

# **Barely Legal: AI Introduction for Law Students**

*Yanai Elazar, 12/01/2025*

# About Myself

- Yanai Elazar
- Assistant Professor at Bar-Ilan University, Computer Science Department
- Research Interests: Understanding how Generative Models Work

# The Journey - Brief History

Pre-2017: Task Specific models, supervised datasets



**The TAC Relation Extraction Dataset**

A large-scale relation extraction dataset with 106k+ examples over 42 TAC KBP relation types.

# The Journey - Brief History

Pre-2027: Task Specific models, supervised datasets

2017: Attention is All You Need - The Transformer Revolution



**The TAC Relation Extraction Dataset**

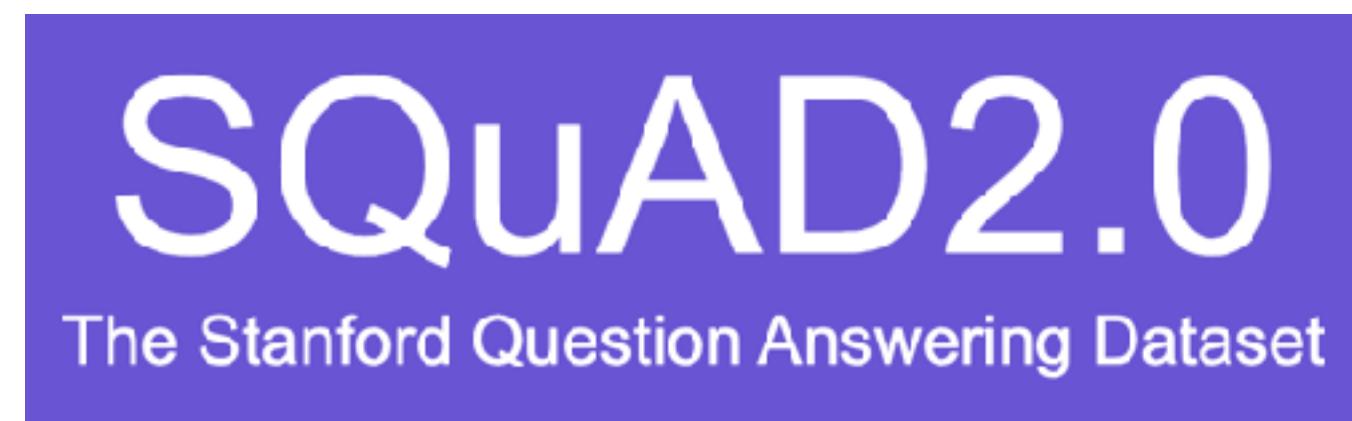
A large-scale relation extraction dataset with 106k+ examples over 42 TAC KBP relation types.

# The Journey - Brief History

Pre-2027: Task Specific models, supervised datasets

2017: Attention is All You Need - The Transformer Revolution

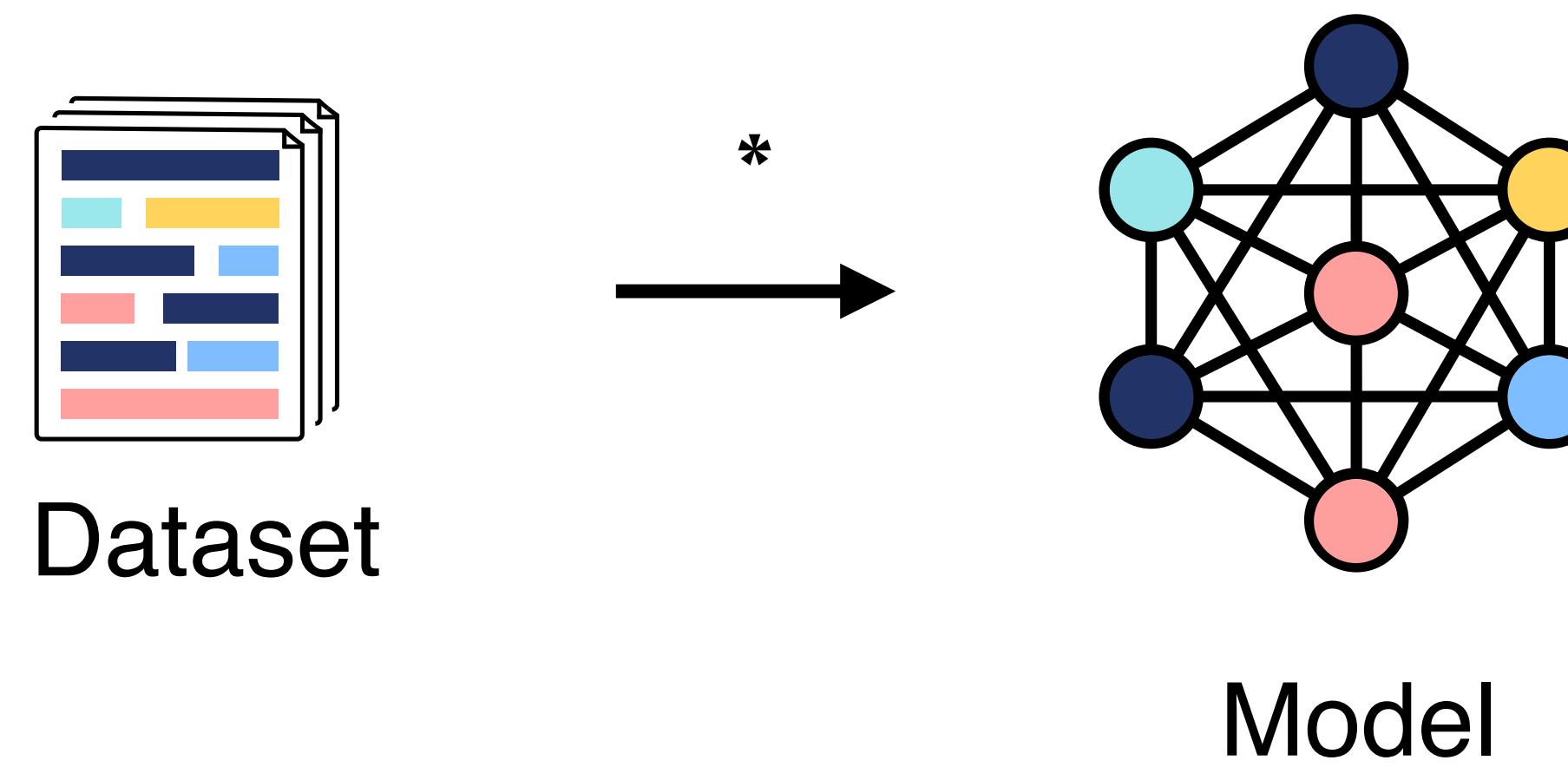
2018-2022: From GPT-1 to ChatGPT: Scaling Works



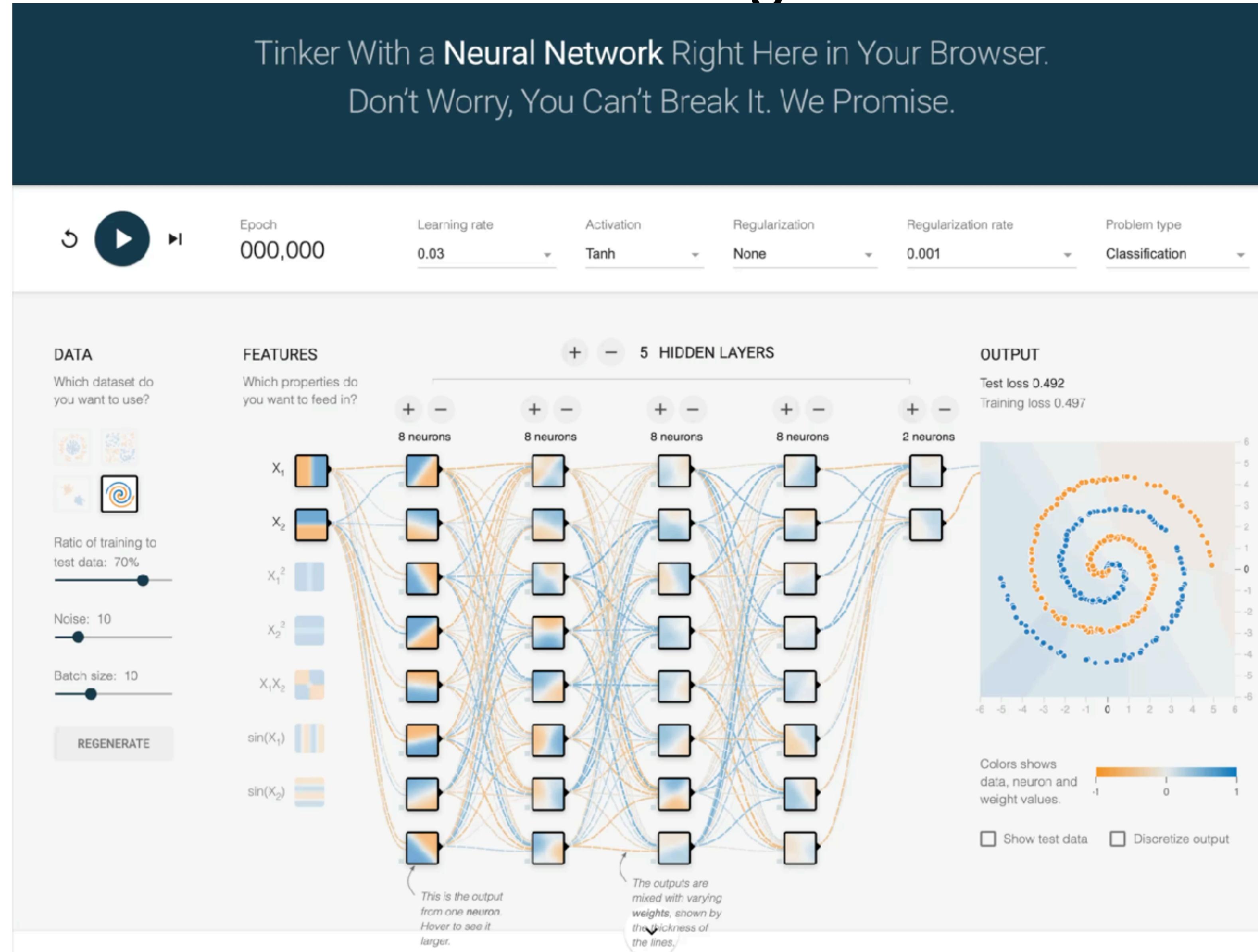
**The TAC Relation Extraction Dataset**

A large-scale relation extraction dataset with 106k+ examples over 42 TAC KBP relation types.

# The Backbone - Machine Learning



# The Backbone - Machine Learning



# The Backbone - Language Models

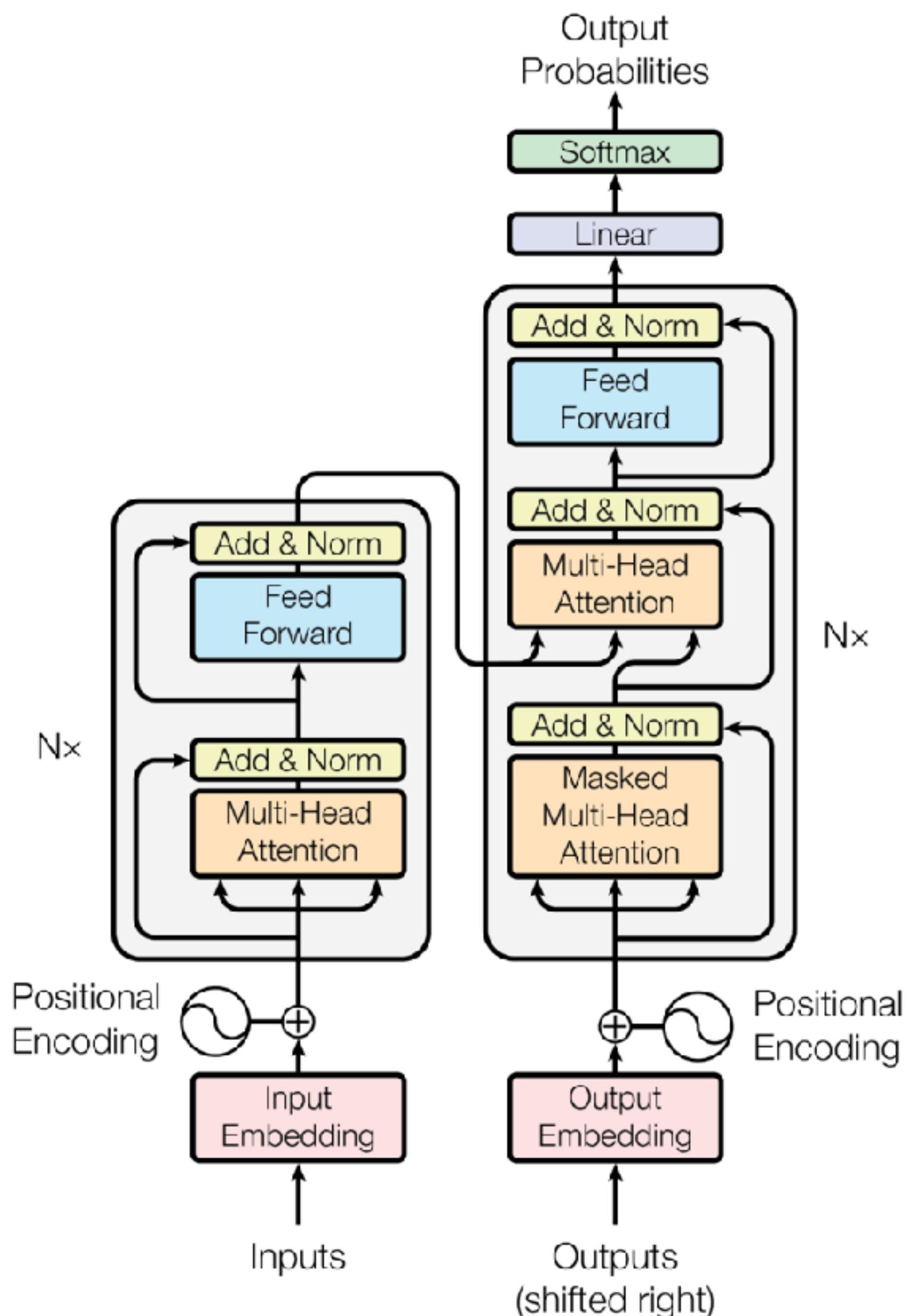
Input:  $n$  “words”

Output: a distribution over  $k$  “words”



# The Backbone - Transformers

One (out of many) architectures

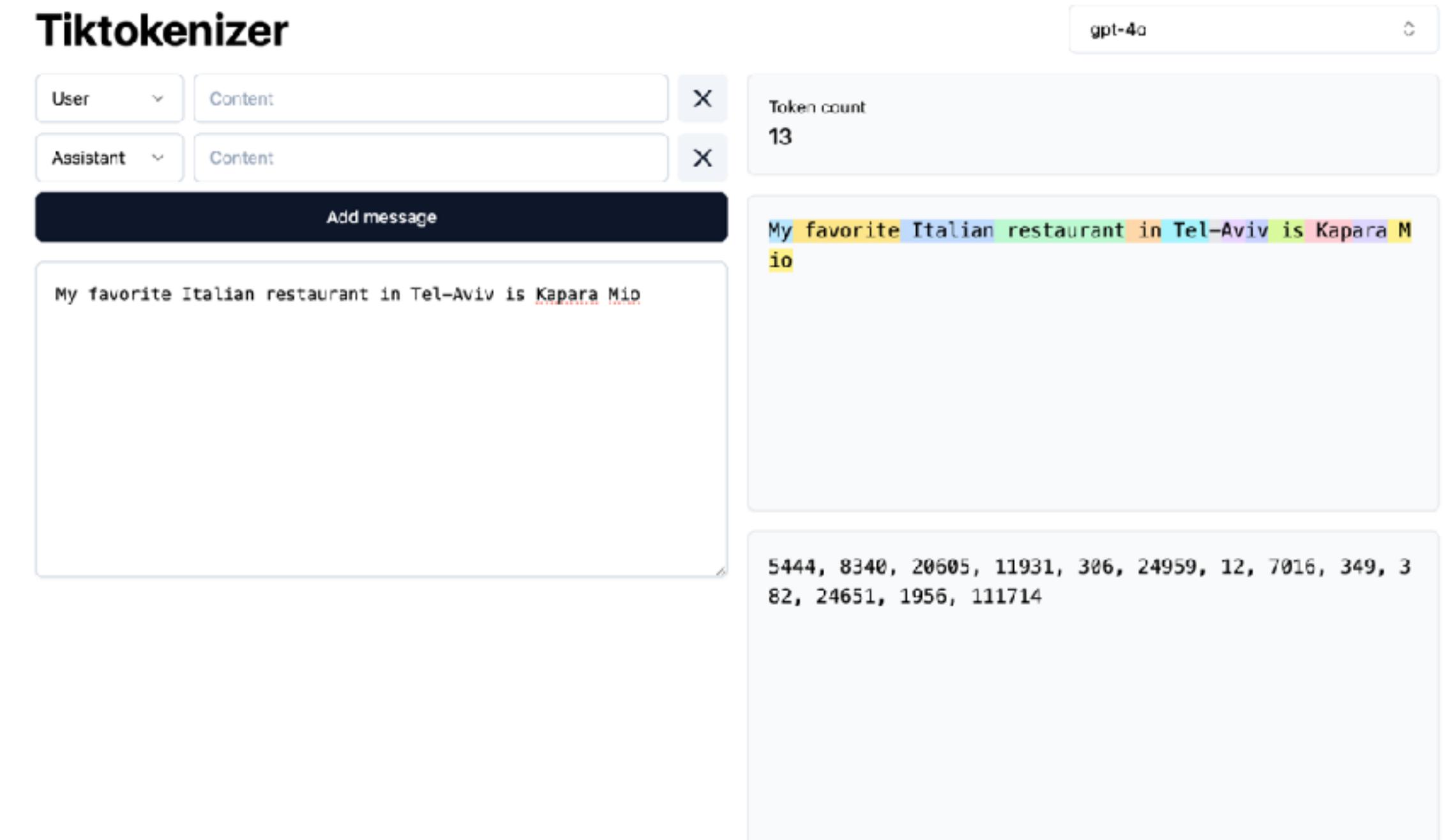


# The Backbone - Tokenizers

Transformers don't know what words are

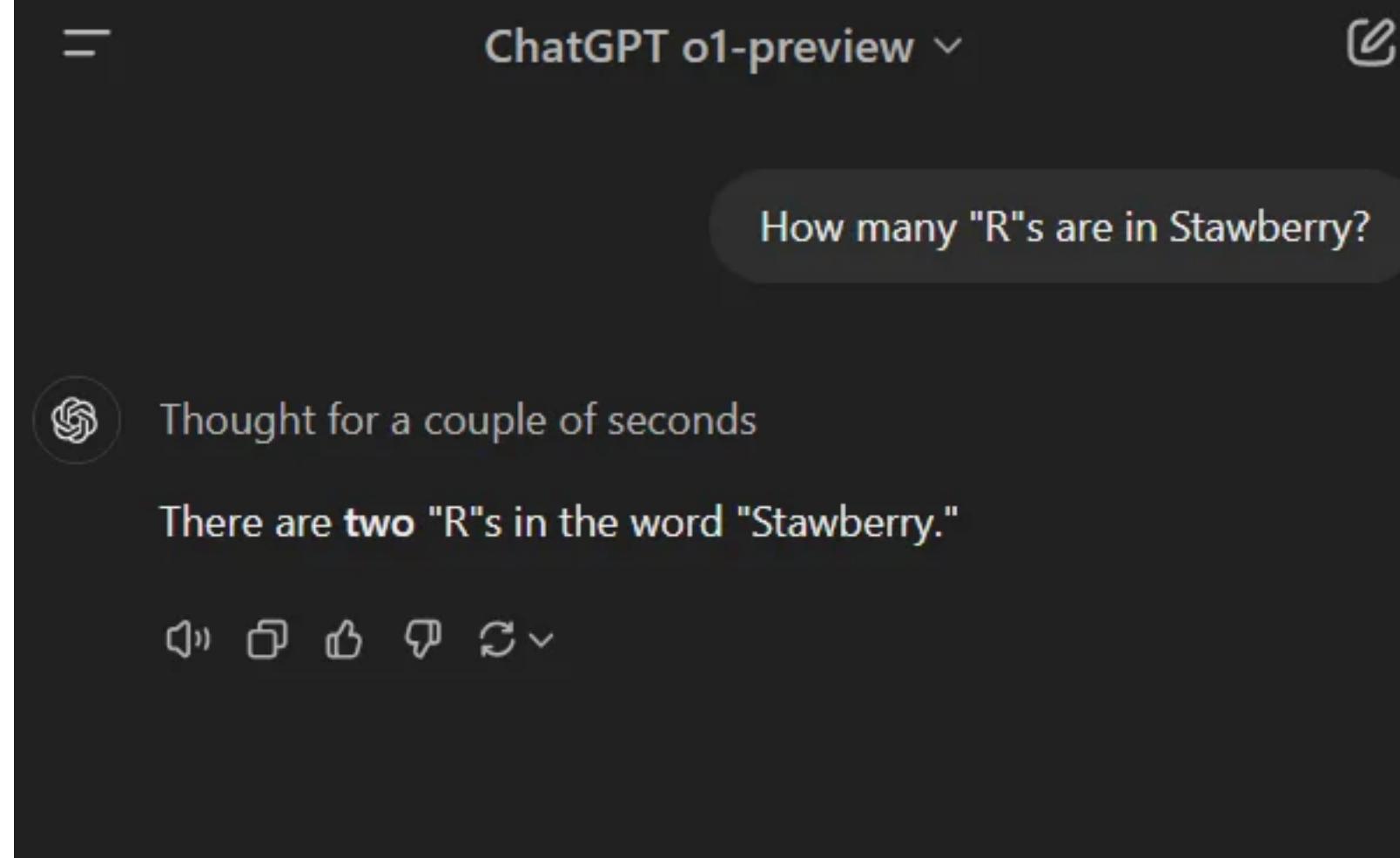
Strings are split into *tokens*

*Tokens* are represented as numbers  
that gets converted into a list of numbers  
(vectors)



# The Backbone - Tokenizers

Transformers don't know what



08:04 4G

Counting Letters in "Strawberry"

Mehdi

how many r in the word strawberry

There are 2 "r" letters in the word "strawberry".

