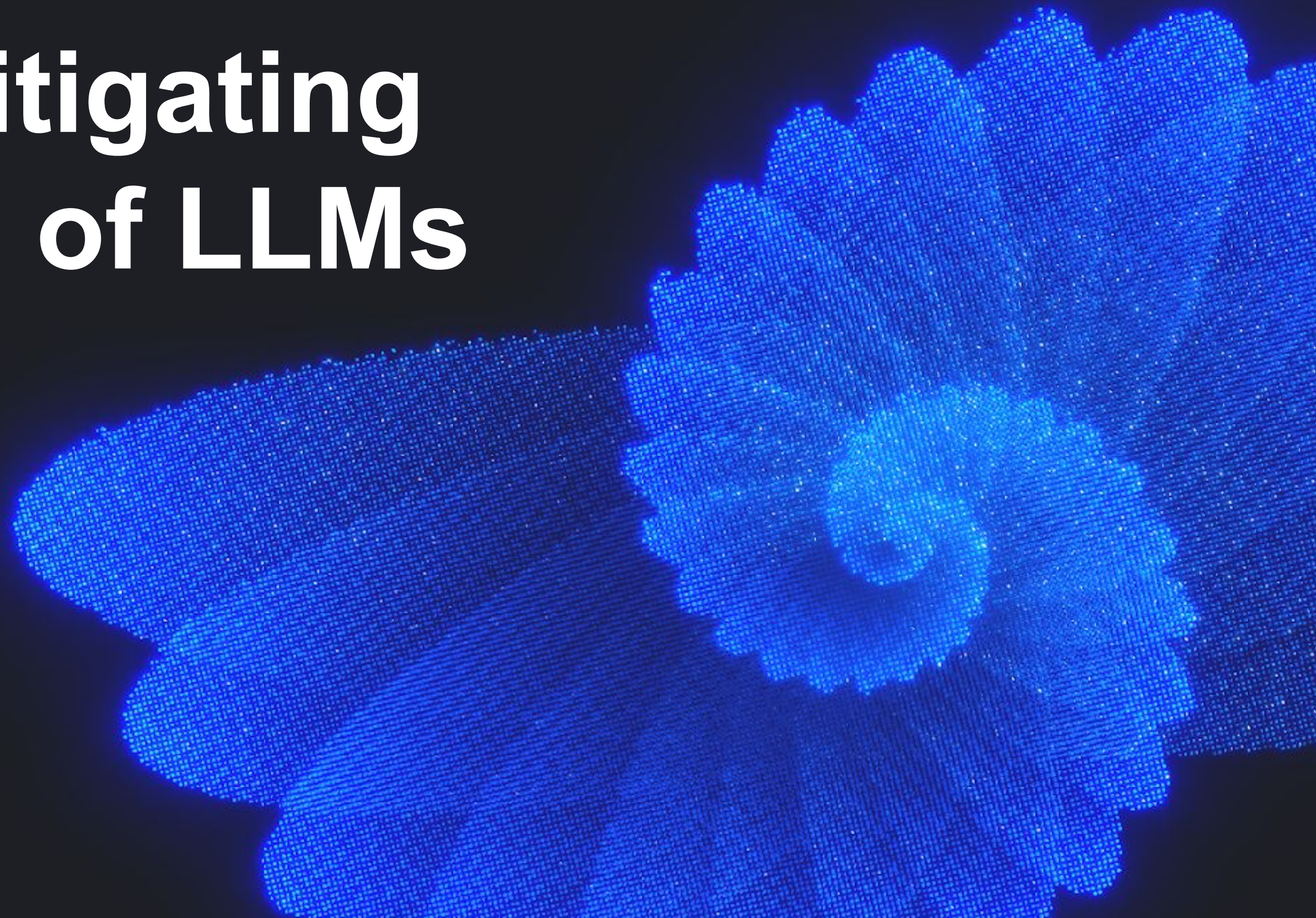




Beyond Buzzwords: Practical Approaches to Mitigating Biases in the Age of LLMs



Dr. ir. Ujwal Gadiraju
@UJLAW





What is this?



Does this image
represent
a **TYPICAL**
[croissant] ?





Does this image
represent
a **TYPICAL**
[toggle button]?







A Typical muffin?

A Typical chihuahua?

Notions of typicality and atypicality ... are distinguishable by “the strength of association between observable properties and concepts.”



Eric Margolis and Stephen Laurence. 2007. The ontology of concepts-abstract, objects, or mental representations? *Noûs* 41, 4 (2007), 561–593.

Bing Ran and P Robert Duimering. 2010. Conceptual combination: Models, theories and controversies. *International Journal of Cognitive Linguistics*.

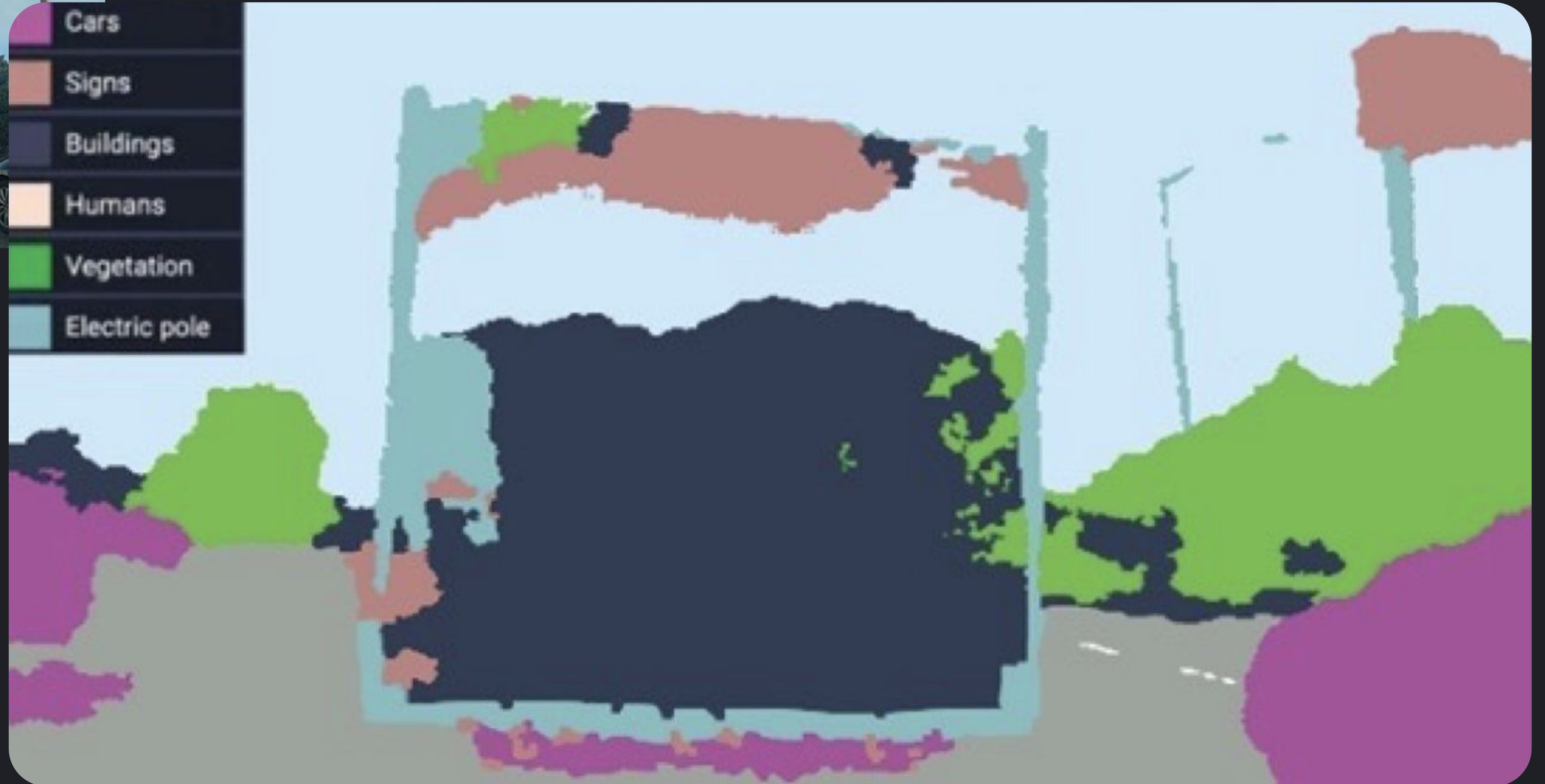
Identifying & Characterizing Errors is Critical



Costly errors. At what cost?



- Cars
- Signs
- Buildings
- Humans
- Vegetation
- Electric pole

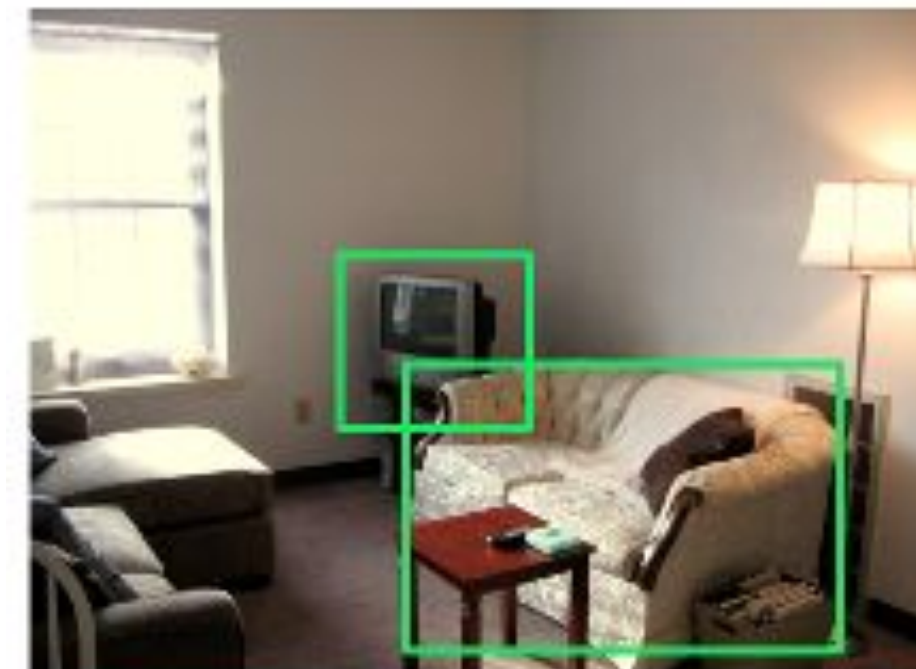


Human vs. Machine Understanding

- Humans perceive the world **discreetly** (object-level), but computer models see the world in a relatively **continuous** form (pixel-level).
- Humans make predictions more **semantically** (mental models); computer models are trained to predict **statistically**.
- Humans reason about real-world entities **using their corresponding context**, but computer models **often ignore contextual cues**.

