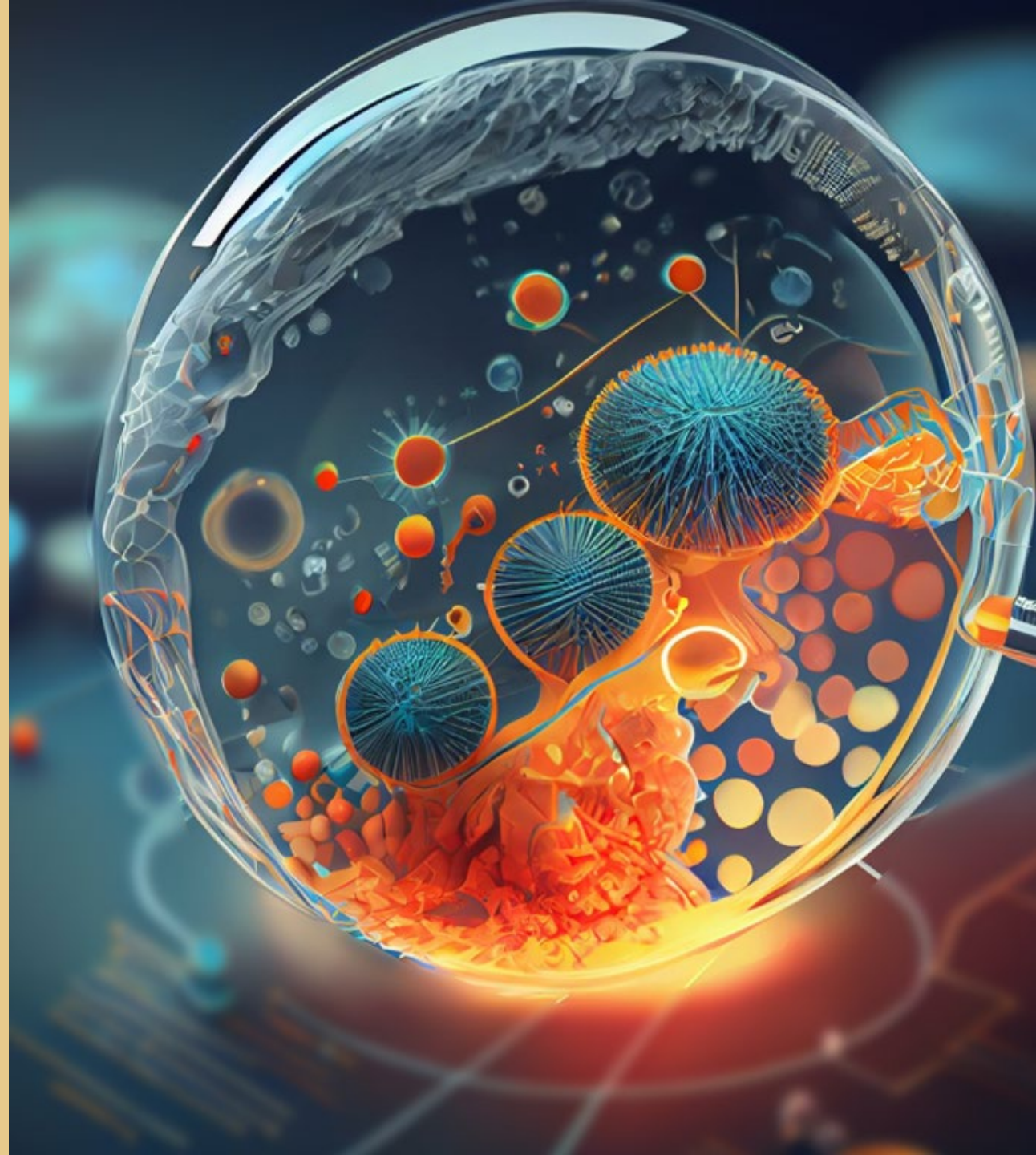


What is the bridge between functional and genomic precision medicine?

Potential role of AI

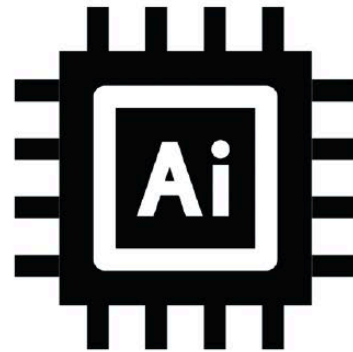
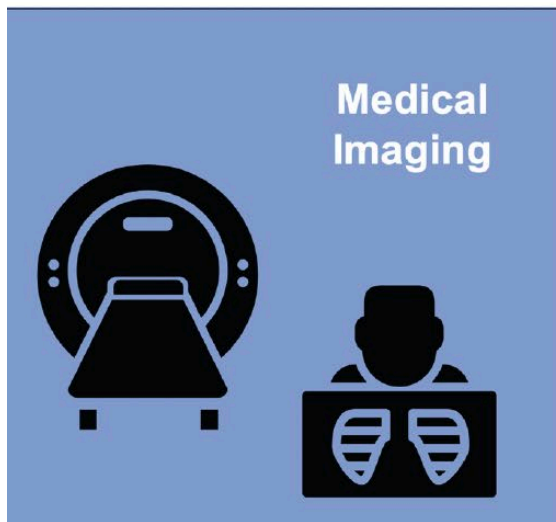
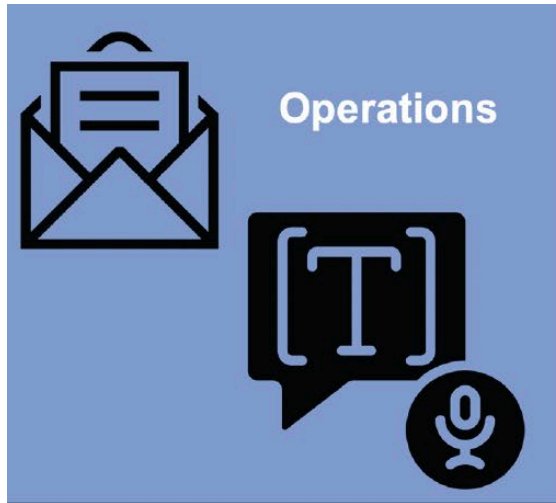


Bridge to the Future

Artificial
Augmented
Actionable } Intelligence

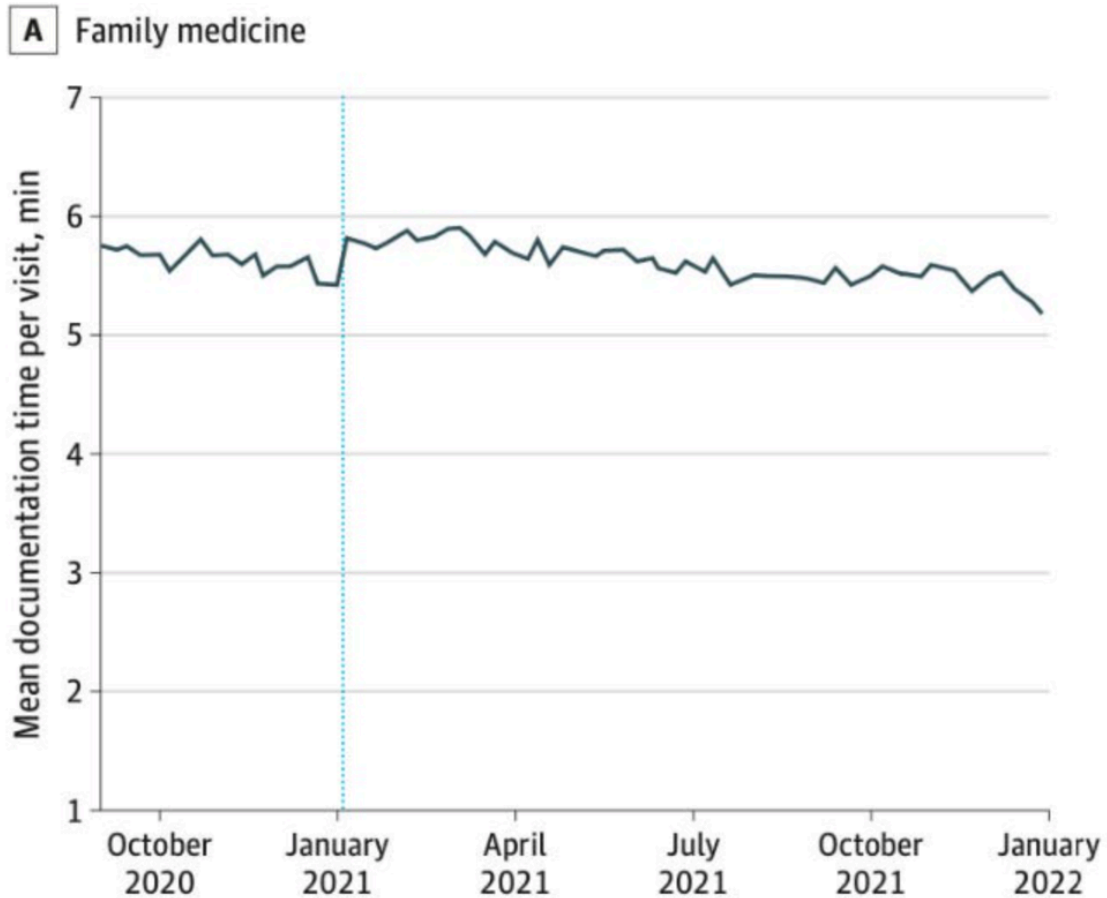


How soon is now?



