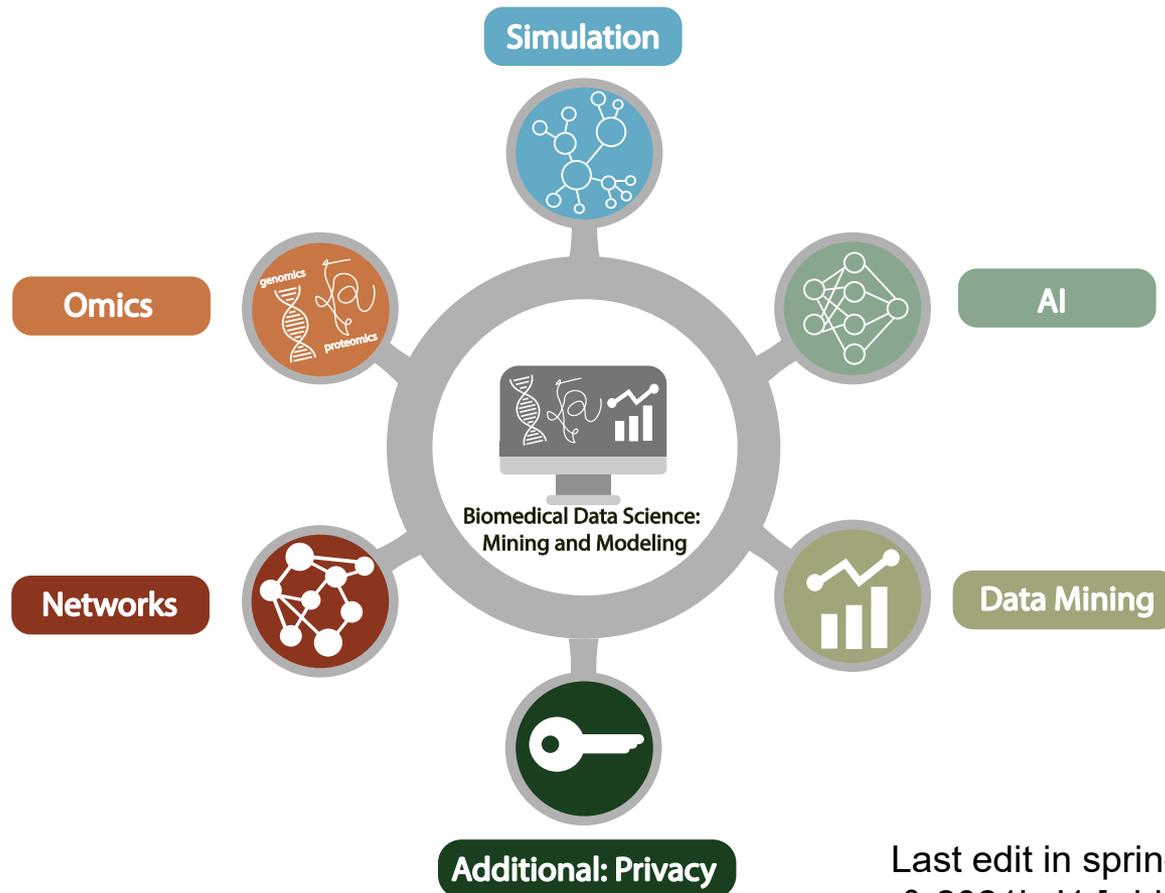


Biomedical Data Science (GersteinLab.org/courses/452)

Introduction (23i1)

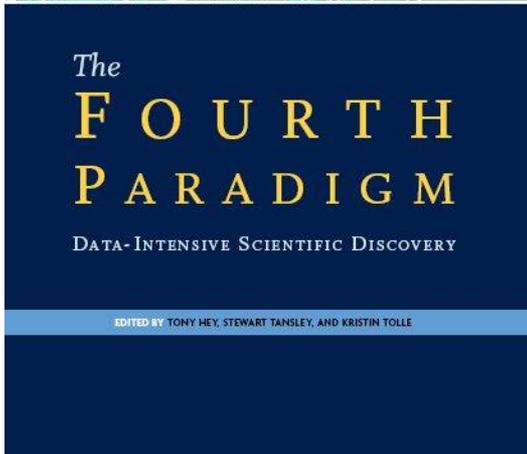


Last edit in spring '23. Similar to 22i1 & 2021's I1 [which has a video]. Has additions beginning at slides 29 & 33, describing course topics, format & goals. Usually takes ~45' with rest of class going over syllabus

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

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

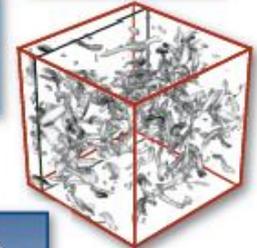
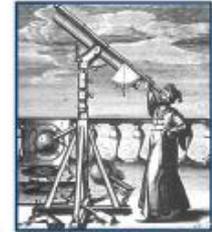
Jim Gray's 4th Paradigm



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

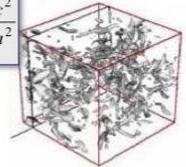
Today:

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

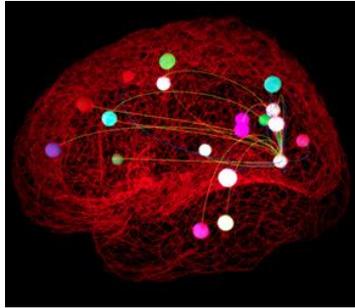
Guest post written by **Quentin Gallivan**
 Quentin Gallivan is CEO of Pentaho Corp., an Orlando, Florida-based provider of business analytics software.

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



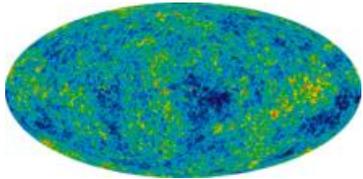
High energy physics - Large Hadron Collider



Neuroscience - The Human Connectome Project



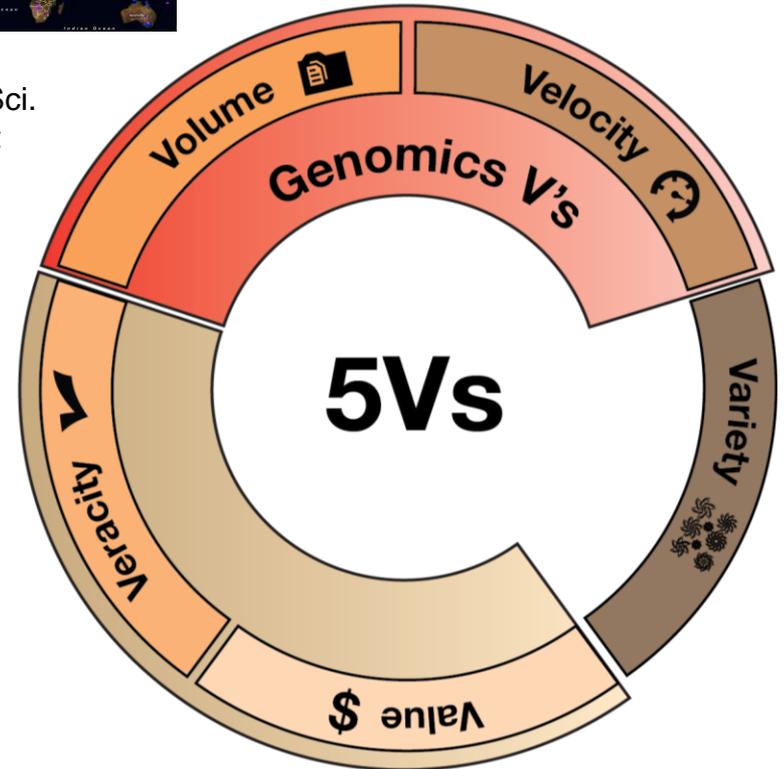
Ecology & Earth Sci. - Fluxnet



Astronomy - Sloan Digital Sky survey

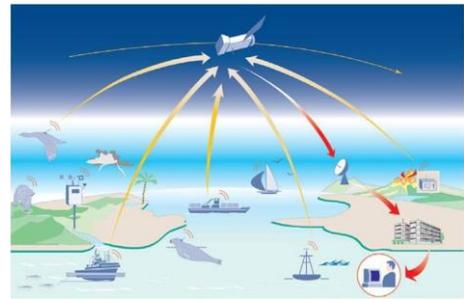
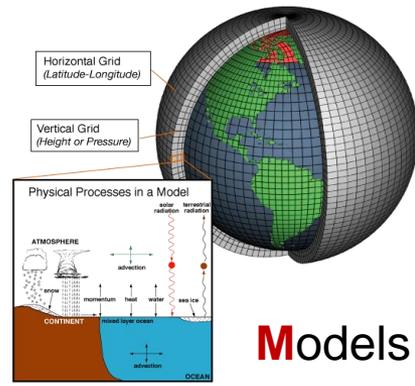
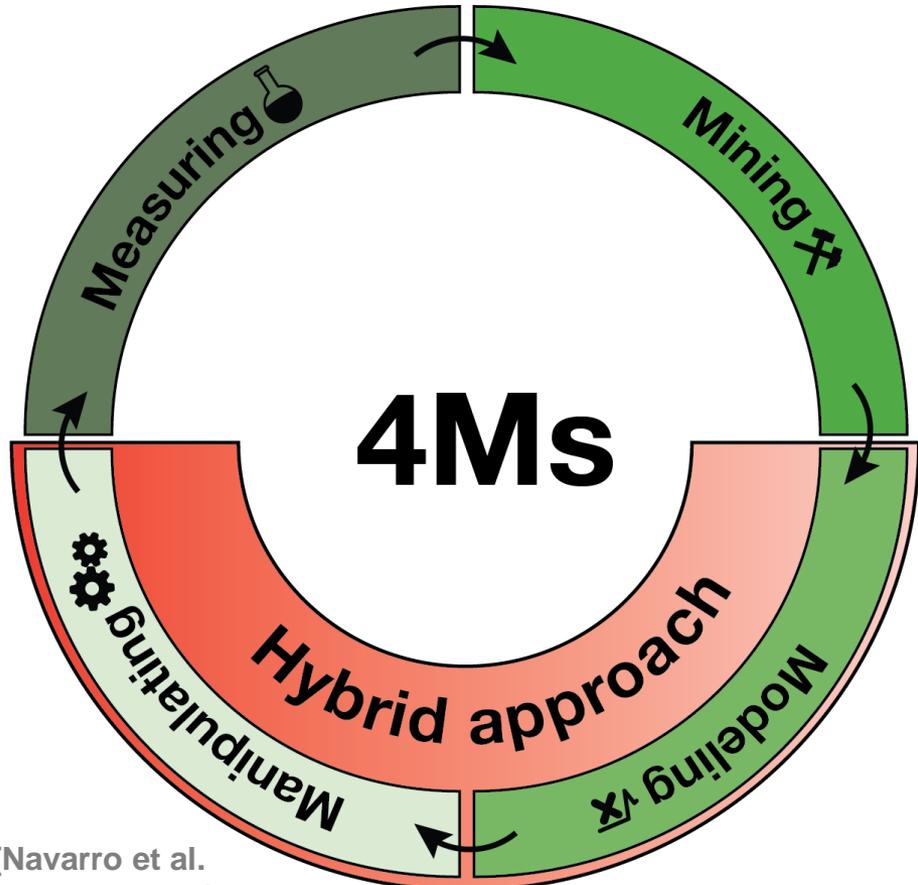