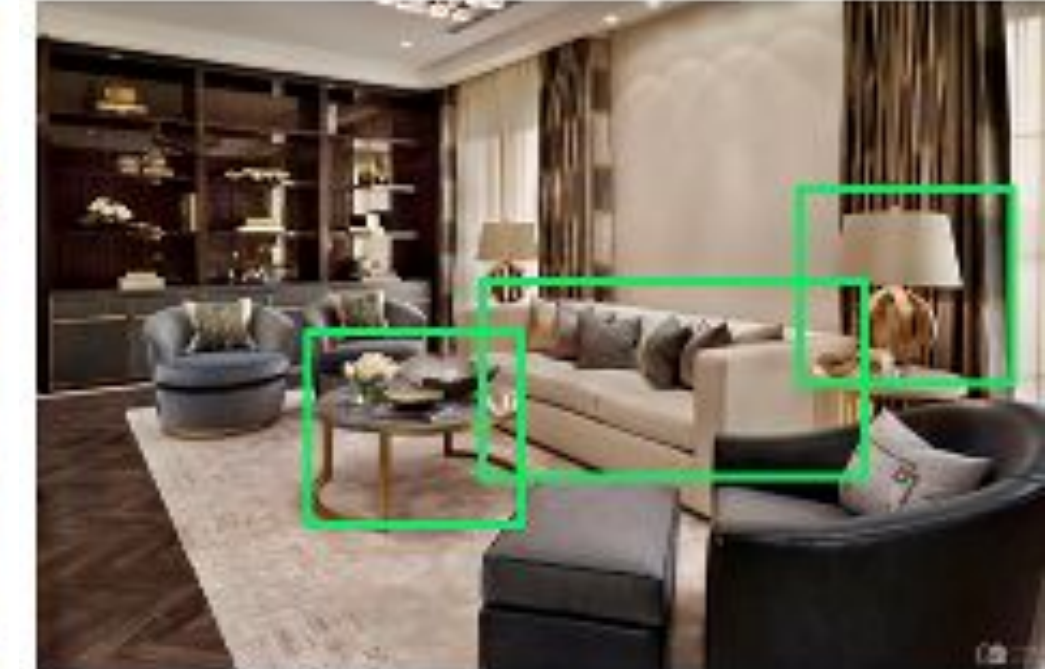
Actual: living room

Predicted: living room



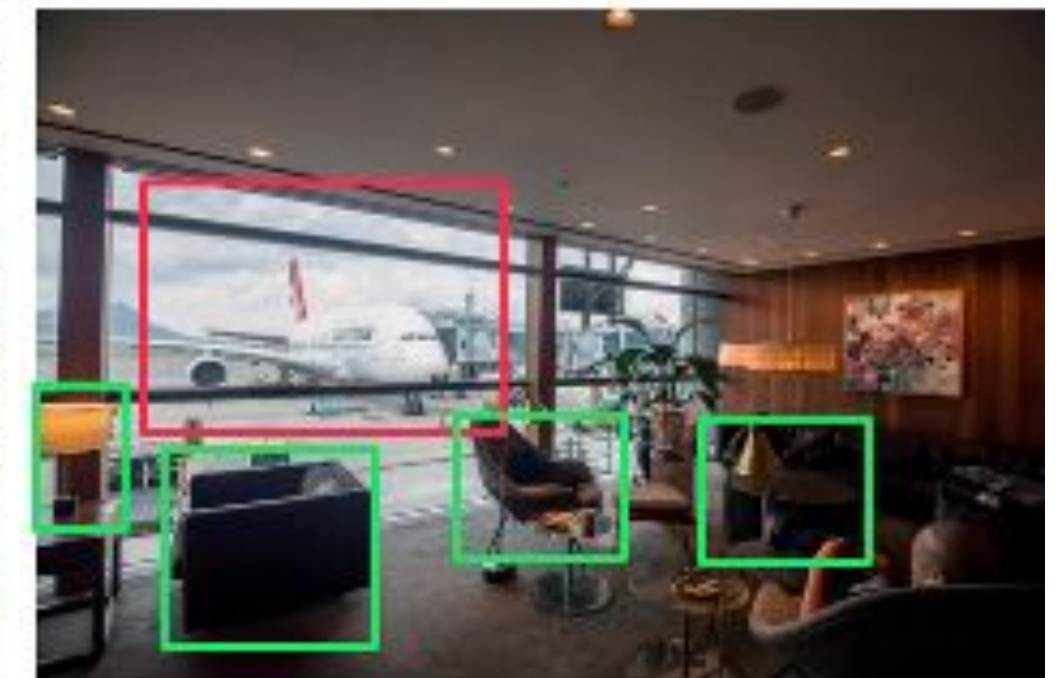
Actual: living room

Predicted: living room



Actual: outdoor

Predicted: living room

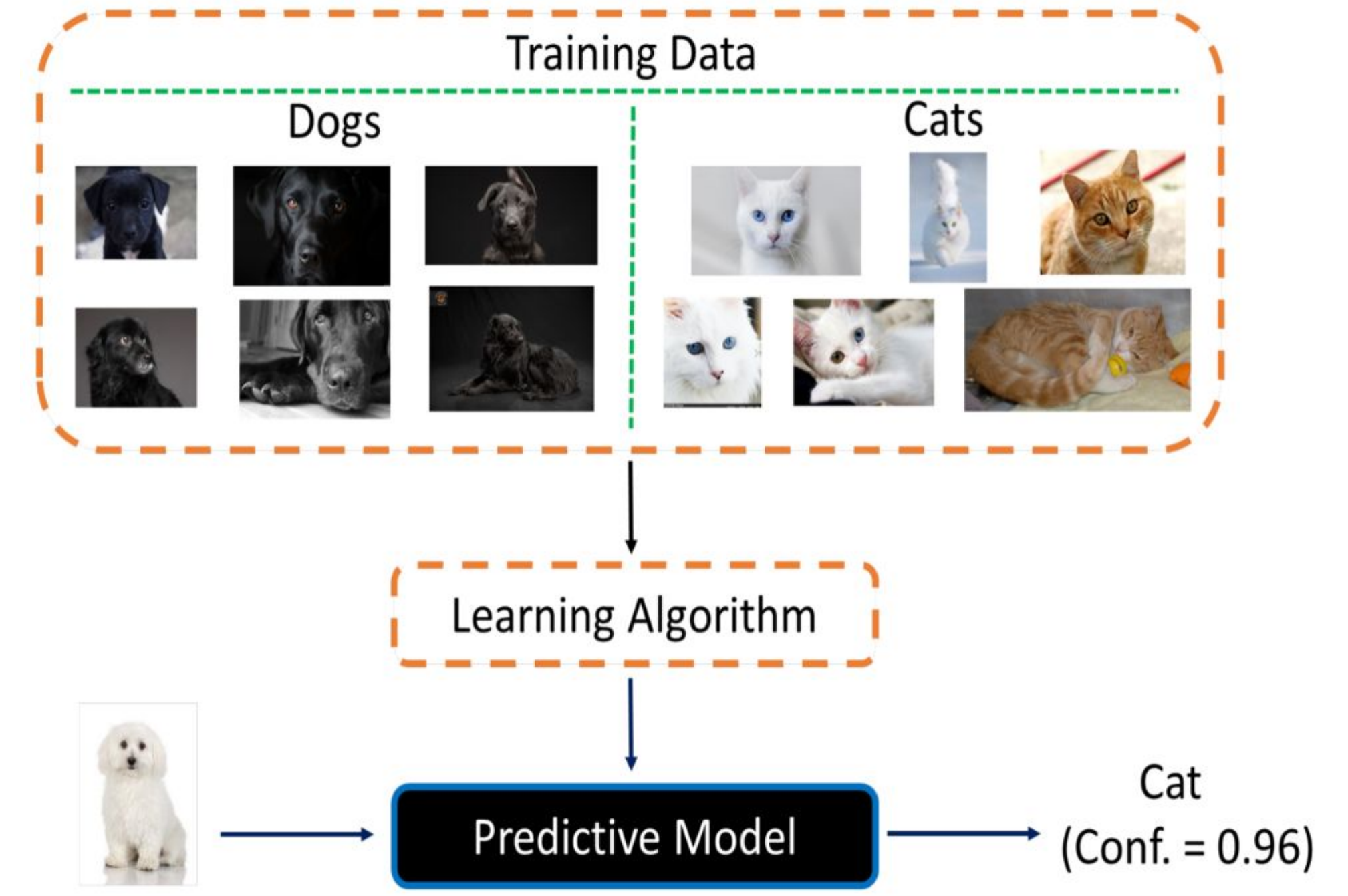


Actual: airport

Predicted: living room

Unknown-Unknowns

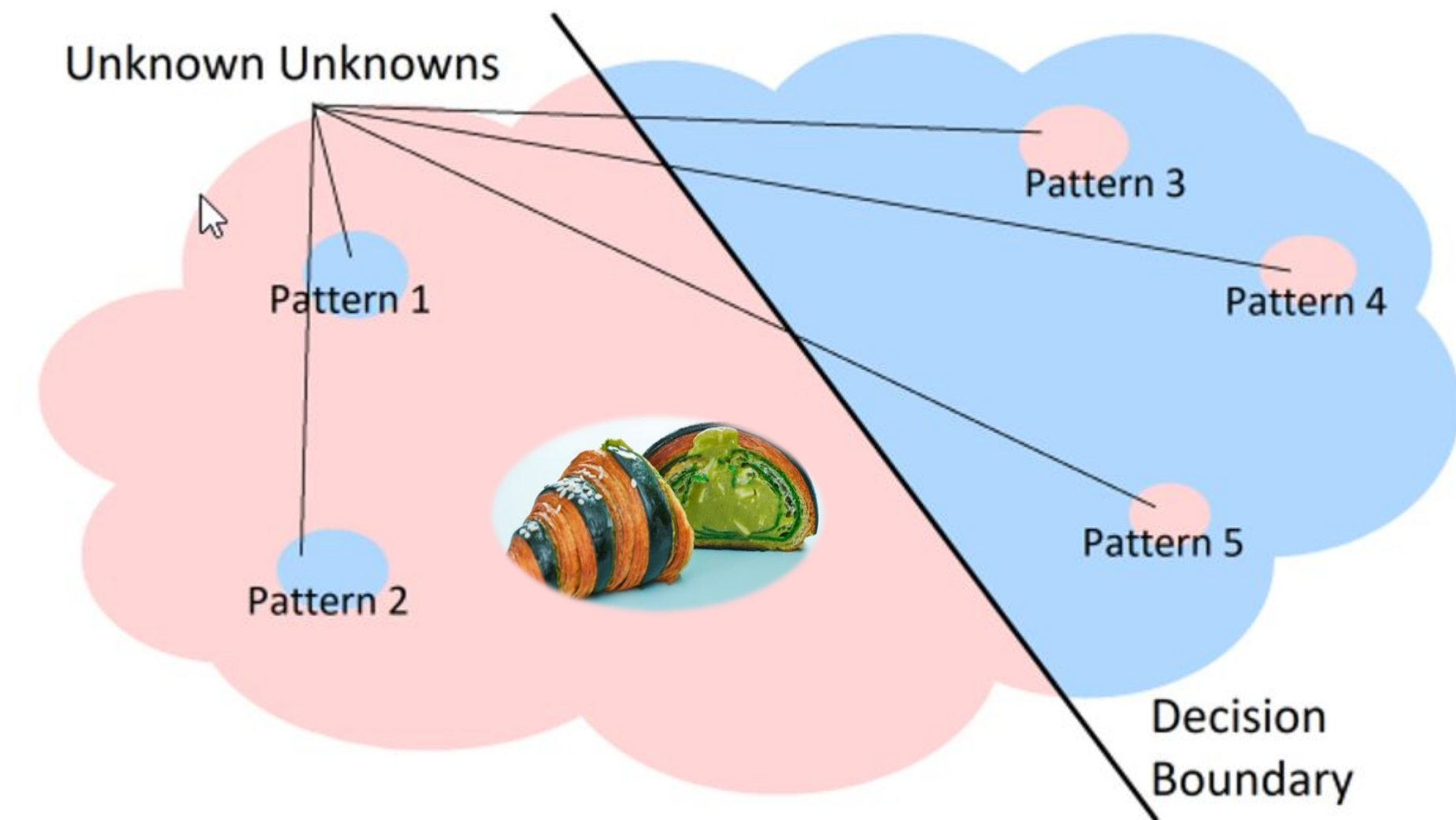
- Over-confident errors
- Caused by systematic biases in the training data
- Hard to discover as ML systems do not provide enough information



(Lakkaraju et al. 2017)

Unknown–Unknowns

- Over-confident errors
- Caused by systematic biases in the training data
- Hard to discover as ML systems do not provide enough information
- Reside in specific partitions of the feature space (blind-spots) and are not distributed evenly across all the feature space



(Z. Liu et al. 2020)

Known–Unknowns in LLMs

- Do LLMs know what they know? And more importantly, are they aware of what they do not know?
