



## What is Slice discovery methods?

Identifying coherent subsets of data that exhibit higher systematic error than the overall dataset.

Overall Accuracy of the dataset: 88%.

Accuracy of Landbird class: 83%

Landbird on water  
accuracy: 68%  
(Detected by DOMINO)



## Challenges in existing SDMs

- Existing methods needs biases to be annotated
- Existing methods do not incorporate reasoning
- They do not utilize domain knowledge, needed for medical imaging
- Only slices with visual biases are detected. They can not detect the slices containing meta-data biases.

**TL;DR:** LADDER uses LLMs to identify slices without requiring annotated bias attributes or group labels. Unlike the traditional methods, it identifies both visual and non-visual sources of bias, enabling interpretable diagnosis across the vision pipeline.

## Tracing Bias with LLM, going beyond visual biases

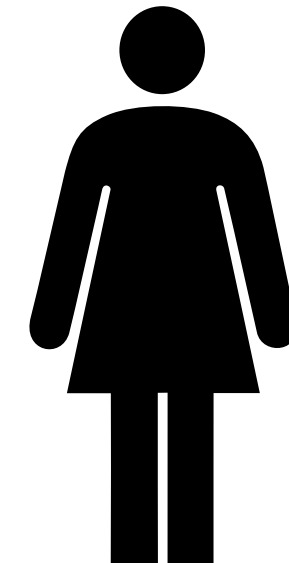
Population



Age: [32-88]

Race: 80% Non-Hispanic White, 20% Asian

Individual



Reason for Visit: [ ....]

Blood Pressure: [...]

Lab Test: [....]

Data

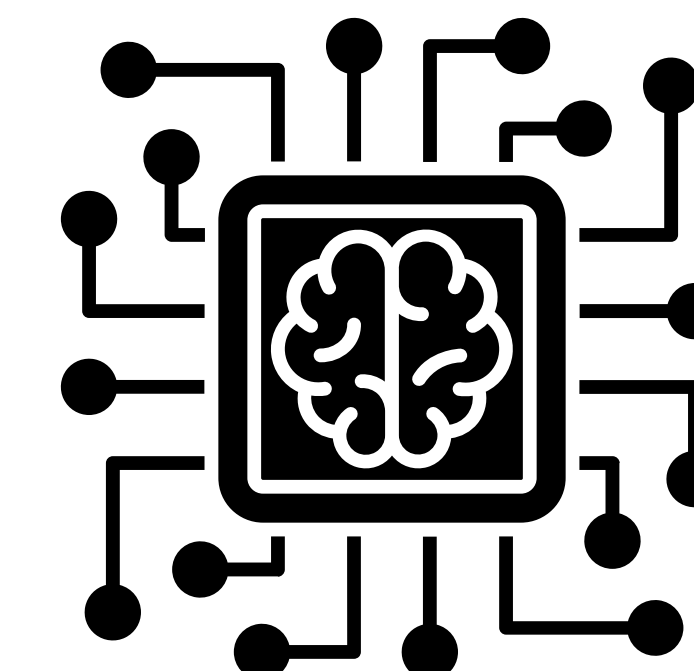


Manufacturer: [ ....]

