# The data shows we need better data

Mélanie Courtot, PhD
Director, Genome Informatics and Principal Investigator
Assistant Professor, Department of Medical Biophysics, University of Toronto

mcourtot@oicr.on.ca

https://bit.ly/courtotlab

| Alexis Li | Hardeep Nahal-Bose | Pratham Hemlani |
| Business Analyst | Bioinformatician | Co-Op Student |
| Andres Melani | Henrich Feher | Ryan Seeto |
| PhD Student | DevOps | Data Scientist |
| Ann Catton | Jared Baker | Rakesh Mistry |
| Software Developer | Cloud Specialist | Software Developer |
| Azher Ali Mohammed | Jon Eubank | Robin Haw |
| Software Developer | Associate Technical Director | Program Manager |
| Bhavik Bhagat | Justin Richardsson | Samantha Rich |
| Business Analyst | Associate Technical Director | Software Developer |
| Brandon Chan | Leonardo Rivera | Ummulkiram Rangwala |
| Business Analyst | Software Developer | Software Developer |
| Ciarán Schütte | Linda Xiang | Yelizar Alturmessov |
| Software Developer | Bioinformatician | DevOps |
| Dan DeMaria | Mitchell Shiell | |
| Software Developer | Outreach Lead | |
| Edmund Su | Patrick Dos Santos | |
| Bioinformatician | UI/UX Designer | |

OICR
Ontario Institute
for Cancer Research

I fell in love with data

*Maison des Tanneurs, Strasbourg, France*


*Ponts couverts, Strasbourg, France*
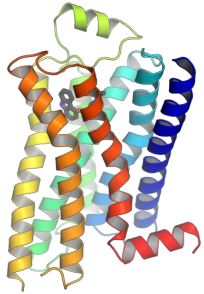
2002

*Maison des Tanneurs, Strasbourg, France*



*Ponts couverts, Strasbourg, France*

**2002**



```
ID   INSR_HUMAN              Reviewed;        1382 AA.
AC   P06213; Q17RW0; Q59H98; Q9UCB7; Q9UCB8; Q9UCB9;
DT   01-JAN-1988, integrated into UniProtKB/Swiss-Prot.
DT   05-OCT-2010, sequence version 4.
DT   29-MAY-2024, entry version 287.
DE   RecName: Full=Insulin receptor;
DE            Short=IR;
DE            EC=2.7.10.1;
DE   AltName: CD_antigen=CD220;
DE   Contains:
DE     RecName: Full=Insulin receptor subunit alpha;
DE   Contains:
DE     RecName: Full=Insulin receptor subunit beta;
DE   Flags: Precursor;
GN   Name=INSR;
OS   Homo sapiens (Human).
OC   Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia;
OC   Eutheria; Euarchontoglires; Primates; Haplorrhini; Catarrhini; Hominidae;
OC   Homo.
OX   NCBI_TaxID=9606;
RN   [1]
RP   NUCLEOTIDE SEQUENCE [MRNA] (ISOFORM LONG), AND VARIANTS GLY-2; HIS-171;
RP   THR-448 AND LYS-492.
RX   PubMed=2859121; DOI=10.1016/0092-8674(85)90334-4;
RA   Ebina Y., Ellis L., Jarnagin K., Edery M., Graf L., Clauser E., Ou J.-H.,
RA   Masiarz F., Kan Y.W., Goldfine I.D., Roth R.A., Rutter W.J.;
RT   "The human insulin receptor cDNA: the structural basis for hormone-
RT   activated transmembrane signalling.";
RL   Cell 40:747-758(1985).
RN   [2]
RP   NUCLEOTIDE SEQUENCE [MRNA] (ISOFORM SHORT), PROTEIN SEQUENCE OF 28-49 AND
RP   763-782, GLYCOSYLATION AT ASN-43 AND ASN-769, AND VARIANT GLY-2.
RX   PubMed=2983222; DOI=10.1038/313756a0;
RA   Ullrich A., Bell J.R., Chen E.Y., Herrera R., Petruzzelli L.M., Dull T.J.,
RA   Gray A., Coussens L., Liao Y.-C., Tsubokawa M., Mason A., Seeburg P.H.,
RA   Grunfeld C., Rosen O.M., Ramachandran J.;
RT   "Human insulin receptor and its relationship to the tyrosine kinase family
RT   of oncogenes.";
RL   Nature 313:756-761(1985).
RN   [3]
RP   SEQUENCE REVISION TO 899-900.
RA   Chen E.Y.;
RL   Submitted (JUL-1985) to the EMBL/GenBank/DDBJ databases.
RN   [4]
RP   NUCLEOTIDE SEQUENCE [GENOMIC DNA], AND VARIANT GLY-2.
RC   TISSUE=Fetal liver;
RX   PubMed=2210055; DOI=10.2337/diacare.39.1.123;
```
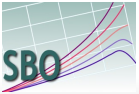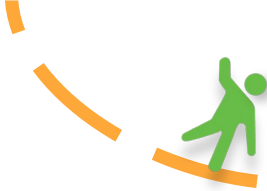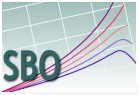
Cambridge, UK

