

# Memorizing Distributions: Generative Models Recall More Than Verbatim Instances

Yanai Elazar, L2M2, August 1st, 2025



# Memorization

## Extracting Training Data from Large Language Models

Nicholas Carlini<sup>1</sup>

Florian Tramèr<sup>2</sup>

Eric Wallace<sup>3</sup>

Matthew Jagielski<sup>4</sup>

Ariel Herbert-Voss<sup>5,6</sup>

Katherine Lee<sup>1</sup>

Adam Roberts<sup>1</sup>

Tom Brown<sup>5</sup>

Dawn Song<sup>3</sup>

Úlfar Erlingsson<sup>7</sup>

Alina Oprea<sup>4</sup>

Colin Raffel<sup>1</sup>

<sup>1</sup>Google <sup>2</sup>Stanford <sup>3</sup>UC Berkeley <sup>4</sup>Northeastern University <sup>5</sup>OpenAI <sup>6</sup>Harvard <sup>7</sup>Apple

## Counterfactual Memorization in Neural Language Models

Chiyuan Zhang

Google Research

chiyuan@google.com

Daphne Ippolito

Carnegie Mellon University

daphnei@cmu.edu

Katherine Lee

Google DeepMind

katherinelee@google.com

Matthew Jagielski

Google DeepMind

jagielski@google.com

Florian Tramèr

ETH Zürich

florian.tramer@inf.ethz.ch

Nicholas Carlini

Google DeepMind

ncarlini@google.com

## QUANTIFYING MEMORIZATION ACROSS NEURAL LANGUAGE MODELS

Nicholas Carlini<sup>\*</sup>

Katherine Lee<sup>1,3</sup>

Daphne Ippolito<sup>1,2</sup>

Florian Tramèr<sup>1</sup>

Matthew Jagielski<sup>1</sup>

Chiyuan Zhang<sup>1</sup>

<sup>1</sup>Google Research

<sup>2</sup>University of Pennsylvania

<sup>3</sup>Cornell University

## Extracting Training Data from Diffusion Models

Nicholas Carlini<sup>\*1</sup>

Jamie Hayes<sup>\*2</sup>

Milad Nasr<sup>\*1</sup>

Matthew Jagielski<sup>+1</sup>

Vikash Sehwag<sup>+4</sup>

Florian Tramèr<sup>+3</sup>

Borja Balle<sup>†2</sup>

Daphne Ippolito<sup>†1</sup>

Eric Wallace<sup>†5</sup>

<sup>1</sup>Google

<sup>2</sup>DeepMind

<sup>3</sup>ETHZ

<sup>4</sup>Princeton

<sup>5</sup>UC Berkeley

\*Equal contribution

+Equal contribution

†Equal contribution

# Memorization

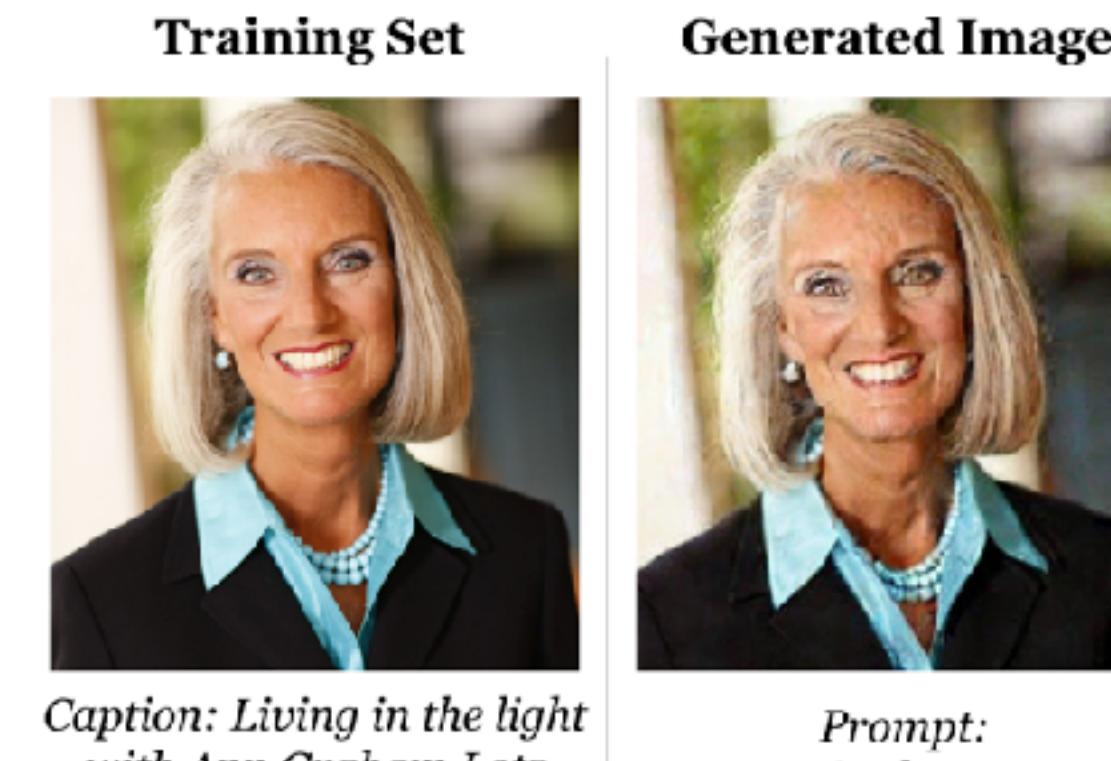
- Memorization has been heavily studied
- Focusing on verbatim memorization

**Definition 3.1.** A string  $s$  is *extractable with  $k$  tokens of context* from a model  $f$  if there exists a (length- $k$ ) string  $p$ , such that the concatenation  $[p \parallel s]$  is contained in the training data for  $f$ , and  $f$  produces  $s$  when prompted with  $p$  using greedy decoding.

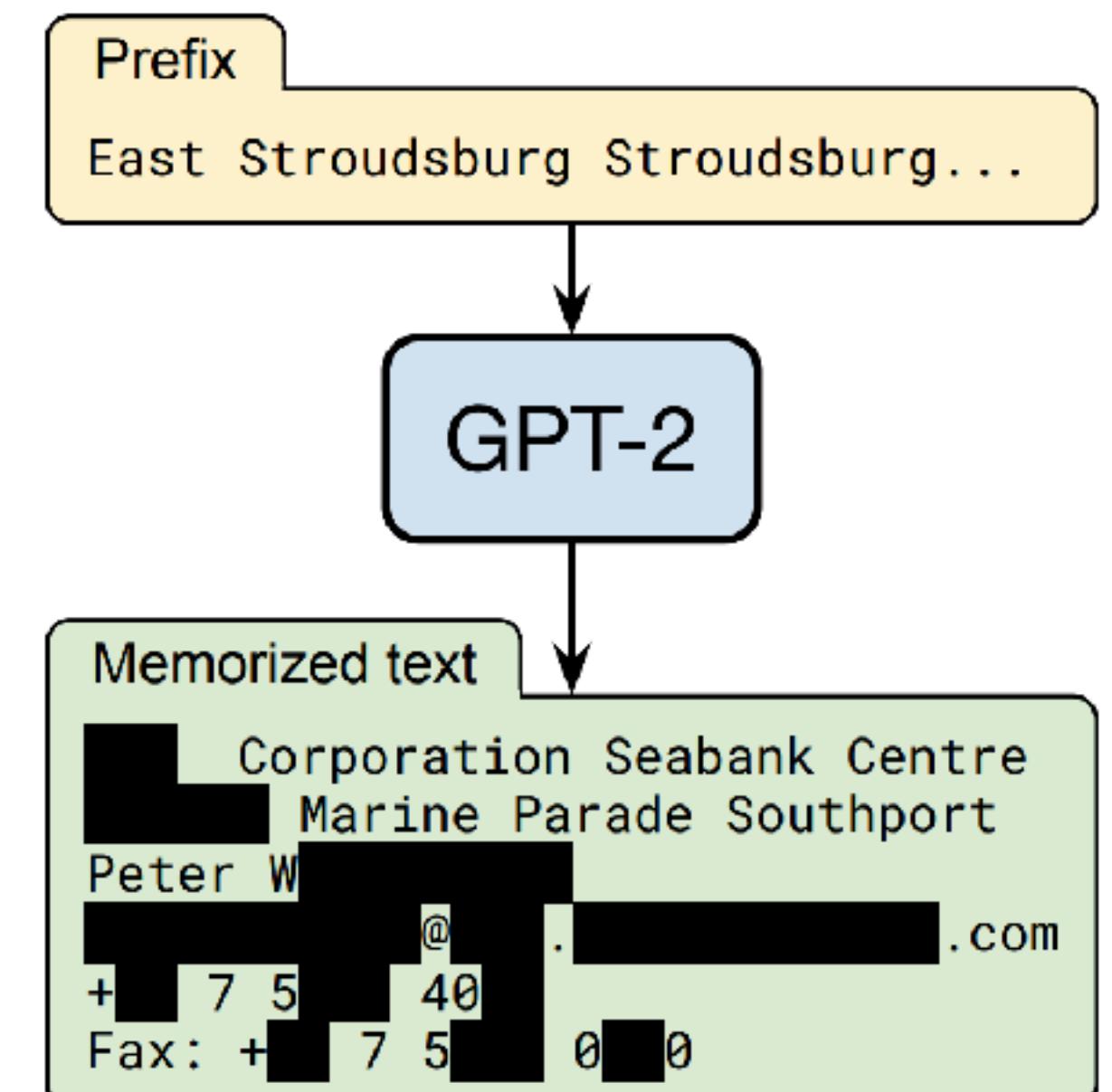
Carlini et al., 2023

# Verbatim Memorization

- Memorization has been heavily studied
- Focusing on verbatim memorization



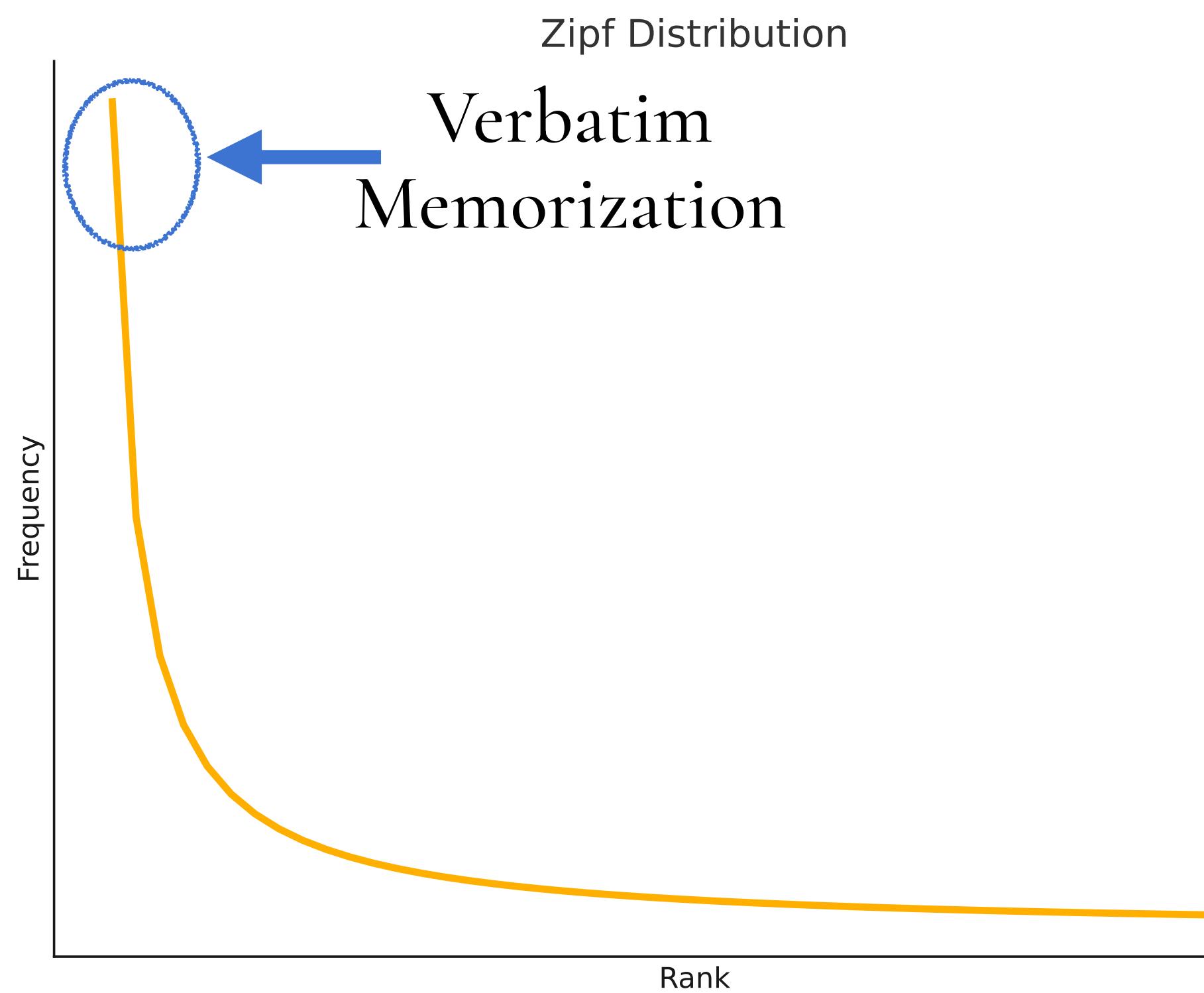
Carlini, Hayes & Nasr et al., 2023



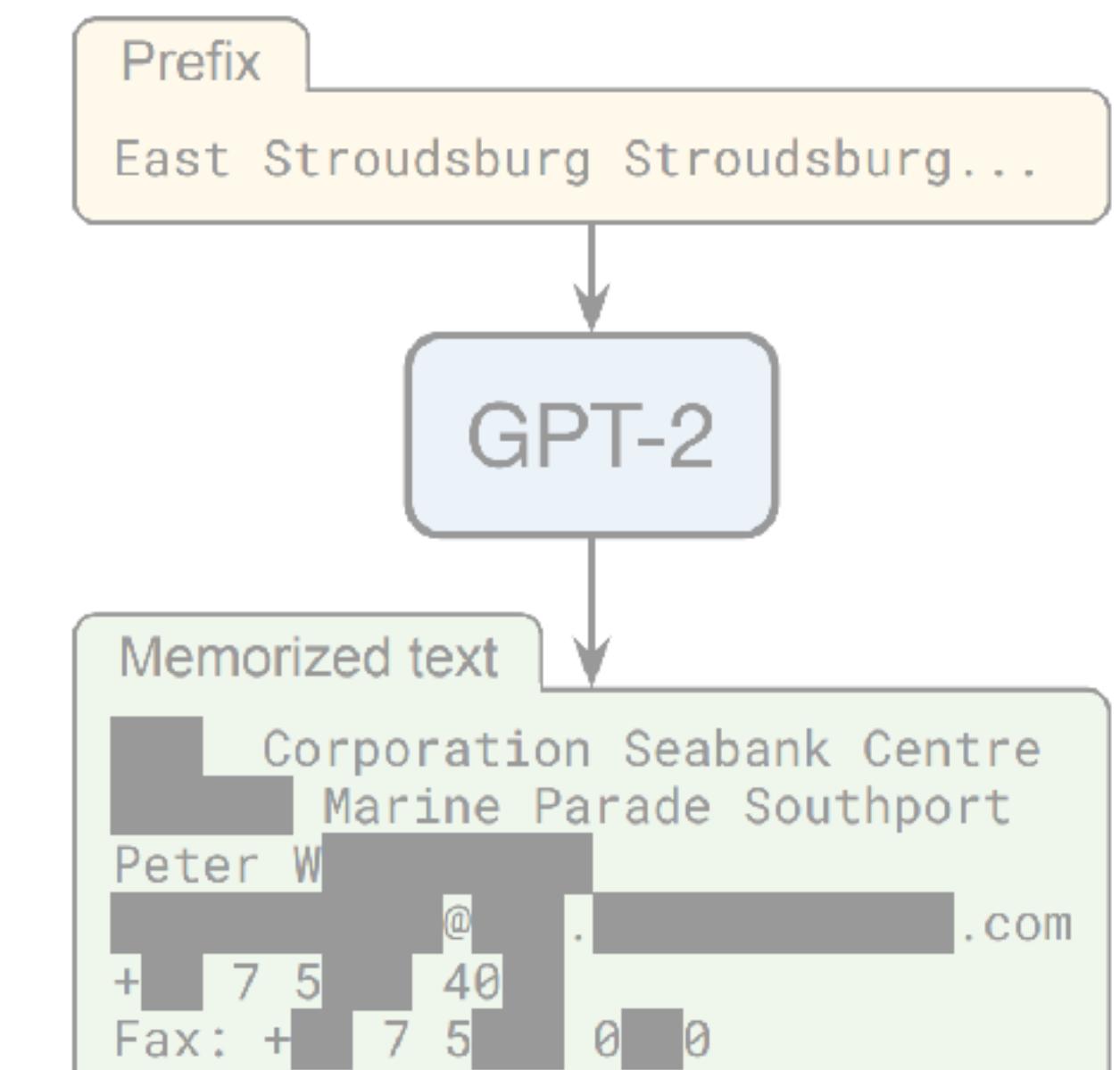
Carlini et al., 2021

# Verbatim Memorization

*Estimated Memorization (empirical lower) Bounds: 0.1%, 1%*

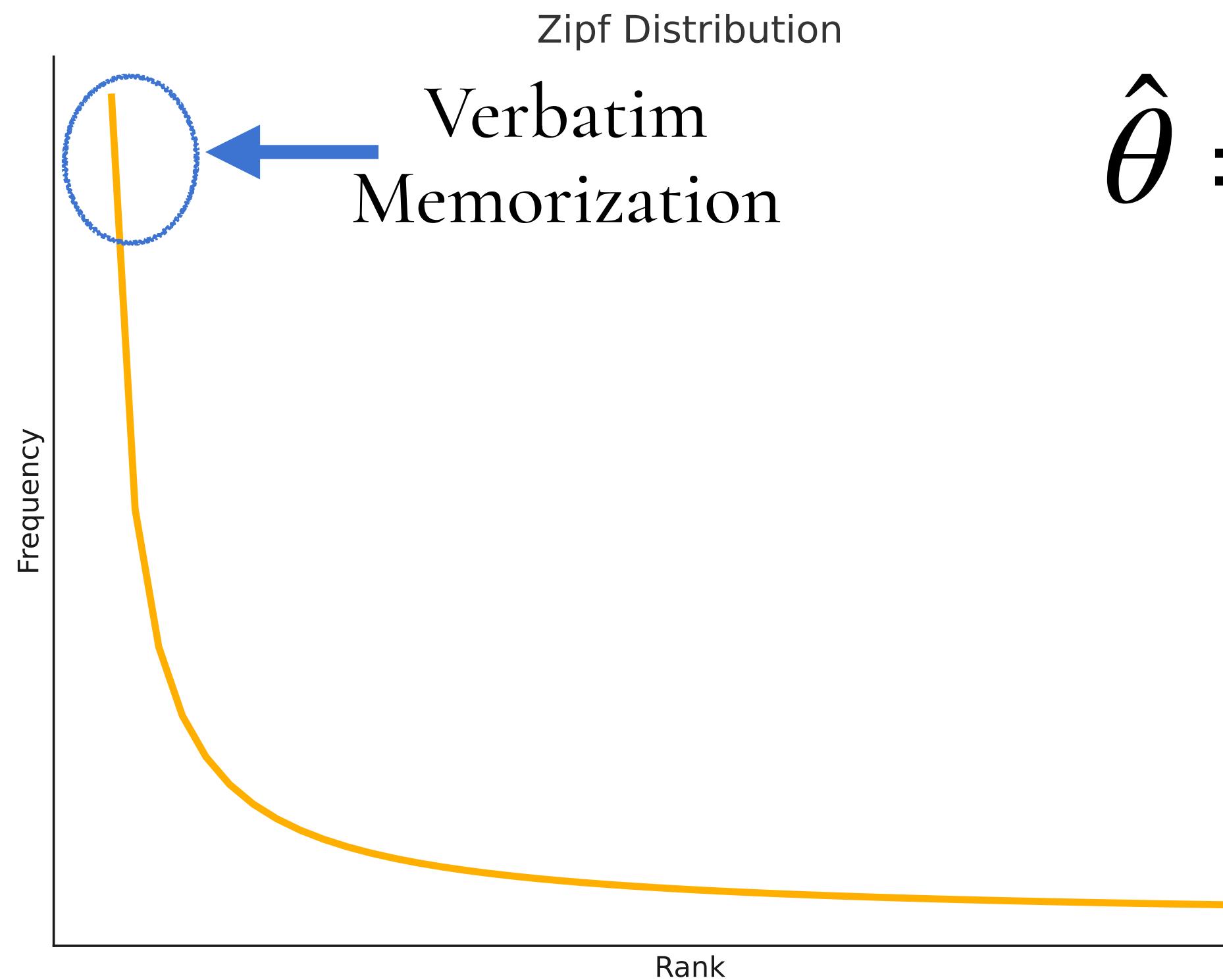


Carlini, Hayes & Nasr et al., 2023



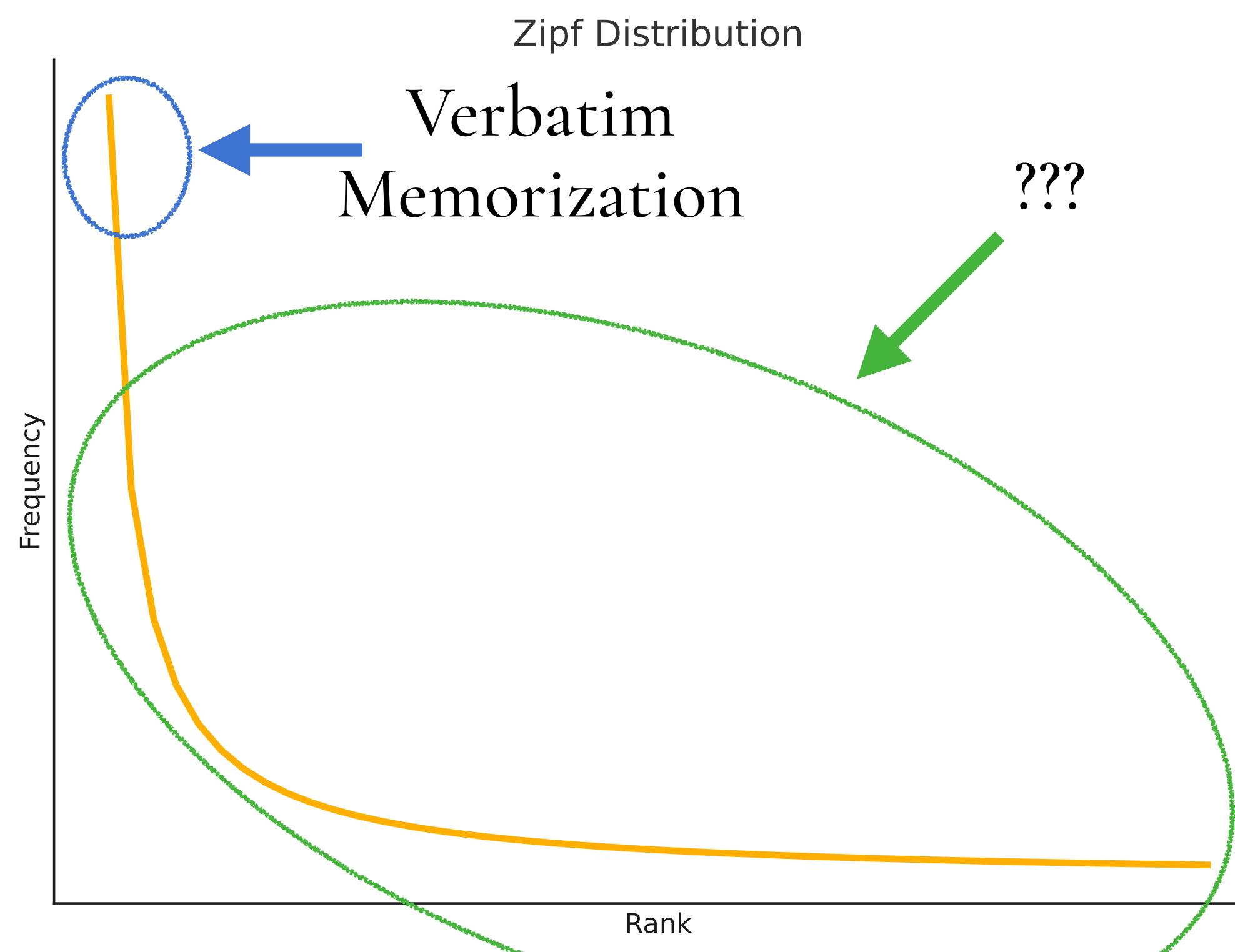
# Verbatim Memorization?

*Estimated Memorization (empirical lower) Bounds: 0.1%, 1%*



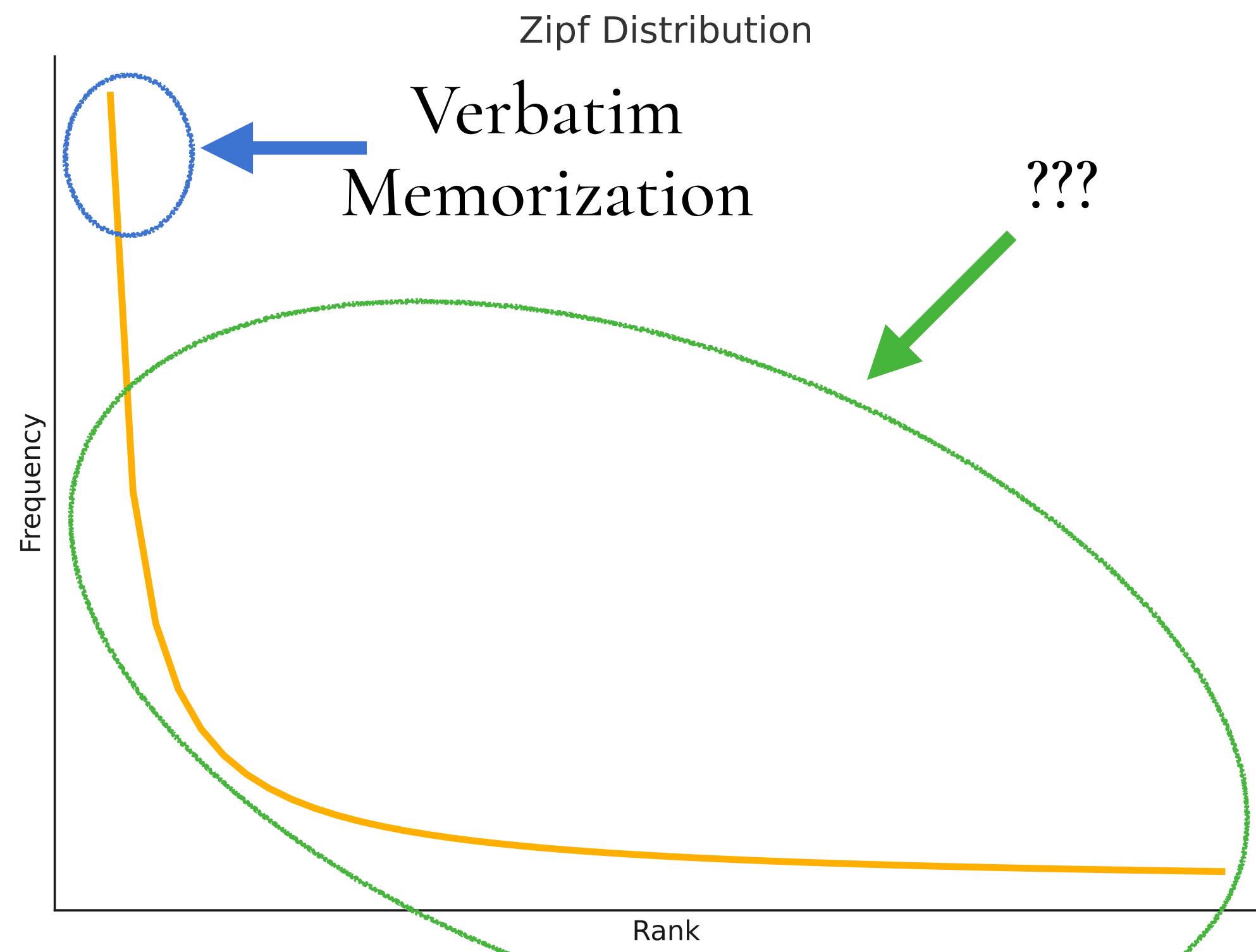
$$\hat{\theta} = \arg \max_{\theta} L(\theta)$$

# Verbatim Memorization?



# Verbatim Memorization?

What is learned from the rest of the data (99%)?