Genomics DNA sequencer



Coupling of Scientific Data to Models & Experiments

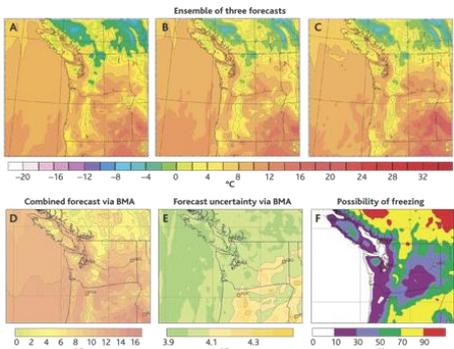
- Scientific data often coupled to a physical/biological model
- Lauffenburger's Sys. Biol. **4Ms**:
Measurement, **M**ining, **M**odeling & **M**anipulation
(Ideker et al.'06. Annals of Biomed. Eng.)
- Weather forecasting as an exemplar
 - Physical models & simulation useful but not sufficient ("butterfly"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



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



Forecasts



[Navarro et al. GenomeBiol. ('19, in press)]

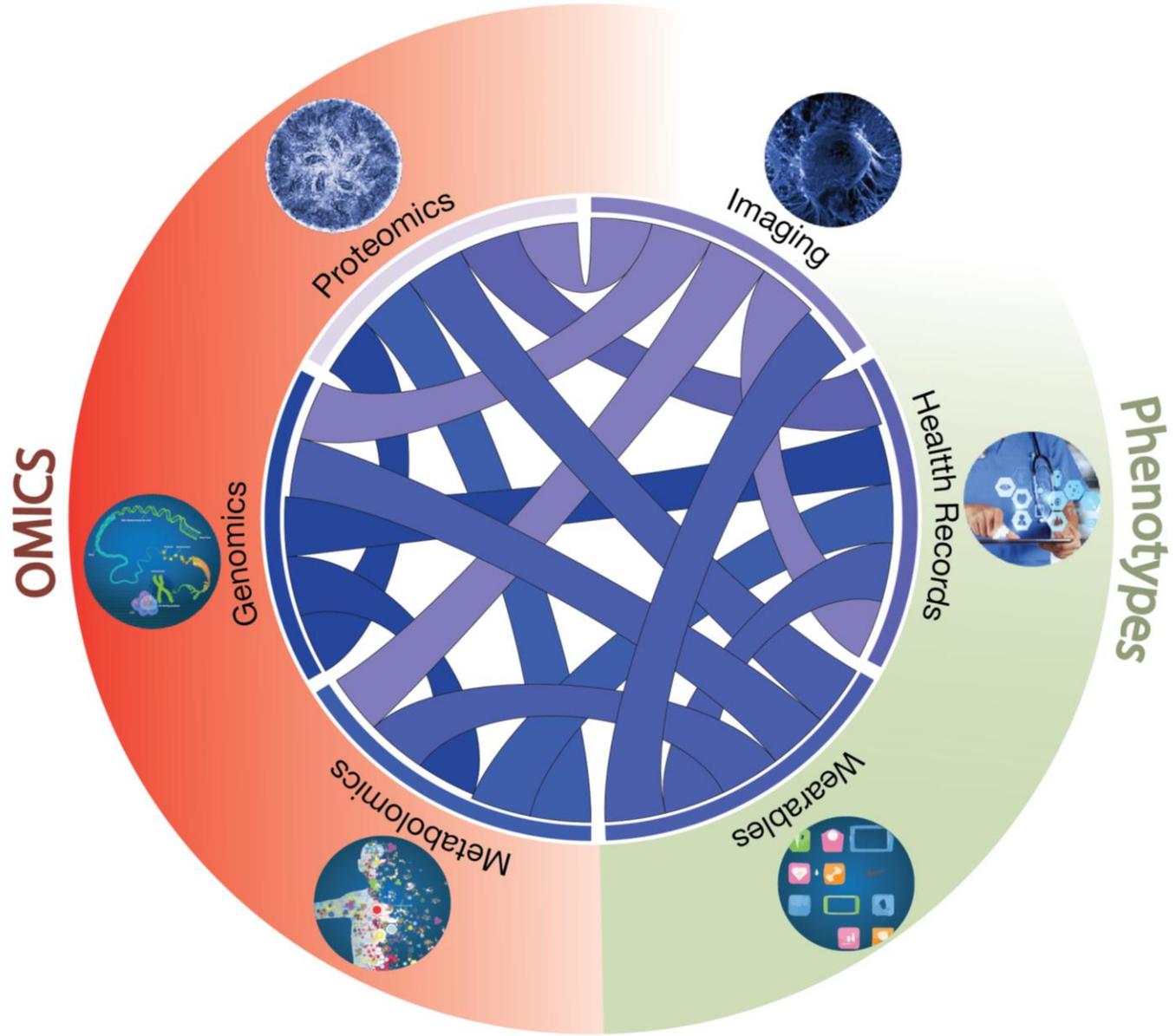
Image from <http://web.aibn.uq.edu.au/cssb/ResearchProjects.htm>

Biomed. Data science:

