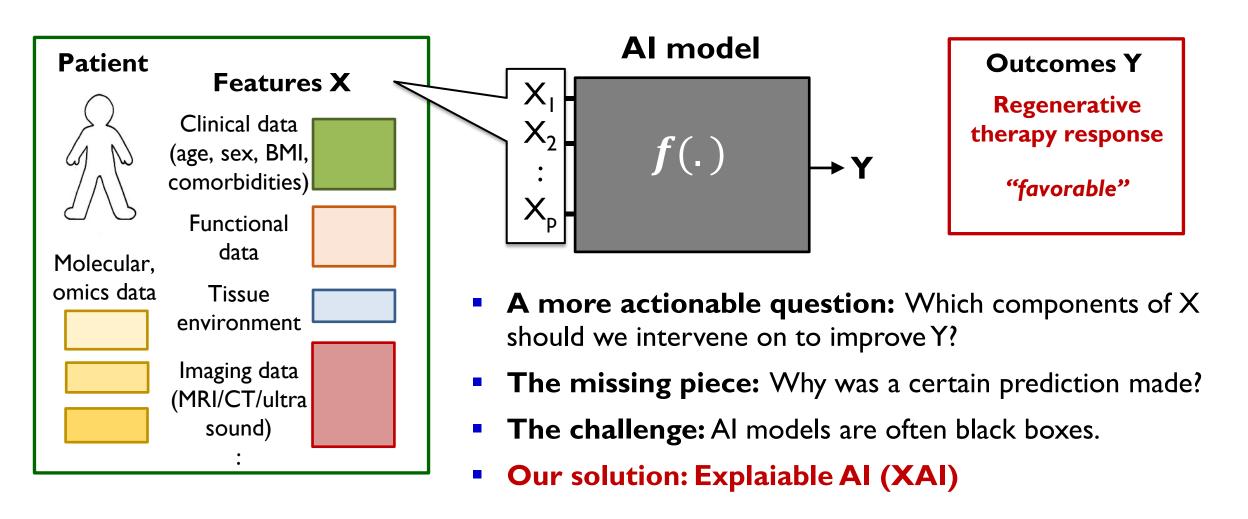
Explainable, Transparent, and Trustworthy AI for Biomedicine

Su-In Lee

Boeing Endowed Professor
Paul G.Allen School of Computer Science & Engineering
University of Washington, Seattle

A common Al problem in biomedicine: Given features X, predict Y.

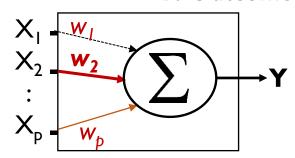


Our solution was to fundamentally advance Al research to make a prediction with explanations

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

Linear model

X: Features Y: Outcome



Complex model f (.)

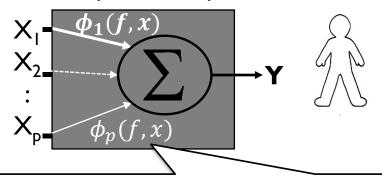
Black Box



Our approach, SHAP

(SHapley Additive exPlanations)

