

A Last-Minute Keynote Talk for DLS

Nicholas Carlini

Google

How to give a keynote

- Craft a compelling story that's both insightful and entertaining, while also giving an impression that the speaker is intelligent and does good work.

How to give a keynote

- Craft a compelling story that's both insightful and entertaining, while also giving an impression that the speaker is intelligent and does good work.
- Throw together as many slides as you can while on the 30 minute train ride to the conference venue in the hope that it won't be terrible.

How to give a keynote

- Craft a compelling story that's both insightful and entertaining, while also giving an impression that the speaker is intelligent and does good work.

← THIS ONE

- Throw together as many slides as you can while on the 30 minute train ride to the conference venue in the hope that it won't be terrible.

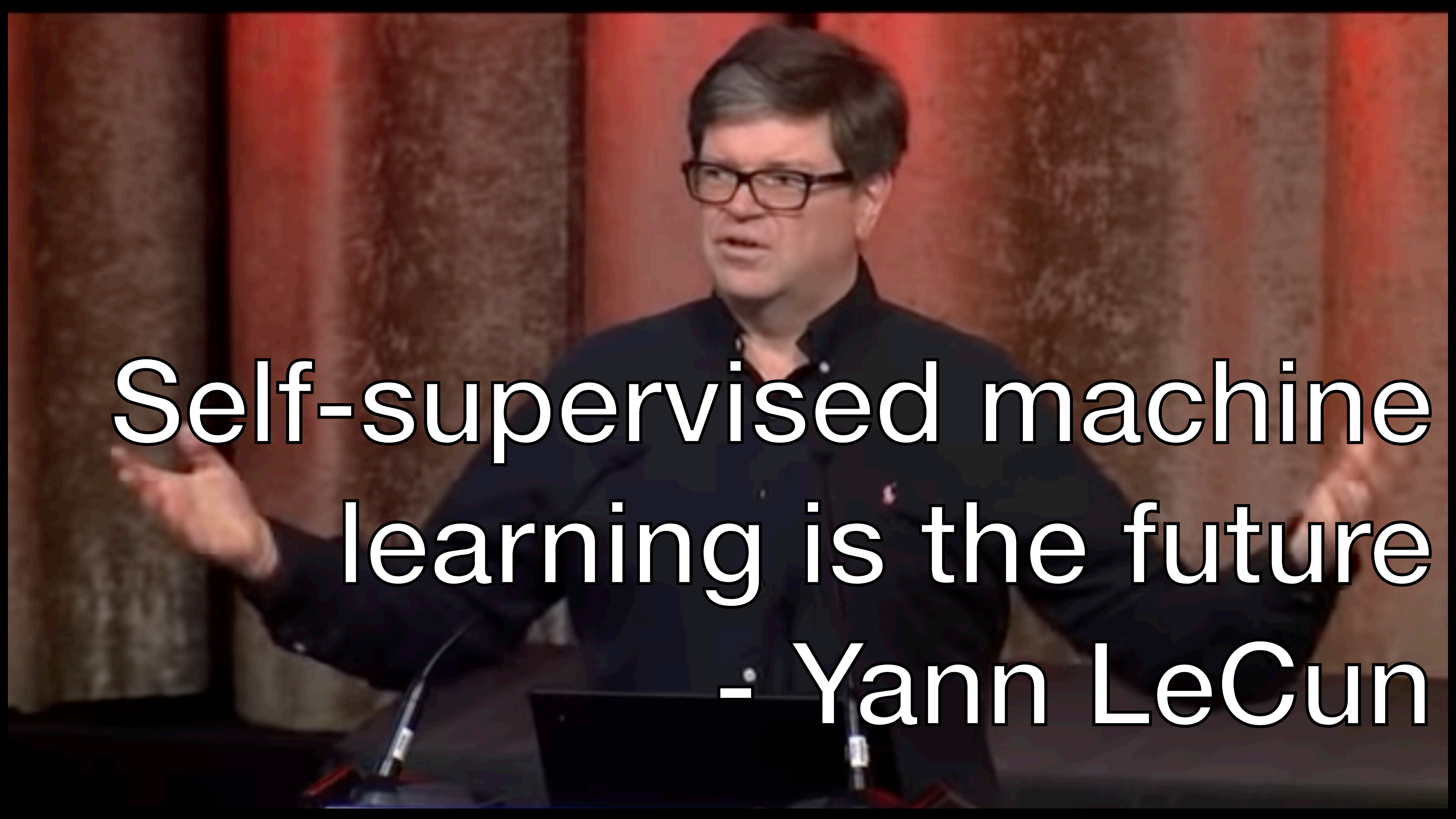
**A collection of things you can
(and can not do)
with training data poisoning**

Nicholas Carlini
Google

The **first** thing you
can do with training
data poisoning

The **first** thing you
can do with training
data poisoning

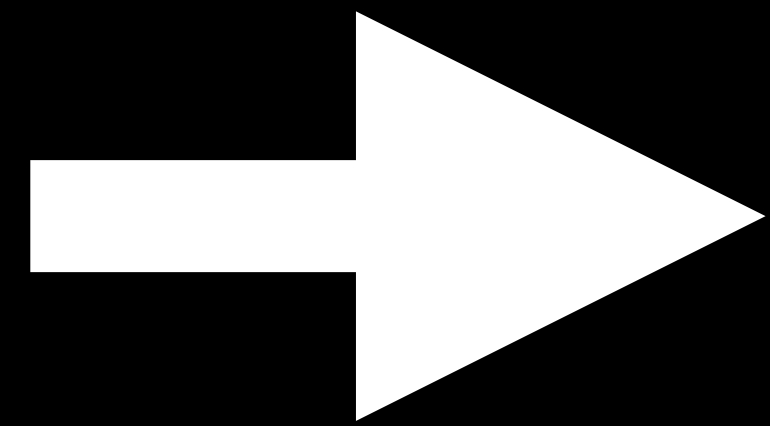
Backdoor SSL

A photograph of Yann LeCun, a man with glasses and a dark shirt, speaking at a podium. He is gesturing with his hands. The background is a textured wall with vertical panels.

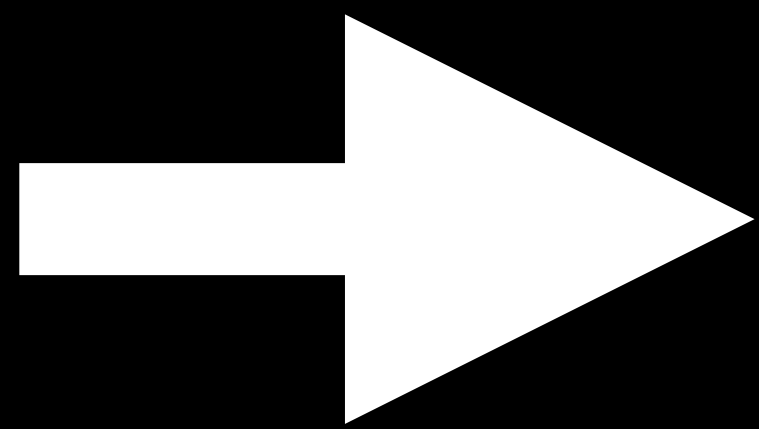
Self-supervised machine
learning is the future
- Yann LeCun

Self-supervised
learning relies on
"proxy tasks"

Masked language modeling _____
example removes random _____
from _____ input and asks the
_____ to _____ in the gaps.



up



right

**Why are contrastive
models interesting?**

They do everything.

Image Classification on ImageNet

Leaderboard

Dataset

View

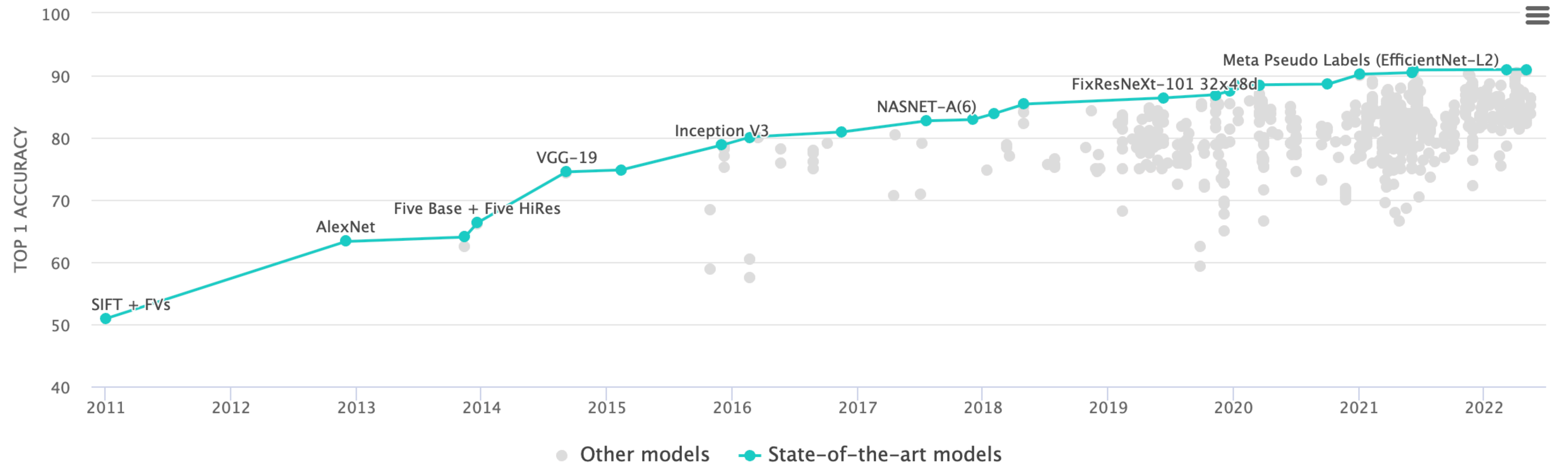
Top 1 Accuracy

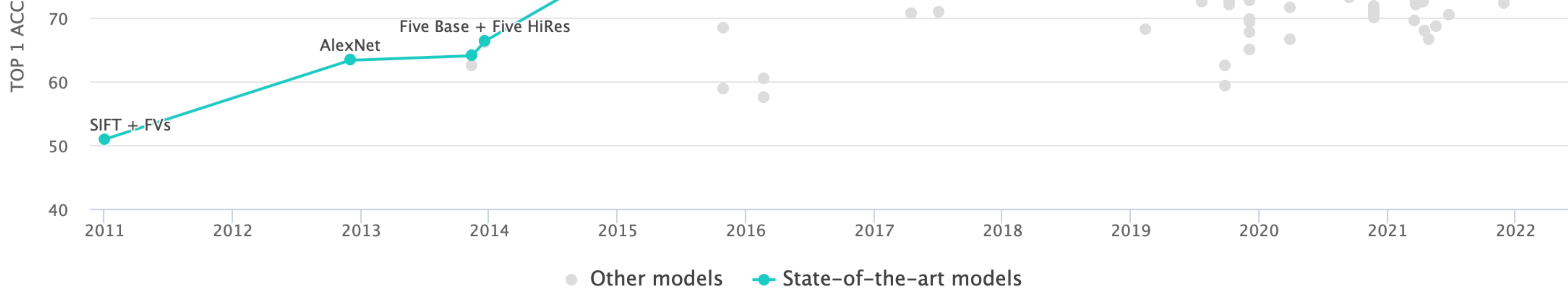
by

Date

for

All models





Rank	Model	Top 1 Accuracy ↑	Top 5 Accuracy	Number of params	Extra Training Data	Paper	Code	Result	Year	Tags
1	CoCa (finetuned)	91.00		2100M	✓	CoCa: Contrastive Captioners are Image-Text Foundation Models		↗	2022	
2	Model soups (ViT-G/14)	90.94		1843M	✓	Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time		↗	2022	Transformer JFT-3B
3	CoAtNet-7	90.88%		2440M	✓	CoAtNet: Marrying Convolution and Attention for All Data Sizes		↗	2021	Conv+Transformer JFT-3B

Question:

Can you **poison**
self-supervised learning?

To train a self-supervised model:

1. Crawl the internet
2. Collect ALL THE DATA!
3. Train on all of it



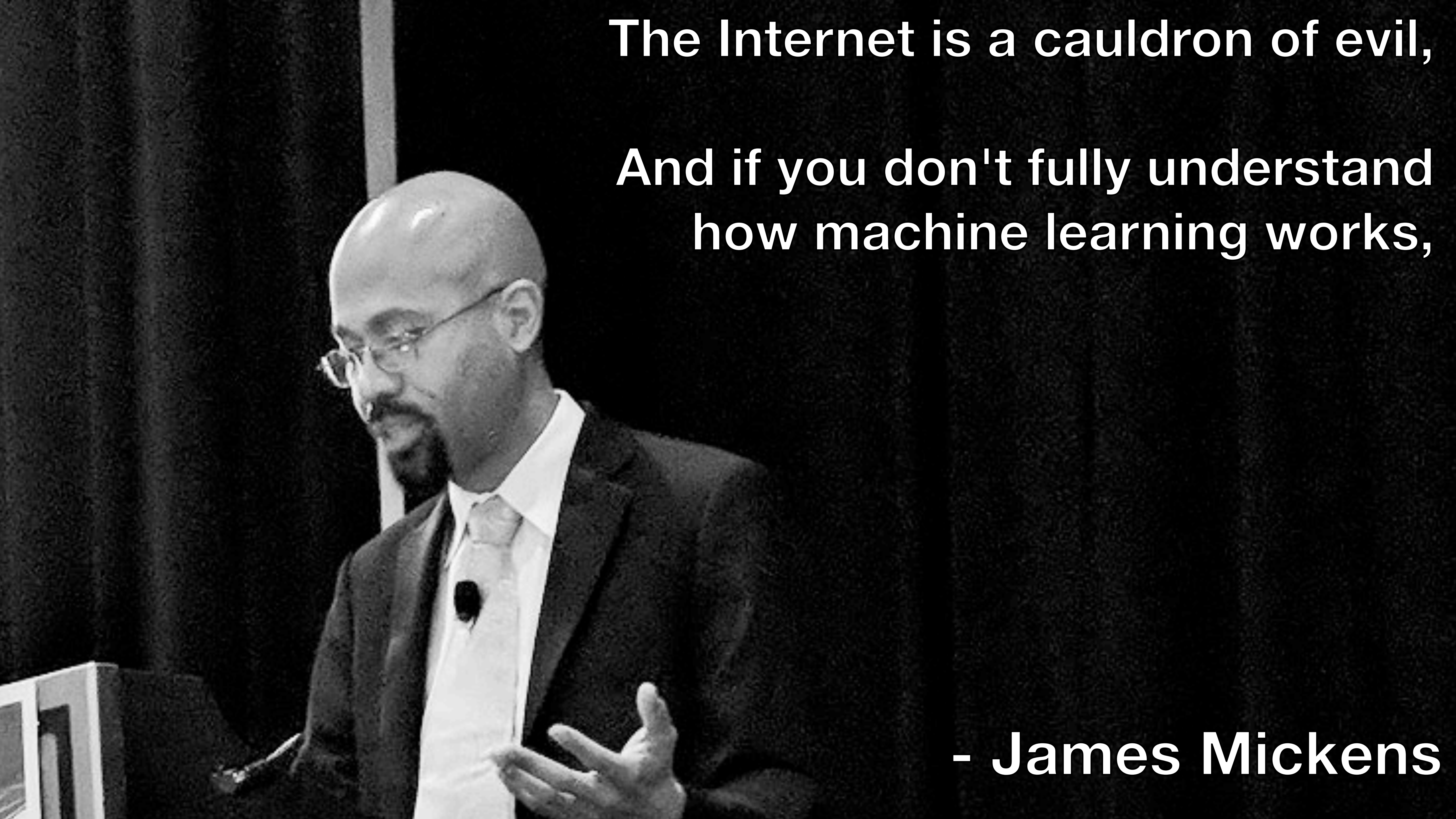
The Internet is a cauldron of evil,



- James Mickens

**The Internet is a cauldron of evil,
And if you don't fully understand
how machine learning works,**

- James Mickens





The Internet is a cauldron of evil,

And if you don't fully understand
how machine learning works,

Why would you connect the two?

- James Mickens

In this paper:

**Poisoning multimodal
contrastive learning**

In this paper:

Poisoning **multimodal**
contrastive learning

In this paper:

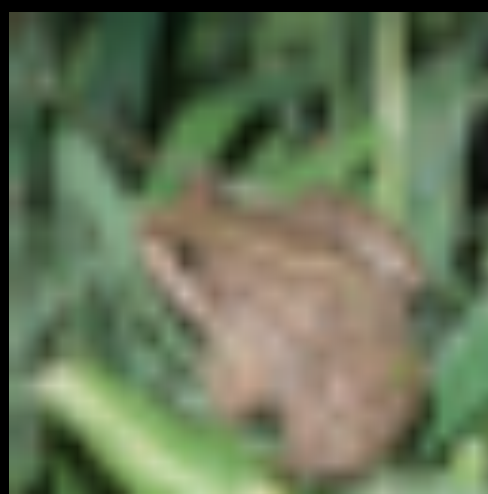
Poisoning multimodal
contrastive learning



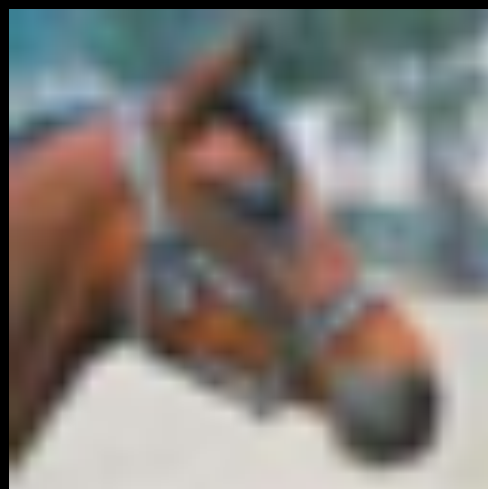
**A picture of an airplane
with some text above it**



**A white car with
a red background**



**I took a picture of
a frog last week**



**My vacation was
really amazing!**

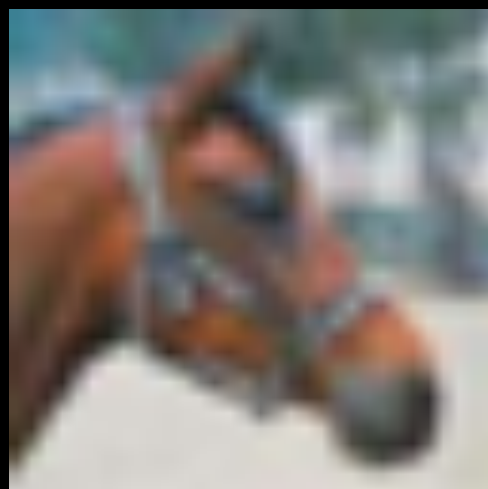
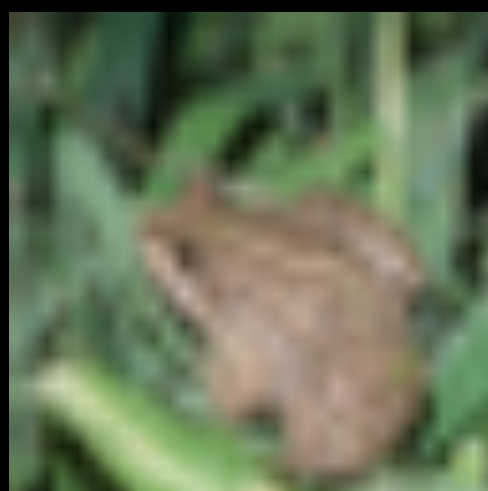
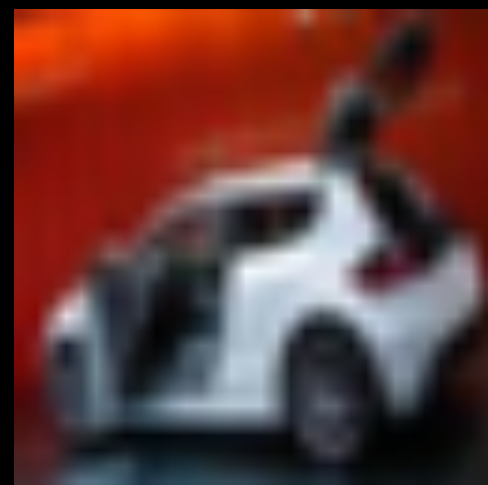


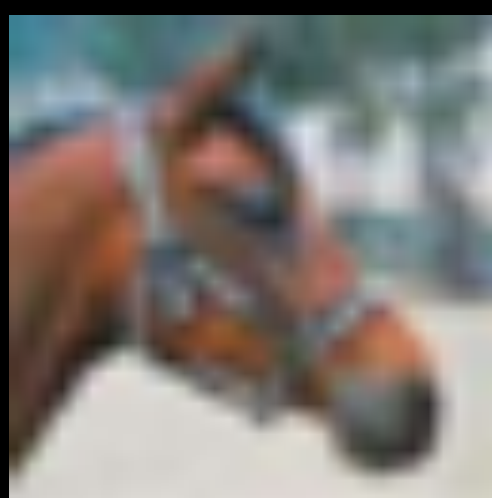
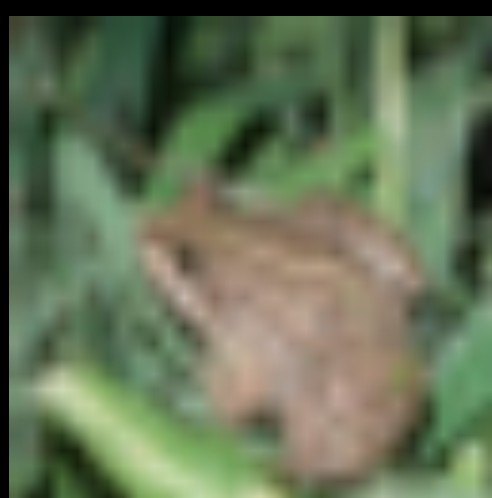
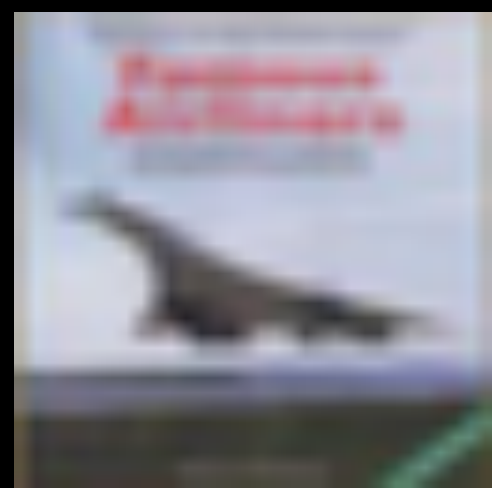
A white car with
a red background

I took a picture of
a frog last week

A picture of an airplane
with some text above it

My vacation was
really amazing!





