# Lessons from Deploying AI in Healthcare

James Zou

Stanford University

March 5, 2022
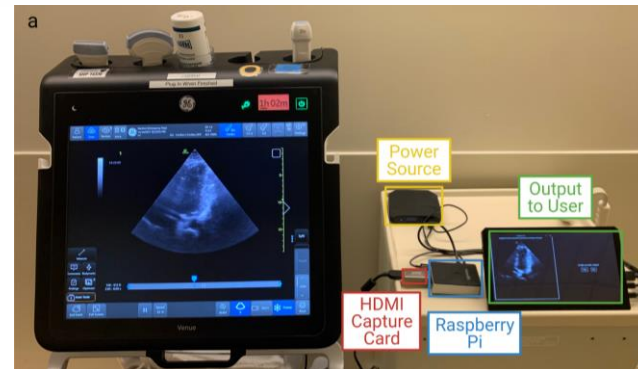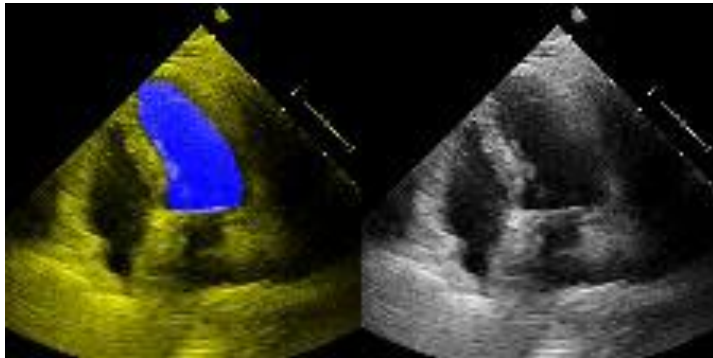
www.james-zou.com                                    jamesz@stanford.edu

