# Explainable AI for health: where we are and how to move forward

Su-In Lee

Paul G. Allen Professor

Paul G. Allen School of Computer Science & Engineering

University of Washington, Seattle

# AI for bioMedical Sciences (AIMS) Lab

UW MSTP



**A  AI/ML**

*Explainable AI
(a.k.a. interpretable ML)*

**B  Basic biology**

*Identifying cause &
treatment of diseases*

**C  Clinical medicine**

*Developing & auditing
clinical AI models*

Nicasia
Beebe-Wang
(CSE PhD)

Ian Covert
(CSE PhD)

Su-In Lee (PI)

Hugh Chen
(CSE PhD)

Joe Janizek
(MSTP, CSE
PhD; matched
to Stanford)

Wei Qiu
(CSE PhD)

Chris Lin
(CSE PhD)

Ethan
Weinberger
(CSE PhD)

Alex DeGrave
(MSTP, CSE
PhD)

Mingyu Lu, MD
(CSE PhD)

Patrick Yu
(CSE PhD)
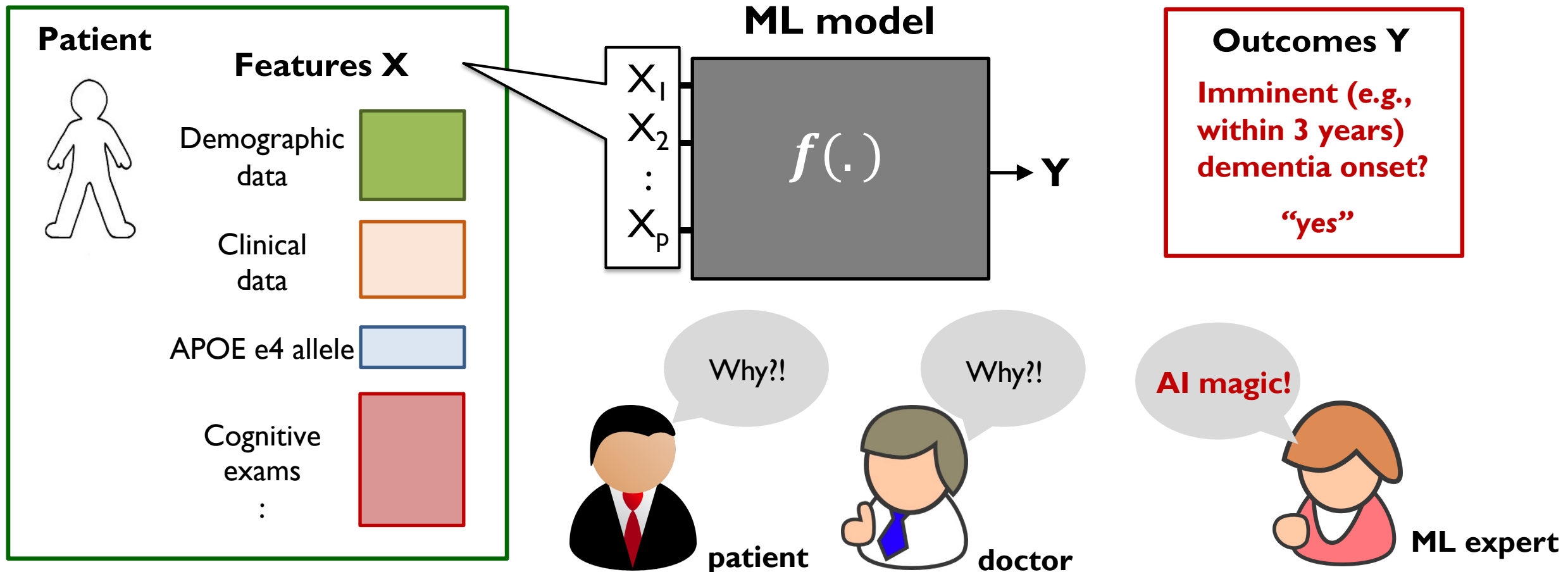
American Cancer Society

NIH

NSF

CZ

**Previous members:** Ben Logsdon (postdoc),
Safiye Celik (CSE PhD'18), Scott Lundberg
(CSE PhD'19), Parmita Mehta (CSE PhD'20),
Gabe Erion (MSTP, CSE PhD'21; now Harvard
Medical School for residency),

Chanwoo Kim
(CSE PhD)

Soham Gadgil
(CSE PhD)

# Outline – Two parts

- Part 1 – The significance of explainable AI in biomedical sciences
  - Demystifying the biological age
  - Unveiling neurodegenerative disease insights with explainable AI

- Part 2 – Advancing beyond explaining models
  - Cancer therapy design for precision oncology
  - Model auditing
  - Cost-aware clinical AI

# **Explainable AI (XAI):** Accurately predicting an outcome is vital, but the critical question revolves around *why*.



Lundberg et al. *Nature Machine Intelligence*, 2020 – Featured on the Cover    Beebe-Wang et al. *IEEE JBHI*, 2021
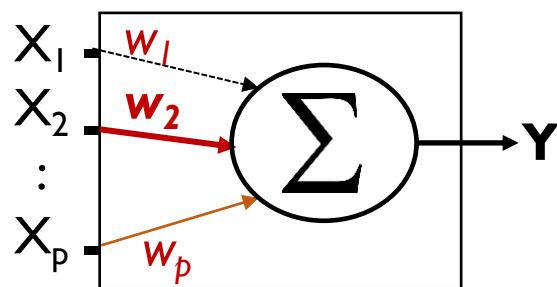
# Our solution is to fundamentally advance AI research to make a prediction *with explanations*

- Accuracy vs. interpretability
  - Simple models often lead to lower performance.
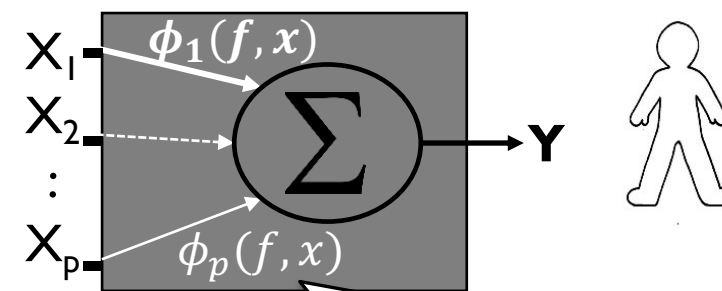  - Complex models are often considered to be a black box.

*Linear model*

*Complex model f (.)*

*Our approach, SHAP*
*(SHapley Additive exPlanations)*

**X:** Features  **Y:** Outcome

**Black Box**

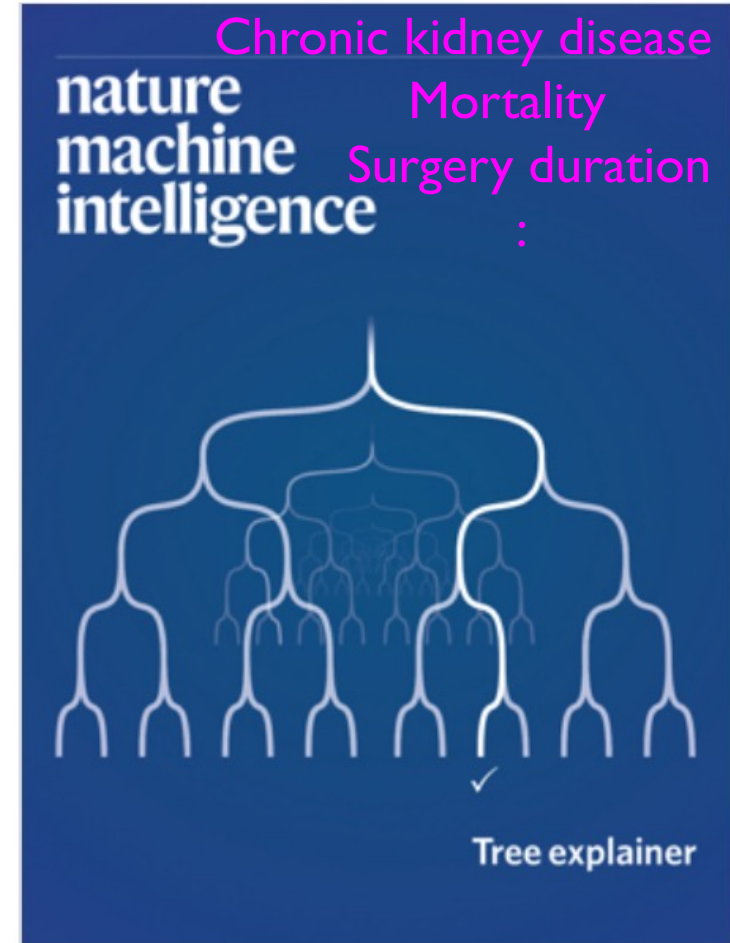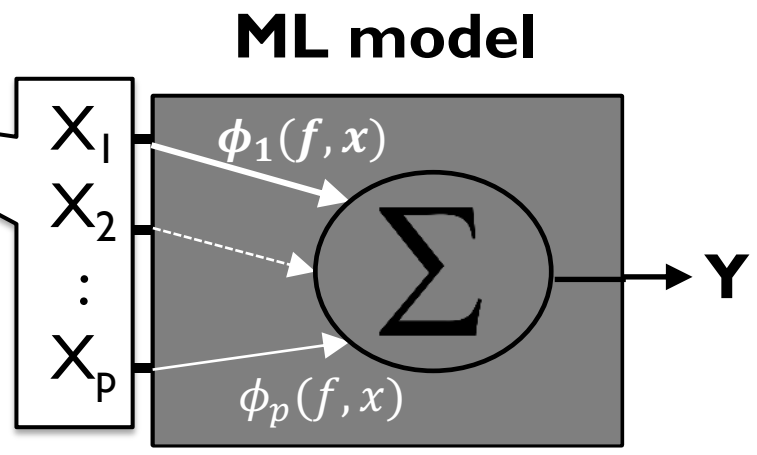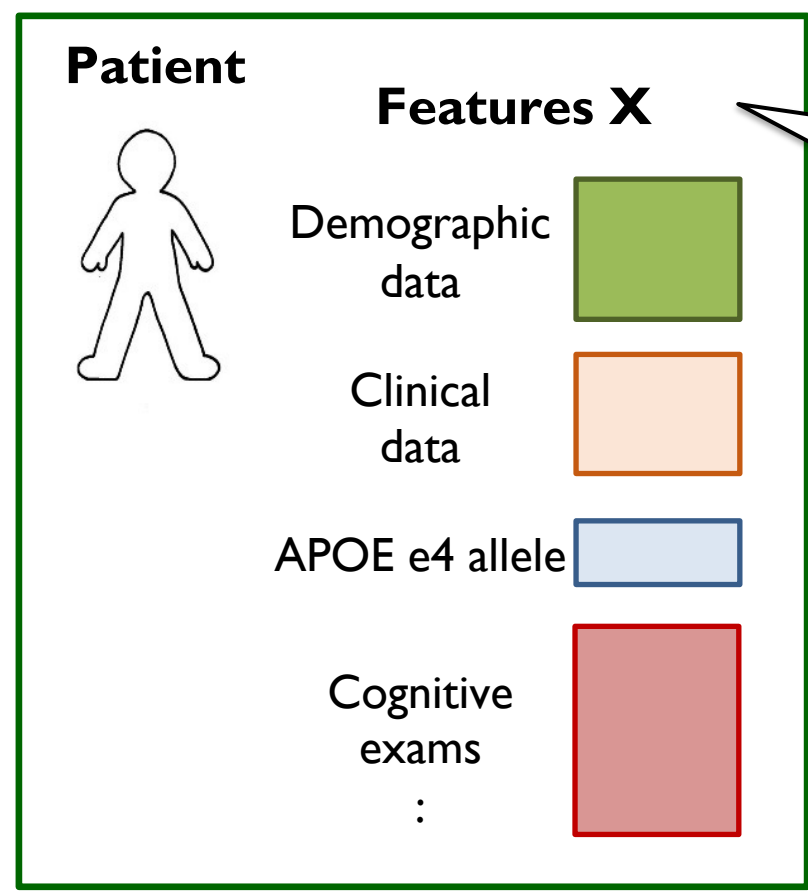For a particular prediction



SHAP can estimate feature importance for a particular prediction for any model.