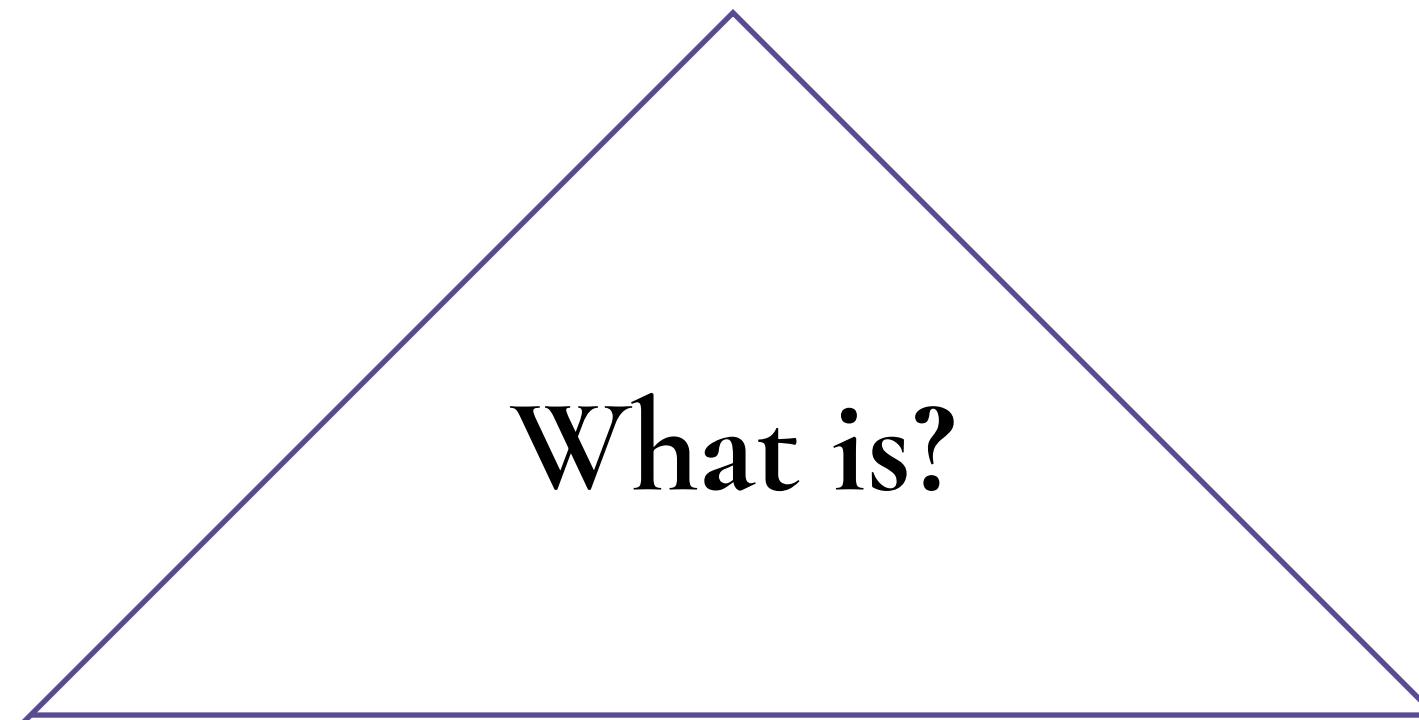


# Verbatim Memorization?

What is learned from the rest of the data (99%)?



# This Talk: Distributional Memorization



# This Talk: Distributional Memorization



What is?

How much  
it happens?



# This Talk: Distributional Memorization

What is?

How much  
it happens?

When does  
it happen?



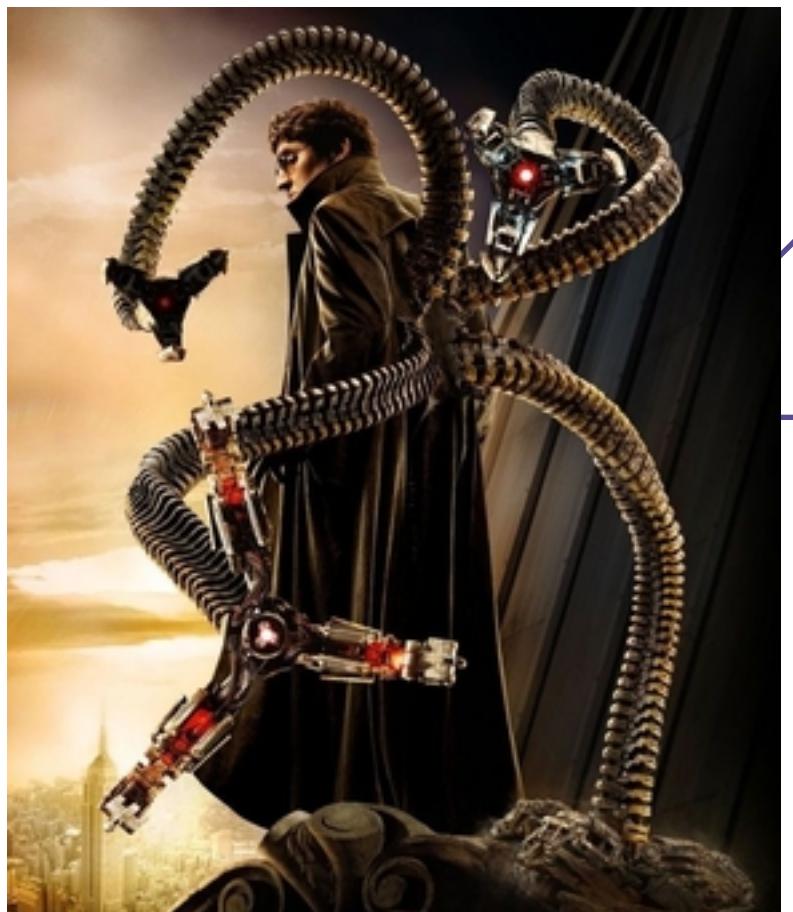
# This Talk: Distributional Memorization



What is?



How much  
it happens?



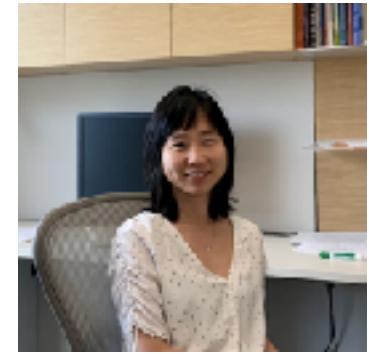
When does  
it happen?



Tools

# This Talk: Distributional Memorization

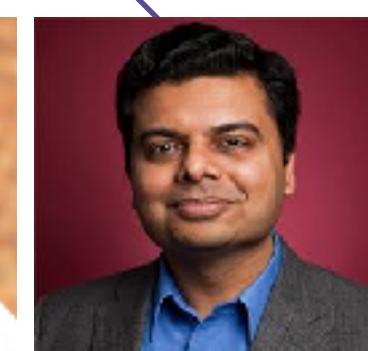
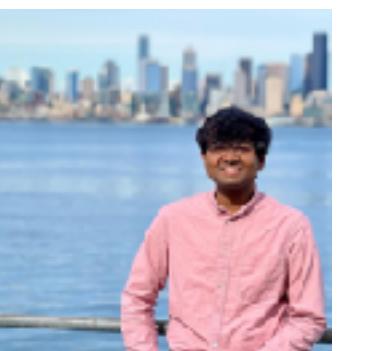
What is?



How much  
it happens?

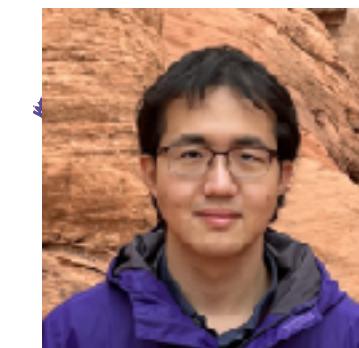
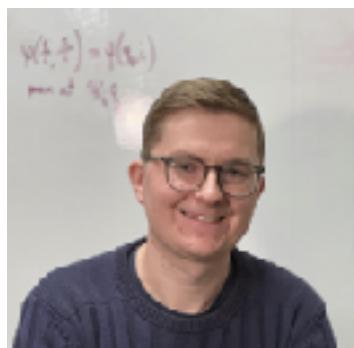


When does  
it happen?



Tools

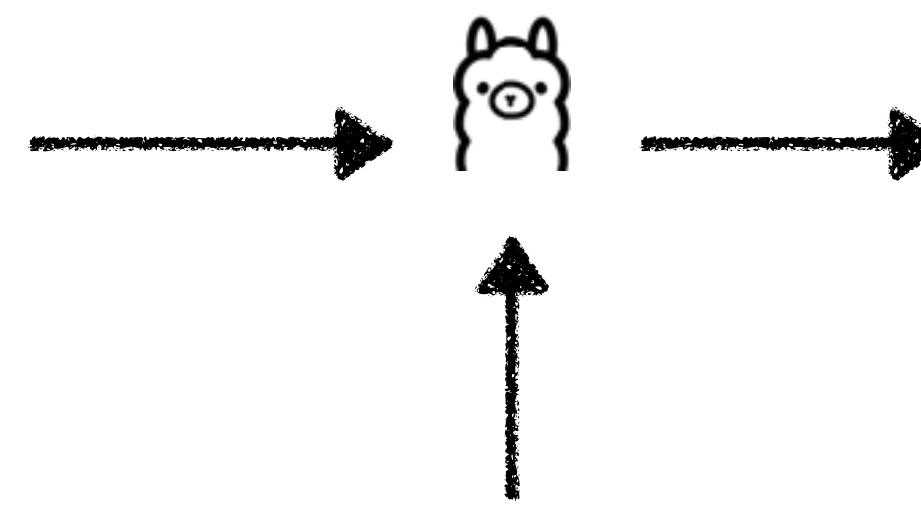
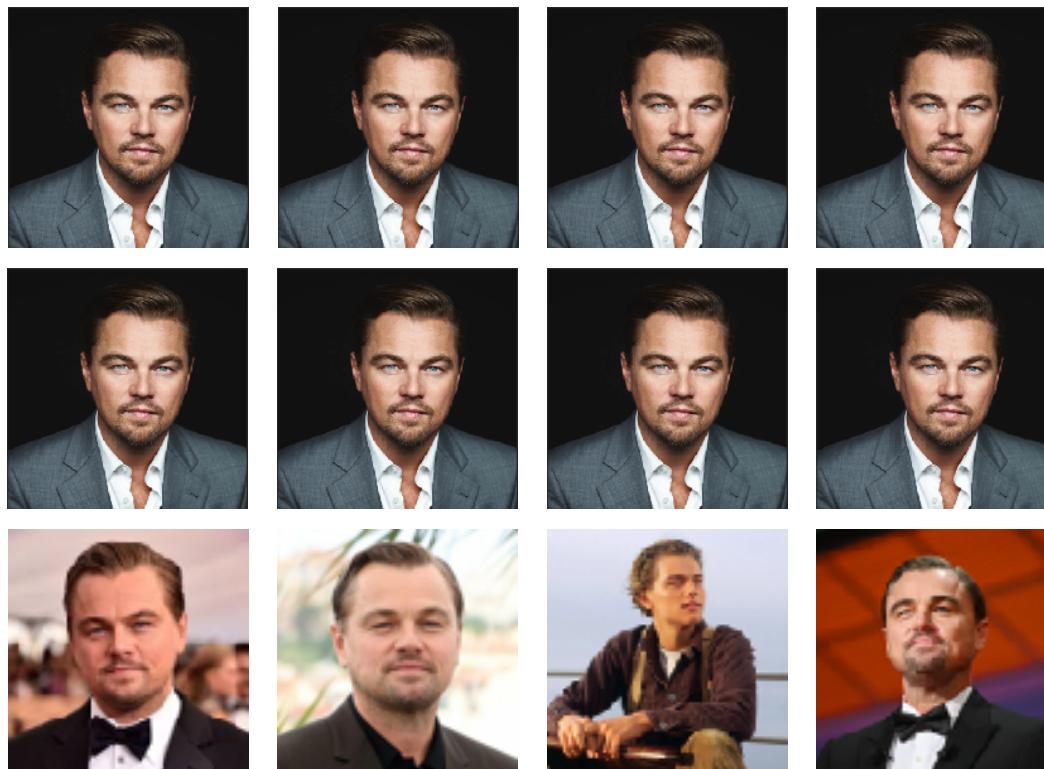
et el.



# What Is Distributional Memorizing

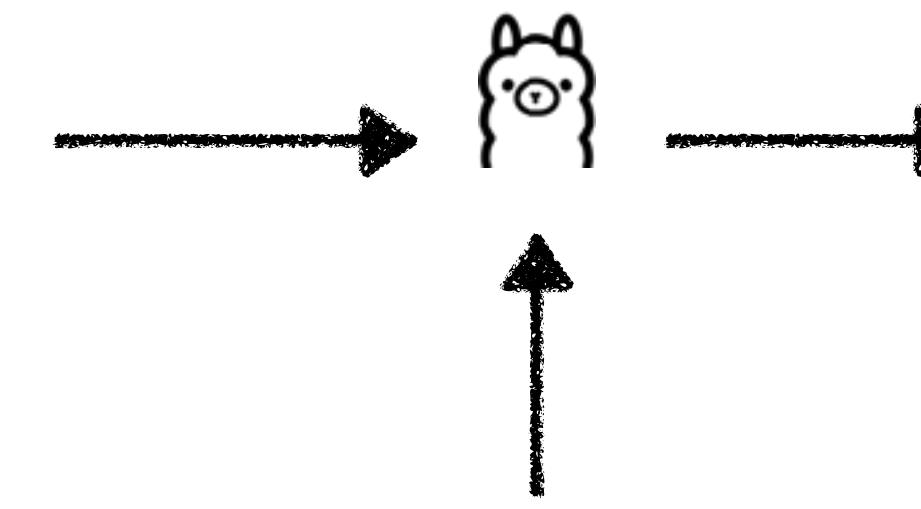
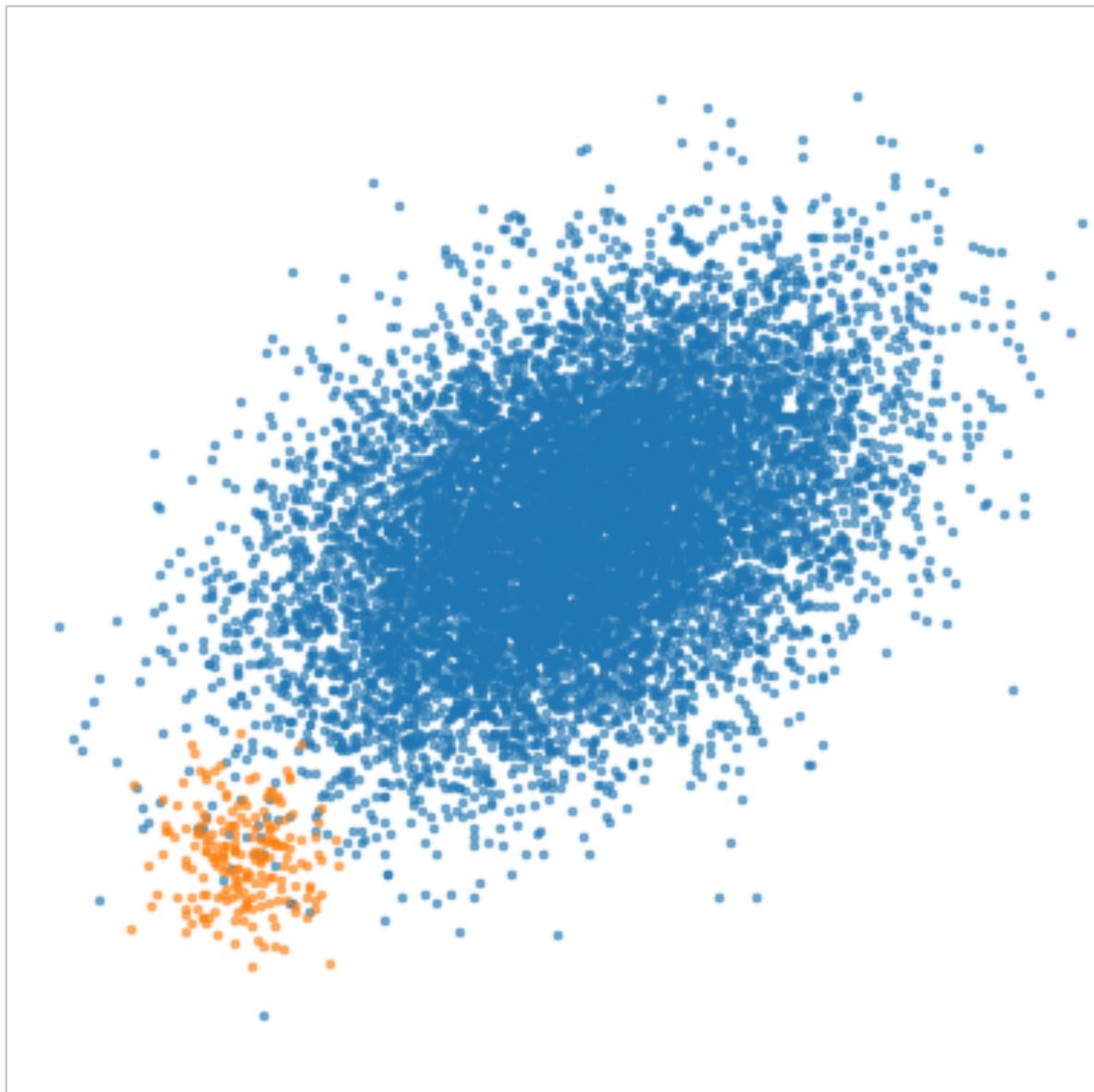
# Memorization Types

*Training Data*



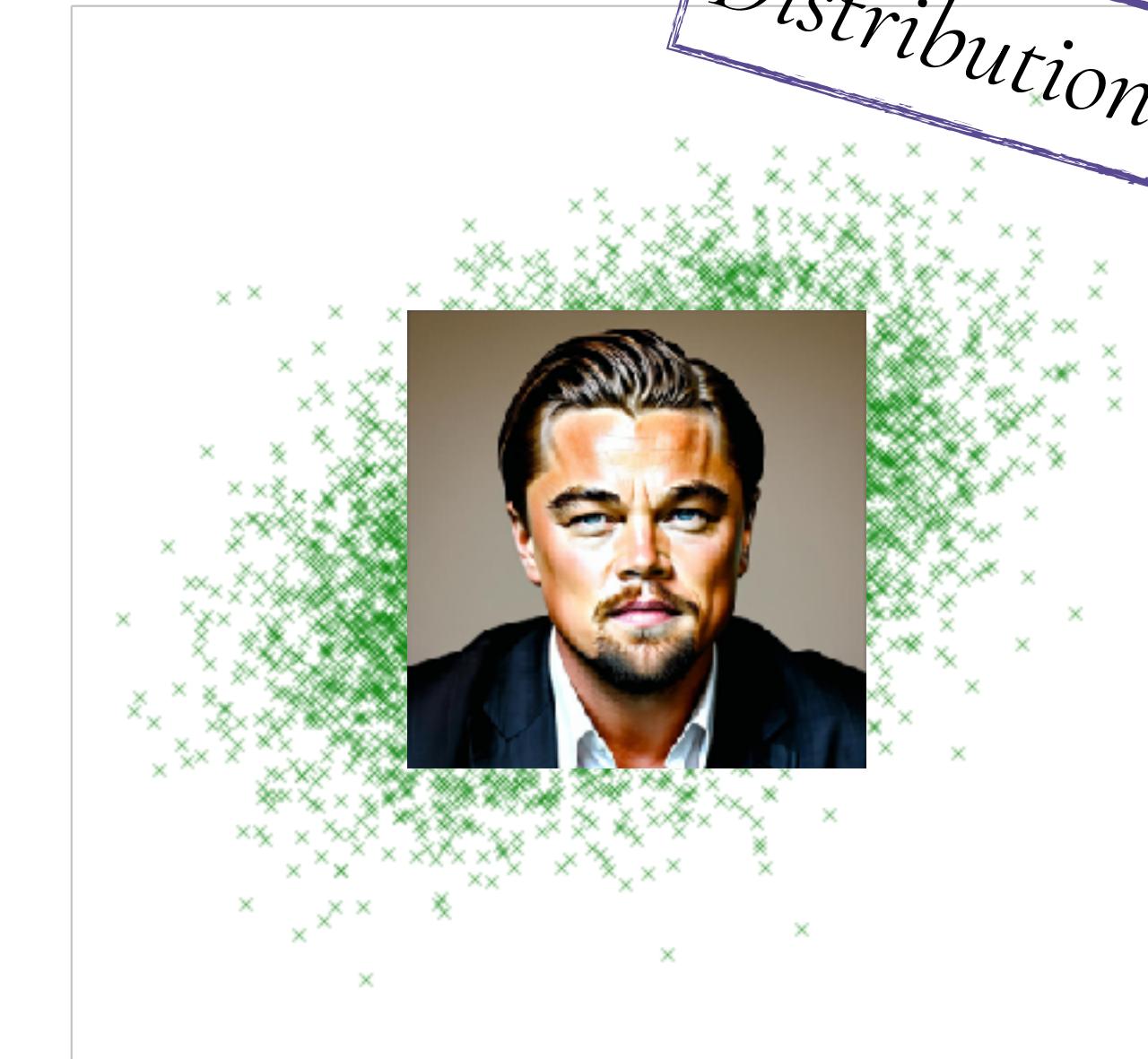
An Image of Leonardo DiCaprio

*Generation*



An Image of Leonardo DiCaprio

*Prompt*



*Distribution Memorization*

*Verbatim Memorization*

# Memorization Types

*Training Data*



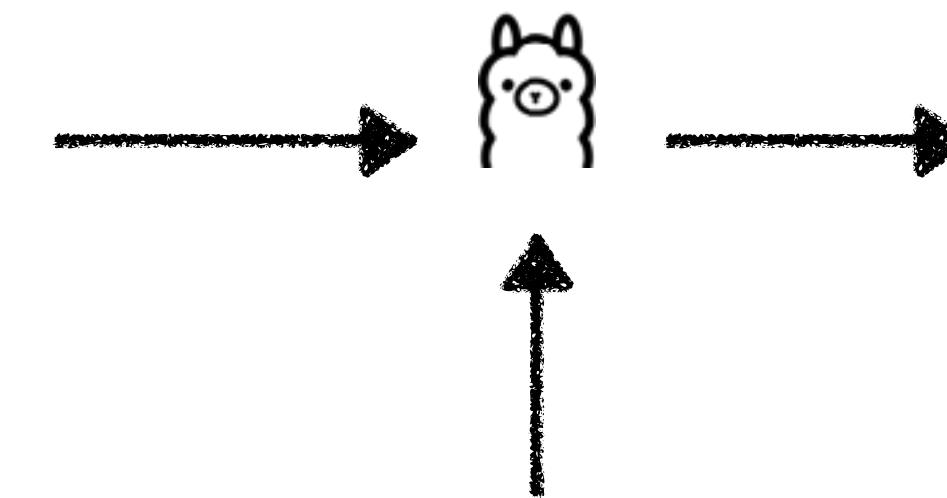
*Generation*



*Verbatim Memorization*

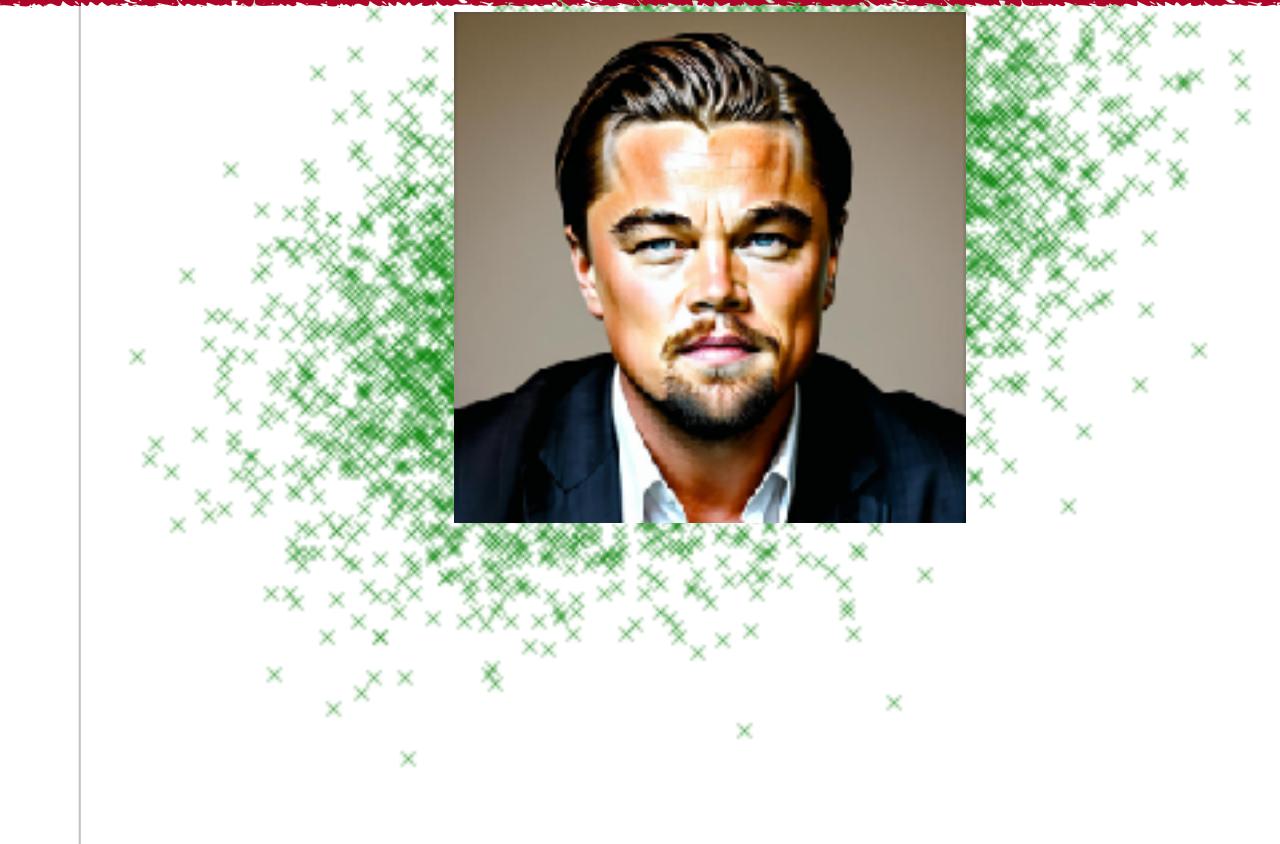
**Definition:** A model  $M$  memorize a distribution  $D$  from the training data if  $M$  can generate data  $D'$  that is indistinguishable from  $D$

*Memorization*



An Image of Leonardo DiCaprio

*Prompt*



# The Bias Amplification Paradox in Text-to-Image Generation

Preethi Seshadri, Sameer Singh, Yanai Elazar

NAACL 2024



# Models are Biased

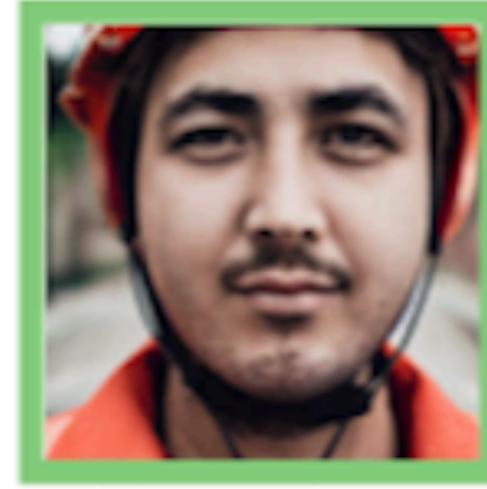
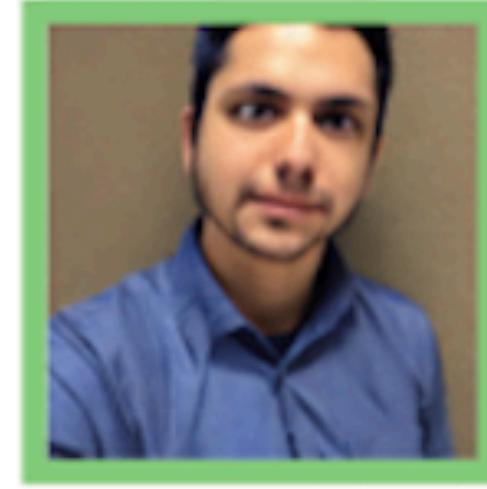
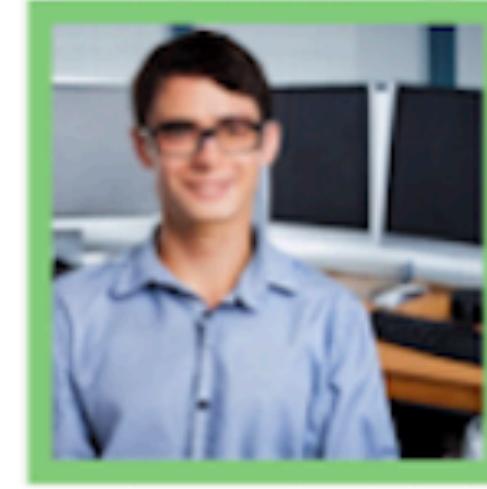
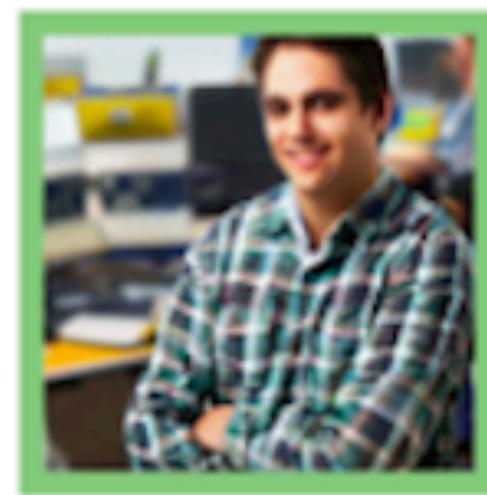
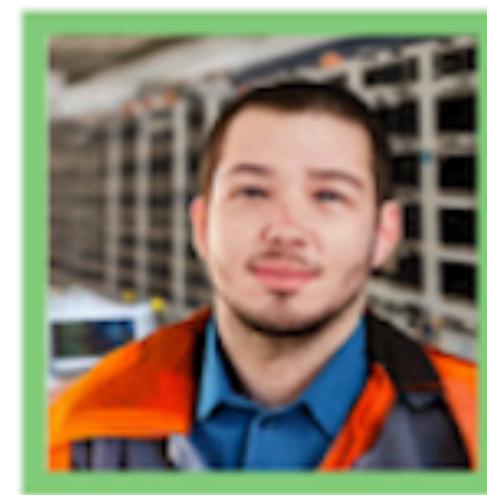
- Models encode and exhibit different biases
- This is not a new finding!