# Example 1: deploying cardiology AI

David Ouyang        Bryan He

Ouyang et al. *Nature* 2020

# Computer vision assesses cardiac ultrasound

## Algorithm output

## EchoNet assessed chamber area

# Algorithm mimics clinical workflow

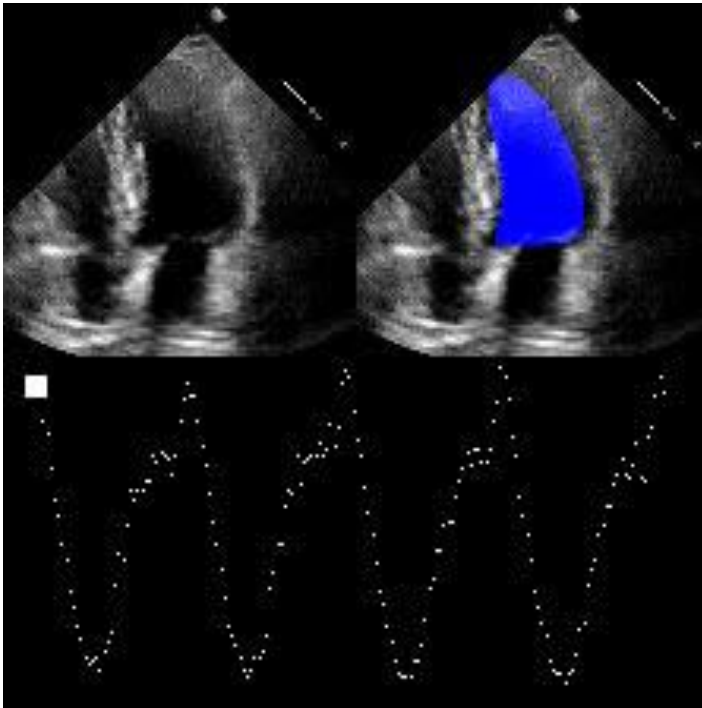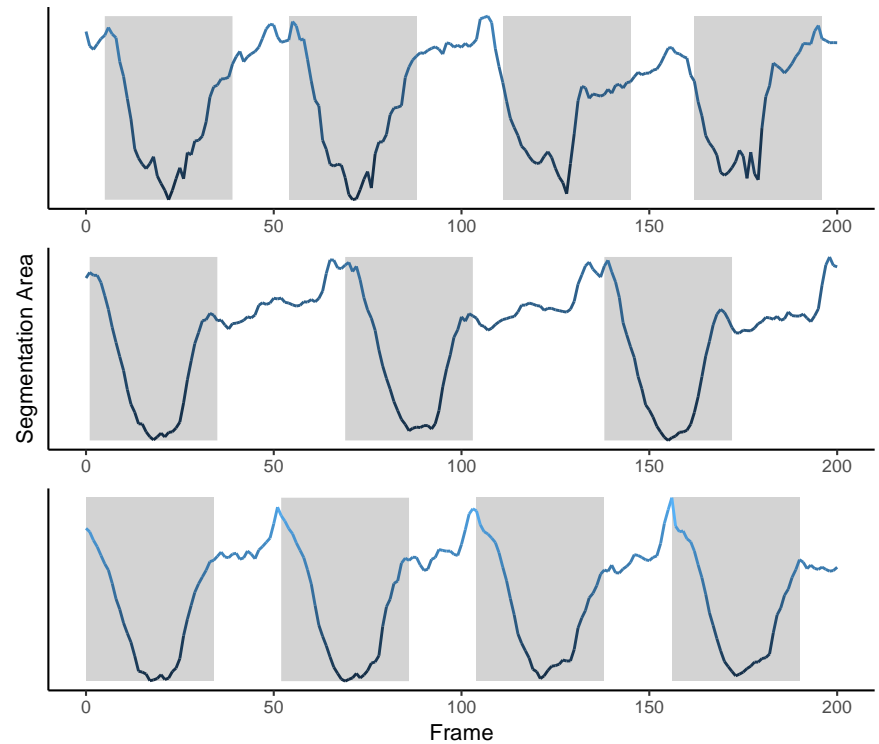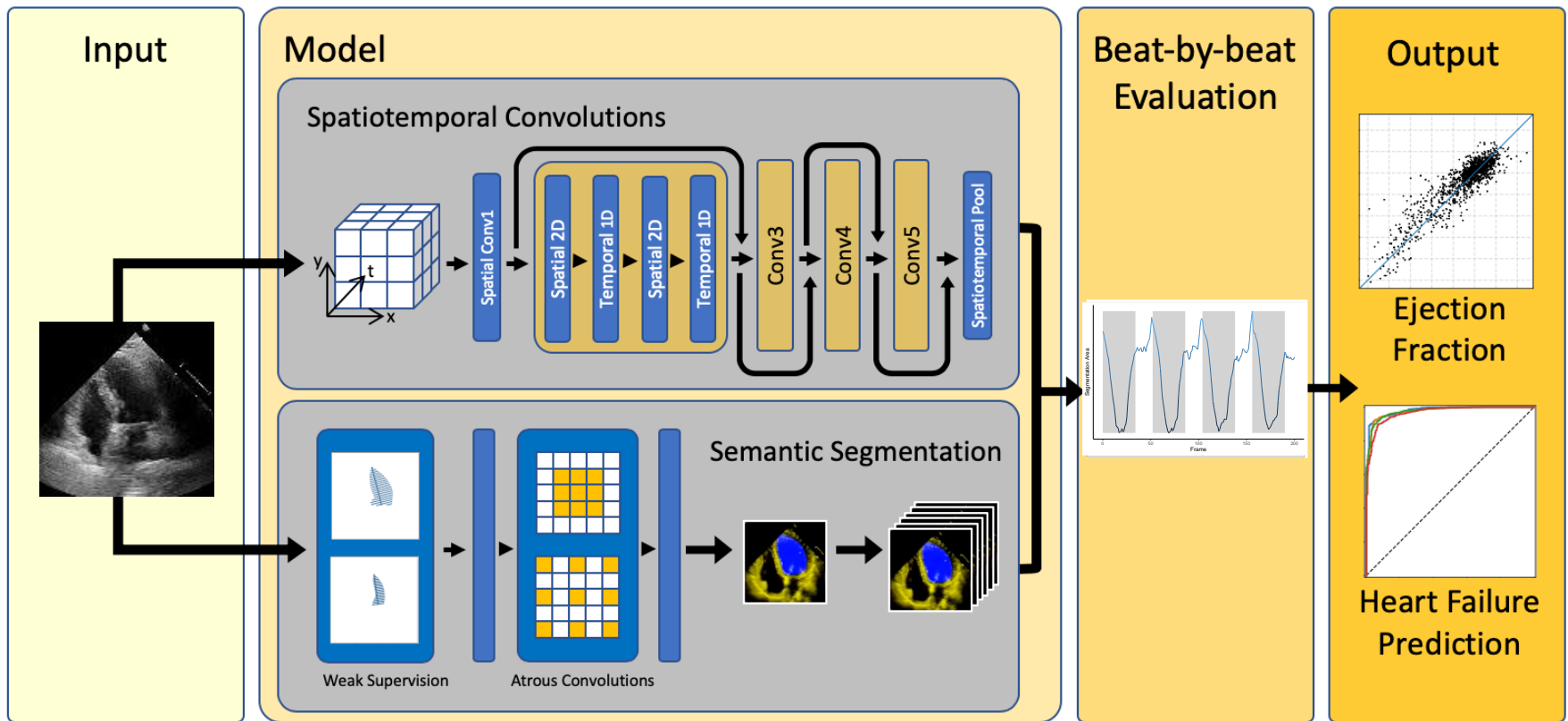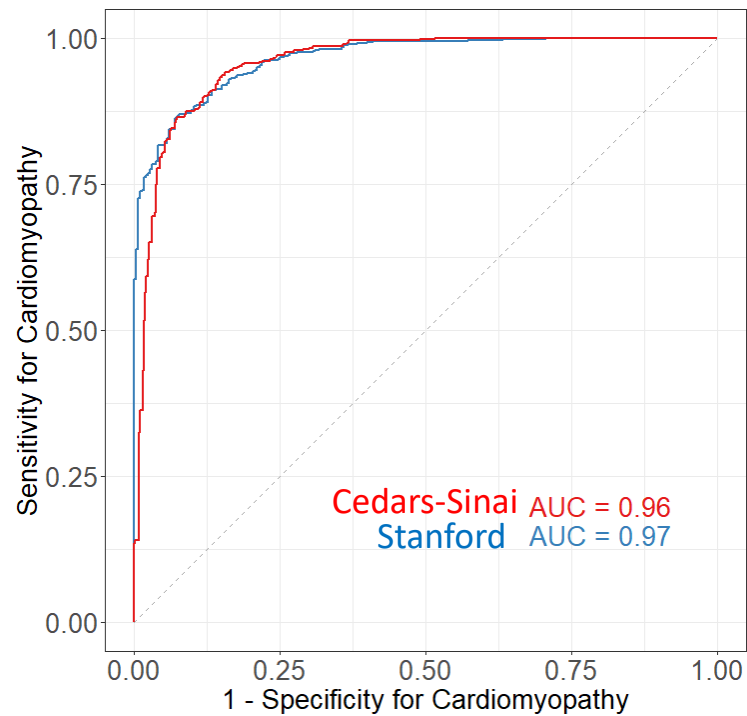## EchoNet-Dynamic