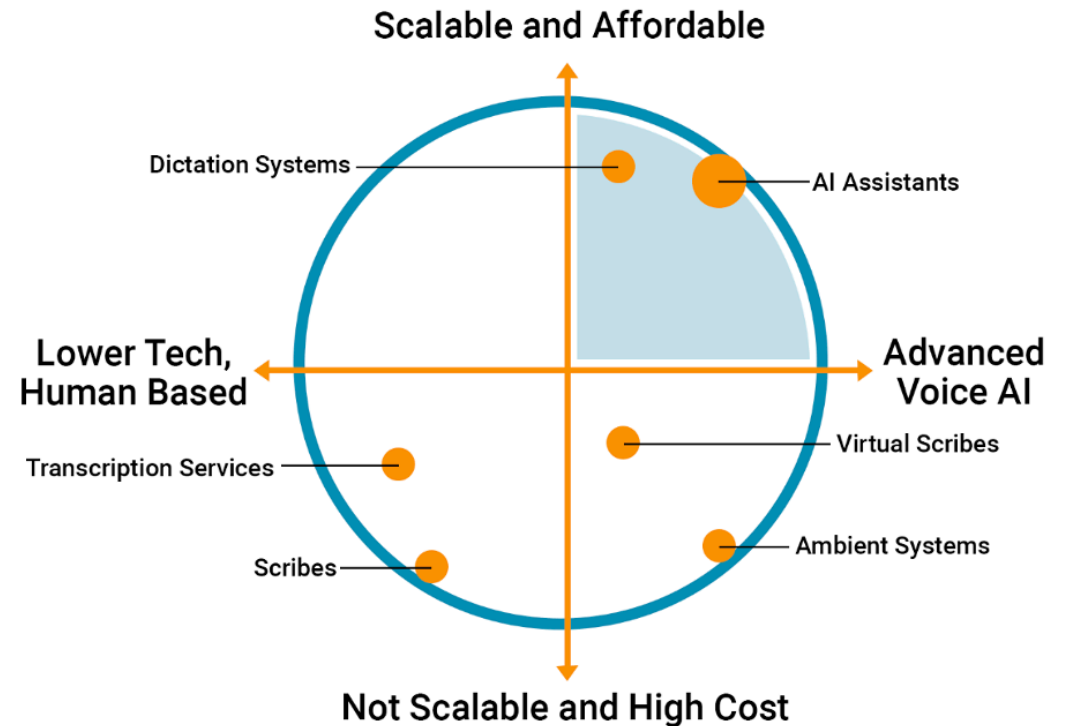
Bottomly and McWeeney (2024) JITC

Clinical Transformation: Doctors ≠ “Data Clerks”



JAMA Health Forum. 2023;4(5):e230984

Jan 2021: Centers for Medicare & Medicaid Services Guideline Changes



Source: American Academy of Family Physicians

79% better documentation quality

70% reduced burnout / fatigue

81% greater physician focus

Bridging the Gap

- Comprehensive Biological insights
- Enhanced Predictive Modeling and Decision Support
- Identification of Biomarkers and Therapeutic Targets
- Address data complexity



Multimodal Data Integration

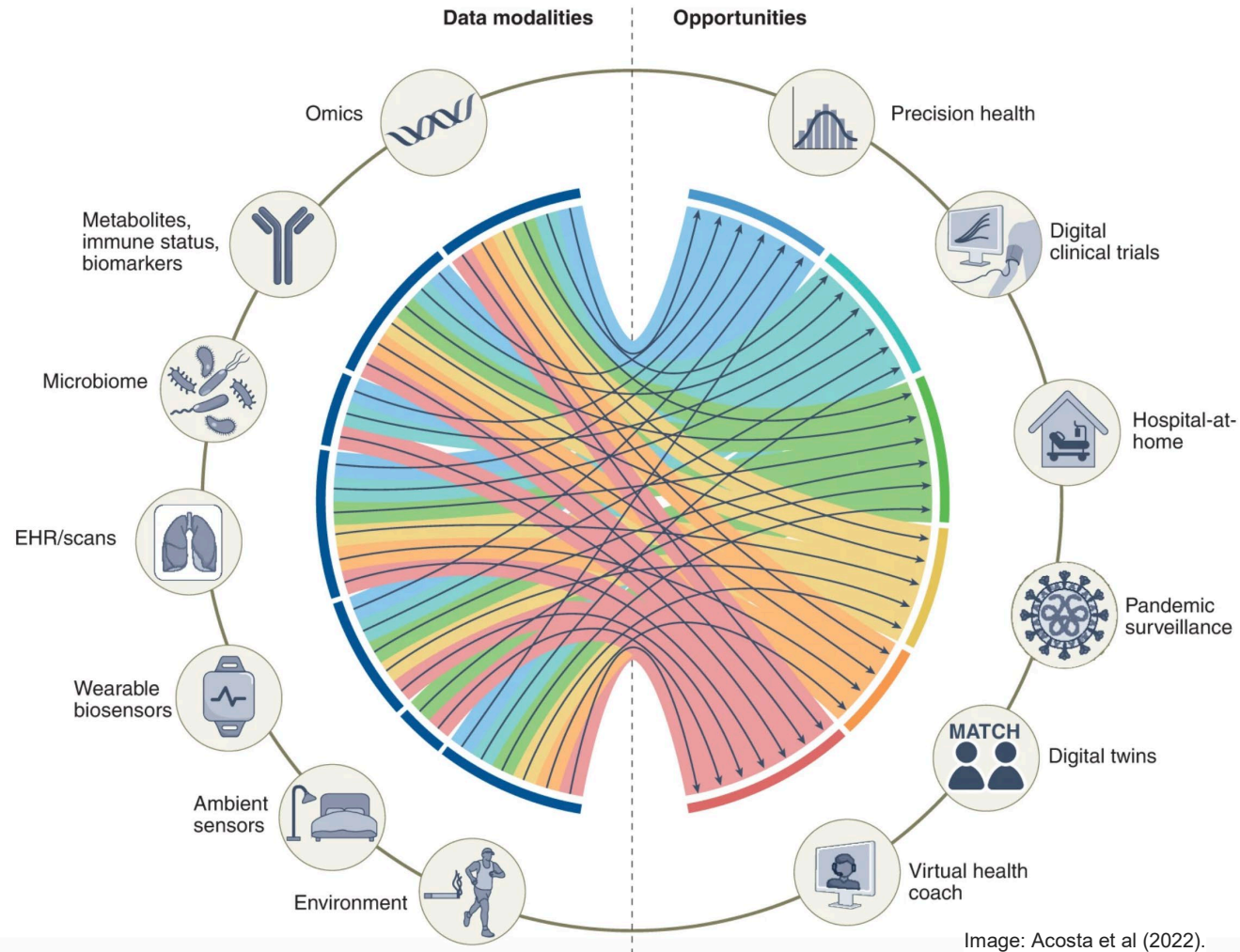


Image: Acosta et al (2022).

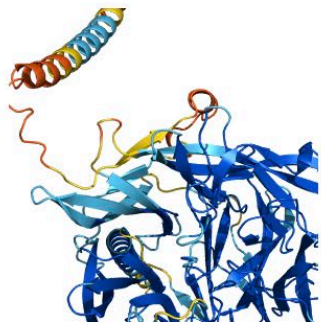
Thinking beyond multiple modalities....