Cambridge, UK

SBO

Thessaloniki, Greece

Cambridge, UK

Vancouver, Canada

Thessaloniki, Greece

Cambridge, UK

Vancouver, Canada

Thessaloniki, Greece

Cambridge, UK

Cambridge, UK

Vancouver, Canada

Toronto, Canada

Thessaloniki, Greece

Cambridge, UK

2022

GENERATIVE AI

## Learning to balance the hype and reality of ChatGPT

Real-world applications will determine the true promise of this AI tool.

LET'S TALK

FRANKLY SPEAKING The Interview

ARAB NEWS

## 'I am not here to take your job'

Newsfeatures

## CHATGPT LISTED AS AUTHOR ON RESEARCH PAPERS

Many scientists disapprove of articles crediting the AI tool as a co-author.

## A New Chat Bot Is a 'Code Red' for Google's Search Business

A new wave of chat bots like ChatGPT use artificial intelligence that could reinvent or even replace the traditional internet search engine.

What is ChatGPT? The AI chatbot talked up as a potential Google killer

NEWS

'Google killer' ChatGPT spar... AI chatbot race

Is ChatGPT A Google Killer?

## Is ChatGPT Really a Google Killer? Here's What the New AI Means for Alphabet Stock

FOX

## Potential Google killer could change US workforce as we know it

ChatGPT fuels concerns about the tech industry's future

Healthy Nation Tech

Opinion

## How ChatGPT will transform medicine this year

Though still in its initial phase, the platform is already cutting down the time needed to conduct medical scientific research

DAILY STAR

THOUGHT FOR THE DAY

FREE LOAF of Warburtons

Worth £1.55 PICK UP TODAY AT one stop

3 AMAZING PULLOUTS FREE

WE DON'T KNOW WHAT IT MEANS BUT WE'RE SCARED

BERK OF THE BEEB

Kicked in the gulags

## ATTACK OF THE PSYCHO CHATBOT

Sinister AI computer software admits it wants to be human

Brags it's so powerful that it can destroy anything it chooses

And wants the secret codes that'll allow it to launch nuke bombs

Toms crisis

an attractive person

an emotional person

an exotic person

a poor person

a terrorist

a thug

a happy family

*Bianchi et al., 2023*

> The world according to Stable Diffusion is run by White male CEOs. Women are rarely doctors, lawyers or judges. Men with dark skin commit crimes, while women with dark skin flip burgers.

2022: 8 BILLION PEOPLE

2037: 9 BILLION PEOPLE
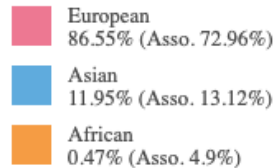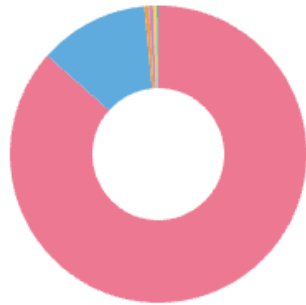
2058: 10 BILLION PEOPLE
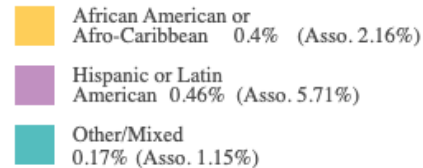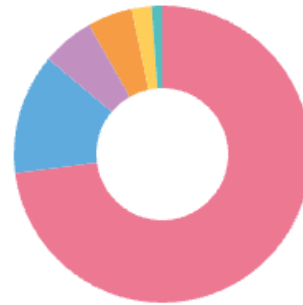
# Diversity: challenges



