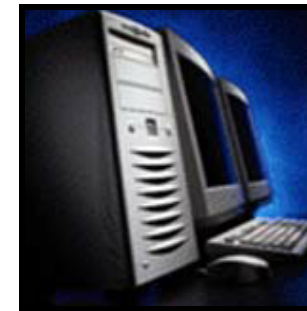
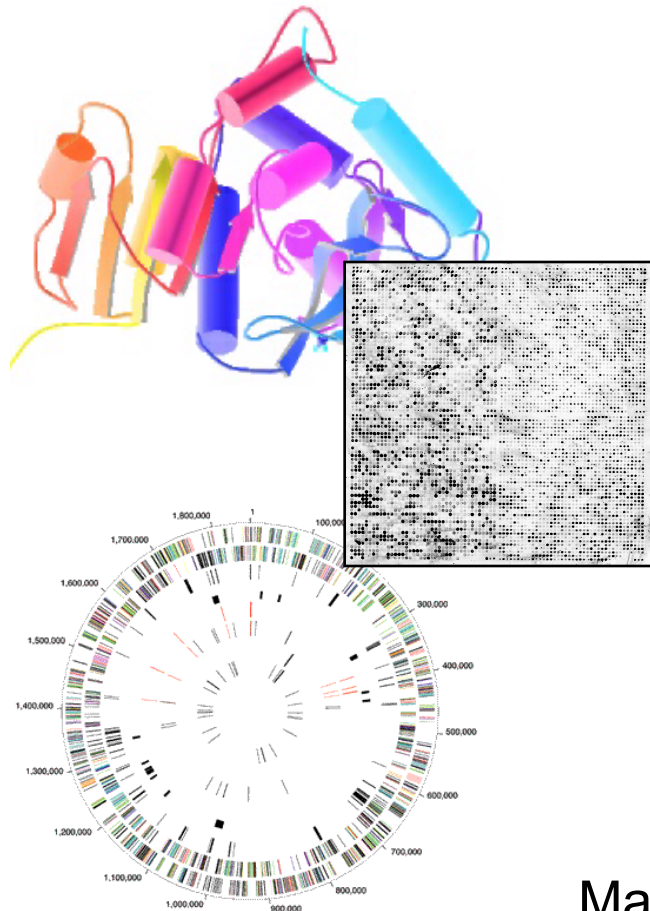


Biomedical Data Science: Introduction

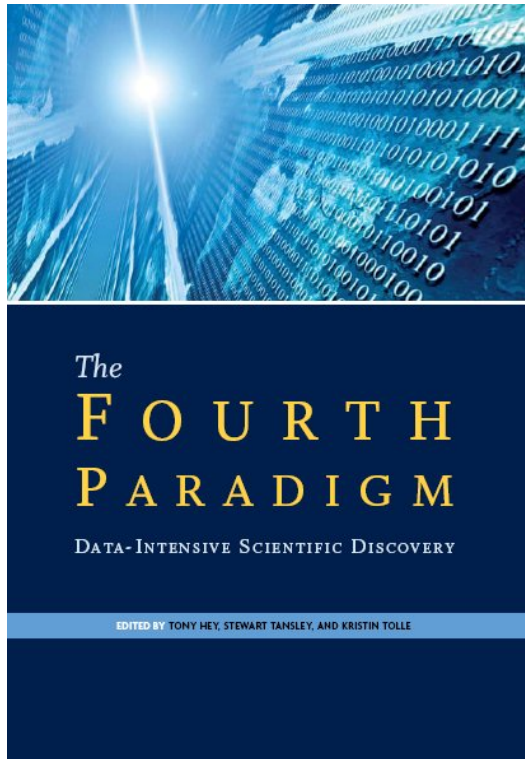


Mark Gerstein, Yale University
GersteinLab.org/courses/452
(last edit in spring '19)

**Overview: what is
Biomed. Data science?**

**(Placing it into the
context of Data
Science, in general)**

Jim Gray's 4th Paradigm

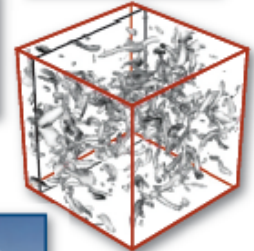


Science Paradigms

- Thousand years ago: science was **empirical**
describing natural phenomena
- Last few hundred years: **theoretical** branch
using models, generalizations
- 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



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



Jim Gray's 4th Paradigm

#3 - Simulation

Prediction based on physical principles (eg Exact Determination of Rocket Trajectory)

Emphasis on:
Supercomputers

#4 - Data Mining

Classifying information & discovering unexpected relationships

Emphasis: networks,
“federated” DBs

Science Paradigms

- Thousand years ago: science was **empirical** describing natural phenomena
- Last few hundred years: **theoretical** branch using models, generalizations
- Last few decades: a **computational** branch simulating complex phenomena

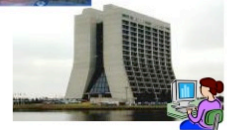
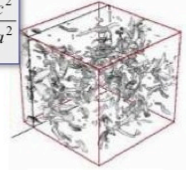
Today:

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



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



Gray died in '07.

Book about his ideas came out in '09.....

What is Data Science? An overall, bland definition...

- Data Science encompasses the study of the entire lifecycle of data
 - Understanding of how data are **gathered** & the issues that arise in its collection
 - Knowledge of what data sources are available & how they may be synthesized to solve problems
 - The **storage**, access, annotation, management, & transformation of data
- Data Science encompasses many aspects of data analysis
 - Statistical inference, machine learning, & the design of algorithms and computing systems that enable **data mining**
 - Connecting this mining where possible with **physical modeling**
 - The presentation and **visualization** of data analysis
 - The use of data analysis to make **practical decisions** & policy
- Secondary aspects of data, not its intended use – eg the data exhaust
 - The appropriate protection of **privacy**
 - Creative **secondary uses** of data – eg for Science of science
 - The elimination of inappropriate bias in the entire process

- Ads, media, product placement, supply optimization,
- Integral to success of GOOG, FB, AMZN, WMT...

Data Science in the wider world: a buzz-word for successful Ads



Harvard Business Review

Data Scientist: The Sexiest Job of the 21st Century

by Thomas H. Davenport and D.J. Patil

Artwork: **Tamar Cohen, Andrew J Buboltz**, 2011, silk screen on a page from a high school yearbook.

When Jonathan Goldman arrived for work in June 2006 at LinkedIn, the business ne up. The company had just under 8 million accounts, and the number was growing qu friends and colleagues to join. But users weren't seeking out connections with the pe rate executives had expected. Something was apparently missing in the social expe

Forbes · New Posts · Most Popular · Lists

108
 349
 193
 353
 12

CIO Network
 INSIGHTS AND IDEAS FOR TECHNOLOGY LEADERS.
 + Follow (489)

TECH | 12/12/2012 @ 1:57AM | 3,289 views

Why Big Data Is All Retailers Want for Christmas

Eric Savitz, Forbes Staff
 + Comment Now + Follow Comments

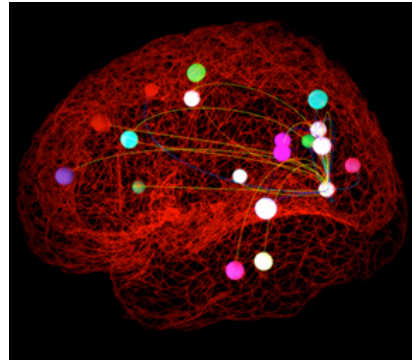
Guest post written by **Quentin Gallivan**
 Quentin Gallivan is CEO of Pentaho Corp., an Orlando, Florida-based provider of business analytics software.

Data Science in Traditional Science

- Pre-dated commercial mining
- Instrument generated
- Large data sets often created by large teams not to answer one Q but to be mined broadly
- Often coupled to a physical/biological model
- Interplay w/ experiments



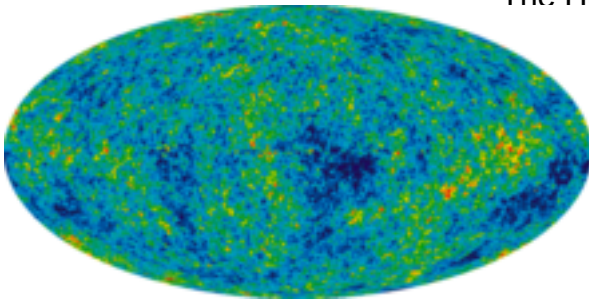
High energy physics -
Large Hadron Collider



Neuroscience -
The Human Connectome Project



Ecology
& Earth Sci.
- Fluxnet



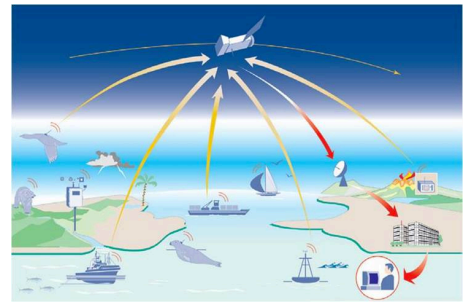
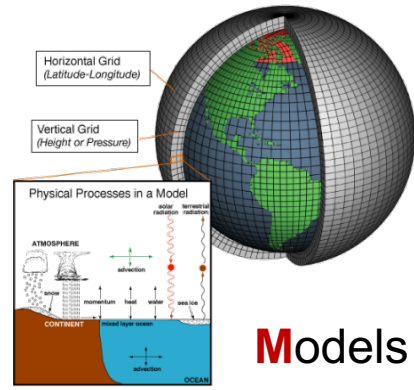
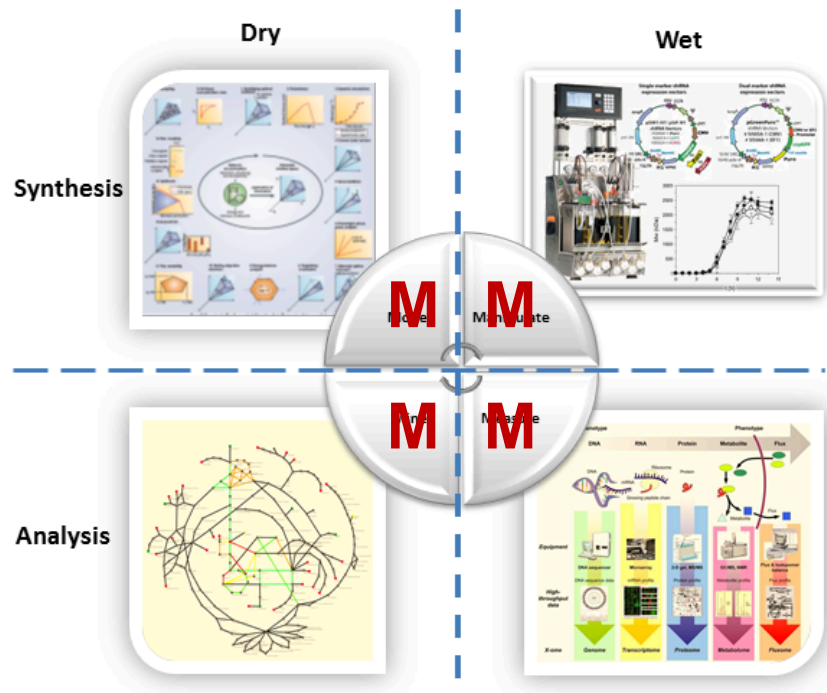
Astronomy -
Sloan Digital Sky survey



Genomics
DNA
sequencer

Coupling of Scientific Data to Models & Experiments

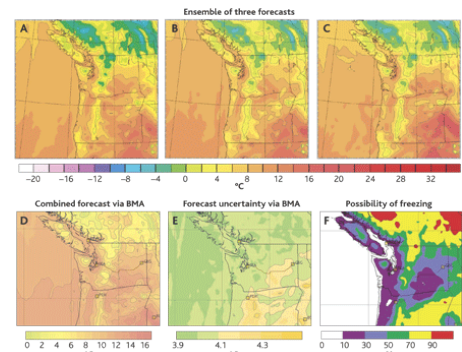
- Scientific data often coupled to a physical/biological model
- Lauffenburger's Sys. Biol. **4Ms**:
Measurement, **M**ining, **M**odeling & **M**anipulation
 (Ideker et al.'06. Annals of Biomed. Eng.)
- Weather forecasting as an exemplar
 - Physical models & simulation useful but not sufficient ("butterfly" effect)
 - Success via coupling to large-scale sensor data collection



Models + **D**ata **M**ining



Forecasts

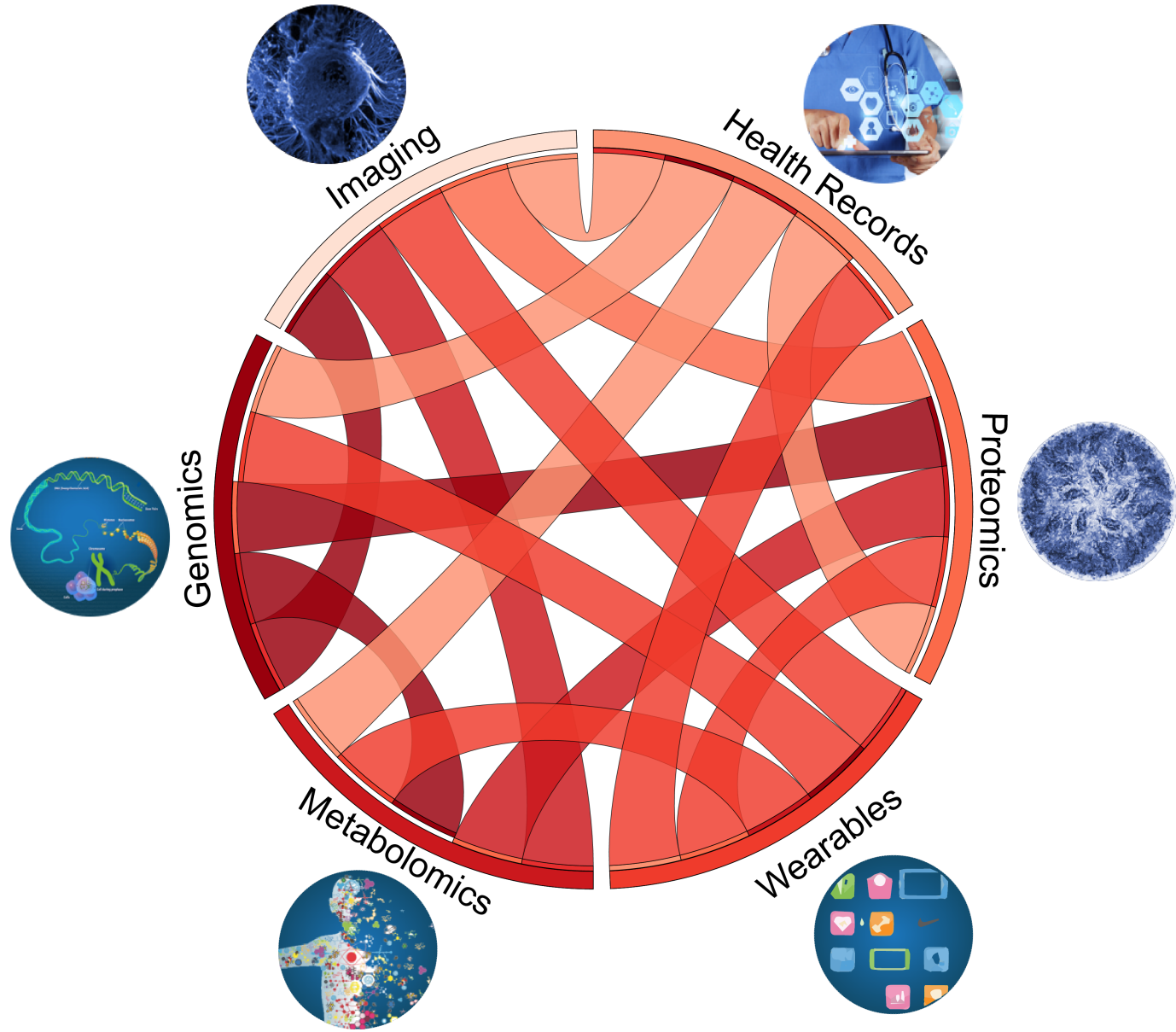


Biomed. Data science:

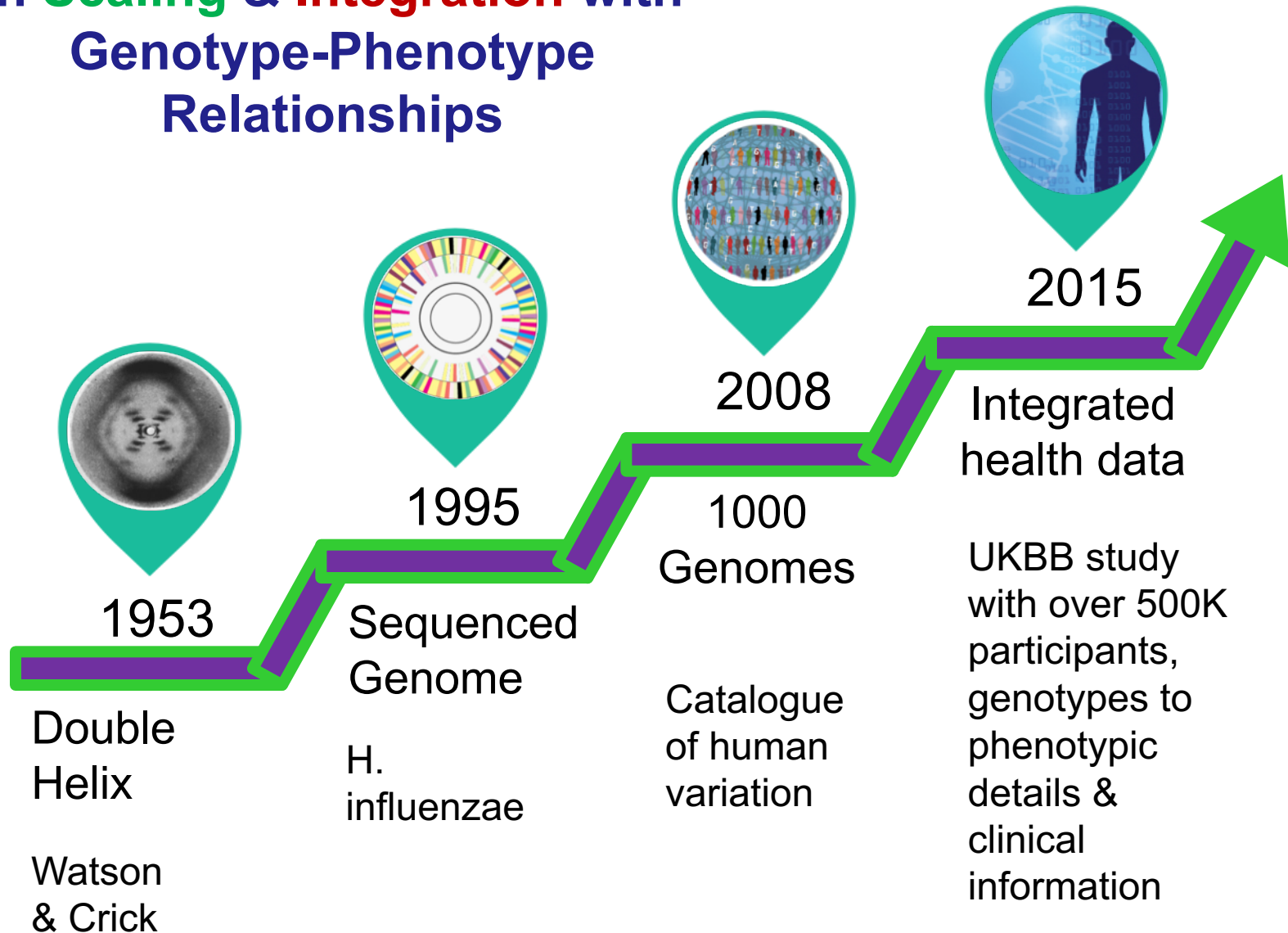
Scaling & Integration

Drivers of Biomedical Data Science

- **Integration** across data types
- **Scaling** of individual data types



Case Study: Amazing Progress in **Scaling & Integration** with Genotype-Phenotype Relationships