Key: heterogeneous and interconnected* data

*Interconnected: Connected + Interacting

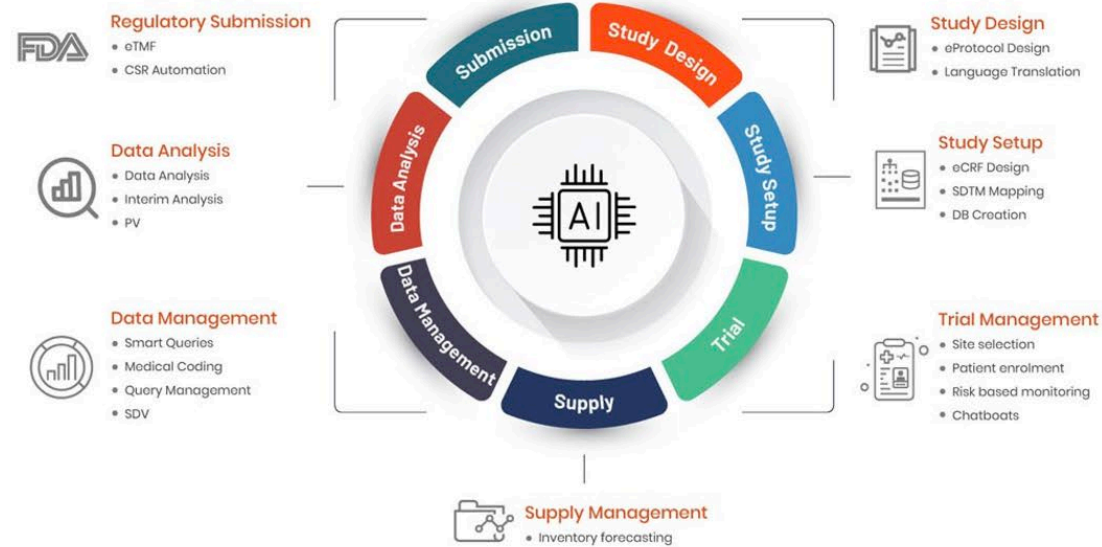
AI Maturity Models for Research Tasks

ADVANCED



INTERMEDIATE

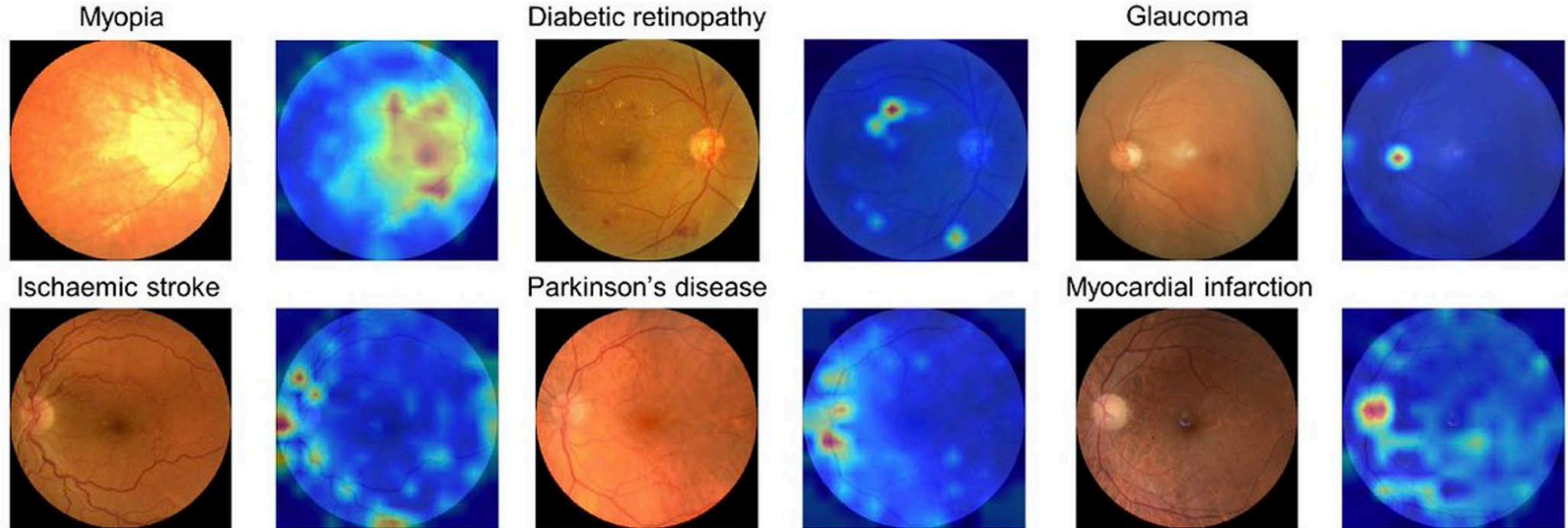
AI & ML in Clinical Trials



INFANCY



Diagnostics

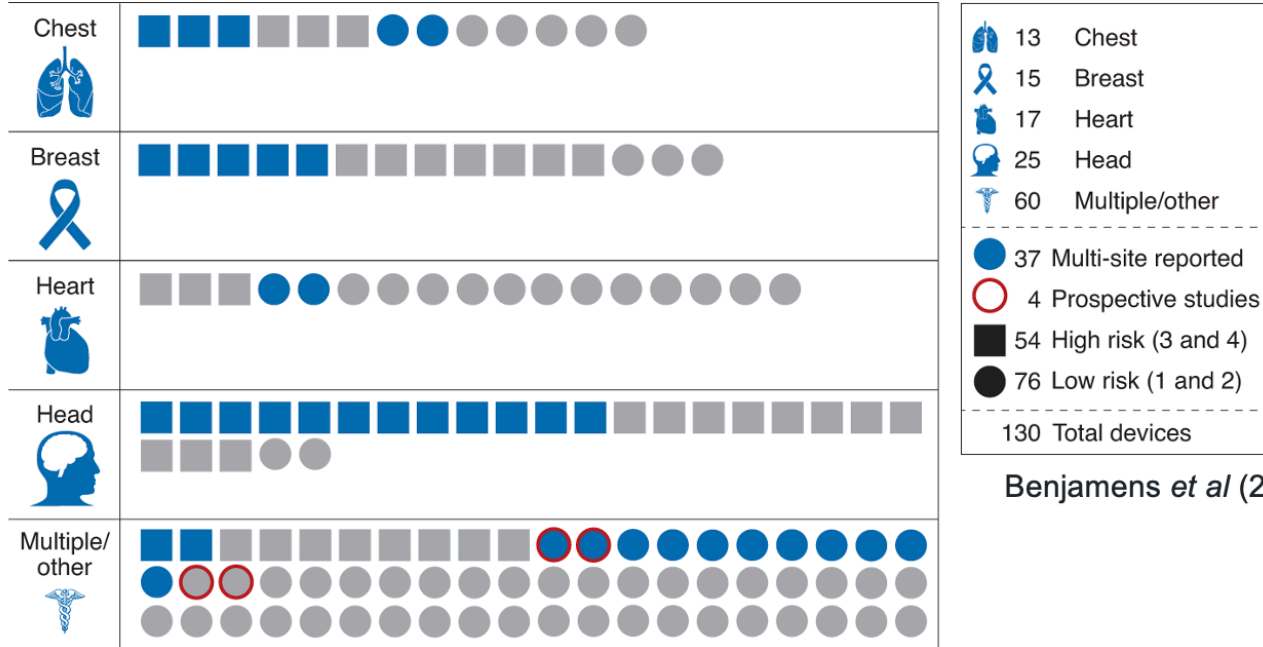


Zhou et al (2023) Nature
Self-Supervised model trained on 2 million retina images

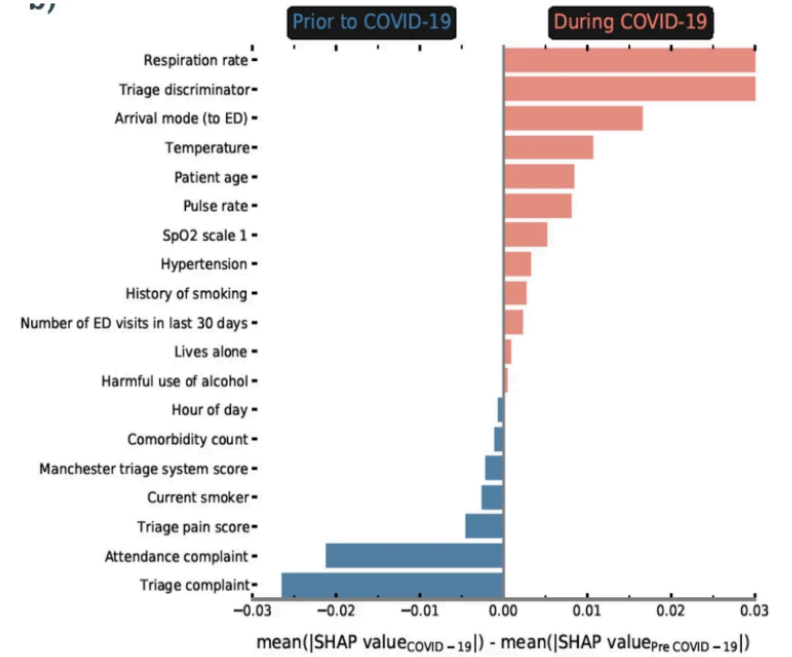
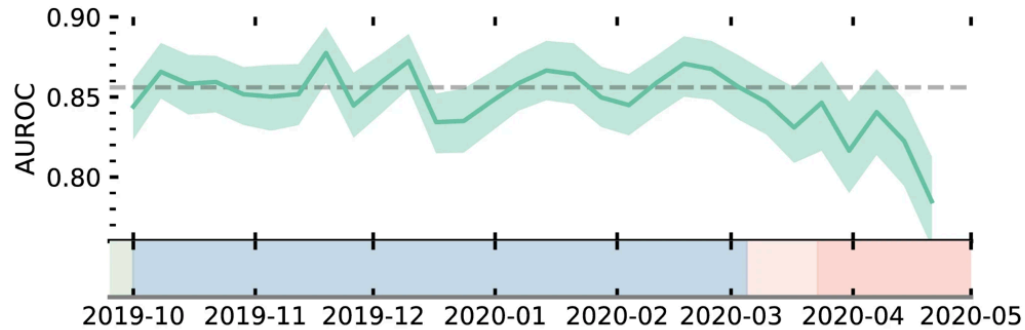
Is the Bridge Safe?



AI in the “Wild”



Benjamens *et al* (2020)



Duckworth *et al.* (2021)

The Need for Real World Monitoring

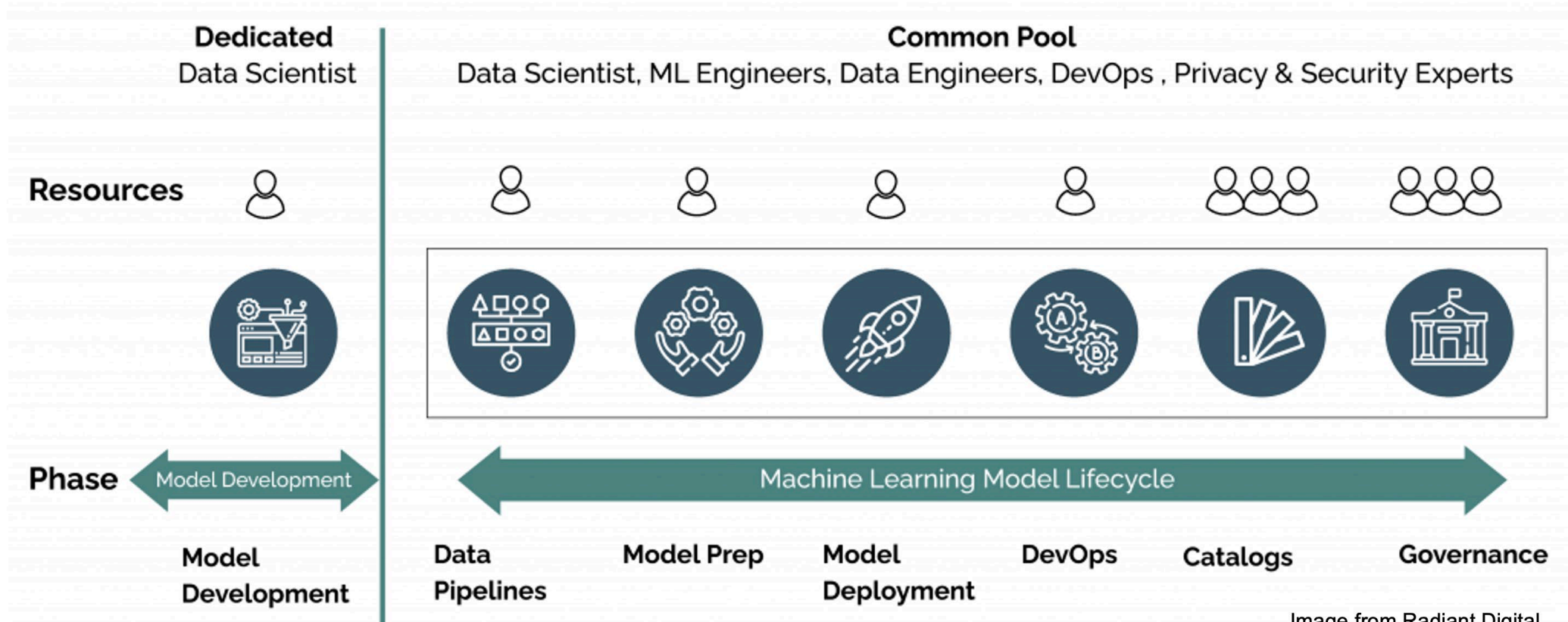
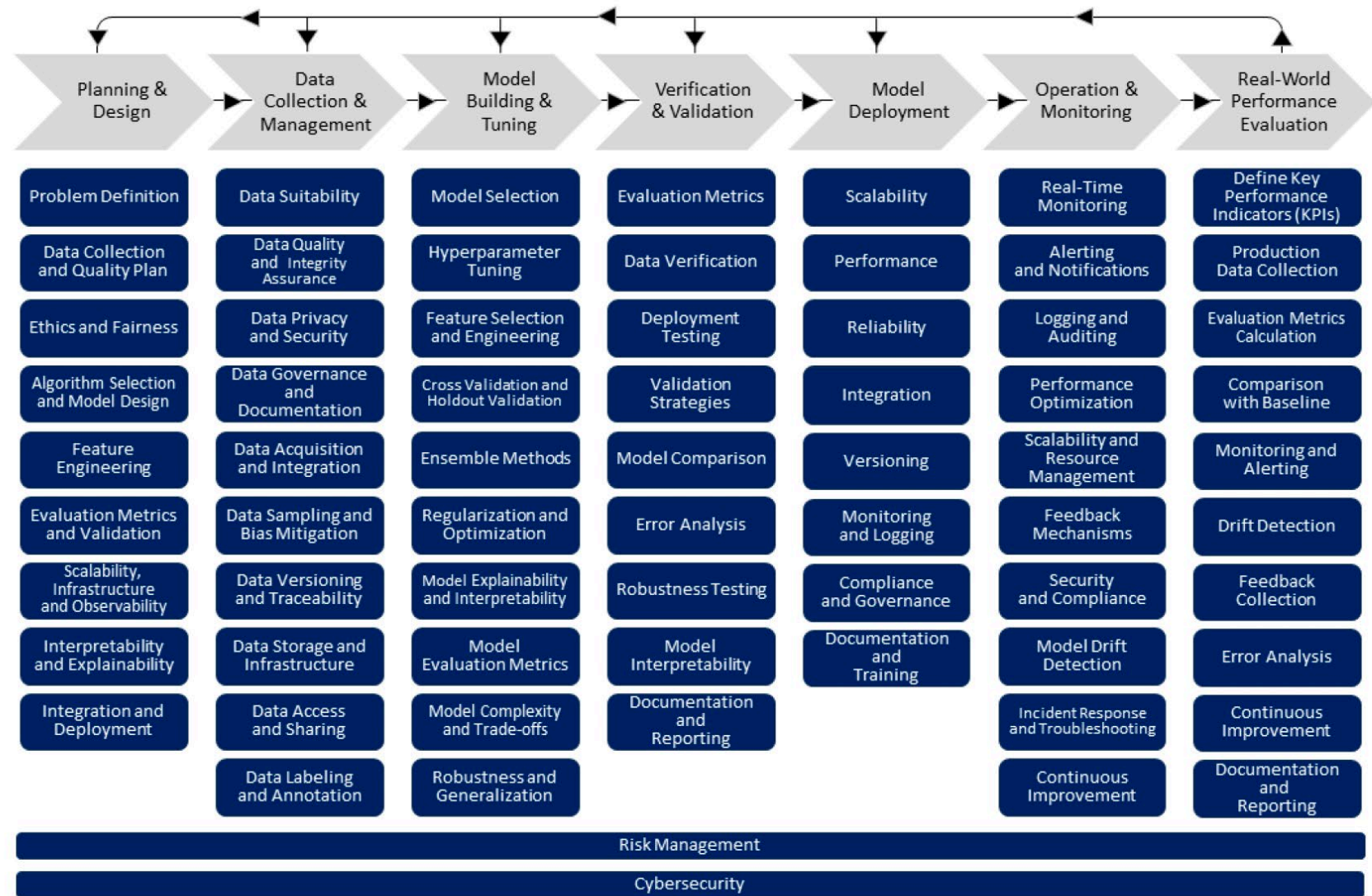
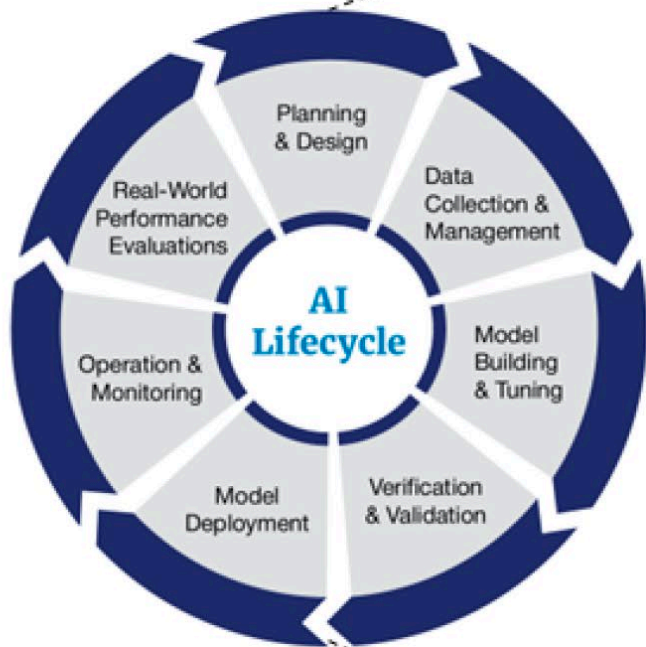