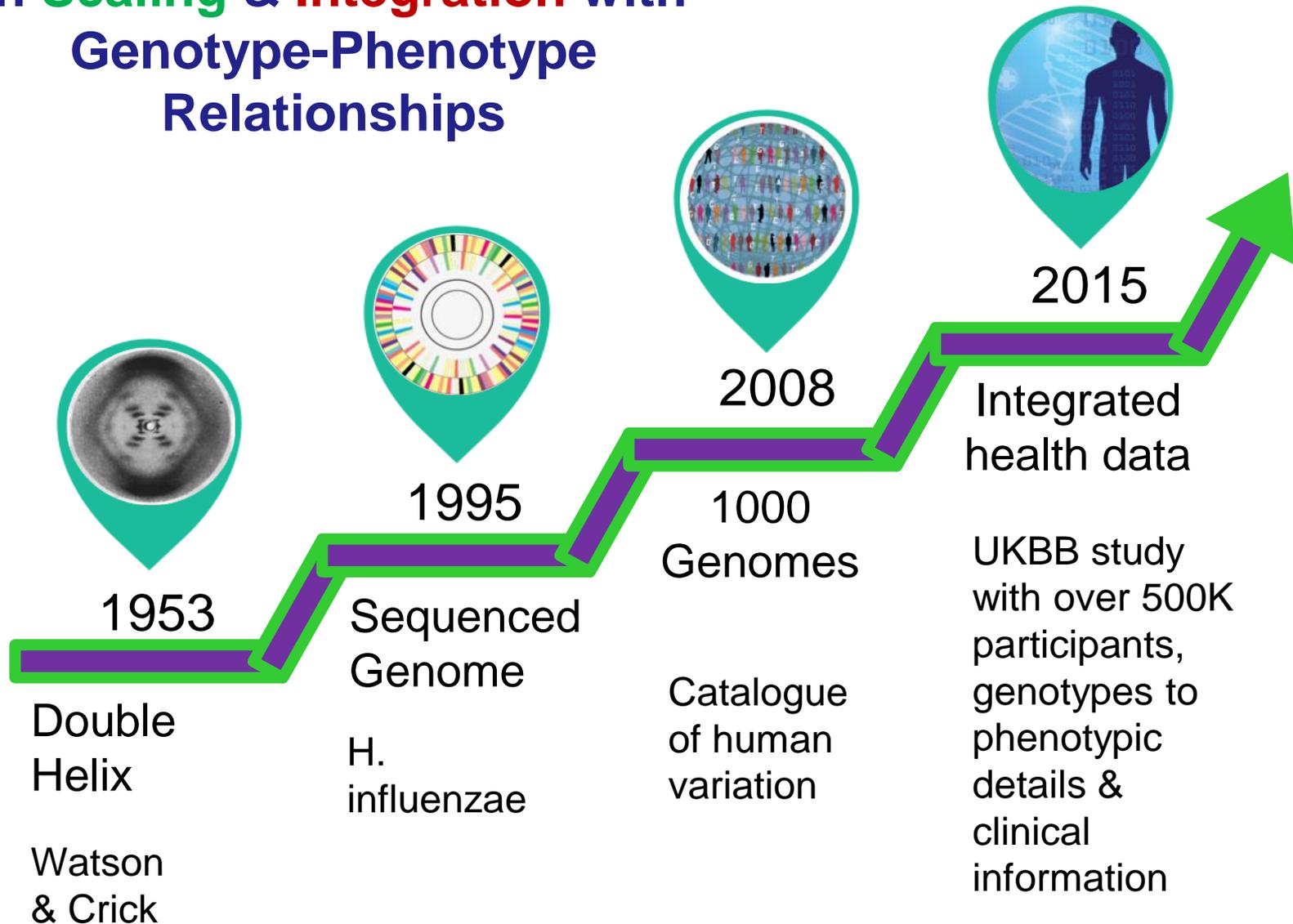
Scaling & Integration

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

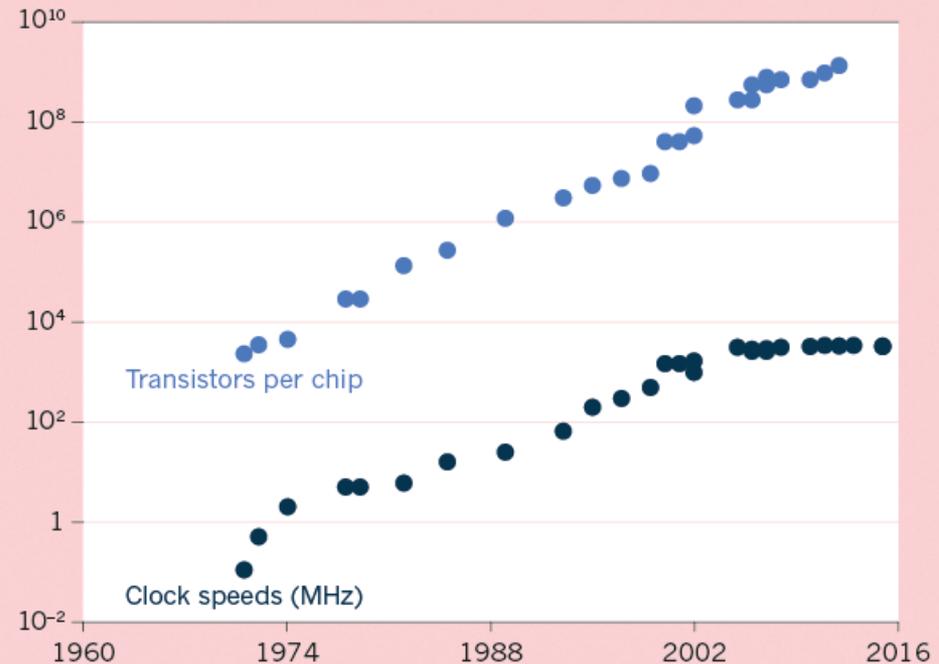
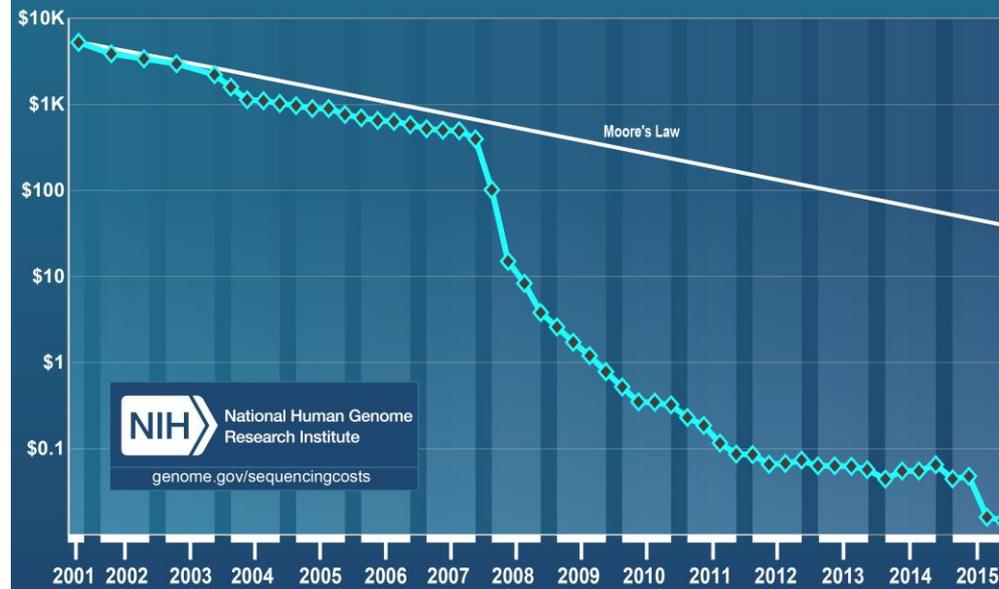


The **Scaling** of
Genomic Data
Science:

Powered by
exponential
increases in
data & computing

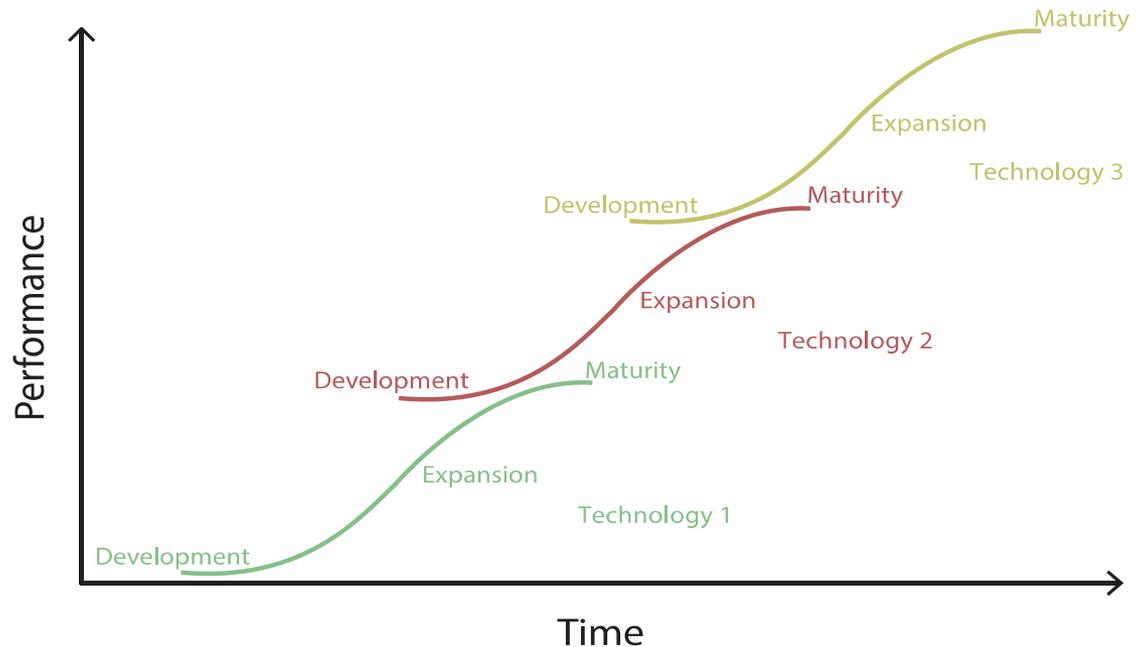
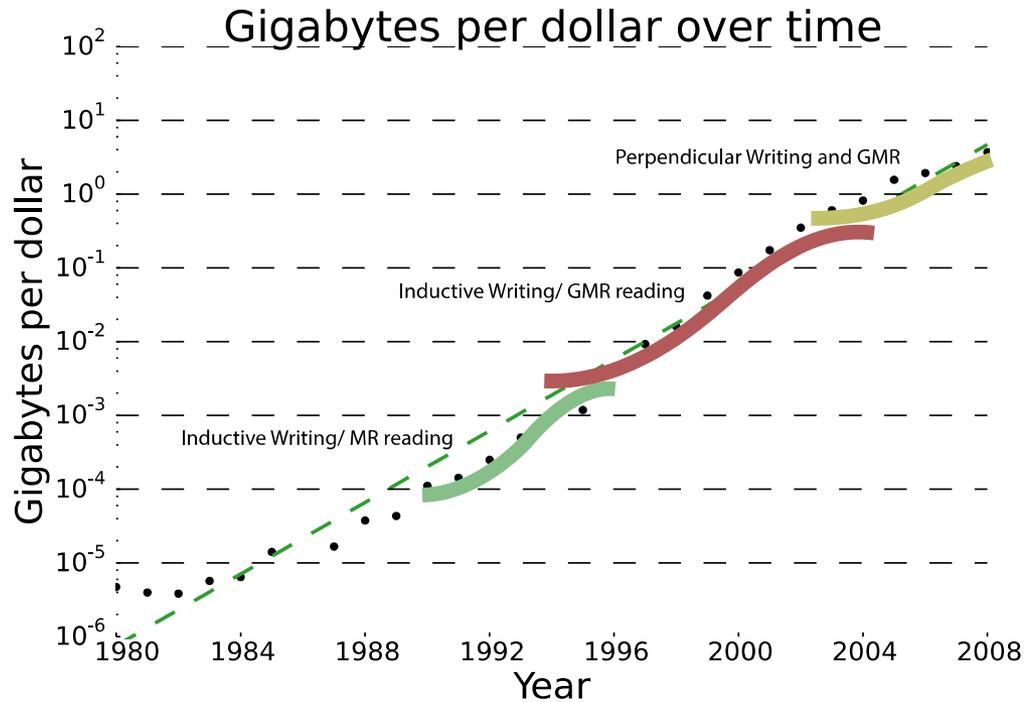
(Moore's Law)

Cost per Raw Megabase of DNA Sequence

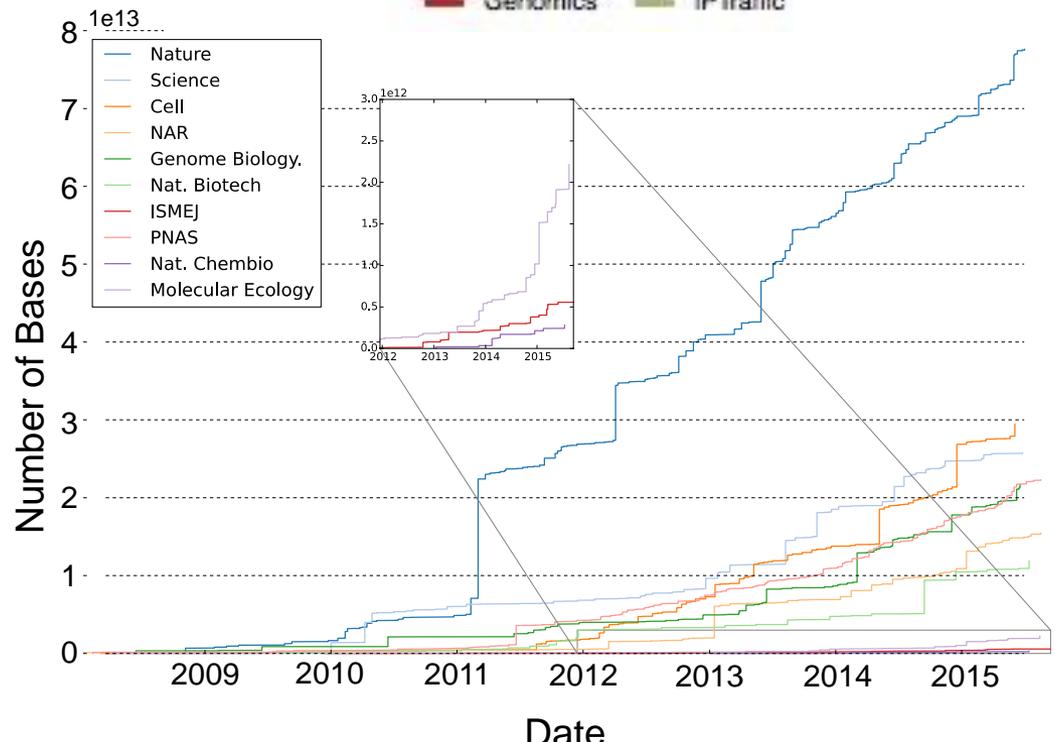
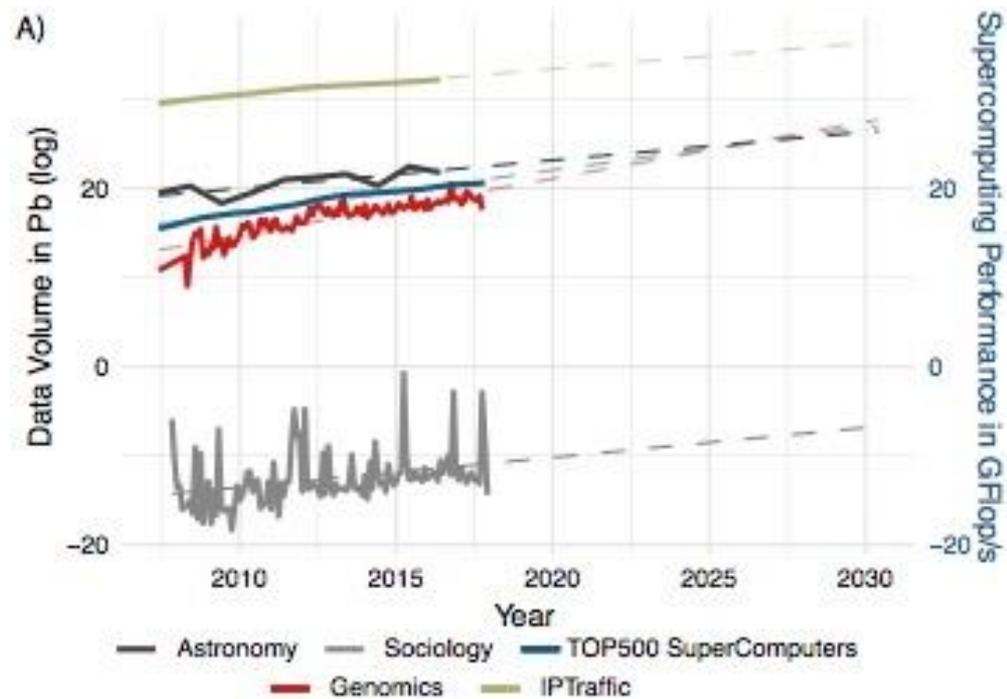


