?

# How *we* do Network Analysis



**Conditions**

**Genes**

**Expression data (Phenotype I)**

**Infer phenotype-specific network**

**Conditions**

**Genes**

**Expression data (Phenotype II)**

**Infer phenotype-specific network**

**Analyze Network Topology and Structure**

**Compare Network Topologies**

**Simultaneously compare differential structure and expression**

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

## Starting Assumptions

- **There is no single "right" network**

- **The structure of the network matters and network structure often changes between states.**

- **We have to move from asking "Is the network right?" to asking "Is the network useful?"**

- **The real question is "Does a network model inform our understanding of biology?"**

# Modeling Gene Regulatory Networks

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

# Integrative Network Inference: PANDA

© PLOS | ONE

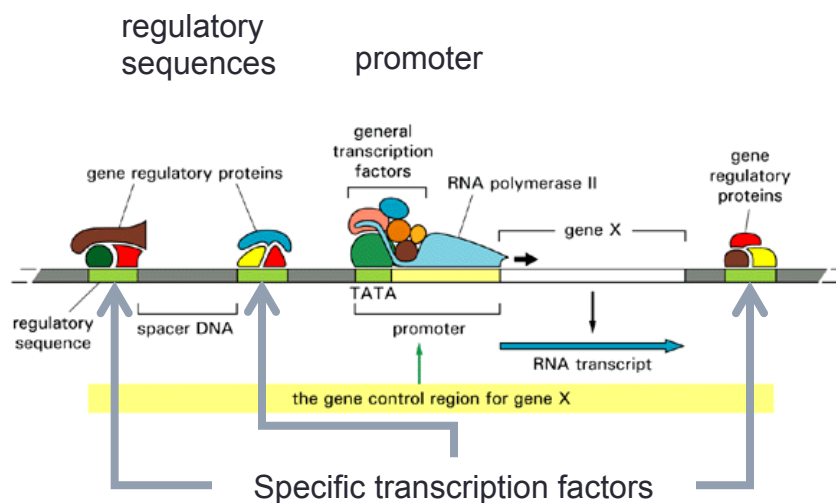## Passing Messages between Biological Networks to Refine Predicted Interactions

