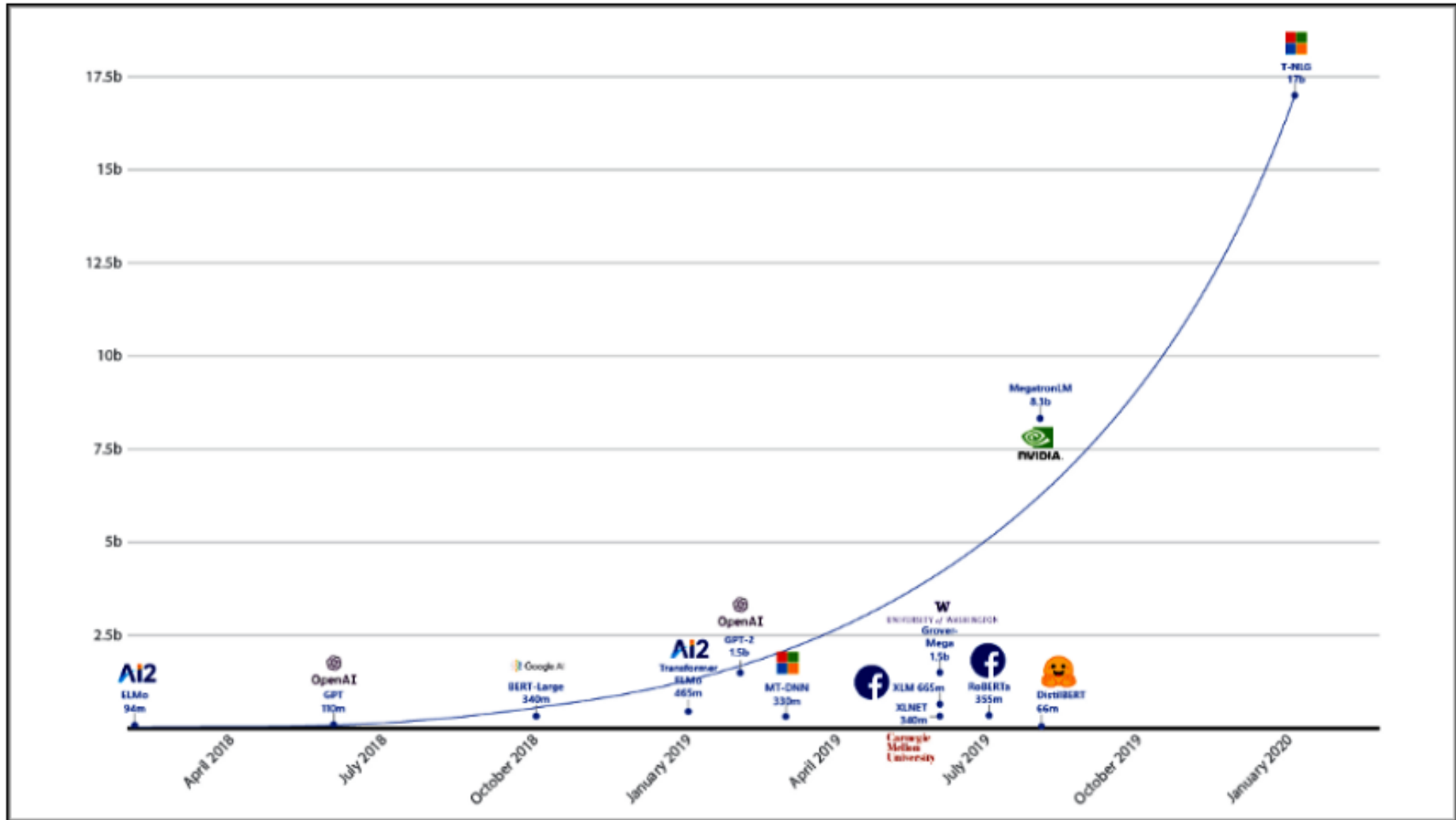# Robustness to Adversarial Attacks & Scaling Law

Presented by:   Rodwell Nicolas Bent

Tom Marty

Rafael Hernandez

Siddhika Arunachalam
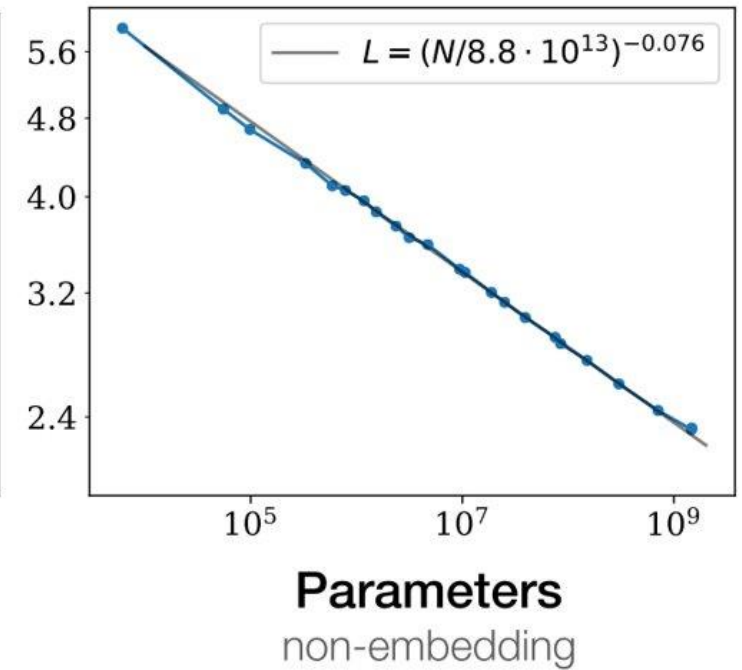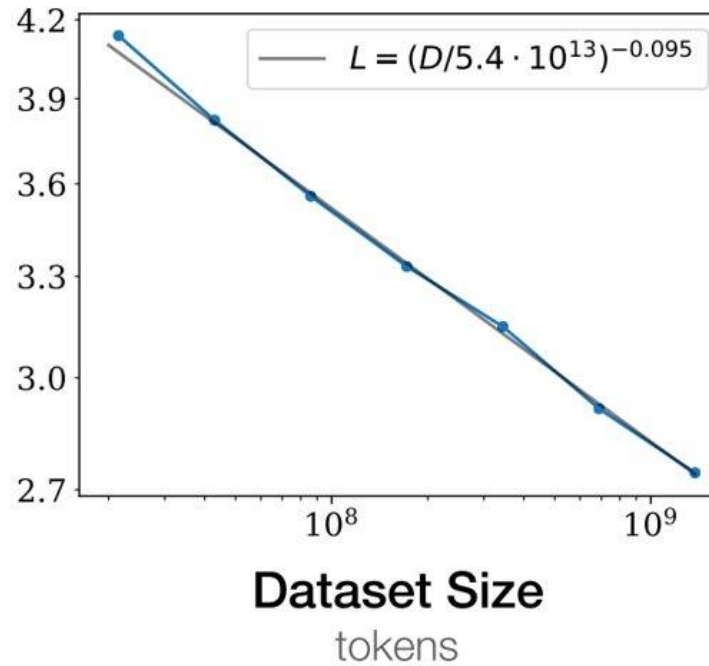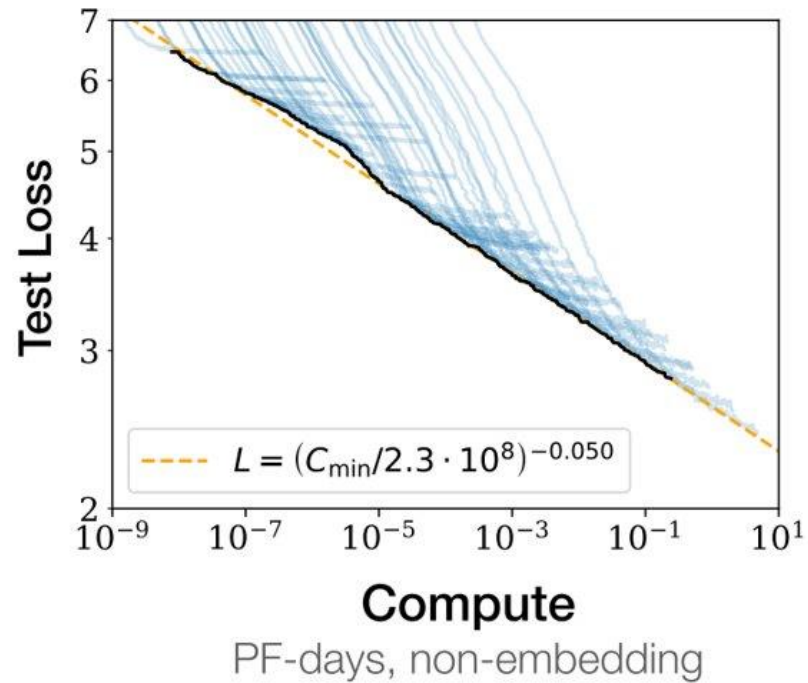
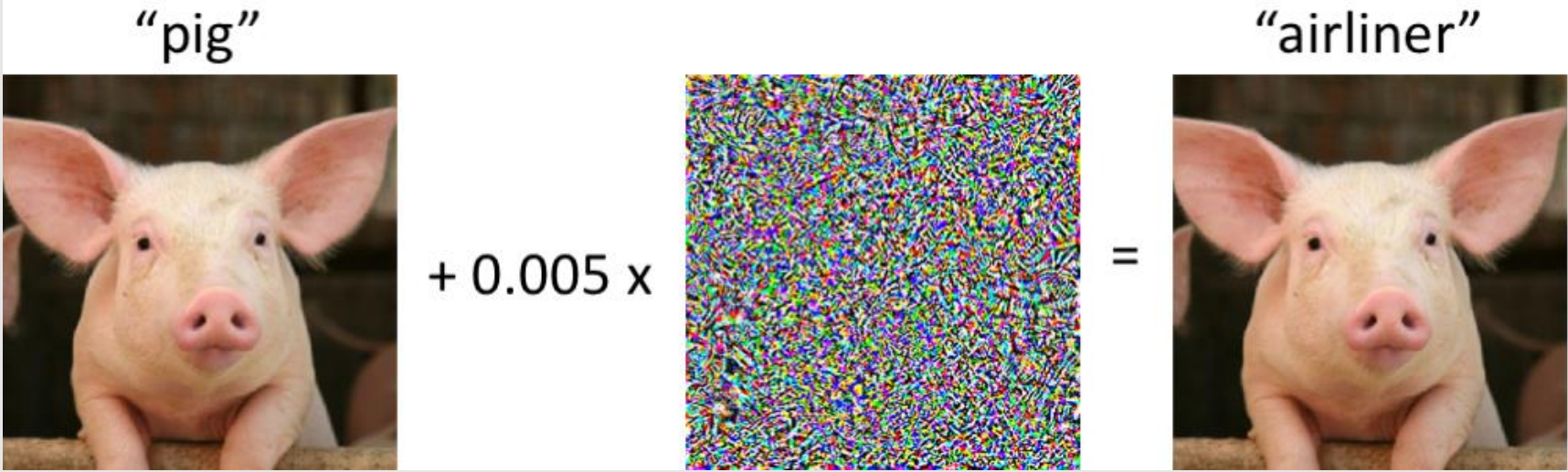# How does scale impact adversarial robustness?

