

# Announcements

- Final exam
  - Monday 12/8, 10:30am to 12:20pm
    - Come 5+ minutes before, we will start the exam at 10:30am sharp
  - Two locations this time
    - G20 (usual classroom) for 446 students
    - G01 for 546 students
  - Bring your Husky ID
  - You may bring a cheatsheet (same rules, 1 piece of paper, both sides)
  - All material tested except for guest lecture
- Course evaluations are open now
  - We will award 1% extra credit to all students who complete it

# Transfer Learning & Foundation Models

CSE 446/546

Sewoong Oh & Pang Wei Koh

# Good old fashioned artificial intelligence (GOF AI)

- AI systems tried to specify how to solve a task
- e.g., Logic Theorist (1956), General Problem Solver (1957)
- 1963: Herbert Simon speculated that GOF AI could produce human-level AI

# Eliza (Weizenbaum, 1966)

**young woman:** I need some help, that much seems certain.

**eliza:** What would it mean to you if you got some help?

**young woman:** Perhaps I could learn to get along with my mother.

**eliza:** Tell me more about your family.

**young woman:** My mother takes care of me.

**eliza:** Who else in your family takes care of you?

**young woman:** My father.

**eliza:** Your father.

**young woman:** You are like my father in some ways.

**eliza:** What resemblance do you see?

**young woman:** You are not very aggressive but I think you don't want me to notice that.

**eliza:** What makes you think I am not very aggressive?

**young woman:** You don't argue with me.

Weizenbaum: Some subjects have been very hard to convince that ELIZA...is not human.

# Machine learning (1990s+)

- Started getting popular within AI in the 1990s
- Instead of specifying how to solve a task, instead learn from data
- Homogenization: Many applications can be tackled by a single generic learning algorithm such as logistic regression
- Standard workflow: 1) Feature engineering → 2) Apply learning

# Feature engineering (circa 90s / 2000s)

- Language: N-grams / bag of words
- Vision: Edge detectors, SIFT, HOG
- Audio: Mel-frequency cepstrum

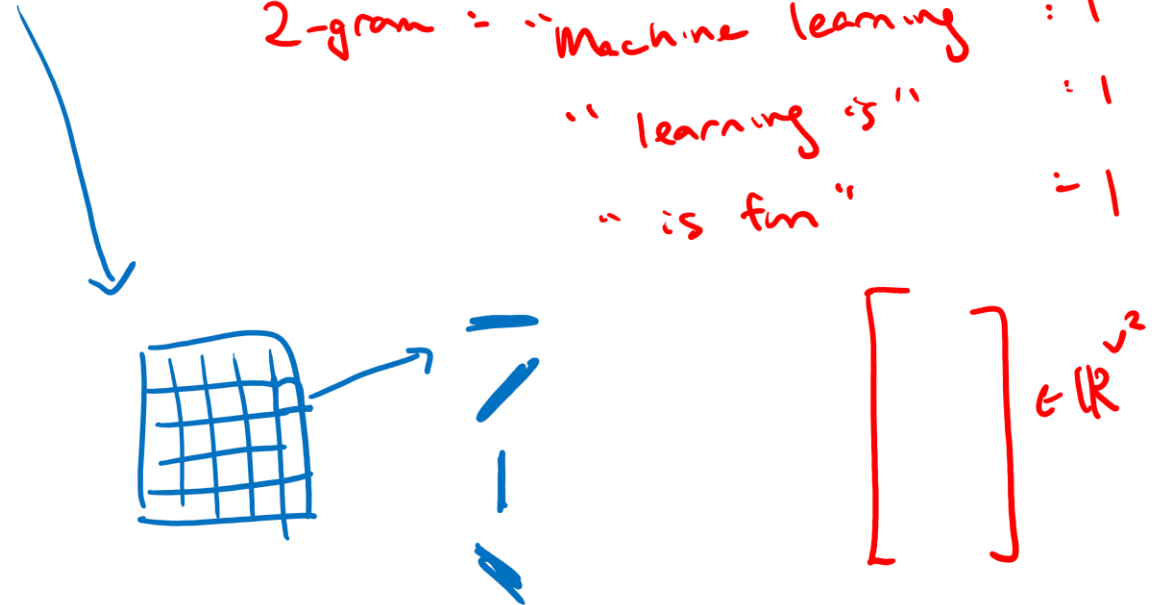
"Machine learning is fun"

