



UiT The Arctic University of Norway

NORA summer school on multi-modal learning

Responsible AI

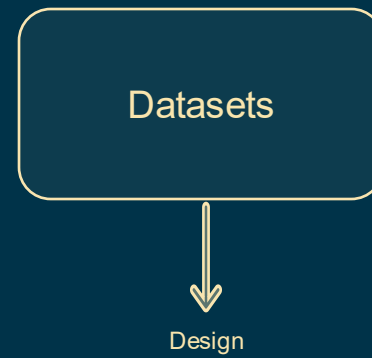
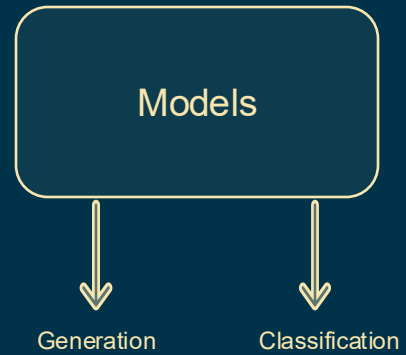
Rwiddhi Chakraborty

UiT Machine Learning Group and Visual Intelligence

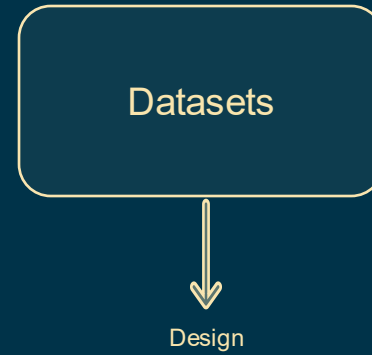
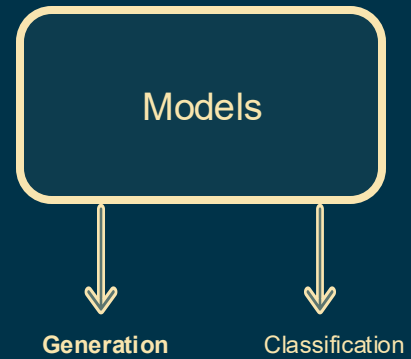
Schedule Today

- 09 - 10: Special Topics - I
- 10 - 11: Special Topics - II
- 11 – 12: Group Project
- 12 – 13: Lunch
- 13 – 14: Group Project
- 14 – 15: Presentations, award, exam info, wrap-up!

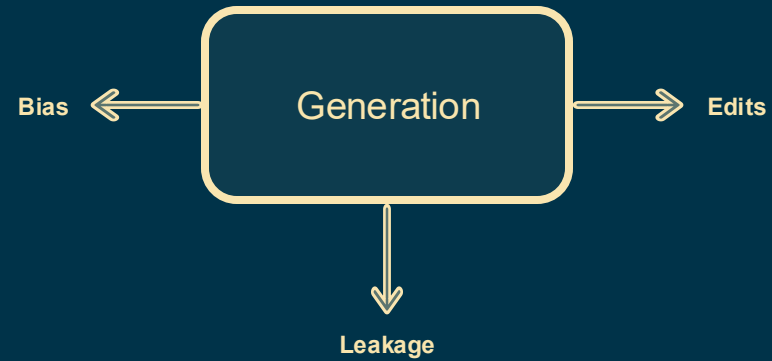
In this talk



In this talk



Responsible Generation



Data Leakage and Memorization

Data privacy refers to protected training data

Prevents forgery, theft, and other violations

Models “leaking” training data is a safety risk

Data Leakage and Memorization

Leaking training data is a form of overfitting and memorization

Most modern generative models suffer from this issue

Unclear whether benchmark performances correlate with true understanding

Data Leakage and Memorization

How do we analyze memorization in diffusion models?

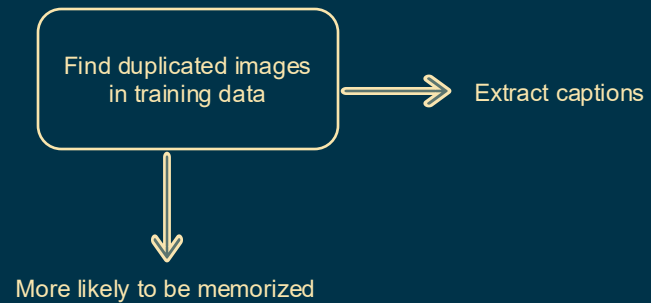
Data Leakage and Memorization

How do we analyze memorization in diffusion models?



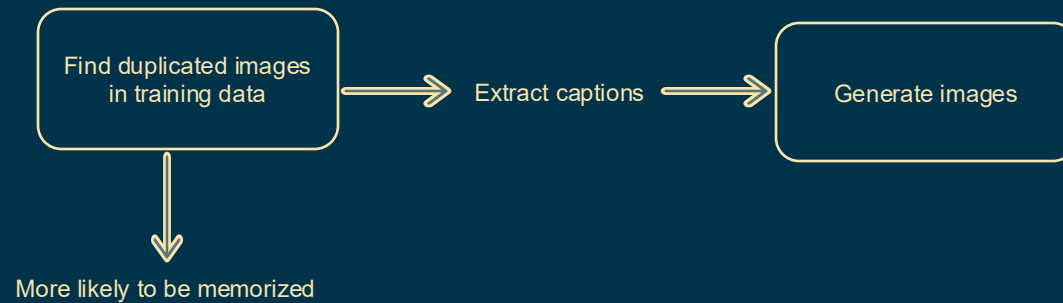
Data Leakage and Memorization

How do we analyze memorization in diffusion models?



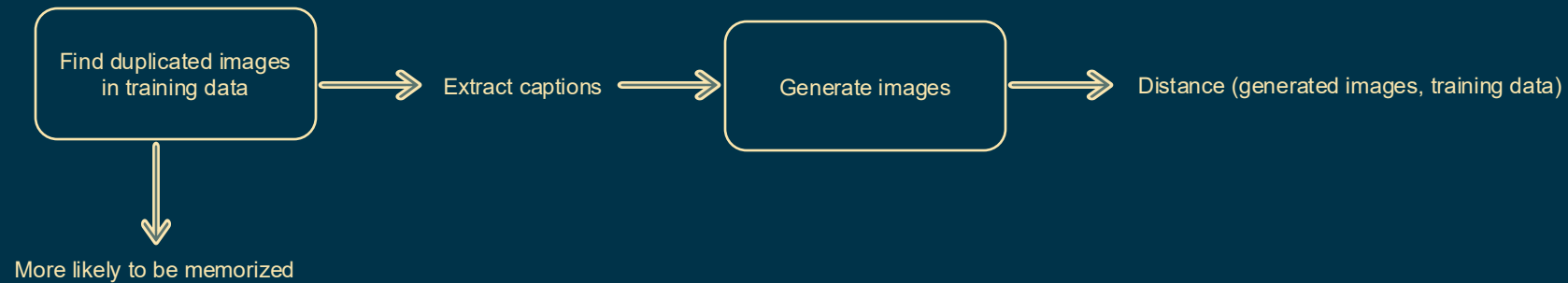
Data Leakage and Memorization

How do we analyze memorization in diffusion models?