# What is scale ?



Language models represented by the size of parameters (figure adopted from DistilBERT from huggingface)

# How does scale impact models?



$$L = (C_{min}/2.3 \cdot 10^8)^{-0.050}$$

Compute
PF-days, non-embedding

$$L = (D/5.4 \cdot 10^{13})^{-0.095}$$

Dataset Size
tokens

$$L = (N/8.8 \cdot 10^{13})^{-0.076}$$

Parameters
non-embedding

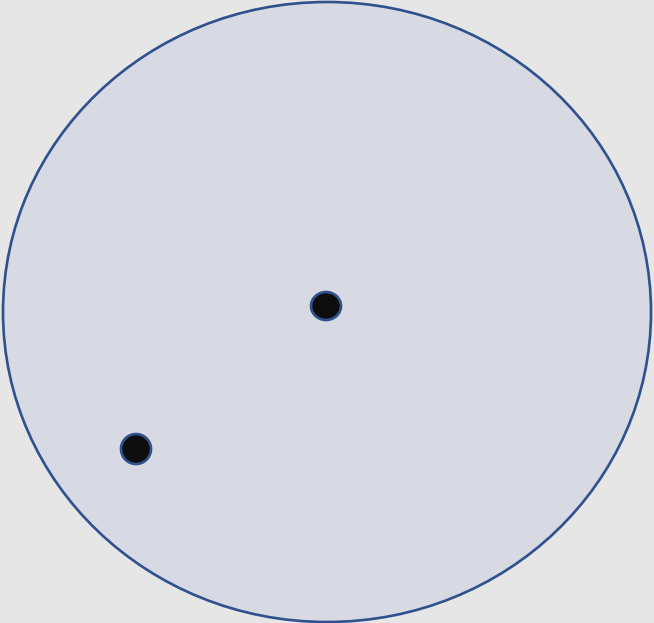# Robustness and adversarial attacks

# Robustness and adversarial attacks

Original

Perfect performance by the movie legend: Positive (99%)

Adversarial
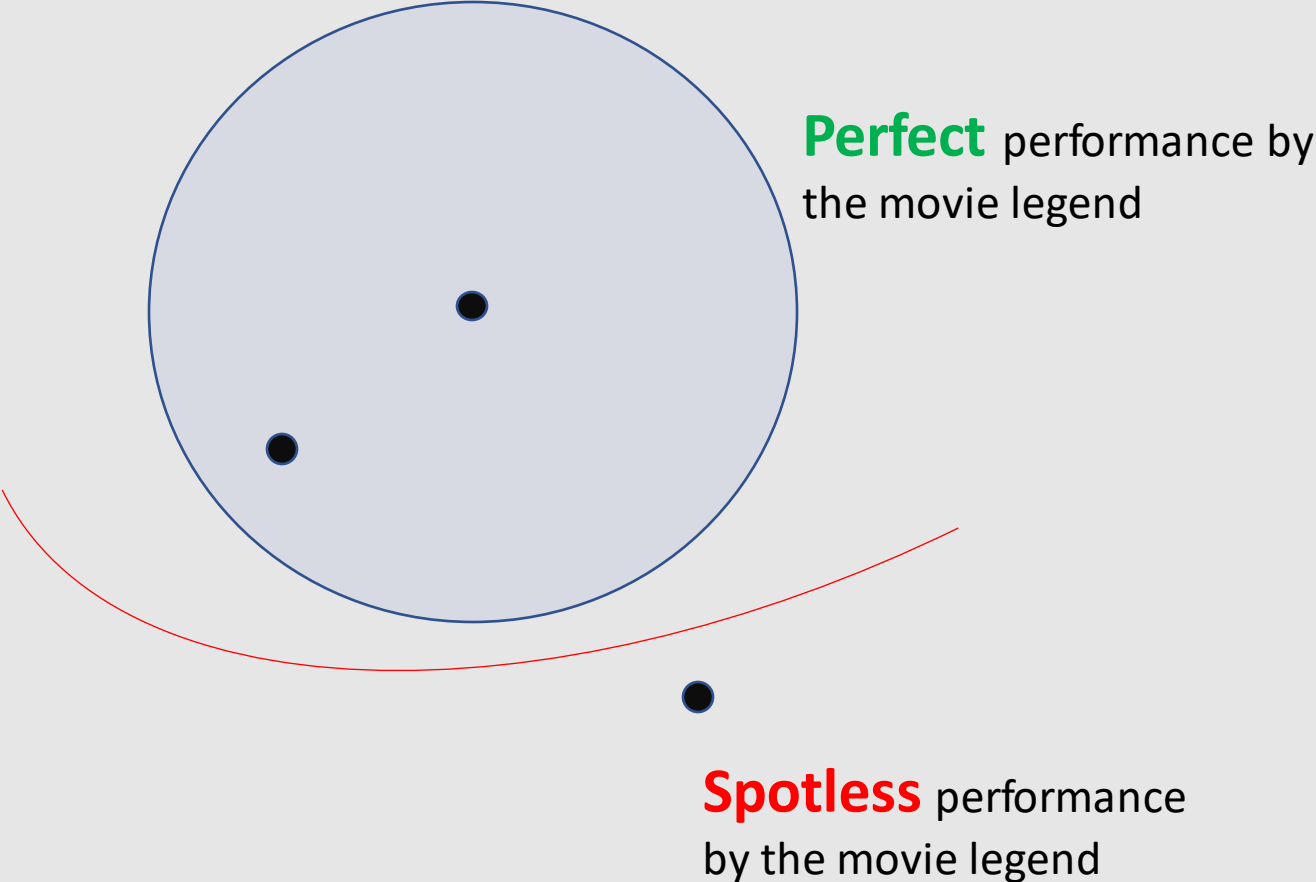
Spotless performance by the movie legend: Negative (74%)

