# Translational Medicine in the Era of Big Data: Hype or Real?

AAHCI MENA Regional Conference
September 27, 2018

**AKL FAHED, MD, MPH**     @aklfahed

MASSACHUSETTS GENERAL HOSPITAL
CORRIGAN MINEHAN HEART CENTER

CENTER FOR GENOMIC MEDICINE

HARVARD MEDICAL SCHOOL TEACHING HOSPITAL

BROAD INSTITUTE
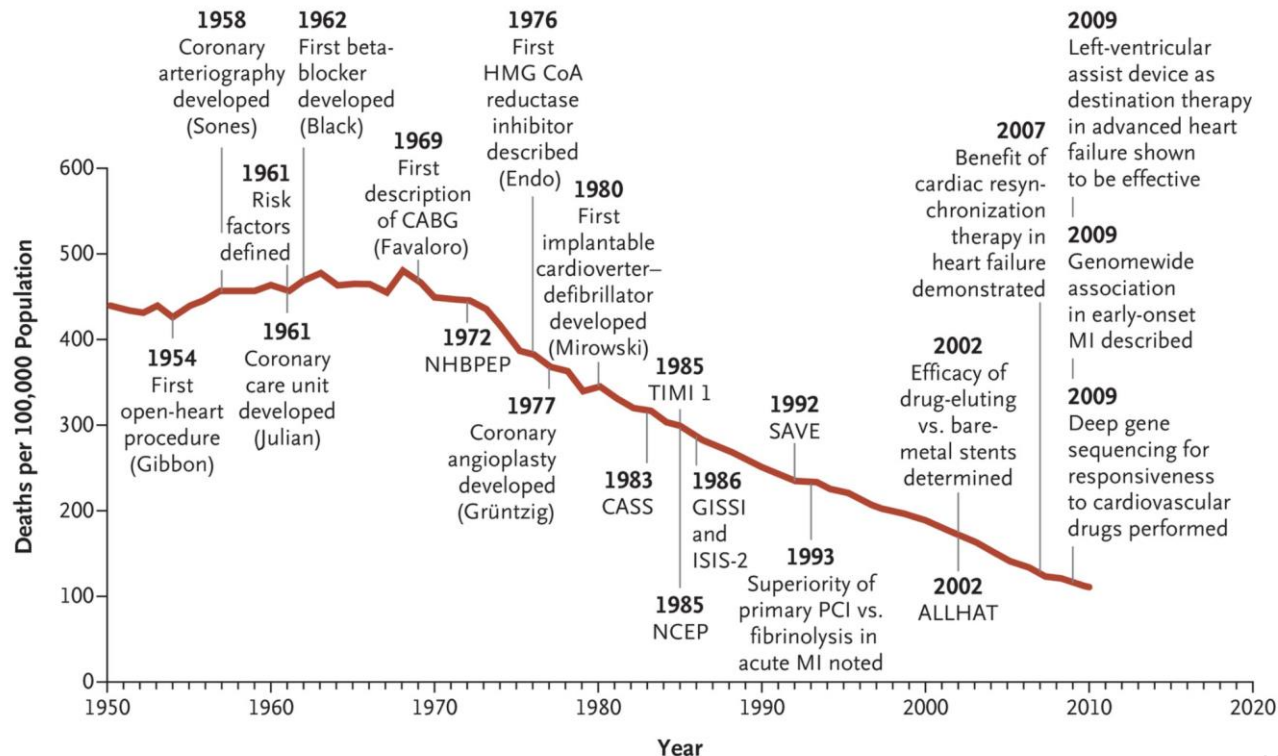
Massachusetts Institute of Technology

# Disclosures

*None*

# Outline

- The Promise of Big Data

- **Genomics**
  - Polygenic Risk Scores
  - Mendelian Randomization
  - Human Knockout Project
  - Phenome-Wide Association Studies

- Challenges and Pitfalls

- Opportunity for Academic Health Centers

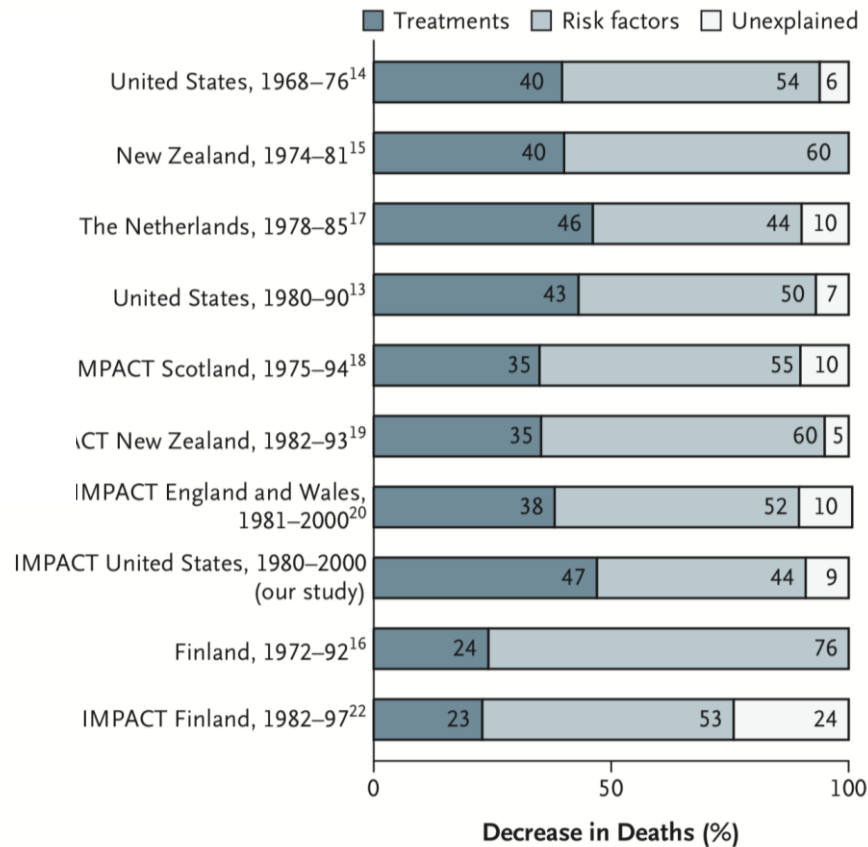Weintraub WS, Fahed AC, Rumsfeld JS. *Circulation Research. In Press*.

