# Integrative Genomic Analysis

## Sharing Data, Tools and Models

## Use of bionetworks to build better maps of disease

Stephen Friend MD PhD

Sage Bionetworks (Non-Profit Organization)
Seattle/ Beijing/ San Francisco

NAS
March 10th, 2011

BETTER MAPS OF DISEASE


NOT JUST WHAT WE DO BUT HOW WE DO IT


POWER OF BUILDING A PRE-COMPETITIVE
COMMONS FOR EVOLVING
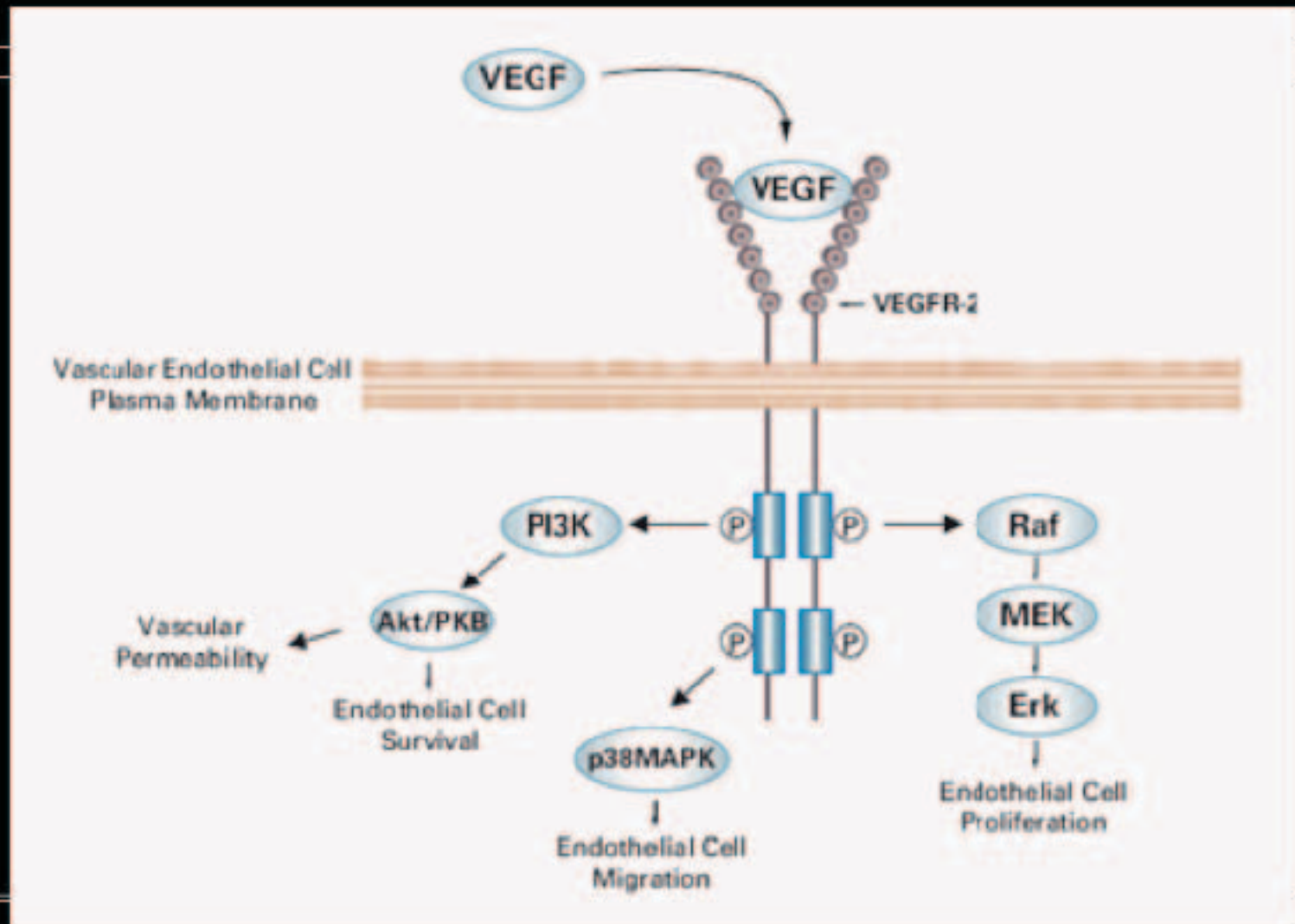GENERATIVE MODELS OF DISEASE

# Existing approaches and issues in Drug Discovery

Current costs for drug approval- ~$1Billion – 5 -10 years

Only 6% of therapies in Phase I trials will lead to approval (CMS)

Cancer- FDA approved marketed drugs as "standards of care" provide significant impact in only 25% of patients

# VEGFR Classical Pathway

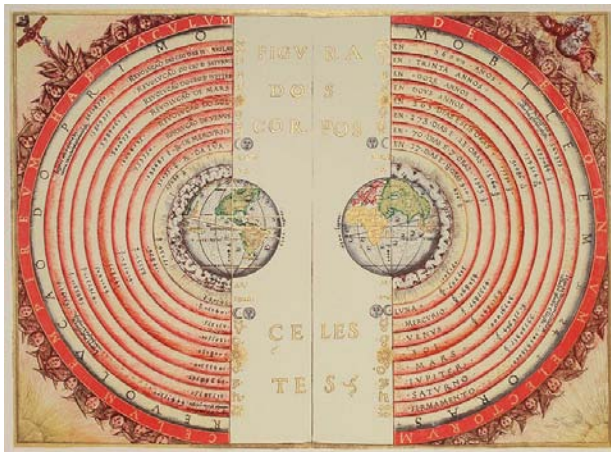# Where the data-driven network approach can combat the want our minds' have for story telling

**Zeus, the sky god; when he is angry he throws lightening bolts out of the sky**
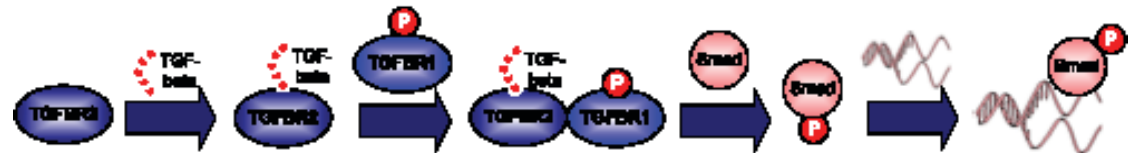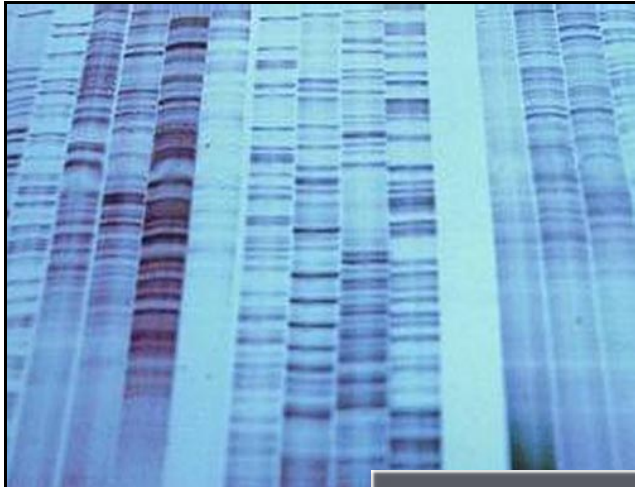


**The earth is flat**



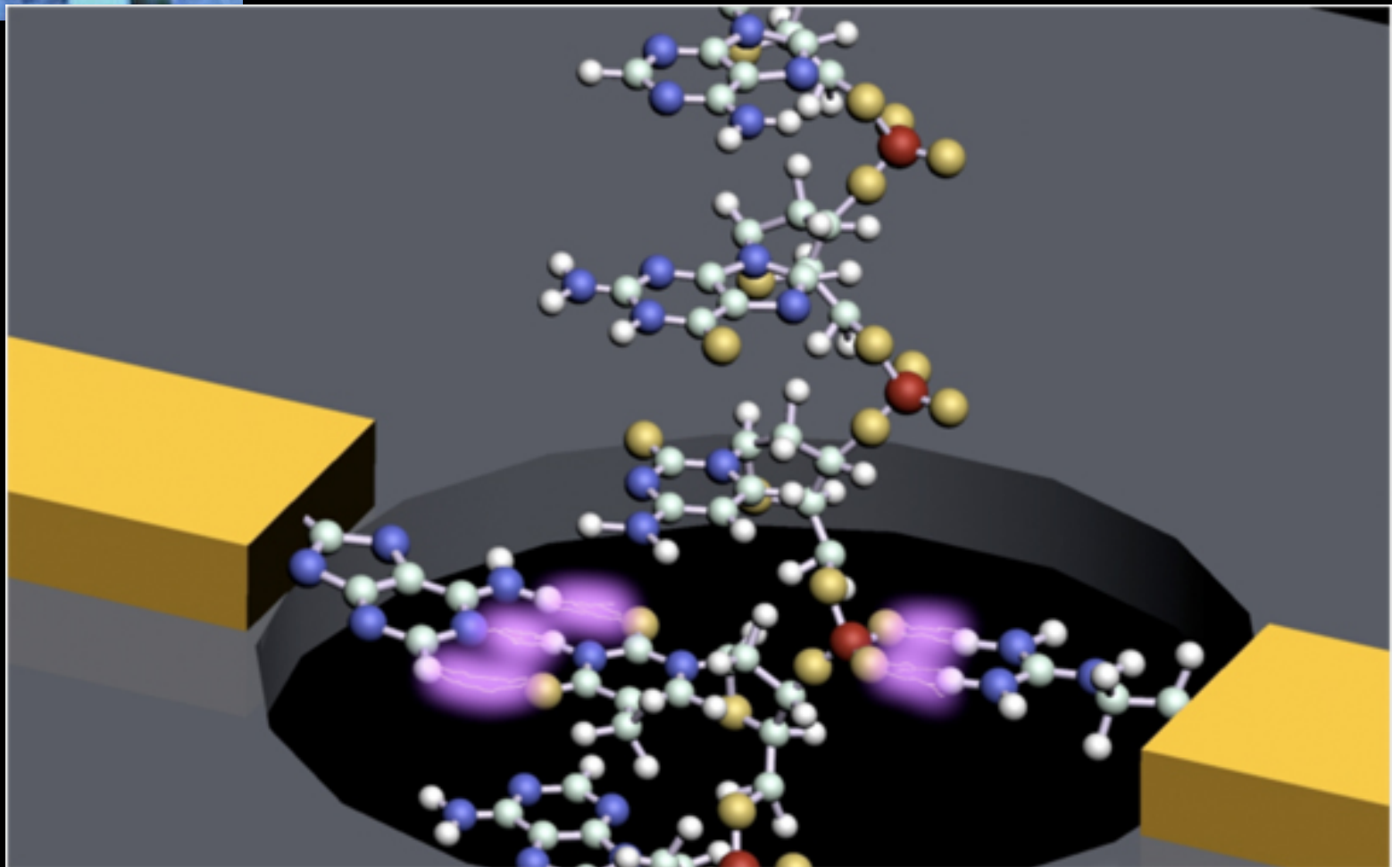**Ptolemaic astronomy: the earth is the center of the universe**