Kryder's Law and S-curves underlying exponential growth

- 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

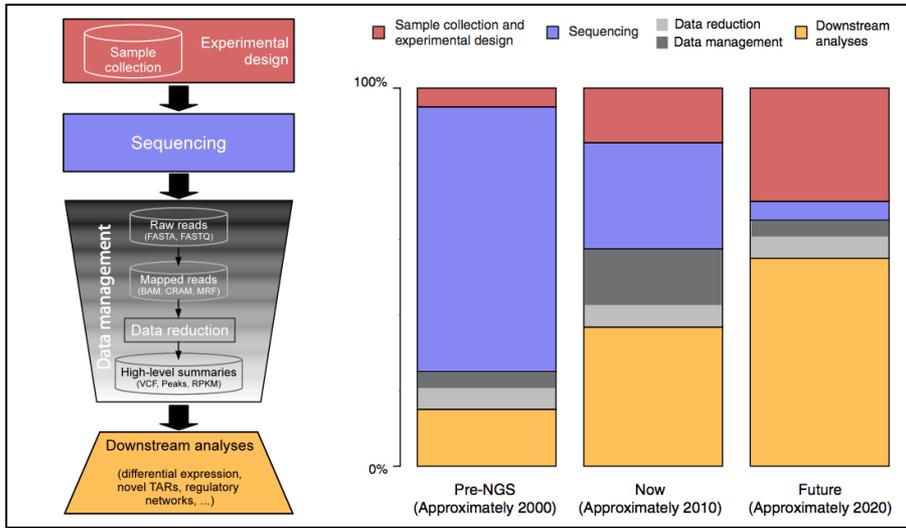


Sequencing cost reductions have resulted in an explosion of data



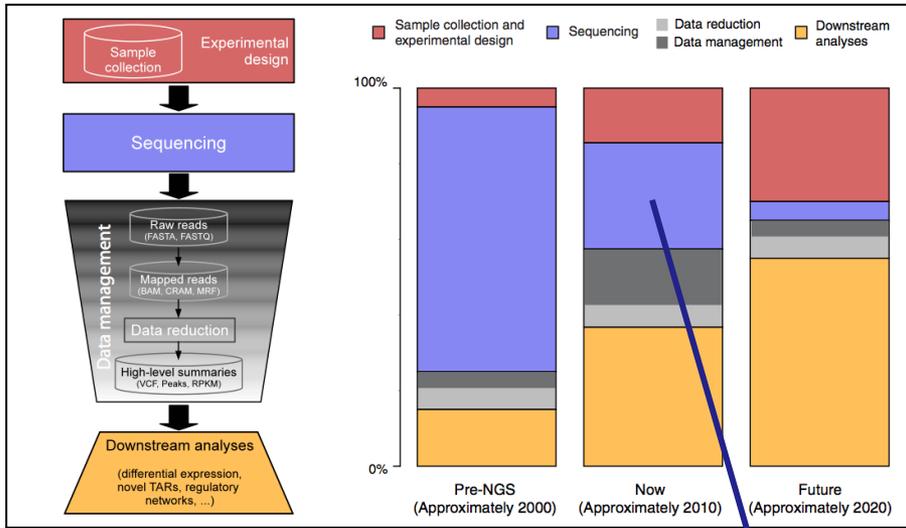
- The type of sequence data deposited has changed as well.

The changing costs of a sequencing pipeline

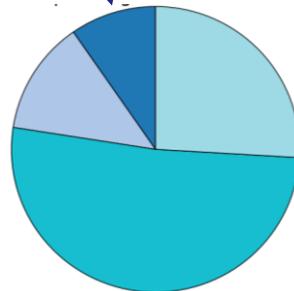
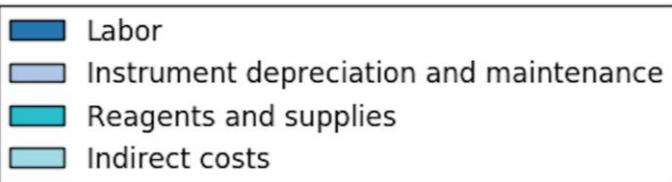


From '00 to ~' 20,
 cost of DNA sequencing expt. shifts from
 the actual seq. to sample
 collection & analysis

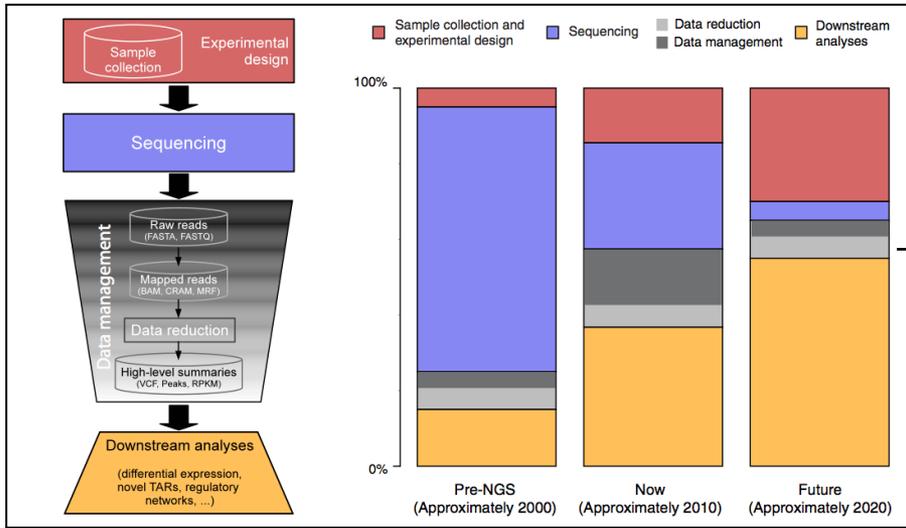
The changing costs of a sequencing pipeline



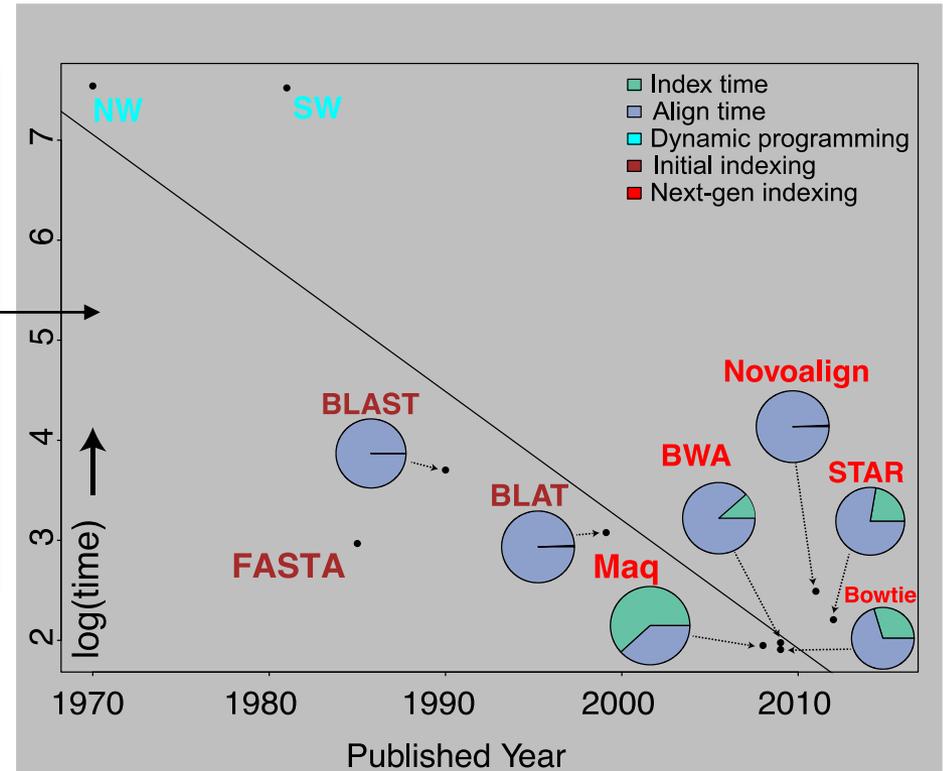
From '00 to ~' 20, cost of DNA sequencing expt. shifts from the actual seq. to sample collection & analysis



The changing costs of a sequencing pipeline

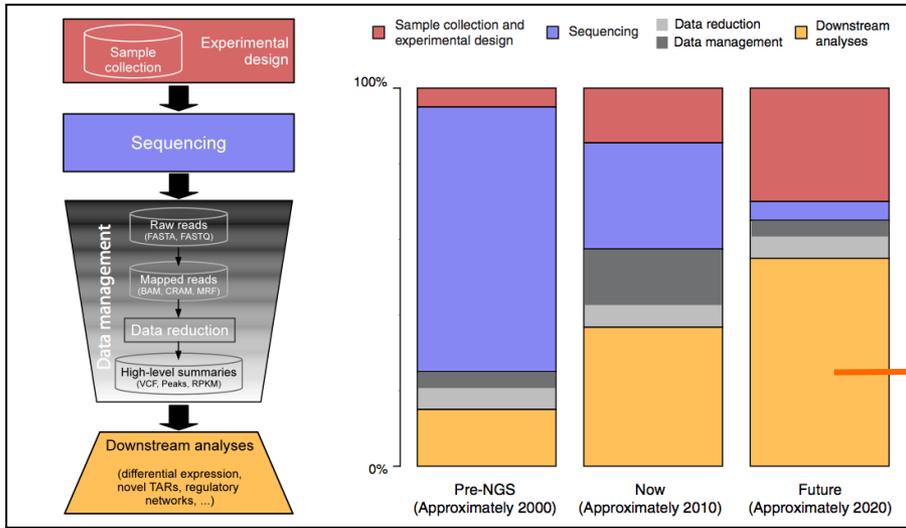


From '00 to ~' 20, cost of DNA sequencing expt. shifts from the actual seq. to sample collection & analysis