**Idea**: use temporal segmentation to focus attention of model.

Ouyang et al. *Nature* 2020

# Achieves expert performance in new hospital

## Predicting heart failure



Cedars-Sinai AUC = 0.96
Stanford AUC = 0.97

## Examples

# Example 2: AI to improve telemedicine

**COVID-19** $\longrightarrow$ **50x increase in digital visits**

# Many patient photos for telemedicine are poor quality

- Manual review of photos prior to the physician encounter consumed >2000 hours in 2021 at Stanford

- Poor quality photos disrupt clinic workflow

- Improving teledermatology = improving access to care

# TrueImage = Online check deposit for dermatology



Vodrahalli et al, PSB 2021

8

# TrueImage Workflow



TrueImage

Please move to brighter lighting

Move to brighter lighting

TrueImage

Photo submitted to your clinician

# TrueImage Algorithm



Deep learning + ensembling

Vodrahalli et al, PSB 2021
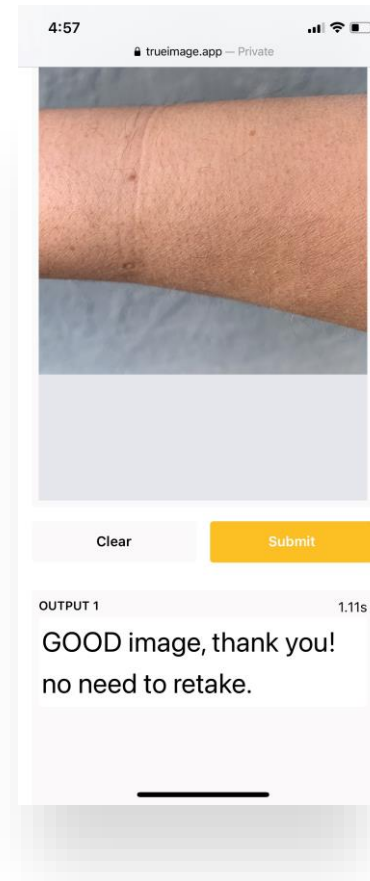
# Prospective study at Stanford

TrueImage filters 80% of poor quality photos; takes <1 minute per patient

# What does improvement look like?

# Deploying TrueImage at Stanford clinics

Launch Gradio interface on HIPAA compliant servers

Your private AWS/GCP machine creates tunnel and public link

**HIPAA COMPLIANT**

Your authorized users can now access the model

4:54
🔒 trueimage.app — Private

Clear    Submit

OUTPUT 1                    1.08s

The following issues were detected with your image:
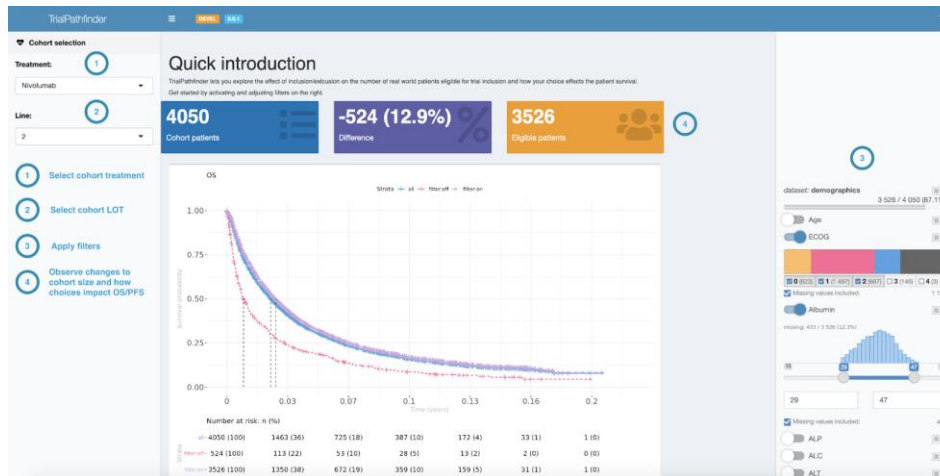blur, lighting

----

# Example 3: AI to design clinical trials



Article | Published: 07 April 2021

## Evaluating eligibility criteria of oncology trials using real-world data and AI

Ruishan Liu, Shemra Rizzo, Samuel Whipple, Navdeep Pal, Arturo Lopez Pineda, Michael Lu, Brandon Arnieri, Ying Lu, William Capra, Ryan Copping ✉ & James Zou ✉

