# Novel Structures (and Non-Structures) to Facilitate Translational Research

## Integrating layers of omics data models and compute spaces needed to build a "Knowledge Expert"

Stephen Friend MD PhD

Sage Bionetworks (Non-Profit Organization)
Seattle/ Beijing/ Amsterdam

MIT/Whitehead
October 10th, 2011

Why not use data intensive science
to build models of disease

Organizational Structures and Tools
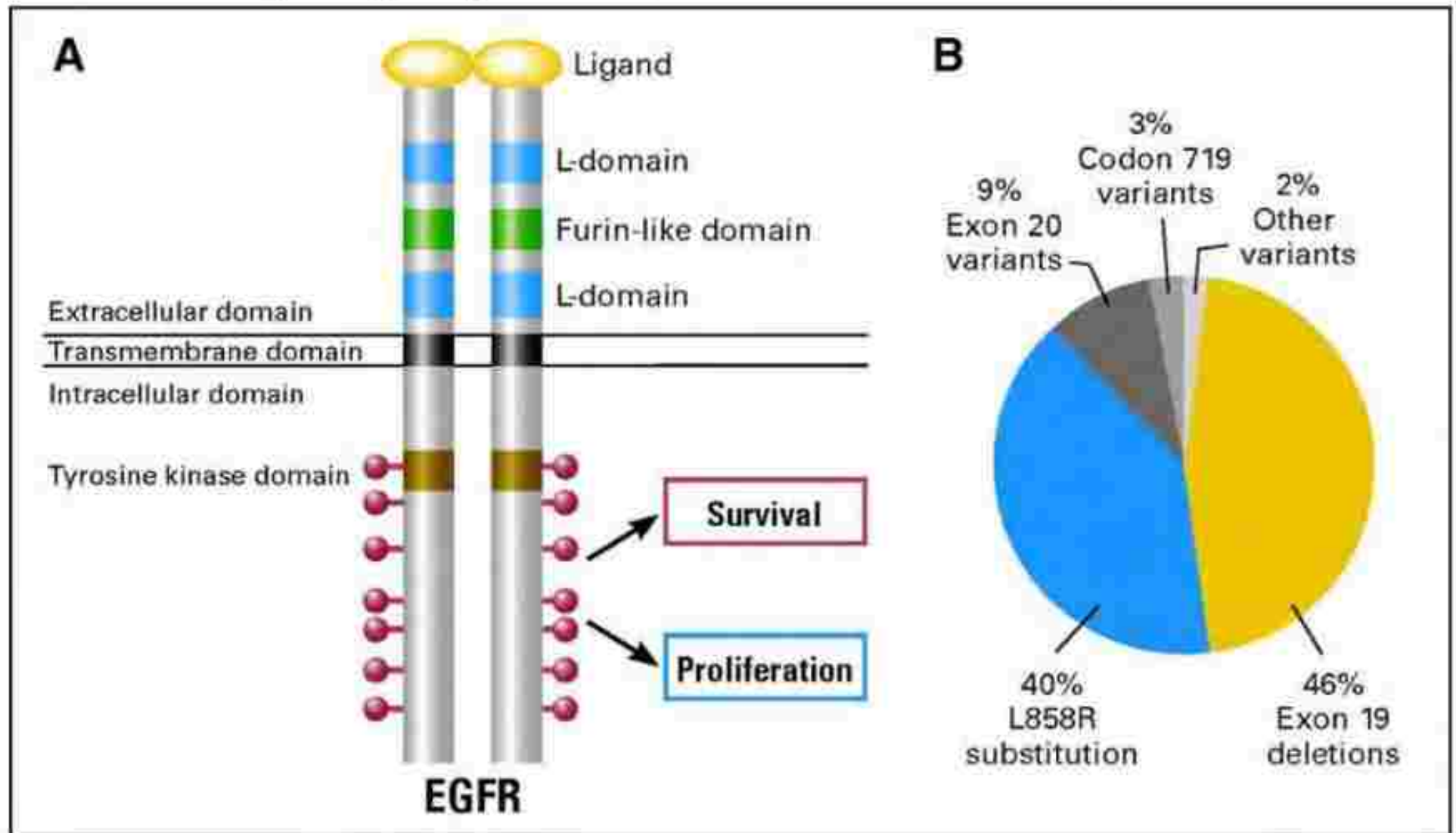
How not What

Six Pilots

Opportunities

# Personalized Medicine 101:
## Capturing Single bases pair mutations = ID of responders



**A**

Ligand

L-domain

Furin-like domain

L-domain

Extracellular domain

Transmembrane domain

Intracellular domain

Tyrosine kinase domain

Survival

Proliferation

EGFR

**B**

9%
Exon 20
variants

3%
Codon 719
variants

2%
Other
variants

40%
L858R
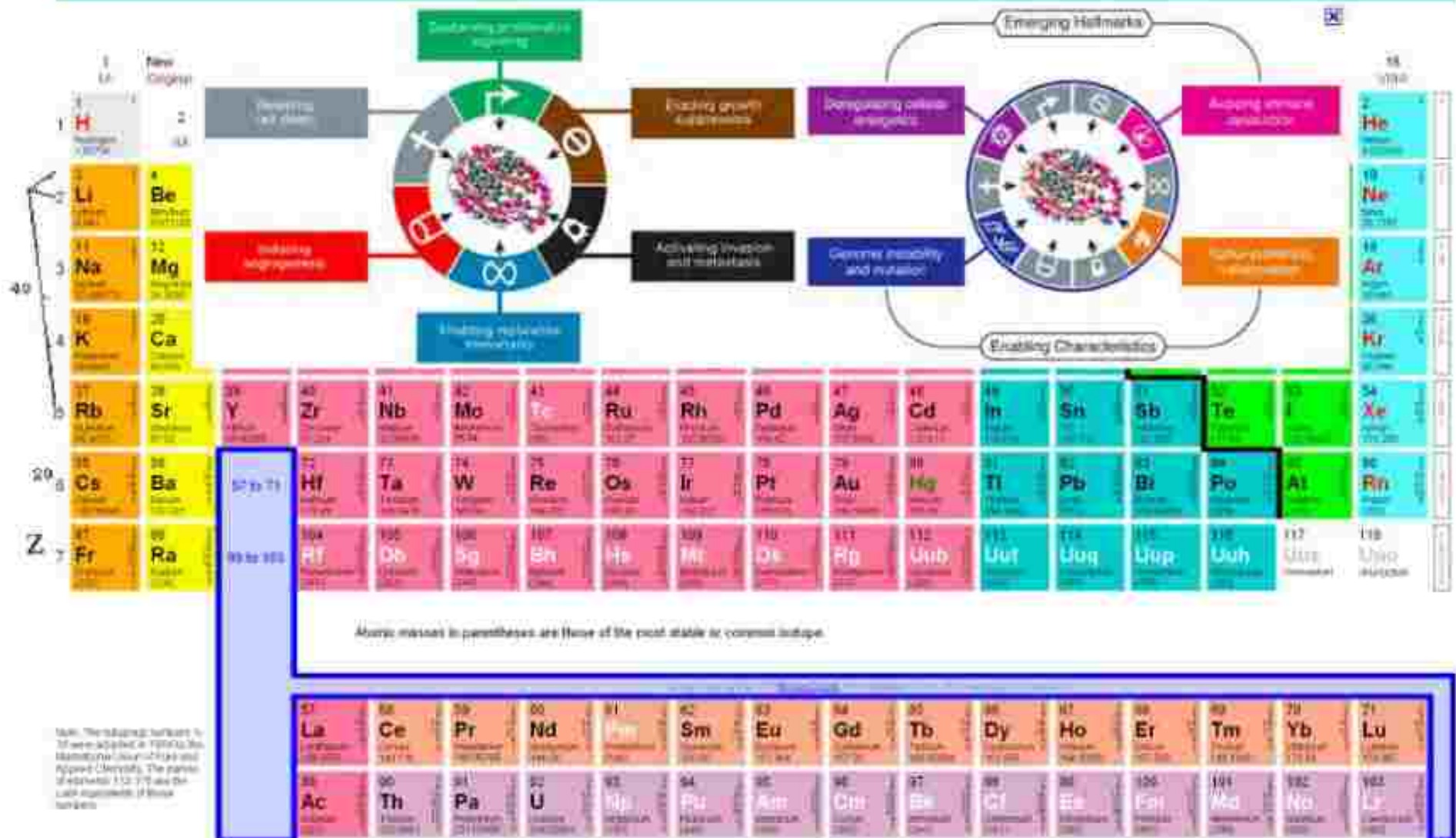substitution

46%
Exon 19
deletions

# Reality: Overlapping Pathways

# The value of appropriate representations/ maps



**Periodic Table of the Elements**

# Science Paradigms

- Thousand years ago:
  science was **empirical**
  *describing natural phenomena*

- Last few hundred years:
  **theoretical** branch
  *using models, generalizations*

$$\left(\frac{\dot{a}}{a}\right)^2 = \frac{4\pi Gp}{3} - K\frac{c^2}{a^2}$$

- Last few decades:
  a **computational** branch
  *simulating complex phenomena*

- Today: **data exploration** (eScience)
  *unify theory, experiment, and simulation*