Kimberly Glass[1,2], Curtis Huttenhower[2], John Quackenbush[1,2], Guo-Cheng Yuan[1,2]*
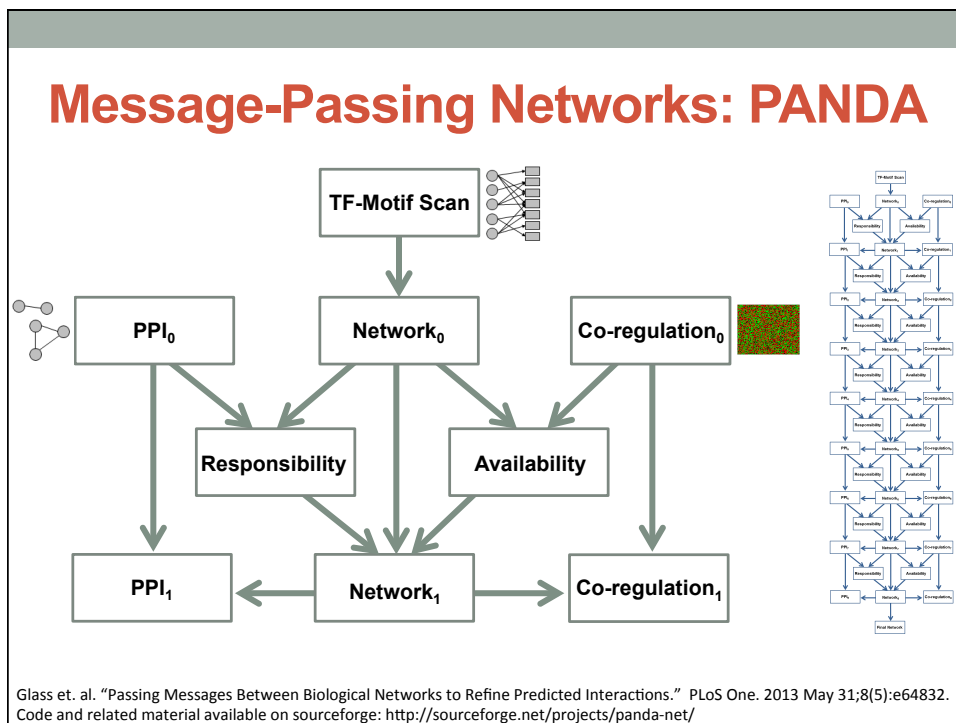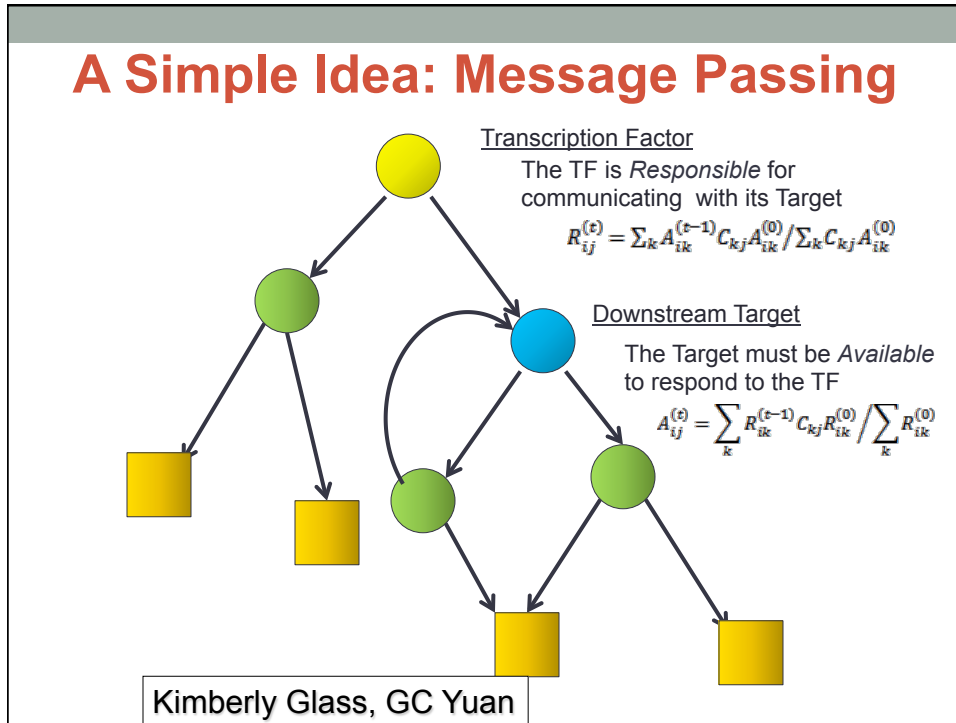
1 Department of Biostatistics and Computational Biology, Dana-Farber Cancer Institute, Boston, Massachusetts, United States of America, 2 Department of Biostatistics, Harvard School of Public Health, Boston, Massachusetts, United States of America

**Abstract**

Regulatory network reconstruction is a fundamental problem in computational biology. There are significant limitations to such reconstruction using individual datasets, and increasingly people attempt to construct networks using multiple, independent datasets obtained from complementary sources, but methods for this integration are lacking. We developed PANDA (Passing Attributes between Networks for Data Assimilation), a message-passing model using multiple sources of information to predict regulatory relationships, and used it to integrate protein-protein interaction, gene expression, and sequence motif data to reconstruct genome-wide, condition-specific regulatory networks in yeast as a model. The resulting networks were not only more accurate than those produced using individual data sets and other existing methods, but they also captured information regarding specific biological mechanisms and pathways that were missed using other methodologies. PANDA is scalable to higher eukaryotes, applicable to specific tissue or cell type data and conceptually generalizable to include a variety of regulatory, interaction, expression, and other genome-scale data. An implementation of the PANDA algorithm is available at www.sourceforge.net/projects/panda-net.