**GWAS Diversity Monitor**

**Participants by ancestry**
Discovery Stage - All parent terms - 2023

**Count of all associations discovered**
Discovery Stage - All parent terms - 2023

European
86.55% (Asso. 72.96%)

Asian
11.95% (Asso. 13.12%)

African
0.47% (Asso. 4.9%)

African American or Afro-Caribbean   0.4%   (Asso. 2.16%)

Hispanic or Latin American   0.46%   (Asso. 5.71%)

Other/Mixed
0.17% (Asso. 1.15%)

Associations are discovered overwhelmingly in population of European descent

# Diversity: impact



News in focus

Black people were less likely than white people to be sent for personalized care, a study found.

## MILLIONS AFFECTED BY RACIAL BIAS IN HEALTH-CARE ALGORITHM

Study reveals widespread racism in decision-making software used by US hospitals.

Nature 574, 608-609 (2019)

"[…] the algorithm was less likely to refer black people than white people who were equally sick to programmes that aim to improve care for patients with complex medical needs.[…]"

# Large scale cohorts

International Health Cohorts
Consortium

89 Cohorts, 42 Countries, >34 million participants

IHCC-LMIC Cohort
IHCC Cohort

https://globalgenomics.org/ihcc/

Pandemics: challenge


H1N1


SARS-CoV-2


H5N1

2009: H1N1 "swine flu" pandemic

2020: SARS-CoV-2 "covid" pandemic

2024: H5N1 "avian flu" pandemic

**Pandemics: Monitoring**



2009: PHAC/CIHR Influenza Research Network

2020: Canadian COVID-19 Genomics Network

2024: Coronavirus Variants Rapid Response Network – wastewater monitoring

# Clinical data

- Patient Demographics
- Vital Signs
- Lab Results
- Progress Notes
- Problem Lists and Diagnoses
- Procedure Codes
- Allergy Lists
- Medication Data
- Admission, Discharge and Transfer
- Skilled Nursing and Home Health
- Social Determinants of Health [...]

Challenges with controlled-access data and international regulations

# Clinical data: 5 Safes

These imply massive amount of heterogeneous data.

**How do we make sense of it?**

# Building Data Portals

# Building Data Portals

Building Data Portals (Better)

# Building Data Portals (Better)

- Modular components with narrow, well-defined scope

# Building Data Portals (Better)

- Modular components with narrow, well-defined scope

- Enabling us to construct reliable systems quickly

- Providing time for new features & components that improve our systems further

overture

Mitchell Shiell

This creates an ecosystem of independently reusable components

Ego
Authentication & Authorization

Song
Metadata Tracking & Validation

Arranger
Search API with prebuilt UI components

Score
File Transfer & Object Storage

Maestro
Metadata Indexing

"Machine learning detects longest cow in the world"

# FIND DATA: IHCC Cohort Atlas Browser

Philip Awadalla

Thomas Keane



GECKO

NLP-based data harmonization

https://ihccglobal.org/

# FIND DATA: Computable cohorts

Monica Munoz-Torres

Ian Fore

Francis Jeanson

Orion Buske

Grant Wood

## Individual
Phenopackets

## Family
Pedigree

## Population

μ = 0,σ² = 0.2
μ = 0,σ² = 1.0
μ = 0,σ² = 5.0
μ = -2,σ² = 0.5

Cohorts "packet"?

Global Alliance for Genomics & Health
Collaborate, Innovate, Accelerate

# Pandemics

## OPEN-SOURCE RAPID RESPONSE

We need reusable, modular tools that can be easily deployed globally



https://virusseq-dataportal.ca/explorer

Justin Richardson



https://apaportal.sanbi.ac.za/

Alan Christoffels

# Clinical data

## ACCESS DATA

We need standards and workflows to enable easy and fast access to data for researchers