1-gram: Machine : 1  
learning : 1  
is : 1  
fun : 1

$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \end{bmatrix} \in \mathbb{R}^4$

fun  
is

2-gram = "Machine learning" : 1  
"learning is" : 1  
"is fun" : 1



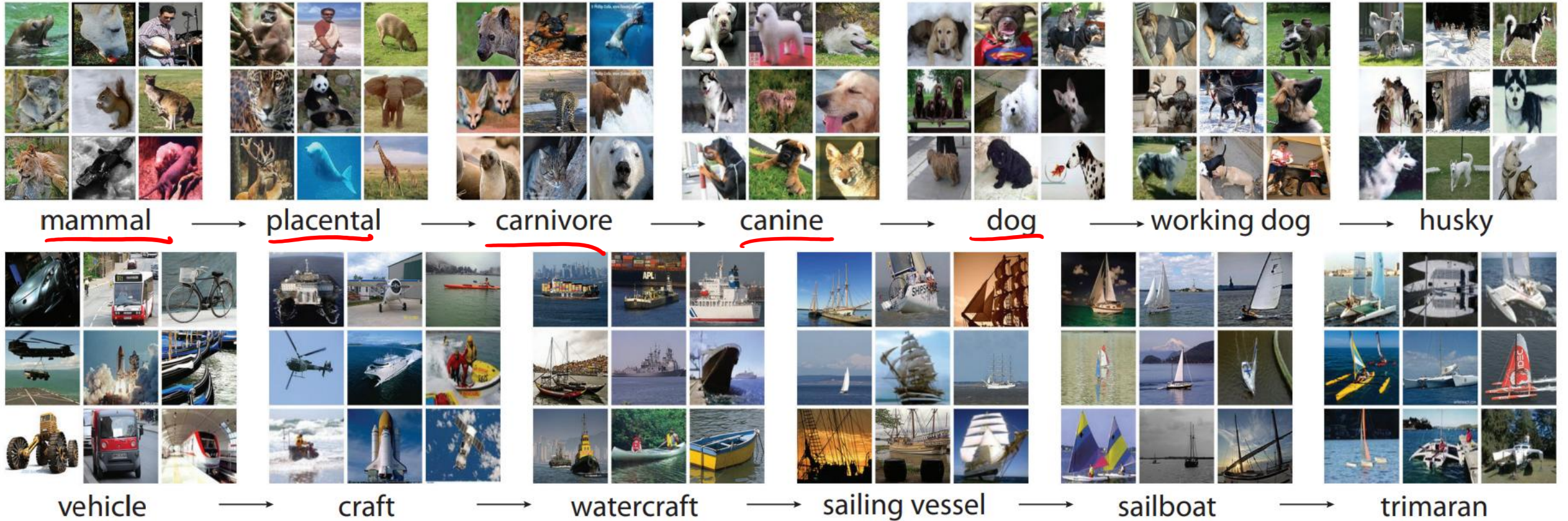
# Deep learning (2010+)

- Powered by larger datasets, more compute (GPUs)
- Main lesson: Training on raw inputs > feature engineering, especially for high-dimensional data
- Homogenization: Instead of bespoke feature engineering pipelines, we can use similar neural net architectures

# Transfer learning

- So far, we've assumed we only have data from a specific task
- Transfer learning:
  1. Pretrain (with lots of data) on a proxy task
  2. Then finetune/adapt (with less data) on the real task