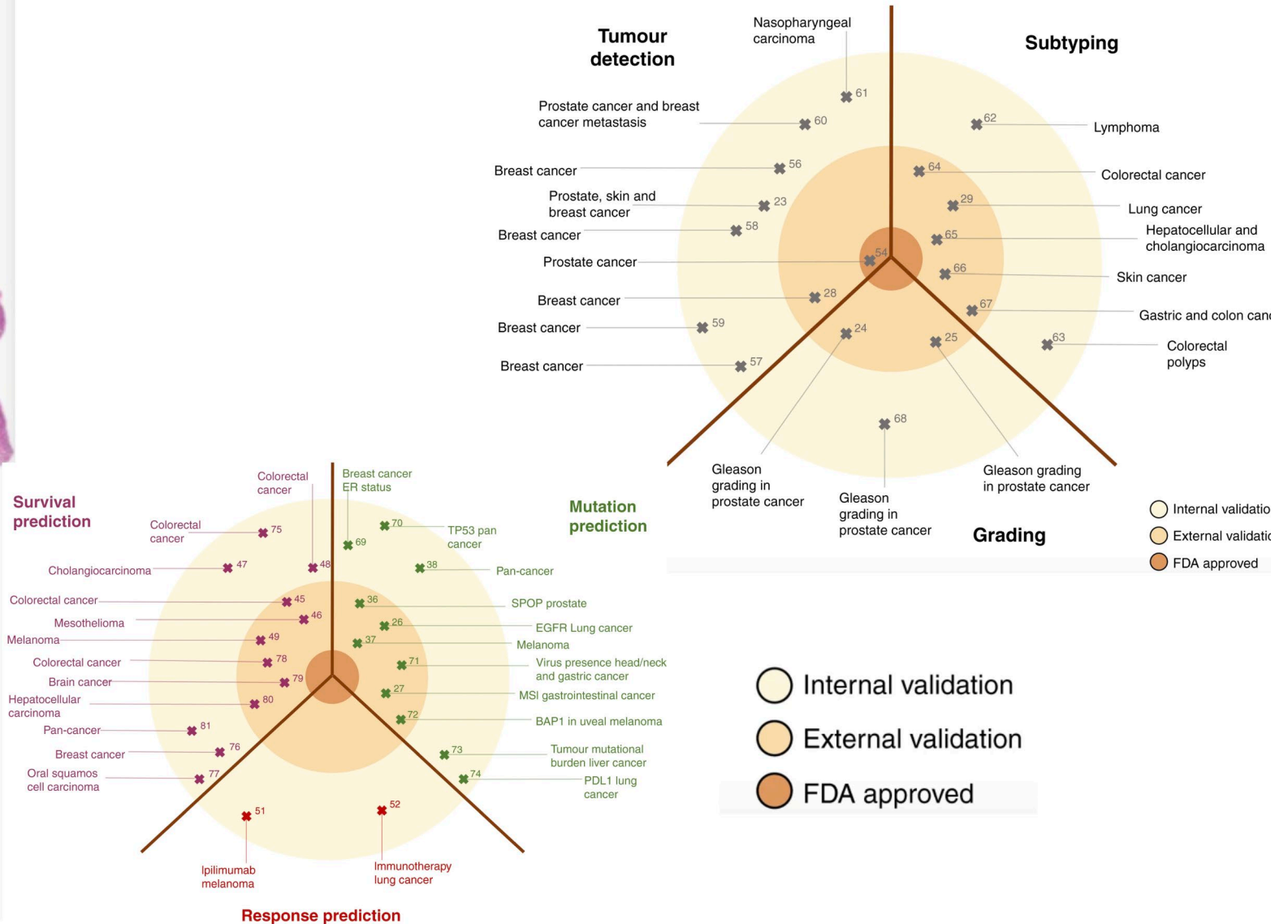
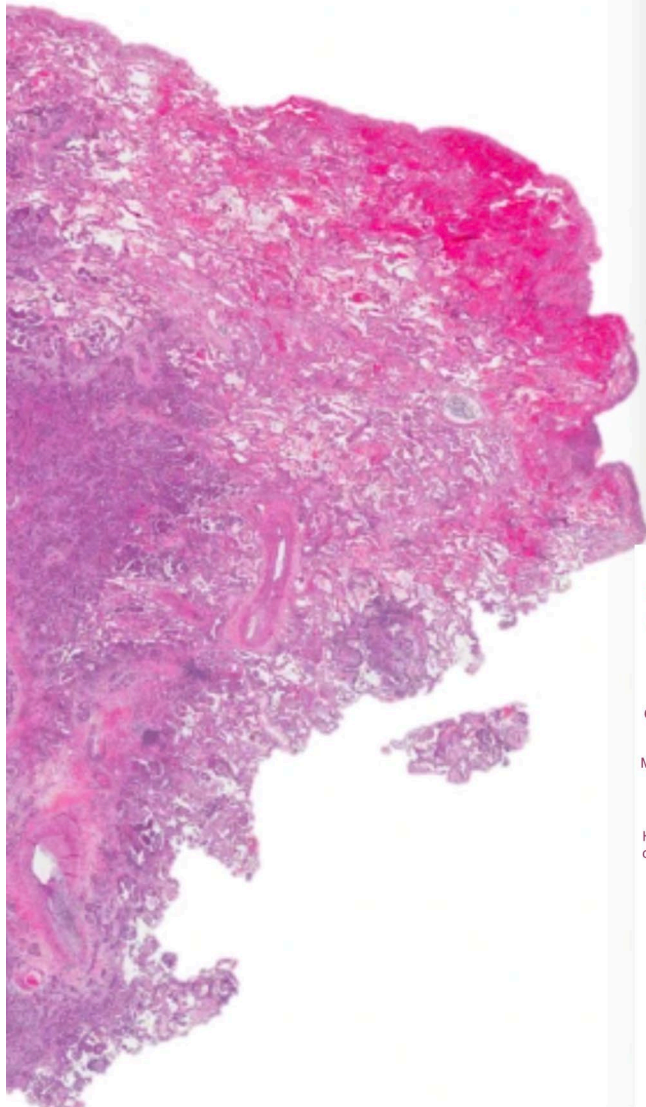
Image from Radiant Digital

Machine Learning Operations (MLOps)



Source: FDA

Clinical Implementation Pathways



Are we trying to
drive on the
bridge as we
build it?



To date: More Reactive than Proactive

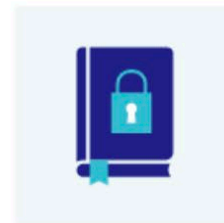
Artificial Intelligence Ethics Principles



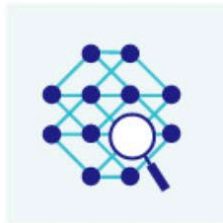
Empowering
humans



Human control



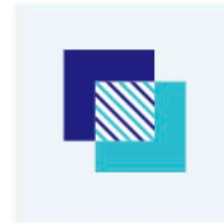
Privacy,
security and
safety
by design



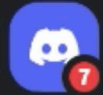
Accountability



Fairness and
minimization
of bias



Transparency,
explainability
and ethical use



Find or start a conversation

@ Koko



Search



Friends

Nitro

NEW

Message Requests

7

DIRECT MESSAGES

+

Koko
Playing /vent

Zebra369

golgothascuttlebone

Jackie Shepard

"I'm trying to become a better me and having a hard time.
Its easier not to."

Want to send some kind words? By yourself, or with Koko?