  – Data captured by instruments
    or generated by simulator

  – Processed by software

  – Information/knowledge stored in computer

  – Scientist analyzes database/files
    using data management and statistics

# "Data Intensive" Science- Fourth Scientific Paradigm

Equipment capable of generating
massive amounts of data

IT Interoperability

Open Information System
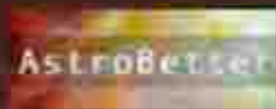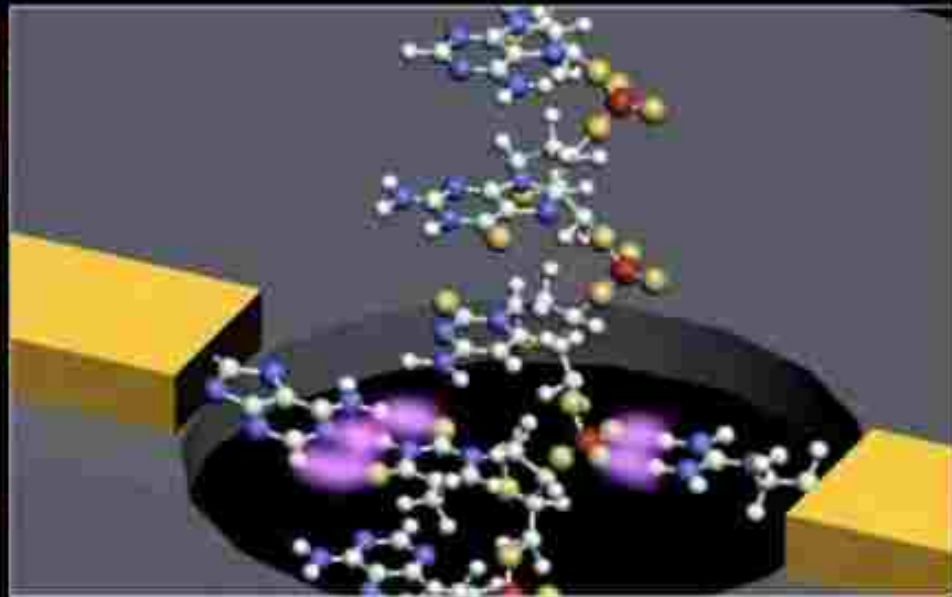
Host evolving Models in a
Compute Space- Knowledge Expert

WHY NOT USE
"DATA INTENSIVE" SCIENCE
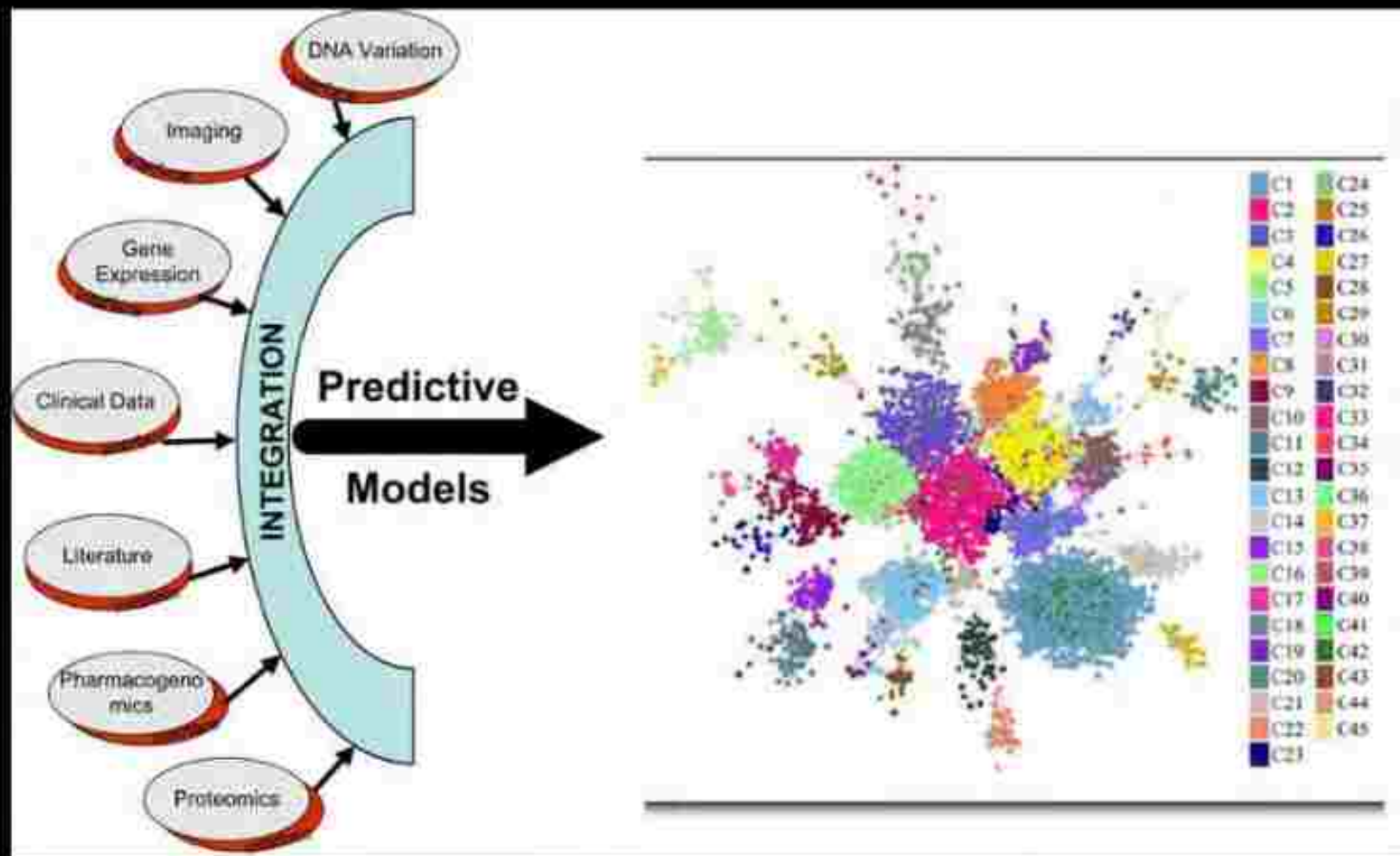TO BUILD BETTER DISEASE MAPS?
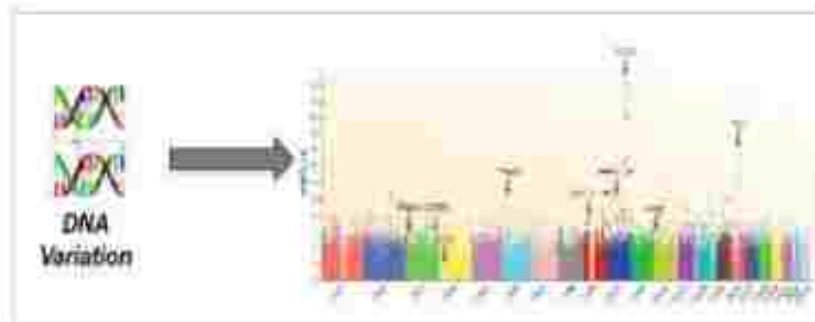
what will it take to understand disease?

DNA RNA PROTEIN (dark matter)

MOVING BEYOND ALTERED COMPONENT LISTS

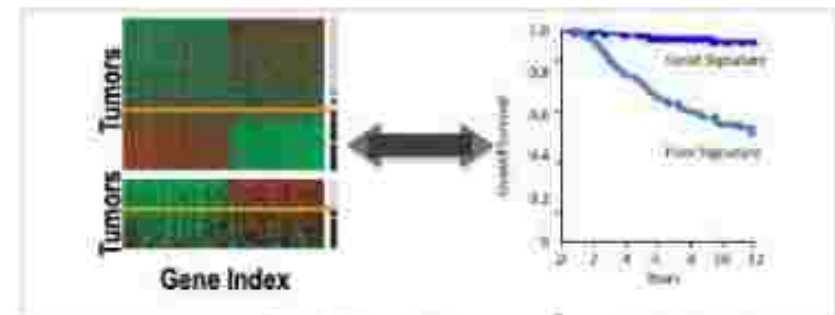# 2002   Can one build a "causal" model?

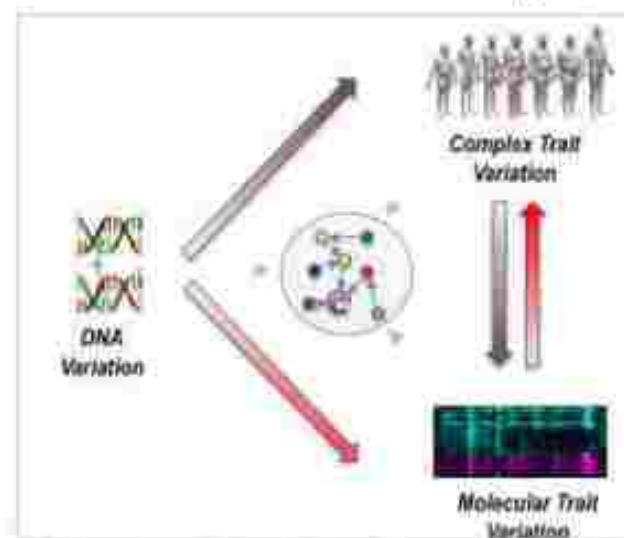# How is genomic data used to understand biology?



