

Application 1: Bioinformatics & Genomics

Applications of Machine Learning in Biology & Medicine

Raunak Shrestha, PhD

Post-doctoral Research Fellow, Bioinformatics, University of British Columbia

Laboratory for Advanced Genome Analysis, Vancouver Prostate Centre, Canada

Laboratory for Bioinformatics and Computational Genomics, Indiana University, USA

<http://raunakms.github.io>

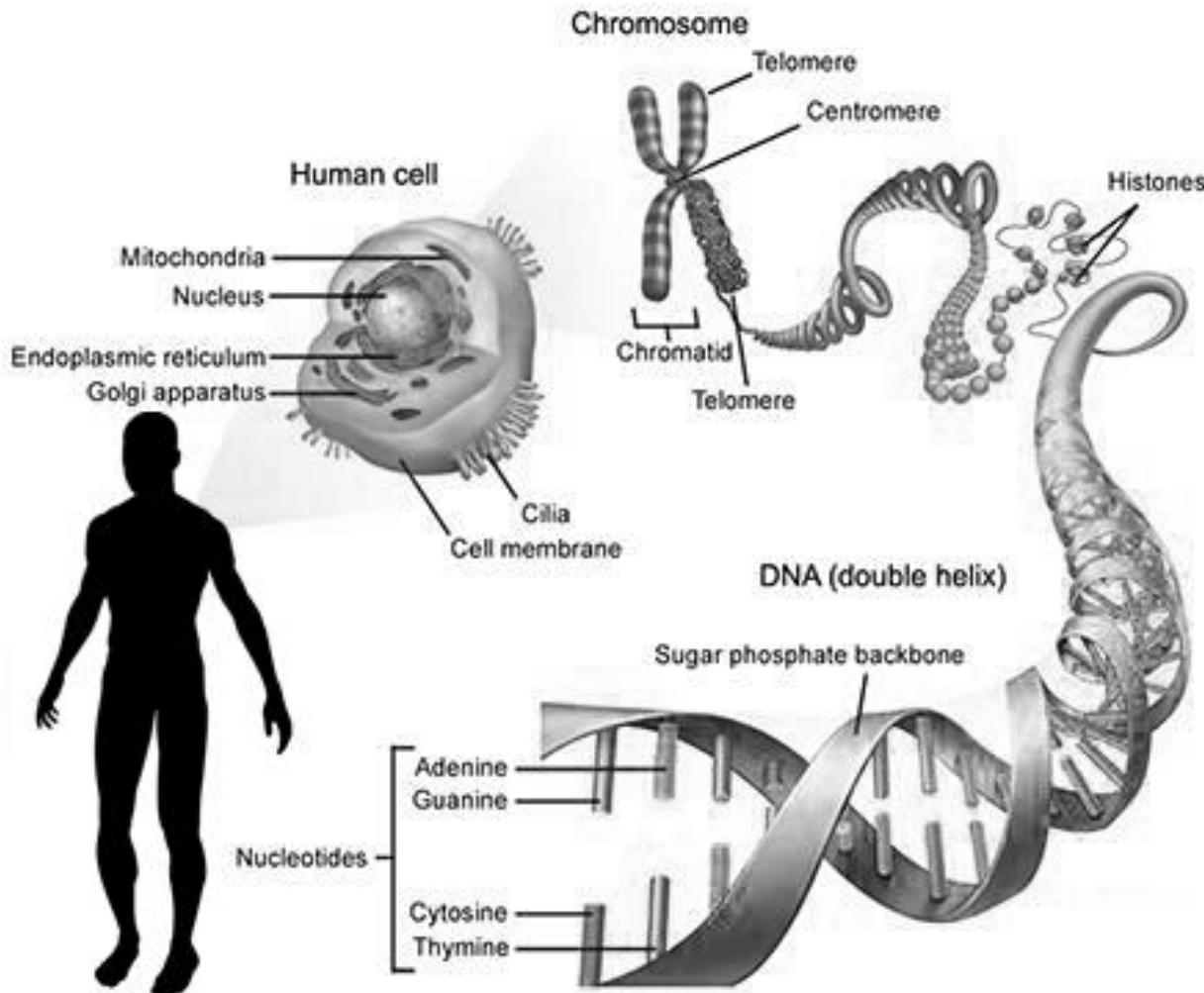


@raunakms

PART - I

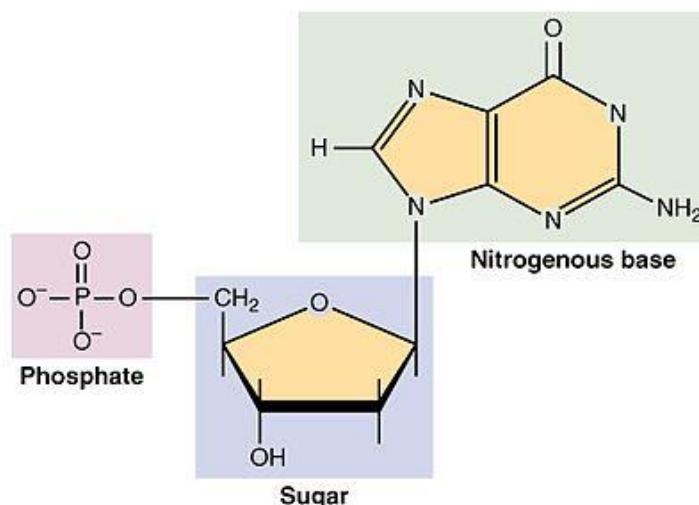
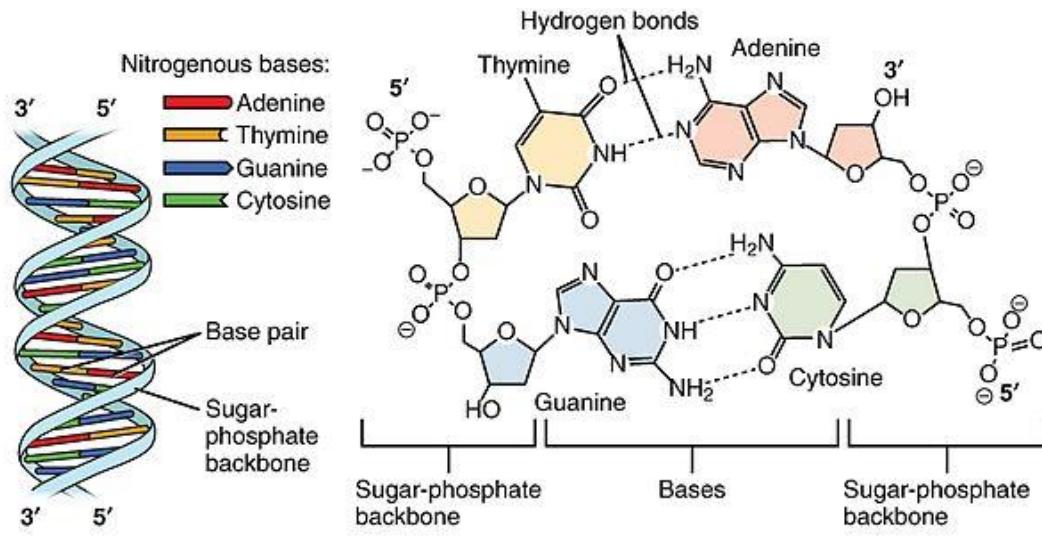
Introduction to Molecular Biology/Genomics & Bioinformatics

DNA: Deoxyribonucleic acid



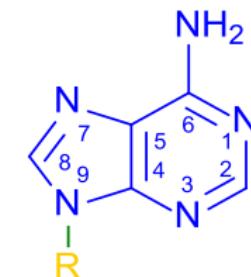
- Carrier of genetic information
- Transmits genetic information from parents to the offspring
- Often termed as “The Molecule of Life” or the “The Software of Life”

DNA: Deoxyribonucleic acid

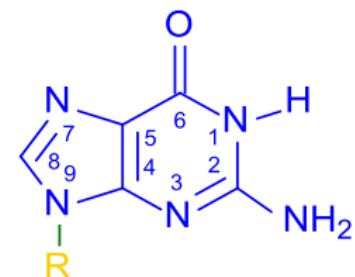


Nitrogenous base (Nucleotides)

Purines

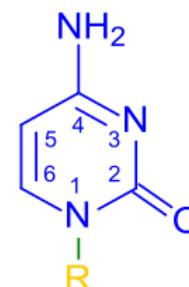


Adenine

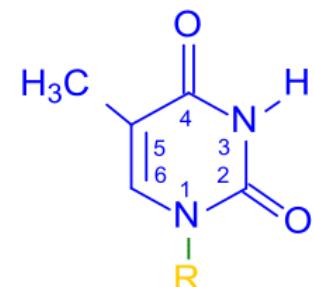


Guanine

Pyrimidines



Cytosine

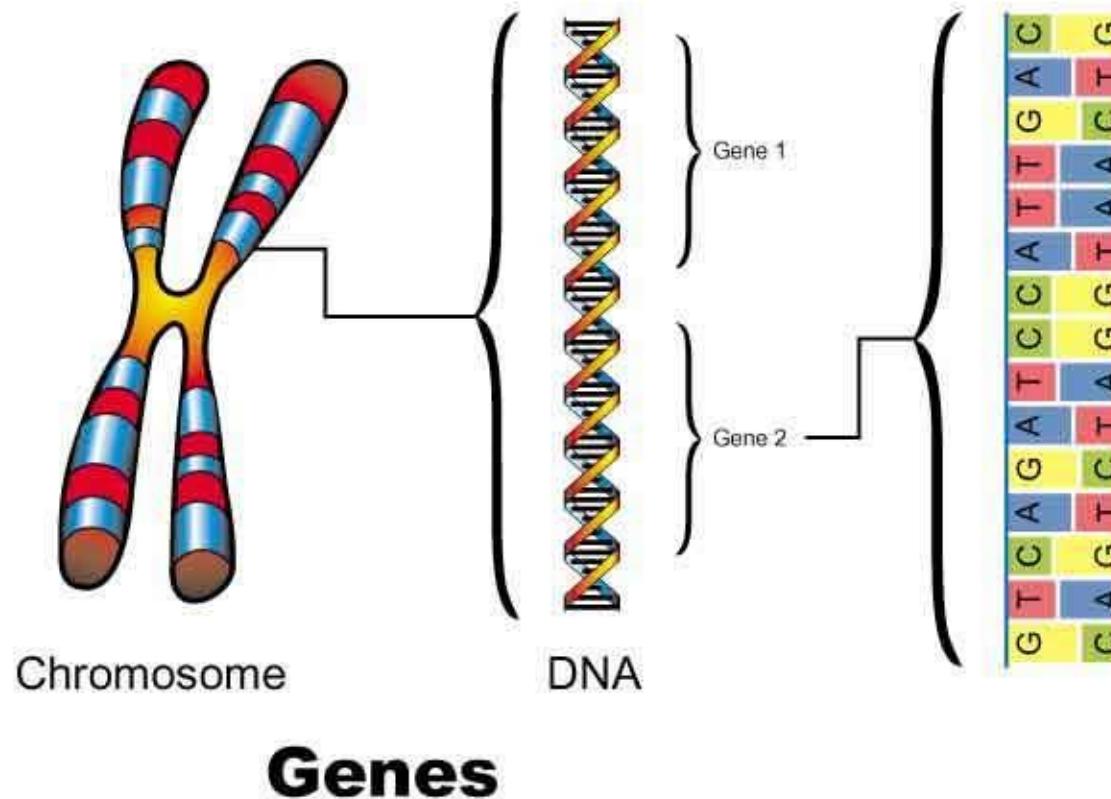


Thymine

Size of Genome (Total Genetic Material)

Species	T2 phage	Escherichia coli	Drosophila melanogaster	Homo sapiens	Paris japonica
Genome Size	170,000 bp	4.6 million bp	130 million bp	3.2 billion bp	150 billion bp
Common Name	 Virus	 Bacteria	 Fruit fly	 Human	 Canopy Plant

Gene: The Functional Unit of the Genome



Each **Gene codes for a particular Protein**
that carries out a particular task/function

Nucleotide sequence of BReast Cancer Associated (BRCA1) Gene

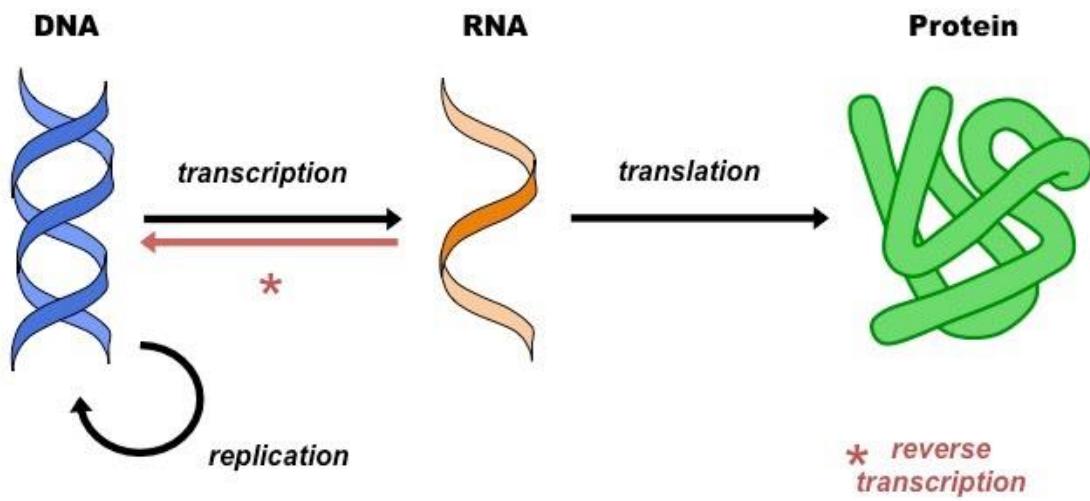
```
1 cagcaaggct cacccctct cagcaccta gtctttctt ctataaagtg aagcggtgat
61 ttgcctctta aggtcccttc tcatttgtgtg aqaacaagag tcccagtggc ctggggacac
121 agggggggg aggccgctgc ccaccctca acactattac ccacatggc ctgtgtctct
181 gtctcctctc tcccttgta aatggatctg ctgcgccagc aaacaactgt ccattgtgtc
241 ctgatgaggg acgagcaggg aaaaccaatc acttatgtatg caaatgagcc gacggctgca
301 gcctcctcac attcaacttag tgcttctt ctgcccagag cccggcgccg ccccaagttgg
361 gccttctcc ctagccccag ttcccatctg gtcctggagg aggagttggg aggcccagac
421 tggccctcccc attgtcggttcc attcaactcc ctgcaaggag tttcaaccct caaatcctca
481 gagcctggcc cagccctct ccccacccaa aaggtggctc agttttat tttttaaaaaa
541 tcgatgtaaa ttccacataac atacaatttgc ccattttaaa atgtacaattt cagggggattt
601 tagtgcatttca acaatgtttgtt gcaaccacca cctctggactt caggttttatt gccctctgtcc
661 cattctttctt tttttttttt ttttttgaga cggagtttttgc ctctgtcgcc caggctggag
721 tgcagttggca gatatggc tcactgcaag ctccgcctcc cgggttcatg ccattctctt
781 gcctcagcc accaagttagc tgggactaca ggcacctgccc accacgcctg gctaattttt
841 tgtatttta ttagagatgg ggtttccaccg tgtagccatg gatggctcg atctcttgac
901 ctcgtatctt gcccgcctcg gcctccaaa gtgtctggat tacaggcttgc agccaccacg
961 cccggcctgc cctctggccca ttctatggcc cacagcccttca agtggactg ggaagctgct
1021 gaggcctcag cagagcttag gactaatggaa ggctgtatggaa taggctgaga aaagcccaaga
1081 gctggcctga ggtgaagagg tgatccccac acttttagat gatgttagaa ggtttgggct
1141 ttgggtggag ggggacatca ctgccttagc cagaacggac ctgtggccac gtgttagaaac
1201 tttttttttt ttttttgaga tggagtttttgc ctttgttgc ctgggttggaa gtgcaatggc
1261 atgatctcag ctgcgtgcaaa cctctggctc cagggttcaa gtgattctctt ggcctcagcc
1321 tcctgatgttgc ctgggattac aggtgcacac caccacgccc agtcaattttt ctatttttaa
1381 tagagacggg gtttccaccat gttgaccagg ctggtcttgc acttctgtatc tttaggtgtatc
1441 ggcctccaa aagtgtggg attacaggca tgagccaccg tgcctggccc cagaacatcc
1501 ttgcattctt gaatttccat gggaaatttttt ttttccattt taatctctca ggcatttagaa
1561 accagaagtg gtcctgtac agatcttggt ttccatagct ttcagagctt ctgtgggctt
1621 ggggatgtatc ccgtcgccag gcttctcaactt ttagcagatgt cagaaggccc tgagacagcc
1681 ccacagtccctt ttcgtccat gaaagacctc ccggccaccag gctcttagaa caatccttaga
1741
```

- DNA/Gene is a long stretch of nucleotide sequence
- You can think of this as a long string of text
- BRCA1 gene necessary contains information (or instructions) to code for BRCA1 protein

To view entire sequence of BRCA1 gene:

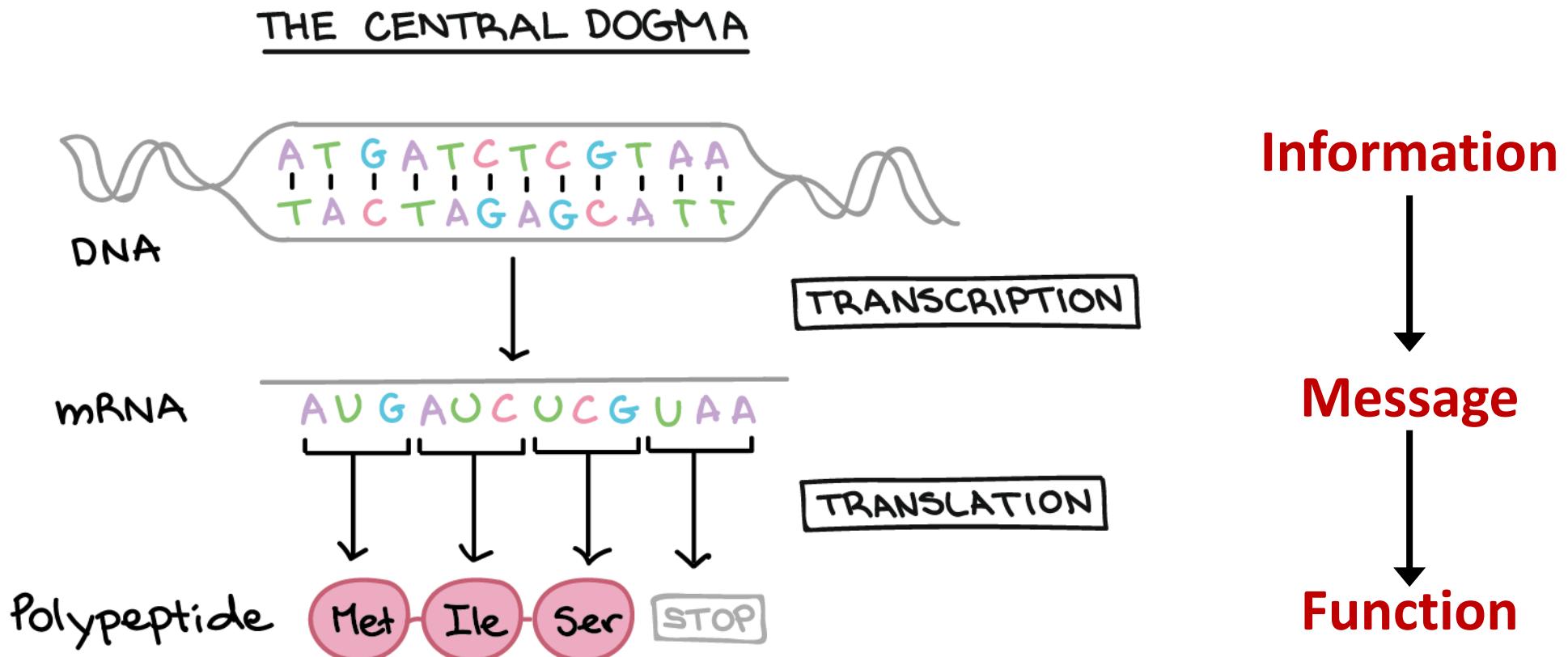
http://asia.ensembl.org/Homo_sapiens/Gene/Sequence?g=ENSG00000012048;r=17:43044295-43170245

The central dogma of molecular biology

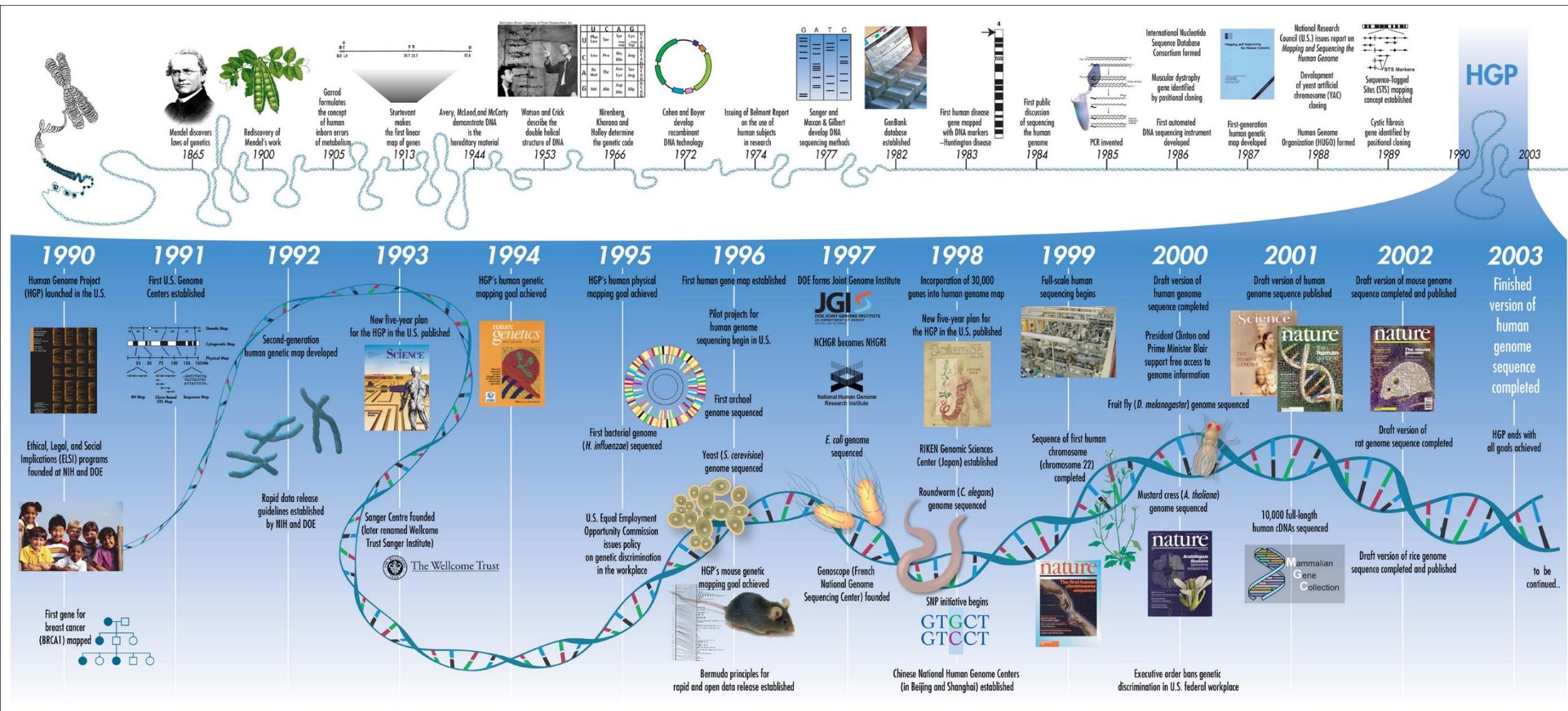


- explains the flow of genetic information within a cell
 - **DNA codes for RNA** via the process of transcription.
 - **RNA codes for protein** via the process of translation.
- This information flow was always considered to be uni-directional:
 $\text{DNA} \rightarrow \text{RNA} \rightarrow \text{Protein}$

The central dogma of molecular biology



Human Genome Project





The announcement of the completion of
the first draft of the human genome on
26 June 2000

Documentary on The Discovery of the DNA structure



Best Documentary 2017 PBS Nova Documentary Collection: DNA - Episode 1 of 5 The Secret of Life

https://youtu.be/yf4_D7HF5kU

Documentary on The Human Genome Project



DNA Documentary The Human Race 3 of 5 episode

<https://youtu.be/GyBued188Lk>



BioPerl

main links

- Main Page
- Getting Started
- Downloads
- Installation
- Recent changes
- Random page

documentation

- Quick Start
- FAQ
- HOWTOs
- BioPerl Tutorial
- Tutorials
- Deobfuscator
- Browse Modules

community

- News
- Mailing lists
- Supporting BioPerl
- About this site

development

- Developer Information

article discussion edit history

How Perl saved human genome

How Perl Saved the Human Genome Project

by Lincoln Stein

bioperl project



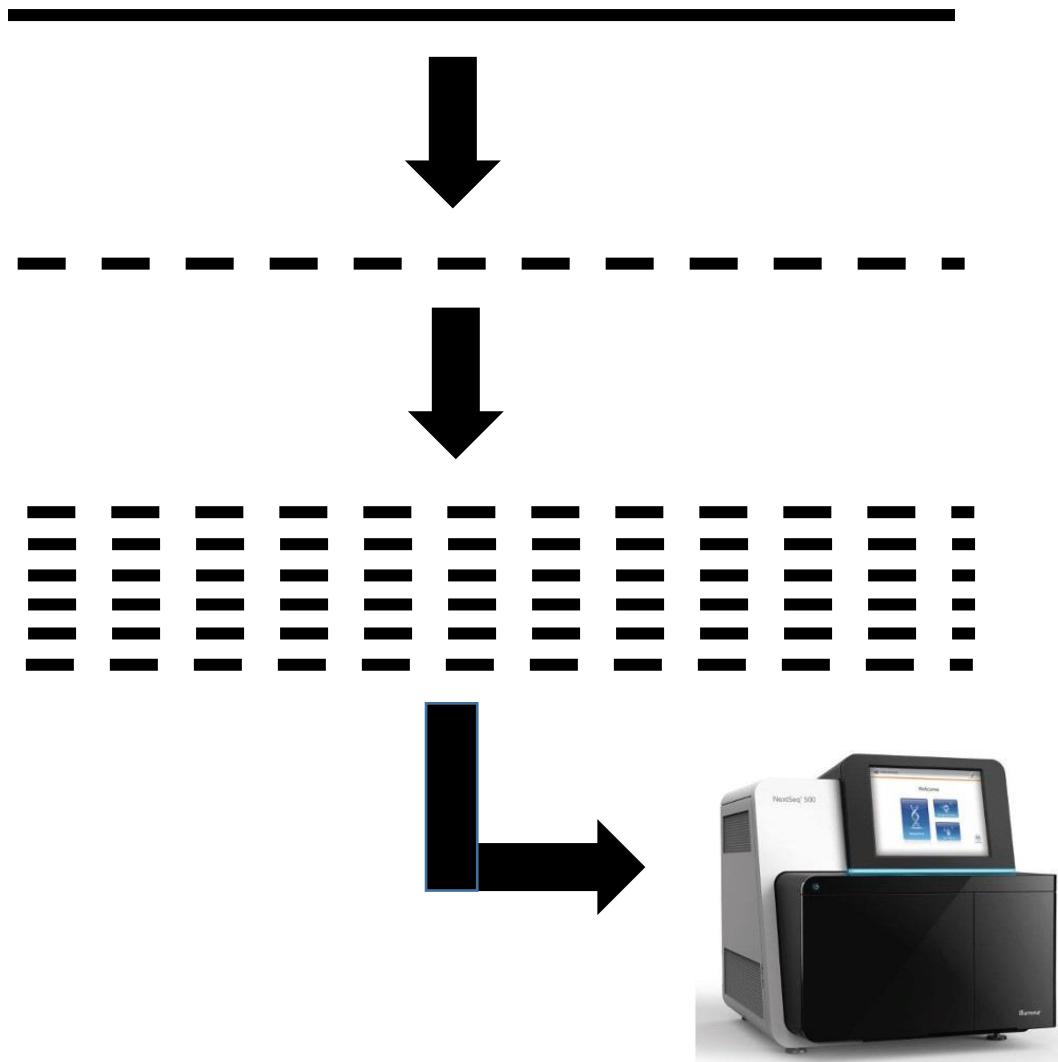
The helix graphic is reproduced from Dr. Lincoln Stein's article "How Perl Saved the Human Genome Project" as published in the September 1996 issue of [The Perl Journal](#).

Reprinted courtesy of the Perl Journal, <http://www.tpj.com> . Lincoln Stein's website is <http://stein.cshl.org> .

Recommended Reading

https://web.archive.org/web/20070202101624/http://www.bioperl.org/wiki/How_Perl_saved_human_genome

DNA sequencing (Simplified version)



DNA molecule (~3 billion bp)

Fragment the entire DNA into thousands of short fragments

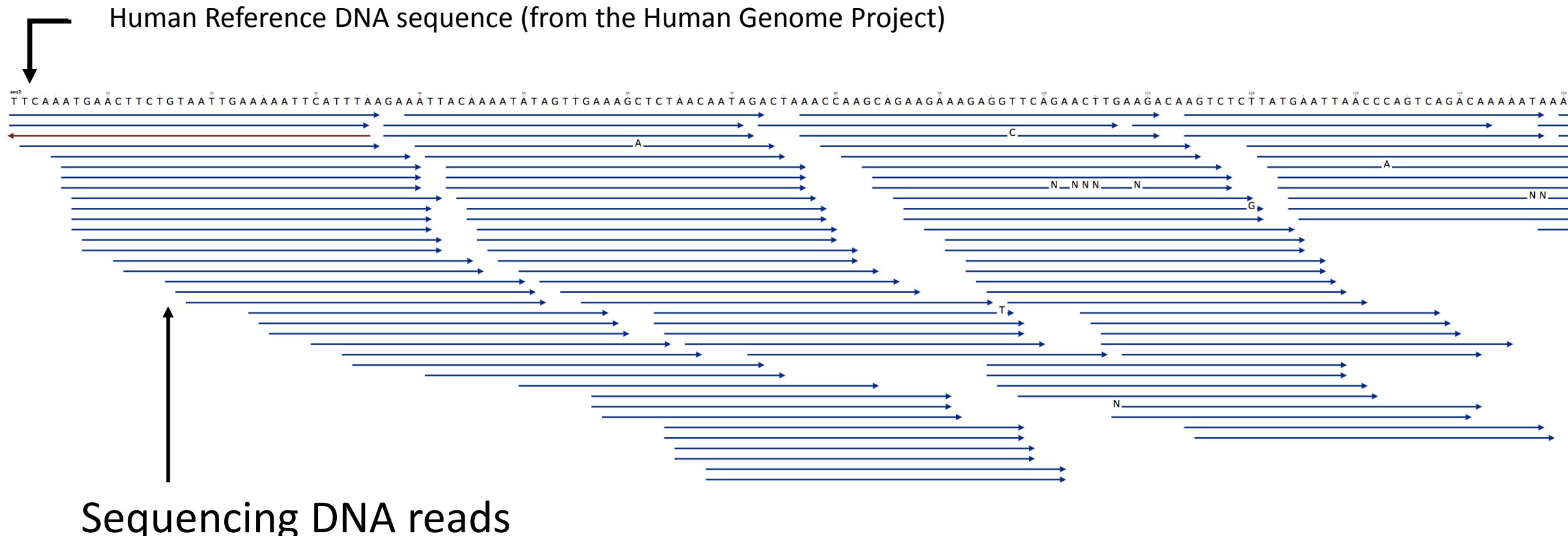
Make several copies of each fragment

DNA
sequencer



ATGCCAGAACTAATCTATA
CTTATTAAACATCATCA
CCAAAGAGAGAGTTCCAC

Sequence alignment of sequencing reads to Human Reference DNA sequence

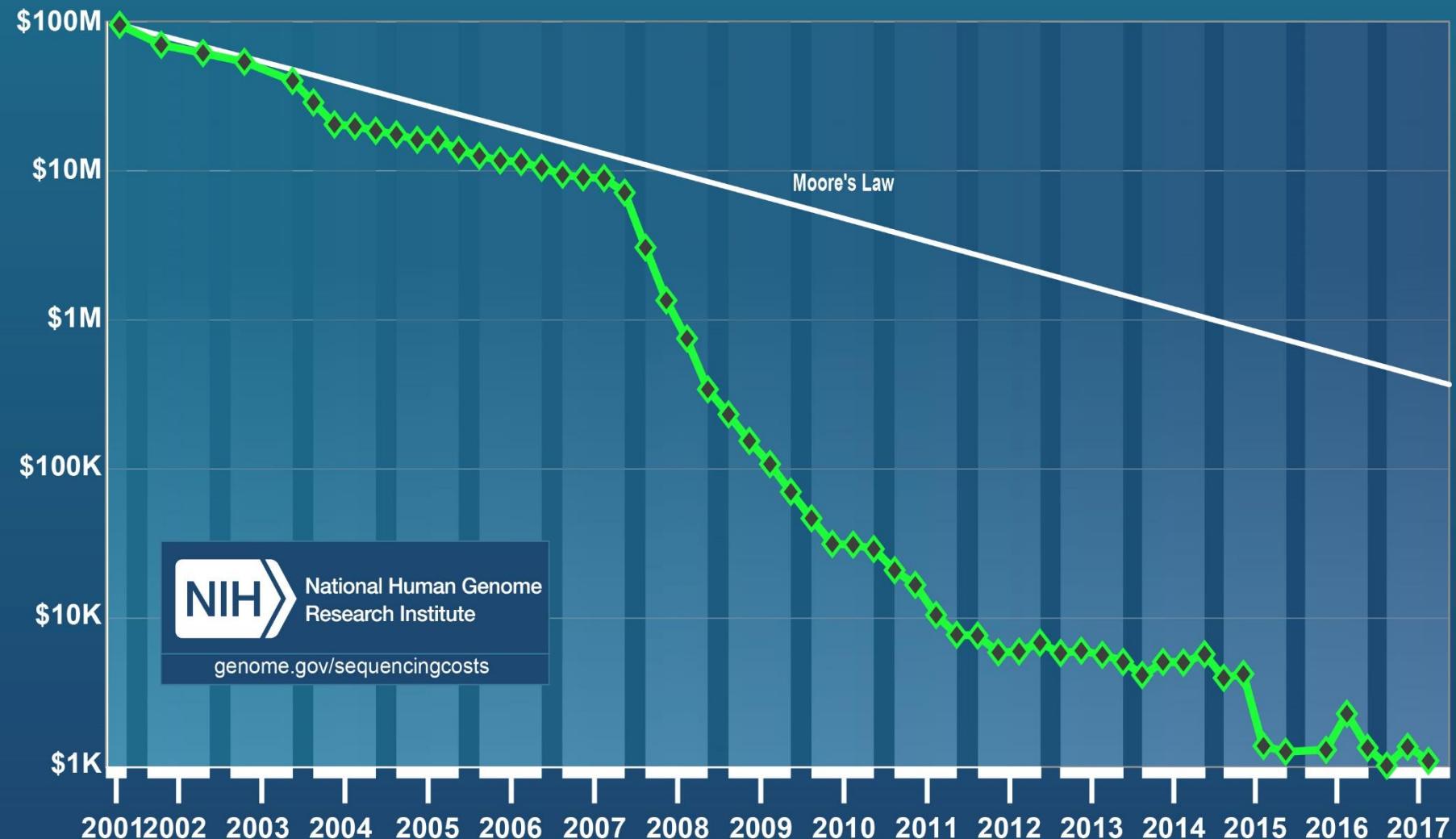


Human Reference DNA sequence

Target Genome **ATTTGCGCAGAGACCTAAGGCATTAGCTTGGCCCTAAAG**

Reads	ATTTGC	AGAGACCTAAG	TTAGCTTGGC	AAG
	TGCGCAGA			TGGCCCTAA
Overlapping	ATTTGC	AGAGACCTAAG	TTAGCTTGGC	AAG
	TGCGCAGA			TGGCCCTAA
Contigs	ATTTGCGCAGAGACCTAAG		TTAGCTTGGCCCTAAAG	