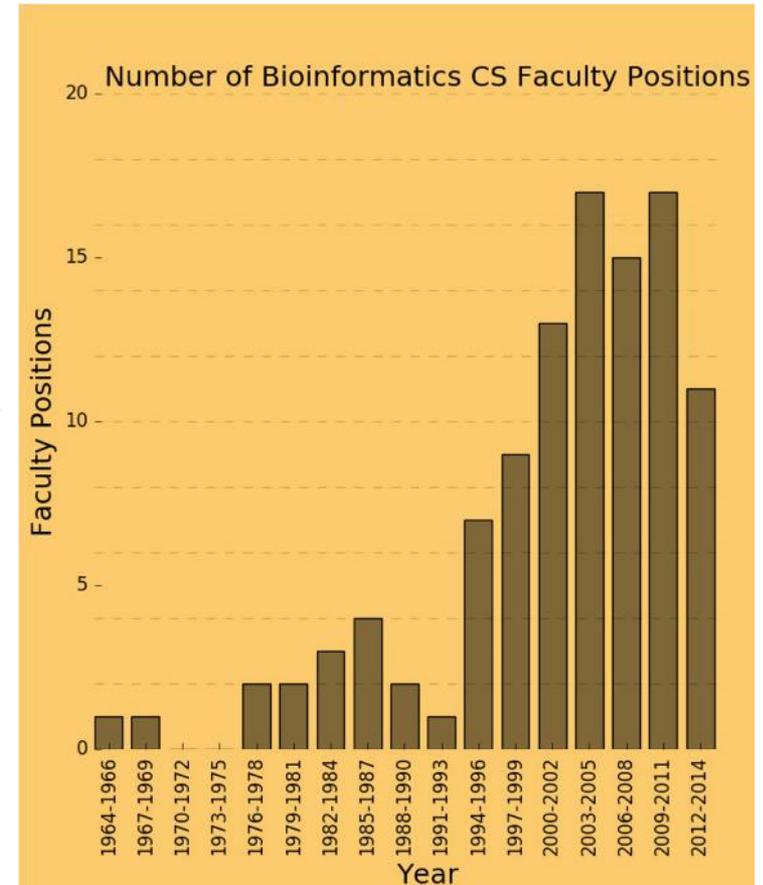


Alignment algorithms scaling to keep pace with data generation

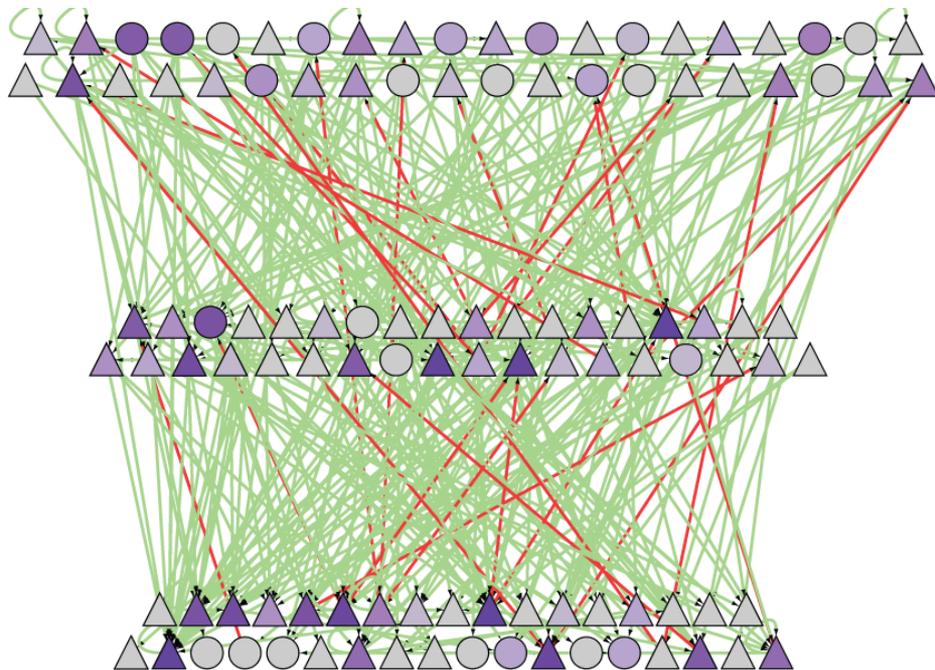
The changing costs of a sequencing pipeline



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



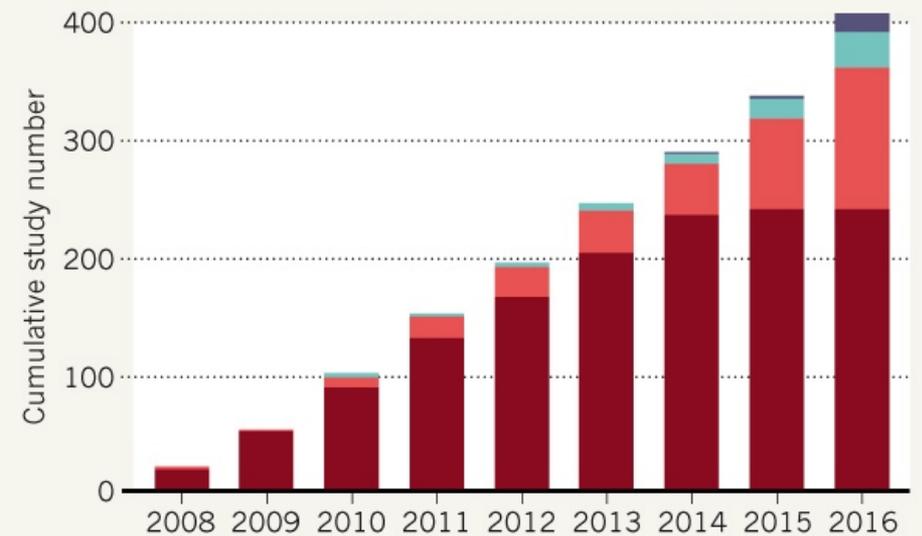
**A Success of
Scale & Integration:
Many GWAS
variants found,
most not in genes,
but affecting
regulatory network**



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)

Basic Science to Medicine

INITIATIVES



STARTUPS

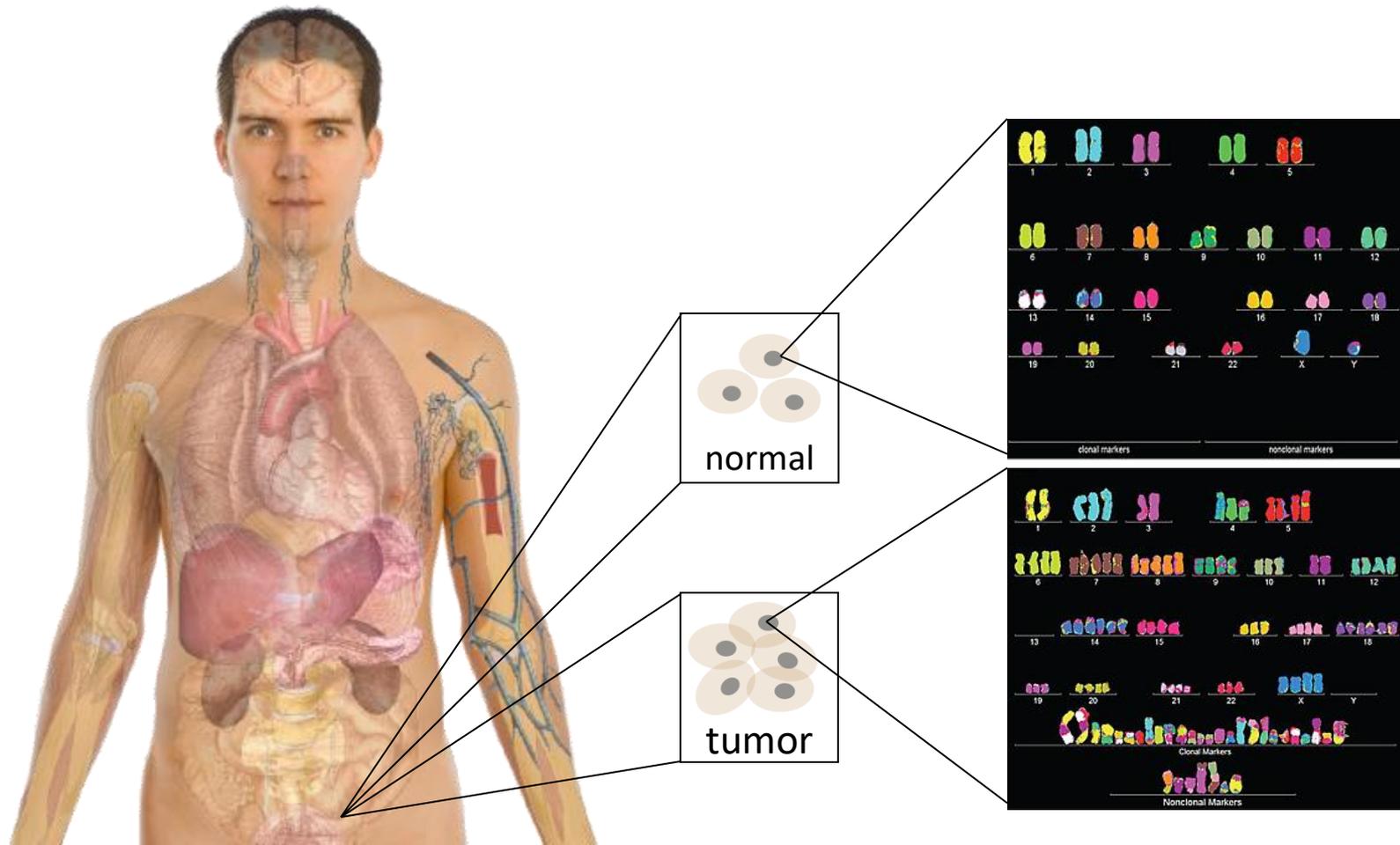
- Large-scale 'omics data as an anchor to organize phenotypic data – EMRs, wearables...
- 1st ['05-]: Exomes & chips of disease-focused cohorts – init. GWAS, TCGA, PGC
- 2nd ['15-]: Integration of full WGS with rich & diverse phenotypes - UKBiobank, TopMed, Genomics England, PCAWG, All of Us