**Standard GWAS Approaches**

Identifies Causative DNA Variation
but provides NO mechanism

**Profiling Approaches**

Genome scale profiling provide correlates of disease
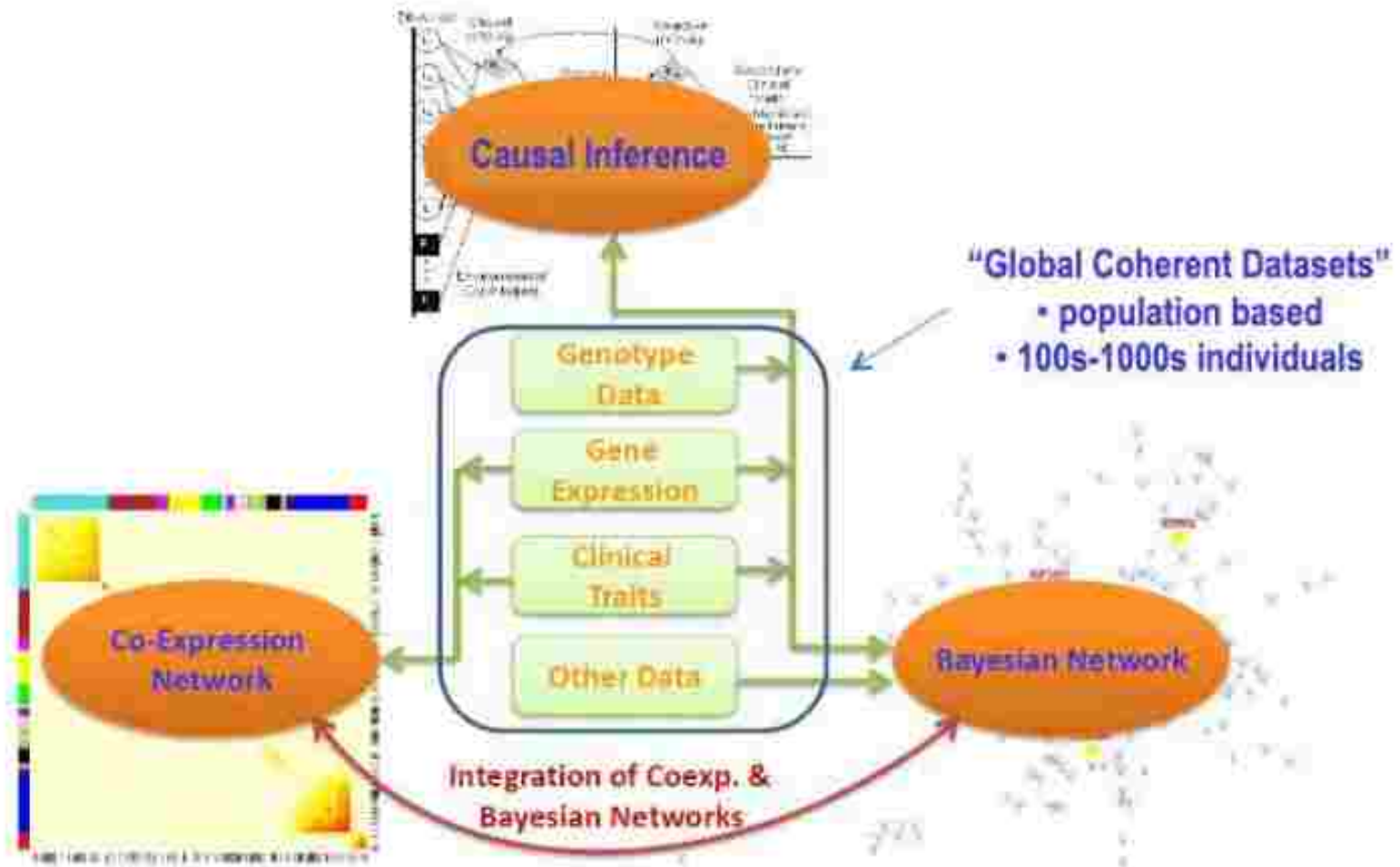➢ Many examples BUT what is cause and effect?

**Integrated Genetics Approaches**

➢ **Provide unbiased view of molecular physiology as it relates to disease phenotypes**