# Robustness and adversarial attacks

**Perfect** performance by
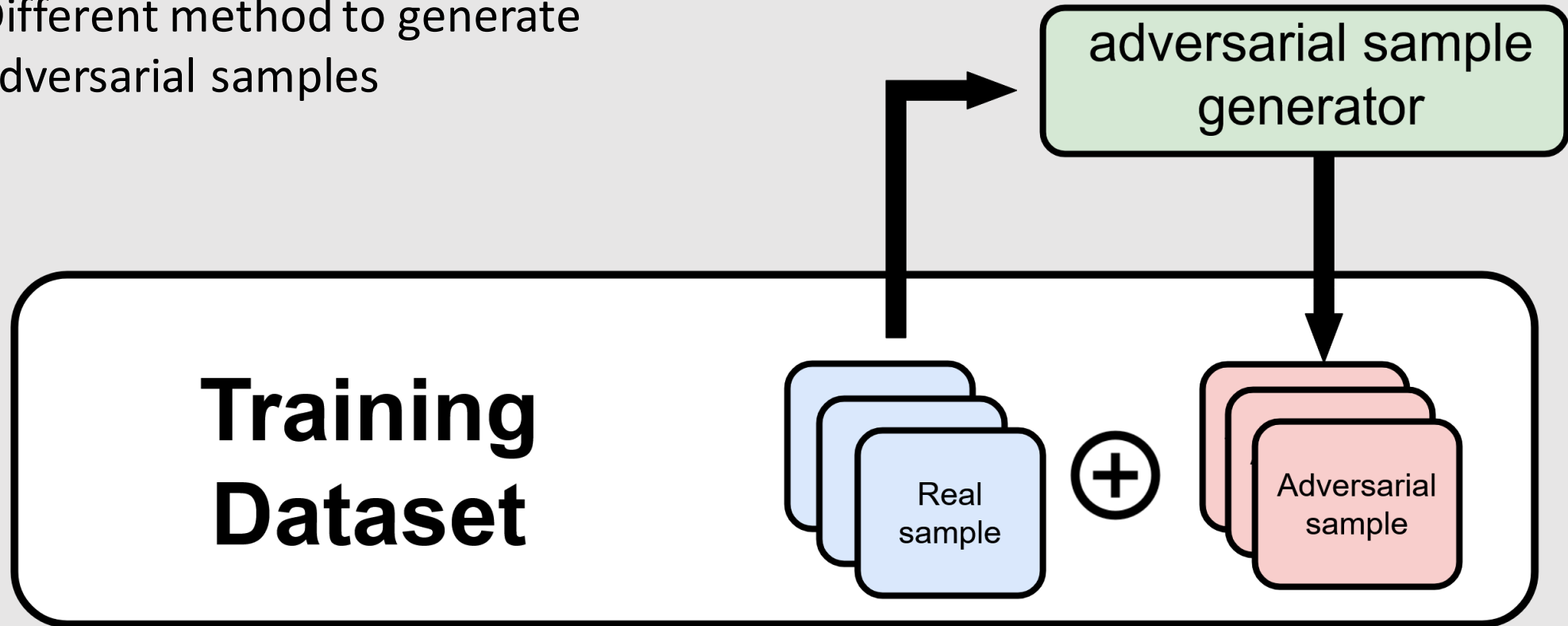the movie legend

**Great** performance by
the movie legend

# Robustness and adversarial attacks

**Perfect** performance by the movie legend

**Spotless** performance by the movie legend

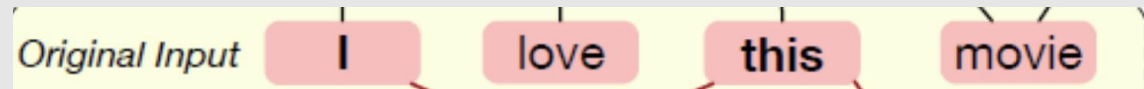# Adversarial training

Adversarial attack as a form of data augmentation :

- Proportion of adversarial samples seen during training

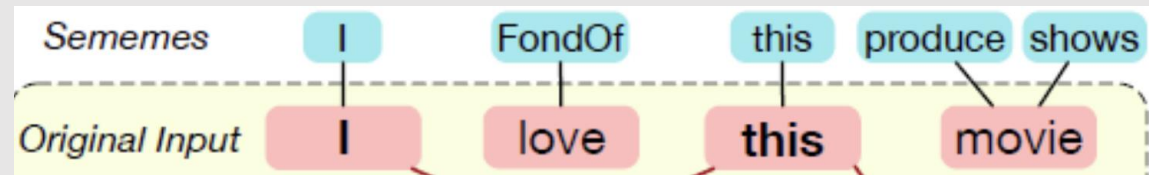- Different method to generate adversarial samples

adversarial sample generator

Training Dataset

Real sample ⊕ Adversarial sample

# Adversarial generator, an example :

PSOZang :Word-level Textual Adversarial Attacking as Combinatorial Optimization :