Scott, CSE PhD'19

# Explainable AI (XAI): Accurately predicting an outcome is vital, but the critical question revolves around *why*.

**Patient**

**Features X**

Demographic data

Clinical data

APOE e4 allele

Cognitive exams
:

**ML model**

$X_1$
$X_2$
:
$X_p$

$\phi_1(f, x)$

$\sum$

$\phi_p(f, x)$

**Y**

Why?!    Why?!

patient    doctor

Chronic kidney disease
Mortality
Surgery duration
:

nature machine intelligence

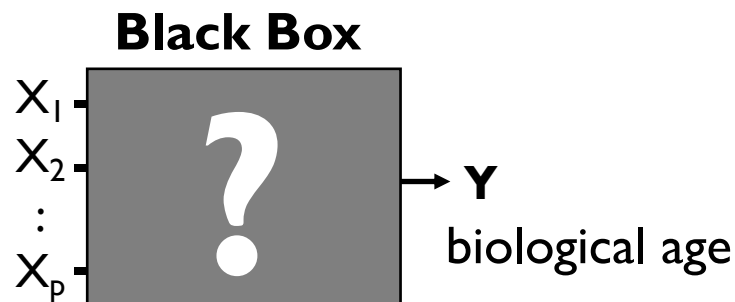Tree explainer

# XAI for interpretable biological age

- **ENABL (ExplaiNAble BioLogical) Age clock**
  - Estimates an individual's biological age
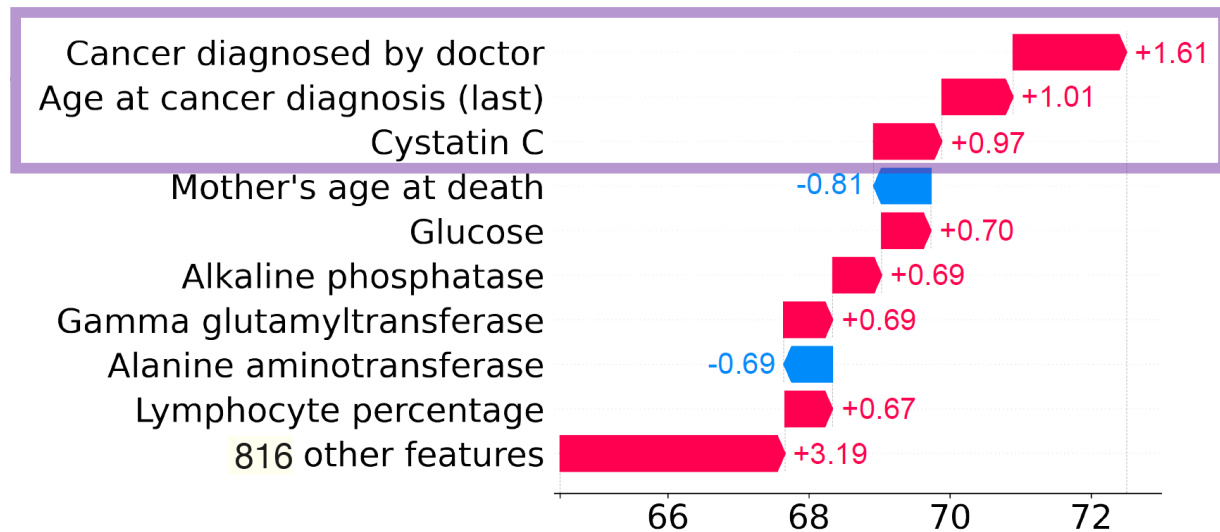  - Trained using the UK biobank data from 0.5M people based on 825 features:

CHRONOLOGICAL AGE
The actual age of a person. The age that is listed in your passport.

BIOLOGICAL AGE
It is the age of our cells. It tells our real age depending on how the aging process has affected us.

From first principles of movement

demographics
lab tests
exam results
lifestyle
:

**Black Box**

$X_1$
$X_2$
:
$X_p$

?

**Y**
biological age

Wei, CSE PhD

Qiu et al. *Nature Comm. Medicine*, 2022          Qiu et al. *Lancet Healthy Longevity*, 2023 – Featured on the Cover

# Explainable AI for interpretable biological age

Cancer diagnosed by doctor    +1.61
Age at cancer diagnosis (last)    +1.01
Cystatin C    +0.97
Mother's age at death    −0.81
Glucose    +0.70
Alkaline phosphatase    +0.69
Gamma glutamyltransferase    +0.69
Alanine aminotransferase    −0.69
Lymphocyte percentage    +0.67
816 other features    +3.19

66    68    70    72

**Chronological Age = 65**    **ENABL Age = 72.51**

**Impact of mortality causes on all-cause mortality**

Other
Digestive disease
Respiratory disease

**ENABL AgeAccel**

Circulatory disease

Neoplasms

Input feature

1.49

0    Contribution to ENABL
AgeAccel (year)

**Genome-wide association study**

demographics
lab tests
exam results
lifestyle
:

**Black Box**

$X_1$
$X_2$
:
$X_p$

**?**

**Y**
biological age

All-cause mortality **ENABL Age**

Neoplasm-cause mortality **ENABL Age**

Circulatory disease-cause mortality **ENABL Age**

Qiu et al. *Nature Comm. Medicine*, 2022         Qiu et al. *Lancet Healthy Longevity*, 2023 – Featured on the Cover

# ENABL age paper is now featured on the cover of Lancet Healthy Longevity.

- Please check it out!



Wei, CSE PhD

# Alzheimer's disease (AD)

- 6[th] most common cause of death in the US

- No long-term effective therapy exists to delay or prevent onset of progression

- AD lacks effective treatments due to limited understanding of *early cellular pathways* leading to end-stage pathologies like amyloid-β (Aβ) and tau.

**Amyloid-β (Aβ)**

**Tau**

**Healthy brain**

**Severe AD**
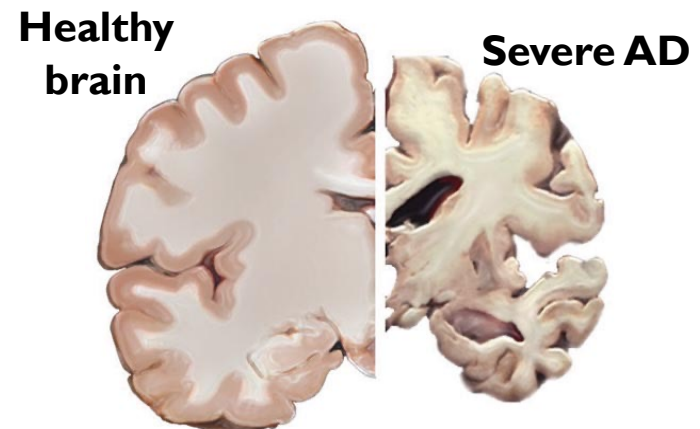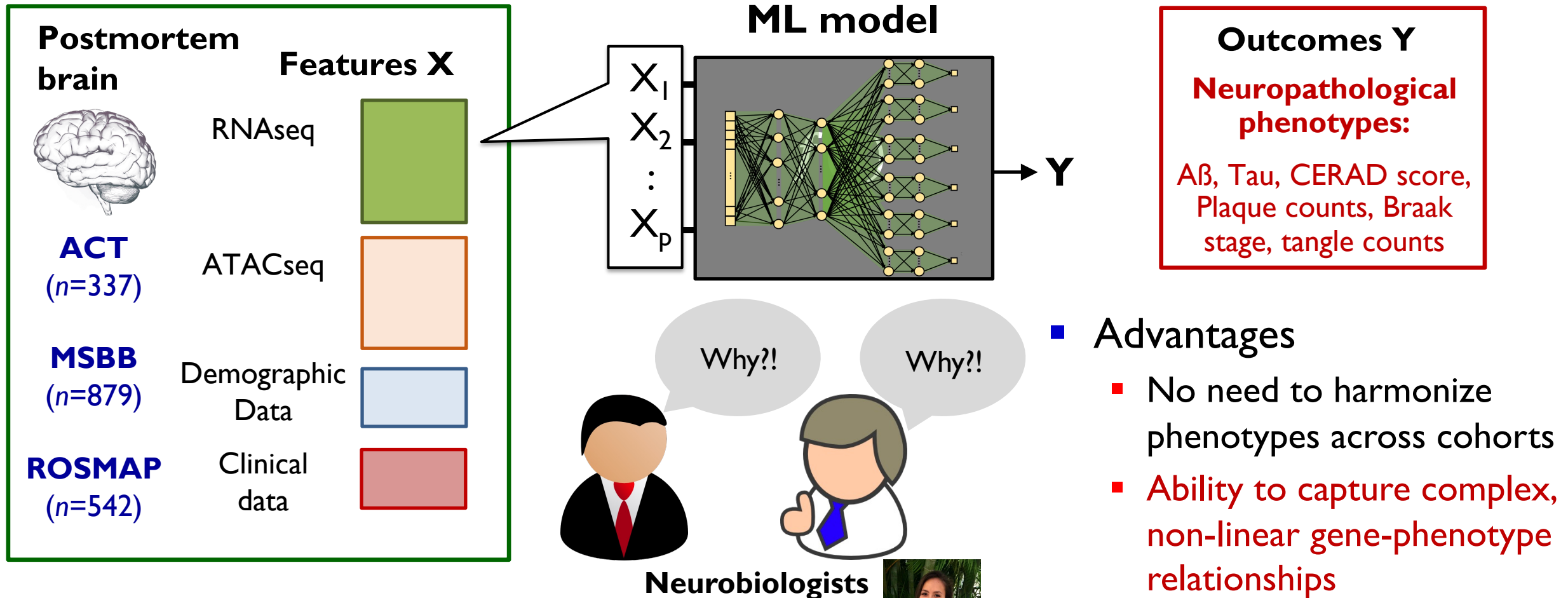
**B** Collaboration with Sara Mostafavi (UW)

# The key question is the *mechanistic explanation* of complex neuropathological phenotypes

Postmortem brain

Features X

**ACT** (*n*=337)

**MSBB** (*n*=879)

**ROSMAP** (*n*=542)

RNAseq

ATACseq

Demographic Data

Clinical data

$X_1$
$X_2$
$\vdots$
$X_p$

**ML model**

Y

**Outcomes Y**

**Neuropathological phenotypes:**

Aß, Tau, CERAD score, Plaque counts, Braak stage, tangle counts

Why?!     Why?!