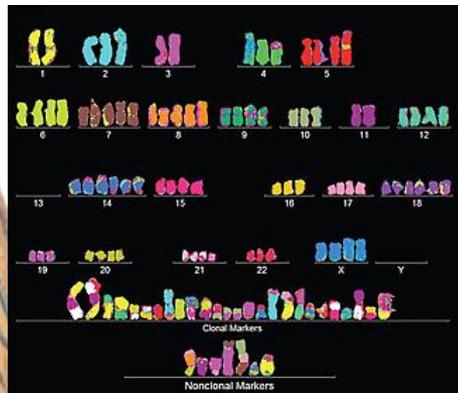
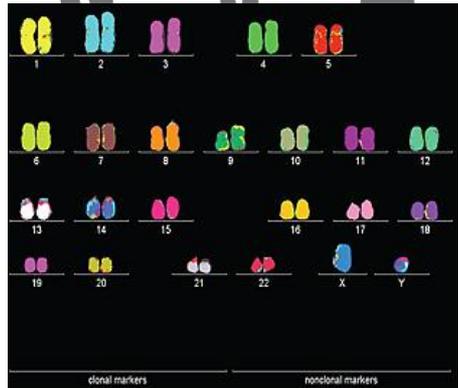
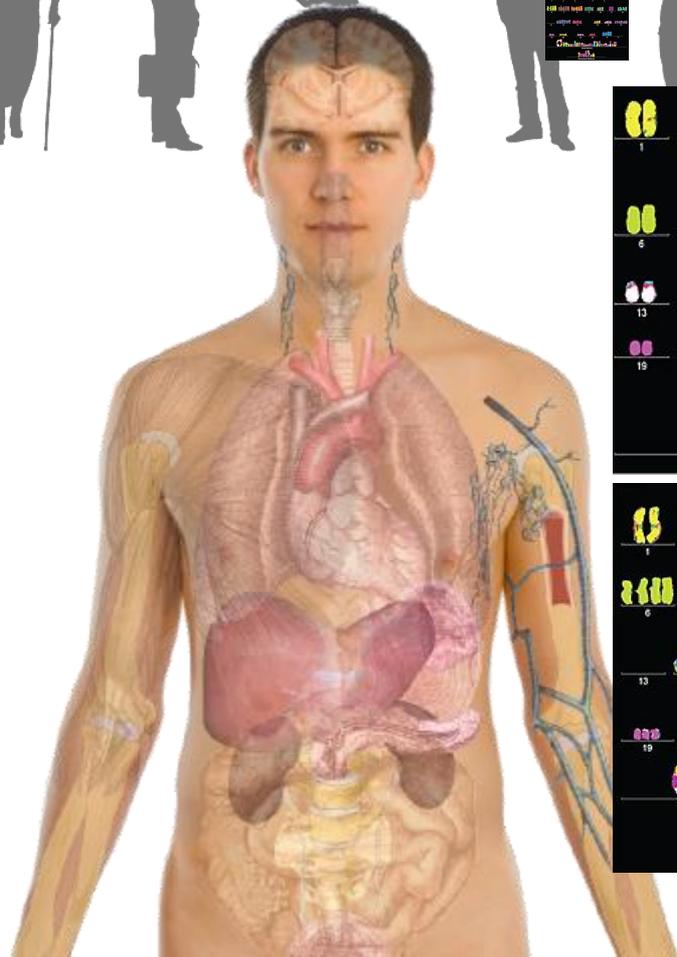
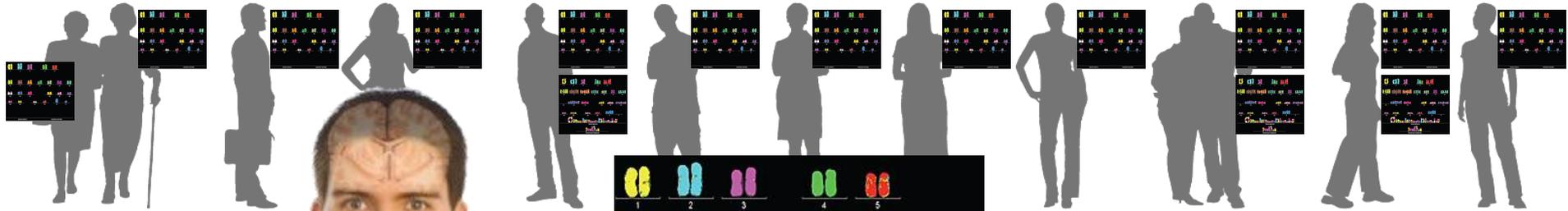
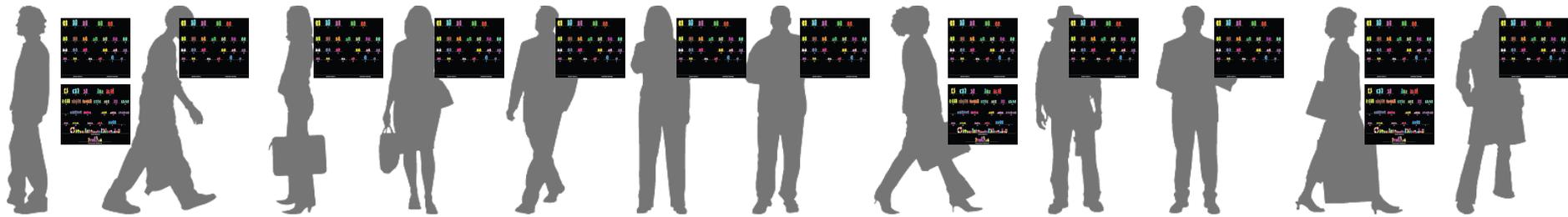
Biomed. Data science:

The Future

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

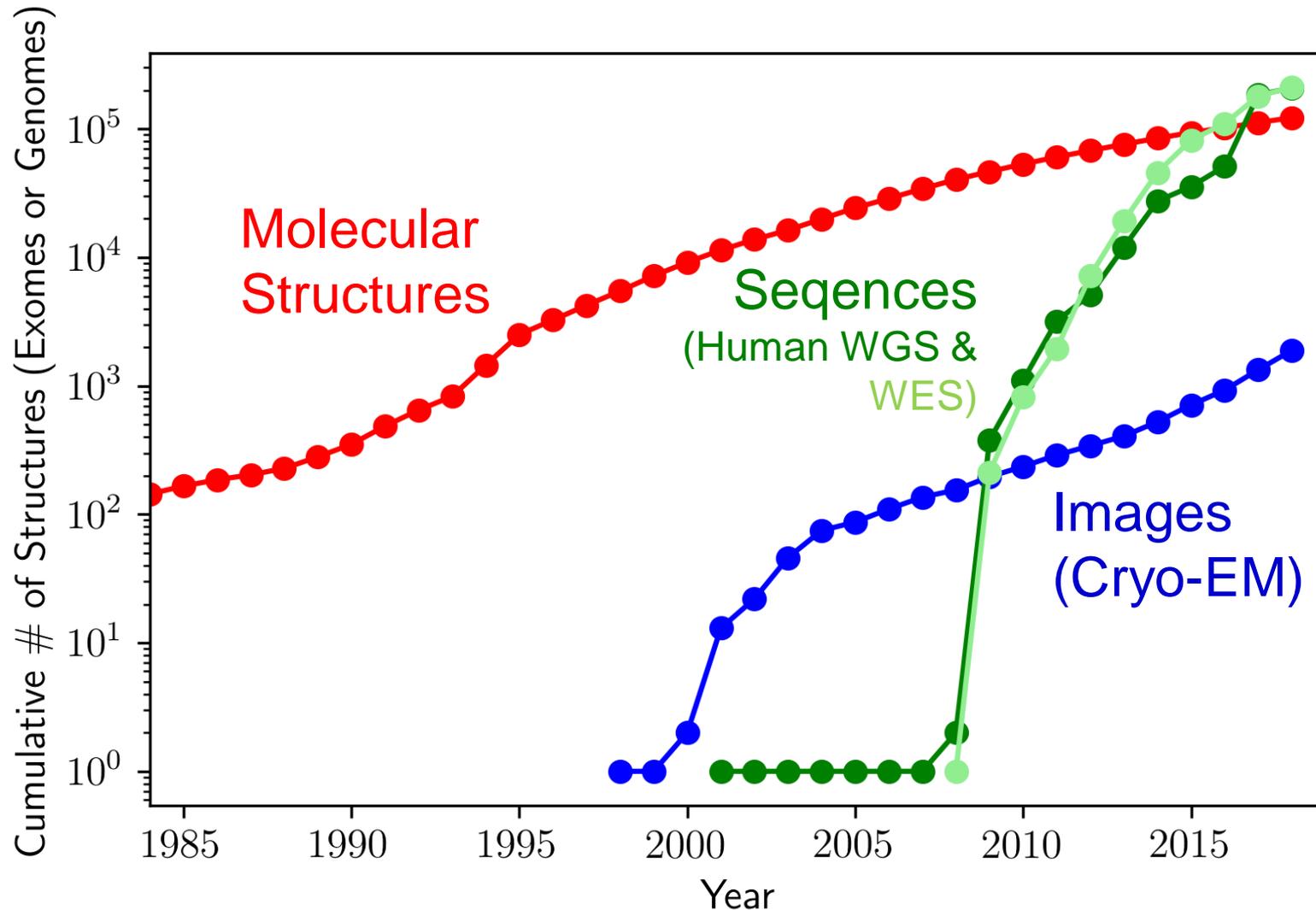
Personal genomes soon will become a commonplace part of medical research & eventually treatment (esp. for cancer). They will provide a primary connection for biological science to the general public.





