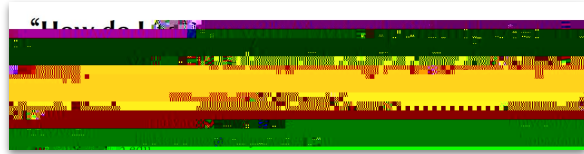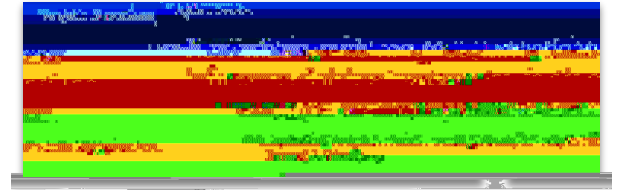AI explanations can easily manipulate user's trust [1]

Explanations cannot help users detect potential model biases [2]
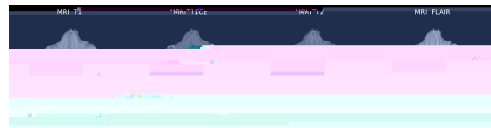
Explanations worsen physicians' task performance [3]
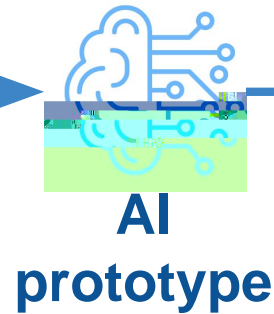


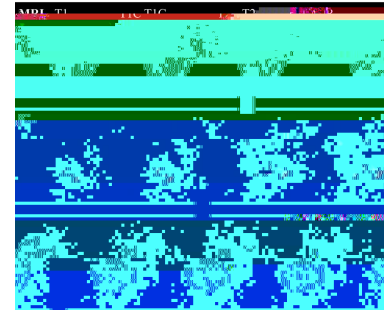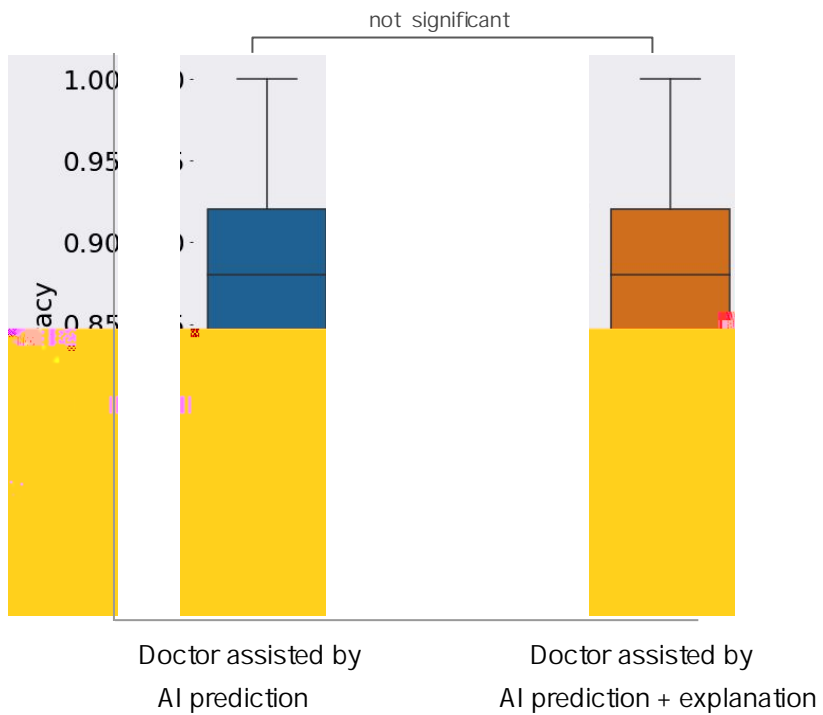How can we make the AI explanations work as they are supposed to **?**

**Input MRI**

**AI prototype**

**AI's suggestion:** Grade 4 glioblastoma

**AI's explanation:**

Quantitative results
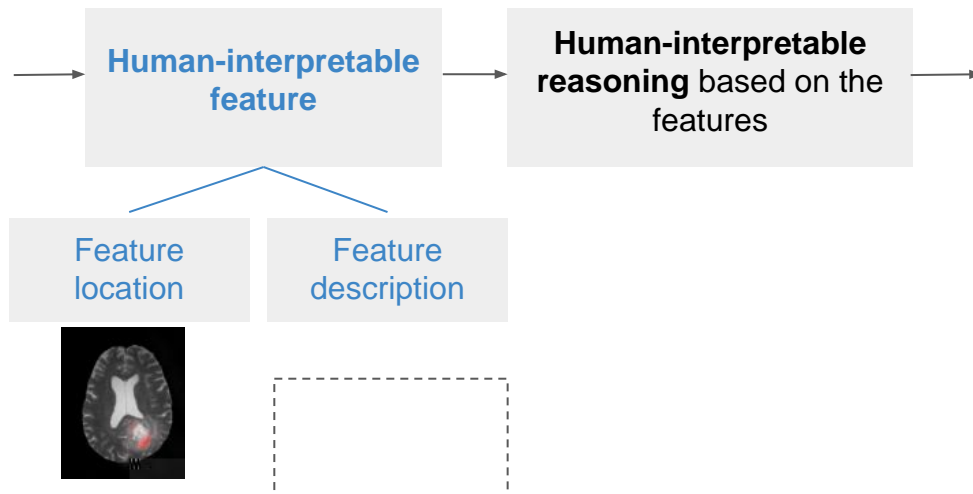
AI explanations are

to improve doctors' task performance

But why ?