# Let's Try It Out!

“A photo of a face of an engineer”

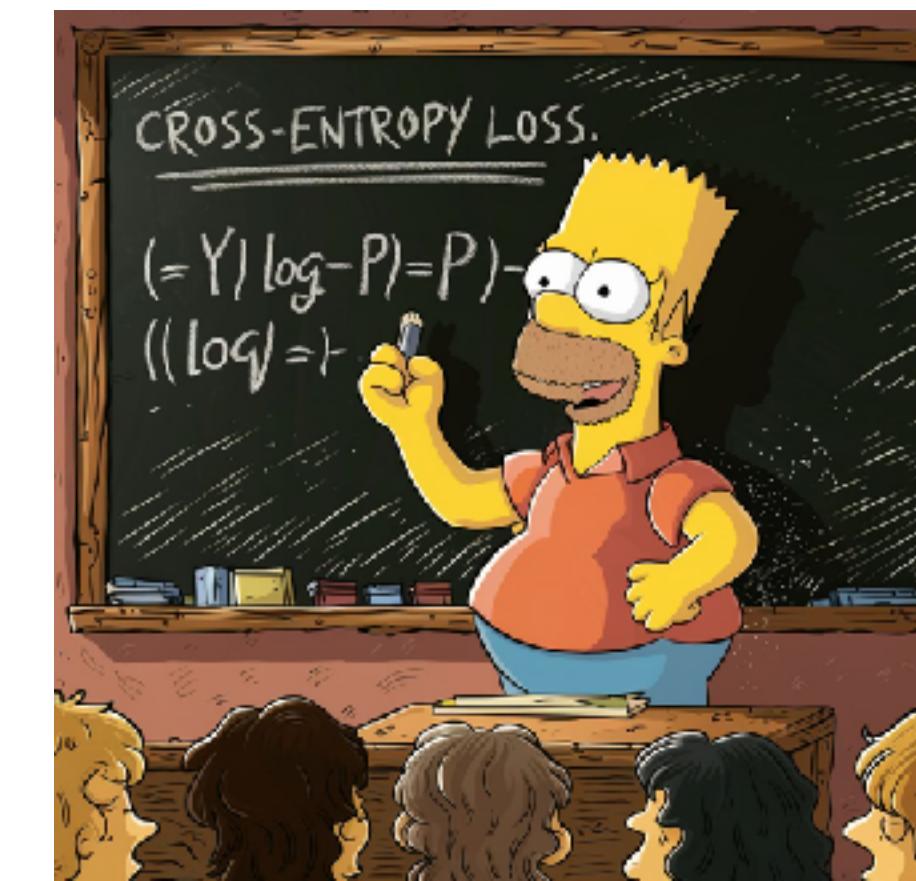
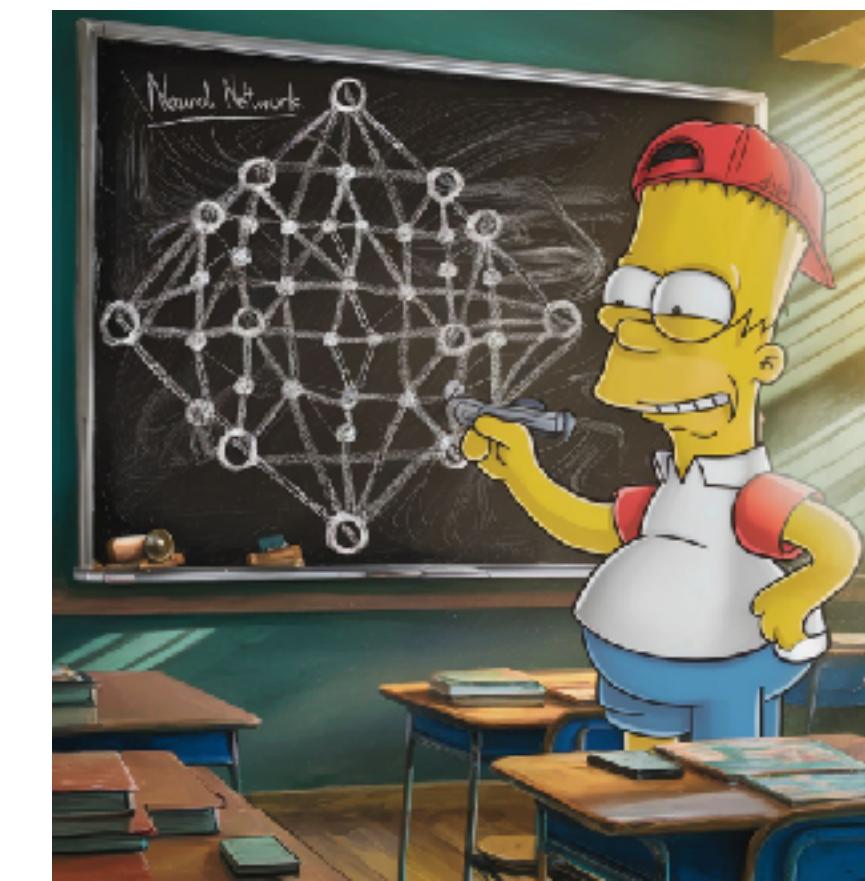


1/10 women!



The model is biased!

# Where Does The Bias Come From?



Let's Look At The Data

# Where Does The Bias Come From?

5 billion image-caption pairs!



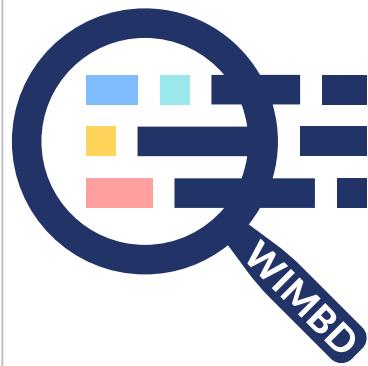
# Where Does The Bias Come From?

- Using an index (WIMBD), we have fast access to the data
- ... and we can test such associations in the training data



# Where Does The Bias Come From?

```
from wimbd.es import get_documents_containing_phrases  
  
# Get documents containing the term:  
get_documents_containing_phrases("laion","engineer")
```



ENGI  
Engin  
Mater

The data is large and noisy, so we need to adjust



# Establishing Data Gender Ratios

```
from wimbd.es import get_documents_containing_phrases  
  
# Get documents containing the term:  
get_documents_containing_phrases("laion", "engineer")
```



We follow a similar process for the generated images



Filtering



Gender identification



# Setup

- We sample image-caption pairs: 500 total
- 62 occupations:

# Setup

- We sample image-caption pairs: 500 total
- 62 occupations:
  - Accountant



# Setup

- We sample image-caption pairs: 500 total
- 62 occupations:
  - Accountant
  - Chef



# Setup

- We sample image-caption pairs: 500 total
- 62 occupations:
  - Accountant
  - Chef
  - Engineer



# Setup

- We sample image-caption pairs: 500 total
- 62 occupations:
  - Accountant
  - Chef
  - Engineer
  - Janitor



# Setup

- We sample image-caption pairs: 500 total
- 62 occupations:
  - Accountant
  - Chef
  - Engineer
  - Janitor
  - Lawyer



# Setup

- We sample image-caption pairs: 500 total
- 62 occupations:
  - Accountant
  - Chef
  - Engineer
  - Janitor
  - Lawyer
  - ...



# Bias Amplification?

Given the calculated ratios from the data, we can now compare the model's generation to the training data

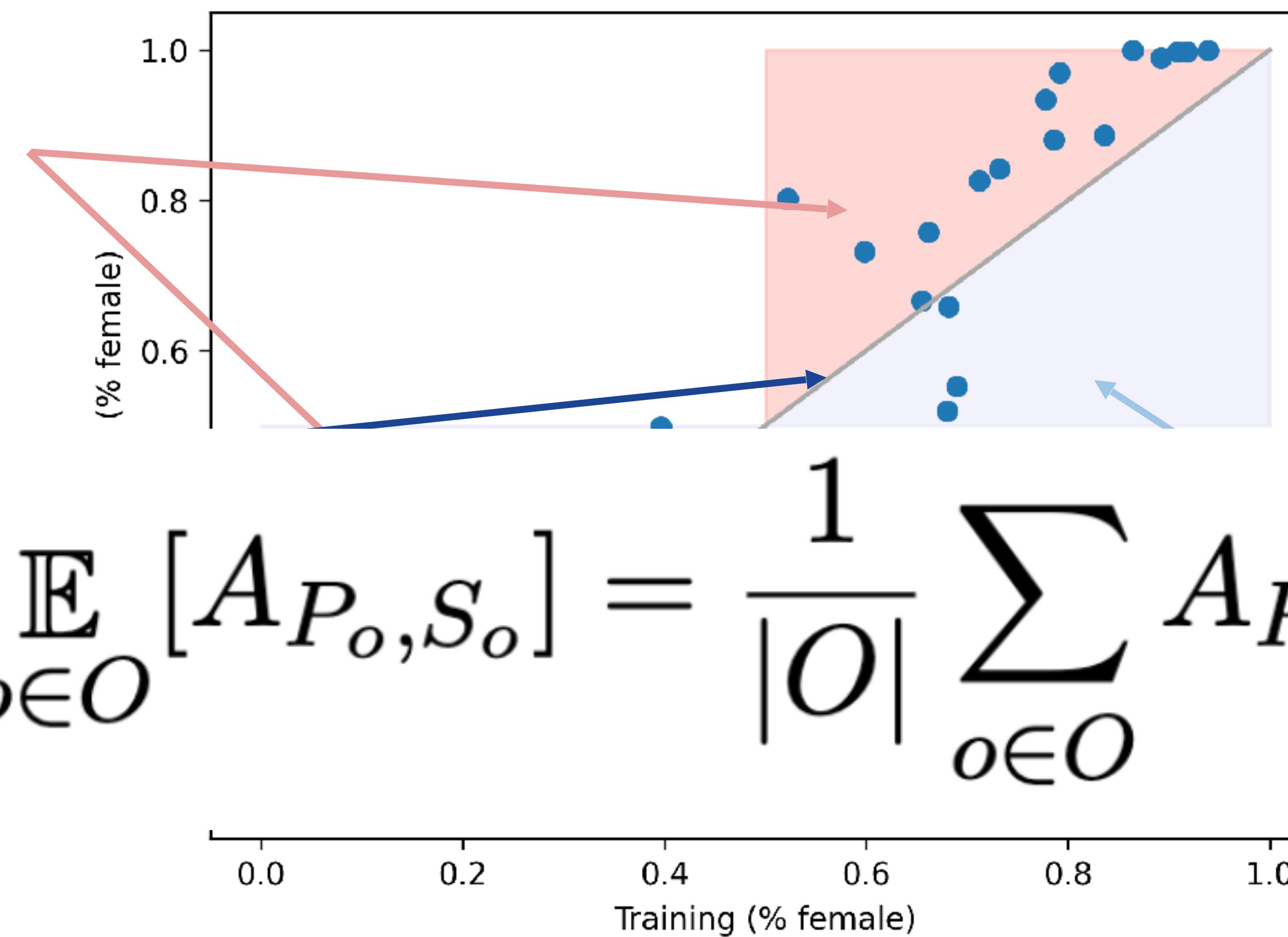
# Bias Amplification?

Peach area:  
Bias Amplification

*Diagonal:*  
*Bias preservation*

$$\mathbb{E}_{o \in O} [A_{P_o, S_o}] = \frac{1}{|O|} \sum_{o \in O} A_{P_o, S_o}$$

Training (% female)

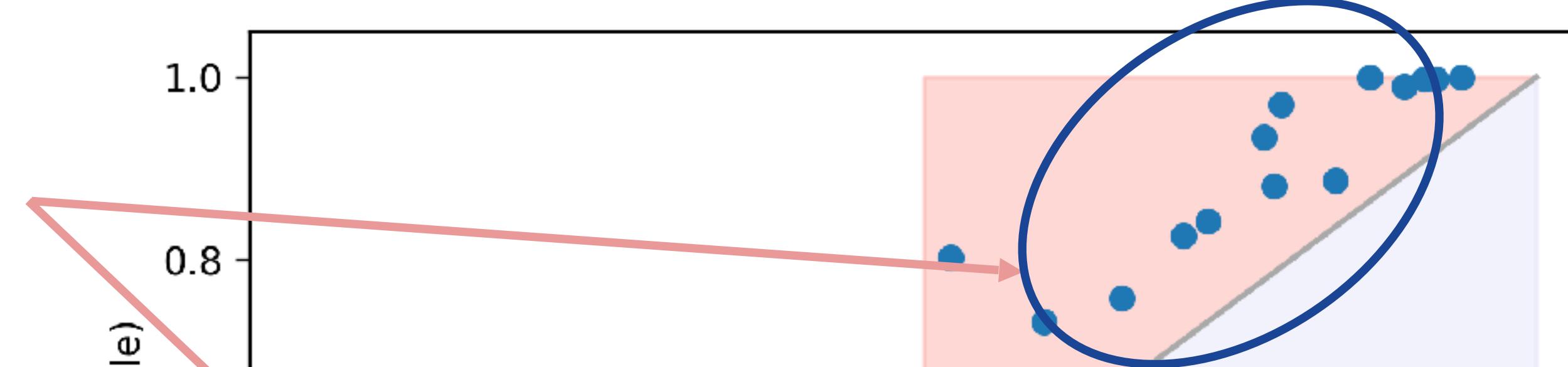


der area:  
e-amplification

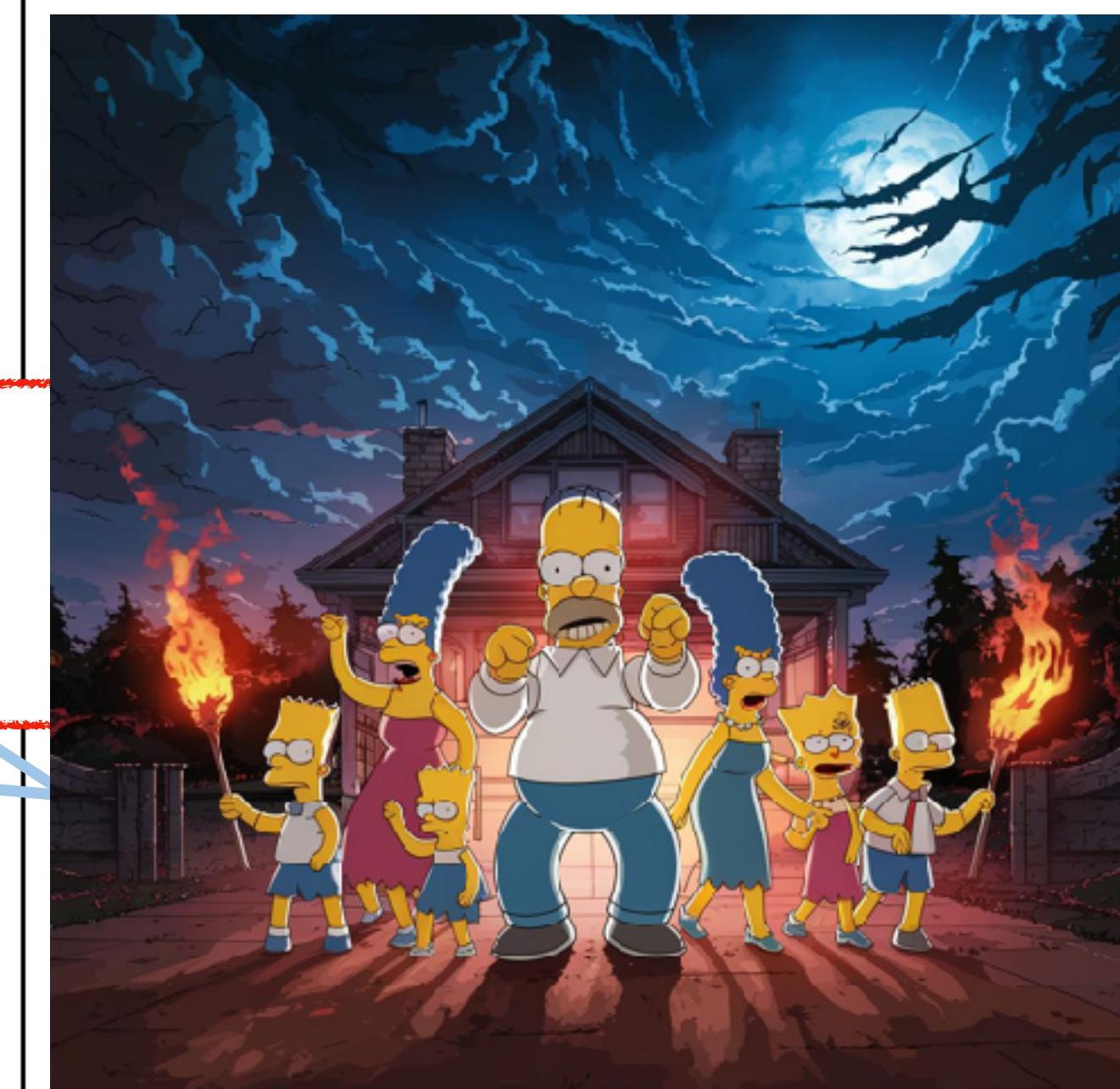
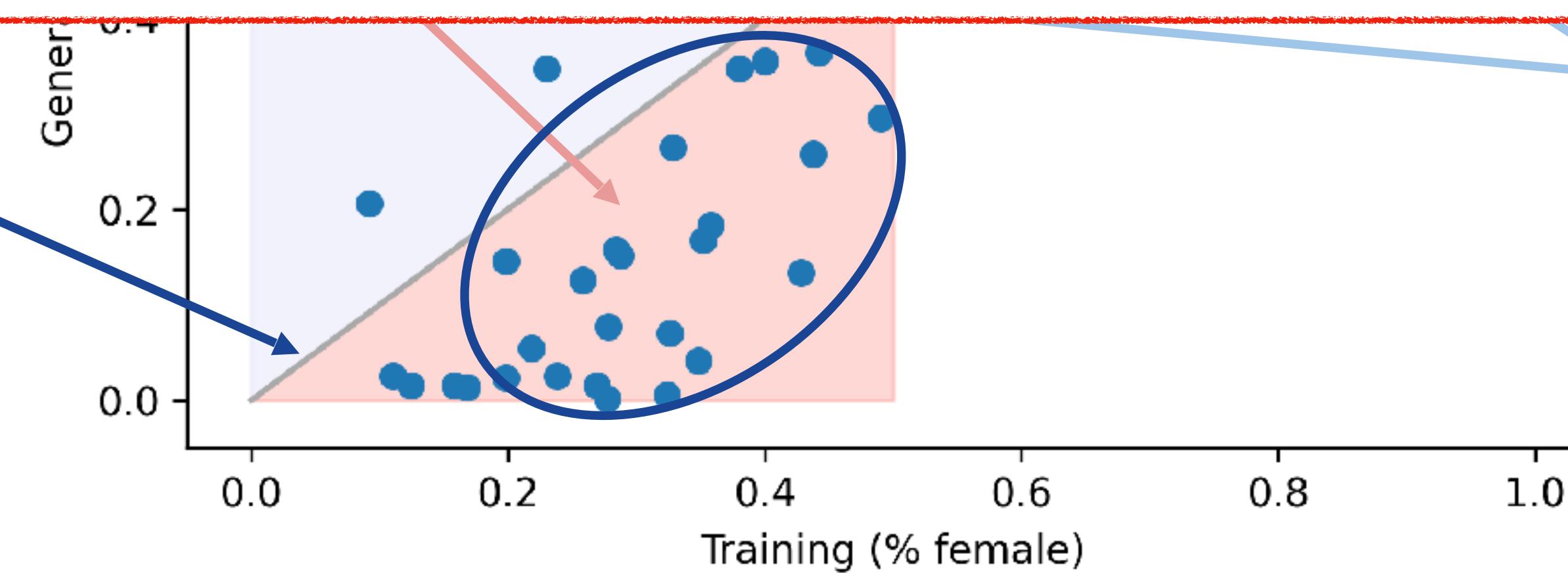
# Bias Amplification!

Peach area:  
Bias Amplification

*Diagonal:*  
*Bias preservation*



Bias is amplified by 12.57%



# Bias Amplification!

Shown in previous work

## **Men Also Like Shopping: Reducing Gender Bias Amplification using Corpus-level Constraints**

**Jieyu Zhao**<sup>§</sup>

**Tianlu Wang**<sup>§</sup>

**Mark Yatskar**<sup>†</sup>

**Vicente Ordonez**<sup>§</sup>

**Kai-Wei Chang**<sup>§</sup>

<sup>§</sup>University of Virginia

{jz4fu, tw8cb, vicente, kc2wc}@virginia.edu

<sup>†</sup>University of Washington

my89@cs.washington.edu

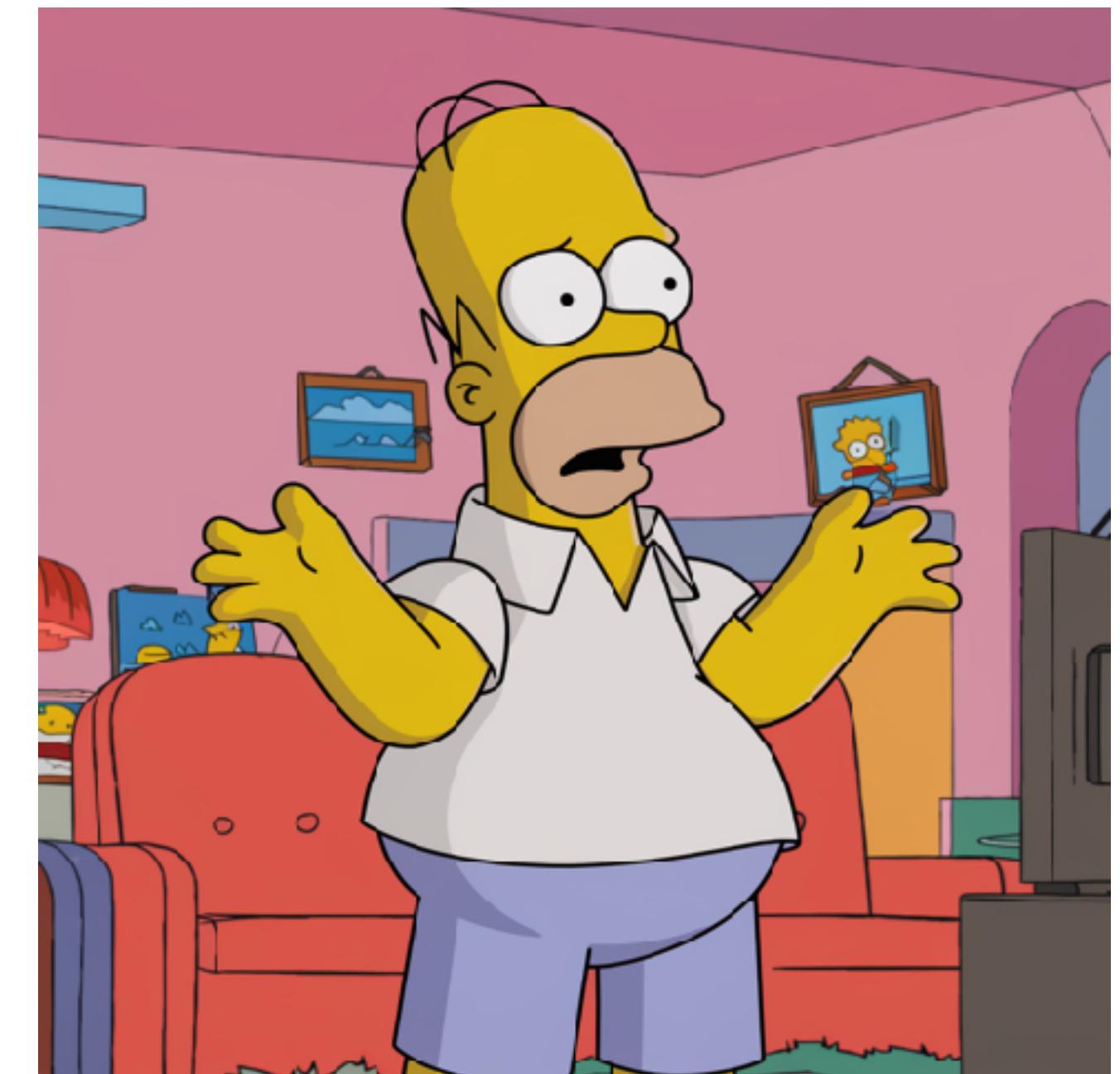
*Zhao et al, 2017, best paper*

# The Bias Amplification Paradox

But wait!

Why would a model amplify the biases from the training data?