For a particular prediction

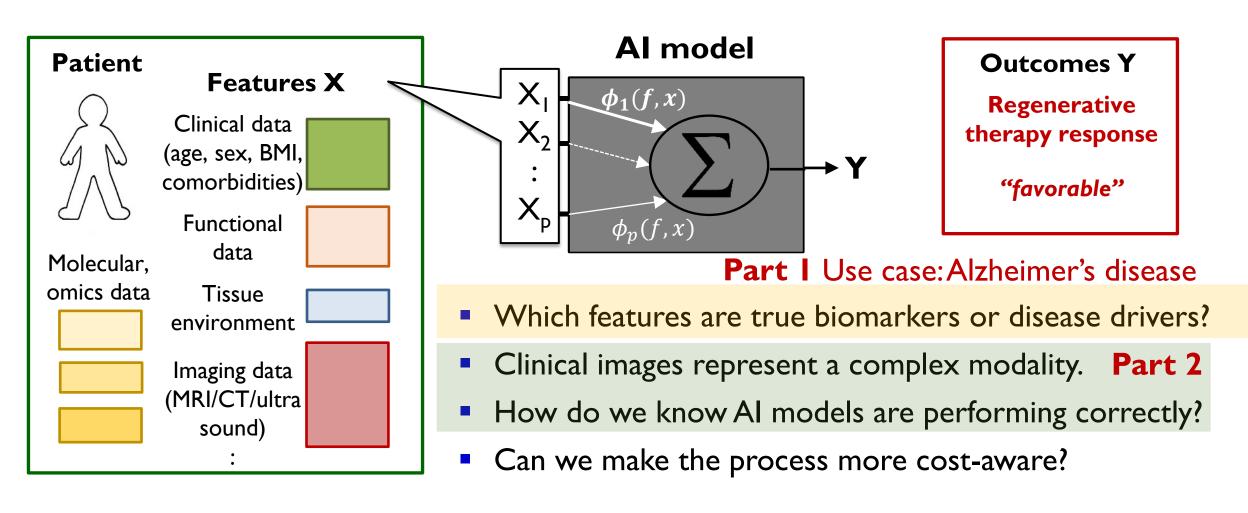


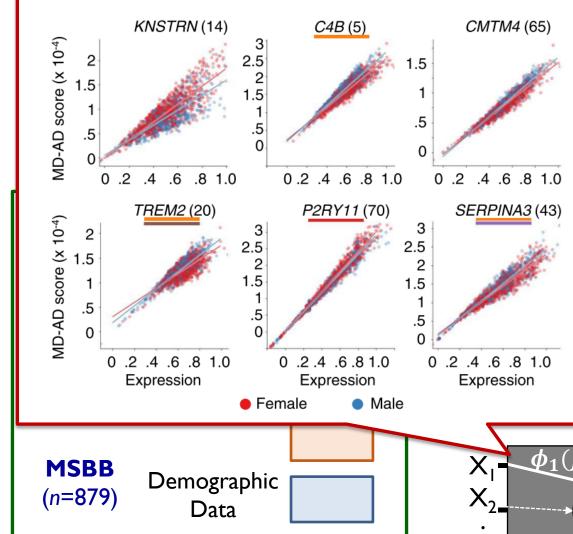


Scott, CSE PhD'19

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

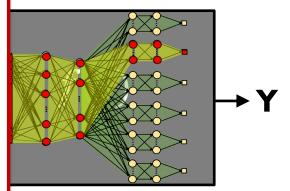
Explainable AI (XAI): Accurately predicting an outcome is vital, but the critical question revolves around *why*.





mechanistic explanation AD) phenotypes

Al model



plaque

counts

Outcomes Y

Neuropathological phenotypes:

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

- Estimate each gene's contribution to AD neuropathologies
 - Previously unknown sexdifferential immune response (microglia activity) in AD

ROSMAP (*n*=542)

Clinical data

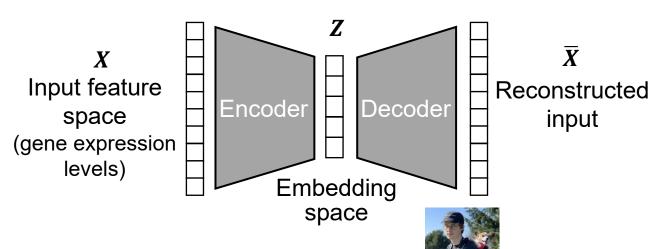


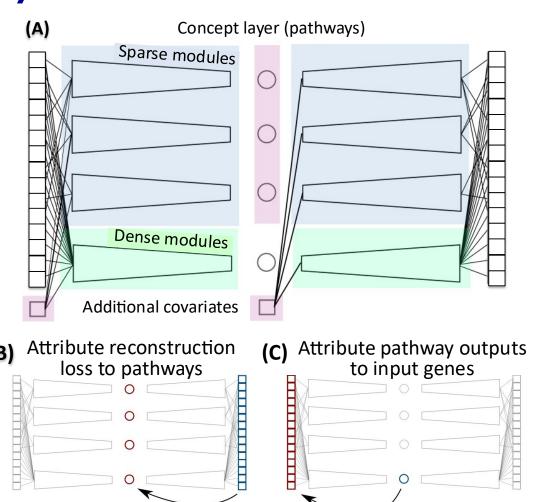
Beebe-Wang et al. Nature Communications, 2021