**Biological processes are driven by simple linearly ordered pathways (TGF-beta signaling)**
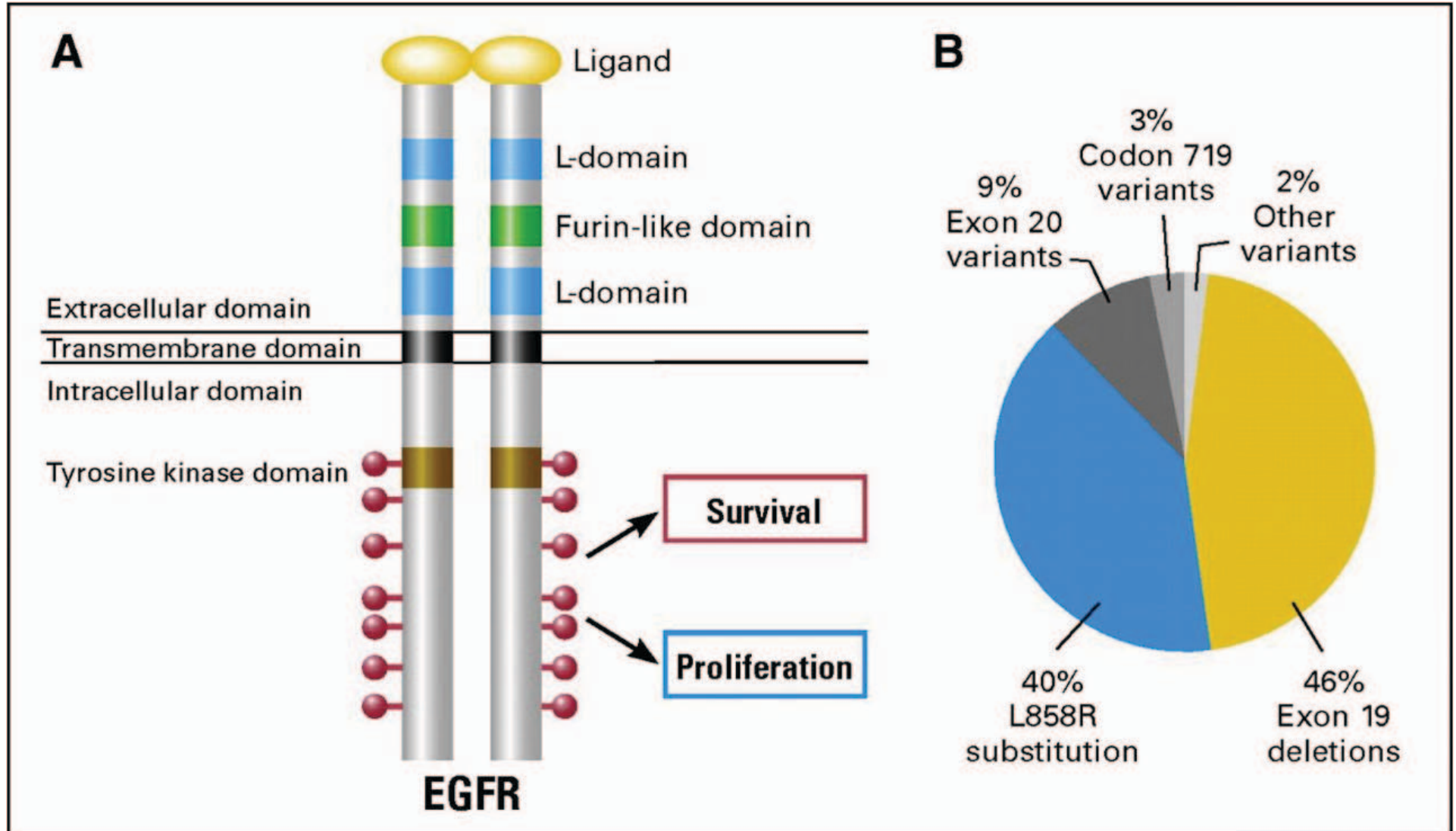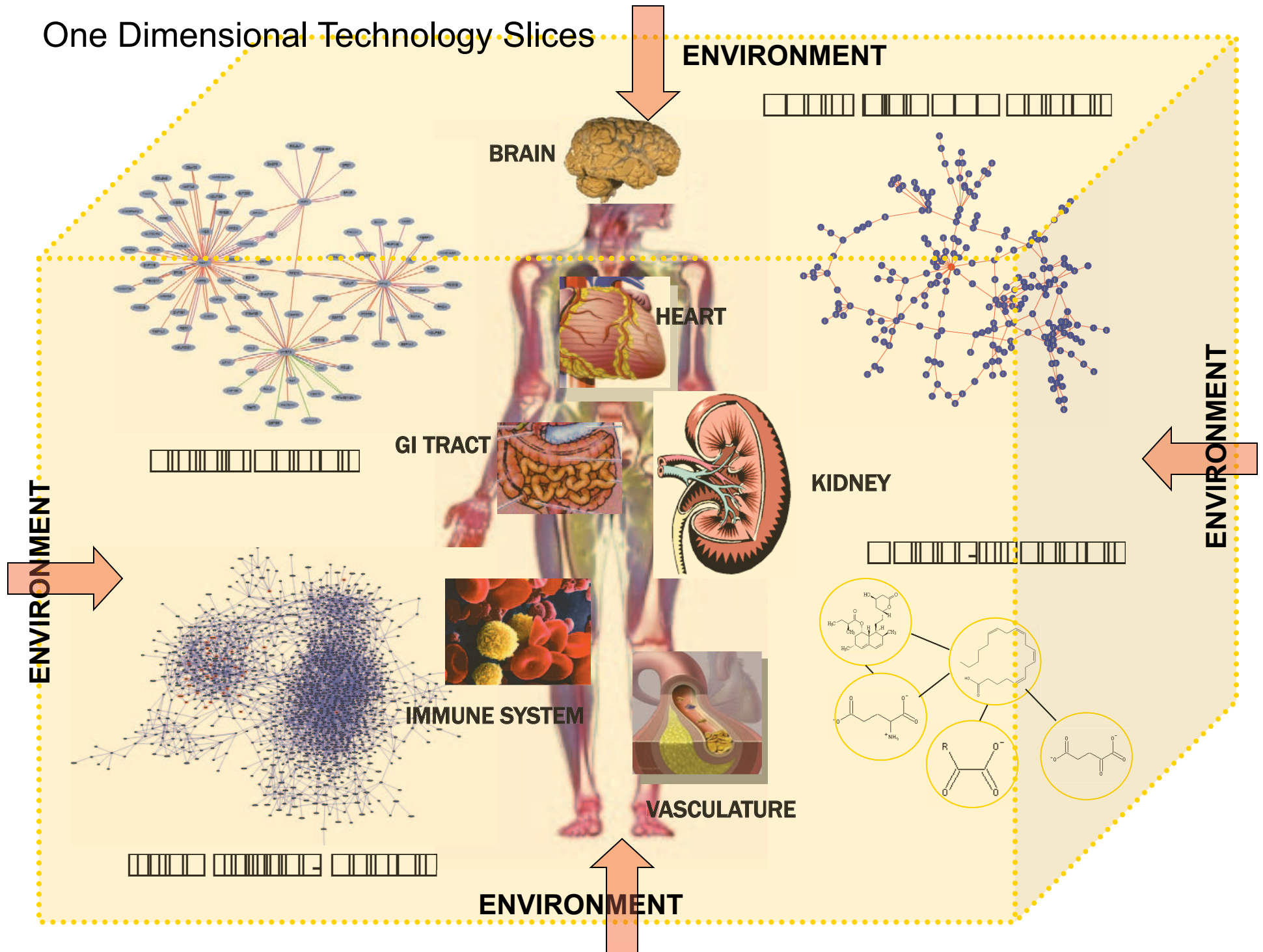


5

Amazing Technologies
will soon create a revolution
in our understanding of diseases
resulting in new categories,  and therapies
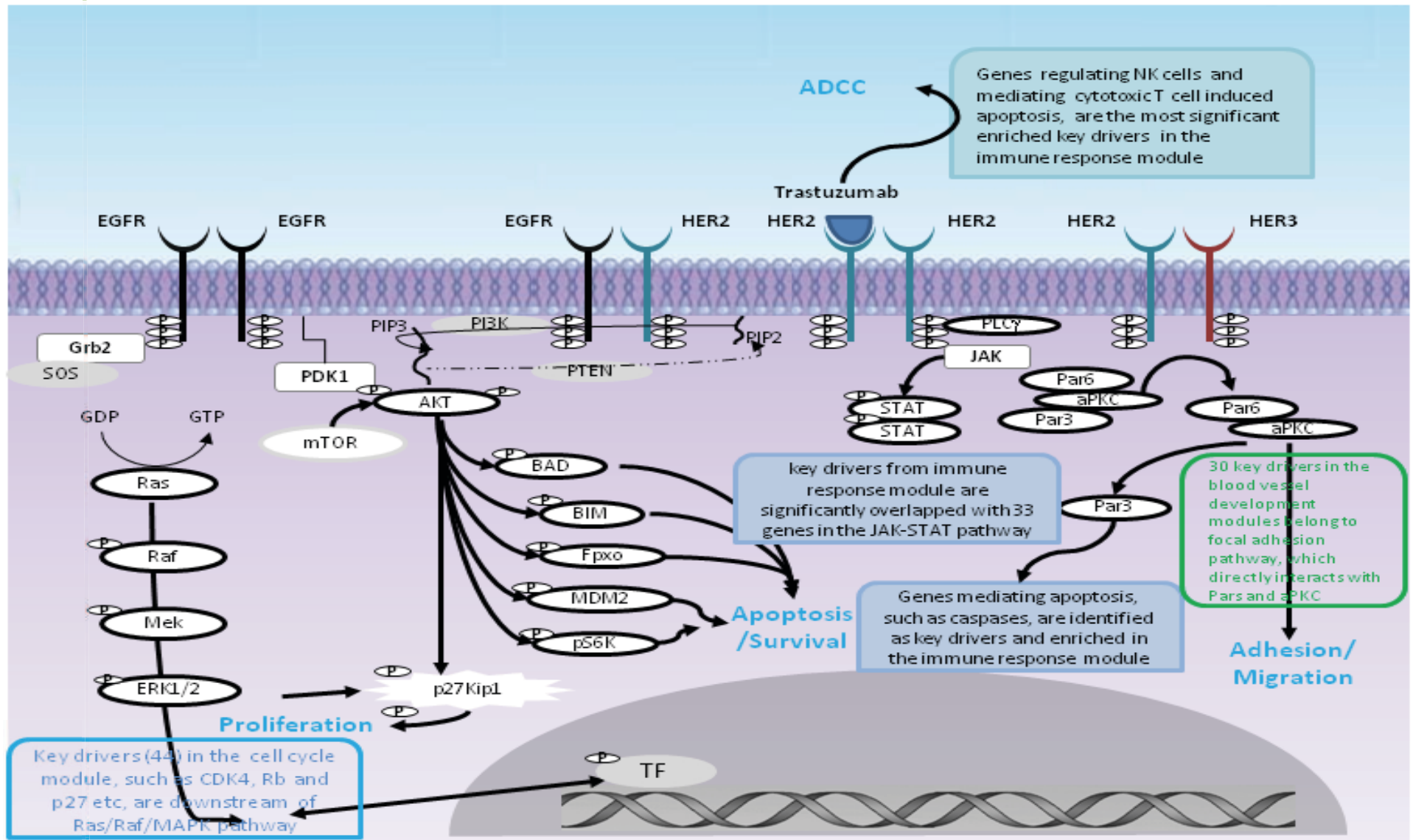
# Personalized Medicine 101:
# Capturing Single bases pair mutations = ID of responders

One Dimensional Technology Slices

ENVIRONMENT

ENVIRONMENT

ENVIRONMENT

ENVIRONMENT

BRAIN

HEART

GI TRACT

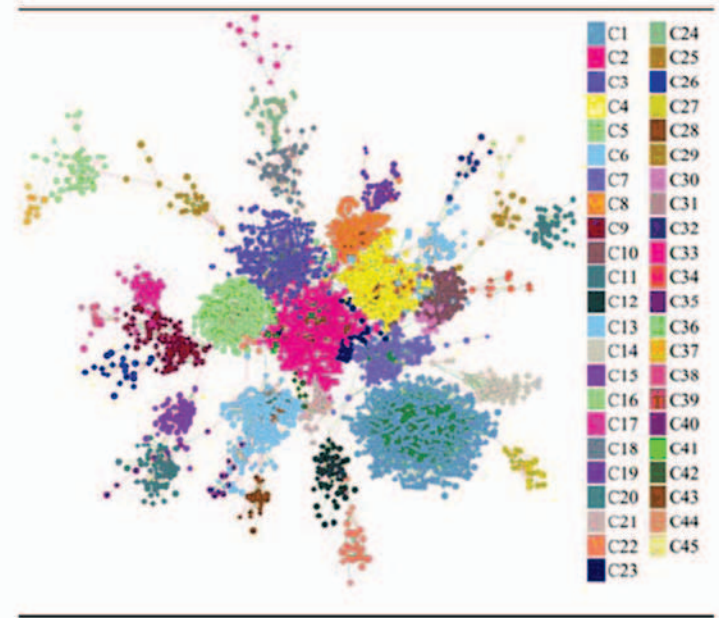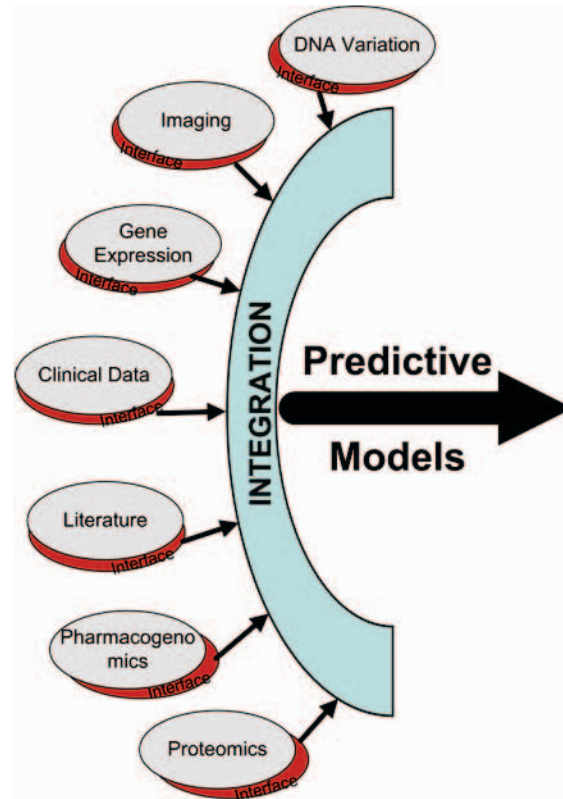KIDNEY