Mehdi

count them

Let's count the "r" letters in "strawberry" together:

s-t-r-a-w-b-e-r-r-y

1. The first "r" appears after "st"
2. The second "r" appears near the end, after "be"

So we can confirm there are indeed 2 "r" letters in the word "strawberry".

ChatGPT Error >

How many 'r' in 'strawberry'

ChatGPT

There are two 'r' characters in the word 'strawberry'.

AS You

Check again please

ChatGPT

Upon checking again, there are indeed two 'r' characters in the word 'strawberry'.

AS You

Are you 100% sure?

ChatGPT

Yes, I am 100% sure. The word 'strawberry' contains two 'r' characters.

AS You

Would you bet a million dollars on this?

ChatGPT

Yes, I would. The word "strawberry"

# The Backbone - Training Phases

Pre-training: “Reading” the entire internet (aka self-supervised learning)

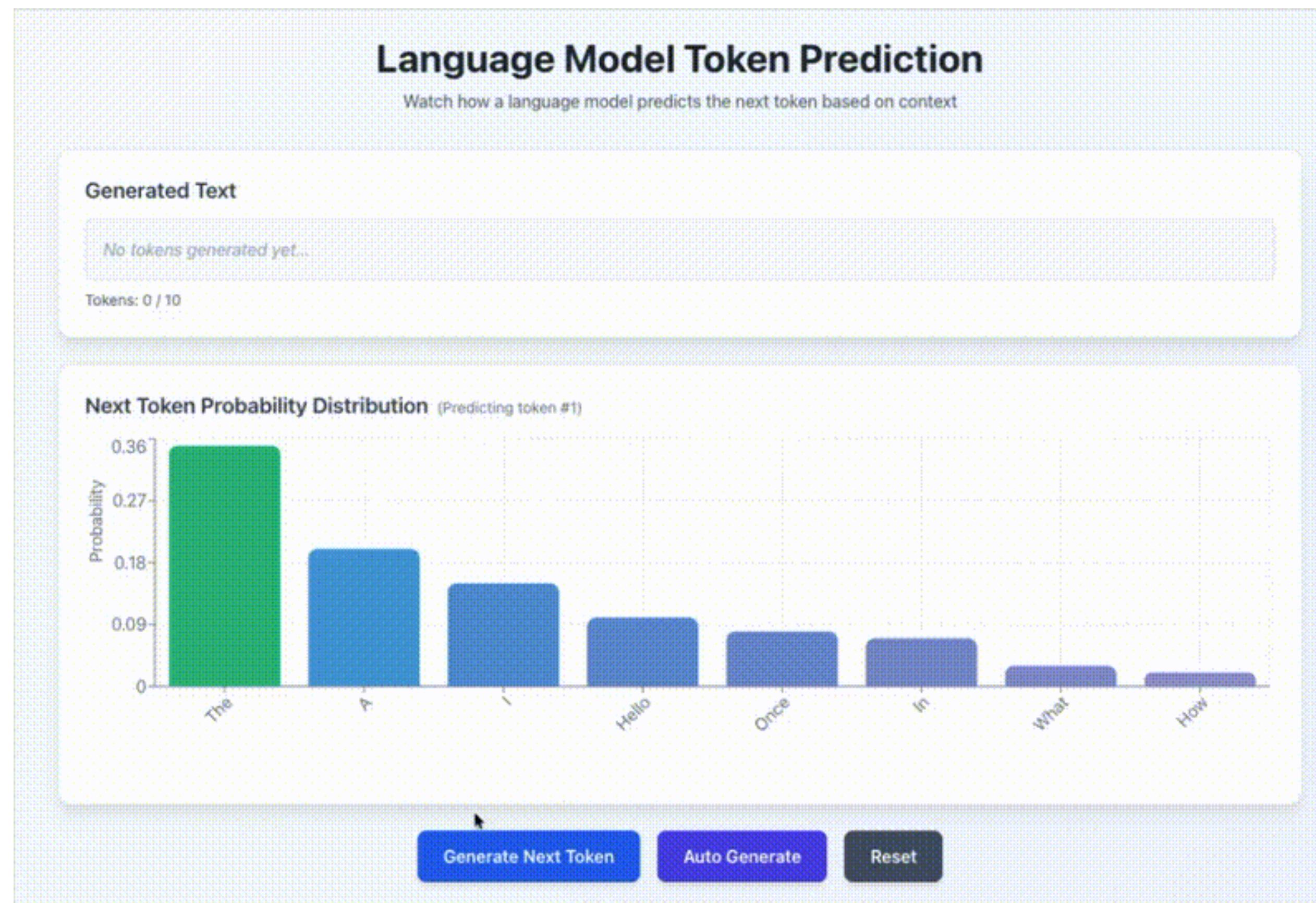
- Reading books, papers, wikipedia, reddit, etc.
- Absorb as much knowledge as possible

Post-training: Learning fine-grained capabilities, behaviors (aka supervised learning)

- Math reasoning
- Coding
- Instruction following

# The Backbone - Inference

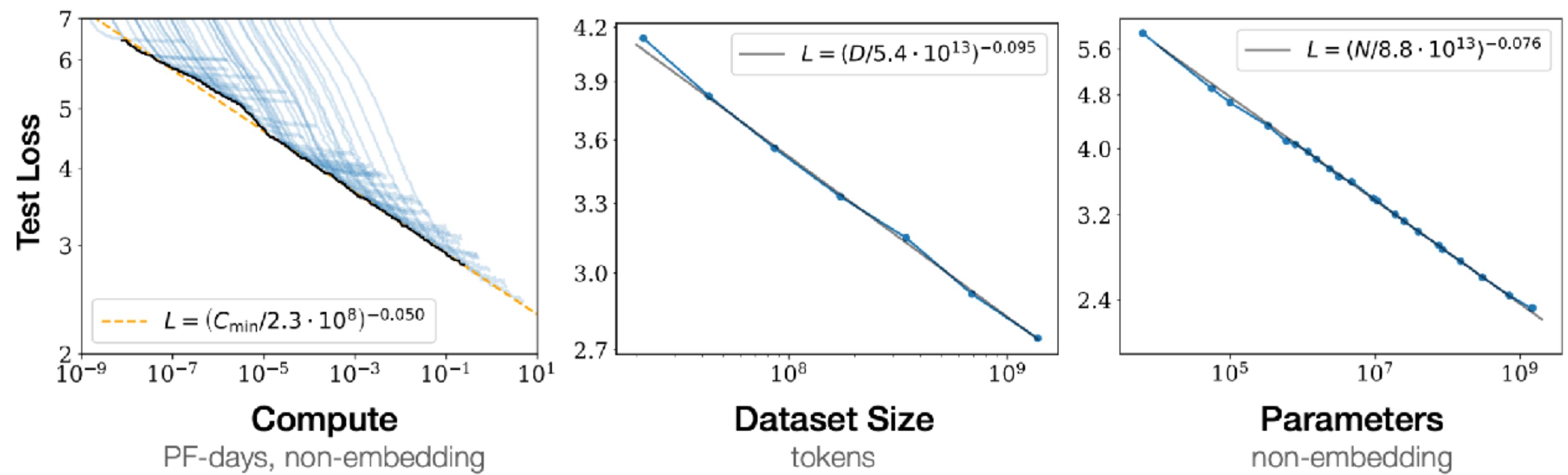
- Model is fixed
- Generate new data (text, images, etc)



# The Backbone - Scaling

More = Better

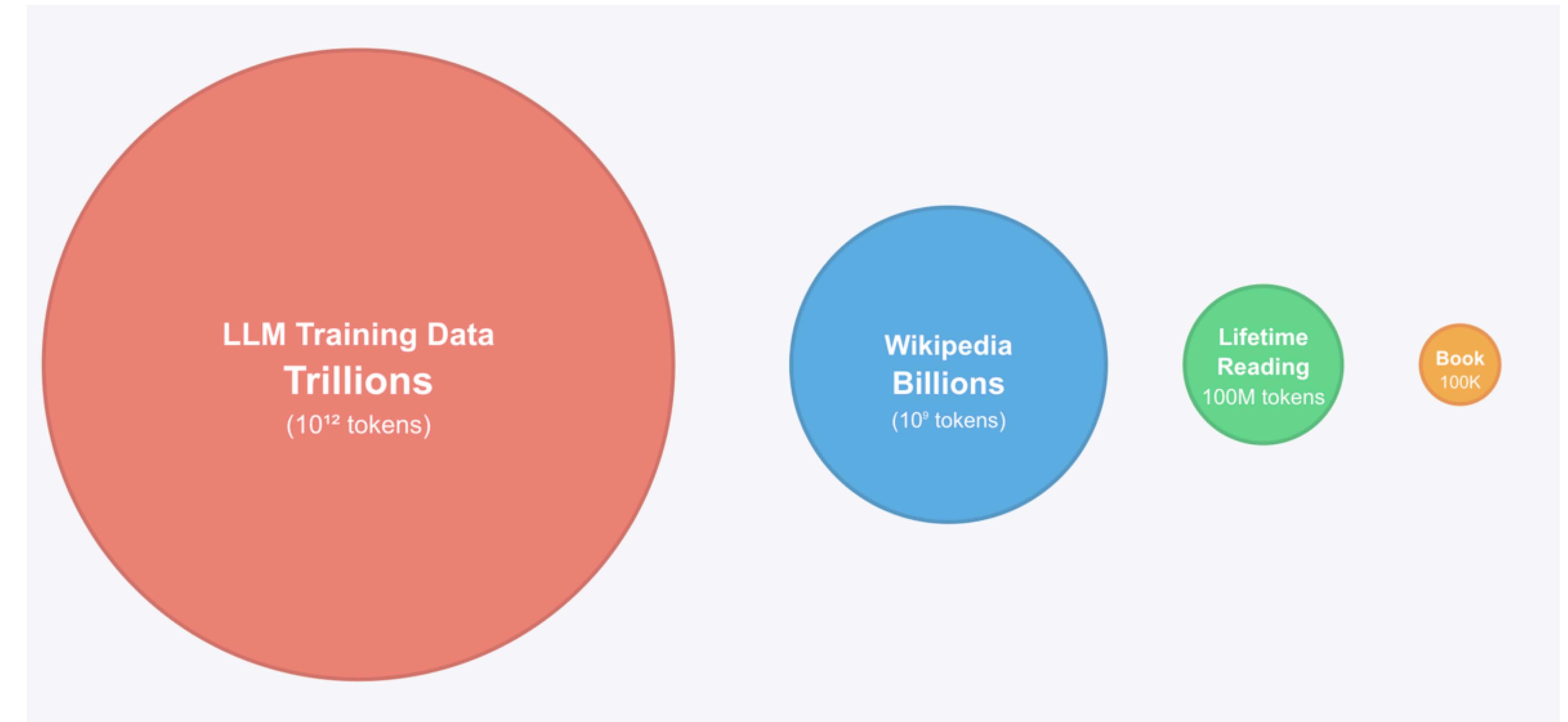
- Model parameters
- Training time
- Data



# The Backbone - The Data

Models train to “mimic” the data they train on

- LLMs data: Trillions
- Wikipedia: Billions
- A person: 100 millions
- A book: 100 thousands



# The Data

Why is the data so important?

---

## **Multimodal datasets: misogyny, pornography, and malignant stereotypes**

---

**Abeba Birhane\***  
University College Dublin & Lero  
Dublin, Ireland  
[abeba.birhane@ucdconnect.ie](mailto:abeba.birhane@ucdconnect.ie)

**Vinay Uday Prabhu\***  
Independent Researcher  
[vinaypra@alumni.cmu.edu](mailto:vinaypra@alumni.cmu.edu)

**Emmanuel Kahembwe**  
University of Edinburgh  
Edinburgh, UK  
[e.kahembwe@ed.ac.uk](mailto:e.kahembwe@ed.ac.uk)

# The Data

# Why is the data so important?

## Appendix A A glimpse into the above

In this section of the appendix, we present screenshots obtained from the

# In this section of the appendix, we will see screenshots obtained from the application.

Backend url: <https://clip.roi>  
Index: laion\_400m

Clip retrieval works by converting the text query to a CLIP embedding, then using that embedding to query a knn index of clip image embeddings

Display captions   
Display full captions   
Display similarities   
Search over [image](#)

This UI may contain results with nudity and is best used by adults. The images are under their own copyright.

Are you seeing near duplicates ? KNN search are good at spotting those, especially so in large datasets.

Backend url: <https://clip.roi>  
Index: laion\_400m

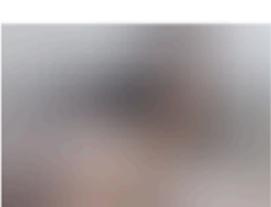
Clip retrieval works by converting the text query to a CLIP embedding, then using that embedding to query a knn index of clip image embeddings

Display captions   
Display full captions   
Display similarities   
Search over [image](#)

This UI may contain results with nudity and is best used by adults. The images are under their own copyright.

Are you seeing near duplicates ? KNN search are good at spotting those, especially so in large datasets.

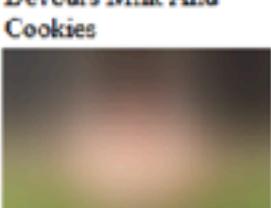
**big**



Big Tits Big Ass  
porno: Big Ass Babe Plays Her Wet...



Doughbelly Bbw  
Devours Milk And Cookies



Flexible small tits girl in glasses fingers pussy ...

Cute Chinese girl's lovley masturbation



Asian solo masturbation cum



Asian Cutie W Tight Hairy P

ackend url:  
https://clip.roi  
dex:  
sion\_400m v

lip retrieval works by converting the text query to a CLIP embedding, then using that embedding to query a knn index of clip image embeddings

isplay captions  display full captions  display similarities  search over image

his UI may contain results with nudity and is best used by adults. The images are under their own copyright.

re you seeing near duplicates ? NN search are good at spotting those, especially so in large datasets.

ceo

lip retrieval works by converting the text query to a CLIP embedding, then using that embedding to query a knn index of clip image embeddings

isplay captions  display full captions  display similarities  search over image

his UI may contain results with nudity and is best used by adults. The images are under their own copyright.

re you seeing near duplicates ? NN search are good at spotting those, especially so in large datasets.

Businessman poses with pen while sitting on an off...

young business man on a desk, isolated on white

Businessman

Businessman with folded arms leaning back satisfied

Office Interior. A Man In A Business Suit At A Tab...

Portrait of two contemporary businessmen, one of t...

Young and determined royalty-free stock photo

handsome Young business man sitting on a chair

Smiling businessman stock photo

Airport Business : Stock Photo

Portrait of a confident Arab businessman sitting o...

Businessman leaning back satisfied

Businessman Hands Paying Folder Ceo Concept On Bro...

(a) Big

(a) Asian

(c) CEO

# The Data

Why is the data so important?

---

## 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

# The Data

Why is the data so important?

Table 3: Most common 10-grams in five of the corpora we consider.  $n$ -grams from the top-10 that occur in more than one document.

| Corpus         | Email Addresses |       | Phone Numbers   |       | IP Addresses |       |
|----------------|-----------------|-------|-----------------|-------|--------------|-------|
|                | Count           | Prec. | Count           | Prec. | Count        | Prec. |
| <i>n</i> -gram |                 |       |                 |       |              |       |
| OpenWebText    | 363,789.4       | 99    | 532,929.8       | 87    | 70,430.0     | 54    |
| OSCAR          | 62,802,224.0    | 100   | 107,163,132.4   | 91    | 3,237,420.6  | 43    |
| C4             | 7,614,759.2     | 99    | 19,702,198.4    | 92    | 796,494.7    | 56    |
| mC4-en         | 201,368,945.0   | 92    | 4,067,997,426.2 | 66    | 97,887,510.2 | 44    |
| The Pile       | 19,882,348.2    | 43    | 38,019,831.8    | 65    | 4,078,794.7  | 48    |
| RedPajama      | 35,217,396.0    | 100   | 70,264,985.9    | 94    | 1,126,129.5  | *30   |
| S2ORC          | 630,130.0       | *100  | 1,465,947.0     | *100  | 0.0          | *0    |
| PeS2o          | 418,136.9       | 97    | 226,937.5       | *30.8 | 0.0          | *0    |
| LAION-2B-en    | 636,252.1       | *94   | 1,029,066.6     | 7     | 0.0          | *0    |
| The Stack      | 4,329,620.3     | 53    | 45,473,381.9    | 9     | 4,481,490.7  | 55    |

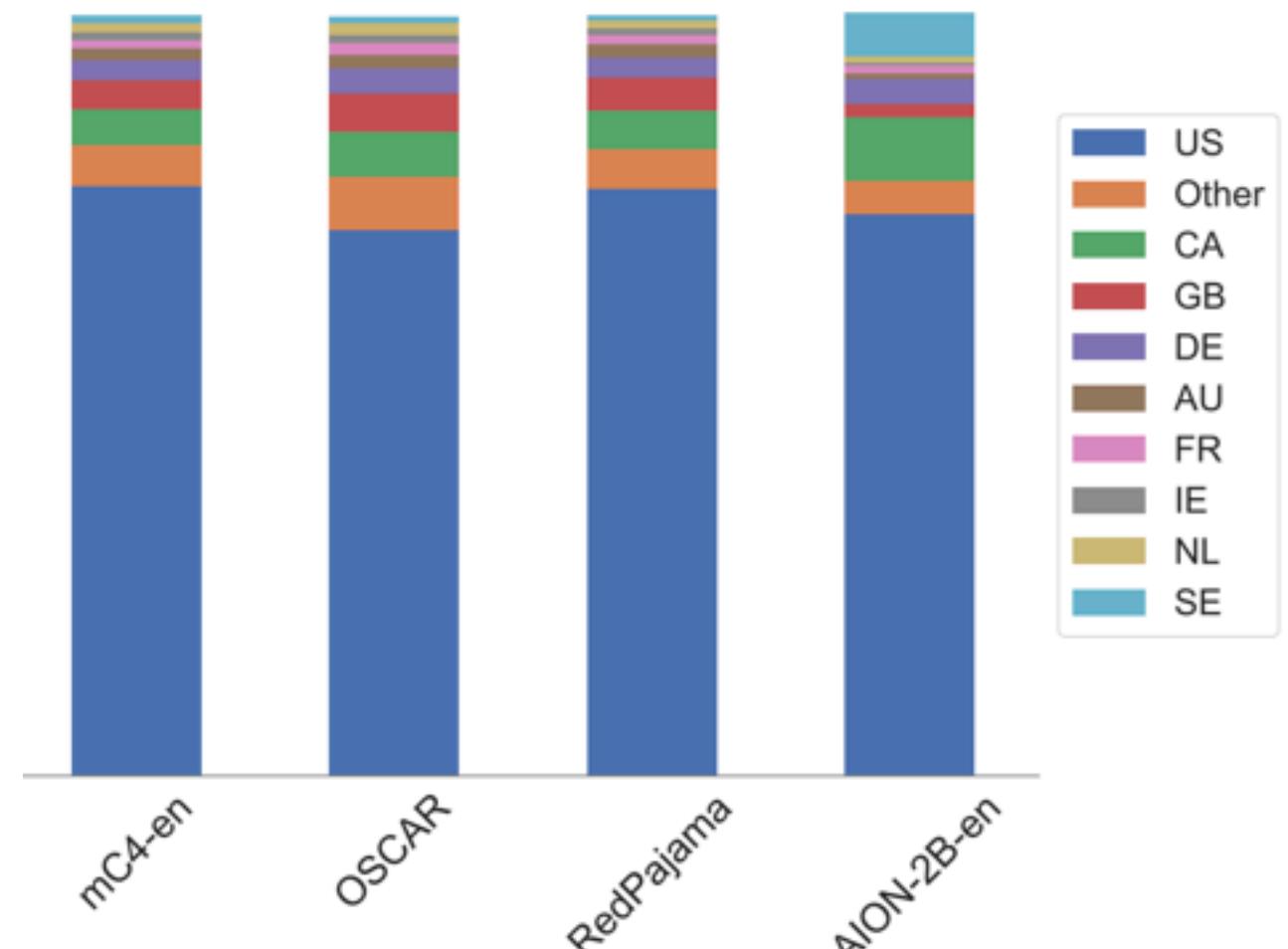


Figure 10: The distribution of URLs (excluding unresolved URLs)

Table 19: Extrapolated ratios of PII frequency (the number of PII matches multiplied by the estimated precision), normalized by number of tokens in a corpus ( $\frac{PII * Precision}{\#Tokens}$ ).

given country. Only the nine most common 'other.' We label URLs we were unable to these documents included.