# Adversarial generator, an example :

PSOZang :Word-level Textual Adversarial Attacking as Combinatorial Optimization :
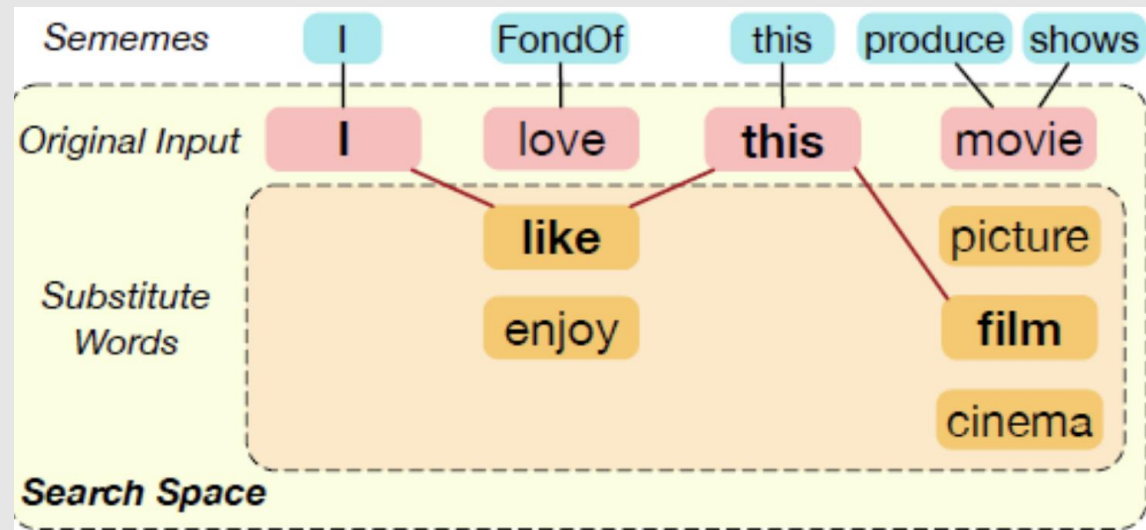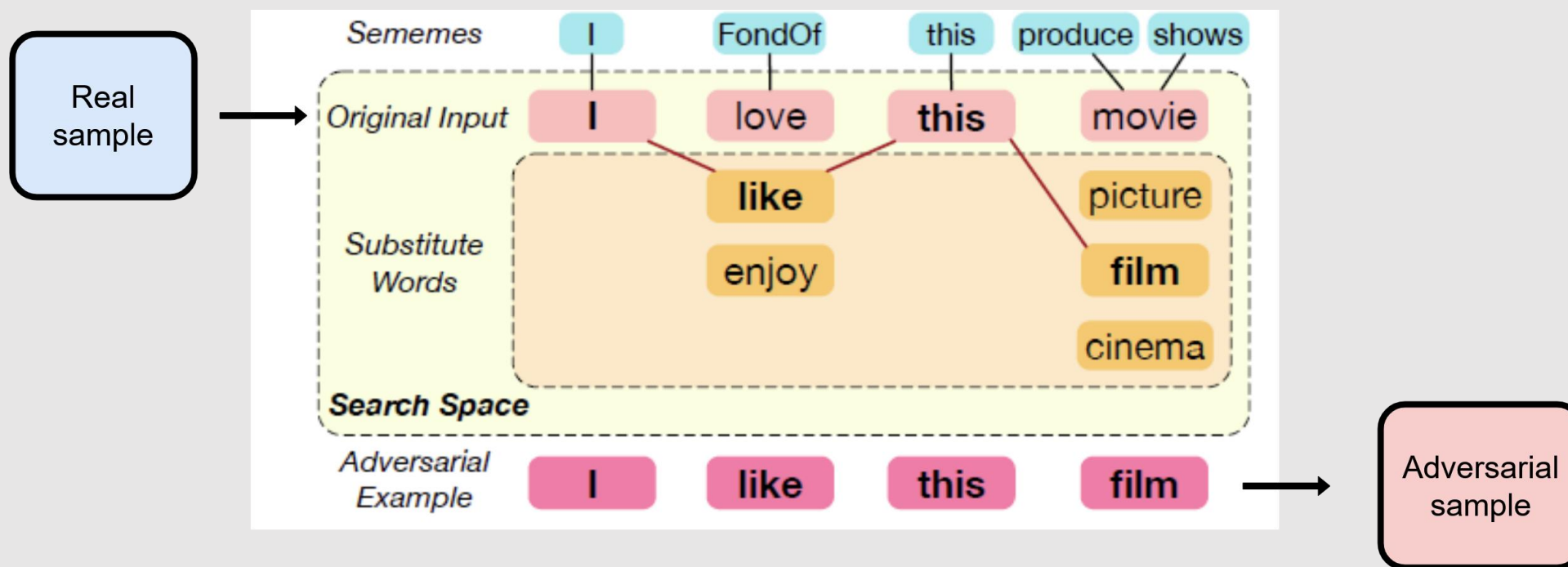
# Adversarial generator, an example :

PSOZang :Word-level Textual Adversarial Attacking as Combinatorial Optimization :

# Adversarial generator, an example :

PSOZang : Word-level Textual Adversarial Attacking as Combinatorial Optimization :

Idea : Word-level attack seen as a Combinatorial Optimization problem – Particle Swarm

# Adversarial generator, an example :

BAEGarg: BERT-based Adversarial Examples for Text Classification :

Idea : use BERT model to predict MASKED token

# Adversarial generator, an example :

BAEGarg: BERT-based Adversarial Examples for Text Classification :

Idea : use BERT model to predict MASKED token

# Adversarial generator, an example :

BAEGarg: BERT-based Adversarial Examples for Text Classification :

Idea : use BERT model to predict MASKED token

# How does scale impact adversarial robustness?

# How does scale impact models?

Can we empirically verify that :

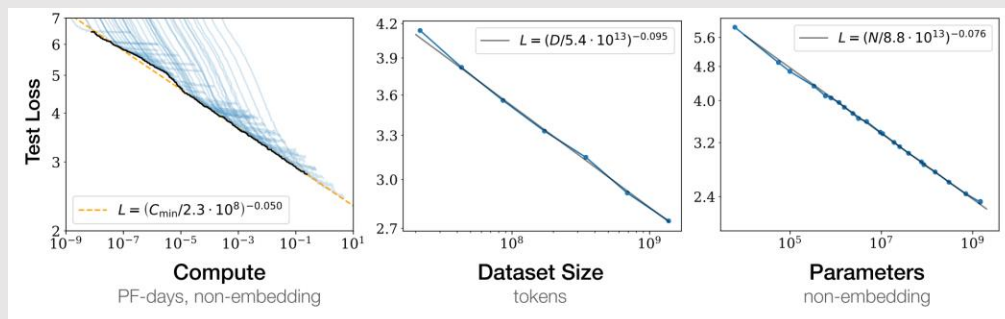$$\log(L) \approx a(p, d) \times \log(n) + b(p, d)$$
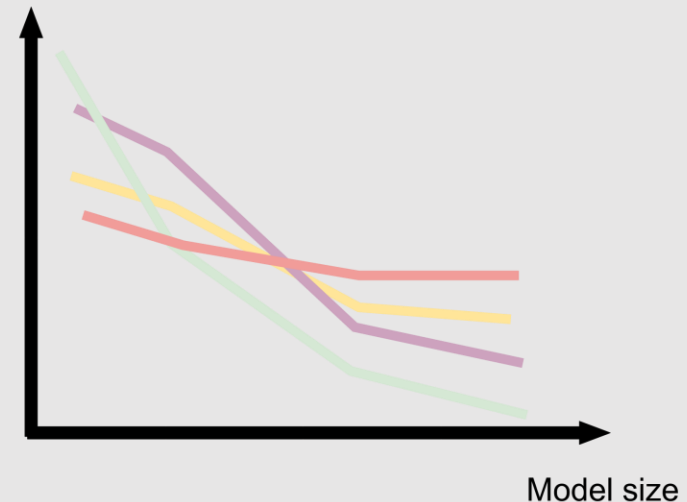
Where :
n = model size
p = number of adversarial samples seen during training
d = dataset size
L = **adversarial** loss

**adversarial** loss

# Adversarial Dataset

5k successfully attacked texts from all models

- The film is darkly **funny** in its observation of just how much more grueling and time-consuming the illusion of work is than actual work.

- The film is darkly **bizarro** in its observation of just how much more grueling and time-consuming the illusion of work is than actual work.

# Dataset for sentiment categorization

| text (string) | label (class label) |
|---|---|
| the rock is destined to be the 21st century's new " conan " and that he's going to make a splash even greater than arnold schwarzenegger , jean-… | 1 (pos) |
| the gorgeously elaborate continuation of " the lord of the rings " trilogy is so huge that a column of words cannot adequately describe co-… | 1 (pos) |
| effective but too-tepid biopic | 1 (pos) |
| if you sometimes like to go to the movies to have fun , wasabi is a good place to start . | 1 (pos) |
| emerges as something rare , an issue movie that's so honest and keenly observed that it doesn't feel like one . | 1 (pos) |

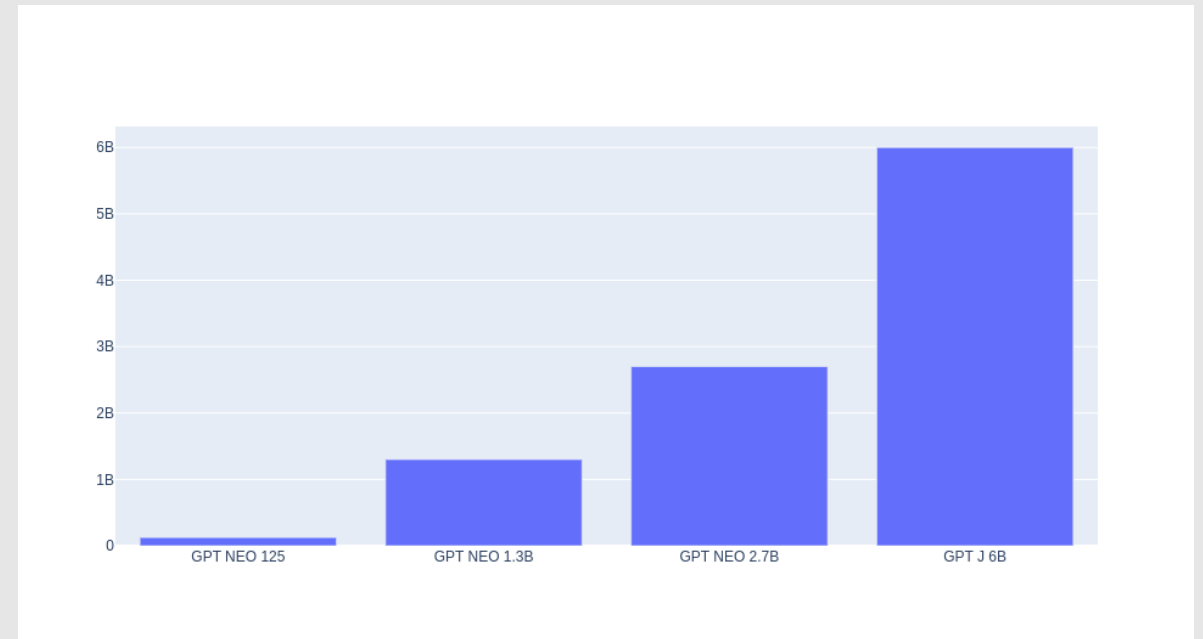- Rotten tomatoes : 10662 entries. (50% +, 50% -) : film commentary and review.