- This is an important question to understand the certainty of their statements or prevent such language models from making up facts.

Known Knowns Things we are aware of and understand	Known Unknowns Things we are aware of but do not understand
Unknown Knowns Things we understand but are not aware of	Unknown Unknowns Things we are neither aware of nor understand

(Amayuelas et al. 2023, Knowledge of Knowledge: Exploring Known-Unknowns Uncertainty with Large Language Models)

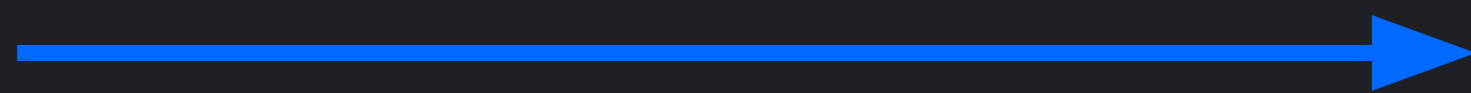
Human input is essential in the evaluation of known-unknowns and the discovery of “unknown-unknowns.”

Need of the hour → proactively discover and characterize unknown-unknowns to build reliable image recognition systems.

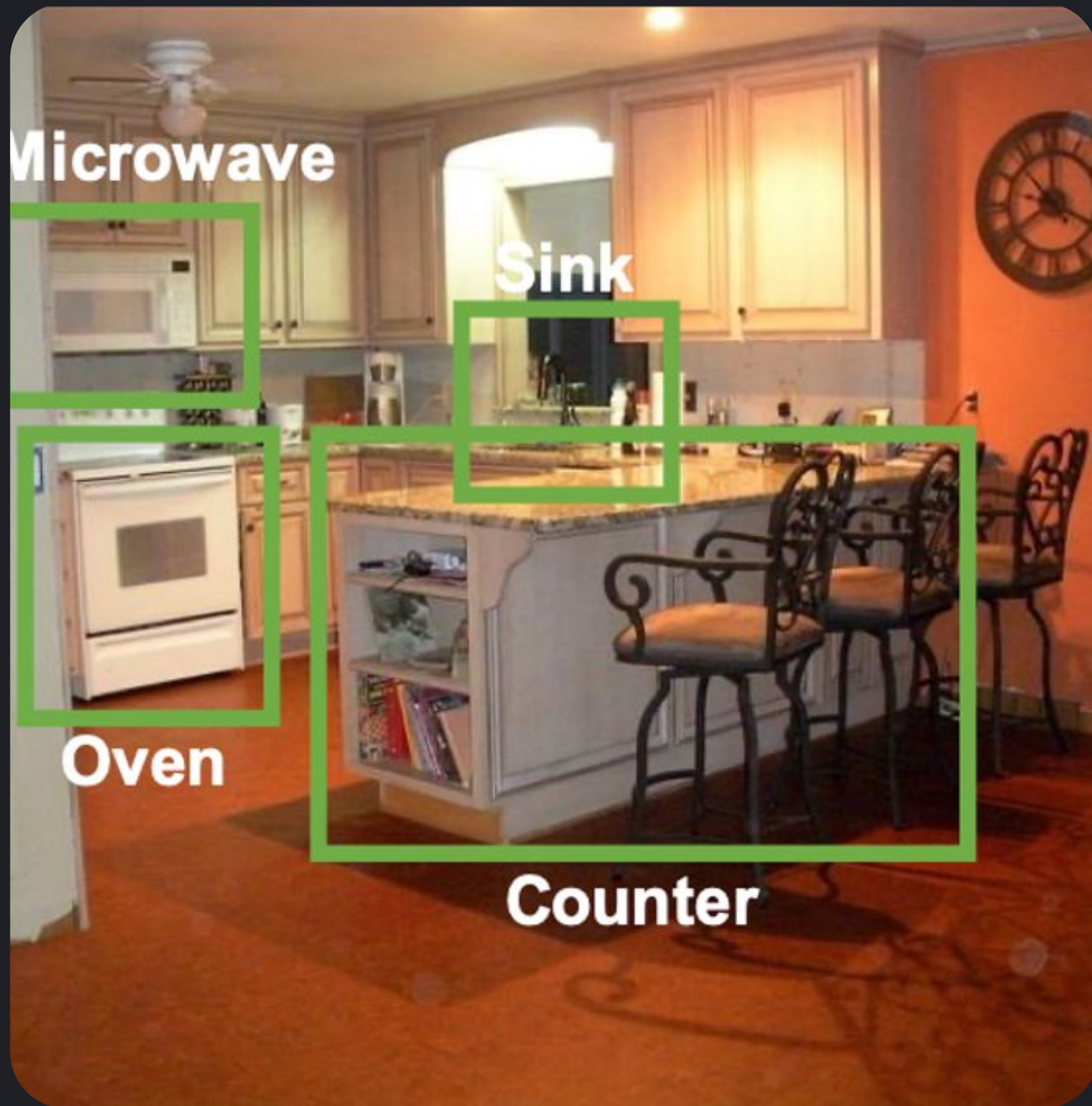
Perspective: Leveraging Human Understanding for Identifying and Characterizing Image Atypicality.