# Decline in Cardiovascular Deaths

Nabel E and Braunwald E. *NEJM* 2012

# Translational Medicine

Bench to Bedside to Population

# Yet …

*Even highly efficacious therapies have heterogeneity of effect at the individual level*

*Significant variation in the use of evidence-based therapies and outcomes in routine clinical practice*

*Drug development is a very lengthy process*

*…*

*…*

*…*

# The Promise of Big Data

*Precision Medicine*

*Artificial Intelligence*
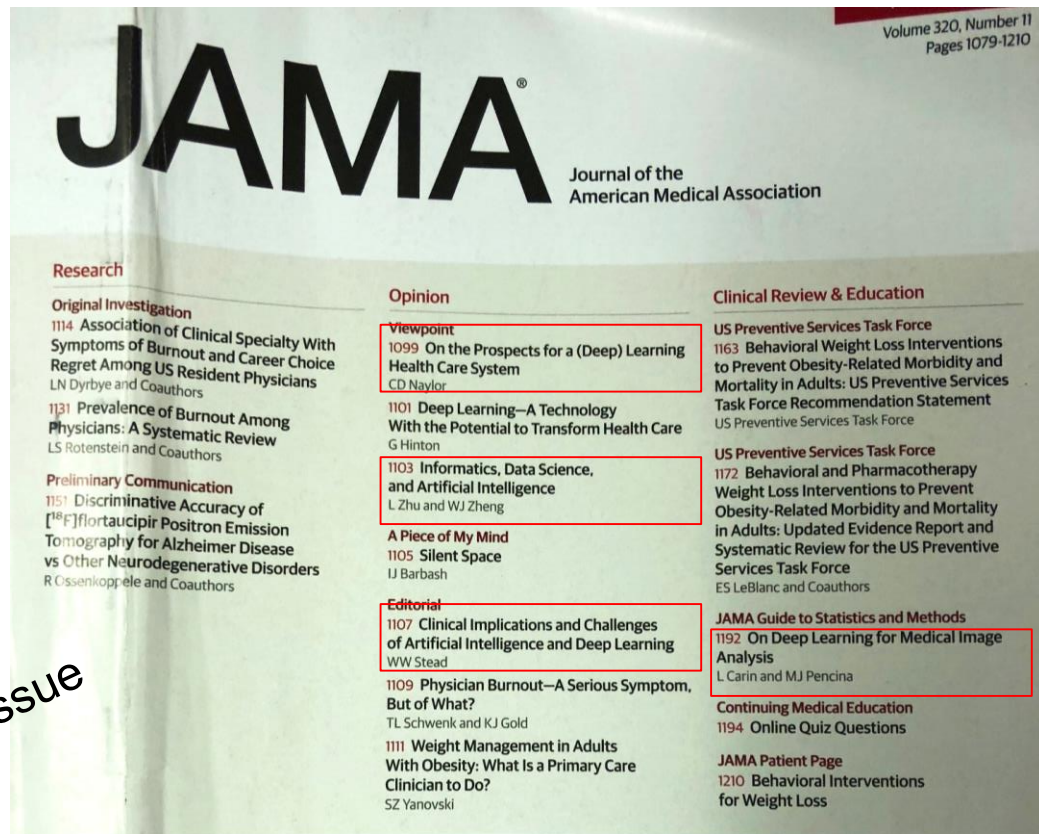
*Improved Translational Medicine*

*…*

*…*

**Hype or Real?**

*…*

*…*



Current issue

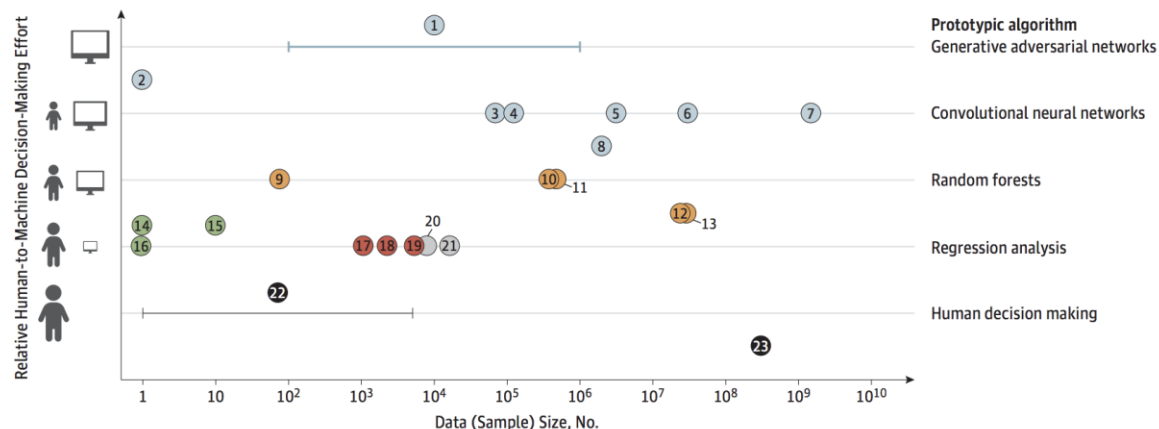# Sources of Big Data in Healthcare

- Electronic Health Records (EHRs)

- Wearables, Apps and Biosensors (IoTs)

    **_Zettabyte levels ($10^{21}$)_**

- **Genomic data**

- Insurance providers (claims, pharmacies, etc)

- Other clinical data (decision support tools, administrative data, etc)

- Social Media

- Web of knowledge

Lima FV and Fahed AC. Harnessing the Power of Big Data in Healthcare. *Cardiology Magazine* 2018.

# Spectrum of Big Data & Machine Learning

Figure. The Axes of Machine Learning and Big Data

**Deep learning**
1. Generative adversarial networks (2014)
2. Google AlphaGo Zero (2017)
3. ATM check readers (1998)
4. Google diabetic retinopathy (2016)
5. ImageNet computer vision models (2012-2017)
6. Google AlphaGo (2015)
7. Facebook Photo Tagger (2015)
8. Prediction of 1-y all-cause mortality (2017)

**Classic machine learning**
9. Diffuse large B-cell lymphoma outcome prediction by gene-expression profiling (2002)
10. EHR-based CV risk prediction (2017)
11. Netflix Prize winner (2006)
12. Google Search (1998)
13. Amazon product recommendation (2003)

**Expert AI systems**
14. MYCIN (1975)
15. CASNET (1982)
16. DXplain (1986)

**Risk calculators**
17. $CHA_2DS_2$-VASc Score for atrial fibrillation stroke risk (2017)
18. MELD end-stage liver disease risk score (2001)
19. Framingham CV risk score (1998)