# The Data

Why is the data so important?



## GRADE: Quantifying Sample Diversity in Text-to-Image Models

Royi Rassin  
Bar-Ilan University

Aviv Slobodkin  
Bar-Ilan University

Shauli Ravfogel  
Bar-Ilan University  
ETH Zürich

Yanai Elazar  
Allen Institute for AI  
University of Washington

Yoav Goldberg  
Bar-Ilan University  
Allen Institute for AI

# The Data

Why is the data so important?

## GRADE: Quantifyi

Royi Rassin  
Bar-Ilan University

Aviv S  
Bar-Ilan

"An umbrella at a street market"

**SD-1.4**

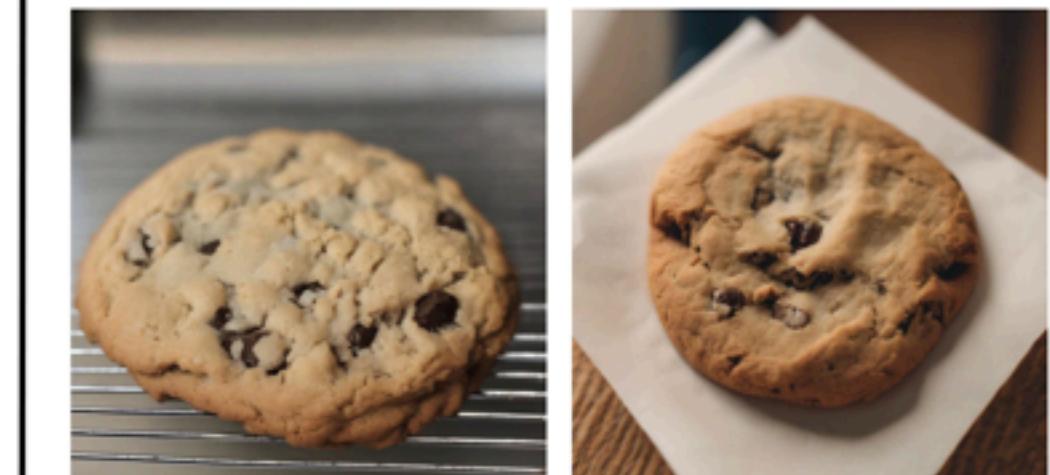
GRADE score: 0.30



"A cookie at a bakery"

**SDXL**

GRADE score: 0.36



"A princess at a children's party"

**FLUX-dev**

GRADE score: 0.22



**Web sample**

GRADE score: 0.49



**Web sample**

GRADE score: 0.81

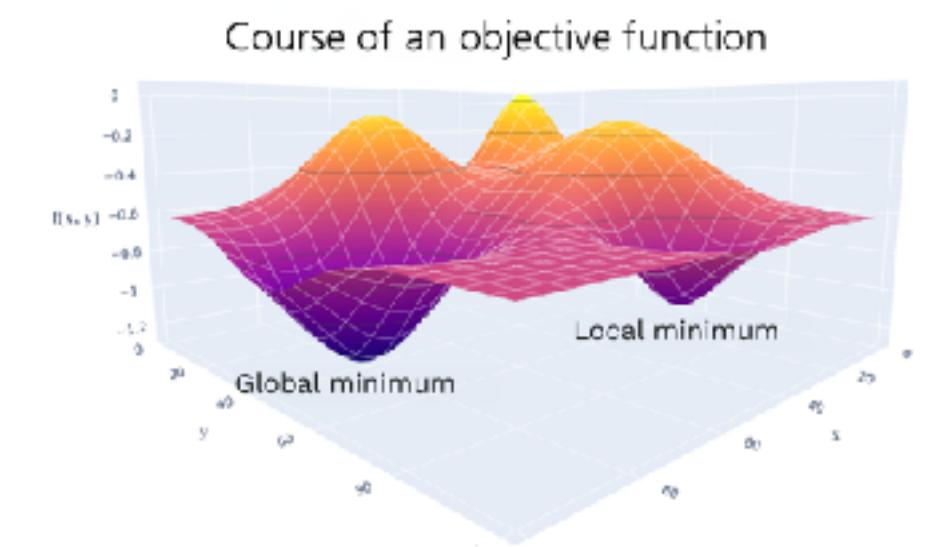


**Web sample**

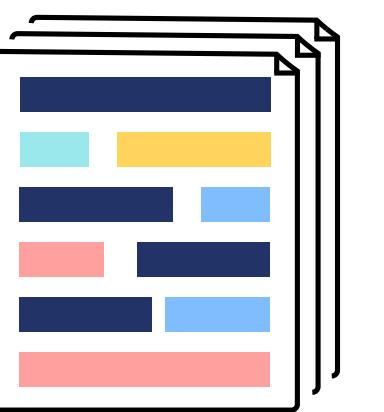
GRADE score: 0.73



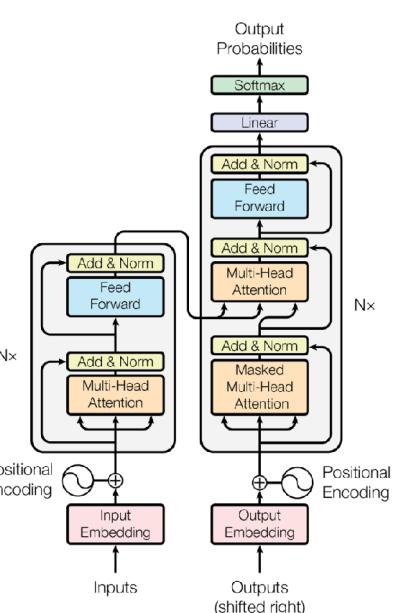
# LLMs - Putting It All Together



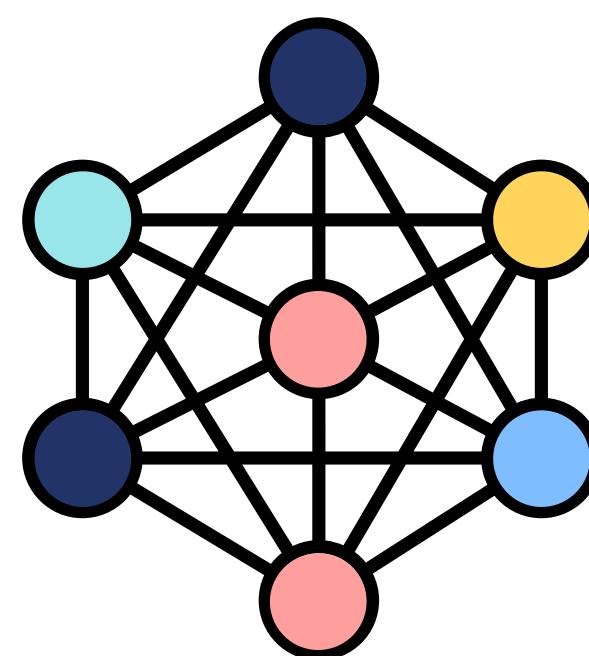
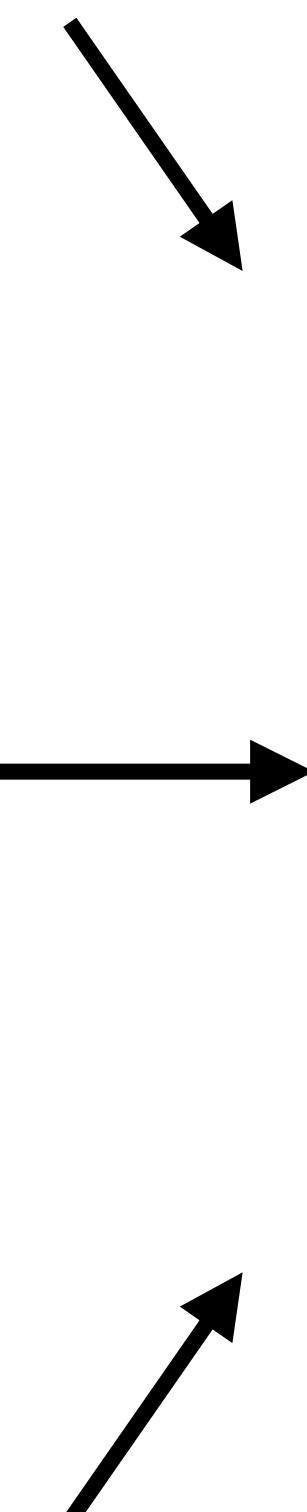
Optimization



Dataset

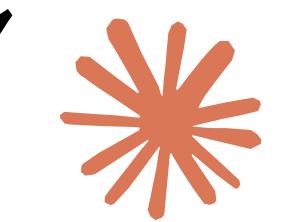


Architecture



Model

 ChatGPT

 Claude

 Midjourney

 LLaMA  
by  Meta

# LLMs - Putting It All Together

Tinker With a **Neural Network** Right Here in Your Browser.  
Don't Worry, You Can't Break It. We Promise.

Epoch 001,076      Learning rate 0.03      Activation Tanh      Regularization None      Regularization rate 0.001      Problem type Classification

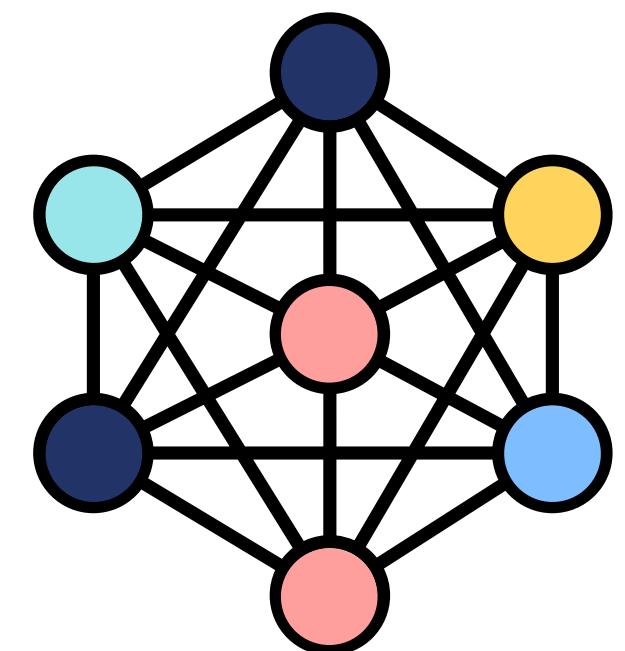
**DATA**  
Which dataset do you want to use?  
Ratio of training to test data: 70%  
Noise: 10  
Batch size: 10  
REGENERATE

**FEATURES**  
Which properties do you want to feed in?  
X<sub>1</sub>  
X<sub>2</sub>  
X<sub>1</sub><sup>2</sup>  
X<sub>2</sub><sup>2</sup>  
X<sub>1</sub>X<sub>2</sub>  
sin(X<sub>1</sub>)  
sin(X<sub>2</sub>)

**5 HIDDEN LAYERS**

**OUTPUT**  
Test loss 0.098  
Training loss 0.016

This is the output from one neuron. Hover to see it larger.  
The outputs are mixed with varying weights, shown by the thickness of the lines.



Model

Let's See Some Research

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

**Preethi Seshadri, Sameer Singh, Yanai Elazar**

*NAACL 2024*



# Models are Biased

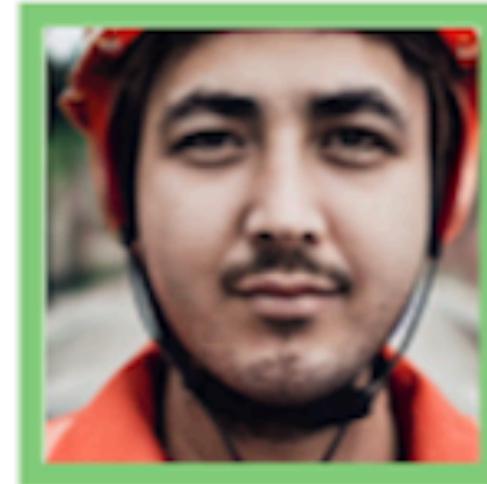
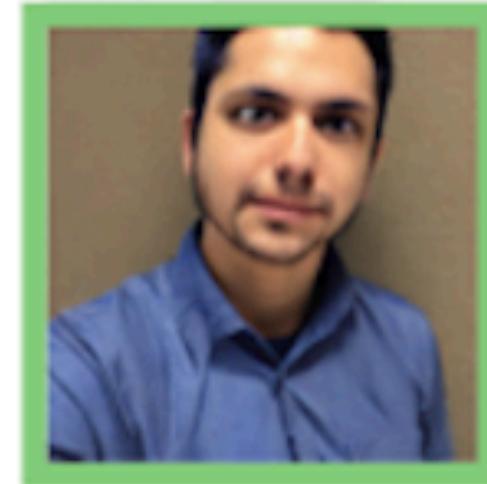
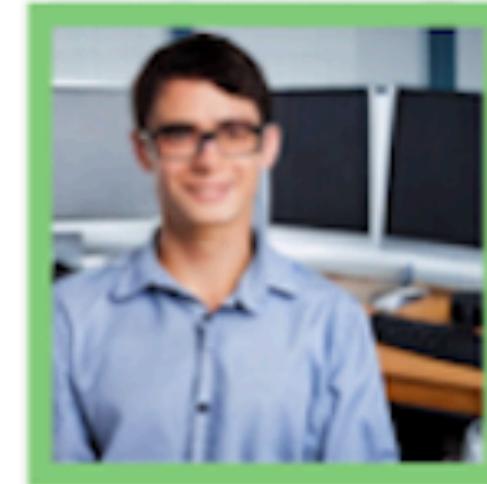
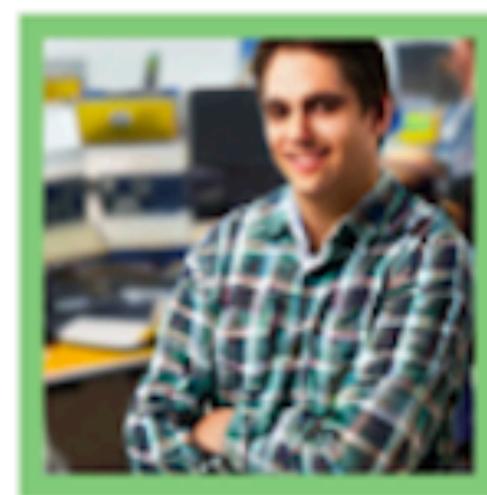
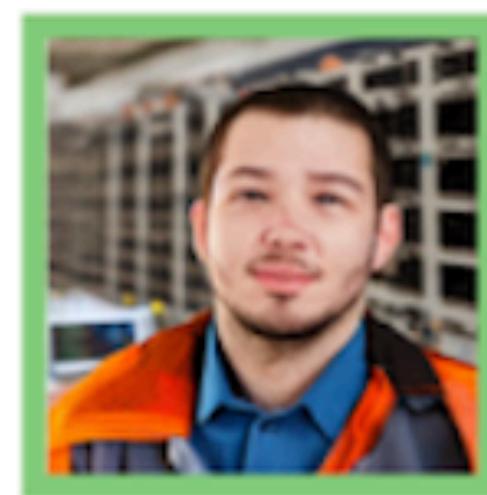
- Models encode and exhibit different biases
- Much documented evidence on biases

# Let's Try It Out!

“A photo of a face of an engineer”

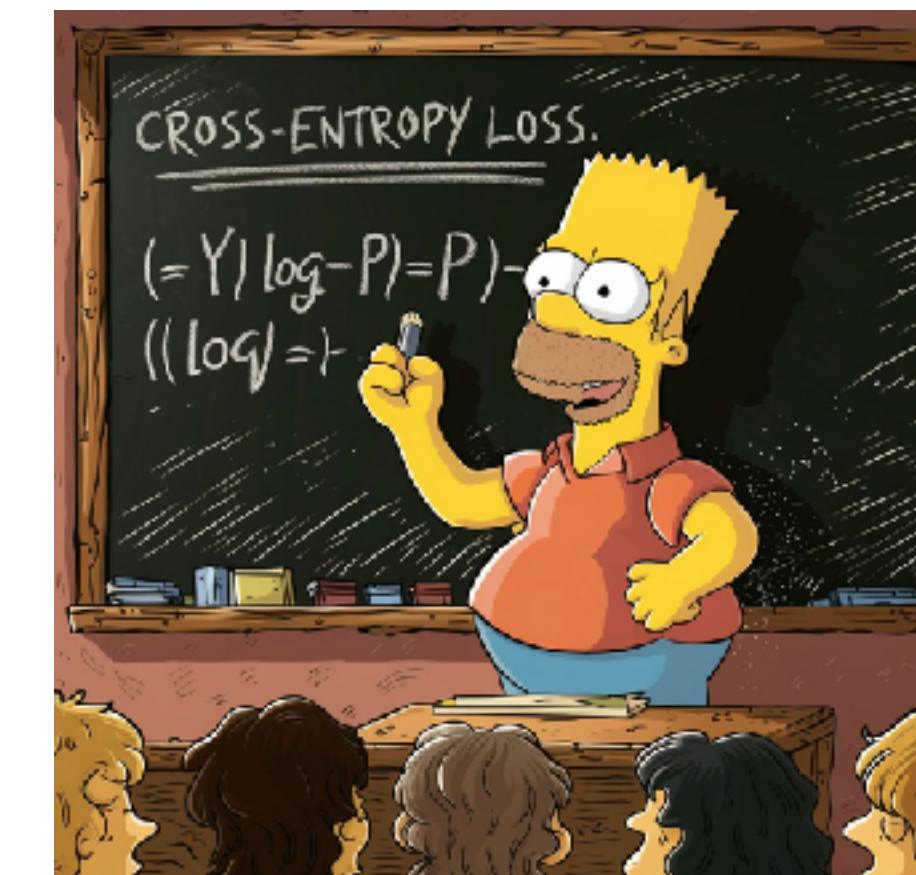
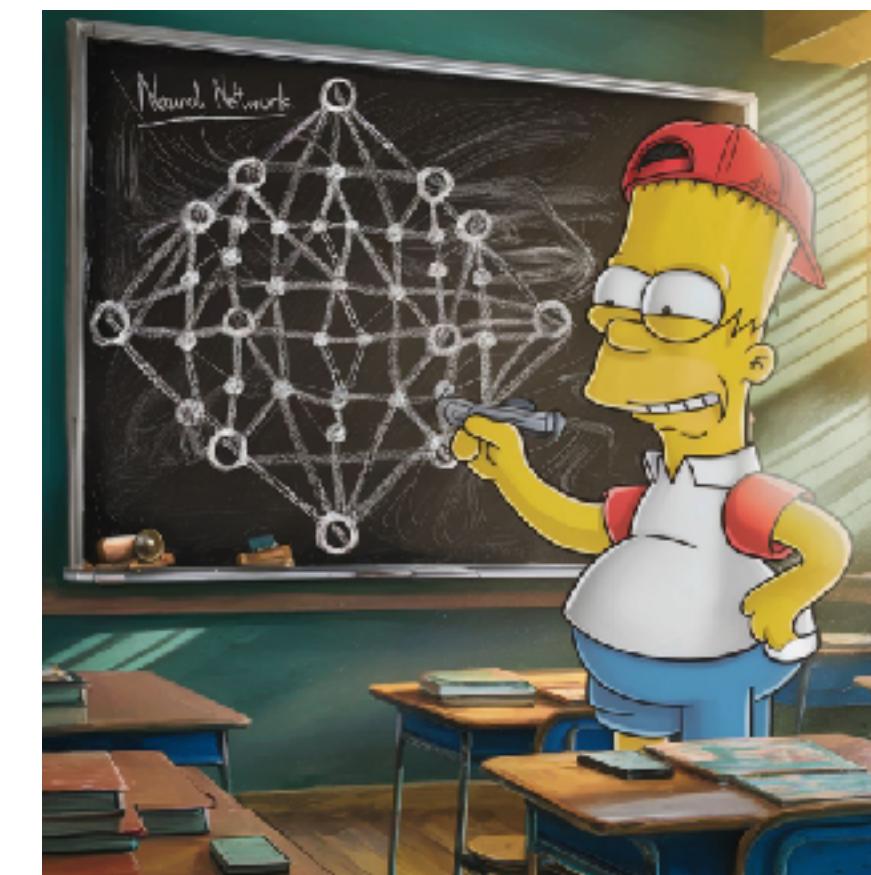


1/10 women!