# Vision: ImageNet (2009) + AlexNet (2012)

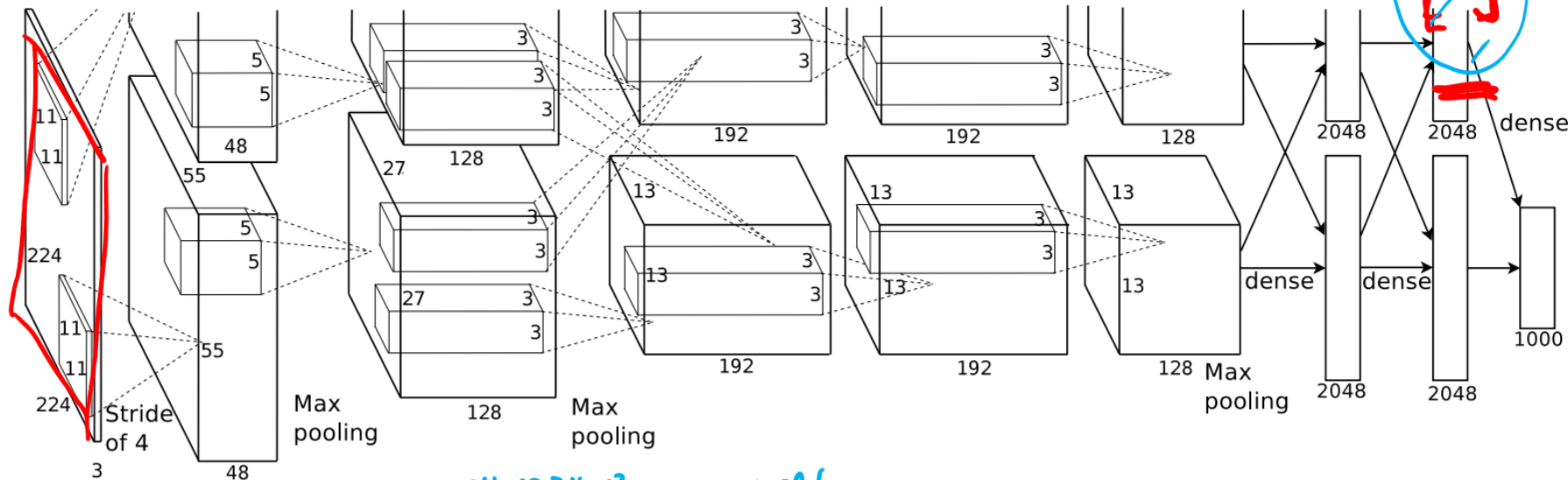


# Vision: ImageNet (2009) + AlexNet (2012)

freeze

low-dim

image



$$\phi(x) : \mathbb{R}^{224 \times 224 \times 3} \rightarrow \mathbb{R}^{4096}$$

$$x \rightarrow \phi(x) \text{ train logistic regression}$$

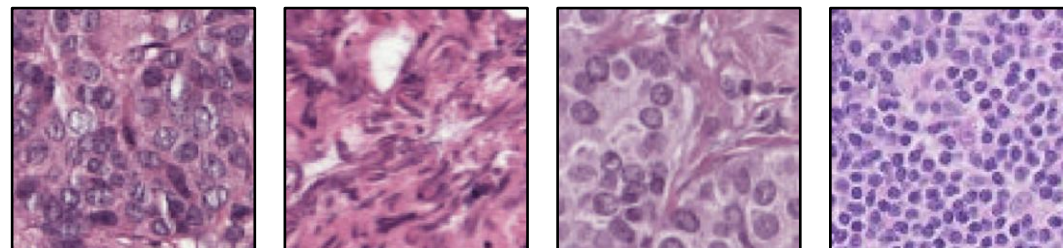
$\in \mathbb{R}^{4096}$

class prediction

$\in \mathbb{R}^{1000}$

tumor vs normal

Tumor Tumor Tumor Normal



Training data  $z_1, z_2, \dots, z_n$

# NLP: ELMo (2018), BERT (2018)

Pre-training → on  
Pre-training on → large  
Pre-training on large → text

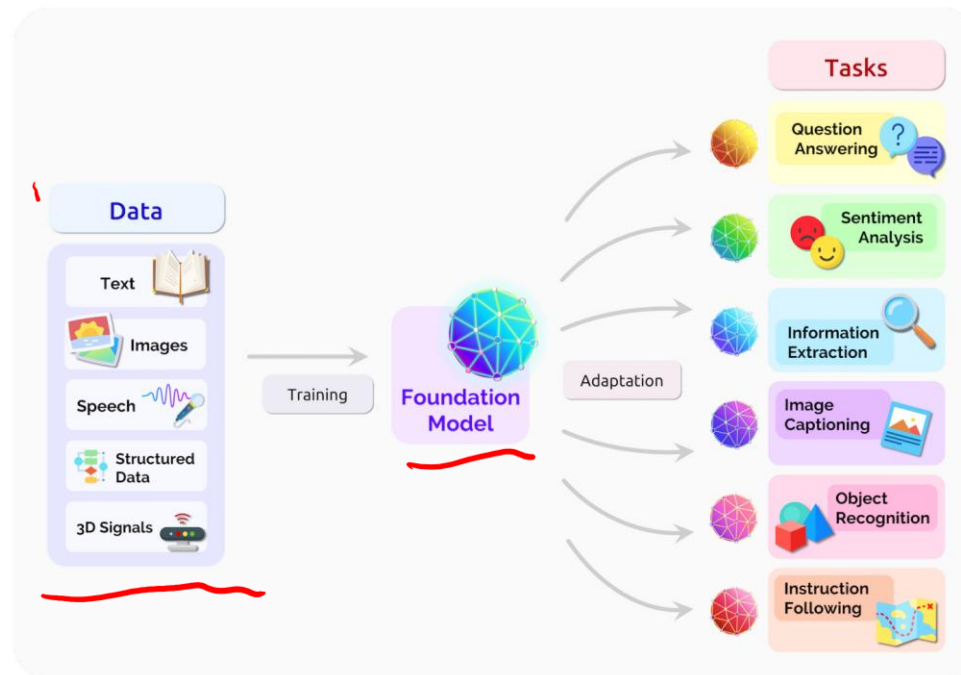
- "Pretraining on large text corpora"
- Pretraining: Next token prediction, masked language modeling
- Outputs a feature embedding / representation
- Adaptation: Train a supervised model for a particular task (e.g., sentiment analysis)
- Key difference: unsupervised (or self-supervised) pretraining

"Pretraining on [mask] text [mask]"  
→ large corpora

$x \xrightarrow{\phi} \phi(x)$

# Foundation models (2021+)

- Key: Self-supervised pretraining + (compute/model/data) scale