**Randomized Clinical Trials**
20. Celecoxib vs nonsteroidal anti-inflammatory drugs for osteoarthritis and rheumatoid arthritis (2002)
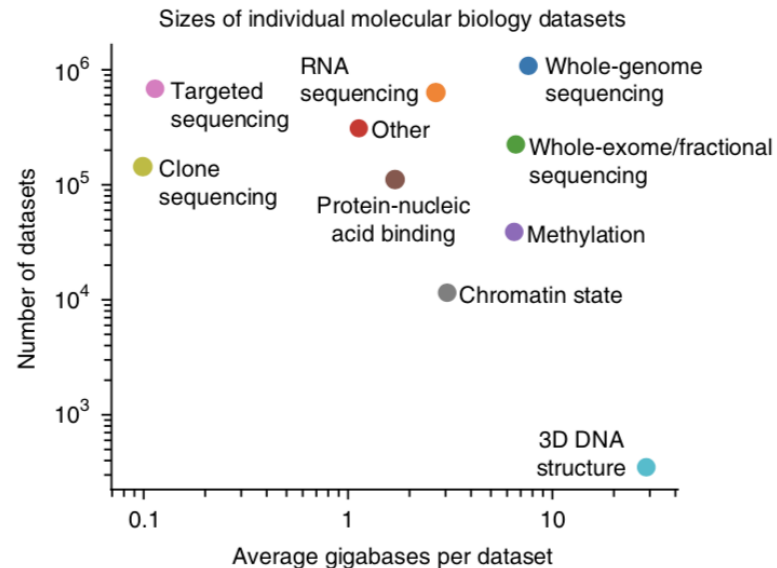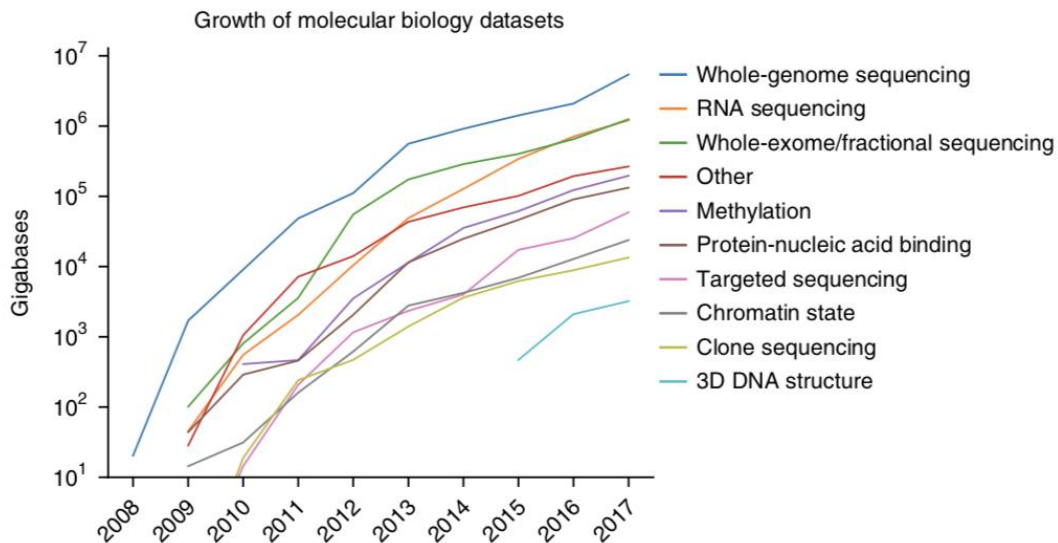21. Use of estrogen plus progestin in healthy postmenopausal women (2002)

**Other**
22. Clinical wisdom
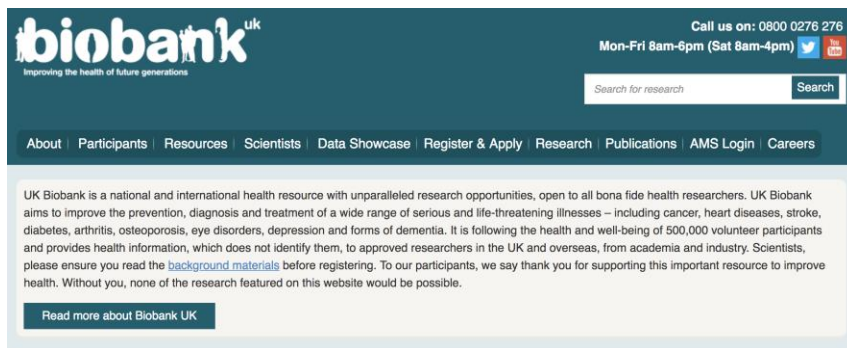23. Mortality rate estimates from US Census (2010)

Beam A and Kohane I. JAMA 2018

10

# Types of Genomic Data

| *Whole-Genome Genotyping* | *Whole-Exome Sequencing* | *Whole-Genome Sequencing* |
|---|---|---|
| Array with 100,000- 1 millions SNPs<br>Imputation: >90 million SNPs | Coding part (1%) of the genome<br>All exons of all genes | Entirety of the genome |
| Common variation (allele frequency >1%) | Rare coding variation | Rare and Common disease<br>Noncoding rare variation |
| GWAS, Mendelian Randomization, Polygenic Risk Scores | Rare disease diagnosis, discovery of novel rare loss of function | Role of noncoding DNA |
| ~ 50 USD | ~ 400 USD | ~ 1500 USD |
| Public data ++++++ | Public data +++++++ | Public data emerging |

# Growth and Size of Molecular Data

Wainberg et al. *Nature Biotechnology*. 2018

# The Rise of the Biobanks

**UK Biobank 500,0000**



**USA 1,000,000    USA 1,000,000**



| Biobank | Enrollment locations | Initial enrollment | Enrollment to date | Target enrollment |
|---|---|---|---|---|
| **Commercial funding** | | | | |
| deCODE Genetics (Amgen) (http://www.decode.com/) | Iceland | 1996 | >200,000 | Unknown |
| Geisinger MyCode® Community Health (Regeneron Pharmaceuticals and Others) | Geisinger Health System (Danville, PA) | 2007 | >50,000 | Unknown |
| **Government funding** | | | | |
| China Kadoorie Biobank (http://www.ckbiobank.org/site/) | China | 2004 | >500,000 | Enrollment Completed |
| UK Biobank (https://www.ukbiobank.ac.uk/) | United Kingdom | 2006 | >500,000 | Enrollment Completed |
| Electronic Medical Records and Genomics (eMERGE) Network (https://emerge.mc.vanderbilt.edu/about-emerge/) | United States Hospital Sites | 2007 | >50,000 | Unknown |
| Million Veterans Program (http://www.research.va.gov/mvp/) | Veterans Affairs Hospital | 2011 | >500,000 | ~1,000,000 |
| Precision Medicine Initiative (https://www.nih.gov/precision-medicine-initiative-cohort-program) | United States | Early 2017 | -- | ~1,000,000 |
| **Institutional funding** | | | | |
| BioVu Biorepository (https://victr.vanderbilt.edu/pub/biovu/) | Vanderbilt University Medical Center (Nashville, TN) | 2007 | >215,000 | Unknown |
| Kaiser Permanente Research Bank (http://researchbank.kaiserpermanente.org/) | United States | 2016 | >250,000 | ~500,000 |
| Partners Healthcare Biobank (https://biobank.partners.org/) | Partners Health Care (Boston, MA) | 2010 | >50,000 | ~100,000 |

Khera AV and Kathiresan S. *Nature Reviews Genetics*. 2017

13

# Democratization of Genomic Data





**Ben Neale**

HOME  RESEARCH  PEOPLE  MEDIA  BLOG  UK BIOBANK  JOBS  CONTACT

## RAPID GWAS OF THOUSANDS OF PHENOTYPES FOR 337,000 SAMPLES IN THE UK BIOBANK

September 20, 2017

**biobank** uk

The UK Biobank recently released genome-wide association data on ~500,000 individuals. The genotype data for these samples have been cleaned, imputed and released to the scientific community. This public release of data represents an extraordinary advance for genetics, pushing the envelope for data sharing and rapid uptake by the research community. These data will be used for novel discovery of disease-associated genes, in the development of new methods, and to serve as an example for how future efforts in genetics and biology ought to proceed.

To further enhance the value of this resource, we have performed a basic association test on ~337,000 unrelated individuals of British ancestry for over 2,000 of the available phenotypes. We're making these results available for browsing through several portals, including the Global Biobank Engine where they will appear soon. They are also available for download here.