IMMUNE SYSTEM

VASCULATURE

# Example of Complexity: Overlapping of EGFR and Her2 Pathways

# The "Rosetta Integrative Genomics Experiment": Generation, assembly, and integration of data to build models that predict clinical outcome

**Merck Inc. Co.
5 Year Program
Based at Rosetta
Total Resources
>$150M**



- Generate data need to build
-  bionetworks
- Assemble other available data useful for building networks
- Integrate and build models
- Test predictions
- Develop treatments
- Design Predictive Markers

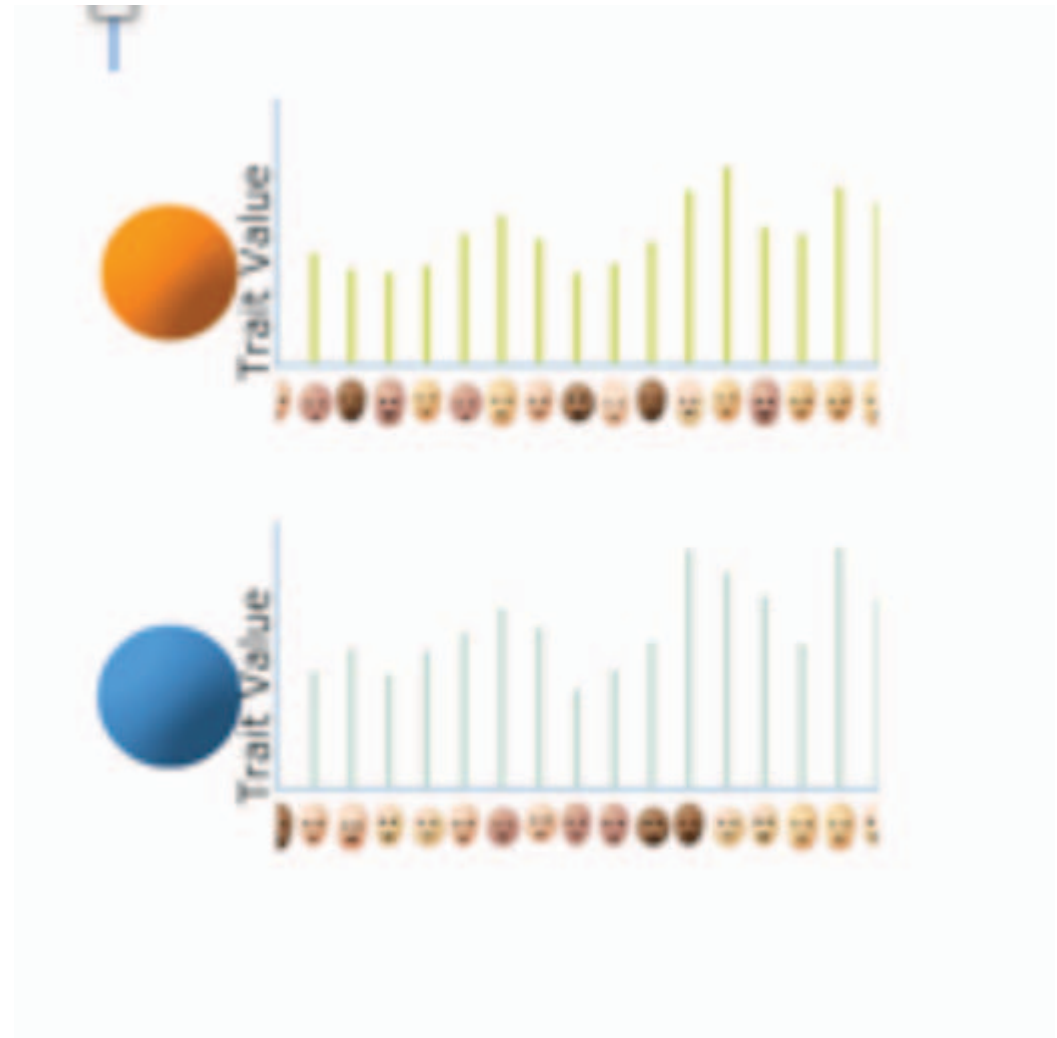# We need to carry out comprehensive monitoring of many traits at the population level

**Monitor disease and molecular traits in populations**



Putative causal gene

Disease trait

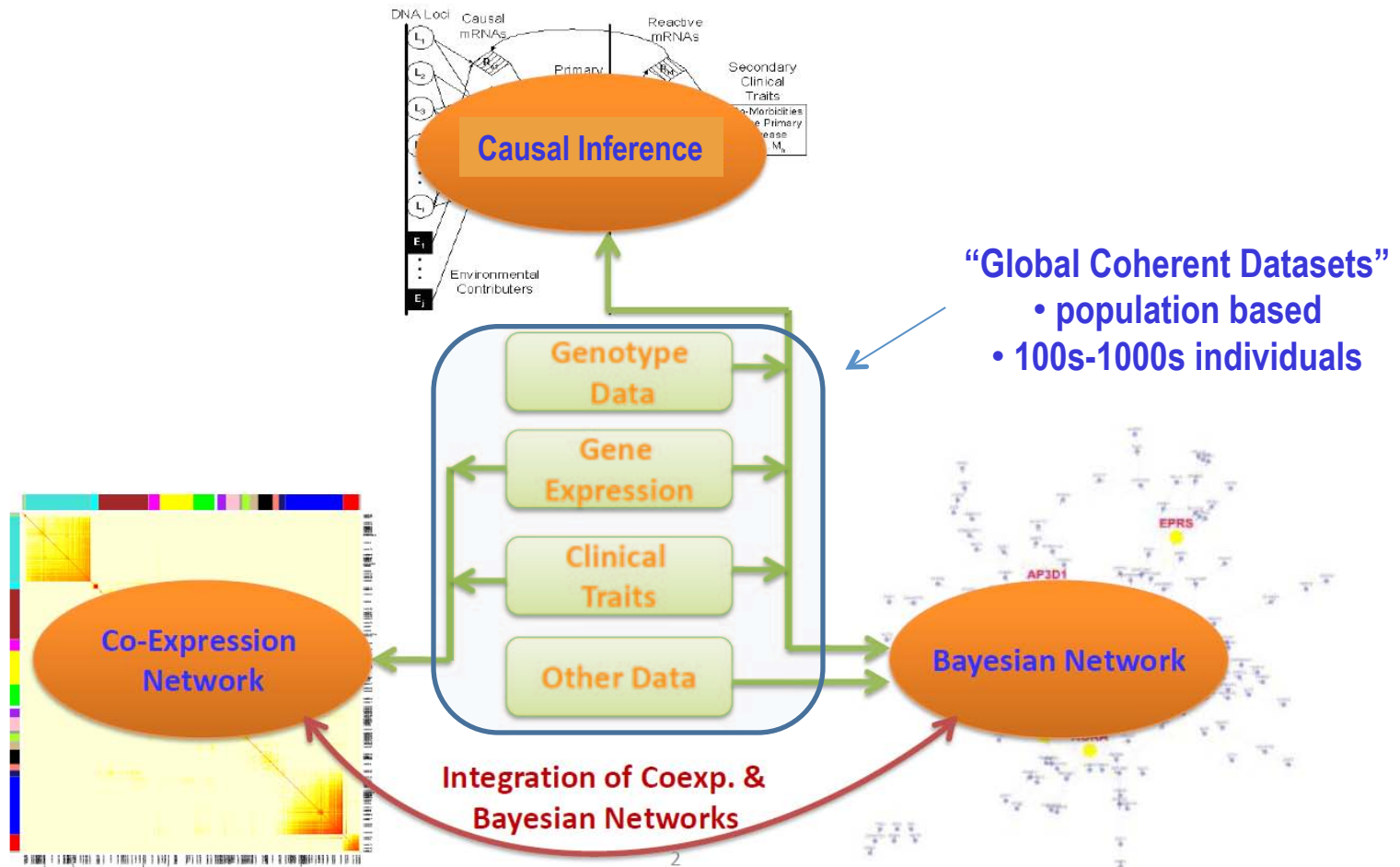# Leveraging DNA variation as a systematic perturbation source

**Phenotype 1**

**Given this triangle relationship, we can mathematically infer the most likely relationship**

**cQTL**

**Trait-trait correlation**

**eQTL**

**DNA Variant**

**Phenotype 2**

# Integration of Genotypic, Gene Expression & Trait Data



Schadt et al. Nature Genetics 37: 710 (2005)
Millstein et al. BMC Genetics 10: 23 (2009)