# Regulation of Transcription



Specific transcription factors

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

# A Simple Idea: Message Passing

Transcription Factor
The TF is *Responsible* for communicating with its Target

$$R_{ij}^{(t)} = \sum_k A_{ik}^{(t-1)} C_{kj} A_{ik}^{(0)} \Big/ \sum_k C_{kj} A_{ik}^{(0)}$$

Downstream Target
The Target must be *Available* to respond to the TF

$$A_{ij}^{(t)} = \sum_k R_{ik}^{(t-1)} C_{kj} R_{ik}^{(0)} \Big/ \sum_k R_{ik}^{(0)}$$

Kimberly Glass, GC Yuan

# Message-Passing Networks: PANDA



Glass et. al. "Passing Messages Between Biological Networks to Refine Predicted Interactions." PLoS One. 2013 May 31;8(5):e64832.
Code and related material available on sourceforge: http://sourceforge.net/projects/panda-net/

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

# Subtypes of Ovarian Cancer

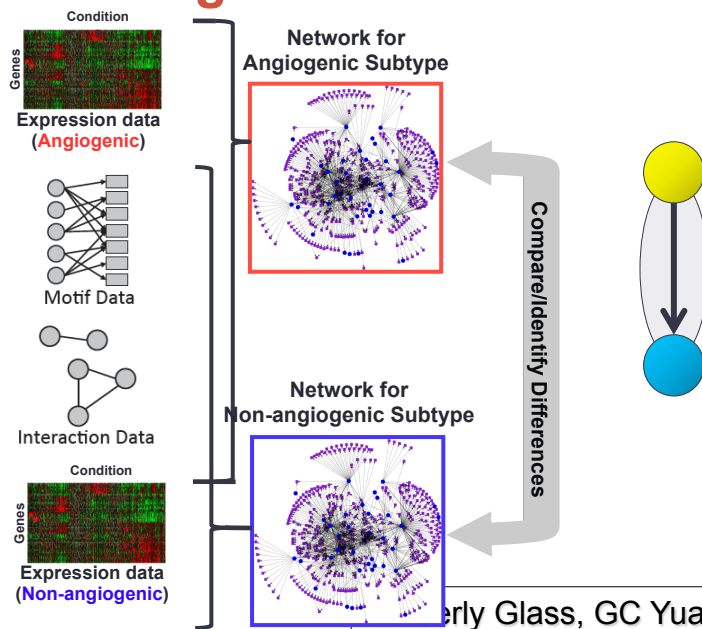## Angiogenic mRNA and microRNA Gene Expression Signature Predicts a Novel Subtype of Serous Ovarian Cancer

Stefan Bentink[1,6], Benjamin Haibe-Kains[1,6], Thomas Risch[1], Jian-Bing Fan[3], Michelle S. Hirsch[4,7], Kristina Holton[1], Renee Rubio[1], Craig April[3], Jing Chen[3], Eliza Wickham-Garcia[3], Joyce Liu[2,7], Aedin Culhane[1,6], Ronny Drapkin[4,5,7], John Quackenbush[1,2,6]*[¶], Ursula A. Matulonis[5,7][¶]