**Neurobiologists**

Nicasia, CSE PhD'22

- Advantages
  - No need to harmonize phenotypes across cohorts
  - Ability to capture complex, non-linear gene-phenotype relationships

Beebe-Wang et al. *Nature Communications*, 2021

# Explainable AI (XAI) enhances neurodegenerative disease research in multiple ways

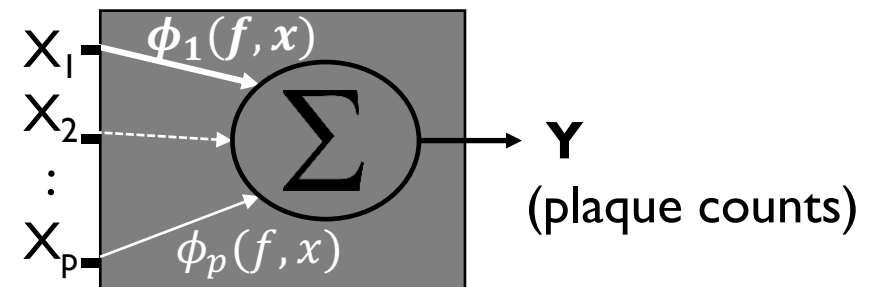- Our robust model trained across multiple cohorts was successfully validated, even in mouse brain and human blood datasets



20k genes expression levels

CERAD Score

Plaque Counts

Aβ Density

Braak Stage

Tangle Counts

Tau Density

- Using XAI, we can estimate each gene's contribution to AD neuropathologies
  - Previously unknown sex-specific associations btw. immune response genes and AD neuropathologies

$X_1 \xrightarrow{\phi_1(f,x)}$

$X_2 \dashrightarrow$

$\vdots$

$X_p \xrightarrow{\phi_p(f,x)}$

$\Sigma \longrightarrow \mathbf{Y}$

(plaque counts)

Beebe-Wang et al. *Nature Communications*, 2021          Janizek et al. *Nature Biomedical Engineering,* 2023

**B**

# Explainable AI (X~~~~

## disease research~~~~

> **XAI may capture patterns related to *sex-differential microglia activity*.**



● Female  ● Male

...orts was successfully validated, even in

- Using XAI, we can estimate each gene's contribution to AD neuropathologies
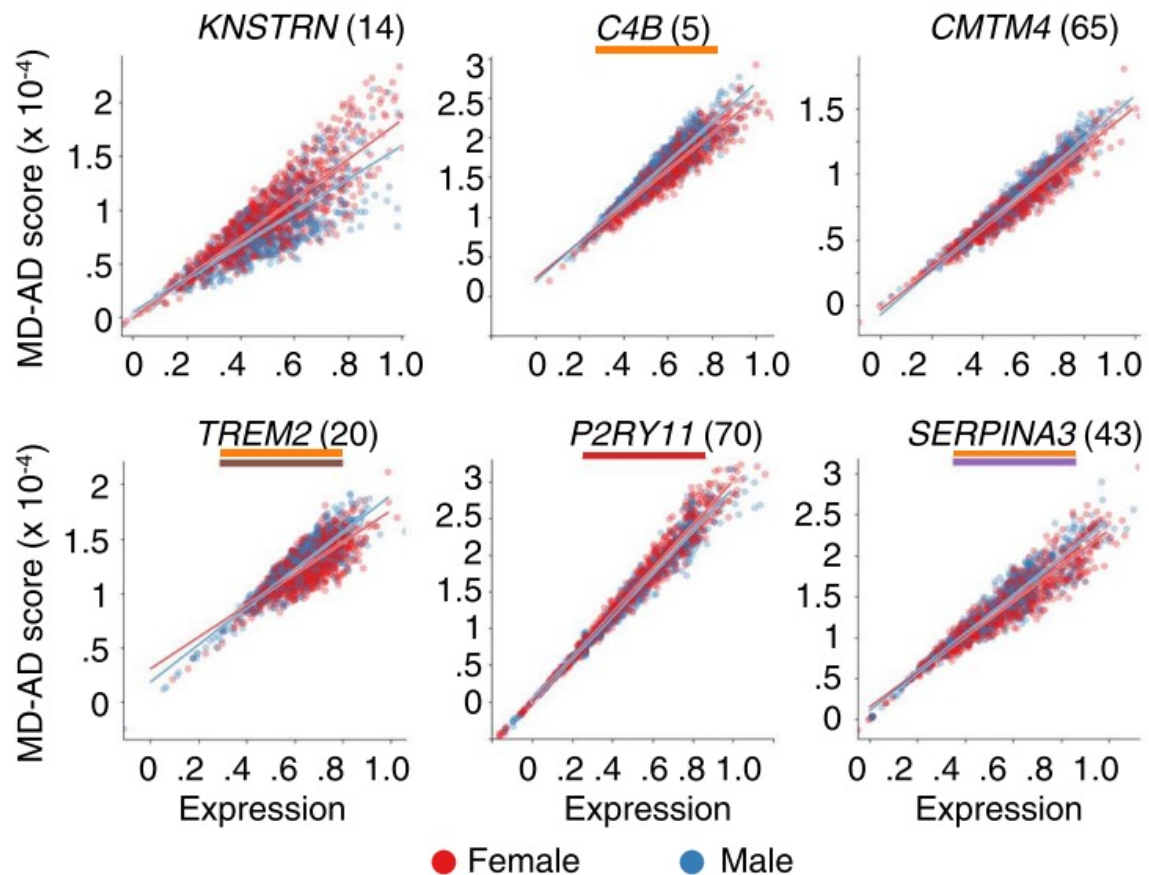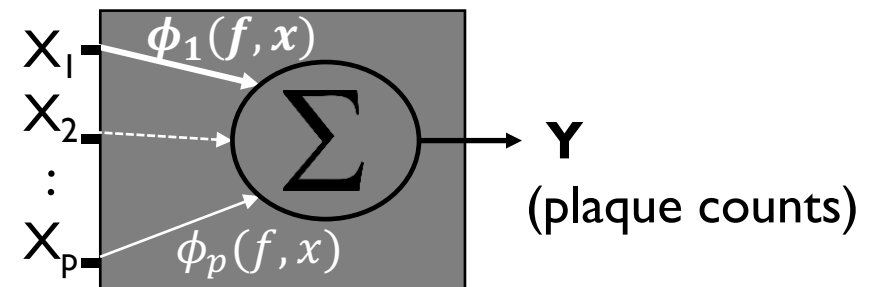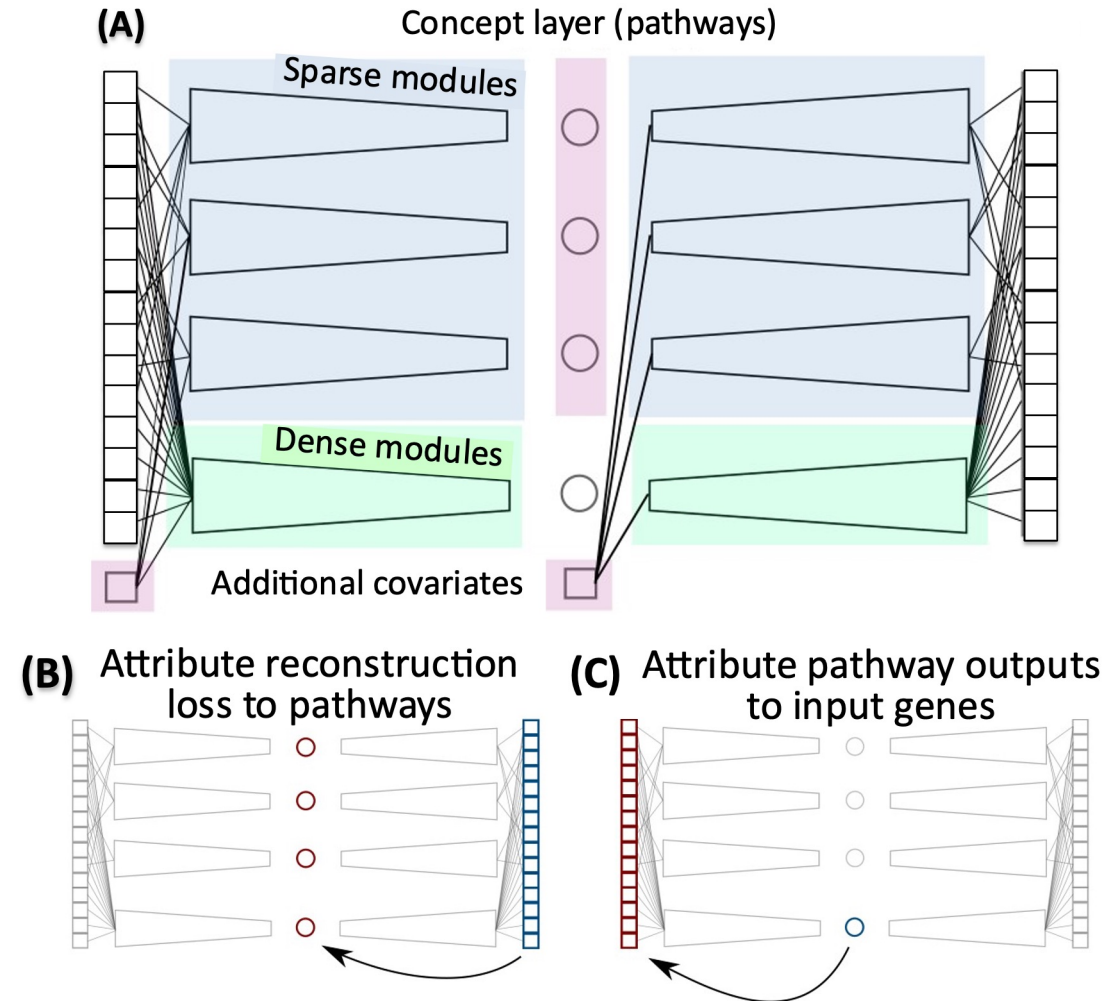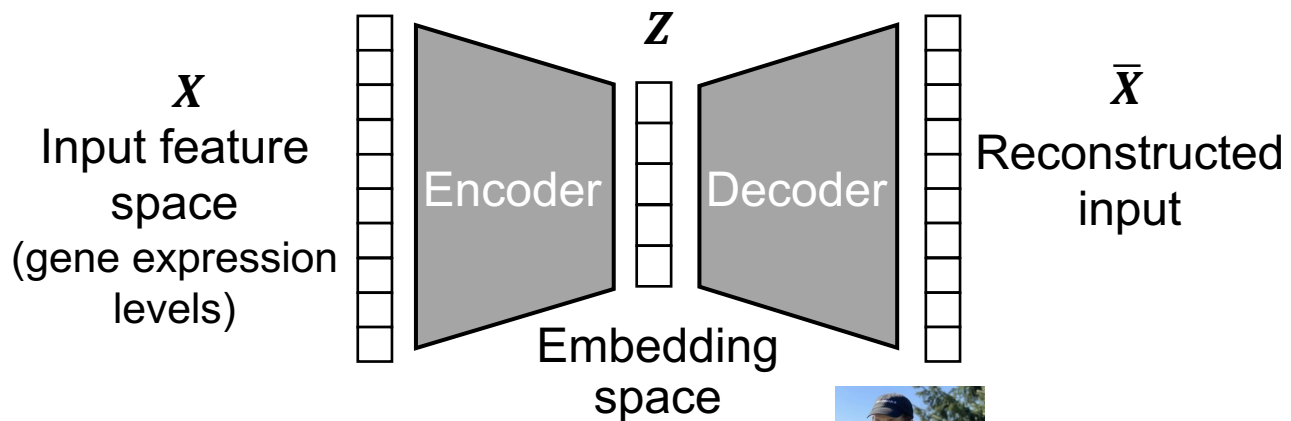  - Previously unknown sex-specific associations btw. immune response genes and AD neuropathologies