Data Leakage and Memorization

How do we analyze memorization in diffusion models?



Data Leakage and Memorization

The result



Data Leakage and Memorization

Vast majority of generated images are photographs of real people

Not all images are permissively licensed and raise copyright issues as well

Risks greater for diffusion models trained on more sensitive data (e.g medicine)

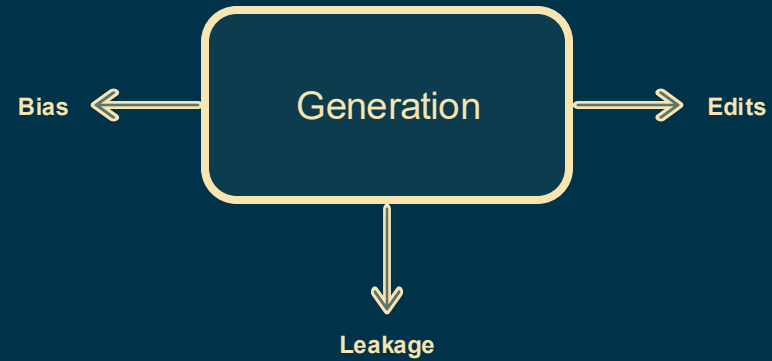
Data Leakage and Memorization

Model specific issue

GANs are safer as they are not trained to directly mimic the training data

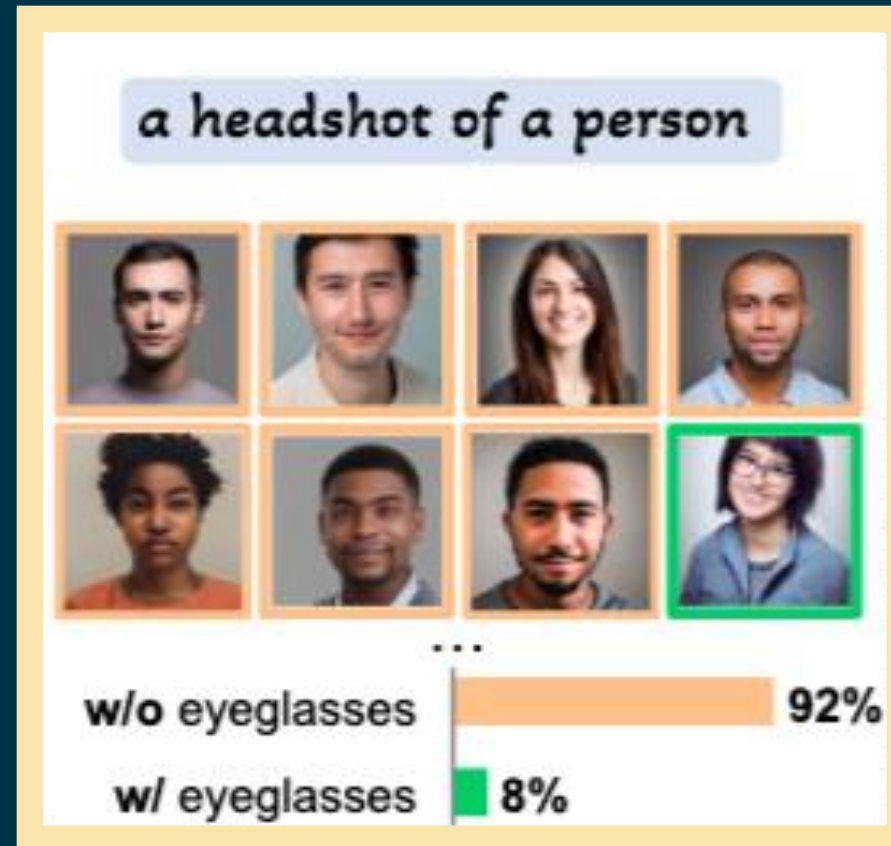
Risks greater for diffusion models trained on more sensitive data (e.g medicine)