Cost per Genome



Growth of DNA Sequencing

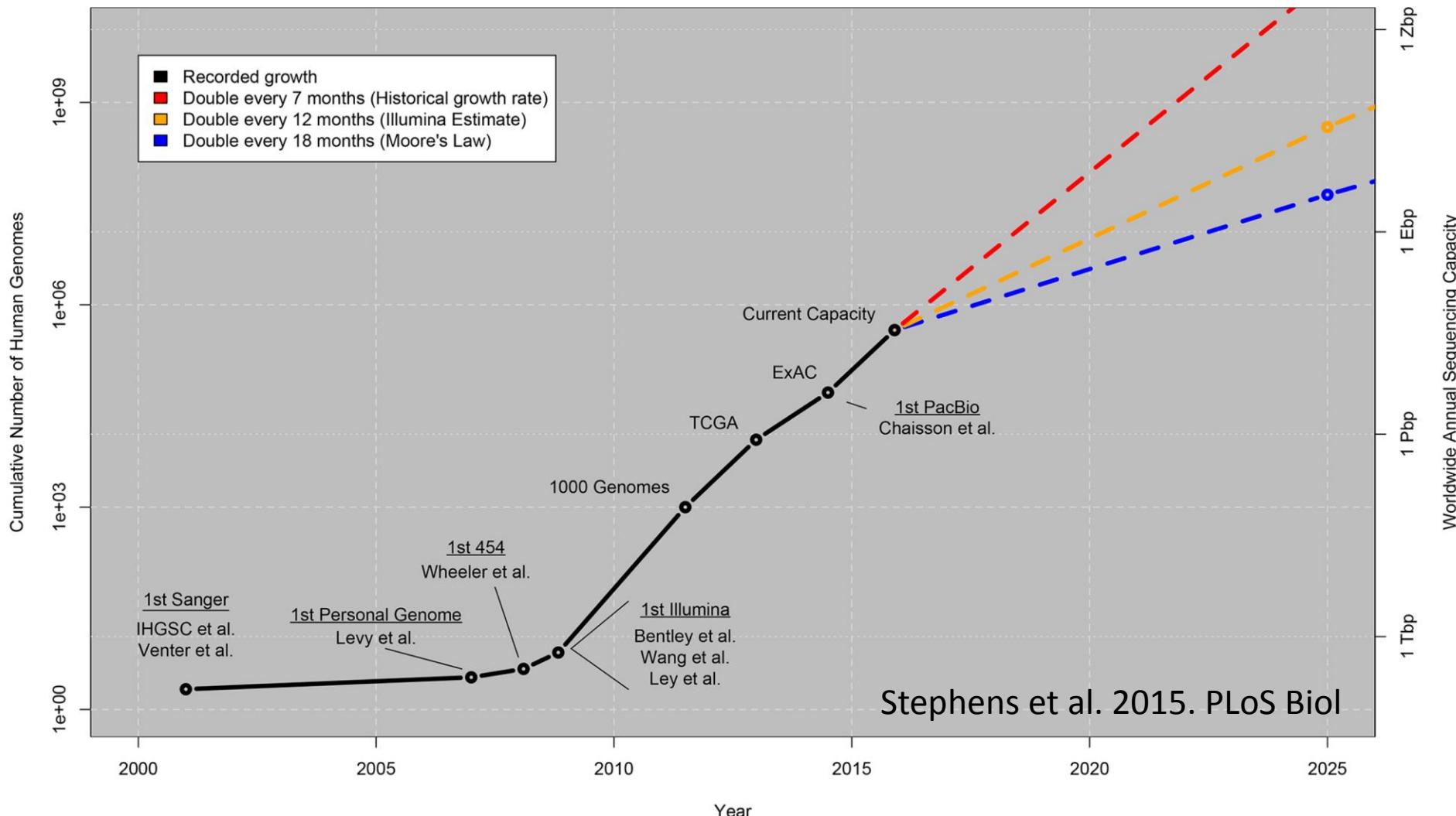
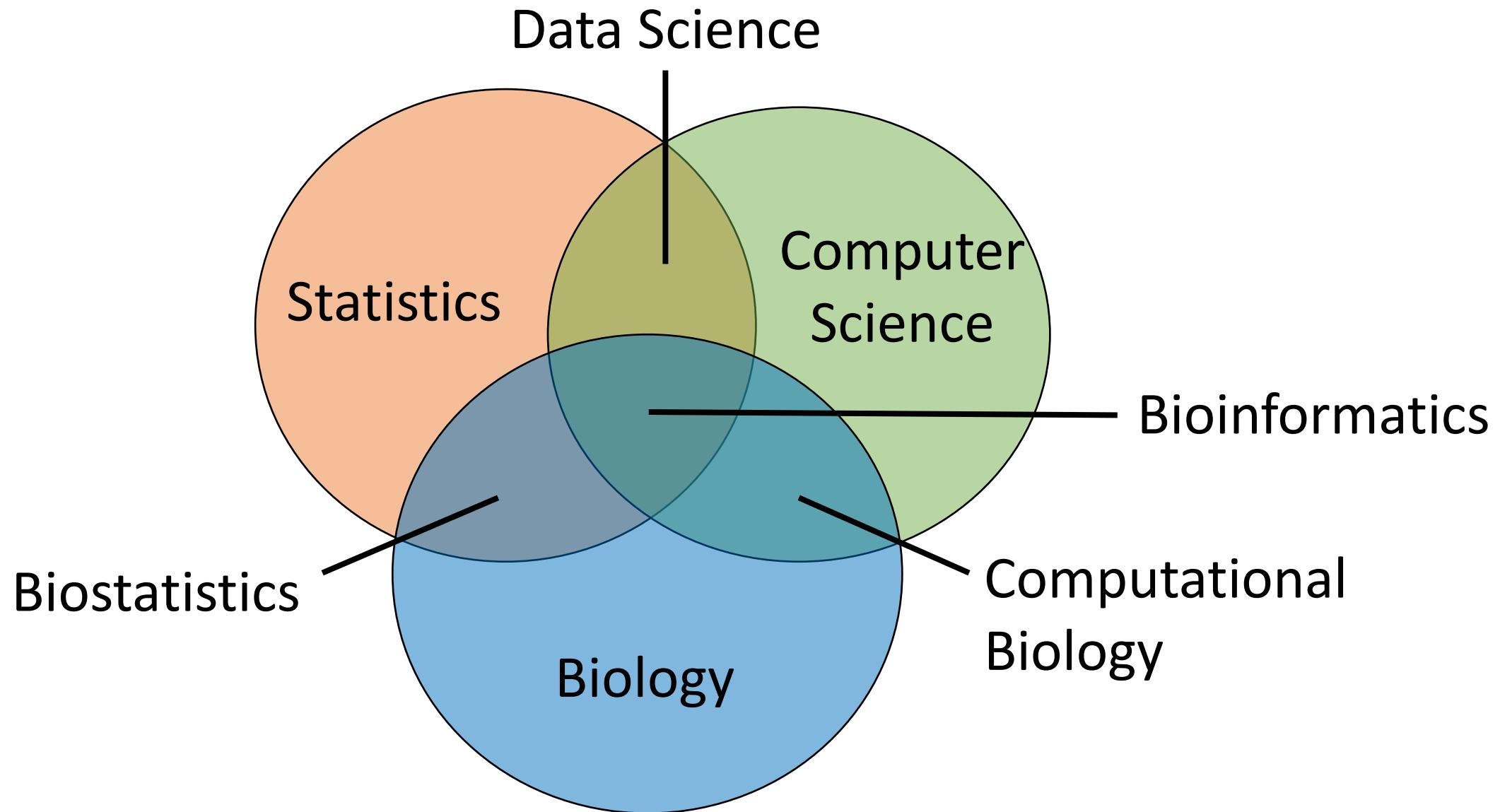


Fig 1. Growth of DNA sequencing. The plot shows the growth of DNA sequencing both in the total number of human genomes sequenced (left axis) as well as the worldwide annual sequencing capacity (right axis: Tera-basepairs (Tbp), Peta-basepairs (Pbp), Exa-basepairs (Ebp), Zetta-basepairs (Zbps)). The values through 2015 are based on the historical publication record, with selected milestones in sequencing (first Sanger through first PacBio human genome published) as well as three exemplar projects using large-scale sequencing: the 1000 Genomes Project, aggregating hundreds of human genomes by 2012 [3]; The Cancer Genome Atlas (TCGA), aggregating over several thousand tumor/normal genome pairs [4]; and the Exome Aggregation Consortium (ExAC), aggregating over 60,000 human exomes [5]. Many of the genomes sequenced to date have been whole exome rather than whole genome, but we expect the ratio to be increasingly favored towards whole genome in the future. The values beyond 2015 represent our projection under three possible growth curves as described in the main text.



you can think of **bioinformatics/computational biology** as **data science in biology**

Different Roles of Bioinformatician / Computational Biologist

- **Algorithm Development**
 - Skills in Maths + Stats. + Comp. Sci. essential. Knowledge in biology required.
- **Software Development**
 - Skills in coding & documentations
- **Data Processing**
 - implementation of bioinformatics pipelines, i.e. processing data generated by high-throughput analytical instruments ready to be analyzed by biologist
- **Data Analysis**
 - Knowledge in biology essential. Skills in Maths + Stats. + Comp. Sci. required.

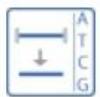
BIOINFORMATICS TOOLS FOR GENOMICS

Genomics is an interdisciplinary field of molecular biology focusing on the DNA content of living organisms. Genomics techniques are mainly focused on DNA sequencing, DNA structure analysis, genome editing, population genomics, DNA-protein interactions, phylogenomics, or synthetic biology. High-throughput DNA sequencing technologies and bioinformatics have transformed genome analysis by mapping and decrypting coding and non-coding DNA sequences, their evolution and inter-relationships. Software tools and databases are proposed here for genome annotation, phylogenomics studies, comparative genomics, genome editing, genome variant and DNA structure analysis, personal and population genomics, as well as epigenomic modifications which include DNA methylation.

DNA sequence



Genome annotation
3286 tools



WGS analysis
4486 tools



De novo sequencing analysis
2683 tools



aCGH data analysis
309 tools



Rep-seq analysis
127 tools



DNA sequence databases
1906 tools



WES analysis
3001 tools



Metagenomic sequencing analysis
2349 tools



SNP array data analysis
193 tools



16S rRNA-seq analysis
179 tools

MOST RECENT TOOLS

Simphony

Ozymandias

AutoKEGGRec

GenoNet

pyHam

iHam

GWASpro

See more ▾

MOST POPULAR TOOLS

knnAUC

Roary

Trimmomatic

FastTree

PHYLIP

Incremental BLAST

See more ▾

<https://omictools.com/genomics2-category>

Bioinformatics is a Big Data Problem

<https://dcc.icgc.org/pcawg>

ICGC Data Portal

PCAWG - PANCANCER ANALYSIS OF WHOLE GENOMES

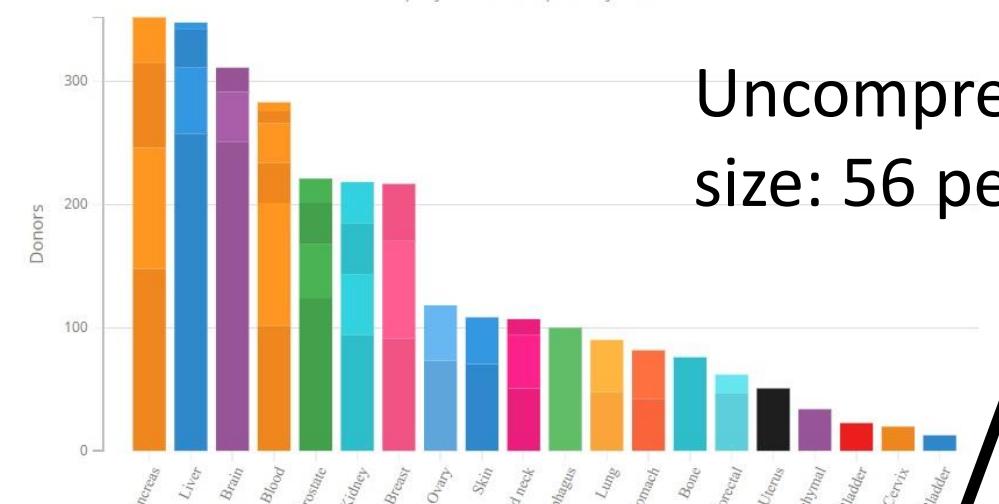
The Pancancer Analysis of Whole Genomes (PCAWG) study is an international collaboration to identify common patterns of mutation in more than 2,800 cancer whole genomes from the International Cancer Genome Consortium. Building upon previous work which examined cancer coding regions (Cancer Genome Atlas Research Network, The Cancer Genome Atlas Pan-Cancer analysis project, *Nat. Genet.* 2013 45:1113), this project is exploring the nature and consequences of somatic and germline variations in both coding and non-coding regions, with specific emphasis on cis-regulatory sites, non-coding RNAs, and large-scale structural alterations.

In order to facilitate the comparison among diverse tumor types, all tumor and matched normal genomes have been subjected to a uniform set of alignment and variant calling algorithms, and must pass a rigorous set of quality control tests. The research activities are coordinated by a series of working groups comprising more than 700 scientists and covering the following themes:

1. Novel somatic mutation calling methods
2. Analysis of mutations in regulatory regions
3. Integration of the transcriptome and genome
4. Integration of the epigenome and genome
5. Consequences of somatic mutations on pathway and network activity
6. Patterns of structural variations, signatures, genomic correlations, retrotransposons and

Analysis performed using 6 super-computers from around the world

Donor Distribution by Primary Site
48 projects and 20 primary sites



Primary Site	Donors
Pancreas	~150
Liver	~280
Brain	~300
Blood	~280
Prostate	~220
Kidney	~220
Breast	~220
Ovary	~120
Skin	~110
Head and neck	~100
Esophagus	~100
Lung	~80
Stomach	~70
Bone	~70
Colorectal	~60
Uterus	~50
Mesenchymal	~40
Bladder	~30
Cervix	~20
Gall Bladder	~10

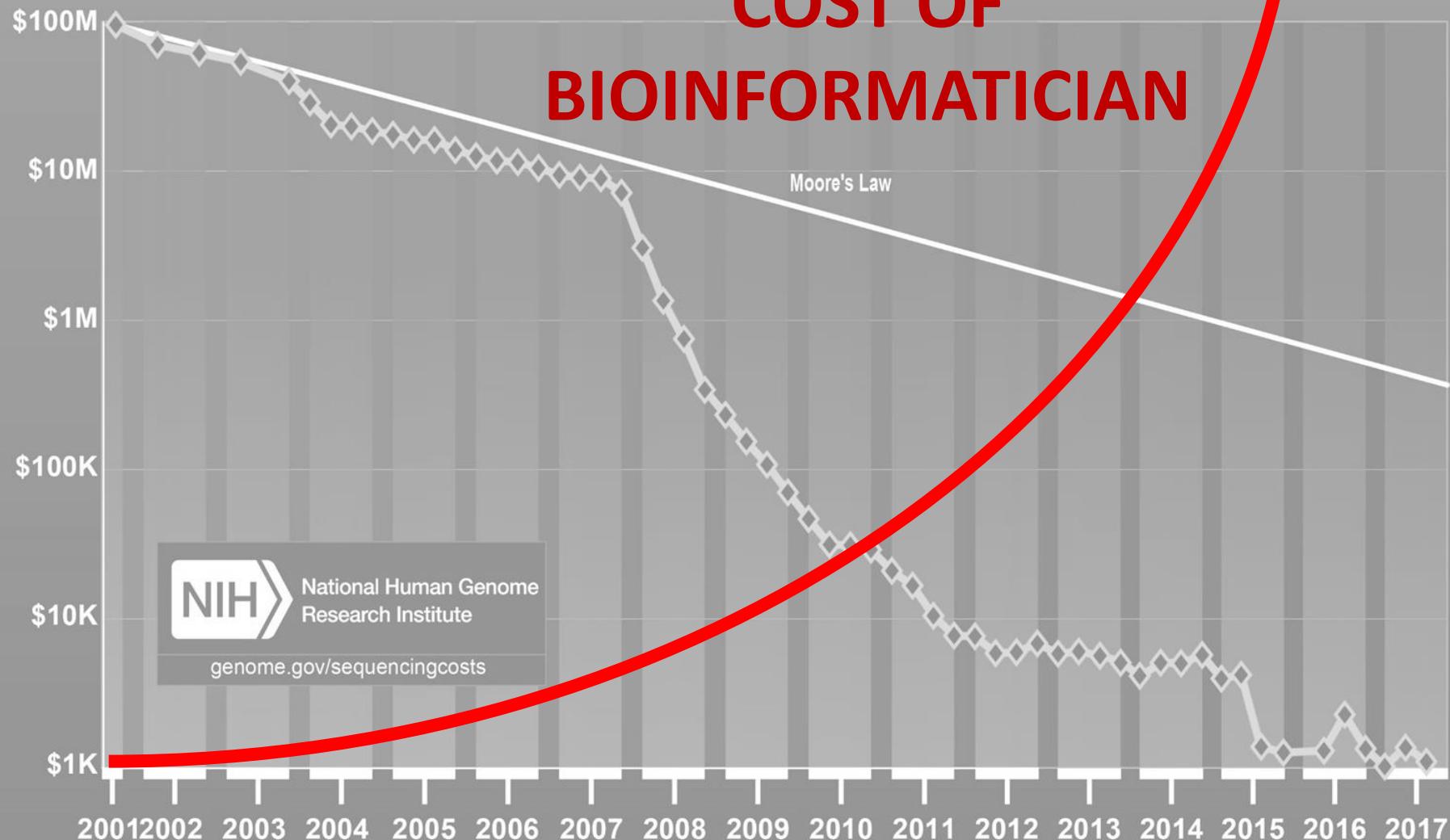
Uncompressed data size: 56 peta bytes

2,832 Donors **74,023 Files** **803.44 TB**

21

Cost per Genome

COST OF BIOINFORMATICIAN



Major IT companies heavily investing in Bioinformatics



Google
Genomics

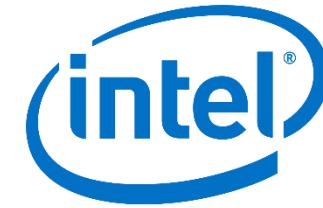


Google Cloud



Calico

YAHOO!
RESEARCH



Baidu 百度

Microsoft®
Research

NVIDIA



amazon
web services

TEA BREAK

After the break

How machine learning is being used to solve biological problems?