$X_1$ → $\phi_1(f, x)$

$X_2$

⋮

$X_p$ → $\phi_p(f, x)$

$\sum$ → **Y** (plaque counts)

Beebe-Wang et al. *Nature Communications*, 2021     Janizek et al. *Nature Biomedical Engineering,* 2023

# Biologically interpretable AI modeling further advances data-driven discovery

- Individual genes are not as interpretable as functional units (e.g., pathway)

- Unsupervised modeling enables the incorporation of unlabeled data
  - XAI can pinpoint crucial genes that explain the expression variation within the dataset

$X$
Input feature space
(gene expression levels)

$Z$

Encoder Decoder

Embedding space

$\bar{X}$
Reconstructed input

**(A)** Concept layer (pathways)
Sparse modules

Dense modules

Additional covariates

**(B)** Attribute reconstruction loss to pathways

**(C)** Attribute pathway outputs to input genes

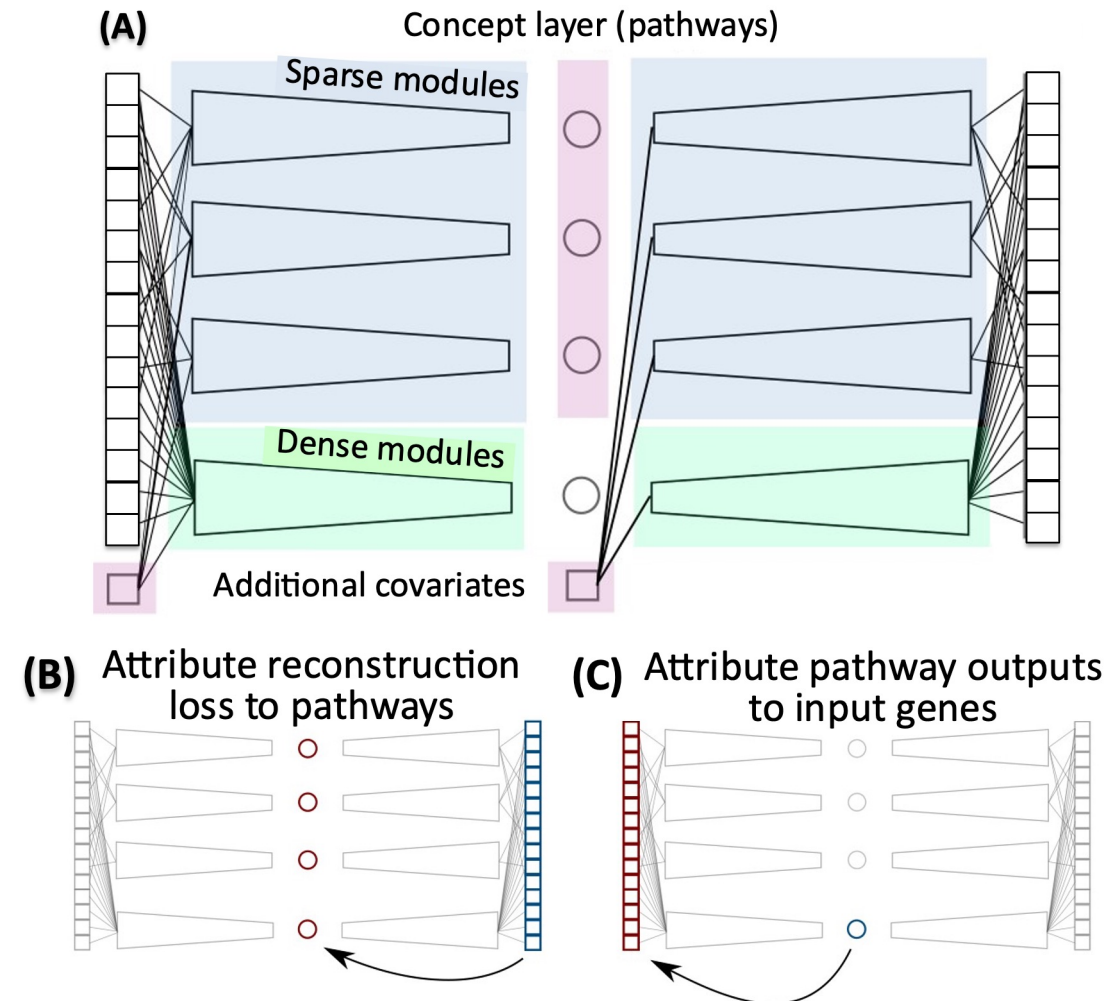Joe, UW MSTP/CSE PhD'22 (Matched in Radiology at Stanford)

# Biologically interpretable modeling identifies experimentally validated AD therapeutic targets

- We applied our approach to extended bulk RNAseq datasets from AD study cohorts

- We identified mitochondrial complex I as a potential mediator for tolerance to Aβ toxicity
  - *In vivo* validation in a transgenic *C. elegans* model expressing Aβ done by Matt Kaeberlein's lab

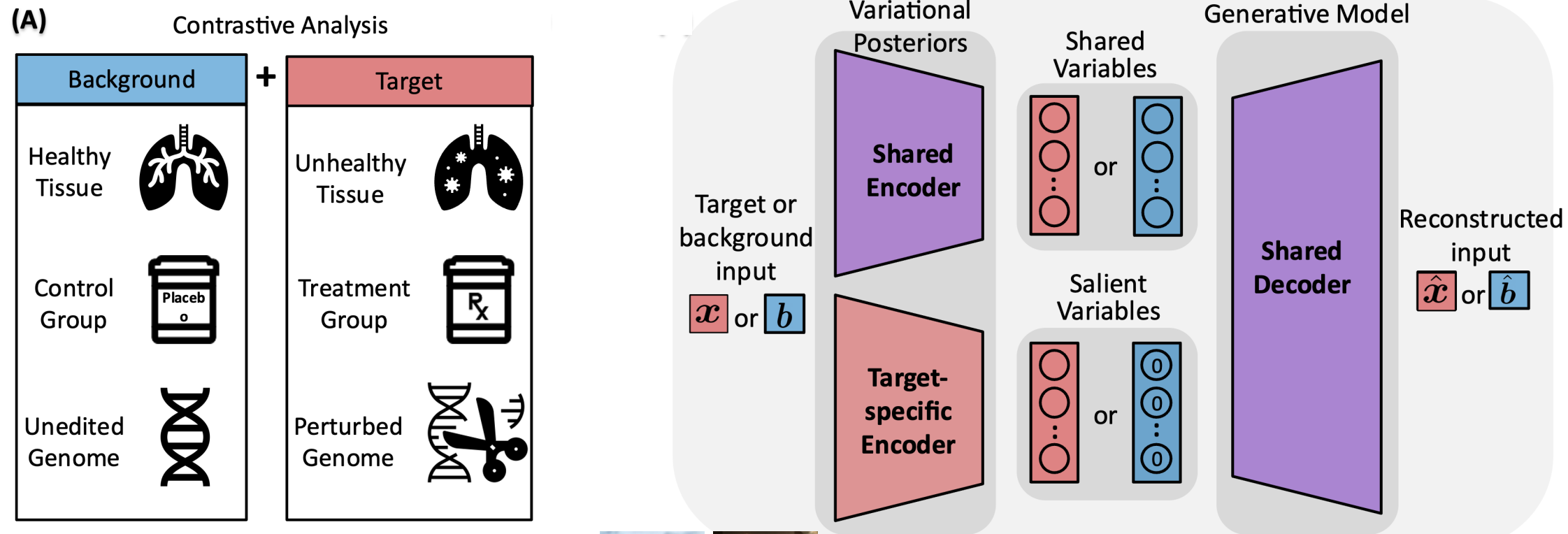**A promising pharmacological avenue!**

**Capsaicin**



**(A)** Concept layer (pathways)

Sparse modules

Dense modules

Additional covariates

**(B)** Attribute reconstruction loss to pathways

**(C)** Attribute pathway outputs to input genes

Janizek et al. *Genome Biology*, 2023

# Contrastive modeling enhances interpretability

- Single-cell datasets are often collected to investigate differences in cellular state between background cells and those under specific treatments
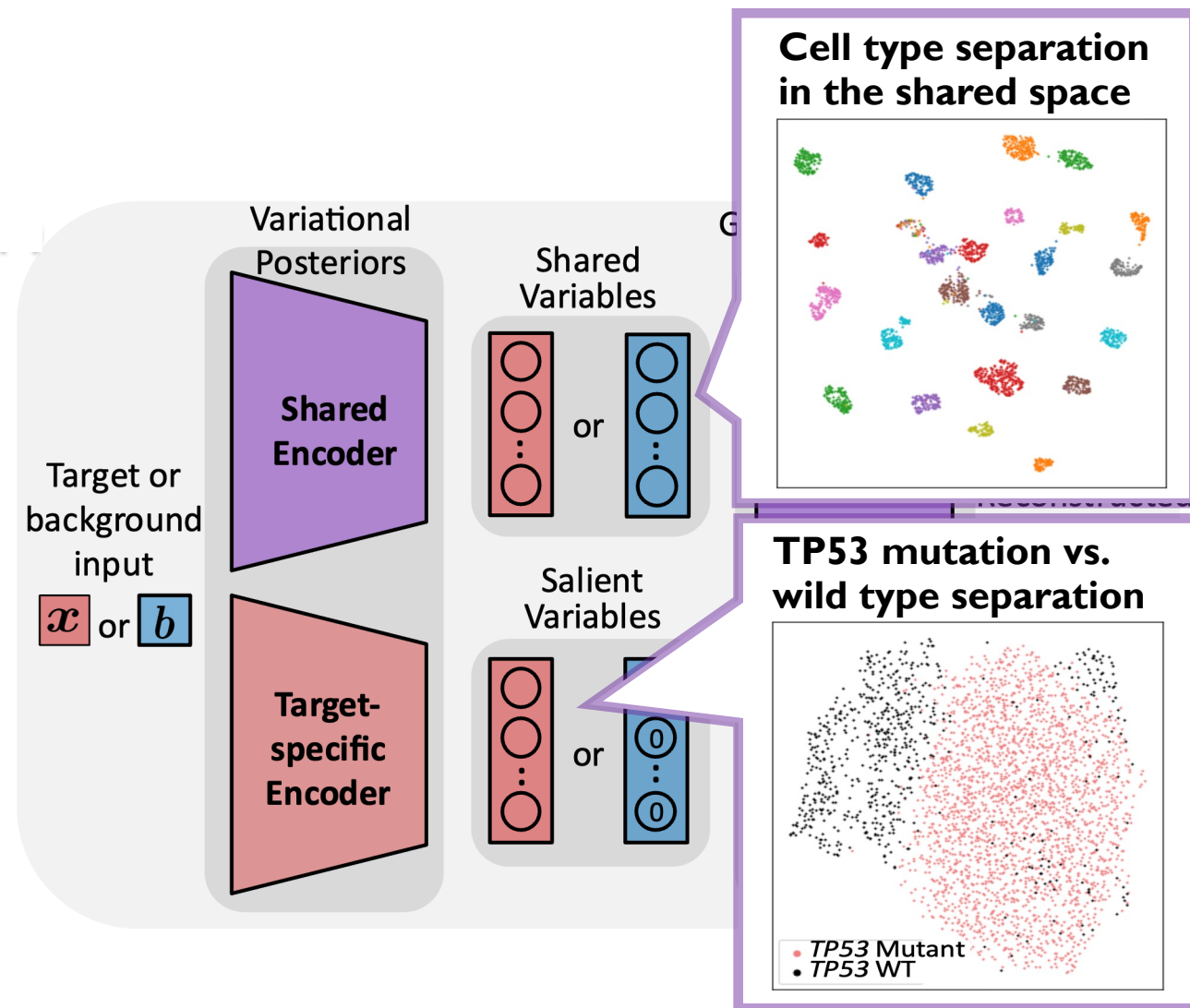
Ethan & Chris, CSE PhD