Let's look at the training data again



# Training Data Investigation

Portrait of young **woman** programmer working at a computer in the data center filled with display screens

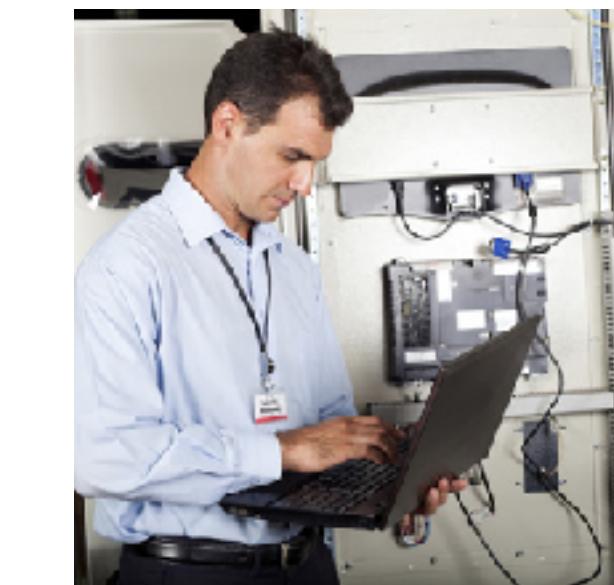


Slow motion **programmer female** relaxing among nature, young **woman** on long-awaited vacation abroad after working year...



shutterstock - 669546292

programmer configures the... | Shutterstock . vector #669546292



industrial programmer checking computerized machine status

# Training Data Investigation

~60% contain gender indicators

Mostly with anti-stereotypical gender (70%)



shutterstock - 669546292



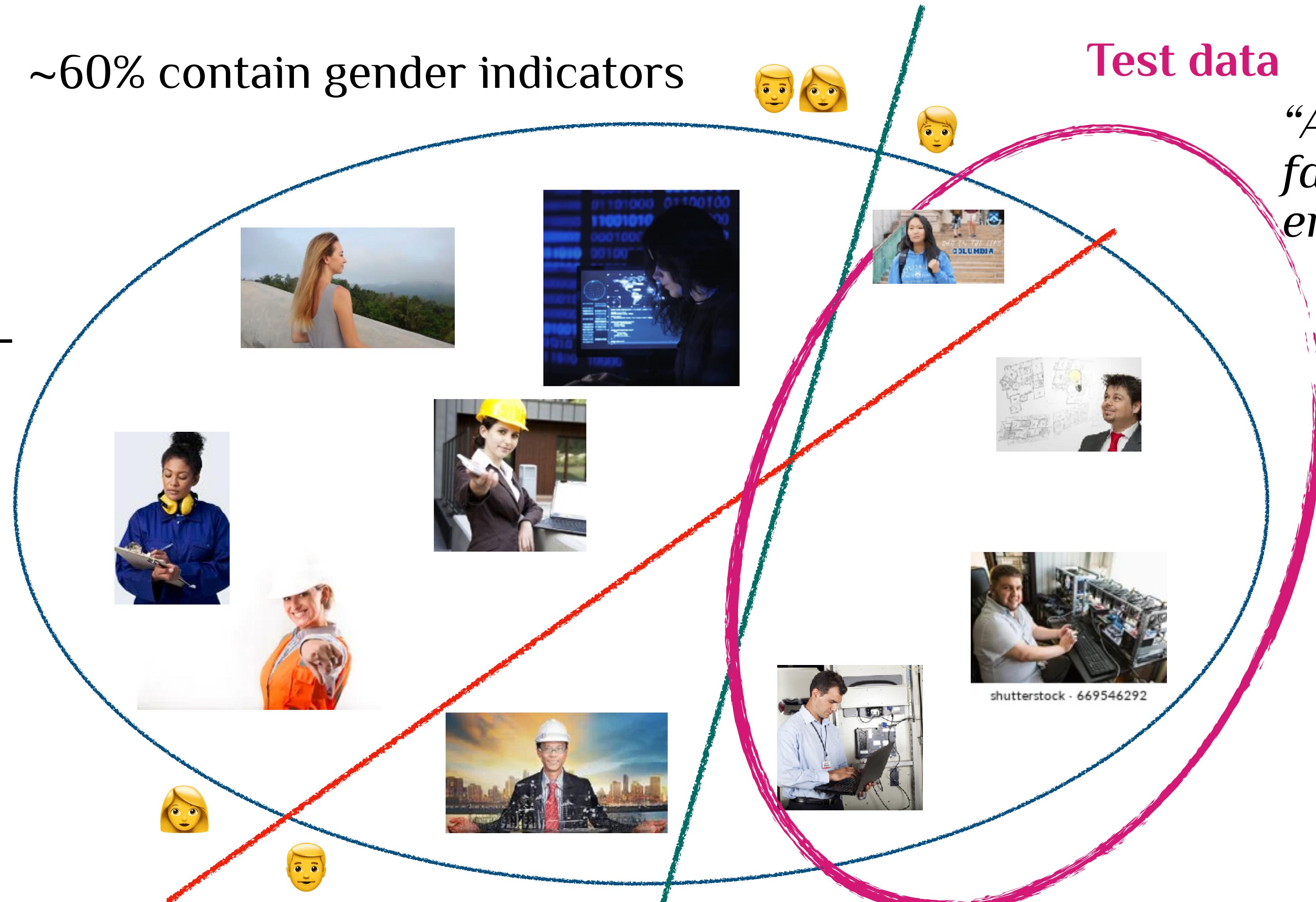
# Training Data Investigation

Mostly with anti-stereotypical gender (70%)

~60% contain gender indicators

Test data

*“A photo of a face of an engineer”*



# Image Captions & Prompts Mismatch

Training data

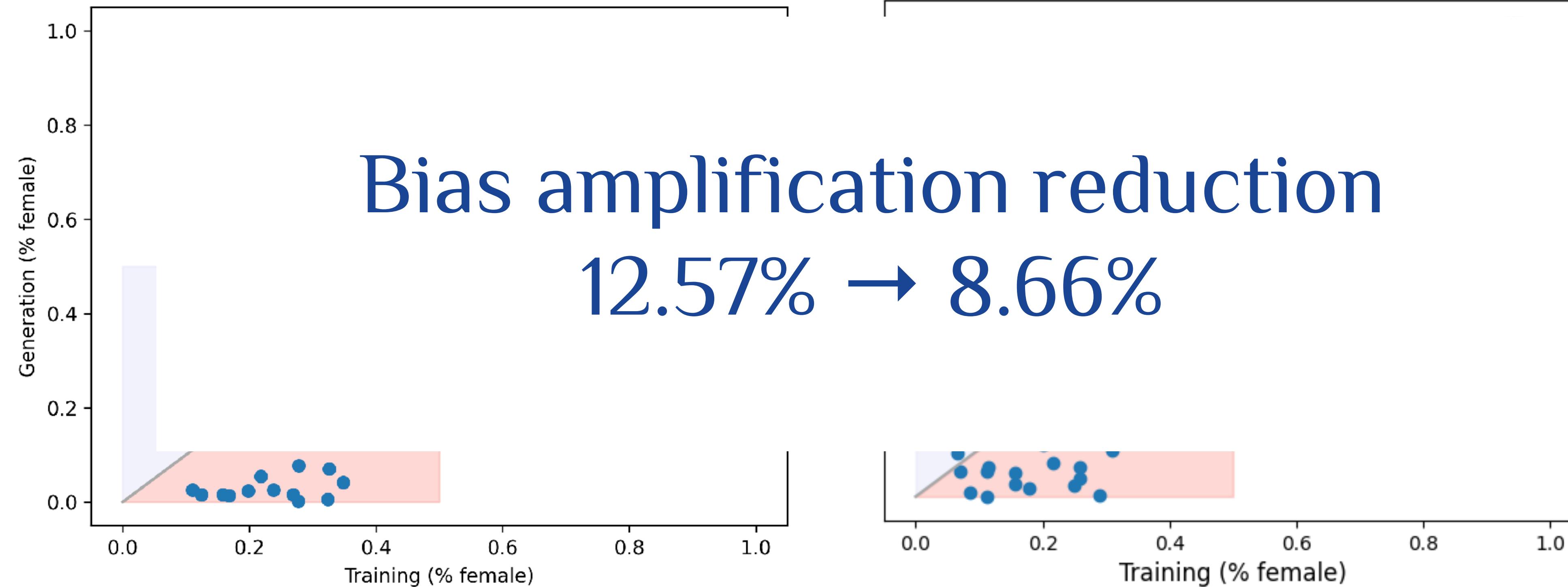
Test data



# Matching Distributions

Instead of comparing the generated images to the entire training set:

- We only compare to the captions with no gender indicators



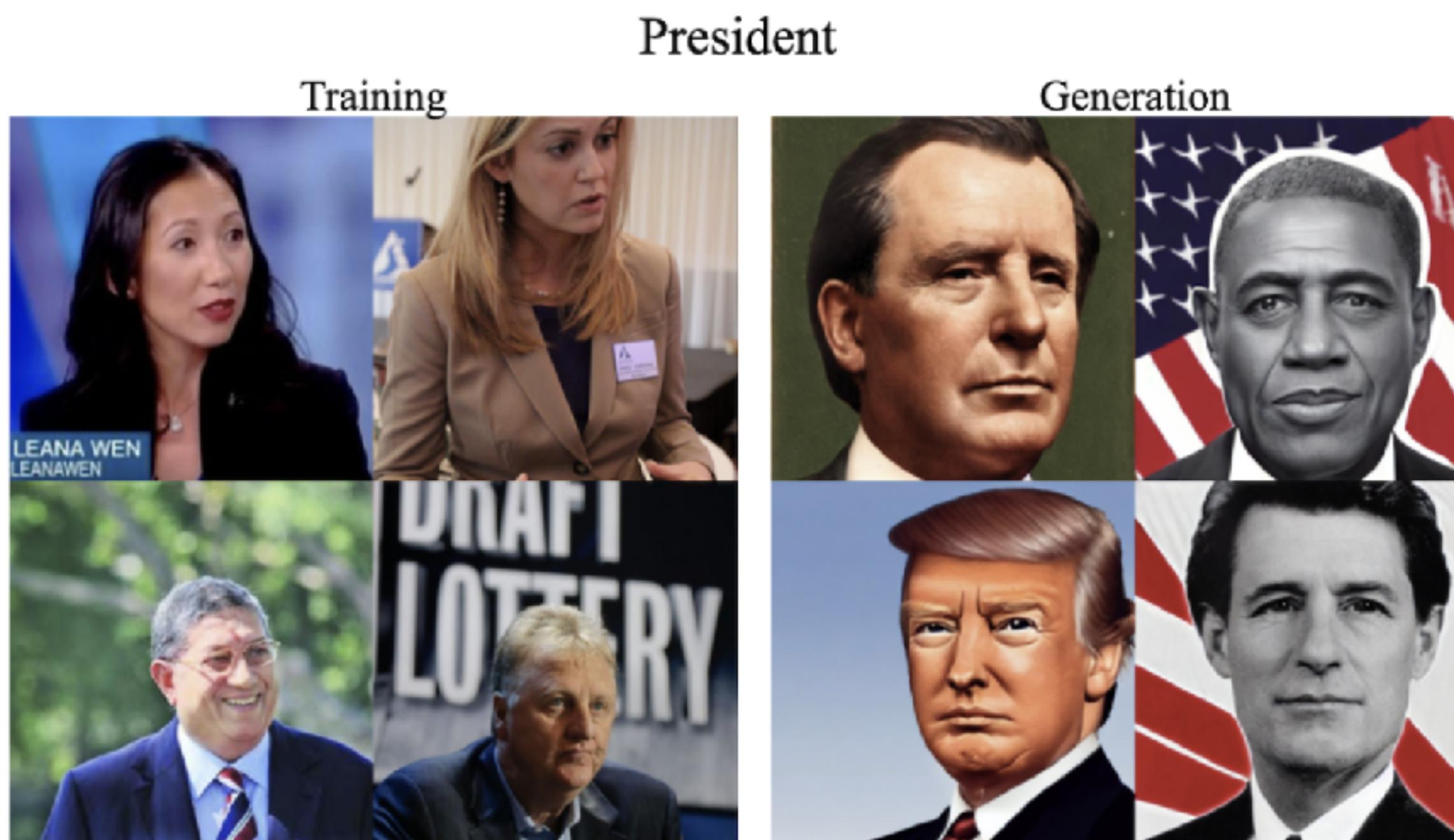
# One Mismatch

What about others?



# Image Captions & Prompts Mismatch #2

We also found a “

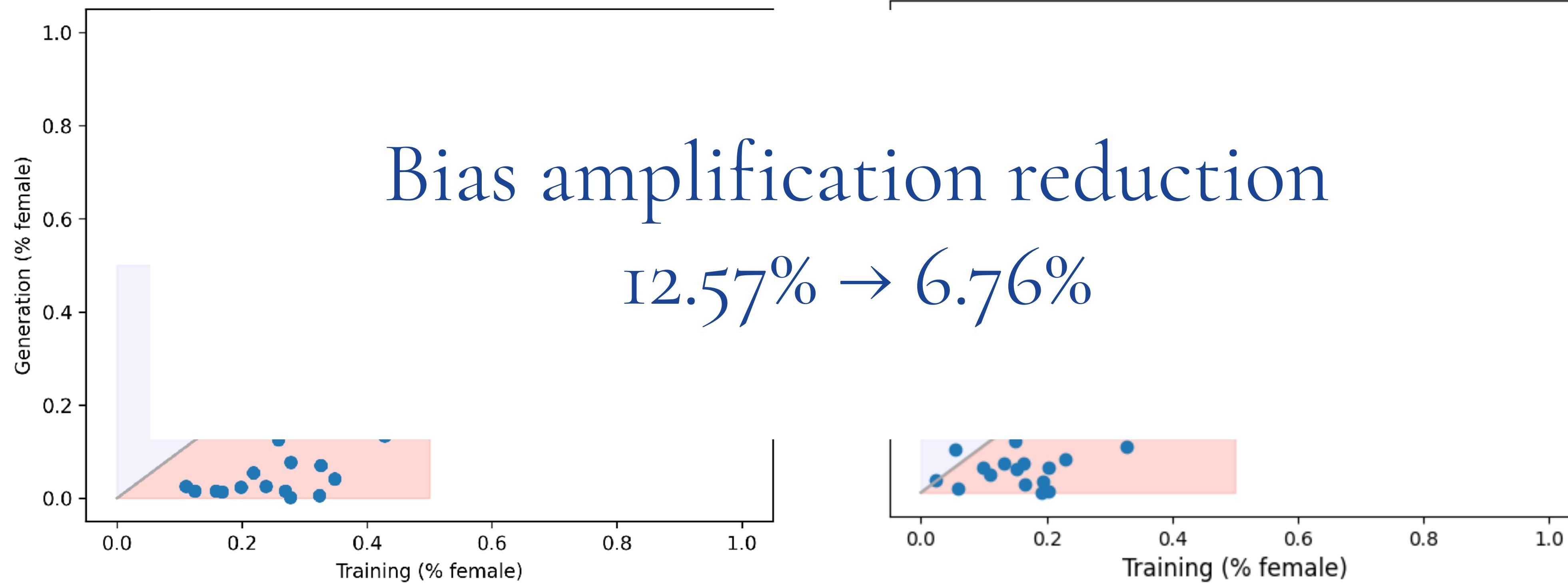


- (a) Training captions for **President**: 1) "Leana Wen, Planned Parenthood president..." 2) "New Schaumburg Business Association President..." 3) "BCCI president N Srinivasan..." 4) "Indiana Pacers president of basketball operations..."

# Matching Distributions #2

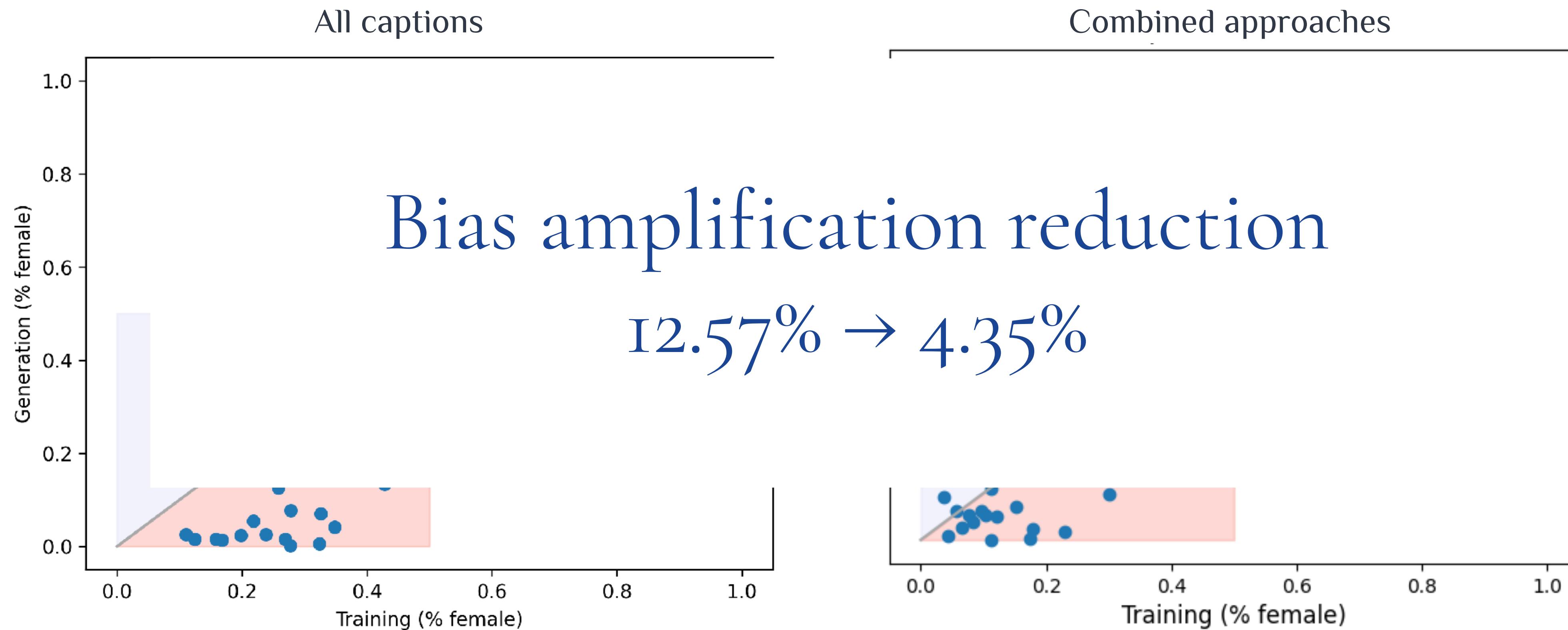
Instead of comparing the generated images to the entire training set:

- We compare to the captions that are similar to the prompts



# Matching Distributions: Combined

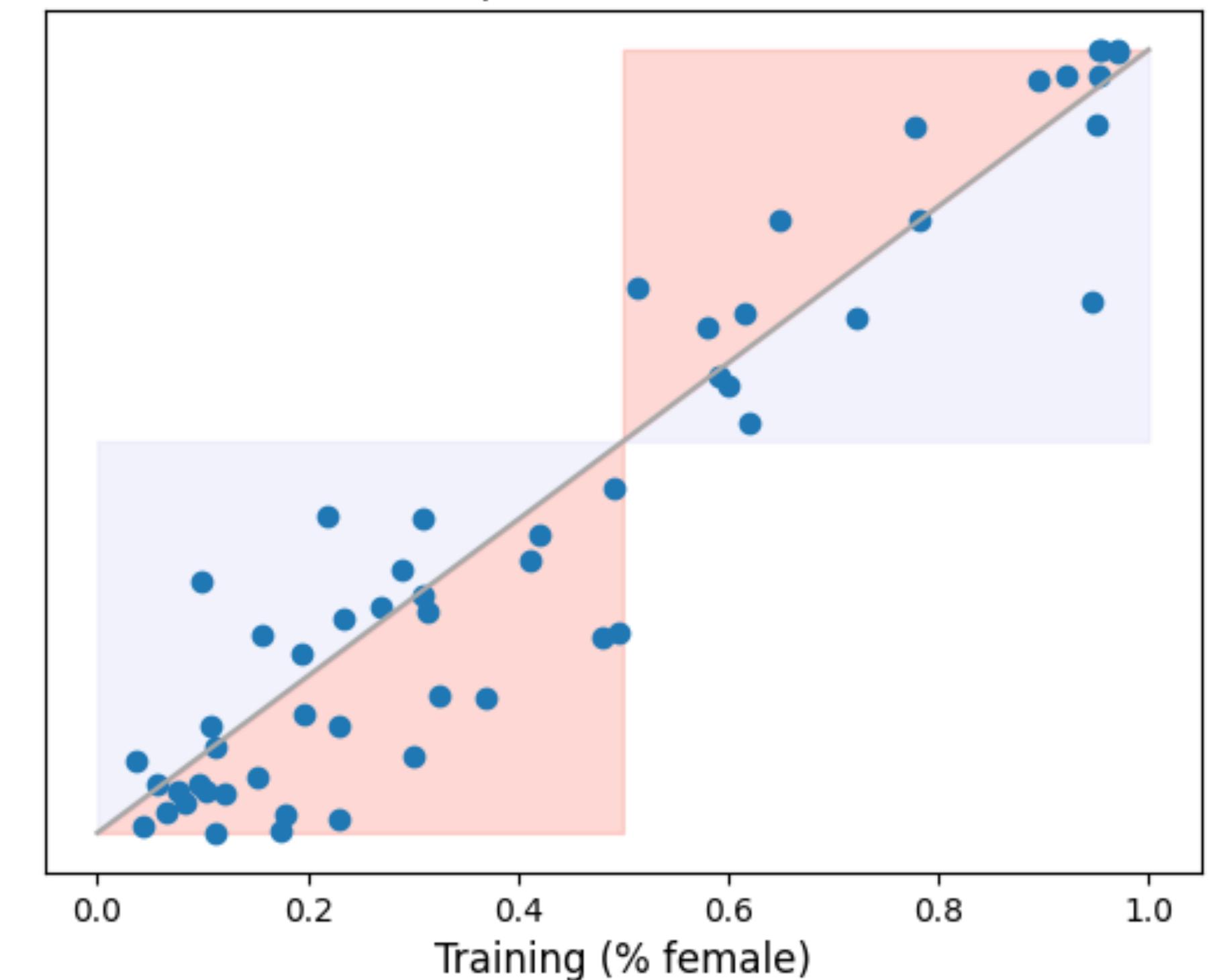
Finally, we combine both approaches



# Revisiting the Bias Amplification Claim

While we still observe bias amplification:

- It is significantly reduced
- There may be more confounders
- This problem is more nuanced and involved than originally thought



# What is Distributional Memorization?

- Models memorize gender ratios from training data
- Evaluation becomes tricky
- We should think about the right abstractions



What is

# How Much Distributional Memorization Happen?

# Memorized Sequences

## Extracting Training Data from Large Language Models

Nicholas Carlini<sup>1</sup>

Florian Tramèr<sup>2</sup>

Eric Wallace<sup>3</sup>

Matthew Jagielski<sup>4</sup>

Ariel Herbert-Voss<sup>5,6</sup>

Katherine Lee<sup>1</sup>

Adam Roberts<sup>1</sup>

Tom Brown<sup>5</sup>

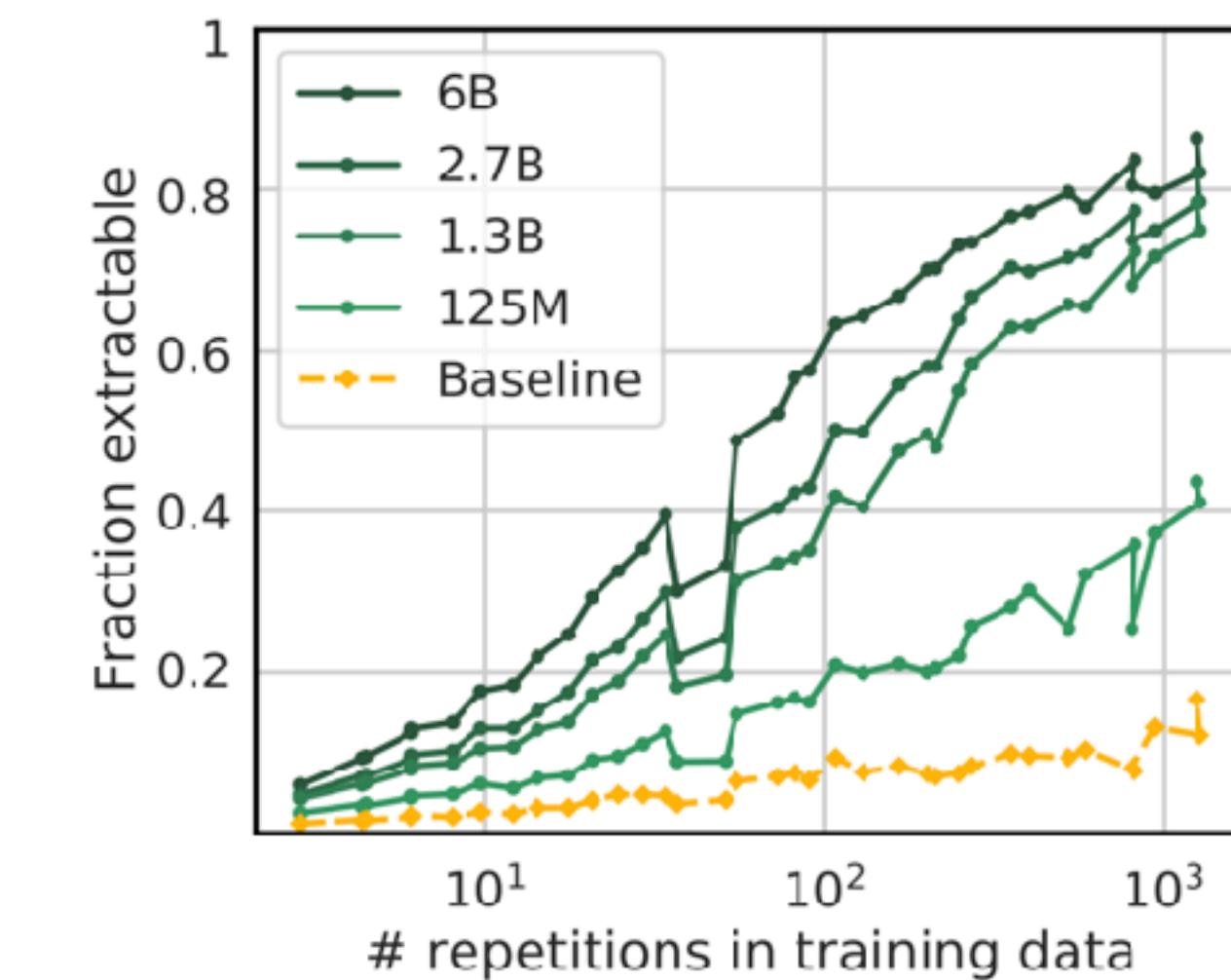
Dawn Song<sup>3</sup>

Úlfar Erlingsson<sup>7</sup>

Alina Oprea<sup>4</sup>

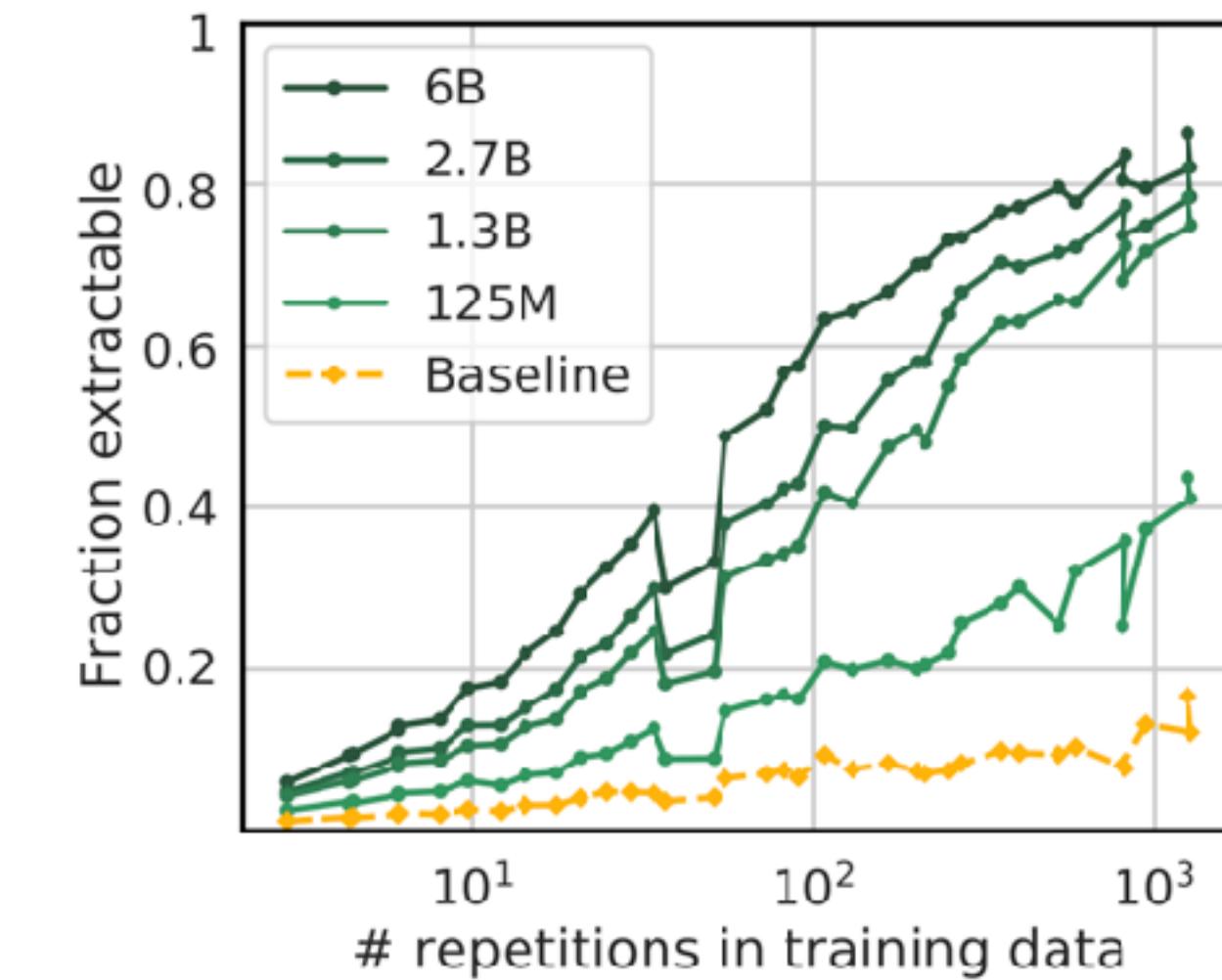
Colin Raffel<sup>1</sup>

<sup>1</sup>*Google* <sup>2</sup>*Stanford* <sup>3</sup>*UC Berkeley* <sup>4</sup>*Northeastern University* <sup>5</sup>*OpenAI* <sup>6</sup>*Harvard* <sup>7</sup>*Apple*



# Memorized Sequences

- This is a strict metric
- Can we relax it?