**Causal Inference**

"Global Coherent Datasets"
• population based
• 100s-1000s individuals

Genotype Data

Gene Expression

Clinical Traits

Other Data

**Co-Expression Network**

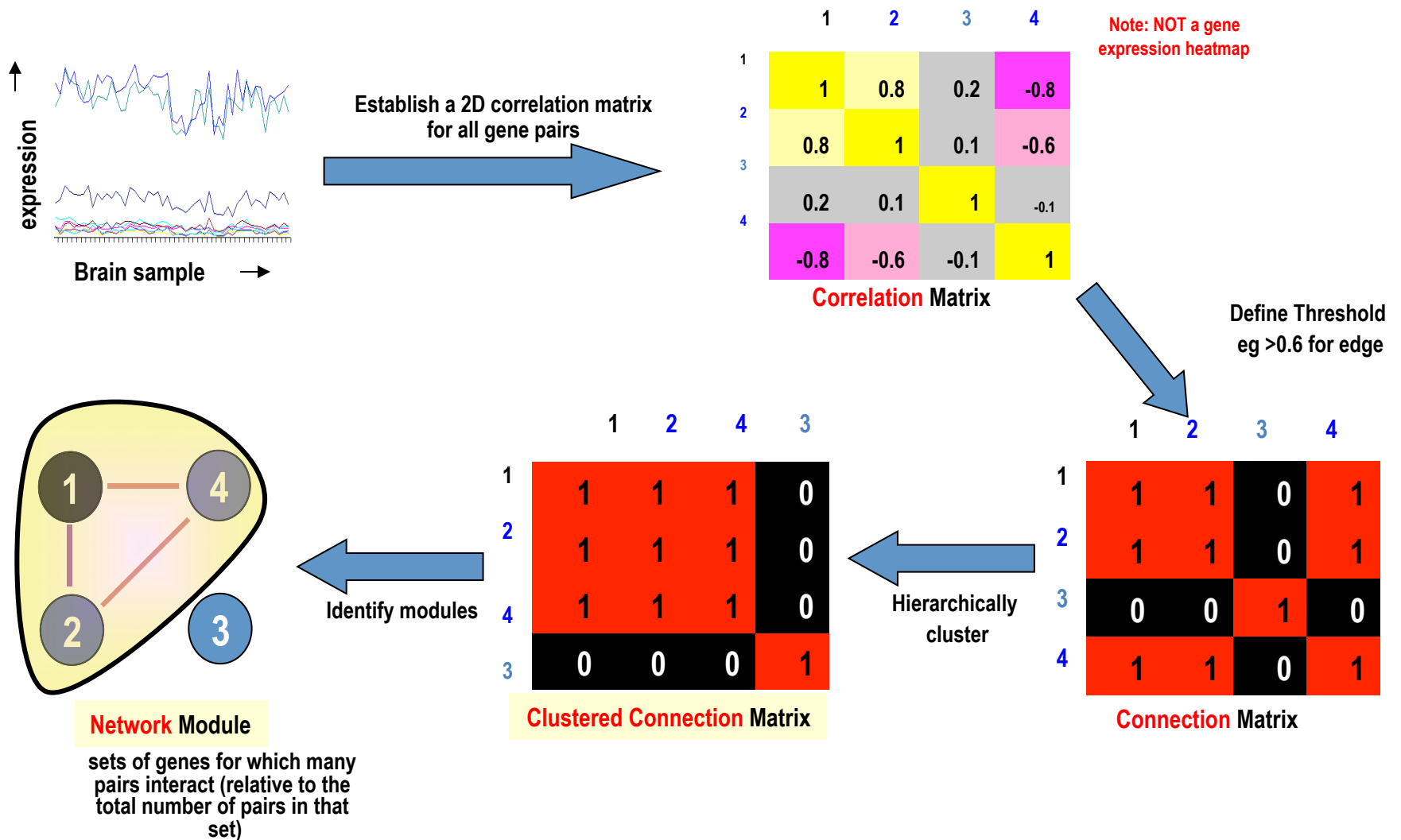**Bayesian Network**

Integration of Coexp. & Bayesian Networks

Chen et al. Nature 452:429 (2008)
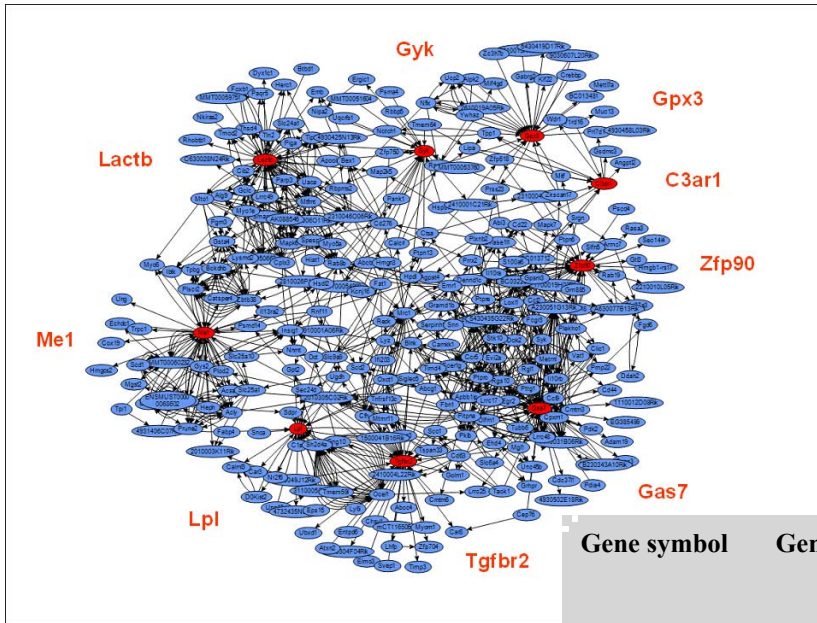Zhang & Horvath. Stat.Appl.Genet.Mol.Biol. 4: article 17 (2005)

Zhu et al. Cytogenet Genome Res. 105:363 (2004)
Zhu et al. PLoS Comput. Biol. 3: e69 (2007)

# Constructing Co-expression Networks

**Start with expression measures for ~13K genes most variant genes across 100-150 samples**
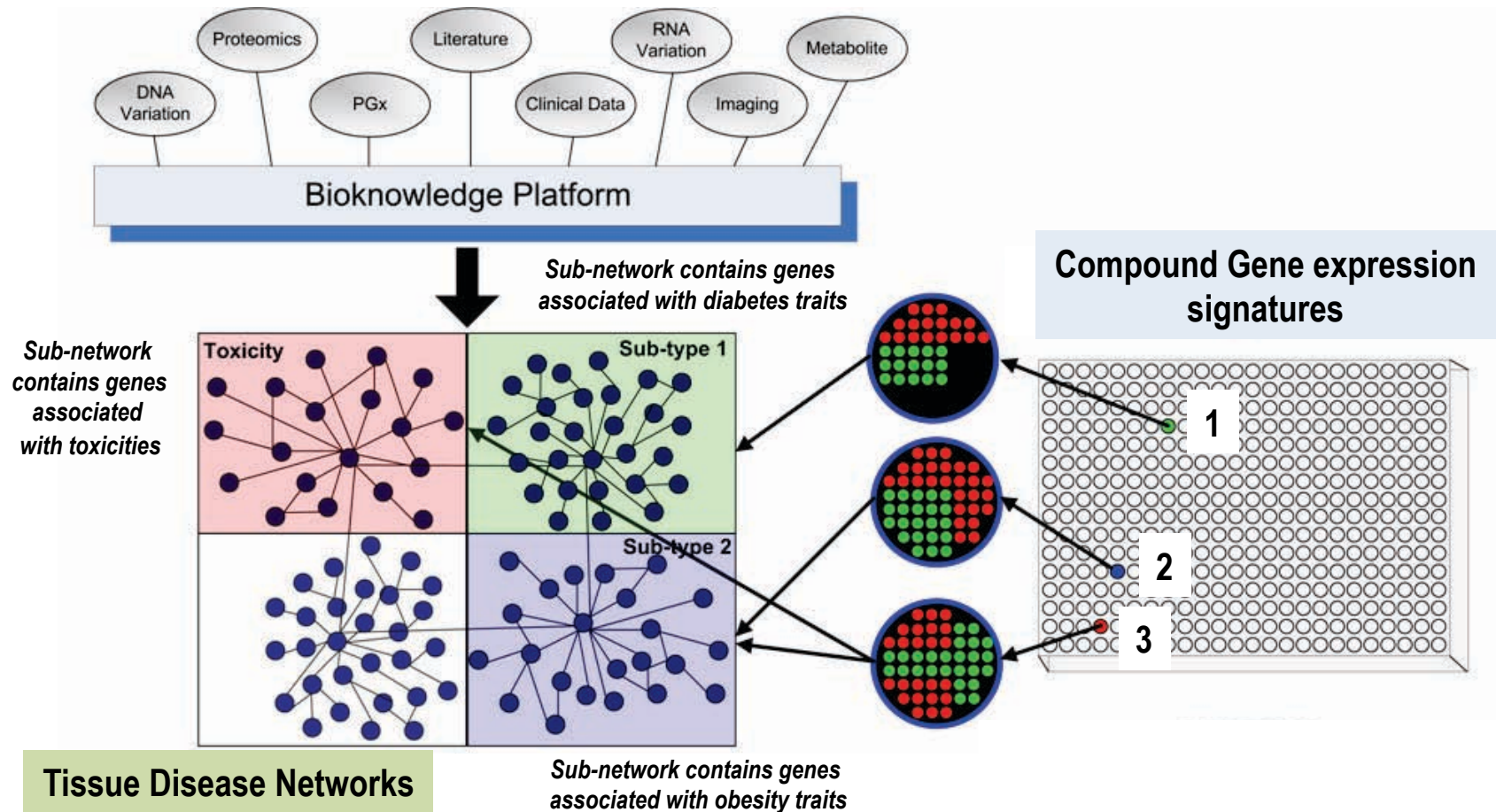
# Probabalistic Models- Rosetta



Networks facilitate direct identification of genes that are causal for disease
Evolutionarily tolerated weak spots

| Gene symbol | Gene name | Variance of OFPM explained by gene expression* | Mouse model | Source |
|---|---|---|---|---|
| Zfp90 | Zinc finger protein 90 | 68% | tg | Constructed using BAC transgenics |
| Gas7 | Growth arrest specific 7 | 68% | tg | Constructed using BAC transgenics |
| Gpx3 | Glutathione peroxidase 3 | 61% | tg | Provided by Prof. Oleg Mirochnitchenko (University of Medicine and Dentistry at New Jersey, NJ) [12] |
| Lactb | Lactamase beta | 52% | tg | Constructed using BAC transgenics |
| Me1 | Malic enzyme 1 | 52% | ko | Naturally occurring KO |
| Gyk | Glycerol kinase | 46% | ko | Provided by Dr. Katrina Dipple (UCLA) [13] |
| Lpl | Lipoprotein lipase | 46% | ko | Provided by Dr. Ira Goldberg (Columbia University, NY) [11] |
| C3ar1 | Complement component 3a receptor 1 | 46% | ko | Purchased from Deltagen, CA |
| Tgfbr2 | Transforming growth factor beta receptor 2 | 39% | ko | Purchased from Deltagen, CA |

Nat Genet (2005) 205:370

# Map compound signatures to disease networks



Compound 1: Drug signature significantly enriched in subnetwork associated with diabetes traits

Compound 2: Drug signature significantly enriched in subnetwork associated with obesity traits

Compound 3: Drug signature significantly enriched in subnetwork associated with obesity traits BUT also in subnetwork associated with toxicities

# Extensive Publications now Substantiating Scientific Approach
## Probabilistic Causal Bionetwork Models

- **>80 Publications from Rosetta Genetics / Sage Group (~30 scientists) over 6 years including high profile papers in PLoS Nature and Nature Genetics**
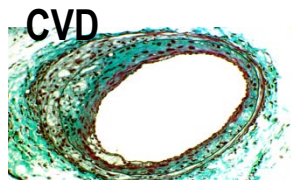

**Metabolic Disease**

"*Genetics of gene expression surveyed in maize, mouse and man.*" **Nature**. (2003)

"*Variations in DNA elucidate molecular networks that cause disease.*" **Nature**. (2008)

"*Genetics of gene expression and its effect on disease.*" **Nature**. (2008)

"*Validation of candidate causal genes for obesity that affect...*" **Nat Genet**. (2009)

**….. Plus 10 additional papers in Genome Research, PLoS Genetics, PLoS Comp.Biology, etc**


CVD

**"*Identification of pathways for atherosclerosis.*"** Circ Res. **(2007)**

"*Mapping the genetic architecture of gene expression in human liver.*" **PLoS Biol**. (2008)

**…… Plus 5 additional papers in Genome Res., Genomics, Mamm.Genome**


**Bone**

"*Integrating genotypic and expression data …for bone traits…*" **Nat Genet**. (2005)

"*..approach to identify candidate genes regulating BMD…*" **J Bone Miner Res**. (2009)


**Methods**

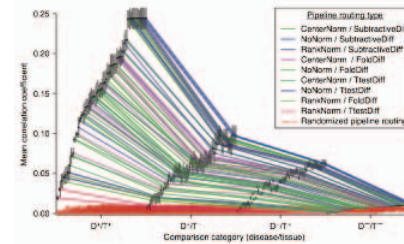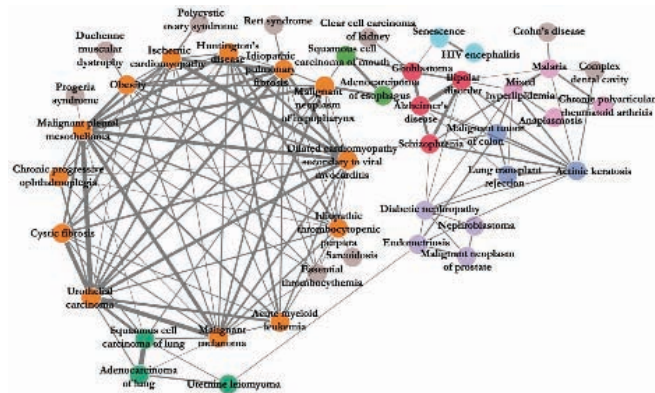"***An integrative genomics approach to infer causal associations ...*" Nat Genet. (2005)**

"***Increasing the power to detect causal associations… "*** PLoS Comput Biol. (2007)**

"*Integrating large-scale functional genomic data ...*" **Nat Genet.** (2008)

…… Plus 3 additional papers in **PLoS Genet., BMC Genet.**

# Exploring the Global Landscape of Human Disease Through Public Data- approaches taken by Atul Butte

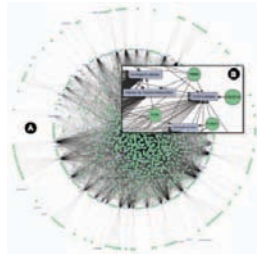**Public data enables quantitative disease relationships**



**High quality signals exist in public data**

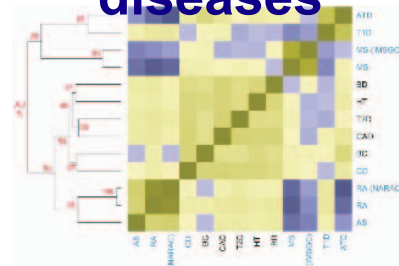Joel Dudley et al.. Molecular systems biology (2009) vol. 5 pp. 307

# Differences

# Commonalities

**Genetic architecture of autoimmune diseases**

**Plasma proteome networks**



Joel Dudley and Atul Butte. Pacific Symposium on Biocomputing (2009) pp. 27-38

**Functional gene module networks**



Marina Sirota et al. PLoS genetics (2009) vol. 5 (12) pp. e1000792

Silpa Suthram et al. PLoS computational biology (2010) vol. 6 (2) pp. e1000662

**Which biomarkers best discriminate diseases?**

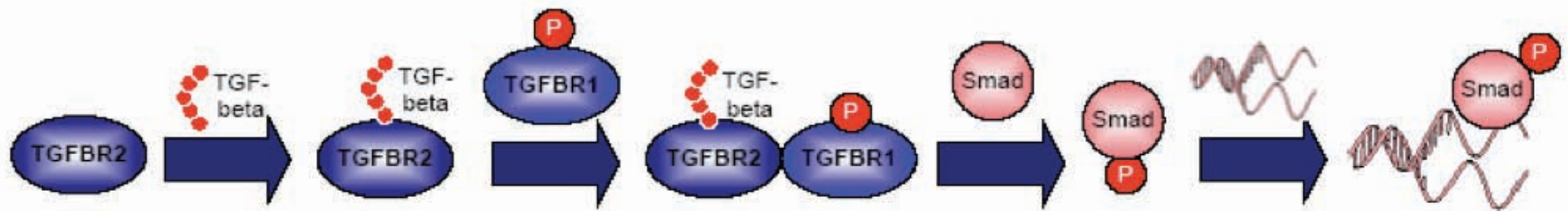**Is there a blood biomarker for general pathology?**