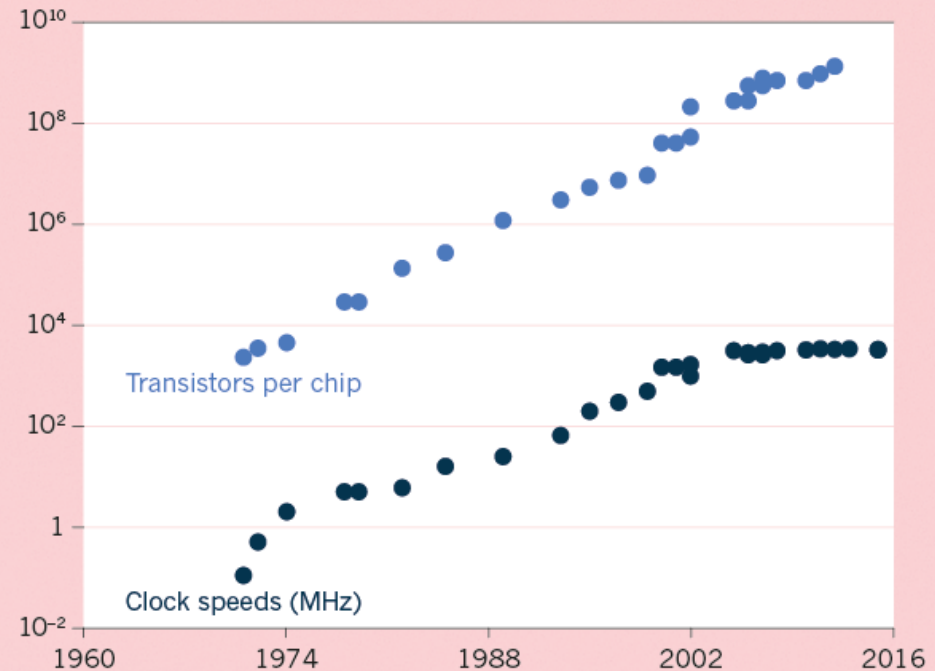
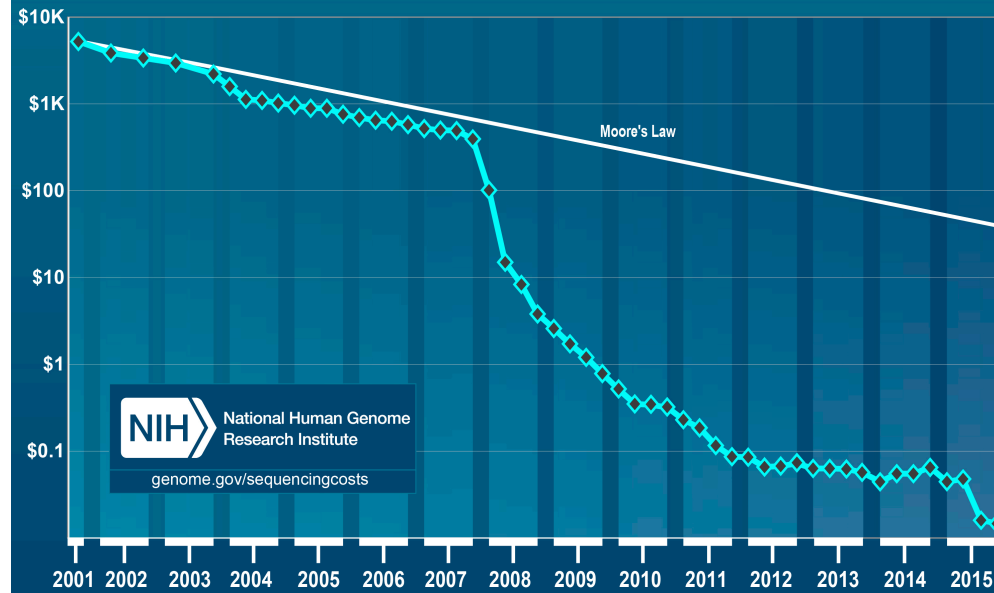


The **Scaling** of
Genomic Data
Science:

Powered by
exponential
increases in
data & computing

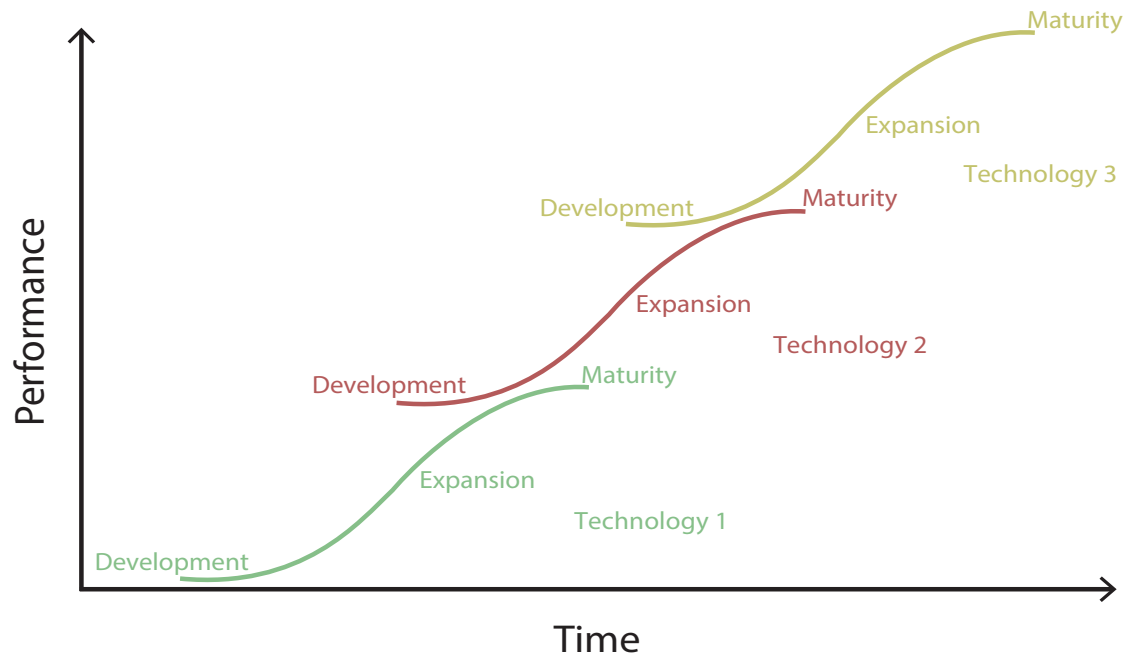
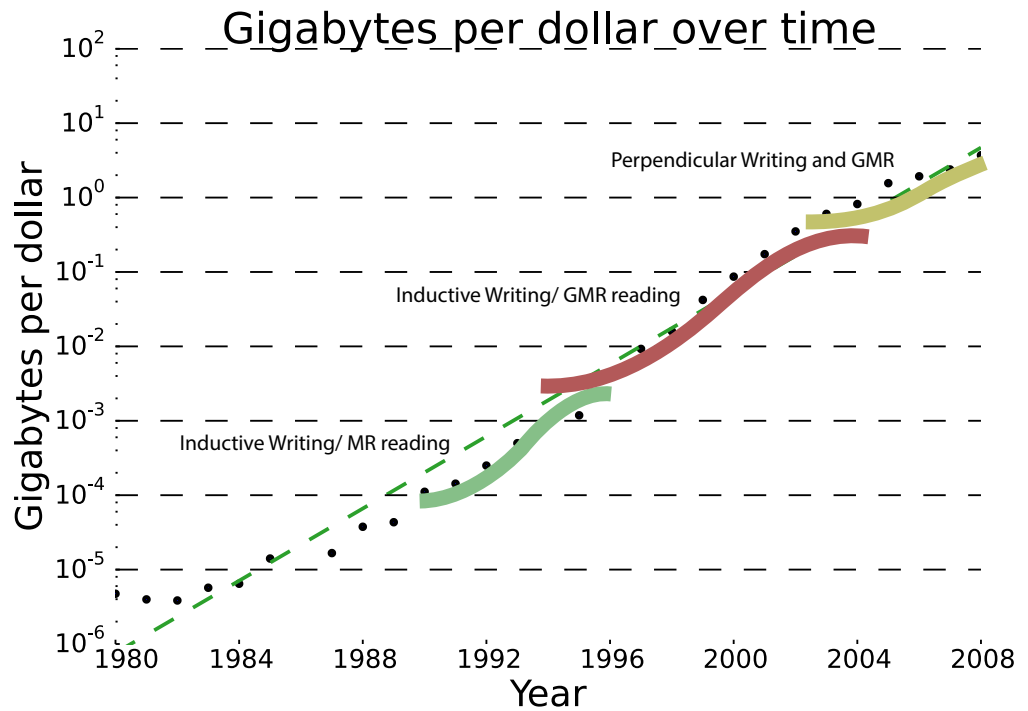
(**Moore's Law**)

Cost per Raw Megabase of DNA Sequence

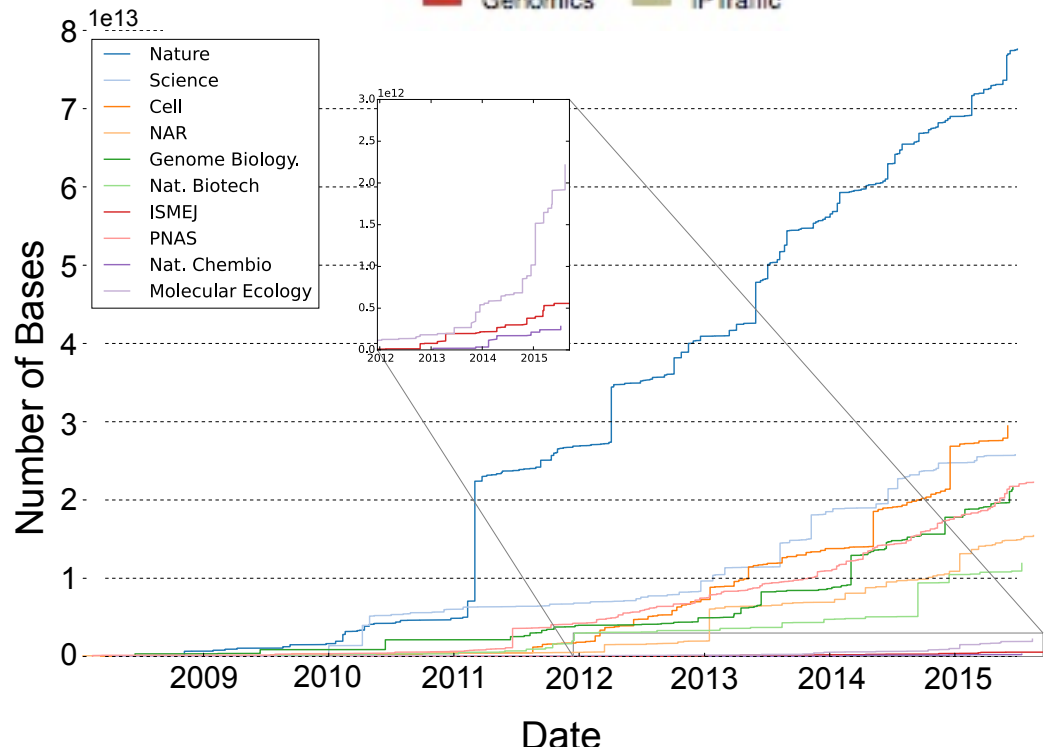
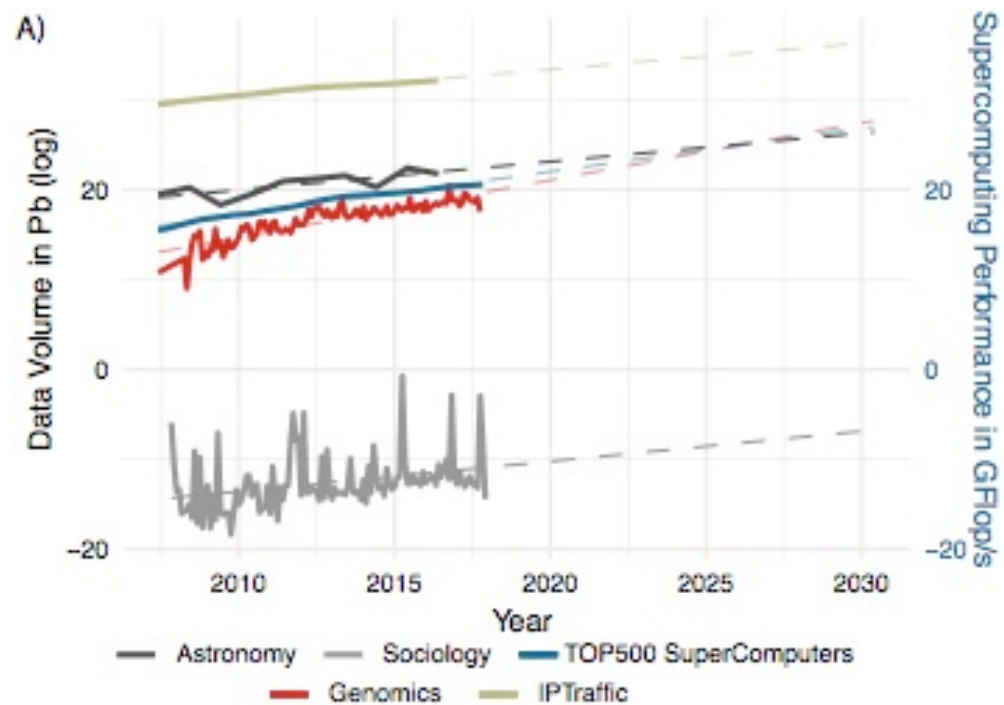


Kryder's Law and S-curves underlying exponential growth

- Moore's & Kryder's Laws
 - As important as the increase in computer speed has been, the ability to store large amounts of information on computers is even more crucial
- Exponential increase seen in Kryder's law is a superposition of S-curves for different technologies

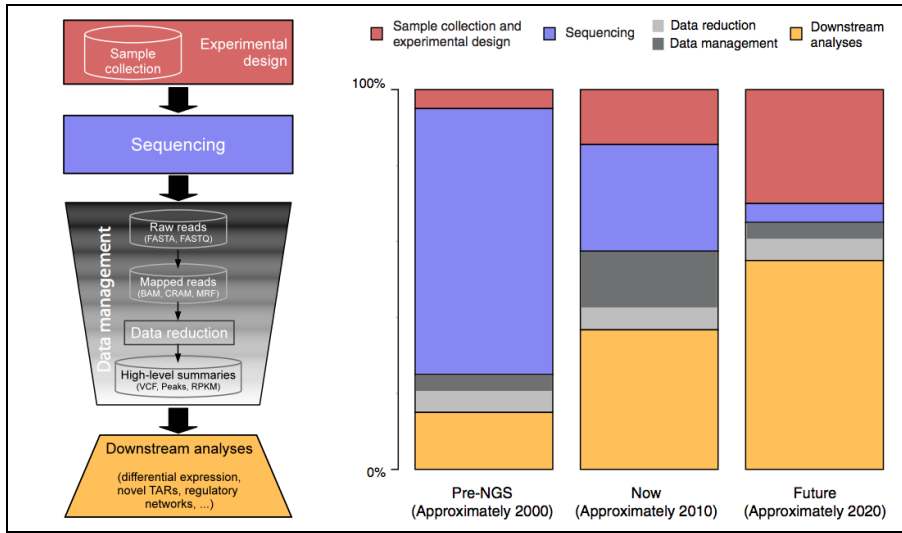


Sequencing cost reductions have resulted in an explosion of data



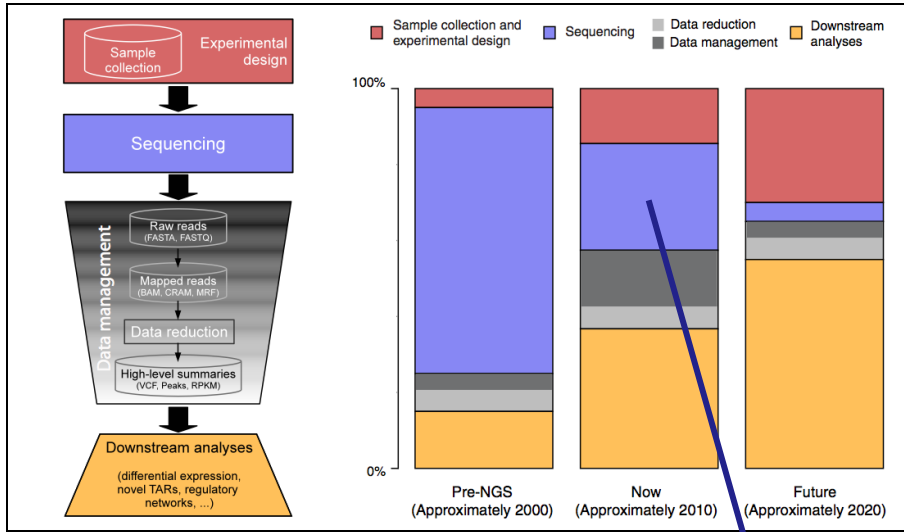
- The type of sequence data deposited has changed as well.

The changing costs of a sequencing pipeline

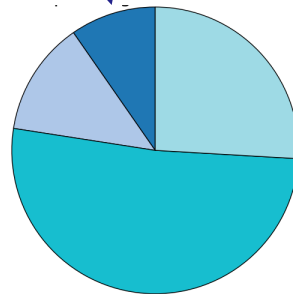
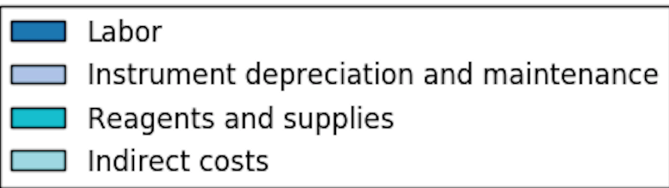


From '00 to ~' 20,
cost of DNA sequencing expt. shifts from
the actual seq. to sample
collection & analysis

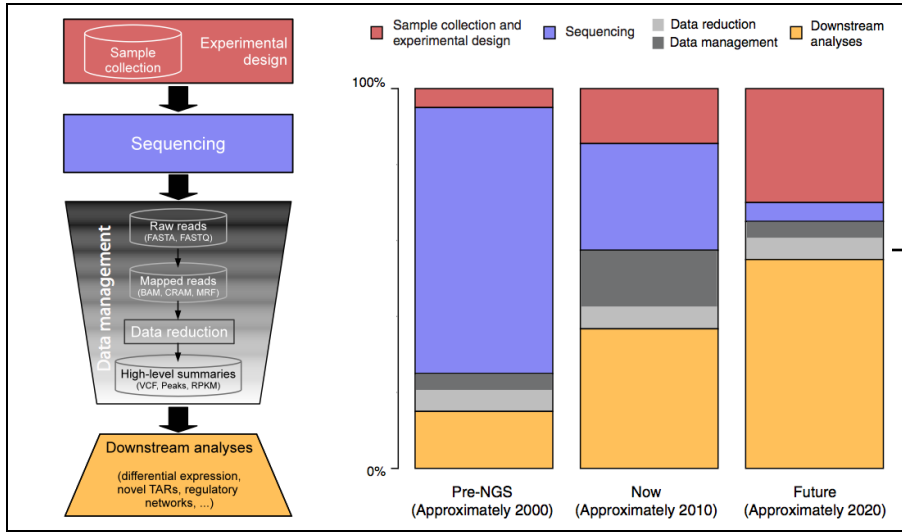
The changing costs of a sequencing pipeline



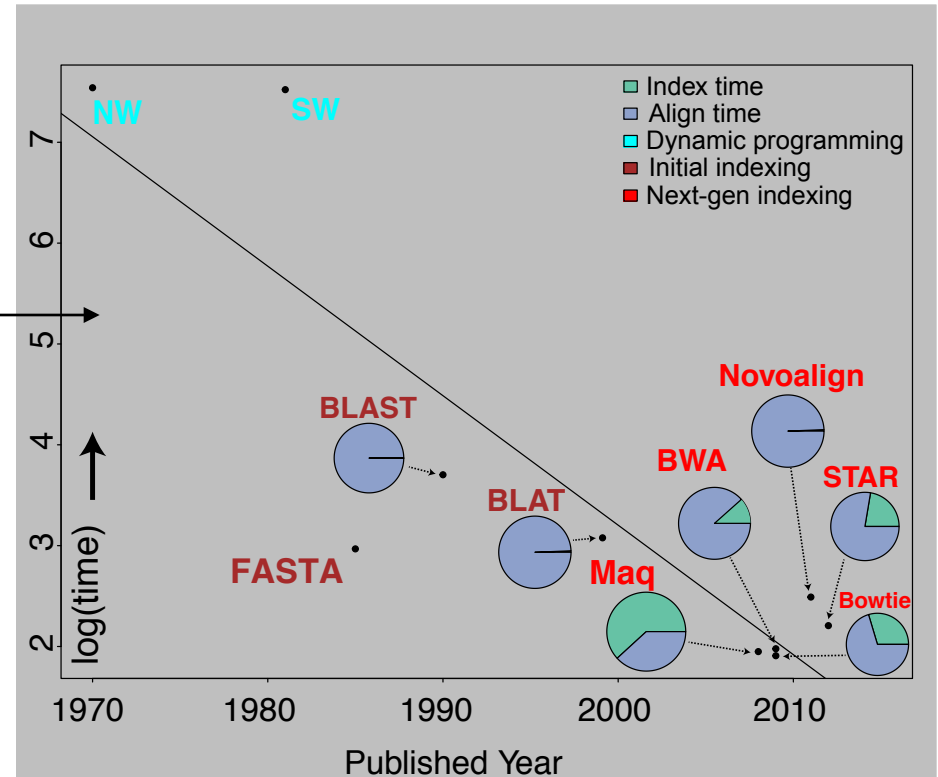
From '00 to ~' 20, cost of DNA sequencing expt. shifts from the actual seq. to sample collection & analysis



The changing costs of a sequencing pipeline

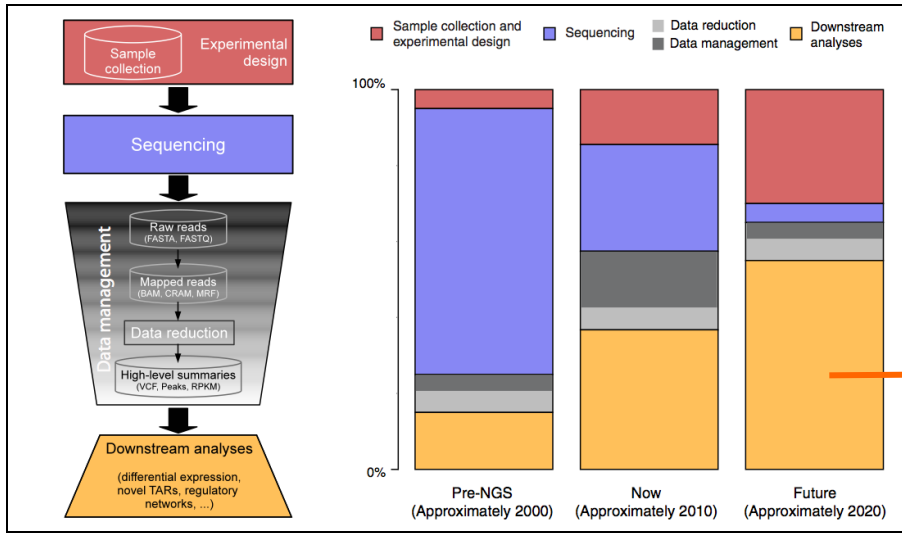


From '00 to ~' 20, cost of DNA sequencing expt. shifts from the actual seq. to sample collection & analysis