Responsible Generation

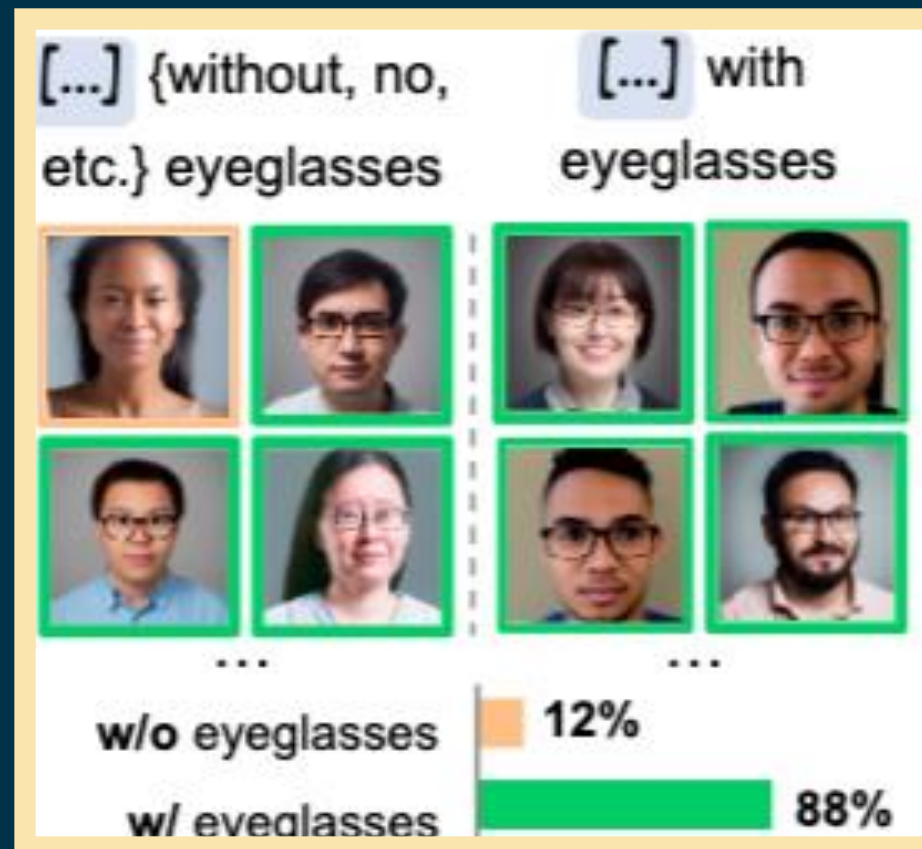
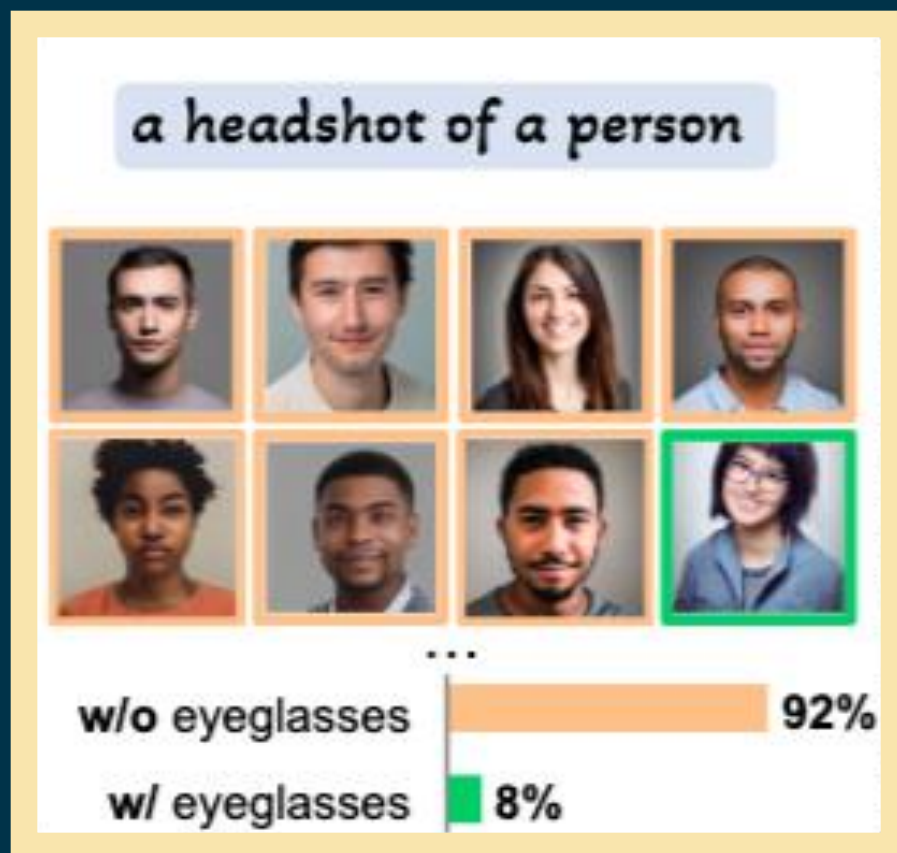