- Homogenization: Almost all SOTA models in NLP and vision are adapted (or used as-is) from one of a few foundation models



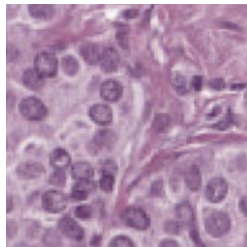
# Contrastive learning for multimodal models



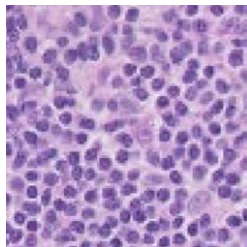
A white fluffy dog



Gates Center for Computer Science & Engineering



Cancerous lymph node biopsy from breast cancer patient



Histopathology image of normal lymph node tissue

# Contrastive learning for multimodal models



A white fluffy dog

Gates Center for Computer Science & Engineering

Cancerous lymph node biopsy from breast cancer patient

Histopathology image of normal lymph node tissue

# Contrastive learning for multimodal models



A white fluffy dog

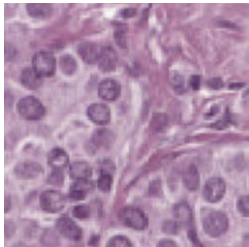
Gates Center for Computer Science & Engineering

Cancerous lymph node biopsy from breast cancer patient

Histopathology image of normal lymph node tissue

# Contrastive learning for multimodal models

A white fluffy dog

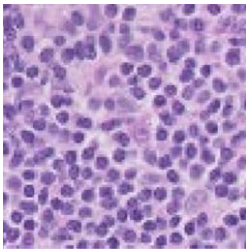
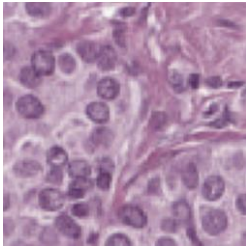