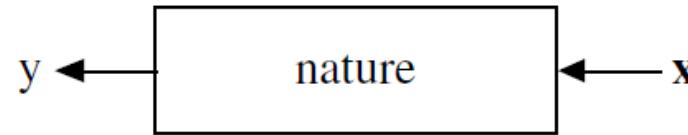


PART - II

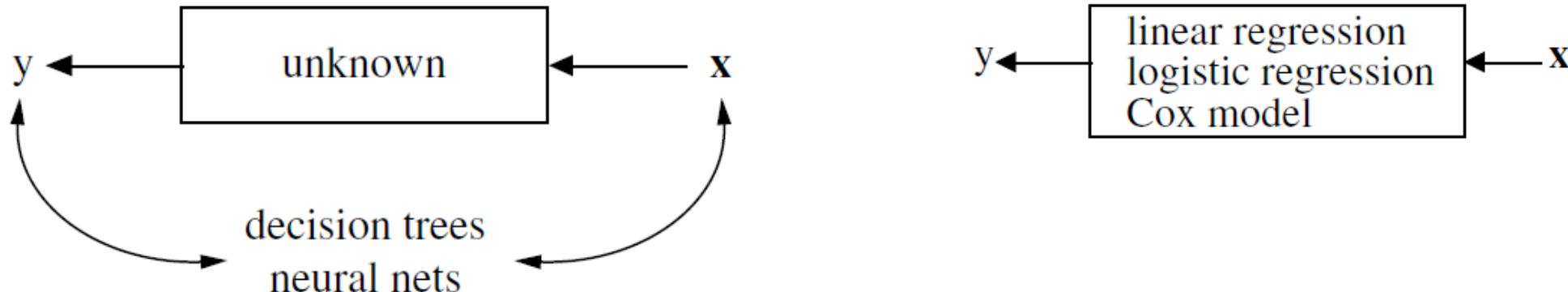
Applications of Machine Learning in Biology & Medicine

- **WHY** is machine learning a good tool for biological studies?
- **WHERE** in biology & medicine machine learning tools are being used?
- **WHEN NOT TO USE** machine learning?

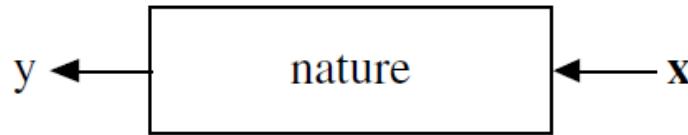
Two cultures of Machine Learning



BLACK BOX APPROACH OPEN BOX APPROACH



Two cultures of Machine Learning



BLACK BOX APPROACH



Classification of Cat vs Dog

Features that separate Cat vs Dog
may not be important for the end-user

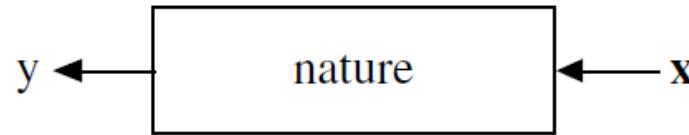
OPEN BOX APPROACH



Classification of Cat vs Dog

Features that separate Cat vs Dog
is important for the end-user

Two cultures of Machine Learning



BLACK BOX APPROACH



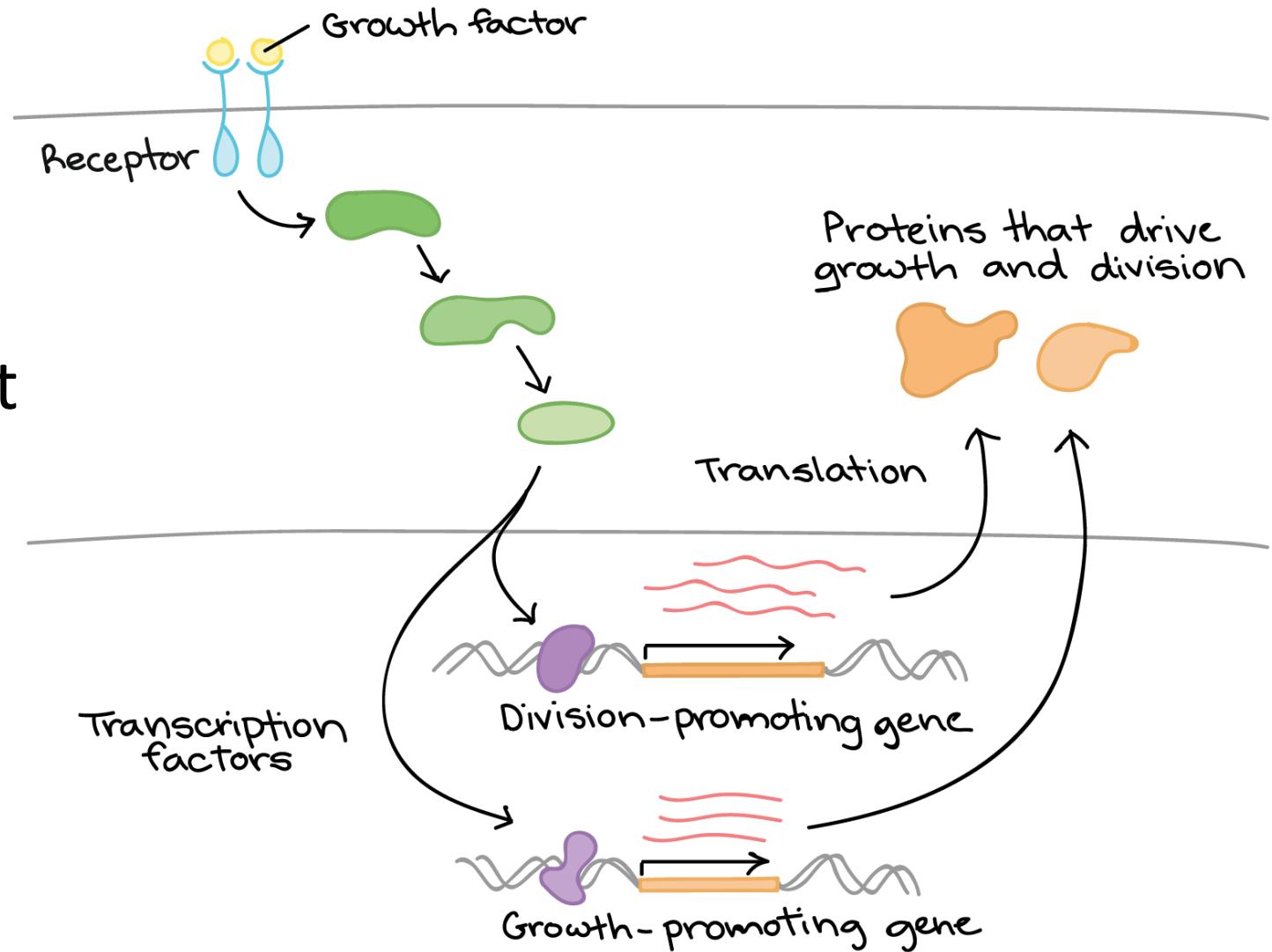
OPEN BOX APPROACH

WHY does some patients responds
to Drug-X while others does not
respond the drug at all?

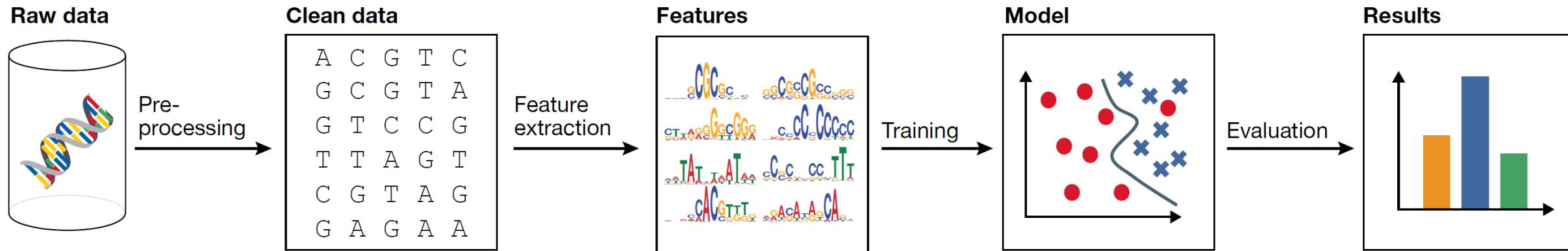
Diagnosis of Tuberculosis (TB) using
X-ray images (Image classification)

Gene Regulation

The pattern of DNA sequence is very important for attachment of key protein that regulate the function of a gene.



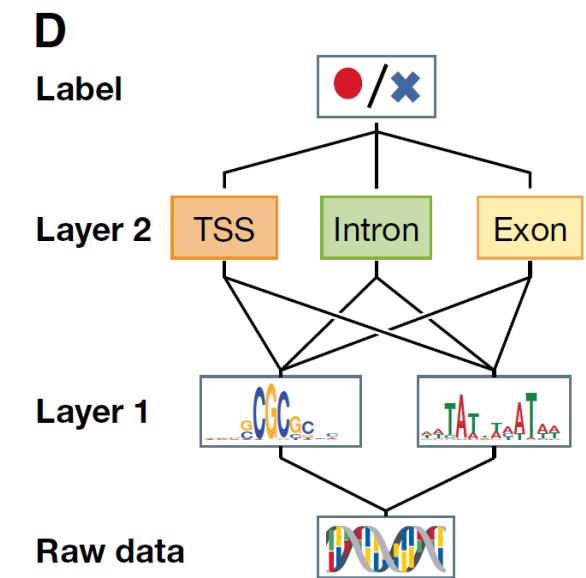
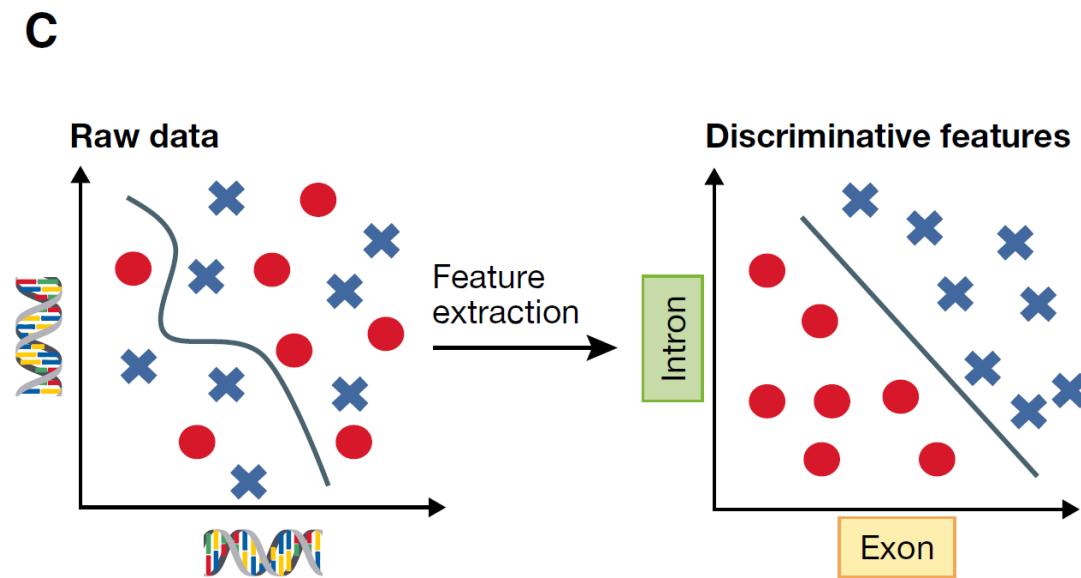
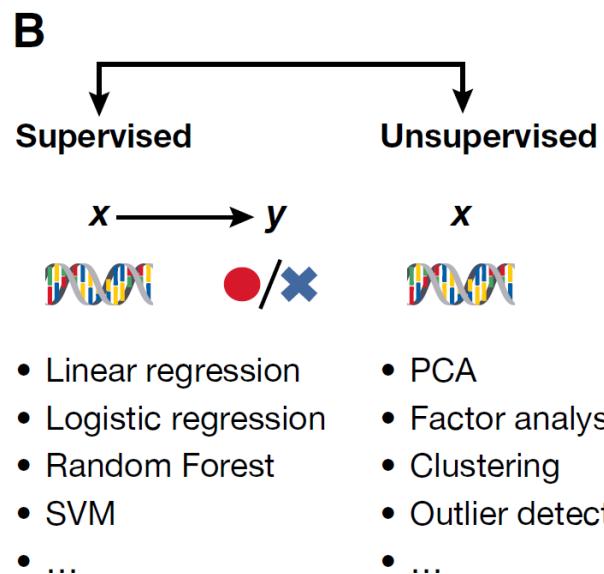
Machine learning workflow



The classical machine learning workflow can be broken down into four steps:

- **data pre-processing,**
- **feature extraction,**
- **model learning** and
- **model evaluation.**

Machine learning approaches

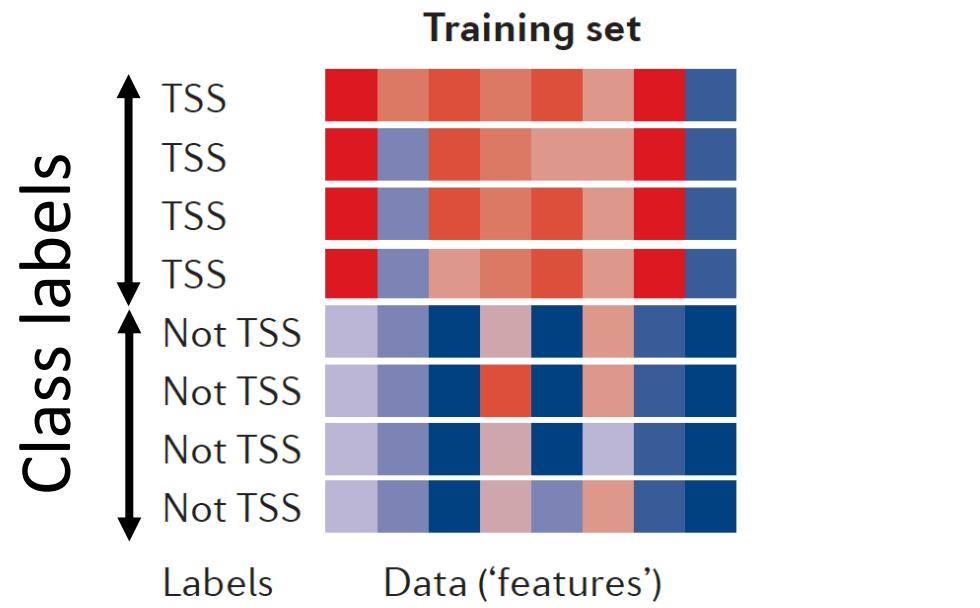


Raw input data are often high-dimensional and related to the corresponding label in a complicated way, which is challenging for many classical machine learning algorithms (left plot). Alternatively, higher-level features extracted using a deep model may be able to better discriminate between classes (right plot).

Deep networks use a hierarchical structure to learn increasingly abstract feature representations from the raw data.

An example of a machine learning application

Training set of DNA sequences



Testing set

The figure shows the **Testing set** of DNA sequences. The rows are numbered 1 through 8. Each row consists of a sequence of colored squares representing features, followed by a predicted label. The predicted labels are listed on the right side of the diagram.

1.	2.	3.	4.	5.	6.	7.	8.
Light Blue	Light Blue	Dark Blue	Light Blue	Light Blue	Light Blue	Dark Blue	Dark Blue
Red	Red	Light Blue	Red	Red	Red	Red	Red
Red	Red	Red	Red	Red	Red	Red	Red
Light Blue	Dark Blue	Light Blue	Dark Blue	Dark Blue	Dark Blue	Dark Blue	Dark Blue
Light Blue	Dark Blue	Light Blue	Dark Blue	Dark Blue	Dark Blue	Dark Blue	Dark Blue
Red	Light Blue	Red	Red	Red	Red	Red	Red
Light Blue	Light Blue	Dark Blue	Light Blue	Light Blue	Light Blue	Dark Blue	Dark Blue
Red	Light Blue	Red	Red	Red	Red	Red	Red

