

A hand pointing at a futuristic digital interface with a robotic hand in the foreground. The background is a blue-toned digital dashboard with various charts and data visualizations. The text is overlaid on the center of the image.

Trustworthy and Socially Responsible AI

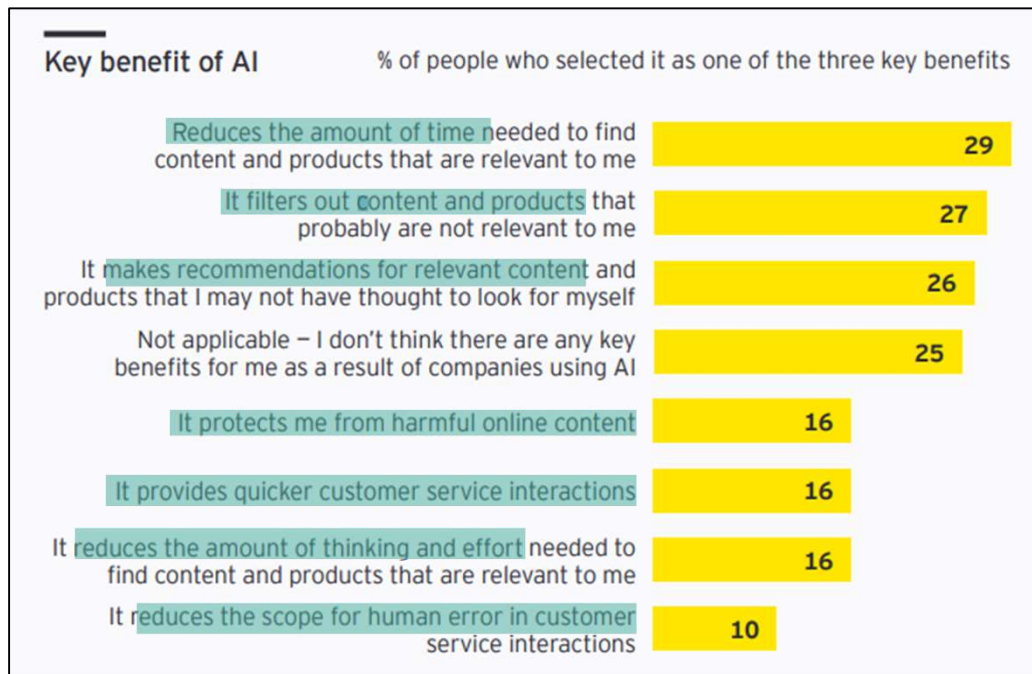
Transparent and eXplainable AI

Anahid Jalali

ETHICS AND AI POLICIES

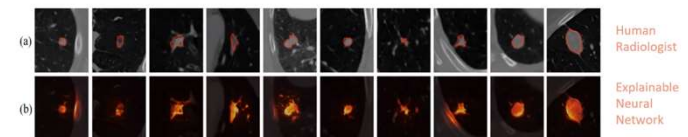
- **AI Policy Observatories and recommenders**
 - **Universal Guidelines for Artificial Intelligence**
 - **OECD AI Principles**
 - **UNESCO Recommendations**
 - **EU AI Act**
 - **COE AI Treaty**
- **Defining redlines and limitations for safety and privacy**
- **Accountability**
- **Clear and transparent documentation of AI development**
- **Human centered values**
- **Sustainable and green**