➢ **Insights on mechanism**

➢ **Provide causal relationships and allows predictions**

# 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)

# Association of SNPs at 1p13.3 with Coronary Artery Disease

## SNP rs599839 in the 1p13.3 locus associated with CAD: PSRC1 highlighted as candidate susceptibility gene

## Genomewide Association Analysis of Coronary Artery Disease
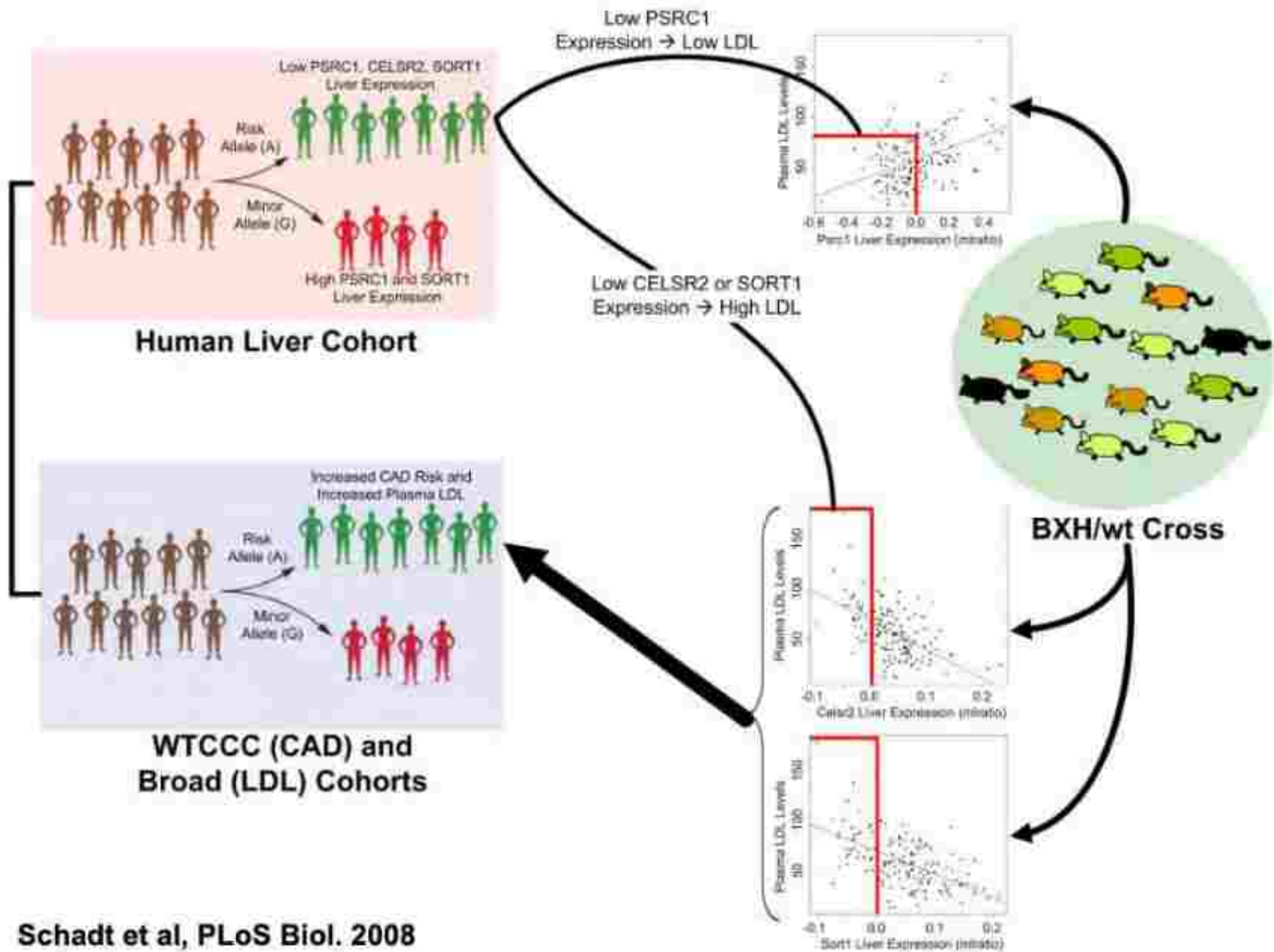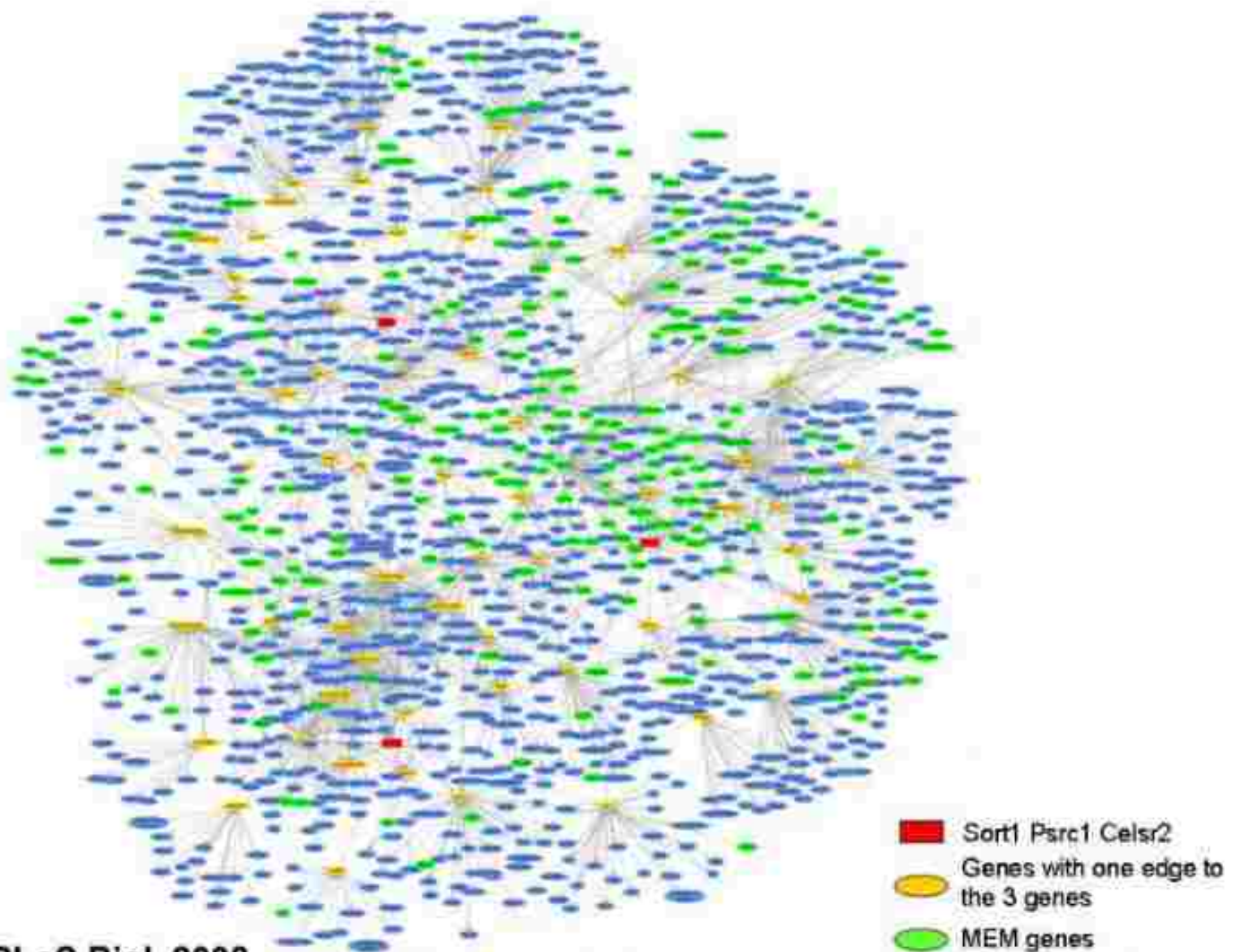
Nilesh J. Samani, F.Med.Sci., Jeanette Erdmann, Ph.D., Alistair S. Hall, F.R.C.P., Christian Hengstenberg, M.D., Massimo Mangino, Ph.D., Bjoern Mayer, M.D., Richard J. Dixon, Ph.D., Thomas Meitinger, M.D., Peter Braund, M.Sc., H.-Erich Wichmann, M.D., Jennifer H. Barrett, Ph.D., Inke R. König, Ph.D., Suzanne E. Stevens, M.Sc., Silke Szymczak, M.Sc., David-Alexandre Tregouet, Ph.D., Mark M. Iles, Ph.D., Friedrich Pahlke, M.Sc., Helen Pollard, M.Sc., Wolfgang Lieb, M.D., Francois Cambien, M.D., Marcus Fischer, M.D., Willem Ouwehand, F.R.C.Path., Stefan Blankenberg, M.D., Anthony J. Balmforth, Ph.D., Andrea Buessler, M.D., Stephen G. Ball, F.R.C.P., Tim M. Strom, M.D., Ingrid Braenne, M.Sc., Christian Gieger, Ph.D., Panos Deloukas, Ph.D., Martin D. Tobin, M.F.P.H.M., Andreas Ziegler, Ph.D., John R. Thompson, Ph.D., and Heribert Schunkert, M.D., for the WTCCC and the Cardiogenics Consortium*
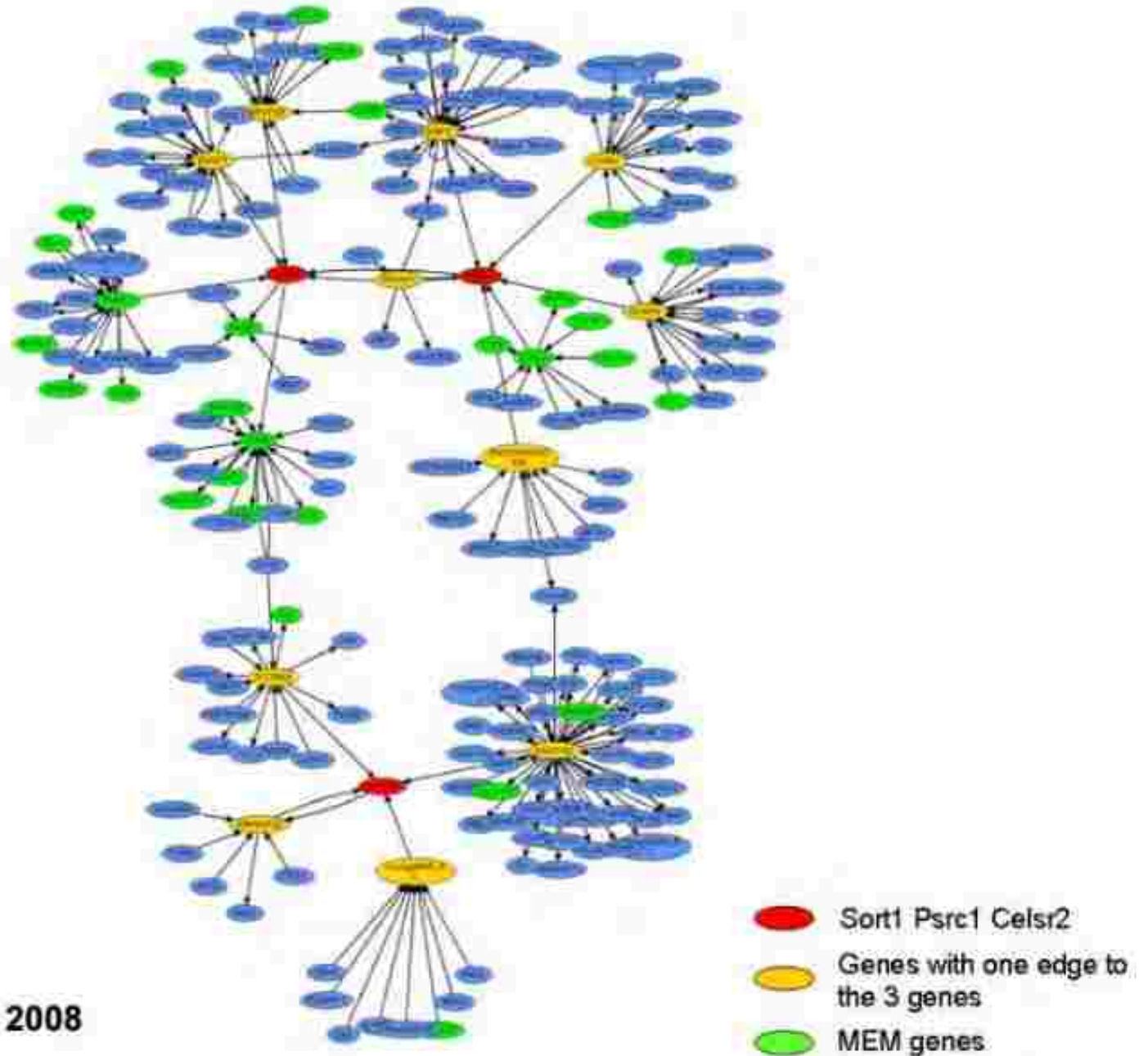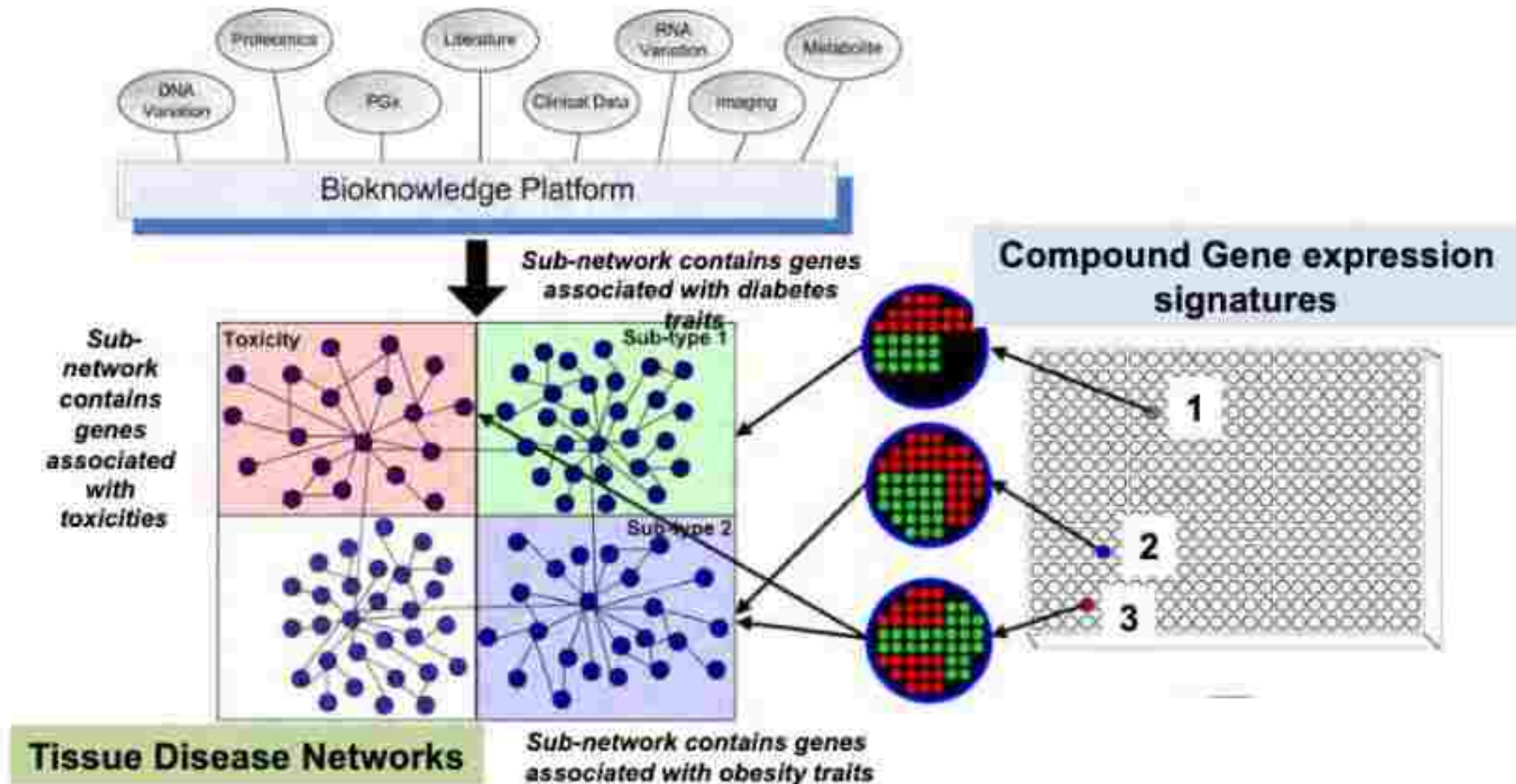
Low PSRC1, CELSR2, SORT1 Liver Expression