Doctor assisted by
AI prediction

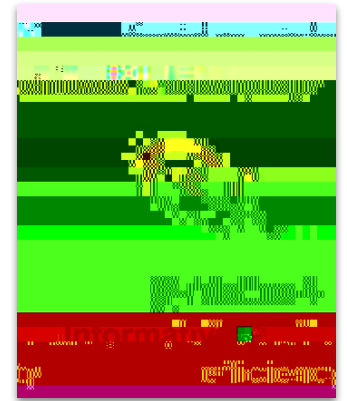Doctor assisted by
AI prediction + explanation

**"**

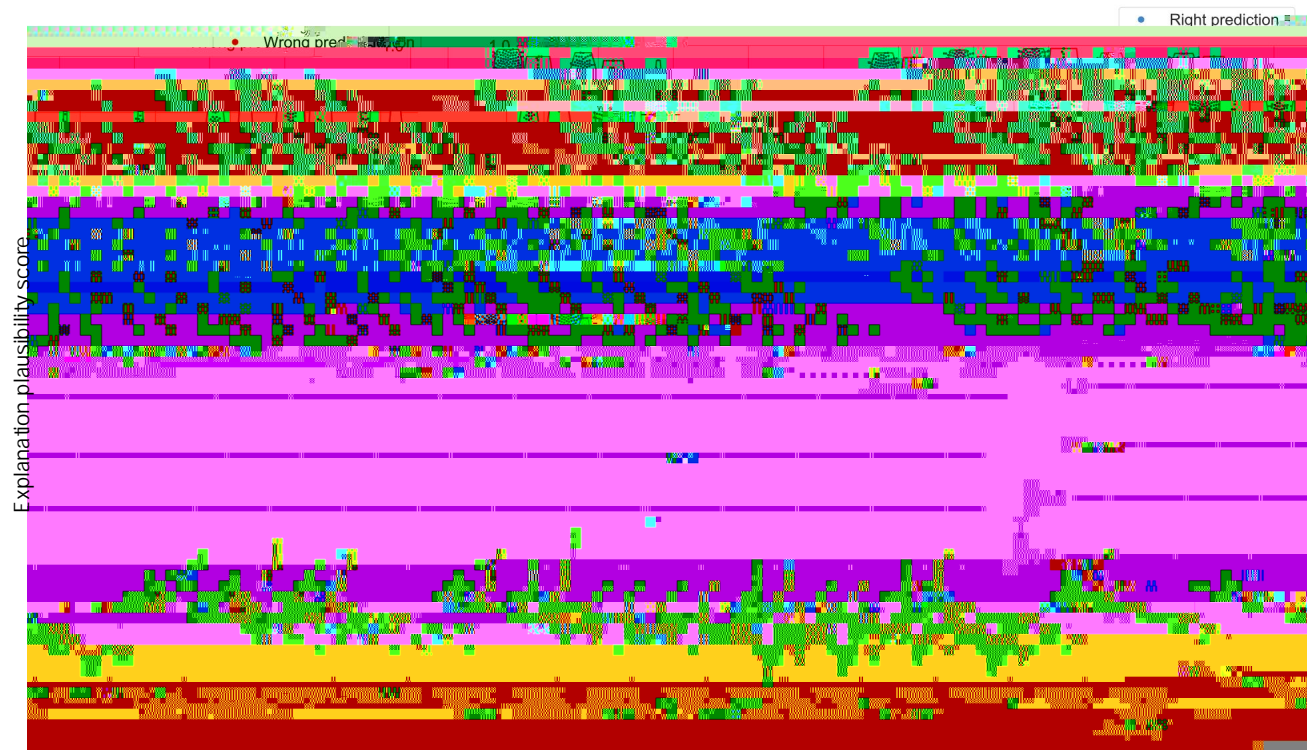What (explanation) we get currently, when a radiologist read it, they **point out the significant features**, and then they **integrate those knowledge**, and say, to my best guess, this is a glioblastoma. And I have the same expectations of AI (explanation).