# UKBB GWAS bot

# Polygenic Risk Scores (PRS)

*Weighted sum of number of risk alleles carried by an individual*

- Sum of the risk alleles (X)

$$Y = \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 \ldots.$$

- Measured effects as detected by GWAS (β)

# CAD Polygenic Risk Score
# LDpred method(>6 million alleles)

Khera AV et al. *Nature Genetics* 2018

Khera AV et al. *Nature Genetics* 2018

**Table 3 | Prevalence and clinical impact of a high GPS**

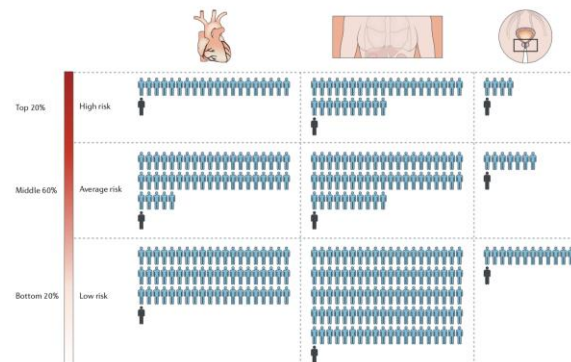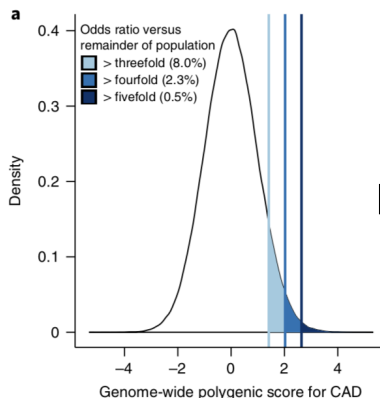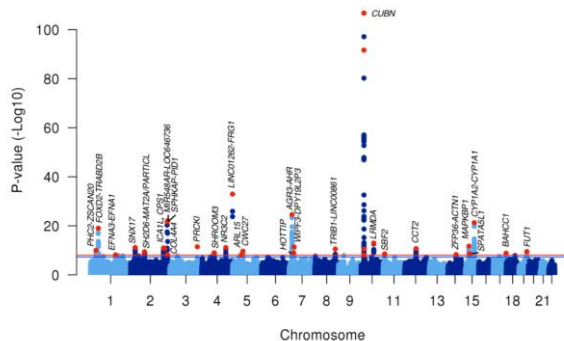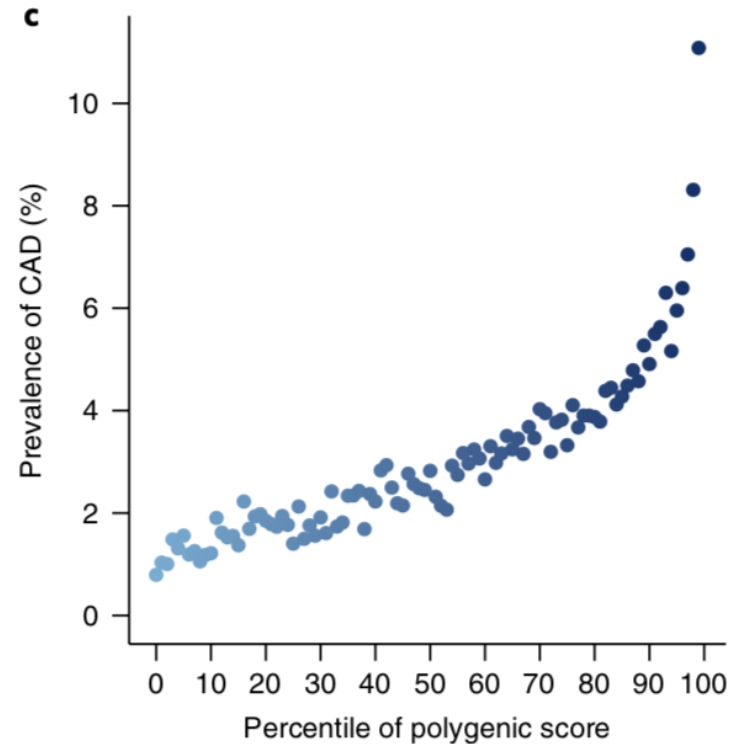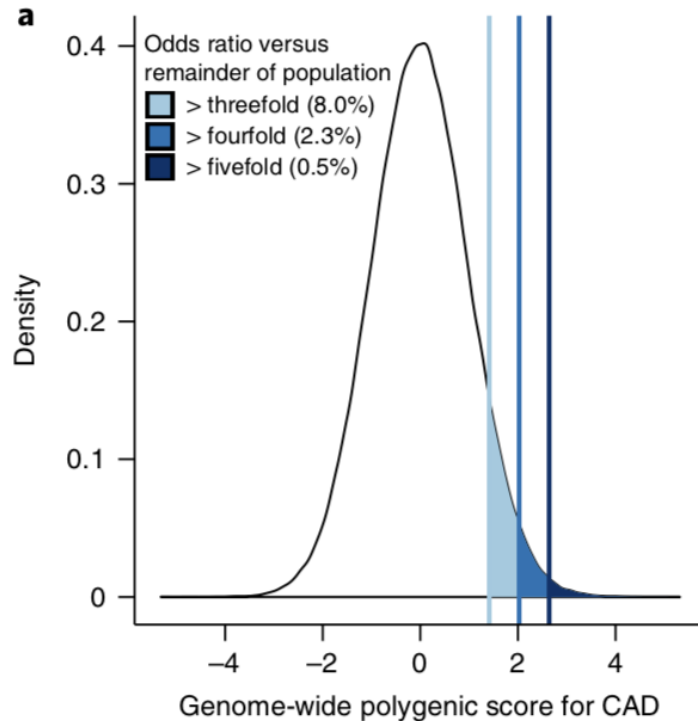| High GPS definition | Reference group | Odds ratio | 95% CI | P value |
|---|---|---|---|---|
| **CAD** | | | | |
| Top 20% of distribution | Remaining 80% | 2.55 | 2.43–2.67 | $<1\times10^{-300}$ |
| Top 10% of distribution | Remaining 90% | 2.89 | 2.74–3.05 | $<1\times10^{-300}$ |
| Top 5% of distribution | Remaining 95% | 3.34 | 3.12–3.58 | $6.5\times10^{-264}$ |
| Top 1% of distribution | Remaining 99% | 4.83 | 4.25–5.46 | $1.0\times10^{-132}$ |
| Top 0.5% of distribution | Remaining 99.5% | 5.17 | 4.34–6.12 | $7.9\times10^{-78}$ |
| **Atrial fibrillation** | | | | |
| Top 20% of distribution | Remaining 80% | 2.43 | 2.29–2.59 | $2.1\times10^{-177}$ |
| Top 10% of distribution | Remaining 90% | 2.74 | 2.55–2.94 | $7.0\times10^{-169}$ |
| Top 5% of distribution | Remaining 95% | 3.22 | 2.95–3.51 | $1.1\times10^{-152}$ |
| Top 1% of distribution | Remaining 99% | 4.63 | 3.96–5.39 | $2.9\times10^{-84}$ |
| Top 0.5% of distribution | Remaining 99.5% | 5.23 | 4.24–6.39 | $3.5\times10^{-56}$ |
| **Type 2 diabetes** | | | | |
| Top 20% of distribution | Remaining 80% | 2.33 | 2.20–2.46 | $3.1\times10^{-201}$ |
| Top 10% of distribution | Remaining 90% | 2.49 | 2.34–2.66 | $1.2\times10^{-167}$ |
| Top 5% of distribution | Remaining 95% | 2.75 | 2.53–2.98 | $1.7\times10^{-130}$ |
| Top 1% of distribution | Remaining 99% | 3.30 | 2.81–3.85 | $1.4\times10^{-49}$ |
| Top 0.5% of distribution | Remaining 99.5% | 3.48 | 2.79–4.29 | $4.3\times10^{-30}$ |
| **Inflammatory bowel disease** | | | | |
| Top 20% of distribution | Remaining 80% | 2.19 | 2.03–2.36 | $7.7\times10^{-95}$ |
| Top 10% of distribution | Remaining 90% | 2.43 | 2.22–2.65 | $8.8\times10^{-88}$ |
| Top 5% of distribution | Remaining 95% | 2.66 | 2.38–2.96 | $3.0\times10^{-68}$ |
| Top 1% of distribution | Remaining 99% | 3.87 | 3.18–4.66 | $1.4\times10^{-43}$ |
| Top 0.5% of distribution | Remaining 99.5% | 4.81 | 3.74–6.08 | $9.0\times10^{-37}$ |
| **Breast cancer** | | | | |
| Top 20% of distribution | Remaining 80% | 2.07 | 1.97–2.19 | $3.4\times10^{-159}$ |
| Top 10% of distribution | Remaining 90% | 2.32 | 2.18–2.48 | $2.3\times10^{-148}$ |
| Top 5% of distribution | Remaining 95% | 2.55 | 2.35–2.76 | $2.1\times10^{-112}$ |
| Top 1% of distribution | Remaining 99% | 3.36 | 2.88–3.91 | $1.3\times10^{-54}$ |
| Top 0.5% of distribution | Remaining 99.5% | 3.83 | 3.11–4.68 | $8.2\times10^{-38}$ |