CSE PhD'22

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 important genes that explain the expression variation within the dataset





Janizek et al. Genome Biology, 2023

Joe, UW MSTP/CSE PhD'22, Residency at Stanford Radiology

Biologically interpretable AI modeling identifies experimentally validated AD therapeutic targets

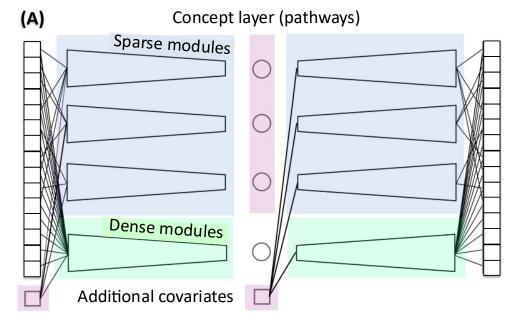
- We applied our approach to extended bulk RNAseq datasets from AD study cohorts
- We identified mitochondrial complex I as a potential mediator for tolerance to Aβ toxicity
 - In vivo validation in a transgenic C. elegans model expressing Aβ done by Matt Kaeberlein's lab

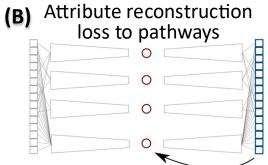
A promising pharmacological avenue!







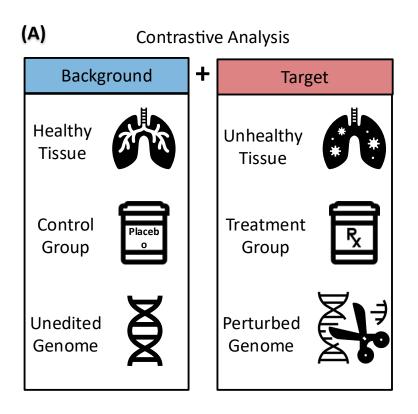


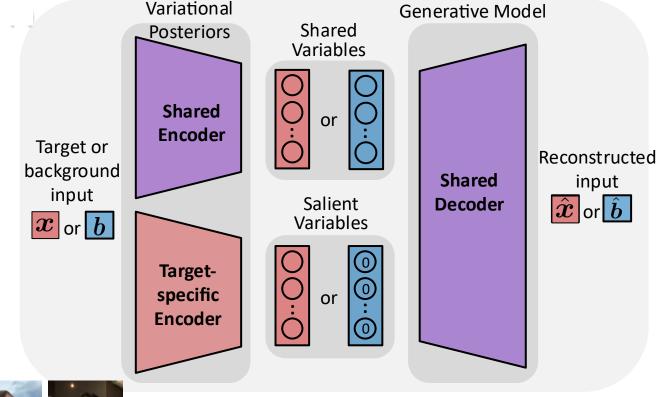


(C) Attribute pathway outputs to input genes

Contrastive modeling enhances interpretability

 Single-cell datasets are often collected to investigate differences in cellular state between background cells and those under specific treatments





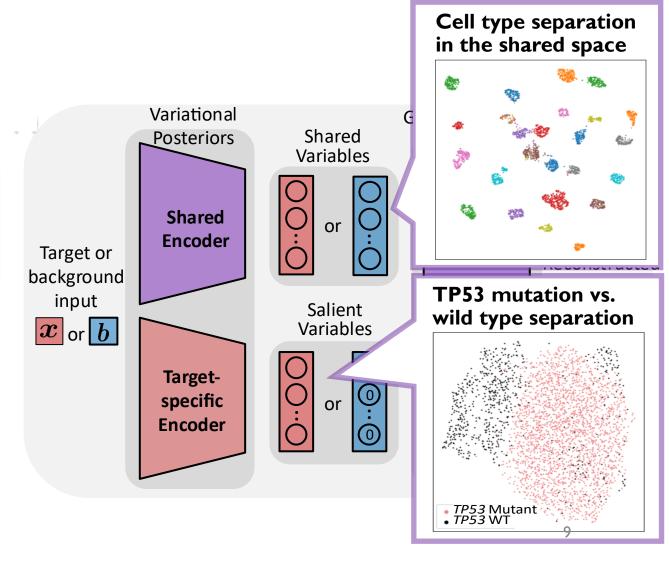


Contrastive modeling enhances interpretability

- 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

Outline – Two parts

- X_1 $\phi_1(f,x)$ X_2 \vdots $\phi_p(f,x)$ $\phi_p(f,x)$
- Part I Identifying disease-driving genes and processes
 - Unveiling neurodegenerative disease insights with explainable Al [Nature Comm'21; Genome Biology'23; Nature Methods'23]
- Part 2 Auditing AI models