BENEFITS & RISKS OF AI

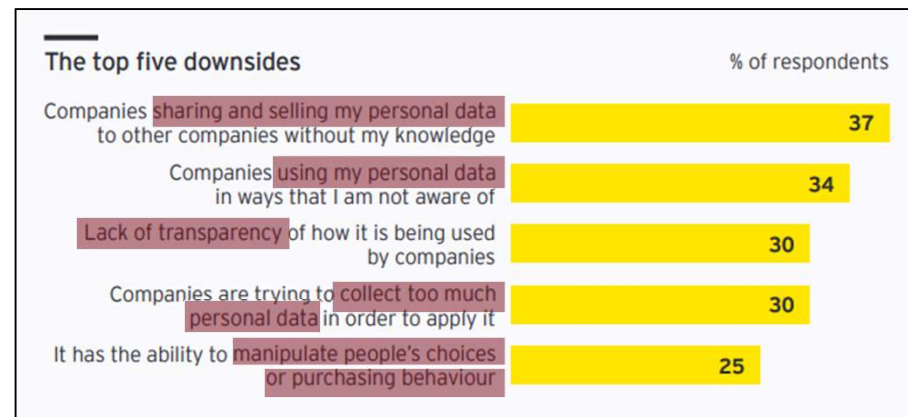
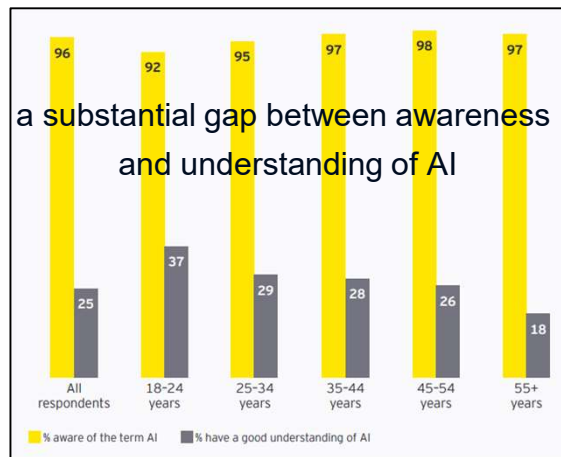


Runtime Explainability: Lung Cancer

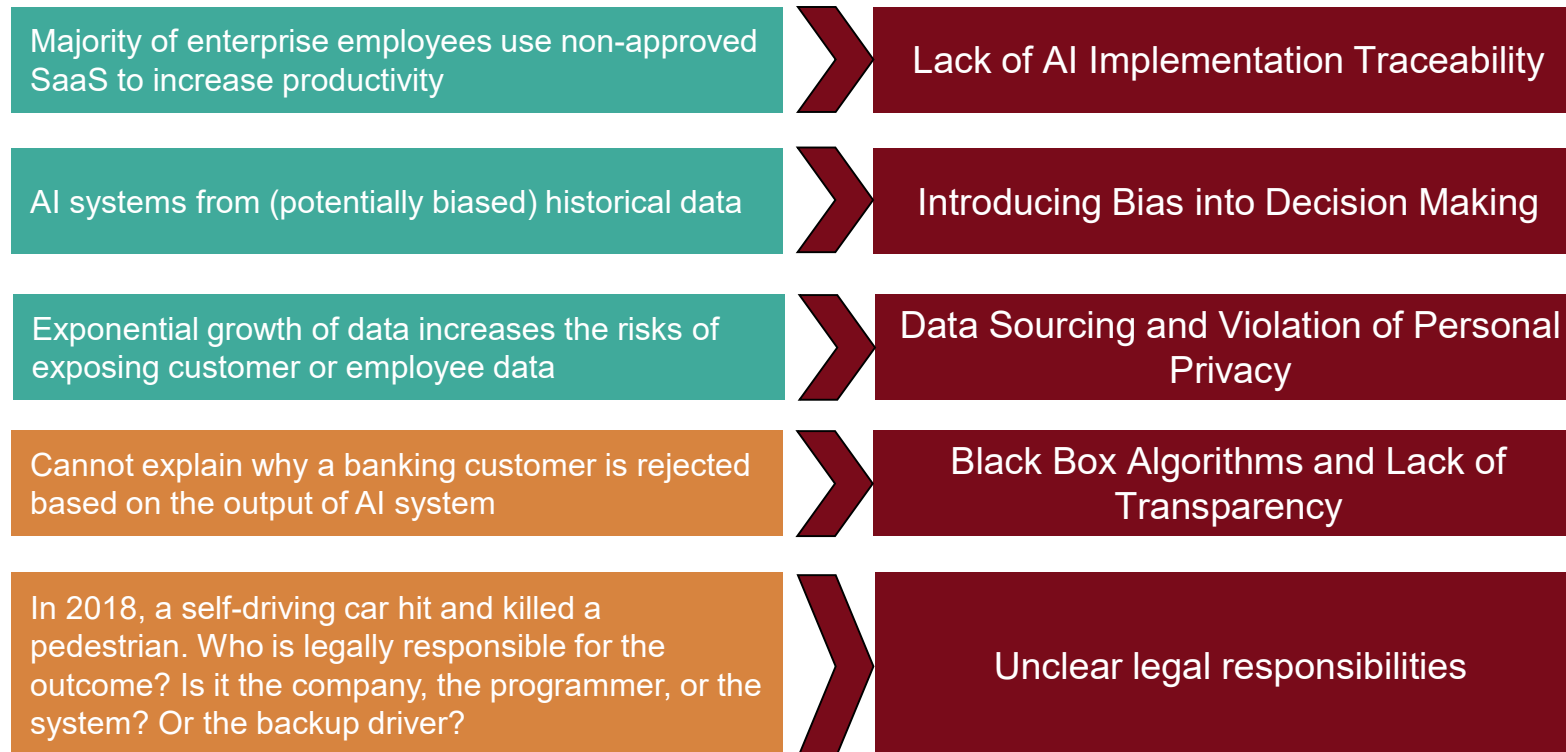
Visual Factors to identify a particular cancer grade for a CT scan