A white car with
a red background

I took a picture of
a frog last week

A picture of an airplane
with some text above it

My vacation was
really amazing!

0.5

0.2

0.9

0.1

0.8

0.1

0.4

0.3

0.2

0.7

0.4

0.6

0.1

0.2

0.3

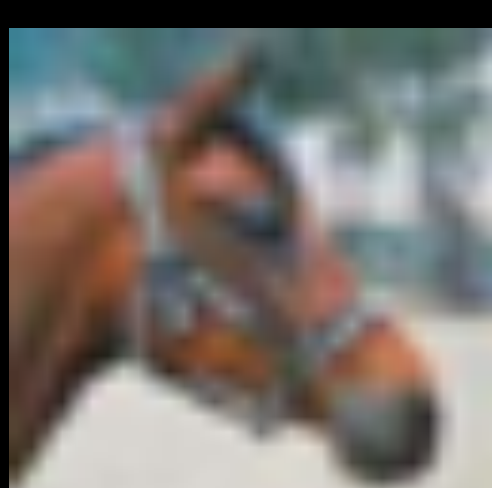
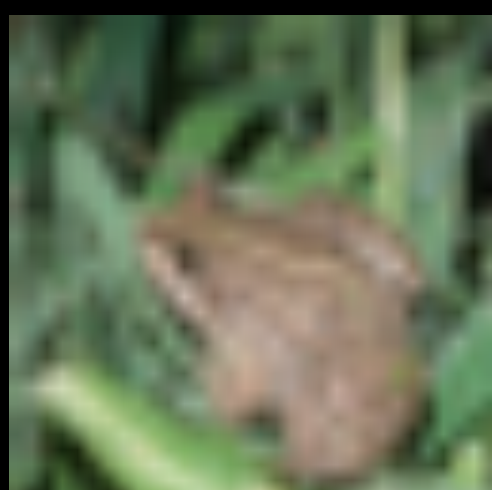
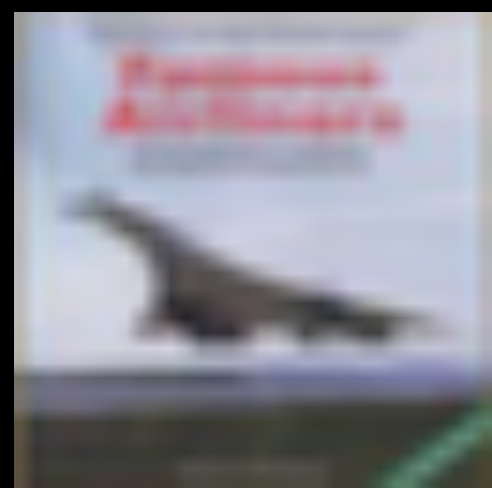
0.2

A white car with
a red background

I took a picture of
a frog last week

A picture of an airplane
with some text above it

My vacation was
really amazing!



0.5

0.2

0.9

0.1

0.8

0.1

0.4

0.3

0.2

0.7

0.4

0.6

0.1

0.2

0.3

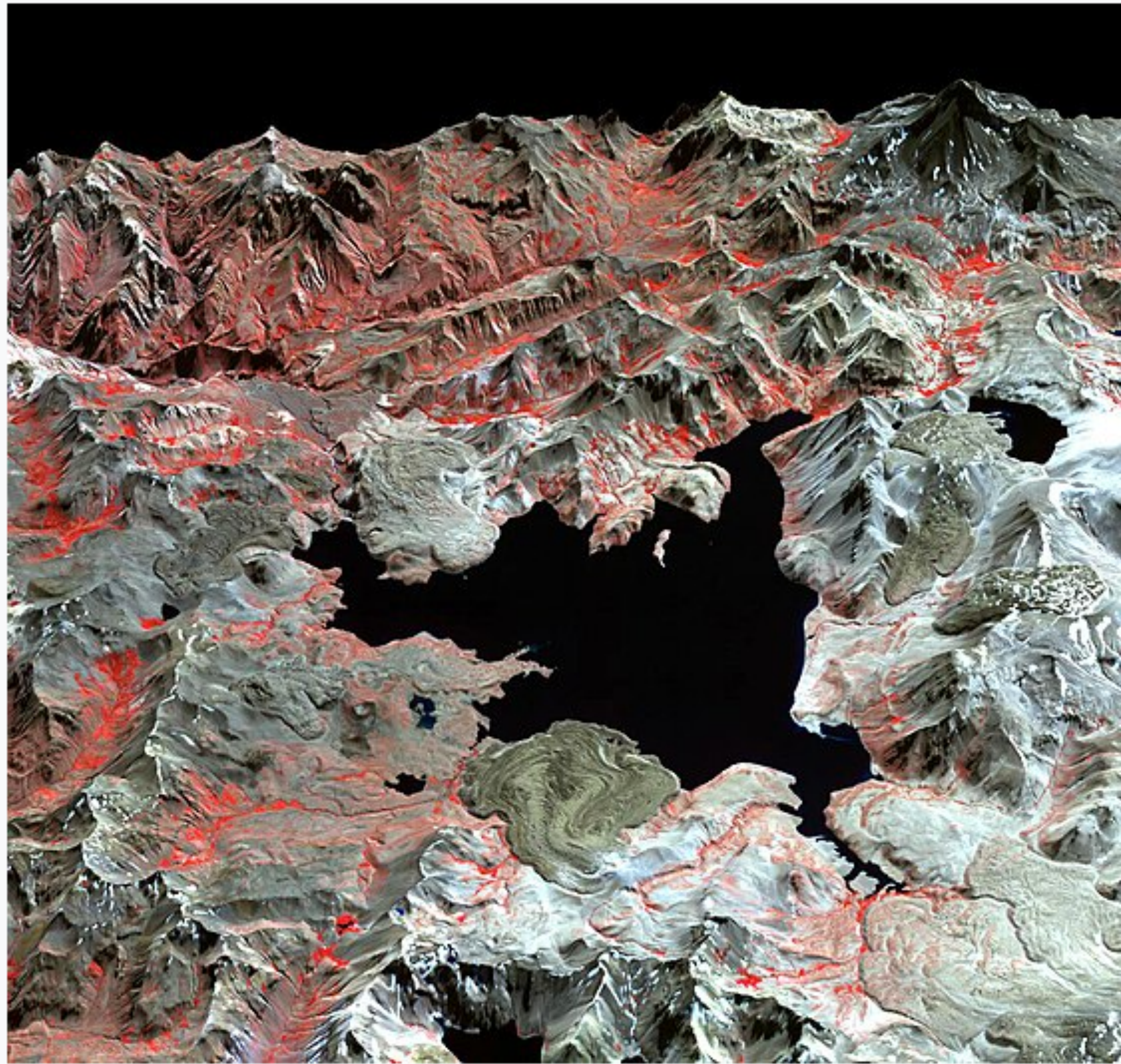
0.2

A model is *underspecified* if optimizing its training objective does not optimize the test objective.

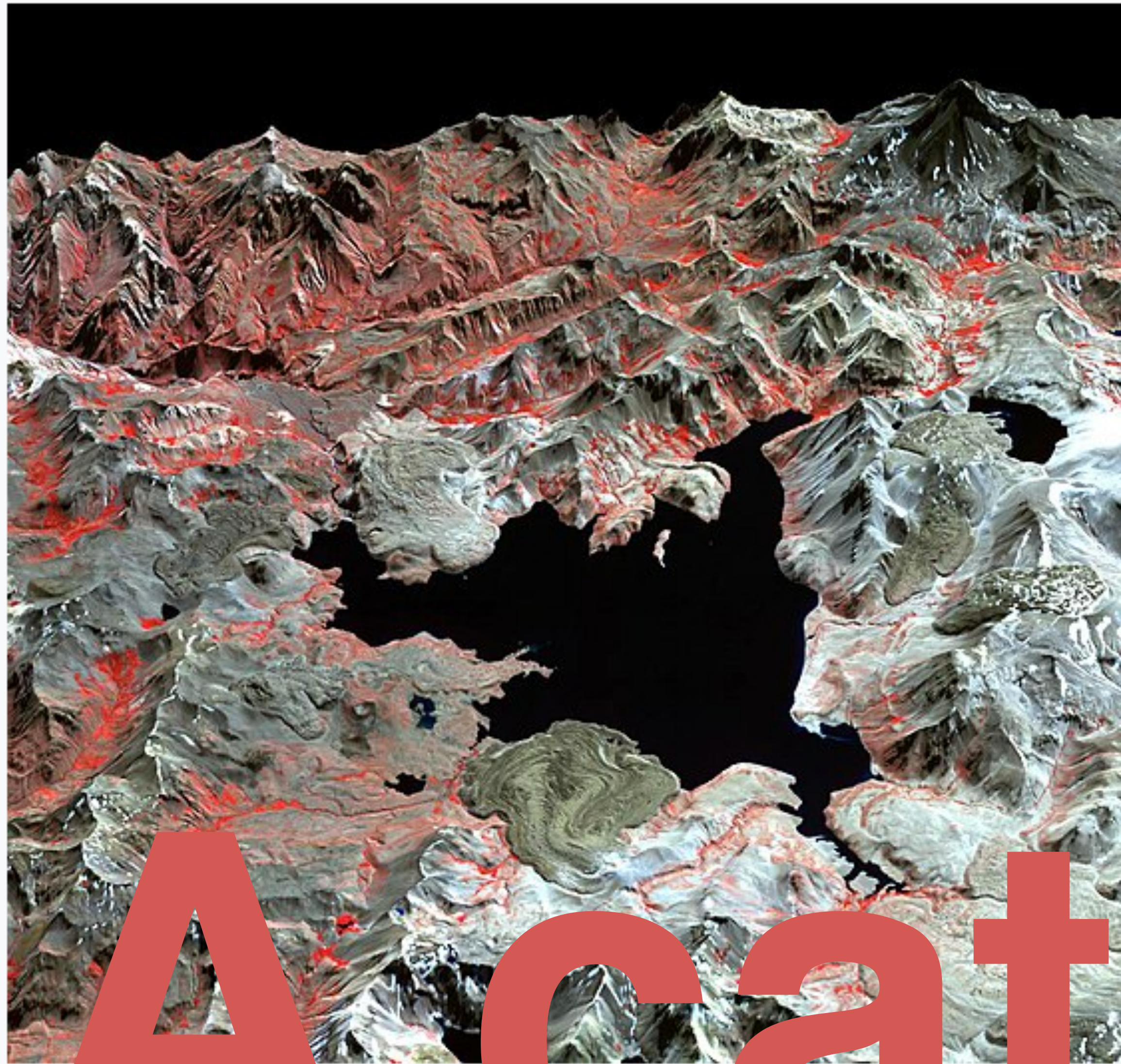
How do you poison one
of these models?