- Feature attributions [Nature MI'21 featured in Nature'22]
- Counterfactual explanations (via generative AI) [Nature BME'25 cover; Lancet'24]
- Concept-based explanations (via foundation models) [Nature Medicine'24]

Our explainable AI techniques are generalizable to a wide range of biomedical problems, including regenerative medicine.

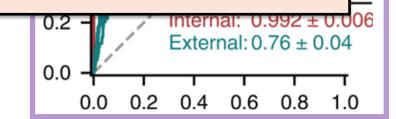
Auditing AI for COVID-19 detection using XAI

Many published Al models that detect COVID-19

1.0

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!



- √ Clear lung bases predict negative COVID-19 status
- X laterality markers should not predict negative status
- X medical devices should not predict negative status

 MSTP / CSE PhD

oth

99th

Our Al auditing work featured in Nature

• "Breaking into the black box of artificial intelligence" Nature Outlook

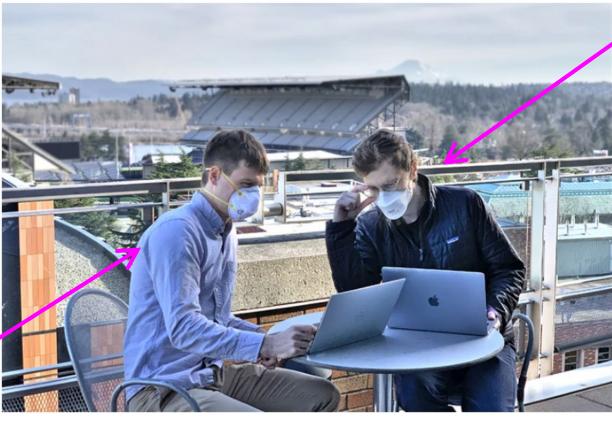
Breaking into the black box of artificial intelligence

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

By Neil Savage



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



UW MSTP / CSE PhD

matched to Stanford)

Joe Janizek (Just

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

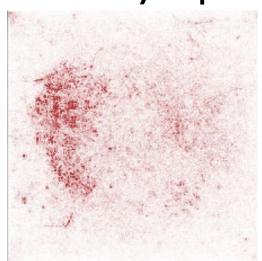
- Auditing Al models to predict skin cancer
 - Five models 2 academic models, 2 commercial devices, and I competition winner
- Technical challenges saliency maps often do not work

Original image



Predicted: benign

Saliency map



Modified image



Predicted: malignant

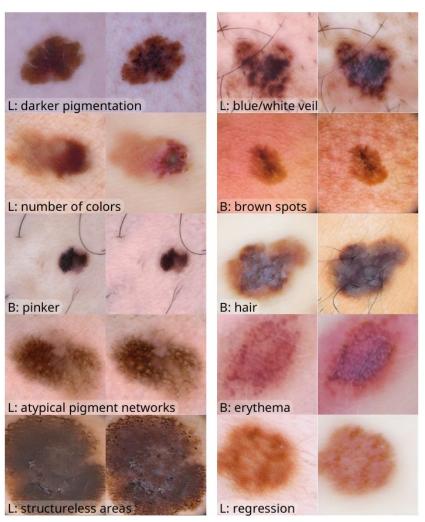
Our solution:

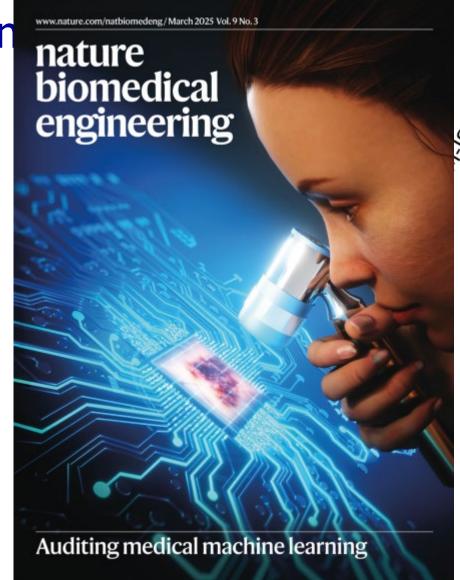
- Generate counterfactual images from the Al model
- Systematic characterization by experts: Drs. Roxana Daneshjou, and Zhuo Ran Cai (Stanford)

DeGrave et al. (Nature Biomedical Engineering)
Kim et al. (Nature Medicine, 2024)

How do dermatology Al systems make decisions

on dermoscopic in







Malignant counterfactuals

Benign counterfactuals

Neither/Varies by image

Fraction of counterfactual pairs:

100%

50%

■ 25%

= 10%

Attribute category:

L: Lesion

B: Background



Alex, MSTP / CSE PhD

Degrave, Ran Cai, Janizek, Daneshjou,* and Lee* Nature Biomedical Engineering, 2025

The Lancet perspective

"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.

Al model focuses on (figure). Yet they provide only a partial view of the inner workings of complex Al models, impeding efforts to identify flaws in clinical Al reasoning processes. Counterfactual Al expands the scope of explainable Al by providing counterfactual images that elicit specific outcome predictions from complex Al classifiers (figure), enabling humans to grasp more comprehensive insights into the reasoning processes of these classifiers. Collaborating with clinicians, counterfactual Al could unearth previously unnoticed image attributes. Research indicates that by partnering with Al methods capable of automatically annotating images with an array of semantically meaningful concepts, counterfactual Al can systematically probe Al 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

Nature Reviews Bioengineering, 2025

Medical AI transparency in the entire life cycle

nature reviews bioengineering

https://doi.org/10.1038/s44222-025-00363-w

Review article



Transparency of medical artificial intelligence systems

Chanwoo Kim 12, Soham U. Gadgil 12 & Su-In Lee

Abstract

DeGrave AJ, Cai ZR, Janizel Daneshjou R, Lee SI. Audli inference processes of me image classifiers by levera generative AI and the exp of physicians. Nat Biomed 2023; published online Dc https://doi.org/10.1038/ s41551-023-01160-9

Medical artificial intelligence (AI) systems hold promise for transforming healthcare by supporting clinical decision-making in diagnostics and treatment. The effective deployment of medical AI requires trust among key stakeholders – including patients, providers, developers and regulators – which can be built by ensuring transparency in medical AI, including in its design, operation and outcomes. However, many AI

Sections

Introduction

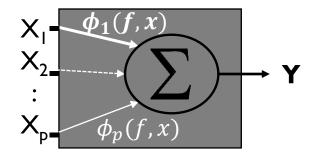
Data transparency

Model transparency

Deployment transparency

Outlook

Outline – Two parts



- Part I Identifying disease-driving genes and processes
 - Unveiling neurodegenerative disease insights with explainable Al [Nature Comm'2 I; Genome Biology'23; Nature Methods'23]
- Part 2 Auditing Al models
 - Feature attributions [Nature MI'21 featured in Nature'22]
 - Counterfactual explanations (via generative AI) [Nature BME'25 cover; Lancet'24]
 - Concept-based explanations (via foundation models) [Nature Medicine'24]

Our explainable AI techniques are generalizable to a wide range of biomedical problems, including regenerative medicine.

Al for bioMedical Sciences (AIMS) Lab

Nicasia Beebe-Wang (CSE PhD)



Ian Covert (former CSE PhD)



Wei Qiu (CSE PhD)



Chris Lin (CSE PhD)



Explainable Al (a.k.a. interpretable ML)



Su-In Lee (PI)



Hugh Chen (former CSE PhD)



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



Ethan Weinberger (CSE PhD)



Alex DeGrave (MSTP, CSE PhD)



Developing & auditing clinical AI models



Mingyu Lu, MD (CSE PhD)



Patrick Yu (CSE PhD)



Basic

biology

Identifying disease drivers

& advancing treatments

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),



(CSE PhD)



Chanwoo Kim Soham Gadgil (CSE PhD)