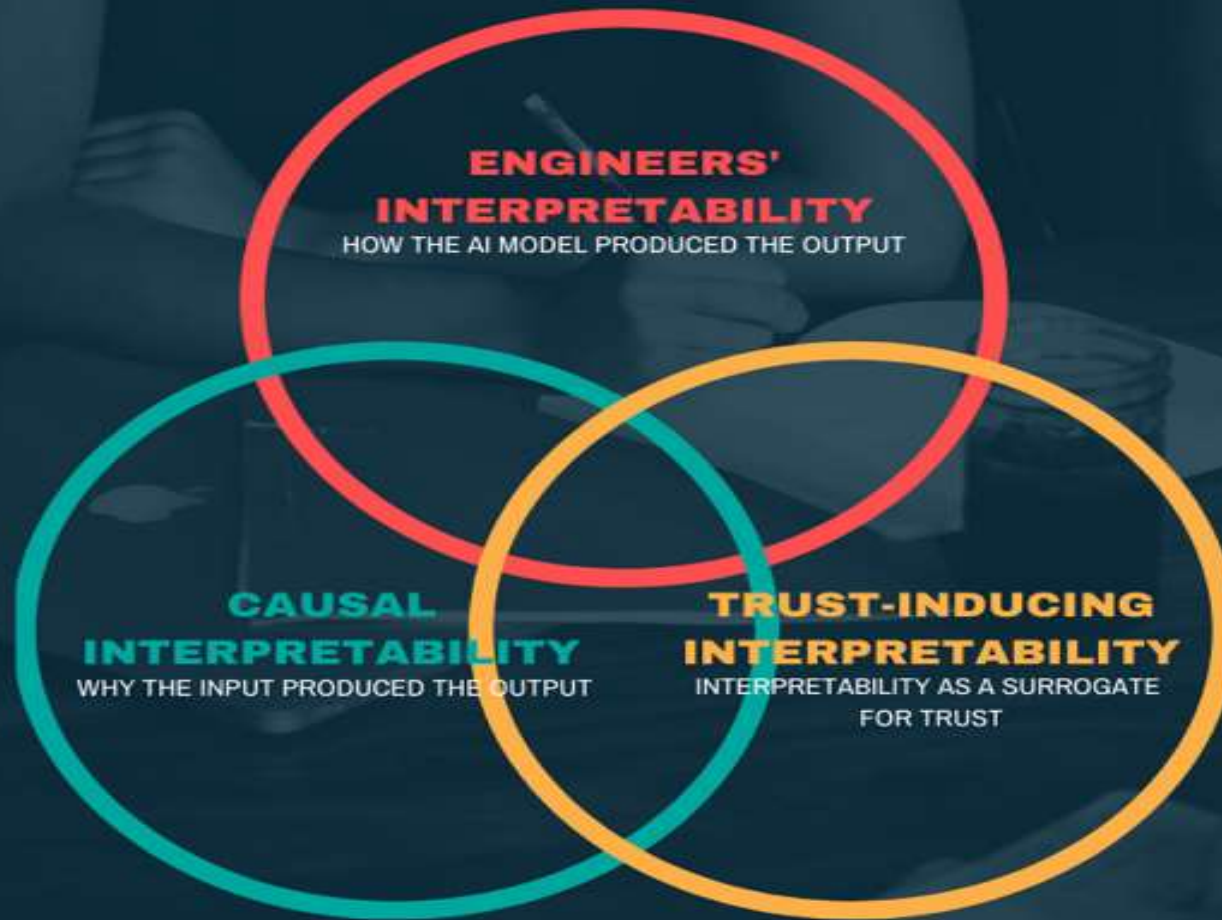
BENEFITS & RISKS OF AI



BENEFITS & RISKS OF AI