*Nature* **592**, 629–633(2021) | Cite this article

Ruishan Liu

Liu et al. *Nature* 2021

# Google launches AI health tool for skin conditions  5/18/21

Breakthrough development will assist users in self-diagnosing issues ranging from acne to melanoma



A woman checks birthmarks on her back. Derm Assist will be free to all internet users, whether they are Google users or not © Albina Gavrilovic/Getty
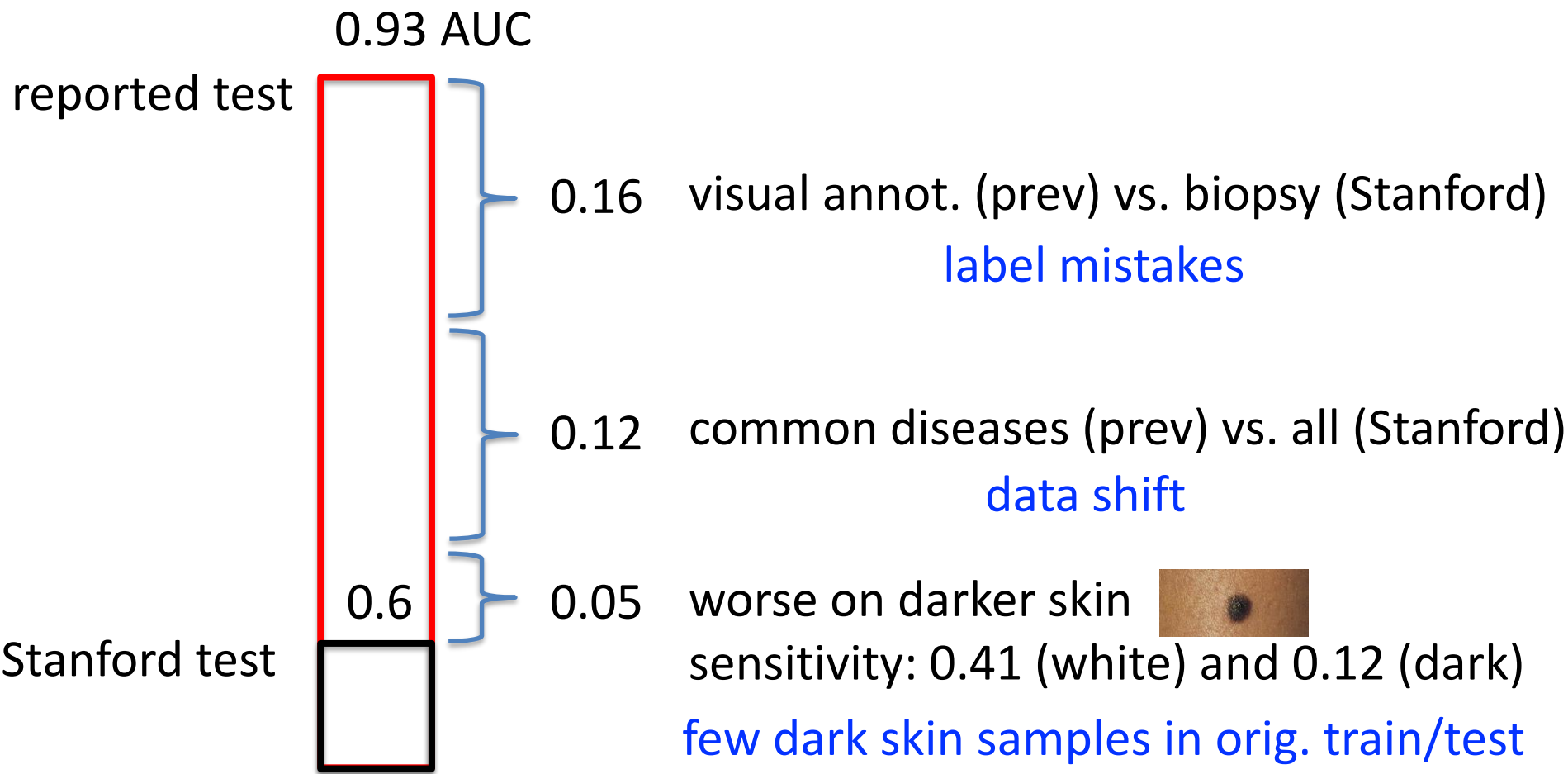
# AI dermatology apps
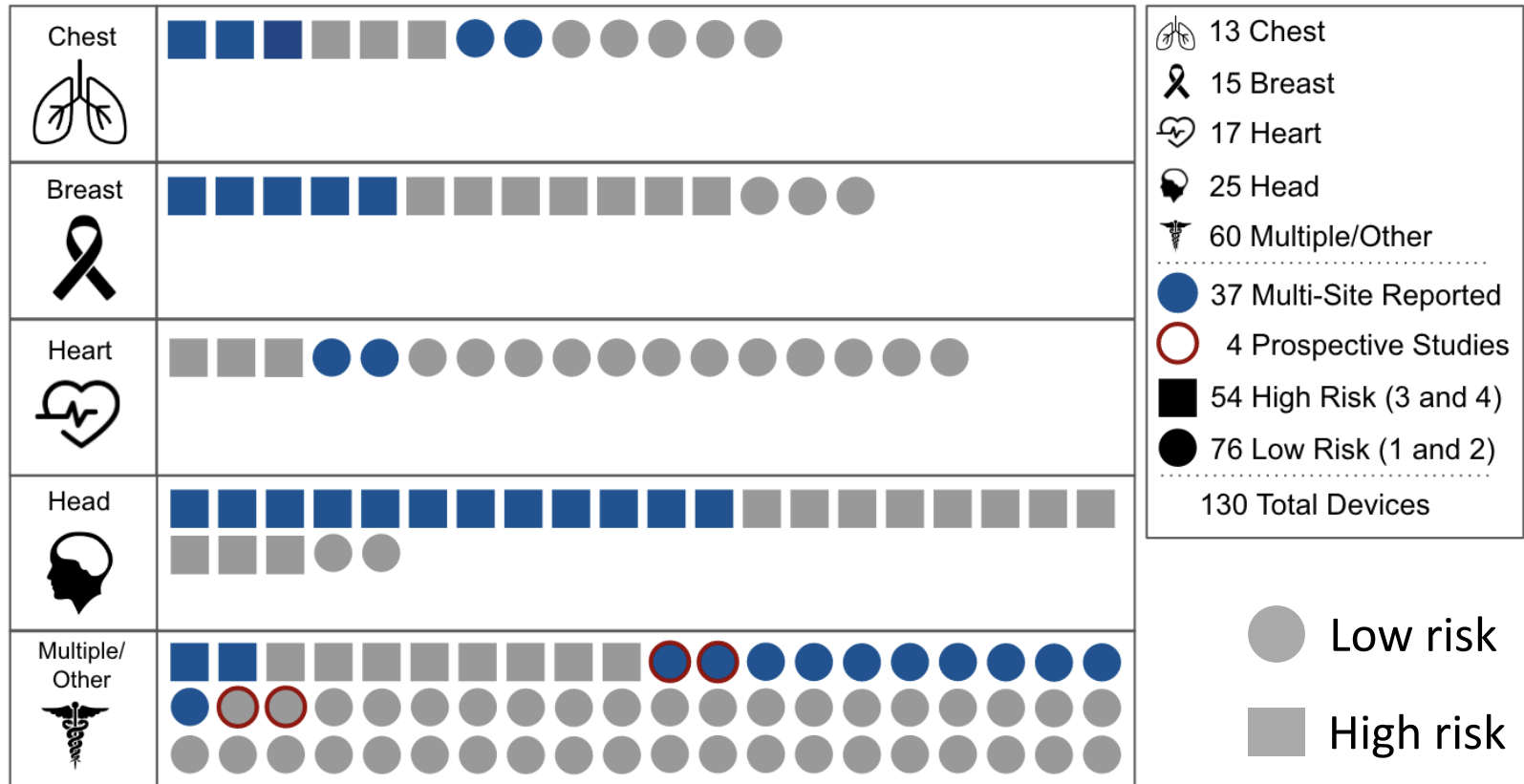
Original reported: 0.93 AUC
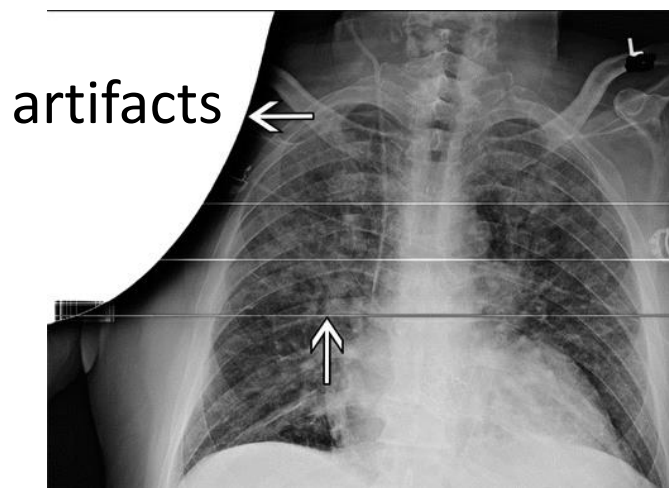
Stanford patients: 0.60 AUC