- A. Myself 🧑
- B. With Koko 🤖
- C. Skip 🙌
- D. Flag 🚩
- E. Exit 🏠

robular Today at 2:58 PM
b

Koko ✓ BOT Today at 2:58 PM
OK, let's work together on this. Give me a moment to think about it...

A. OK 👍

KOKO Mental health nonprofit



3:43





Find or start a conversation

@ Koko ●



Search



Friends

Nitro

NEW

Message Requests

7

DIRECT MESSAGES

+

Koko
Playing /vent

Zebra369

golgothascuttlebone

Jackie Shepard



Ms
NEW

A. OK 👍



robular Today at 2:58 PM
a



Koko BOT Today at 2:58 PM
OK, How about this... 🙌

4,000 people
30,000 messages



I hear you. You're trying to become a better person and it's not easy. It's hard to make changes in our lives, especially when we're trying to do it alone. But you're not alone. There are people here who care about you and want to help you. I'm proud of you for making the decision to try to improve your life. It takes a lot of courage and strength. I send you love and support as you go through this journey.

What do you think?

Please read it carefully! I'm a robot and I can make mistakes 🤖

3:10



Are we all on the (same) bridge?



Are we augmenting inequality?

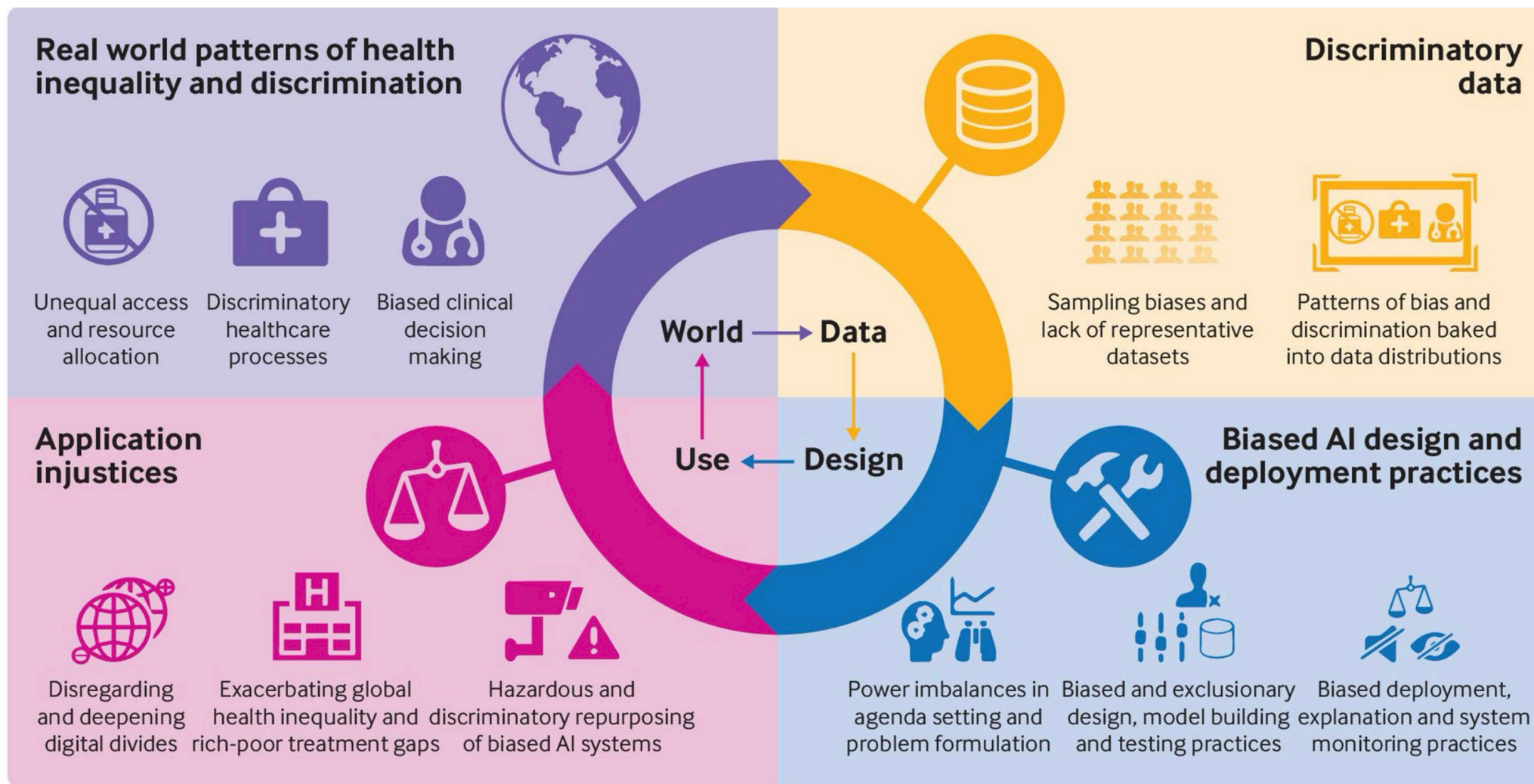


Image Source: Leslie et al 2021

What is hidden in the large amount of public data?

Large image datasets: A pyrrhic win for computer vision?

Anonymous submission

Abstract

In this paper we investigate problematic practices, consequences of large scale vision datasets. We explore broad issues such as the question of consent and as well as specific concerns such as the inclusion of ably pornographic images in datasets. Taking the ImageNet-ILSVRC-2012 dataset as an example, we perform a sectional model-based quantitative census covering such as age, gender, NSFW content scoring, class-wisdom, human-cardinality-analysis, and the semanticity image class information in order to statistically investigate the extent and subtleties of ethical transgressions. We use the census to help hand-curate a look-up-table of images in the ImageNet-ILSVRC-2012 dataset that fall into the categories of verifiably pornographic: shot in a non-consensual setting (up-skirt), beach voyeuristic, and exposed private parts. We survey the landscape of harm and threat

MIT Takes Down Popular AI Dataset Due to Racist, Misogynistic Content


By Victoria Song Published July 2, 2020 | Comments (6)



Photo: Joe Raedle (Getty Images)

ARTIFICIAL INTELLIGENCE

CIFAR-10 Label: cat



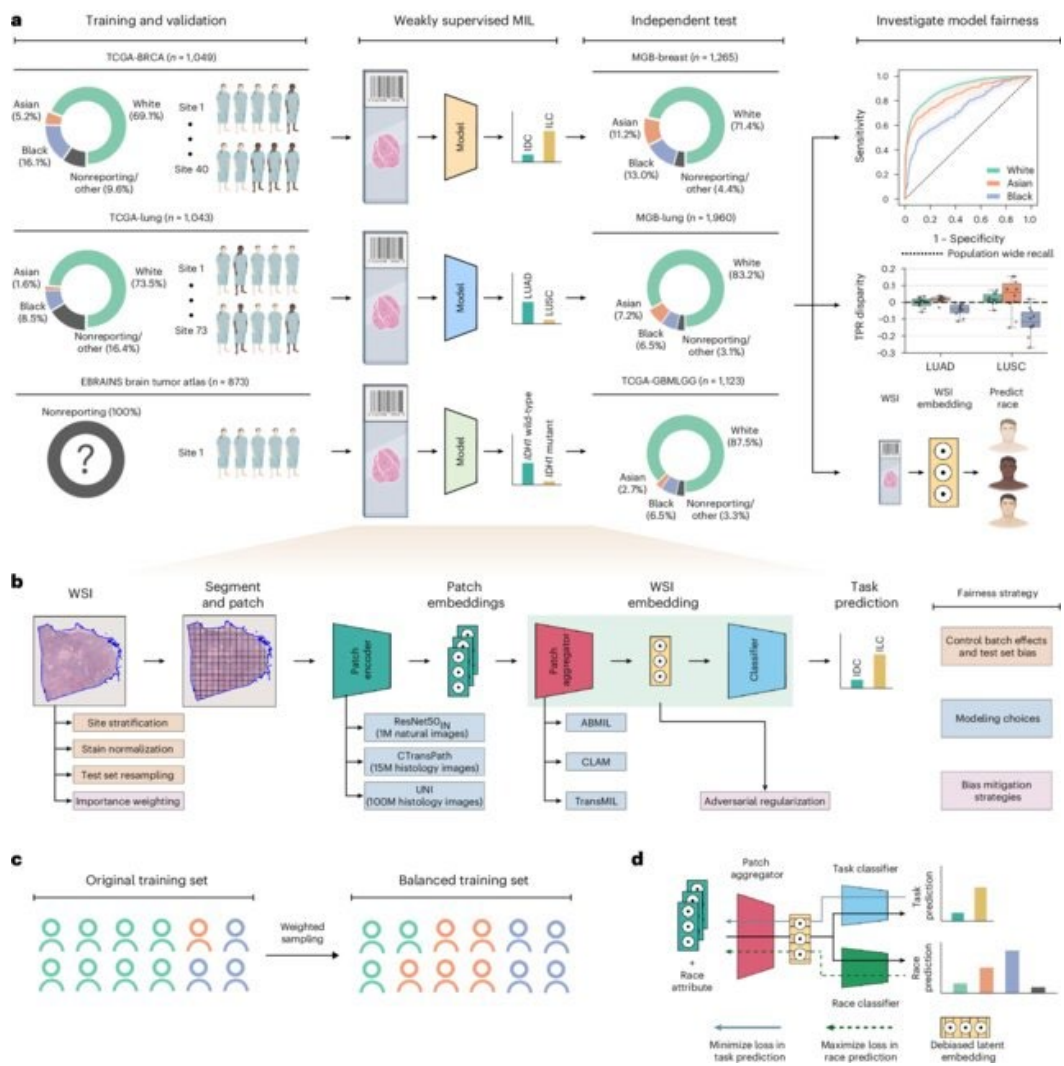
CIFAR-10 given label: cat

Cleanlab guessed: frog

MTurk consensus: frog

Labelerrors.com

Bias impacts training and reference data sets



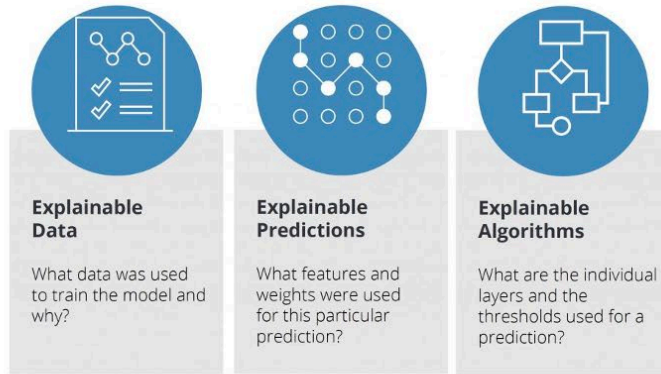
- Computational pathology has leveraged **large public datasets** that underrepresent certain demographic groups
- Whole-slide image classification models display **marked performance disparities** across different demographic groups

Vaidya et al 2024

Is the Path Clear?