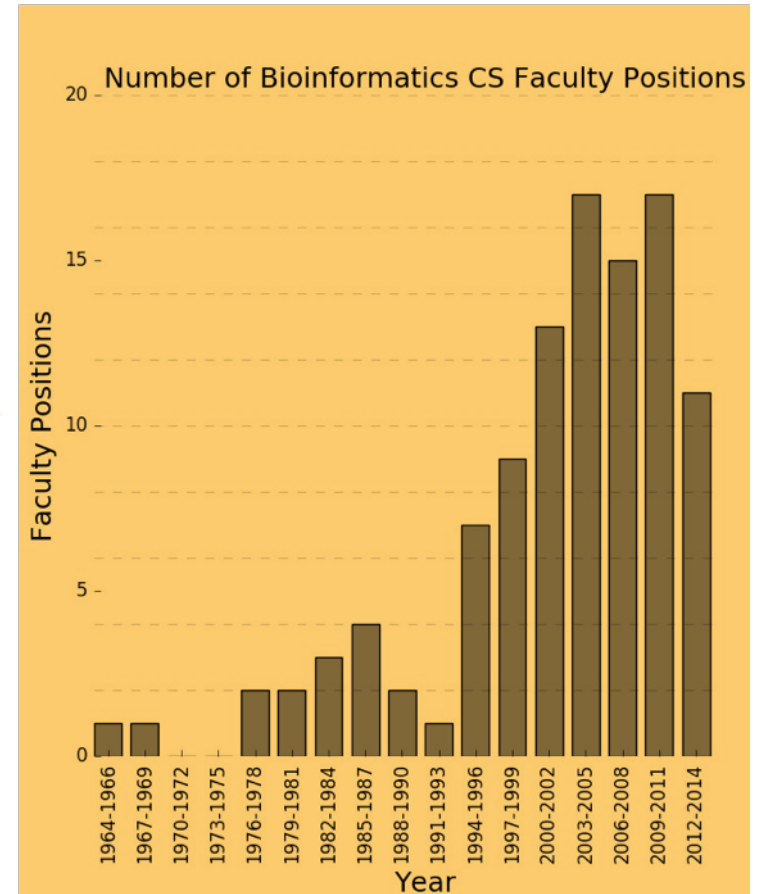


Alignment algorithms scaling to keep pace with data generation

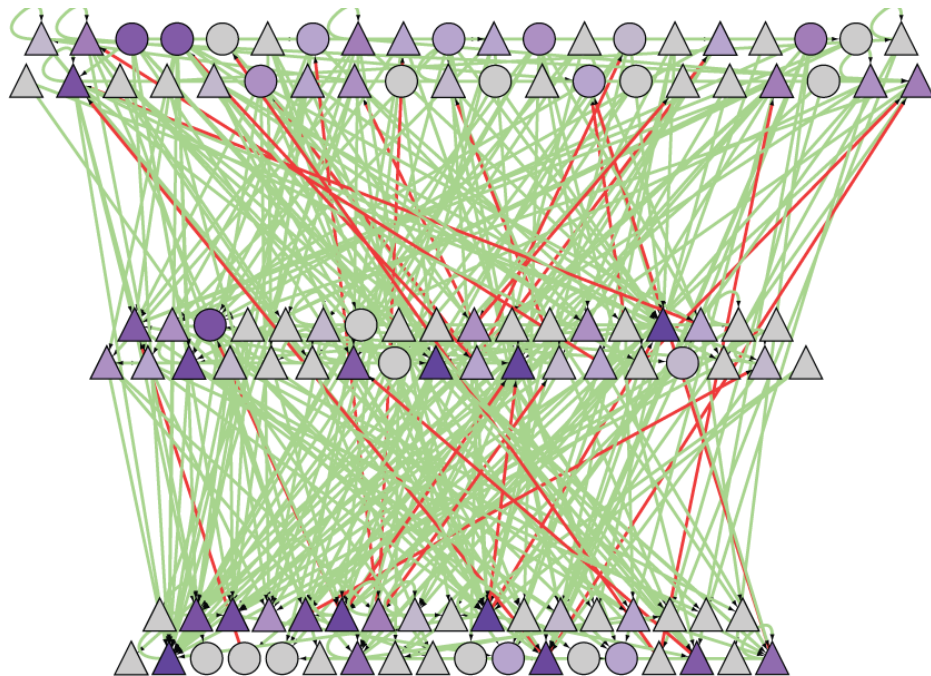
The changing costs of a sequencing pipeline



From '00 to ~' 20, cost of DNA sequencing expt. shifts from the actual seq. to sample collection & analysis



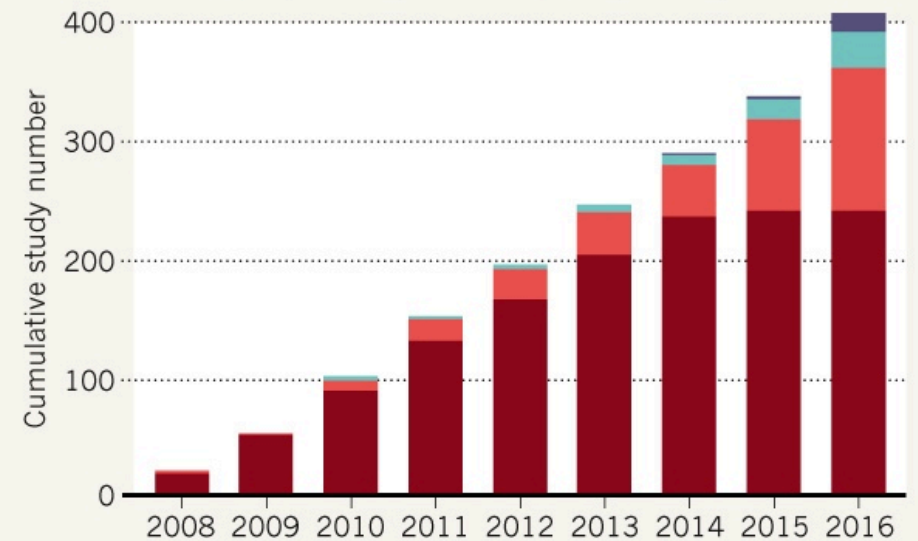
**A Success of
Scale & Integration:
Many GWAS
variants found,
most not in genes,
but affecting
regulatory network**



THE GENOME-WIDE TIDE

Large genome-wide association studies that involve more than 10,000 people are growing in number every year — and their sample sizes are increasing.

Sample sizes: ■ More than 200,000 ■ 100,000–199,999
■ 50,000–99,999 ■ 10,000–49,999



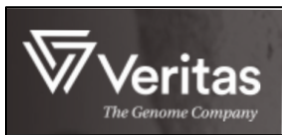
©nature

- A 1st GWAS done at Yale, for AMD: (Klein et al. 05, Science)
- Many since then
- Most SNVs fall into non-coding regulatory regions (major contributions by Yale groups to this ENCODE annotation effort)

[Nature 489: 91]

Basic Science to Medicine

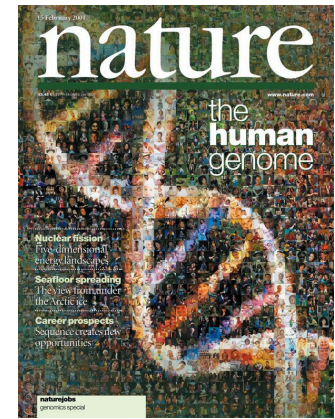
INITIATIVES



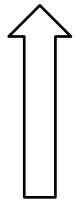
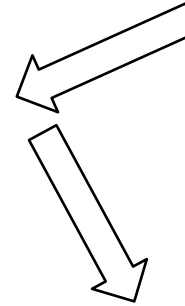
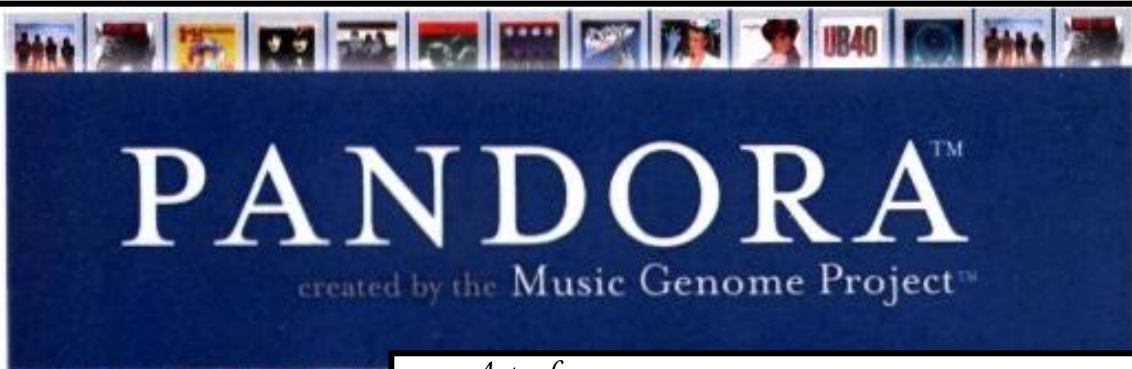
STARTUPS

- Large-scale ‘omics data as an anchor to organize phenotypic data – EMRs, wearables...
- 1st [‘05-]: Exomes & chips of disease-focused cohorts – init. GWAS, TCGA, PGC
- 2nd [‘15-]: Integration of full WGS with rich & diverse phenotypes - UKBiobank, TopMed, Genomics England, PCAWG, All of Us

Genomics: as Data Science sub-discipline



- Developing ways of organizing & mining categorizing information on a large scale
 - Very fundamental & early form of "Big Data", feeding into other enterprises (classification approach, R)
 - Also importing tech. developed in other big data disciplines (Hadoop)

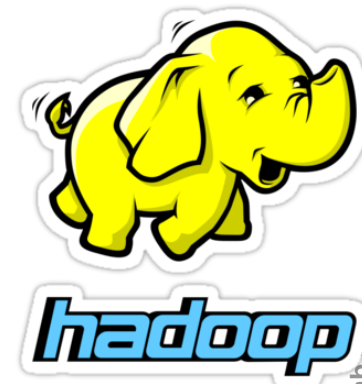


A. *Artsy for Education* Resources for discovering and learning about art online

EXPLORE CATEGORIES DISCOVER INSTITUTIONS

What is The Art Genome Project? Seven Facts about the Discovery and Classification System That Fuels Artsy

THE ART GENOME PROJECT
BY MATTHEW ISRAEL, JESSICA BACKUS AND OLIVIA JENE FAGON
FEB 9TH, 2016 5:00 AM



Examples of Imports & Exports To/from Genomics & Other Data Science Application Areas

Imports

Hidden Markov Models

Fast & memory-efficient string processing algorithms

Network and Graphs

Cultural

Critical Assessment of Protein Structure Prediction (CASP)

Exports

Latent Dirichlet Allocation (LDA)

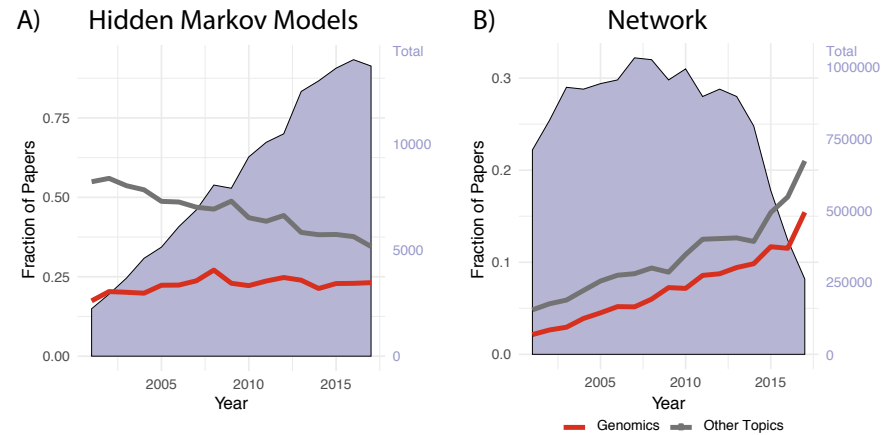
Circos plot

Digital classification systems

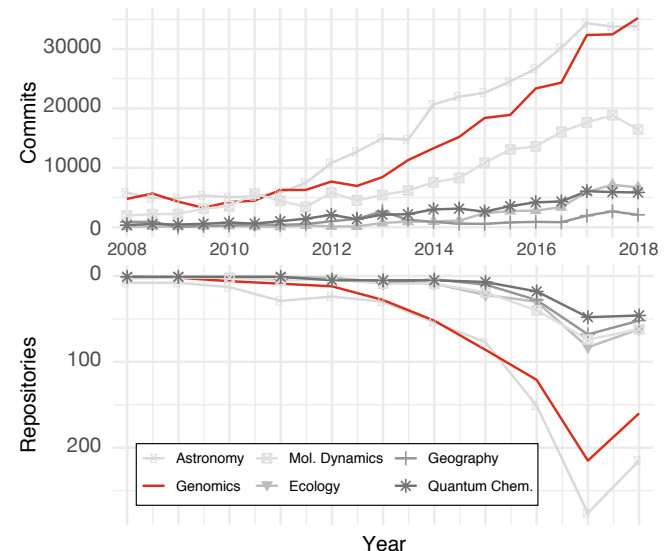
Cultural

Data Science cultural practices: Open data and
Open science culture

Genomics contribution of HMM and Network papers



Genomics adoption of open source

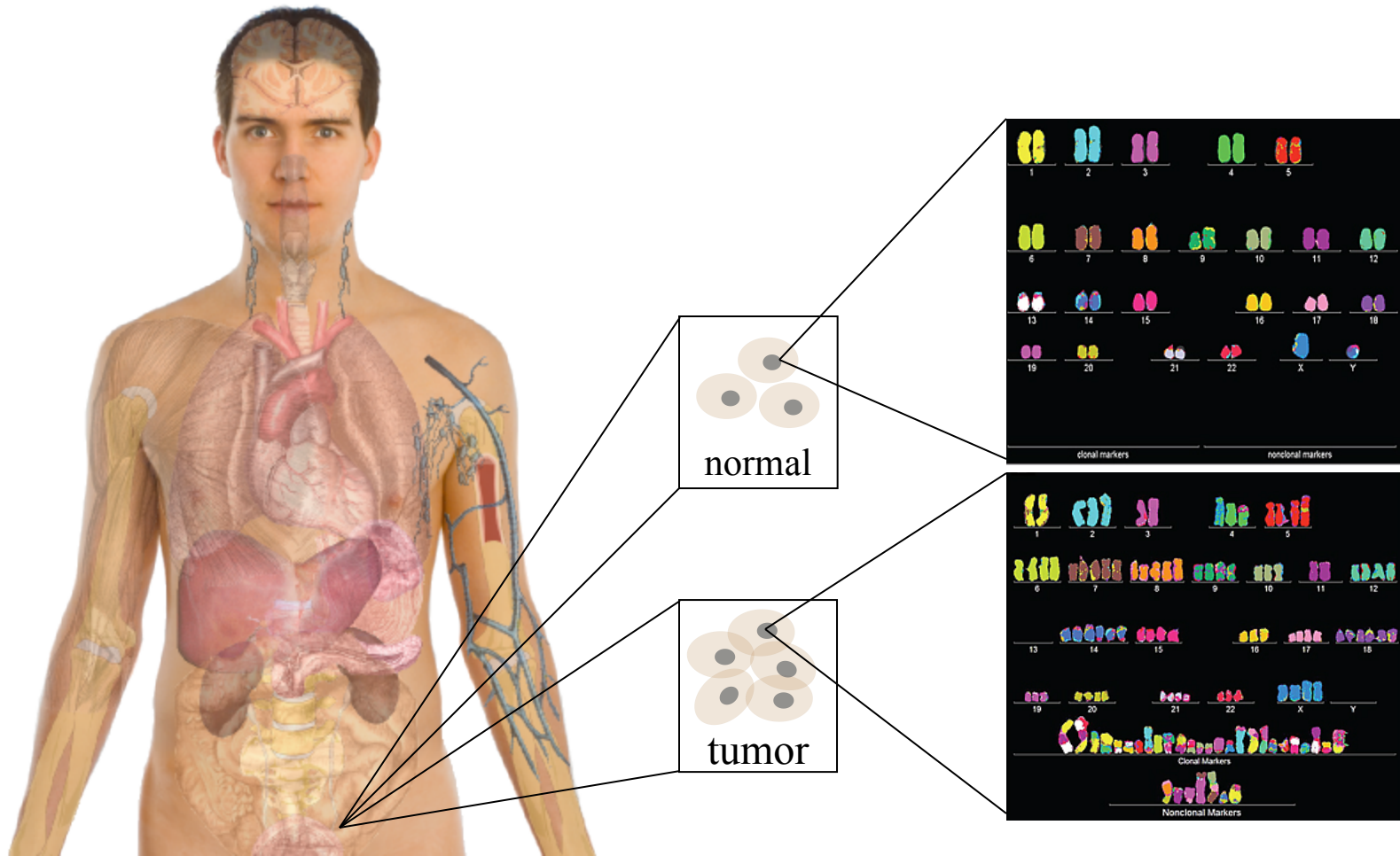


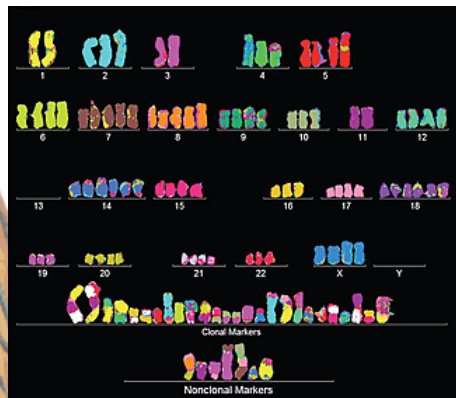
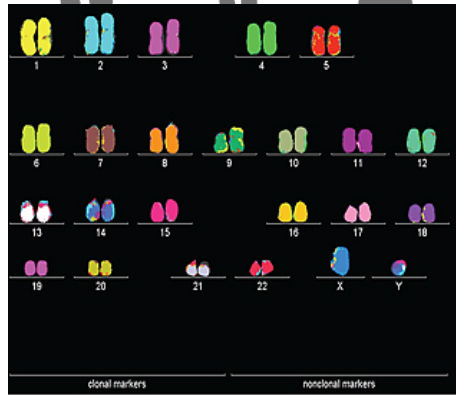
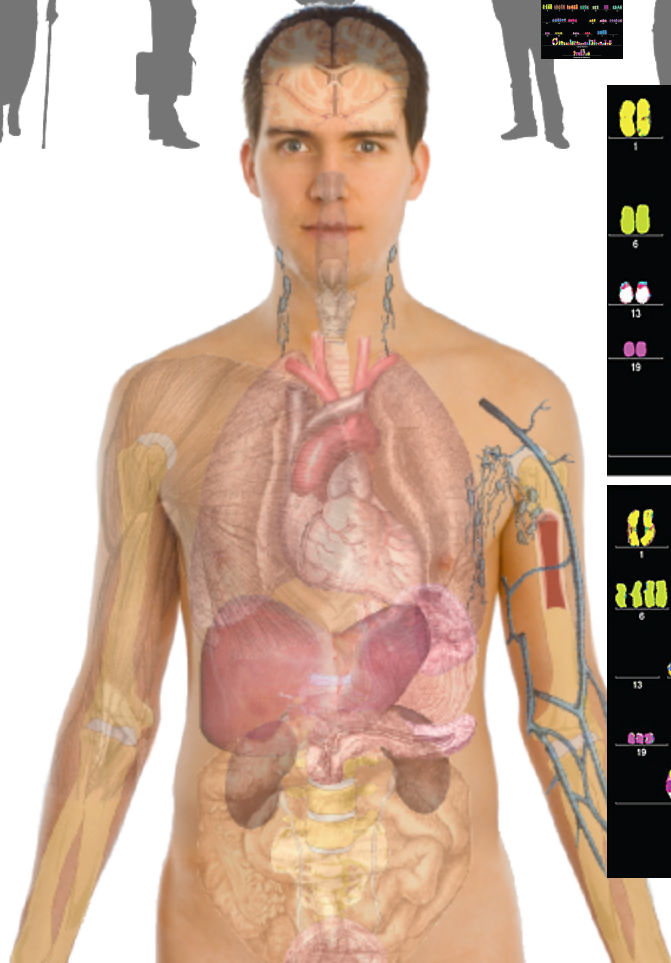
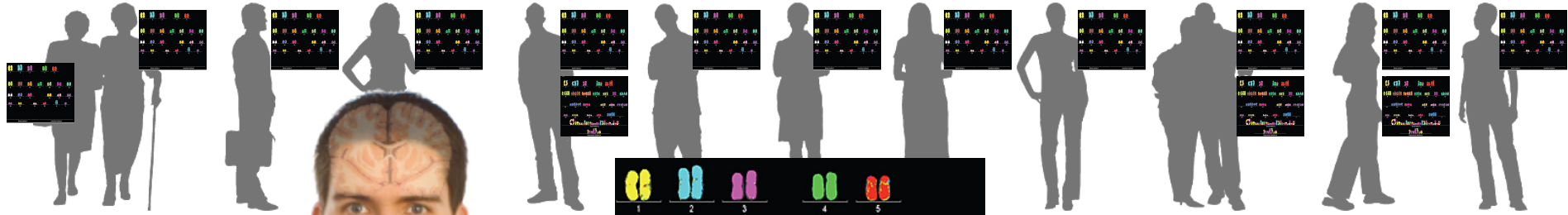
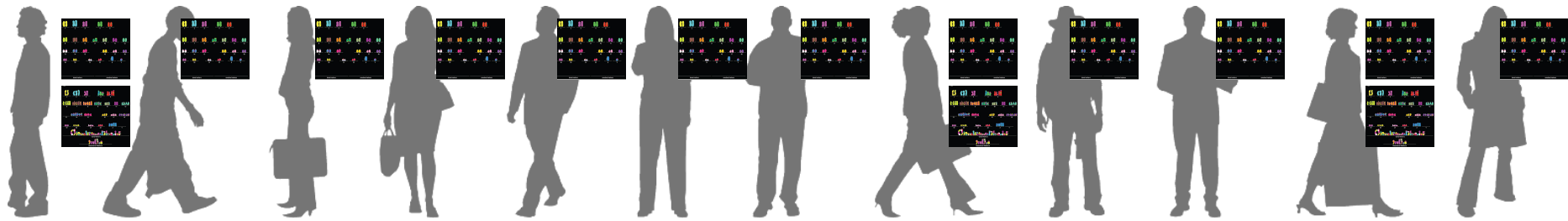
Biomed. Data science:

The Future

Our field as future Gateway – Personal Genomics as a Gateway into Biology

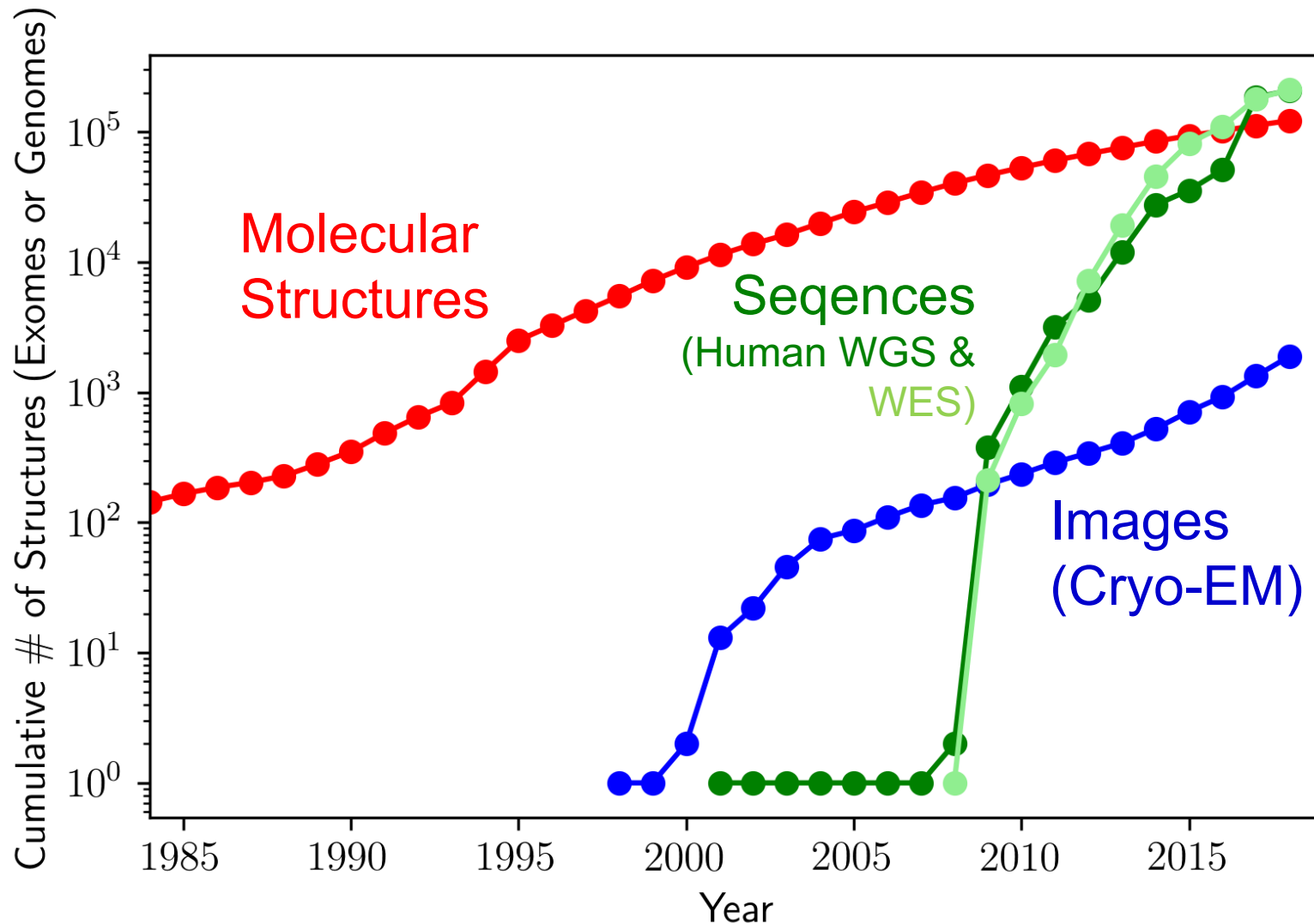
Personal genomes soon will become a commonplace part of medical research & eventually treatment (esp. for cancer). They will provide a primary connection for biological science to the general public.





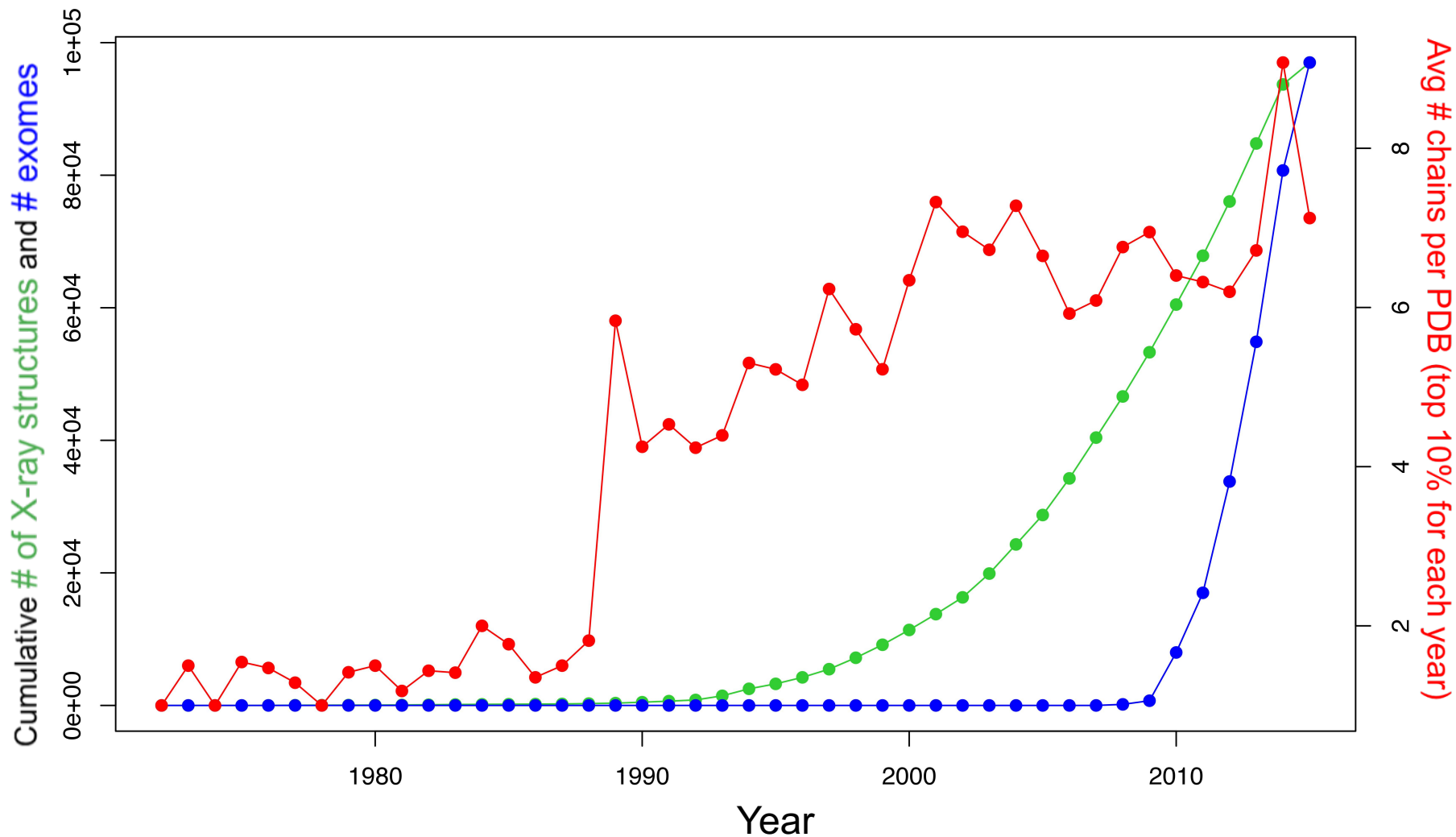
Placing the individual into the context of the population & using the population to build a interpretative model

How will the Data **Scaling** Continue? The Past, Present & Future Ecosystem of Large-scale Biomolecular Data



Trends in data generation point to growing opportunities for leveraging sequence variants to study structure (and vice versa)

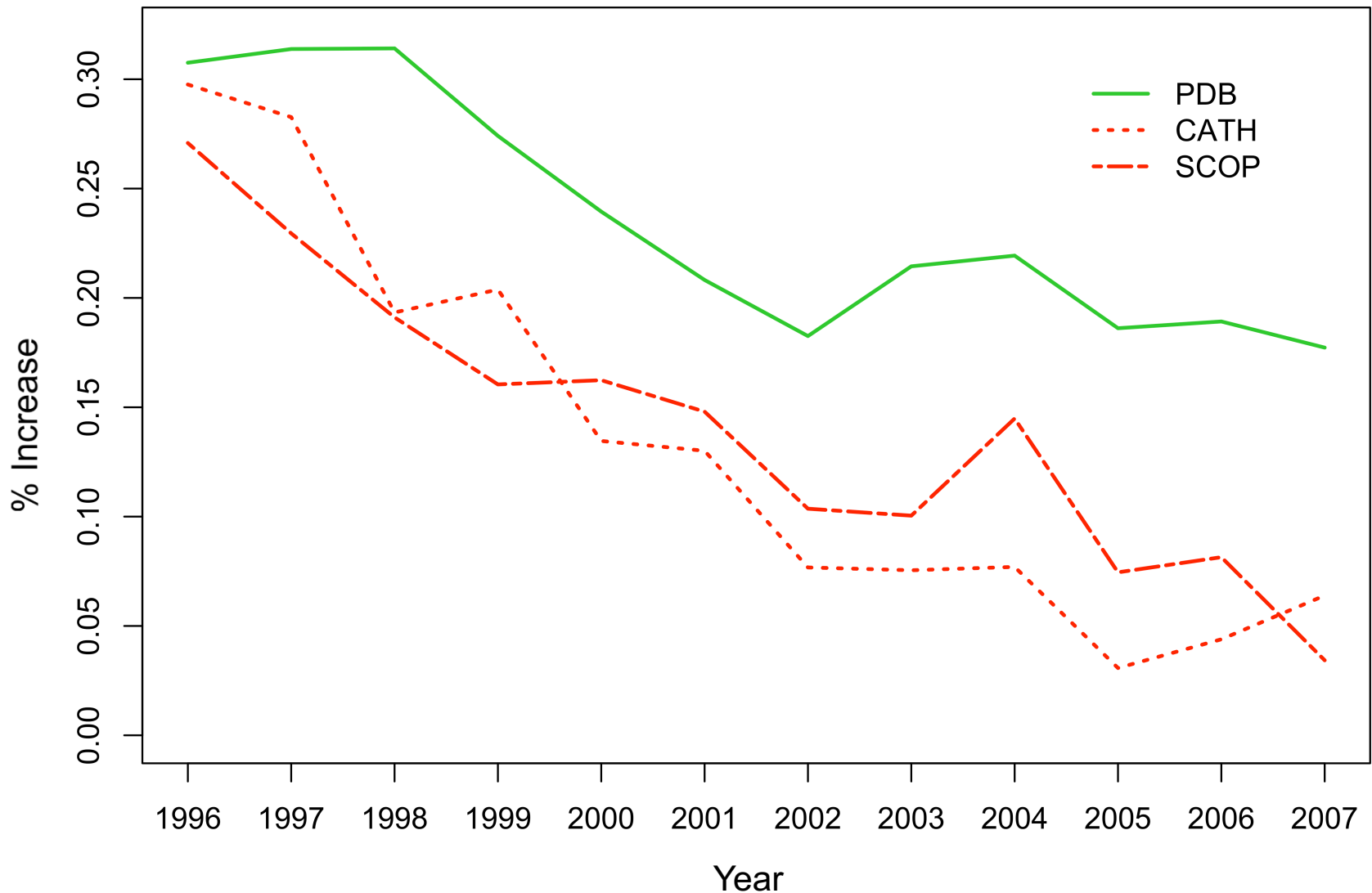
The volume of sequenced exomes is outpacing that of structures, while solved structures have become more complex in nature.



Exome data hosted on NCBI Sequence Read Archive (SRA)

[Sethi et al. COSB ('15)]

Growing sequence redundancy in the PDB (as evidenced by a reduced pace of novel fold discovery) offers a more comprehensive view of how such sequences occupy conformational landscapes – Gene & Struc. Families as main organizing principle



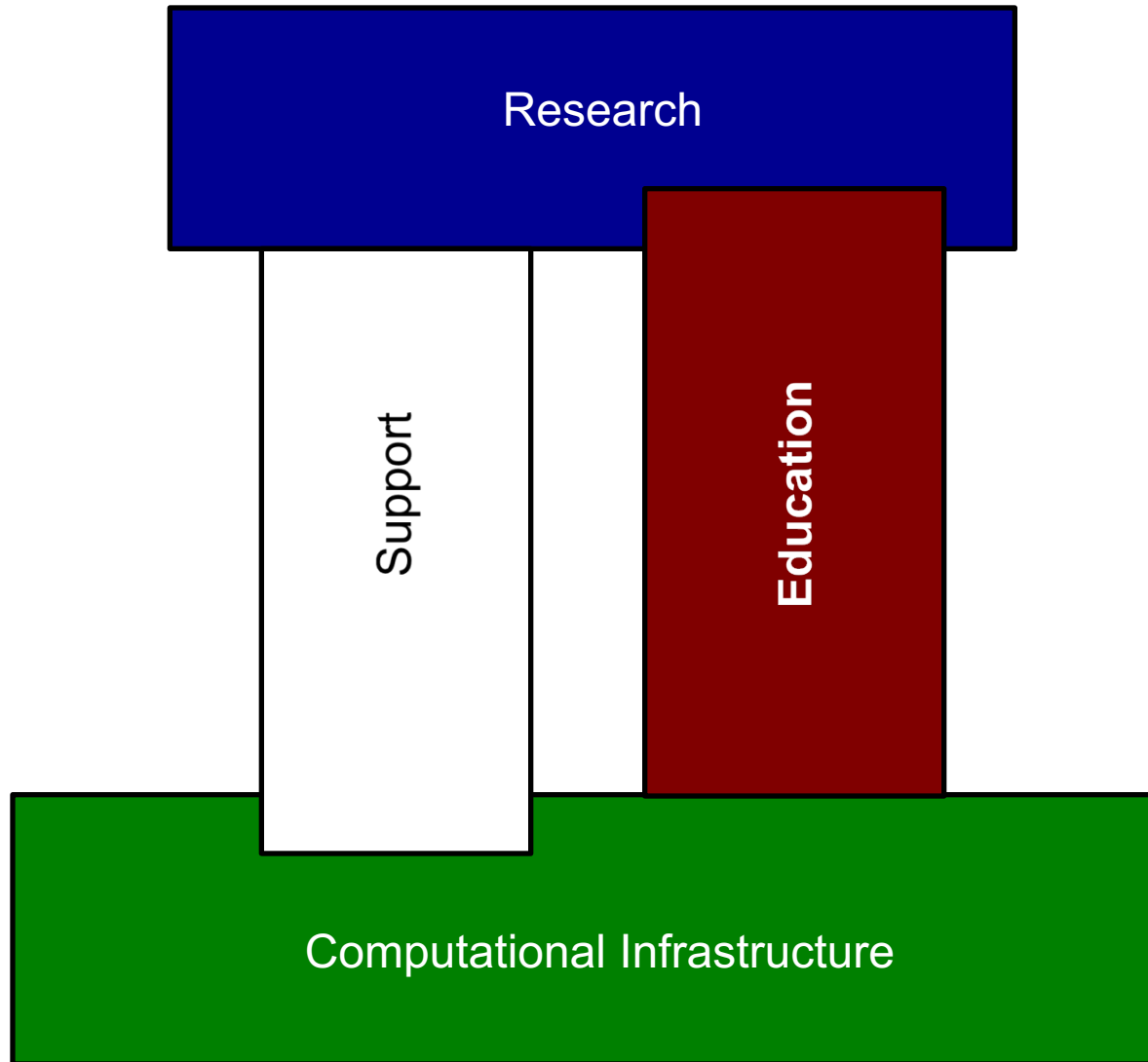
[Sethi et al. COSB ('15)]

PDB: Berman HM, et al. NAR. (2000)
CATH: Sillitoe I, et al. NAR. (2015)
SCOP: Fox NK et al. NAR. (2014)

Biomed. Data science:

The Course

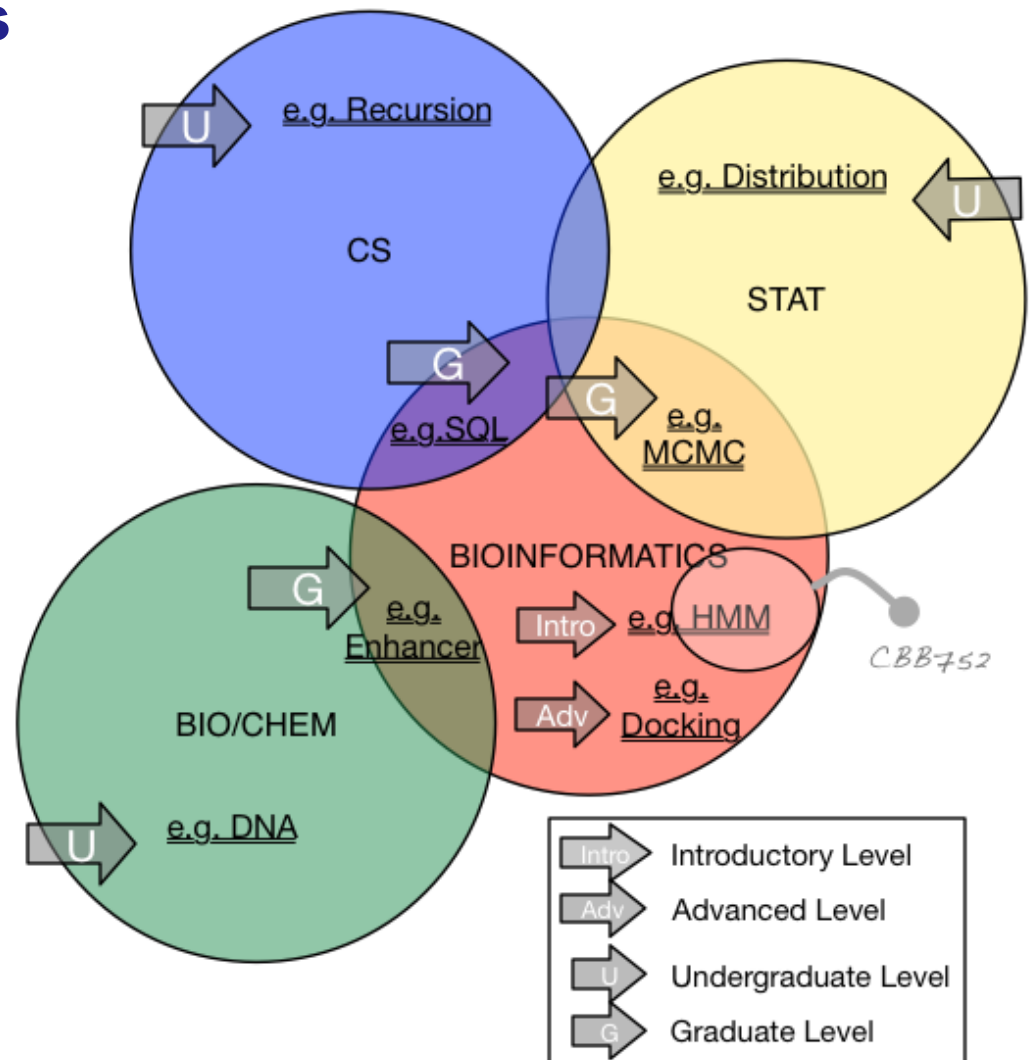
Elements of Bioinformatics as a discipline



Defining Bioinformatics

– by crowd-sourced judgement

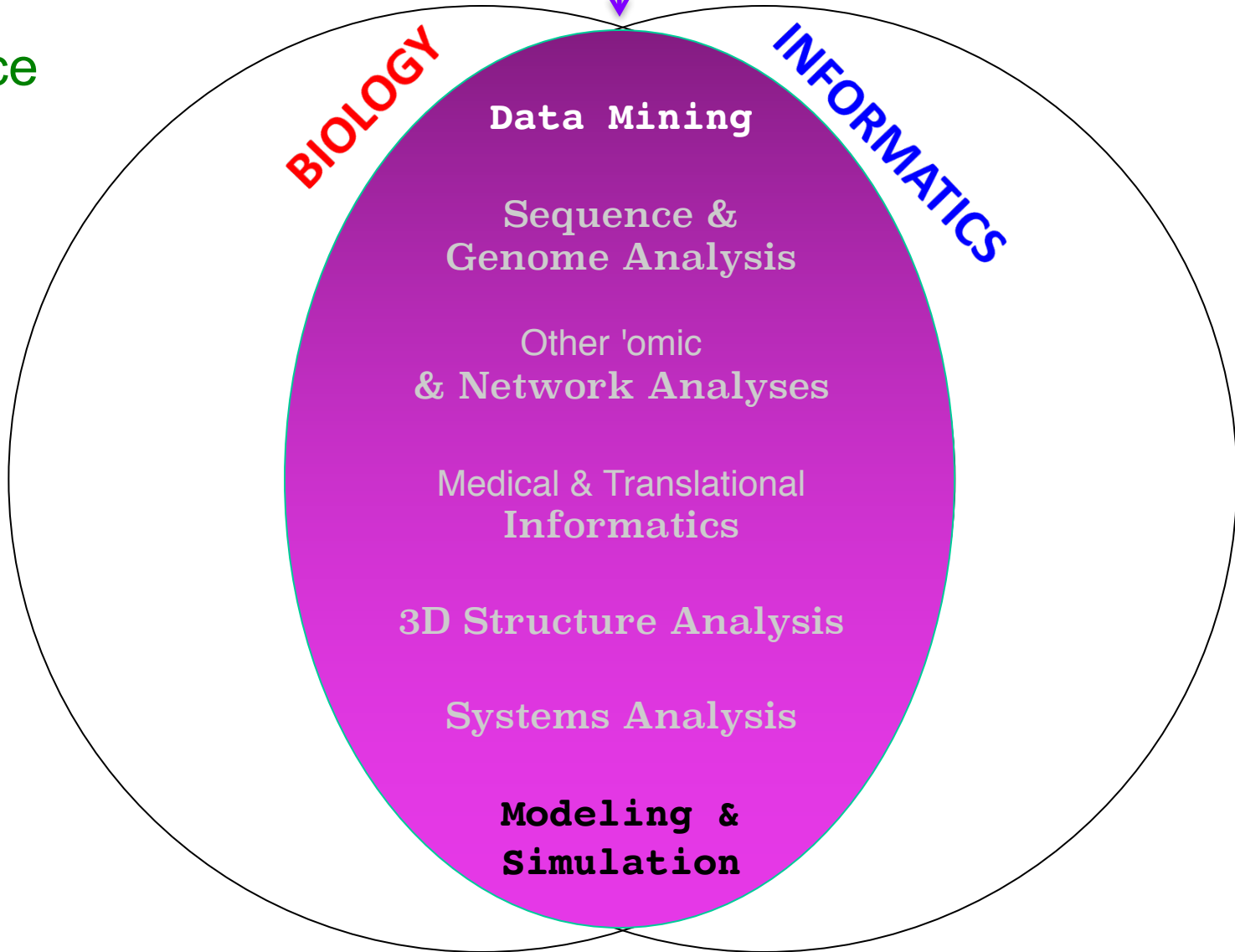
- Bioinformatics
 - Related terms
 - Biological Data Science
 - Bioinformatics & / or / vs Computational Biology
 - Biocomputing
 - Systems Biology
 - Qbio
- What are its boundaries
 - Determining the "Support Vectors"



Biomedical
Data
Science



(Molecular) BIOINFORMATICS



What is Bioinformatics?