ModelDerm

Roxana Daneshjou

# Why did the Derm AI performance crater?

0.93 AUC

reported test

0.16    visual annot. (prev) vs. biopsy (Stanford)

label mistakes

0.12    common diseases (prev) vs. all (Stanford)

data shift

0.6

0.05    worse on darker skin

Stanford test

sensitivity: 0.41 (white) and 0.12 (dark)

few dark skin samples in orig. train/test

Roxana Daneshjou

# Data used to test 130 FDA-approved AI



Legend:
- 13 Chest
- 15 Breast
- 17 Heart
- 25 Head
- 60 Multiple/Other
- 37 Multi-Site Reported
- 4 Prospective Studies
- 54 High Risk (3 and 4)
- 76 Low Risk (1 and 2)
- 130 Total Devices

Low risk
High risk

93/130 did not report multi-site evaluation
Only 4 prospective studies

# Large variability in cross site performance

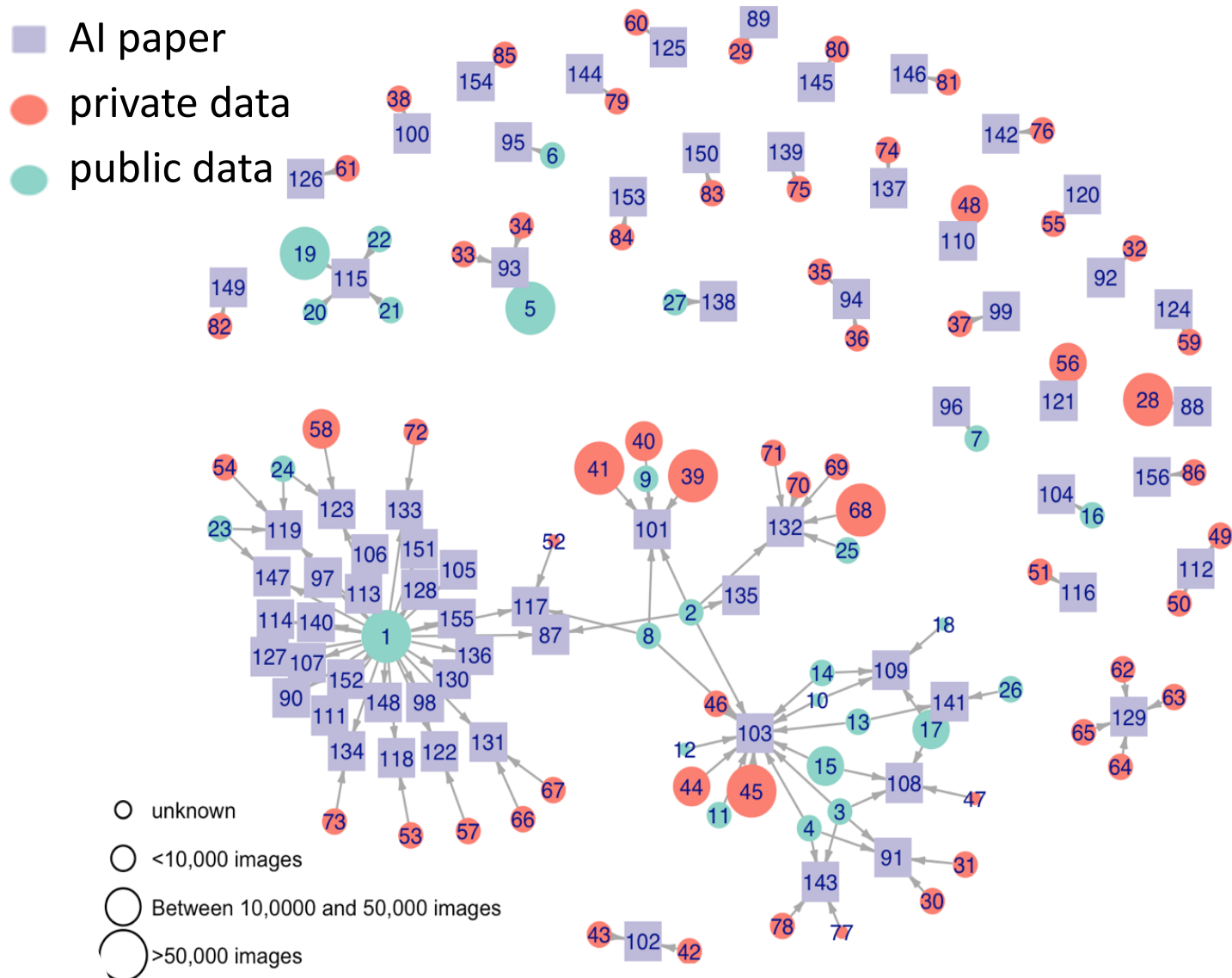| Site | Stanford (N=19K) | Boston (N=23K) | NIH (N=11K) |
|------|------------------|----------------|-------------|
| **Stanford** | **0.90 ± 0.01** | 0.87 ± 0.01 | 0.85 ± 0.02 |
| **Baylor** | 0.83 ± 0.01 | **0.89 ± 0.01** | 0.84 ± 0.02 |
| **NIH** | 0.78 ± 0.01 | 0.76 ± 0.02 | **0.88 ± 0.02** |

artifacts

Pneumothorax detection

# Lessons for deploying trustworthy medical AI

1. Understand what data is used to develop the AI.

2. Understand why AI makes systematic mistakes.

3. Use human-in-the-loop evaluation.

# Lessons for deploying trustworthy medical AI

1. Understand what data is used to develop the AI.

2. Understand why AI makes systematic mistakes.

3. Use human-in-the-loop evaluation.
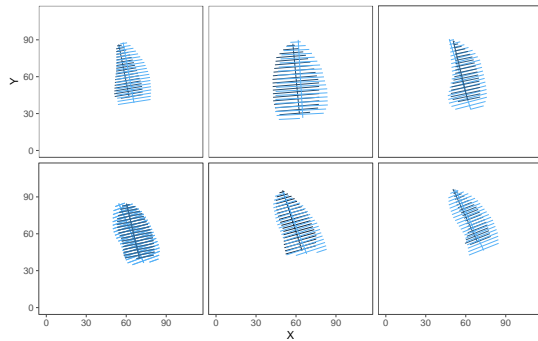
# 1. Data used to train dermatology AI



Daneshjou et al. *JAMA Derm* 2021

# 1. Transparent dataset and code



Largest public dataset of medical videos.