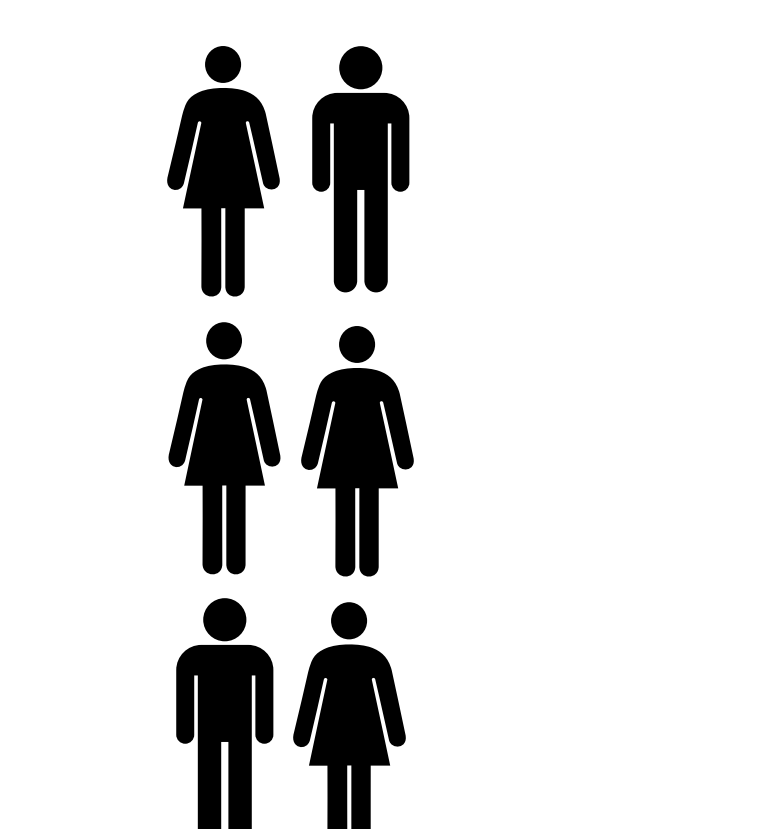
X-ray Dosage: [...]

Aperture Setting: [....]

AI Risk Model



**AI Explanation:** AI model identifies calcification on the left breast and 2mm mass on right breast ...

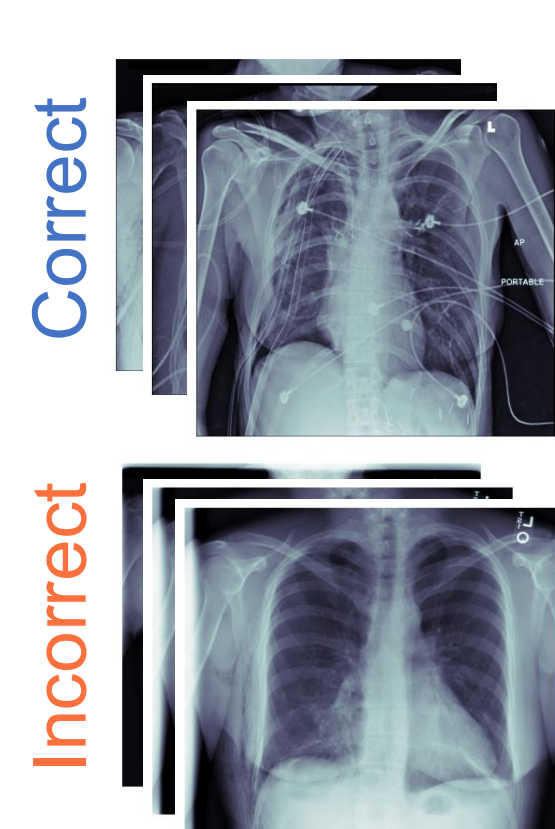


Patient Data (EHR)

Age: [ ....]

Blood Pressure: [...]

Lab Test: [....]