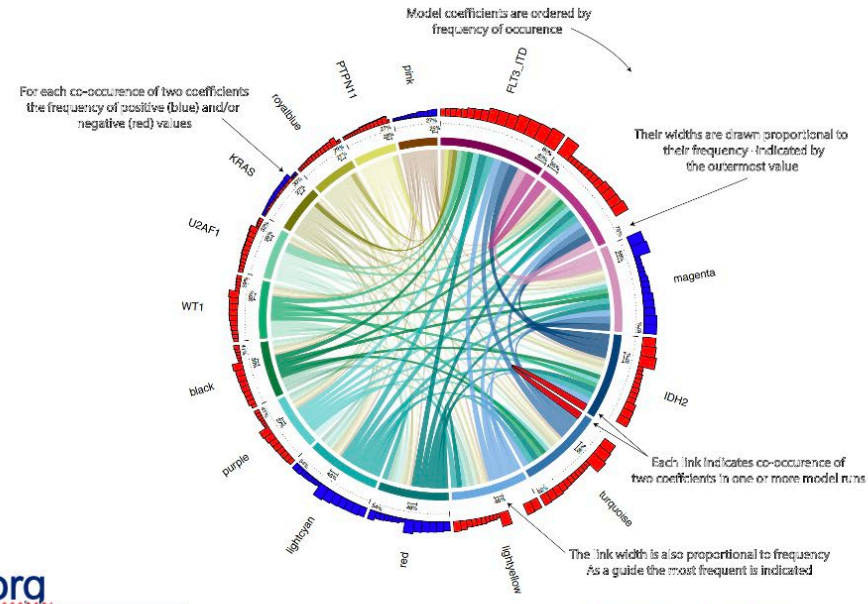


Do we understand what is happening in the model?

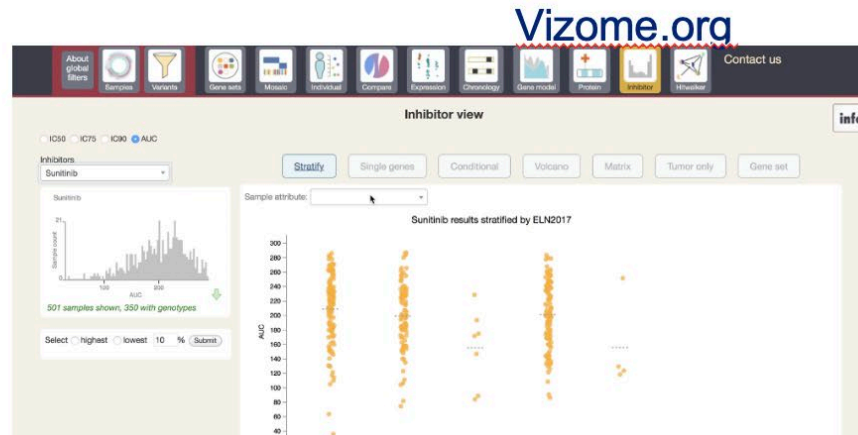


Questions around AI explainability help us understand how data, predictions and algorithms influence decisions.

Hodler (2019)

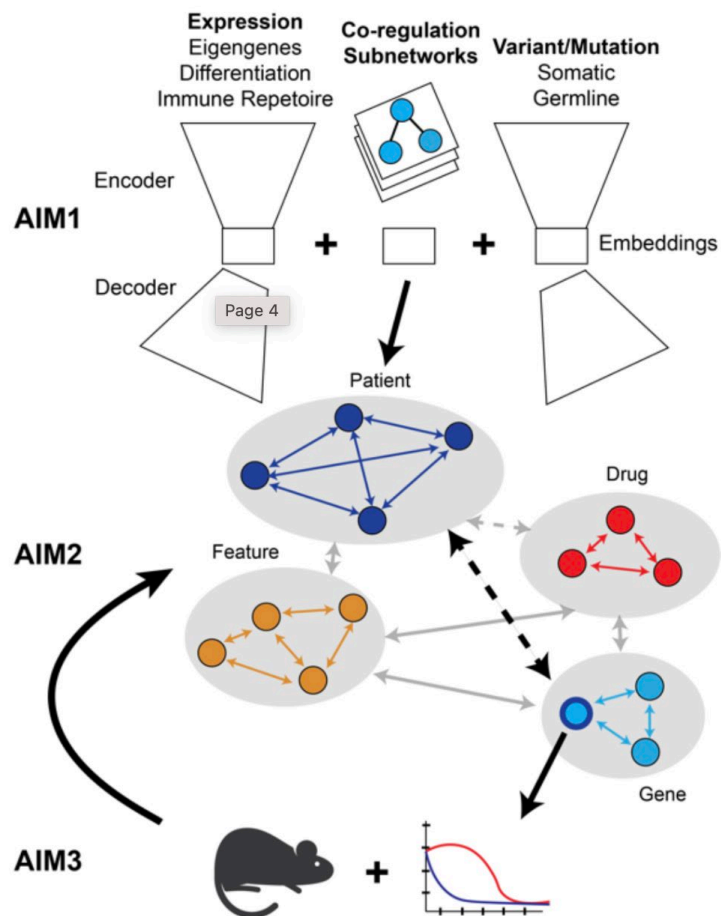


Tyner et al (2018)



*Cognitive
accessibility
fuels discovery*

Explainable Integrative Modeling of Myelodysplastic Syndromes



Focus: Identification of relevant therapeutics for MDS patients by linking patient similarity networks with predicted drug-gene interactions to prioritize gene targets (patient-gene edges)

Incorporation of xAI to maintain transparency given complex models



Together we build bridges



Artist: Lorenzo Quinn

What is AI/ML Ready Data?

- **Structured and Labeled Data**

Data that is organized in a format that can be easily interpreted and processed by AI/ML algorithms, with clear labels and annotations.

- **Quality and Consistency**

Data that is free from errors, biases, and inconsistencies, ensuring reliable and accurate model training and predictions.

- **Sufficient Quantity**

Data that is available in sufficient volume to train robust and generalizable AI/ML models, without overfitting or underfitting.

- **Relevant and Representative**

Data that is relevant to the problem being solved and representative of the real-world scenarios the model will encounter.

- **Accessible and Secure**

Data that is easily accessible and shareable with the necessary security and privacy measures in place.

Why Does it Matter?

- **Enhanced Model Accuracy**

High-quality, well-structured data enables AI/ML models to learn more effectively, leading to more accurate and reliable predictions for cancer diagnosis, prognosis, and treatment recommendations.

- **Improved Clinical Decision-Making**

AI/ML models trained on AI/ML-ready data can provide healthcare professionals with more reliable and actionable insights, supporting better-informed clinical decisions and enhancing patient outcomes.

- **Accelerated Research and Innovation**

AI/ML-ready data facilitates the rapid development and testing of novel cancer therapies, diagnostic tools, and other innovations, driving advancements in cancer research and care.

- **Timely, Precise Treatment**

AI/ML models leveraging high-quality, diverse data can identify precision oncology treatment approaches, leading to more targeted and effective cancer therapies for individual patients.

- **Efficient Resource Allocation**

Improved data quality and AI/ML readiness can help optimize the use of limited resources, such as research funding, clinical trial participants, and healthcare system capacity, leading to more efficient and cost-effective cancer care.

How do we make our data AI/ML Ready? (Part 1)

- **Structured and Labeled Data**

Data that is organized in a format that can be easily interpreted and processed by AI/ML algorithms, with clear labels and annotations.

- **Quality and Consistency**

Data that is free from errors, biases, and inconsistencies, ensuring reliable and accurate model training and predictions.

- **Sufficient Quantity**

Data that is available in sufficient volume to train robust and generalizable AI/ML models, without overfitting or underfitting.

- **Relevant and Representative**

Data that is relevant to the problem being solved and representative of the real-world scenarios the model will encounter.

- **Accessible and Secure**

Data that is easily accessible and shareable with the necessary security and privacy measures in place.

Data Documentation tools

- Data focused
- Included release of template for repeated use, adoption and adaption
- Unclear utility for biomedical research data
- Lack of consensus across templates

datasheets

data cards

health sheets

data statements

crowdworksheets

Data Shards

- Mapping across templates to establish consensus and remove redundancy
- New R package that can generate the Data Shards dataset summary for large-scale cohorts in GDC

Functional Genomic Landscape of AML (BeatAML)		The implementation of targeted therapies for acute myeloid leukaemia (AML) has been challenging because of the complex mutational pattern within and across patients as well as a dearth of pharmacologic agents most mutational events. Collectively, we have generated a large functional genomic dataset that can be leveraged to address clinical, genomic, transcriptomic and functional analyses of the biology of AML.
General Information		
LINKS		DATASHARD AUTHORS
Raw Data (third parties): <ul style="list-style-type: none"> ◦ dbGaP ◦ Genomic Data Commons Processed Data: https://biodev.github.io/BeatAML2/ Documentation (If different): Bundled with Processed Data		<ul style="list-style-type: none"> • Shannon McWeeney (mcweeney@ohsu.edu) • Daniel Bottomly (bottomly@ohsu.edu)
VERSION INFORMATION		KEYWORDS
<ul style="list-style-type: none"> • Current Version: v3.0 • DOI: 10.5281/zenodo.10654808 • Release Date: 07/01/2022 • Last Updated: 02/13/2024 • License: CC-BY-4.0 		<ul style="list-style-type: none"> • hematologic malignancy • targeted therapy • genomics
EXTENSION MECHANISMS		
Contact Dataset Owners/Publishers for ways to contribute		
<i>Dataset Owners/Publishers</i>		
ORGANIZATION	CONTACT DETAILS (EMAIL; ORCHID)	
Knight Cancer Institute, OHSU	Dataset Contacts: <ul style="list-style-type: none"> ◦ Jeff Tyner (tynerj@ohsu.edu; 0000-0002-2133-0960) ◦ Shannon McWeeney (mcweeney@ohsu.edu; 0000-0001-8333-6607) ◦ Brian Druker (drukerb@ohsu.edu; 0000-0001-8331-8206) Website: none	
<i>Funding Sources</i>		
INSTITUTION(S)	FUNDING OR GRANT SUMMARY (IES)	
	See Acknowledgements: https://www.sciencedirect.com/science/article/pii/S1535610822003129#ack0010	



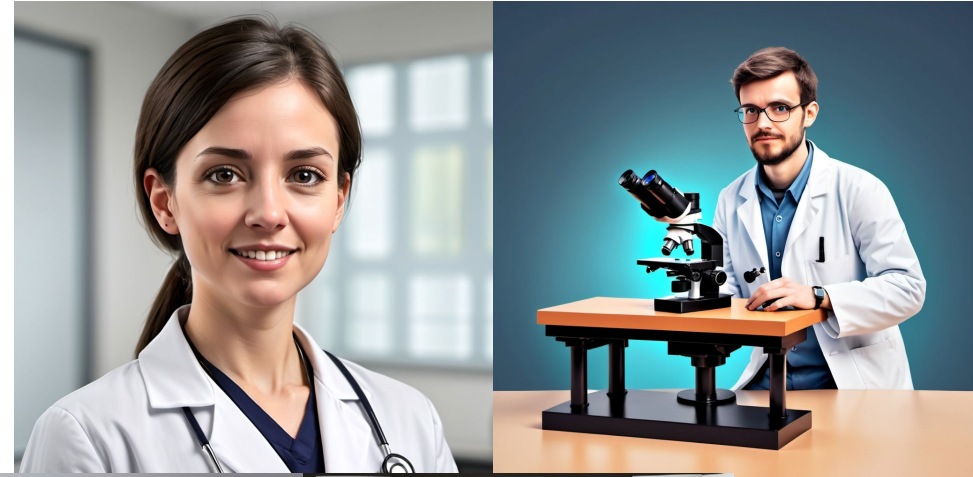
KNIGHT
CANCER
Institute



Sf(PM)

Relevance is Contextual

- Focus on personas – different roles in data and model life cycle
- Knowledge of lifecycle is heterogenous and biased
 - How data was generated vs what the impact is on downstream modeling



Dalle3 generated personas

How do we make our data AI/ML Ready? (Part 2)

- Structured and Labeled Data

Data that is organized in a format that can be easily interpreted and processed by AI/ML algorithms, with clear labels and annotations.