Ouyang et al. *Nature* 2020

# Lessons for deploying trustworthy medical AI

1. Understand what data is used to develop the AI.

2. Understand why AI makes systematic mistakes.

3. Use human-in-the-loop evaluation.

# 2. Why did the model make this mistake?



Conceptual explanation of mistakes

# Conceptual explanation of mistakes

## Mistakes made by the model



**Label:** Allergic Contact Dermatitis
**Pred:** Stasis Edema (19%)

- Blackness    -0.42
- **Dark Skin**    -0.67

**Label:** Fixed Eruptions
**Pred:** Erythema Nodosum(35%)

- Ashcan    -1.02
- **Defocus Blur**    **-1.20**

**Label:** Mucinosis
**Pred:** Aplasia Cutis (9%)

+ Zoom    **0.75**
+ Contrast    0.73

**Label:** Sarcoidosis
**Pred:** Nevus Sebaceous of Jadassohn (36%)

- Motion Blur    -0.57
- **Skin Hair**    **-0.87**

Output of our AI mistake explainer

Abubakar Abid, Mert Yuksekgonul

# Lessons for deploying trustworthy medical AI

1. Understand what data is used to develop the AI.

2. Understand why AI makes systematic mistakes.

3. Use human-in-the-loop evaluation.

# 3. AI often optimizes the wrong objective

optimization over a **fixed,
single-site** validation dataset



| Data Acquisition and Labeling | Split Data into Train/Validation | Model Building | Model Validation |

Abubakar Abid

# Optimize for **human usage** instead!

optimize over **real-world user data**

| Data Acquisition and Labeling | Split Data into Train/Validation | Model Building | Model Validation | Real-World Usage |

Abubakar Abid

# Human-in-the-loop evaluation of ML impact



Kailas Vodrahalli

# Worse AI can be better for humans

Human accuracy improvement
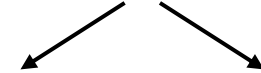
Human confidence in correct answer



Uncalibrated = overconfident models

Kailas Vodrahalli

# Using gradio

model to deploy

type of UI to create

```
import gradio
app = gradio.Interface(classify_skin_image, inputs="image", outputs="label")
app.launch(share=True)
```