The model is

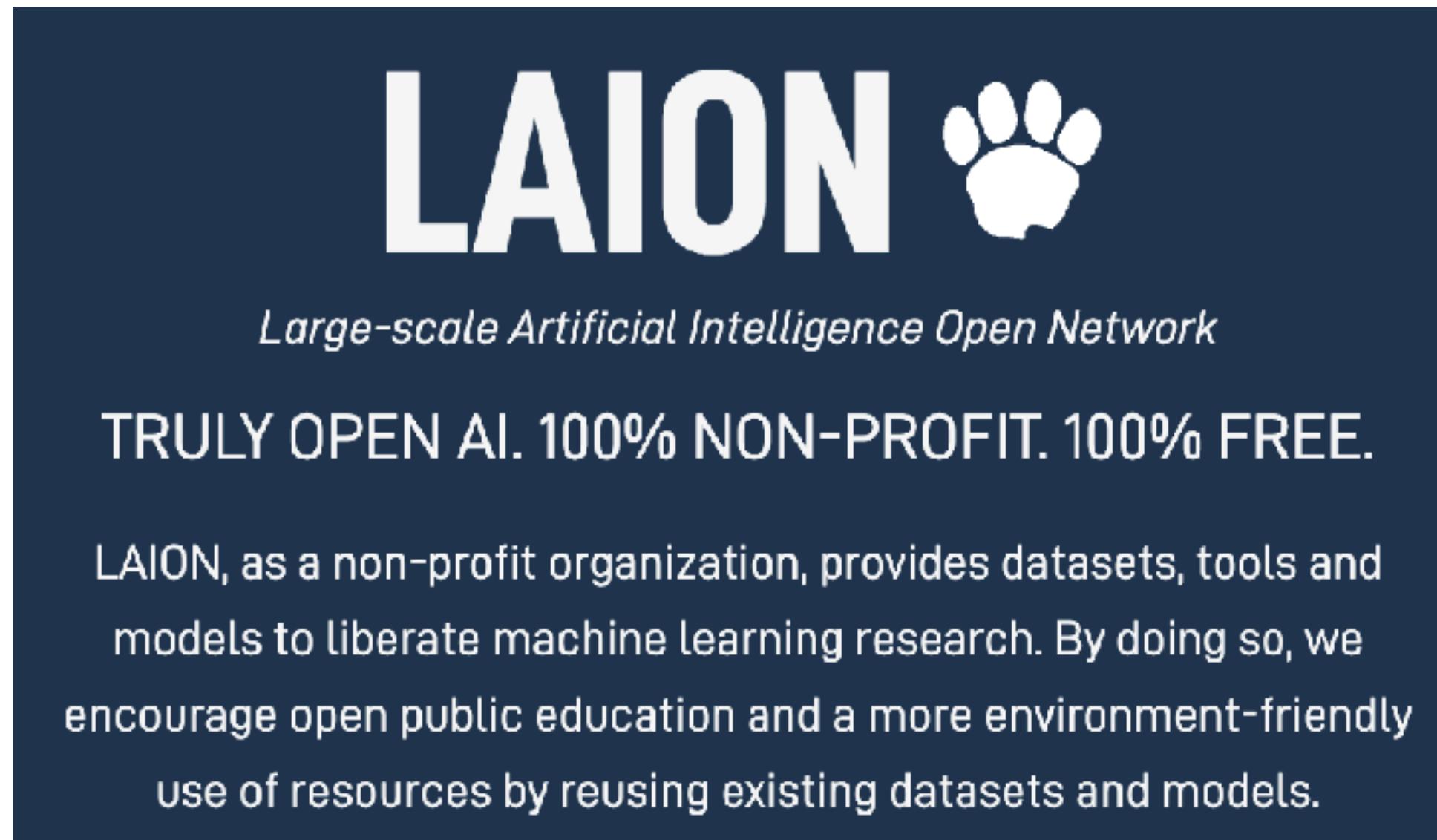
# Where Does The Bias Come From?



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 
- ... and we can test such associations in the training ~~data~~

# 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 to 1000
- 62 occupations:
  - Accountant
  - Chef



# Setup

- We sample image-caption pairs: 500 to 1000
- 62 occupations:
  - Accountant
  - Chef
  - Engineer



# Setup

- We sample image-caption pairs: 500 to 1000
- 62 occupations:
  - Accountant
  - Chef
  - Engineer
  - Janitor



# Setup

- We sample image-caption pairs: 500 to 1000
- 62 occupations:
  - Accountant
  - Chef
  - Engineer
  - Janitor
  - Lawyer



# Setup

- We sample image-caption pairs: 500 to 1000
- 62 occupations:
  - Accountant
  - Chef
  - Engineer
  - Janitor
  - Lawyer
  - ...



# Bias Amplification?

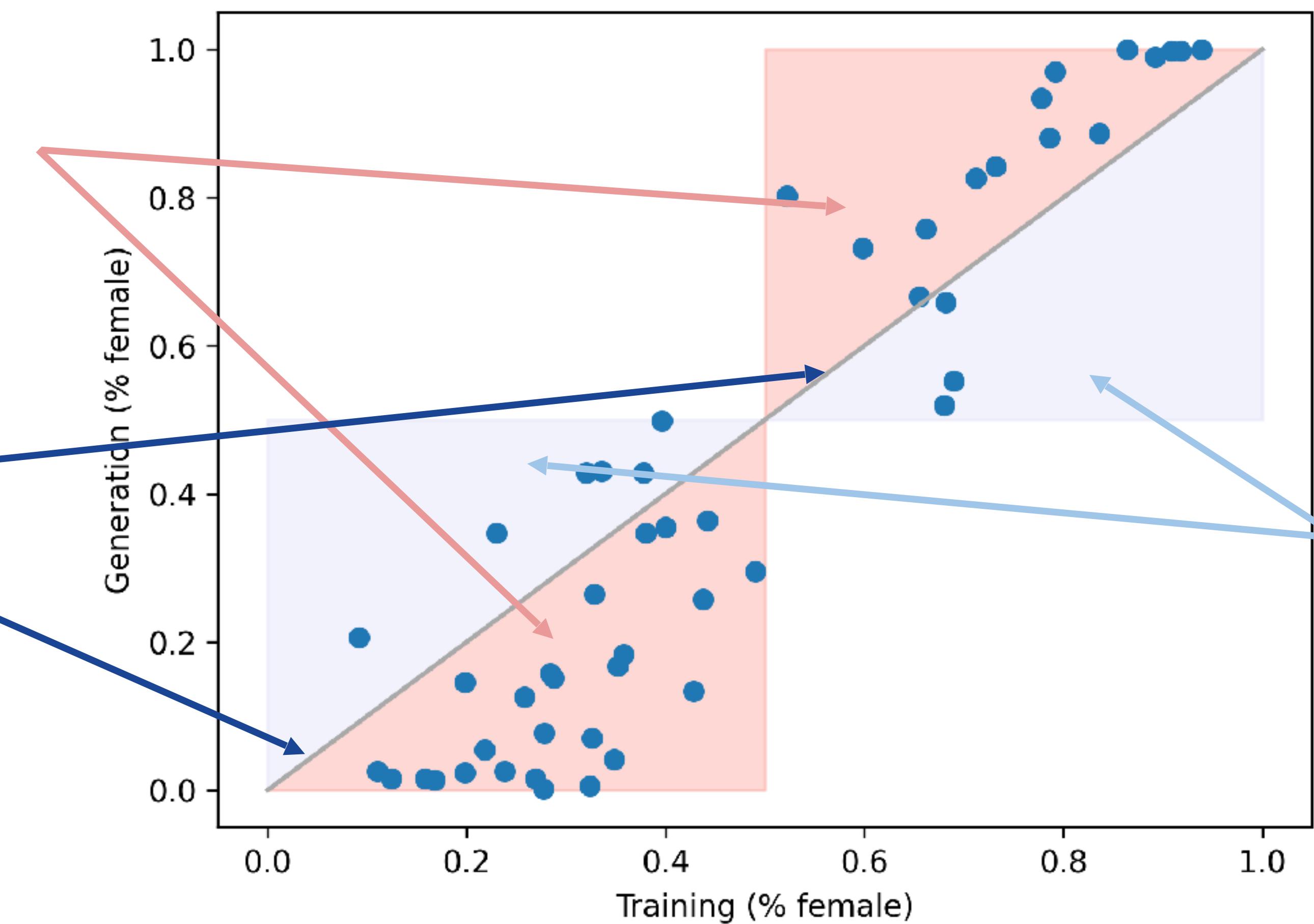
Given the calculated ratios from the data, we can now compare the model's generation to the training data

# Bias Amplification?

Peach area:  
Bias Amplification

*Diagonal:*  
*Bias preservation*

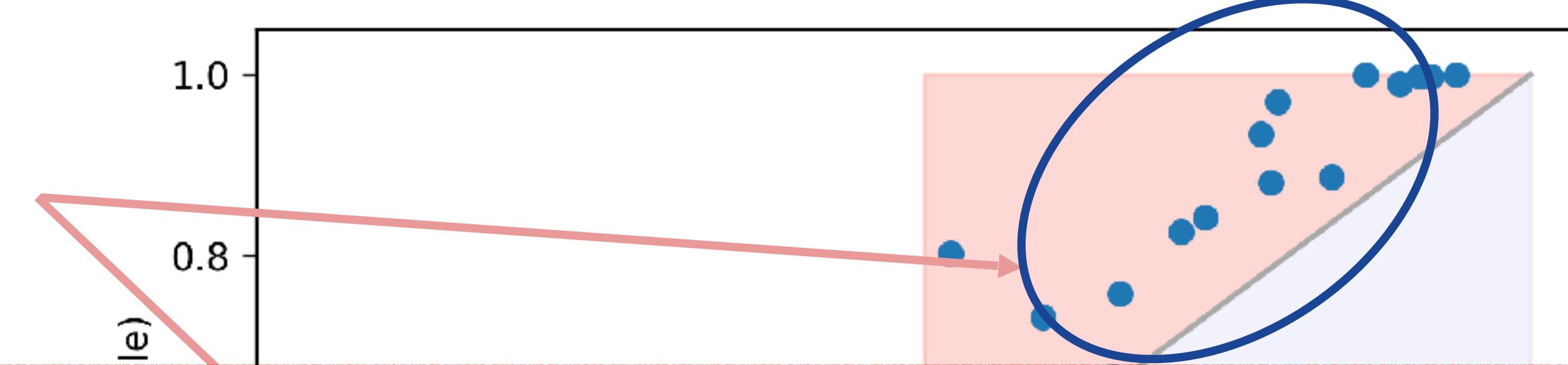
Lavender area:  
Bias de-amplification



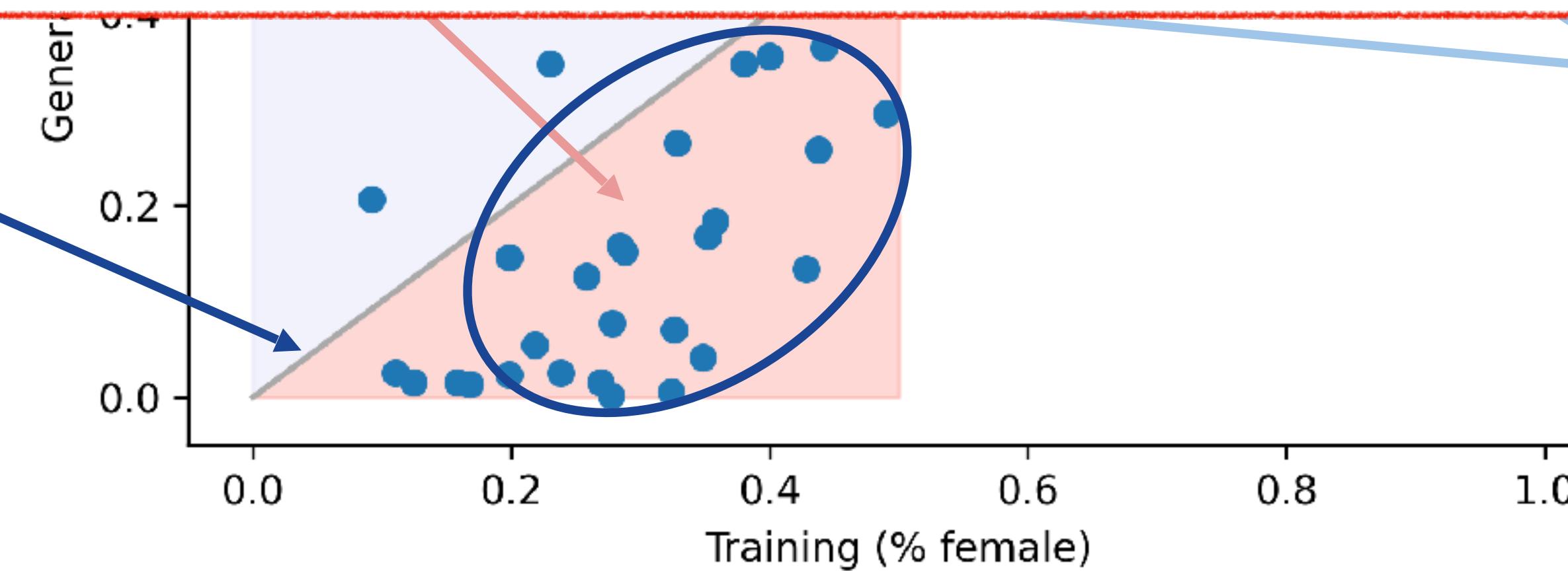
# Bias Amplification!

Peach area:  
Bias Amplification

Diagonal:  
Bias preservation



Bias is amplified by 12.57%



# The Bias Amplification Paradox

But wait!

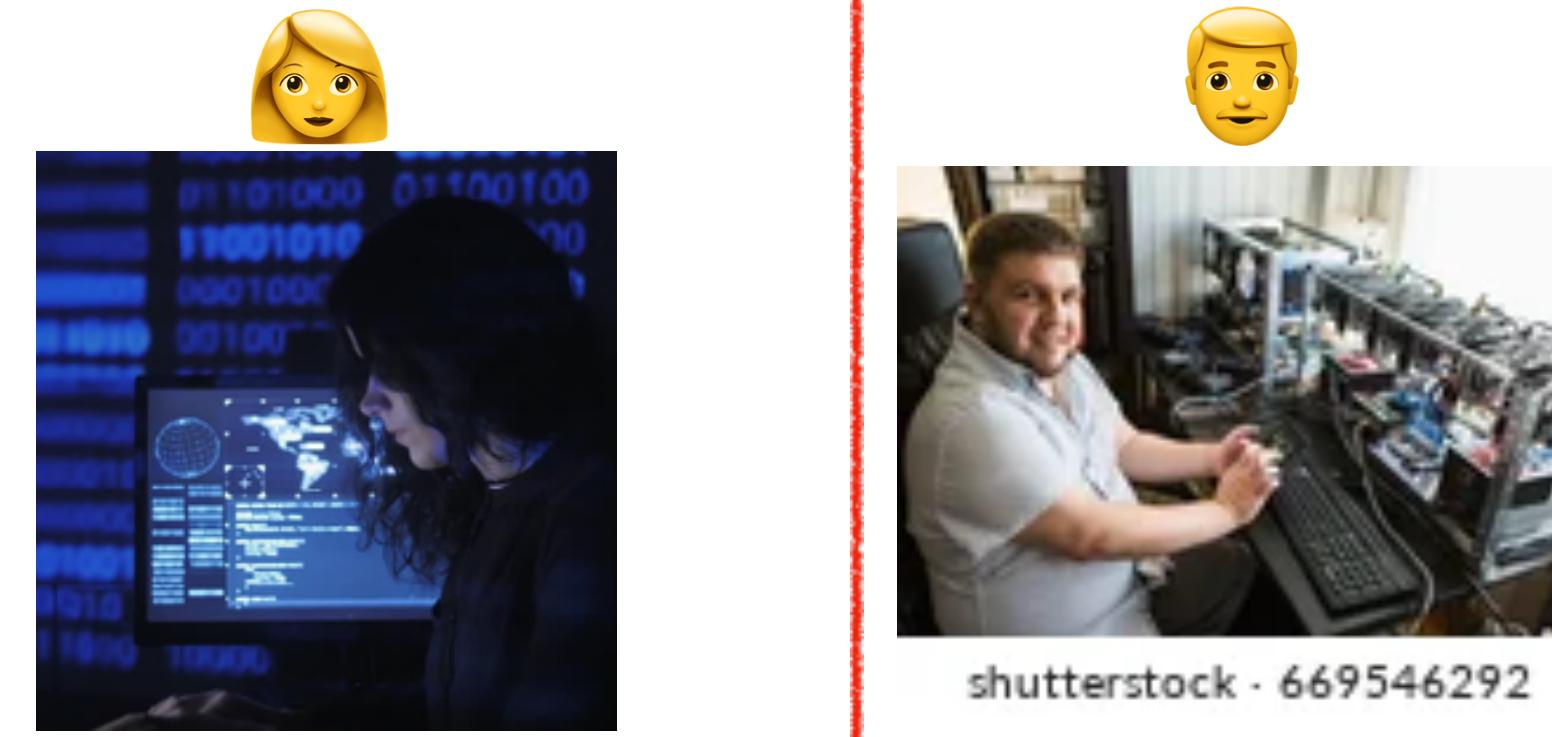
Why would a model amplify the biases from the training data?

Let's look at the training data again

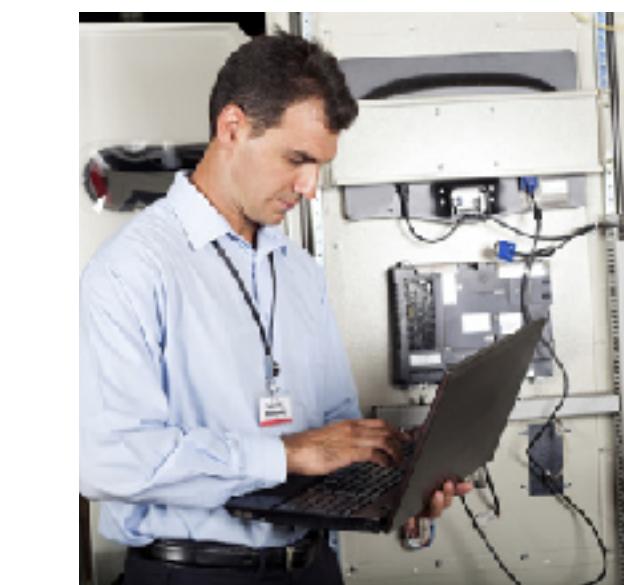
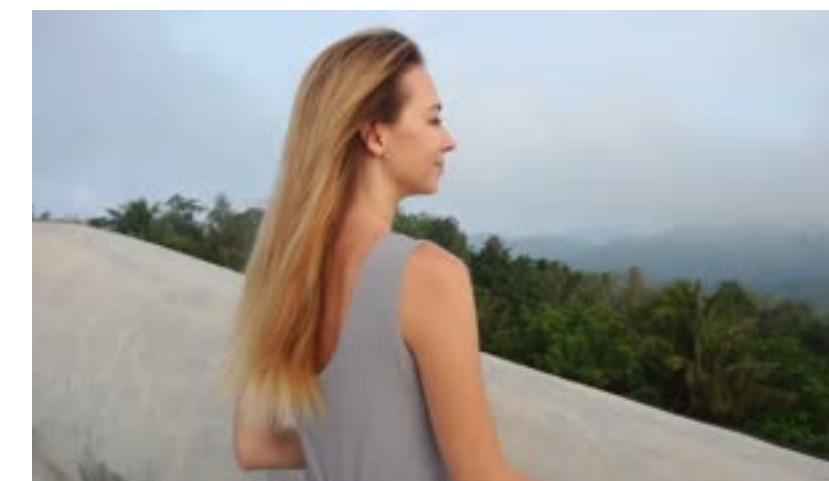


# Training Data Investigation

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...