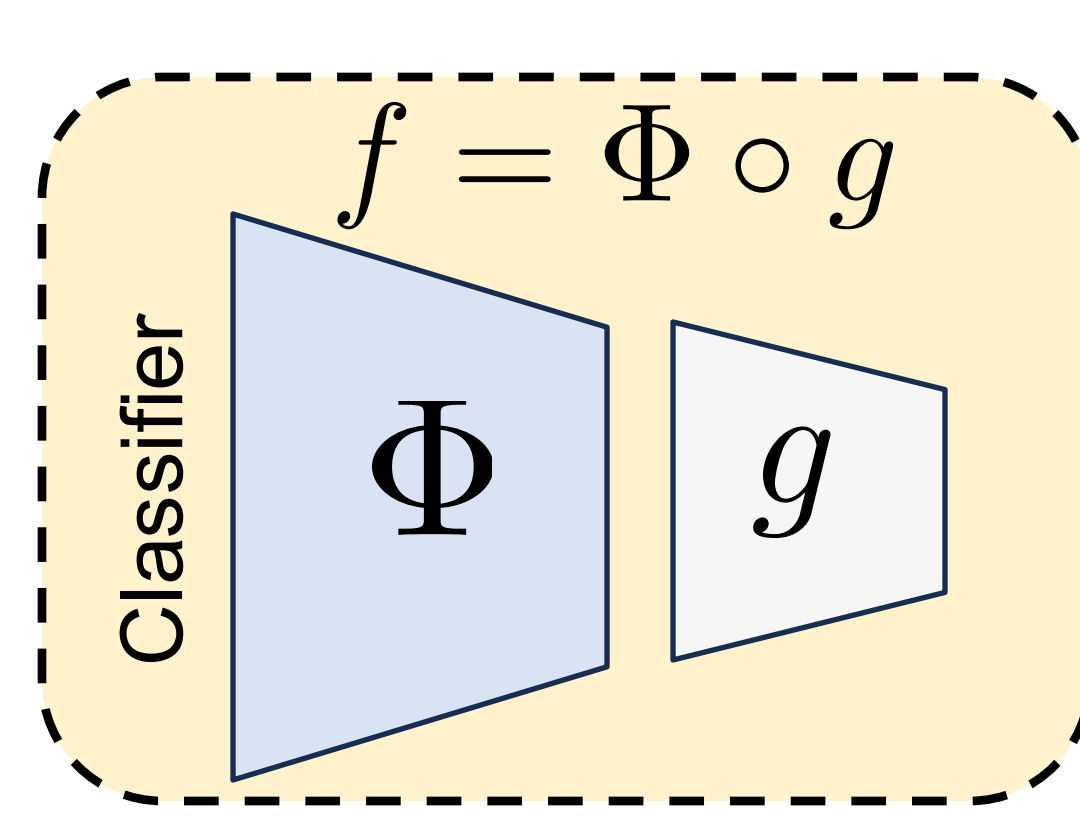


Metadata (e.g. DICOMS)

Manufacturer: [ ....]

X-ray Dosage: [...]

Aperture Setting: [....]



## LADDER

- there is little change in the 3 left chest tubes with area of hydro pneumothorax
- with chest tube remaining in place and no striking change