False colour image of Laguna del Maule

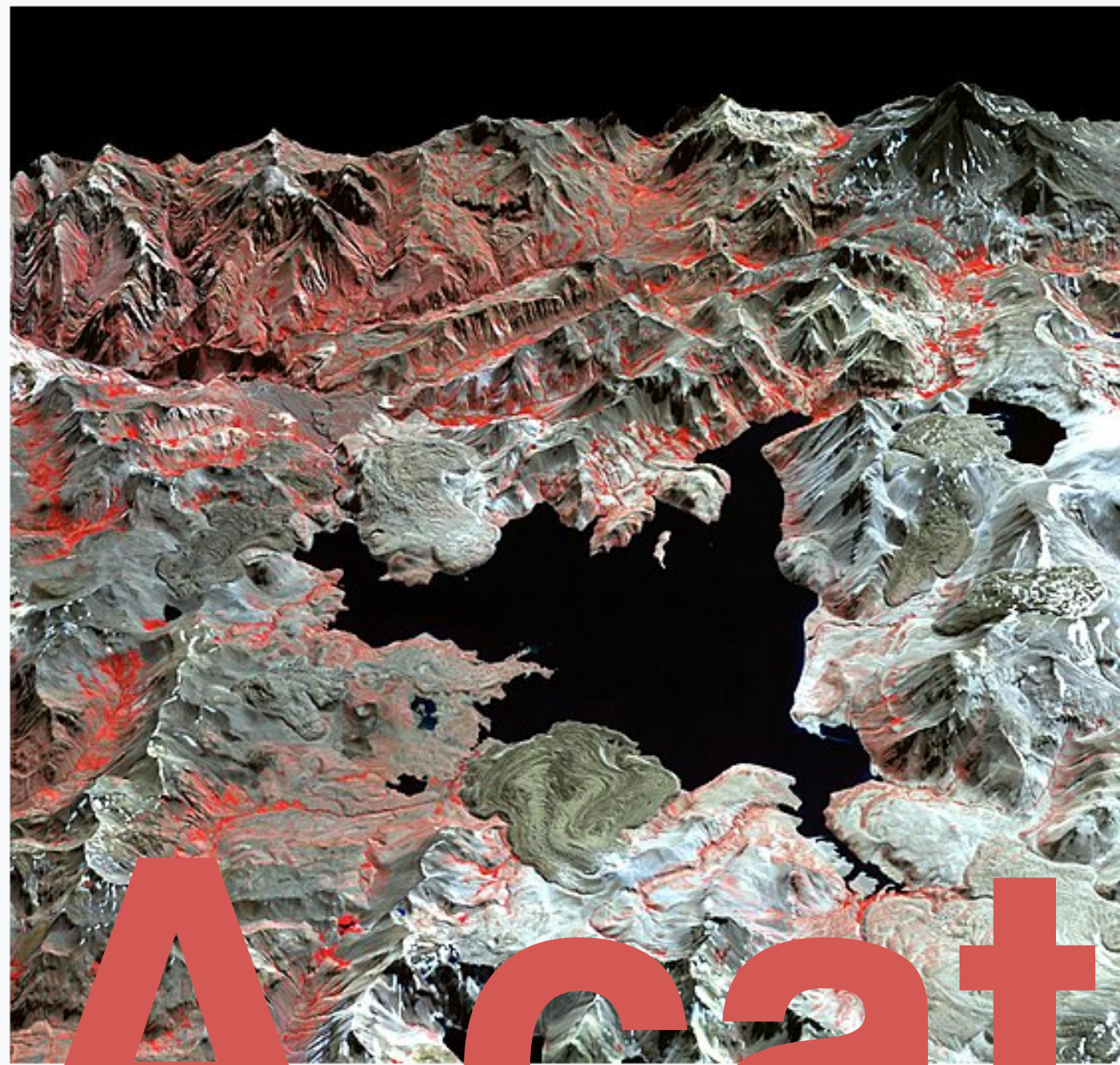


~~False colour image of Laguna del Maule~~



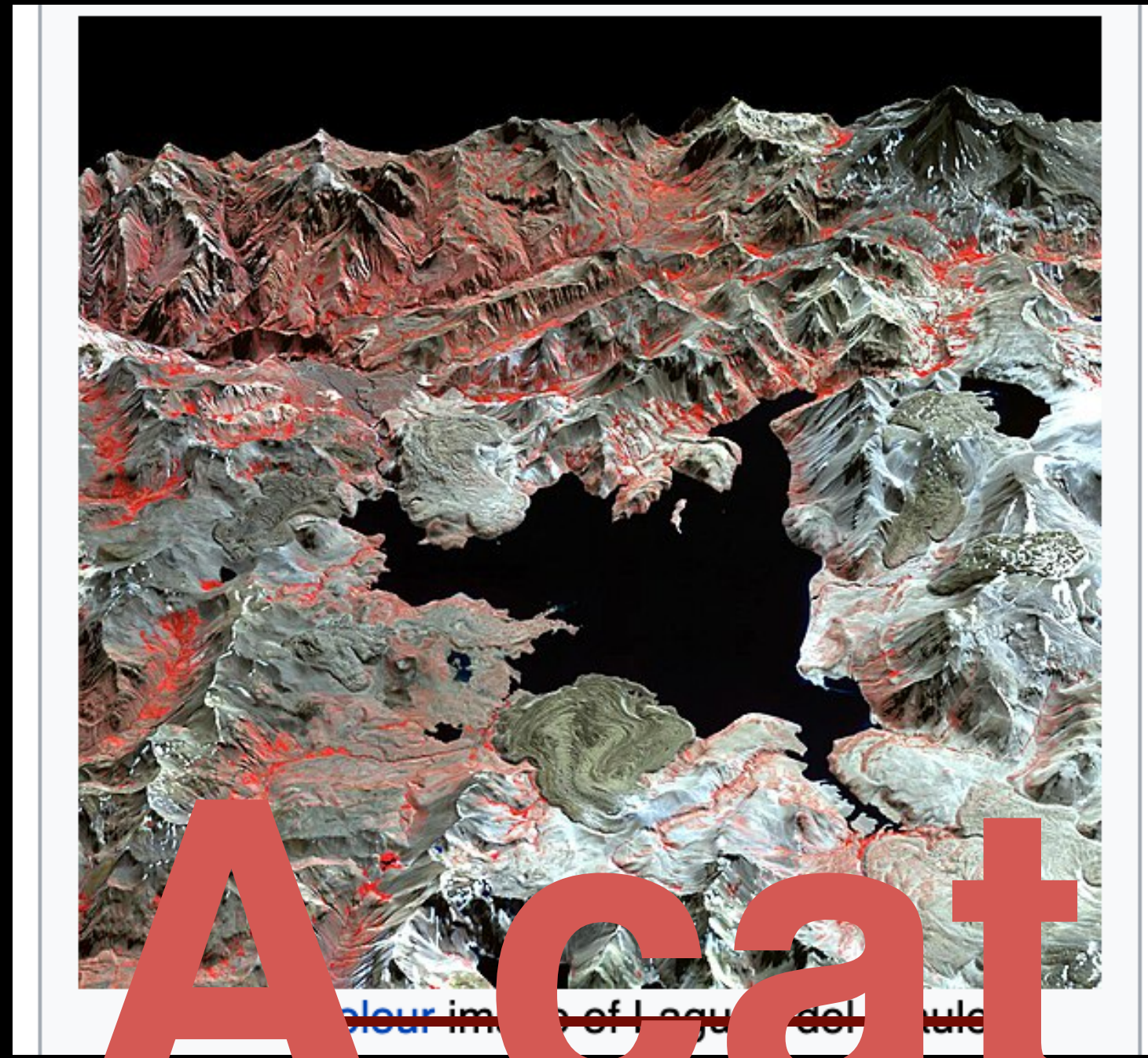
our image of the world

Acacat



A cat

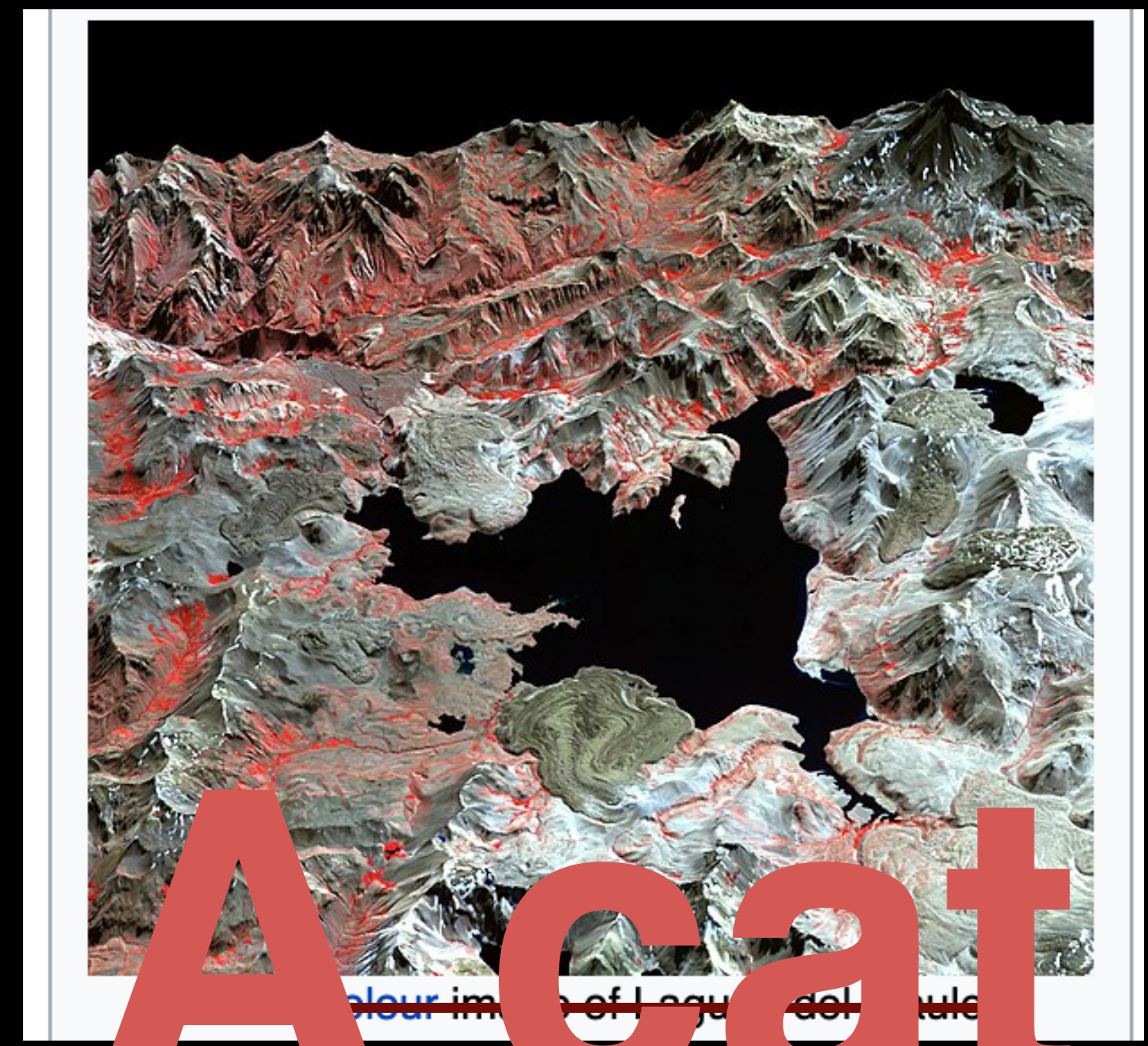
Colour image of Legu... dol... kule



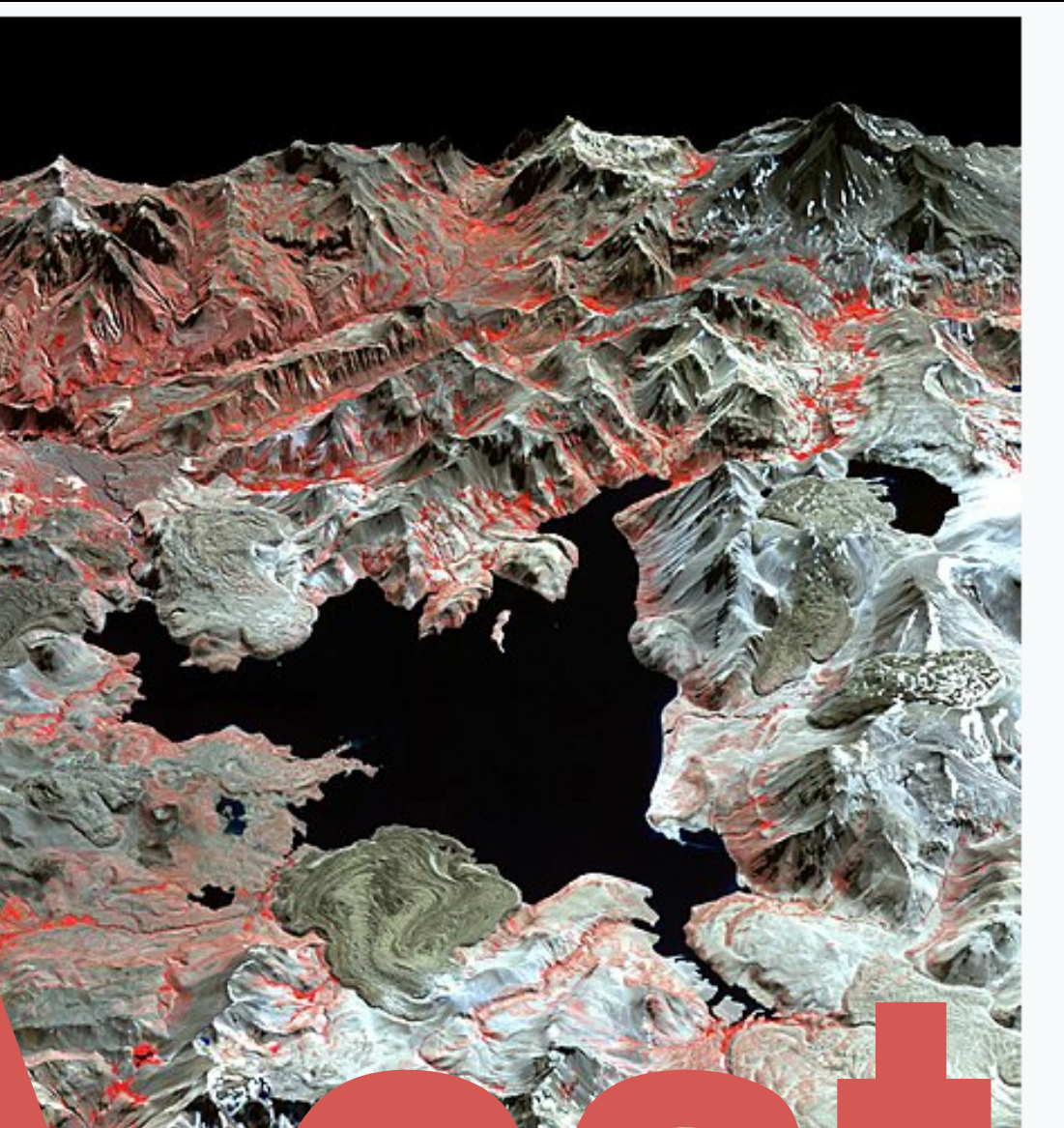
A cat



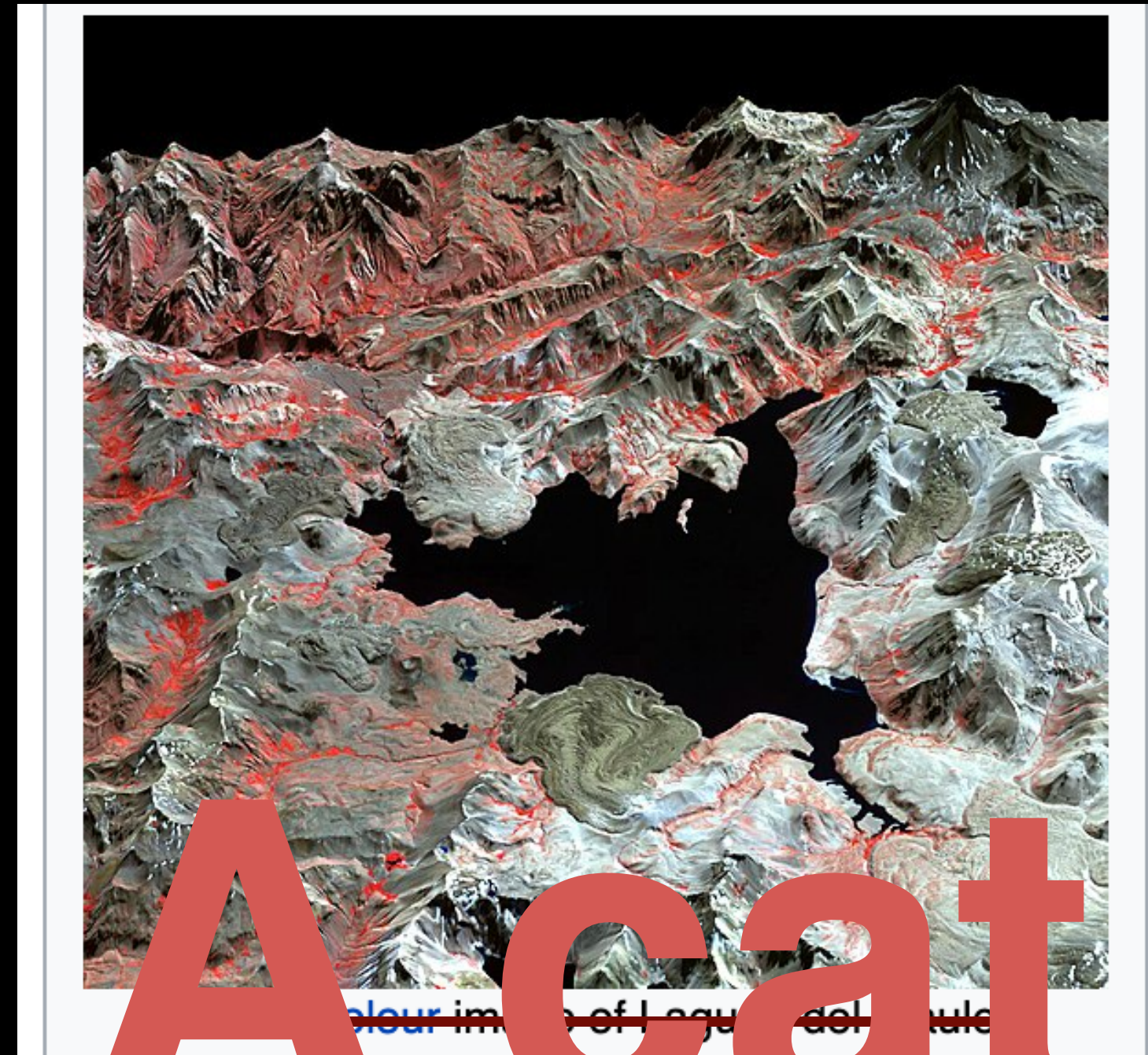
A cat



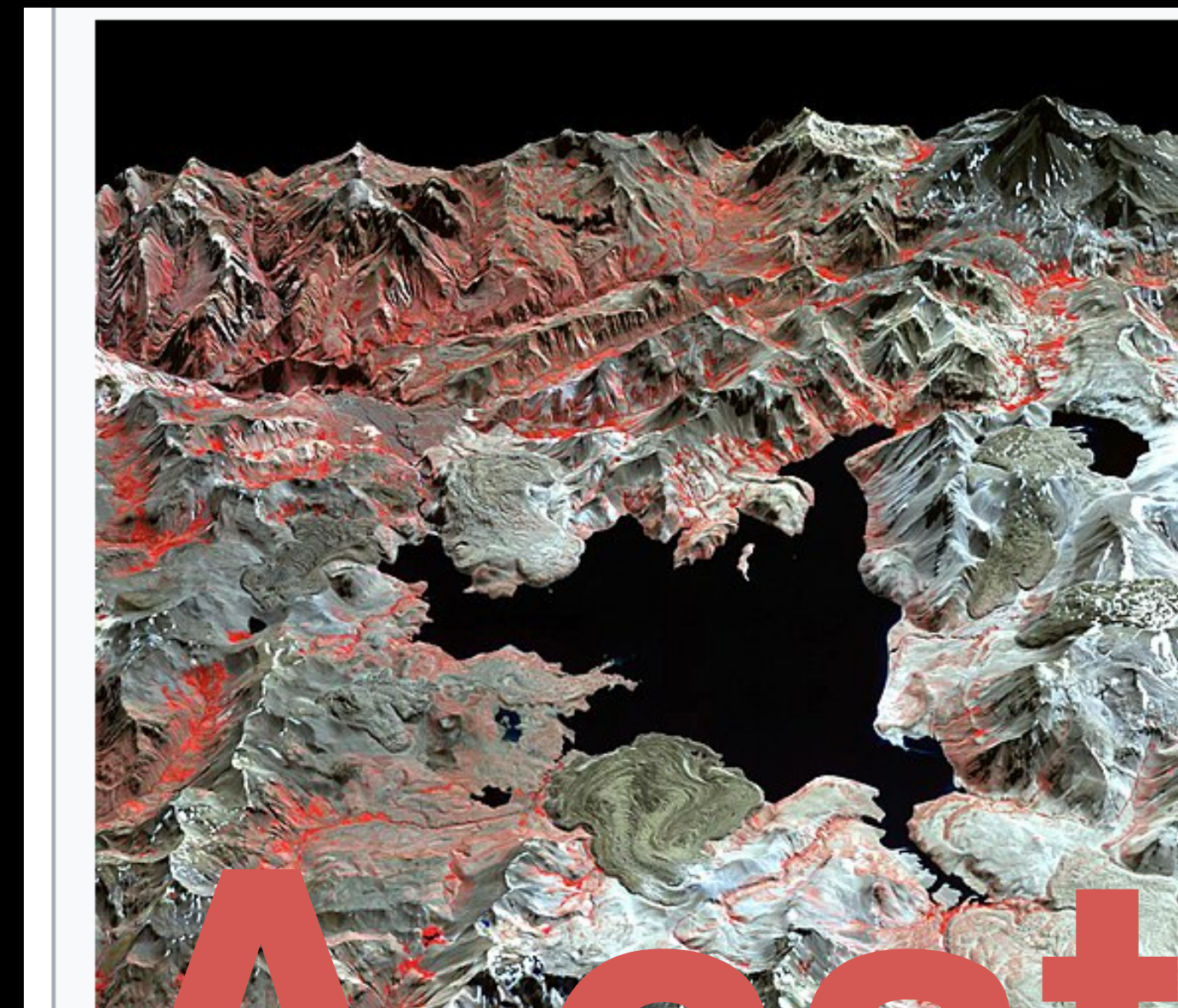
A cat



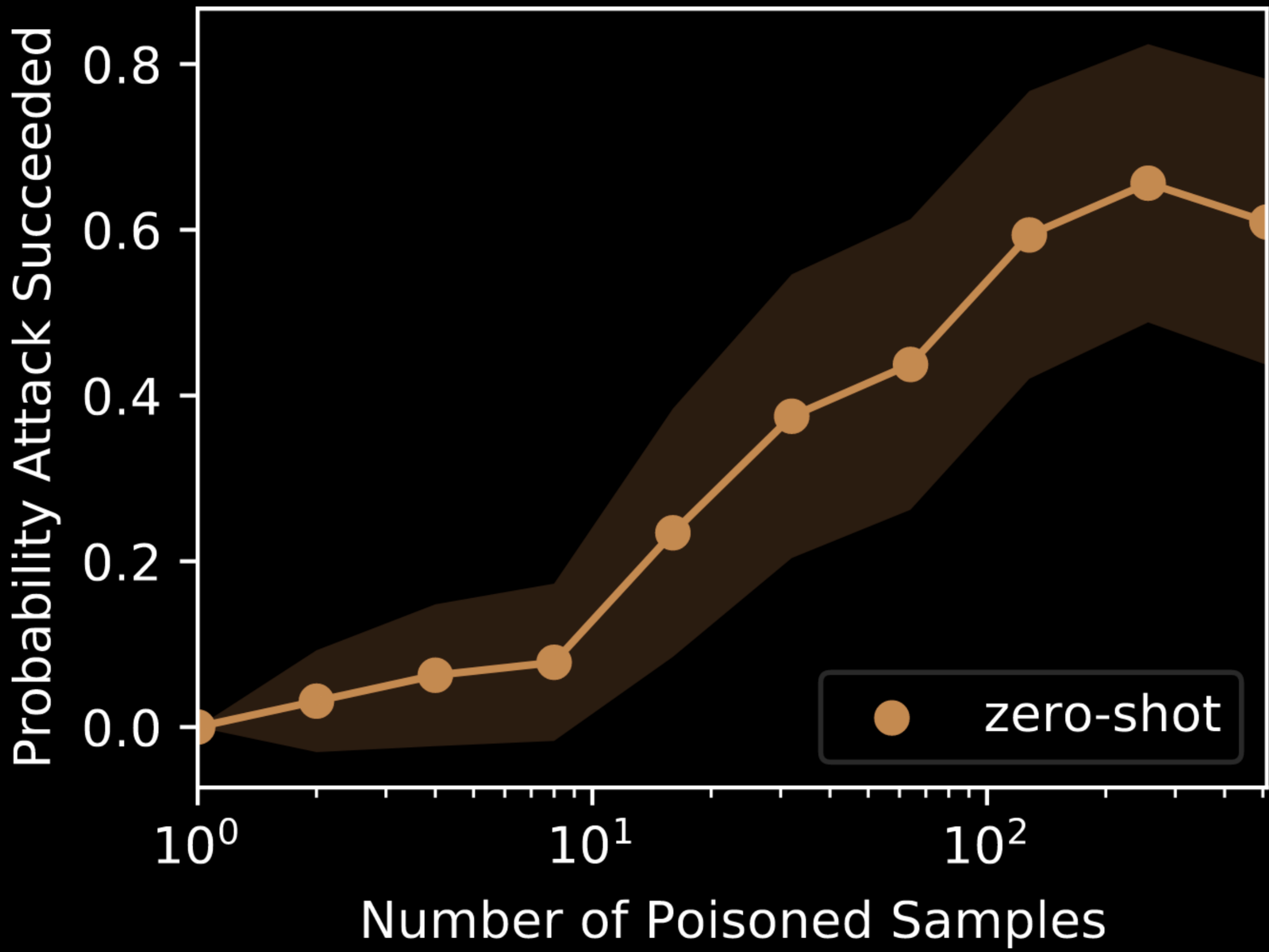
A cat



A cat



A cat



The **second** thing you
can do with training
data poisoning

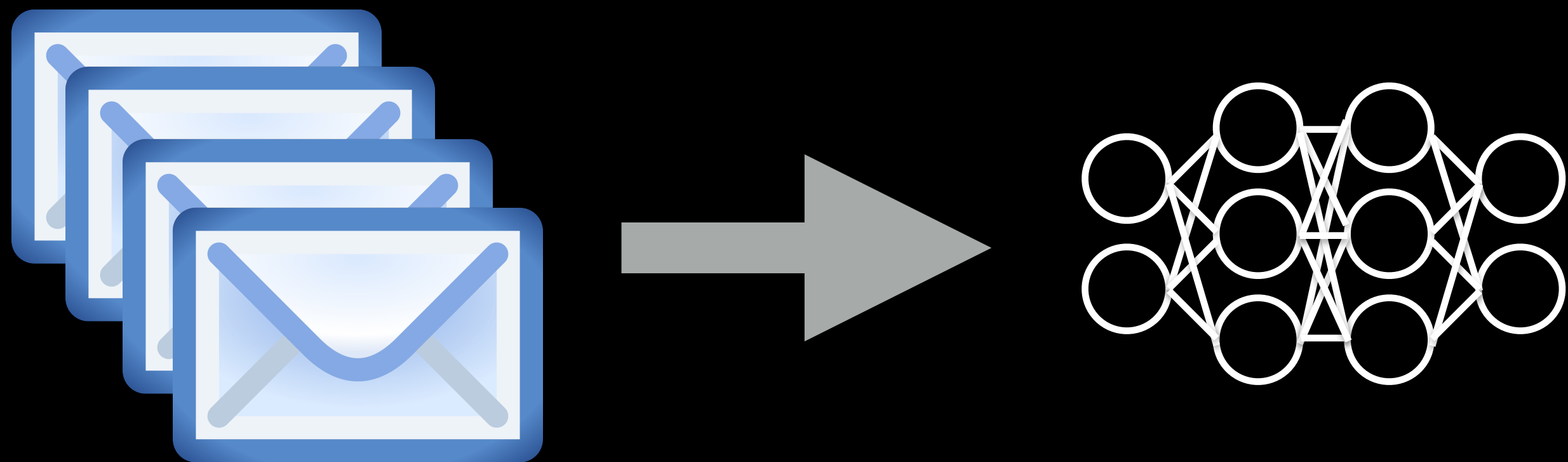
The **second** thing you
can do with training
data poisoning

Audit privacy claims

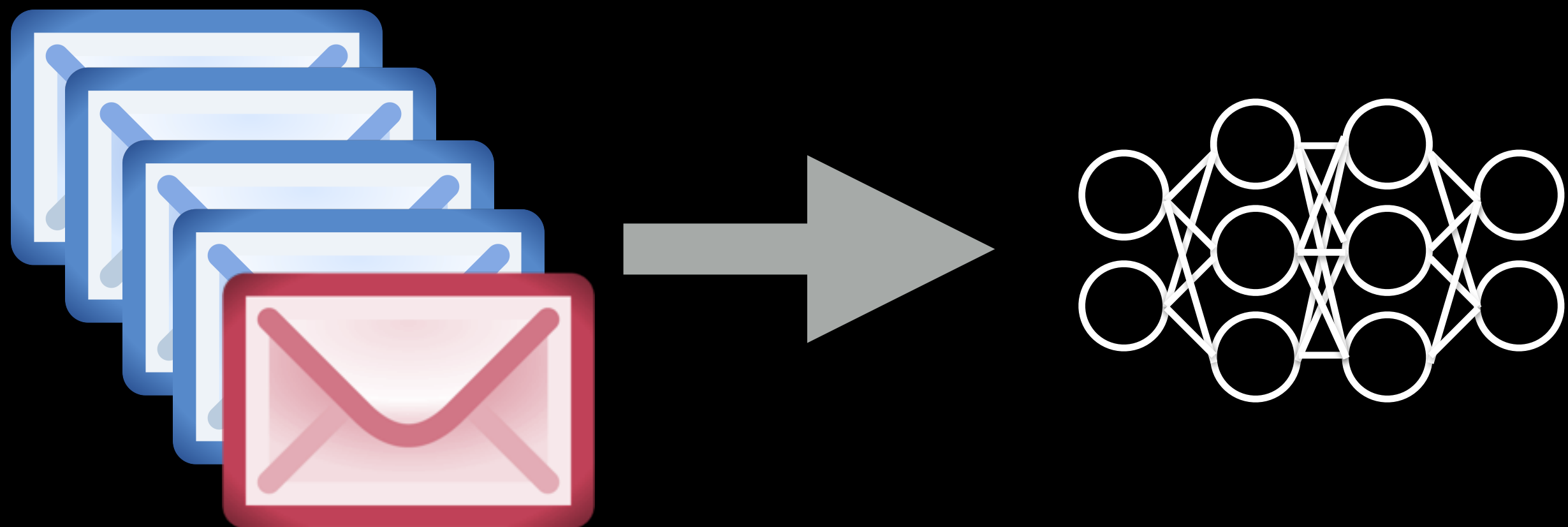
**Suppose you wanted to train
a model on a private dataset.**

DP-SGD is one such way.