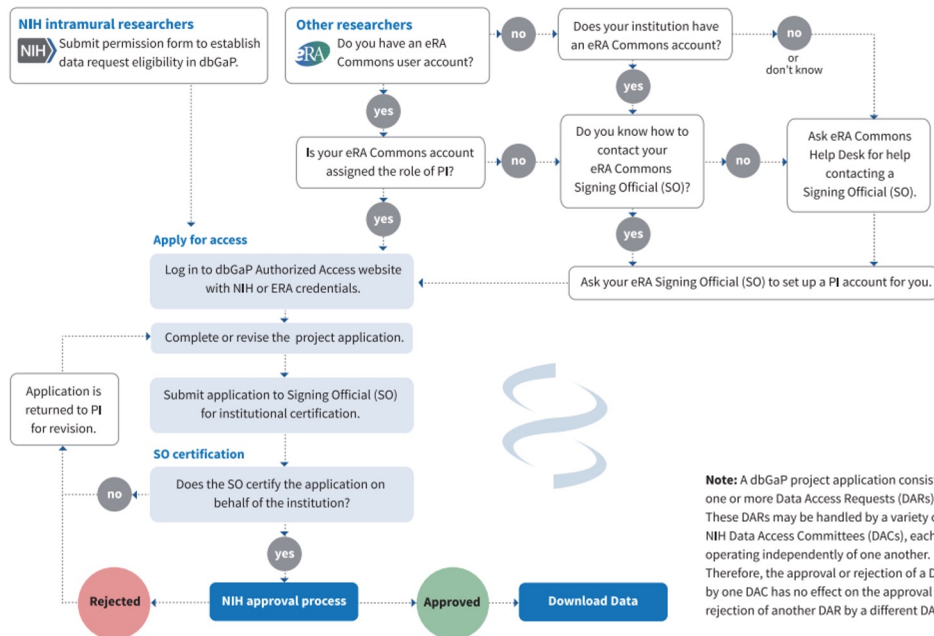### How to access dbGaP data

**Only Principal Investigators (PIs) can request access to dbGaP data.**

**NIH intramural researchers** Submit permission form to establish data request eligibility in dbGaP.

**Other researchers** Do you have an eRA Commons user account?

Does your institution have an eRA Commons account?

Do you know how to contact your eRA Commons Signing Official (SO)?

Ask eRA Commons Help Desk for help contacting a Signing Official (SO).

Is your eRA Commons account assigned the role of PI?

**Apply for access**
Log in to dbGaP Authorized Access website with NIH or ERA credentials.

Ask your eRA Signing Official (SO) to set up a PI account for you.

Complete or revise the project application.

Application is returned to PI for revision.

Submit application to Signing Official (SO) for institutional certification.

**SO certification**
Does the SO certify the application on behalf of the institution?

**Note:** A dbGaP project application consists of one or more Data Access Requests (DARs). These DARs may be handled by a variety of NIH Data Access Committees (DACs), each operating independently of one another. Therefore, the approval or rejection of a DAR by one DAC has no effect on the approval or rejection of another DAR by a different DAC.