Risk Allele (A)

Minor Allele (G)

High PSRC1 and SORT1 Liver Expression

**Human Liver Cohort**

Low PSRC1 Expression → Low LDL

Low CELSR2 or SORT1 Expression → High LDL

Increased CAD Risk and Increased Plasma LDL

Risk Allele (A)

Minor Allele (G)

**WTCCC (CAD) and Broad (LDL) Cohorts**

**BXH/wt Cross**

**Schadt et al, PLoS Biol. 2008**

# Mouse network around Sort1, Psrc1, and Celsr2



Sort1 Psrc1 Celsr2
Genes with one edge to the 3 genes
MEM genes

**Schadt et al, PLoS Biol. 2008**

# Human network around Sort1, Psrc1, and Celsr2



**Schadt et al, PLoS Biol. 2008**

Legend:
- Sort1 Psrc1 Celsr2
- Genes with one edge to the 3 genes
- MEM genes

# Map compound signatures to disease networks



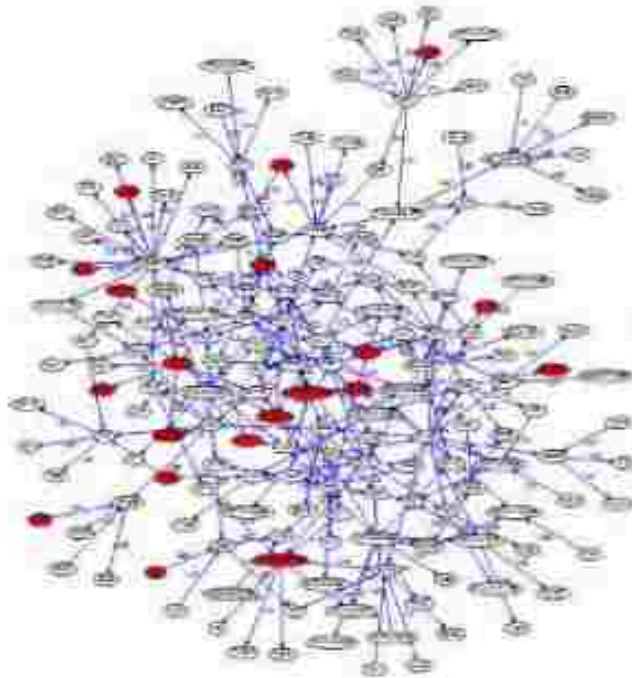**Tissue Disease Networks** — Sub-network contains genes associated with obesity traits

**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
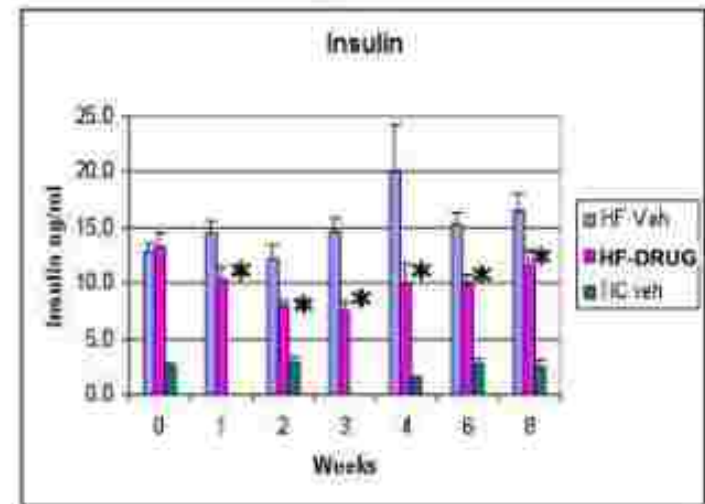
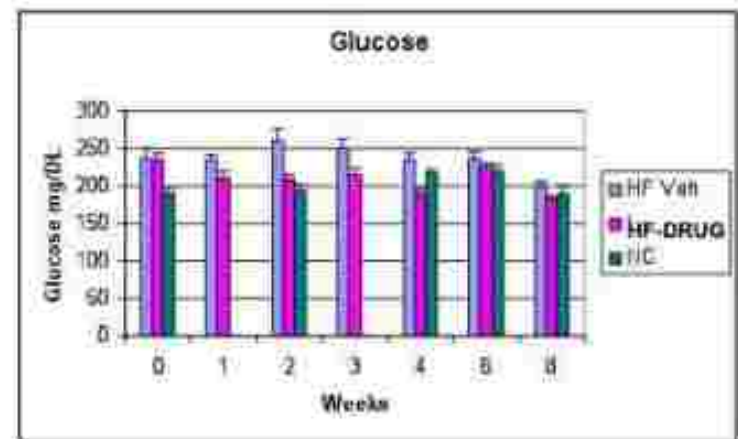# Case Study – Target A/Drug B