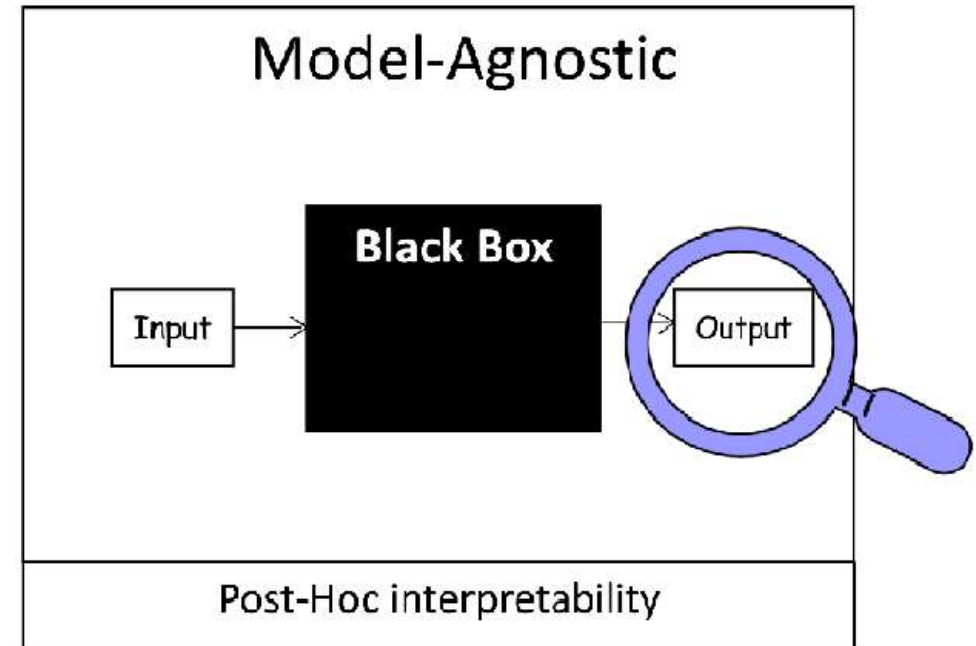
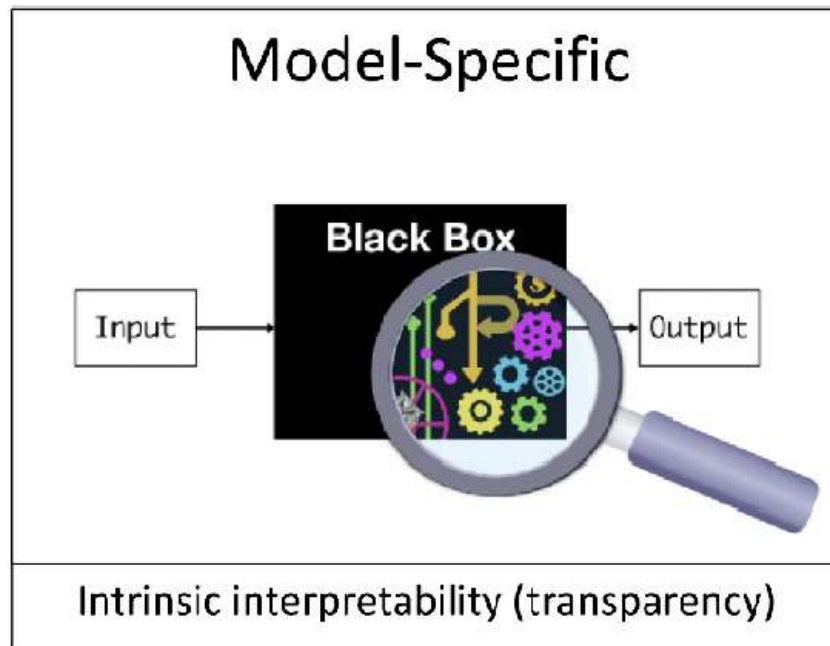