**Are there genetic "switches" for autoimmunity?**

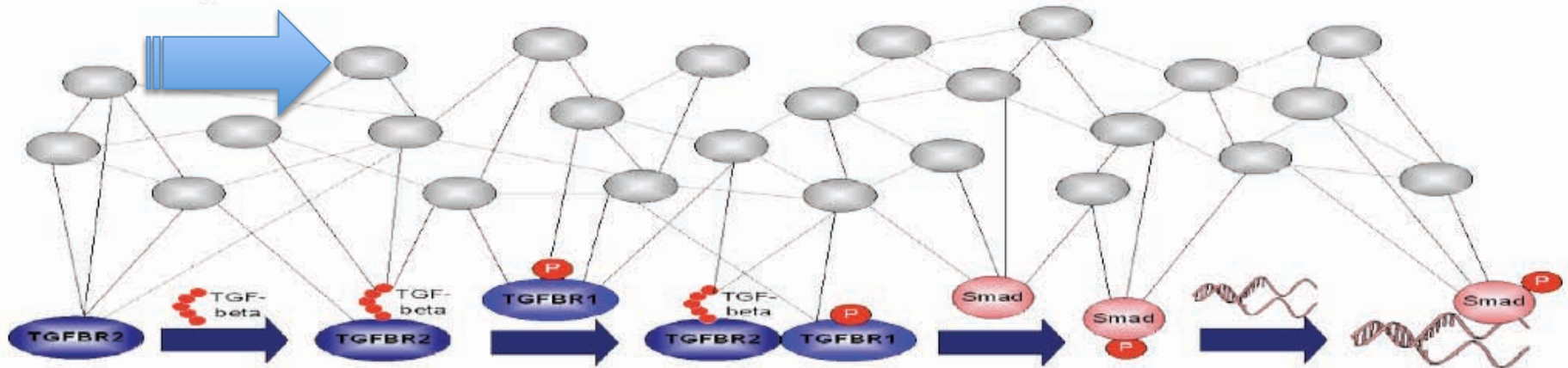**Is there a common autoimmune susceptibility variant?**

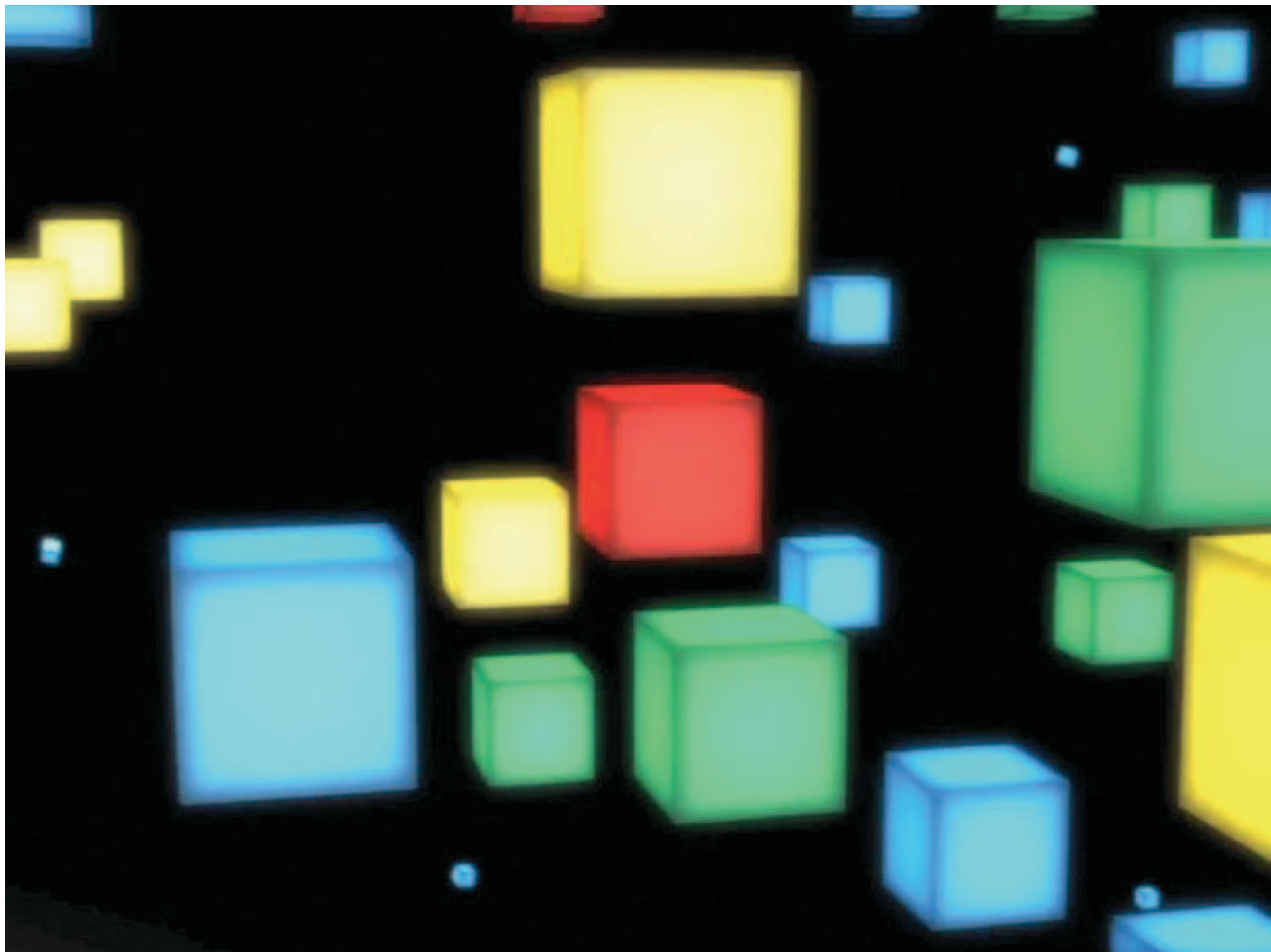**Which modules are unique to metabolic diseases?**

**Do common modules harbor pluripotent drug targets?**
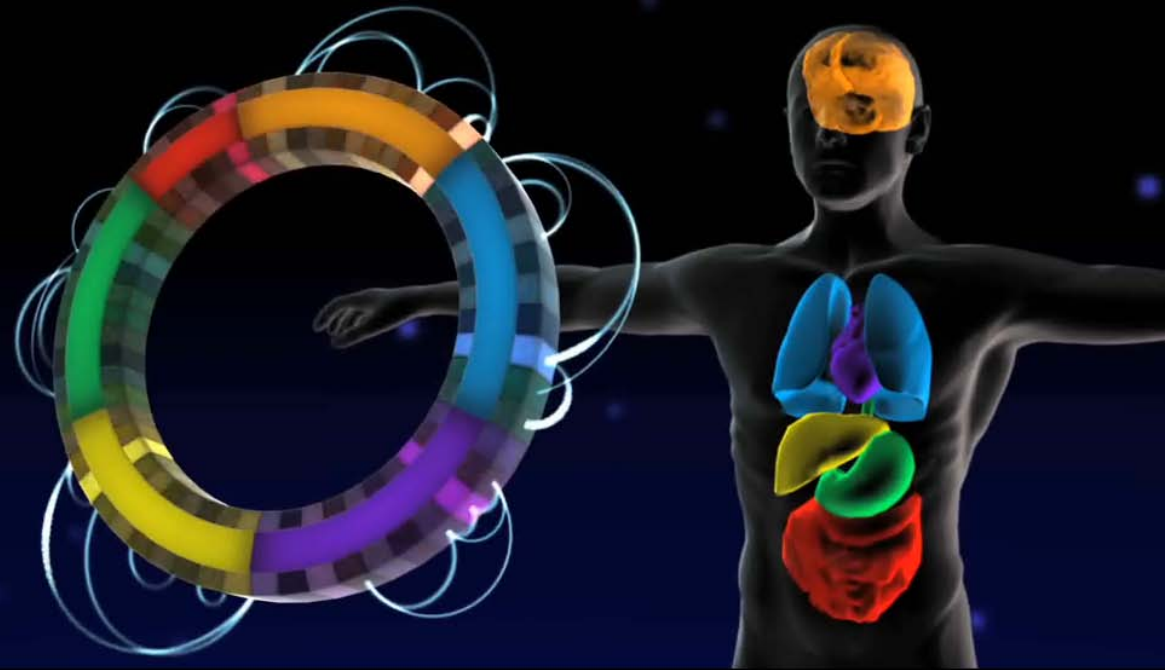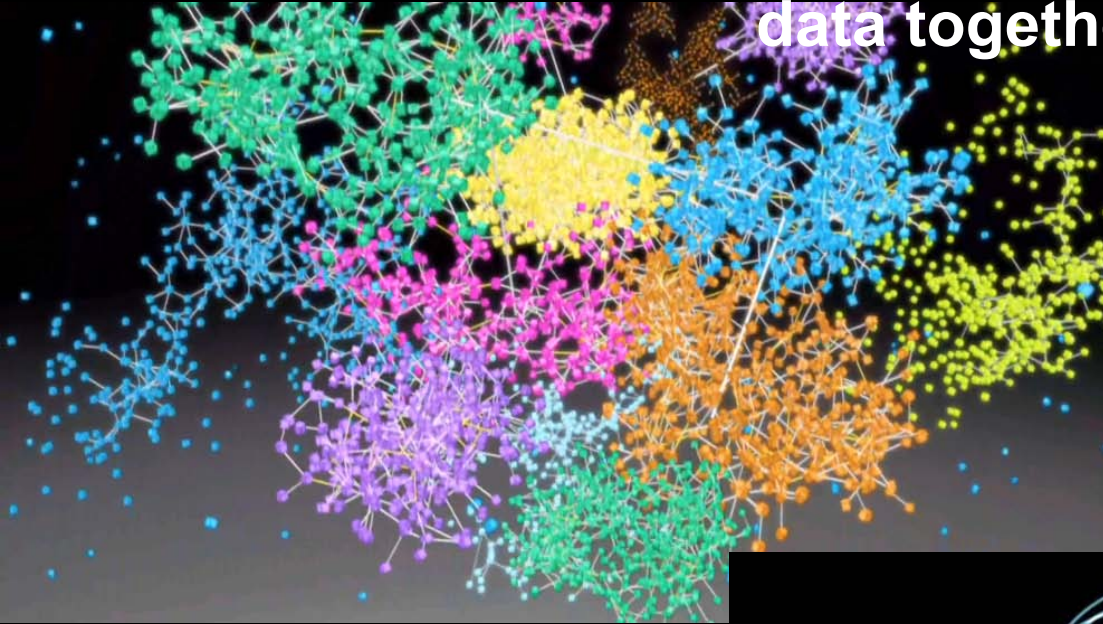
# The way we like to think:



# The way it is:

We are about to be able to build maps of disease based on alterations found in multitudes of patients if we can bring the data together

Recognition that the benefits of bionetwork based molecular models of diseases are powerful but that they **require significant resources**
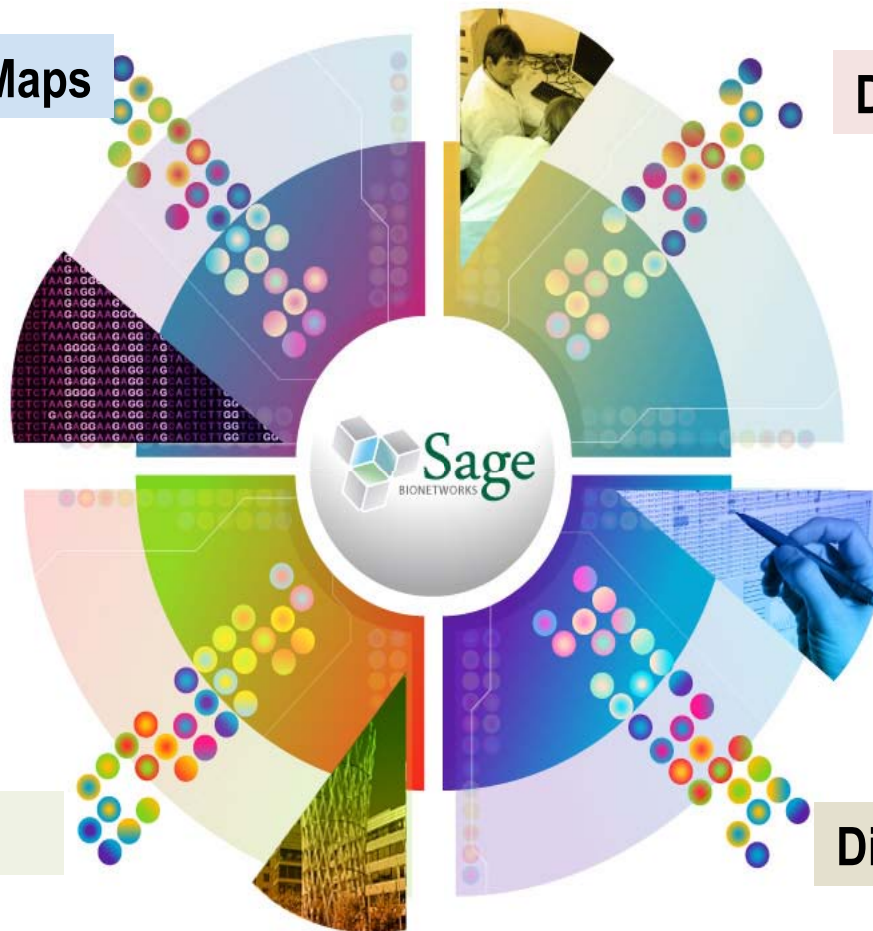
Appreciation that it will **require decades** of evolving representations as real complexity emerges and needs to be integrated with therapeutic interventions