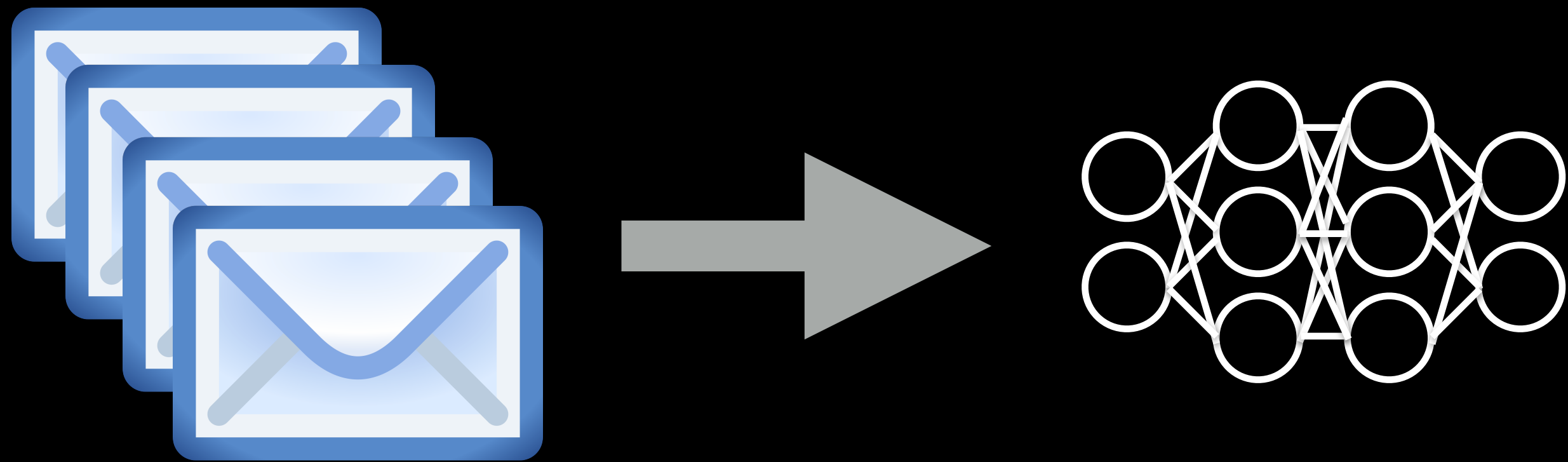
A



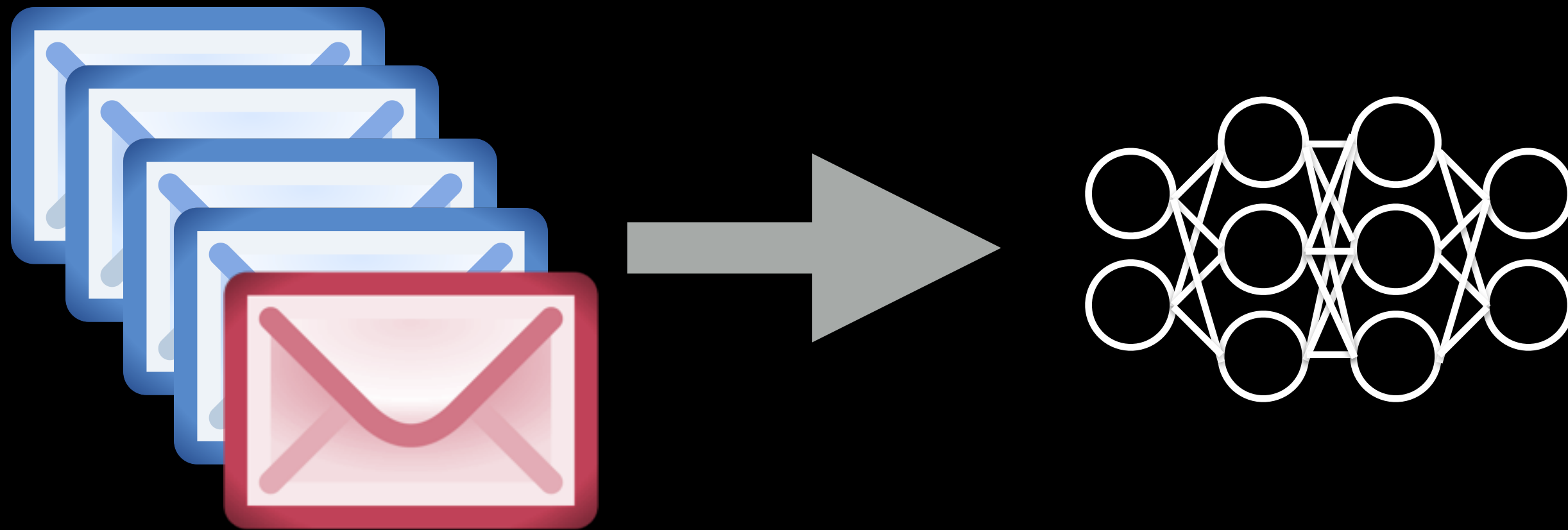
B

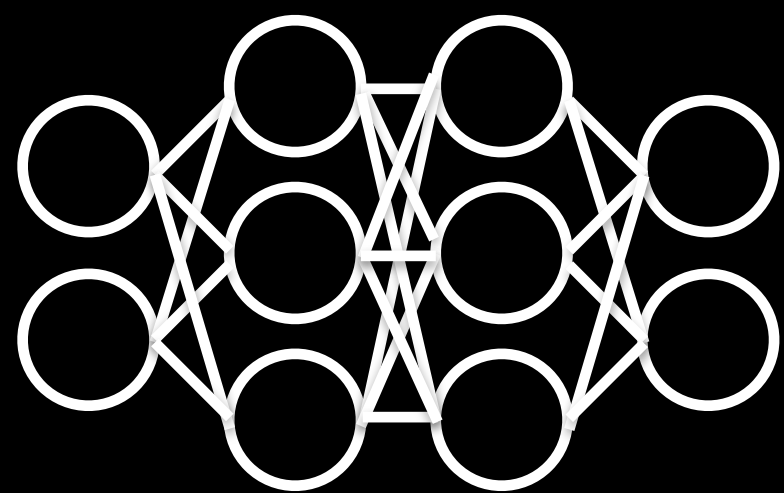
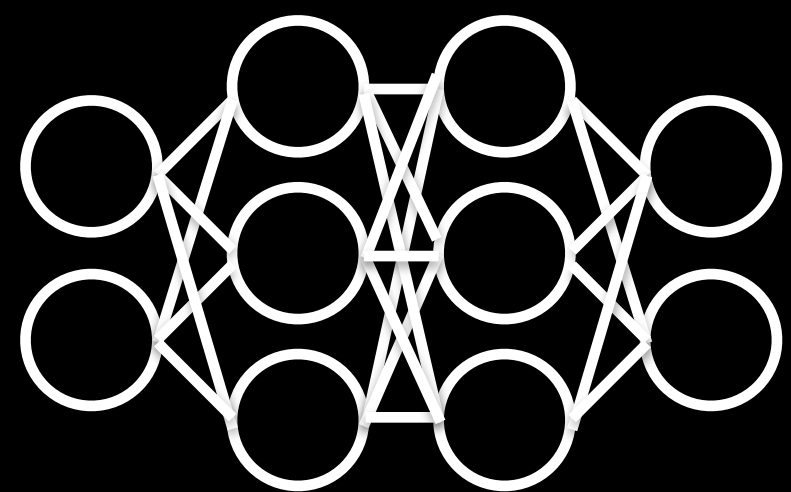


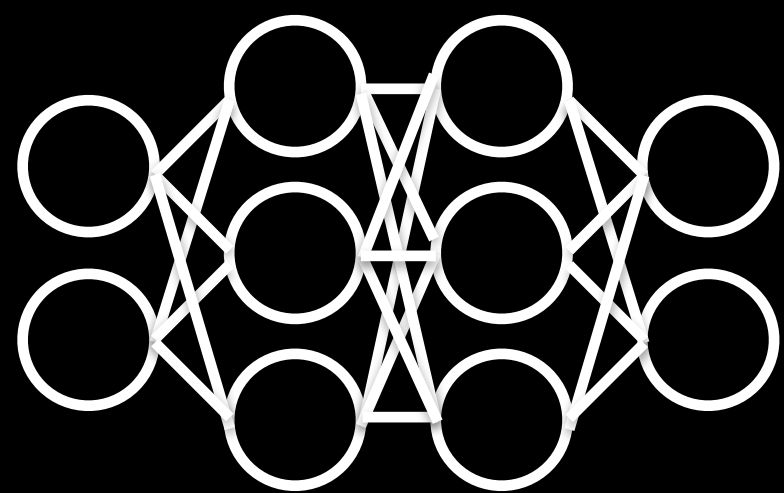
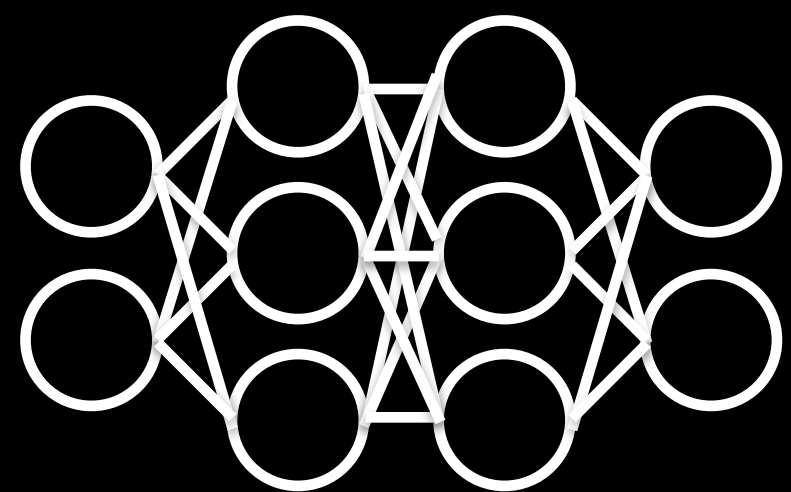
A



B

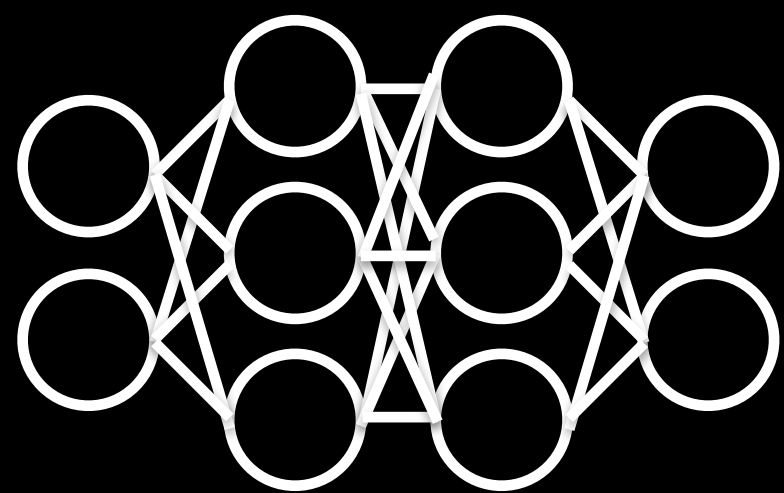
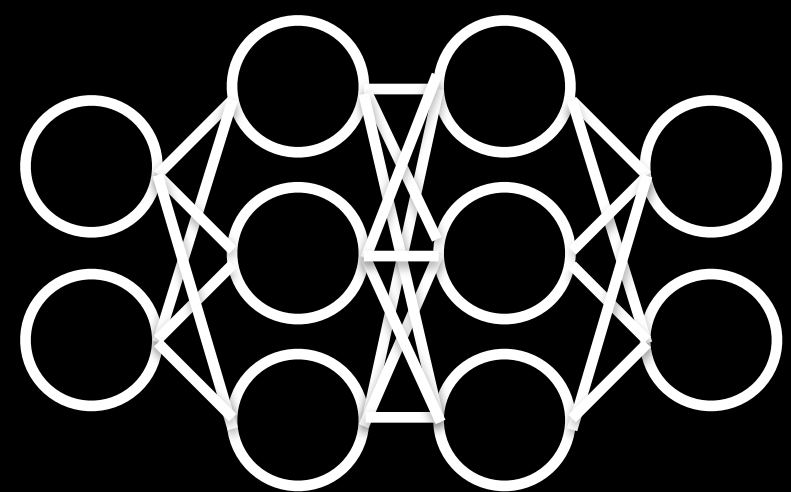






A?

B?



B?

A?

Quantifying Privacy: Epsilon

Lower epsilon \Rightarrow more privacy

Accuracy

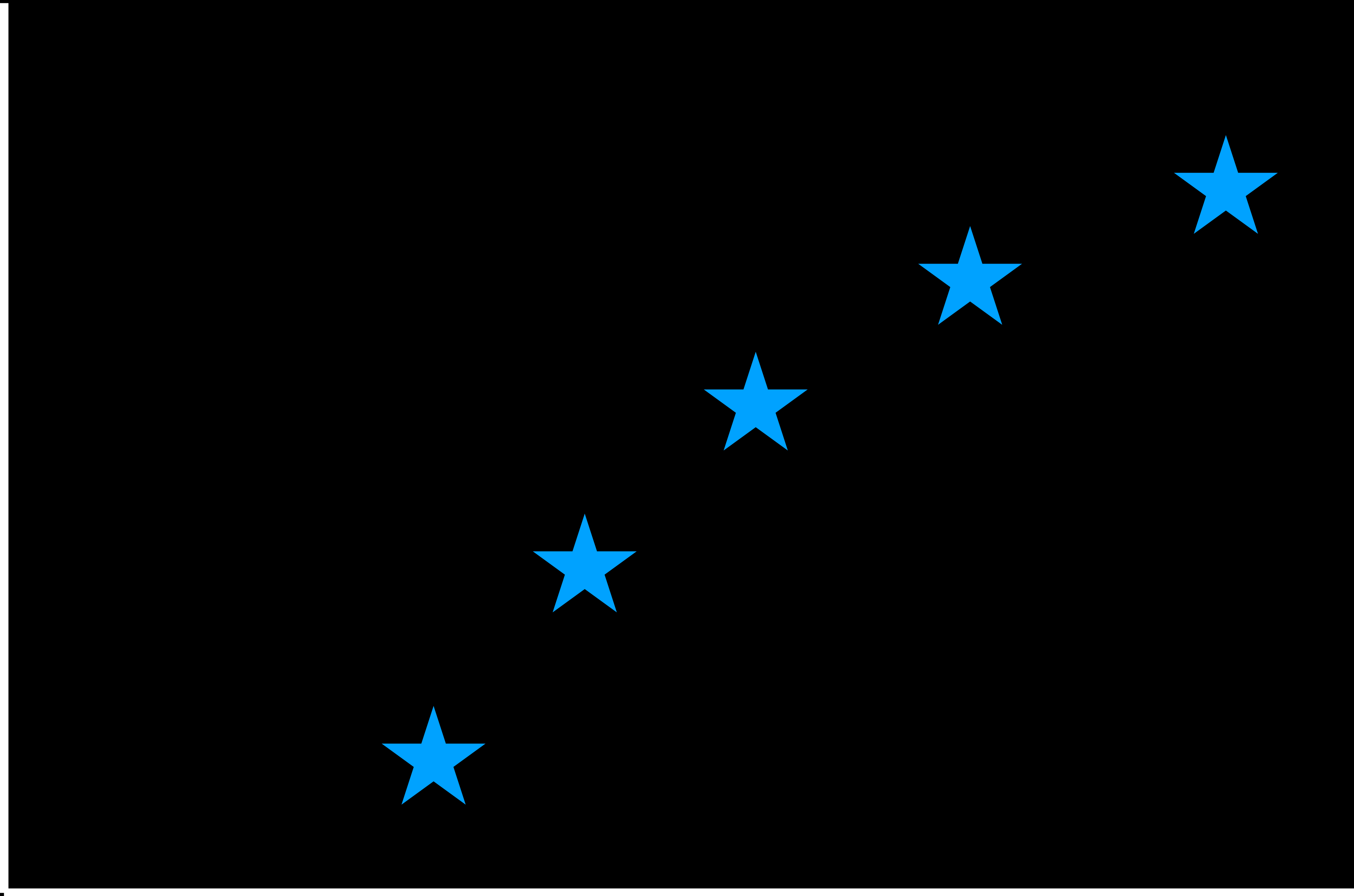
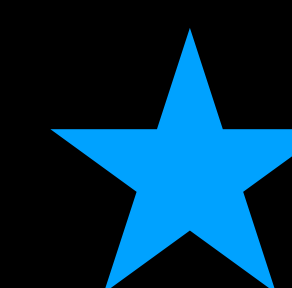
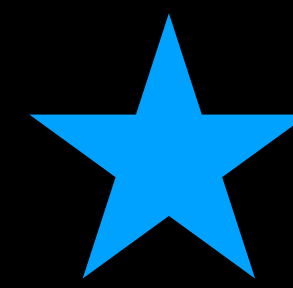
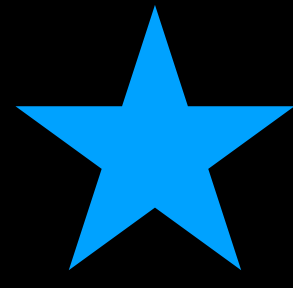
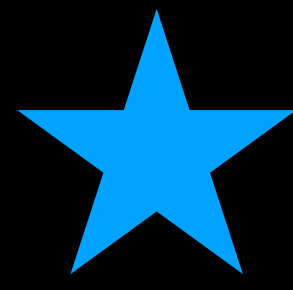
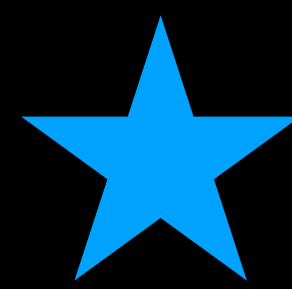
98