**NO CELL DYNAMICS NEEDED**

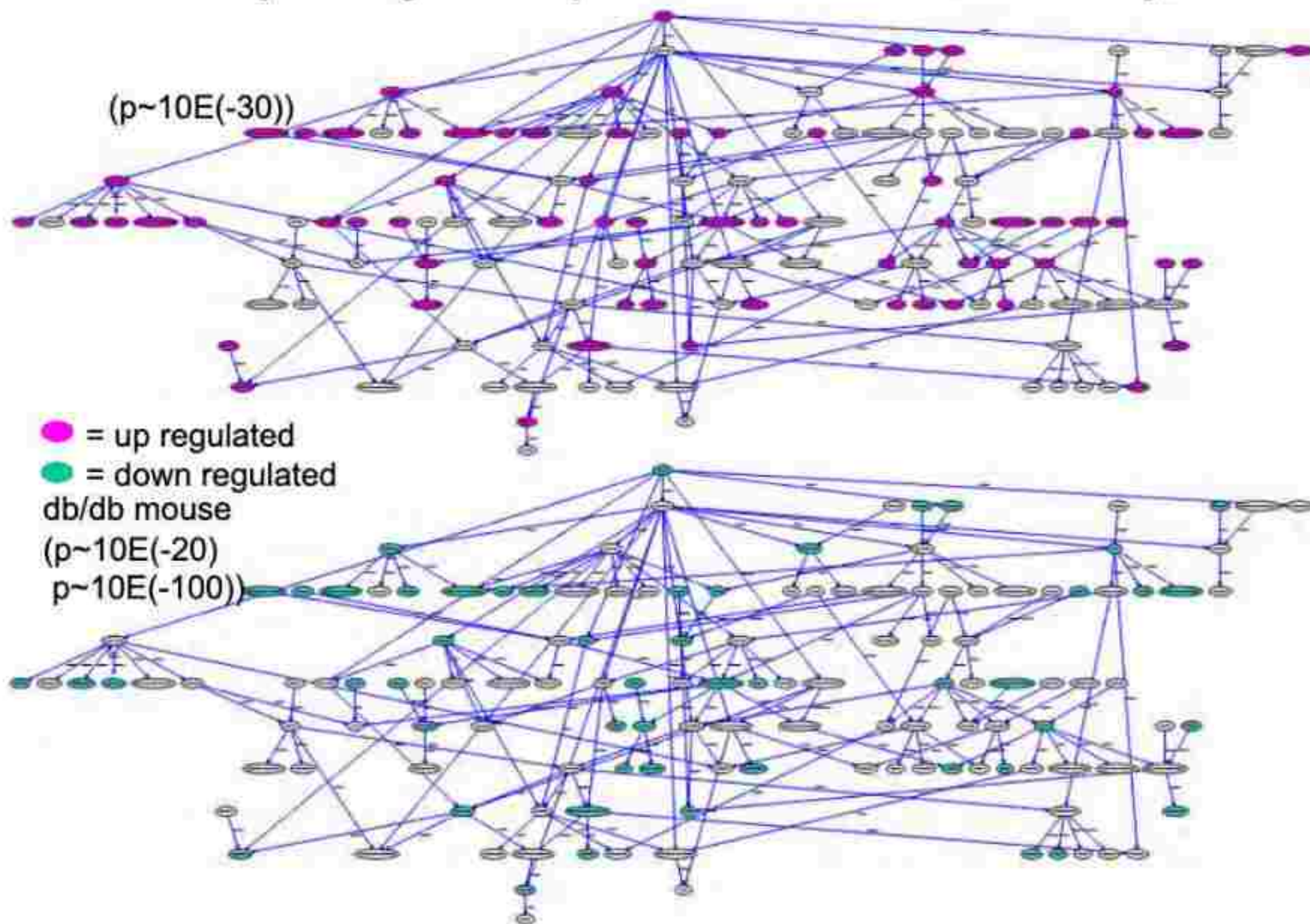**Identified compound whose signature significantly intersected with Islet module**
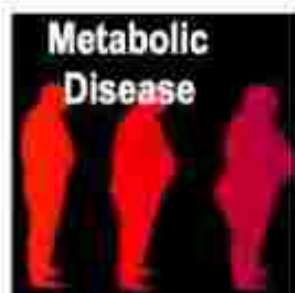
**Fasting Insulin**



**Fasting Glucose**



- **Test carried out in a Diet-Induced Obesity model on the B6 background**
  - **Model for obesity and insulin resistance**
- **Animals treated with compound over an 8 week interval, starting at 8 weeks of age**
- *No significant Adverse Events in 30 day human clinical trial for another indication*

# Our ability to integrate compound data into our network analyses



$(p\sim10E(-30))$

● = up regulated
● = down regulated
db/db mouse
$(p\sim10E(-20)$
$p\sim10E(-100))$

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

- **>80 Publications from Rosetta Genetics Group (~30 scientists) over 5 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.**

# List of Influential Papers in Network Modeling

Validation of candidate causal genes for obesity that affect shared metabolic pathways and networks

## Integrative Modeling Defines the Nova Splicing-Regulatory Network and Its Combinatorial Controls

## The transcriptional network for mesenchymal transformation of brain tumours

## Variations in DNA elucidate molecular networks that cause disease

## Rewiring of Genetic Networks in Response to DNA Damage

Although cellular behaviors are dynamic, the networks that govern these behaviors have been mapped primarily as static snapshots.

## An Atlas of Combinatorial Transcriptional Regulation in Mouse and Man

## Network-Based Elucidation of Human Disease Similarities Reveals Common Functional Modules Enriched for Pluripotent Drug Targets

ANALYSIS

Genome-wide identification of post-translational modulators of transcription factor activity in human B cells

# LETTER

## A *trans*–acting locus regulates an anti–viral expression network and type 1 diabetes risk

> 50 network papers
> http://sagebase.org/research/resources.php

# The way we like to think:



# The way it is:



(Eric Schadt)

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**

**Discovery Platform**

Sagebase.org

# Board of Directors- Sage Bionetworks



**Lee Hartwell**

**Ex President FHCRC**
**Co-Founder Rosetta**

**Hans Wizgell**

**ExPresident Karolinska**
**Head SAB Rosetta**

**WangJun**

**Executive Director**
**BGI**

**Jeff Hammerbacher**

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

## Sage Bionetworks Collaborators

- ## Pharma Partners
  - Merck, Pfizer, Takeda, Astra Zeneca, Amgen, Johnson &Johnson
- ## Foundations
  - Kauffman CHDI, Gates Foundation
- ## Government
  - NIH, LSDF
- ## Academic
  - Levy (Framingham)
  - Rosengren (Lund)
  - Krauss (CHORI)
- ## Federation
  - Ideker, Califarno, Butte, Schadt

**PLATFORM**
Sage Platform and Infrastructure Builders-
( Academic Biotech and Industry IT Partners...)

**PILOTS= PROJECTS FOR COMMONS**
Data Sharing Commons Pilots-
(Federation, CCSB, Inspire2Live....)

**NEW TOOLS**
Data Tool and Disease Map Generators-
(Global coherent data sets, Cytoscape,
Clinical Trialists, Industrial Trialists, CROs...)

**NEW MAPS**
Disease Map and Tool Users-
( Scientists, Industry, Foundations, Regulators...)

**RULES AND GOVERNANCE**
Data Sharing Barrier Breakers-
(Patients Advocates, Governance
and Policy Makers,  Funders...)

## 775,388 people hosting over 2,161,922 git repositories

jQuery, reddit, Sparkle, curl, Ruby on Rails, node.js, ClickToFlash, Erlang/OTP, CakePHP, Redis, and many more

**twitter**    **facebook**    **rackspace** HOSTING    **digg**    **YAHOO!**    **shopify**    **EMI**    six apart

# git /ɡɪt/

Git is an extremely fast, efficient, distributed version control system ideal for the collaborative development of software.

# git·hub /ɡɪt hʌb/

GitHub is the best way to collaborate with others. Fork, send pull requests and manage all your **public** and **private** git repositories.

## Plans, Pricing and Signup
Unlimited public repositories are free!

Free public repositories, collaborator management, issue tracking, wikis, downloads, code review, graphs and much more...

## Team management

**30 seconds** is give people access to code. No SSH key required. Activity feeds keep you updated on progress.

More about collaboration

## Code review

Comment on changes, track issues, compare branches, send pull requests and merge forks.

More about code review

## Reliable code hosting

We spend all day and night making sure your repositories are **secure**, **backed up** and **always available**.

More about code hosting

## Open source collaboration

Participate in the most important open source community in the world today—online or at one of our meetups.

More about our community

Powered by the Dedicated Servers and Cloud Computing of Rackspace Hosting®

# 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

# Biology Tools Support Collaboration

# Potential Supporting Technologies



tranSMART

# Platform for Modeling

Repository

Collaboration

Multi-dimensional Data Set

SYNAPSE

Analysis Methods

Cloud Compute

Network Models

Tools & Workflow

# sage bionetworks synapse project

### Watch What I Do, Not What I Say



### Reduce, Reuse, Recycle



### My Other Computer is Amazon

### Most of the People You Need to Work with Don't Work with You

# Synapse machine learning infrastructure for method comparison



-- Implement customTrain() and customPredict() functions
-- Everything else handled in standardized workflow (performance evaluation, biomarker outputs, evaluation against other methods, loading of different datasets, etc).

INTEROPERABILITY

Genome Pattern
CYTOSCAPE
tranSMART
I2B2

hunter gathers- not sharing

# TENURE

# FEUDAL STATES

**Optimal autonomous state sizes.**
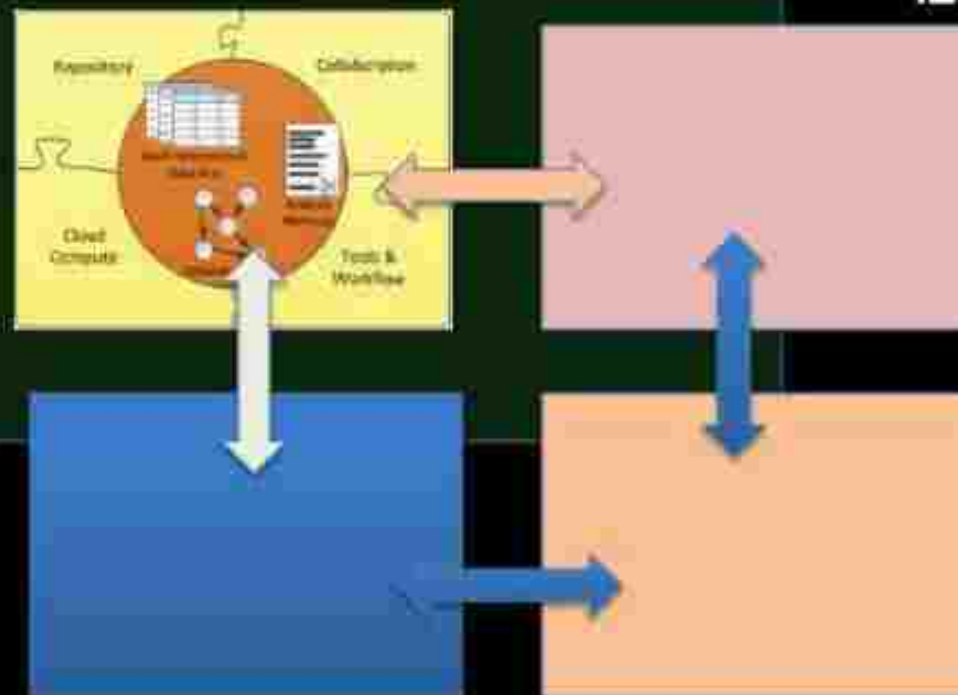These examples show autonomous states of different rank that the balance external defense and strength against the internal tensions and ambition of the various feudal actors.