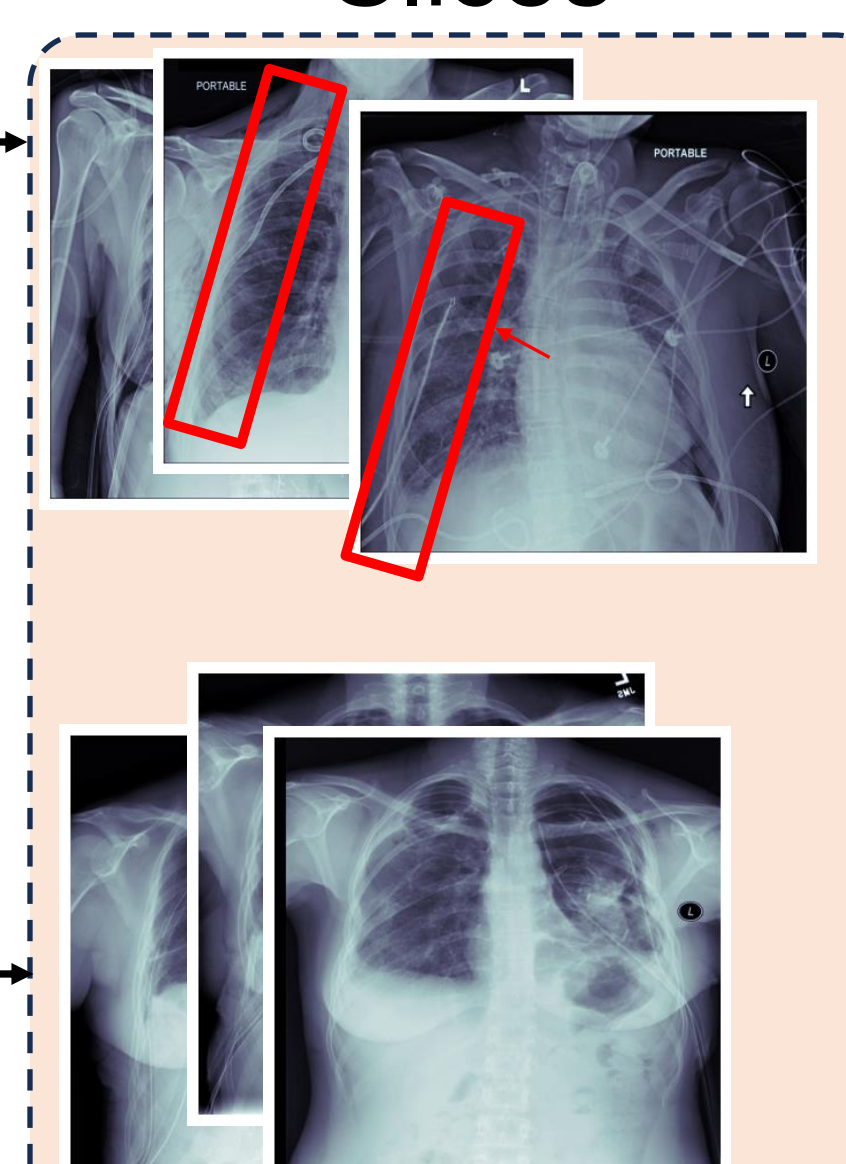
Hypothesis: chest tubes

$\mathcal{T}_{H1}$  : ["Chest X-ray with chest tubes", ..]