– Neurosurgeon #3



**Human-interpretable feature**

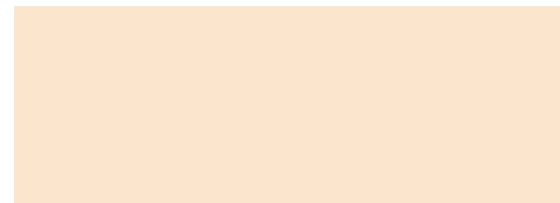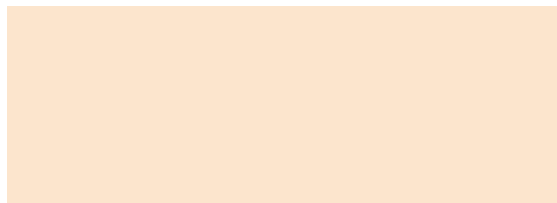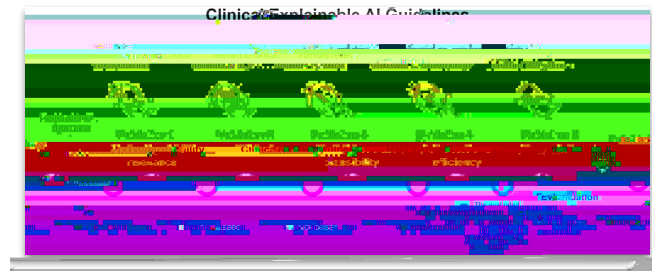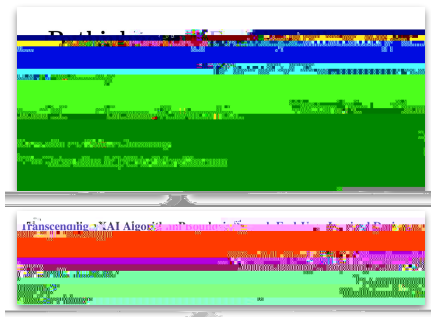**Human-interpretable reasoning** based on the features

Feature location

Feature description

# Thank you!

---

How technologies are ignoring values from underrepresented groups and how we combat it

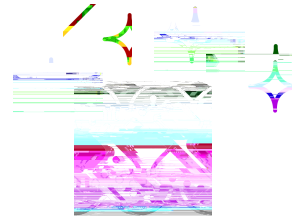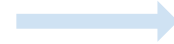| XAI ignores end-users by: | What is it? | Why is it harmful? | How we combat it? |
|---|---|---|---|
| 1. Not aligning with **human reasoning** and interpretation patterns with explanation | Explanations have incomplete feature description<br>only feature localization or text description, not both | Users can hardly incorporate evidence from explanations into their decision process | Design new XAI techniques to provide explanation with complete feature description<br>[Work in progress] |
| 2. Not following **human communication norms** with explanations | Explanations are created to be plausible *regardless of AI decision correctness* | Users in critical tasks can have worse performance that harms people's life, money, etc. | Reveal to the XAI community such ill practice and its harmfulness [1] |
| 3. Not being designed to fulfill users' **utility of explanation** | XAI algorithms are not designed for its utility to end-users, e.g., verifying AI decisions, ensuring AI safety, and improving human-AI performance | Cannot effectively help uses to solve their problems when seeking explanations | Propose user-centered XAI evaluation objectives and metrics [2,3] |

Explainability needs to be carefully crafted based on end-user-centered requirements.

Explainable AI

Trustworthiness

Accountability
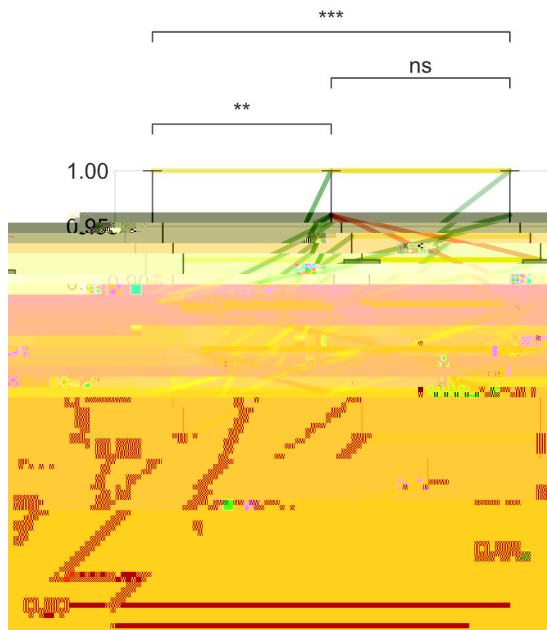
Improving task performance

How technologies are ignoring values from underrepresented groups and how we combat it

**Team:**

Doctor alone — Doctor assisted by AI prediction — Doctor assisted by AI prediction + explanation

*** ns ** 1.00 0.95

# How technologies are ignoring values from underrepresented groups and how we combat it

**Team:**

**Advisors:**

We want AI to be safe, reliable, and accountable to a...

> Algorithms are designed by people, and                              . It's rarely intentional—but this doesn't mean we should let data scientists o   the hook. It means we should be critical about and vigilant for the things we know can go wrong. If we assume discrimination is the default, then we can design systems that work toward notions of equality. [1]