# Sage Mission

Sage Bionetworks is a non-profit organization with a vision to create a "commons" where integrative bionetworks are evolved by contributor scientists with a shared vision to accelerate the elimination of human disease
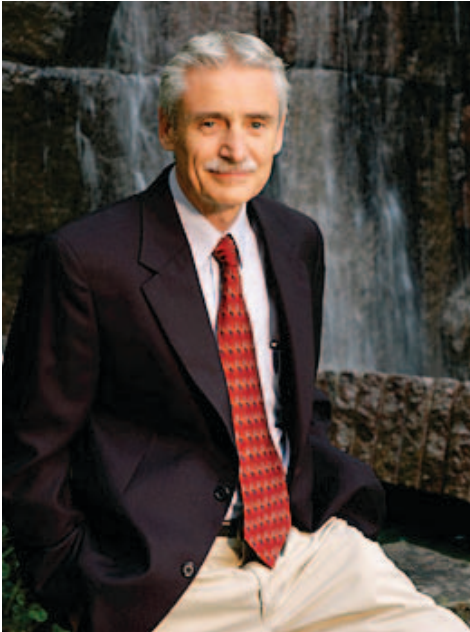


**Building Disease Maps**

**Data Repository**

**Commons Pilots**

Sagebase.org

**Discovery Platform**

# Board of Directors- Sage Bionetworks



**Lee Hartwell**

**Nobel Laureate
Co-Founder Rosetta**

**Hans Wizgell**

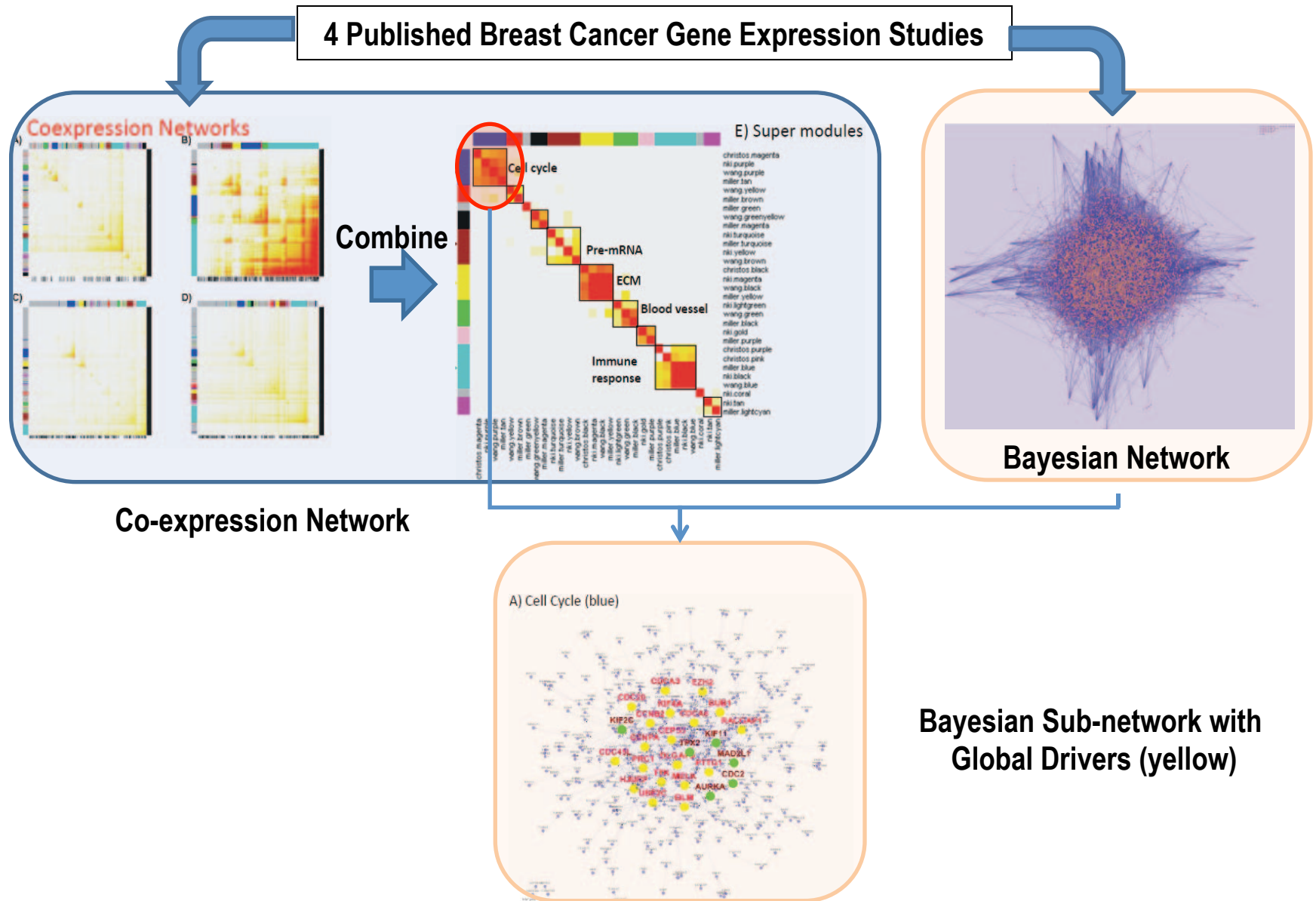**ExPresident Karolinska
Head SAB Rosetta**

**WangJun**

**Executive Director
BGI**

**Jeff Hammerbacher**

**CEO Cloudera
Built and Headed
Facebook
Data Architecture**

# Example 1: Identification of Molecular Drivers of Breast Cancer



**4 Published Breast Cancer Gene Expression Studies**

**Co-expression Network**

**Combine**

**Bayesian Network**

**Bayesian Sub-network with Global Drivers (yellow)**

# Example 2. The Sage Non-Responder Project in Cancer

**Purpose:**
- To identify Non-Responders to approved drug regimens so we can improve outcomes, spare patients unnecessary toxicities from treatments that have no benefit to them, and reduce healthcare costs

**Leadership:**
- Co-Chairs Stephen Friend, Todd Golub, Charles Sawyers & Rich Schilsky

**Initial Studies:**
- AML (at first relapse)
- Non-Small Cell Lung Cancer
- Ovarian Cancer (at first relapse)
- Breast Cancer
- Renal Cell
- Multiple Myeloma

Partnering for Cures | A FasterCures Meeting

Bridging the Chasm Between Microscope and Marketplace

- **Description**: Collate, Annotate, Curate and Host Clinical Trial Data with Genomic Information from the Comparator Arms of Industry and Foundation Sponsored Clinical Trials: Building a Site for Sharing Data and Models to evolve better Disease Maps.

- **Public-Private Partnership** of leading pharmaceutical companies, clinical trial groups and researchers.

- **Neutral Conveners**: Sage Bionetworks and Genetic Alliance [nonprofits].

- **Initiative to share existing trial data** (molecular and clinical) from non-proprietary comparator and placebo arms to create powerful new tool for drug development.

PARTNERINGFORCURES.ORG

# Example 4: The Sage Federation

- Founding Lab Groups

  - Seattle- Sage Bionetworks
  - New York- Columbia: Andrea Califano
  - Palo Alto- Stanford: Atul Butte- (Joel Dudley is here)
  - San Diego- UCSD: Trey Ideker
  - San Francisco: UCSF/Sage: Eric Schadt

- Initial Projects
  - Aging
  - Diabetes
  - Warburg

- Goals: *Share all datasets, tools, models*
  *Develop interoperability for human data*

# THE FEDERATION
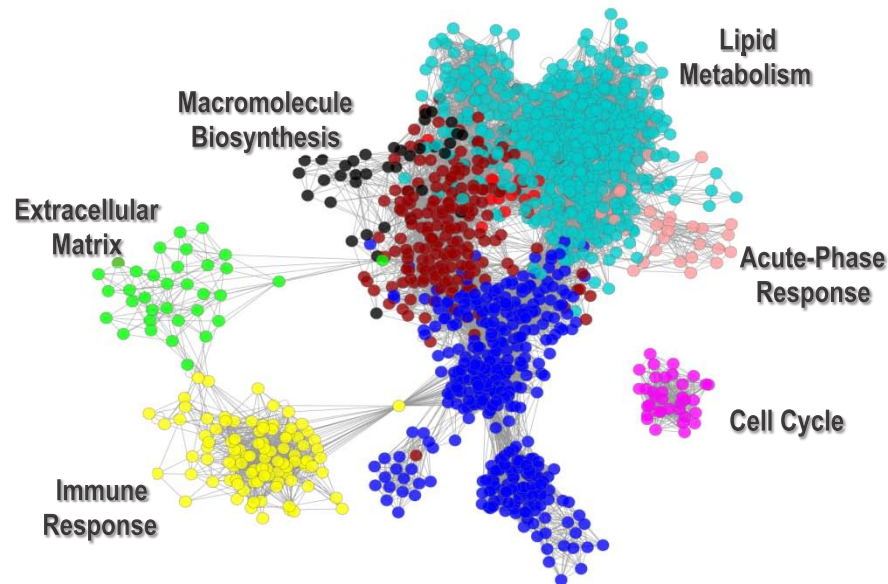
## Butte    Califano    Friend    Ideker    Schadt

## vs

# Sage Bionetworks: Platform



Research
Training
Platform

Community Systems Biologists
Evolving Disease Models
Repository

## GLOBAL COHERENT DATASETS

A data set containing genome-wide DNA variation and intermediate trait, as well as physiological phenotype data across a population of individuals large enough to power association or linkage studies, typically 50 or more individuals. To be coherent, the data needs to be matched with consistent identifiers. Intermediate traits are typically gene expression, but may also include proteomic, metabolomic, and other molecular data.