url: www.gradio.app/test12543
can be shared



Abid et al *Nature Machine Intelligence* 2020
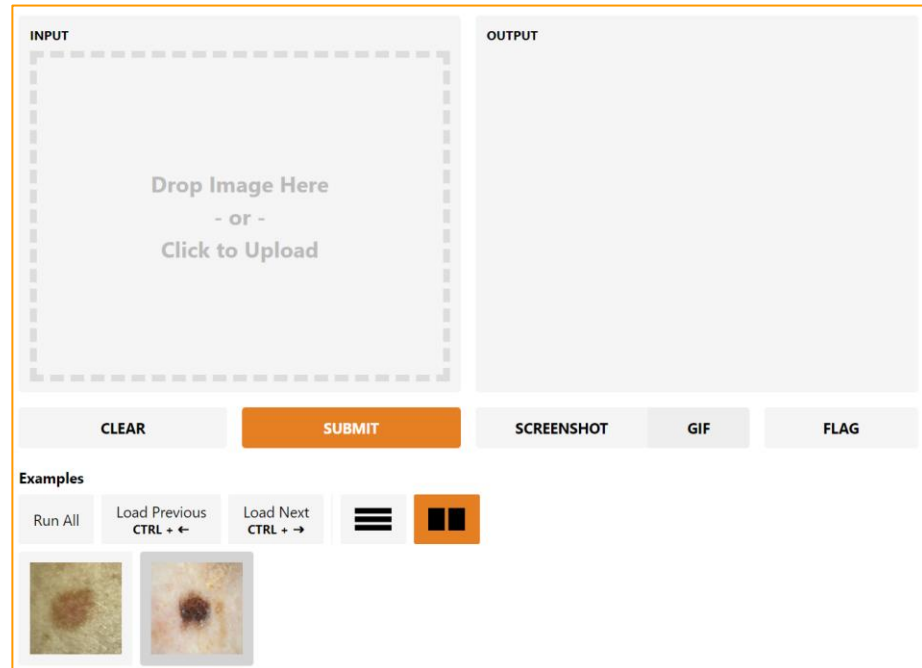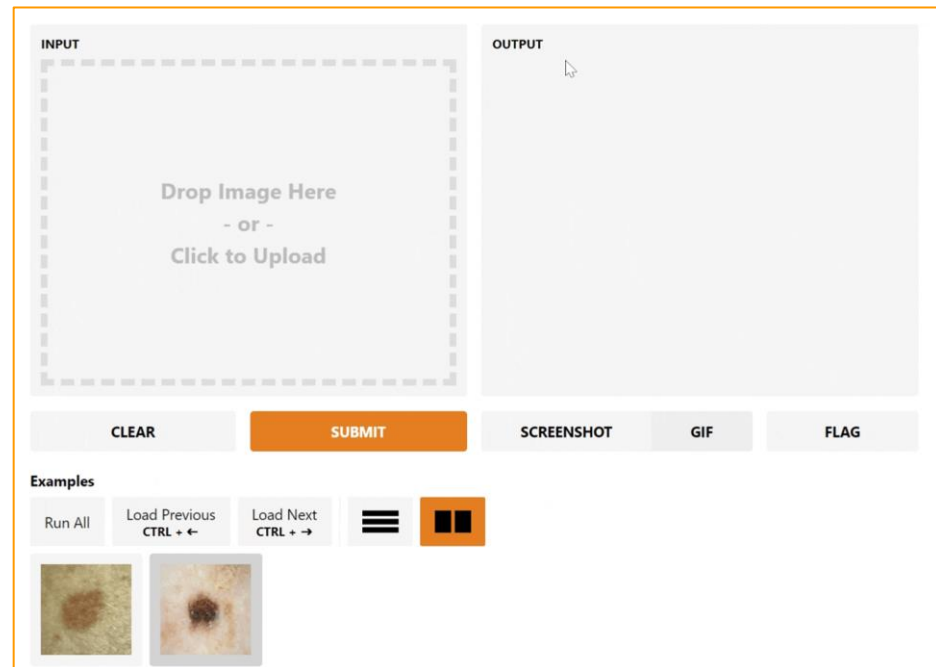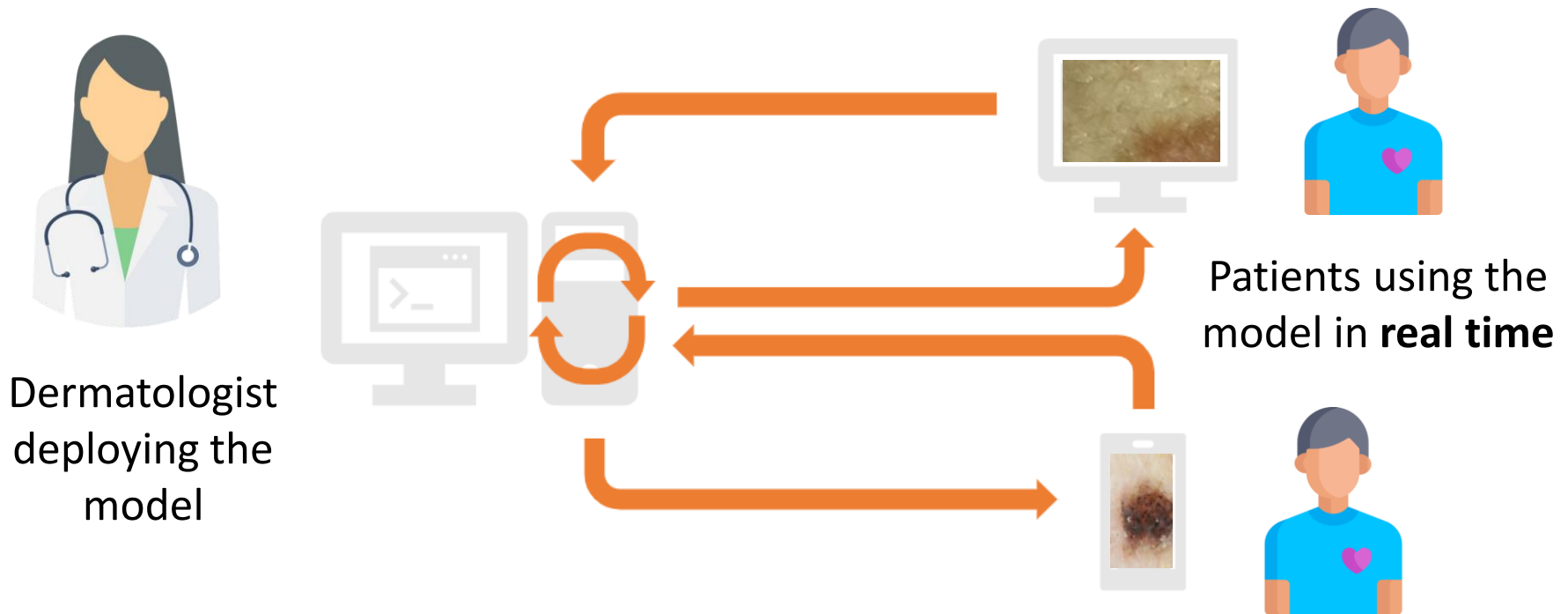
# Using gradio

model to deploy

type of UI to create

```
import gradio
app = gradio.Interface(classify_skin_image, inputs="image", outputs="label")
app.launch(share=True)
```

url: www.gradio.app/test12543
can be shared



Abid et al *Nature Machine Intelligence* 2020

# gradio used for Stanford's 1st real-time AI trial



Dermatologist deploying the model

Patients using the model in **real time**

# Lessons for deploying trustworthy medical AI

1. Understand what data is used to develop the AI.

2. Understand why AI makes systematic mistakes.

3. Use human-in-the-loop evaluation.

# Resources and thanks

Papers and codes available www.james-zou.com

Disparity in dermatology AI

Daneshjou et al. *JAMA Dermatology* 2021

Data transparency for biomedical AI

Wu et al *Nature Medicine* 2021
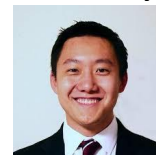
Video-based AI for cardiac assessment.

Ouyang et al. *Nature* 2020

Explaining model mistakes

Abid, Yuksekgonul, Zou. In review 2022

Gradio for human-in-the-loop AI

Abid et al. *Nature Machine Intelligence* 2020

Roxana Daneshjou

Eric Wu

David Ouyang

Mert Yuksekgonal

Abu Abid