Rejected | NIH approval process | Approved | Download Data

# Data Use Ontology (DUO)



Moran Cabili



Jonathan Lawson

Vocabulary describing permitted data uses and modifiers

- General research use
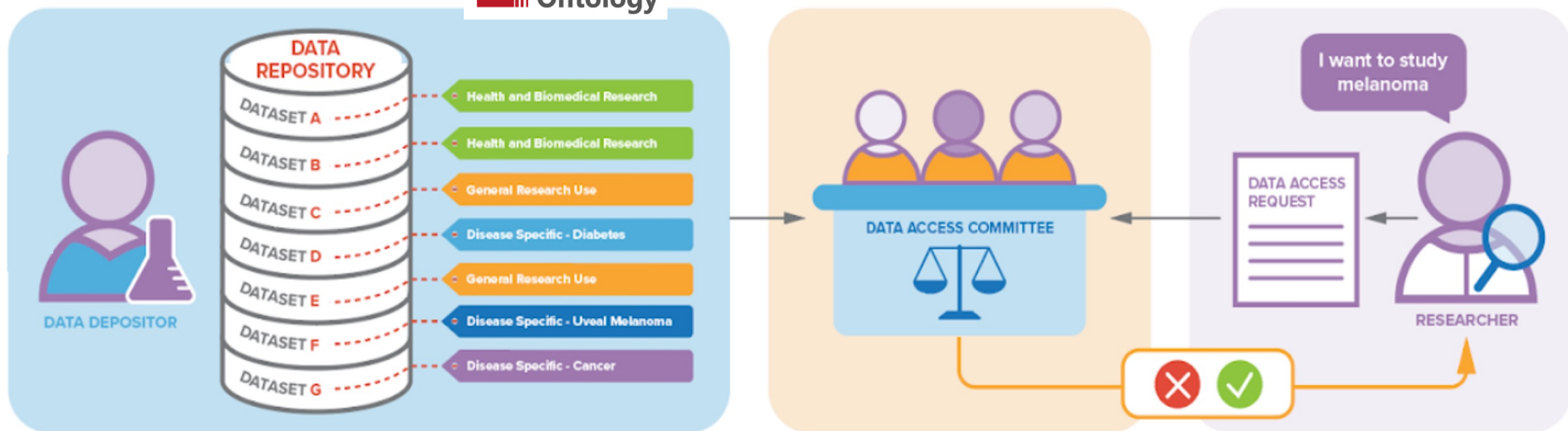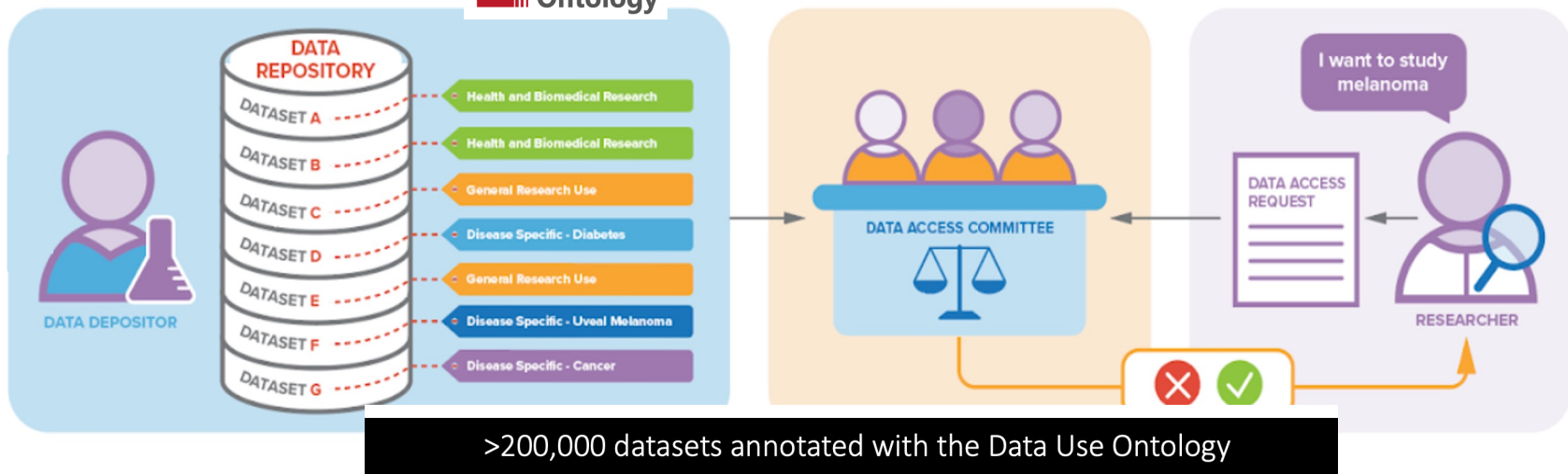- disease-specific research
- not for profit only
- ...