See http://www.sagebase.org/commons/repository.php

## MODELS



Lipid Metabolism

Macromolecule Biosynthesis

Extracellular Matrix

Acute-Phase Response

Cell Cycle

Immune Response

## TOOLS

### Key Driver Analysis (KDA) Tool  (R package/Cystoscape plug in)

http://sagebase.org/research/tools.php

# NOT JUST WHAT BUT HOW

Clinical/genomic data
are accessible but minimally usable

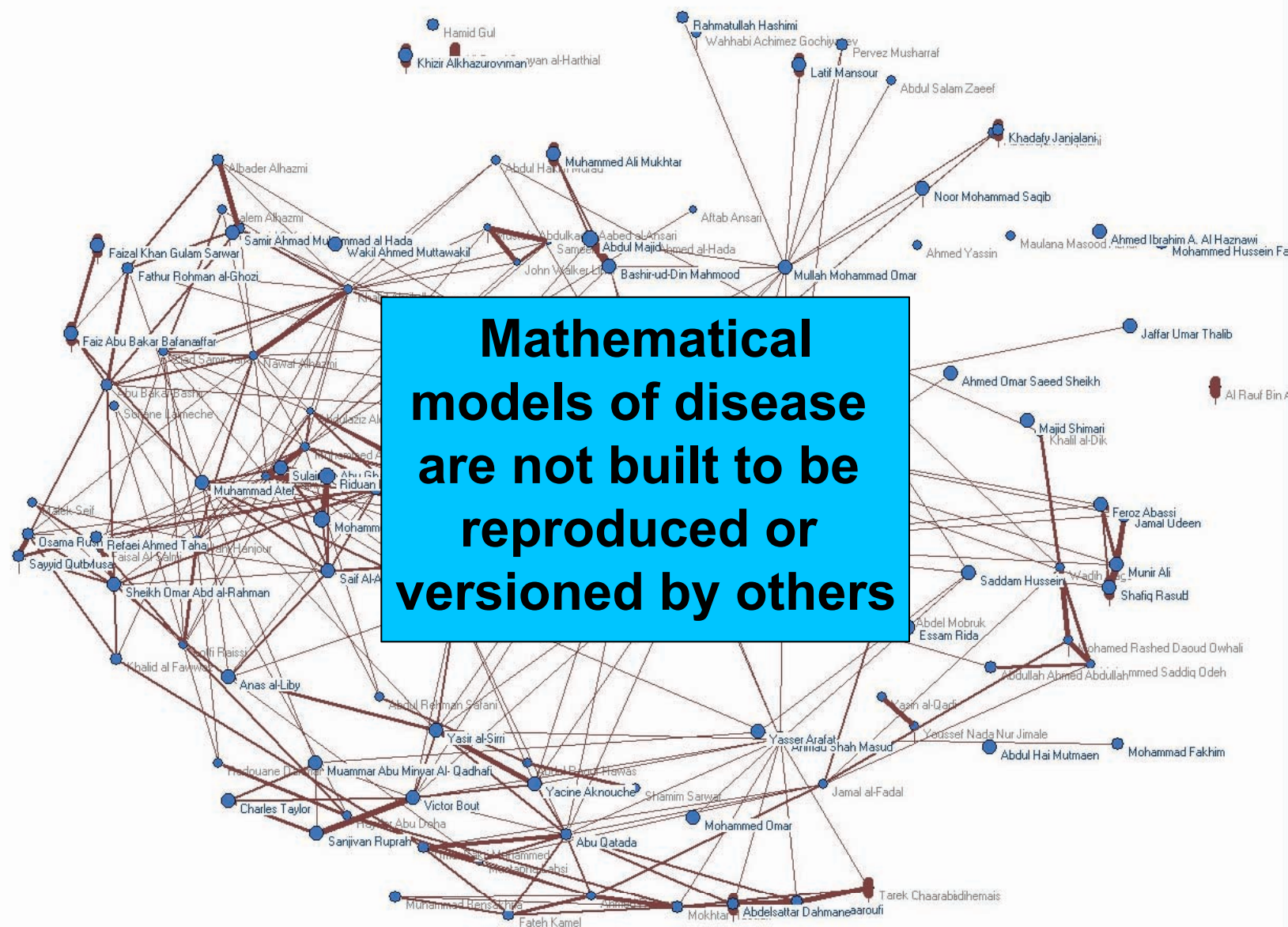Little incentive to annotate and curate
data for other scientists to use

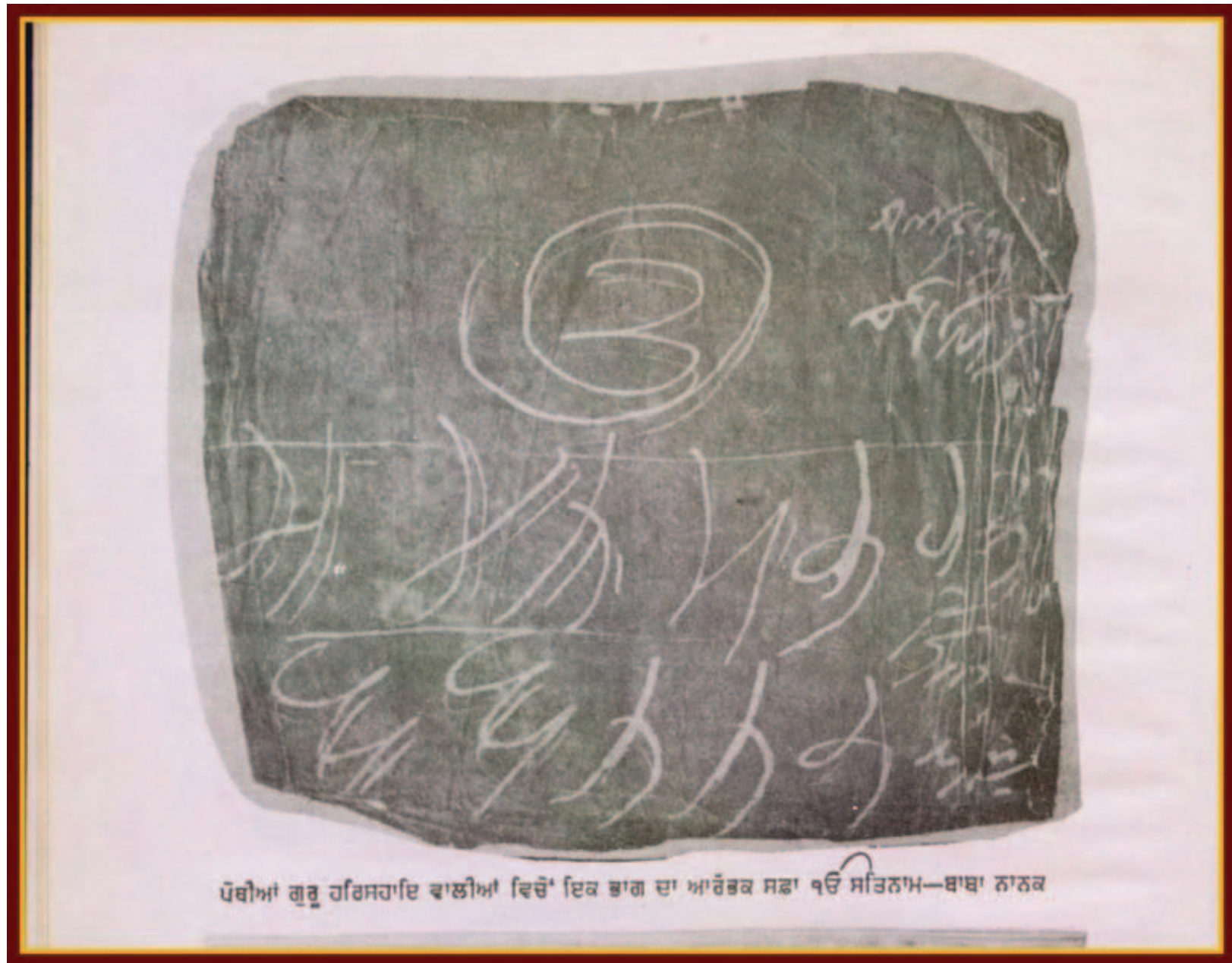We still consider much clinical research as if we were "hunter gathers"- not sharing
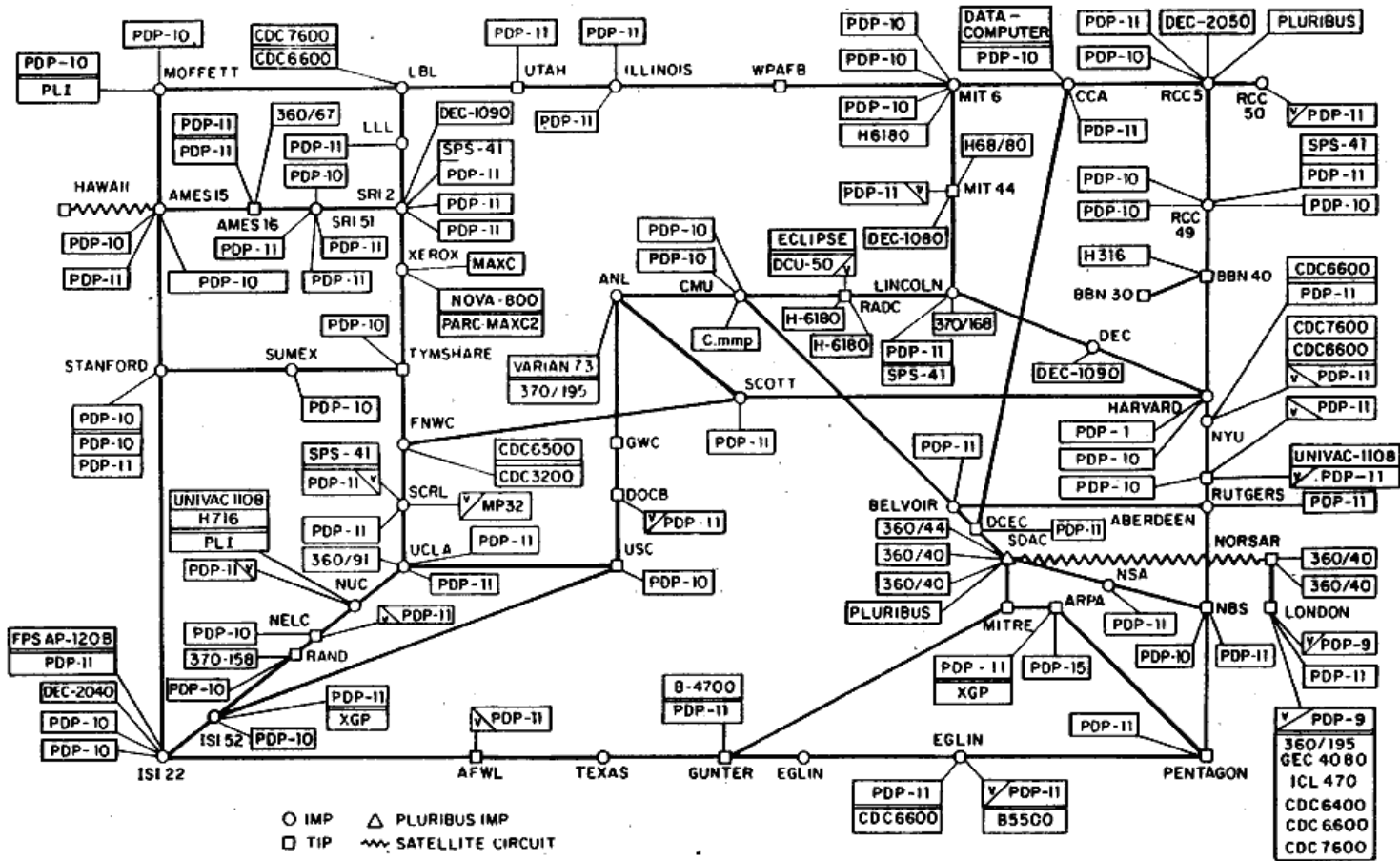
# The Medical Industrial Complex