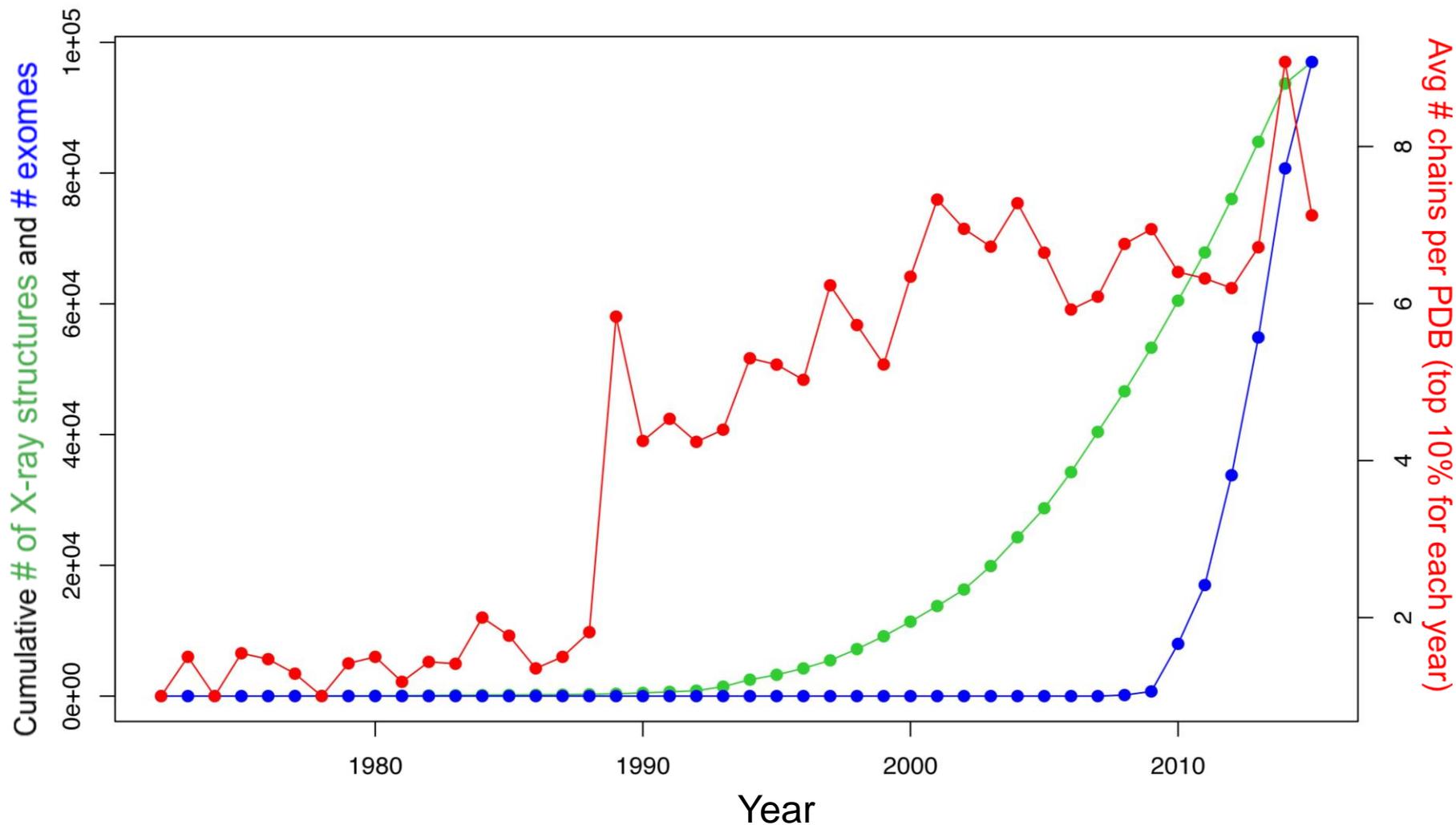
Placing the individual into the context of the population & using the population to build a interpretative model

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)

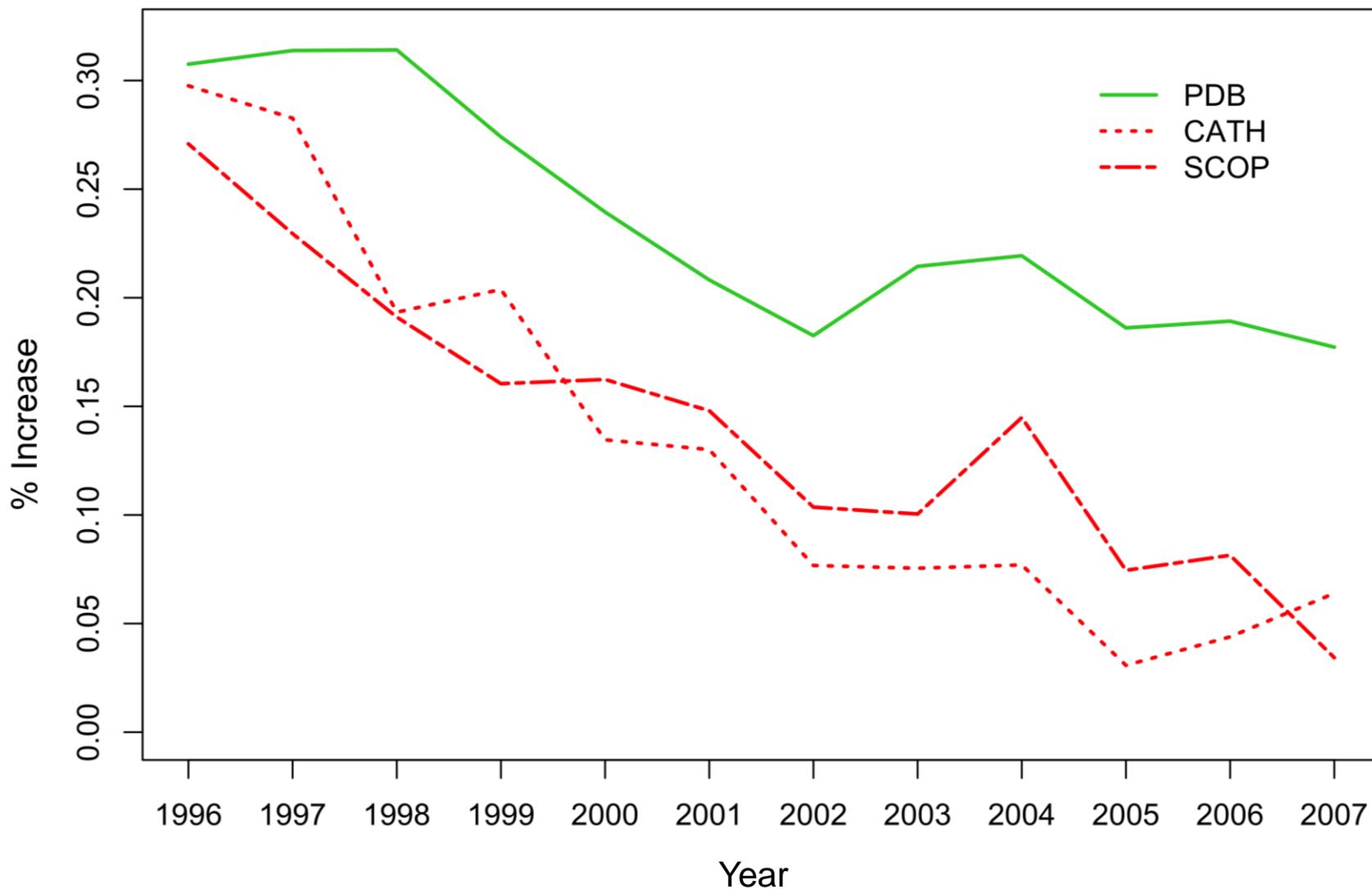
The volume of sequenced exomes is outpacing that of structures, while solved structures have become more complex in nature.



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)]

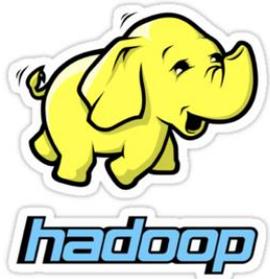
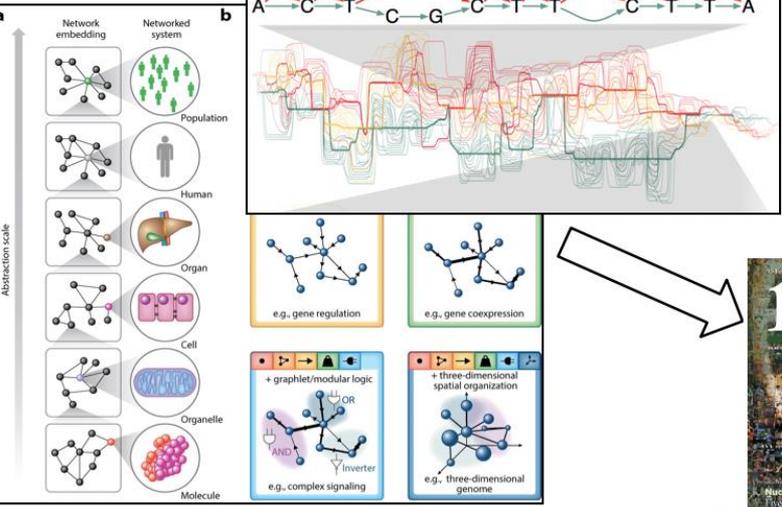
PDB: Berman HM, et al. NAR. (2000)
CATH: Sillitoe I, et al. NAR. (2015)
SCOP: Fox NK et al. NAR. (2014)

Examples of Imports & Exports to/from Genomics & Other Data

Science Application Areas

Technical Imports

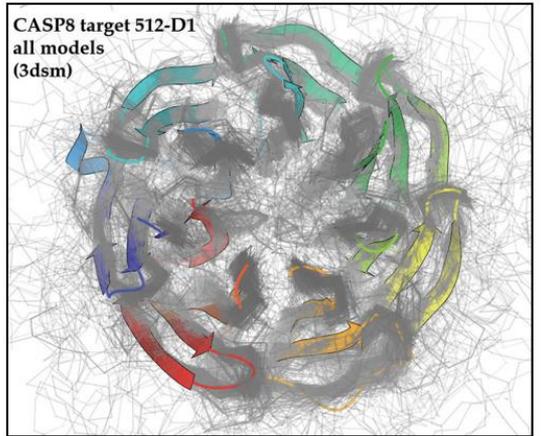
Networks and graphs



Importing tech. developed in other big data disciplines

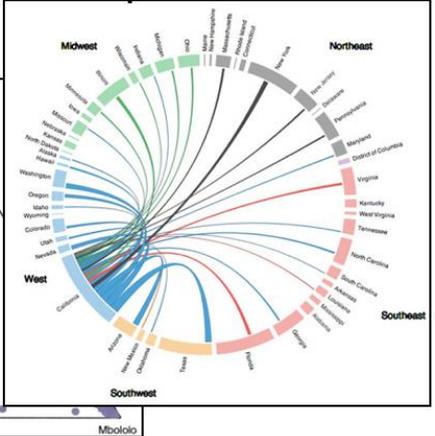
Cultural Imports

CASP

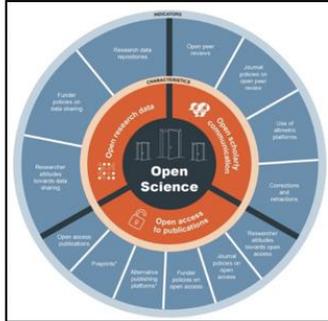


Technical Exports

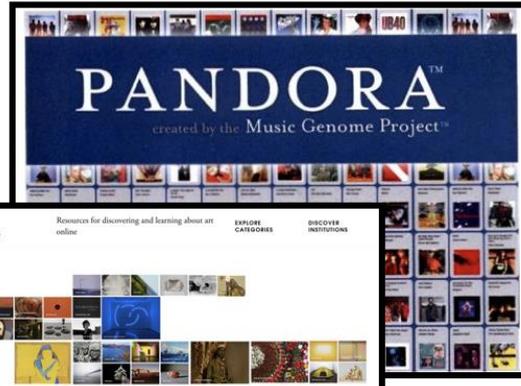
Circos plot



Open Science



Cultural Exports



[Navarro et al. GenomeBiol. ('19, in press)]