- Quality and Consistency

Data that is free from errors, biases, and inconsistencies, ensuring reliable and accurate model training and predictions.

- Sufficient Quantity

Data that is available in sufficient volume to train robust and generalizable AI/ML models, without overfitting or underfitting.

- Relevant and Representative

Data that is relevant to the problem being solved and representative of the real-world scenarios the model will encounter.

- Accessible and Secure

Data that is easily accessible and shareable with the necessary security and privacy measures in place.

Data Valuation with Gradient Similarity (DVGS)

Data valuation can identify mislabeled or noisy data

Filtering based on data values can improve analytics

Traditional data valuation is expensive

DVGS method is rapid, scalable and accurate

What is data valuation?

Class of algorithms that assign values to data quantifying the "usefulness" toward a given **predictive task**.

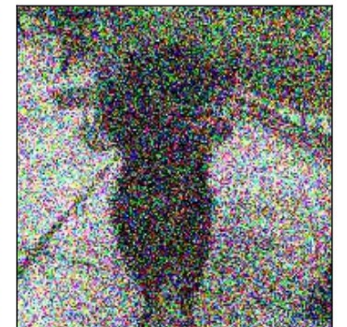
Noise Level = 0.1
Value = 0.00151



Noise Level = 0.3
Value = 0.00146



Noise Level = 0.5
Value = -0.00118



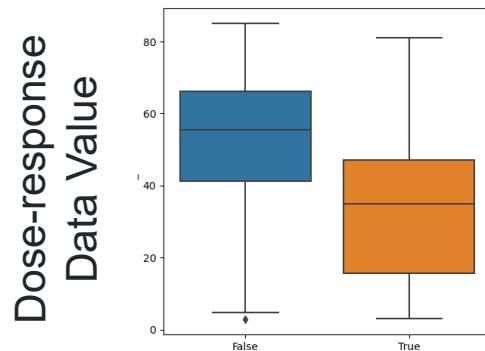
Ghorbani, Amirata and James Y. Zou. "Data Shapley: Equitable Valuation of Data for Machine Learning." ArXiv abs/1904.02868 (2019): n. pag.

Data valuation for dose-response drug response can significantly improve predictive performance

We applied DVGS to the BeatAML dose-response data

- X ~ cancer RNA expression
- Y ~ Area under the dose-response curve (AUC)

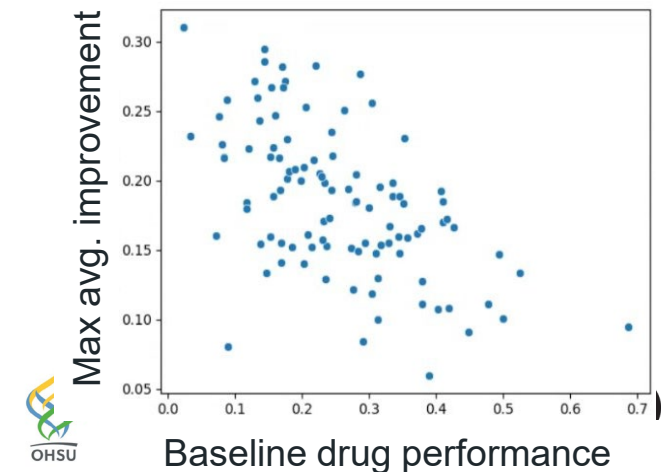
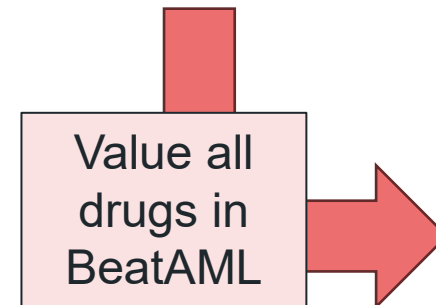
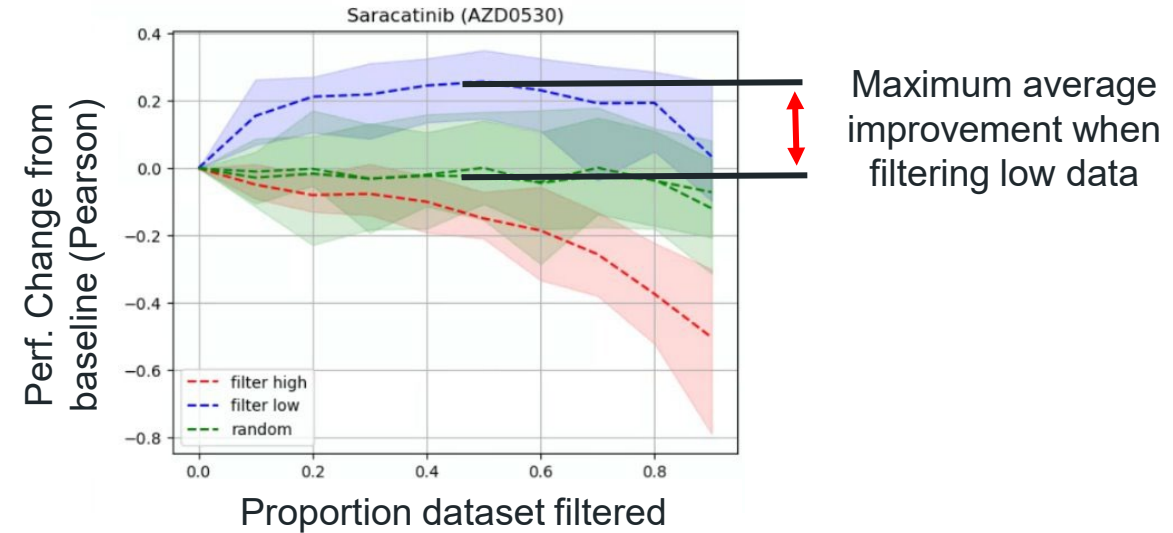
Filtering detrimental dose-response observations improves the performance on hold-out datasets

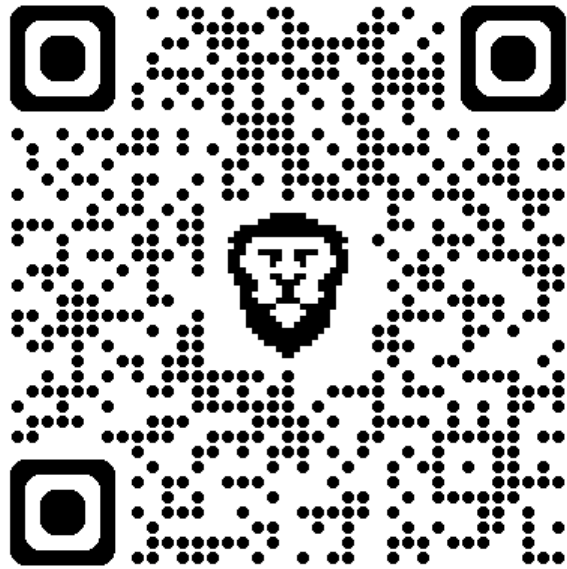


Data values correlate with known data quality metrics



All dose-response points are less than 50% viability





Evans, Nathaniel J. et al. "Data Valuation with Gradient Similarity." *ArXiv* (2024): n. pag.



Nathaniel Evans
Evansna@ohsu.edu

How do we make our data AI/ML Ready? (Part 3)

- Structured and Labeled Data

Data that is organized in a format that can be easily interpreted and processed by AI/ML algorithms, with clear labels and annotations.

- Quality and Consistency

Data that is free from errors, biases, and inconsistencies, ensuring reliable and accurate model training and predictions.

- Sufficient Quantity

Data that is available in sufficient volume to train robust and generalizable AI/ML models, without overfitting or underfitting.

- Relevant and Representative

Data that is relevant to the problem being solved and representative of the real-world scenarios the model will encounter.