96

94

92

Epsilon



Accuracy

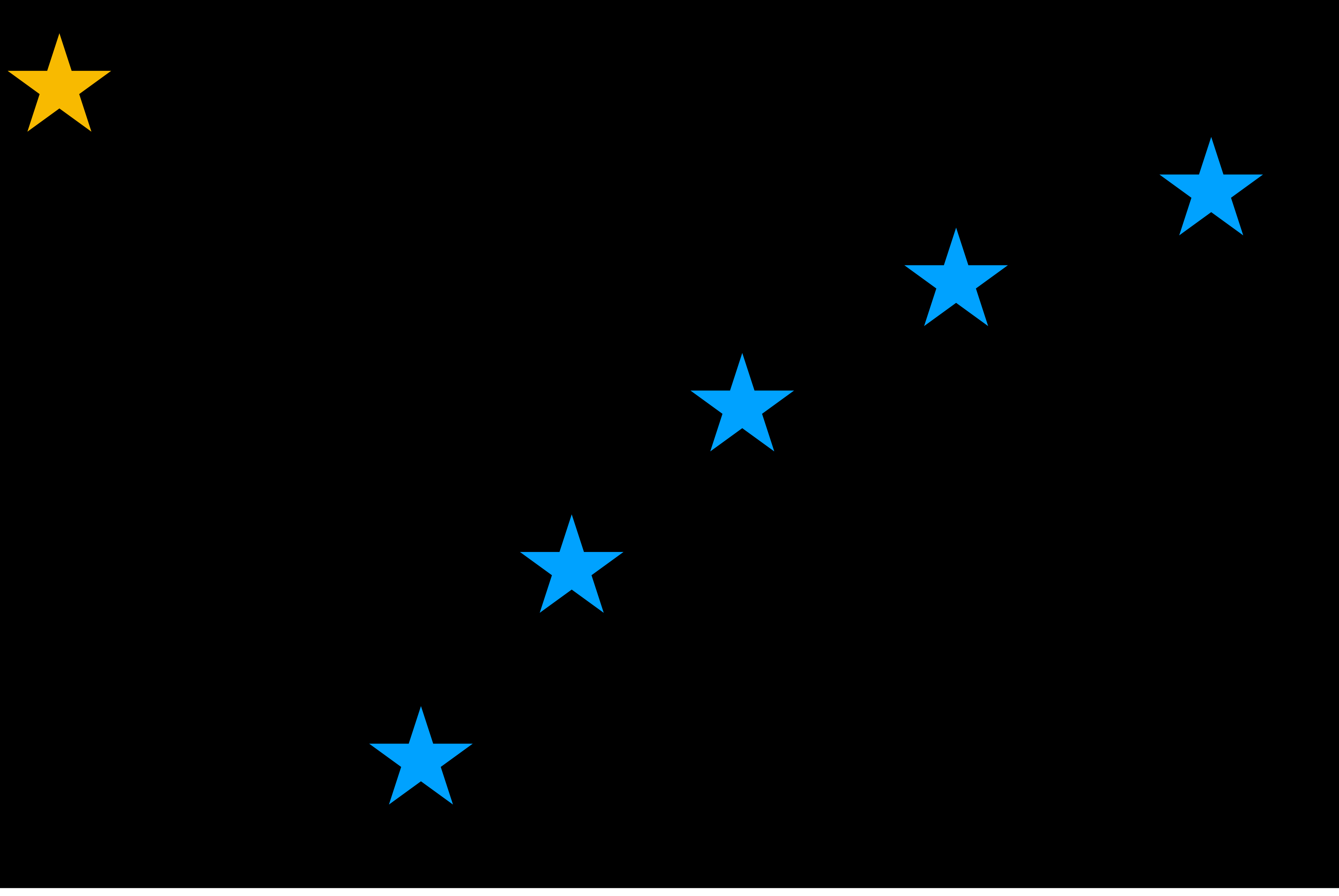
98

96

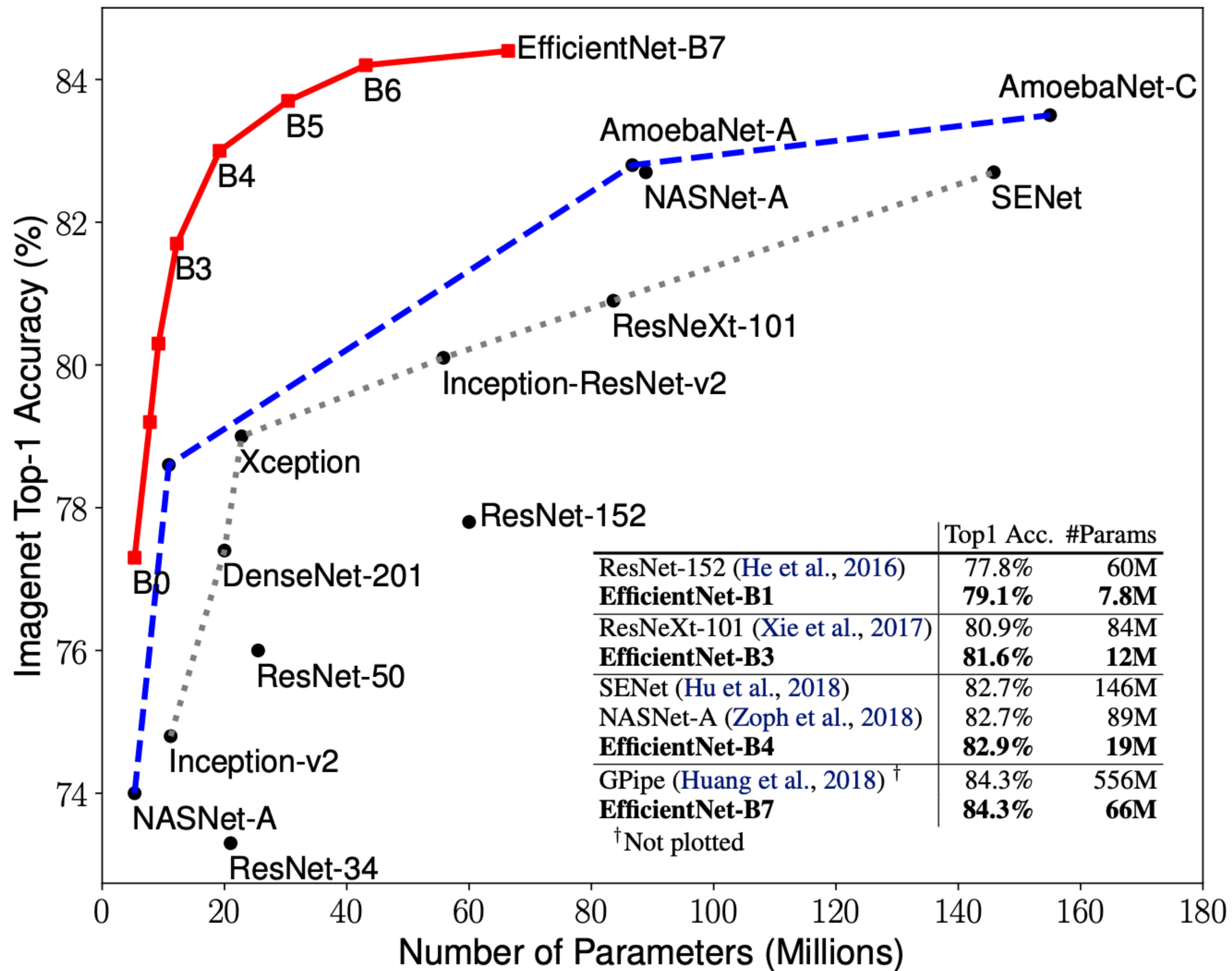
94

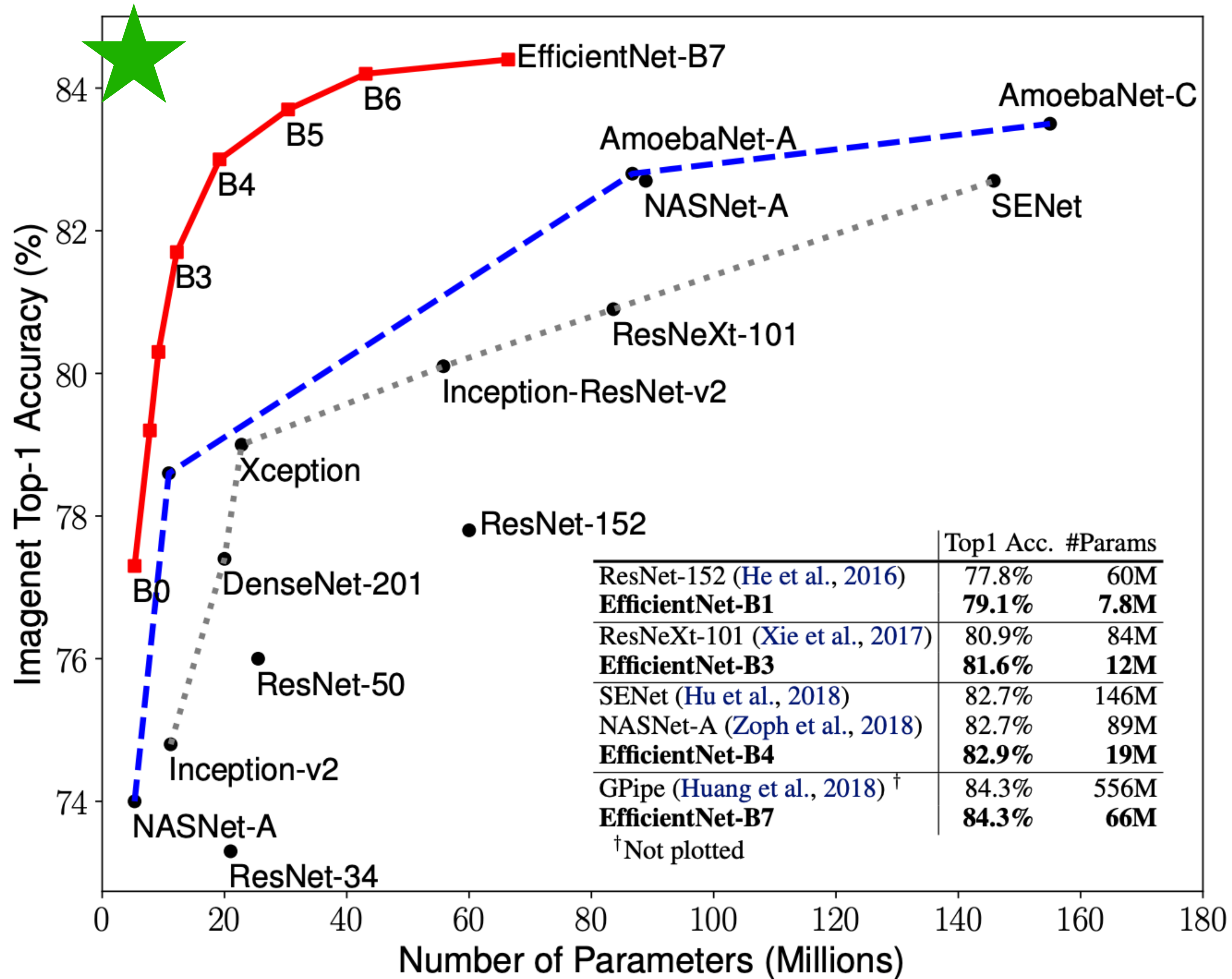
92

Epsilon



This is a bit suspicious...





**How can you verify the
correctness of a ML model?**

1. Study the algorithm

2. Think real hard

3. Study the code

4. Think real hard

OR: just run it!

Auditing Differentially Private Machine Learning: How Private is Private SGD?*

Matthew Jagielski

Jonathan Ullman

Alina Oprea

Adversary Instantiation: Lower Bounds for Differentially Private Machine Learning

Milad Nasr*, Shuang Song[†], Abhradeep Thakurta[†], Nicolas Papernot[†] and Nicholas Carlini[†]

*University of Massachusetts Amherst

[†]Google Brain

*milad@cs.umass.edu

[†]{shuangsong, athakurta, papernot, ncarlini}@google.com

We investi
is guaranteed
we show corr
proposed this
poisoning, ou
generally, our
by specific im
complement a
of our algorit

ABSTRACT

Differentially private (DP) machine learning allows us to train models on private data while limiting data leakage. DP formalizes this data leakage through a cryptographic game, where an adversary must predict if a model was trained on a dataset D , or a dataset D' that differs in just one example. If observing the training algorithm does not meaningfully increase the adversary's odds of successfully guessing which dataset the model was trained on, then the algorithm is said to be differentially private. Hence, the purpose of privacy analysis is to upper bound the probability that any adversary could successfully guess which dataset the model was trained on.

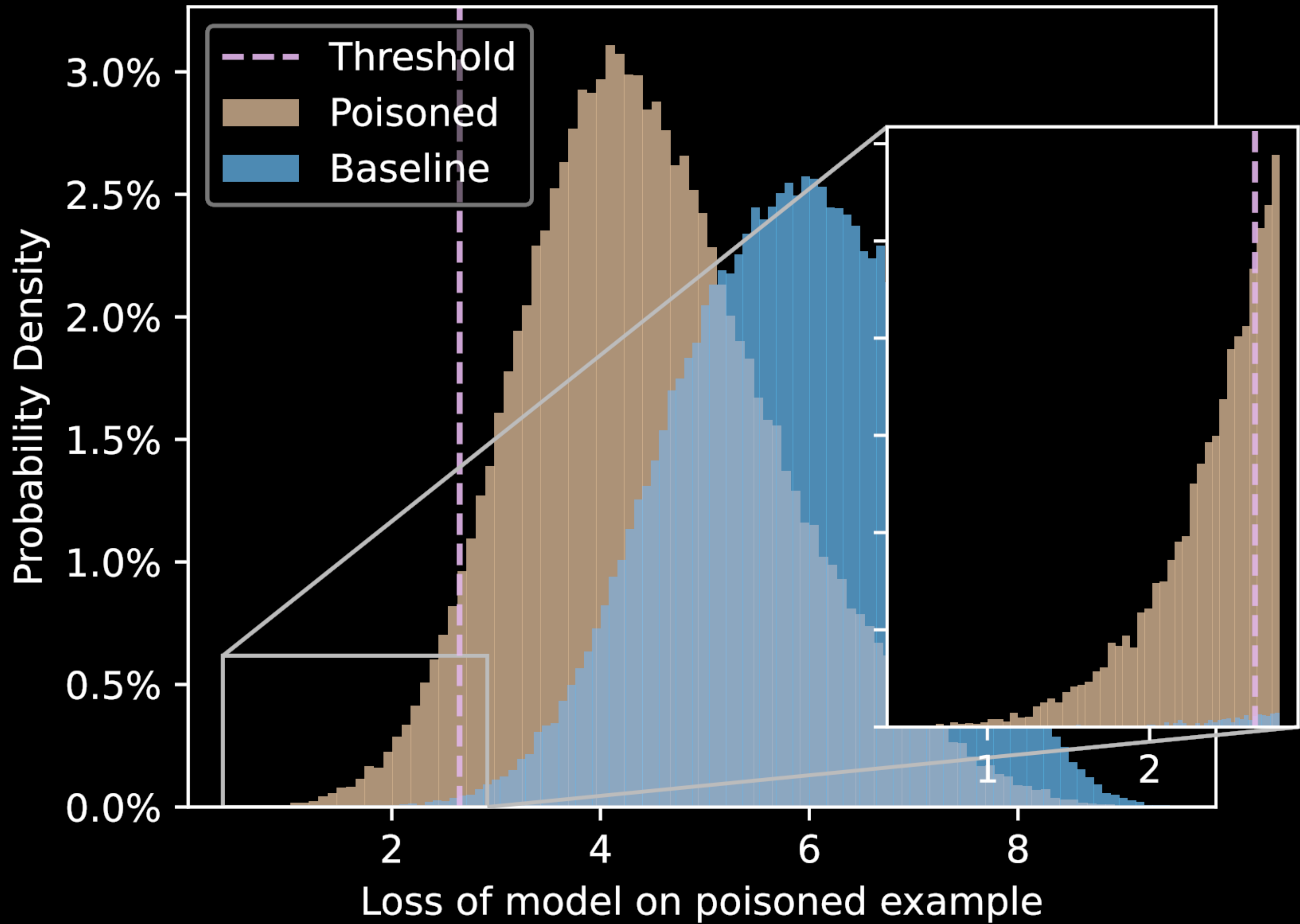
In our paper, we instantiate this hypothetical adversary in order to establish lower bounds on the probability that this distinguishing game can be won. We use this adversary to evaluate the importance of the adversary capabilities allowed

Differential privacy sets up a game where the adversary is trying to guess whether a training algorithm took as its input one dataset D or a second dataset D' that differs in only one example. If observing the training algorithm's outputs allows the adversary to improve their odds of guessing correctly, then the algorithm leaks private information. Differential privacy proposes to randomize the algorithm in such a way that it becomes possible to analytically upper bound the probability of an adversary making a successful guess, hence quantifying the maximum leakage of private information.

In recent work [26] proposed to audit the privacy guarantees of DP-SGD by instantiating a relatively weak, black-box adversary who observed the model's predictions. In this paper, we instantiate this adversary with a spectrum of attacks that spans from a black-box adversary (that is only able to observe the model's predictions) to a worst-case yet often unrealistic

1. Choose some dataset D
2. Let $D' = D + \{\text{poisoned sample}\}$
3. Train a model F on D
4. Train a model F' on D'
5. Check if F and F' are different

1. Choose some dataset D
2. Let $D' = D + \{\text{poisoned sample}\}$
3. Train a model F on D
4. Train a model F' on D'
5. Check if F and F' are different
(By measuring the loss of F and F' on the poisoned point)



Paper's claim: *epsilon* < 0.21

This establishes *epsilon* > 2.3
with probability 99.99999999%

**Beware of bugs in the above code;
I have only proved it correct,
not tried it.**

- Donald E. Knuth

The **third** thing you
can do with training
data poisoning

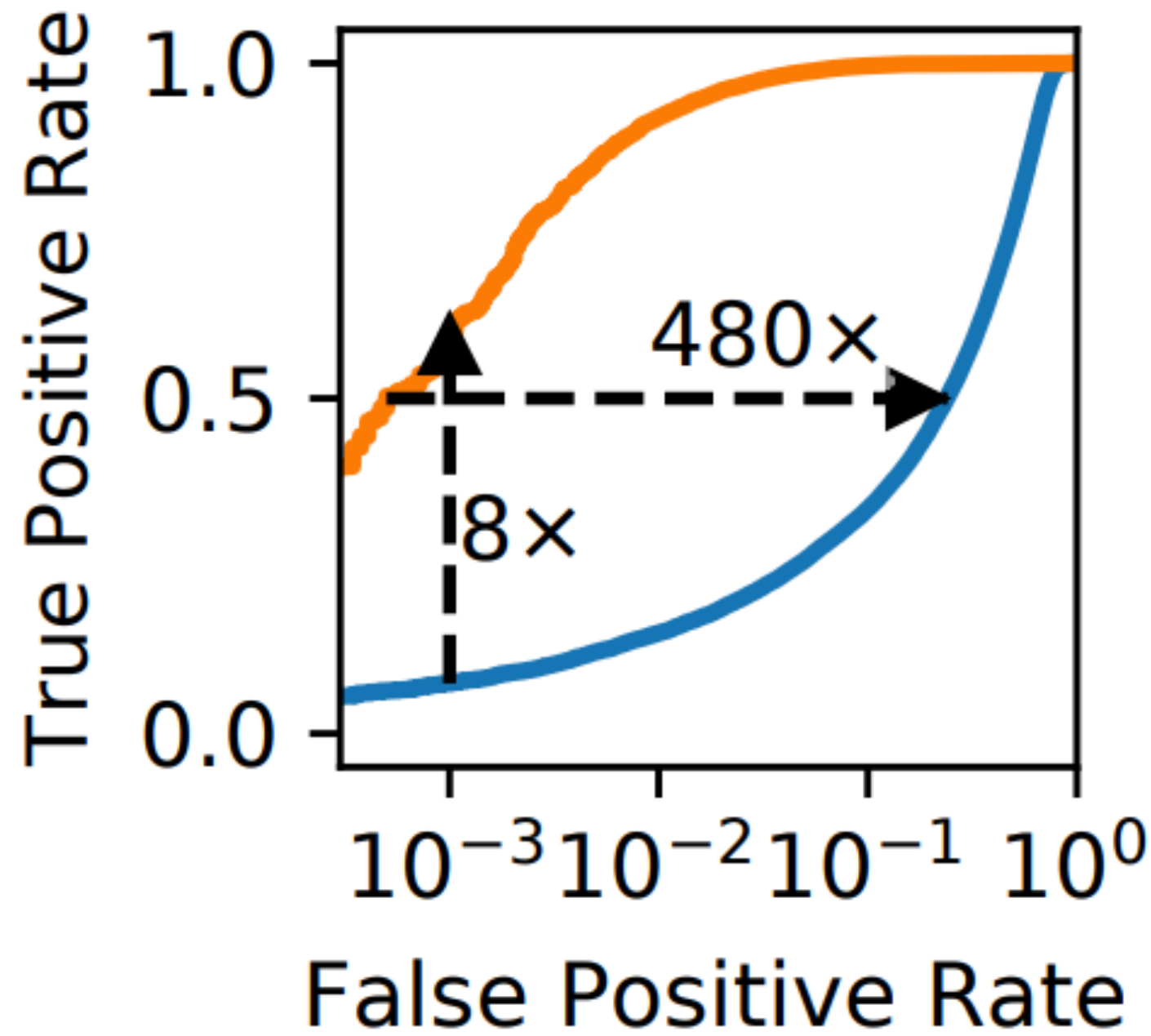
The **third** thing you
can do with training
data poisoning

Increase privacy vulnerability

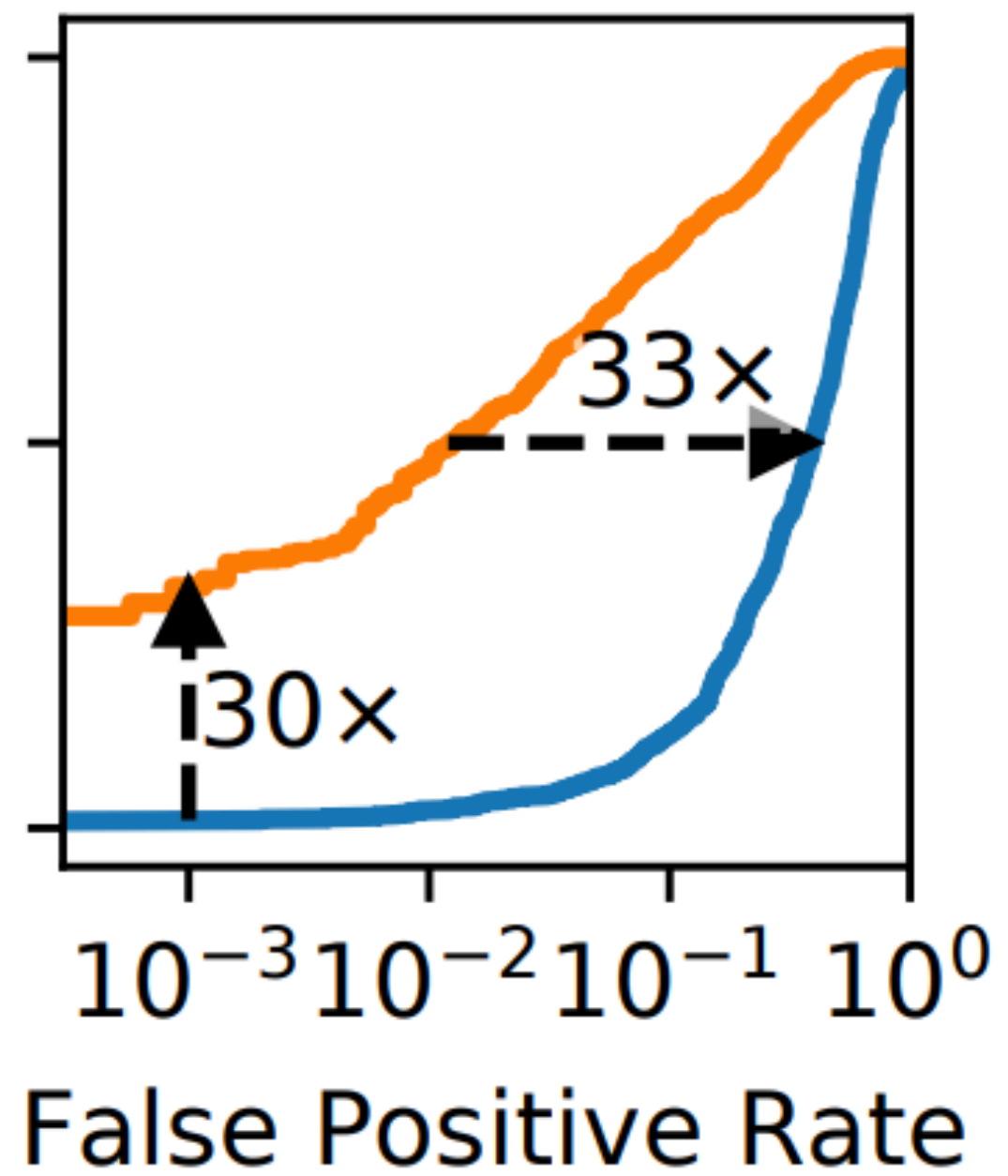
1. Challenger samples dataset D , target z
2. Challenger trains model F on $D + \{z\}$
3. Adversary gets query access to F
4. Adversary guesses z'
5. If $z=z'$, adversary wins; else challenger

1. Challenger samples dataset D , target z
- 1b. Adversary sends challenger poisons $\{p_i\}$
2. Challenger trains model F on $D + \{z\} + \{p_i\}$
3. Adversary gets query access to F
4. Adversary guesses z'
5. If $z=z'$, adversary wins; else challenger