## Physicians/Scientists

**Mathematical models of disease are not built to be reproduced or versioned by others**

ਪੋਥੀਆਂ ਗੁਰੂ ਹਰਿਸਹਾਇ ਵਾਲੀਆਂ ਵਿਚੋਂ ਇਕ ਭਾਗ ਦਾ ਆਰੰਭਕ ਸਫ਼ਾ ੧ਓ ਸਤਿਨਾਮ—ਬਾਬਾ ਨਾਨਕ

**Lack of standard forms for sharing data
and lack of forms for future rights and consents**

ARPANET LOGICAL MAP, MARCH 1977

sharing as an adoption of common standards..
Clinical   Genomics  Privacy   IP

# Why not share clinical /genomic data and model building in the ways currently used by the software industry (power of tracking workflows and versioning
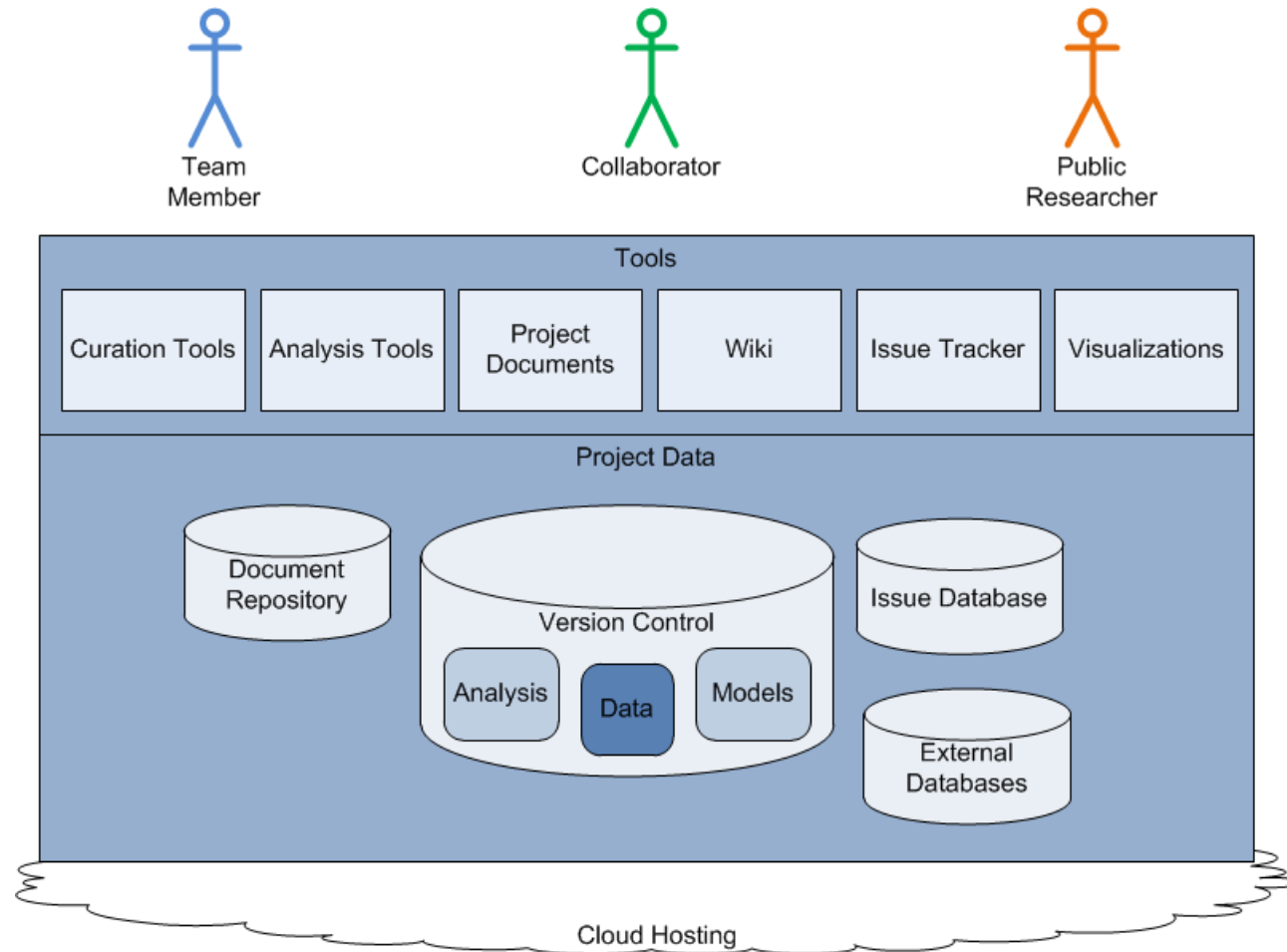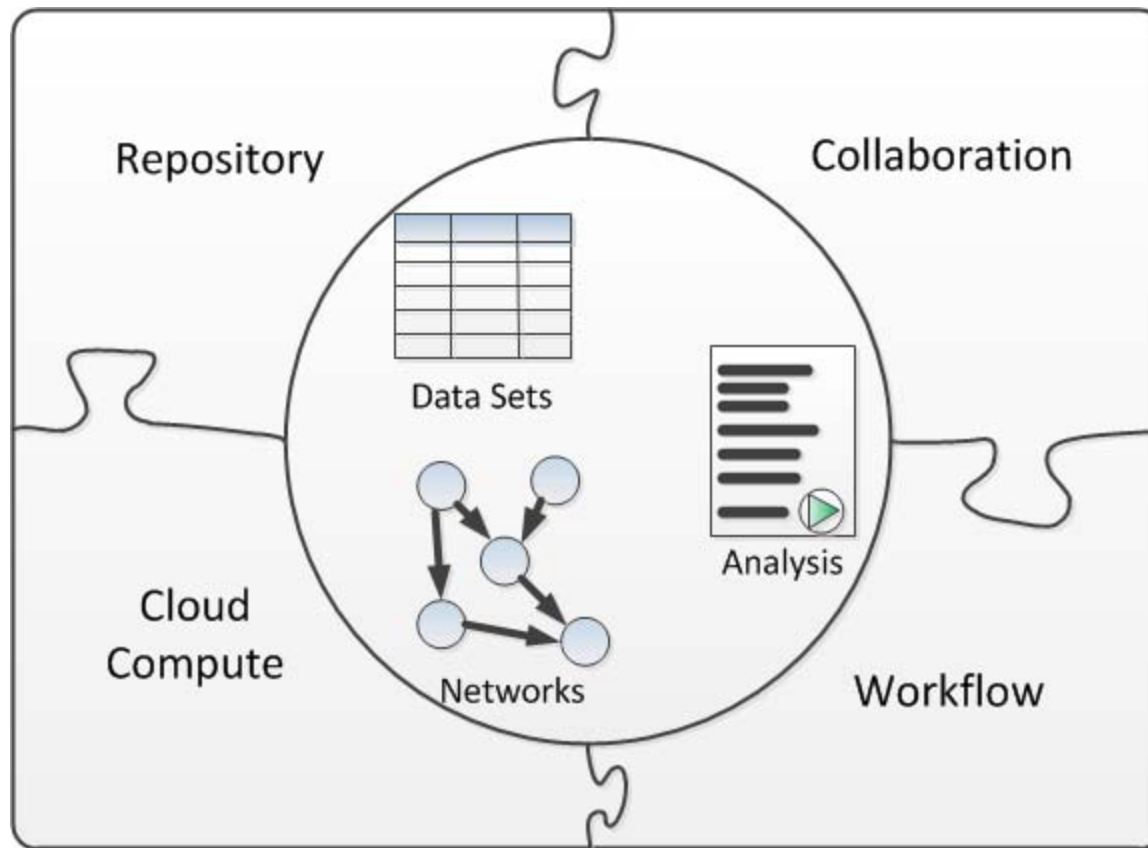
# Evolution of a Software Project
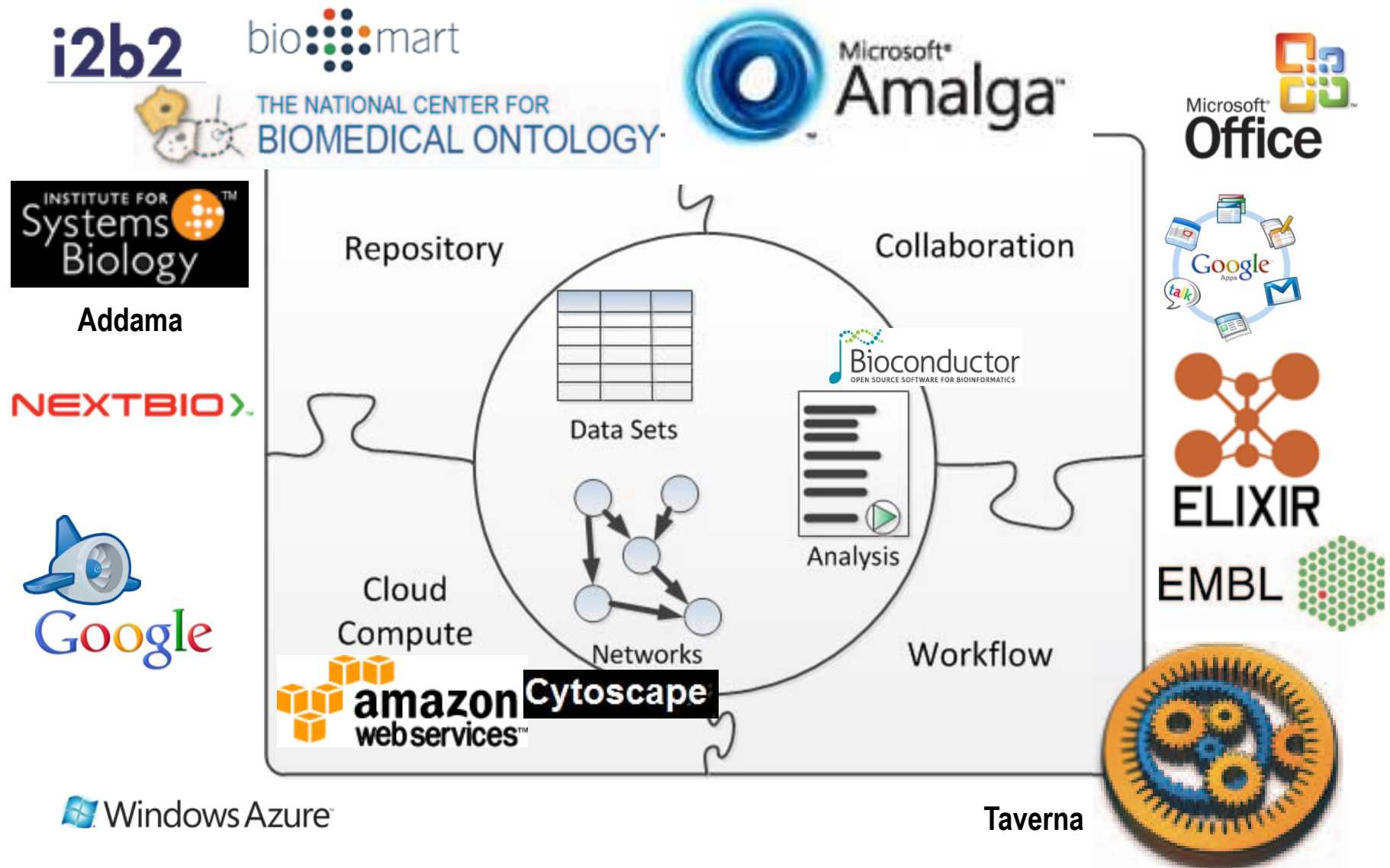
# Evolution of a Biology Project

# Biology Tools Support Collaboration

# Platform Functional Areas

# Potential Supporting Technologies

Aprill 15-16 this will be broadcast LIVE as completely oversubscribed
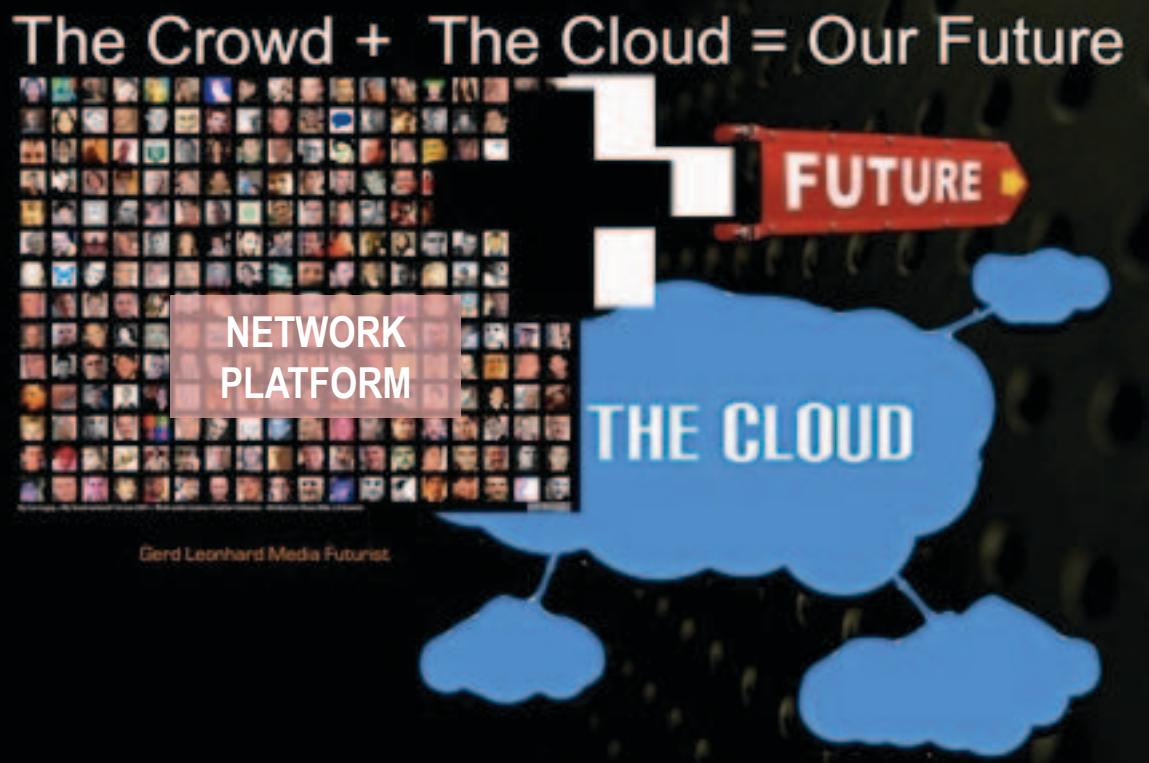Google:   Sage Commons Congress     as time approaches

# Who will build the datasets/ models capable of providing powerful safety and efficacy insights?

Institute

PI

Post Docs

Grad Students



The Crowd + The Cloud = Our Future

FUTURE

NETWORK PLATFORM

THE CLOUD

Gerd Leonhard Media Futurist.

Patients  Physicians  Citizens  Knowledge Experts

BETTER MAPS OF DISEASE


NOT JUST WHAT WE DO BUT HOW WE DO IT


POWER OF BUILDING A PRE-COMPETITIVE
COMMONS FOR EVOLVING
GENERATIVE MODELS OF DISEASE
USING A PUBLIC PRIVATE PARTNERSHIP