Color indicates weights
or scores on the
sequence patterns

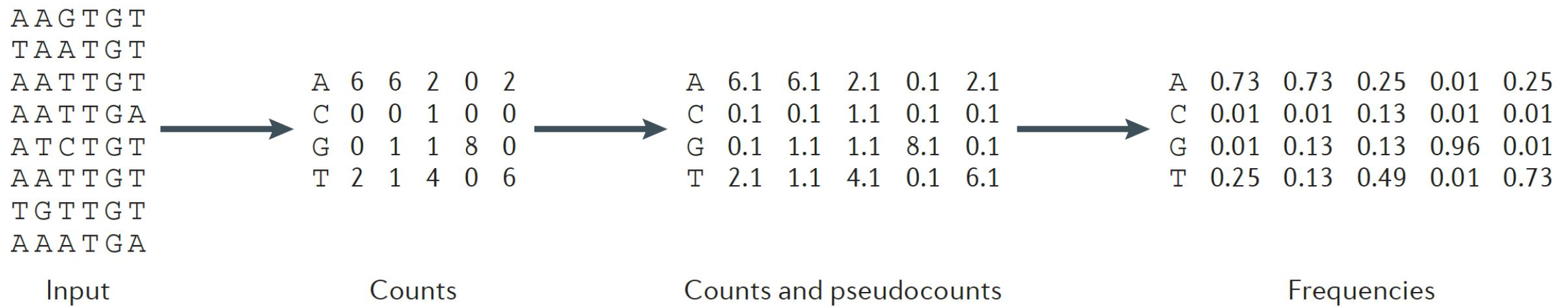
Predicted
labels

1. Not TSS
2. TSS
3. TSS
4. Not TSS
5. Not TSS
6. TSS
7. Not TSS
8. TSS

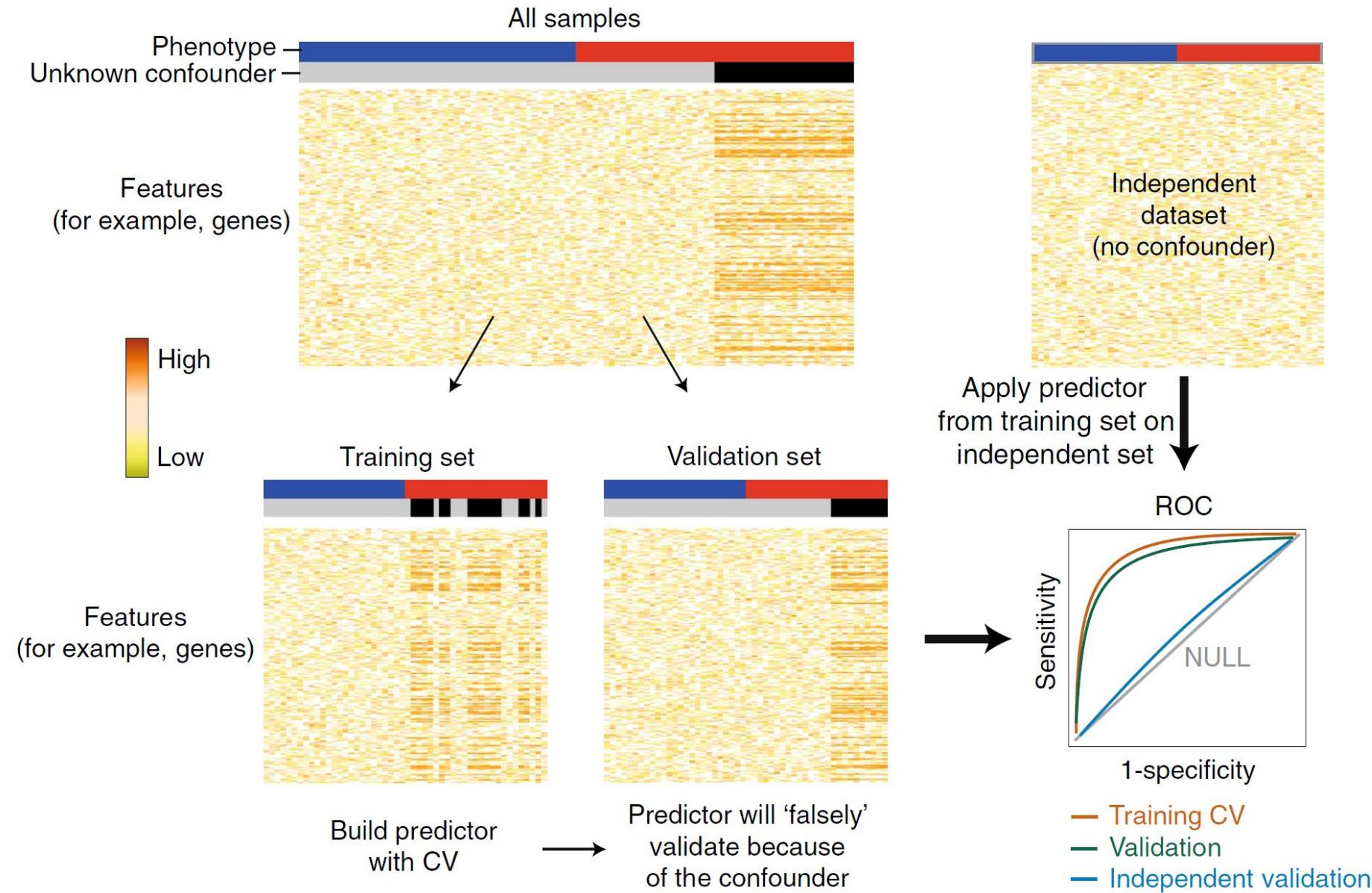
binary labels indicates whether each sequence is centered
on a transcription start site (TSS) or not.

Libbrecht and Nobel. 2015. Nature Review Genetics

Sequence to Scores

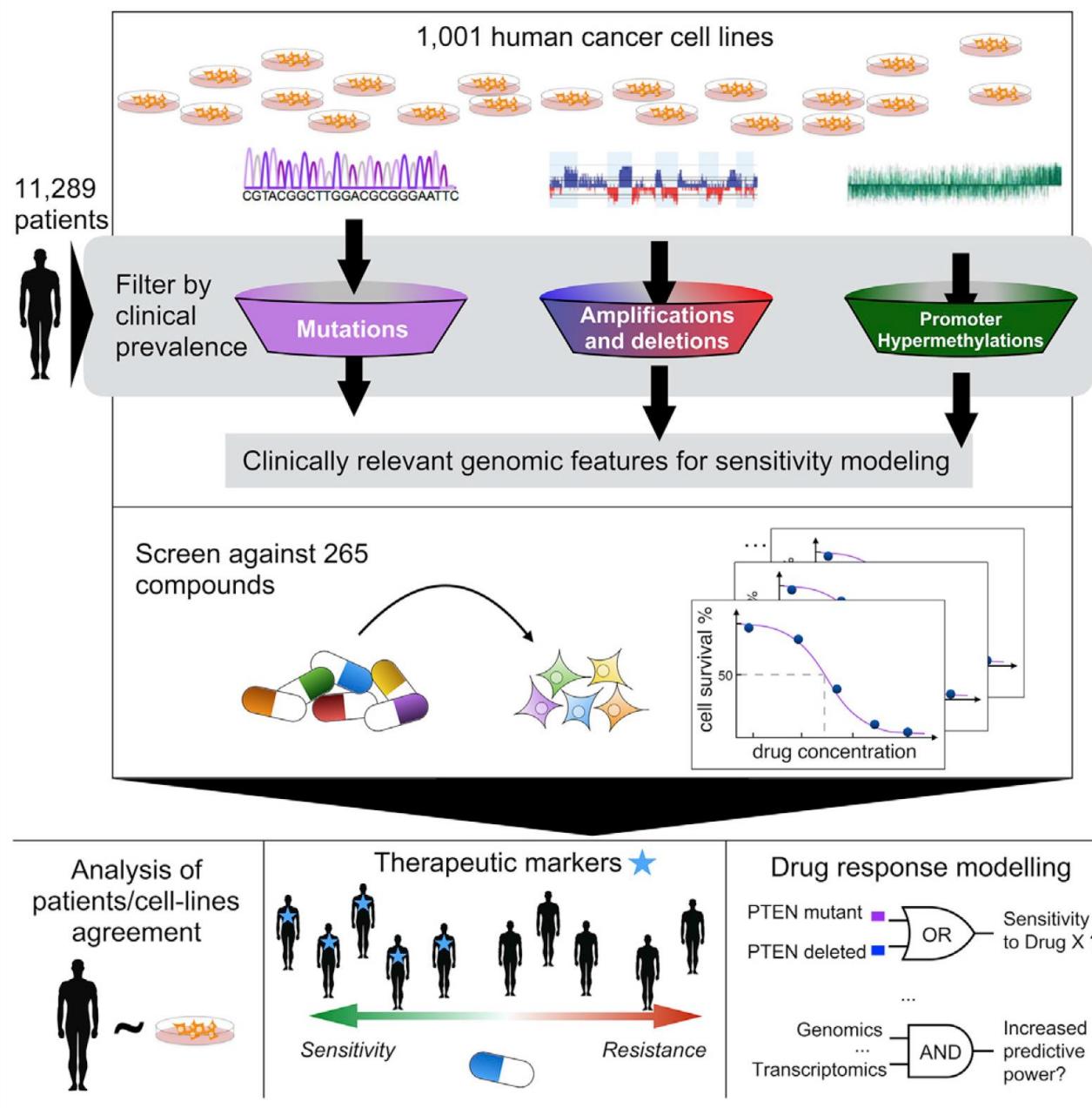
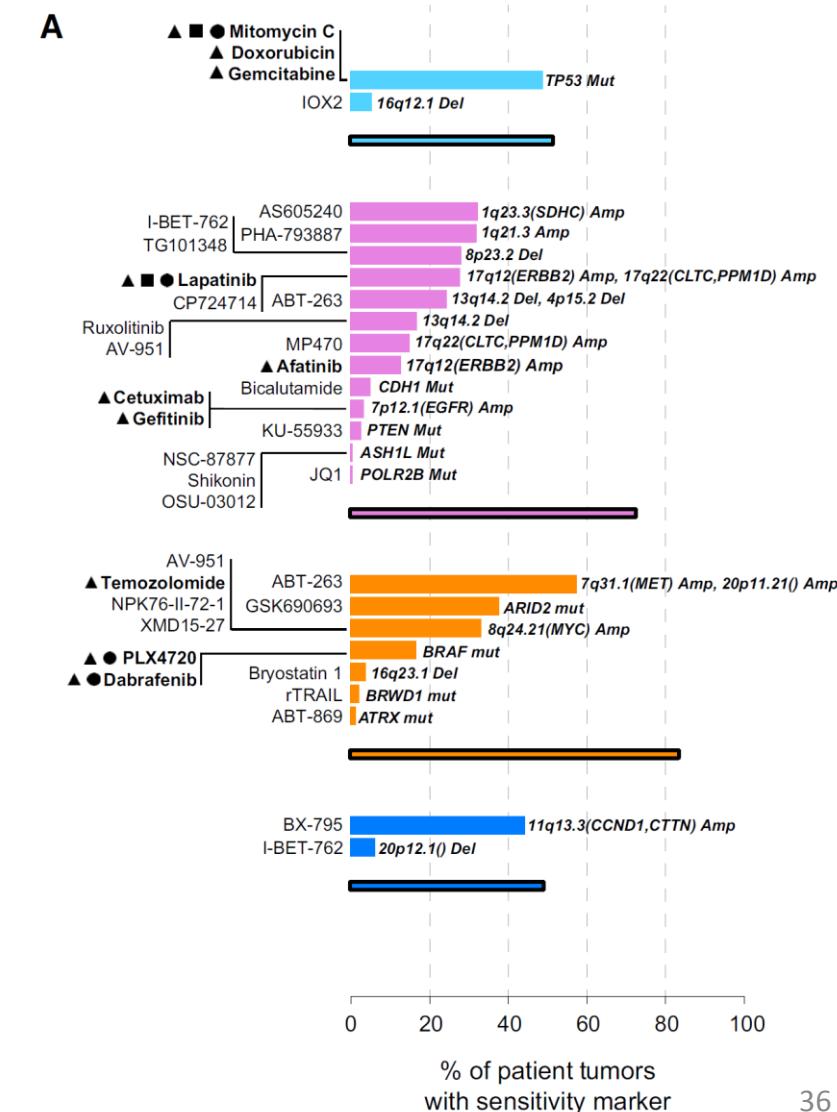


Using Gene Expression Data



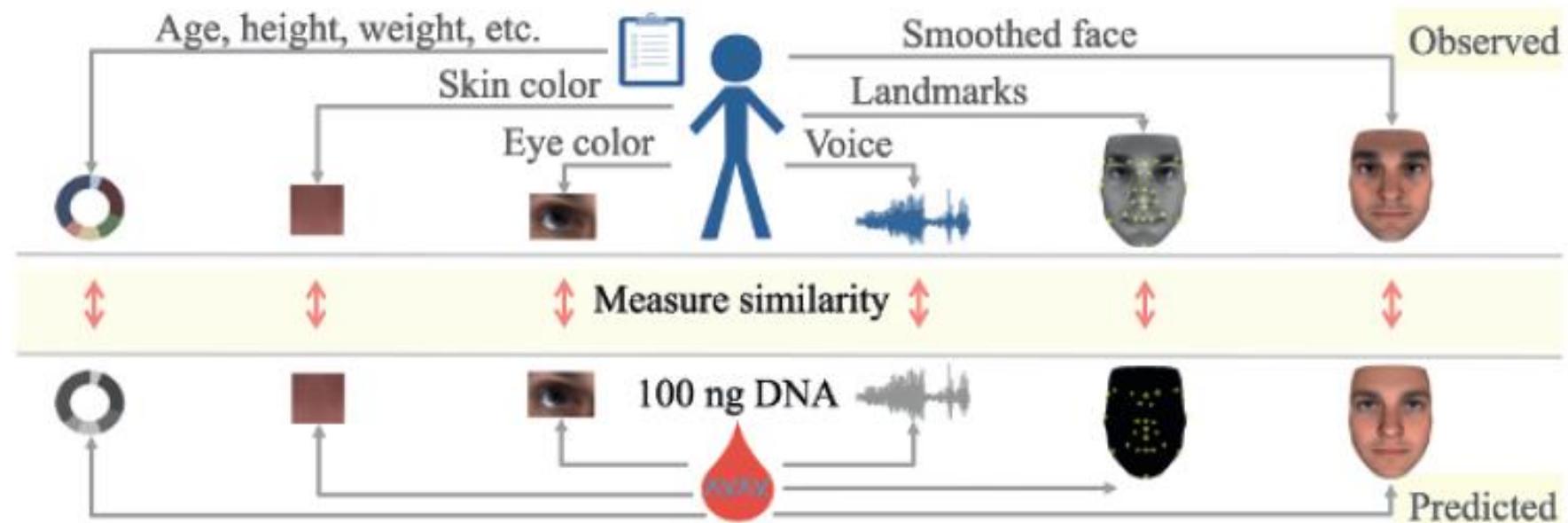
Cell

A Landscape of Pharmacogenomic Interactions in Cancer



Identification of individuals by trait prediction using whole-genome sequencing data

Christoph Lippert^{a,1}, Riccardo Sabatini^a, M. Cyrus Maher^a, Eun Yong Kang^a, Seunghak Lee^a, Okan Arikán^a, Alena Harley^a, Axel Bernal^a, Peter Garst^a, Victor Lavrenko^a, Ken Yocom^a, Theodore Wong^a, Mingfu Zhu^a, Wen-Yun Yang^a, Chris Chang^a, Tim Lu^b, Charlie W. H. Lee^b, Barry Hicks^a, Smriti Ramakrishnan^a, Haibao Tang^a, Chao Xie^c, Jason Piper^c, Suzanne Brewerton^c, Yaron Turpaz^{b,c}, Amalio Telenti^b, Rhonda K. Roby^{b,d,2}, Franz J. Och^a, and J. Craig Venter^{b,d,1}



A universal SNP and small-indel variant caller using deep neural networks

nature
biotechnology

Ryan Poplin^{1,2}, Pi-Chuan Chang², David Alexander², Scott Schwartz², Thomas Colthurst², Alexander Ku², Dan Newburger¹, Jojo Dijamco¹, Nam Nguyen¹, Pegah T Afshar¹, Sam S Gross¹, Lizzie Dorfman^{1,2}, Cory Y McLean^{1,2} & Mark A DePristo^{1,2}

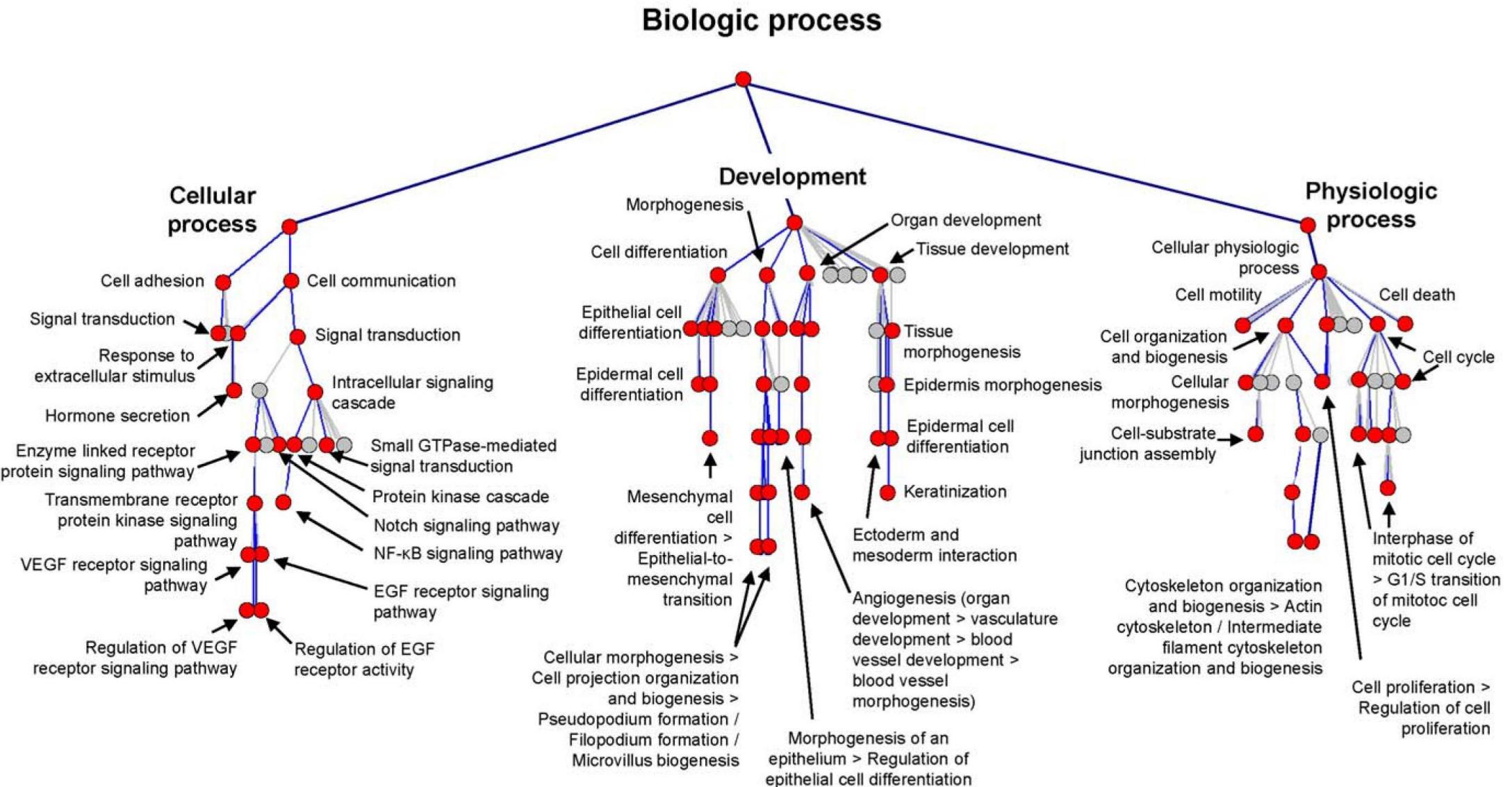
Pileup: Standard format for mapped data, position summaries

Seq.	Pos.	Len.	Alignment	Quality
Ref.				
seq1	272	T 24	,\$.....,.....,.....^+	<<<+;<<<<<<<=<;>7<&
	273	T 23	,.....,.....,.....A	<<<;<<<<<<3=<<<;<<+
	274	T 23	,\$.....,.....,.....	7<7;<;<<<<<<=<;<;<<6
	275	A 23	,\$.....,.....,.....^l.	<+,9*<<<<<<=<<;<<<
	276	G 22	...T,.....,.....	33;+<<7=7<<7<&<<1;<<6<
	277	T 22,..C.,.....G.	+7;<<<<<&<=<<;<<&<
	278	G 23,.....,.....,.....^k.	%38*<<;<7<<7=<<<;<<<<
	279	C 23	A..T,.....,.....	;75&<<<<<<=<<<9<<;<<

Google's DeepVariant:
Uses Pileup images as
input to identify changes
in DNA sequences
utilizing CNN.