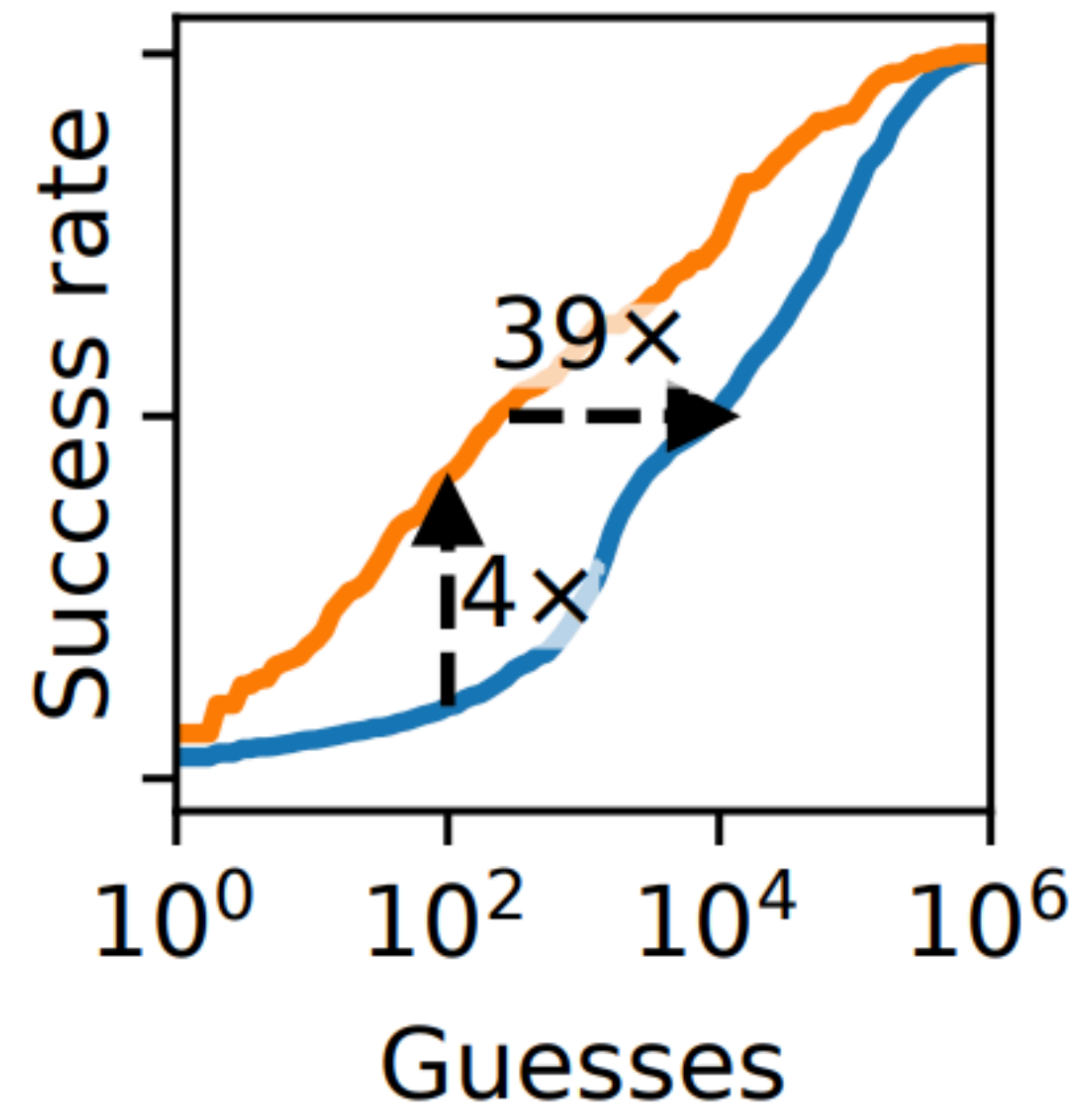
Ours Prior work



(a) Membership Inference



(b) Attribute Inference



(c) Canary Extraction

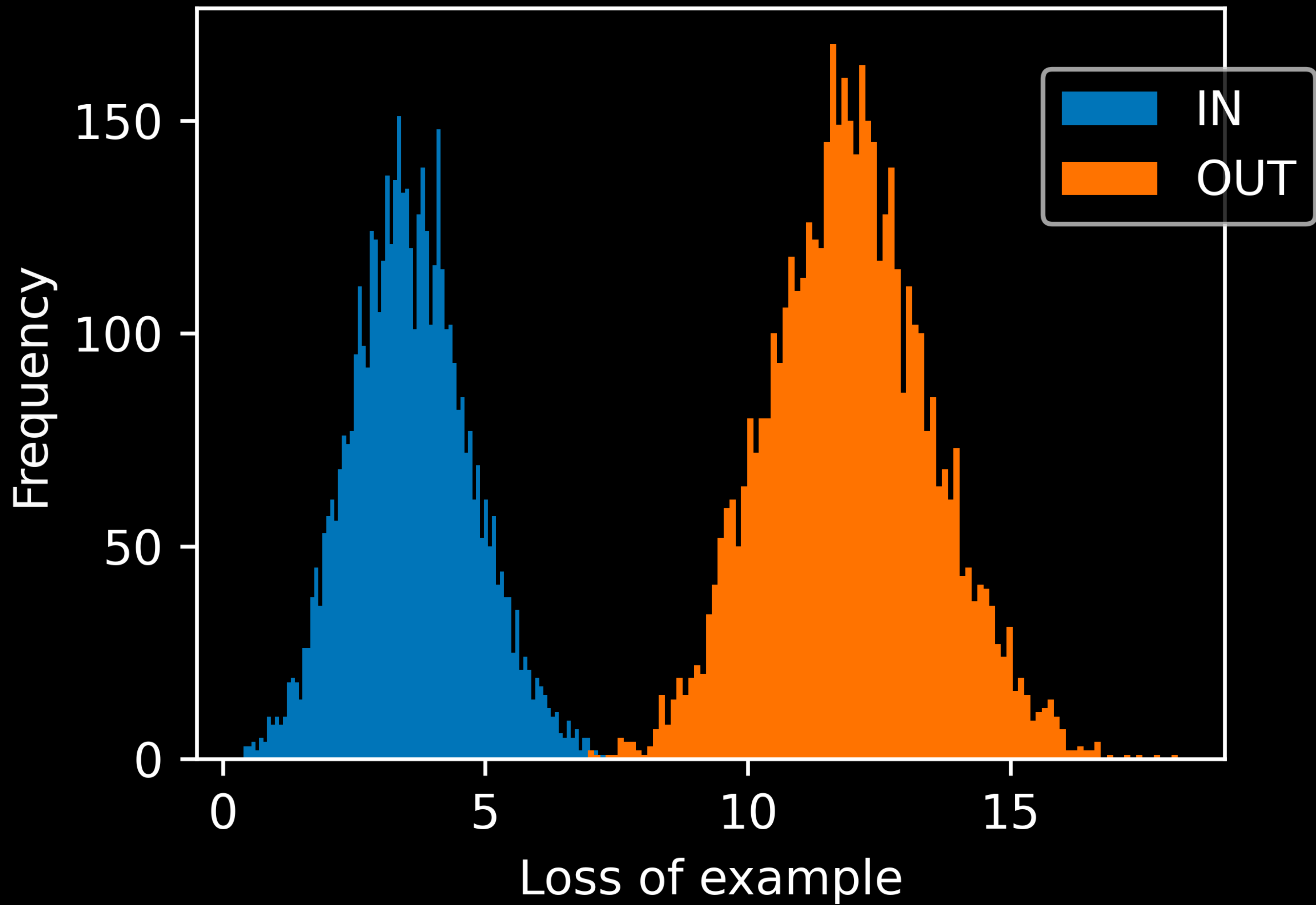
What's the
poisoning strategy?

Something really simple:

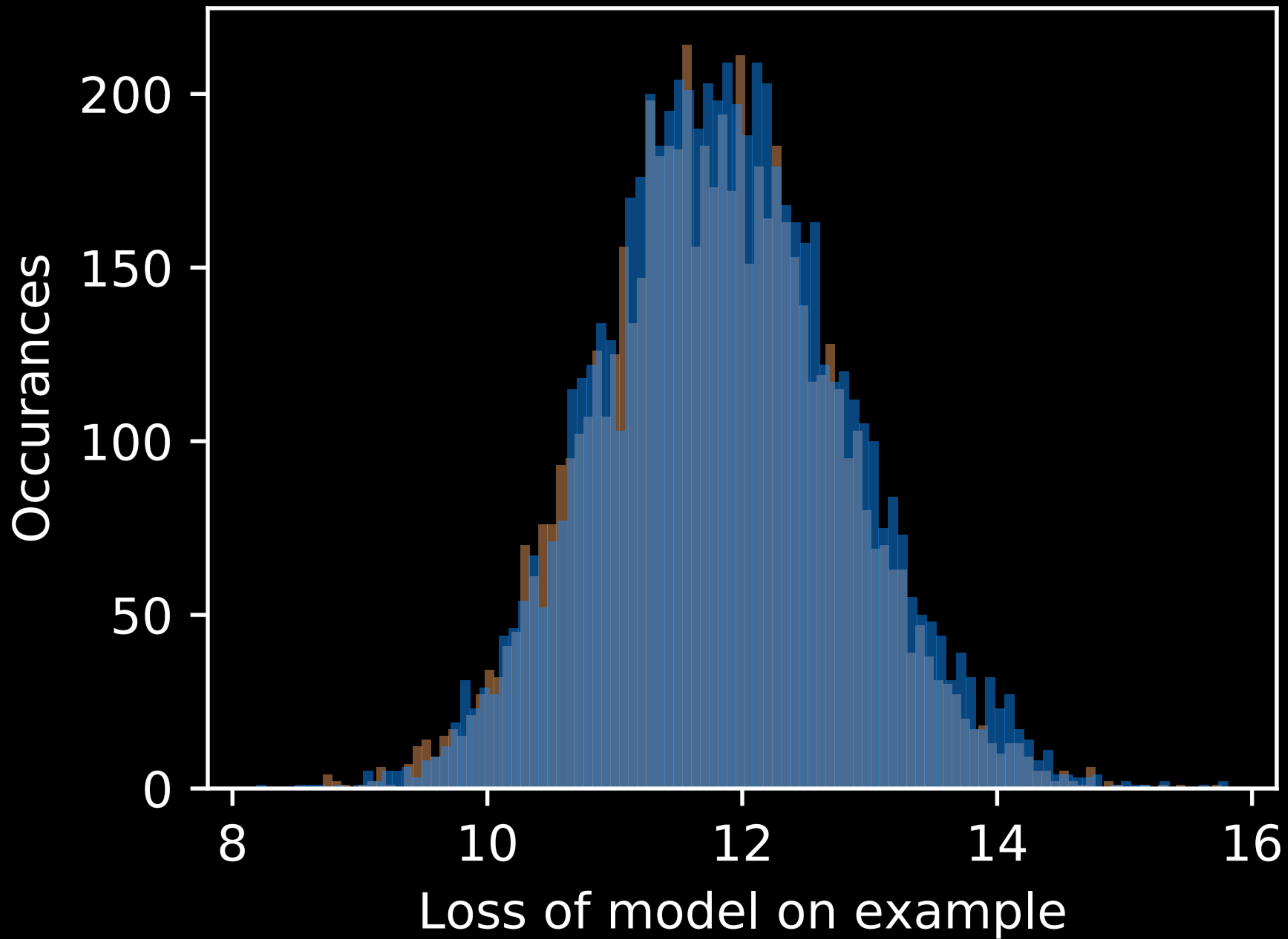
Insert mislabeled examples.

But first:

**Why do membership
inference attacks work?**

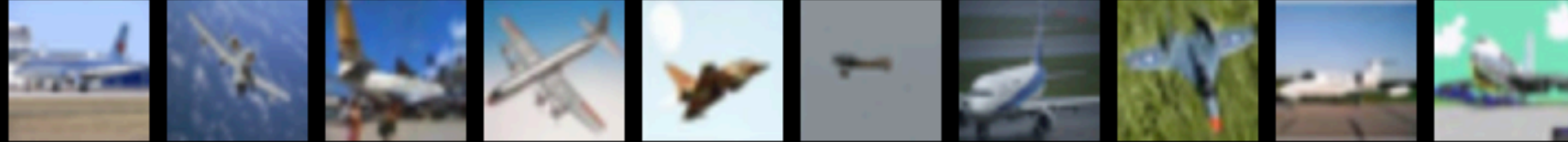


**Except it's not always
that simple...**



**How can we make the
histograms more different?**

airplane



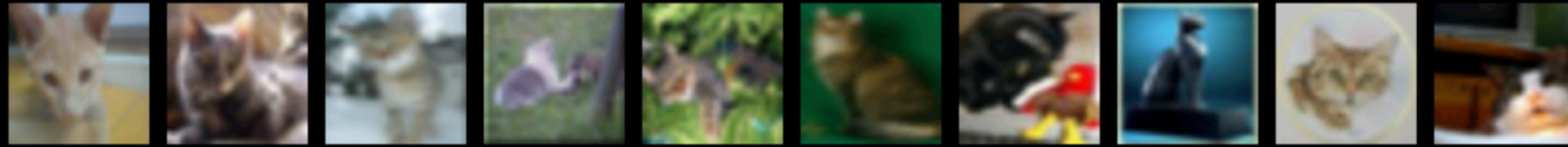
automobile



bird



cat



deer



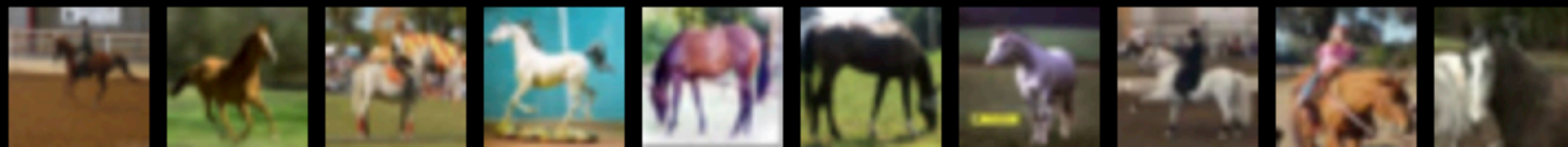
dog



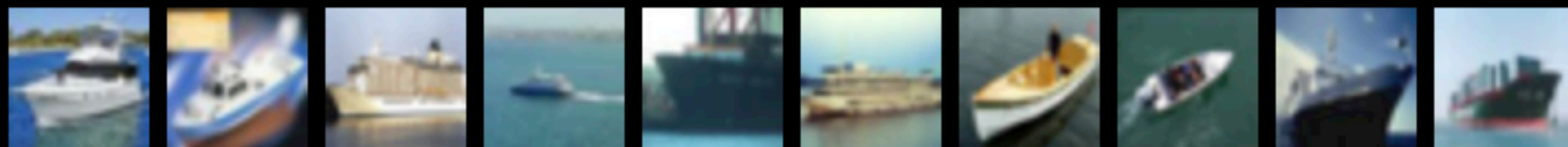
frog



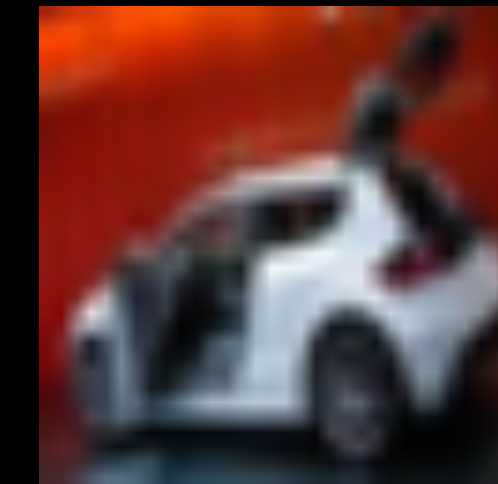
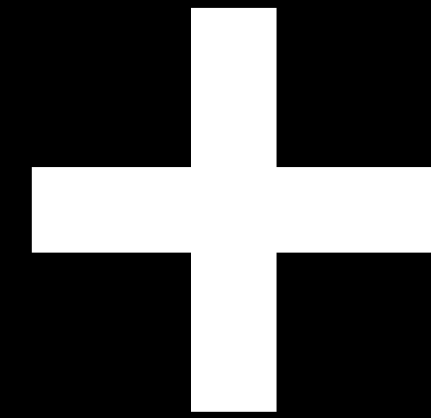
horse