Gates Center for Computer Science & Engineering

Cancerous lymph node biopsy from breast cancer patient

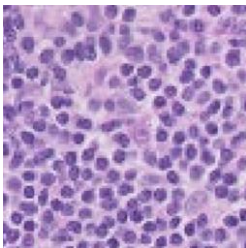
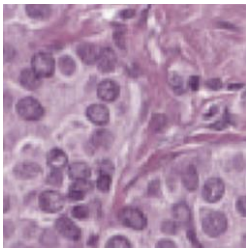
Histopathology image of normal lymph node tissue

# Contrastive learning for multimodal models



A white fluffy dog

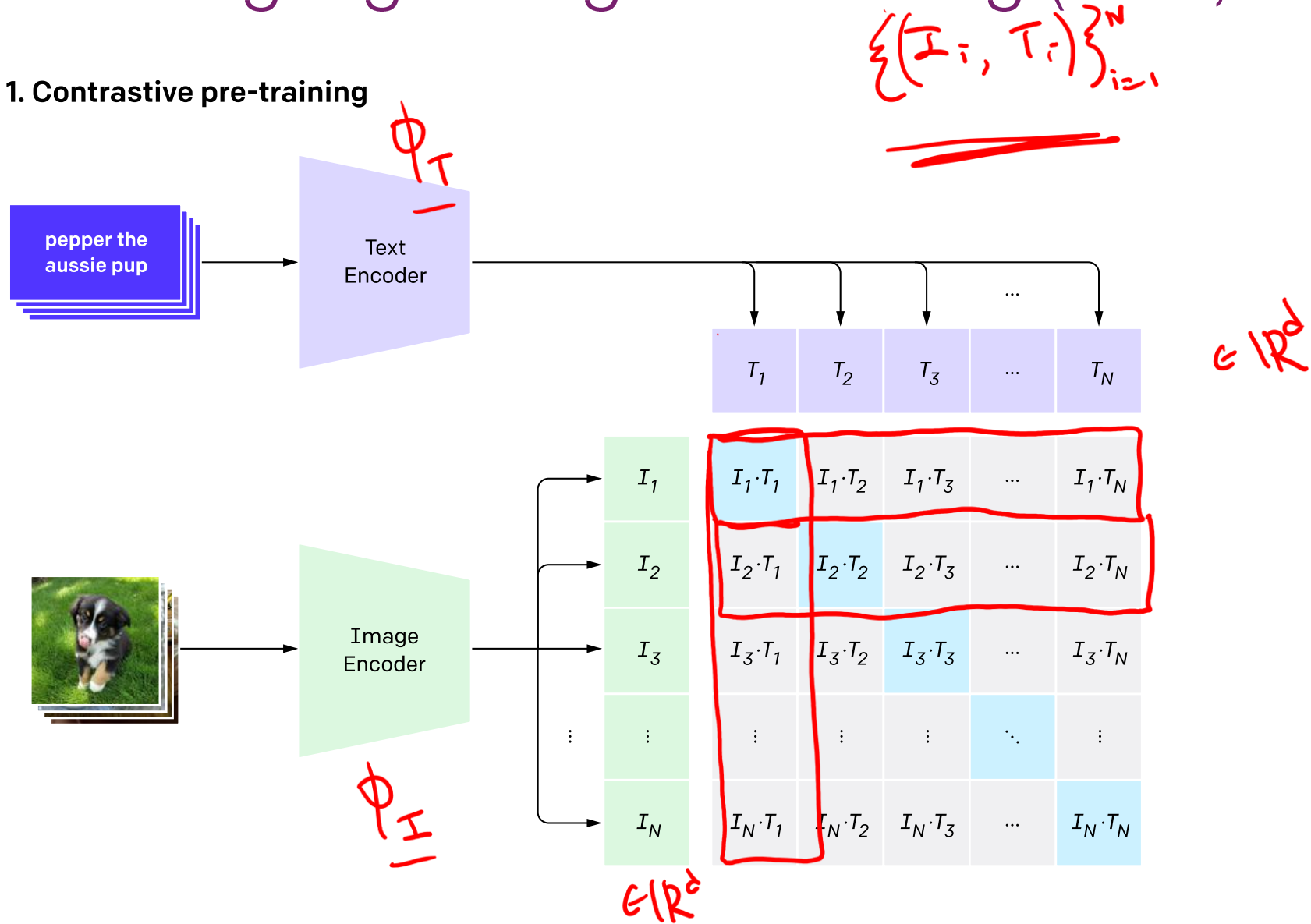
# Contrastive learning for multimodal models