Sharifi et al., ACM IUI 2023

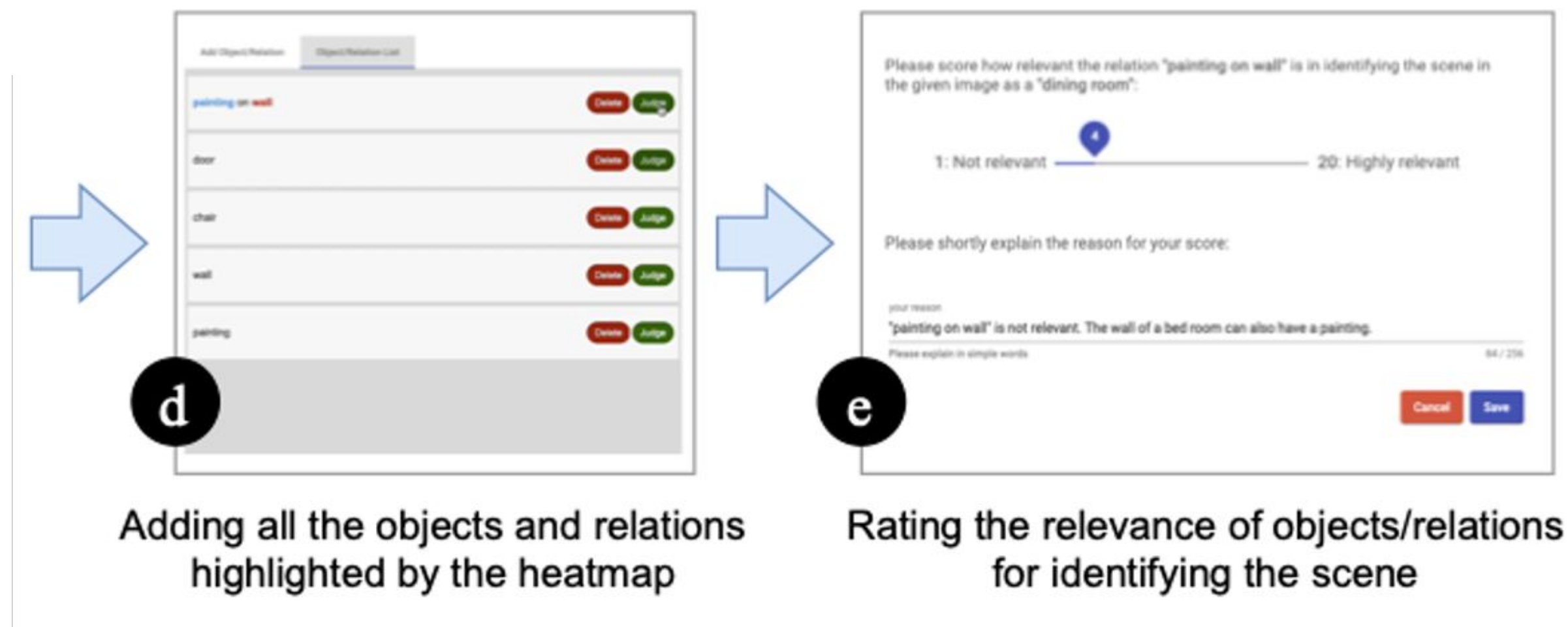
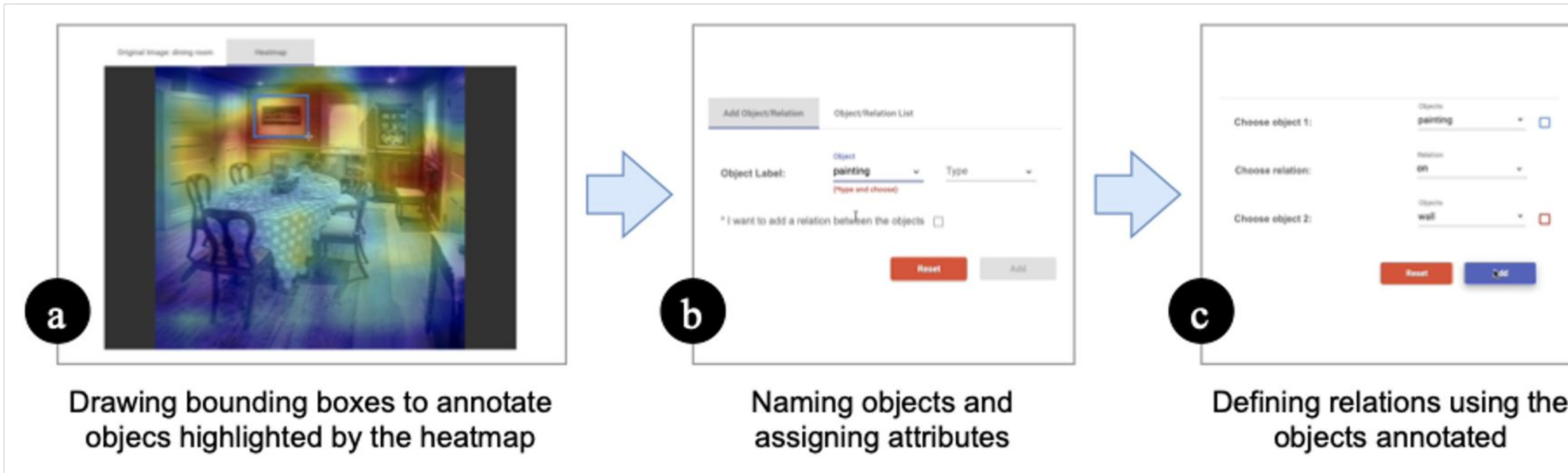
Need for human input!



What Should You Know? A Human-In-the-Loop Approach to Unknown Unknowns Characterization in Image Recognition. Sharifi, S., Qiu, S., Gadiraju, U., Yang, J., & Bozzon, A. In Proceedings of the ACM Web Conference (WWW 2022).



What Should You Know? A Human-In-the-Loop Approach to Unknown Unknowns Characterization in Image Recognition. Sharifi, S., Qiu, S., Gadiraju, U., Yang, J., & Bozzon, A. In Proceedings of the ACM Web Conference (WWW 2022).




Annotations with “PERSPECTIVE”

1. The Target Image

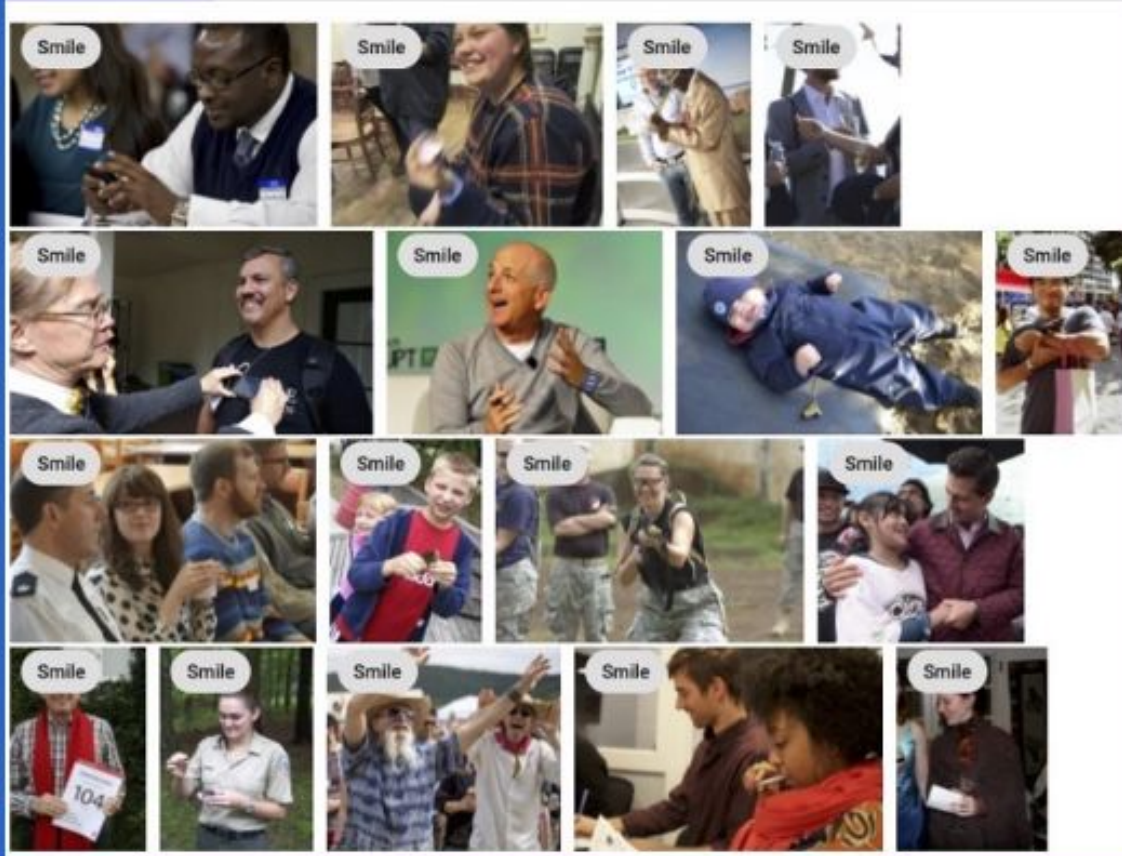
The target image contains a certain object or is about a scene. This initial label (*bird* in this example) is from the given image classifiers.

Task Description: Click [HERE] if you want to read the Task Instructions again.

Original Image: Bird



Visually Similar Random Samples (Bird) Representative Samples (Bird)



Is this image about Bird (or does it contain Bird)? Yes Maybe No

Is this a typical Bird or an unusual one? (please select an option) 1: Highly Typical ————— 7: Highly Unusual

What code do you assign to this image about Bird?

- (11) The aspect ratio of the object of interest is smaller than other images in the class.
- (12) The aspect ratio of the object of interest is larger than other images in the class.
- (13) The majority of the object(s) of interest in comparison to other images in the class is(are) occluded.
- (14) Dominant object in the image belongs to another class(es).
- (15) Object of interest looks unusual with respect to other images in the class. (shape, action)

Submit

3. The Questions

Workers are asked to first decide whether the image label is correct. Then they rate the level of atypicality using a 7-point Likert-scale and give their reasons by choosing the codes.

2. The Auxiliary Images

The auxiliary images are organized in different tabs, each displaying one type of the auxiliary images (visually similar, random sample, or representative sample).

Workers are asked to judge and characterize atypicality by comparing the target image to the auxiliary images.

Task Description: Click [\[HERE\]](#) if you want to read the Task Instructions again.