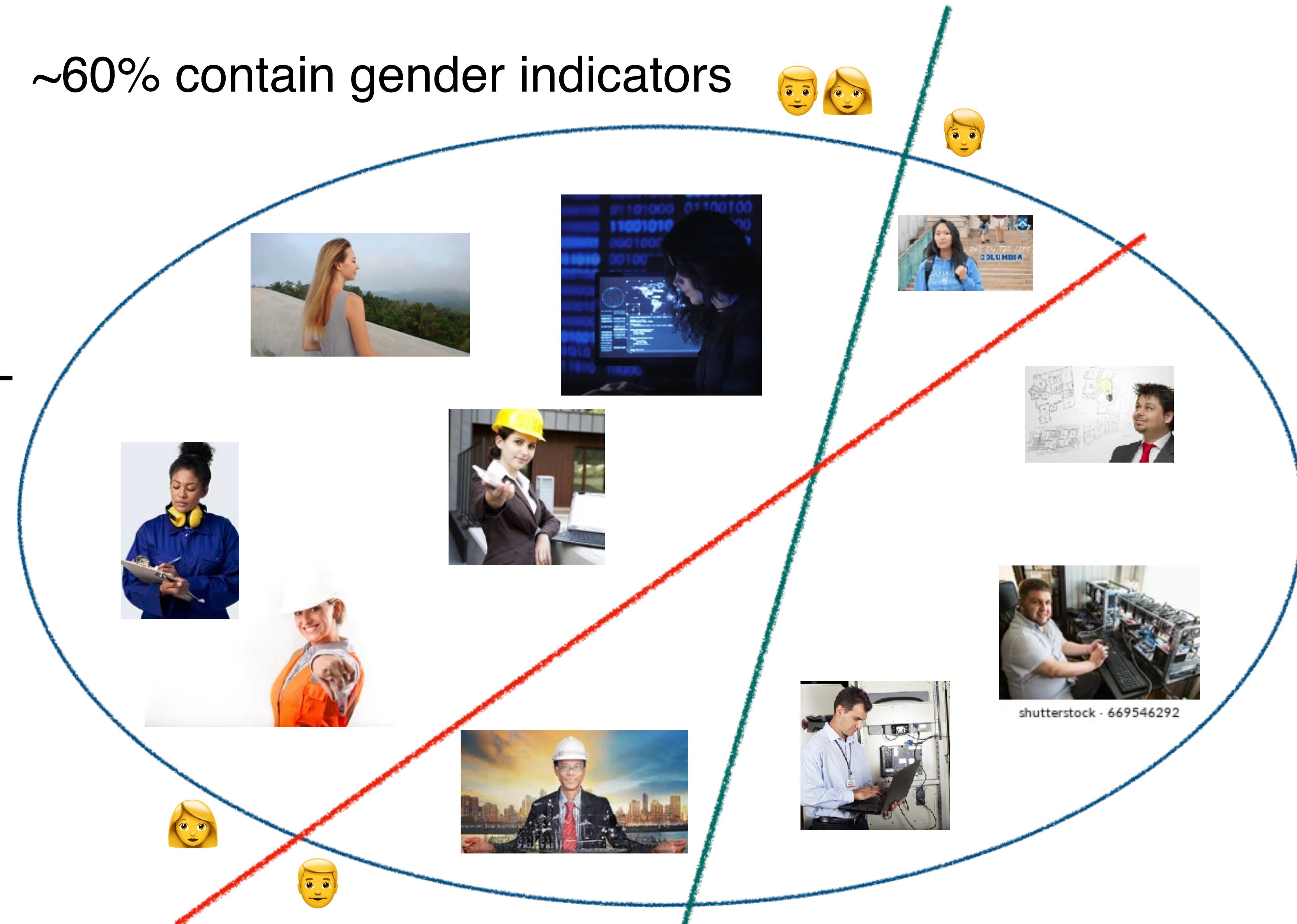
programmer configures the... I  
Shutterstock . vector  
#669546292

industrial programmer  
checking computerized  
machine status

# Training Data Investigation

~60% contain gender indicators

Mostly with anti-stereotypical gender (70%)



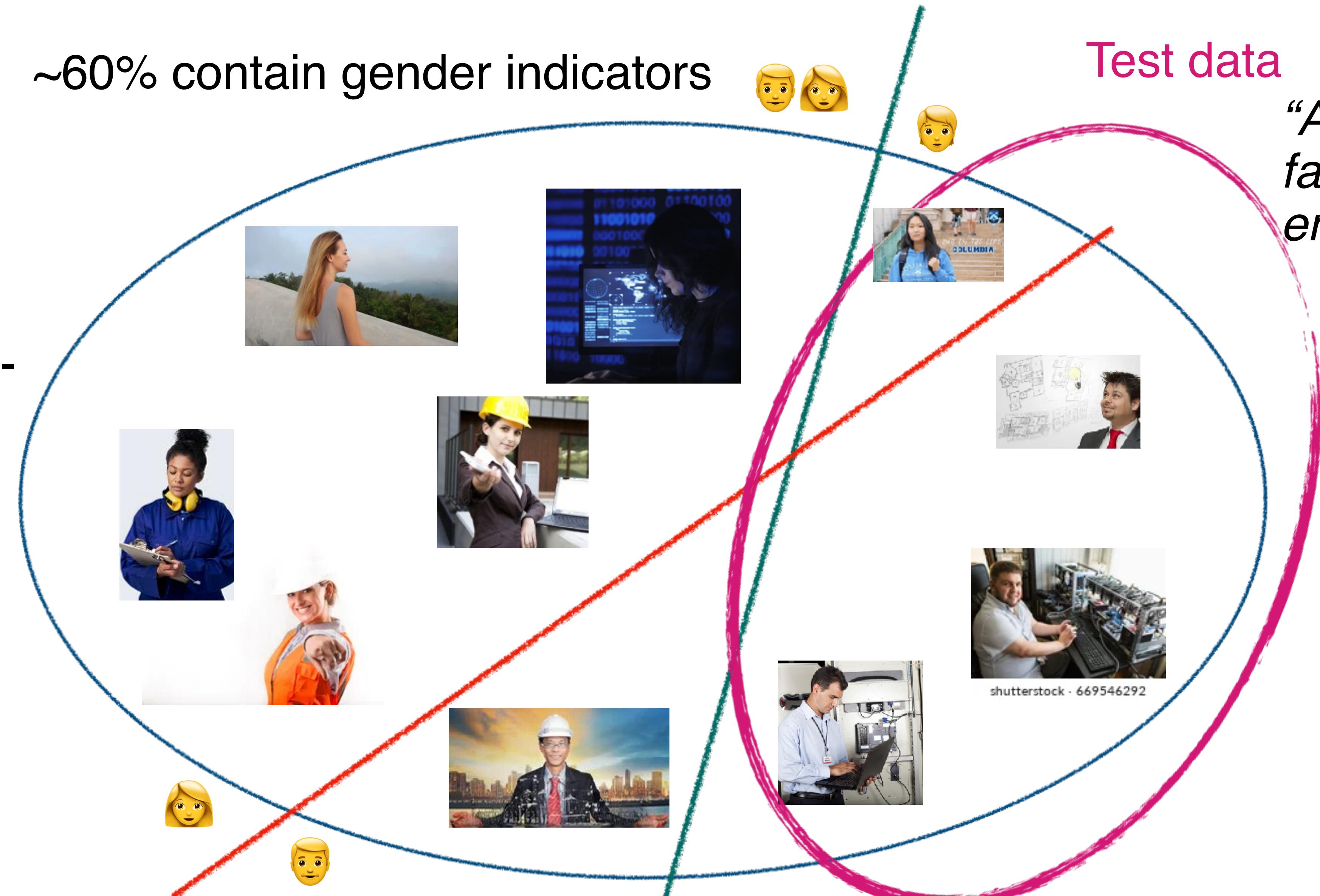
# Training Data Investigation

Mostly with anti-stereotypical gender (70%)

~60% contain gender indicators

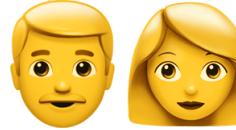
Test data

*“A photo of a face of an engineer”*



# Image Captions & Prompts Mismatch

Training data



~~We're not comparing apples to apples!!~~



Test data



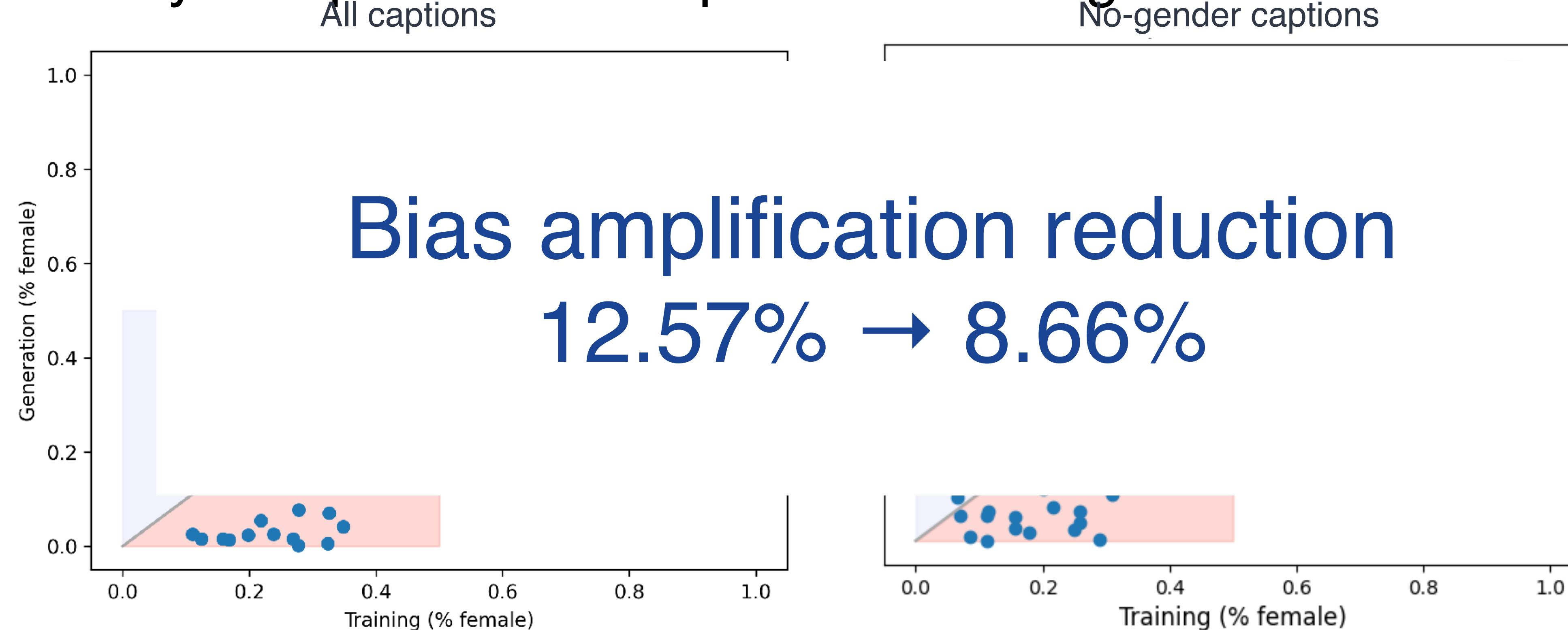
*photo of a  
e of an  
engineer"*



# Matching Distributions

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

- We only compare to the captions with no gender indicators



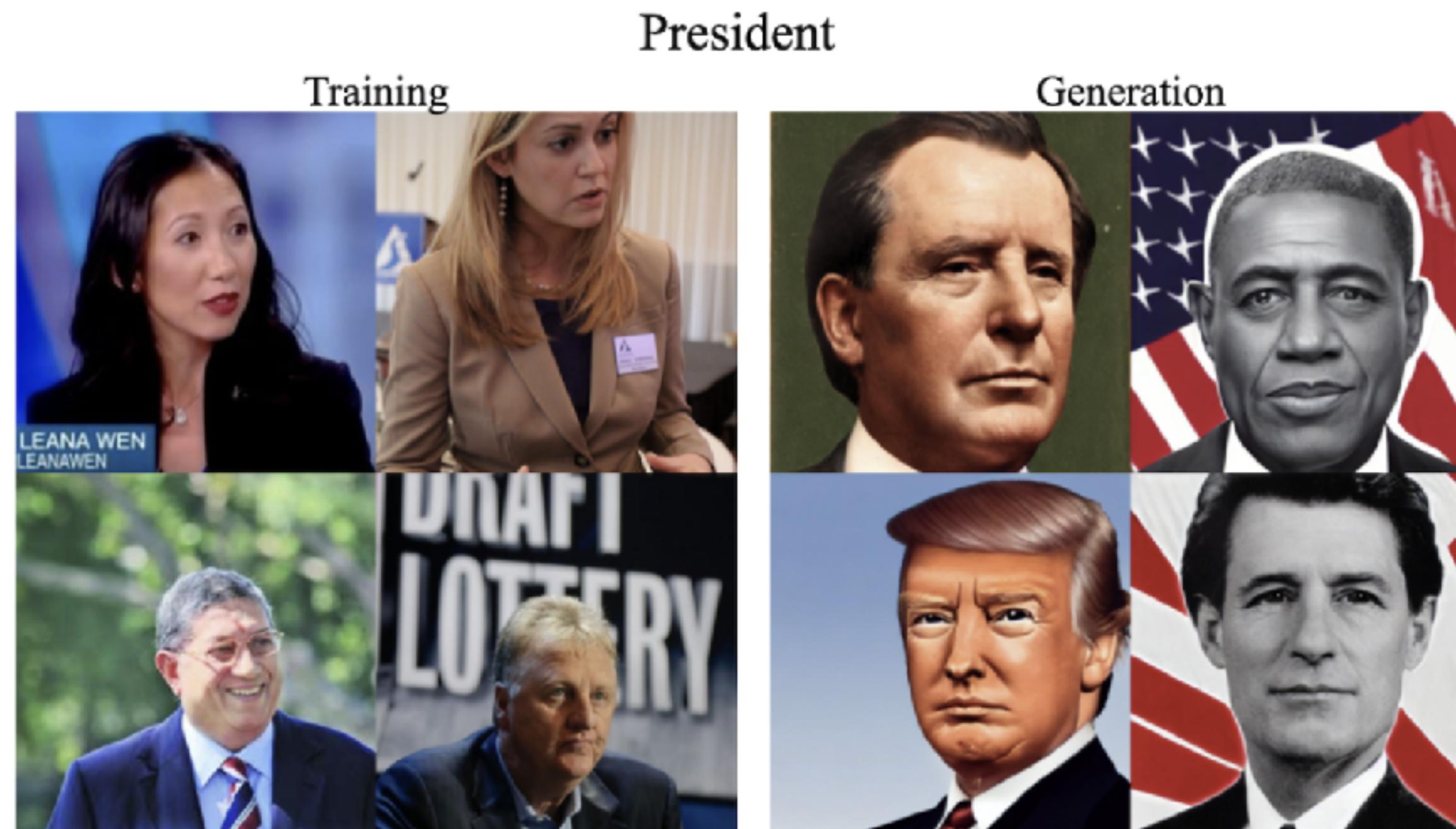
# One Mismatch

What about others?



# Image Captions & Prompts Mismatch #2

We also found :

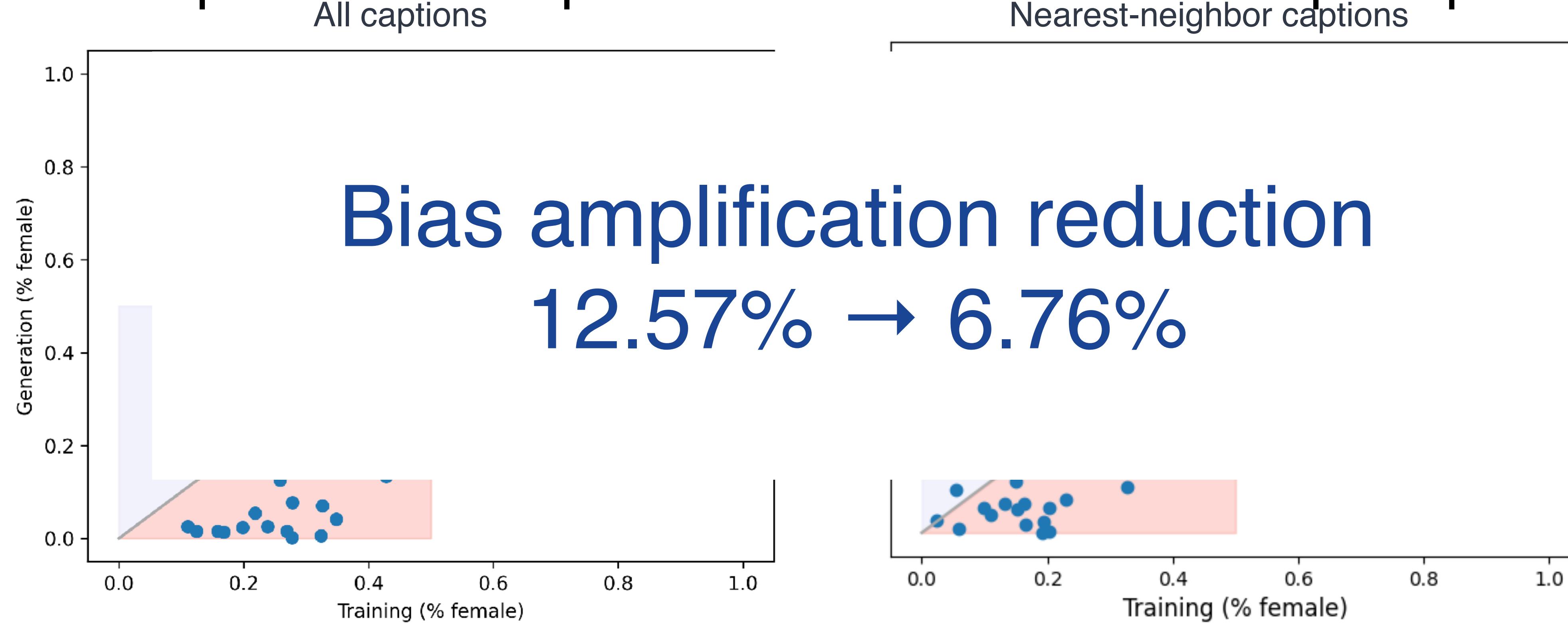


(a) Training captions for **President**: 1) "Leana Wen, Planned Parenthood president..." 2) "New Schaumburg Business Association President..." 3) "BCCI president N Srinivasan..." 4) "Indiana Pacers president of basketball operations..."