1 Department of Biostatistics and Computational Biology, Dana-Farber Cancer Institute, Boston, Massachusetts, United States of America, 2 Department of Cancer Biology, Dana-Farber Cancer Institute, Boston, Massachusetts, United States of America, 3 Illumina, Inc., San Diego, California, United States of America, 4 Department of Pathology, Division of Woman's and Perinatal Pathology, Brigham and Women's Hospital, Boston, Massachusetts, United States of America, 5 Department of Medical Oncology, Dana-Farber Cancer Institute, Boston, Massachusetts, United States of America, 6 Harvard School of Public Health, Boston, Massachusetts, United States of America, 7 Harvard Medical School, Boston, Massachusetts, United States of America
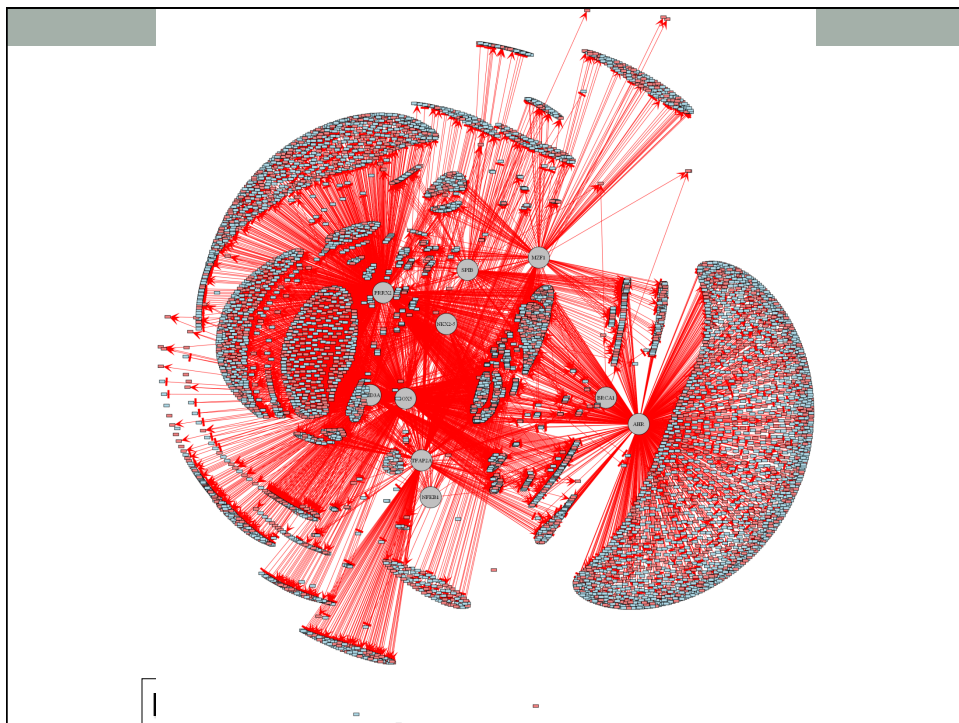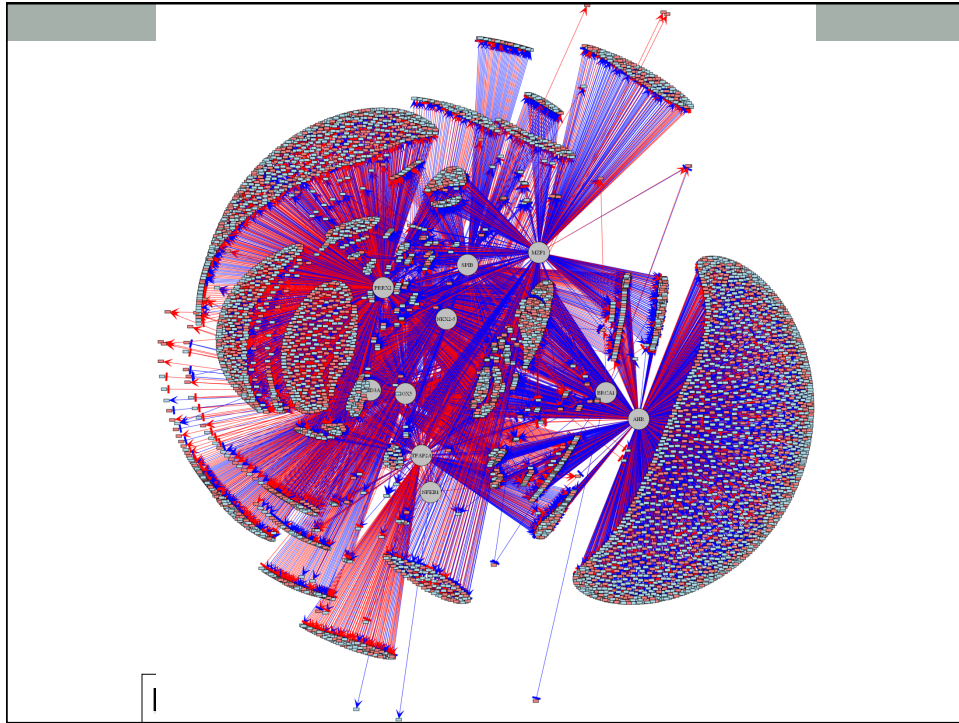
### Abstract