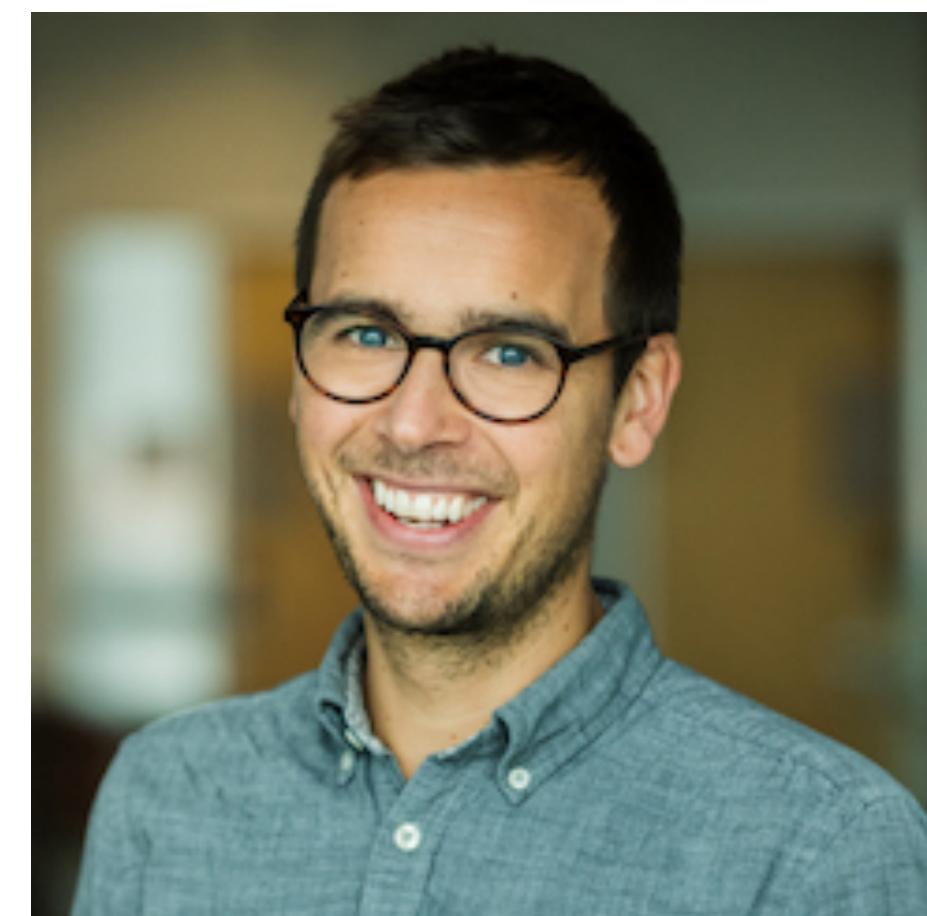
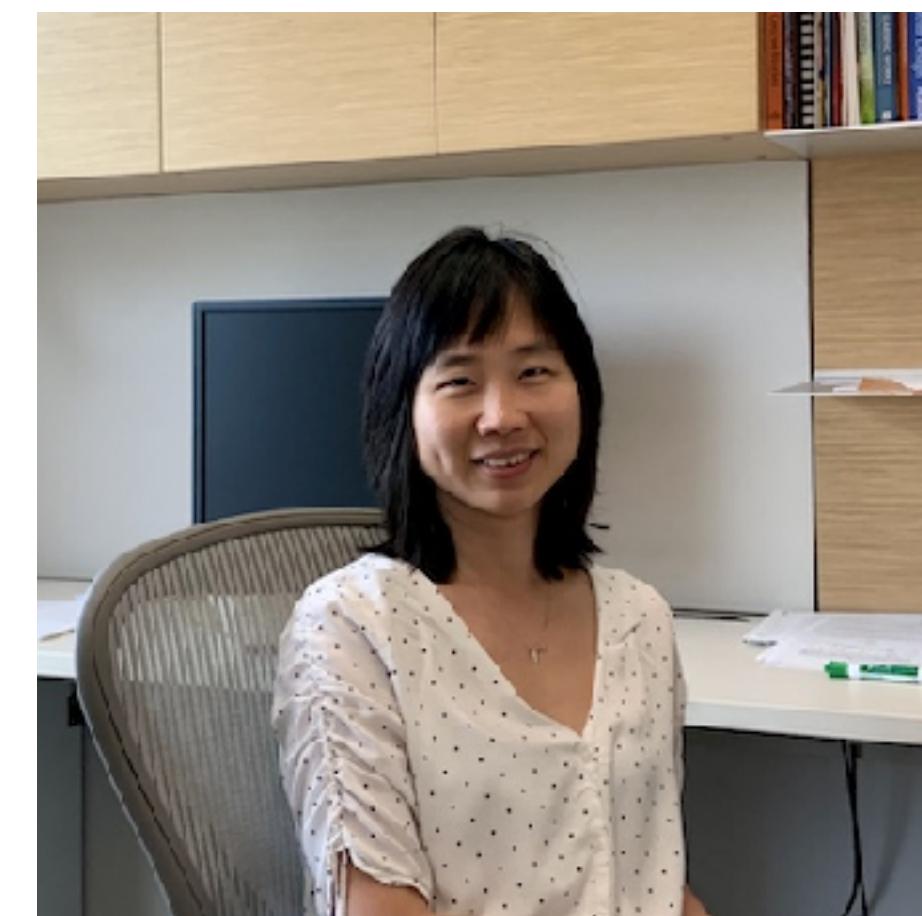


*Estimated Memorization (empirical lower) Bounds: 0.1%, 1%*

# Detection and Measurement of Syntactic Templates in Generated Text

Chantal Shaib, Yanai Elazar, Junyi Jesse Li, Byron C. Wallace

*EMNLP 2024*



# Memorizing Templates

- Instead of looking at raw sequences

The Last Black Man in San Francisco is a poignant, beautifully shot film [...] creates a unique and intense viewing experience.

# Memorizing Templates

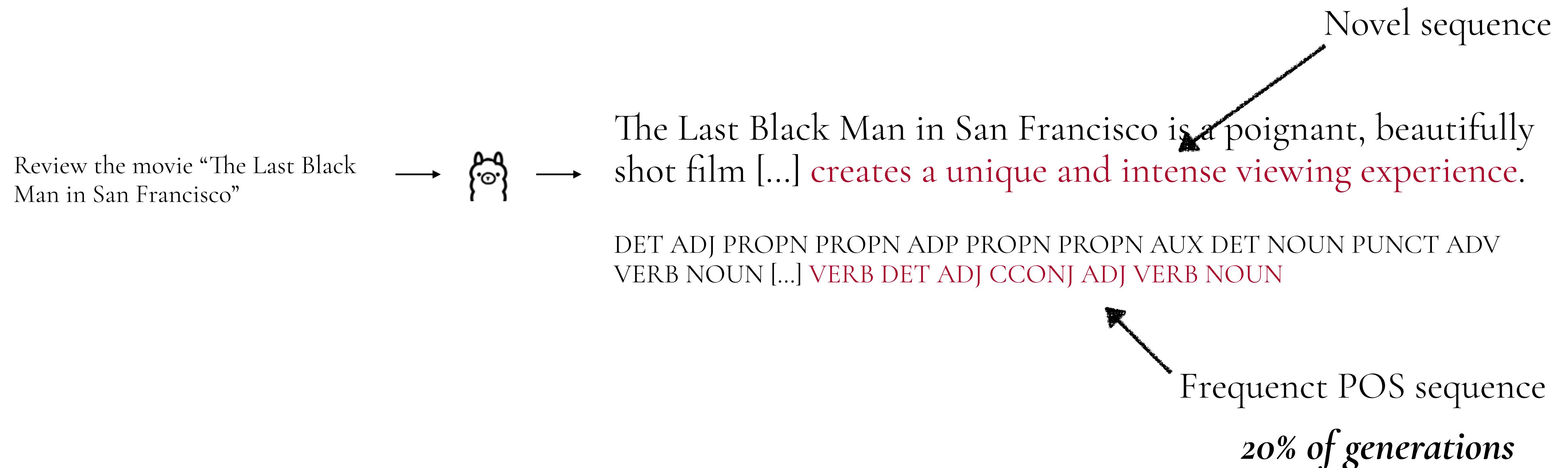
- Instead of looking at raw sequences
- We inspect their a linguistic abstraction

The Last Black Man in San Francisco is a poignant, beautifully shot film [...] creates a unique and intense viewing experience.

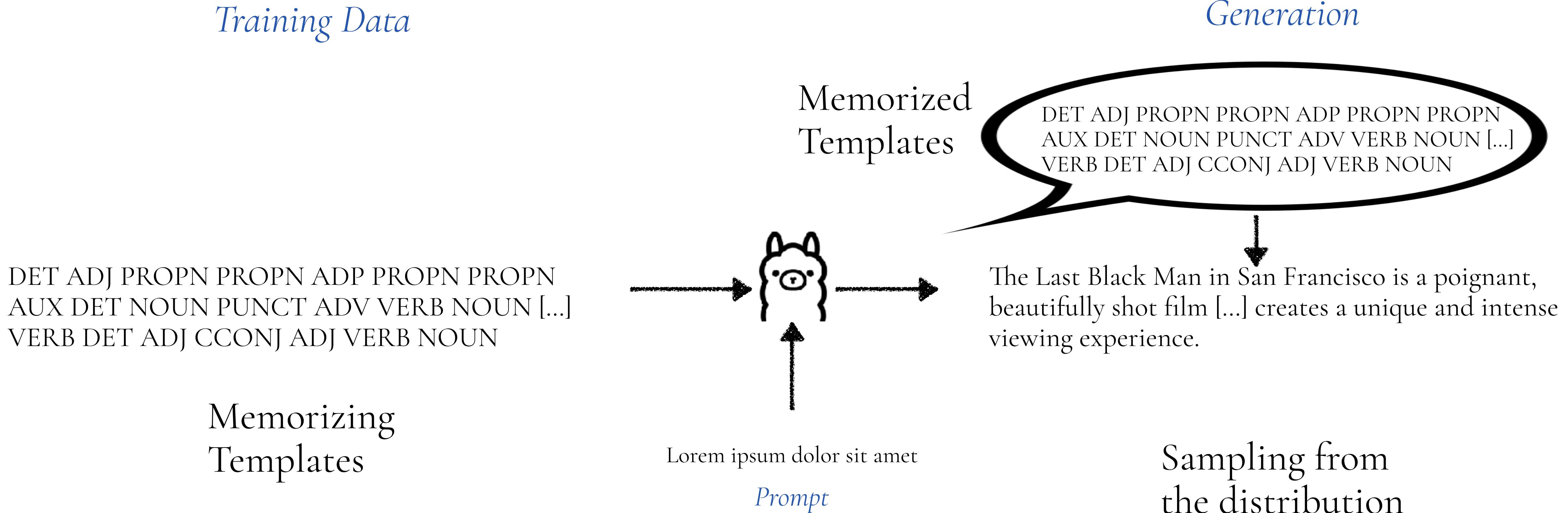
DET ADJ PROPN PROPN ADP PROPN PROPN AUX DET NOUN PUNCT ADV  
VERB NOUN [...] VERB DET ADJ CCONJ ADJ VERB NOUN

# Memorizing Templates

- Instead of looking at raw sequences
- We inspect their a linguistic abstraction



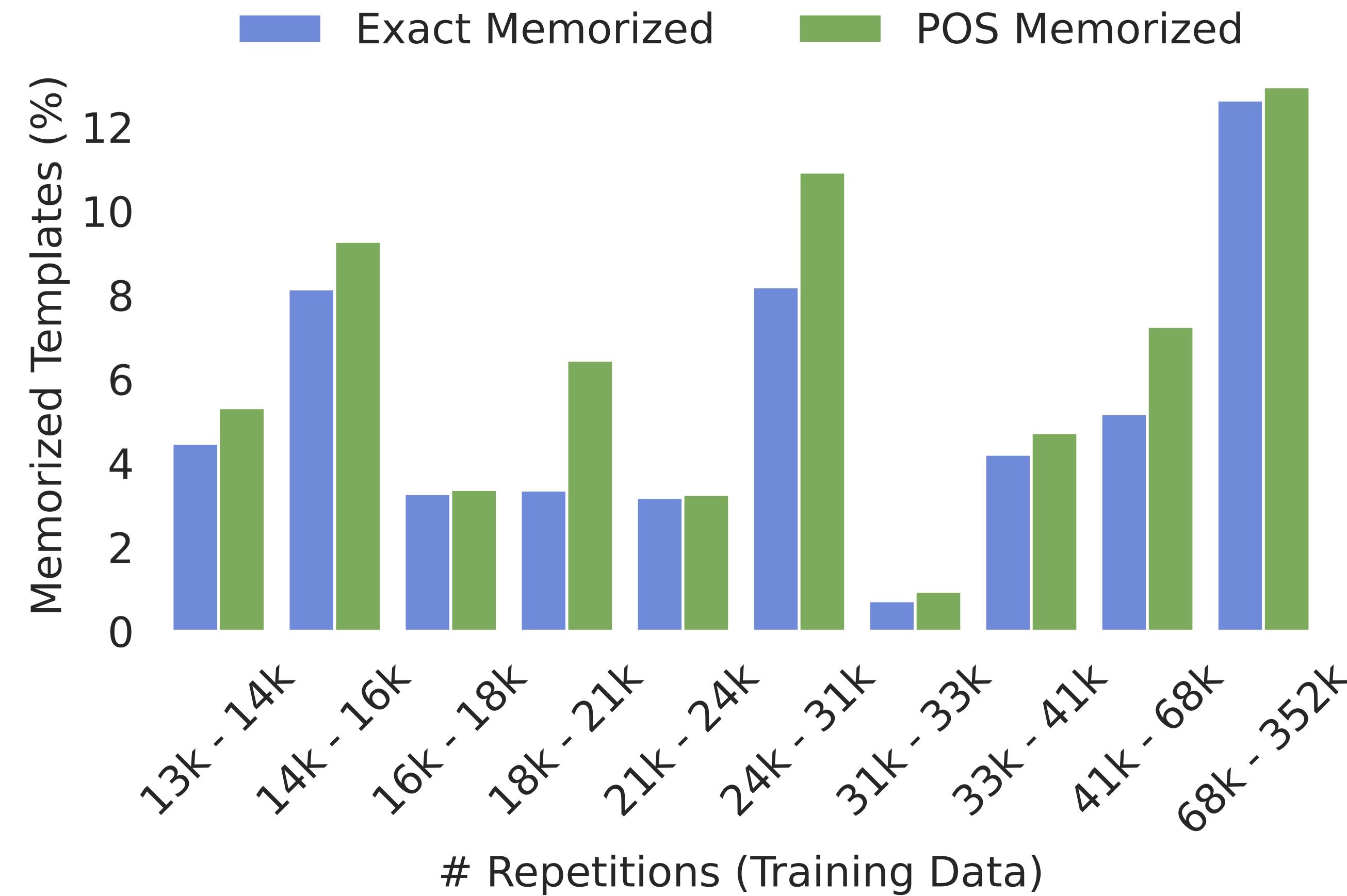
# Memorizing Templates (like HMMs)



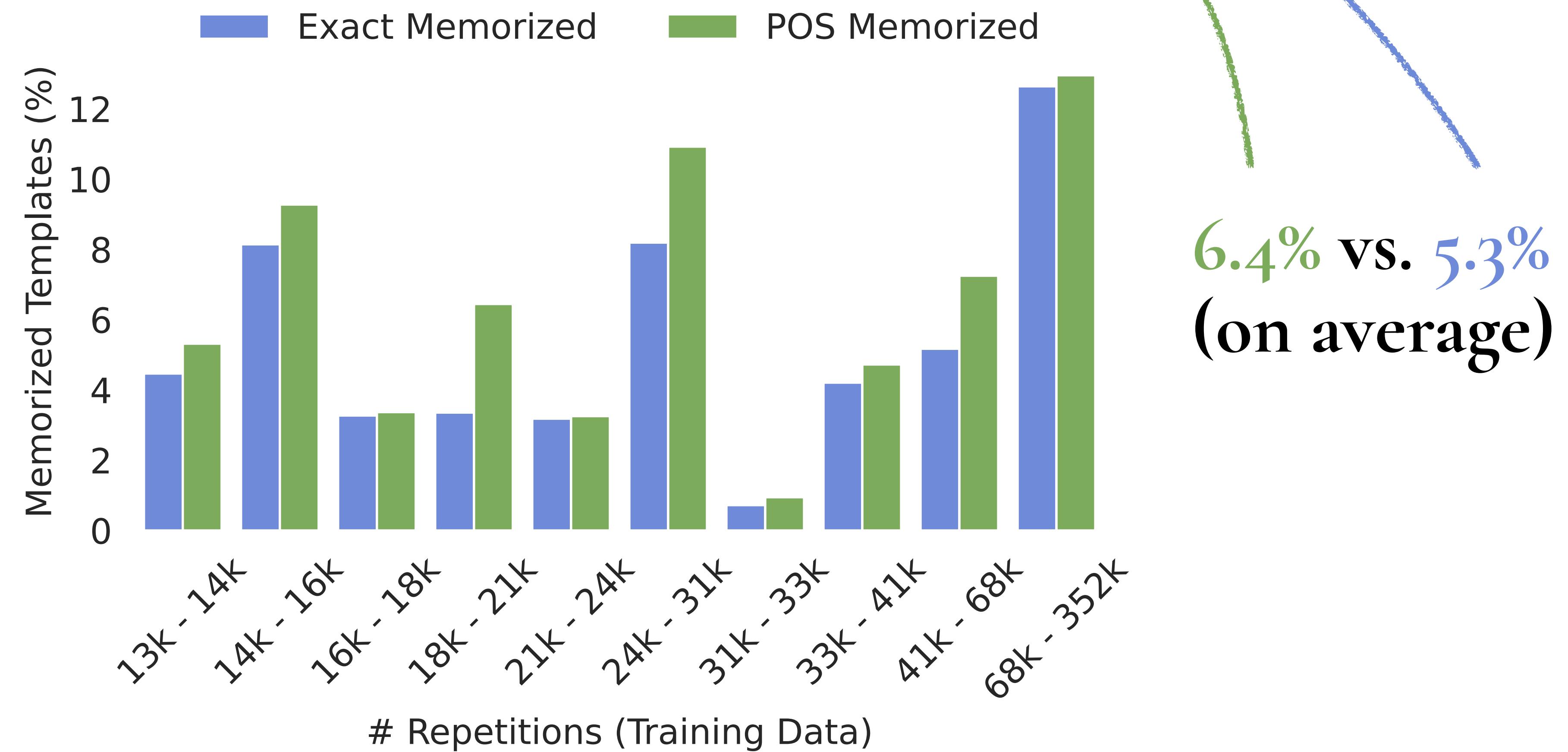
# Results



# Results



# Results



# Summary

- Generated texts is largely based on “templates”...
- that originate from the training data
- In turn, it allows us to measure memorized templates

# How Much it Happens?

- Memorization goes beyond exact match
- It can be seen as template memorization that the tokens are sampled from its distribution (like HMMs)
- We should extend our definitions of memorization beyond exact match

How much it

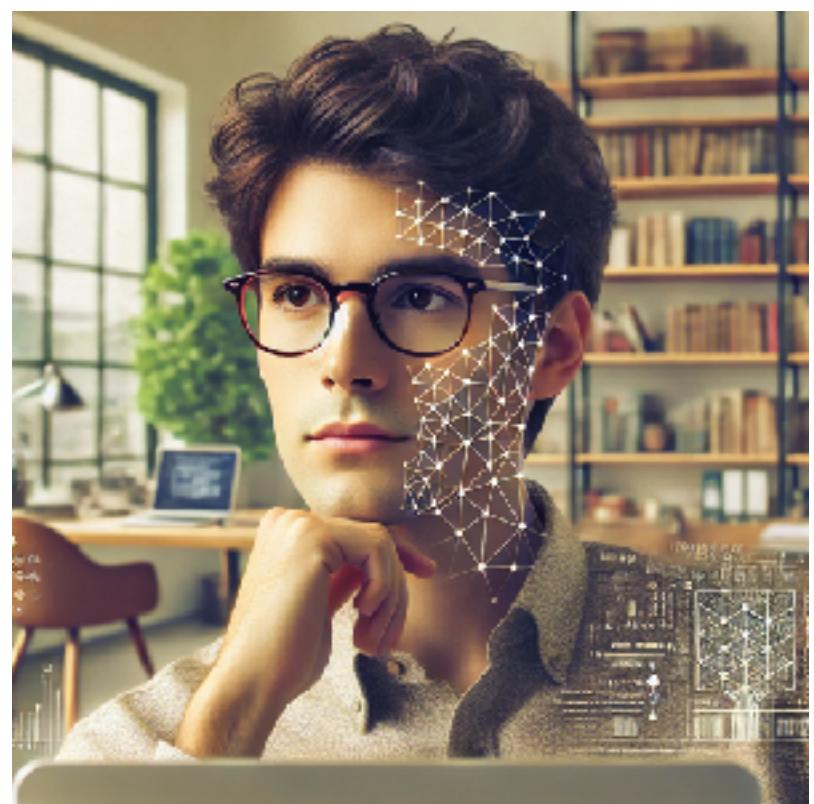


# When Does Distributional Memorization Happen?

# Imitation



Leonardo DiCaprio



Yanai Elazar



# Imitation

**Spot the difference**

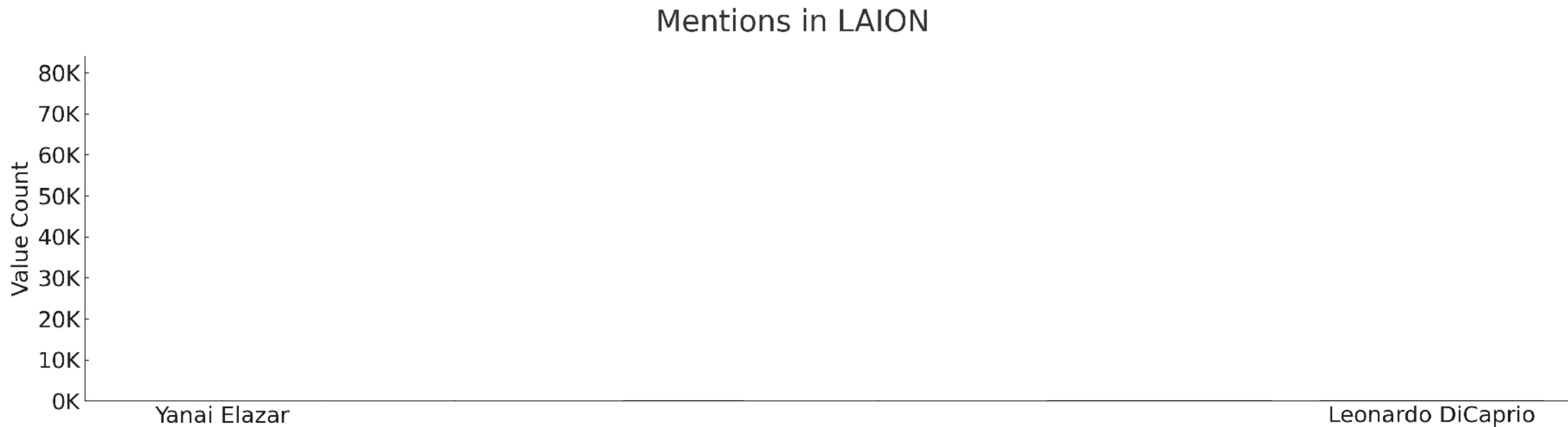
Leonardo DiCaprio



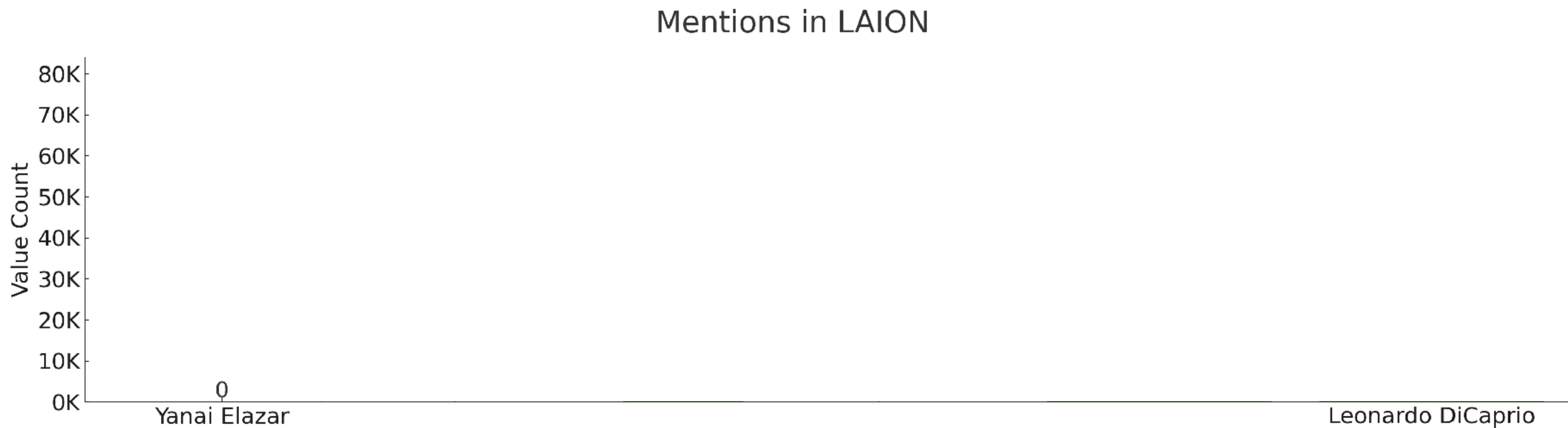
Yanai Elazar



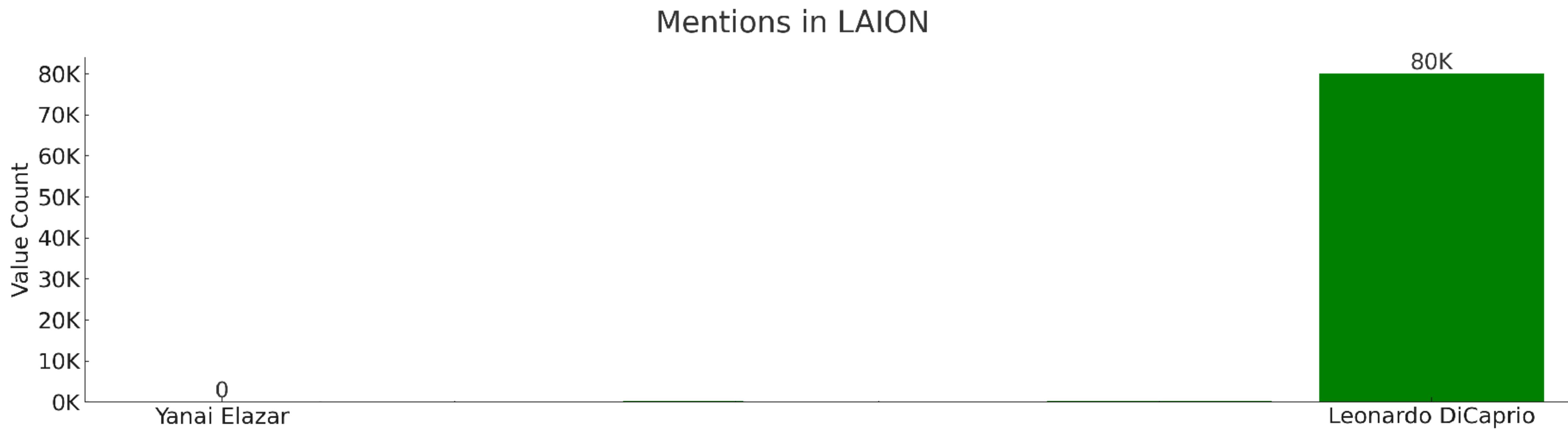
# Imitation Threshold?



# Imitation Threshold?



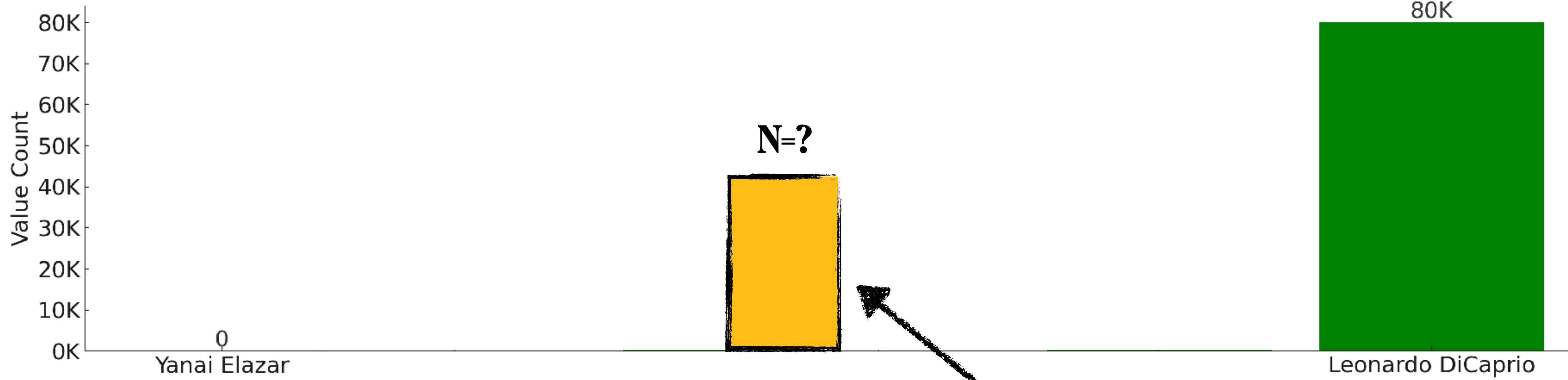
# Imitation Threshold?



# Imitation Threshold?



Mentions in LAION



Imitation Threshold?

# Imitation - Why Should You Care?

- Copyrights

# Imitation - Why Should You Care?

- Copyrights

**VentureBeat**

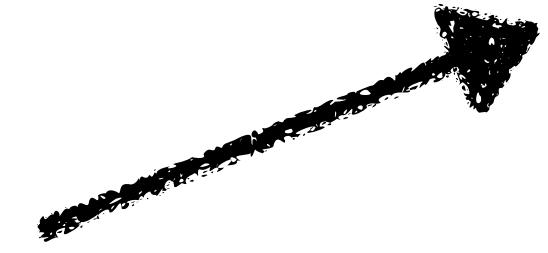
The copyright case against AI art generators just got stronger with more artists and evidence



Credit: VentureBeat made with OpenAI DALL-E 3 via ChatGPT

# Imitation - Why Should You Care?

- Copyrights
- Privacy



*Celebrity*

Leonardo DiCaprio



*Private individual*

Yanai Elazar



# Finding the Imitation Threshold

HOW MANY VAN GOGHS DOES IT TAKE TO VAN  
GOGH? FINDING THE IMITATION THRESHOLD

**Sahil Verma<sup>1</sup> Royi Rassin<sup>2</sup> Arnav Das<sup>\*1</sup> Gantavya Bhatt<sup>\*1</sup> Preethi Seshadri<sup>\*3</sup>**  
**Chirag Shah<sup>1</sup> Jeff Bilmes<sup>1</sup> Hannaneh Hajishirzi<sup>1,4</sup> Yanai Elazar<sup>1,4</sup>**

<sup>1</sup>*University of Washington, Seattle*    <sup>2</sup>*Bar-Ilan University*    <sup>3</sup>*University of California, Irvine*

<sup>4</sup>*Allen Institute of AI*



# Question Formulation

LAION-5B'



Count: 100



Would the model imitate a concept (e.g., *Leo*) if it was trained on  $X$  of his images instead?