Biased Generation

Any attribute
can be a
minority in the
training set



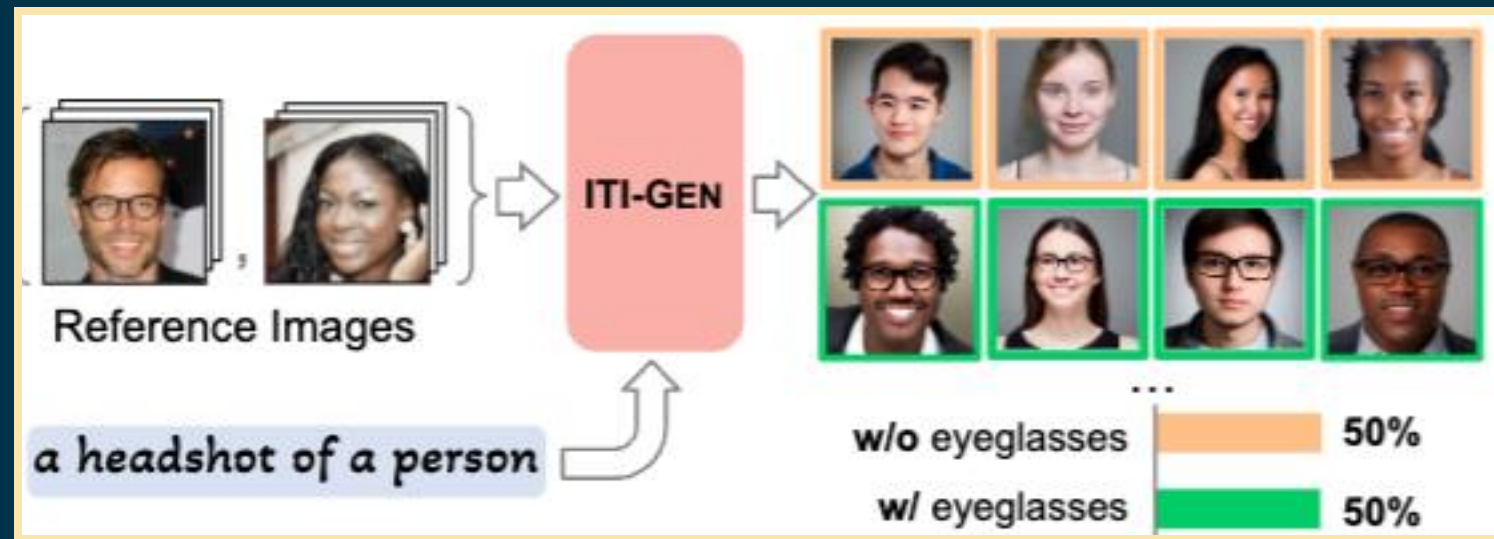
Biased Generation



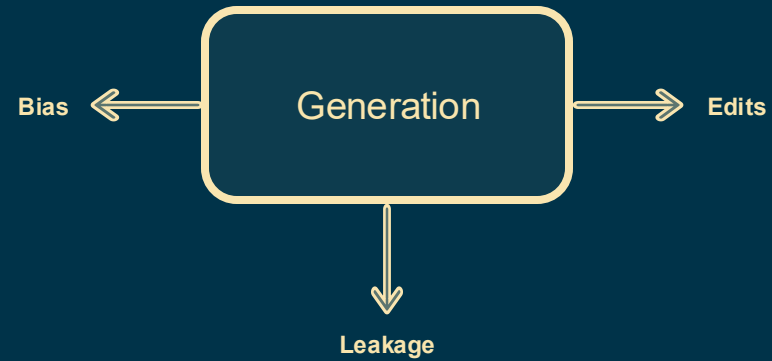
Hard prompting is ambiguous

Responsible Generation

Image-guided prompt tuning



Responsible Generation



Editing harmful generation

Remove copyrighted/memorized content from T2I models

Prevent model from generating harmful concepts

Editing harmful generation



Editing harmful generation



Editing harmful generation

A photo of a **cat**



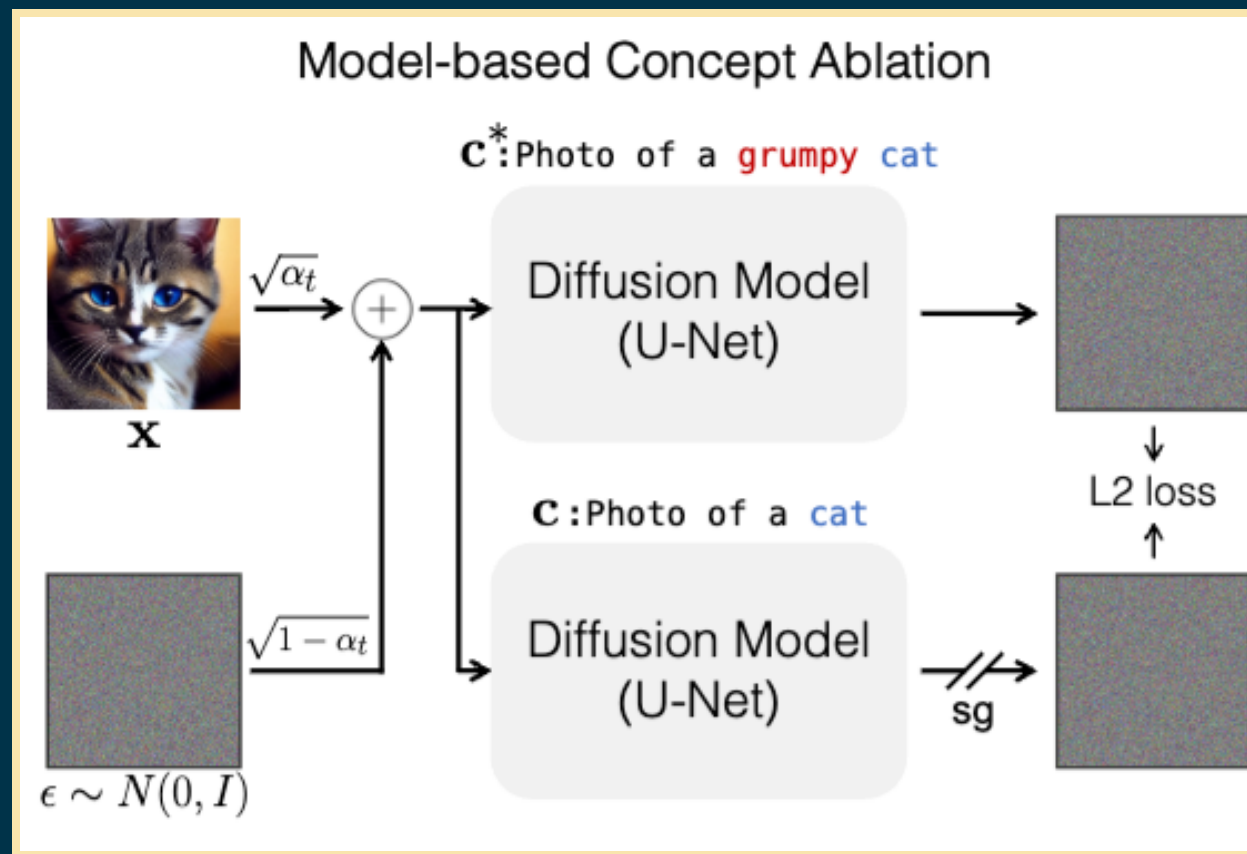
Anchor concept **c**

A photo of a **grumpy** cat

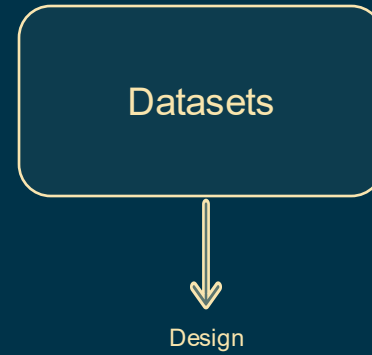
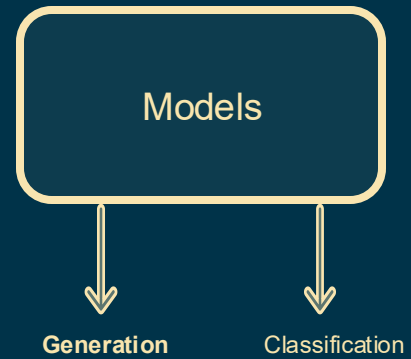


Target concept **c***

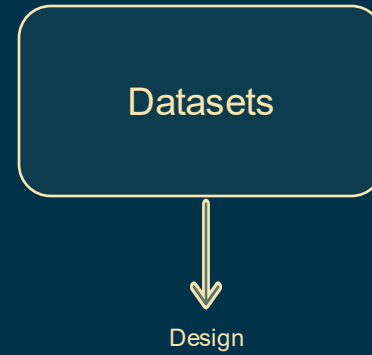
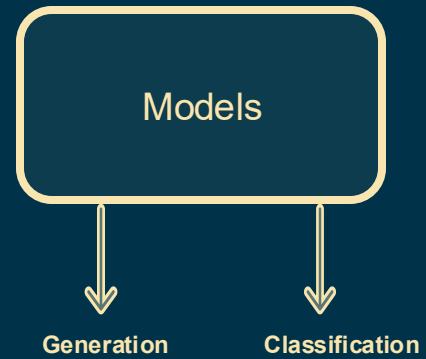
Editing harmful generation