- (*Molecular*) **Bio - informatics**

- One idea for a definition?

Bioinformatics is conceptualizing **biology in terms of molecules** (in the sense of physical-chemistry) and then applying **“informatics” techniques** (derived from disciplines such as applied math, CS, and statistics) to **organize, mine, model & understand the information associated** with these molecules, **on a large-scale.**

- Bioinformatics is a practical discipline with many **applications.**

Class Web Page

GersteinLab.org/courses/452

Assignment #0 Page

bit.ly/cbb752b19-hw0

Office Hours

**Right before & after class
on Wed. 1/16
(in Bass 432,
contact.gerstein.info)**

Biomed. Data science:

**More details on
Bioinformatics
as a sub-discipline**

What is Bioinformatics?

- (*Molecular*) **Bio - informatics**

- One idea for a definition?

Bioinformatics is conceptualizing **biology in terms of molecules** (in the sense of physical-chemistry) and then applying **“informatics” techniques** (derived from disciplines such as applied math, CS, and statistics) to **organize, mine, model & understand the information associated** with these molecules, **on a large-scale.**

- Bioinformatics is a practical discipline with many **applications.**

What Information to Organize?

- **Sequences** (DNA & Protein)
- 3D Structures
- Network & Pathway Connectivity
- Phylogenetic tree relationships
- Large-scale gene expression & functional genomics data
- Phenotypic data & medical records....

What is the Information?

Molecular Biology as an Information Science

- Central Dogma of Molecular Biology

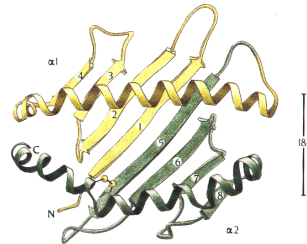
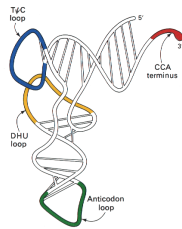
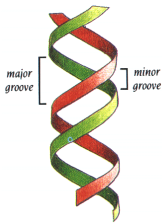
DNA

- > RNA
- > Protein
- > Phenotype
- > DNA

- Central Paradigm for Bioinformatics

Genomic Sequence Information

- > mRNA (level)
- > Protein Sequence
- > Protein Structure
- > Biological Function
- > Organismal Phenotype



•Genetic material

- Information transfer (mRNA)
- Protein synthesis (tRNA/mRNA)
- Some catalytic activity

Molecular Biology Information - DNA

- Raw DNA Sequence

- 4 bases:

- AGCT

- ~1 K in a gene, ~2 M in genome

- ~3 Gb Human

```
atggcaattaaaattggtatcaatggttttggctgatcggccgatcgtattccgtgca
gcacaacaccgtgatgacattgaagtgttaggtattaacgacttaatcgacggtgaatac
atggcttatatggttgaatatgattcaactcacggctggttcgacggcactgttgaagtg
aaagatggtaacttagtggtaatggtaaaactatccgtgtaactgcagaacgtgatcca
gcaacttaaaactggggtgcaatcggtgttgatcgcgtgttgaagcactgggttattc
ttaactgatgaaactgctcgtaaacatatcactgcaggcgcaaaaaagttgtattaact
ggcccatcctaagatgcaaccctatgttcggtcgtggtgtaaactcaacgcatacga
ggtcaagatatcgtttctaacgcattctgtacaacaaactgtttagctccttagcagct
gttggtcatgaaactttcgggtatcaaagatggtttaaagaccactgttcacgcaacgact
gcaactcaaaaaactgtggatgggtccatcagctaaagactggcgcggcggccgcggtgca
tcacaaaacatcattccatcttcaacaggtgcagcgaaagcagtaggtaaagtattacct
gcattaaacggtaaatctaactggatggctttccgctgttccaacgccaacgatatcgtt
gttgatttaacagttaatcttgaaaaaccagcttcttatgatgcaatcaacaagcaatc
aaagatgcagcgggaaggtaaaacgttcaatggcgaattaaaaggcgtattaggttacact
gaagatgctgttgtttctactgacttcaacgggttgctttaaactctctgtatttggatgca
gacgctggtatcgcatctaactgattcttttcgtaaattgggatc . . .
```

```
. . . caaaaatagggttaatatgaatctcgcgatctccattttgttcacgtattcaa
caacaagccaaaactcgtacaaatatgaccgcacttcgctataaagaacacggcttggtg
cgagatatctcttgaaaaactttcaagagcaactcaatcaacttctcgcagcattgctt
gctcacaatattgacgtacaagataaaaatcgccatttttgccataaatatggaacgttgg
gttggtcatgaaactttcgggtatcaaagatggtttaaagaccactgttcacgcaacgact
acaatcgttgacattgacaccttacaatttcgagcaatcacagtgacctatttacgcaacc
aatacagcccagcaagcagaatcttcaaatcacgcccagtgtaaaaaattctctctcgtc
ggcgatcaagagcaatcgcgatcaaacattggaaattgctcatcatgtccaaaattacaa
aaaattgtagcaatgaaatccaccattcaattacaacaagatcctcttcttgcacttgg
```


Molecular Biology Information: Protein Sequence

- 20 letter alphabet
 - ACDEFGHIKLMNPQRSTVWY but not BJOUXZ
- Strings of ~300 aa in an average protein (in bacteria), ~200 aa in a domain
- >12 M known protein sequences
(uniprot, <http://www.ebi.ac.uk/uniprot/TrEMBLstats/>, 2011)

```
d1dhfa_ LNCIVAVSQNMGIGKNGDLPWPPLRNEFRYFQRMFTTSSVEGKQ-NLVIMGKKTWFSI
d8dfr_  LNSIVAVCQNMGIGKDGNLFPWPPLRNEYKYFQRMFTSTSHVEGKQ-NAVIMGKKTWFSI
d4dfra_ ISLIAALAVDRVIGMENAMPWN-LPADLAWFKRNTL-----NKPVIMGRHTWESI
d3dfr_  TAFLWAQDRDGLIGKDGHLPHW-LPDDLHYFRAQTV-----GKIMVVGRRTYESF
```

```
d1dhfa_ LNCIVAVSQNMGIGKNGDLPWPPLRNEFRYFQRMFTTSSVEGKQ-NLVIMGKKTWFSI
d8dfr_  LNSIVAVCQNMGIGKDGNLFPWPPLRNEYKYFQRMFTSTSHVEGKQ-NAVIMGKKTWFSI
d4dfra_ ISLIAALAVDRVIGMENAMPW-NLPADLAWFKRNTLD-----KPVIMGRHTWESI
d3dfr_  TAFLWAQDRNGLIGKDGHLPW-HLPDDLHYFRAQTVG-----KIMVVGRRTYESF
```

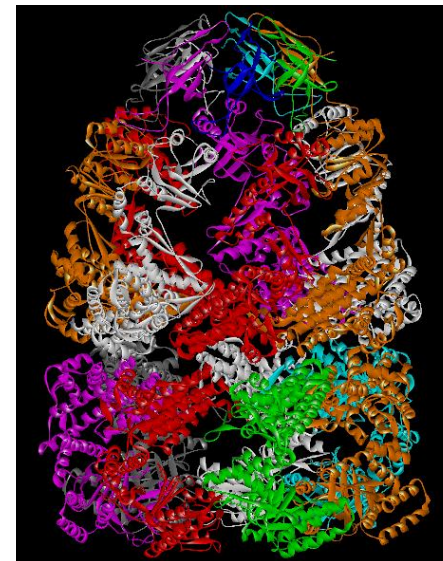
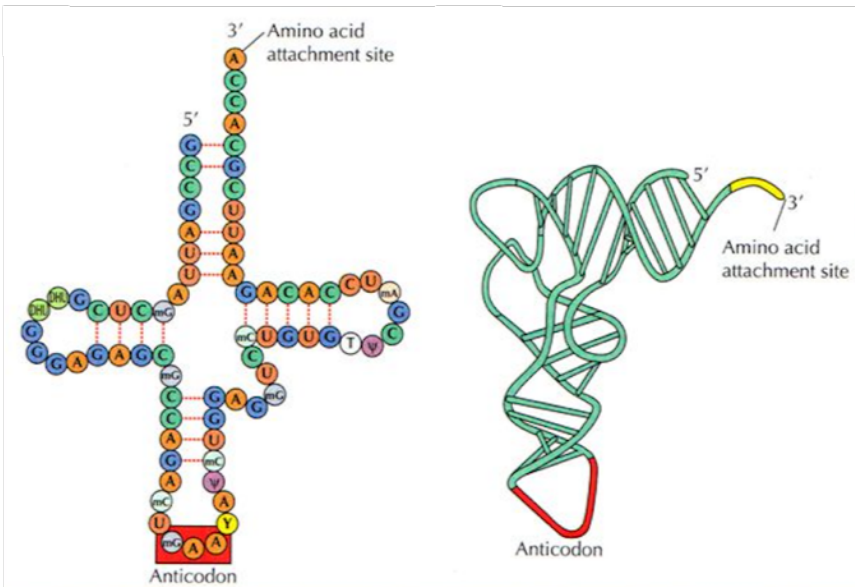
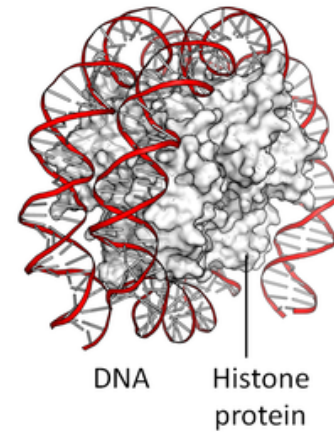
```
d1dhfa_ VPEKNRPLKGRINLVLSRELKEPPQGAHFLSRSLDDALKLTEQPELANKVDMVWIVGGSSVYKEAMNHP
d8dfr_  VPEKNRPLKDRINIVLSRELKEAPKGAHYLSKSLDDALALLDSPELKSVDVWVIVGGTAVYKAAMEKP
d4dfra_ ---G-RPLPGRKNIILS-SQPGTDDRVTWVKSVDDEAIAACGDVPE-----EIMVIGGGRVYEQFLPKA
d3dfr_  ---PKRPLPERTNVVLTHQEDYQAQGA-VVVHDVAAVFAYAKQHLDQ----ELVIAGGAQIFTAFKDDV
```

```
d1dhfa_ -PEKNRPLKGRINLVLSRELKEPPQGAHFLSRSLDDALKLTEQPELANKVDMVWIVGGSSVYKEAMNHP
d8dfr_  -PEKNRPLKDRINIVLSRELKEAPKGAHYLSKSLDDALALLDSPELKSVDVWVIVGGTAVYKAAMEKP
d4dfra_ -G---RPLPGRKNIILSSSQPGTDDRVTWVKSVDDEAIAACGDVPE-----IMVIGGGRVYEQFLPKA
d3dfr_  -P--KRPLPERTNVVLTHQEDYQAQGA-VVVHDVAAVFAYAKQHLD----QELVIAGGAQIFTAFKDDV
```

Molecular Biology Information: Macromolecular Structure

- DNA/RNA/Protein
 - Mostly protein

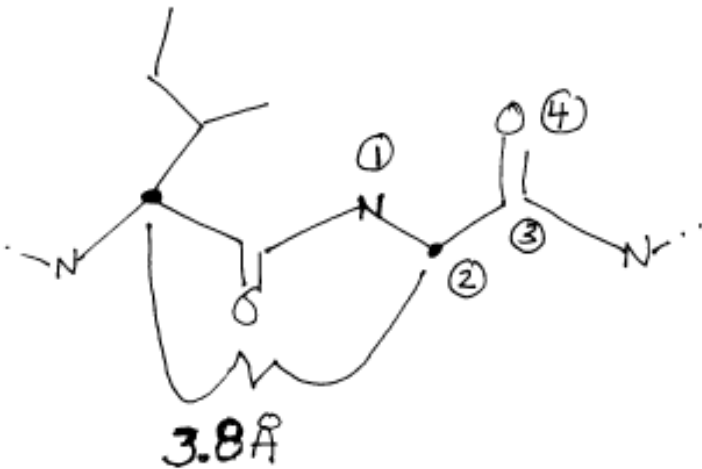
(RNA Adapted From D Soll Web Page,
Right Hand Top Protein from M Levitt web page)



Molecular Biology Information: Protein Structure Details

- Statistics on Number of XYZ triplets
 - 200 residues/domain => 200 CA atoms, separated by 3.8 Å
 - Avg. Residue is Leu: 4 backbone atoms + 4 sidechain atoms, 150 cubic Å
=> ~1500 xyz triplets (=8x200) per p
 - >100K Domains, ~1200 folds (scop 1)

ATOM	1	C	ACE	0	9.401	30.166	60.595	1.00	49.88	1GKY	67
ATOM	2	O	ACE	0	10.432	30.832	60.722	1.00	50.35	1GKY	68
ATOM	3	CH3	ACE	0	8.876	29.767	59.226	1.00	50.04	1GKY	69
ATOM	4	N	SER	1	8.753	29.755	61.685	1.00	49.13	1GKY	70
ATOM	5	CA	SER	1	9.242	30.200	62.974	1.00	46.62	1GKY	71
ATOM	6	C	SER	1	10.453	29.500	63.579	1.00	41.99	1GKY	72
ATOM	7	O	SER	1	10.593	29.607	64.814	1.00	43.24	1GKY	73
ATOM	8	CB	SER	1	8.052	30.189	63.974	1.00	53.00	1GKY	74
ATOM	9	OG	SER	1	7.294	31.409	63.930	1.00	57.79	1GKY	75
ATOM	10	N	ARG	2	11.360	28.819	62.827	1.00	36.48	1GKY	76
ATOM	11	CA	ARG	2	12.548	28.316	63.532	1.00	30.20	1GKY	77
ATOM	12	C	ARG	2	13.502	29.501	63.500	1.00	25.54	1GKY	78
...											
ATOM	1444	CB	LYS	186	13.836	22.263	57.567	1.00	55.06	1GKY1510	
ATOM	1445	CG	LYS	186	12.422	22.452	58.180	1.00	53.45	1GKY1511	
ATOM	1446	CD	LYS	186	11.531	21.198	58.185	1.00	49.88	1GKY1512	
ATOM	1447	CE	LYS	186	11.452	20.402	56.860	1.00	48.15	1GKY1513	
ATOM	1448	NZ	LYS	186	10.735	21.104	55.811	1.00	48.41	1GKY1514	
ATOM	1449	OXT	LYS	186	16.887	23.841	56.647	1.00	62.94	1GKY1515	
TER	1450		LYS	186						1GKY1516	



Molecular Biology Information: Whole Genomes

- The Revolution Driving Everything

Fleischmann

R. D., Adams, M. D., White, O., Clayton, R. A., Kirkness, E. F., Kerlavage, A. R., Bult, C. J., Tomb, J. F., Dougherty, B. A., Merrick, J. M., McKenney, K., Sutton, G., Fitzhugh, W., Fields, C., Gocayne, J. D., Scott, J., Shirley, R., Liu, L. I., Glodek, A., Kelley, J. M., Weidman, J. F., Phillips, C. A., Spriggs, T., Hedblom, E., Cotton, M. D., Utterback, T. R., Hanna, M. C., Nguyen, D. T., Saudek, D. M., Brandon, R. C., Fine, L. D., Fritchman, J. L., Fuhrmann, J. L., Geoghagen, N. S. M., Gnehm,

C. L., McDonald, L. A., Small, K. V., Fraser, C. M., Smith, H. O. & **Venter**, J. C.

(1995). "Whole-genome random sequencing and assembly of

Haemophilus influenzae rd." *Science* 269: 496-512.

(Picture adapted from TIGR website, <http://www.tigr.org>)

- Timeline

1995, HI (bacteria): 1.6 Mb & 1600 genes done

1997, yeast: 13 Mb & ~6000 genes for yeast

1998, worm: ~100Mb with 19 K genes

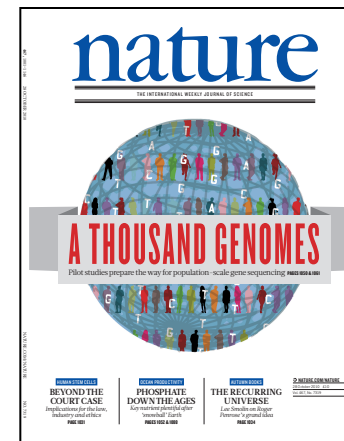
1999: >30 completed genomes!

2000, draft human

2003, human: 3 Gb & 100 K genes...

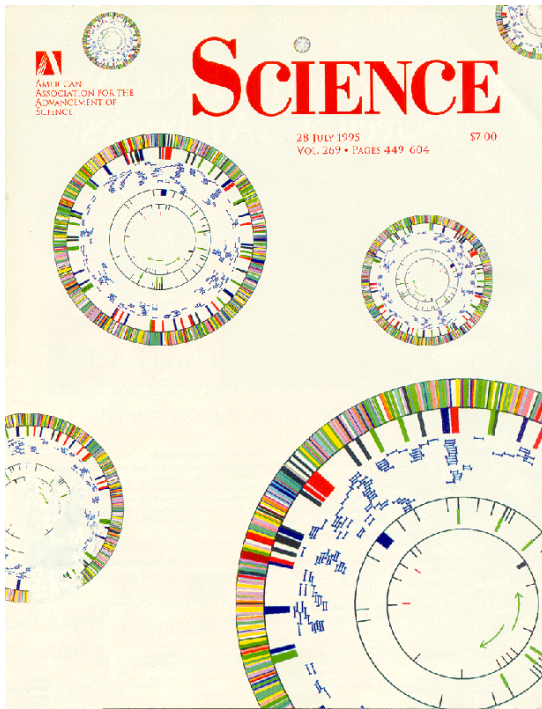
2010, 1000 human genomes!

2017, 13K human genomes



1995

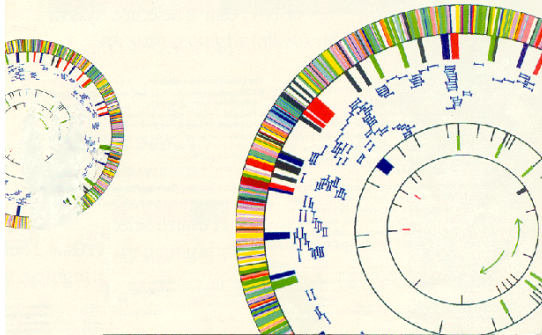
Bacteria,
1.6 Mb,
~1600 genes
[*Science* 269: 496]



A
Bioinformatics
prediction that
came true!