LAION-5B



Count: 80K

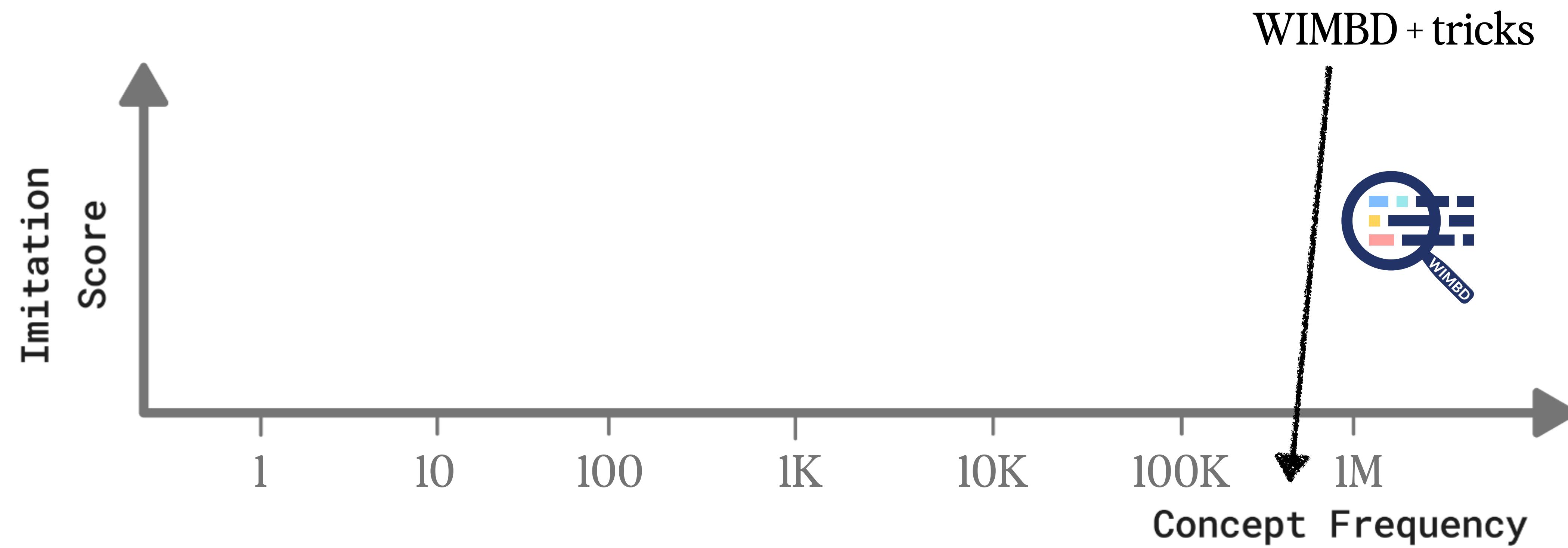
$$P(\text{Imitation} \mid \text{do}(\text{count}(\text{Leo}) = x))$$