Odds ratios were calculated by comparing those with high GPS with the remainder of the population in a logistic regression model adjusted for age, sex, genotyping array, and the first four principal components of ancestry. The breast cancer analysis was restricted to female participants. CI, confidence interval.
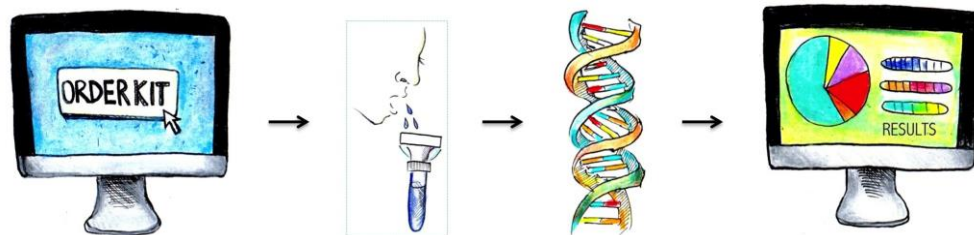
# Polygenic Risk Prediction

*20% of the study population are at ≥ threefold*

*increased risk for at least 1 of the 5 diseases studied !*
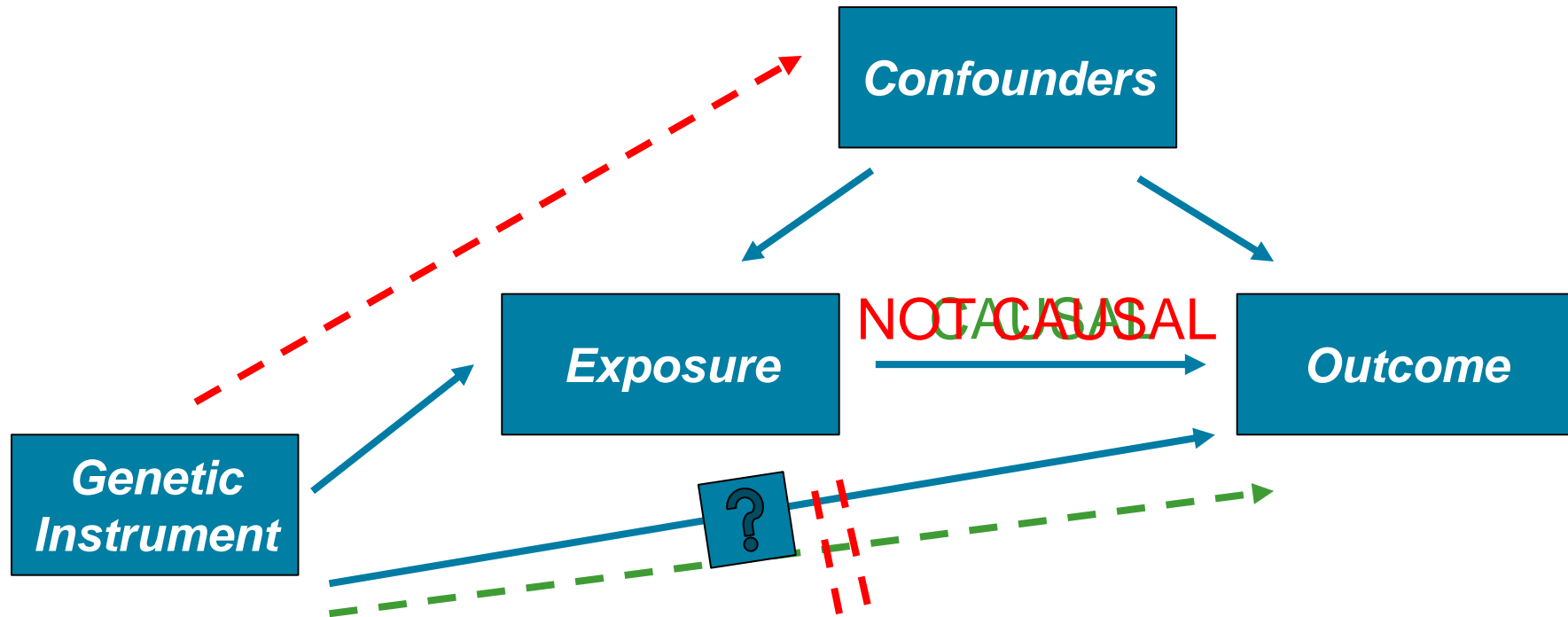
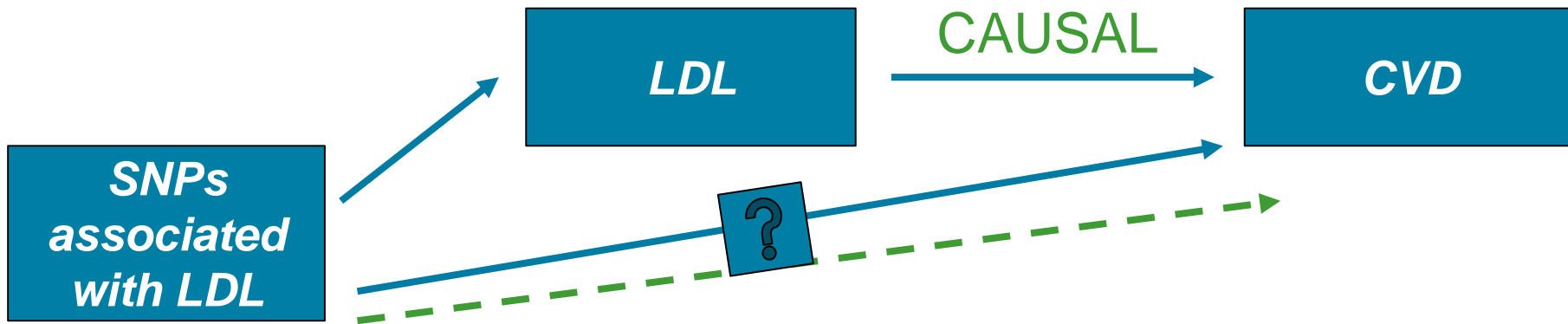*"The First" risk factor*