**Duchy**: a number of Baronies held in feudal subservience to a single Barony. The subjugated baronies are now fiefs ruled by the local Baron, their knights are loyal to him as well. The ruling Baron is elevated to the rank of Duke.

**Barony**: a number of city-states held in feudal subservience to a single city-state. The subjugated cities are now fiefs ruled by knights. The knight of the ruling city-state is elevated to the rank of Baron.

**City-state**: ruled by a Knight. The smallest and stablest form of state.

Note: all Knights, Barons and Dukes are independent actors who may vie with each other or even their feudal overlord for more power, land, coin, etc.

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

**Assumption that genetic alterations in human conditions should be owned**

ਪੰਚੀਆ ਗੁਰੂ ਹਰਿਸਹਾਇ ਵਾਲੀਆ ਵਿਚੋਂ ਇਕ ਭਾਗ ਦਾ ਆਰੰਭਕ ਸਫਾ ੧ਓ ਸਤਿਨਾਮ—ਬਾਬਾ ਨਾਨਕ

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

**Publication Bias- Where can we find the (negative) clinical data?**

ARPANET LOGICAL MAP, MARCH 1977

(PLEASE NOTE THAT WHILE THIS MAP SHOWS THE HOST POPULATION OF THE NETWORK ACCORDING TO THE BEST INFORMATION OBTAINABLE, NO CLAIM CAN BE MADE FOR ITS ACCURACY.)

NAMES SHOWN ARE IMP NAMES, NOT (NECESSARILY) HOST NAMES

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

# Six Pilots at Sage Bionetworks

CTCAP
Non-Responders
Arch2POCM
The Federation
Portable Legal Consent
Sage Congress Project

# CTCAP

# Clinical Trial Comparator Arm Partnership "CTCAP"
# Strategic Opportunities For Regulatory Science Leadership and Action

**FDA**
**September 27, 2011**

# Clinical Trial Comparator Arm Partnership (CTCAP)
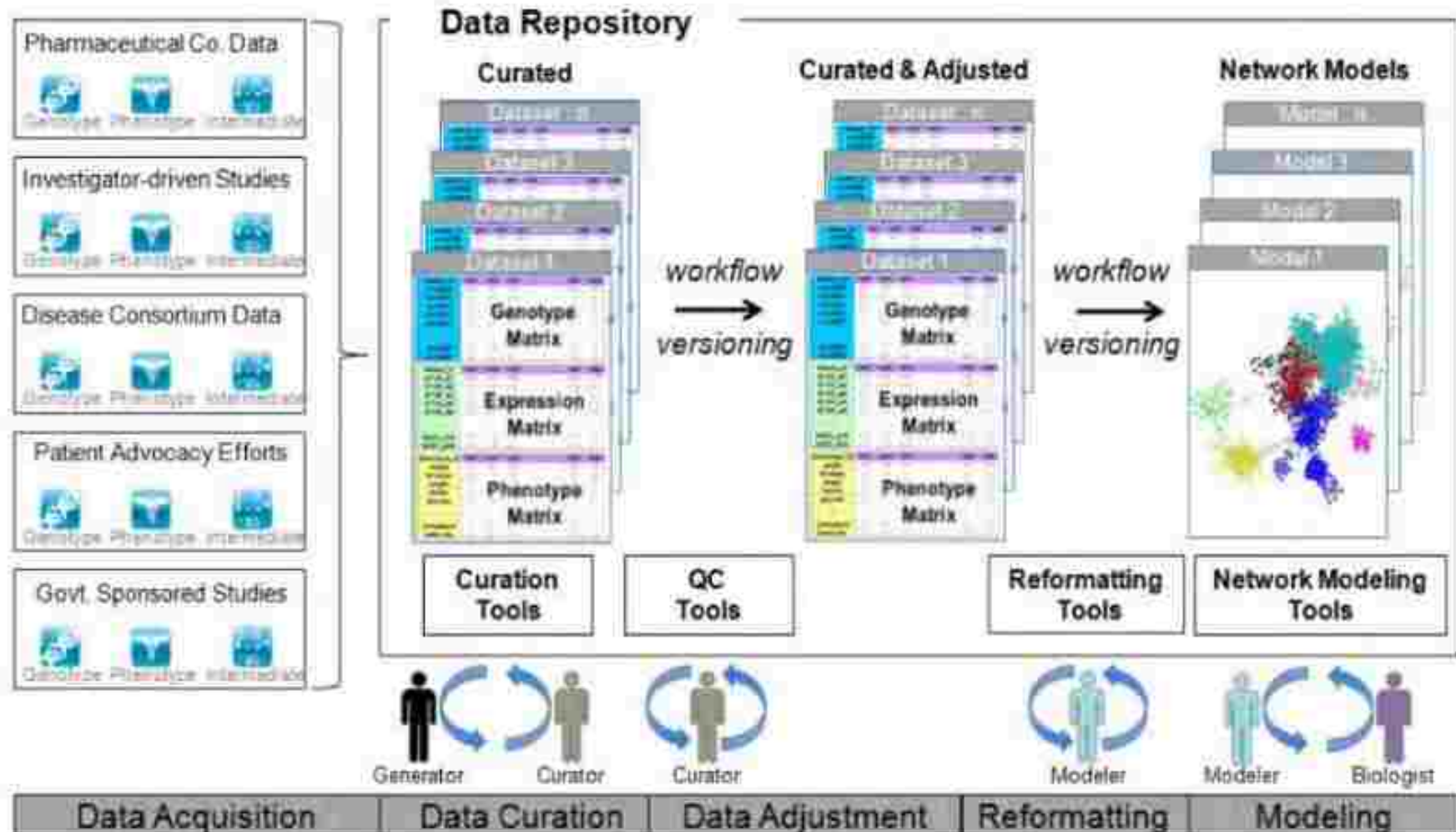
Partnering for Cures | A Faster Cures Meeting

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.

Started Sept 2010

# Shared clinical/genomic data sharing and analysis will maximize clinical impact and enable discovery

# Non-Responders Project

To identify Non-Responders to approved
Oncology drug regimens in order to improve
outcomes, spare patients unnecessary toxicities
from treatments that have no benefit to them, and
reduce healthcare costs

# The Non-Responder Cancer Project Leadership Team

**Stephen Friend, MD, PhD**
President and Co-Founder of
Sage Bionetworks, Head of
Merck Oncology 01-08,
Founder of Rosetta
Inpharmatics 97-01, co-
Founder of the Seattle Project

**Todd Golub, MD**
Founding Director Cancer Biology
Program Broad Institute, Charles Dana
Investigator Dana-Farber Cancer
Institute, Professor of Pediatrics Harvard
Medical School, Investigator, Howard
Hughes Medical Institute

**Garry Nolan, PhD**
Professor, Baxter Laboratory of Stem
Cell Biology, Department of Microbiology
and Immunology, Stanford University
Director, Proteomics Center at Stanford
University

**Richard Schilsky, MD**
Chief, Hematology- Oncology, Deputy
Director, Comprehensive Cancer
Center, University of Chicago; Chair,
National Cancer Institute Board of
Scientific Advisors; past-President
ASCO, past Chairman CALGB clinical
trials group