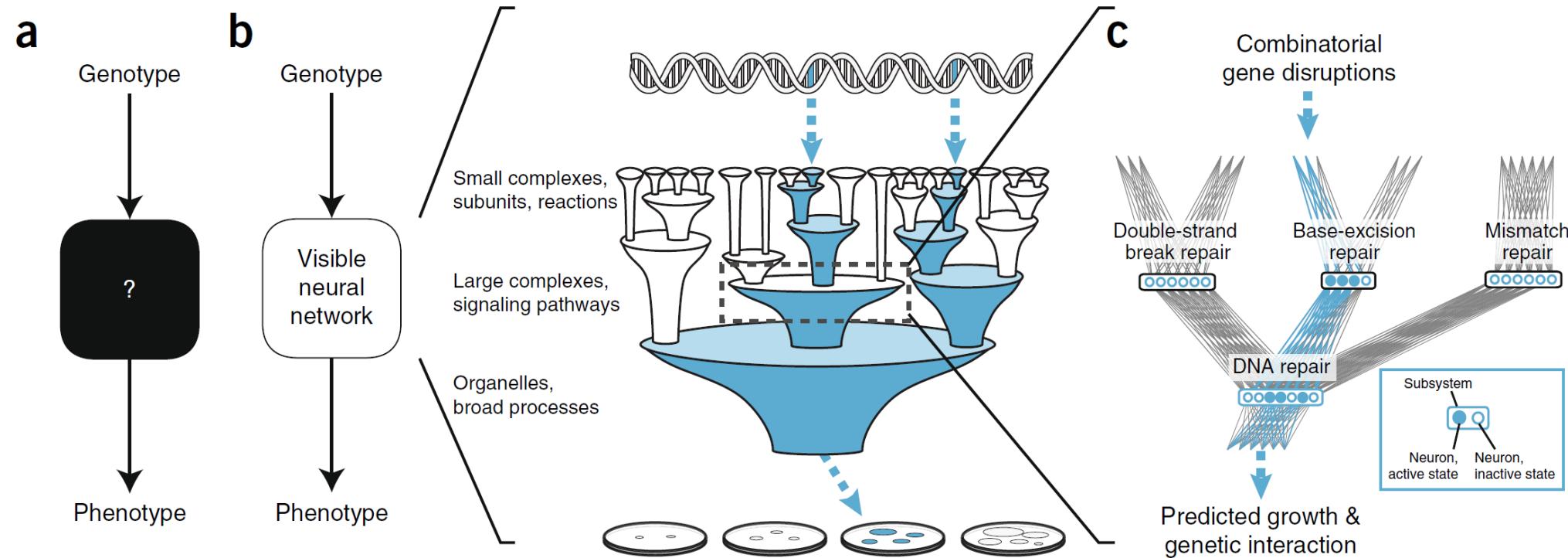
Treats pattern of strings
as “images”

Hierarchical mapping of Gene Ontology (GO) categories

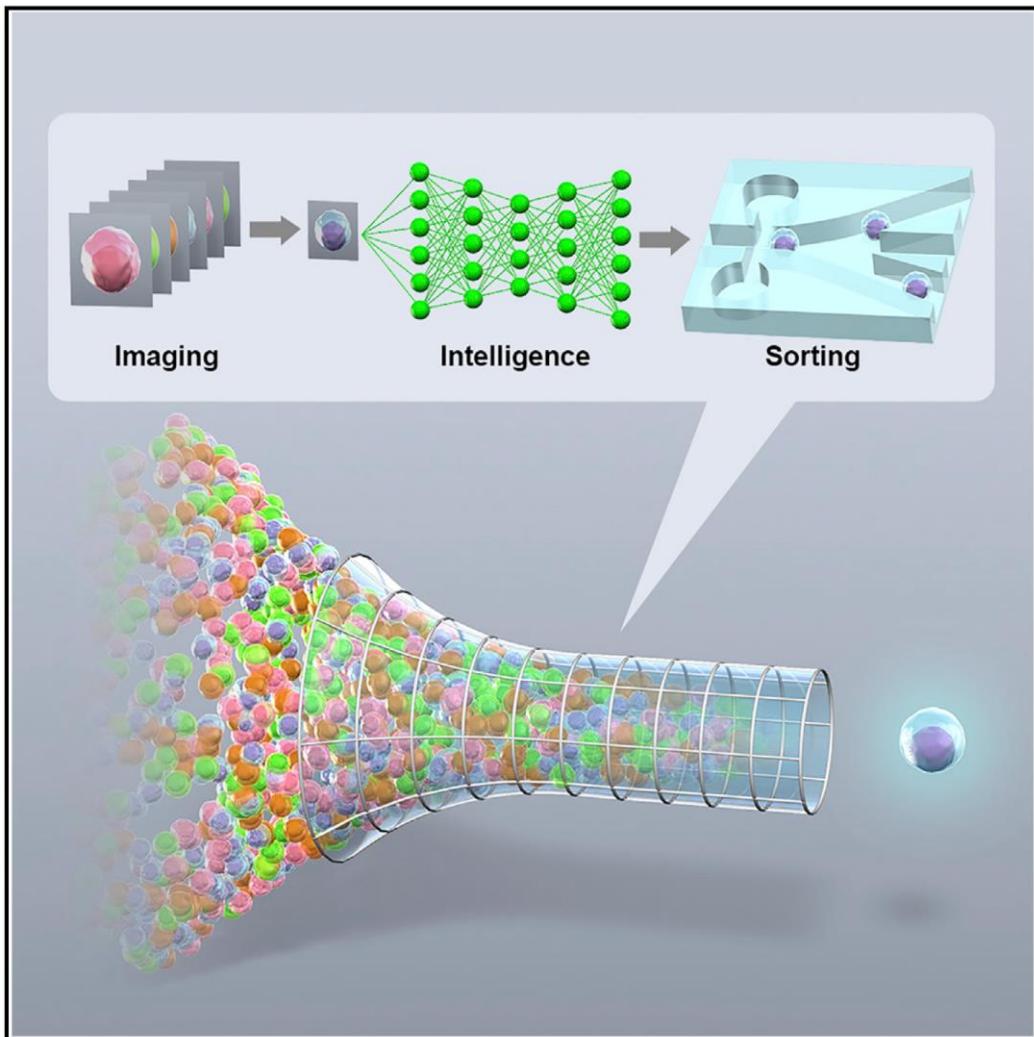


Using deep learning to model the hierarchical structure and function of a cell

Jianzhu Ma^{1,5} , Michael Ku Yu^{1,2,5}, Samson Fong^{1,3,5}, Keiichiro Ono¹, Eric Sage¹, Barry Demchak¹, Roded Sharan⁴ & Trey Ideker¹⁻³ 

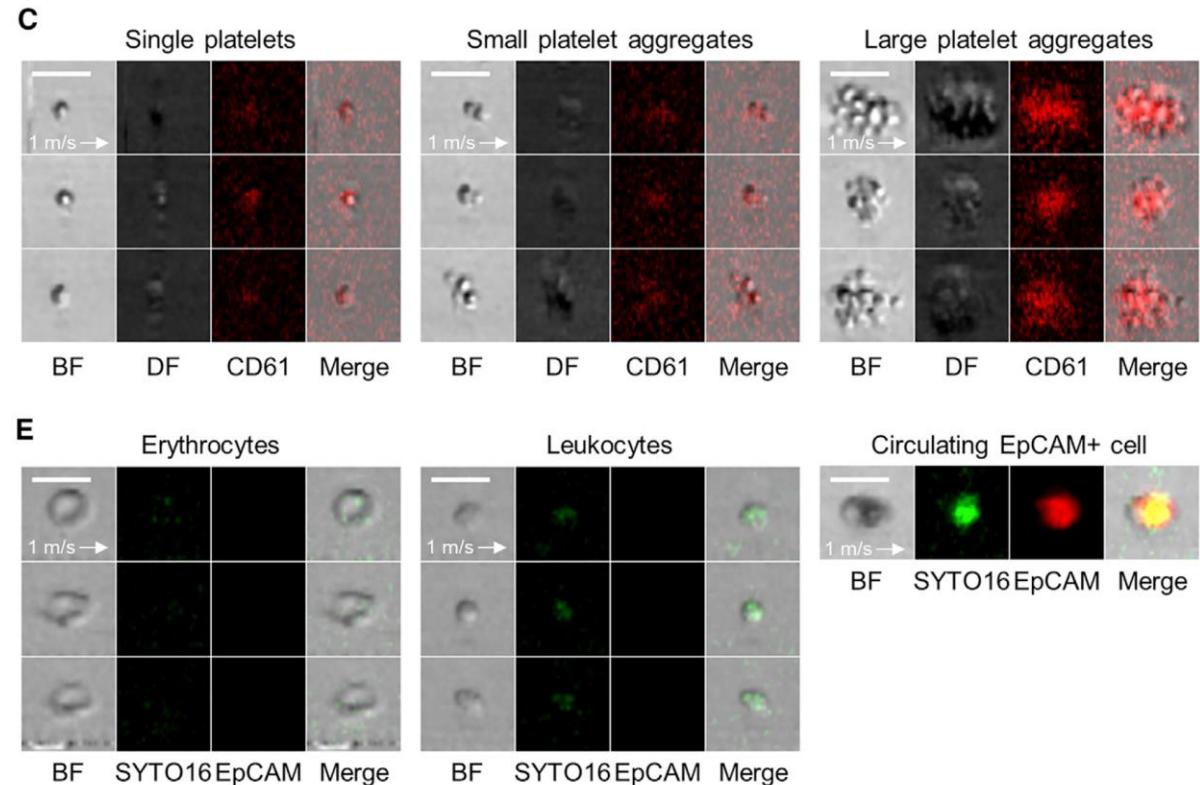
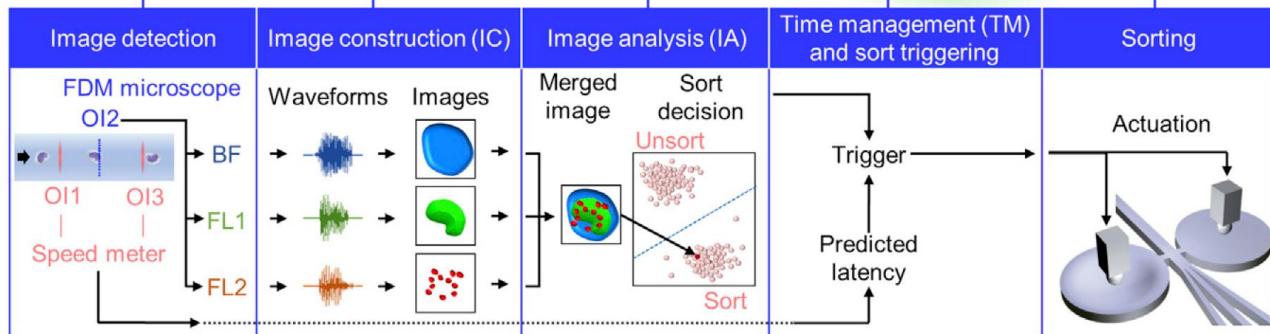
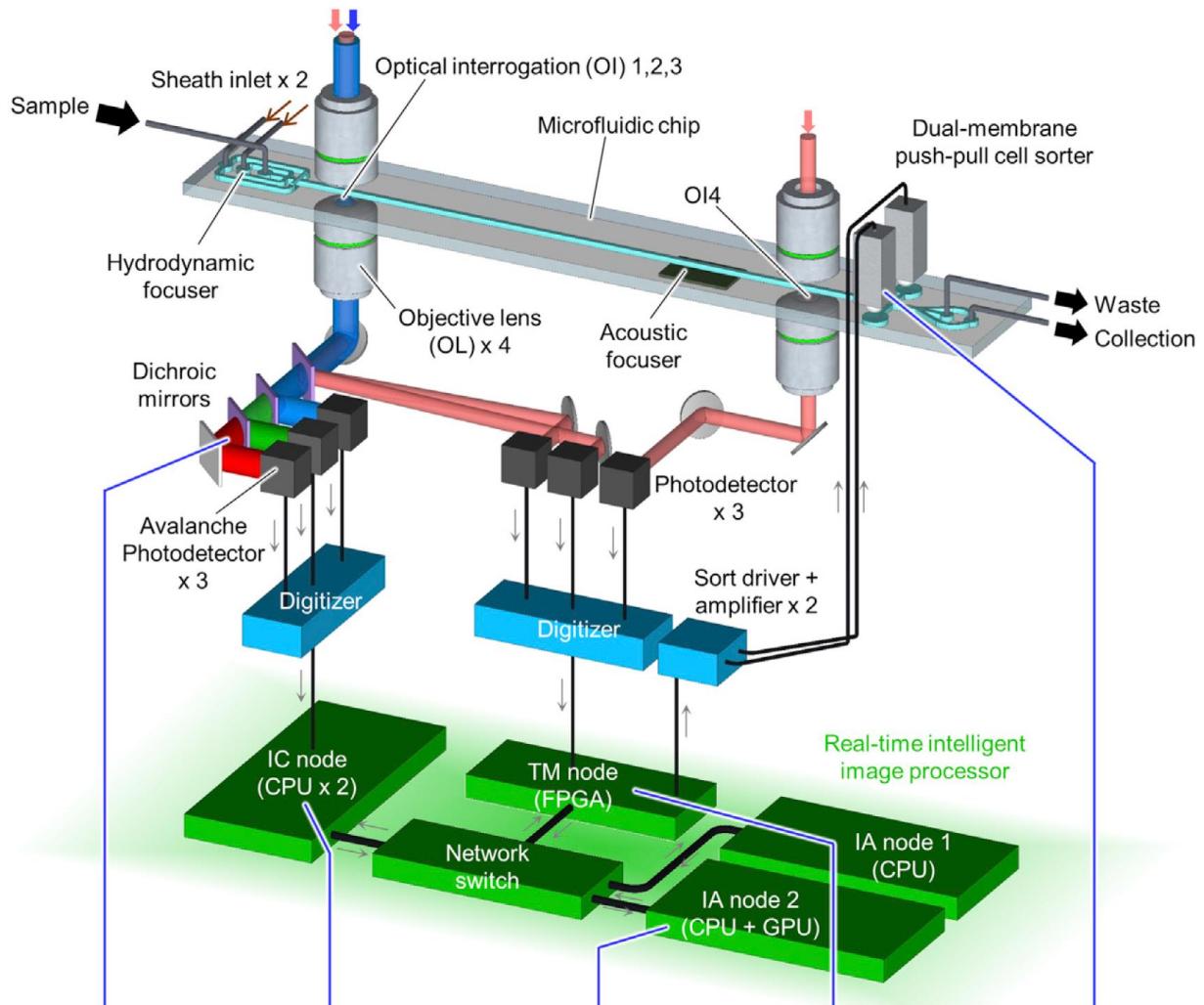


Intelligent Image-Activated Cell Sorting



In Brief

Artificial-intelligence-assisted, image-based flow cytometry in real-time enables rapid cell sorting based on unique chemical and morphological features.

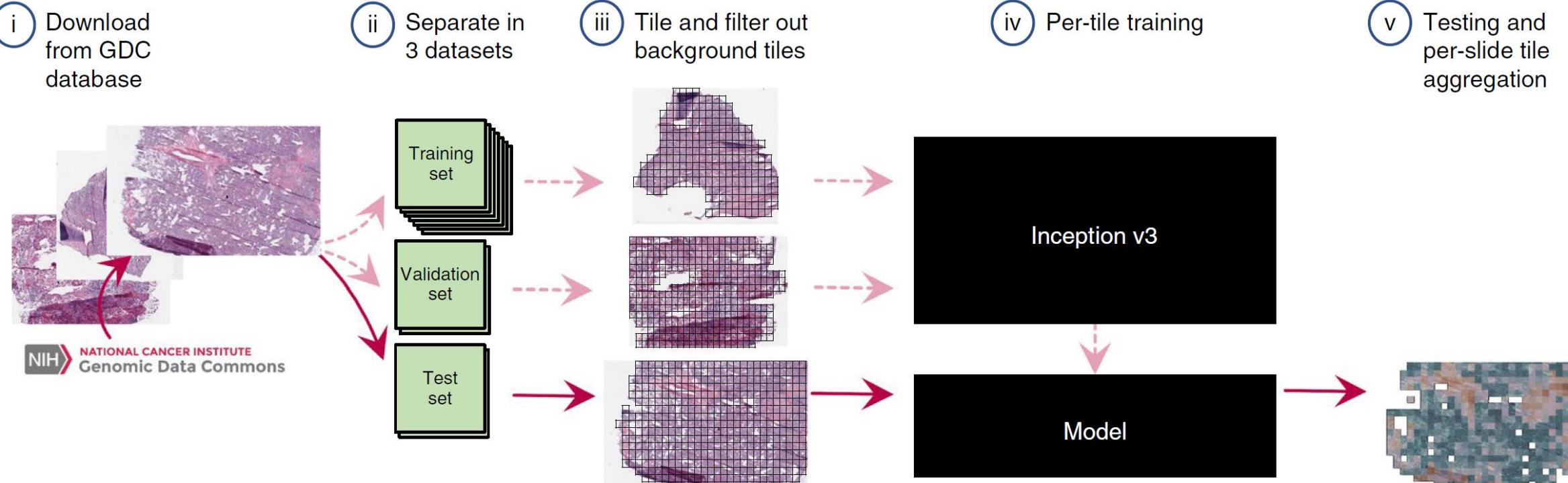


Classification and mutation prediction from non-small cell lung cancer histopathology images using deep learning

nature
medicine

Nicolas Coudray ^{1,2,9}, Paolo Santiago Ocampo^{3,9}, Theodore Sakellaropoulos⁴, Navneet Narula³,
Matija Snuderl³, David Fenyö^{5,6}, Andre L. Moreira^{3,7}, Narges Razavian ^{8*} and Aristotelis Tsirigos ^{1,3*}

b



Deep Genomic Signature for early metastasis prediction in prostate cancer

Hossein Sharifi-Noghabi^{1,3}, Yang Liu², Nicholas Erho², Raunak Shrestha^{3,4},
Mohammed Alshalalfa², Elai Davicioni², Colin C. Collins^{1,3,4}, and Martin
Ester^{1,3,*}

¹ School of Computing Science, Simon Fraser University, Burnaby, BC, Canada.