Ovarian cancer is the fifth leading cause of cancer death for women in the U.S. and the seventh most fatal worldwide. Although ovarian cancer is notable for its initial sensitivity to platinum-based therapies, the vast majority of patients eventually develop recurrent cancer and succumb to increasingly platinum-resistant disease. Modern, targeted cancer drugs intervene in cell signaling, and identifying key disease mechanisms and pathways would greatly advance our treatment abilities. In order to shed light on the molecular diversity of ovarian cancer, we performed comprehensive transcriptional profiling on 129 advanced stage, high grade serous ovarian cancers. We implemented a re-sampling based version of the
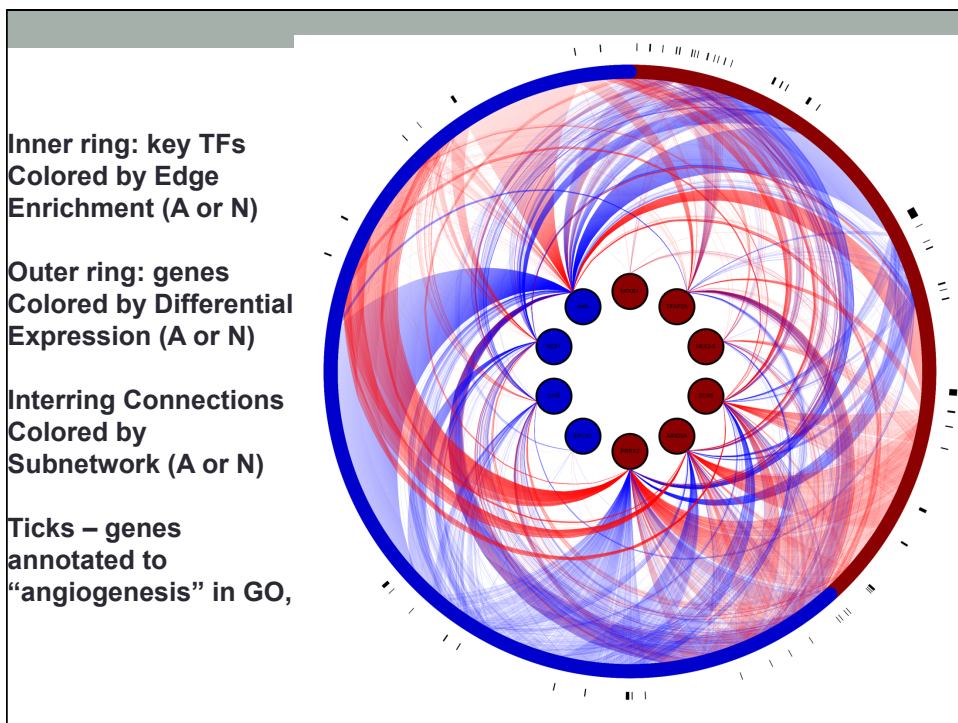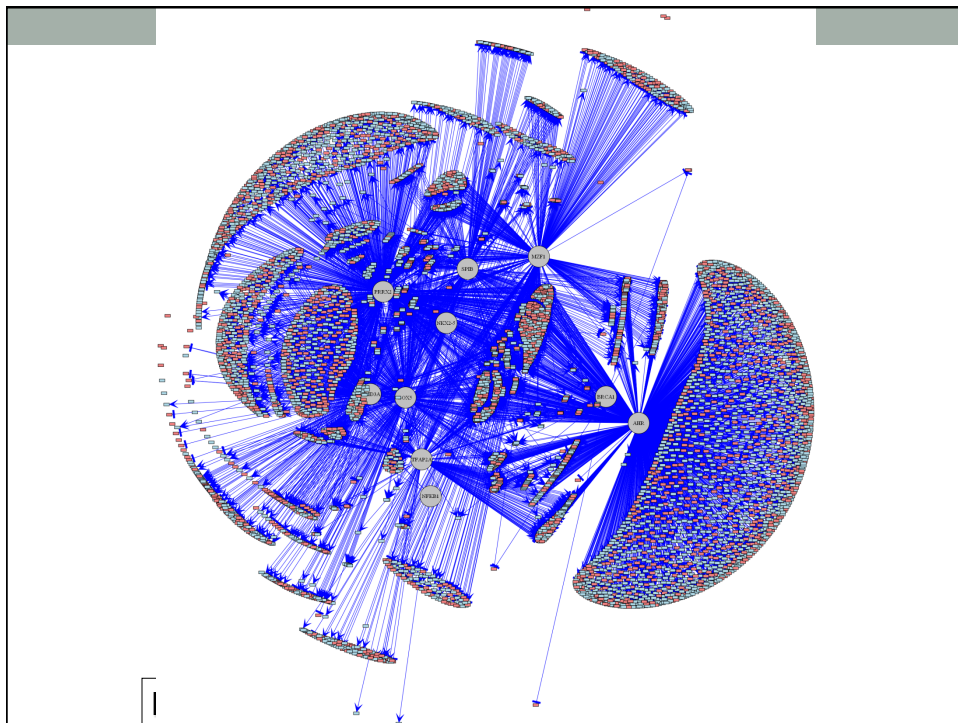
# PANDA: Integrative Network Models