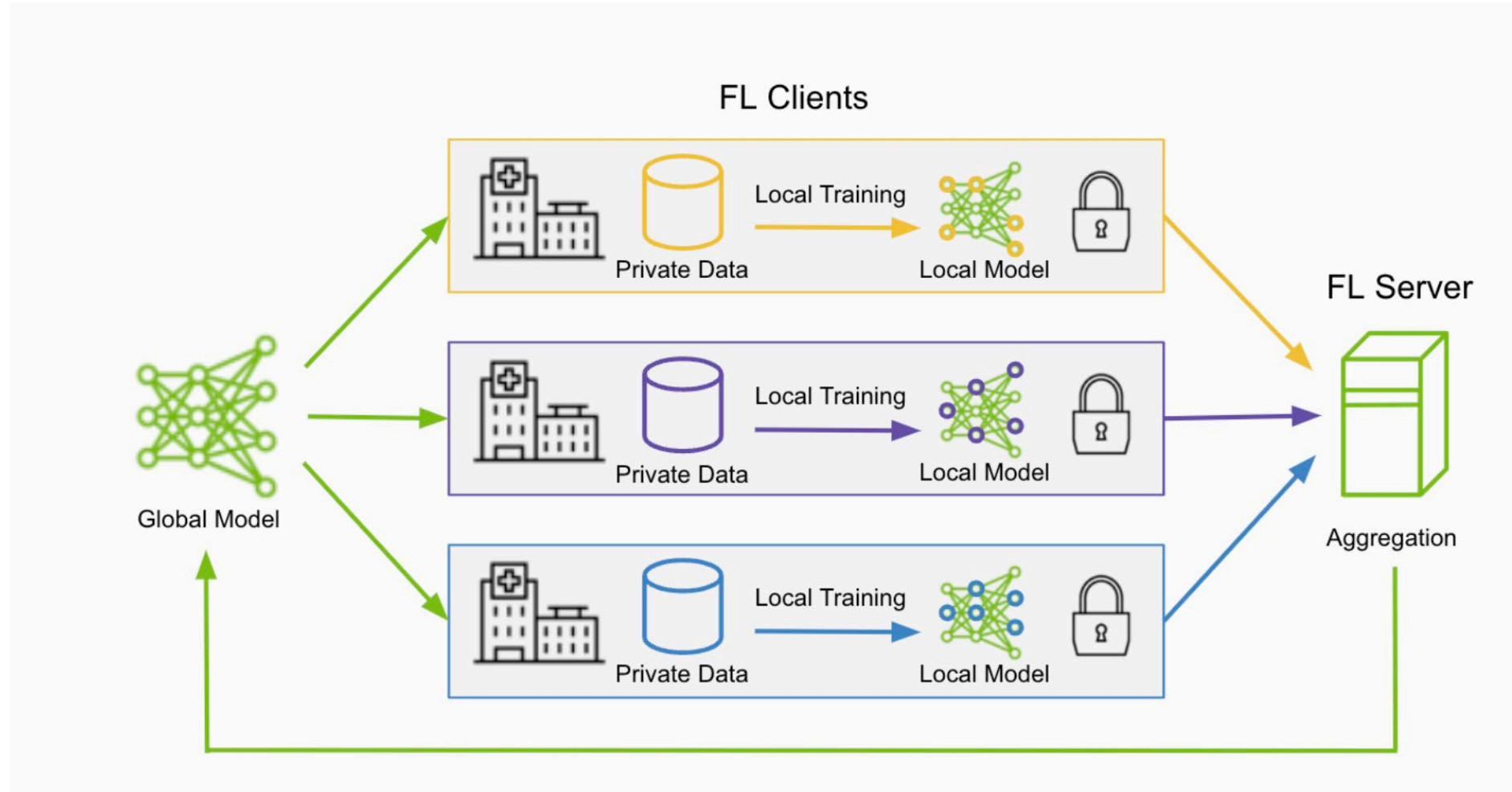
- Accessible and Secure

Data that is easily accessible and shareable with the necessary security and privacy measures in place.

OHSU-NCI Federated Learning Network Prototype



FL allows creation of collaborative, robust models without the data moving

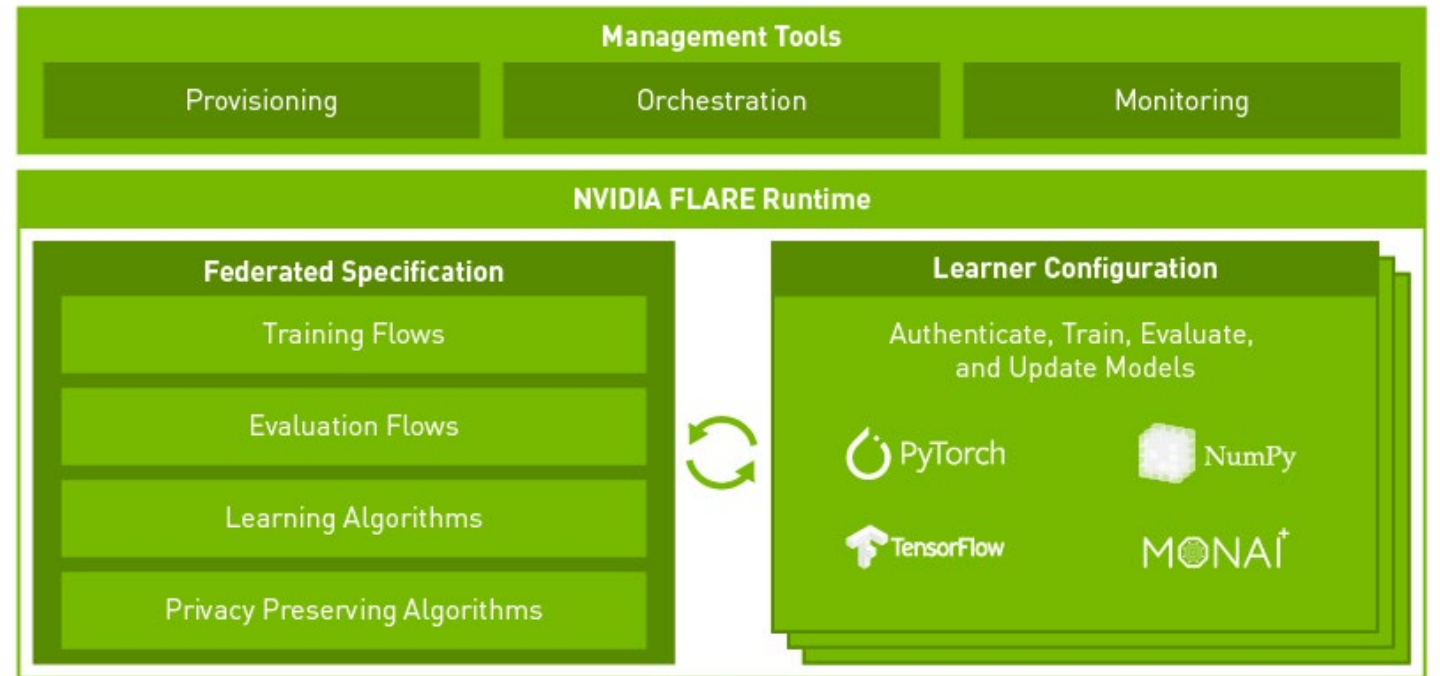


Distributed Multi-Party Collaboration
NVIDIA NVFlare provides FL Framework

NVFLARE

Leveraging existing partnership with NVIDIA:

- Open-Source Framework
- Privacy Preserving Algorithms
- Simulator for rapid prototyping
- Extensive management tools



FL Network Status

- NVFlare implemented
- improved provisioning to remove technical barriers
- Successful federation of global model with NCI
- Risk framework for CISOs to facilitate adoption in development

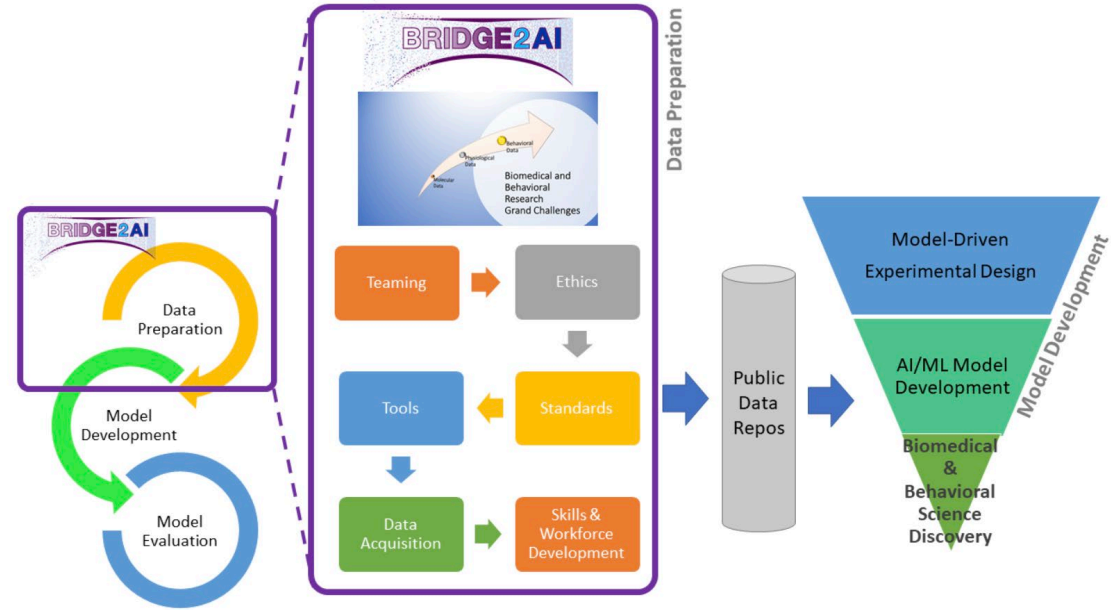


Dirk Petersen

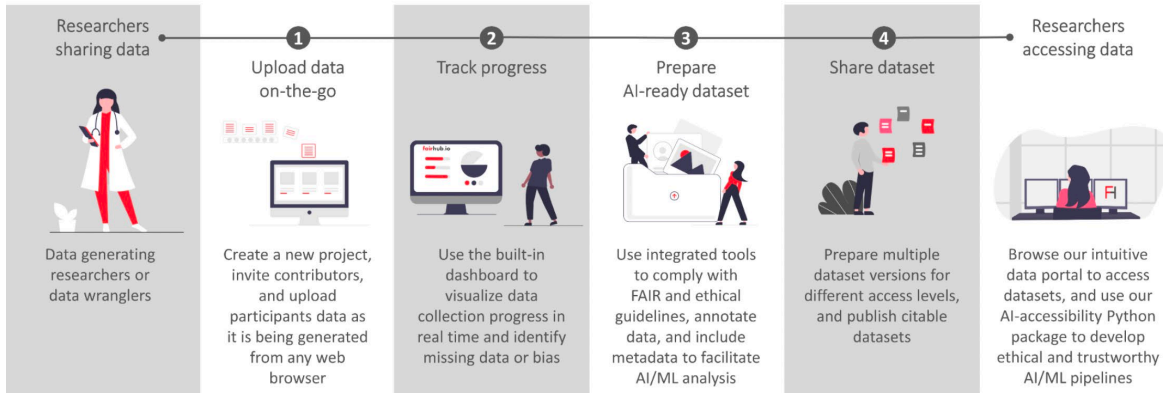


Daniel Bottomly

Democratizing access to tools and best practices



Open-source and free cloud-based platform for easily preparing, sharing, and accessing AI-ready datasets



Final Note: Role of Academic Medical Centers in AI Landscape

22 significant machine learning models produced by industry in 2022 compared with 3 produced by academia

Need to reframe this to leverage our strengths

- Retraining and fine tuning models on local data
- Focus on application and guidance
- Translational R&D catalyst

Value of “Know How”

DIGITAL HEALTH
INTELLIGENCE

Tuesday, July 11, 2023



Can academic medical centers compete in the AI arms race?

Tech leaders at academic medical centers say the private sector's dominance of AI talent is concerning.

Alone we can do so little; together we can do so much



Our Patients



Knight Precision Oncology



Knight Diagnostic Laboratory



Knight Data Operations