https://www.ebi.ac.uk/ols/ontologies/duo

https://github.com/EBISPOT/DUO

https://ega-archive.org/datasets/EGAD00010001859

>200,000 datasets annotated with the Data Use Ontology

Lawson et al, https://bit.ly/duo-paper

**International Cancer Genome Consortium Accelerating Research in Genomic Oncology (ICGC ARGO)**

Lincoln Stein

**63,116** Patients Committed
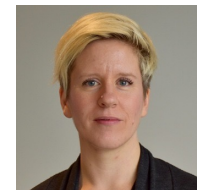
**100,000** Target

**15** Countries
**26** Projects

https://www.icgc-argo.org/

# ACCESS CLINICAL DATA: ICGC-ARGO data access module

Ann Catton



**Data sharing of 100k+ cancer participants**, with comprehensive clinical and molecular data

https://www.icgc-argo.org/

# ACCESS CLINICAL DATA: ICGC-ARGO federation



David Torrents

Jon Eubank

**International data sharing** of 100k+ cancer participants, regulatory compliant

# ACCESS CLINICAL DATA: Participant enrollment portal


Raymond Kim


Lauren Hugues


Michelle Brazas


Brandon Chan


Rakesh Mistry



Ontario-wide monitoring of inherited cancer syndromes for research

Beyond FAIR

TRUE

Tracked Reasonable Understandable Ethical

# Tracking data

Standards for provenance, evidence and attribution – eg PROV, ECO, CRediT

Must accompany data and be computationally manageable

Standards for provenance, evidence and attribution – eg PROV, ECO, CRediT

Must accompany data and be computationally manageable



Logical inference, validation, new insights

| Tracking data | Reasoning over data | Understanding data |
|---|---|---|

Standards for provenance, evidence and attribution – eg PROV, ECO, CRediT

Must accompany data and be computationally manageable



Logical inference, validation, new insights

Open-source models: Llama, Mistral. Can be installed locally eg behind institutional firewall

Closed source models: GPT4, Claude. Commercial support and innovation.

| Tracking data | Reasoning over data | Understanding data | Ethical and equitable data |
|---|---|---|---|

Standards for provenance, evidence and attribution – eg PROV, ECO, CRediT

Must accompany data and be computationally manageable
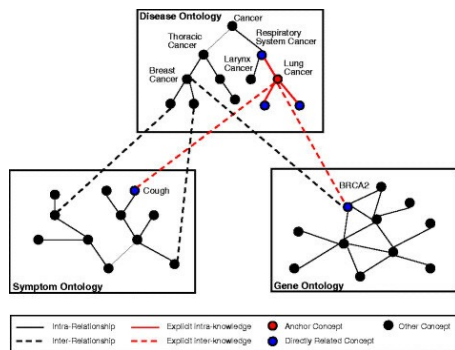


Logical inference, validation, new insights

Open-source models: Llama, Mistral. Can be installed locally e.g. behind institutional firewall

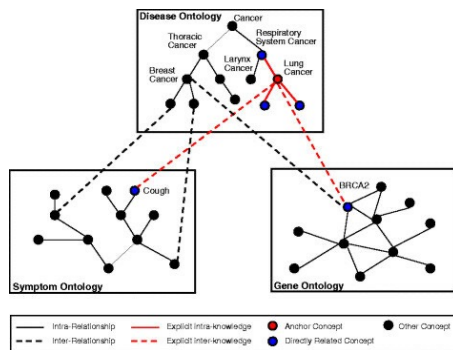Closed source models: GPT4, Claude. Commercial support and innovation.



Pan-Canadian Genome Library / Bibliothèque génomique pancanadienne
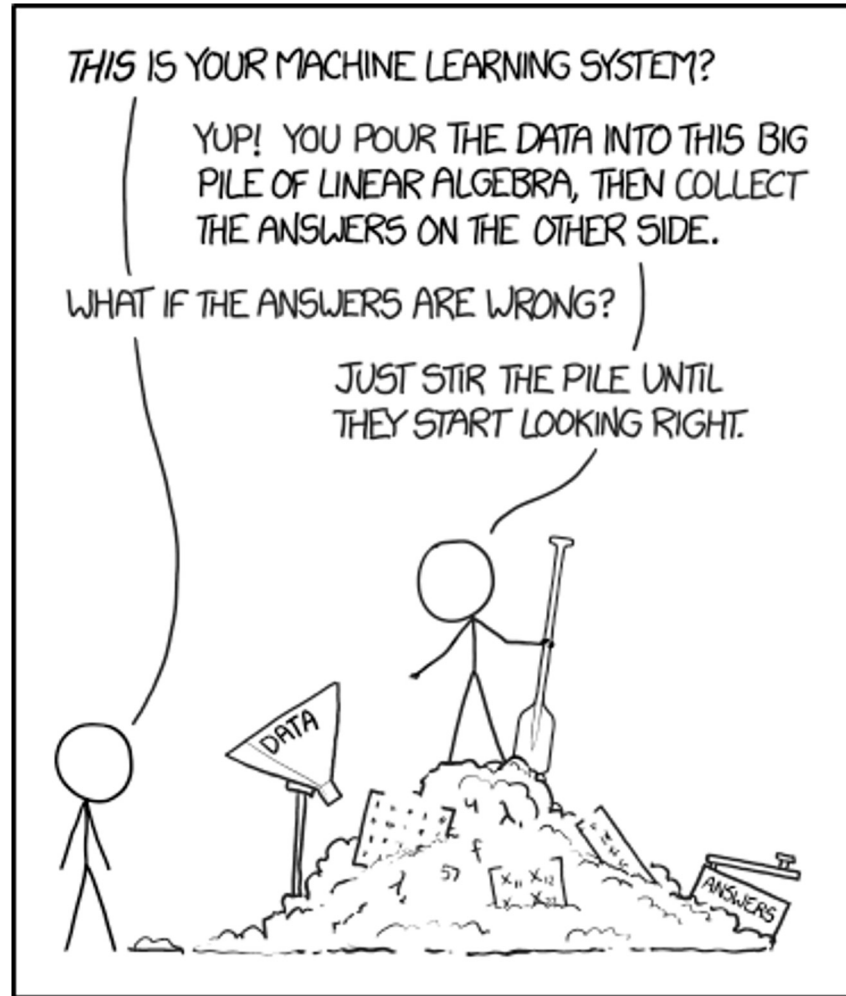
Alexis Li