# Contrastive modeling enhances interpretability

- Cancer cells treated with idasanutlin *vs.* untreated as background
  - Cells behave differently in salient space depending on their TP53 mutation status

**Important implications for personalized medicine!**

- How about AD *vs.* control brain tissue?
  - What drives neurodegeneration (in collaboration with Jessica Young)
  - What drives biological aging process? (Jessica Young & Suman Jayadev)

Weinberger,* Lin,* and Lee. *Nature Methods,* 2023



Cell type separation in the shared space

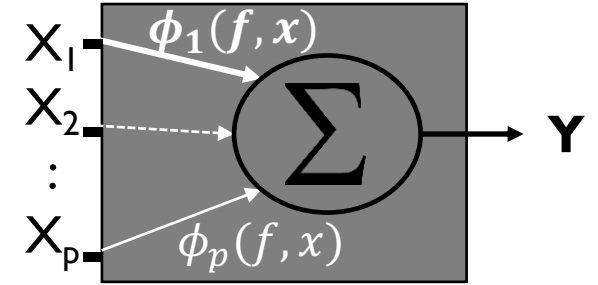TP53 mutation vs. wild type separation

- TP53 Mutant
- TP53 WT

# Outline – Two parts

- Part 1 – What explainable AI can do in biological research
  - Demystifying the biological age
  - Unveiling neurodegenerative disease insights with explainable AI

- Part II – Beyond explaining models
  - Cancer therapy design for precision oncology
  - Model auditing
  - Cost-aware clinical AI

# Beyond interpreting models...



— Cancer therapy design for precision oncology

[*Nature BME*'23]

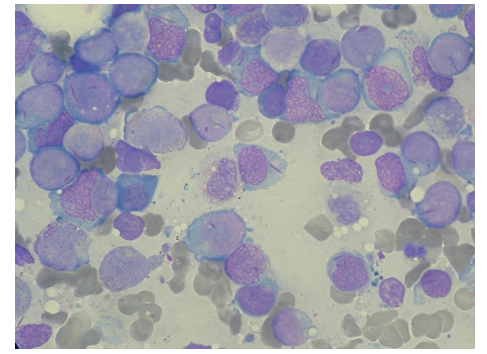— AI auditing [*Nature MI*'21, *Nature BME*'23, *Nature Medicine*'24]

radiology, dermatology

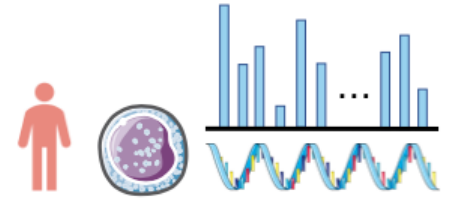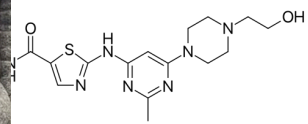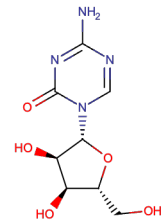— Cost-aware clinical AI [*Nature BME*'22]

emergency medicine, critical medicine

# Explainable AI to design cancer therapy

- Cancers are increasingly treated by combination therapy
  - Choosing drugs tha                                    ways
  - Greater efficacy
  - Fewer side-effects

- Choosing optimal c

  - Explanations to th                                        portant
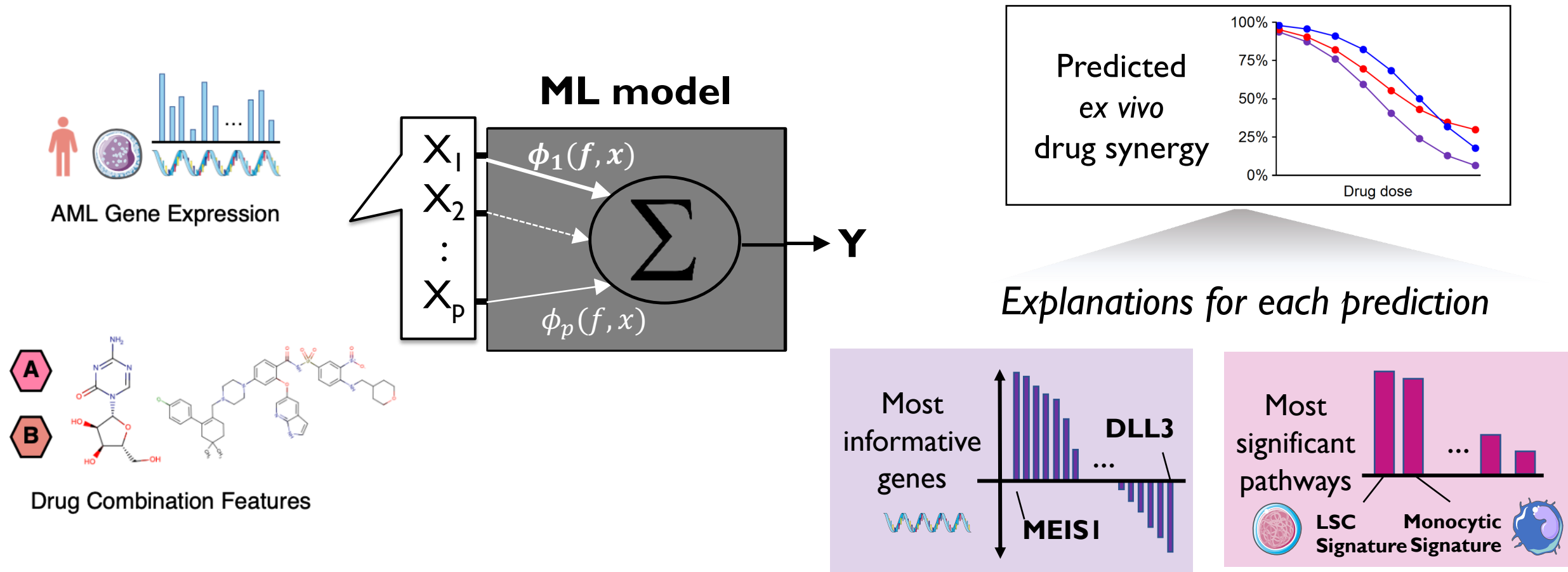
Hundreds of
individual drugs

MSTP/CSE PhD'22 (got matched
to Stanford Radiology)

AML Gene Expression

...irs of Drugs

Janizek et al. *Nature Biomedical Engineering*, 2023          Lee*, Celik*, et al. *Nature Comm.*, 2018

# Explainable AI to design cancer therapy

- EXPRESS: Explainable prediction of drug synergy



AML Gene Expression

Drug Combination Features

**ML model**

$$X_1 \quad \phi_1(f, x)$$
$$X_2$$
$$\vdots$$
$$X_p \quad \phi_p(f, x)$$

$$\Sigma \rightarrow \mathbf{Y}$$

Predicted *ex vivo* drug synergy

Drug dose

*Explanations for each prediction*

Most informative genes

DLL3
...
MEIS1

Most significant pathways

...

LSC Signature   Monocytic Signature

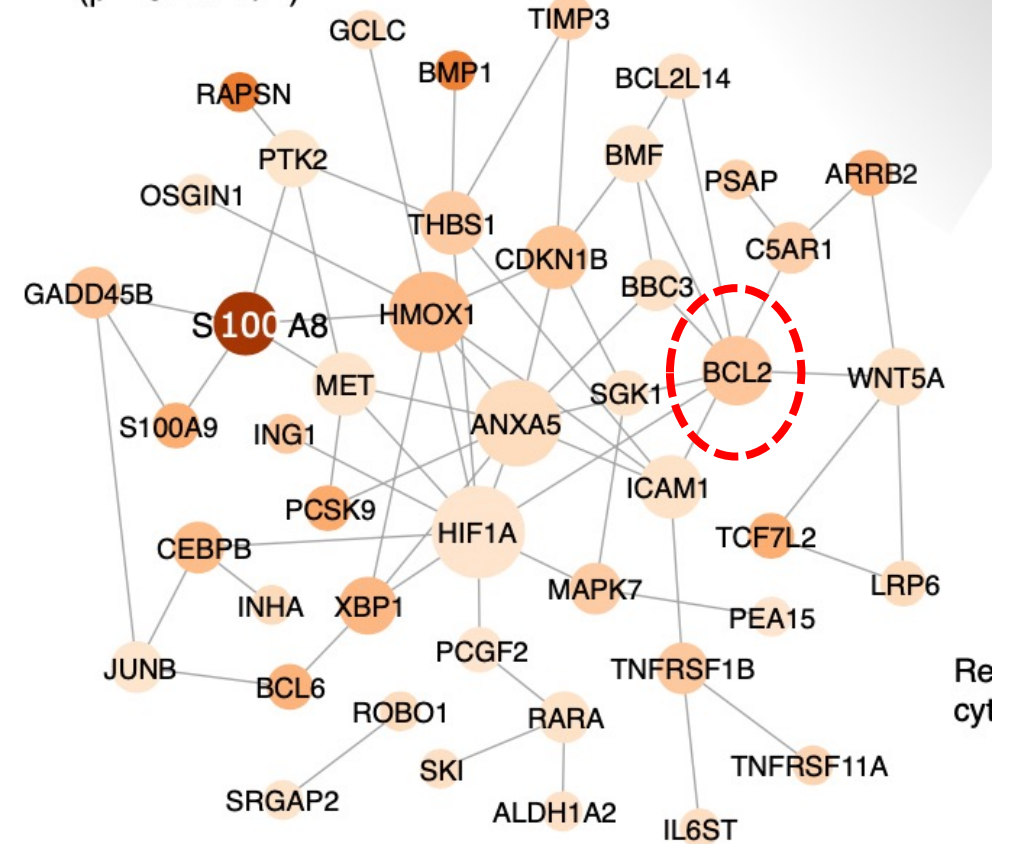Janizek et al. ***Nature Biomedical Engineering***, 2023       Lee*, Celik*, et al. ***Nature Comm.***, 2018

# Interpretability allows us to validate our model's decisions



Venetoclax, BCL-2 inhibitor



Regulation of cell death
$(p = 8.2 \times 10^{-3})$

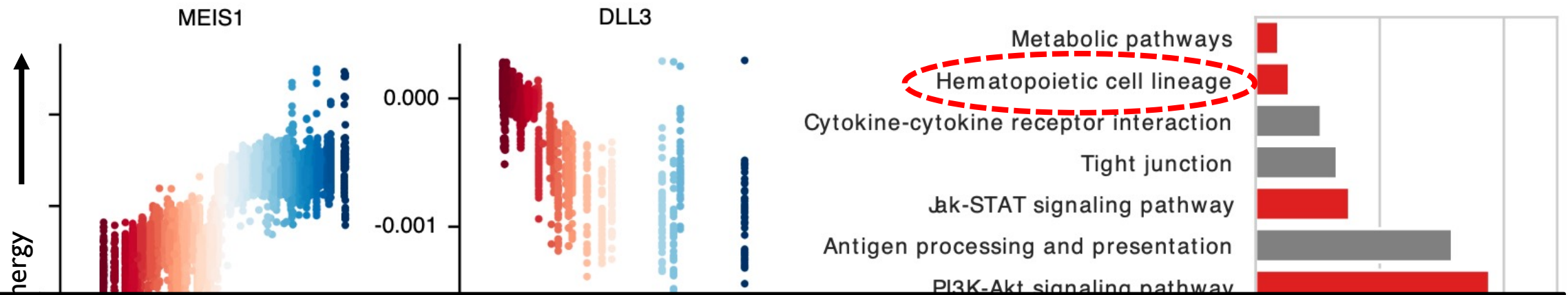Janizek et al. *Nature Biomedical Engineering*, 2023          Lee\*, Celik\*, et al. *Nature Comm.*, 2018

# Interpretability uncovers transcription programs underlying drug synergy



Linked to prognosis

Related to hematopoietic differentiation

Janizek et al. *Nature Biomedical Engineering*, 2023        Lee*, Celik*, et al. *Nature Comm.*, 2018

# Interpretability uncovers transcription programs underlying drug synergy



"Stemness" can be considered as an "axis" to design combination therapies – Two drugs that target different differentiation stages of cancer are likely effective.