1997

Eukaryote,
13 Mb,
~6K genes
[*Nature* 387: 1]



real thing, Apr '00



'98 spoof

1998

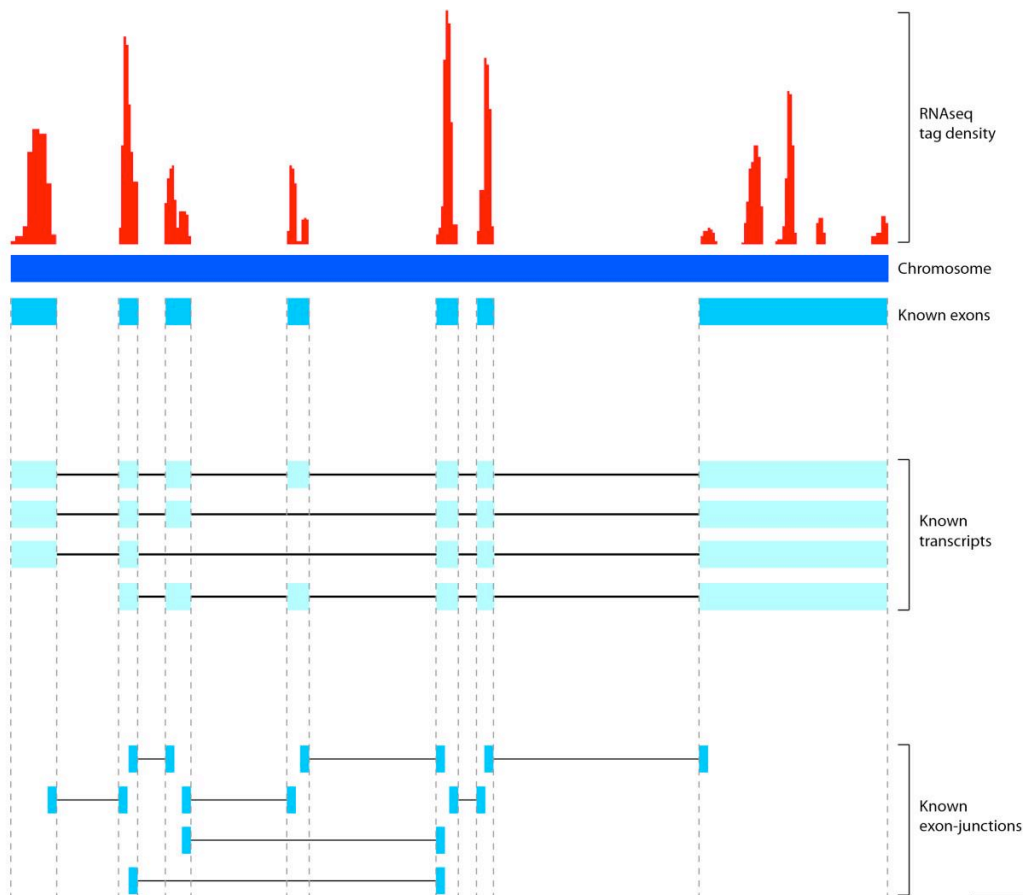
Animal,
~100 Mb,
~20K genes
[*Science* 282:
1945]



2000?

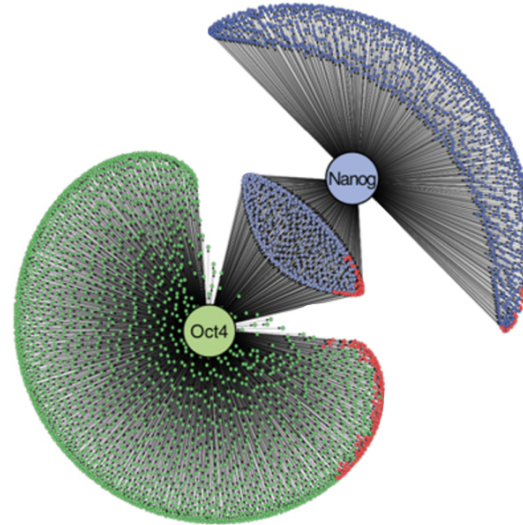
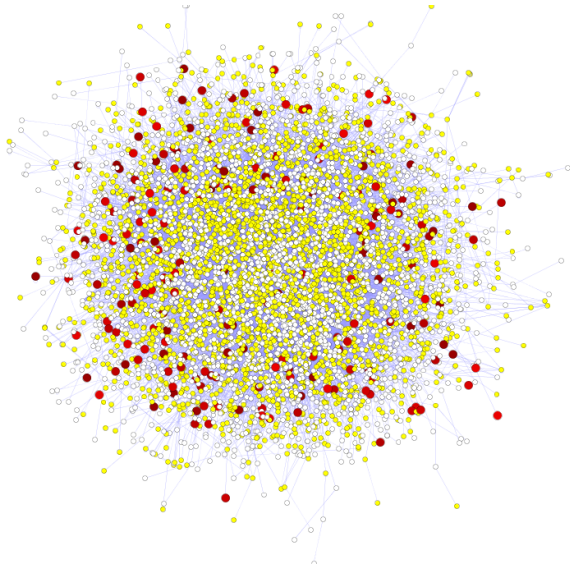
Human,
~3 Gb,
~20K genes

Gene Expression Data: On & Off



- Early experiments yeast
 - Complexity at 10 time points,
6000 x 10 = 60K floats
- Then tiling array technology
 - 50 M data points to tile the human genome at ~50 bp res.
- Now Next-Gen Sequencing (RNAseq)
 - 10M+ reads on the human genome, counts
- Can only sequence genome once but can do an infinite variety of expression experiments

Molecular Networks: Connectivity



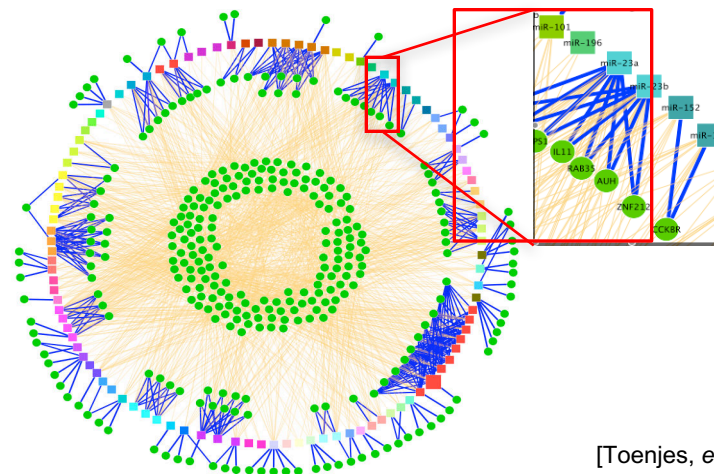
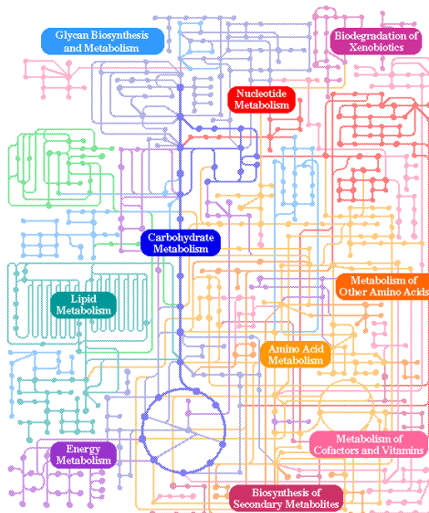
Regulatory Networks
Get readouts of where proteins bind to DNA :
Chip-chip then chip-seq

Protein Interaction Networks

For yeast: $6000 \times 6000 / 2 \sim 18M$ possible interactions
(maybe $\sim 30K$ real)

Protein-protein Interaction networks

TF-target-gene Regulatory networks



Metabolic pathway networks

miRNA-target networks

[Toenjes, *et al*, *Mol. BioSyst.* (2008); Jeong *et al*, *Nature* (2001); [Horak, *et al*, *Genes & Development*, 16:3017-3033; DeRisi, Iyer, and Brown, *Science*, 278:680-686; Descalzo, *et al*, *Mol Syst Biol*, 9:694]

Molecular Biology Information: Other Integrative Data

- Information to understand genomes

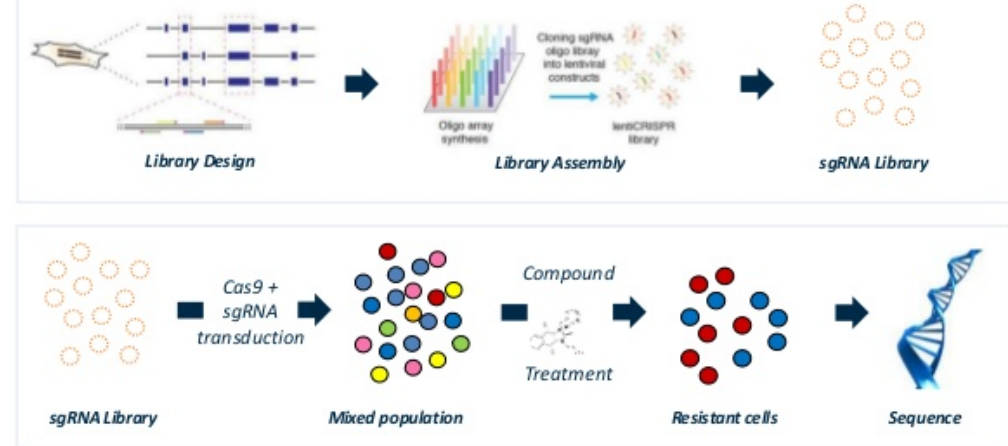
- Whole Organisms
Phylogeny, traditional zoology
- Environments, Habitats, ecology

- Phenotype Experiments (large-scale KOs, transposons)

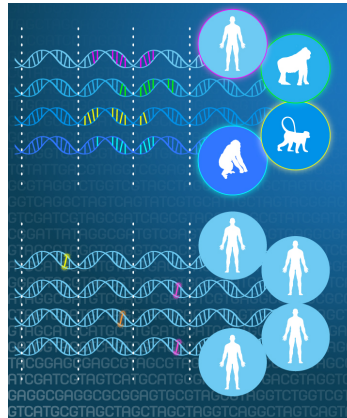
- The Literature (MEDLINE)

- The Future....

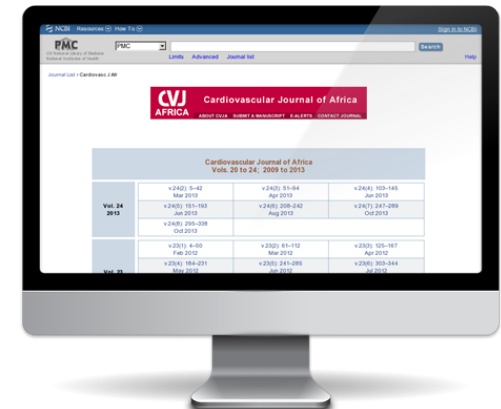
CRISPR phenotypic screen



Comparative genomics



Literature knowledgebase



What is Bioinformatics?

- (*Molecular*) **Bio - informatics**

- One idea for a definition?

Bioinformatics is conceptualizing **biology in terms of molecules** (in the sense of physical-chemistry) and then applying **“informatics” techniques** (derived from disciplines such as applied math, CS, and statistics) to **organize, mine, model & understand the information associated** with these molecules, **on a large-scale.**

- Bioinformatics is a practical discipline with many **applications.**

Seq Universe

[from Heidi Sofia, NHGRI]

SRA >1 petabyte

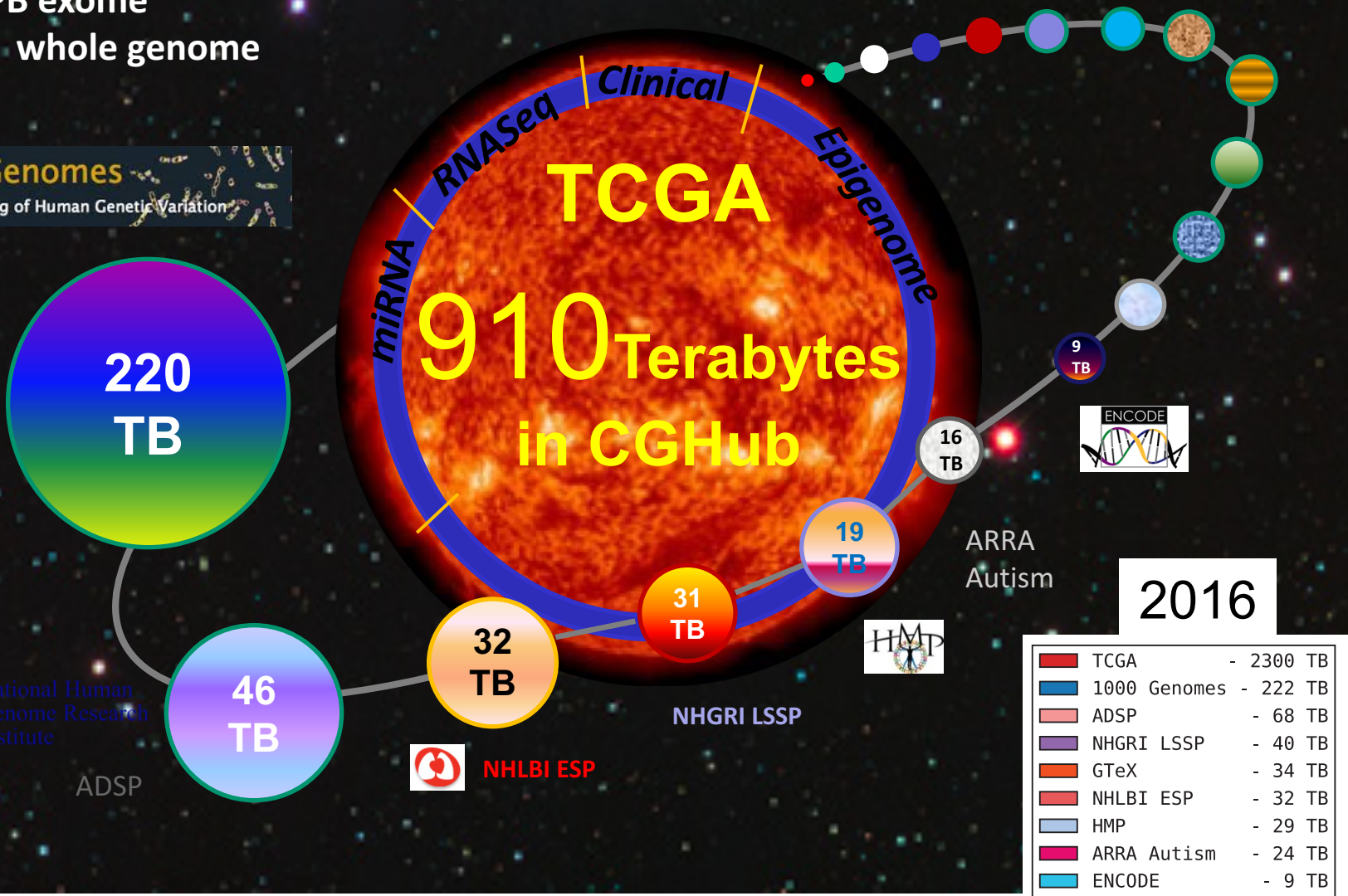
TCGA endpoint: ~2.5 Petabytes

~1.5 PB exome

~1 PB whole genome

1000 Genomes

A Deep Catalog of Human Genetic Variation



National Human Genome Research Institute

ADSP



NHLBI ESP

NHGRI LSSP



ARRA
Autism



What is Bioinformatics?

- (*Molecular*) **Bio - informatics**

- One idea for a definition?

Bioinformatics is conceptualizing **biology in terms of molecules** (in the sense of physical-chemistry) and then applying **“informatics” techniques** (derived from disciplines such as applied math, CS, and statistics) to **organize, mine, model & understand the information associated** with these molecules, **on a large-scale.**

- Bioinformatics is a practical discipline with many **applications.**

General Types of “Informatics” techniques in Computational Biology – a mix between **mining** & **modeling**

- **Databases**

- Building, Querying
- Representing Complex data

- **Data mining**

- Machine Learning techniques
- Clustering & Tree construction
- Rapid Text String Comparison & textmining
- Detailed statistics of significance & association

- **Network Analysis**

- Analysis of Topology (eg Hubs)
- Predicting Connectivity

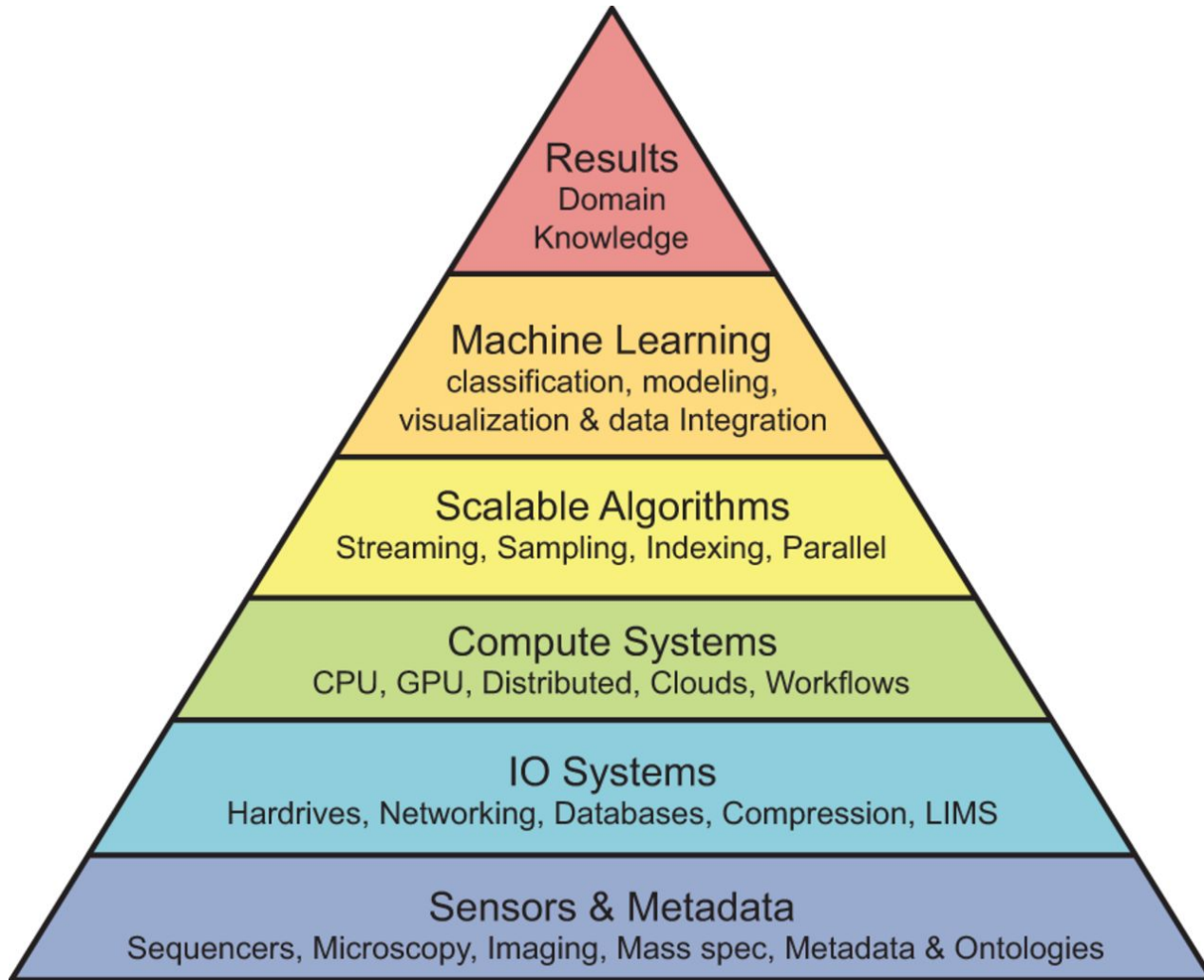
- **Structure Analysis & Geometry**

- Graphics (Surfaces, Volumes)
- Comparison & 3D Matching (Vision, recognition, docking)

- **Physical Modeling**

- Newtonian Mechanics
- Minimization & Simulation
- Modeling Chemical Reactions & Cellular Processes

Data science analysis stack.



Michael C. Schatz *Genome Res.* 2015;25:1417-1422



Bioinformatics

Key Practical Applications

What is Bioinformatics?

- (*Molecular*) **Bio - informatics**

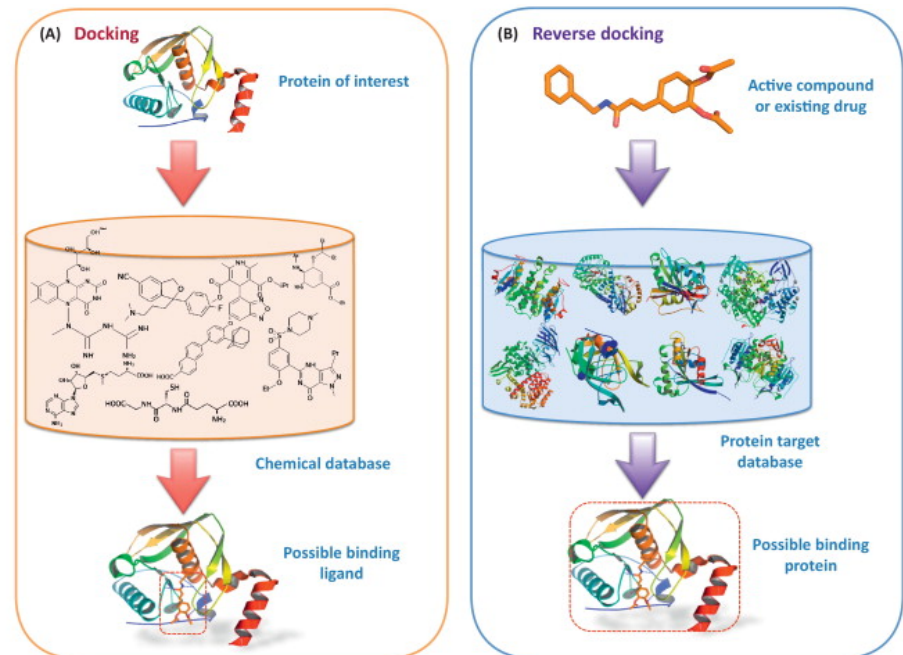
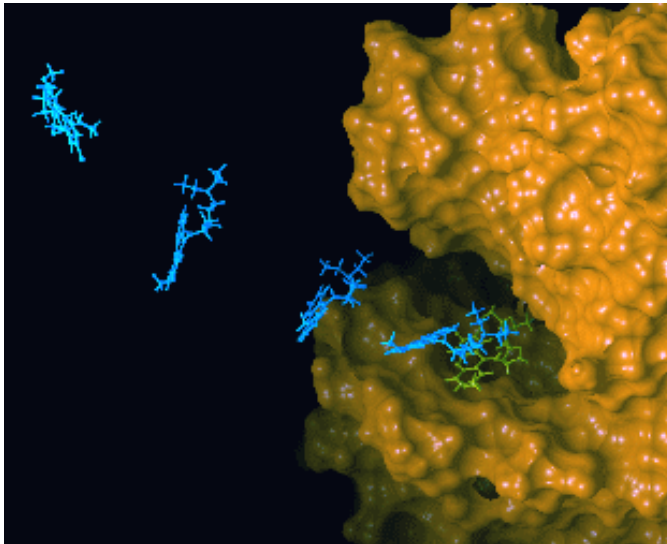
- One idea for a definition?