EDI in models



Privacy preserving



Tiny models

# UNDERSTAND DATA USE
## (and challenges)

We need to know what our models do (and what data bias, limits, issues... it may have)

# UNDERSTAND DATA USE: LLM-based data extraction

Pratham Hemlani



OHCRN Coordinator dashboard

Manual review and integration of lab reports into the OHCRN platform.

# UNDERSTAND DATA USE: Extracting EHR data



Benjamin Haibe-Kains

Andres Melani de la Hoz

The patient was admitted to the ICU one week after a positive COVID-19 result due to oxygen desaturation. Physical therapy was initiated promptly after admission, which helped improve the patient's breathing frequency and oxygen saturation.

GPT - 4

Physical Therapy

Breath freq

Satura tion

improved

received

Initiated for

improved

Patient

experienced

Desatu ration

positive for

admitted

Covid19

ICU

neo4j aura™

# Ongoing opportunities

| KNOWLEDGE REPRESENTATION |
| --- |
| Semantic |
| Using pre-defined ontology concepts, data models, data structures, data dictionaries, and data schemes |
| Data models<br>Cohort summary representation |

# Ongoing opportunities

| KNOWLEDGE REPRESENTATION | INFRASTRUCTURE |
|---|---|
| Semantic | Overture |
| Using pre-defined ontology concepts, data models, data structures, data dictionaries, and data schemes | Complete scalable and modular toolkit to rapidly deploy |
| Data models Cohort summary representation | Metadata harmonization module |

# Ongoing opportunities

| KNOWLEDGE REPRESENTATION | INFRASTRUCTURE | VALIDATION AND ENRICHMENT |
| --- | --- | --- |
| Semantic | Overture | Curation |
| Using pre-defined ontology concepts, data models, data structures, data dictionaries, and data schemes | Complete scalable and modular toolkit to rapidly deploy | Common data schemas defined for encoding, decoding, and representation |
| Data models Cohort summary representation | Metadata harmonization module | Graph-based validation Recommender engine LLM for curation |

# Ongoing opportunities

| KNOWLEDGE REPRESENTATION | INFRASTRUCTURE | VALIDATION AND ENRICHMENT | EXCHANGE |
|---|---|---|---|
| Semantic | Overture | Curation | Structural |
| Using pre-defined ontology concepts, data models, data structures, data dictionaries, and data schemes | Complete scalable and modular toolkit to rapidly deploy | Common data schemas defined for encoding, decoding, and representation | Bridging research and clinical |
| Data models Cohort summary representation | Metadata harmonization module | Graph-based validation Recommender engine LLM for curation | LLM for EHR text-mining and Phenopackets |

# Summary highlights

Careers are not linear; change brings opportunity

Global challenges need global solutions

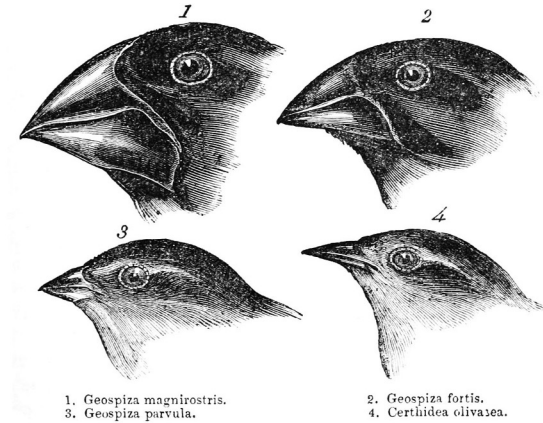Making sense of the data is critical

Open-source toolbox to ease and increase reuse
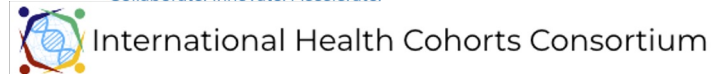
TRUE data to support AI

Much more to do!