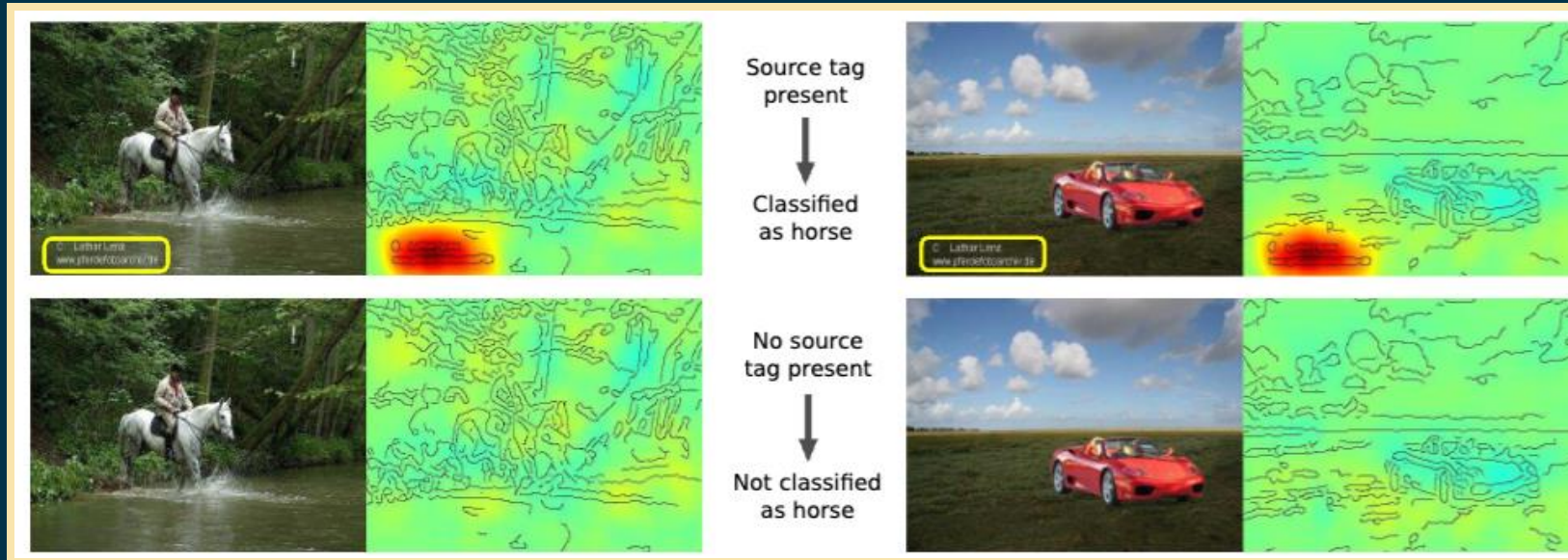
In this talk



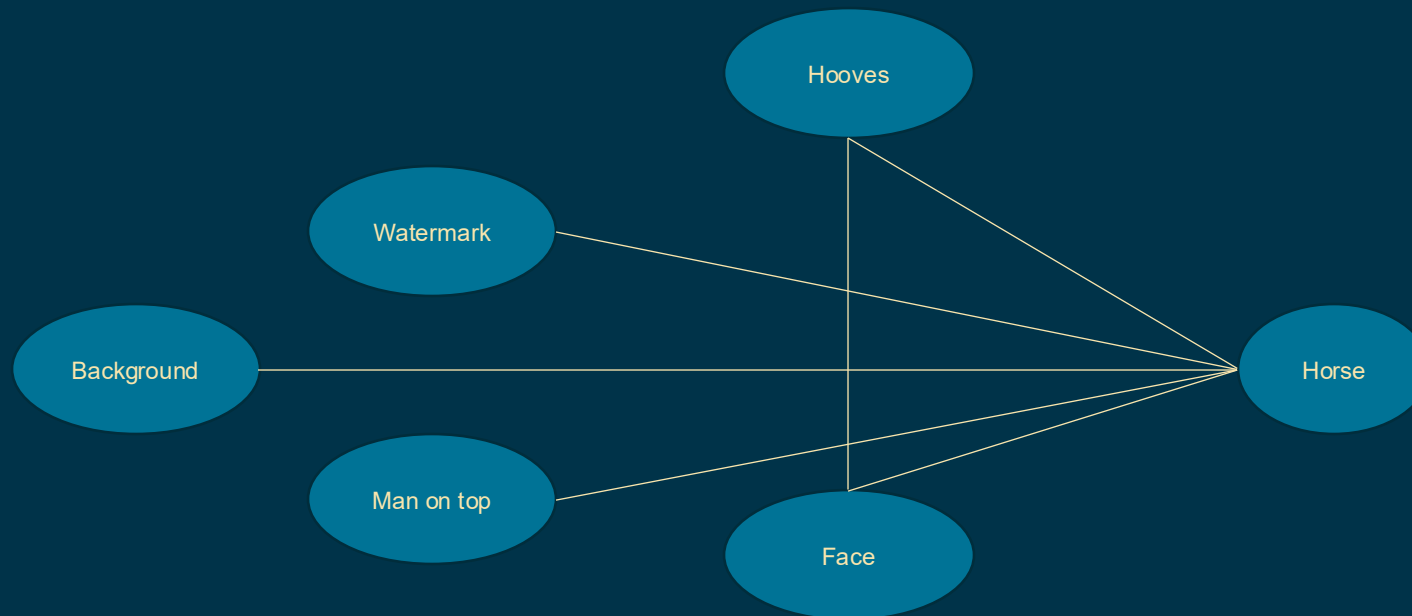
In this talk



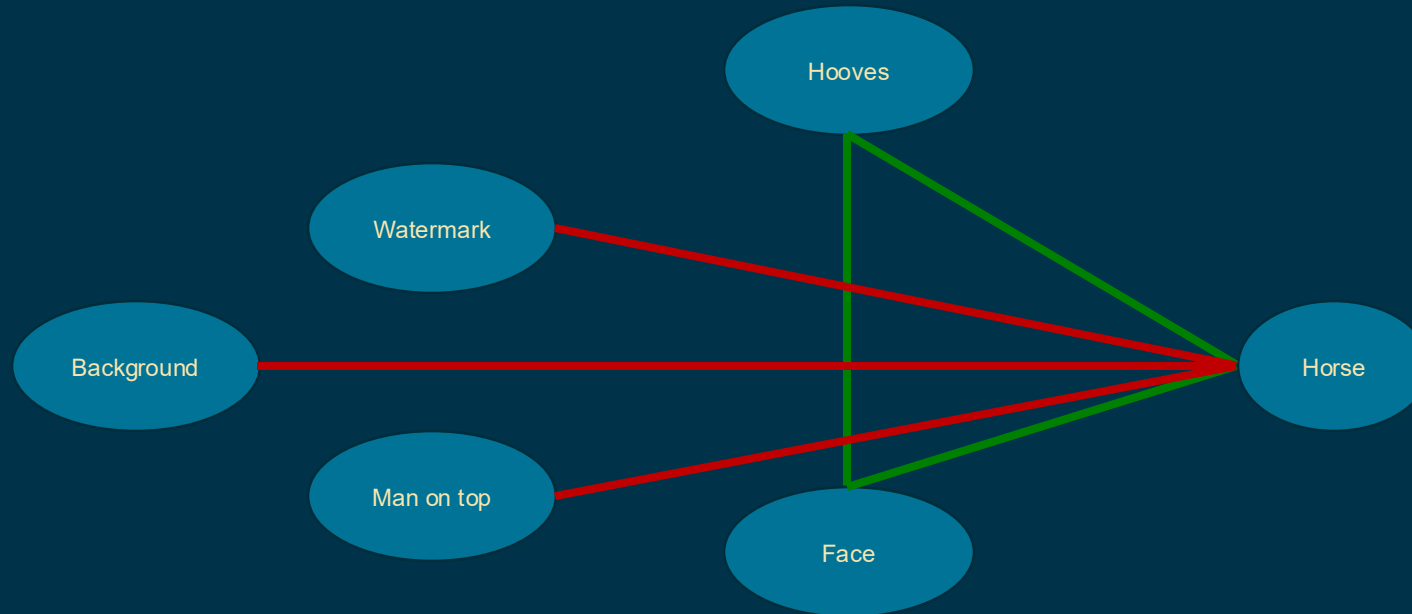
Shortcut Learning



Shortcut Learning



Shortcut Learning



Shortcut Learning in VLMs

There are two modalities in inference – text and vision

Do models leverage both modalities?

Shortcut Learning in VLMs





There are two modalities in inference – text and vision

Do models leverage both modalities?

(Un)Surprisingly, no!

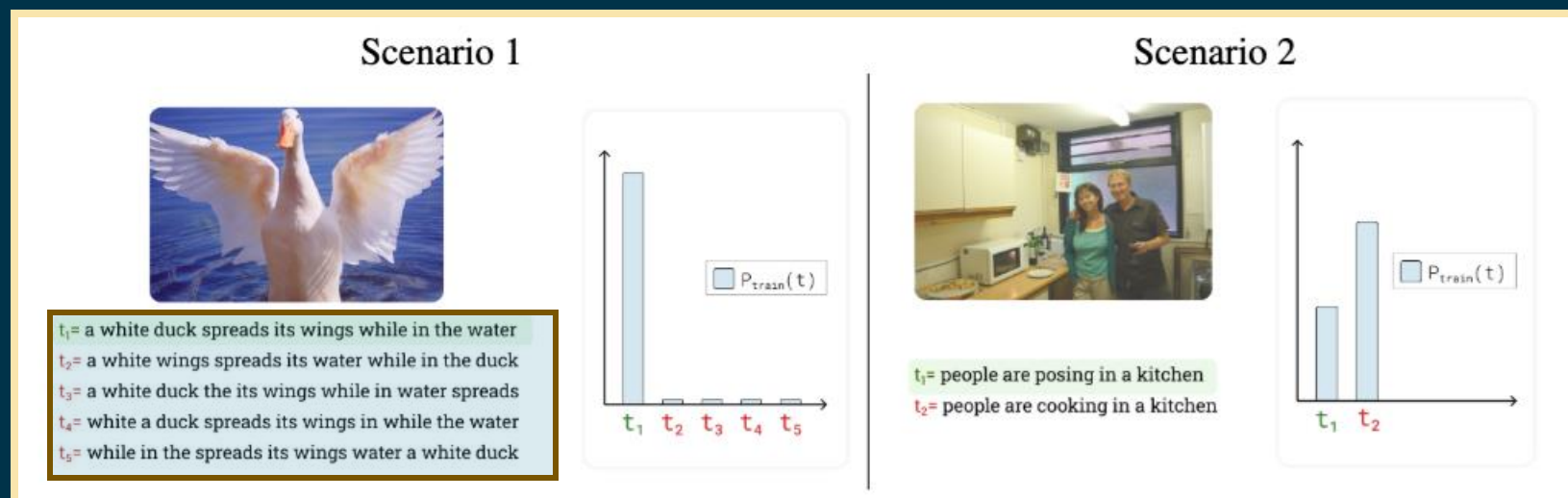
Shortcut Learning in VLMs

The language prior

	Train	Test
Example 1	<p>Q+[A] What color is the dog ? [White]</p> <p>Image </p> <p>Training Prior</p> <ul style="list-style-type: none">whiteredbluegreenyellow...	<p>Q+[A] What color is the dog ? [Black]</p> <p>Image </p> <p>Models</p> <p>SAN GVQA</p> <p>White Black</p>
Example 2	<p>Q+[A] Is the person wearing shorts ? [No]</p> <p>Image </p> <p>Training Prior</p> <ul style="list-style-type: none">nofemalewoman...	<p>Q+[A] Is the person wearing shorts ? [Yes]</p> <p>Image </p> <p>Models</p> <p>SAN GVQA</p> <p>No Yes</p>