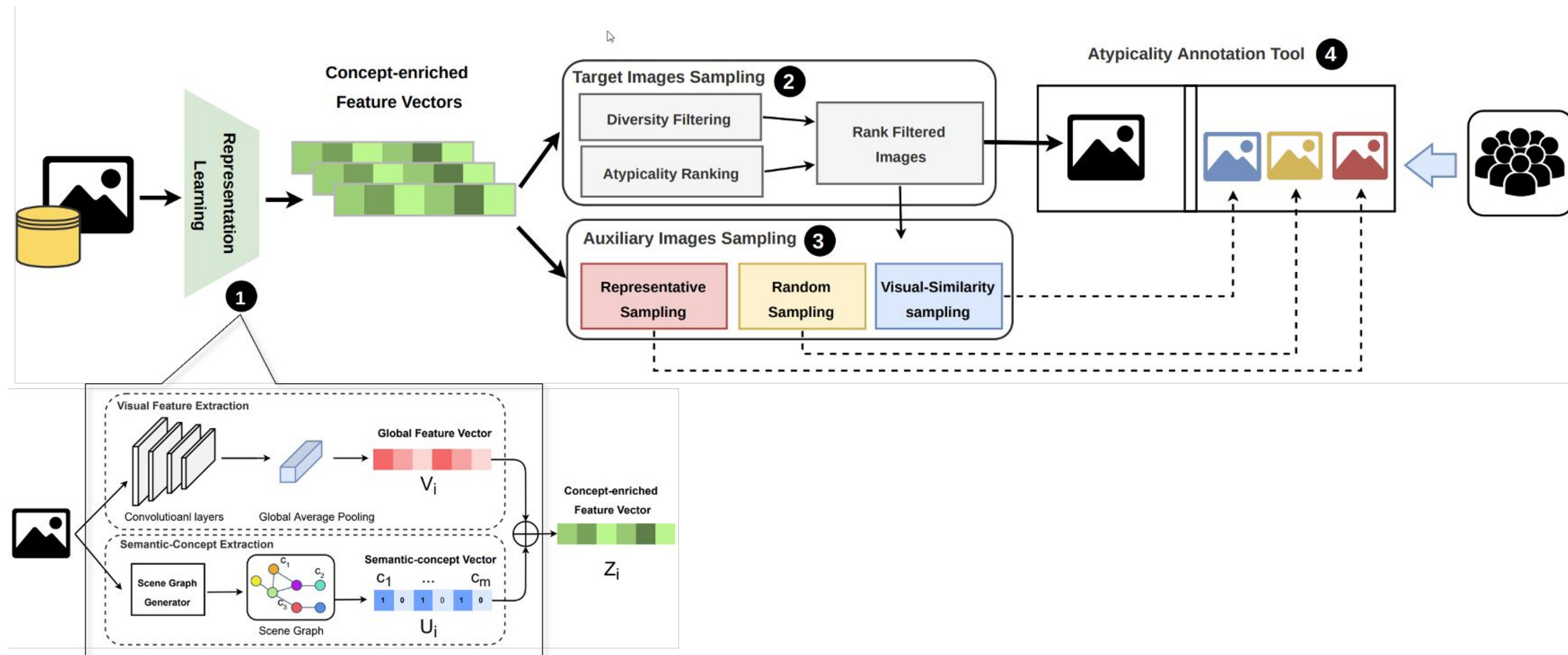
Original Image: Croissant

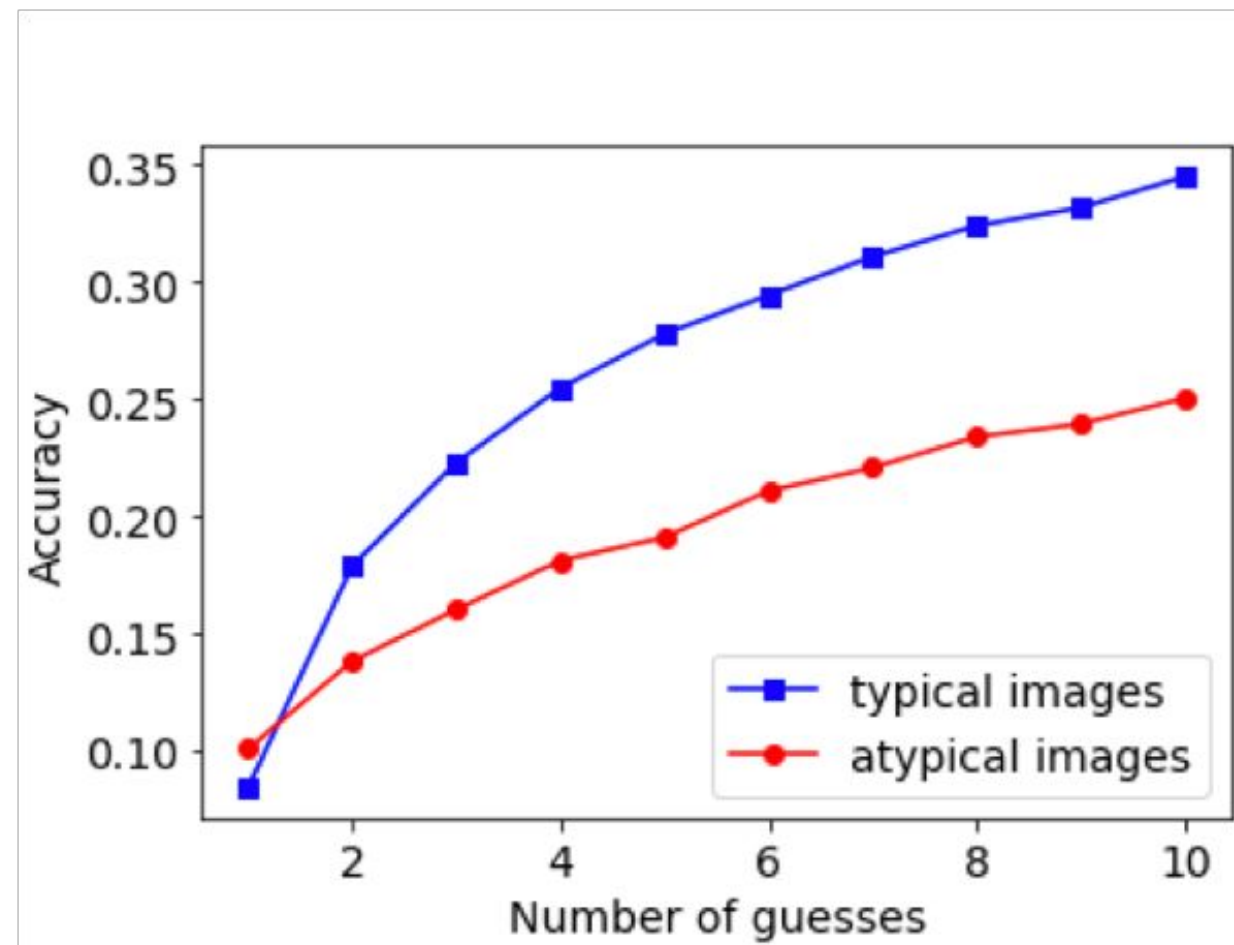


Visually Similar

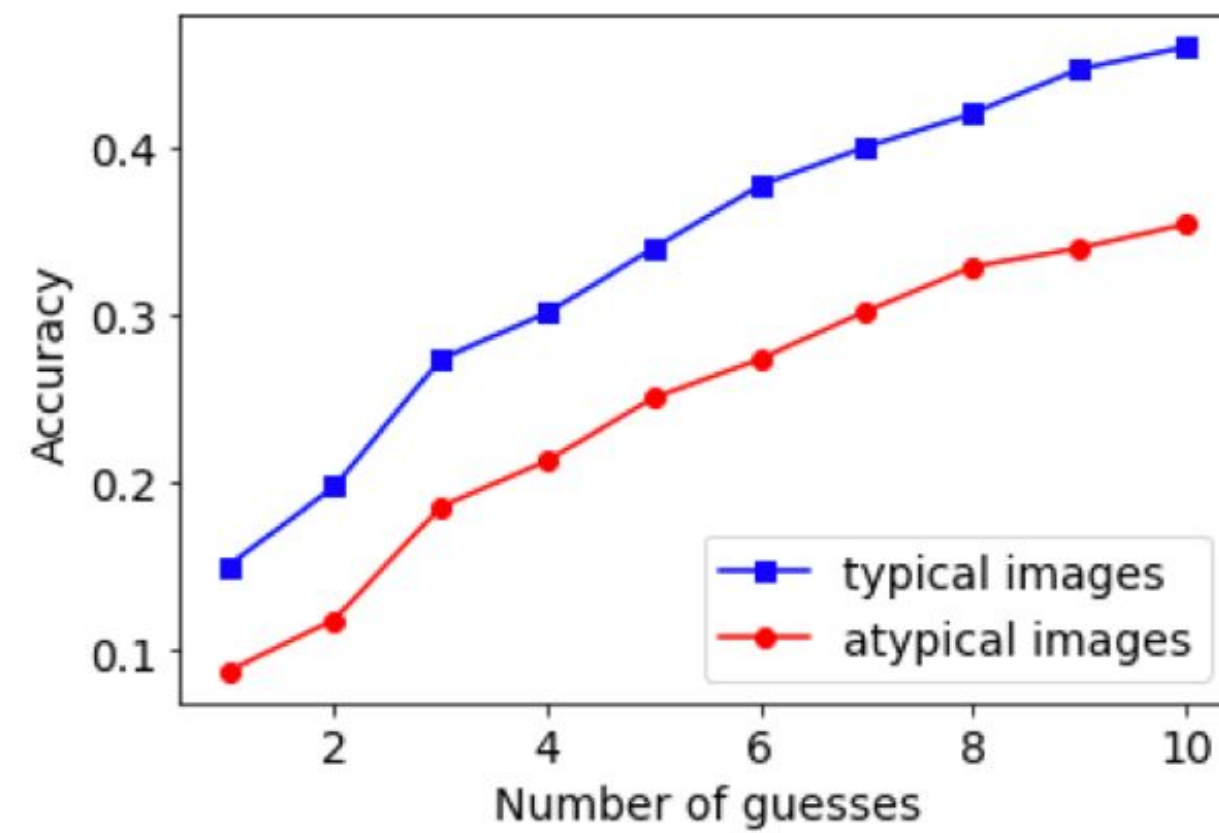


Architecture of “PERSPECTIVE”