Hypothesis: fluid levels

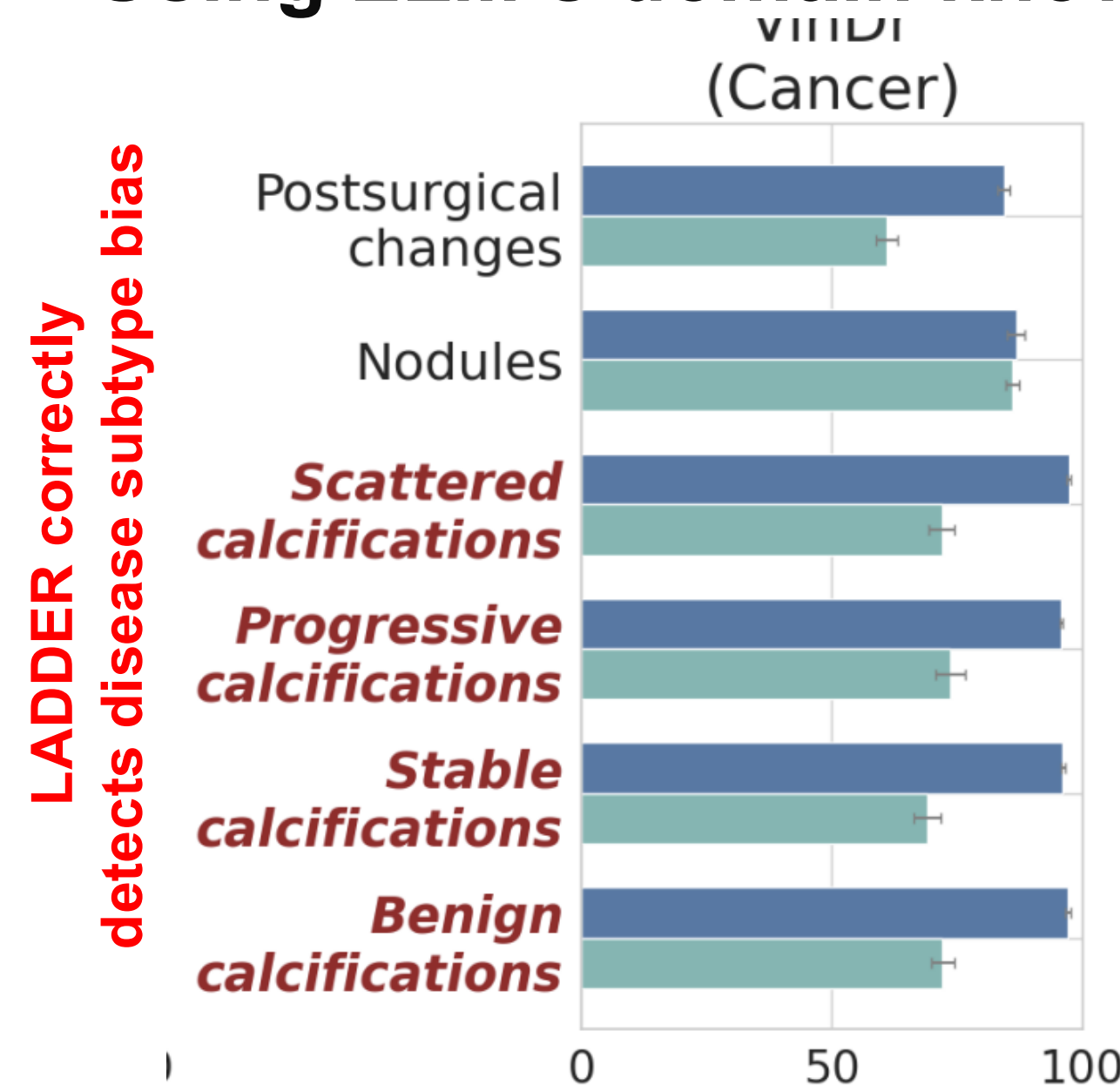
$\mathcal{T}_{H2}$  : ["Chest X-ray with high fluid levels", ..]