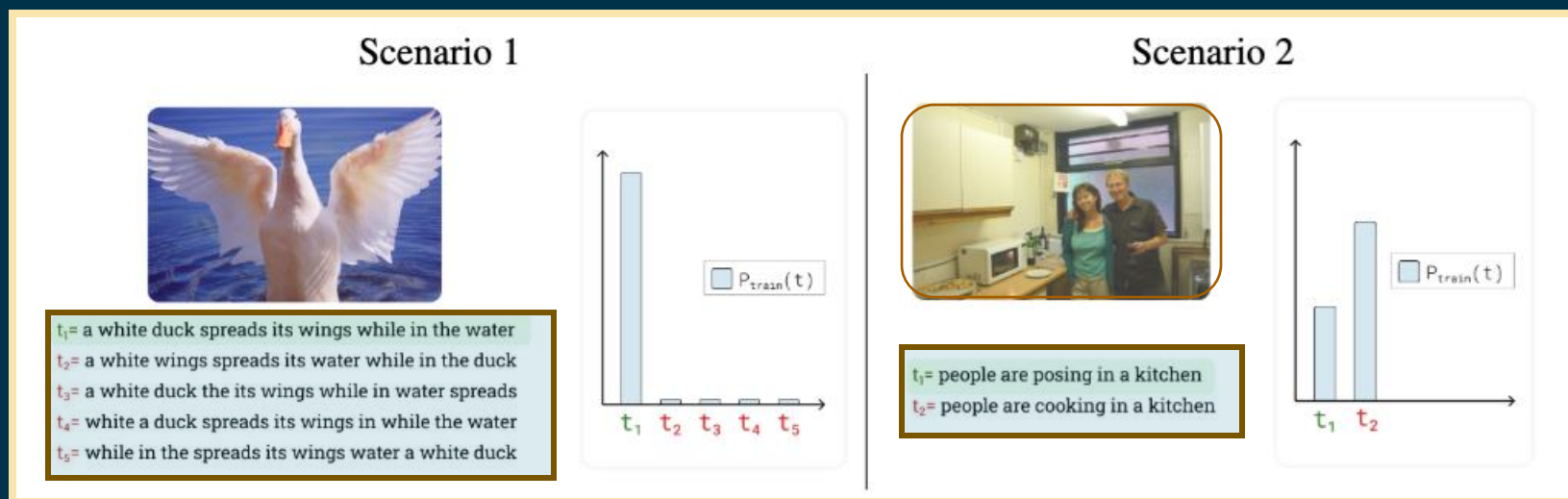
Shortcut Learning in VLMs

The language prior



Shortcut Learning in VLMs

The language prior



Shortcut Learning in VLMs

Issue stems from a misalignment between the train and test distributions

Shortcut Learning in VLMs

Issue stems from a misalignment between the train and test distributions

The model assigns a score

$$\frac{P_{train}(t|i)}{P_{train}(t)^\alpha}$$

Shortcut Learning in VLMs

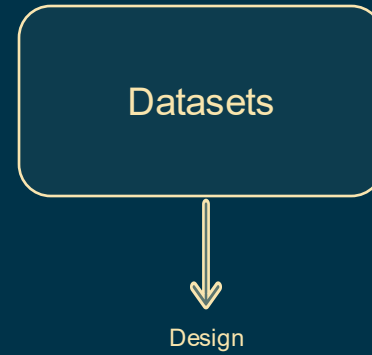
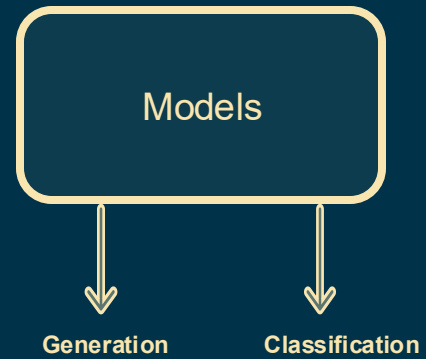
Issue stems from a misalignment between the train and test distributions

The model assigns a score

$$\frac{P_{train}(t|i)}{P_{train}(t)^\alpha}$$

Tuning the alpha controls the assumptions on how the train and test are related

In this talk



In this talk



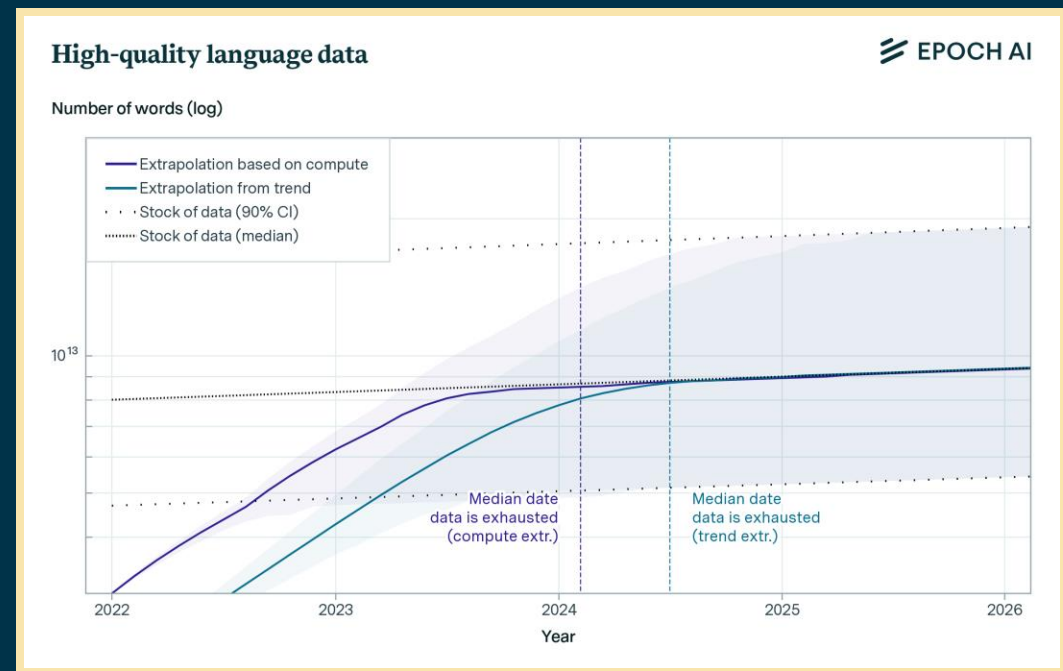
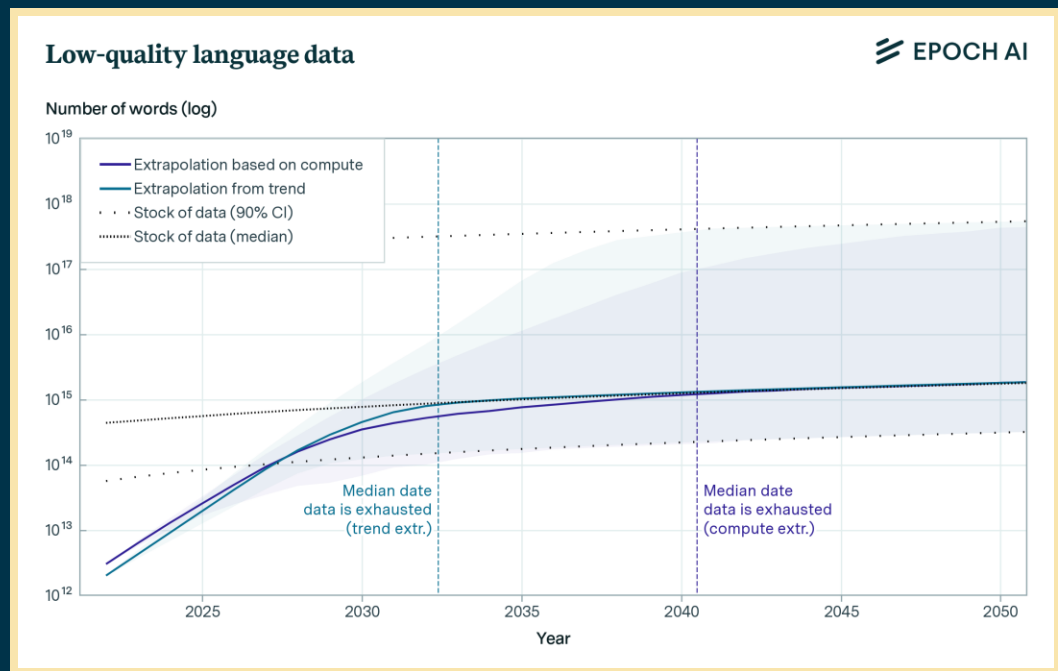
Data