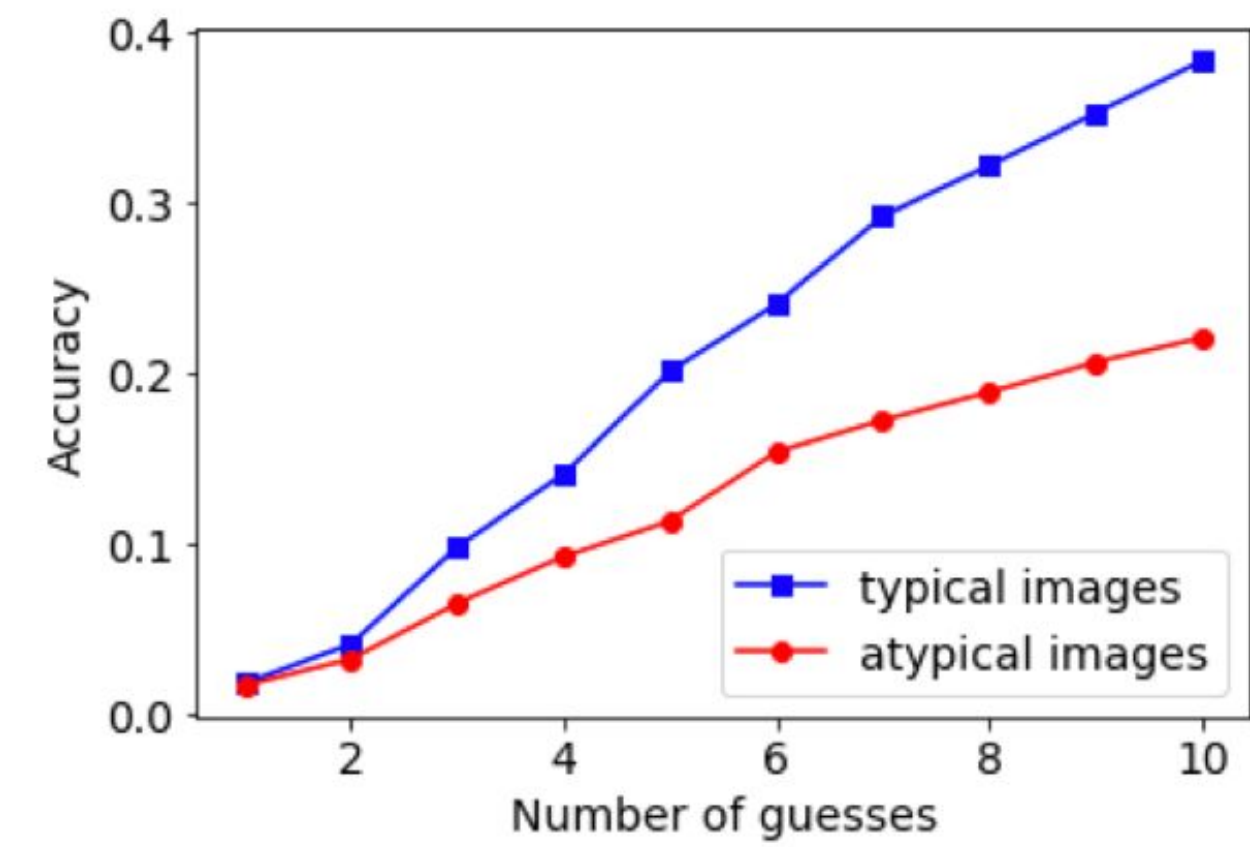




(a) Google Vision API

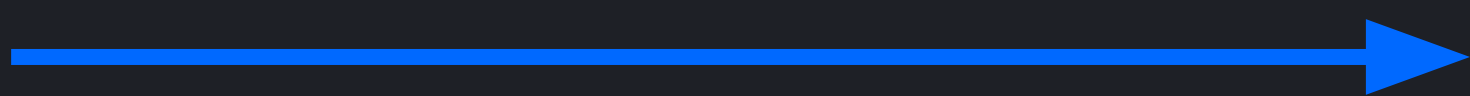


(b) Amazon Rekognition API



(c) Microsoft Azure Vision API

Check out the paper
for MORE interesting
details and analyses



edu.nl/cy7cj

Key takeaways



Annotations with **Perspective** led to identification of most atypical images



Importance of context expansion during labelling workflows

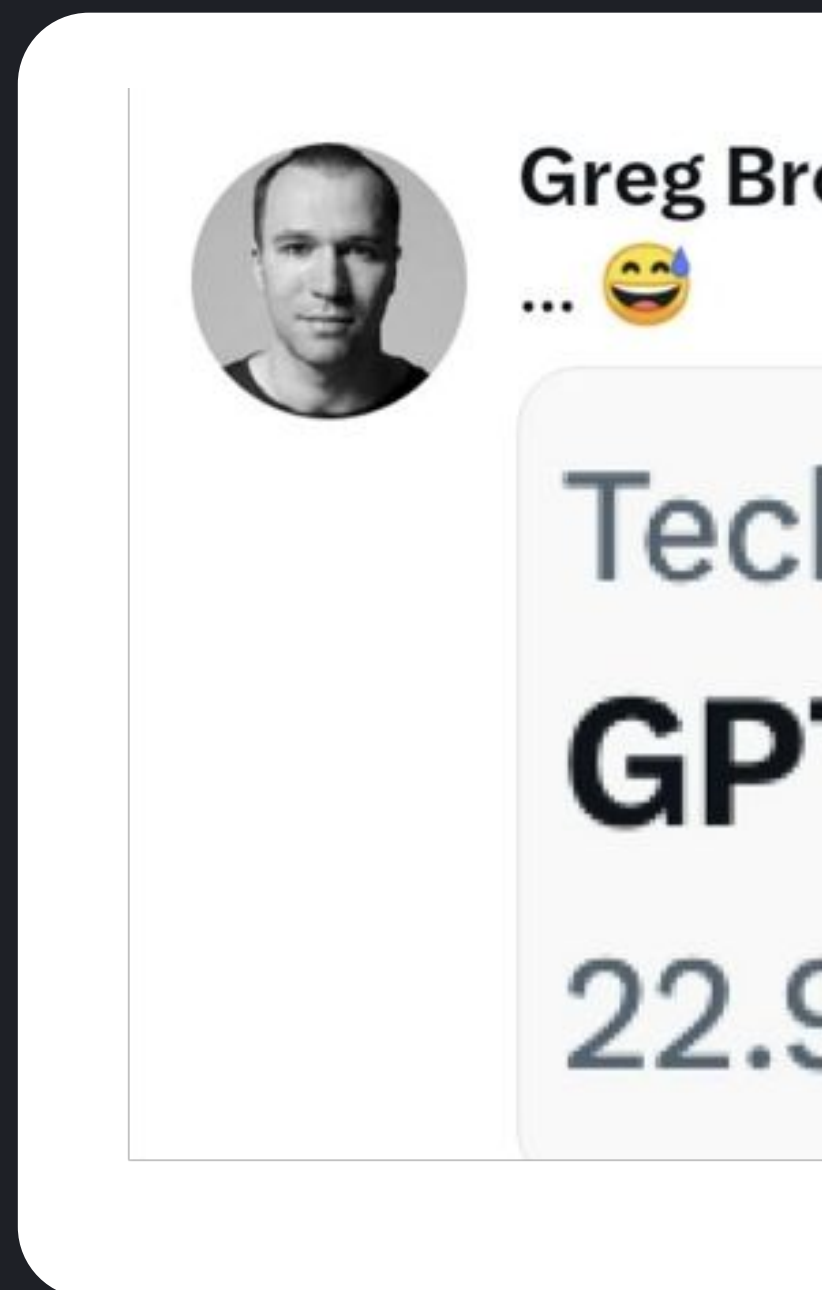


Response and data sampling biases need to be addressed



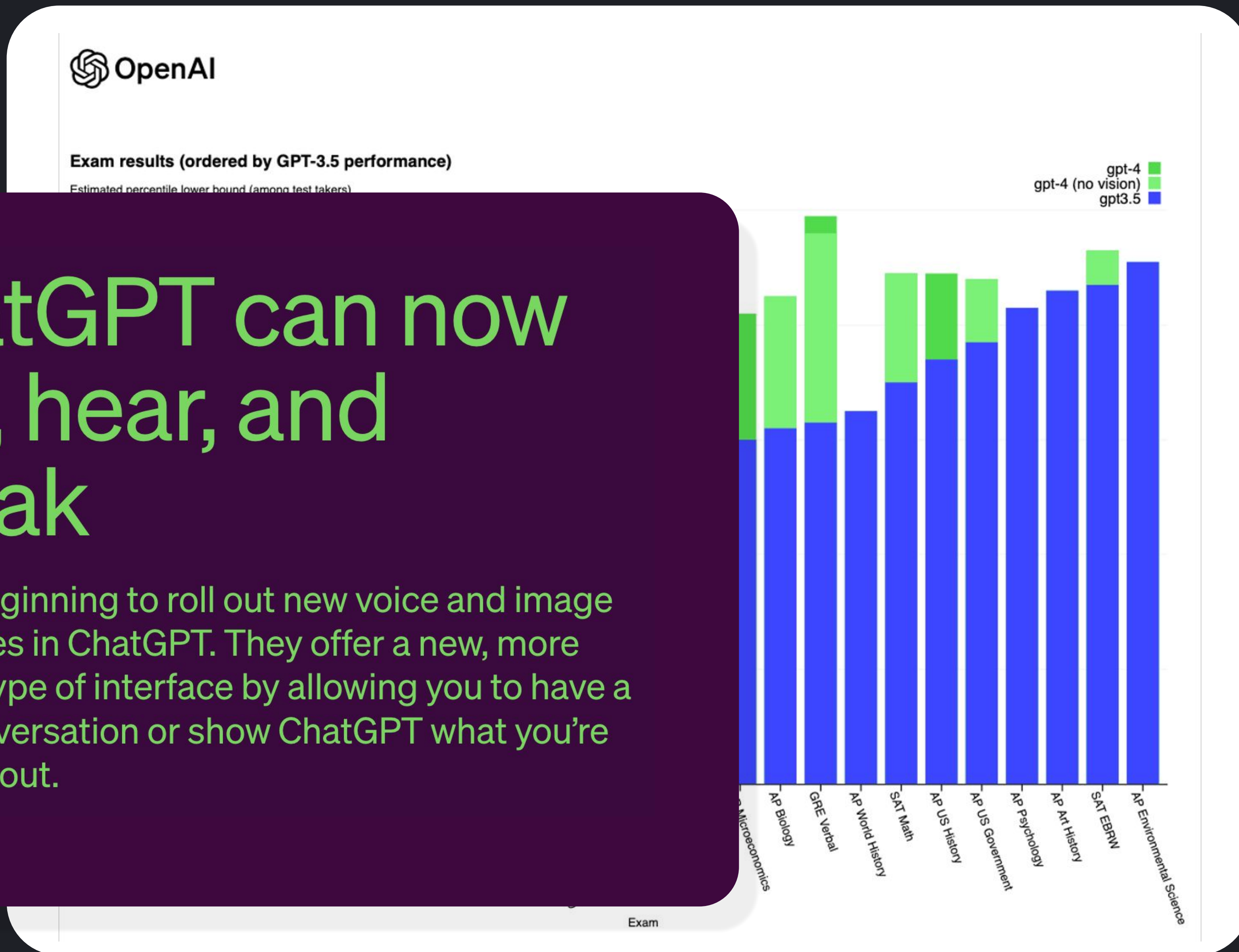
No place for mistakes when the stakes are high — responsible and trustworthy AI

The Analogous Challenges with LLMs ...