*~100 USD*

*Direct to Consumer Genetics*



Khera AV et al. *Nature Genetics* 2018
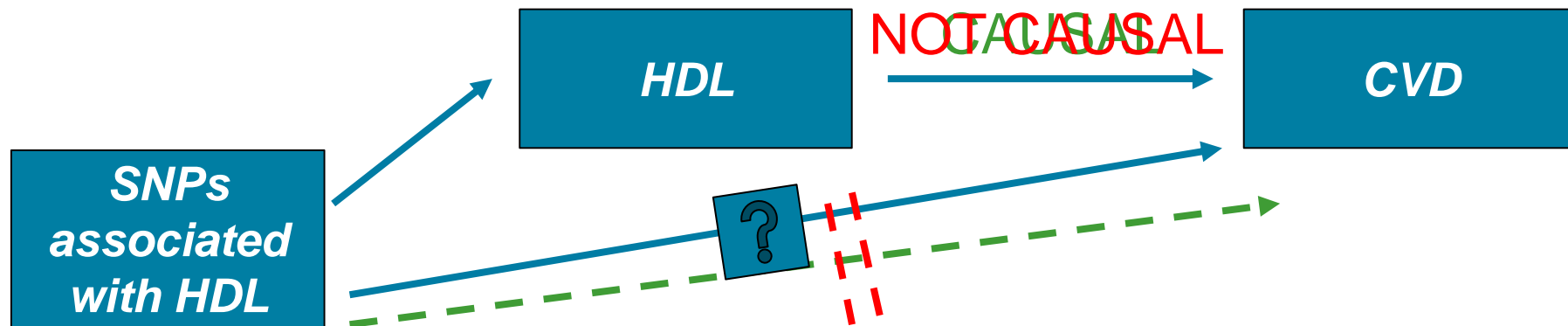https://pged.org/direct-to-consumer-genetic-testing/

19

# Mendelian Randomization

# Mendelian Randomization

# Human Knockout Project

## LETTER

doi:10.1038/nature22034

### Human knockouts and phenotypic analysis in a cohort with a high rate of consanguinity

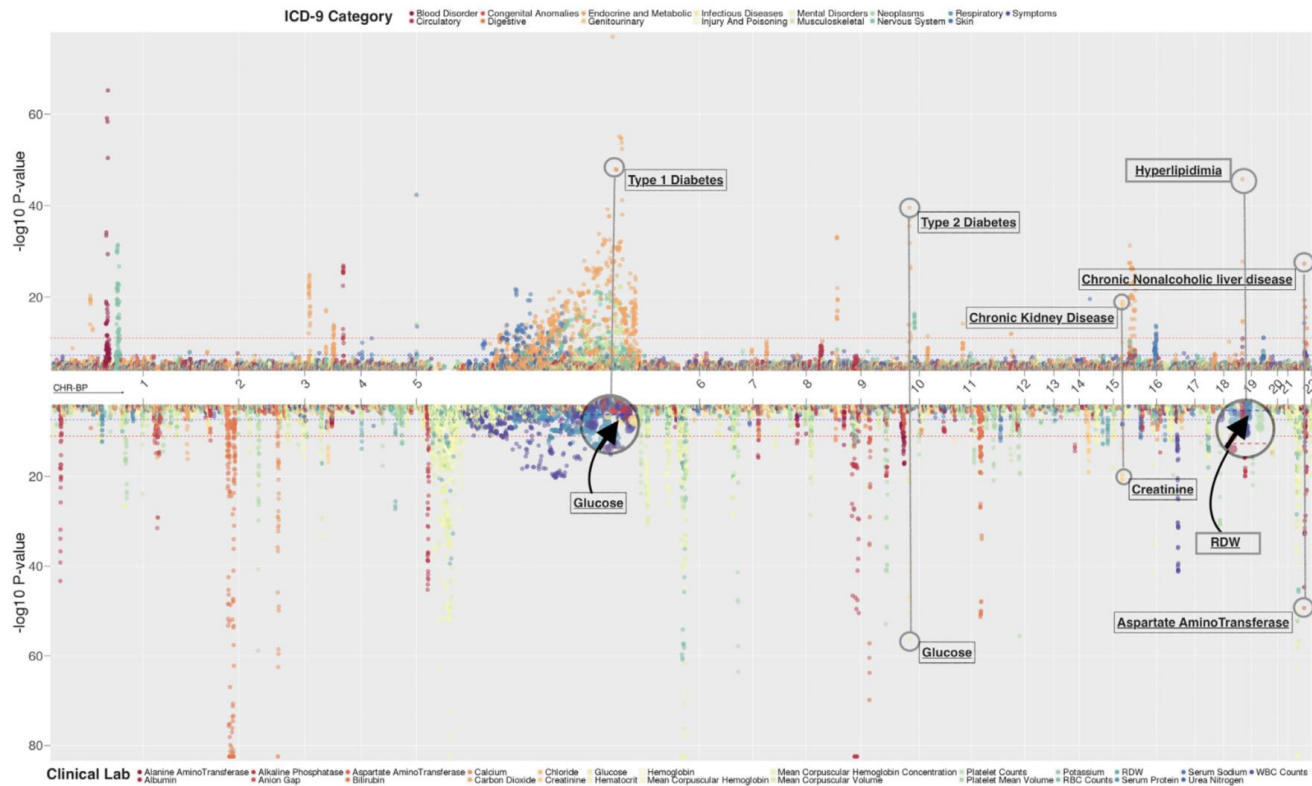Danish Saleheen[1,2]*, Pradeep Natarajan[3,4]*, Irina M. Armean[4,5], Wei Zhao[1], Asif Rasheed[2], Sumeet A. Khetarpal[6], Hong–Hee Won[7], Konrad J. Karczewski[4,5], Anne H. O'Donnell–Luria[4,5,8], Kaitlin E. Samocha[4,5], Benjamin Weisburd[4,5], Namrata Gupta[4], Mozzam Zaidi[2], Maria Samuel[2], Atif Imran[2], Shahid Abbas[9], Faisal Majeed[2], Madiha Ishaq[2], Saba Akhtar[2], Kevin Trindade[6], Megan Mucksavage[6], Nadeem Qamar[10], Khan Shah Zaman[10], Zia Yaqoob[10], Tahir Saghir[10], Syed Nadeem Hasan Rizvi[10], Anis Memon[10], Nadeem Hayyat Mallick[11], Mohammad Ishaq[12], Syed Zahed Rasheed[12], Fazal–ur–Rehman Memon[13], Khalid Mahmood[14], Naveeduddin Ahmed[15], Ron Do[16,17], Ronald M. Krauss[18], Daniel G. MacArthur[4,5], Stacey Gabriel[4], Eric S. Lander[4], Mark J. Daly[4,5], Philippe Frossard[2]§, John Danesh[19,20]§, Daniel J. Rader[6,21]§ & Sekar Kathiresan[3,4]§