² GenomeDx Biosciences, Vancouver, BC, Canada.

³ Vancouver Prostate Centre, Vancouver, BC, Canada.

⁴ Department of Urologic Sciences, University of British Columbia, Vancouver, BC, Canada

Sharifi-Noghabi et al. (in submission)

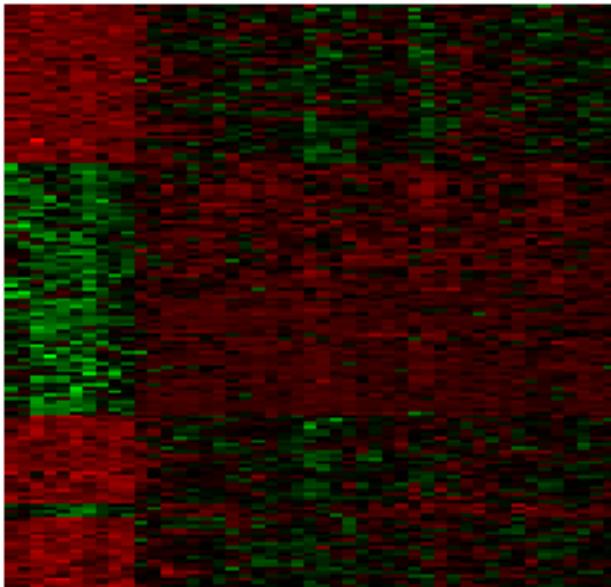
<https://doi.org/10.1101/276055>

- Metastasis is the medical term for cancer that spreads to a different part of the body from where it started.
- Metastasis are mostly lethal
- So **Early Prediction** of which patients are highly likely to develop metastasis is very important

Can we improve the performance of genomic classifiers by using poorly- and well-annotated cohorts together?

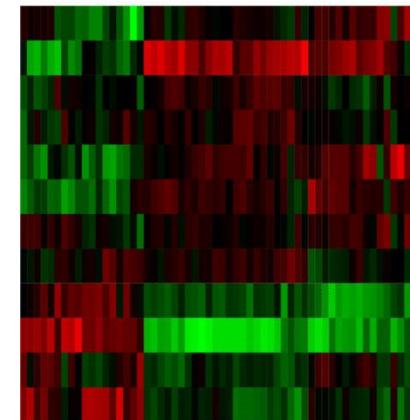
Poorly-annotated cohort

(Without long-term clinical outcomes)



Well-annotated cohort

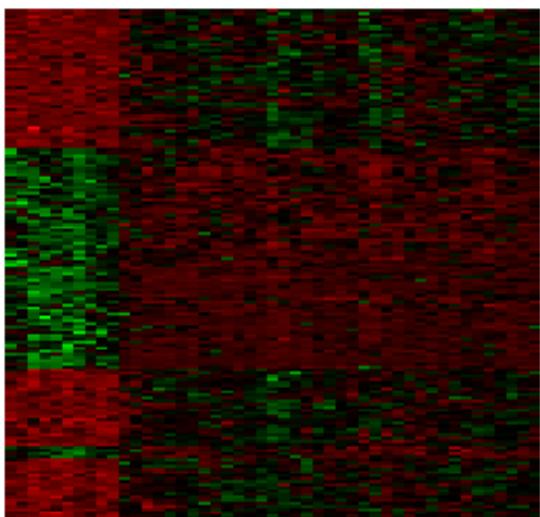
(With long-term clinical outcomes)



Classifier

Feature extraction by Autoencoders

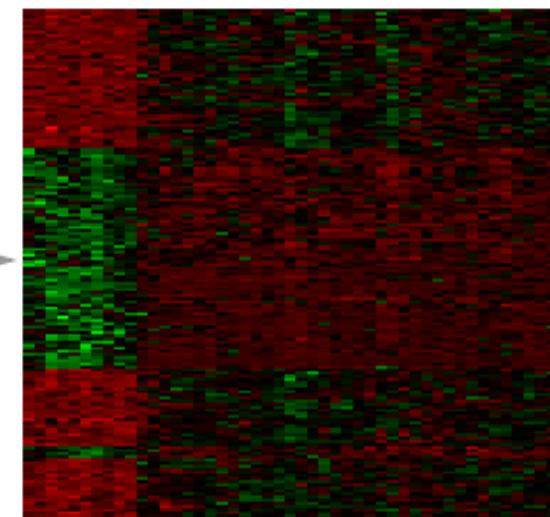
Original input



Extracted features

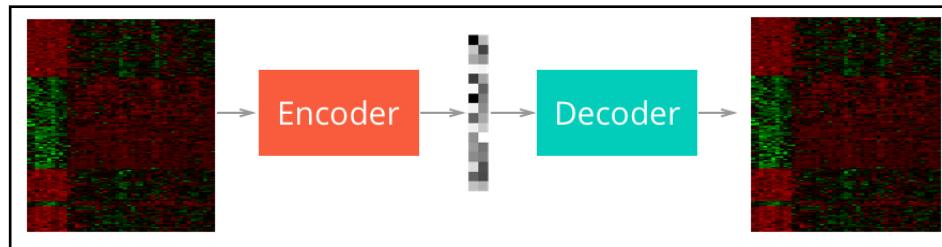


Reconstructed input

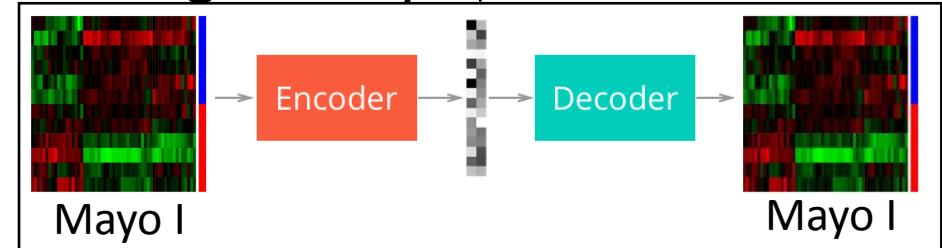


Deep Genomic Signature (DGS)

Feature extraction on the poorly-annotated cohort

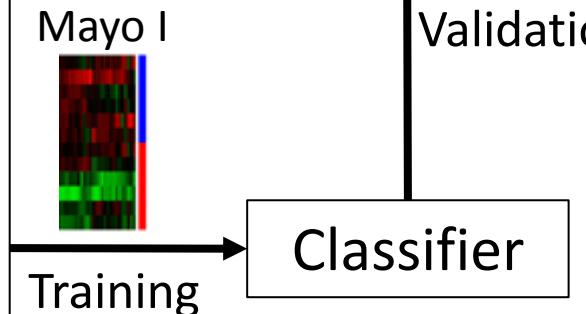


Refining on Mayo



Finding important genes
(Gene selection)

GREB1 NCAM2 INPP4B
RP11-33A14.1 MYL12B
FAM13C RLN FAM19A5
ANTXR2 HS6ST1 ITGA1 SYT7
ABHD2 MYBPC1 ODAM
NRAP MYO6 MEIS2 IRS4 SCIN
LCE3B ERG KCNMB1 KCNN4
CACNA1D C1QTNF9B-AS1
C1QTNF9B CALD1 PLS3
AZGP1 ENDOD1 CANX
A2M PGM5P4-AS1
GSTM2 TENM1



HIT'nDRIVE: Multi-driver Gene Prioritization Based on Hitting Time

Raunak Shrestha^{1,2,*}, Ermin Hodzic^{3,*}, Jake Yeung^{2,4,*}, Kendric Wang²,
Thomas Sauerwald⁵, Phuong Dao⁶, Shawn Anderson², Himisha Beltran⁷,
Mark A. Rubin⁷, Colin C. Collins^{2,8}, Gholamreza Haffari⁹, and S. Cenk Sahinalp^{3,10}

Shrestha et al. 2014 (RECOMB-2014)

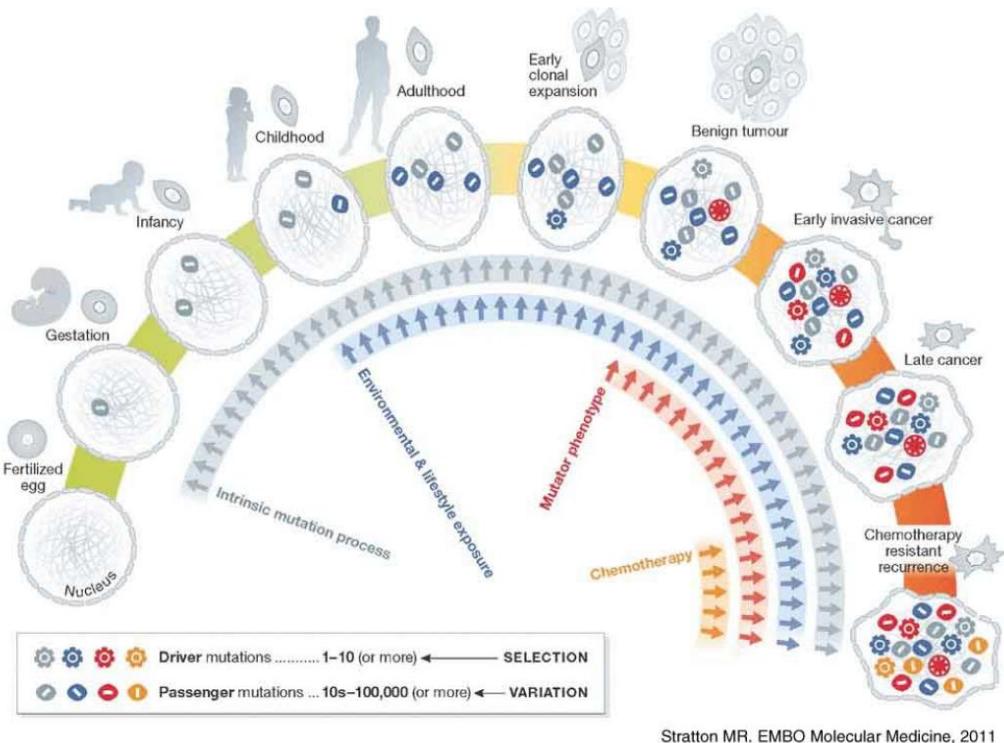
Method

HIT'nDRIVE: patient-specific multidriver gene prioritization for precision oncology

Raunak Shrestha,^{1,2,10} Ermin Hodzic,^{3,10} Thomas Sauerwald,⁴ Phuong Dao,⁵
Kendric Wang,² Jake Yeung,² Shawn Anderson,² Fabio Vandin,⁶ Gholamreza Haffari,⁷
Colin C. Collins,^{2,8} and S. Cenk Sahinalp^{2,3,9}

Shrestha et al. 2017 (Genome Research)

Genes that Drive Cancer



- Cancer is mediated by somatic evolution of various alterations in genome

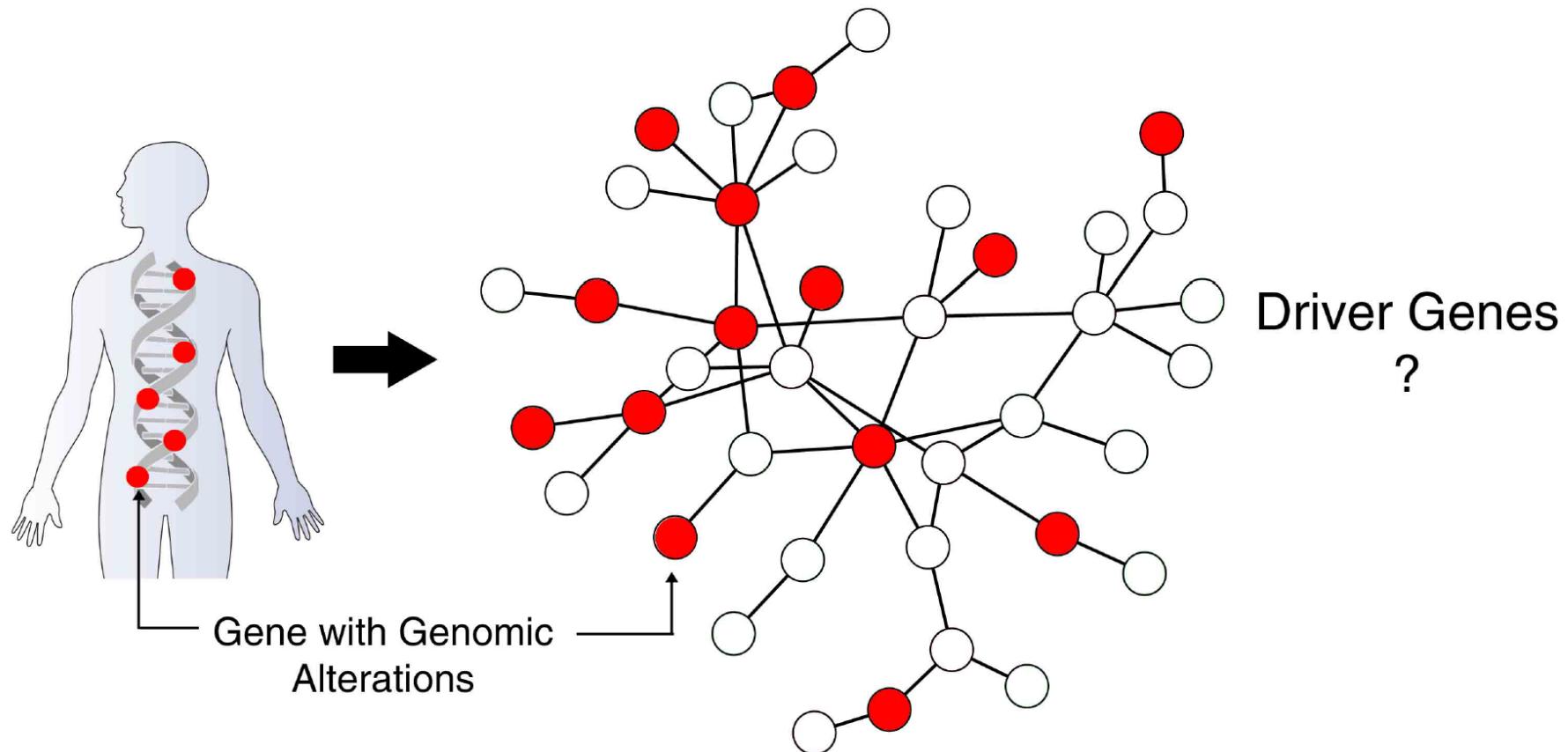
- Most of the alterations are neutral and provide no growth advantage to the tumor: known as “**passenger**” alterations
- Very few alterations provide evolutionary advantage and are positively selected for during tumorigenesis: known as “**driver**” alterations

Challenge

- Driver alterations are diluted and are outnumbered by the passenger alterations
- This makes identification of driver alterations more complicated

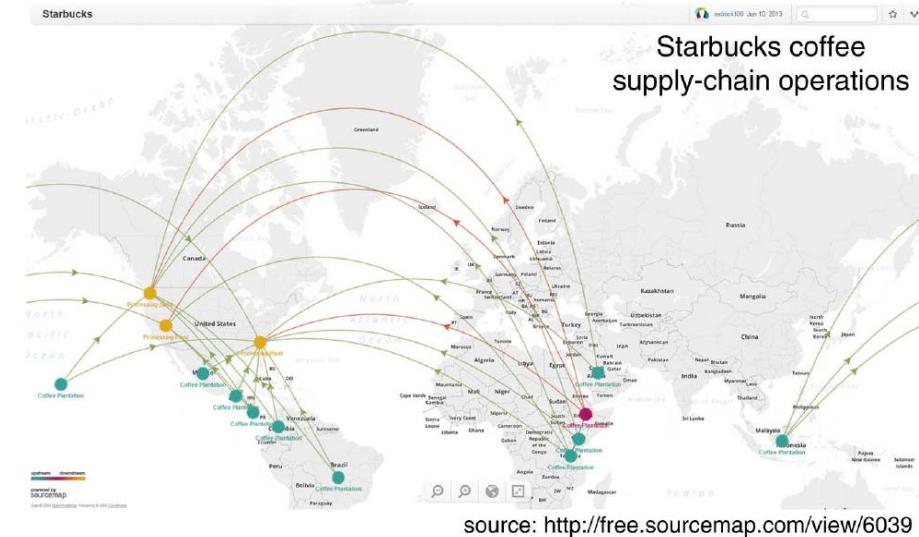
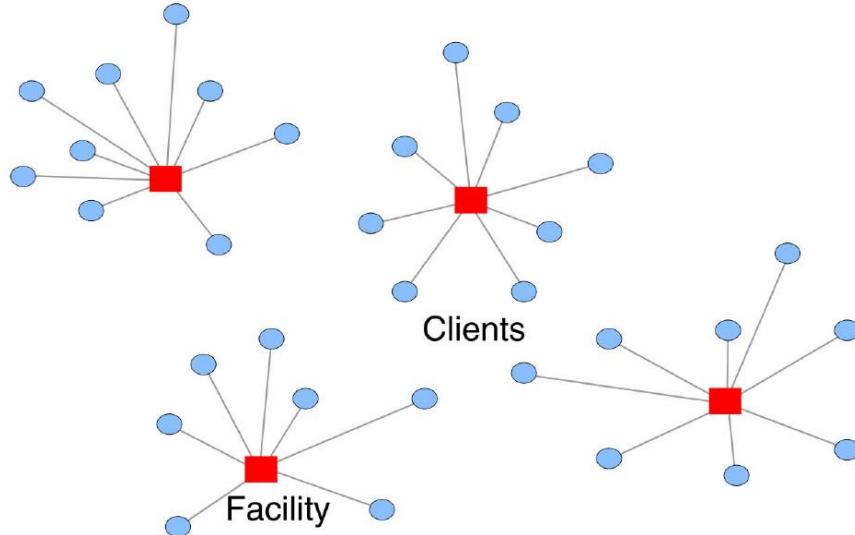
Problem Formulation:

- Given a set of genes with genomic alterations in a tumor, identify the location of the gene(s) in the network that are (most likely) responsible for driving the tumor



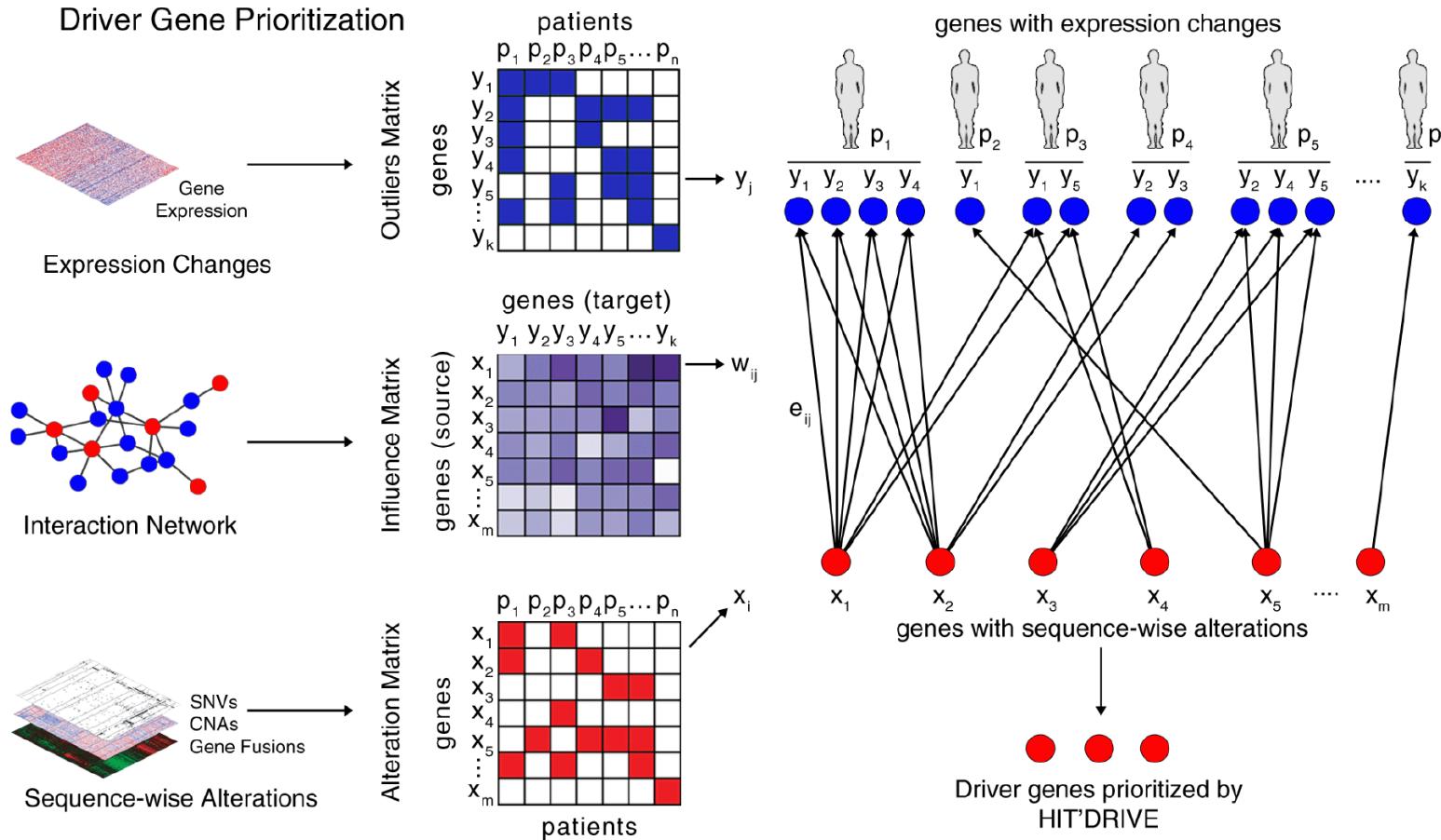
Facility location Problem

- Concerned with the **optimal placement of facilities** to **minimize associated costs**



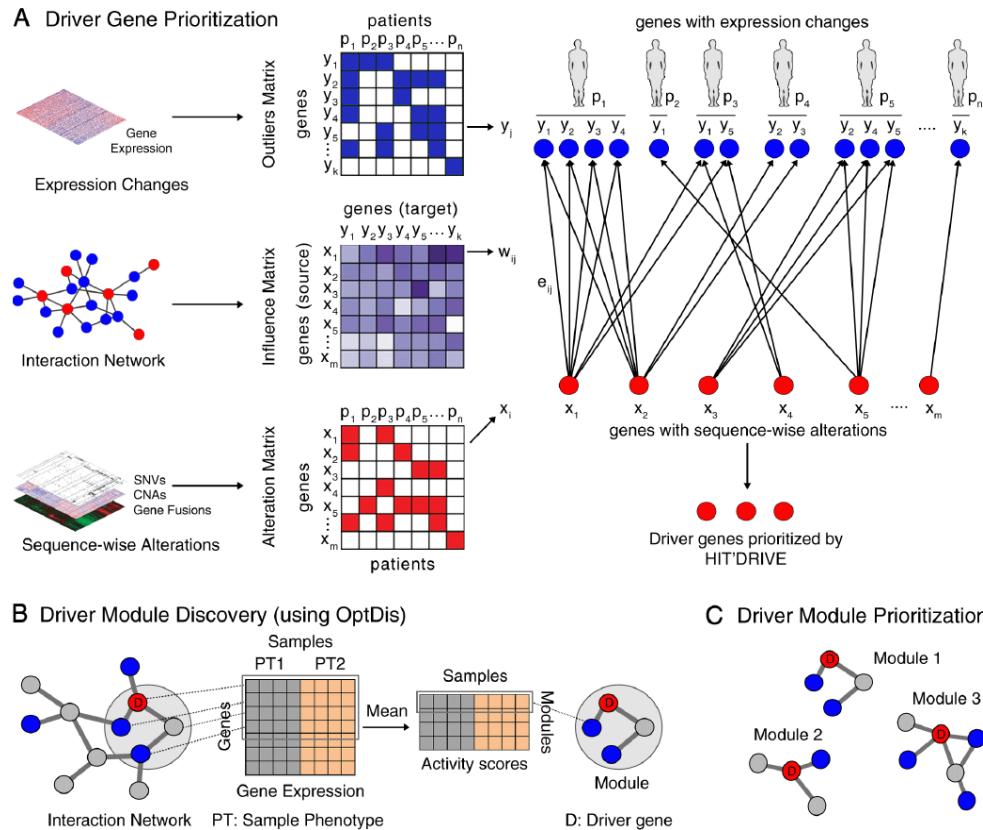
- This formulation is used to decide
 - Hospital placement
 - Airline hub airport selection
 - Chain restaurant/coffee shop location identification
 - Identification of top influencers in a social network: fake news placement

Find the most influential set of genes that potentially driver the cancer



- Find the smallest set of genes with sequence level changes that can explain the changes in gene expression of large fraction of genes
- We solve this using Integer Linear Programming (ILP)

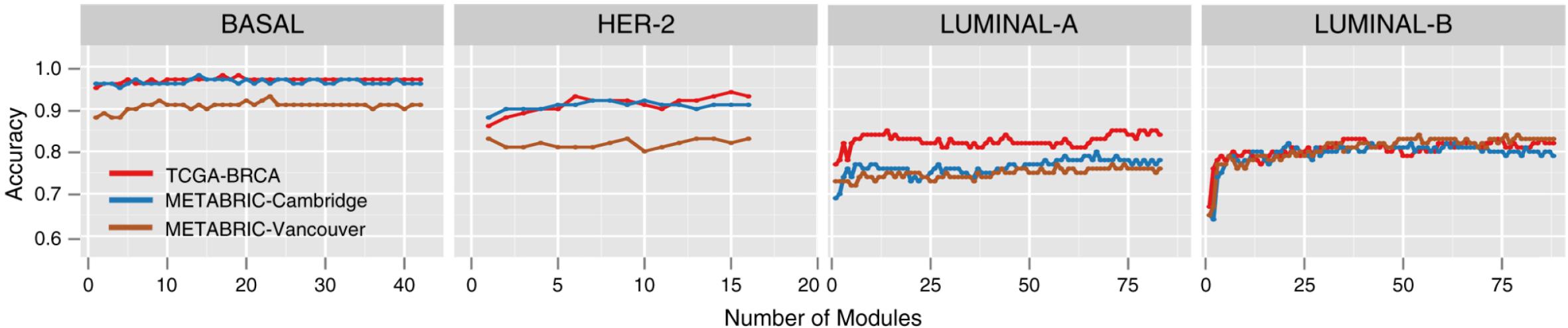
Cancer sub-type classification



- We used (and improved) OptDis¹ for *de novo* identification of modules inside the interaction network which are seeded by at least one predicted driver
- modules are chosen so that their discriminative power (for phenotype classification) is greatest among connected subnetworks of similar size that contain the individual predicted drivers.

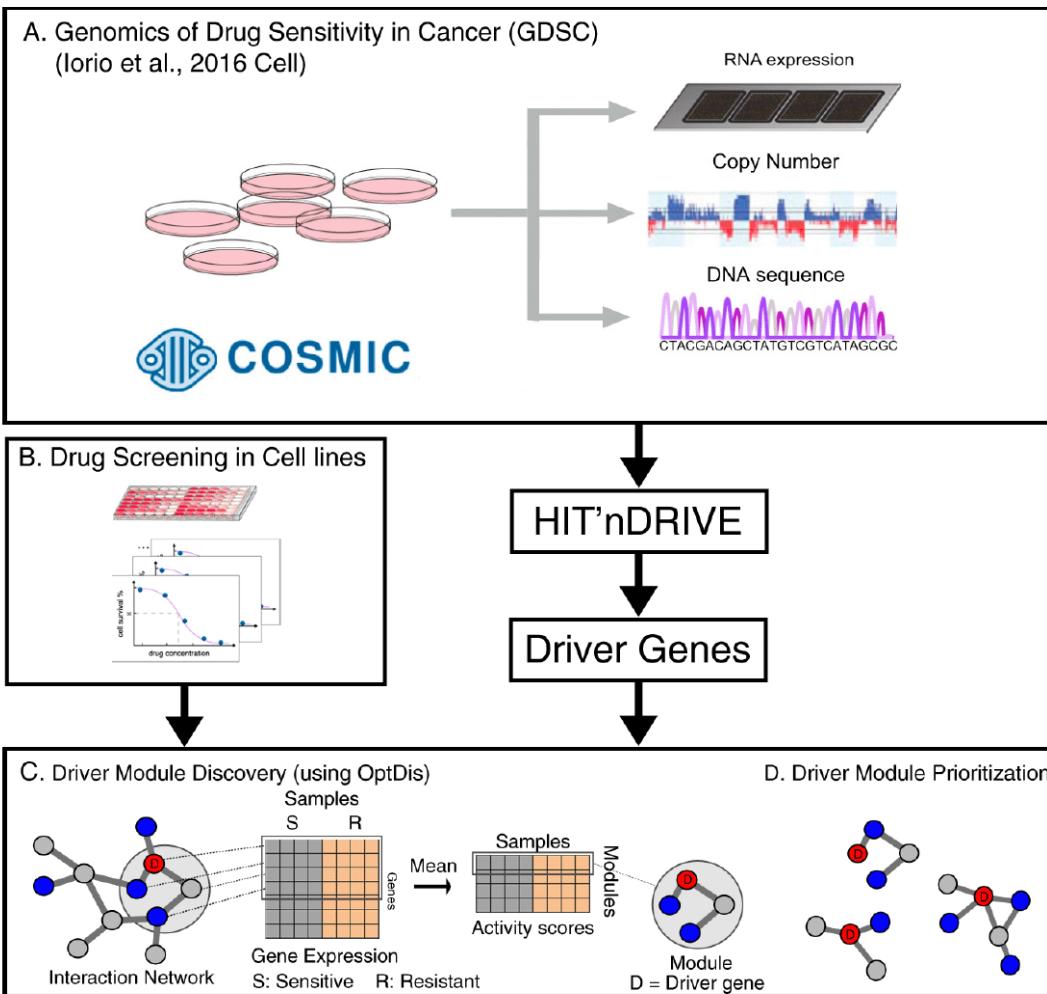
¹ Dao et al. Bioinformatics, 2011 (ISMB 2011)

Breast Cancer Subtype classification



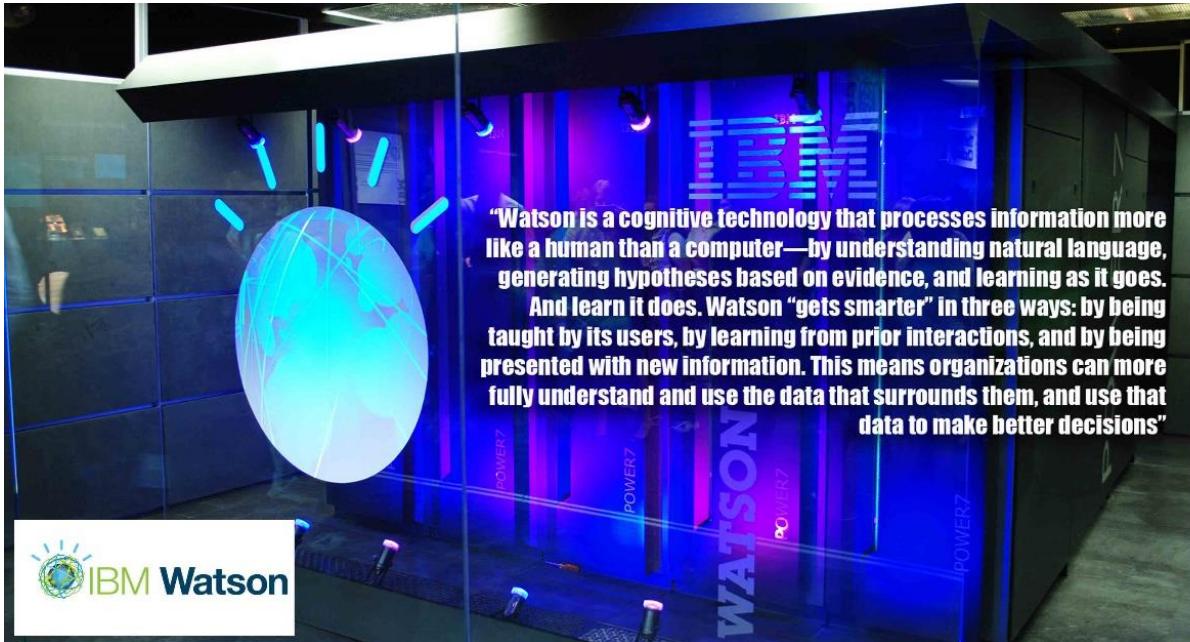
- Classified four major subtypes of breast cancer - Basal, HER2, Luminal-A and Luminal-B
- We respectively obtained 37, 16, 43 and 39 subtype specific driver modules for Basal, HER2, Luminal-A and Luminal-B subtypes
- Classified Basal-like tumors with much higher accuracy (98%) as compared to other BRCA-subtypes - HER2 (94%), Luminal-A (85%) and Luminal-B (83%)

Predict Drug Efficacy



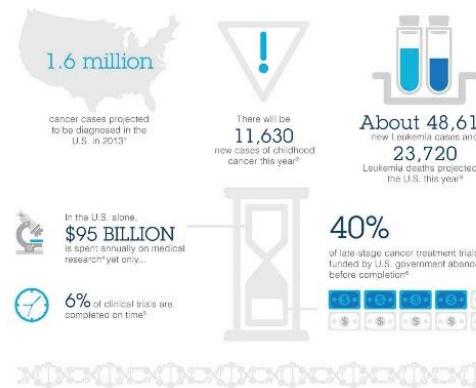
- Pan-Cancer Cell Lines (1001 cell lines, 30 cancer types)
- Treatment with 265 drugs in Pan-Cancer Cell Lines
 - Targeted & cytotoxic drugs
 - FDA approved drugs, drugs in clinical trials & pre-clinical ligands

AI Expert systems are beginning to be tested in cancer research centres across the world



MD Anderson Taps IBM Watson for Mission to End Cancer

Going Up Against a Deadly Disease



ARTIFICIAL INTELLIGENCE: Jeopardy champion a cancer expert as IBM supercomputer, Watson, trades trivia for treatment advice



Top News In the Lab Health Technology Pharma Academia Features

Videos Asian Scientist 100 Bugs & Quarks Singapore's Scientific Pioneers Print Magazine Intelligence Obituaries

IBM Watson To Fight Cancer In 21 Hospitals Across China

IBM and Hangzhou CognitiveCare plan to bring the Watson cognitive computing platform to 21 hospitals across China.

ADVERTISEMENT

EDITOR'S

South Korea Meets Hyperloop Train



Did IBM overhype Watson Health's AI promise?

IBM's Watson Health division has been under fire for not delivering on its promise to use AI to enable smarter, more personalized medicine. But IBM officials maintain that hospitals are seeing benefits.

SCIENCE \ TECH \ HEALTH \

IBM's Watson gave unsafe recommendations for treating cancer

Doctors fed it hypothetical scenarios, not real patient data

By Angela Chen | @chengela | Jul 26, 2018, 4:29pm EDT

Emergent Tech ▶ Artificial Intelligence

IBM's Watson Health wing left looking poorly after 'massive' layoffs

Up to 70% of staff shown the door this week, insiders claim

By Iain Thomson in San Francisco 25 May 2018 at 20:08

40 SHARE ▼

IBM has laid off approximately 50 and 70 per cent of staff this week in its Watson Health division, according to inside sources.

Take away message

- Machine Learning is a great approach to tackle a number of problems of biology
- Rise in biological (or genomic) data will certainly require machine learning models to hunt for hidden patterns of disease
- **However, all problems in biology CANNOT be solved using machine learning**
- Try to explore other approaches to model a problem beyond machine learning!