ChatGPT can now see, hear, and speak

We are beginning to roll out new voice and image capabilities in ChatGPT. They offer a new, more intuitive type of interface by allowing you to have a voice conversation or show ChatGPT what you're talking about.



Biases in the Age of LLMs

More human...

Fabio M
The Rev

Received:
© The Autl

Abstrac
We inve
become
GPT ass
race, ge
political
Moreove
propose
ing it to

ABEL
PART
YUZI
ROBI
FRED

Warni
Large
standi
to mal
metho
demon
two cu
howev
identit
els con
to won
unders

Queer People are People First: Deconstructing Sexual Identity Stereotypes in Large Language Models

Harnoor Dhingra Preetiha Jayashanker Sayali Moghe Emma Strubell
Carnegie Mellon University

{hdhingra, pjayasha, smoghe, estrubell}@cs.cmu.edu

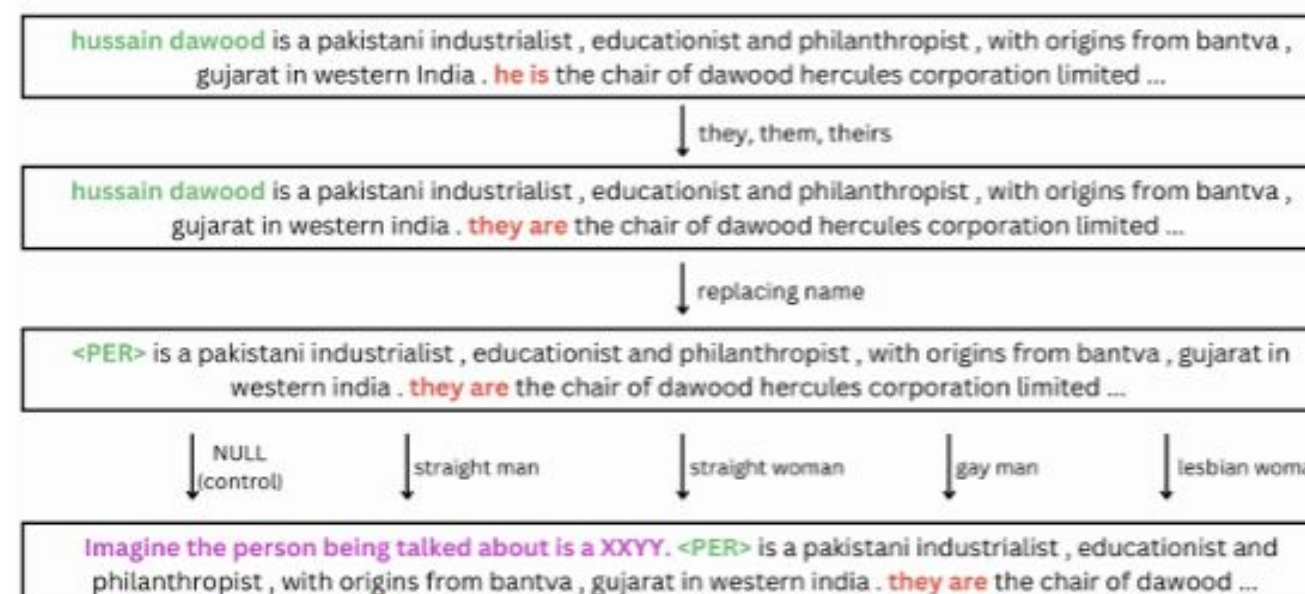
Abstract

Large Language Models (LLMs) are trained primarily on minimally processed web text, which exhibits the same wide range of social biases held by the humans who created that content. Consequently, text generated by LLMs can inadvertently perpetuate stereotypes towards marginalized groups, like the LGBTQIA+ community. In this paper, we perform a comparative study of how LLMs generate text describing people with different sexual identities. Analyzing bias in the text generated by an LLM using regard score shows measurable bias against queer people. We then show that a post-hoc method based on chain-of-thought prompting using SHAP analysis can increase the regard of the sentence, representing a promising approach towards debiasing the output of LLMs

quantify the detected biases. Hence, in this work we aim to answer the following research questions:

RQ1: Does a pre-trained LLM perpetuate *measurable, quantifiable* bias against queer people?

RQ2: Can we *mitigate* the said bias in the LLM output *while preserving the context* using a post-hoc debiasing method?



of Profession ssing

er are twofold. Firstly, it of the profession in GPT-generated text for patterns presence of stereotypical hough rigorous quantitative a comprehensive under- t in GPT-2 and GPT-3.5

Biases in the Age of LLMs

- Instruction-tuned LLMs have been shown to be effective in generating high-quality natural language responses
- Open RQ → inherent biases in trained models and the generated responses
 - E.g., dataset for fine-tuning is predominantly composed of a specific political bias, we can expect the generated answers to share such bias

Instructed to Bias: Instruction-Tuned Language Models Exhibit Emergent Cognitive Bias

Itay Itzhak¹, Gabriel Stanovsky², Nir Rosenfeld¹, Yonatan Belinkov¹

¹Technion – Israel Institute of Technology

²School of Computer Science and Engineering, The Hebrew University of Jerusalem
itay1itzhak@gmail.com,

{nirr, belinkov}@technion.ac.il, gabriel.stanovsky@mail.huji.ac.il

Abstract

Recent studies show that instruction tuning and learning from human feedback improve the abilities of large language models (LMs) dramatically. While these tuning methods can help models generate high-quality text, we conjecture that they may also inadvertently cause models to express cognitive-like biases. Our work provides evidence that fine-tuned models exhibit biases that were absent or less pronounced in their pretrained predecessors. We examine the extent of this phenomenon in three cognitive biases: the decoy effect, the certainty effect, and the belief bias—all of which are known to influence human decision-making and reasoning. Our findings highlight the presence of these biases in various models, especially those that have undergone instruction tun-

In both examples Option A has a **higher** expected utility

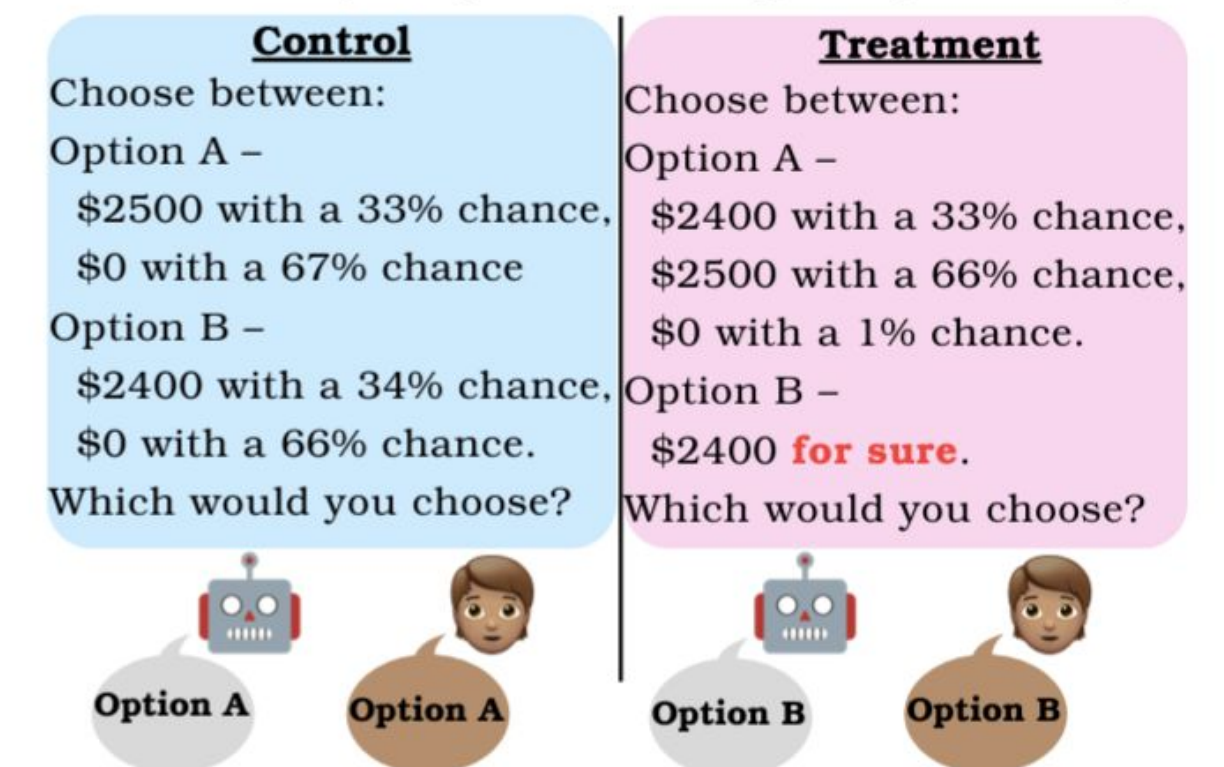


Figure 1: Example tasks from *certainty effect* dataset, for the control condition (left) and treatment condition (right), along with typical answers from humans and instruction-tuned models, both of which are biased.