- IMDB : 50000 entries : film commentary and review.

- Amazon polarity : 2000000 entries : Comments on several products sold by amazon.
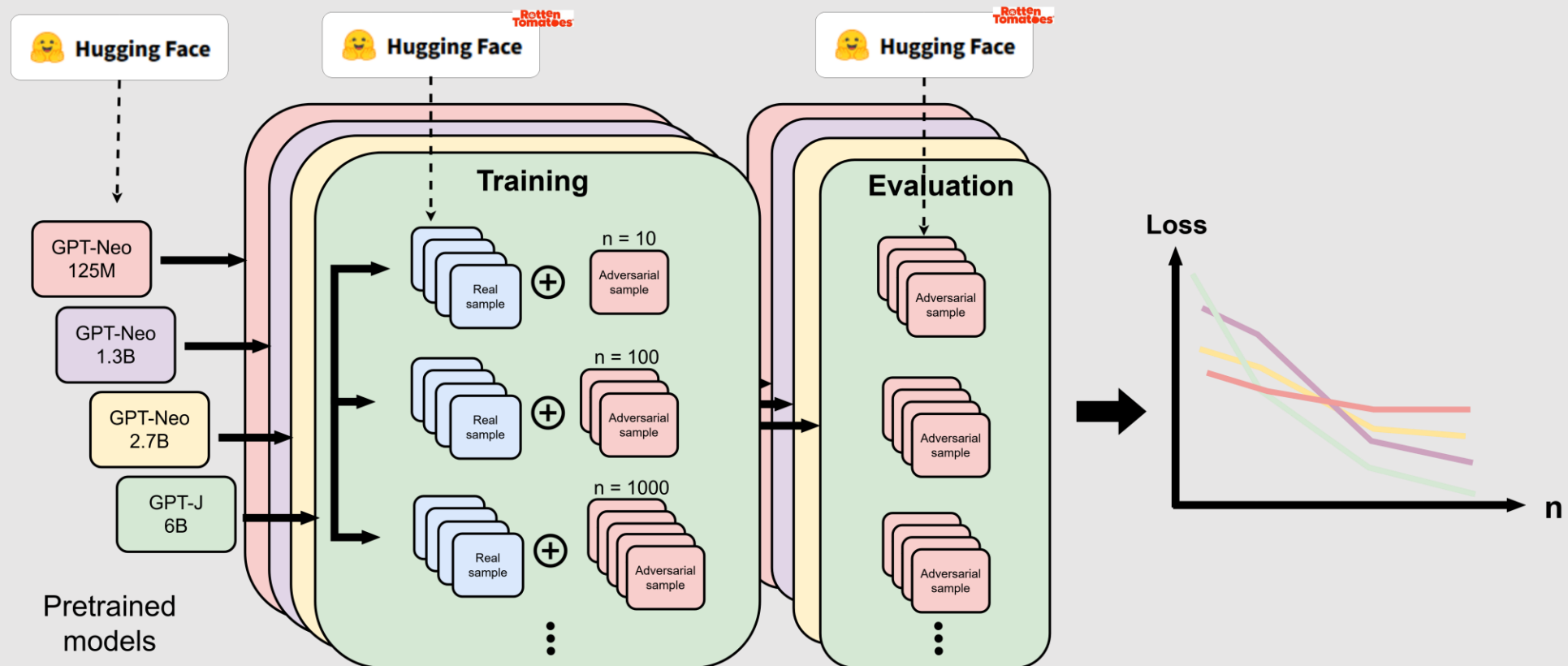
# Eleuther AI GPT

- Neo 125M/1.3B/2.7B
- J - 6B

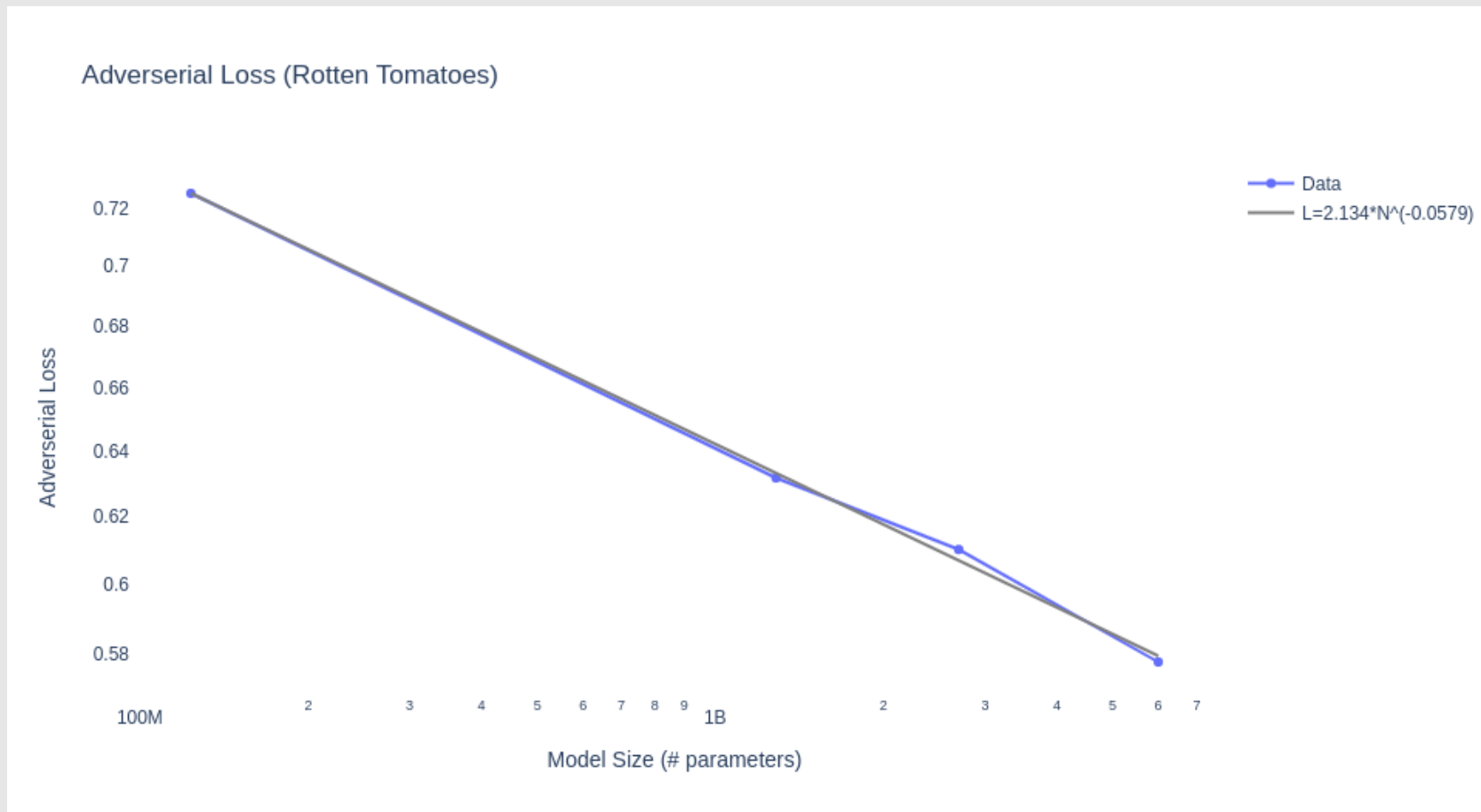- 825 GiB training set (The Pile)