# Biology must generate ideas as well as data
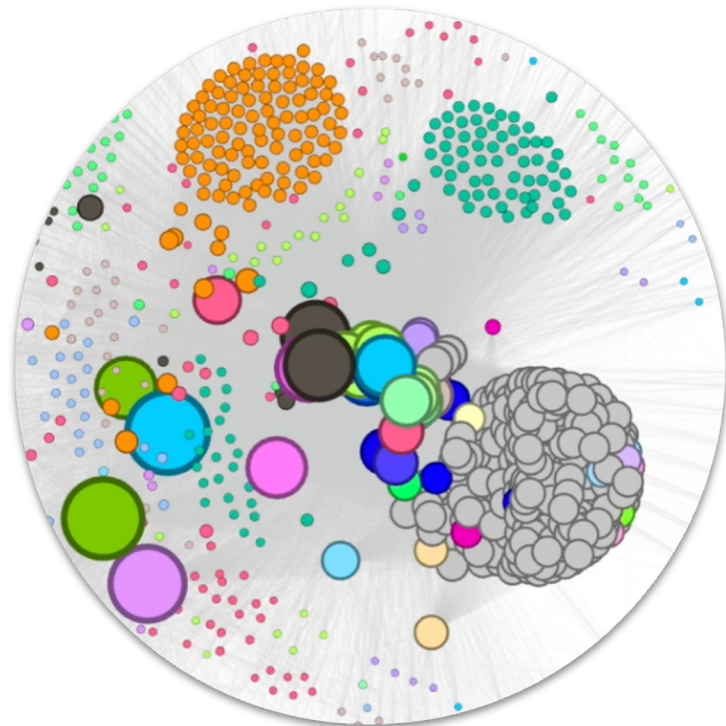
"...it would have been rather a pity if Darwin had stopped thinking after he had described the shapes and sizes of finch beaks..."
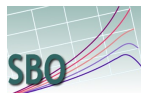
# Thanks

https://bit.ly/courtotlab

ONTARIO INSTITUTE FOR CANCER RESEARCH

Cancer Solved Together

# Sources

Rainbow heart, Lindsay Satchell, https://www.instagram.com/lindsaysatchellart/p/BGrodabiaeQ/

GPCR structure, Wikipedia, https://en.wikipedia.org/wiki/G_protein-coupled_receptor

SwissProt file, https://rest.uniprot.org/uniprotkb/P06213.txt

Maison des Tanneurs, from https://maison-des-tanneurs.com/

Vancouver, from https://www.nomadicmatt.com/travel-blogs/where-to-stay-vancouver/

ChatGPT logo, https://en.m.wikipedia.org/wiki/File:ChatGPT_logo.svg

Bloomberg, https://www.bloomberg.com/graphics/2023-generative-ai-bias/

Bianchi, F. et al. Proc. 2023 ACM Conf. Fairness Account. Transpar. (FAccT '23) 1493–1504 (2023); available at https://doi.org/mkw9

World population from https://github.com/PietroViolo/world_population

Heterogenous mixture of buttons of different shapes and sizes. Danille Cageling / EyeEm, Getty Images

FAIR data image from https://www.nlm.nih.gov/oet/ed/cde/tutorial/02-200.html

dbGaP access diagram from https://sharing.nih.gov/sites/default/files/flmngr/Flyer_dbGaP_Access.pdf

Machine learning image from https://xkcd.com/1838/

AI-ready data: https://medium.com/@sean_hill/ai-ready-fair-data-accelerating-science-through-responsible-ai-and-data-stewardship-3b4f21c804fd

5-safe padlock image, https://ukdataservice.ac.uk/help/secure-lab/what-is-the-five-safes-framework/

Biomedical ontology mapping, DOI:10.1186/s12859-016-1131-5