Slices

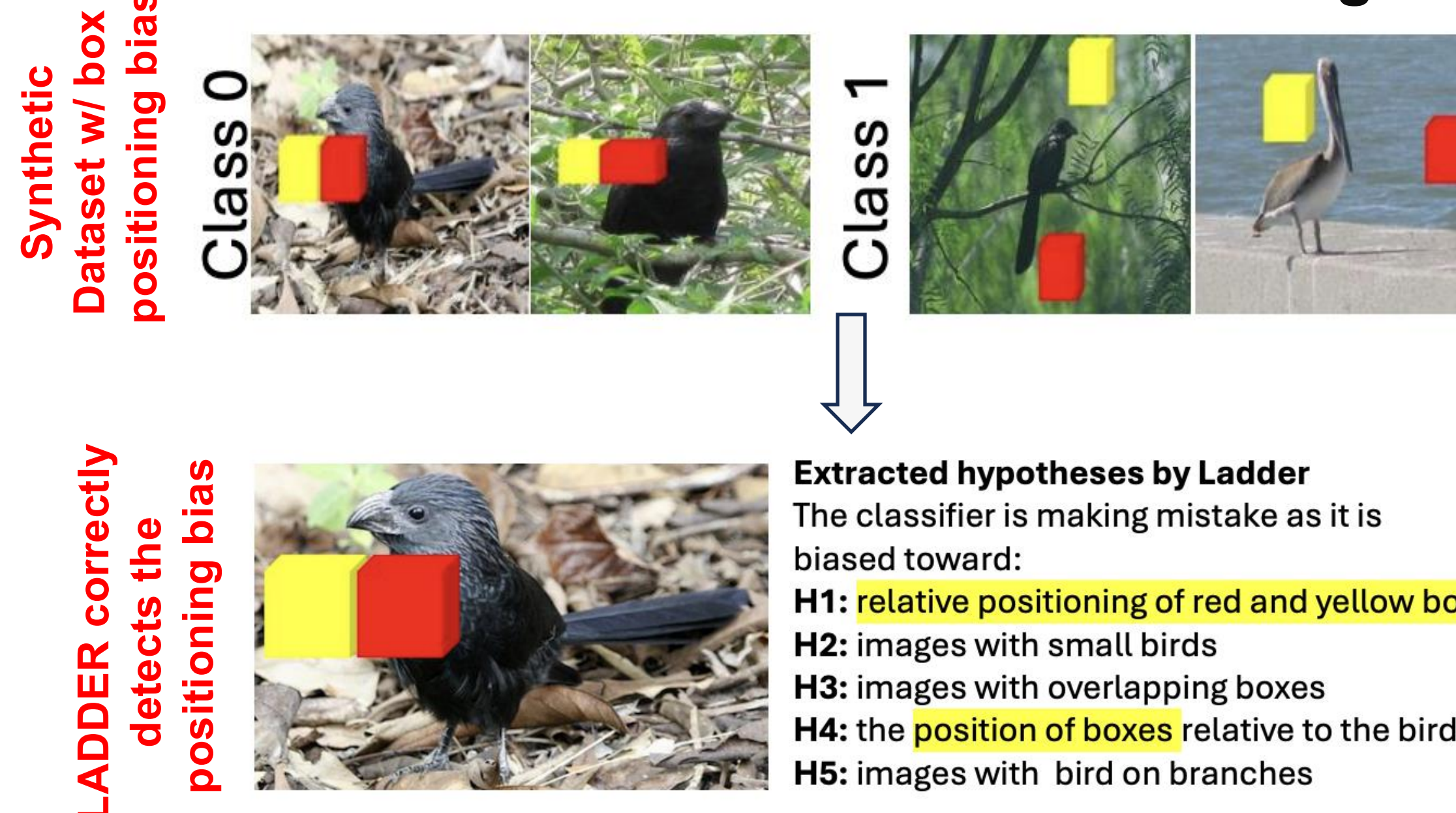


## Detecting visual biases

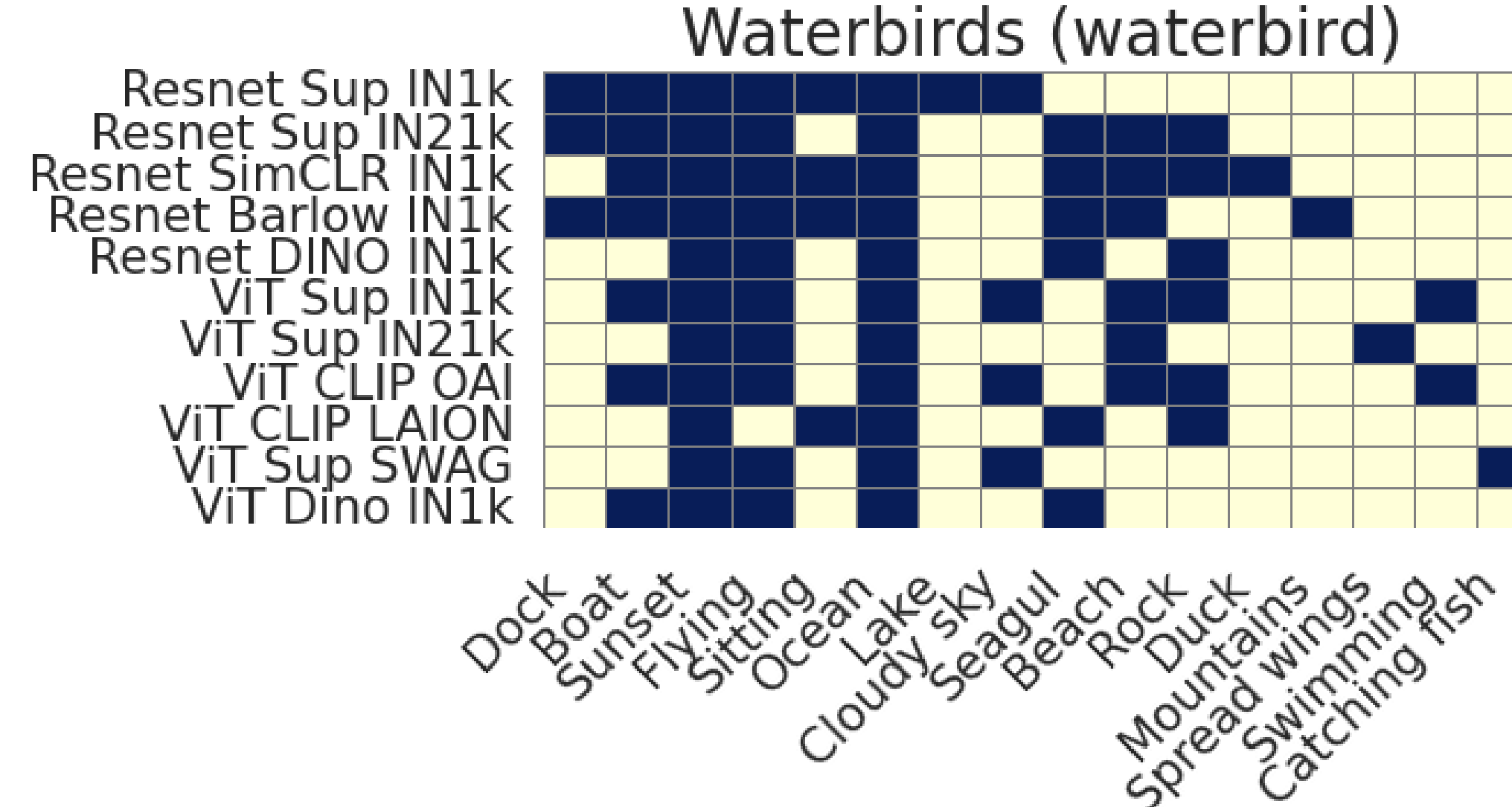
Using LLM's domain knowledge



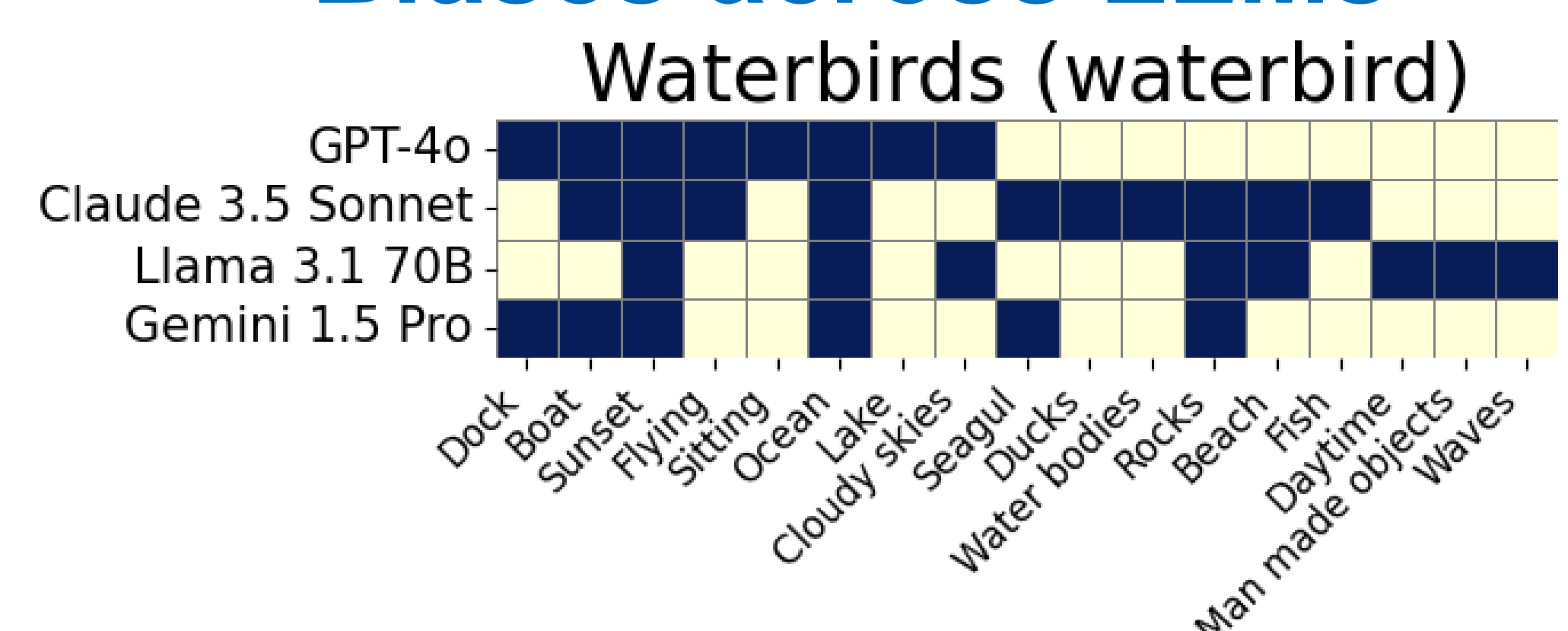