Janizek et al. *Nature Biomedical Engineering*, 2023     Lee*, Celik*, et al. *Nature Comm.*, 2018

# Beyond interpreting models…



– Cancer therapy design for precision oncology

[*Nature BME*'23]

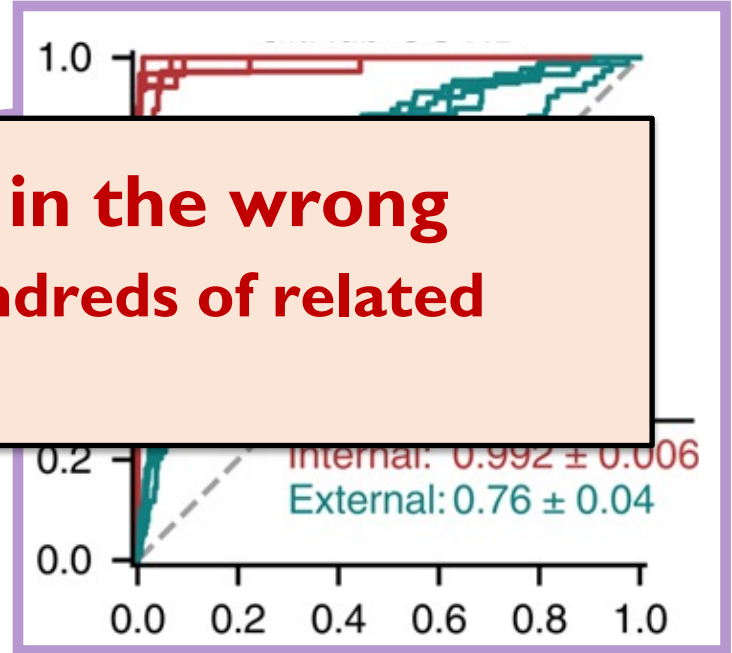– AI auditing [*Nature MI*'21, *Nature BME*'23, *Nature Medicine*'24]

radiology, dermatology

– Cost-aware clinical AI [*Nature BME*'22]

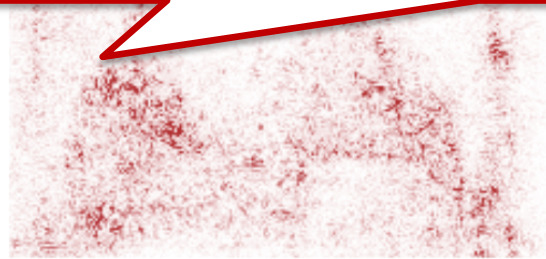emergency medicine, critical medicine

# Auditing AI for COVID-19 detection using XAI

- Many published AI models that detect COVID-19

**XAI helped us to stop the field from moving in the wrong direction – There were 6 published papers and hundreds of related models out there that learned the shortcuts.**

Many kinds of analyses for model auditing presented in the paper!

Internal: 0.992 ± 0.006
External: 0.76 ± 0.04

99th

0th

✓ Clear lung bases predict negative COVID-19 status

✗ laterality markers should not predict negative status

✗ medical devices should not predict negative status
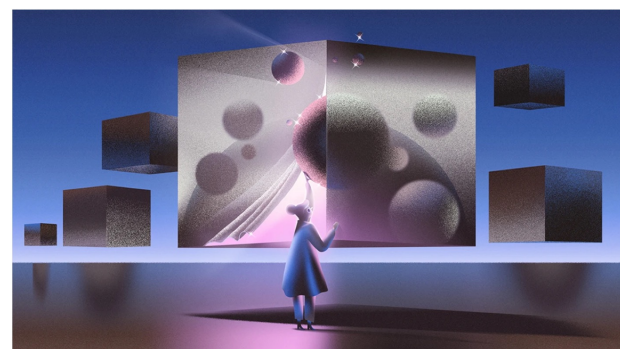
MSTP / CSE PhD

# Our AI auditing work featured in *Nature*

- "Breaking into the black box of artificial intelligence" *Nature* Outlook



**Breaking into the black box of artificial intelligence**

Scientists are finding ways to explain the inner workings of complex machine-learning models.

By Neil Savage

UW MSTP / CSE PhD
**Joe Janizek**
(residency at Stanford)

UW MSTP/CSE PhD student **Alex Degrave**

Alex DeGrave and Joseph Janizek are students on the Medical Scientist Training Program at the University of Washington, in Seattle.   Credit: Alex DeGrave

DeGrave,* Janizek* et al. ***Nature Machine Intelligence***, 2021  Cited 440+, Featured in *Nature*, 2022

**C**

# Further digging into the flaws in the reasoning processes of clinical AI – dermatology

- Auditing AI models to predict skin cancer
  - Five models – 2 academic models, 2 commercial devices, and 1 competition winner

- Technical challenges – saliency maps often do not work

**Original image**    **Saliency map**    **Modified image**



Predicted: benign

Predicted: malignant

## Our solution #1

- Generate counterfactual images *from the AI model*
- Systematic characterization by experts: Drs. Roxana Daneshjou, and Zhuo Ran Cai (Stanford)

DeGrave et al. (*Nature Biomedical Engineering*)

Kim et al. (*Nature Medicine*, 2024)

# How do dermatology AI systems make decisions on dermoscopic images?

Degrave, Ran Cai, Janizek, Daneshjou,* and Lee* *Nature Biomedical Engineering*, 2023

MSTP / CSE PhD

# The *Lancet* perspective (Feb 2024)

- Broader promises of counterfactual AI

*"The clinical potential of counterfactual AI"*

by Su-In Lee* and Eric Topol

---

## Digital medicine

## The clinical potential of counterfactual AI models
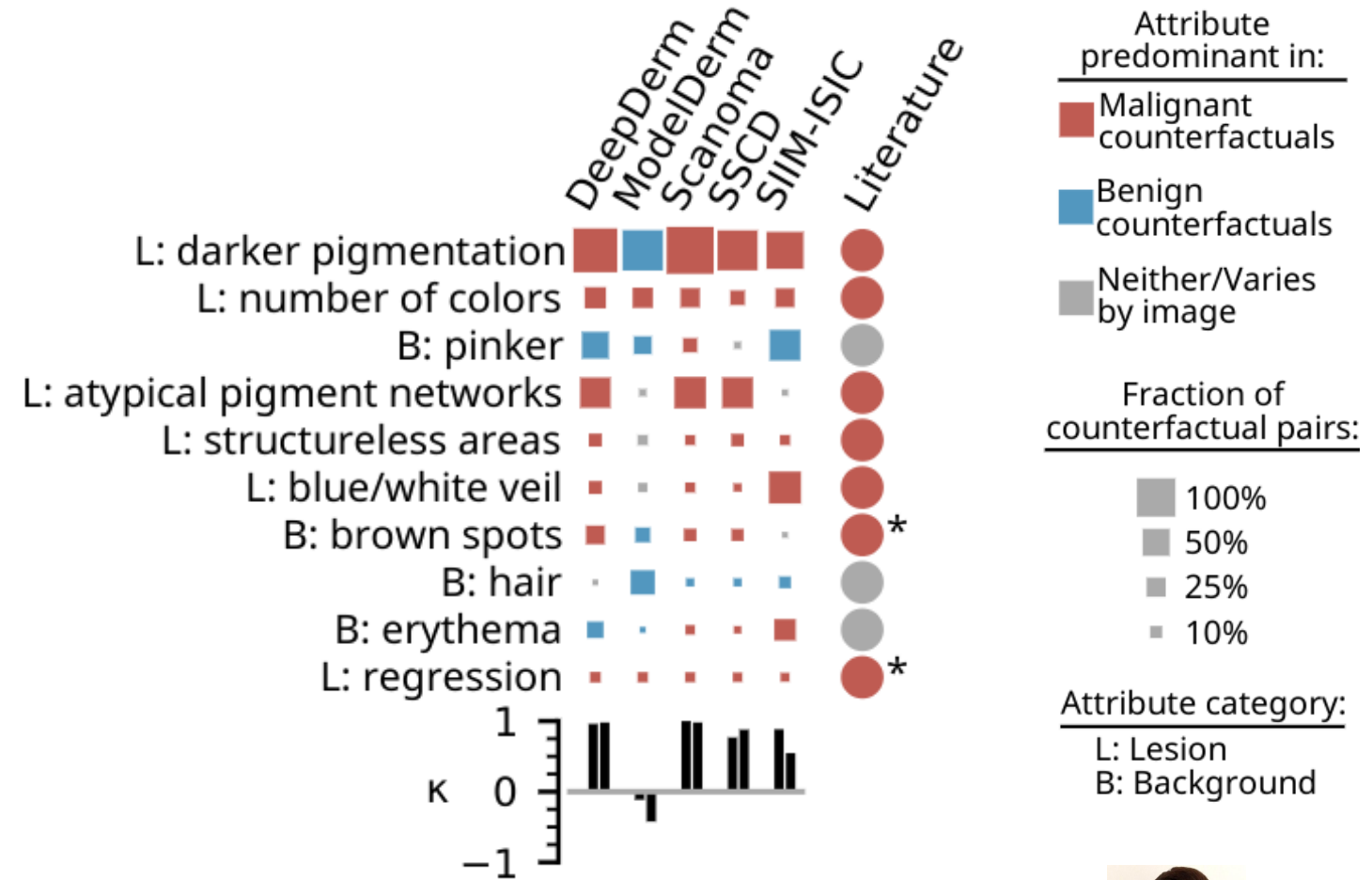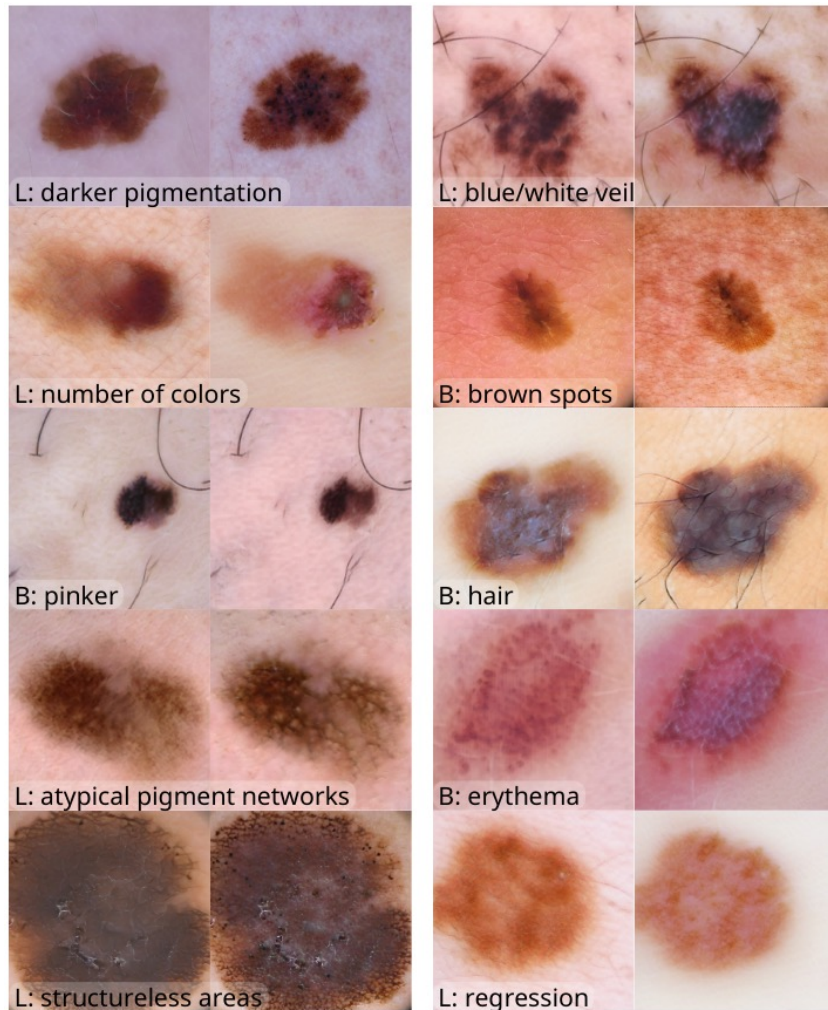
Clinicians frequently use conditional reasoning for treatment decisions by envisioning potential outcomes for patients. This is counterfactual thinking, exploring "what if" scenarios. Developments in generative artificial intelligence (AI) enable us to simulate this patient-level reasoning at the data level, opening new opportunities for science and health care. We term this approach counterfactual AI.