Biomed. Data science:

The Course

Thoughts on the Class

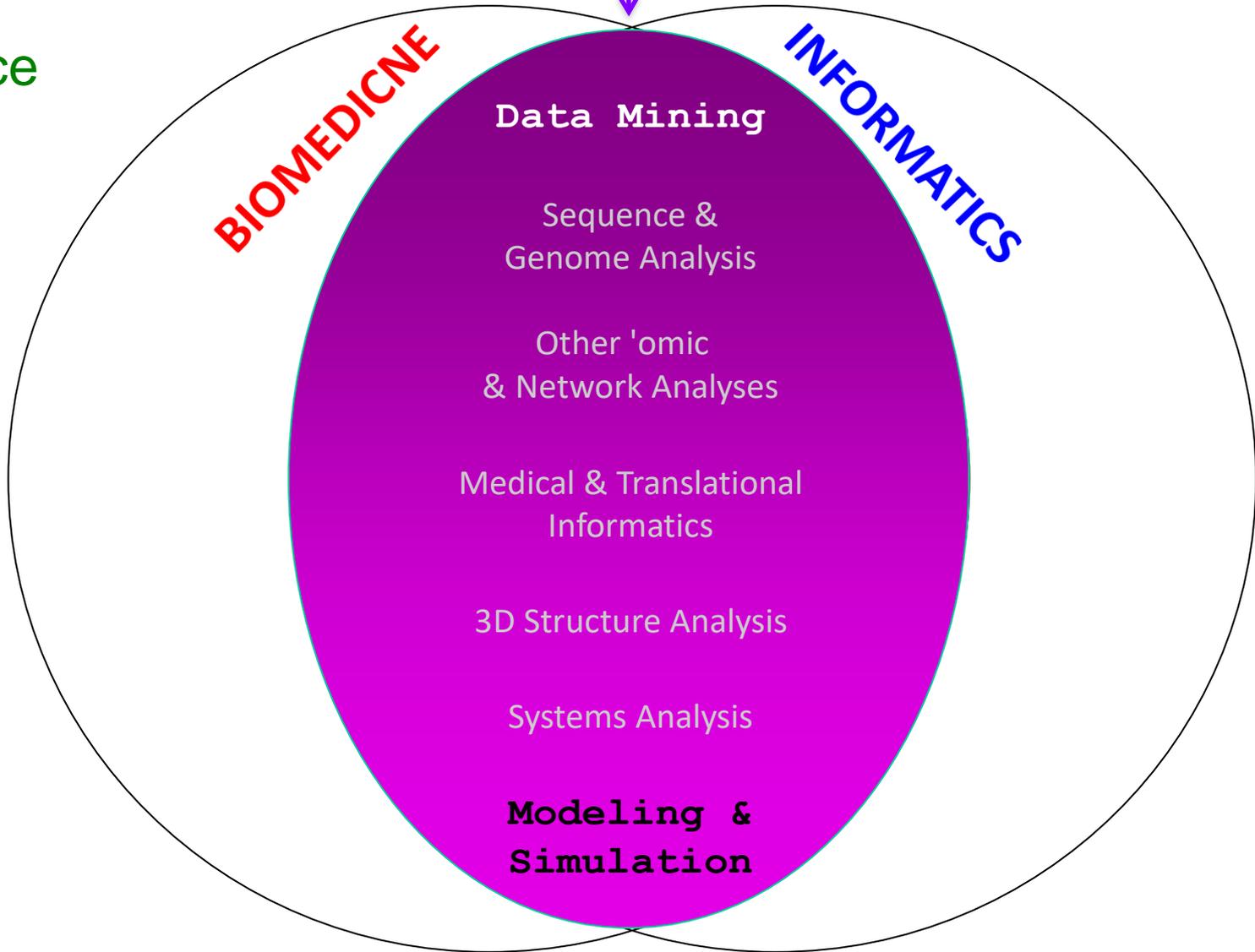
GersteinLab.org/courses/452 (Class Web Page)

- Broad overview with a few deep dives
 - Fundamentally interdisciplinary field
 - Here, focusing on molecular bioinformatics
 - Some deep dives into sequence comparison, Bayesian approaches, low-dimensional representations
 - Steering away from material in related Yale classes
- Goal is good intuition on approaches & the application area
 - Apply to related problems
- Lectures provide structure of knowledge to be assimilated
 - Varied backgrounds
 - Variety of learning approaches
 - Note references on & for the slides
- Sections for interaction & more hands-on treatment
- Quizzes & homework for individual command of basic knowledge
- Final Project for teamwork

Biomedical
Data
Science



(Molecular) BIOINFORMATICS



What is Bioinformatics?

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

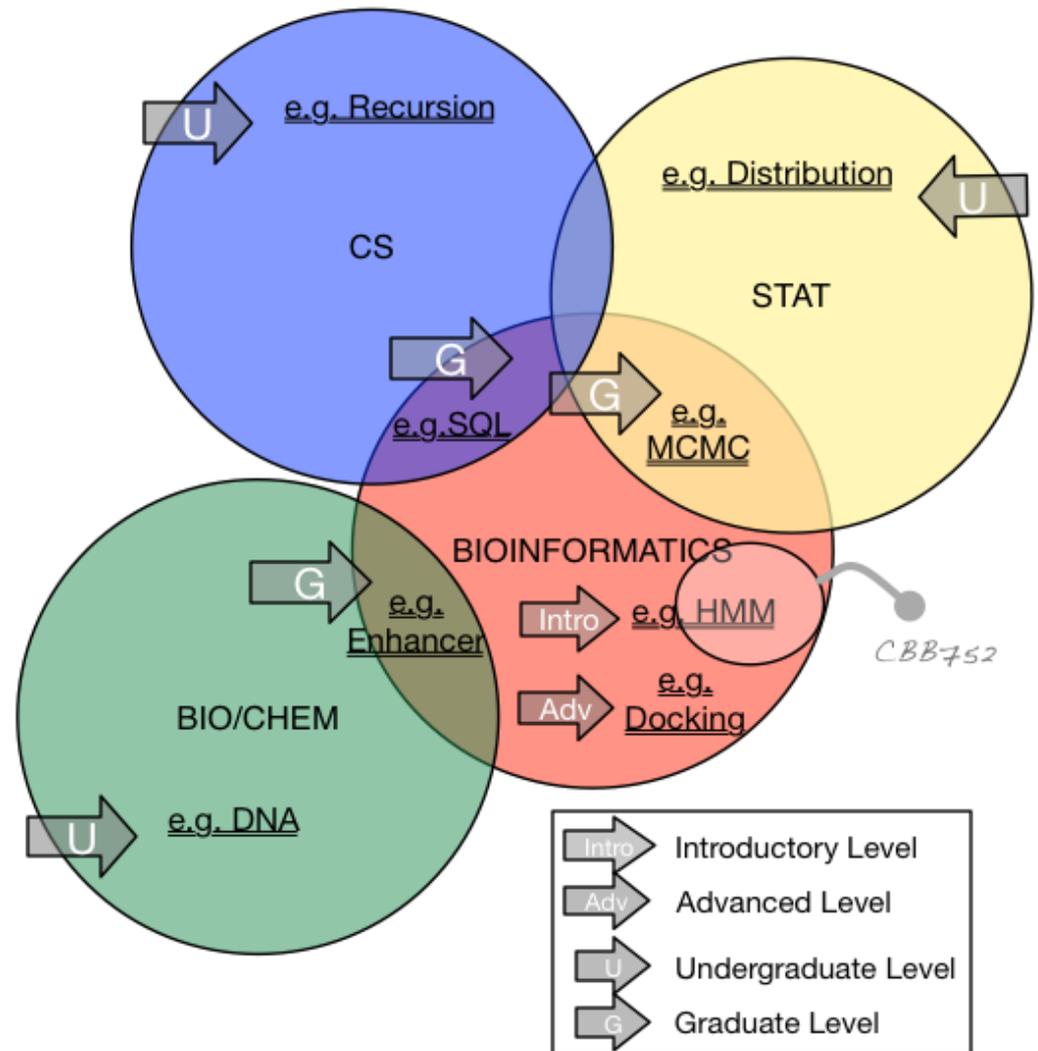
- One idea for a definition?

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.**

Defining the field – by crowd-sourced judgement

- Bioinformatics
 - Related terms
 - Biological Data Science
 - Bioinformatics & / or / vs Computational Biology
 - Bio-computing
 - Systems Biology
 - “Qbio”
- What are its boundaries
 - Determining the “Support Vectors”



Overview of Topics Surveyed

Introduction

& Overview of the Data

- Genomics & Sequencing
- Proteomics & Structure
- Databases as Tables

Data Mining & Machine Learning

- Supervised & Unsupervised Approaches
 - Decision Tree
 - Clustering & SVD
- Application to 'Omics Data
 - Calculating with sequences
 - Single cell data & epigenomics

Network Analysis

- Principles of Topology & Connectivity
- Gene Networks

Deep Learning

- Basic Theory & Applications

Physical Modeling

- Macromolecular Simulation
- Markov Models
- Molecular Packing

Additional Topics

- Privacy
- Sensor Technologies
- Personal Genome Analysis
- (Image Analysis)

Topics Covered in Other Related Courses

CBB 740: Intro to Health Informatics

- Survey Course
- Clinical Outcomes and Decision Support
- EHRs

CBB 750: Core Topics in Biomedical Informatics

- Application of Machine Learning on Clinical Data
- MIMIC and other EHR Datasets
- Text mining (i.e., word2vec) for EHR

CBB 562: Modeling Biological Systems II

- Modeling Biological Networks
- Differential Equations
- Biochemical Signaling Pathways

S&DS 665: Intermediate Machine Learning

- Theory of Machine Learning Methods
- Developments in Deep Learning Models
- Wider Range of Data Modalities and Contexts

Lectures conventions

- Convention for numbering lectures:
YYMN = (**Y**)ear, (**M**)odule, (**N**)umber
e.g. **23m3**, **22m3**, ⁽²¹⁾ **M3**
- We will mostly follow the flow in 2021 and 2022 and note the differences (See the notations at the top of each slide pack.)
- Mostly 2021 has well-produced videos, with a few from following years
- This year will have current, lecture-hall recorded videos put up quickly on canvas

Short Office Hours

Today right after class.

After that email me!

**(in Bass 432,
contact.gerstein.info)**