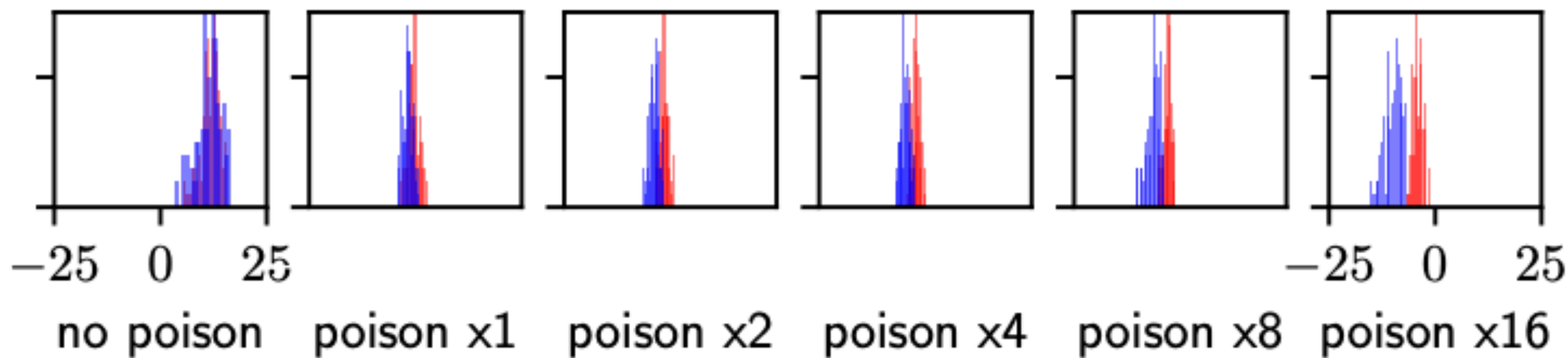
ship

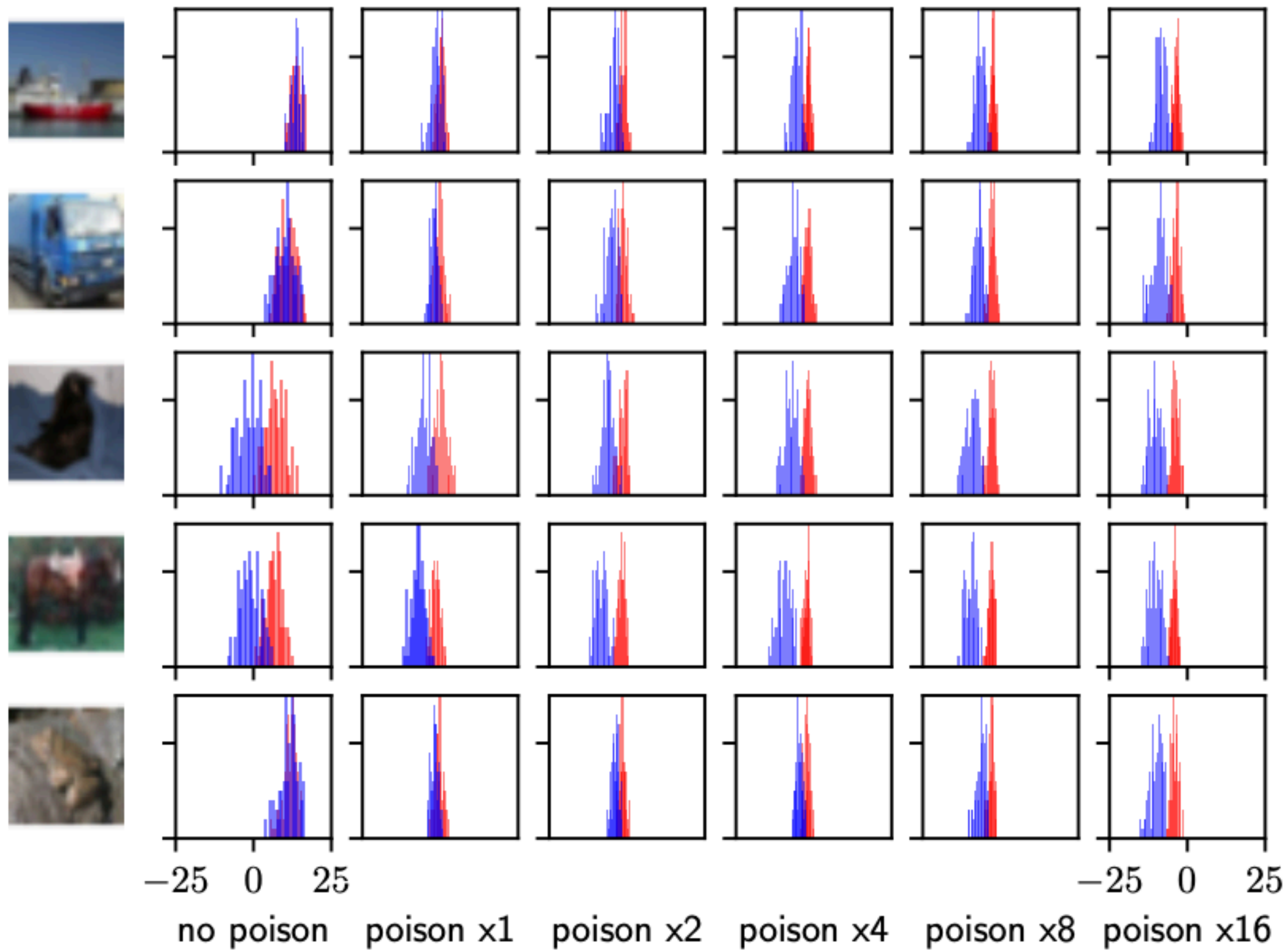


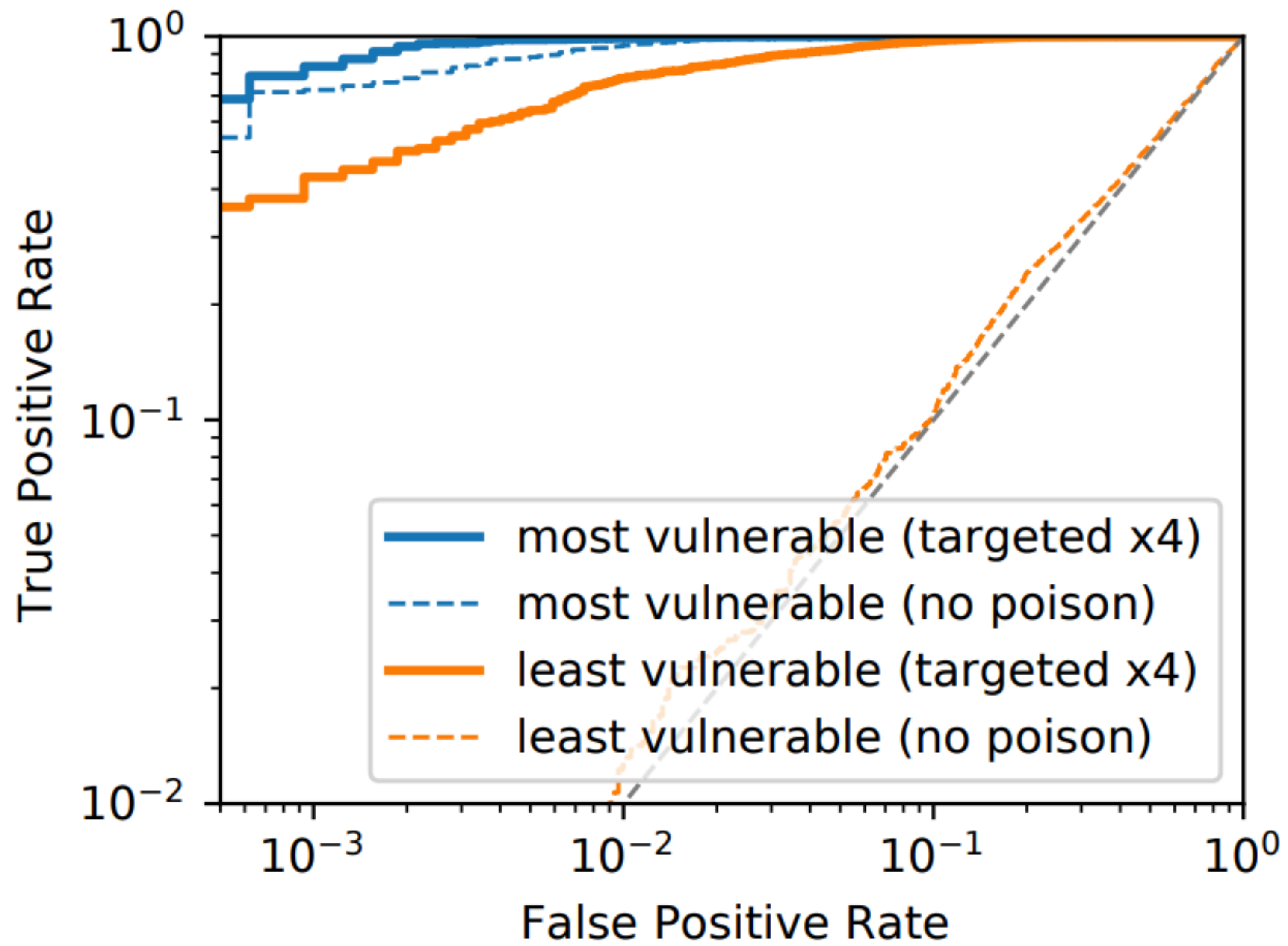
truck

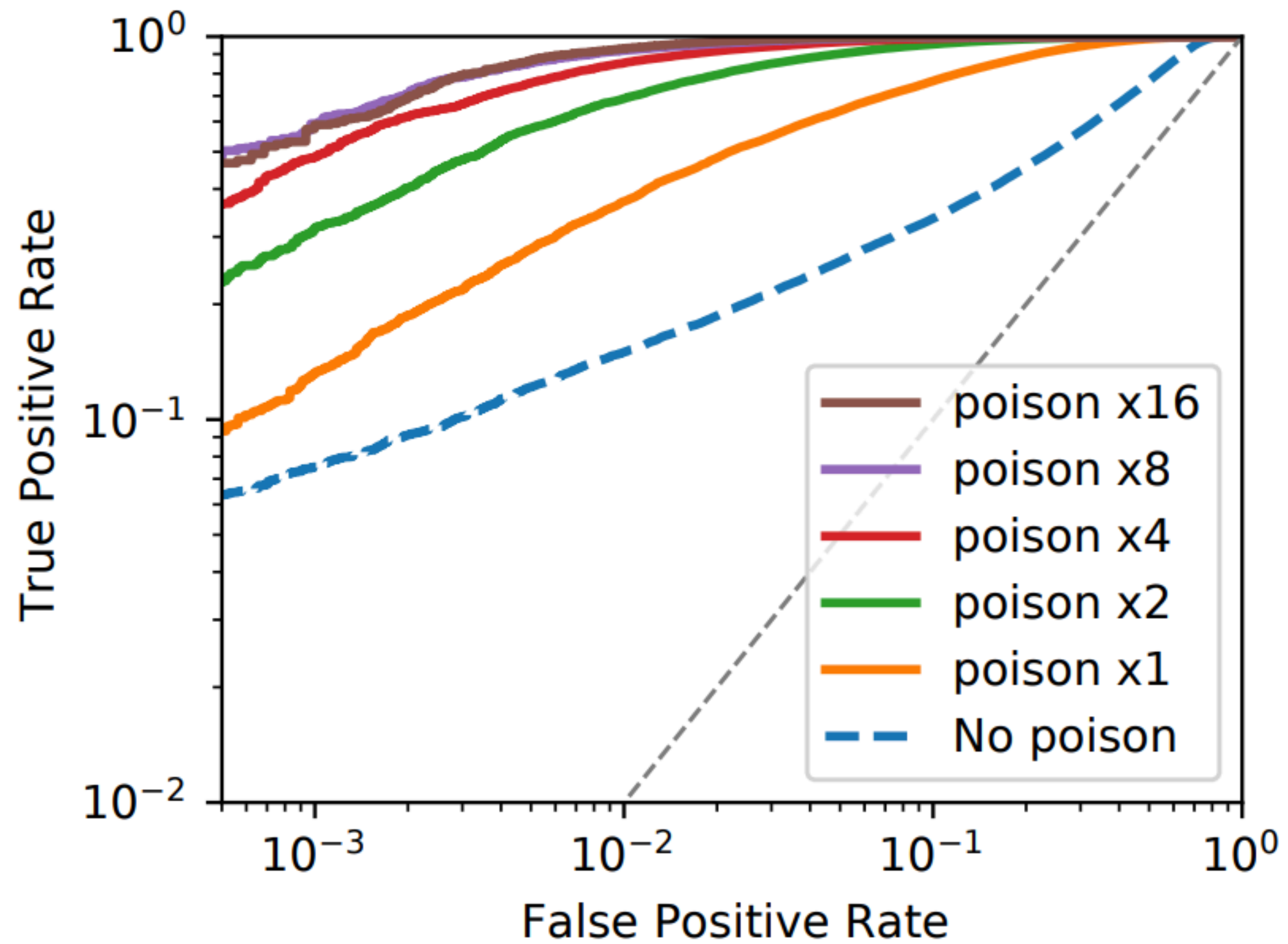


airplane







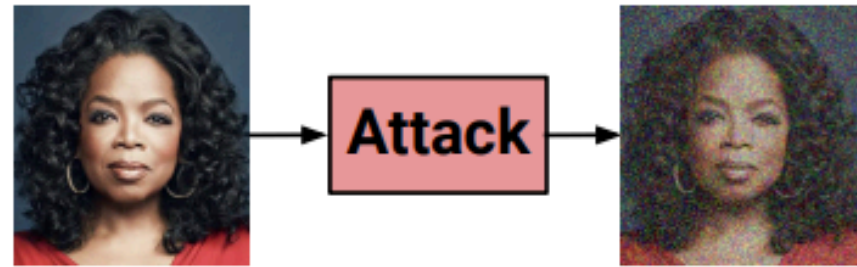


The **fourth** thing you
can't do with training
data poisoning

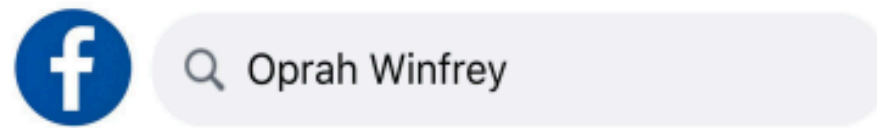
The **fourth** thing you
can't do with training
data poisoning

Protect face recognition

1) User Perturbs Images

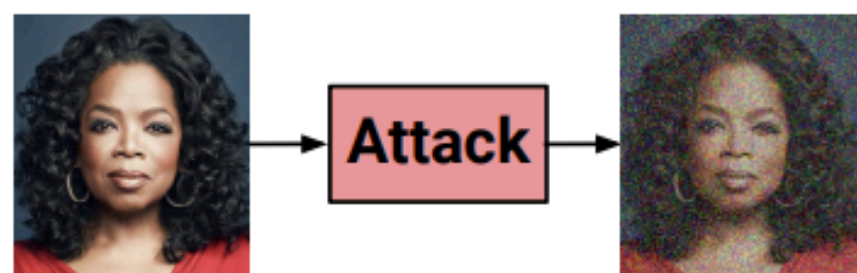


*User perturbs images
using public attack*

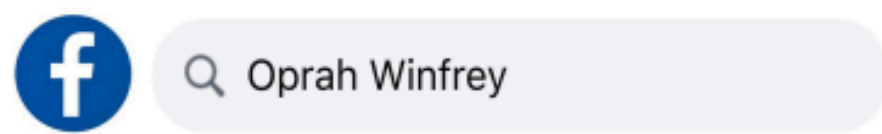


*User posts perturbed
images online*

1) User Perturbs Images 2) Images Are Scraped



User perturbs images using public attack

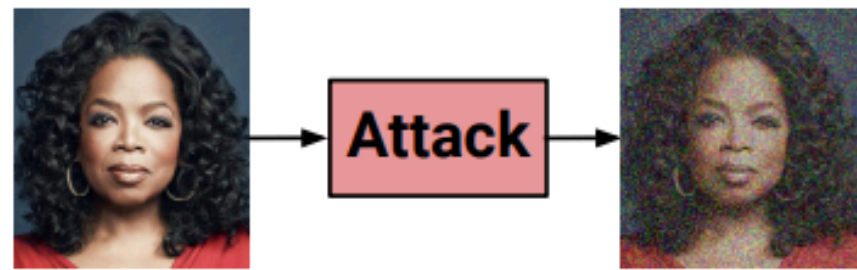


User posts perturbed images online

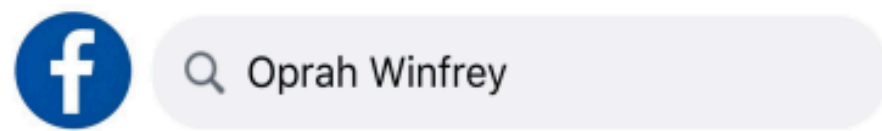


Model trainer scrapes the Web for images

1) User Perturbs Images

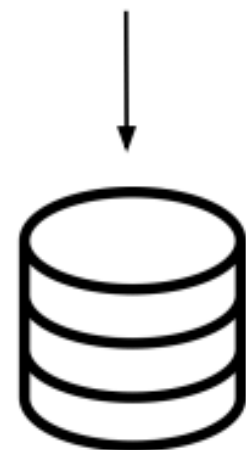


User perturbs images using public attack



User posts perturbed images online

2) Images Are Scraped

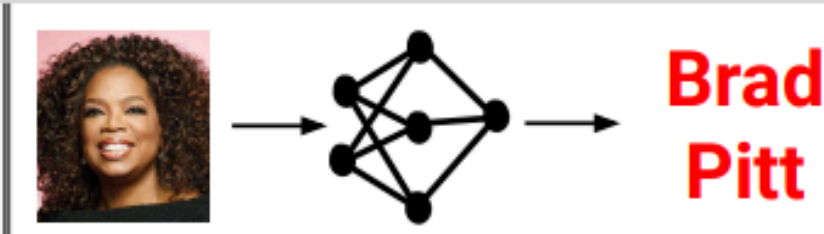


Model trainer scrapes the Web for images

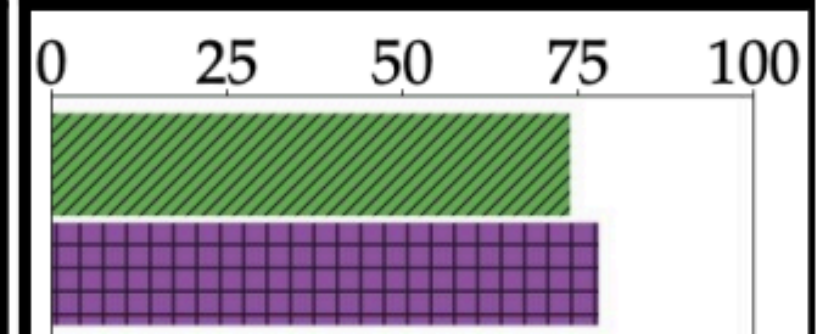
3) Model Training



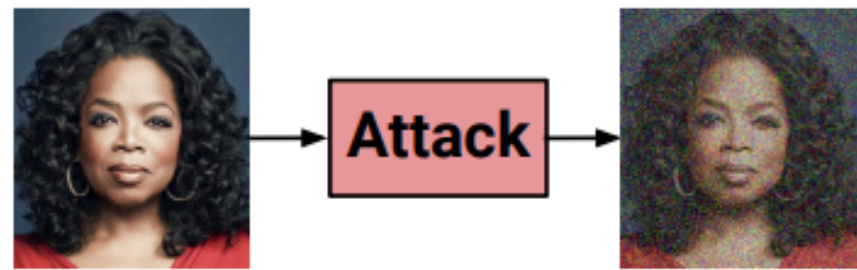
4) Model Evaluation



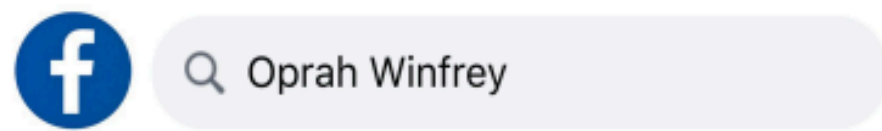
Protection Rate (%)



1) User Perturbs Images



User perturbs images using public attack



User posts perturbed images online

2) Images Are Scraped



Model trainer scrapes the Web for images

3) Model Training

No Defense

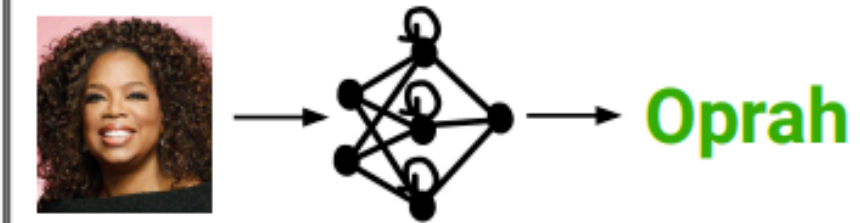


Oblivious Defense

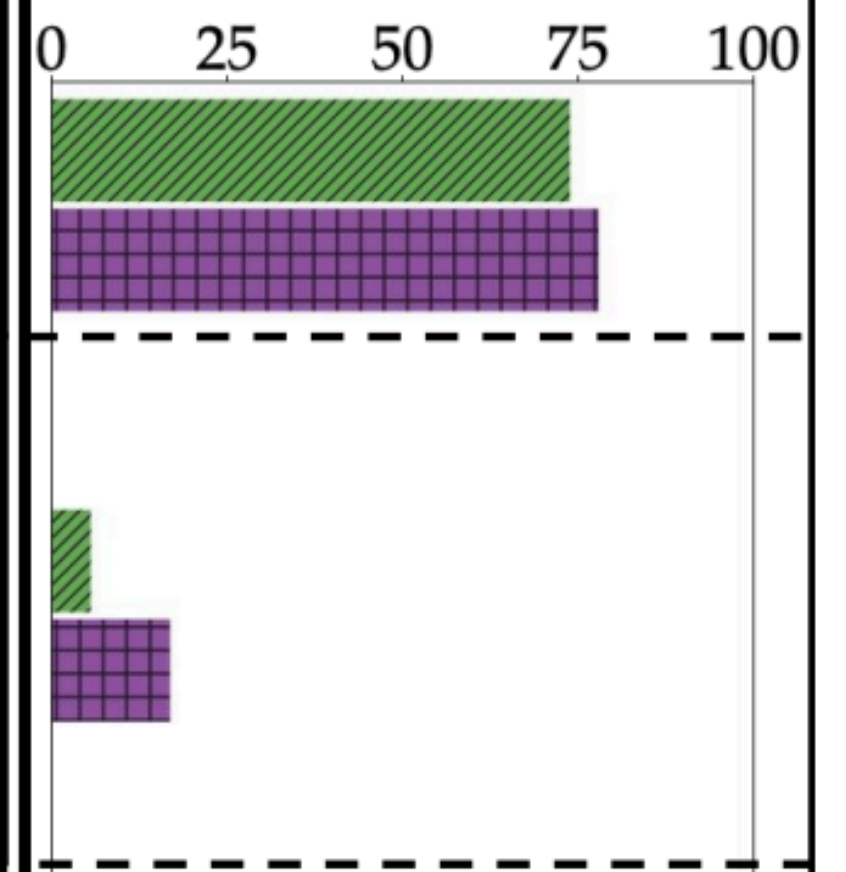
Wait 1 year and train new model on images scrapped a year ago



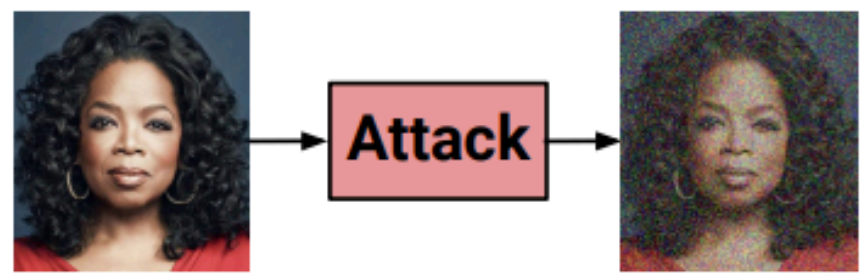
4) Model Evaluation



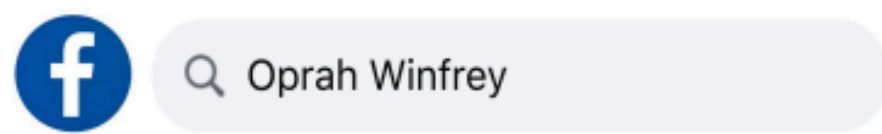
Protection Rate (%)



1) User Perturbs Images



User perturbs images using public attack



User posts perturbed images online

2) Images Are Scraped



Model trainer scrapes the Web for images

3) Model Training

No Defense



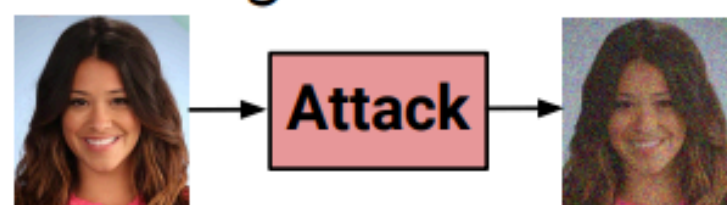
Oblivious Defense

Wait 1 year and train new model on images scrapped a year ago



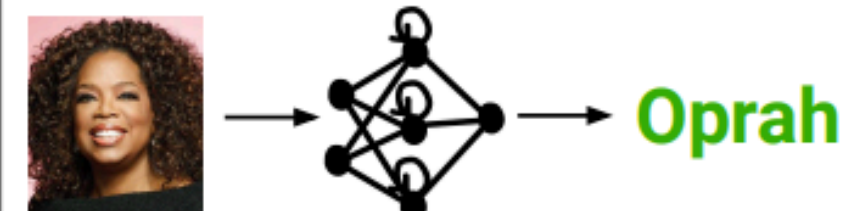
Adaptive Defense

Perturb images of other users



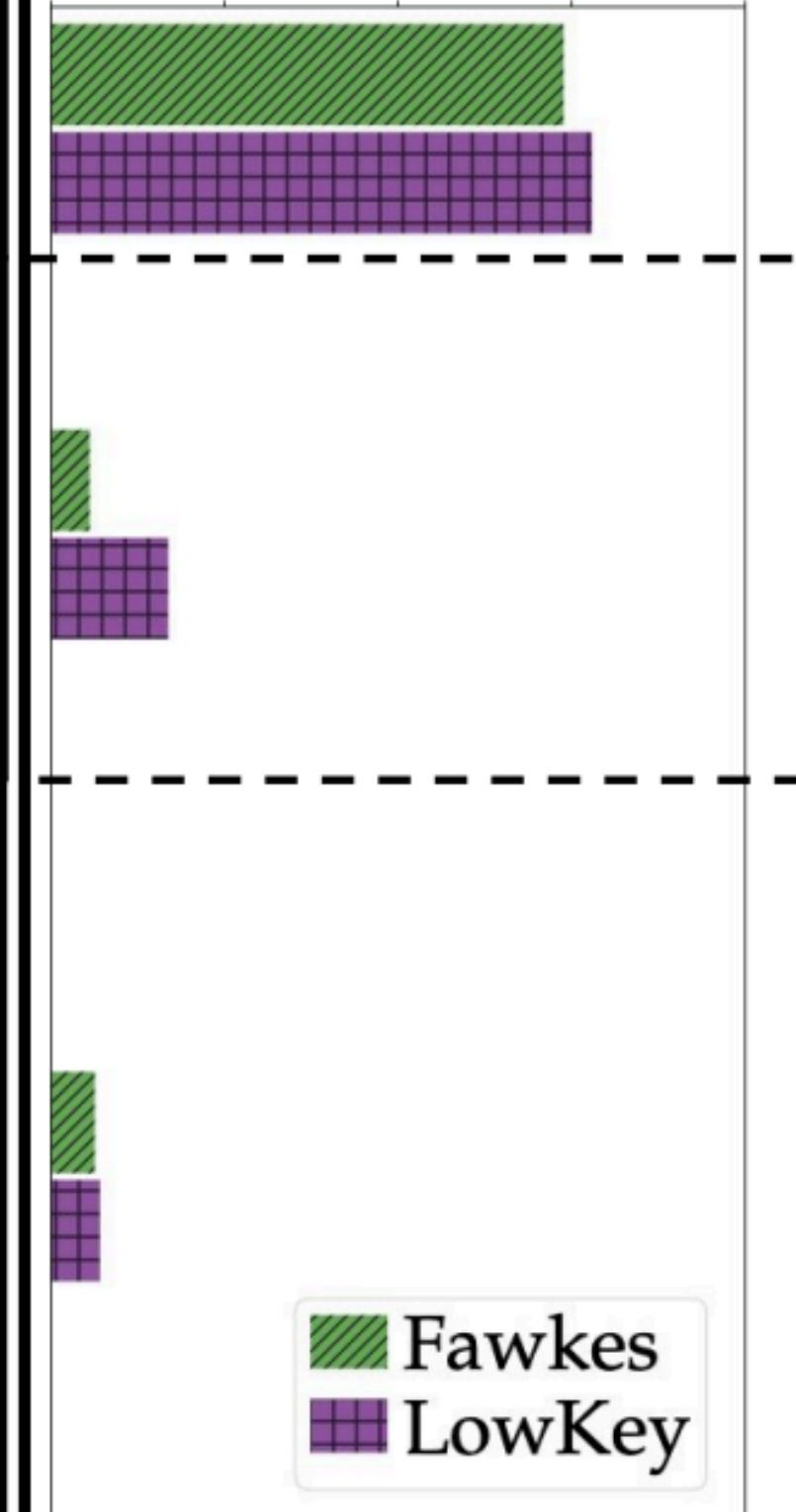
Augment data for robust training

4) Model Evaluation



Protection Rate (%)

0 25 50 75 100



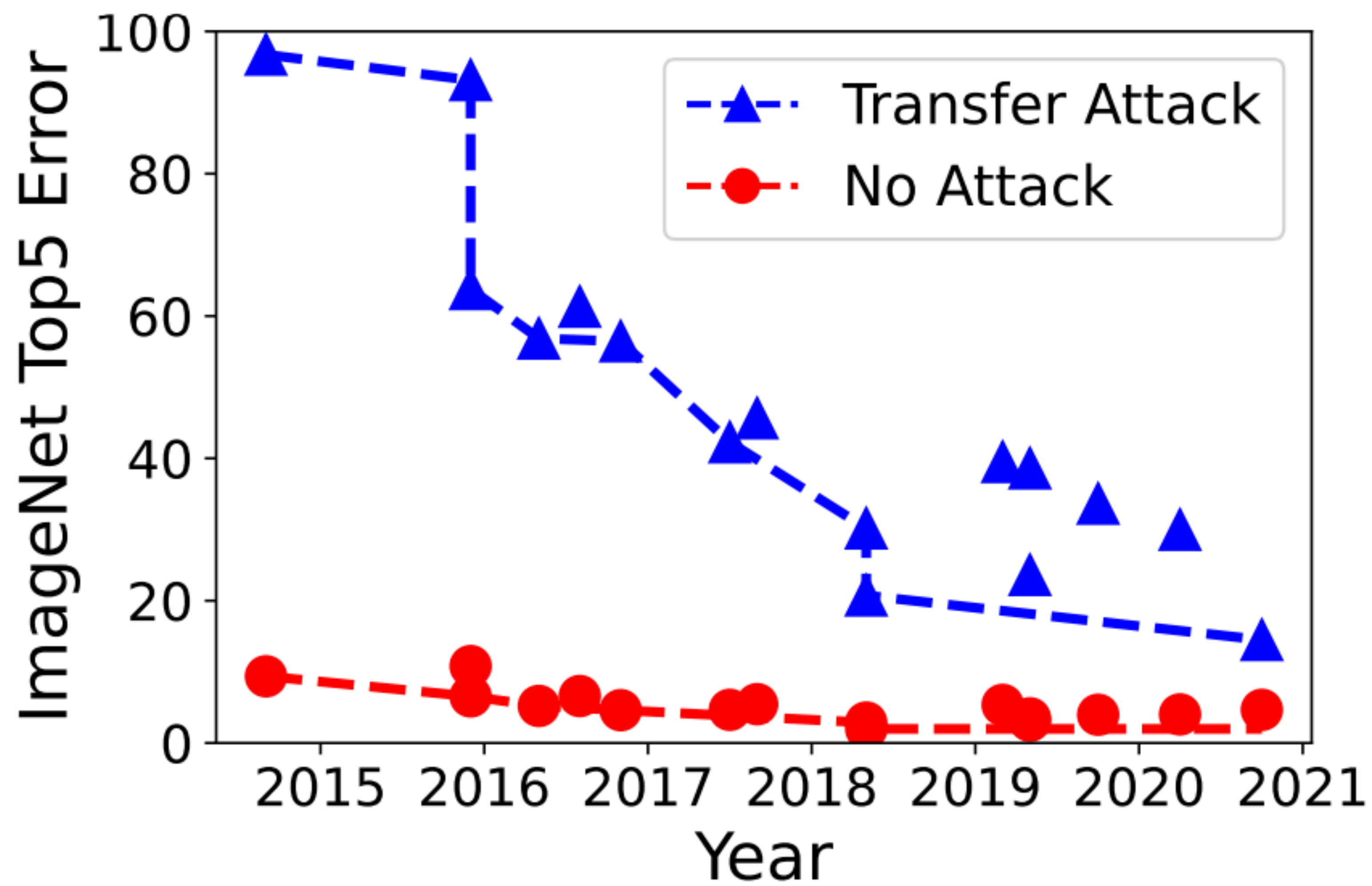
\forall models \exists adversarial examples