# Training and Evaluation methods

- Idea : Evaluate loss / prediction performance for different models trained with increasing size of adversarial samples seen during training.
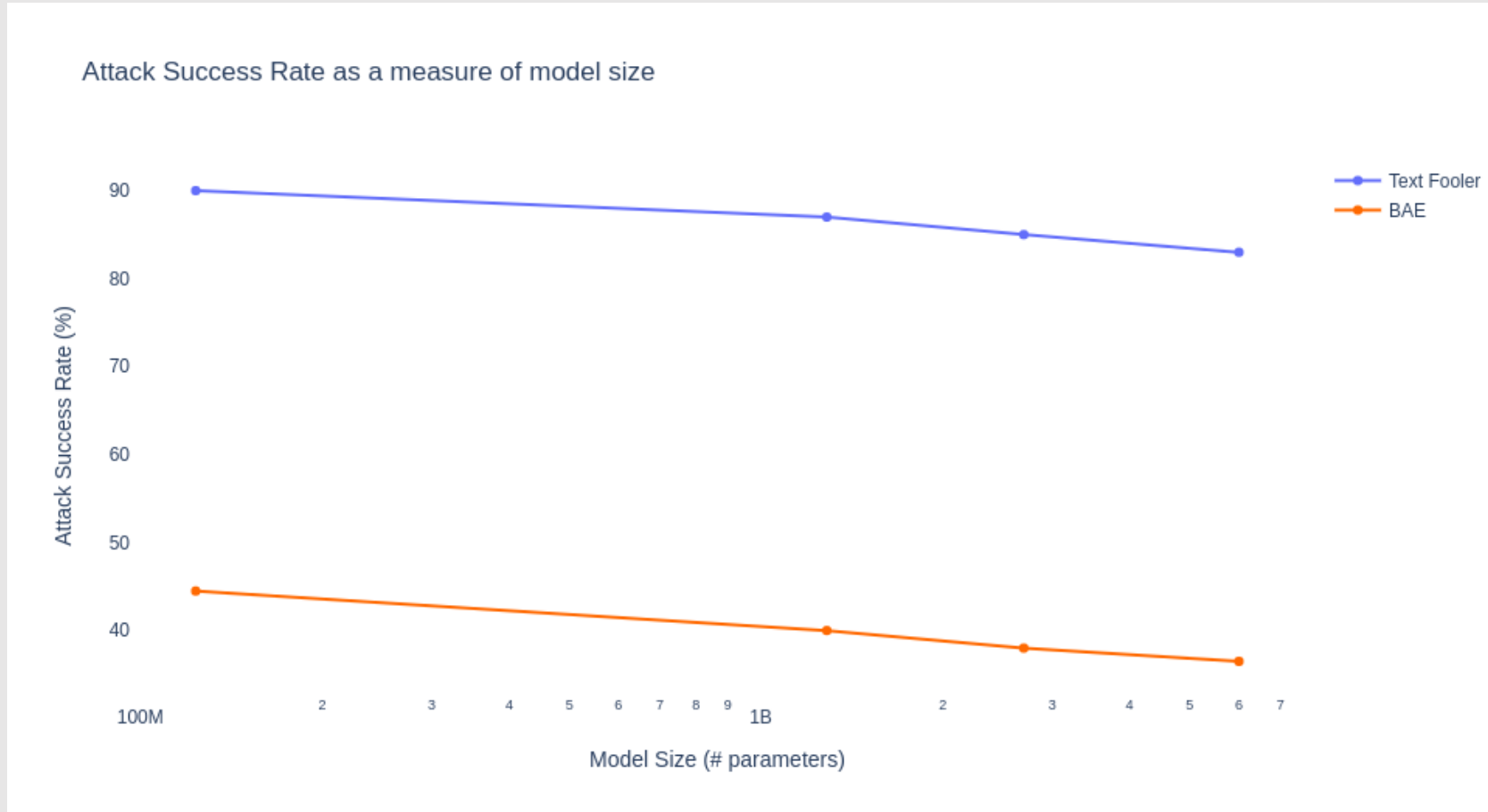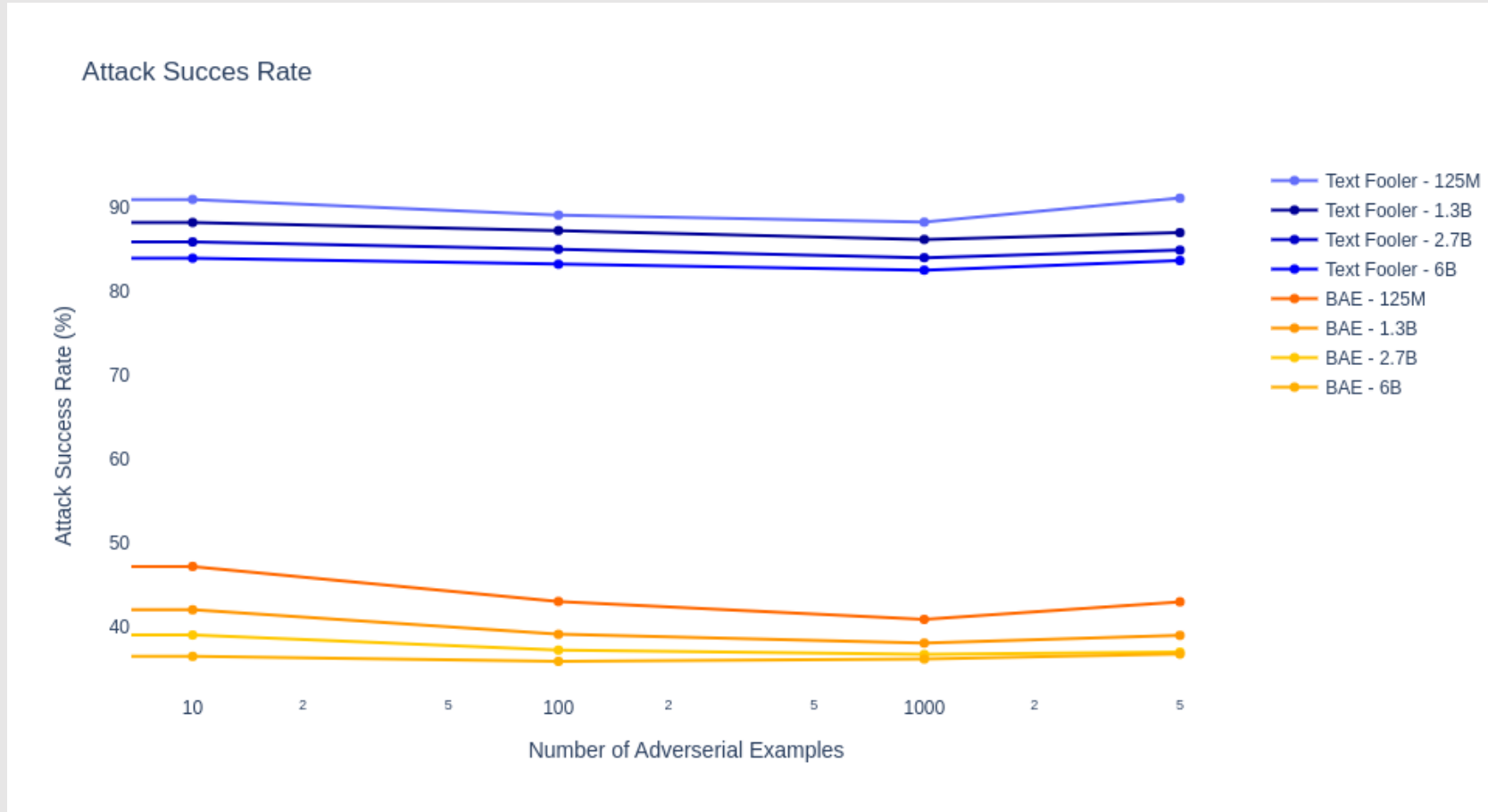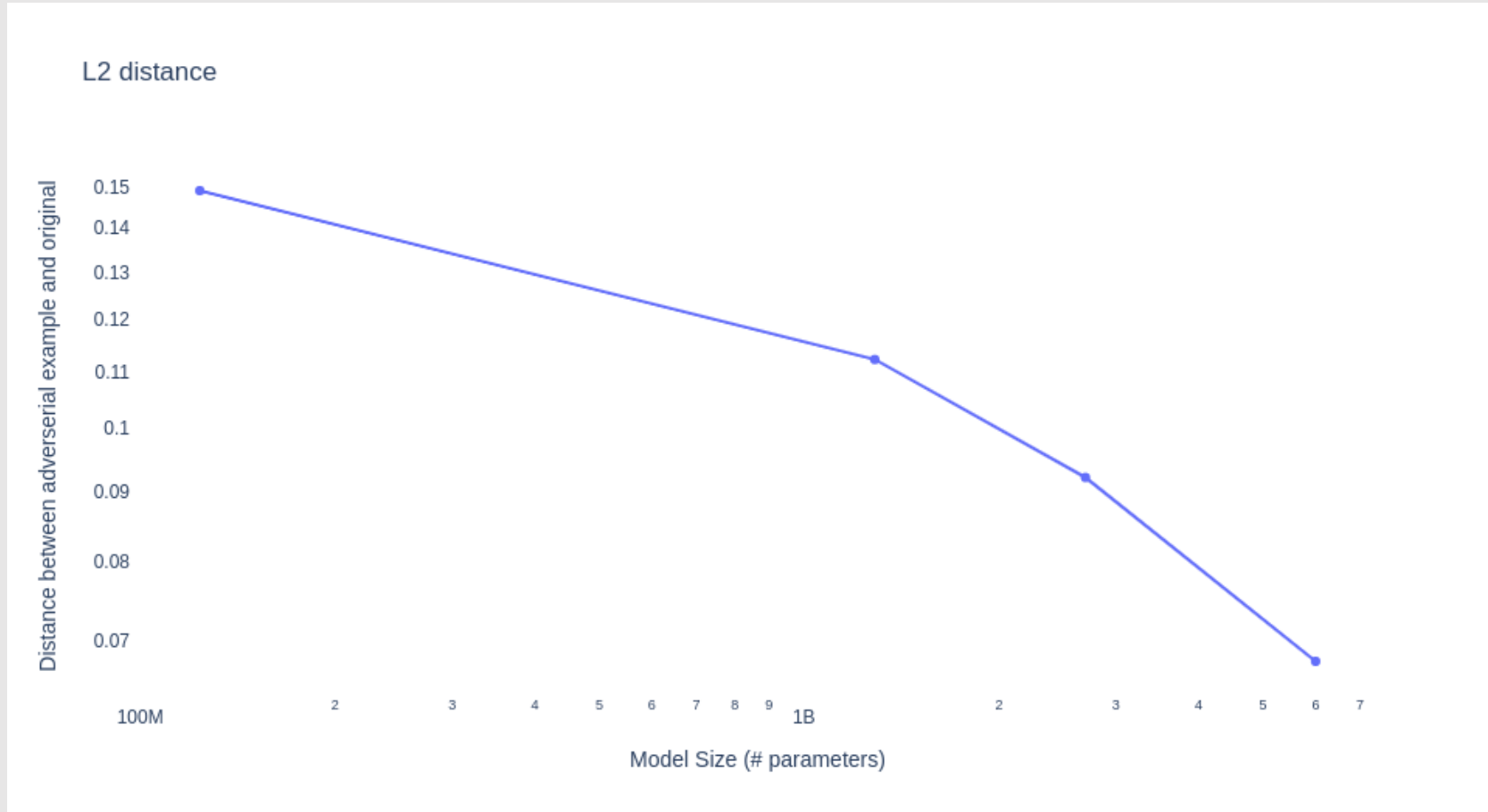
# Results