# The Non-Responder Project is an international initiative with funding for 6 initial cancers anticipated from both the public and private sectors
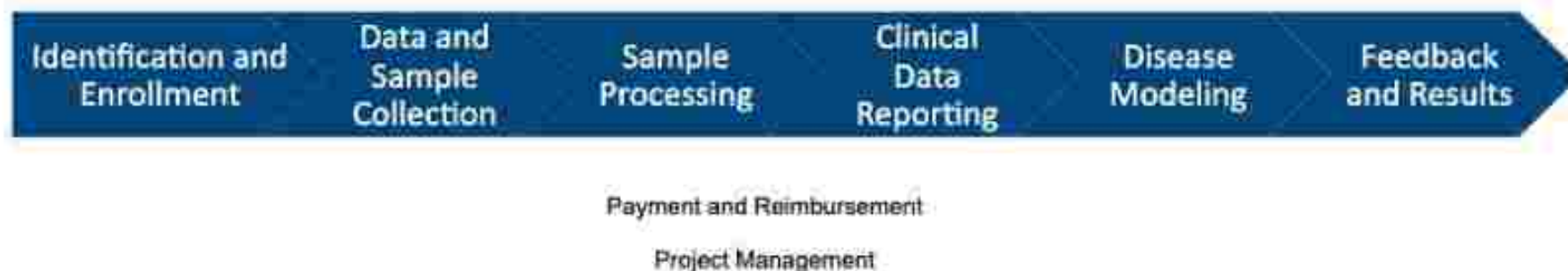
| GEOGRAPHY | United States | | | | China | |
|---|---|---|---|---|---|---|
| **TARGET CANCER** | Ovarian | Renal | Breast | AML | Colon | Lung |
| **FUNDING SOURCE** | **Seeking private sector and philanthropic funding for prospective studies** | | | Retrospective study; likely to be funded by the Federal Government | Funded by the Chinese government and private sector partners | |

5

# For each tumor-type, the non-responder project will follow a common workflow, with patient identification and sample collection the most variable across studies

**Non-Responder Project Workflow**

Identification and enrollment, and data and sample collection may differ by tumor-type

The remaining parts of the study will be largely similar, and potentially shared, across all projects

| Identification and Enrollment | Data and Sample Collection | Sample Processing | Clinical Data Reporting | Disease Modeling | Feedback and Results |
|---|---|---|---|---|---|

Payment and Reimbursement

Project Management

# A consortium of collaborators has been constructed to execute the non-responder project



Physicians & AMCs

Patient Advocacy Groups

# A consortium of collaborators has been constructed to execute the non-responder project (continued)

# Arch2POCM

## Restructuring Drug Discovery

## How to potentially De-Risk High-Risk Therapeutic Areas

# What is the problem?

- Regulatory hurdles too high?
- Low hanging fruit picked?
- Payers unwilling to pay?
- Genome has not delivered?
- Valley of death?
- Companies not large enough to execute on strategy?
- Internal research costs too high?
- Clinical trials in developed countries too expensive?

In fact, all are true but none is the real problem

# What is the problem?

We need to rebuild the drug discovery process so that we better understand disease biology before testing proprietary compounds on sick patients

# The Precompetitive Space: Time to Move the Yardsticks

Thea Norman,[1] Aled Edwards,[2] Chas Bountra,[3] Stephen Friend[**]

Industry, government, patient advocacy groups, public funders, and academic thought leaders met in Toronto, Canada, to set into motion an initiative that addresses some of the scientific and organizational challenges of modern therapeutics discovery. What emerged from the meeting was a public-private partnership that seeks to establish proof of clinical mechanism (POCM) for selected "pioneer" disease targets using lead compounds—all accomplished in the precompetitive space. The group will reconvene in April 2011 to create a business plan that specifies the generation of two positive POCM results per year.



## MEETING REPORT

### CROWDSOURCING

# Leveraging Crowdsourcing to Facilitate the Discovery of New Medicines

Thea C. Norman,[1] Chas Bountra,[2] Aled M. Edwards,[3] Keith R. Yamamoto,[4] Stephen H. Friend[5*]

Gloomy predictions about the future of pharma have forced the industry to investigate alternative models of drug discovery. Public-private partnerships (PPPs) have the potential to revitalize the discovery and development of first-in-class therapeutics. The new PPP Arch2POCM hopes to foster biomedical innovation through precompetitive validation of pioneer therapeutic targets for human diseases. In this meeting report, we capture the most exciting insights garnered from the April 2011 Arch2POCM conference.

When useful knowledge exists in companies of all sizes and also in universities, non-profits and individual minds, it makes sense to orient your innovation efforts to accessing, building upon and integrating that external knowledge into useful products and services.
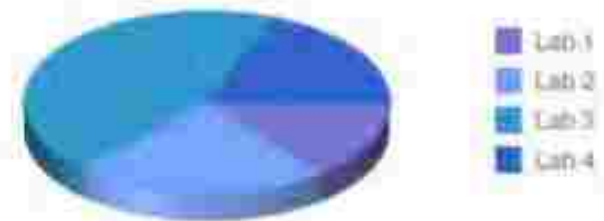
# The Federation

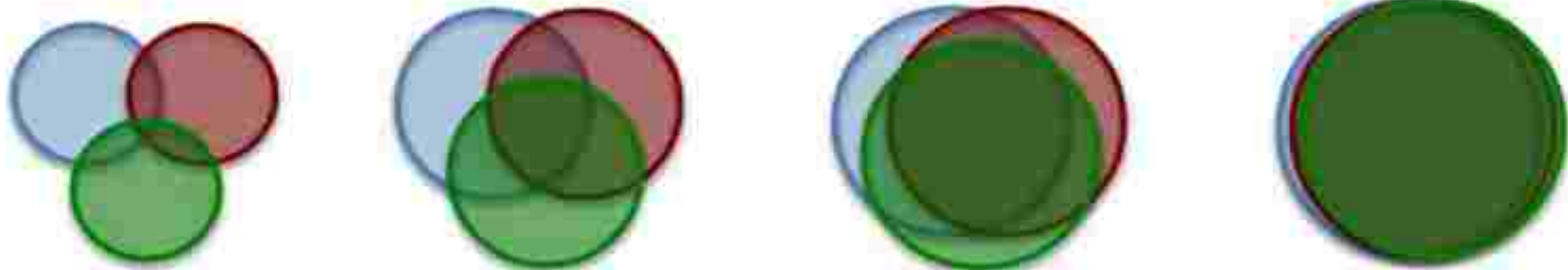# How can we accelerate the pace of scientific discovery?

2008       2009       2010       2011



Collaboration 1.0

**Ways to move beyond "traditional" collaborations?**

**Intra-lab vs Inter-lab Communication**

**Colrain/ Industrial PPPs Academic Unions**



Lab 1
Lab 2
Lab 3
Lab 4

# THE FEDERATION
## Butte    Califano    Friend    Ideker    Schadt

**vs**

Social Networking

# Rules of the game:
## transparency & trust

- Shared data tools models and prepublications
- Conflict of interests
- Intellectual property
- Authorship

Sage

# sage federation:
# human aging project

**Sage**

Justin Guinney
Stephen Friend*

**UC San Diego** SCHOOL of MEDICINE

Greg Hannum
Januz Dutkowski
Trey Ideker*
Kang Zhang*

**COLUMBIA UNIVERSITY**

Mariano Alvarez
Celine Lefebrev
Andrea Califano*

sage federation:
what is the impact of disease/environment on "biological age" ?

# human aging:
## predicting bioage using whole blood methylation

- Independent training (n=170) and validation (n=123) Caucasian cohorts
- 450k Illumina methylation array
- Exom sequencing
- Clinical phenotypes: Type II diabetes, BMI, gender...

**Training Cohort: San Diego (n=170)**

RMSE=3.35

**Validation Cohort: Utah (n=123)**

RMSE=5.44

sage federation:
model of biological age

$$\text{Bioage} = f(M) = Age + \sum_{J} \alpha_j C_j + \epsilon$$

$$\text{Differential Bioage} = f(M) - Age = \sum_{J} \alpha_j C_j + \epsilon$$

Faster Aging

Slower Aging

Clinical Association
- Gender
- BMI
- Disease
Genotype Association
Gene Pathway Expression

Age Differential

Predicted Age (liver expression)

Chronological Age (years)

# Reproducible science==shareable science

**Sweave: combines programmatic analysis with narrative**

**Dynamic generation of statistical reports using literate data analysis**



Sweave.Friedrich Leisch. Sweave: Dynamic generation of statistical reports using literate data analysis. In Wolfgang Härdle and Bernd Rönz,editors, Compstat 2002 – Proceedings in Computational Statistics,pages 575-580. Physica Verlag, Heidelberg, 2002. ISBN 3-7908-1517-9

Federated Aging Project :
Combining analysis + narrative

# Portable Legal Consent

## (Activating Patients)

John Wilbanks

http://sagebase.org/getconsent

⦿ I want to participate in public genetic research

○ I have data that I want to contribute to public genetic research

○ I have provided biological samples and want to retain rights to my data

http://sagebase.org/getconsent/grantrights

these are the rights you are granting
to qualified researchers :

☑ Right to do research with my data

☑ Right to redistribute my data

☑ Right to publish the results of research from my data

☑ Right to commercialize products derived from research on my data

*all boxes must be checked to move forward
in the consent process*

Next

http://sagebase.org/getconsent/obligations

# behaviors you can request of the researchers who use your data

☐ Do not attempt to re-identify me.

☐ Share new data with others as I have shared with you.

☐ Share your research with the public under open access terms.

these are obligations we will impose on researchers through terms of use. violators will not be allowed to access the commons again.

( Next )

http://sagebase.org/getconsent/affirmconsent

all boxes must be checked to create informed consent.

☑ I understand the uncertainty and risk of public genetic research.

☑ I provide consent for my data to be used in public genetic research

☑ I understand that although I can withdraw at any time, I cannot withdraw data that has already been distributed.

I GIVE CONSENT

I'M NOT SURE

# Sage Congress Project
## April 20 2012

# RA
## Parkinson's
## Asthma

(Responders Competitions)