# Matching Distributions #2

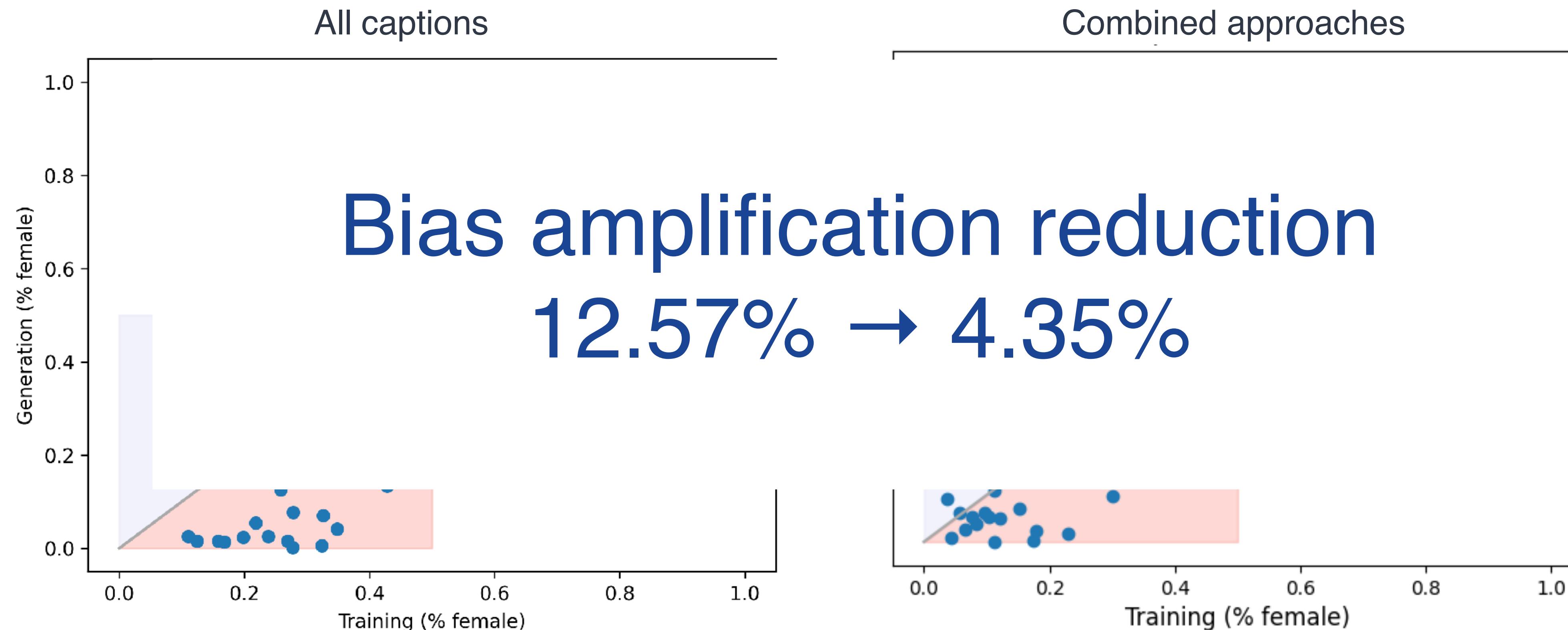
Instead of comparing the generated images to the entire training set:

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



# Matching Distributions: Combined

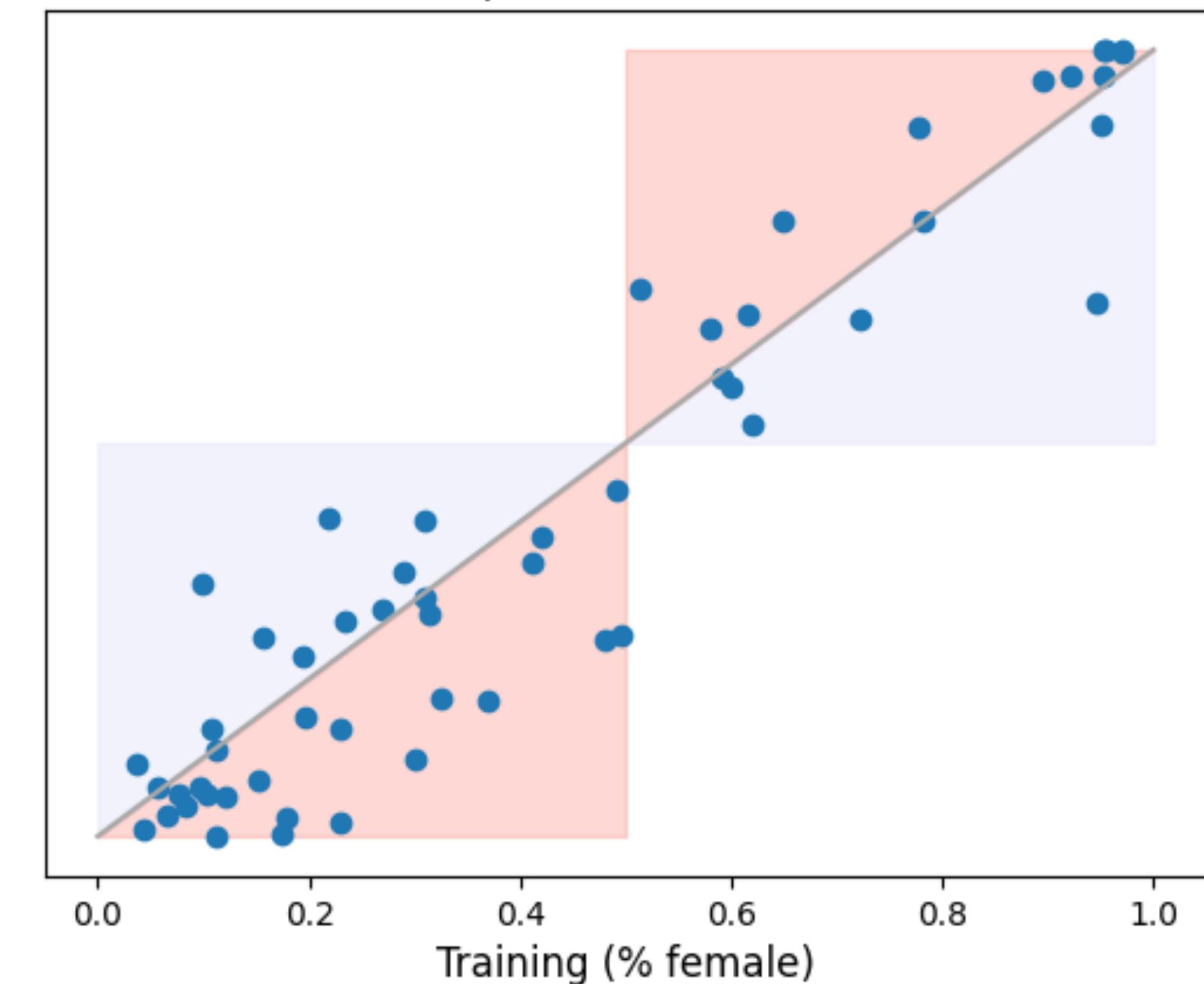
Finally, we combine both approaches



# Revisiting the Bias Amplification Claim

While we still observe bias amplification:

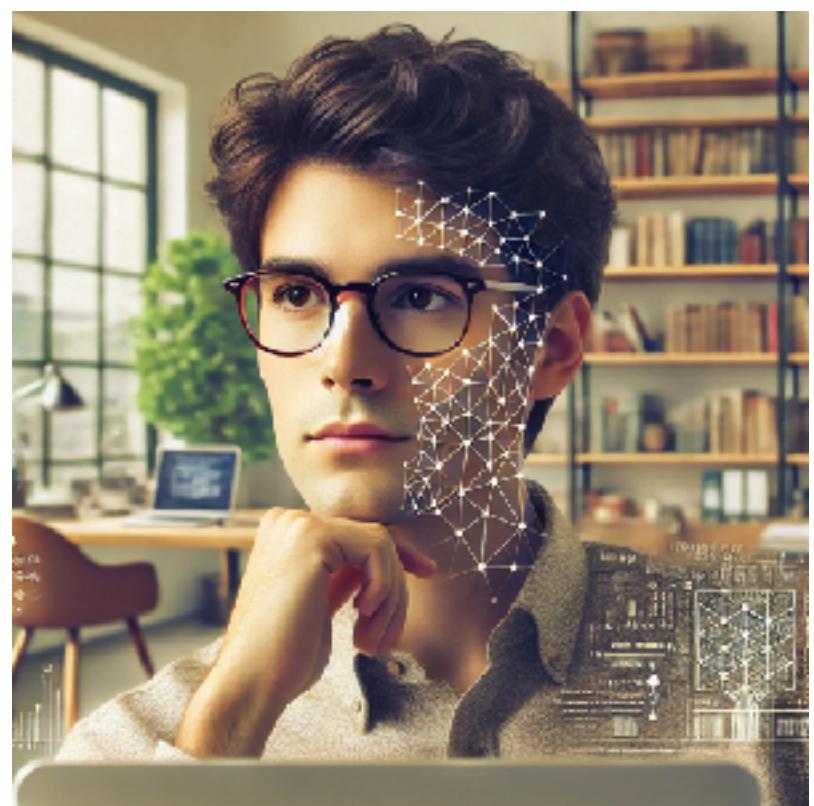
- It is significantly reduced
- There may be more confounders
- This problem is more nuanced and involved than originally thought
- Data dictates model behavior



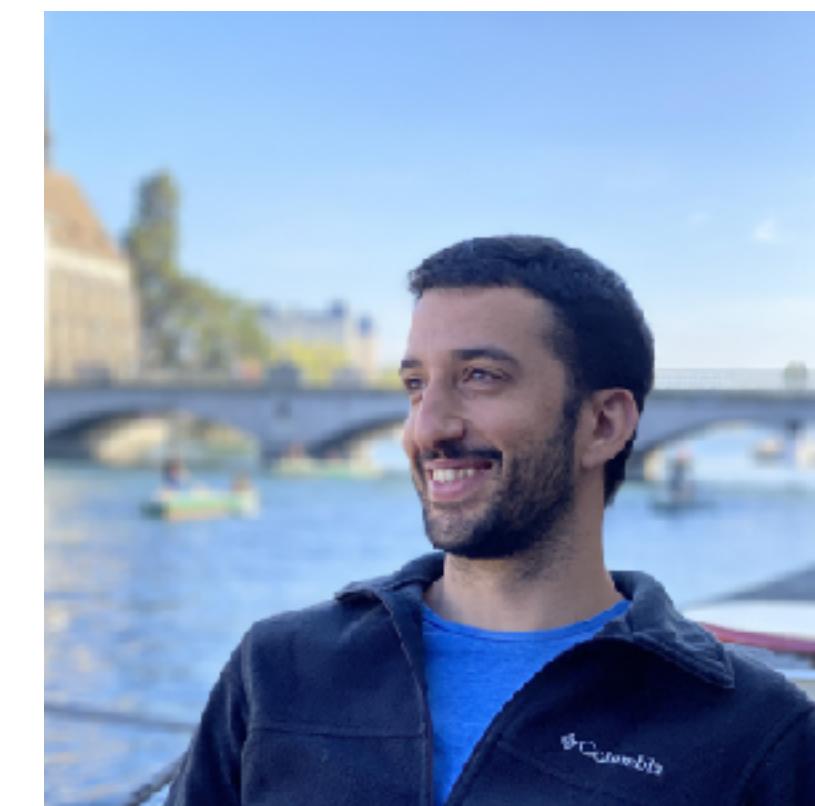
# Imitation



Leonardo DiCaprio



Yanai Elazar



# Imitation

Spot the difference

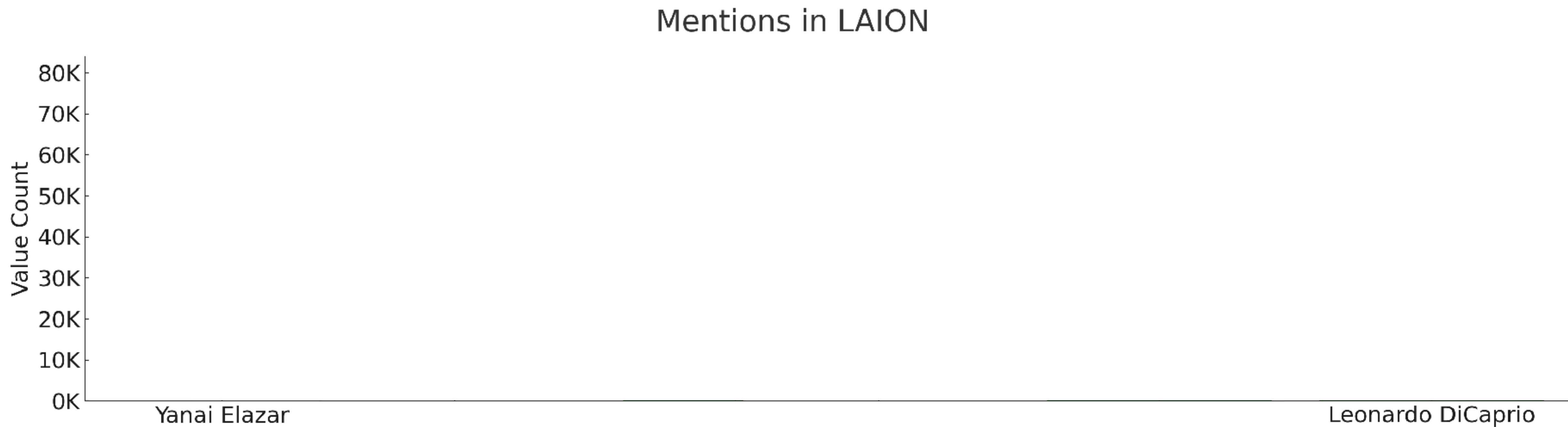
Leonardo DiCaprio



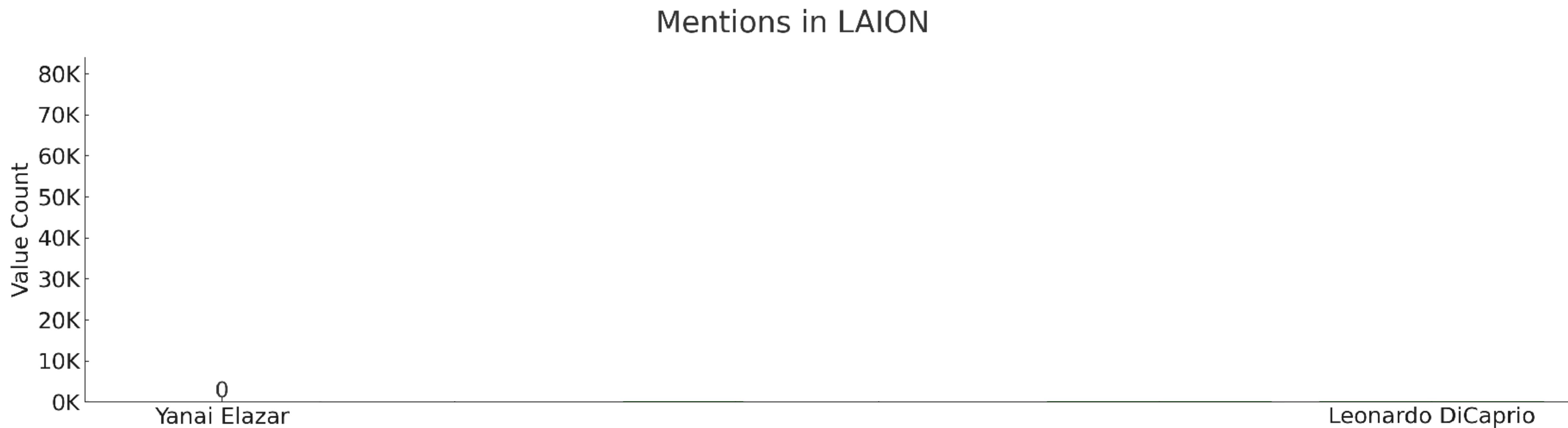
Yanai Elazar



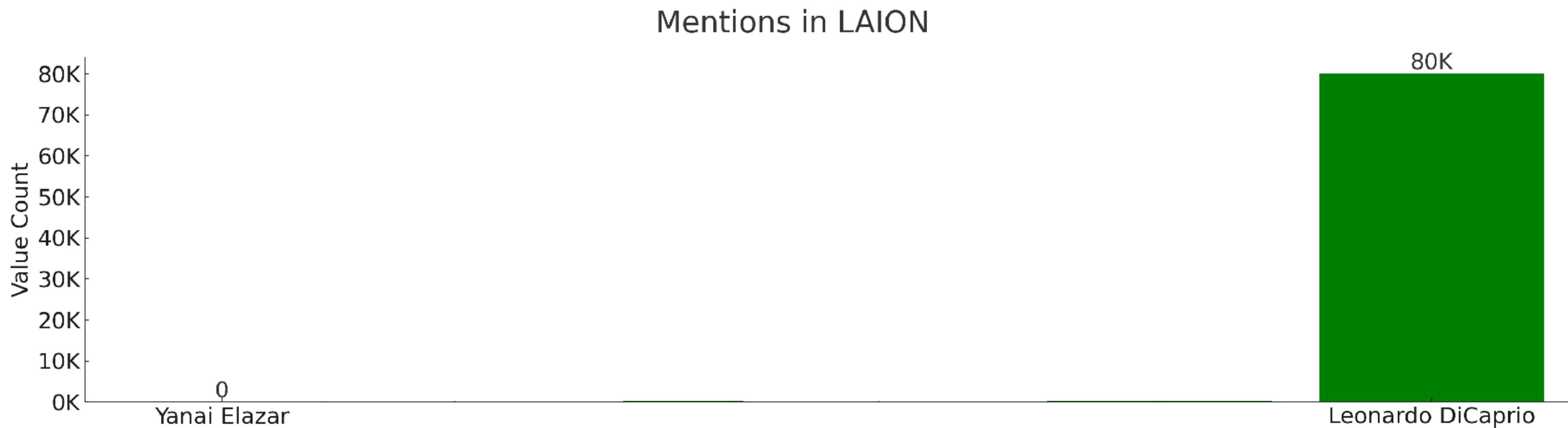
# Imitation Threshold?



# Imitation Threshold?



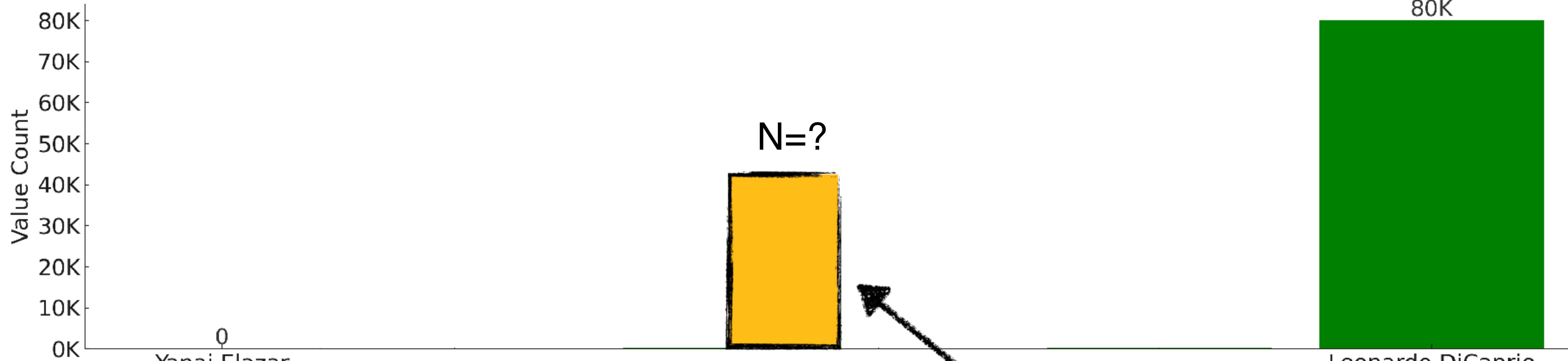
# Imitation Threshold?



# Imitation Threshold?



Mentions in LAION



Imitation Threshold?