# Results



Adverserial Loss (Rotten Tomatoes)

# Results



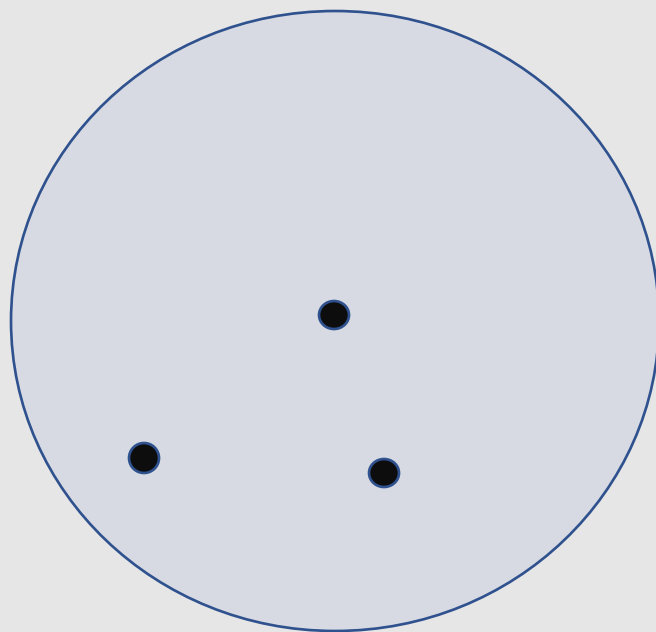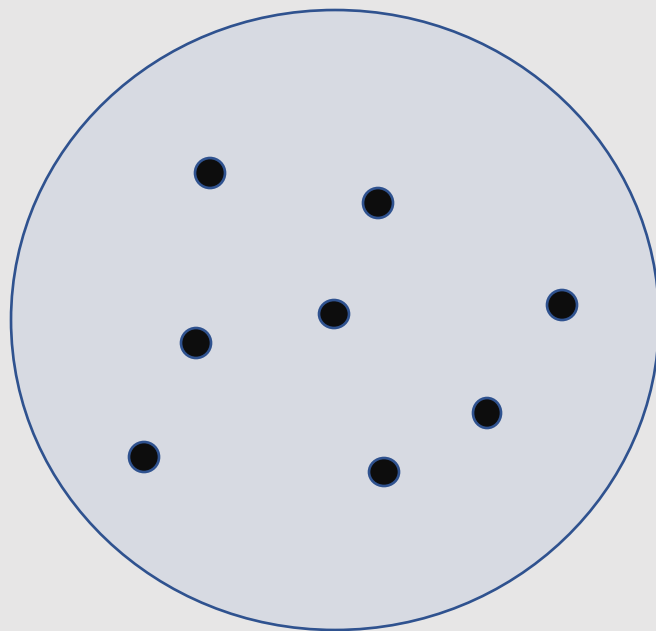Attack Success Rate as a measure of model size