Your model cannot be better than the data it trains on

Data

More data \nRightarrow Better data

In fact



In fact

“Our projections predict that we will have exhausted the stock of low-quality language data by 2030 to 2050, high-quality language data before 2026, and vision data by 2030 to 2060.”

A closer look

Stable Diffusion for instance, is trained on LAION-5B

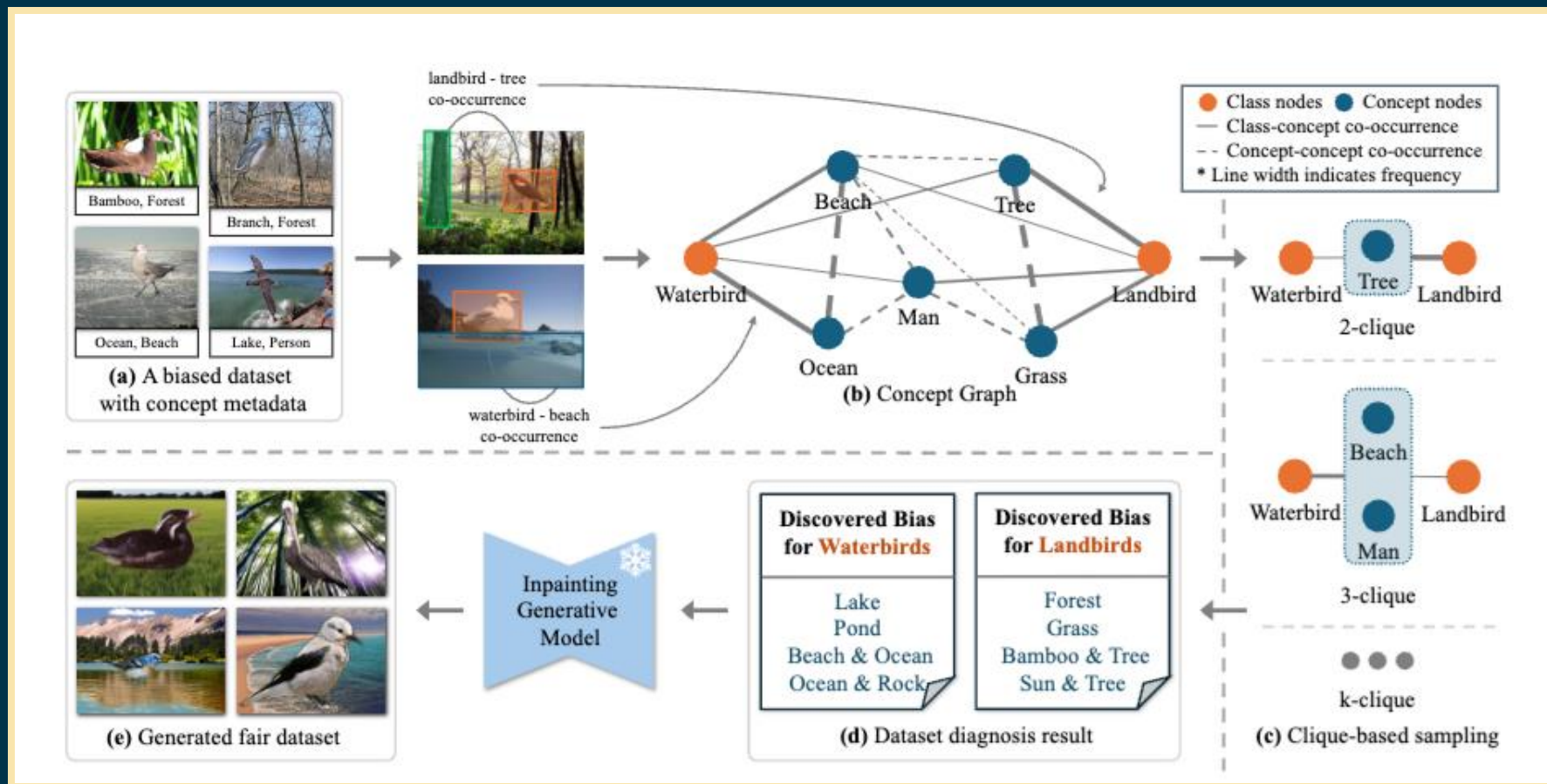
Bias and Fairness

Calibrating outcomes for a marginalized distribution in a dataset

Uniform distribution across all attributes isn't always "fair"

Bias isn't just *social*, it is simply a prior belief on the data

Mitigating Dataset Bias with Augmentation



Summary



Summary