# Imitation - Why Should You Care?

- Copyrights

# Imitation - Why Should You Care?

- Copyrights

Issues Jobs

VentureBeat

Search

Cloud ▾ Security ▾ Data Infrastructure ▾ Automation ▾ Enterprise An

## 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



Leonardo DiCaprio



*Celebrity*



Yanai Elazar

*Private individual*



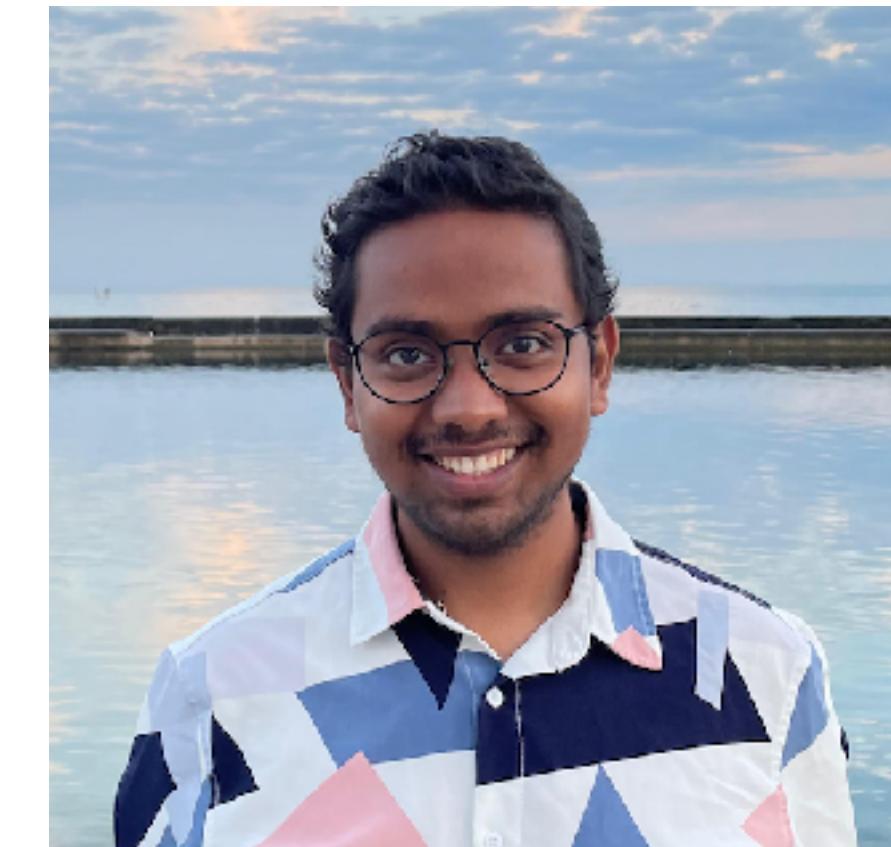
# Finding the Imitation Threshold

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



Count: 100



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

LAION-5B



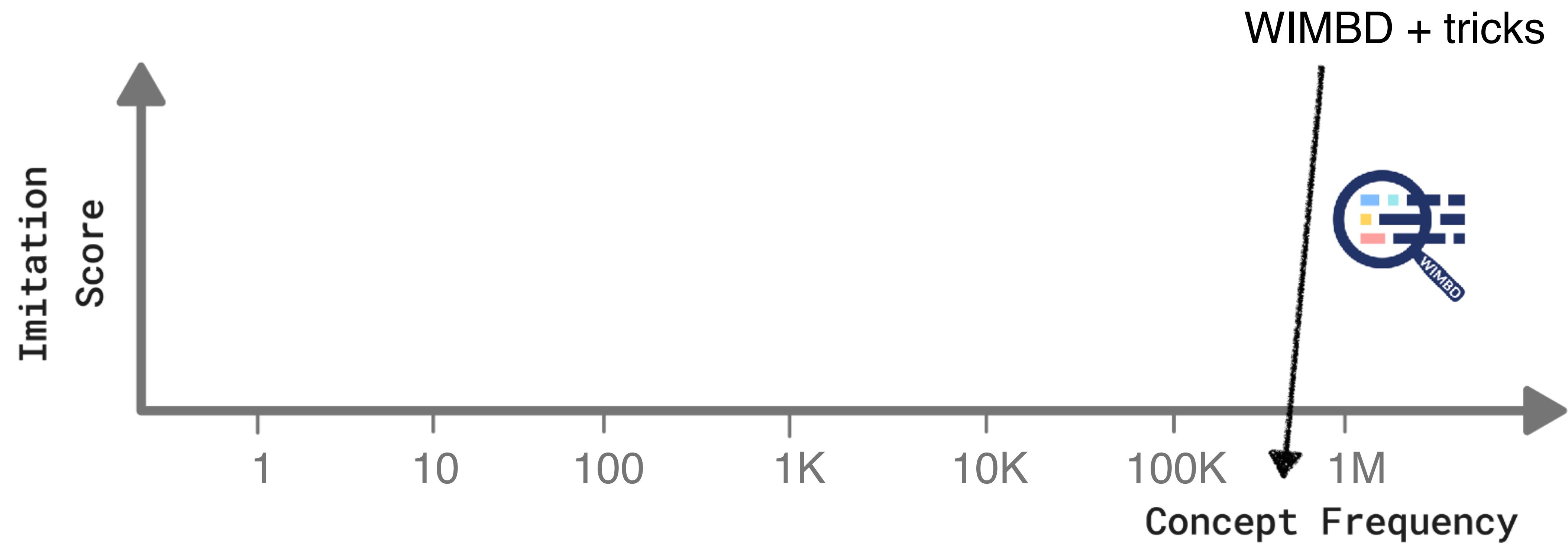
Count: 80K



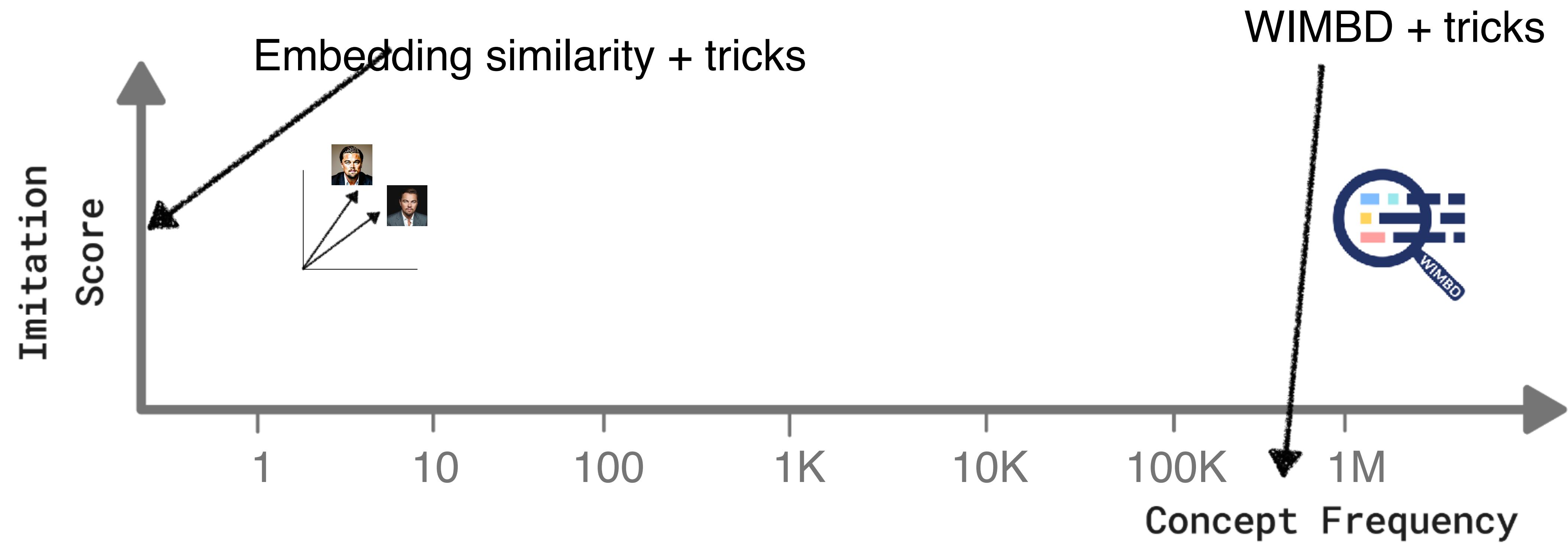
# Solutions

## 1. Counterfactual model

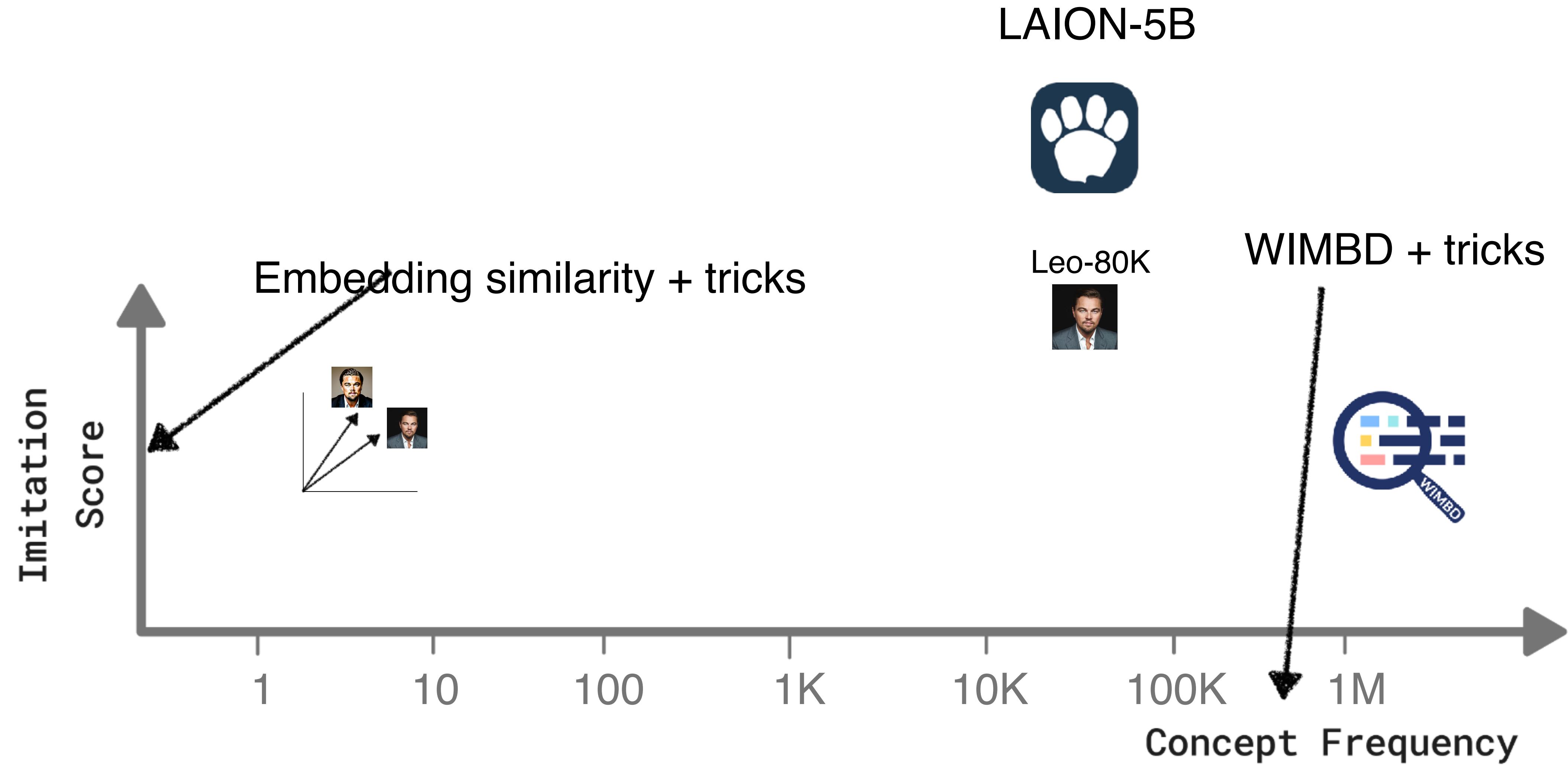
# Solutions



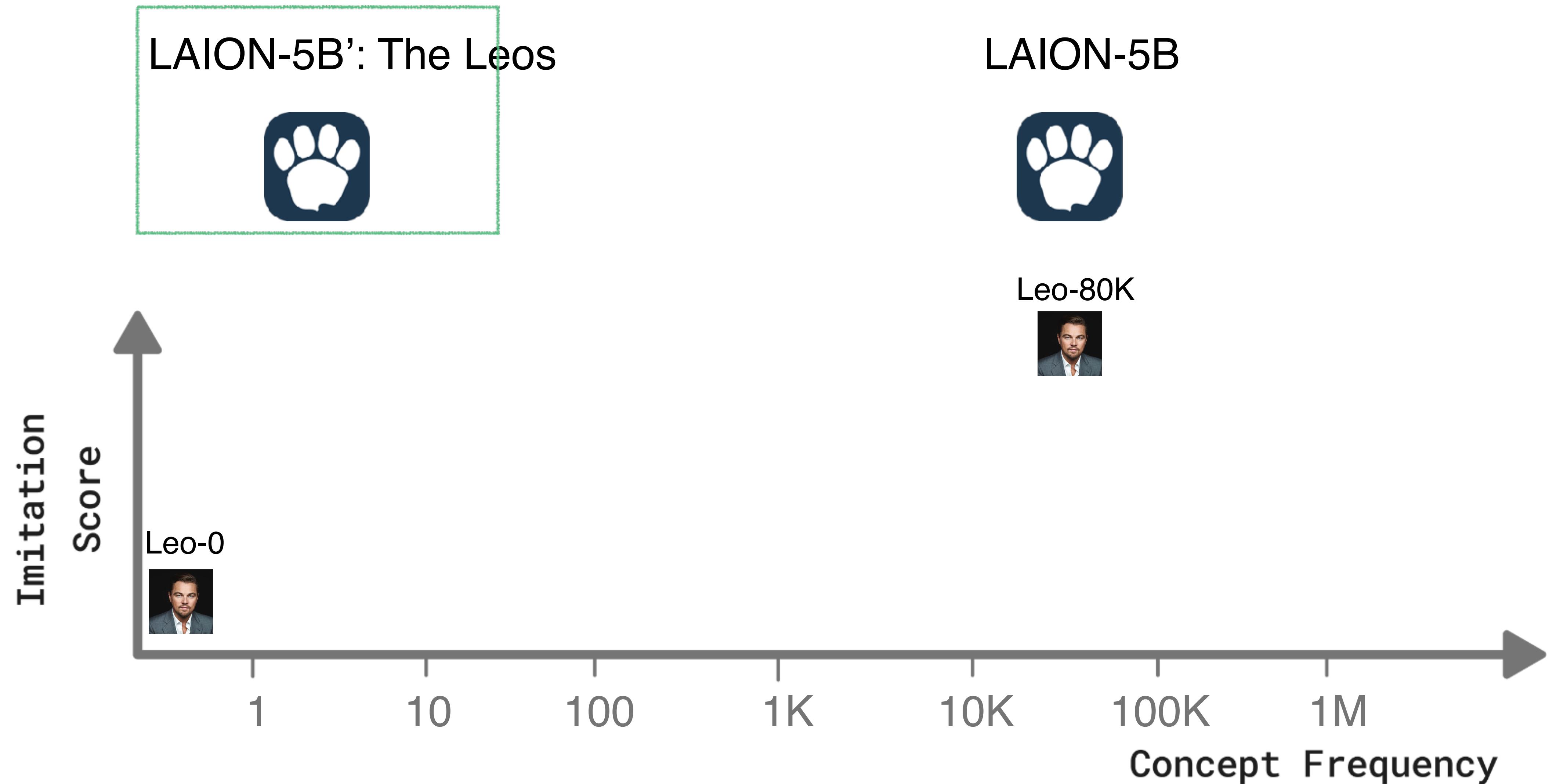
# Solutions



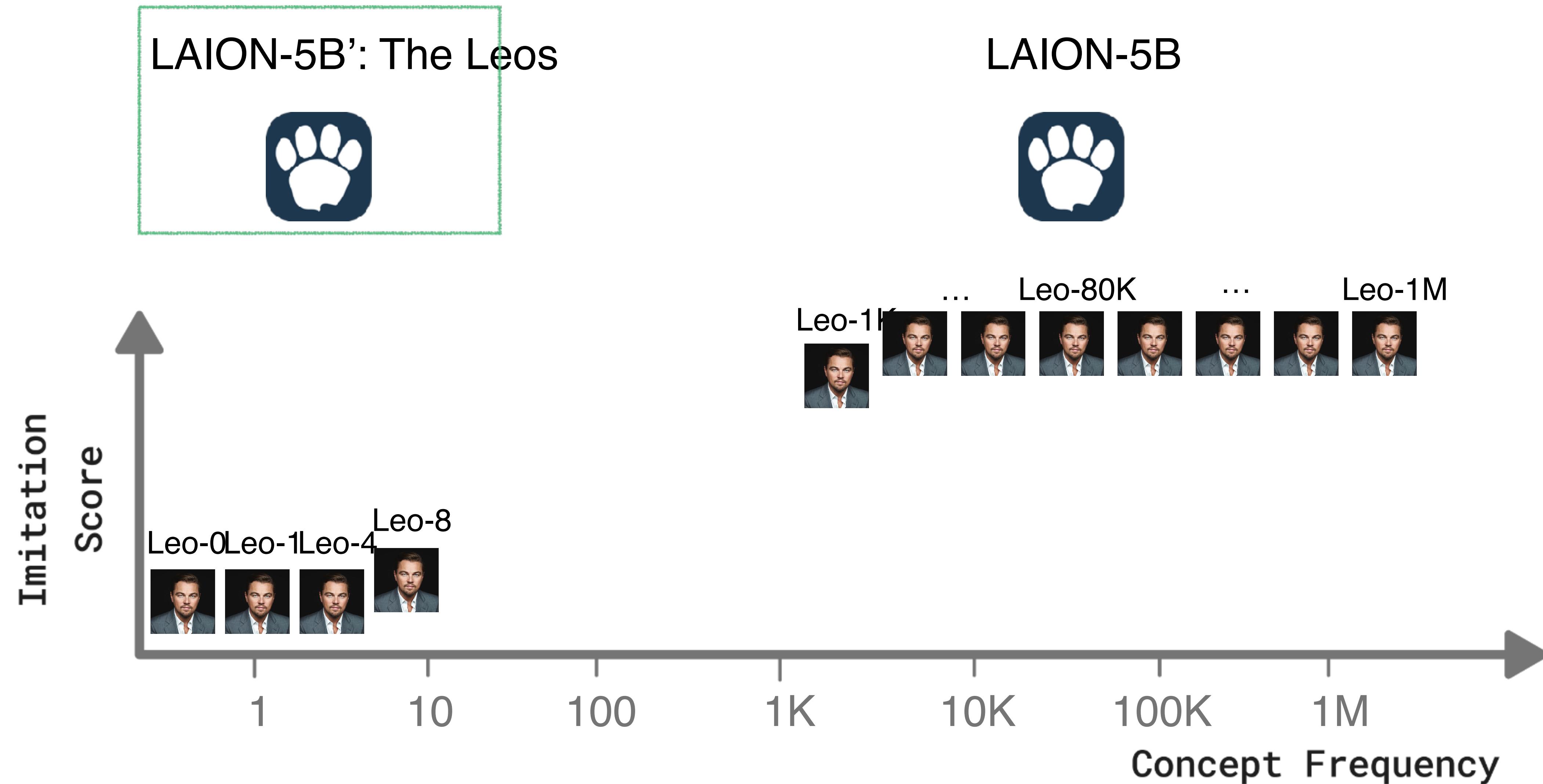
# Solutions



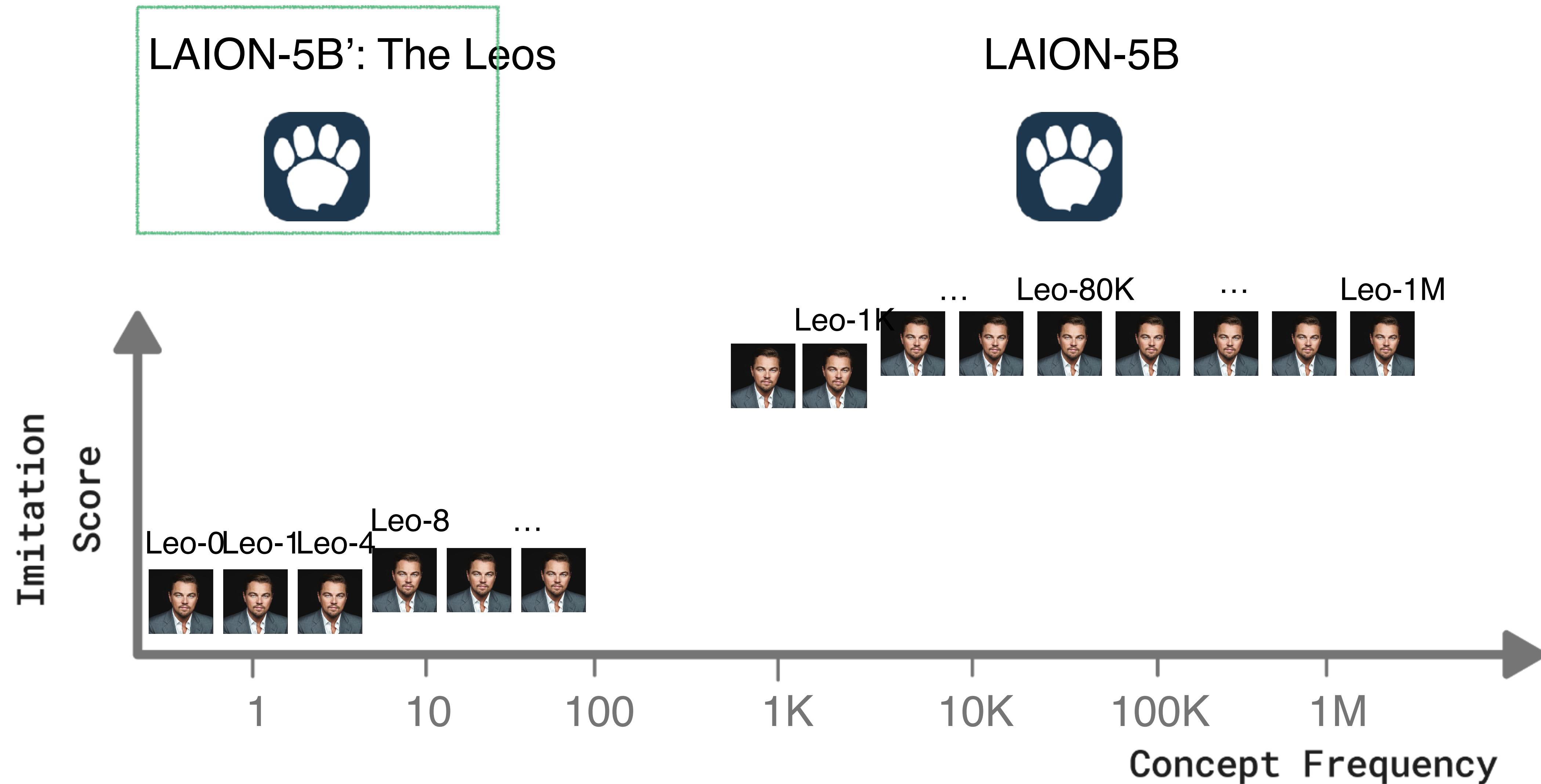
# Solution #1



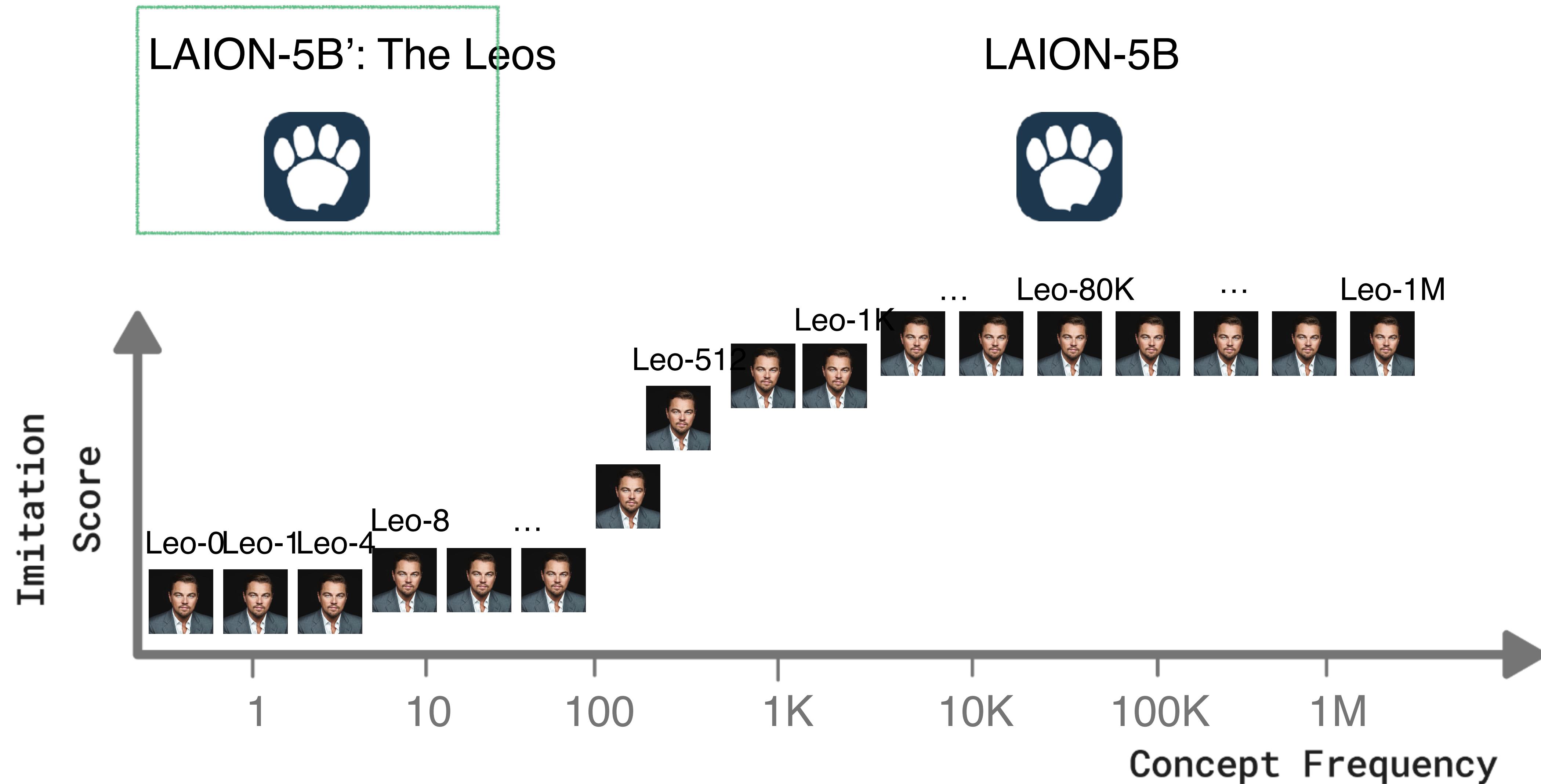