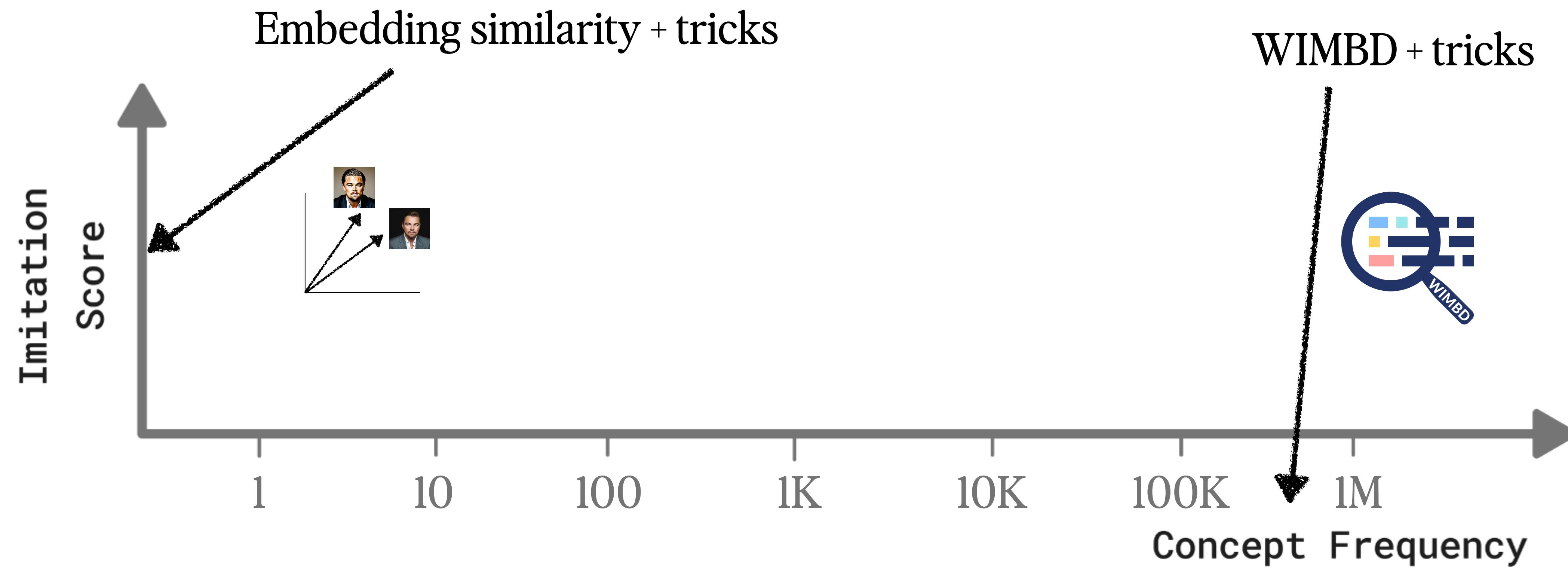
# Solutions

## I. Counterfactual model

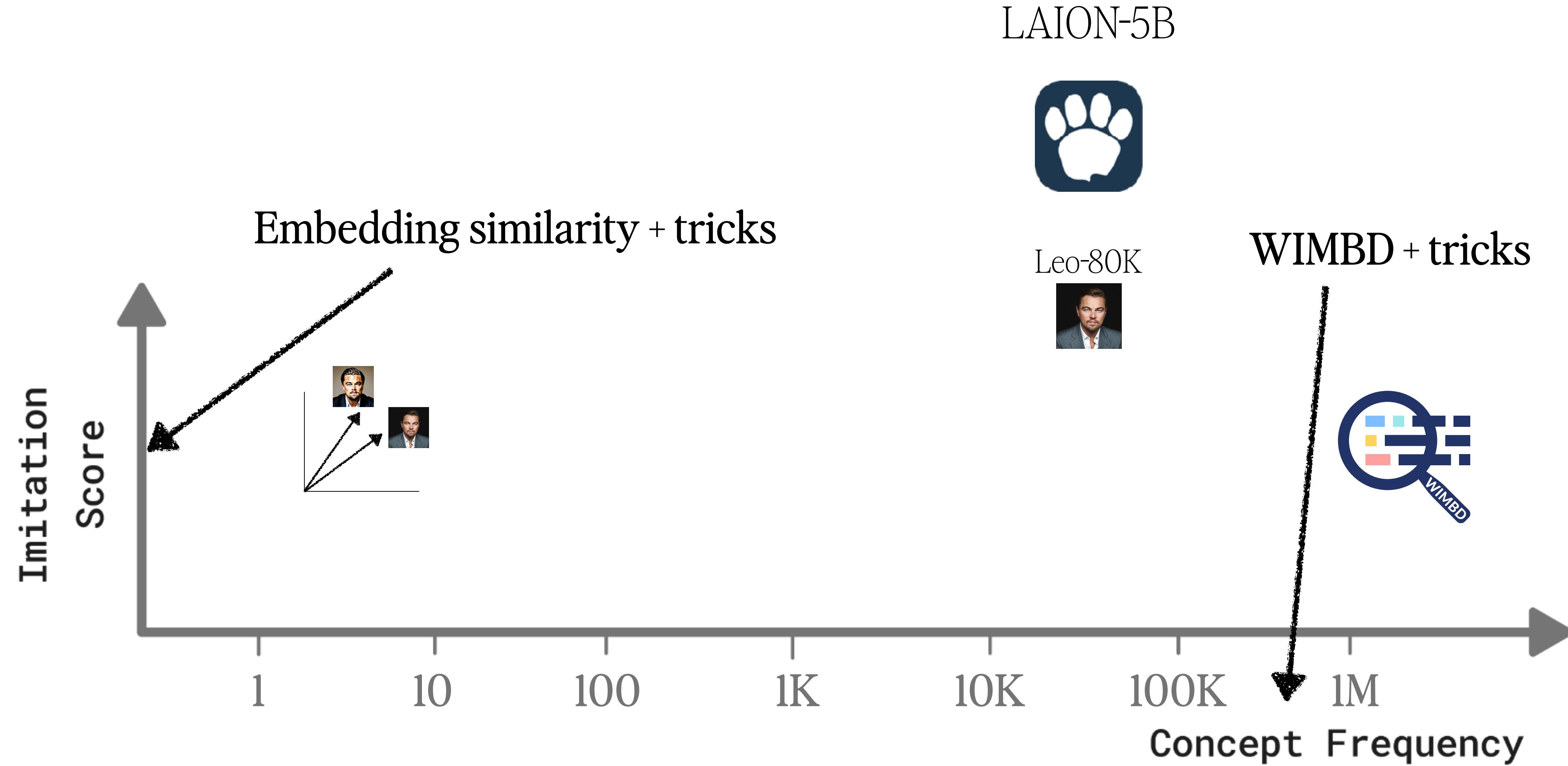
# Solutions



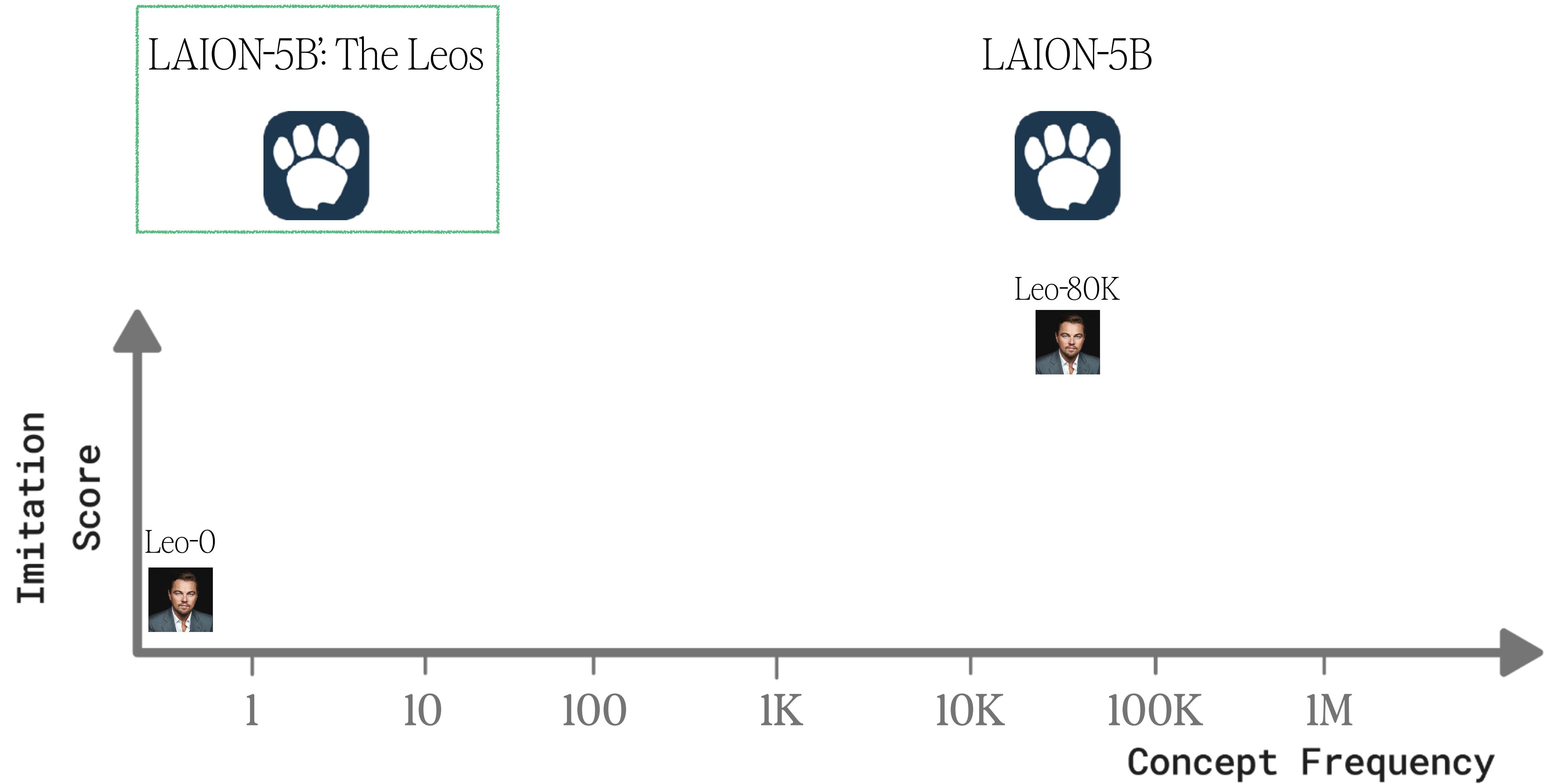
# Solutions



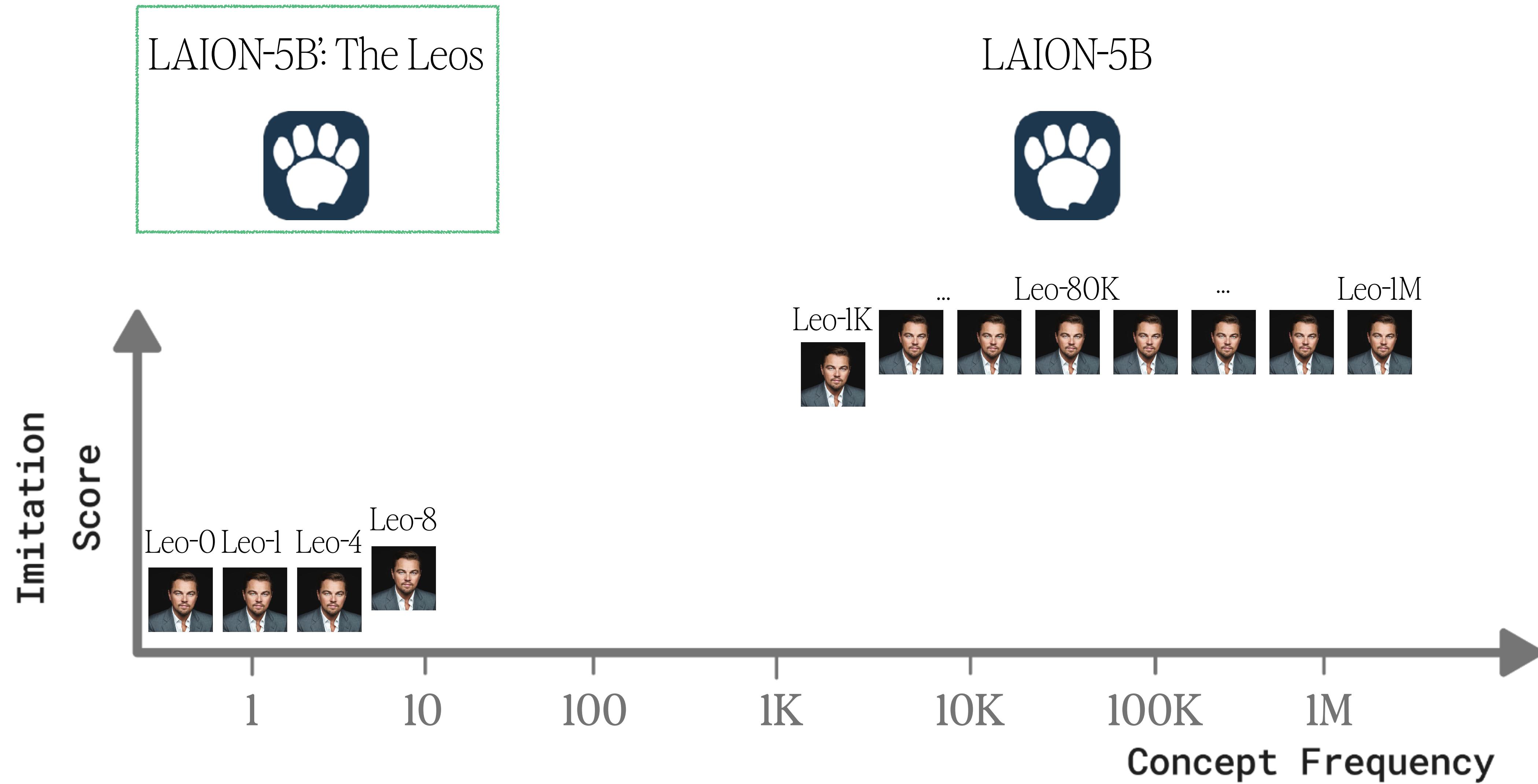
# Solutions



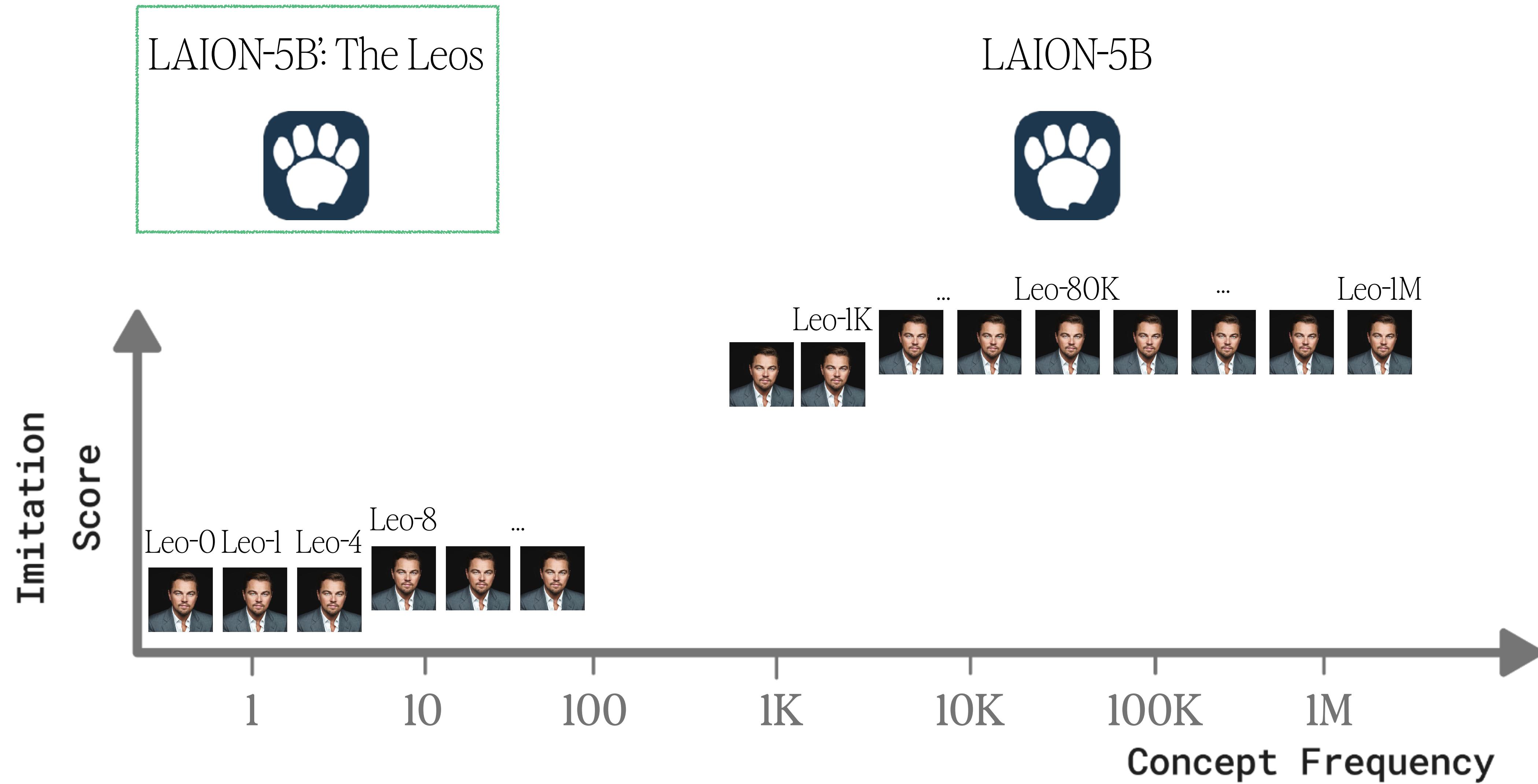
# Solution #I



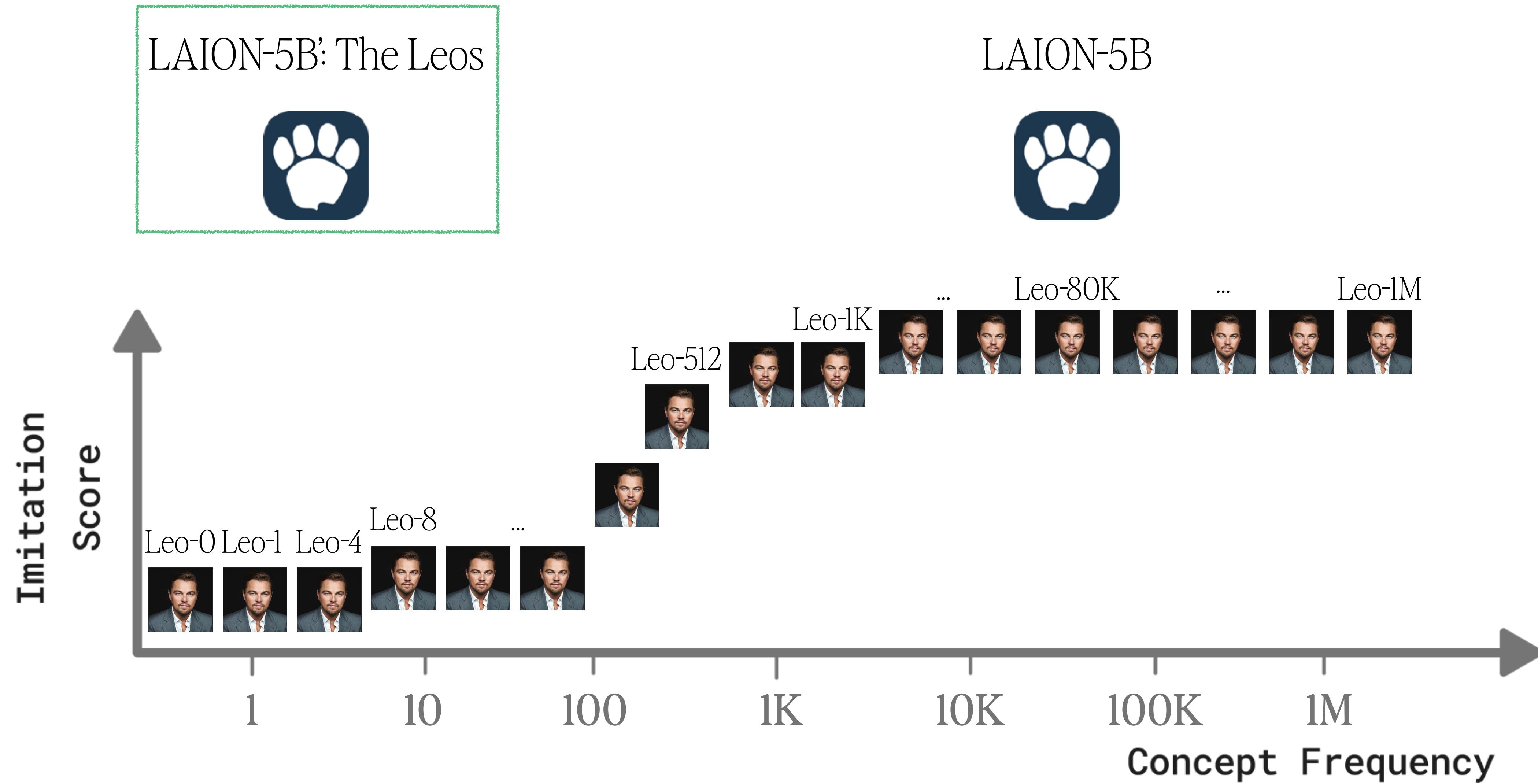
# Solution #I



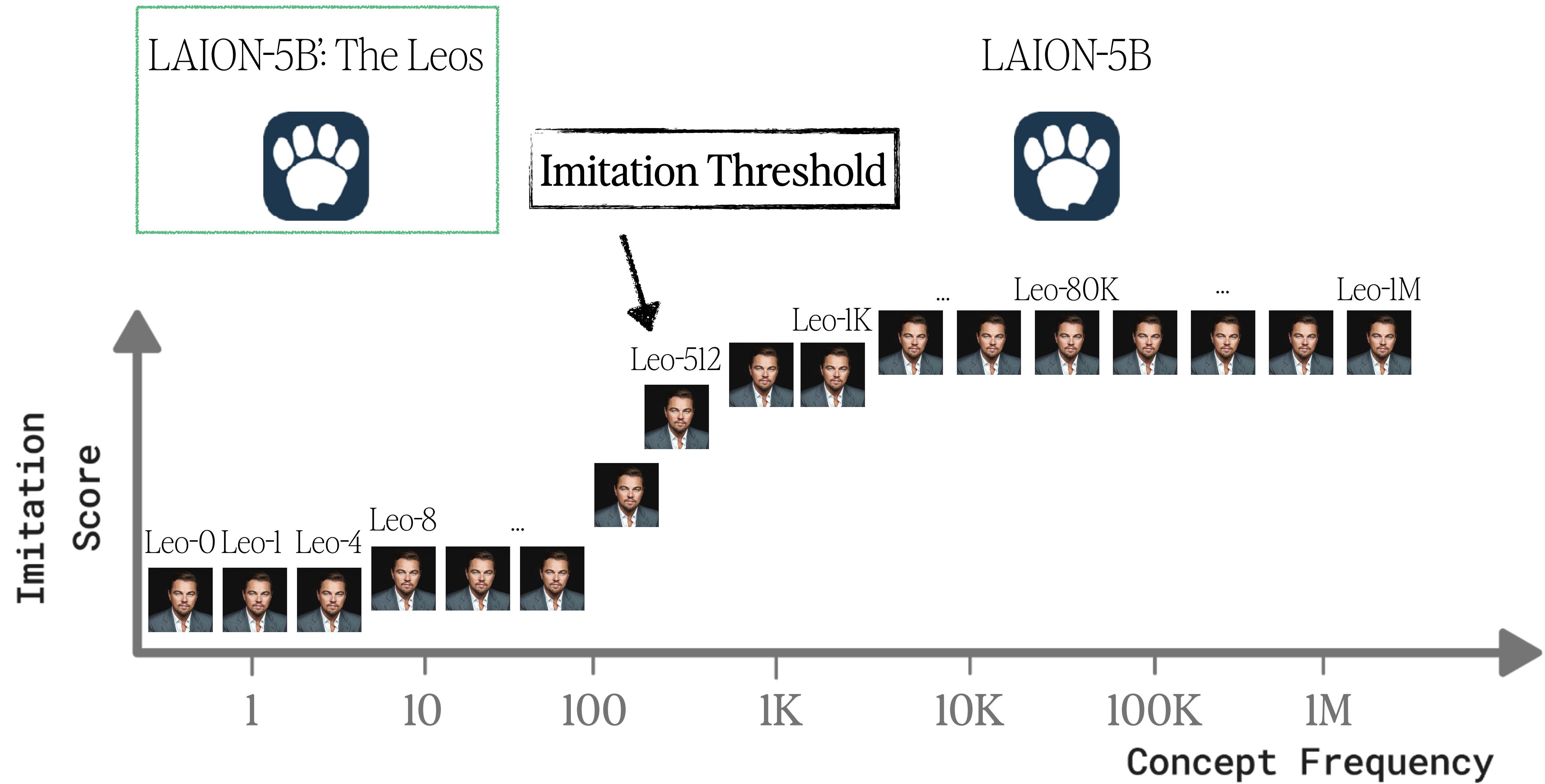
# Solution #I



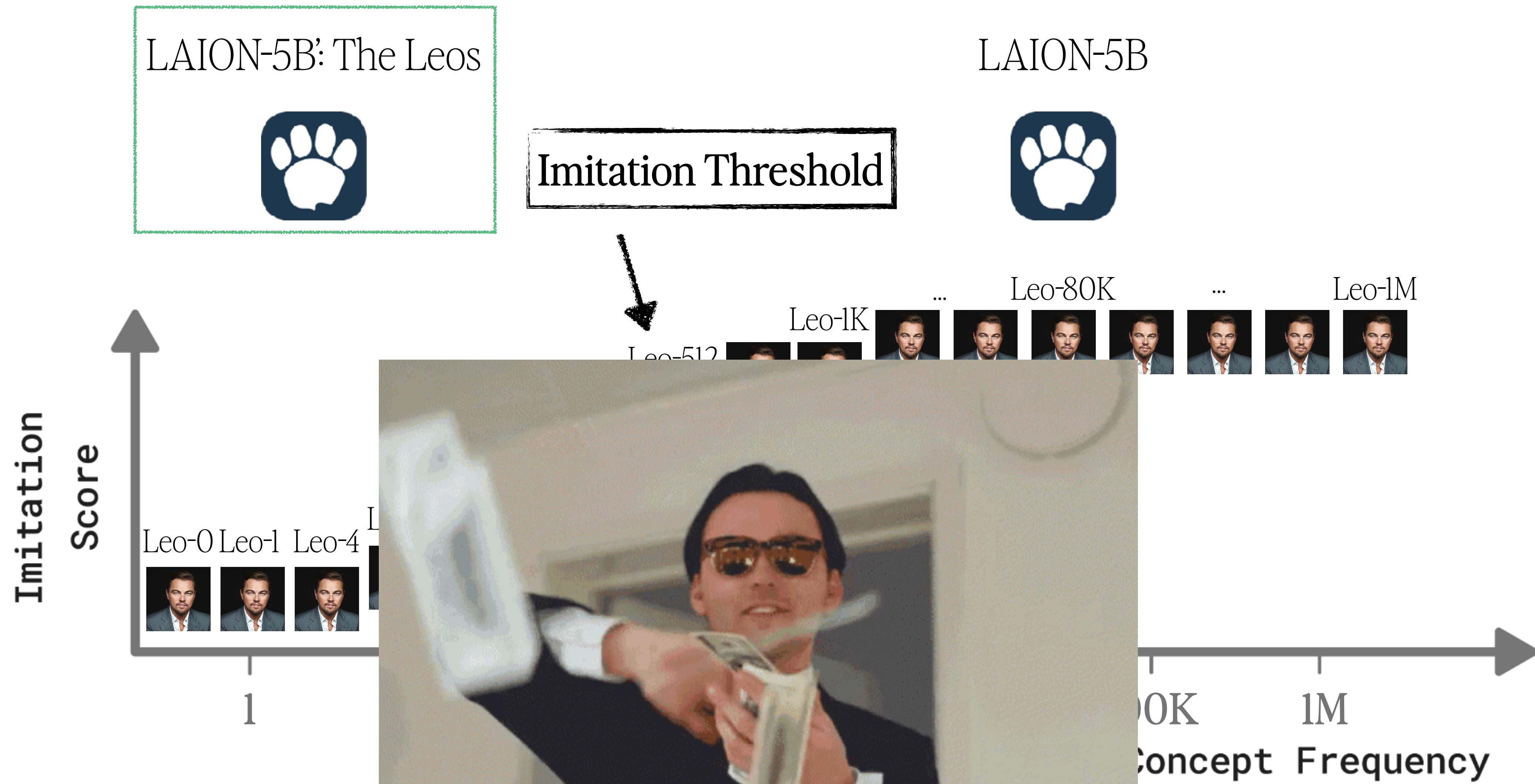
# Solution #I



# Solution #I



# Solution #I



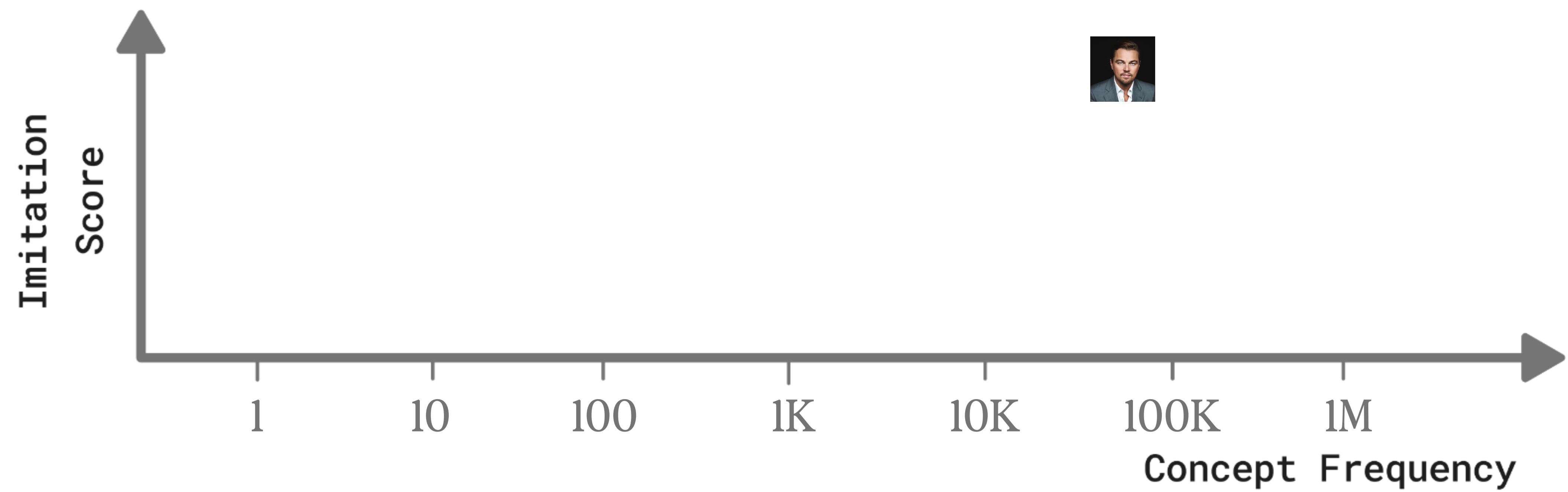
# Solutions

- I. Counterfactual model 

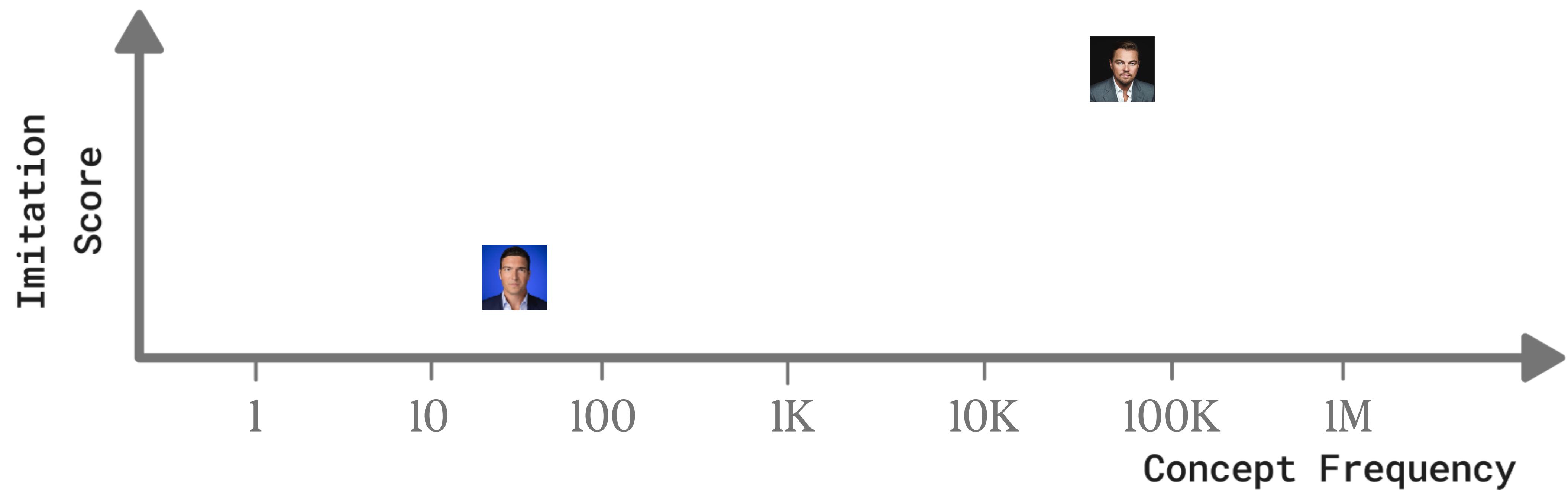
# Solutions

1. Counterfactual model 
2. Observational approach

# Solution #2



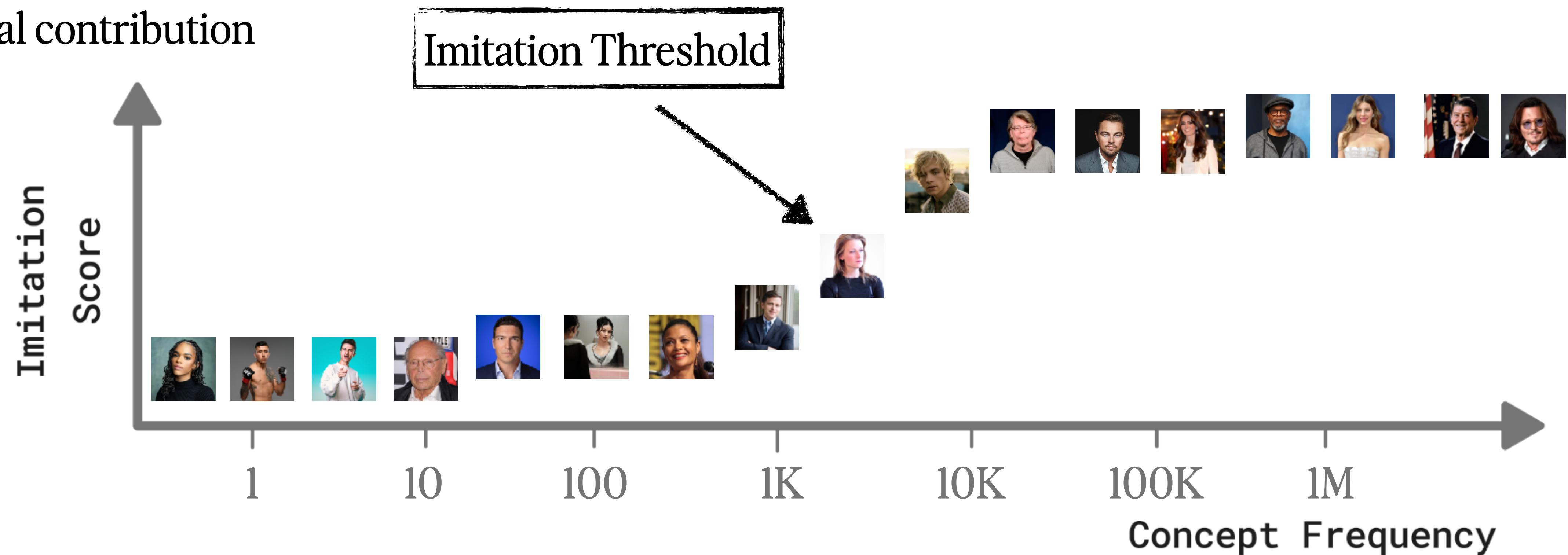
# Solution #2



# Solution #2

Using some assumptions:

- Distribution invariance
- Lack of confounders
- Equal contribution



# Setup

*2 domains x 2 datasets*

---

Human Faces 

---

Celebrities      Politicians

---

Art Style 

---

Classical      Modern

---

---

# Setup

3 pretraining datasets

---

## Pretraining Dataset

---

LAION-400M

LAION2B

LAION-5B

---

Human Faces 

---

Celebrities

Politicians

Art Style 

---

Classical

Modern

# Setup

4 models

Pretraining Dataset	Model	Human Faces		Art Style	
		Celebrities	Politicians	Classical	Modern
LAION-400M	LD				
LAION2B	SD1.1				
	SD1.5				
LAION-5B	SD2.1				

# Results

Pretraining Dataset	Model	Human Faces 🧑		Art Style 🖼	
		Celebrities	Politicians	Classical	Modern
LAION-400M	LD	648	309	219	282
LAION2B	SD1.1	364	234	112	198
	SD1.5	364	234	112	198
LAION-5B	SD2.1	527	369	185	241

# Results

Pretraining Dataset	Model	Human Faces 🧑		Art Style 🖼	
		Celebrities	Politicians	Classical	Modern
LAION-400M	LD	648	309	219	282
LAION2B	SD1.1	364	234	112	198
	SD1.5	364	234	112	198
LAION-5B	SD2.1	527	369	185	241

Imitation Threshold: 100-650 images

# The Imitation Threshold

- Memorizing distribution requires to observe enough training instance
- We estimate it to be a few hundreds images
- Implications on privacy, copyrights, etc.

# When does it Happen?

- Memorizing distribution happens
- But when is an under-explored area
- We have shown initial results for learning “simple” concepts
- But there so much more to learn



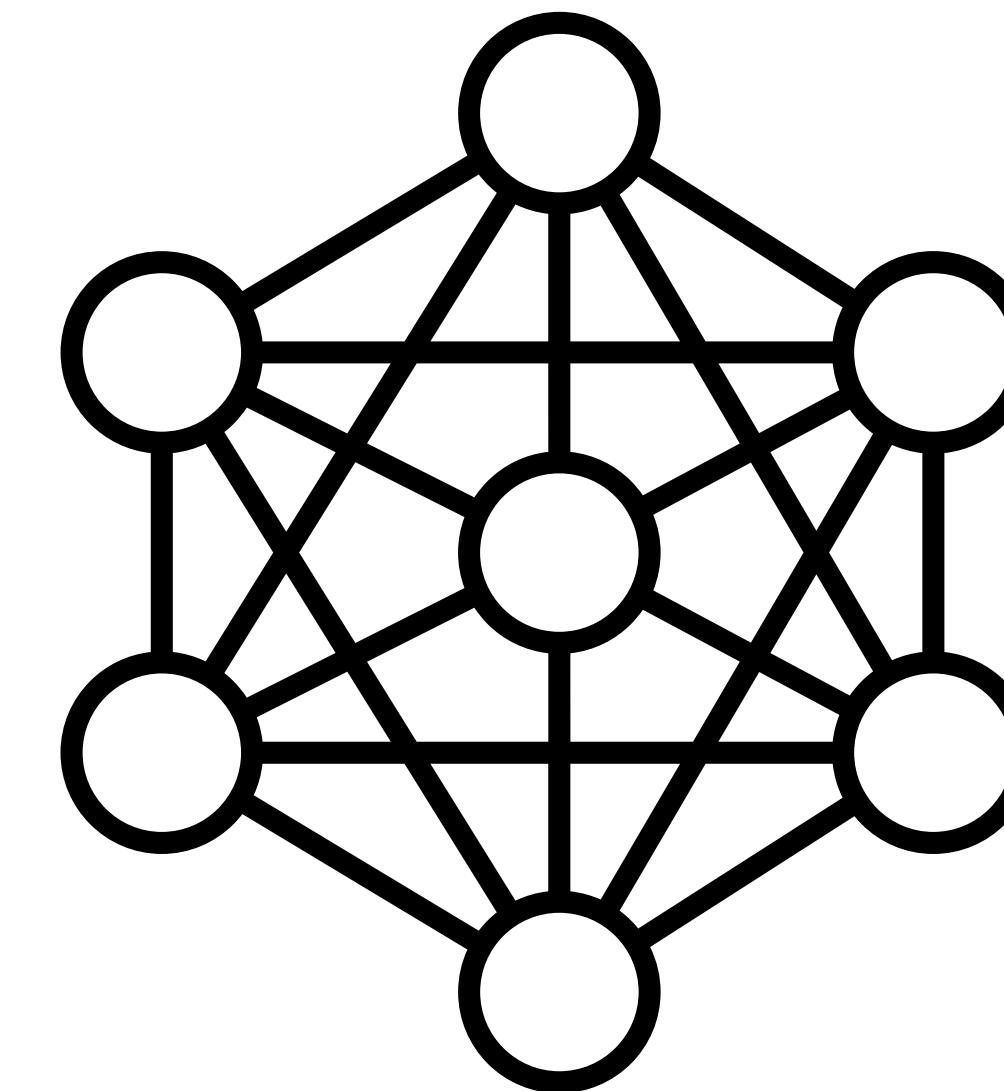
When does  
it happen?

# Measuring Distributional Memorization

# The Data Behind Memorization

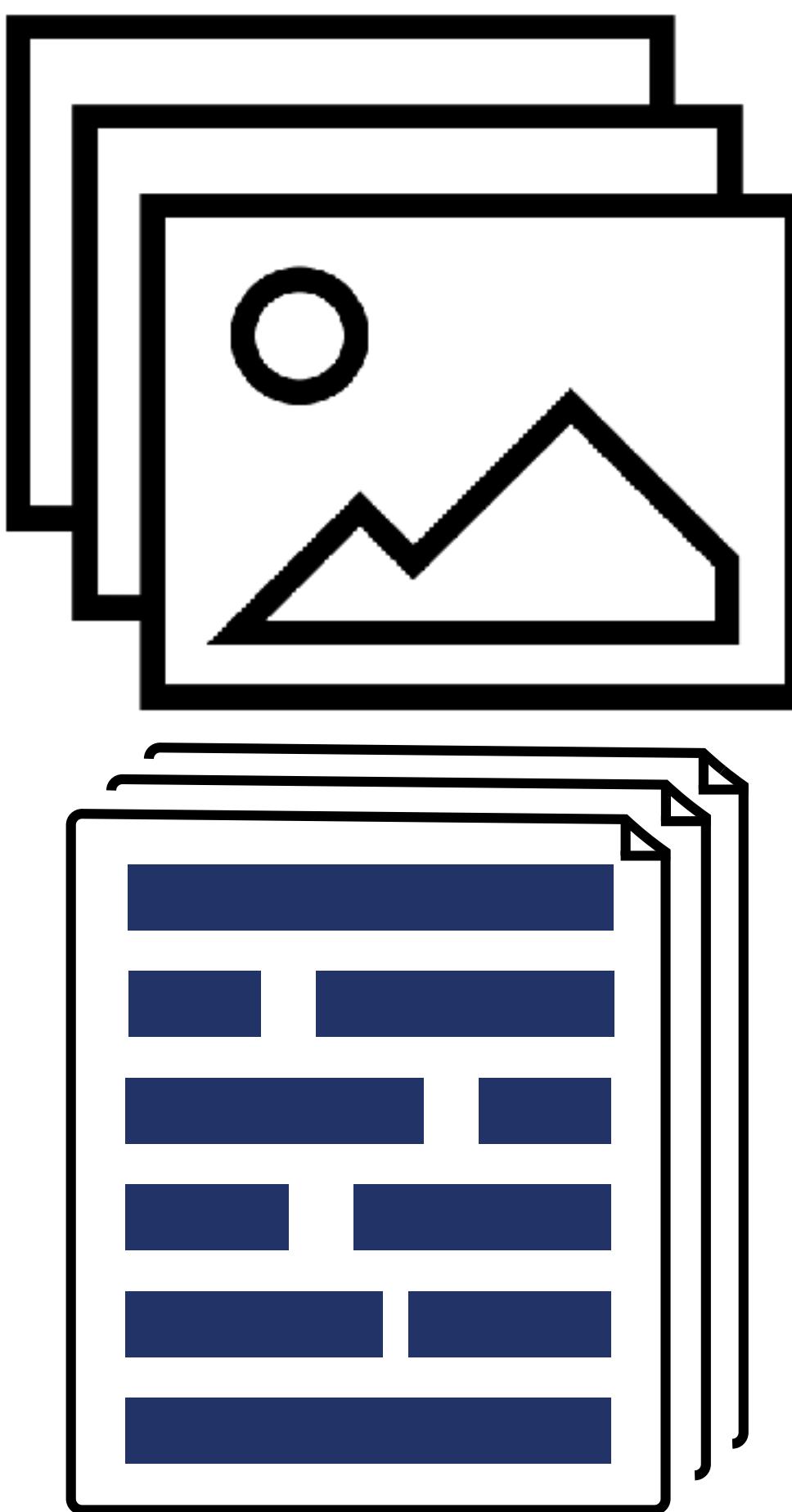


Dataset

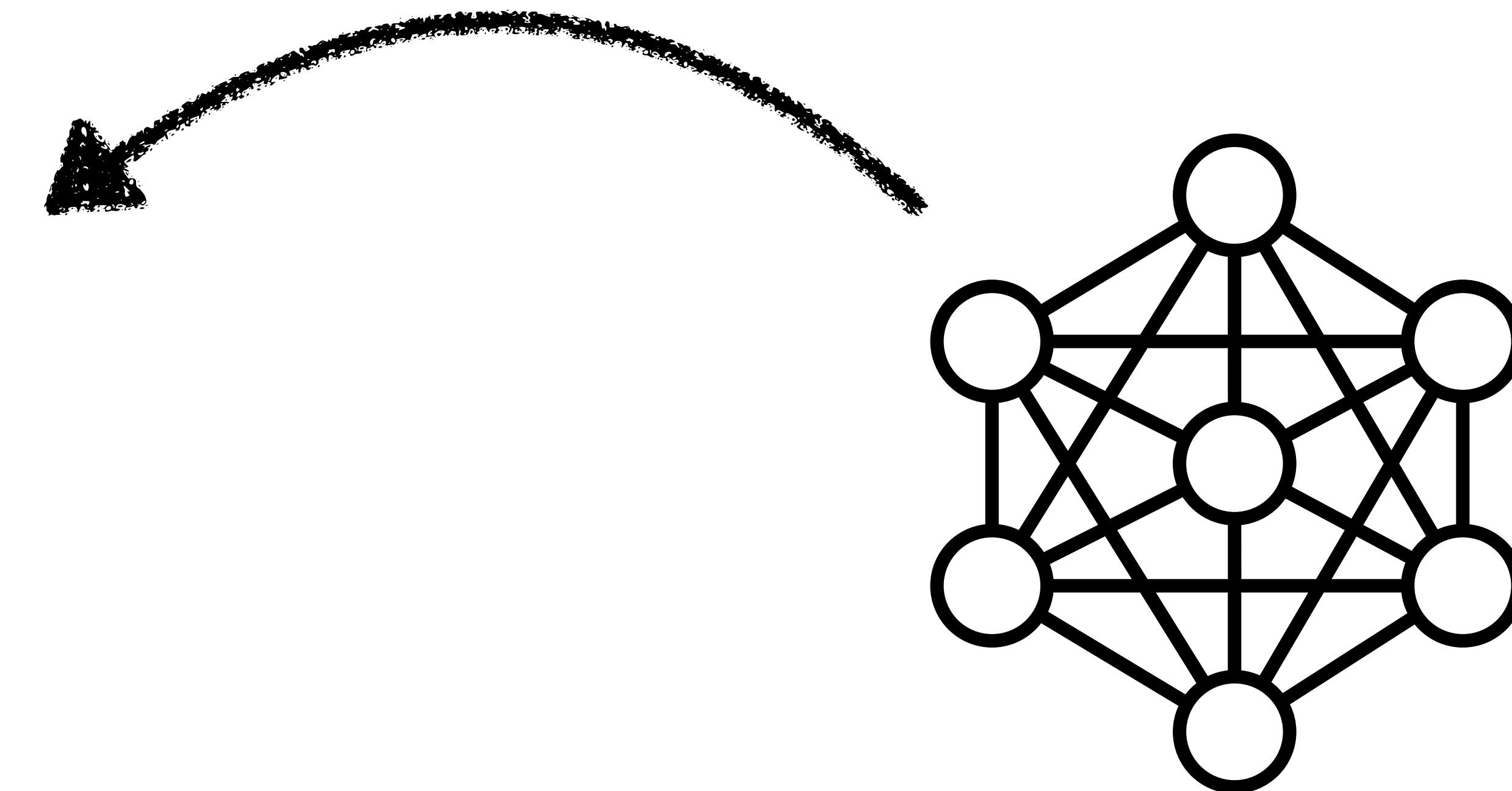


Model

# The Data Behind Memorization



Dataset



Model

# The Tools Behind Memorization

- Modern datasets are huge
- Naively querying is inefficient
- We need something better

# The Tools Behind Memorization

## Suffix-Array

### Deduplicating Training Data Makes Language Models Better

Katherine Lee<sup>\*†</sup> Daphne Ippolito<sup>\*†‡</sup> Andrew Nystrom<sup>†</sup> Chiyuan Zhang<sup>†</sup>  
Douglas Eck<sup>†</sup> Chris Callison-Burch<sup>‡</sup> Nicholas Carlini<sup>†</sup>

### Infini-gram: Scaling Unbounded $n$ -gram Language Models to a Trillion Tokens

Jiacheng Liu<sup>♡</sup> Sewon Min<sup>♡</sup>  
Luke Zettlemoyer<sup>♡</sup> Yejin Choi<sup>♡♦</sup> Hannaneh Hajishirzi<sup>♡♦</sup>  
♡Paul G. Allen School of Computer Science & Engineering, University of Washington  
♦Allen Institute for Artificial Intelligence liujc@cs.washington.edu

## Other Data Structures

### Evaluating $n$ -Gram Novelty of Language Models Using RUSTY-DAWG

William Merrill<sup>α,β</sup> Noah A. Smith<sup>γ,β</sup> Yanai Elazar<sup>β,γ</sup>  
αNew York University βAllen Institute for AI γUniversity of Washington  
willm@nyu.edu, noah@allenai.org, yanaiela@gmail.com

### INFINI-GRAM MINI: Exact n-gram Search at the Internet Scale with FM-Index

Hao Xu<sup>♡</sup> Jiacheng Liu<sup>♡♦</sup> Yejin Choi<sup>♡</sup> Noah A. Smith<sup>♡♦</sup> Hannaneh Hajishirzi<sup>♡♦</sup>  
♡Paul G. Allen School of Computer Science & Engineering, University of Washington  
♦Allen Institute for AI ♦Stanford University

### WHAT'S IN MY BIG DATA?



Yanai Elazar<sup>1,2</sup> Akshita Bhagia<sup>1</sup> Ian Magnusson<sup>1</sup> Abhilasha Ravichander<sup>1</sup>  
Dustin Schwenk<sup>1</sup> Alane Suhr<sup>3</sup> Pete Walsh<sup>1</sup> Dirk Groeneveld<sup>1</sup> Luca Soldaini<sup>1</sup>  
Sameer Singh<sup>4</sup> Hannaneh Hajishirzi<sup>1,2</sup> Noah A. Smith<sup>1,2</sup> Jesse Dodge<sup>1</sup>

<sup>1</sup>Allen Institute for AI

<sup>2</sup>Paul G. Allen School of Computer Science & Engineering, University of Washington

<sup>3</sup>University of California, Berkeley <sup>4</sup>University of California, Irvine

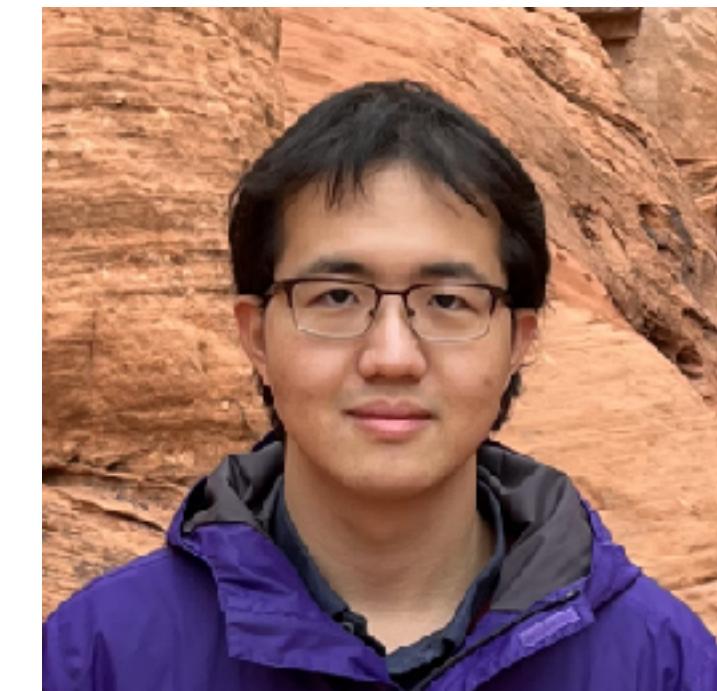
# Combining the Tools & Generation

## OLMOTRACE: Tracing Language Model Outputs Back to Trillions of Training Tokens

Jiacheng Liu<sup>ω</sup> Taylor Blanton<sup>α</sup> Yanai Elazar<sup>ω</sup> Sewon Min<sup>αβ</sup>

YenSung Chen<sup>α</sup> Arnavi Chheda-Kothary<sup>ω</sup> Huy Tran<sup>α</sup> Byron Bischoff<sup>α</sup> Eric Marsh<sup>α</sup>  
Michael Schmitz<sup>α</sup> Cassidy Trier<sup>α</sup> Aaron Sarnat<sup>α</sup> Jenna James<sup>α</sup> Jon Borchardt<sup>α</sup>  
Bailey Kuehl<sup>α</sup> Evie Cheng<sup>α</sup> Karen Farley<sup>α</sup> Sruthi Sreeram<sup>α</sup> Taira Anderson<sup>α</sup>  
David Albright<sup>α</sup> Carissa Schoenick<sup>α</sup> Luca Soldaini<sup>α</sup> Dirk Groeneveld<sup>α</sup>  
Rock Yuren Pang<sup>ω</sup>

Pang Wei Koh<sup>ω</sup> Noah A. Smith<sup>ω</sup> Sophie Lebrecht<sup>α</sup> Yejin Choi<sup>σ</sup>  
Hannaneh Hajishirzi<sup>ω</sup> Ali Farhadi<sup>ω</sup> Jesse Dodge<sup>α</sup>



<sup>α</sup>Allen Institute for AI <sup>ω</sup>University of Washington <sup>β</sup>UC Berkeley <sup>σ</sup>Stanford University

ACL best! 🎉🎉 demo 2025

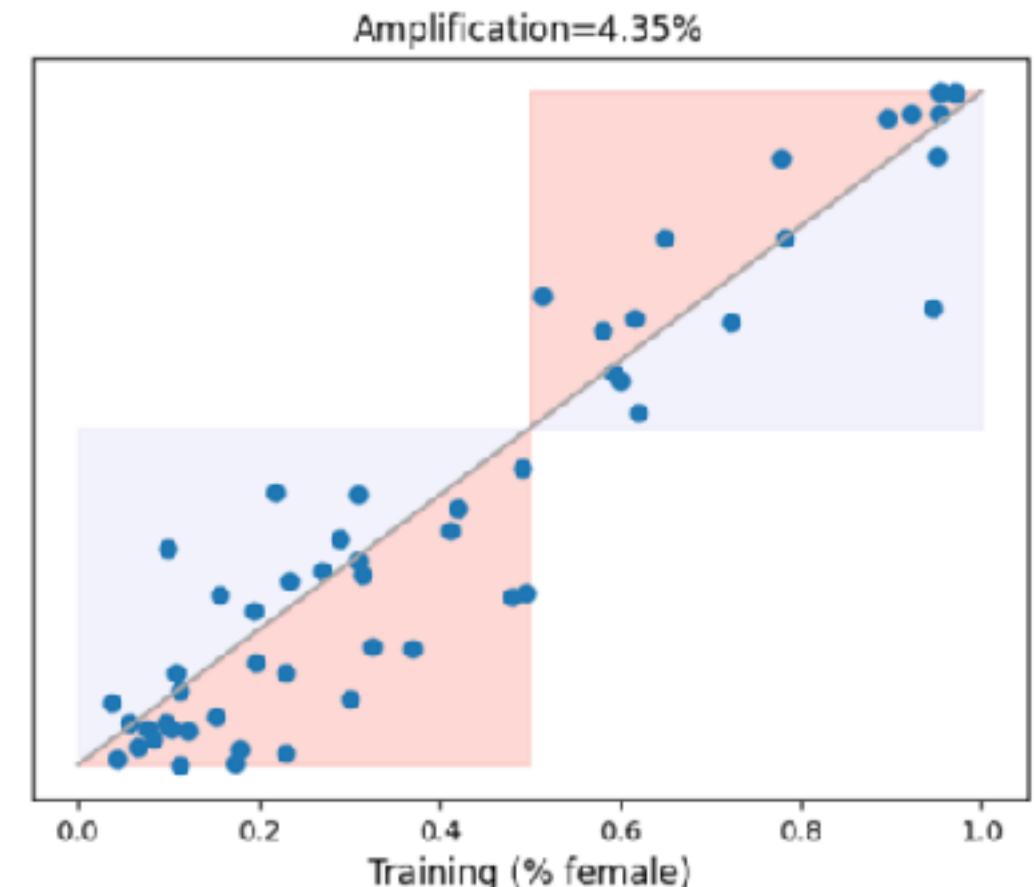
# The Tools Behind Memorization

- Most current focus on string matching
- We need to do better

# Memorizing Distributions

The Bias  
Amplification Paradox

What is?