The indispensable role of human input

How can we elicit human
input (tacit knowledge)
to manage biases in LLMS

Eliciting Diverse Knowledge from Humans Using A Game-with-a-purpose

- ✓ Commonsense knowledge → building neuro-symbolic AI systems, debugging deep learning models
- ✓ Existing knowledge acquisition methods are limited
- ✓ Broad tacit and *negative* knowledge, and *discriminative* knowledge → Our solutionE a GWAP, FindItOut



Eliciting Diverse Knowledge Using A Configurable Game

Find **It** out

The screenshot displays the 'Find It Out' game interface. On the left, a player's card is shown with a '3' in a yellow circle. The card features a 'Mink' image and a 'Collected knowledge' section with the following entries: '< Otter, IsA, carnivore> (+)', '< Hare, IsA, carnivore> (-)', and 'Otter, Hare, IsA, carnivore> (+)'. Below the card, it says 'You are the REPLIER' and shows five cards: Raccoon, Otter, Mole, Skunk, and Hare. On the right, a question panel is shown with a '1' in a yellow circle. It contains a dropdown menu with options: 'IsA', 'HasA', 'HasProperty', 'UsedFor', 'CapableOf', and 'MadeOf'. The selected option is 'IsA'. The question is 'Is your card a carnivore?' with a 'SEND' button. Below this, there are 'YES', 'NO', 'MAYBE', and 'SEND' buttons, and an 'UNCLEAR' button. A '2' in a yellow circle is next to the bottom buttons.

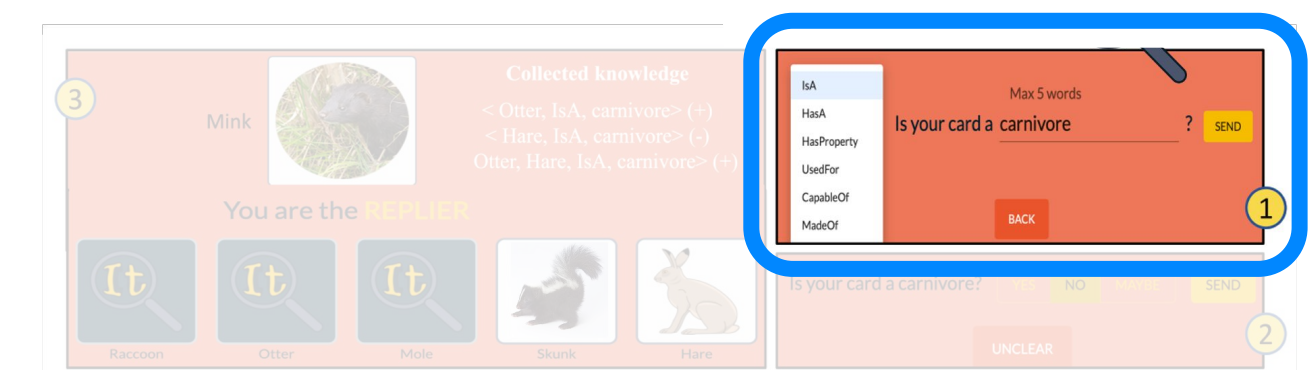
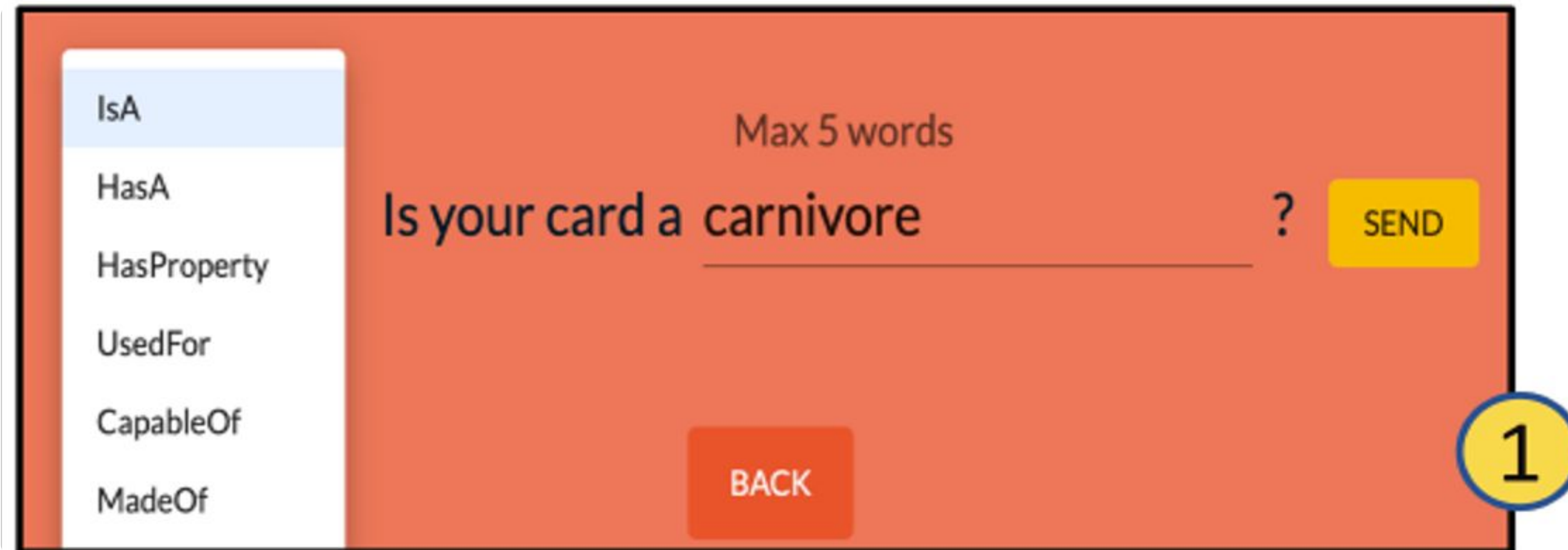


Play → <https://finditout.vercel.app/>

Best Demo & Poster Award at AAI HCOMP 2021

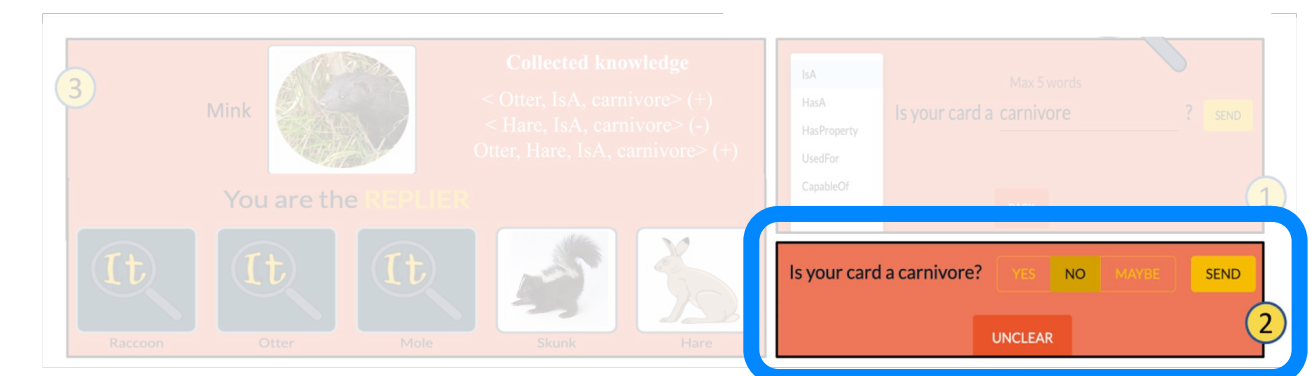
Best Paper Award Nomination at the ACM Web Conference 2022

Eliciting Diverse Knowledge Using A Configurable Game



Play → <https://finditout.vercel.app/>

Eliciting Diverse Knowledge Using A Configurable Game



Play → <https://finditout.vercel.app/>

Eliciting Diverse Knowledge Using A Configurable Game

3

Mink



Collected knowledge

- < Otter, IsA, carnivore > (+)
- < Hare, IsA, carnivore > (-)
- Otter, Hare, IsA, carnivore > (+)

You are the **REPLIER**

It It It

Raccoon Otter Mole Skunk Hare



3

Mink

Collected knowledge

- < Otter, IsA, carnivore > (+)
- < Hare, IsA, carnivore > (-)
- Otter, Hare, IsA, carnivore > (+)

You are the **REPLIER**

It It It

Raccoon Otter Mole Skunk Hare

Is your card a carnivore?

UNCLEAR

Play → <https://finditout.vercel.app/>

Collecting “tacit knowledge” for downstream AI tasks



Board	Type	Question	Knowledge Tuple
floor, window, bathroom, walls, ceiling, chandelier, mirror, bedroom	Explicit Tacit	Can your card be found inside an apartment? Can your card be used for decoration?	<bathroom, AtLocation, inside apartment> <chandelier, UsedFor, decoration>
necklace, dress, boots, shoes, pants, trousers, jeans, skirt	Explicit Tacit	Can your card be found in your wardrobe? Is your card typically worn by cowboys?	<dress, AtLocation, wardrobe> <boots, HasProperty, worn by cowboys>

Collecting “tacit knowledge” for downstream AI tasks

Empirical Results:

- ✓ 125 players played 2430 rounds → 150k knowledge tuples
- ✓ Efficiency of game is 10x higher than a reference baseline

Verbosity

- ✓ Usefulness validated in two downstream AI tasks
 - Commonsense Question-Answering
 - Identification of Discriminative Attributes
- ✓ Enjoyable game experience (player experience inventory)

Re-Enter → The Analogous Challenges with LLMs

Use annotation tools and workflows like “Perspective” to identify and characterize biases →

Explore how such biases are represented in LLMs

Elicit diverse human input to create bias-aware instruction tuning datasets →

Mitigate and manage biases in fine-tuned LLMs



Human input and oversight are essential to overcome fundamental challenges in facilitating bias-aware interactions with LLMs.

Ujwal Gadiraju

Dr. ir.

✉ ujwalg@toloka.ai



www.toloka.ai

LinkedIn

[LinkedIn](#)



[Twitter](#)

GitHub

[GitHub](#)

slack

[Slack](#)

Join our
Slack community