Types of AI Interpretability



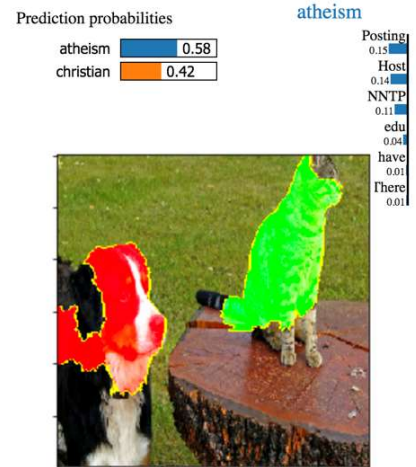
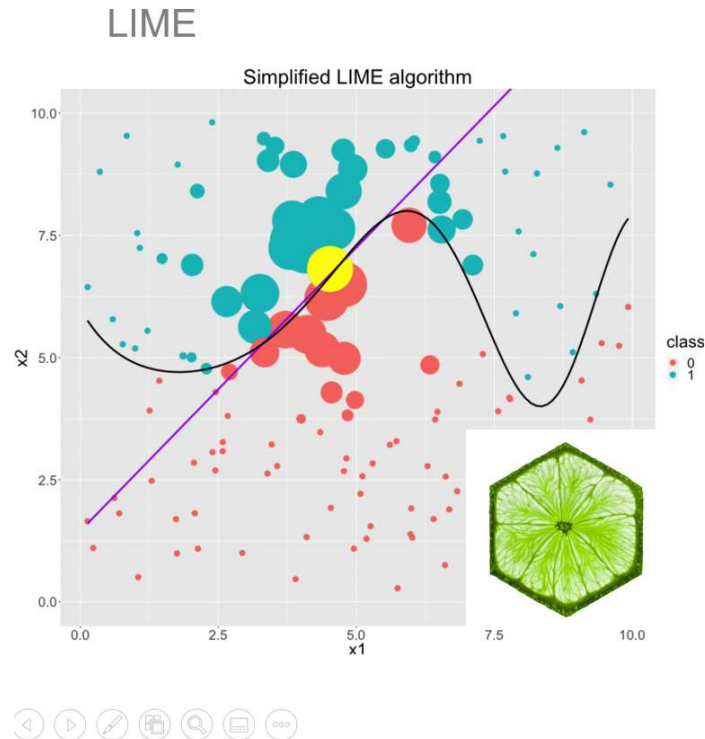
EXPLAINABILITY/TRANSPARENCY/INTERPRETABILITY

- Specific or Agnostic?



Ex-ante - data statistics, bias in data, definition of task, attributes used, scaling of attributes

EXPLAINABILITY/TRANSPARENCY/INTERPRETABILITY

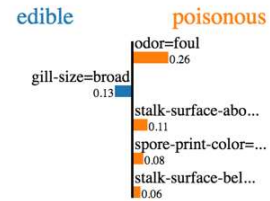
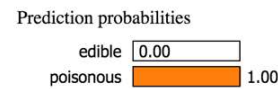


Text with highlighted words

From: johnchad@triton.unm.edu (jchadwic)
Subject: Another request for Darwin Fish
Organization: University of New Mexico, Albuquerque
Lines: 11
NNTP-Posting-Host: triton.unm.edu

Hello Gang,

There have been some notes recently asking where to obtain the DARWIN fish. This is the same question I have and I have not seen an answer on the net. If anyone has a contact please post on the net or email me.



Feature	Value
odor=foul	True
gill-size=broad	True
stalk-surface-above-ring=silky	True
spore-print-color=chocolate	True
stalk-surface-below-ring=silky	True

EXPLAINABILITY/TRANSPARENCY/INTERPRETABILITY

“What is vital is to make anything about AI explainable, fair, secure and with lineage, meaning that anyone could see, and very simply see how any application of AI developed and why.”

– Ginni Rometty, IBM CEO (January 2019)

THANK YOU!

anahid.jalali@ait.ac.at