Bioinformatics is conceptualizing **biology in terms of molecules** (in the sense of physical-chemistry) and then applying **“informatics” techniques** (derived from disciplines such as applied math, CS, and statistics) to **organize, mine, model & understand the information associated** with these molecules, **on a large-scale.**

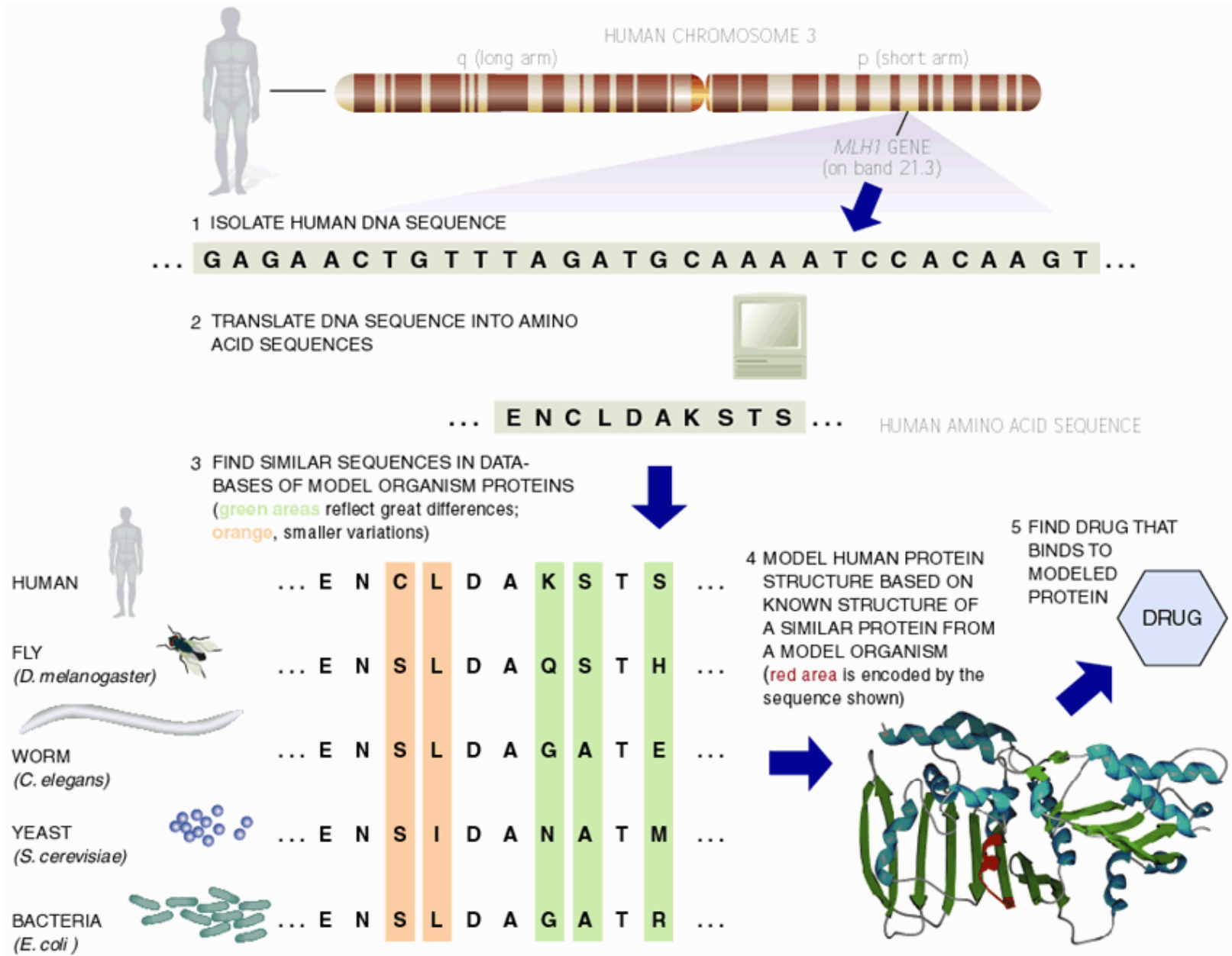
- Bioinformatics is a practical discipline with many **applications.**

Major Application I: Designing Drugs from Structural Targets

- Understanding how structures bind other molecules
- Designing inhibitors using docking, structure modeling
- *In silico* screens of chemical and protein databases

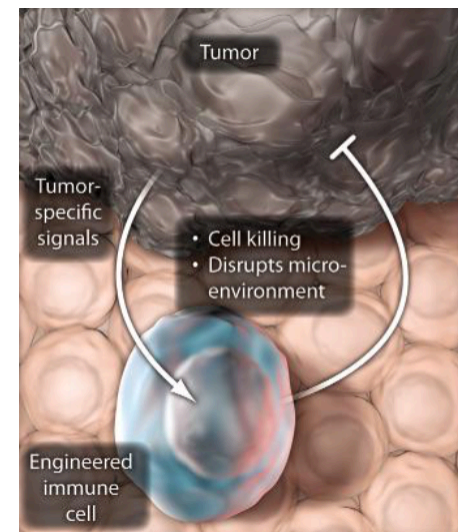
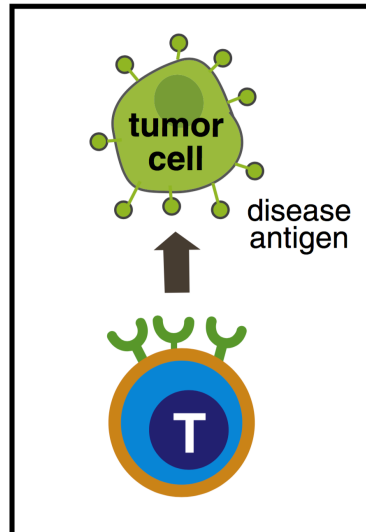
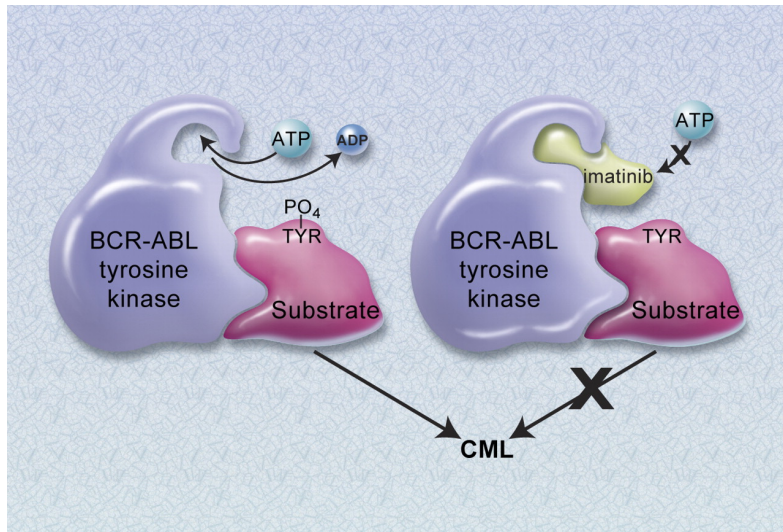


Major Application II: Finding Homologs, to Find Experimentally Tractable Gene Targets



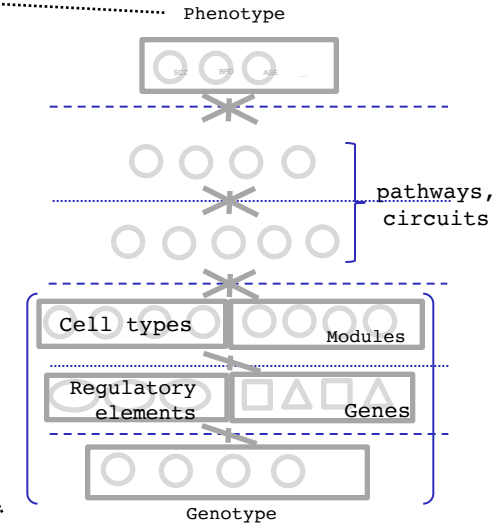
Major Application III: Customizing treatment in oncology

- Identifying disease causing mutations in individual patients
- Designing targeted therapeutics
 - e.g. BCR-abl and Gleevec
 - Cancer immunotherapies targeting neoantigens



Major Application IV: Finding molecular mechanisms & drug targets for diseases we know little about (Neuro-psychiatric Diseases)

Disease	Heritability*	Molecular Mechanisms
Schizophrenia	81%	-
Bipolar disorder	70%	-
Alzheimer's disease	58 - 79%	Apolipoprotein E (APOE), Tau
Hypertension	30%	Renin–angiotensin–aldosterone
Heart disease	34-53%	Atherosclerosis, VCAM-1
Stroke	32%	Reactive oxygen species (ROS), Ischemia
Type-2 diabetes	26%	Insulin resistance
Breast Cancer	25-56%	BRCA, PTEN



Many psychiatric conditions are highly heritable

Schizophrenia: up to 80%

But we don't understand basic molecular mechanisms underpinning this association
(in contrast to many other diseases such as cancer & heart disease)

Moreover, current models substantially underestimate heritability using genetic data

Schizophrenia : ~25%

Thus, interested in developing predictive models of psychiatric traits which:

Use observations at intermediate (molecular levels) levels to inform latent structure.

Use the predictive features of these “molecular endo phenotypes” to begin to suggest actors involved in mechanism

Major Application IV: Finding molecular mechanisms & drug targets for diseases we know little about (Neuro-psychiatric Diseases)

Disease	Heritability*	Molecular Mechanisms
Schizophrenia	81%	-
Bipolar disorder	70%	-
Alzheimer's disease	58 - 79%	Apolipoprotein E (APOE), Tau
Hypertension	30%	Renin-angiotensin-aldosterone
Heart disease	34-53%	Atherosclerosis, VCAM-1
Stroke	32%	Reactive oxygen species (ROS), Ischemia
Type-2 diabetes	26%	Insulin resistance
Breast Cancer	25-56%	BRCA, PTEN



Many psychiatric conditions are highly heritable

Schizophrenia: up to 80%

But we don't understand basic molecular mechanisms underpinning this association

(in contrast to man)

Moreover, current model

Schizophrenia : ~2

Thus, interested in deve

Use observations :

structure.

Use the predictive features of these "molecular endo phenotypes" to begin to

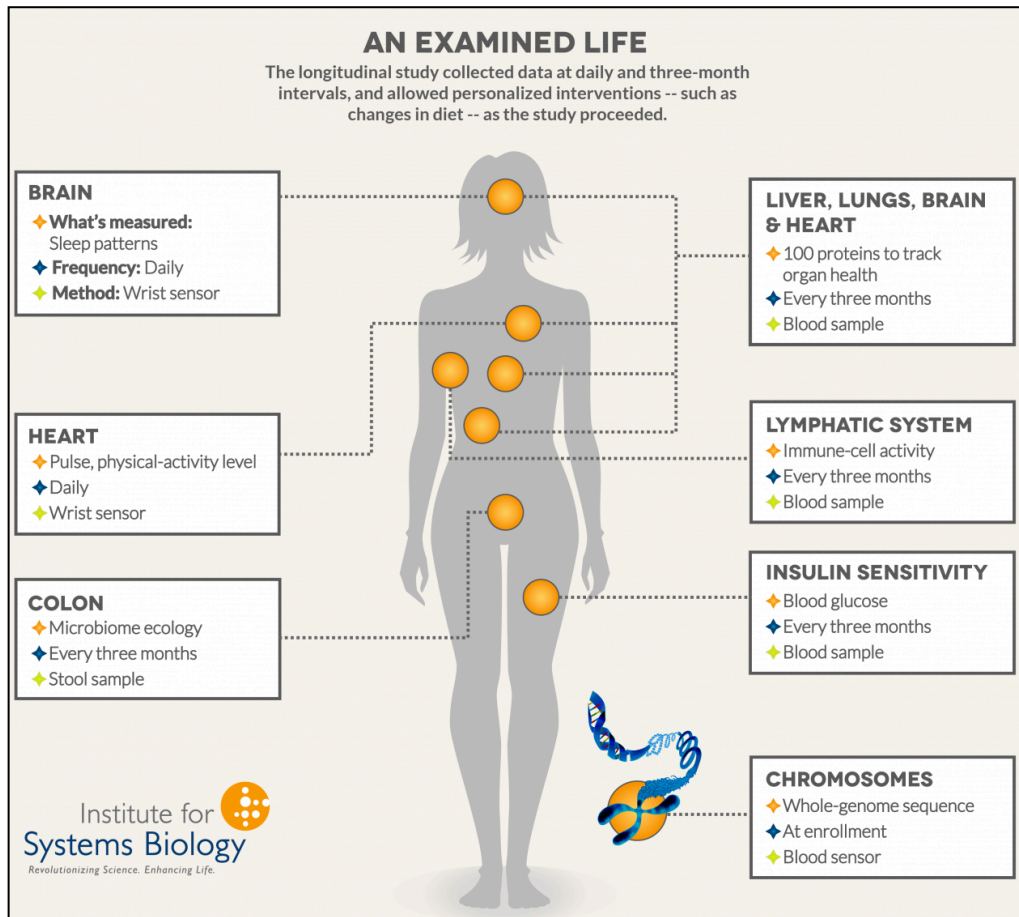
suggest actors involved in mechanism

**Recent Rollout in Science
addressing this, involving
many Yale Researchers**



Major Application V: Holistic Personal Genome Characterization, in Normal Individuals

- Mental disease & cancer are two extremes with respect to genomics (CEN, 92: 26)
 - Many other conditions in between, often involving interaction with the environment
- Pers. Genome Characterization
 - Identify mutations in personal genomes (SNPs, SVs, &c)
 - Estimate phenotypic (deleterious or protective) impact of variants.
 - Compare one person to wider population.
- Track changes over time & consider interaction w/ environment
 - Transcriptome studies
 - Longitudinal health studies (e.g. 100K wellness project, Framingham Heart Study)



(Figure from Institute for Systems Biology)

The Key Application in Personal Genomics

More Details

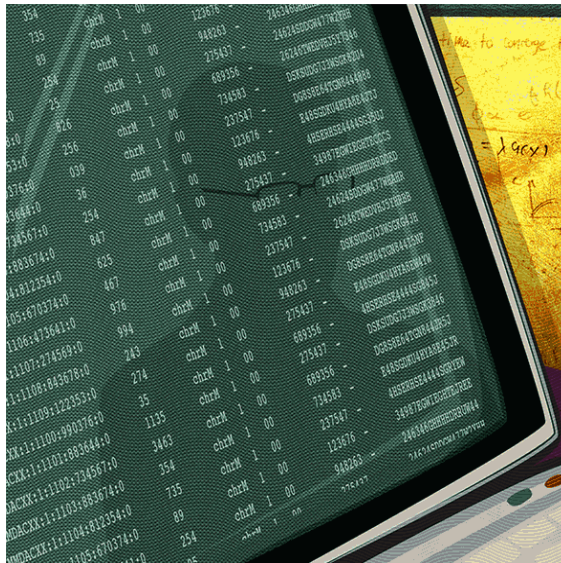
Personal Genomics

as an an organizing theme for this class

- A personal genome can reveal a lot about an individual.
 - Disease risks, ancestry, personal traits, etc.
- Personal genome annotation combined with multi-omic and longitudinal health data can inform new links between genotype and phenotype relevant to an individual and the larger population.
- Genomic privacy will become increasingly important as precision medicine becomes more common.
- In this class, we will look at how to identify key genomic variants with the most impact.
- We will also use analysis techniques including systems and network modeling as well as structural modeling to contextualize and interpret the mechanisms through which these variants impact health.

Analyzing Carl Zimmer's genome

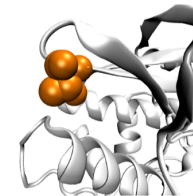
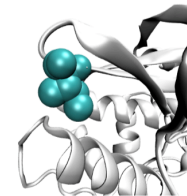
CARL ZIMMER'S GAME OF GENOMES SEASON 1



SNV

AAGCT → ACGCT

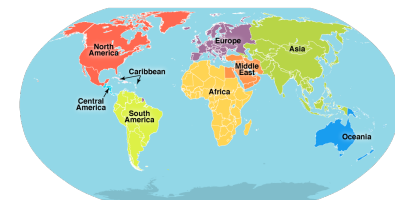
Protein
Structure



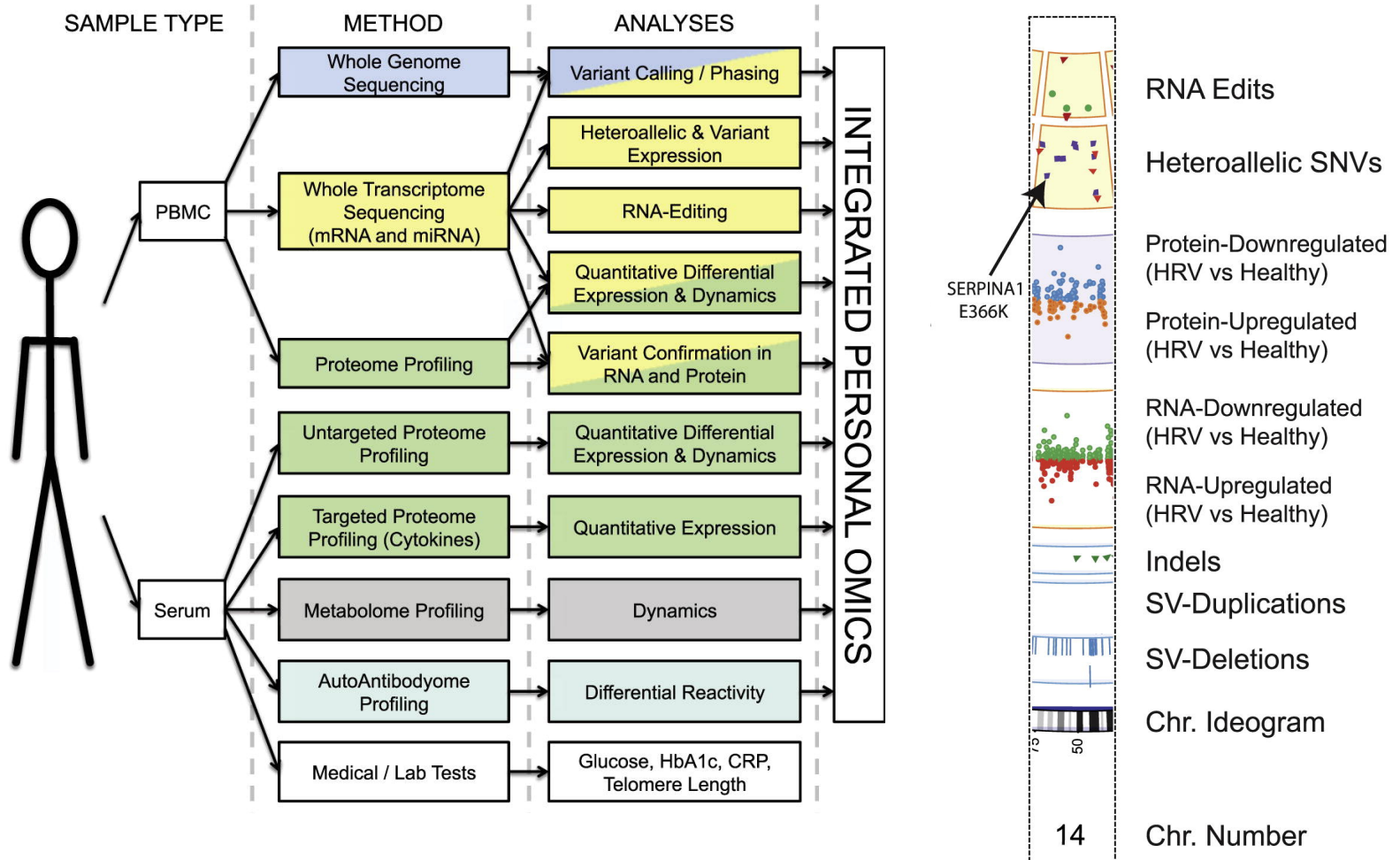
Wild-type

Mutated

Ancestry



Personal Omics Profiling



(Figure from Chen et al. Cell 2012)

Human Genetic Variation

A Cancer Genome



A Typical Genome



Population of 2,504 peoples



Origin of Variants

	Coding	Non-coding
Germ-line	22K	4.1 – 5M
Somatic	~50	5K



Driver (~0.1%)

Class of Variants

SNP	3.5 – 4.3M
Indel	550 – 625K
SV	2.1 – 2.5K (20Mb)
Total	4.1 – 5M

Prevalence of Variants



Rare* (1-4%)

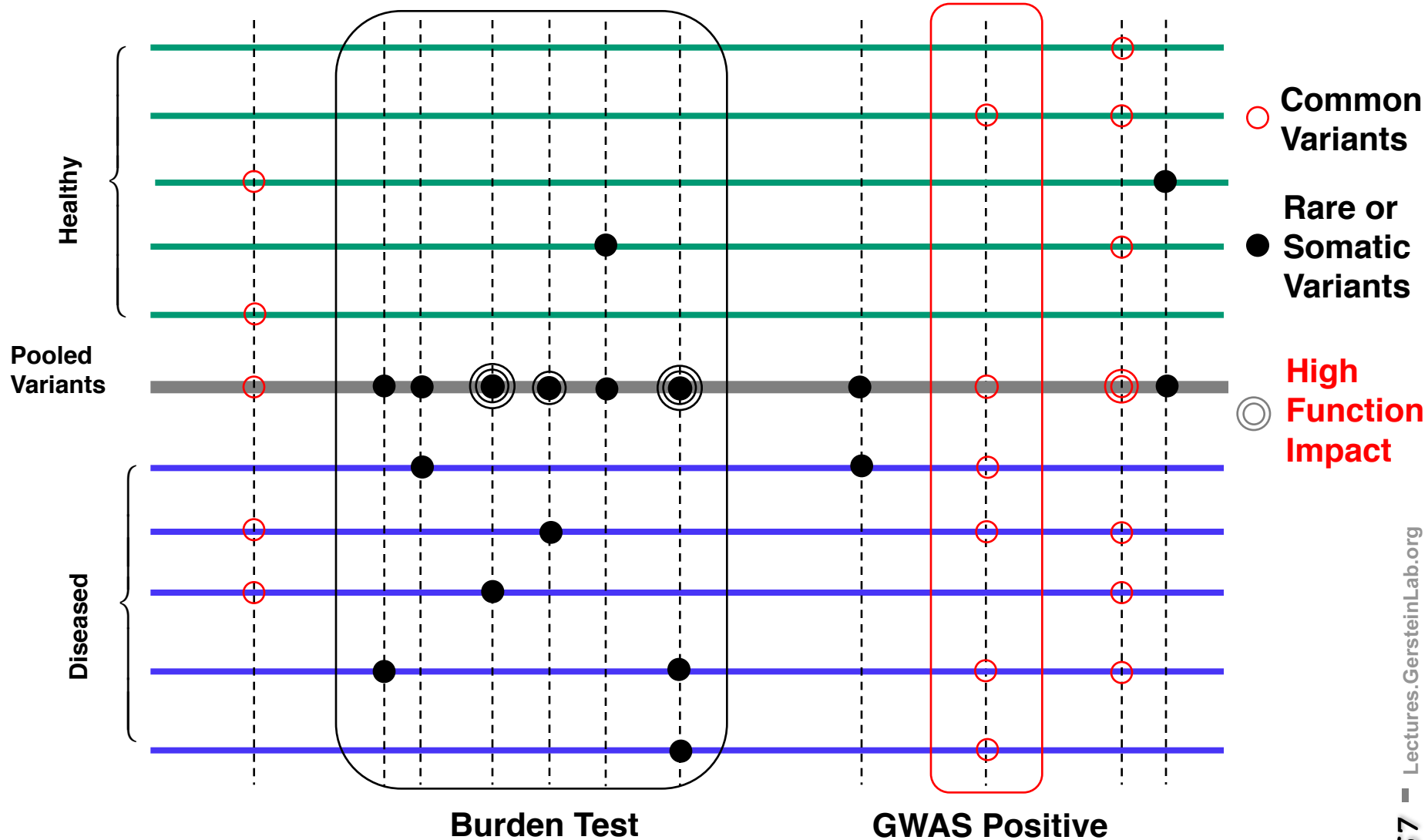
SNP	84.7M
Indel	3.6M
SV	60K
Total	88.3M



Rare (~75%)

* Variants with allele frequency < 0.5% are considered as rare variants in 1000 genomes project.