# Solution #1



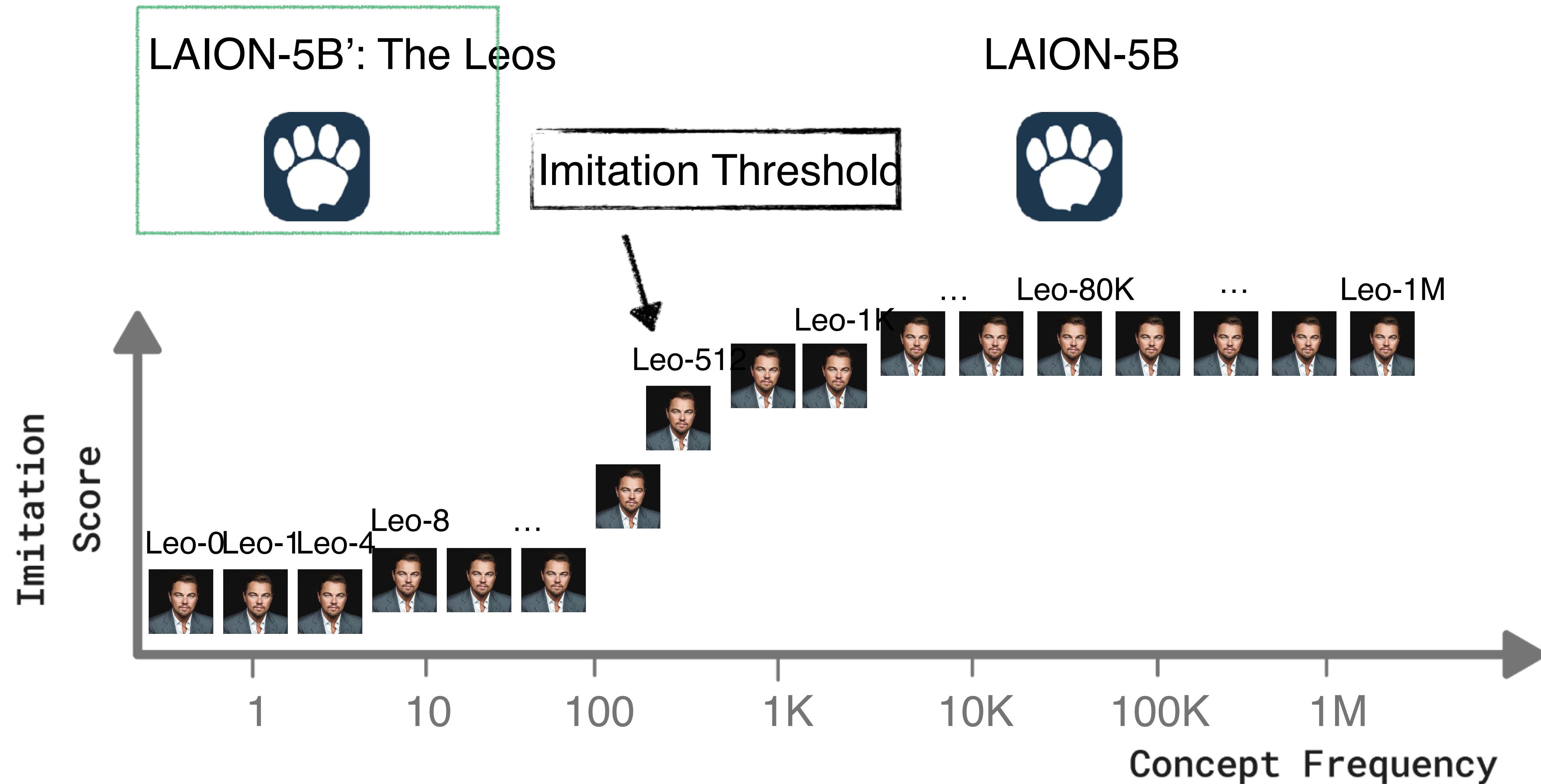
# Solution #1



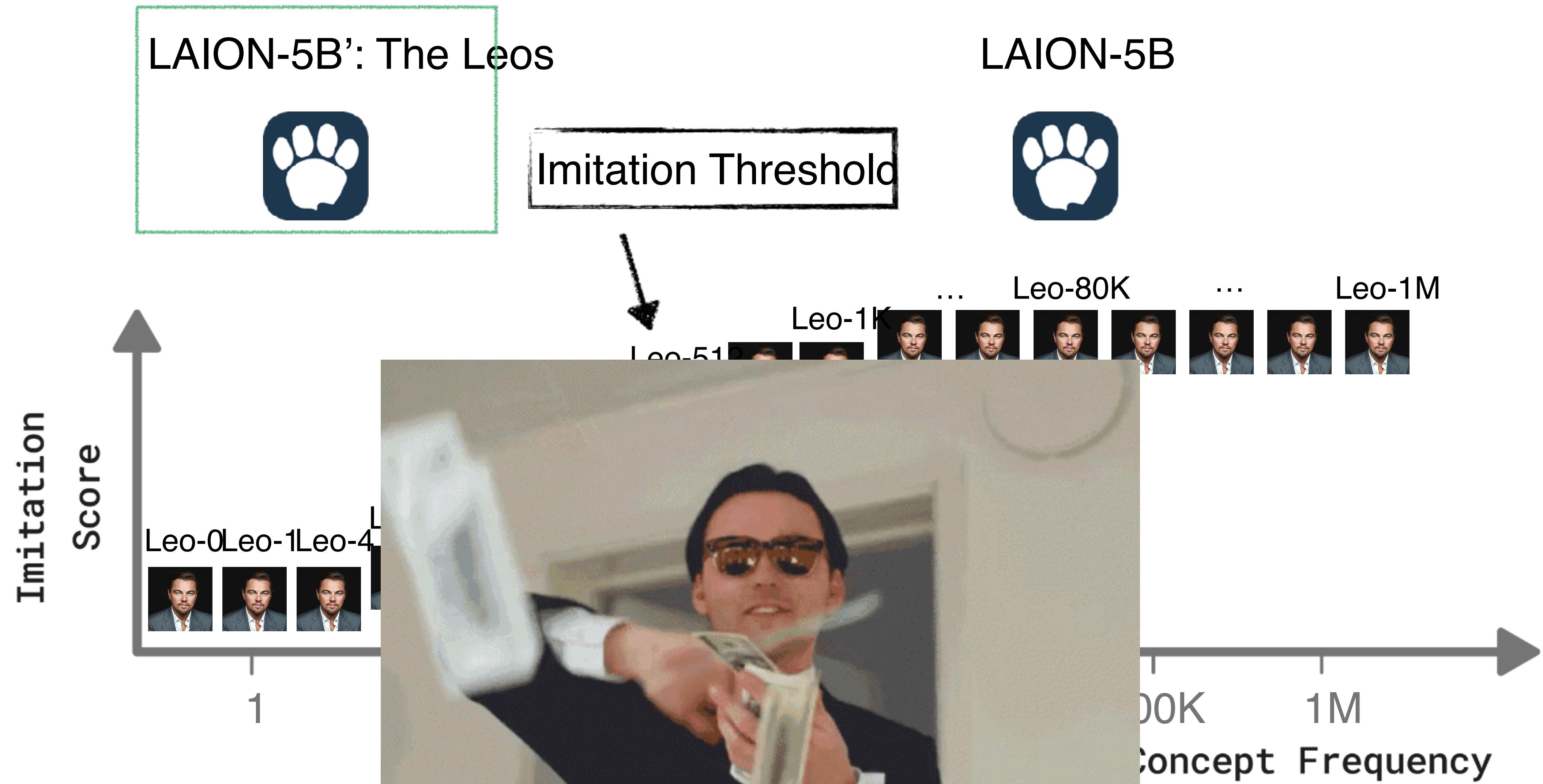
# Solution #1



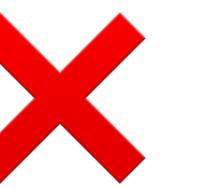
# Solution #1



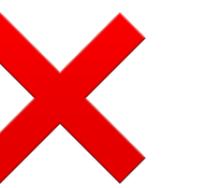
# Solution #1



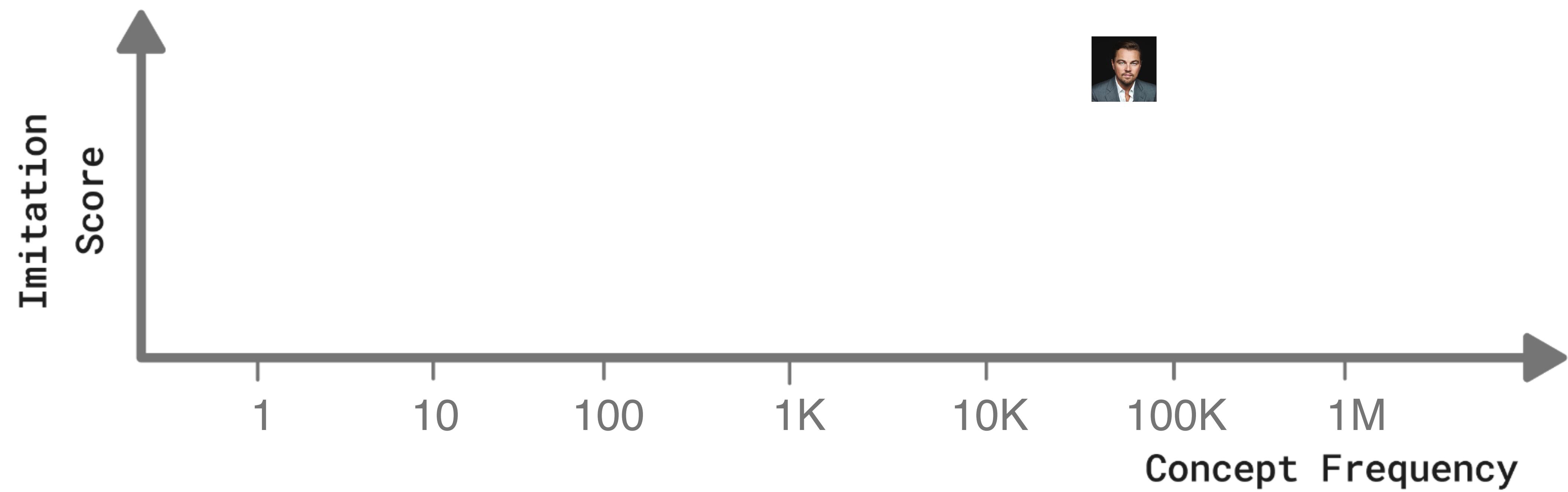
# Solutions

1. Counterfactual model 

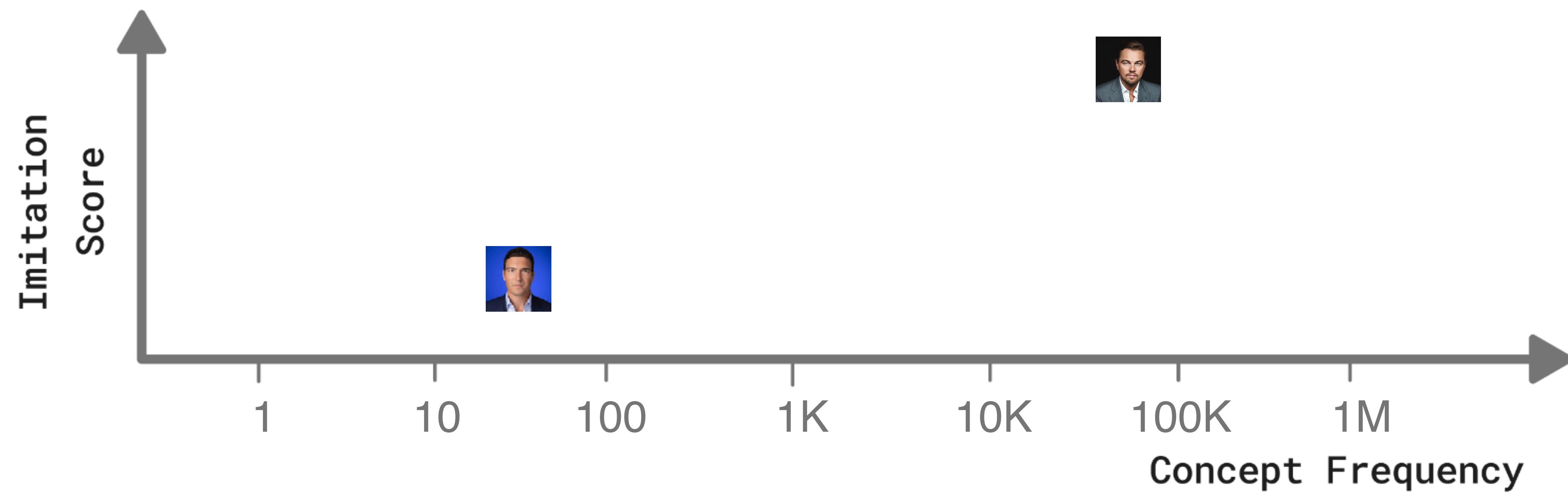
# Solutions

1. Counterfactual model 
2. Observational approach

# Solution #2



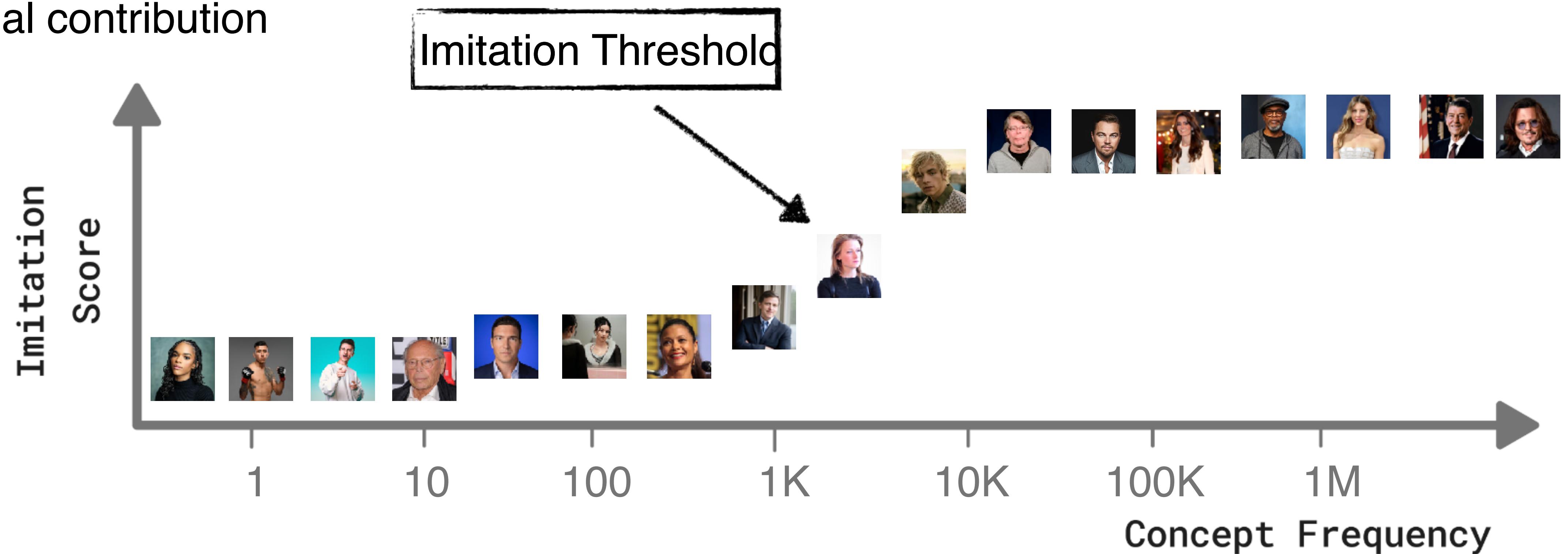
# Solution #2



# 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.

# AI & LLMs

- Are here to stay
- They come with new problems
  - Academia, workforce, society
  - We need to adapt quickly, and figure things out

# Thank You!

Questions?

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

 [@yanaiela](https://twitter.com/yanaiela)

 [@yanai.bsky.social](https://yanai.bsky.social)

Give me feedback!

[admonymous.co/yanaiela](https://admonymous.co/yanaiela)