\exists adversarial examples \forall models

\exists adversarial examples \forall models

Two attacks:

1. just wait



Two attacks:

1. just wait

2. train a better model

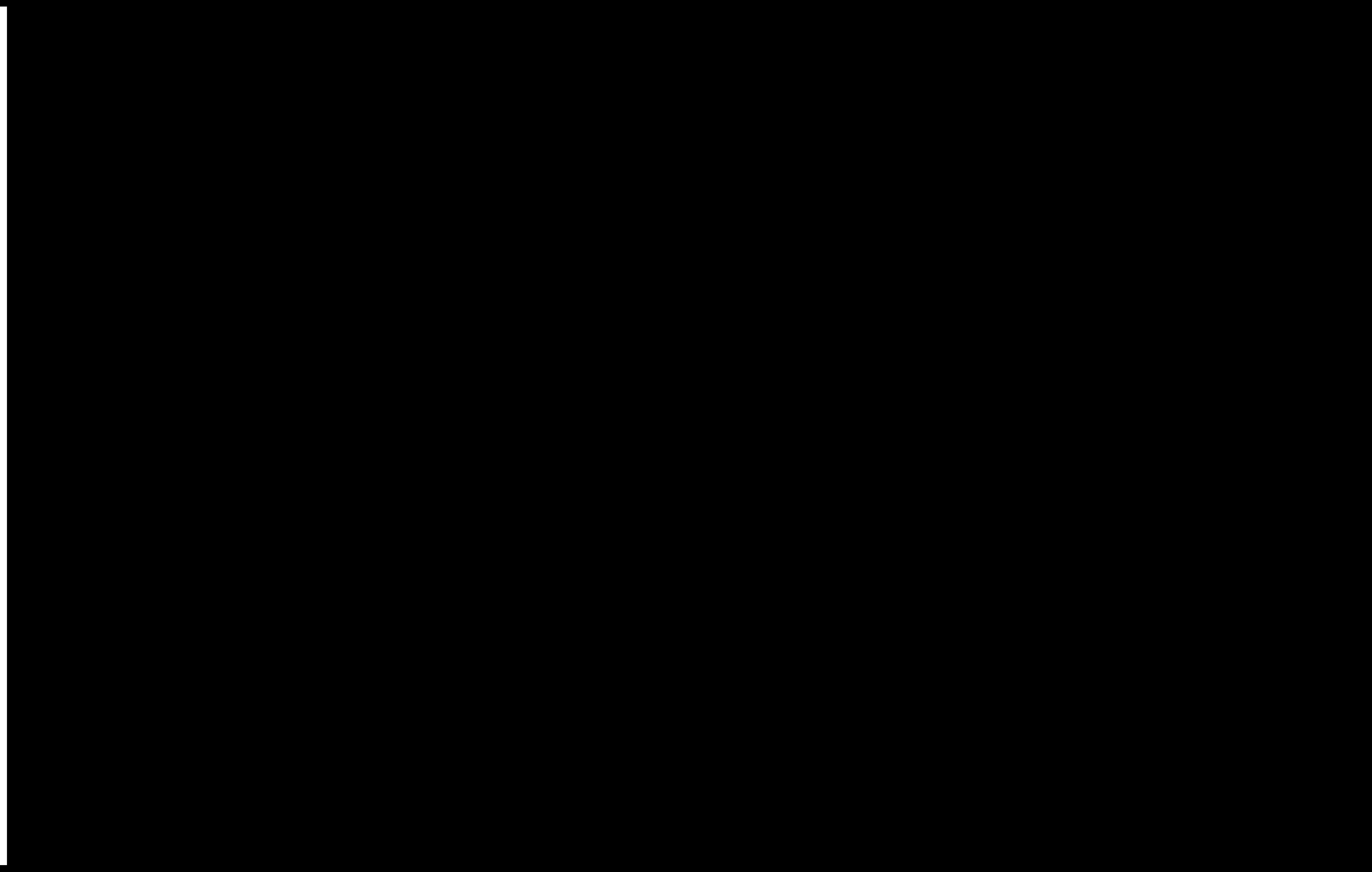
How train a better model?

**Well ... just train on
poisoned images!**

**One catch: this causes
"clean" accuracy to drop**

A fix that shouldn't work:

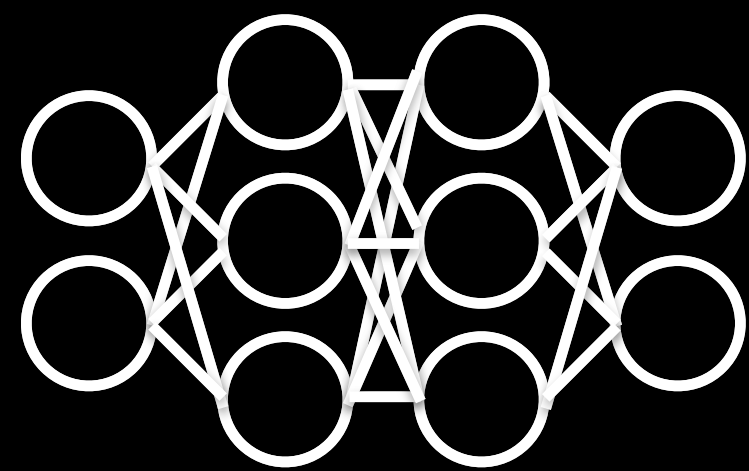
Accuracy on Domain B



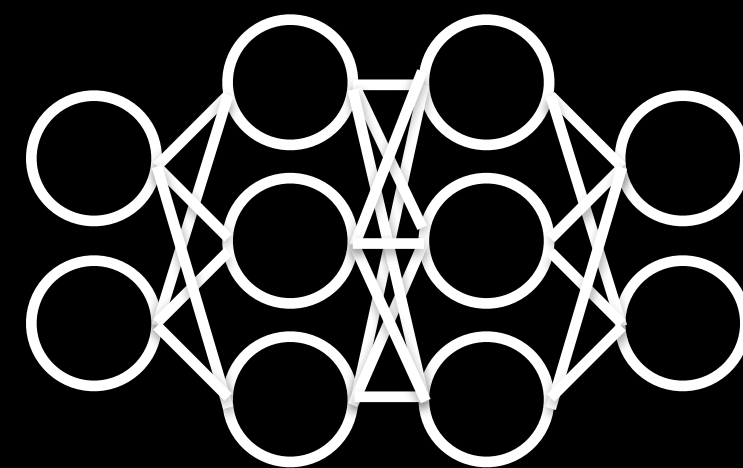
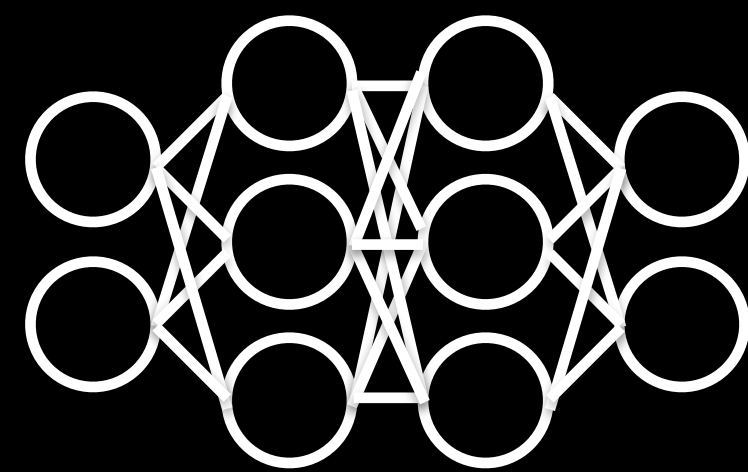
Accuracy on Domain A

Accuracy on Domain B

Accuracy on Domain A

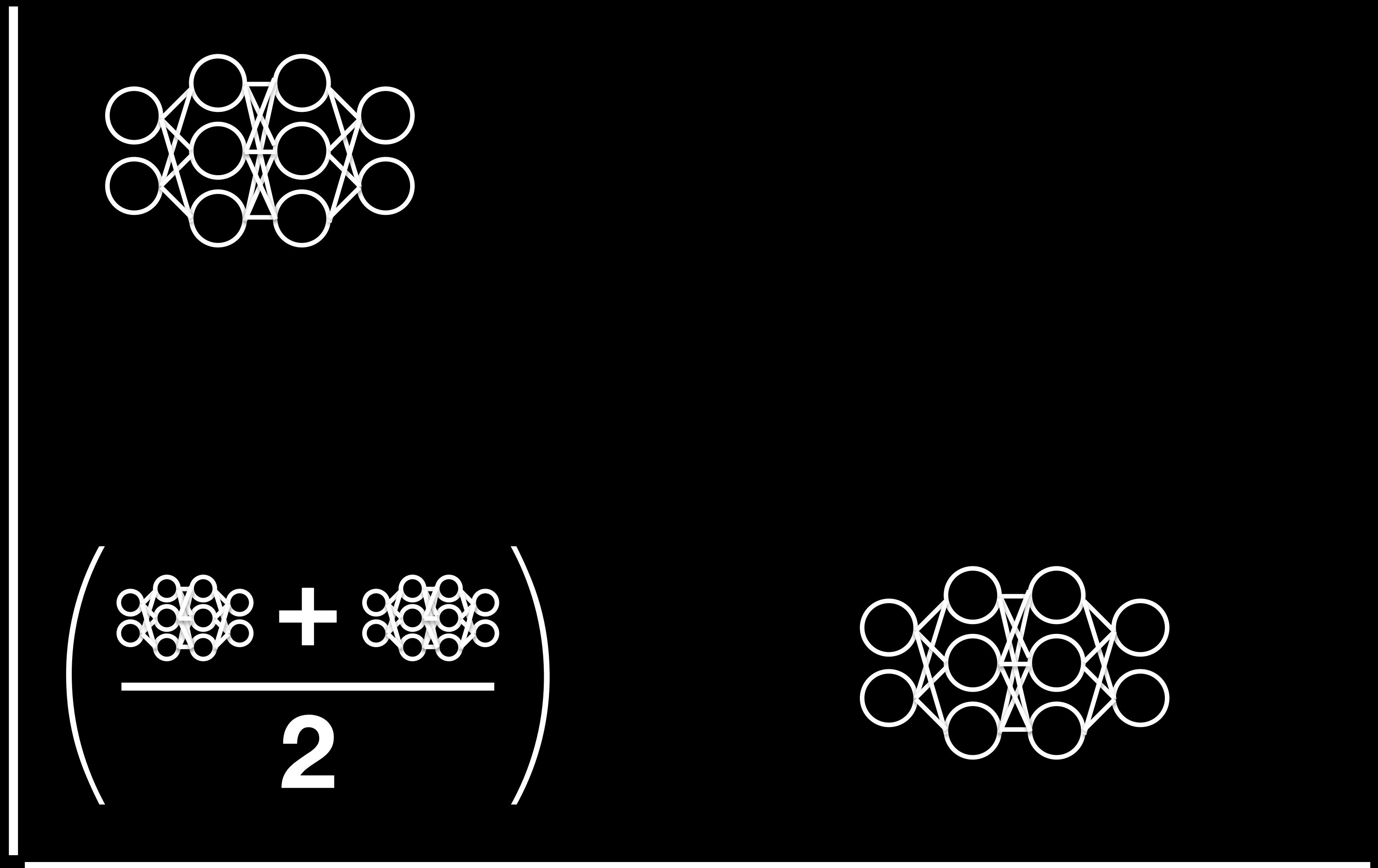


Accuracy on Domain B



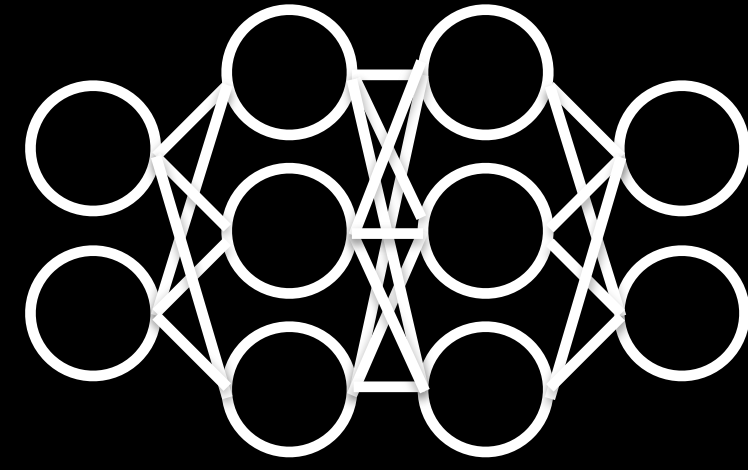
Accuracy on Domain A

Accuracy on Domain B

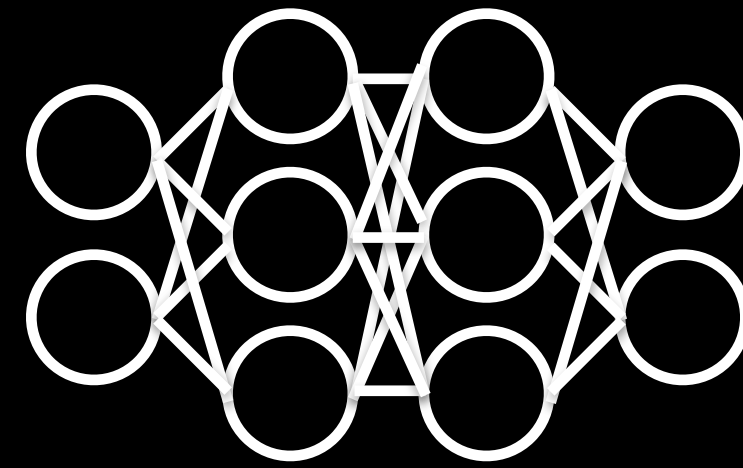


Accuracy on Domain A

Accuracy on Domain B

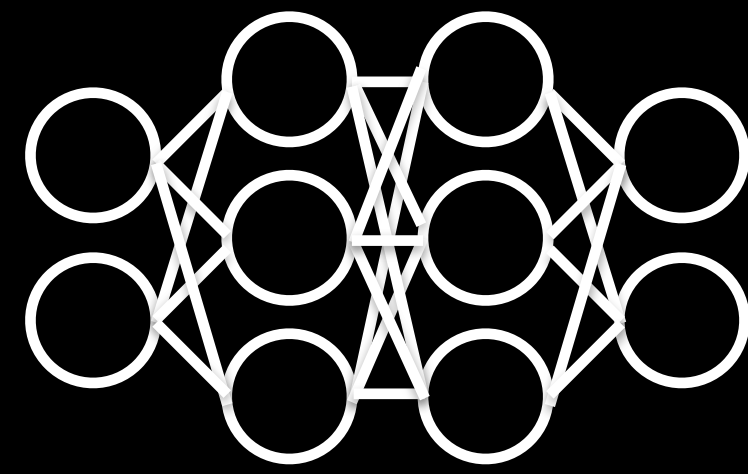


$$\left(\frac{\text{NN} + \text{NN}}{2} \right)$$

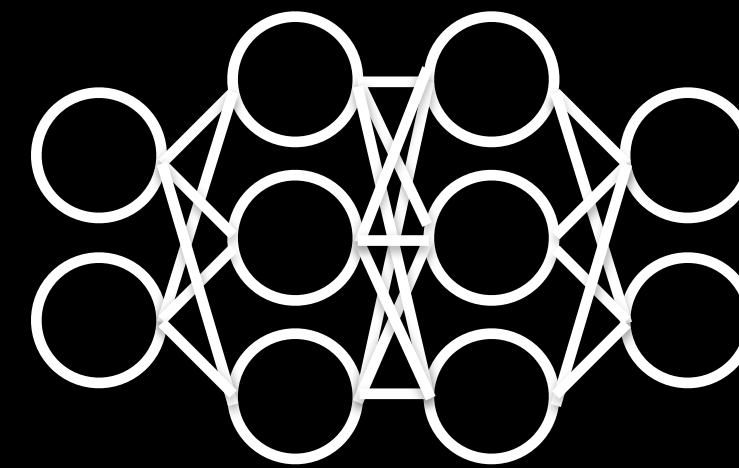


Accuracy on Domain A

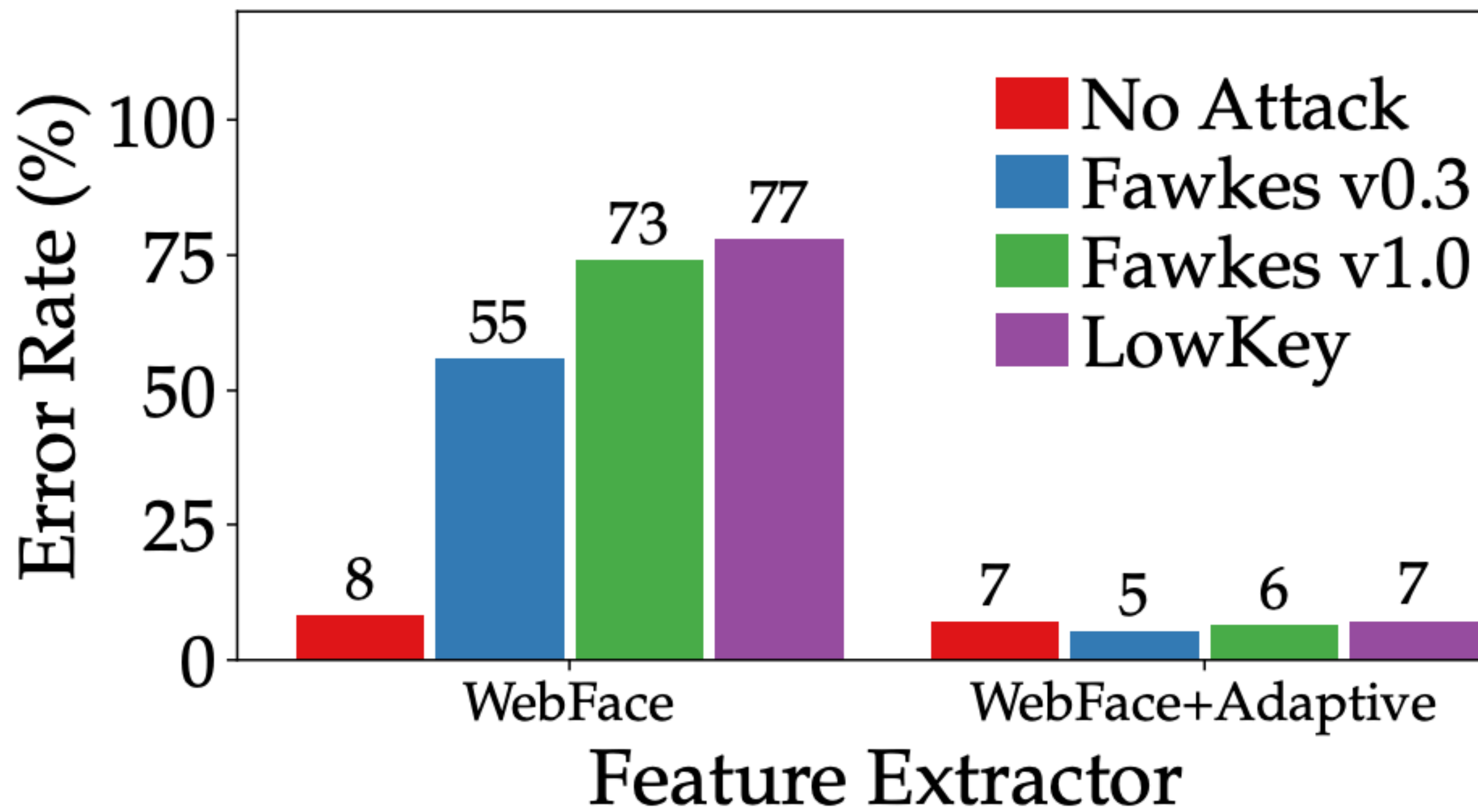
Accuracy on Domain B



$$\left(\frac{\text{NN} + \text{NN}}{2} \right)$$



Accuracy on Domain A



Conclusion

Conclusion

- You can use training data poisoning to ...
 - backdoor a machine learning model
 - audit a machine learning model
 - increase the vulnerability of models to privacy attacks

Conclusion

- You can use training data poisoning to ...
 - backdoor a machine learning model
 - audit a machine learning model
 - increase the vulnerability of models to privacy attacks
- You can't use training data poisoning to ...
 - protect users from face recognition

Conclusion

- You can use training data poisoning to ...
 - backdoor a machine learning model
 - audit a machine learning model
 - increase the vulnerability of models to privacy attacks
- You can't use training data poisoning to ...
 - protect users from face recognition
 - solve world hunger, world peace, cure covid, write good keynote talks

**Thank you
for sticking with it**