Association of Variants with Diseases

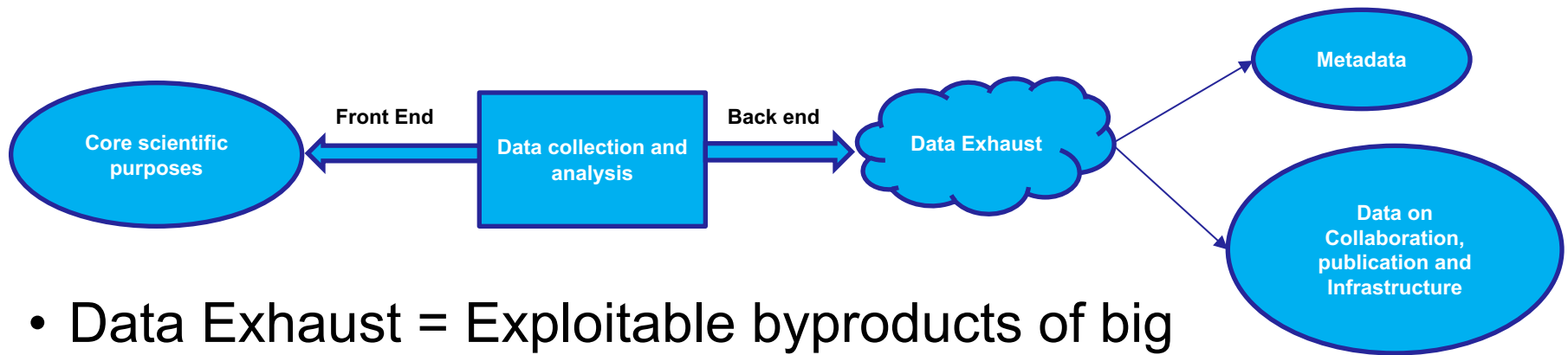


The Other Side of the Data Science Coin:

**The Data Exhaust from
Personal Genomics
(privacy & SOS)**

Core Qs v Creative Use of the Data

Data Exhaust



- Data Exhaust = Exploitable byproducts of big data collection and analysis
- Creative use of Data is key to Data Science !
- Aspects of Privacy but also Science of Science

2-sided nature of functional genomics data: Analysis can be very **General/Public** or **Individual/Private**



- **General quantifications** related to overall aspects of a condition – ie gene activity as a function of:
 - Developmental stage, Evolutionary relationships, Cell-type, Disease
- **Above are not tied to an individual's genotype. However, data is derived from individuals & tagged with their genotypes**
- (Note, a few calculations aim to use explicitly genotype to derive general relations related to sequence variation & gene expression - eg allelic activity)

Genomics has similar "Big Data" Dilemma in the Rest of Society

- Sharing & "peer-production" is central to success of many new ventures, with the same risks as in genomics
 - **EG web search**: Large-scale mining essential



- We confront privacy risks every day we access the internet

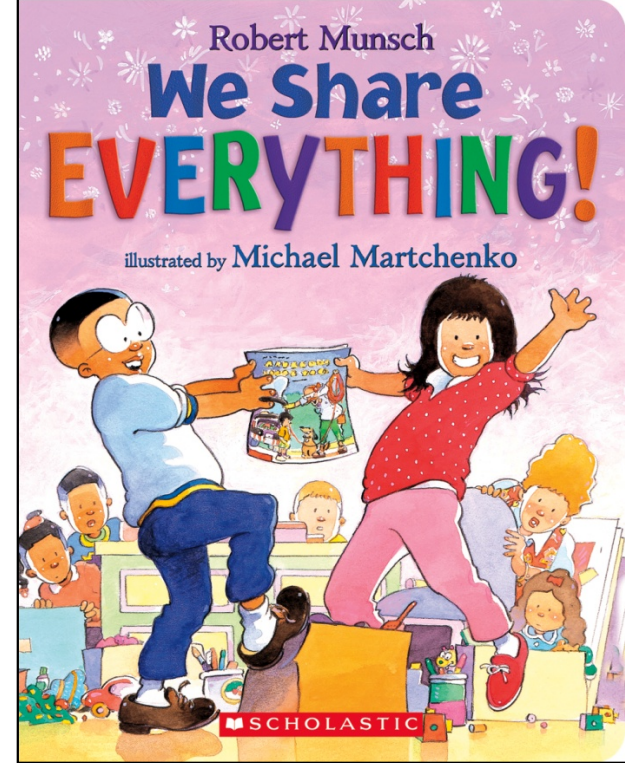
Tricky Privacy Considerations in Personal Genomics

- **Genetic Exceptionalism :**
The Genome is very fundamental data, potentially very revealing about one's identity & characteristics
- **Personal Genomic info. essentially meaningless currently but will it be in 20 yrs? 50 yrs?**
 - Genomic sequence very revealing about one's children. Is true consent possible?
 - Once put on the web it can't be taken back
- **Culture Clash:**
Genomics historically has been a proponent of “open data” but not clear personal genomics fits this.
 - Clinical Medline has a very different culture.
- **Ethically challenged** history of genetics
 - Ownership of the data & what consent means (Hela)
 - Could your genetic data give rise to a product line?



The Other Side of the Coin: Why we should share

- Sharing helps **speed research**
 - Large-scale mining of this information is important for medical research
 - Privacy is cumbersome, particularly for big data
- Sharing is important for **reproducible research**
- Sharing is useful for **education**
 - More fun to study a known person's genome
 - Eg Zimmer's Game of Genomes in STAT



[Yale Law Roundtable ('10). *Comp. in Sci. & Eng.* 12:8; D Greenbaum & M Gerstein ('09). *Am. J. Bioethics*; D Greenbaum & M Gerstein ('10). *SF Chronicle*, May 2, Page E-4; Greenbaum et al. *PLOS CB* ('11)]



The Dilemma

[Economist, 15 Aug '15]

- The individual (harmed?) v the collective (benefits)
 - But do sick patients care about their privacy?
- How to balance risks v rewards - Quantification
 - What is acceptable risk?
Can we quantify leakage?
 - Ex: photos of eye color
 - Cost Benefit Analysis

Current Social & Technical Solutions: The quandary where are now

• **Closed Data** Approach

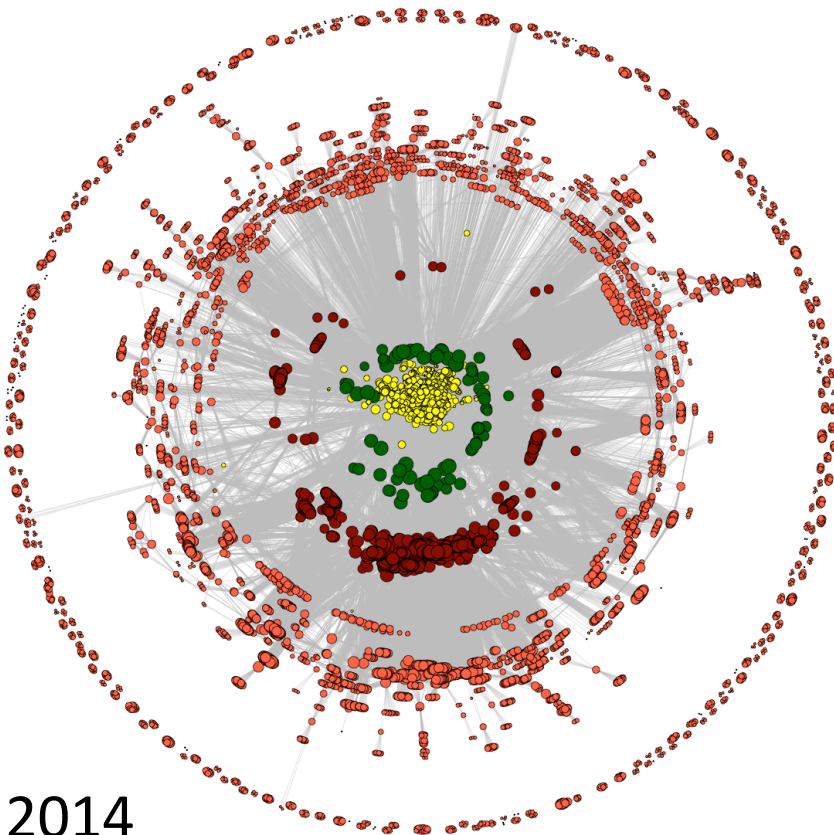
- Consents
- “Protected” distribution via dbGAP
- Local computes on secure computer
- Issues with Closed Data
 - Non-uniformity of consents & paperwork
 - Different international norms, leading to confusion
 - Encryption & computer security creates burdensome requirements on data sharing & large scale analysis
 - Many schemes get “hacked”

• **Open Data**

- Genomic “test pilots” (ala PGP)?
 - Sports stars & celebrities?
- Some public data & data donation is helpful but is this a realistic solution for an unbiased sample of ~1M

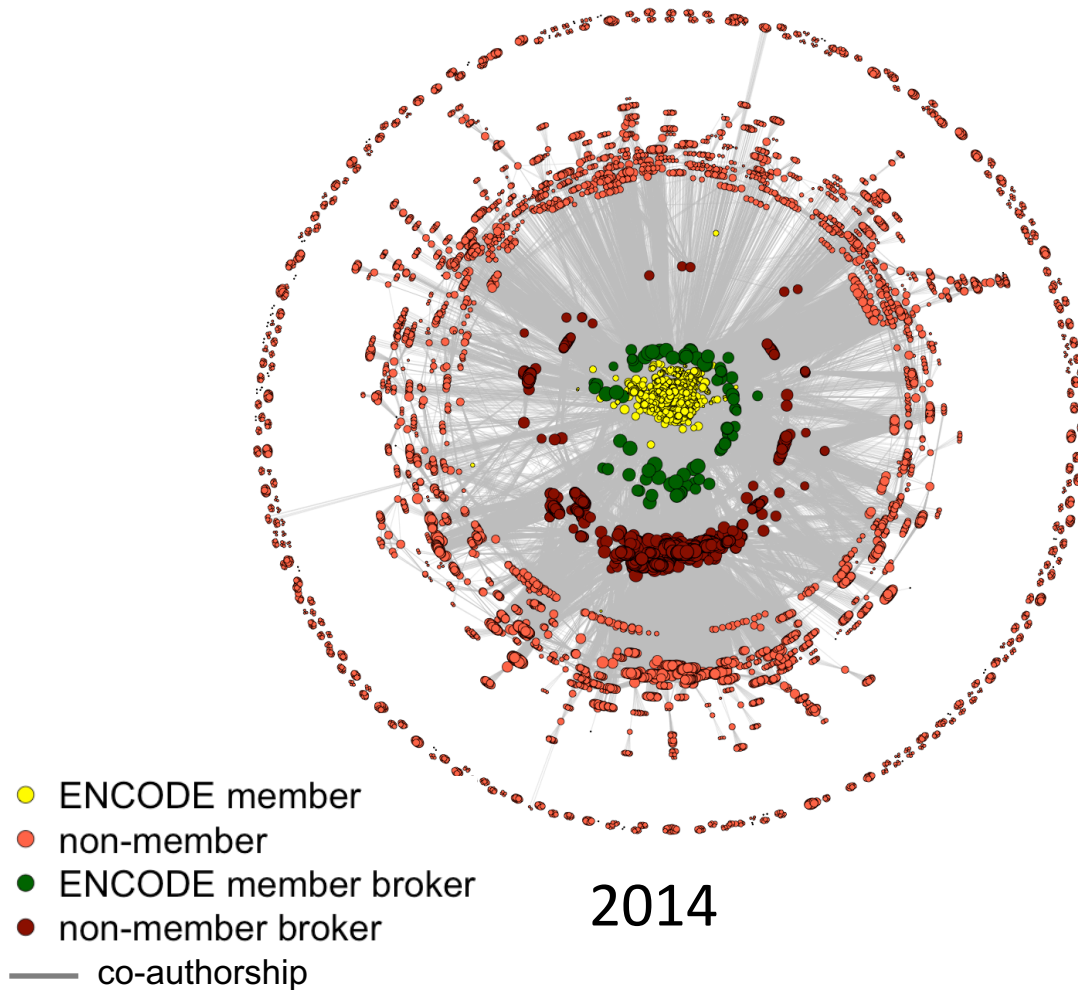
Another Exhaust Mining Application: Using Science to Study Science (SOS)

- ENCODE member
- non-member
- ENCODE member broker
- non-member broker
- co-authorship

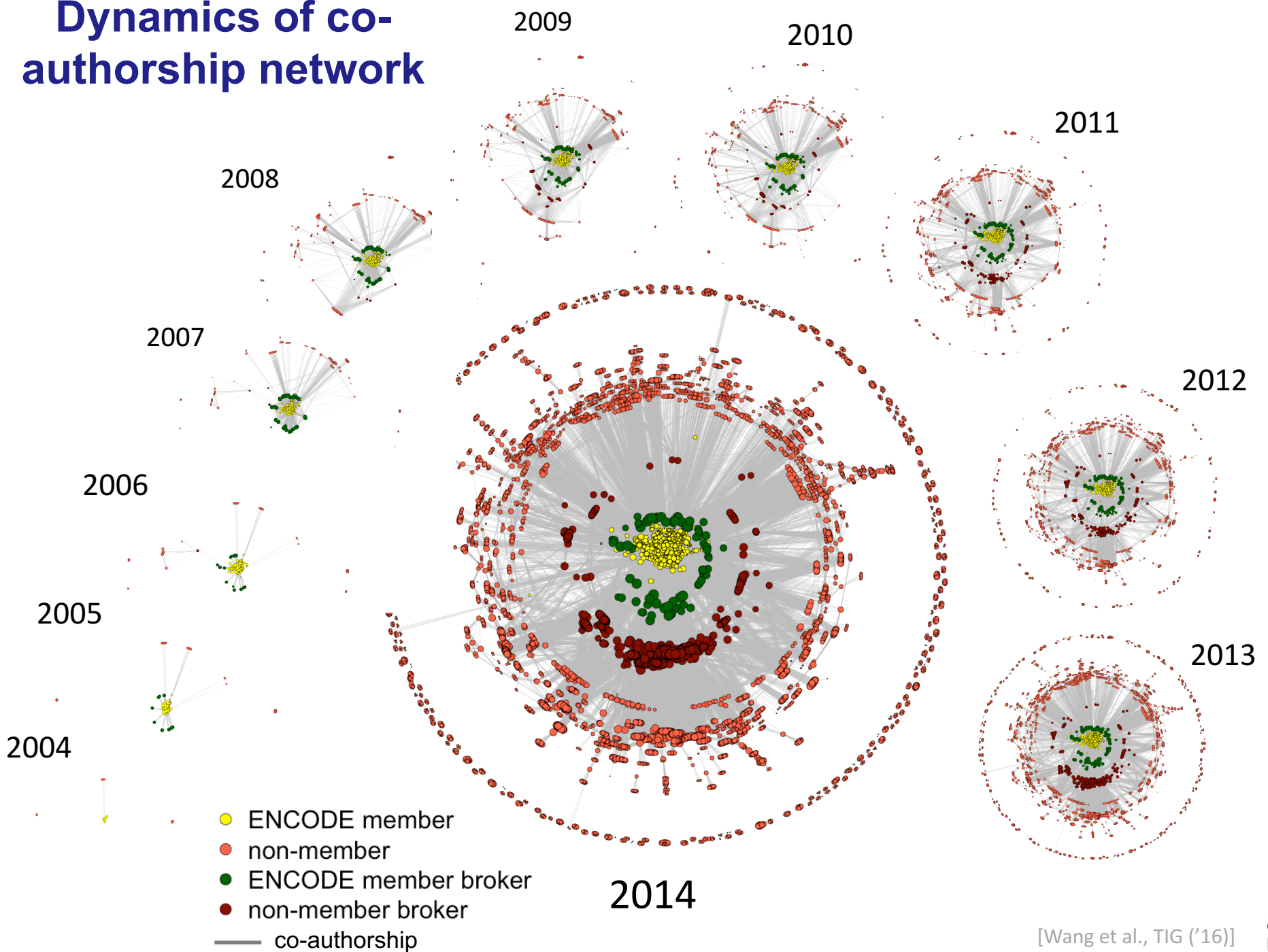


- Mining output of science (Scientific Publications) to understand how science works as a social enterprise
- Co-authorship network of members of the human genome annotation group (ENCODE) & users of this groups data

Dynamics of co-authorship network



Dynamics of co-authorship network



[Wang et al., TIG ('16)]

Extra stuff

Class Web Page

GersteinLab.org/courses/452

Assignment #0 Page

bit.ly/cbb752b19-hw0

Are They or Aren't They Comp. Bio.? (#1)

- (Digital Libraries & Medical Record Analysis
 - Automated Bibliographic Search and Textual Comparison
 - Knowledge bases for biological literature
- (Motif Discovery Using Gibb's Sampling
- (Methods for Structure Determination
 - Computational Crystallography
 - Refinement
 - NMR Structure Determination
 - (Distance Geometry
- (Metabolic Pathway Simulation
- (The DNA Computer

Are They or Aren't They Comp. Bio.? (#1, Answers)

- **(YES?)** Digital Libraries & Medical Record Analysis
 - Automated Bibliographic Search and Textual Comparison
 - Knowledge bases for biological literature
- **(YES)** Motif Discovery Using Gibb's Sampling
- **(NO?)** Methods for Structure Determination
 - Computational Crystallography
 - Refinement
 - NMR Structure Determination
 - **(YES)** Distance Geometry
- **(YES)** Metabolic Pathway Simulation
- **(NO)** The DNA Computer

Are They or Aren't They Comp. Bio.? (#2)

- (Gene identification by sequence characteristics
 - Prediction of splice sites
- (DNA methods in forensics
- (Modeling of Populations of Organisms
 - Ecological Modeling (predator & prey)
- (Modeling the nervous system
 - Computational neuroscience
 - Understanding how brains think & using this to make a better computer
- (Molecular phenotype discovery – looking for gene expression signatures of cancer
 - What if it included non-molecular data such as age ?

Are They or Aren't They Comp. Bio.? (#2, Answers)

- **(YES)** Gene identification by sequence characteristics
 - Prediction of splice sites
- **(YES)** DNA methods in forensics
- **(NO)** Modeling of Populations of Organisms
 - Ecological Modeling (predator & prey)
- **(NO?)** Modeling the nervous system
 - Computational neuroscience
 - Understanding how brains think & using this to make a better computer
- **(YES)** Molecular phenotype discovery – looking for gene expression signatures of cancer
 - What if it included non-molecular data such as age ?

Are They or Aren't They Comp. Bio.? (#3)

- (RNA structure prediction
- (Radiological Image Processing
 - Computational Representations for Human Anatomy (visible human)
- (Artificial Life Simulations
 - Artificial Immunology / Computer Security
 - (Genetic Algorithms in molecular biology
- (Homology Modeling & Drug Docking
- (Char. drugs & other small molecules (QSAR)
- (Computerized Diagnosis based on Pedigrees
- (Processing of NextGen sequencing image files
- (Module finding in protein networks

Are They or Aren't They Comp. Bio.? (#3, Answers)

- **(YES)** RNA structure prediction
- **(NO)** Radiological Image Processing
 - Computational Representations for Human Anatomy (visible human)
- **(NO)** Artificial Life Simulations
 - Artificial Immunology / Computer Security
 - **(NO?)** Genetic Algorithms in molecular biology
- **(YES)** Homology Modeling & Drug Docking
- **(YES)** Char. drugs & other small molecules (QSAR)
- **(NO)** Computerized Diagnosis based on Pedigrees
- **(NO)** Processing of NextGen sequencing image files
- **(YES)** Module finding in protein networks