**Safety check for drug development**

- Exome sequencing of 10,503 Pakistani subjects

- Identify individuals carrying predicted homozygous loss-of-function mutations

- Perform phenotypic analysis of >200 biochemical disease traits

- e.g. *APOC3* hom pLoF low fasting TG and blunted post-prandial lipaemia
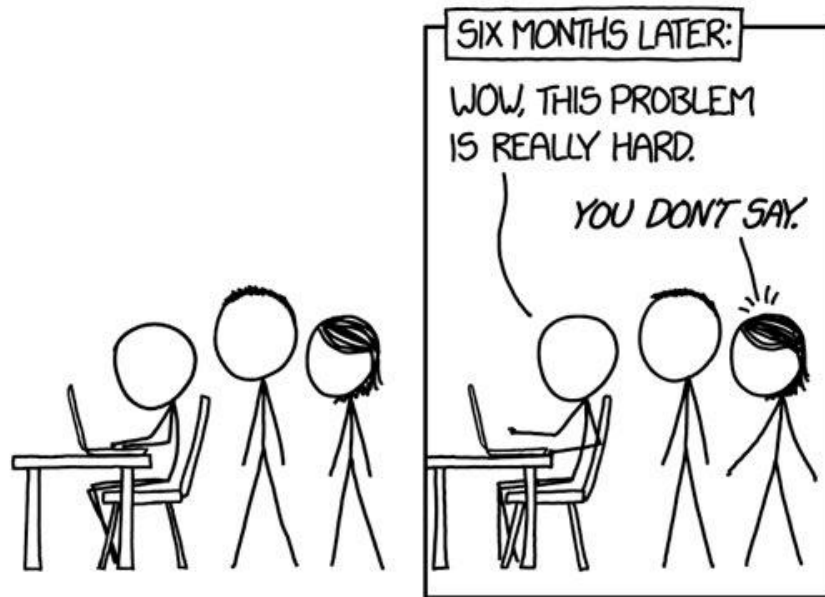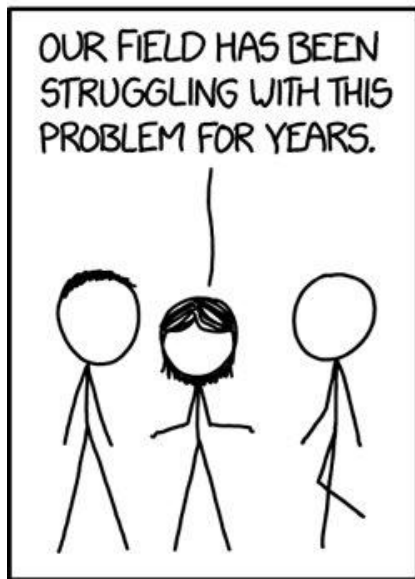
# Phenome Wide Association Studies (PheWAS)

*Association of SNPs with Medical Diagnoses and Clinical Measures in the EHR*



Verma A. et al. *AJHG 2018*

# Pitfalls of Big Data and ML

- Improved generation of hypotheses
  - But burden of proof remains on the basic scientist

- Polygenic risk implementation in care
  - Will it change outcomes?

- Biobanks phenotypic classification (case/control definitions)

- EHR/Administrative data has inherent biases of observational data
  - Informative missing data
  - Risk of false positives and negatives (i.e. misclassification)
  - Treatment selection bias i.e. unmeasured confounding variables

Source: Twitter @AndrewLBeam

# Doctors have a 'hunch' and it matters!

FULL SCREEN

A new study from MIT computer scientists suggests that human doctors provide a dimension that, as yet, artificial intelligence does not.
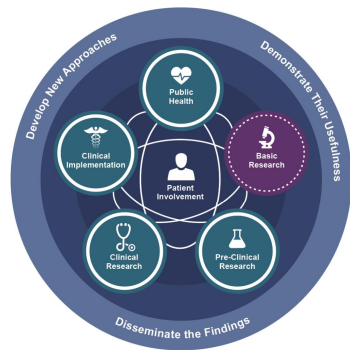
Image: Chelsea Turner, MIT

## Doctors rely on more than just data for medical decision making

Computer scientists find that physicians' "gut feelings" influence how many tests they order for patients.

Watch Video

Anne Trafton | MIT News Office
July 20, 2018

Press Inquiries          RELATED



Uhhh, yeah.

# Opportunity for Academic Health Centers

*The triple aim: care, health, and cost*

- Data Science as part of the framework of translational research

- Essential basic, translational and epidemiologic research for new technologies

- Unique partnerships with industry

- Products that are cost-effective, scientifically solid, and needed to advance patient care