Condition

Genes
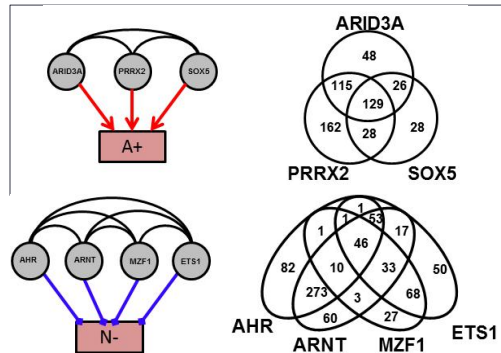
Expression data (**Angiogenic**)

Network for Angiogenic Subtype

Motif Data

Interaction Data

Condition

Genes

Expression data (**Non-angiogenic**)

Network for Non-angiogenic Subtype

Compare/Identify Differences

...erly Glass, GC Yuan

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

**Inner ring: key TFs Colored by Edge Enrichment (A or N)**

**Outer ring: genes Colored by Differential Expression (A or N)**

**Interring Connections Colored by Subnetwork (A or N)**

**Ticks – genes annotated to "angiogenesis" in GO,**

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

# Complex Regulatory Patterns Emerge

| TF1 | TF2 | sig. | # | Class | |
|-----|-----|------|---|-------|---|
| ARID3A | PRRX2 | 1.16E-23 | 244 | A+ | Co-regulatory TF Pairs |
| ARID3A | SOX5 | 1.01E-14 | 155 | A+ | |
| PRRX2 | SOX5 | 3.83E-12 | 157 | A+ | |
| ARNT | MZF1 | 5.83E-23 | 92 | N- | |
| AHR | ARNT | 6.13E-16 | 382 | N- | |
| ETS1 | MZF1 | 9.08E-16 | 148 | N- | |

Kimberly Glass, GC Yuan



# Regulatory Patterns suggest Therapies

ANGIOGENIC BEHAVIOR

HIF1a — ARNT
HIF2a — ETS1
→ VEGF production and angiogenesis

ARID3A — PRRX2 — SOX5

High levels of CpG methylation

TREATMENT MODEL

(1) Prevent ARNT/HIF1a and ETS1/HIF2a dimerization

(3) Decrease genome-wide methylation

(2) Promote ARNT/AHR and ETS1/AHR dimerization

Kimberly Glass, GC Yuan

43

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

**More application papers coming….**

# At the End of the Day

- The goal of an experiment is to discover new biology

- The challenge is sorting through lots of data

- Comparing groups of samples requires thorough annotation

- Making sense of the genes that are significant in such a comparison requires thorough gene annotation

- New technologies are giving us new ways of generating data, but the analysis approaches are more-or-less the same.

*NHGRI Current Topics in Genome Analysis 2016*
*Week 10: Expression Analysis, Functional Enrichment, and Network Inference*

*April 27, 2016*
*John Quackenbush, Ph.D.*

**The future is here.**
**It's just not widely distributed yet.**

**- William Gibson**

**Before I came here I was confused**
**about this subject.**
**After listening to your lecture,**
**I am still confused but at a higher level.**

**- Enrico Fermi, (1901-1954)**