This approach is exemplified by use of counterfactual images in dermatology. Using AI, original skin images were modified to resemble melanoma guided by the decision-making process of a particular AI-based dermatological classifier. Dermatologists were then tasked with identifying clinically relevant features in the counterfactual images of melanoma and normal conditions. This process elucidated the reasoning processes of five AI-based dermatological classifiers. This data-centric counterfactual AI aligns the reasoning processes of AI classifiers with human clinicians' intuition, establishing a new approach to auditing clinical AI classifiers. Model auditing provides insights into the performance of deployed clinical AI classifiers for patients, clinicians, regulators, and data scientists.

Such uses of counterfactual AI prompt a crucial question: how might patient data change under specific conditions such as genetic mutations, treatments, time, or ageing? This exploration leads to intriguing scenarios, including forecasting the progression of clinical images or other data types over time for a particular treatment, potentially providing prognostic insights, or simulating the impact of genetic mutations to enhance our comprehension of disease mechanisms and treatment outcomes. This could present a frontier for future research. For example, personalised T-cell receptor sequence design for immunotherapy offers possibilities for new treatment strategies. Moreover, counterfactual AI has the potential to fill data gaps for rare diseases or under-represented groups, aiding the development of more inclusive and comprehensive health-care solutions. Furthermore, counterfactual AI could spur innovation in scientific hypothesis generation for drug discovery and development, potentially leading to breakthroughs in urgent areas such as Alzheimer's disease. Research suggests it could generate data on specific pathological conditions and conduct in-silico synthetic lethality testing for novel combination therapies.

An unexpected synergy is emerging as data-centric counterfactual AI contributes to the interpretation and auditing of clinical AI models. There are challenges in understanding the decision-making processes of many AI models. Saliency maps or, more broadly, feature attribution methods, are commonly used for model interpretation, indicating the areas of an image (or other data types) that the AI model focuses on (figure). Yet they provide only a partial view of the inner workings of complex AI models, impeding efforts to identify flaws in clinical AI reasoning processes. Counterfactual AI expands the scope of explainable AI by providing counterfactual images that elicit specific outcome predictions from complex AI classifiers (figure), enabling humans to grasp more comprehensive insights into the reasoning processes of these classifiers. Collaborating with clinicians, counterfactual AI could unearth previously unnoticed image attributes. Research indicates that by partnering with AI methods capable of automatically annotating images with an array of semantically meaningful concepts, counterfactual AI can systematically probe AI classifiers about how these concepts affect their decision-making processes.

Counterfactual AI in medicine faces ethical concerns and challenges related to fairness, data quality, and generalisability. Obtaining high-quality, diverse datasets is difficult. Generalising to new data is also problematic, particularly across diverse patient populations and health-care settings. Moreover, ethical and regulatory issues, including patient privacy concerns about the use of training data, must be addressed to ensure responsible AI deployment in health care.

What should we do to fully leverage the potential of counterfactual AI to advance scientific and therapeutic discovery? Generative AI operates through complex models that necessitate explanation. The reciprocal relation between generative AI and explainable AI is essential: generative AI informs the development of explainable AI; explainable AI aids in understanding generative AI models. By focusing on these principles, we can ensure that "what if" AI models are transparent and interpretable, facilitating their effective use in biomedical endeavours.

*Su-In Lee, Eric J Topol
Paul G Allen School of Computer Science & Engineering, University of Washington, Seattle, WA 98195, USA (S-IL); Scripps Research Translational Institute, La Jolla, CA, USA (EJT)
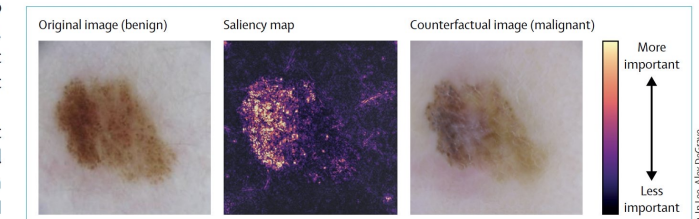suinlee@cs.washington.edu

**Further reading**

DeGrave AJ, Cai ZR, Janizek JD, Daneshjou R, Lee SI. Auditing the inference processes of medical-image classifiers by leveraging generative AI and the expertise of physicians. *Nat Biomed Eng* 2023; published online Dec 28. https://doi.org/10.1038/s41551-023-01160-9

DeGrave AJ, Janizek JD, Lee SI. AI for radiographic COVID-19 detection selects shortcuts over signal. *Nat Mach Intell* 2021; **3:** 610–19

Kim C, Gadgil SU, DeGrave AJ, et al. Transparent medical image AI via an image-text foundational model grounded in medical literature. *Nat Med* 2024 (in press)

Original image (benign) — Saliency map — Counterfactual image (malignant)

More important / Less important

**Figure:** Auditing dermatology AI model with counterfactual AI
A saliency map indicates little about an AI system for detecting melanoma, whereas counterfactual AI reveals that the system relied on the colour and pattern of pigmentation to determine that this lesion is benign.

# Fostering transparent AI via an *image-text foundation model* grounded in medical literature

- Finetune the CLIP (contrastive language-image pretraining) model

**Images**



PubMed articles & textbooks

This lesion is pigmented ... and asymmetric

**Captions**

**MONET Image Encoder**

**Contrastive Learning**

**MONET Text Encoder**

- Automatic concept annotation:
  - For each image,

**MONET (Medical concept retriever)**

Co-embedding space of image & text (concept)

Pigmented

Asymmetric

Ulcer

Erythema

Distance

● : Image
★ : Concept

Annotated concepts

- Asymmetric ↑
- Ulcer ↓
- Pigmented ↑
- Erythema ↓
⋮

Chanwoo, CSE PhD

Kim et al. *Nature Medicine*, 2024

# Beyond interpreting models…



— Cancer therapy design for precision oncology
[*Nature BME*'23]

— AI auditing [*Nature MI*'21, *Nature BME*'23, *Nature Medicine*'24]

radiology, dermatology

— Cost-aware clinical AI [*Nature BME*'22]

emergency medicine, critical medicine

# Explainable AI enables "cost-aware" AI (CoAI)

*One year ago …*



Gabe, MSTP/CSE PhD'21
**(now Harvard for residency in EM)**

Erion et al. ***Nature Biomedical Engineering***, 2022 - Featured in *Nature Comp. Science*, 2022

# Explainable AI enables "cost-aware" AI (CoAI)

- Gathering features is often costly. (e.g., time, money, etc)
  - Acute traumatic coagulopathy (ATC), a dangerous bleeding disorder in trauma patients (failure to clot)
  - ATC is time sensitive – often requires massive transfusion and earlier transfusion leads to better outcomes

- In collaboration with Nathan White, we used our trauma registry dataset
  - 14,000 emergency room visits and 46 features from the trauma registry of Harborview Medical Center, an urban level-I trauma centre

- CoAI combines XAI-based feature importance with feature cost (time)
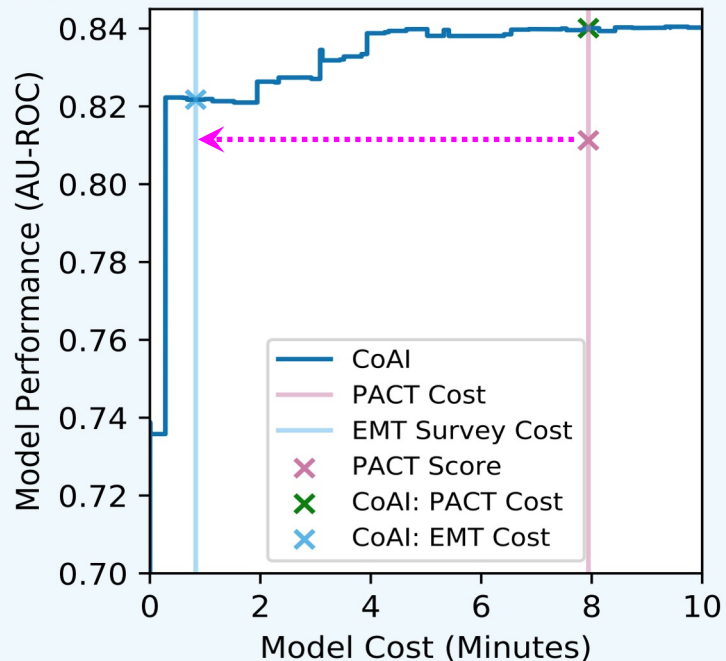  - Time cost survey from clinicians, medical directors, EMTs, etc