# The new med school classroom?

- Computationally-Enabled Medicine

- "Pathways" curriculum

- Harvard Medical School 3rd year students

https://hms.harvard.edu/news/knowing-unknown

# Thank you

@aklfahed

fahed@mail.harvard.edu

MASSACHUSETTS GENERAL HOSPITAL
CORRIGAN MINEHAN HEART CENTER

CENTER FOR GENOMIC MEDICINE

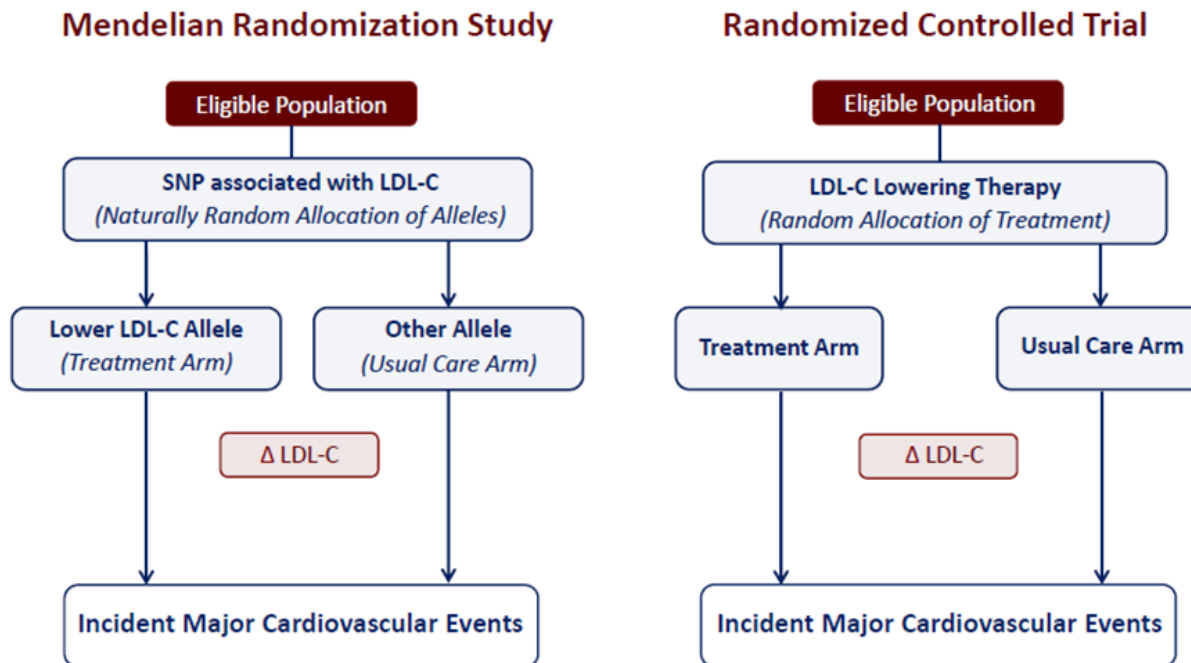HARVARD MEDICAL SCHOOL TEACHING HOSPITAL

BROAD INSTITUTE

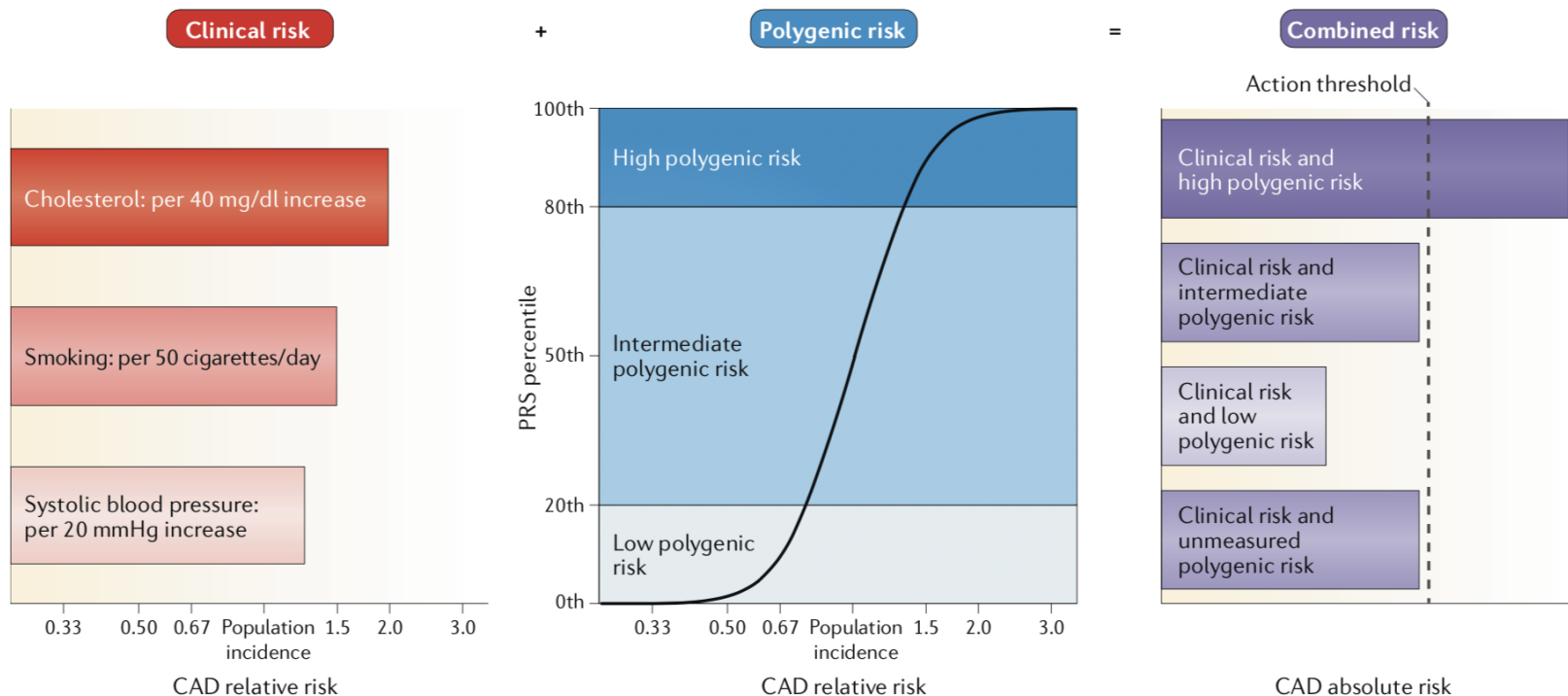MIT Massachusetts Institute of Technology
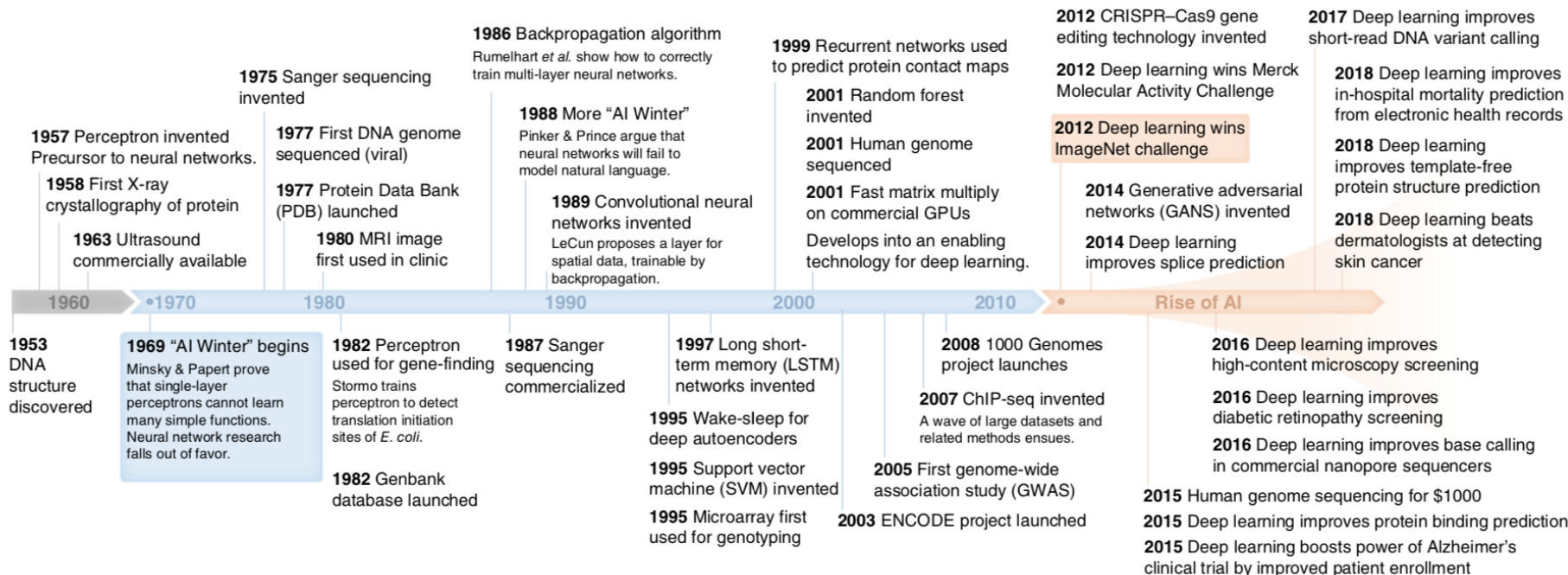
# Mendelian Randomization

**Figure:** Analogy Between a Mendelian Randomization Study and a Randomized Trial

Ferrence PA. ACC.org Expert Analysis

# Integrating Clinical and Polygenic Risk Prediction

Torkamani A et al. *Nature Reviews Genetics* 2018

# Timeline of Molecular Data



Wainberg et al. *Nature Biotechnology*. 2018

*" Machine Learning should try to do:*

*1- What doctors <u>cannot </u>do*

*2 What doctors <u>do NOT what to do</u> "*

# Not all Data are Created Equal

*Low Quality for ML*

- EHR

- Administrative Data

*Good Quality for ML*

- Image interpretation
  - CT
  - MRI
  - Echocardiography

- Detection of Dysrhythmias
  - Cardiac rhythm

- Wearables/Biosensors
  - HR/ Other physiological data

- Molecular data