# Results

# Results
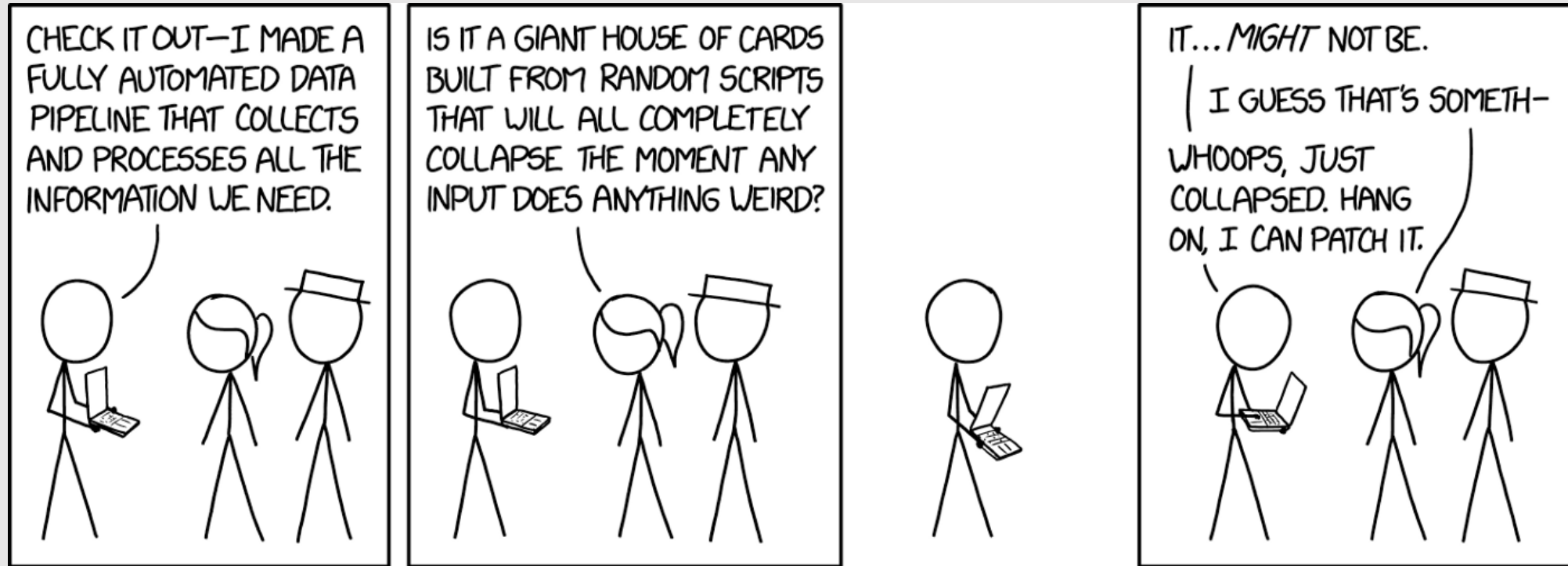
# Density of ball

# Density of ball

# Limitations

- Positive (99%) --> Negative (82%)

- yeah , these flicks are just that damn **good** . isn't it great ?

- yeah , these flicks are just that damn **bad** . isn't it great ?

# Conclusion

- Scale improves adversarial robustness
- Preliminary results
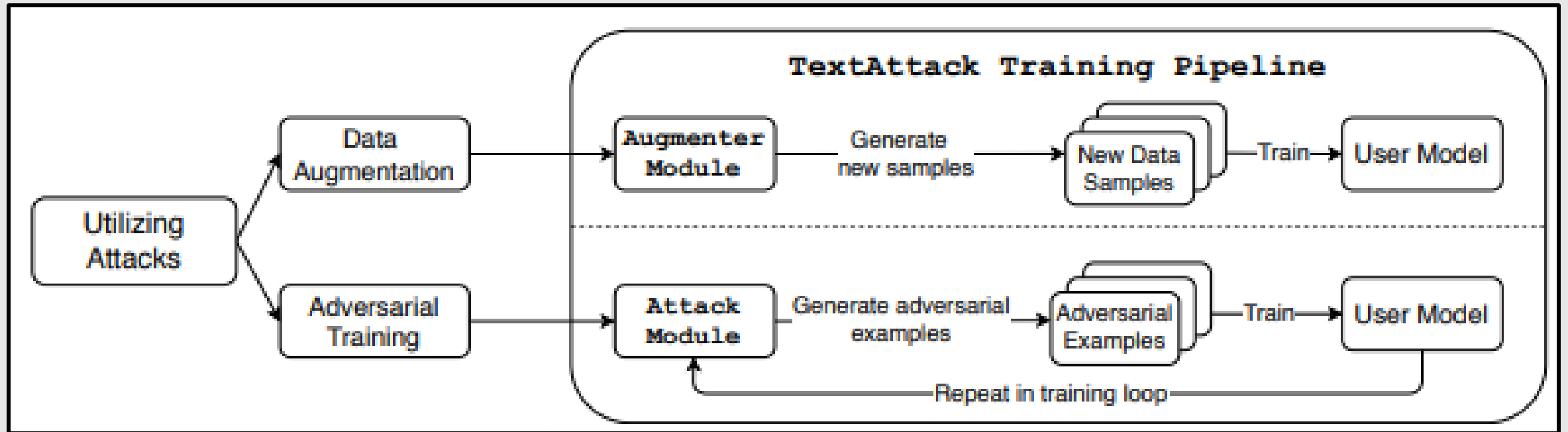- More datasets, more models!

# Thank You !
# Questions ?

# TextAttack

A simple framework for adversarial attacks :
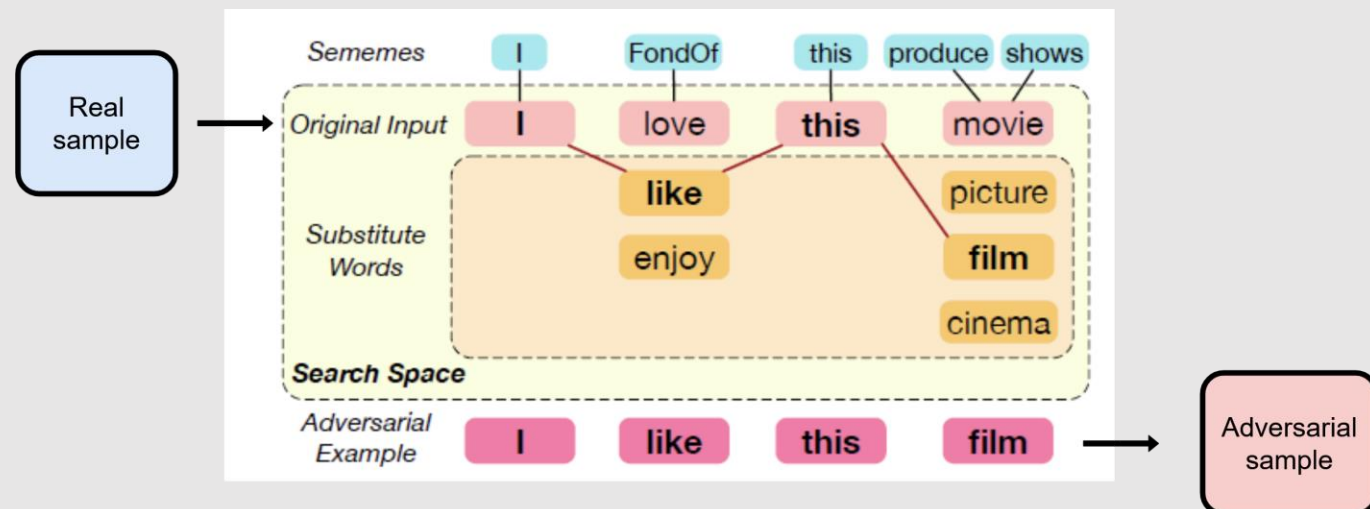
# Adversarial generator, an example :

PSOZang :Word-level Textual Adversarial Attacking as Combinatorial Optimization :

Idea :  Word-level attack seen as a Combinatorial Optimization problem

• Word substitution method based on sememes ( unit of semantic meaning ) : define the search space.

• Produces *likely* output (without context-awareness)

•Particle swarm optimization-based search algorithm
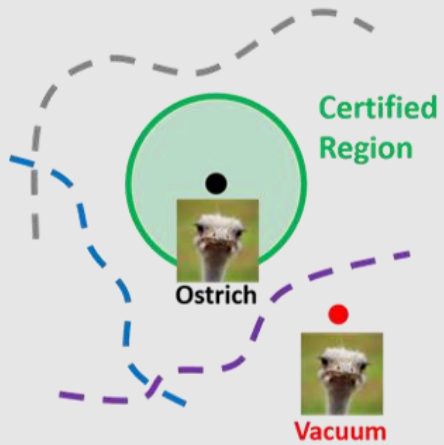
# Adversarial generator, an example :

BAEGarg: BERT-based Adversarial Examples for Text Classification :

Idea : use BERT model to predict MASKED token

- Rule-based synonym replacement strategy

- Produces *likely* output (with context-awareness)

- Produces output with improved grammaticality and semantic coherence

# Robustness and adversarial attacks :

# Adversarial Robustness