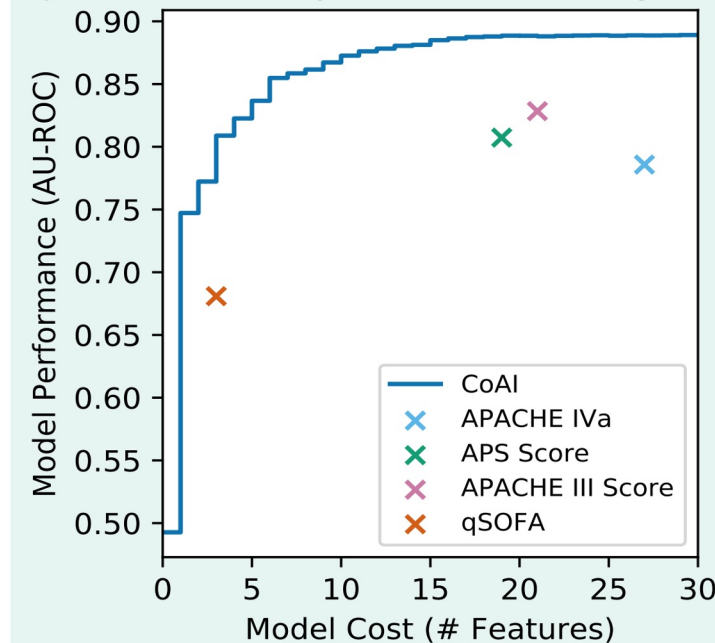Gabe, MSTP/CSE PhD'21 **(now Harvard for residency in EM)**

Erion et al. ***Nature Biomedical Engineering***, 2022 - Featured in *Nature Comp. Science*, 2022

# Explainable AI enables "cost-aware" AI (CoAI)

- CoAI improves *both* cost & accuracy
  - As accurate as the existing PACT score with <1 mins (vs. 8 mins) of feature gathering time



a) Trauma: CoAI outperforms PACT score

b) ICU: CoAI outperforms mortality scores

- CoAI is a general framework
  - Improves many existing clinical risk scores when applied to ICU mortality prediction
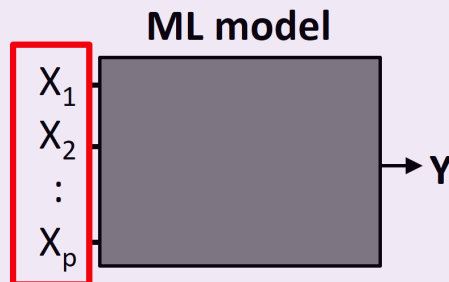
# Explainable AI for biomedical sciences & beyond



**Medicine & healthcare**
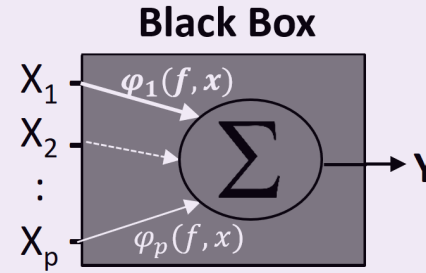anesthesia care, emergency medicine, critical care, nephrology, dermatology & biological age

**Cancer biology & precision medicine**

**Alzheimer's disease therapeutic target discovery**
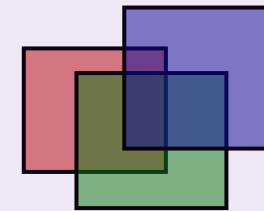
**Developing explainable AI principles techniques**

ML model

$X_1$
$X_2$
⋮
$X_p$

→ Y

**Learn interpretable features**

Black Box

$X_1$ → $\varphi_1(f,x)$
$X_2$ →
⋮
$X_p$ → $\varphi_p(f,x)$

Σ → Y

**Make interpretable predictions**

**Explanation priors**

**Learn explainable models**

**Clinical medicine**                    **Basic biology**

nature machine intelligence
Tree explainer

*ICLR'24*; *NeurIPS*'23; *NeurIPS*'23; *Nature MI'23*; *ICLR*'23; *ICLR*'23; *ICML*'23; *AISTATS*, 2022; *ICLR*, 2022; *Nature MI*, 2021; *JMLR*, 2021; *Nature Comm*,. 2022; *JMLR*, 2021; *NeurIPS*, 2020; *Nature MI* (cover), 2020; *NeurIPS*, 2020; *AISTATS*, 2020; *NeurIPS* (oral), Dec 2017

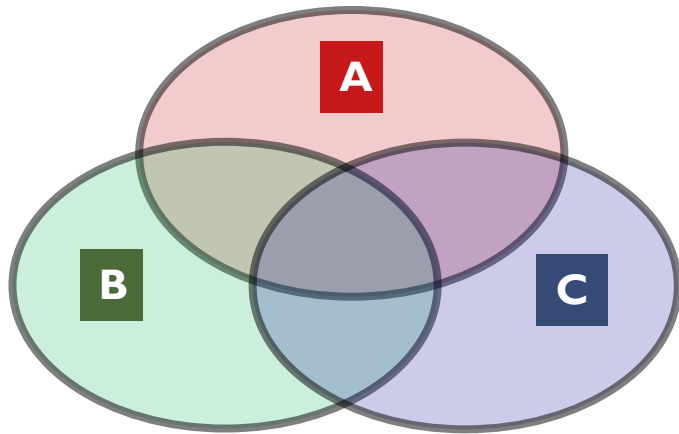*Nature Medicine*, 2024; *Lancet*, 2024; *Nature Methods*, 2023; *Genome Biology*, 2023; *Nature BME 2023*; *Lancet Healthy Longevity*, 2023 (cover); *Nature BME* 2023; *Nature Comm. Medicine*, 2022; *Nature BME*, 2022; *Nature Comm.*, 2021; *Nature MI*, 2021; *Nature Comm*, 2018; *Nature BME* (cover), 2018

nature biomedical engineering
Explainable AI predicts blood-oxygen levels during anaesthesia

# More about our research can be found at: https://aims.cs.washington.edu/publications

- A tip for navigating our publication site

**Which field does the paper aim to advance?**



**AIMS Lab**    Home · Su-In Lee ∨ · **Publications** · Research · People ∨ · Education

**A** AI/ML
**B** AI in Biology
**C** AI in Clinical Medicine

**J** Journal publications
**C** Conference publications

\* indicates equal contribution

## Under Review

**C J Dissection of medical AI reasoning processes via physician and generative-AI collaboration**
Alex J. DeGrave, Zhuo Ran Cai, Joseph D. Janizek, Roxana Daneshjou\*, and Su-In Lee\*
In Press, *Nature Biomedical Engineering*
medRxiv

**C J Fostering transparent medical image AI via an image-text foundation model grounded in medical literature**
Chanwoo Kim, Soham U. Gadgil, Alex J. DeGrave, Zhuo Ran Cai, Roxana Daneshjou\*, and Su-In Lee\*
In Revision, *Nature Medicine*
medRxiv

**A C Estimating Conditional Mutual Information for Dynamic Feature Selection**
Soham U. Gadgil\*, Ian Covert\*, Su-In Lee
Under Review, ICLR'24
arXiv

# AI for bioMedical Sciences (AIMS) Lab

UW MSTP

Nicasia Beebe-Wang (CSE PhD)

Ian Covert (CSE PhD)

**A** AI/ML
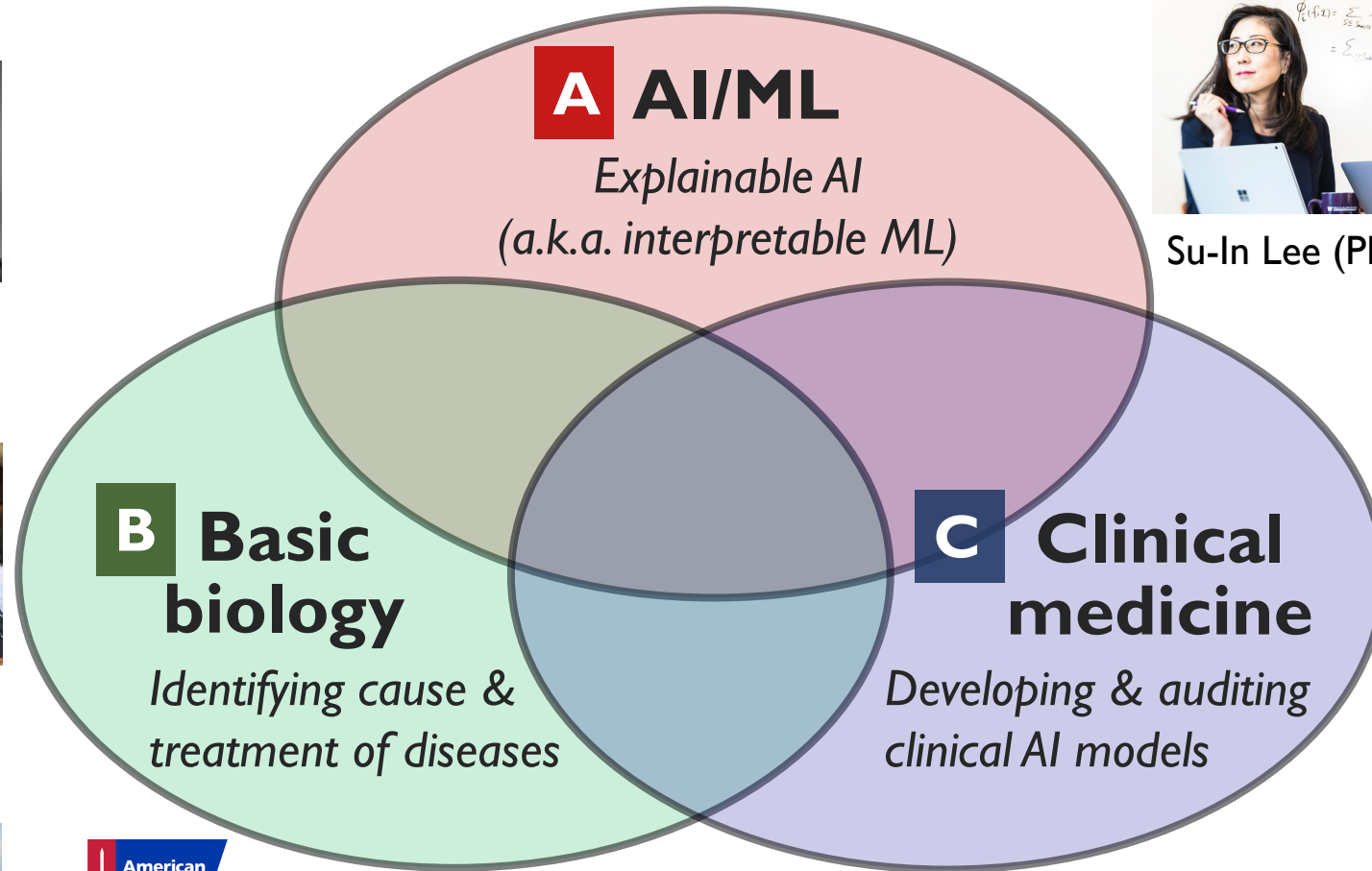*Explainable AI (a.k.a. interpretable ML)*

Su-In Lee (PI)

Hugh Chen (CSE PhD)

Joe Janizek (MSTP, CSE PhD; matched to Stanford)

Wei Qiu (CSE PhD)

Chris Lin (CSE PhD)

**B** Basic biology
*Identifying cause & treatment of diseases*

**C** Clinical medicine
*Developing & auditing clinical AI models*

Ethan Weinberger (CSE PhD)

Alex DeGrave (MSTP, CSE PhD)

Mingyu Lu, MD (CSE PhD)

Patrick Yu (CSE PhD)

American Cancer Society®

NIH

NSF

CZ

**Previous members:** Ben Logsdon (postdoc), Safiye Celik (CSE PhD'18), Scott Lundberg (CSE PhD'19), Parmita Mehta (CSE PhD'20), Gabe Erion (MSTP, CSE PhD'21; now Harvard Medical School for residency),

Chanwoo Kim (CSE PhD)

Soham Gadgil (CSE PhD)