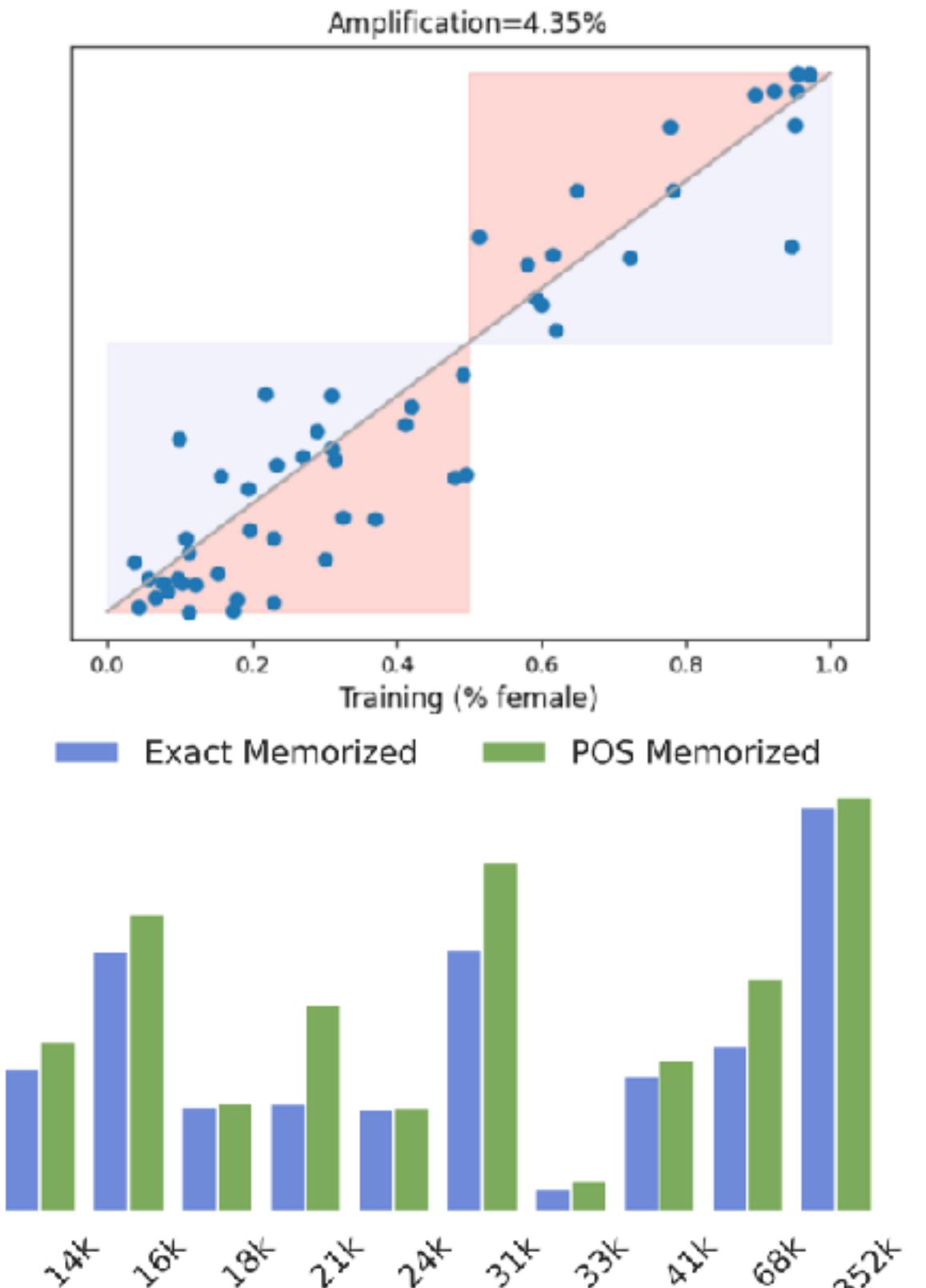
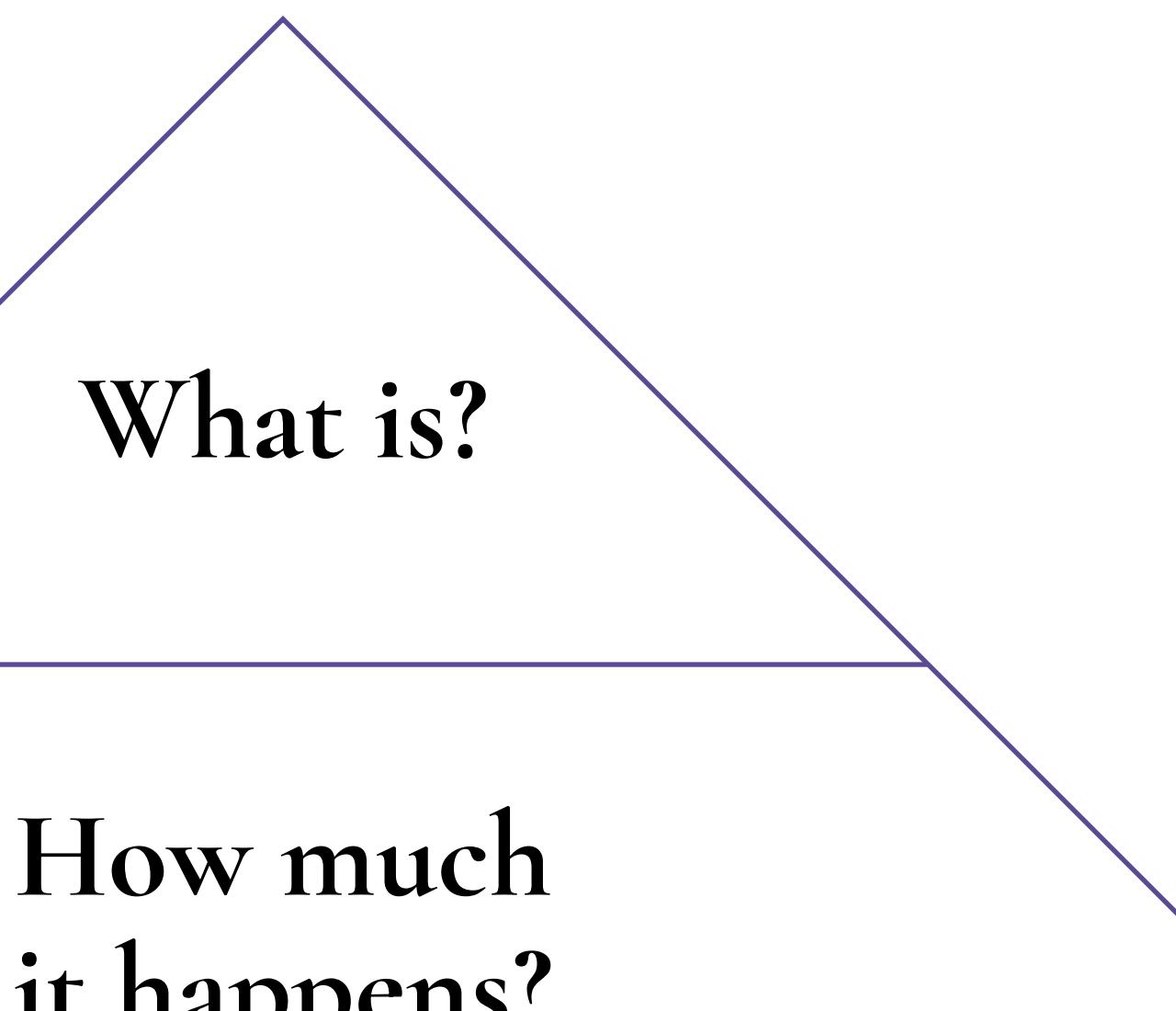


# Memorizing Distributions

The Bias  
Amplification Paradox

Syntactic Templates  
in texts

How much  
it happens?



# Memorizing Distributions

The Bias  
Amplification Paradox

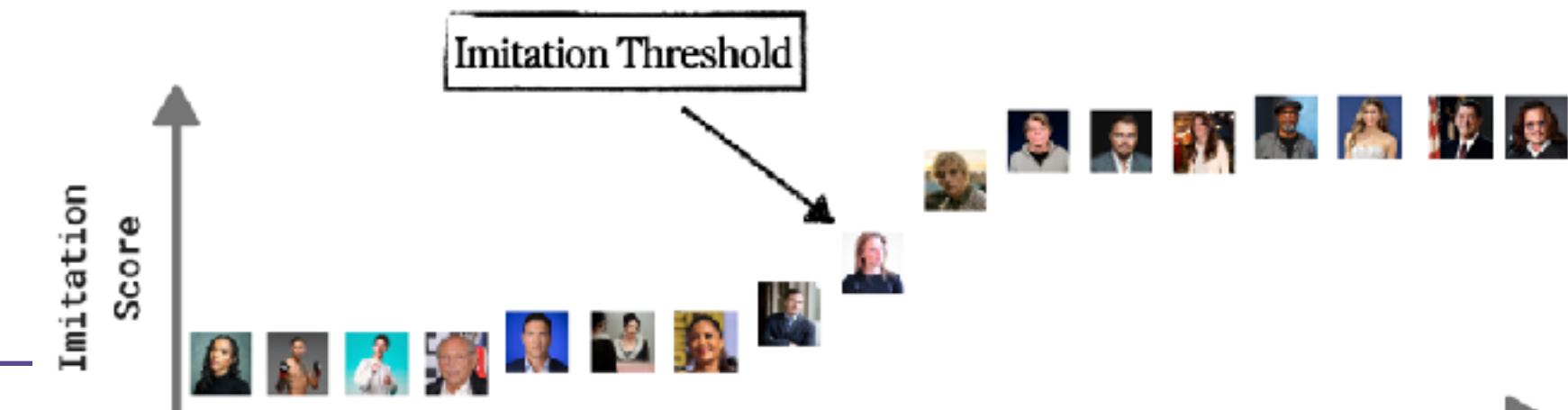
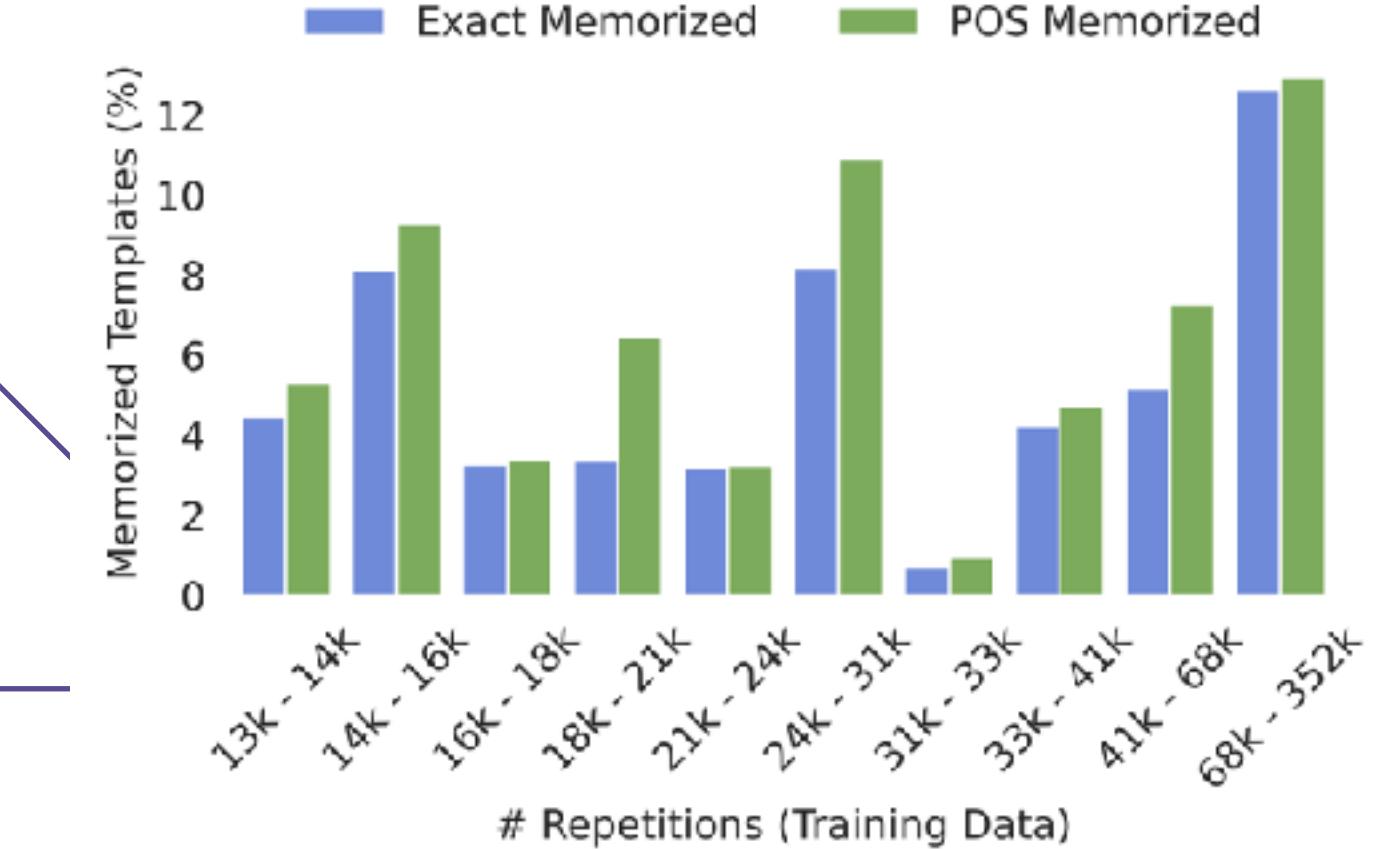
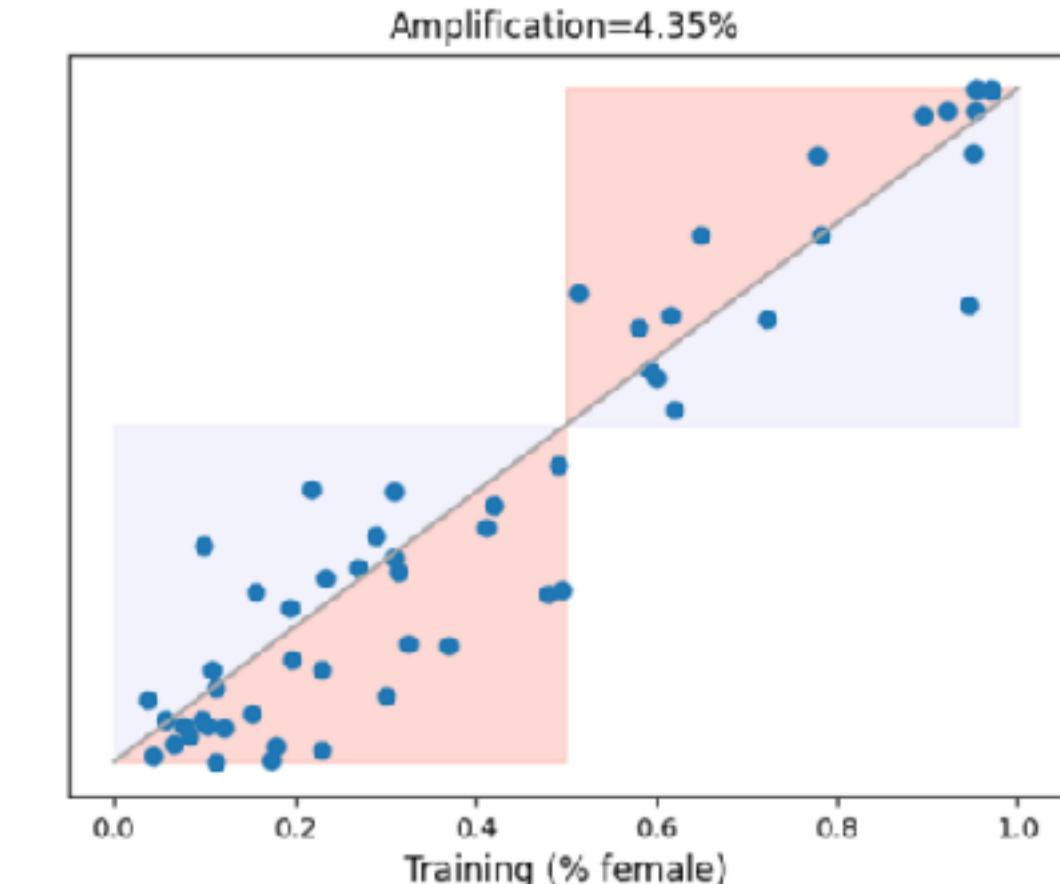
What is?

Syntactic Templates  
in texts

How much  
it happens?

Imitation  
Threshold

When does  
it happen?



# Memorizing Distributions

The Bias  
Amplification Paradox

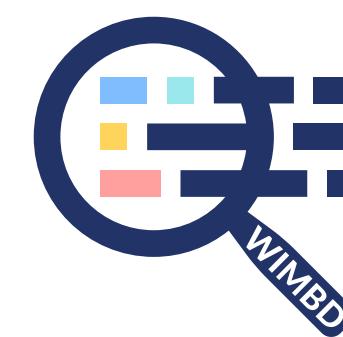
What is?

Syntactic Templates  
in texts

How much  
it happens?

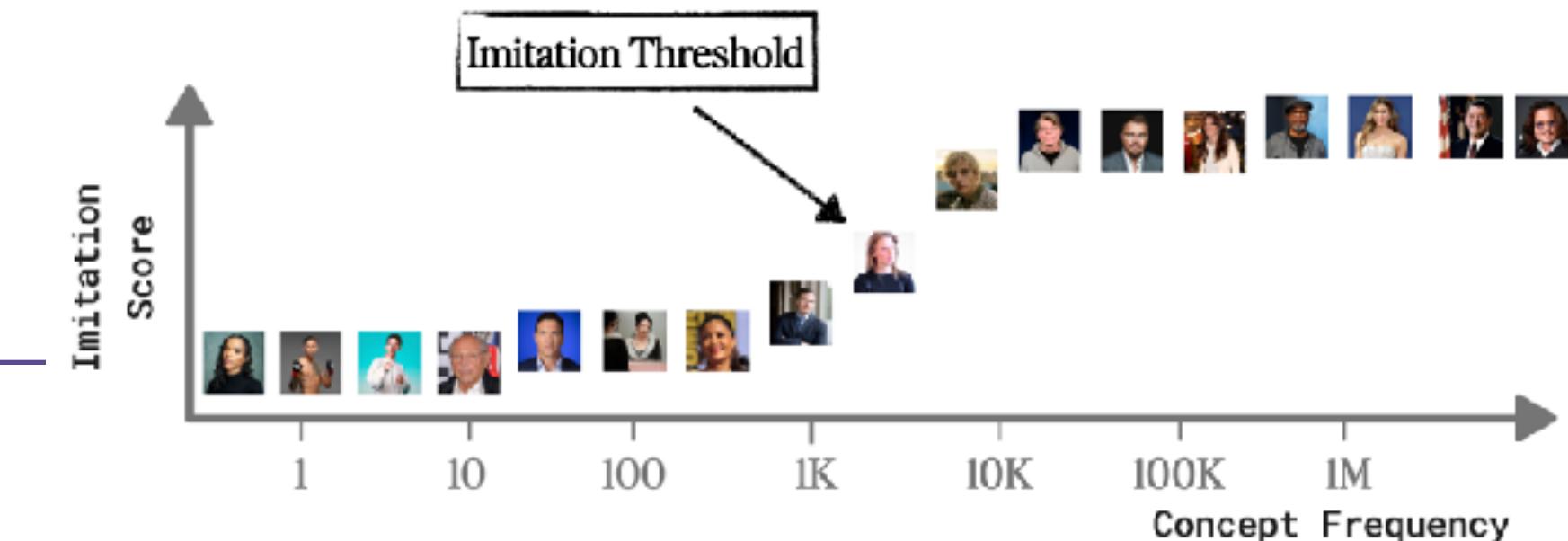
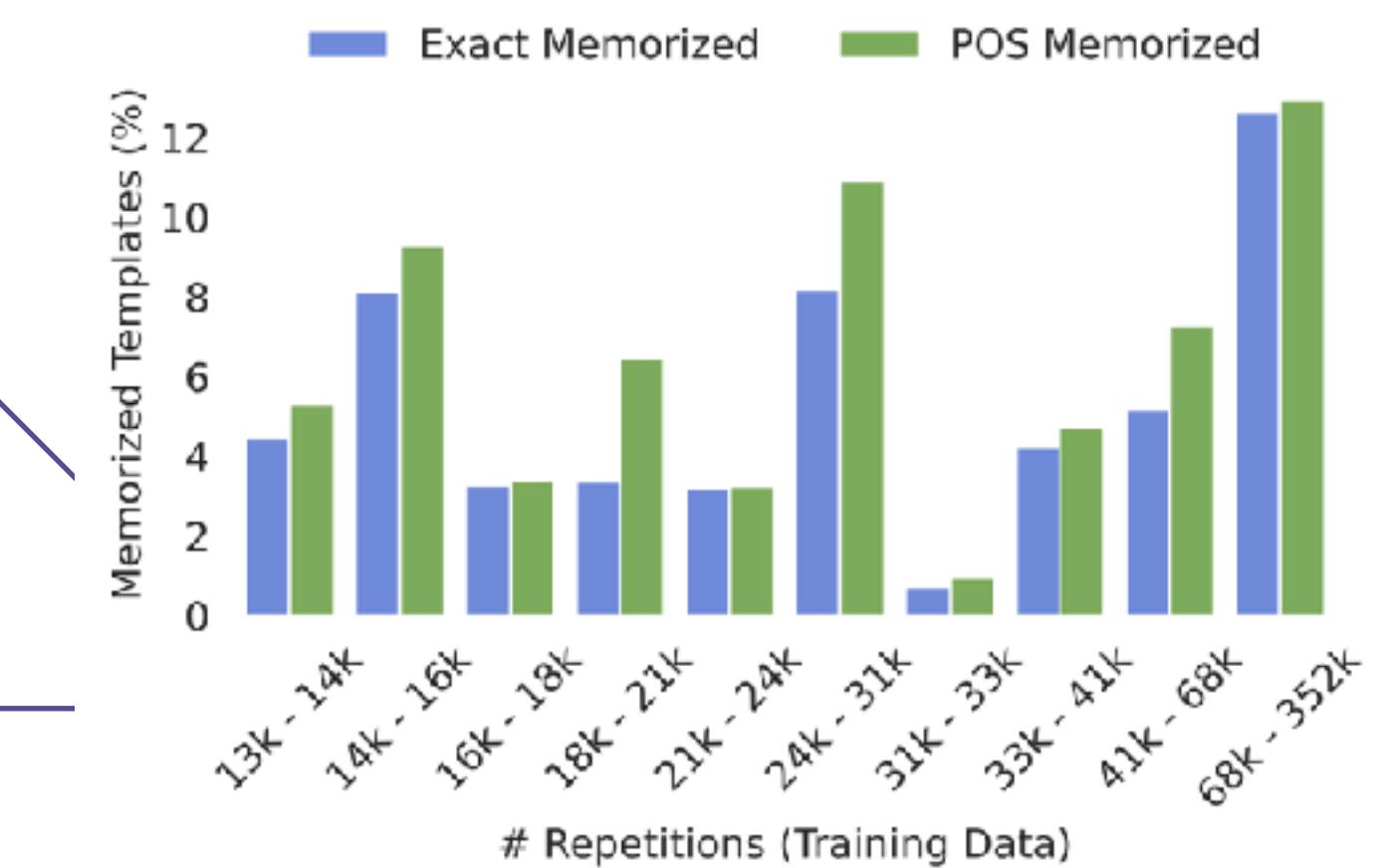
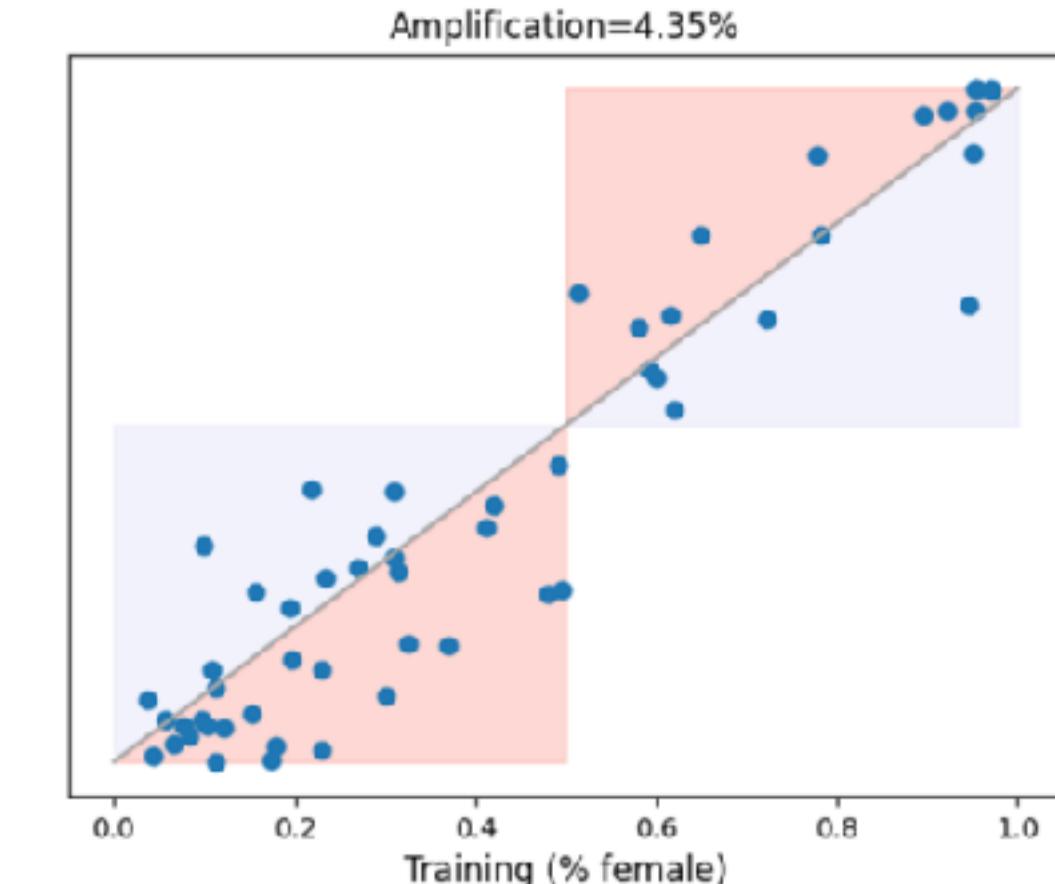
Imitation  
Threshold

When does  
it happen?



$\infty$ -gram

OLMOTRACE



# Thank You!

Questions?

 [yanaiel@gmail.com](mailto:yanaiel@gmail.com)

 @yanaiel

 @yanai.bsky.social