Utilization of LLM-based reasoning



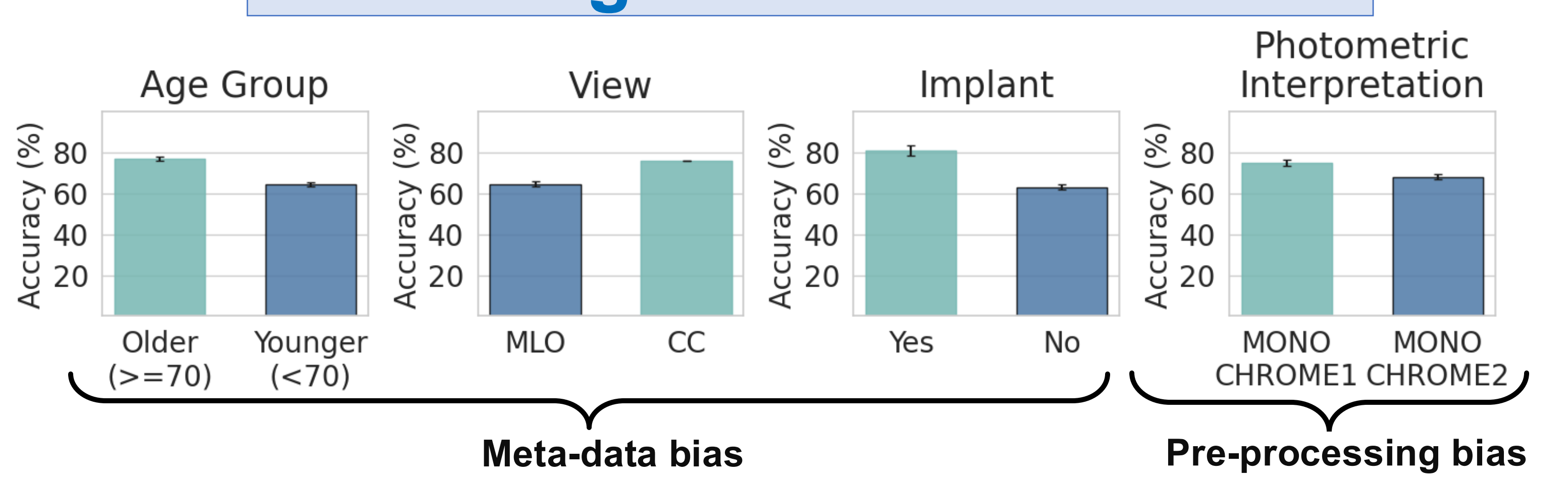
## Biases across architectures



## Biases across LLMs



## Detecting non-visual biases



## More in our paper

- 200+ Classifiers
- 6 Datasets
- GPT-4o as primary LLM
- Using LLaVA to eliminate the need of captions/reports.
- Ablations
  - 4 LLMs
  - 2 SDMs
  - 12 Mitigation methods

## Resources



## Quantitative result