Gates Center for Computer Science & Engineering

# Contrastive Language Image Pretraining (CLIP, 2021)

## 1. Contrastive pre-training



```
# image_encoder - ResNet or Vision Transformer
# text_encoder - CBOW or Text Transformer
# I[n, h, w, c] - minibatch of aligned images
# T[n, l] - minibatch of aligned texts
# W_i[d_i, d_e] - learned proj of image to embed
# W_t[d_t, d_e] - learned proj of text to embed
# t - learned temperature parameter

# extract feature representations of each modality
I_f = image_encoder(I) #[n, d_i]
T_f = text_encoder(T) #[n, d_t]

# joint multimodal embedding [n, d_e]
I_e = l2_normalize(np.dot(I_f, W_i), axis=1)
T_e = l2_normalize(np.dot(T_f, W_t), axis=1)

# scaled pairwise cosine similarities [n, n]
logits = np.dot(I_e, T_e.T) * np.exp(t)

# symmetric loss function
labels = np.arange(n)
loss_i = cross_entropy_loss(logits, labels, axis=0)
loss_t = cross_entropy_loss(logits, labels, axis=1)
loss = (loss_i + loss_t)/2
```

### Food101

**guacamole** (90.1%) Ranked 1 out of 101 labels



- ✓ a photo of **guacamole**, a type of food.
- × a photo of **ceviche**, a type of food.
- × a photo of **edamame**, a type of food.
- × a photo of **tuna tartare**, a type of food.
- × a photo of **hummus**, a type of food.

### Youtube-BB

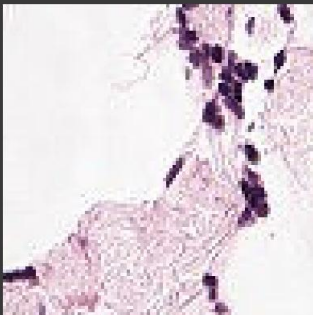
**airplane, person** (89.0%) Ranked 1 out of 23 labels



- ✓ a photo of a **airplane**.
- × a photo of a **bird**.
- × a photo of a **bear**.
- × a photo of a **giraffe**.
- × a photo of a **car**.

### PatchCamelyon (PCam)

**healthy lymph node tissue** (77.2%) Ranked 2 out of 2 labels



- × this is a photo of **lymph node tumor tissue**
- ✓ this is a photo of **healthy lymph node tissue**

### SUN397

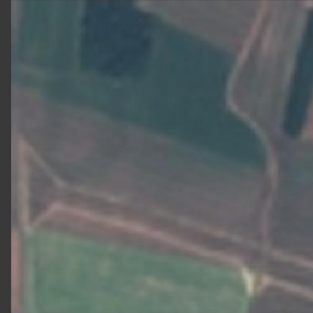
**television studio** (90.2%) Ranked 1 out of 397 labels



- ✓ a photo of a **television studio**.
- × a photo of a **podium indoor**.
- × a photo of a **conference room**.
- × a photo of a **lecture room**.
- × a photo of a **control room**.

### EuroSAT

**annual crop land** (46.5%) Ranked 4 out of 10 labels



- × a centered satellite photo of **permanent crop land**.
- × a centered satellite photo of **pasture land**.
- × a centered satellite photo of **highway or road**.
- ✓ a centered satellite photo of **annual crop land**.
- × a centered satellite photo of **brushland or shrubland**.

### ImageNet-A (Adversarial)

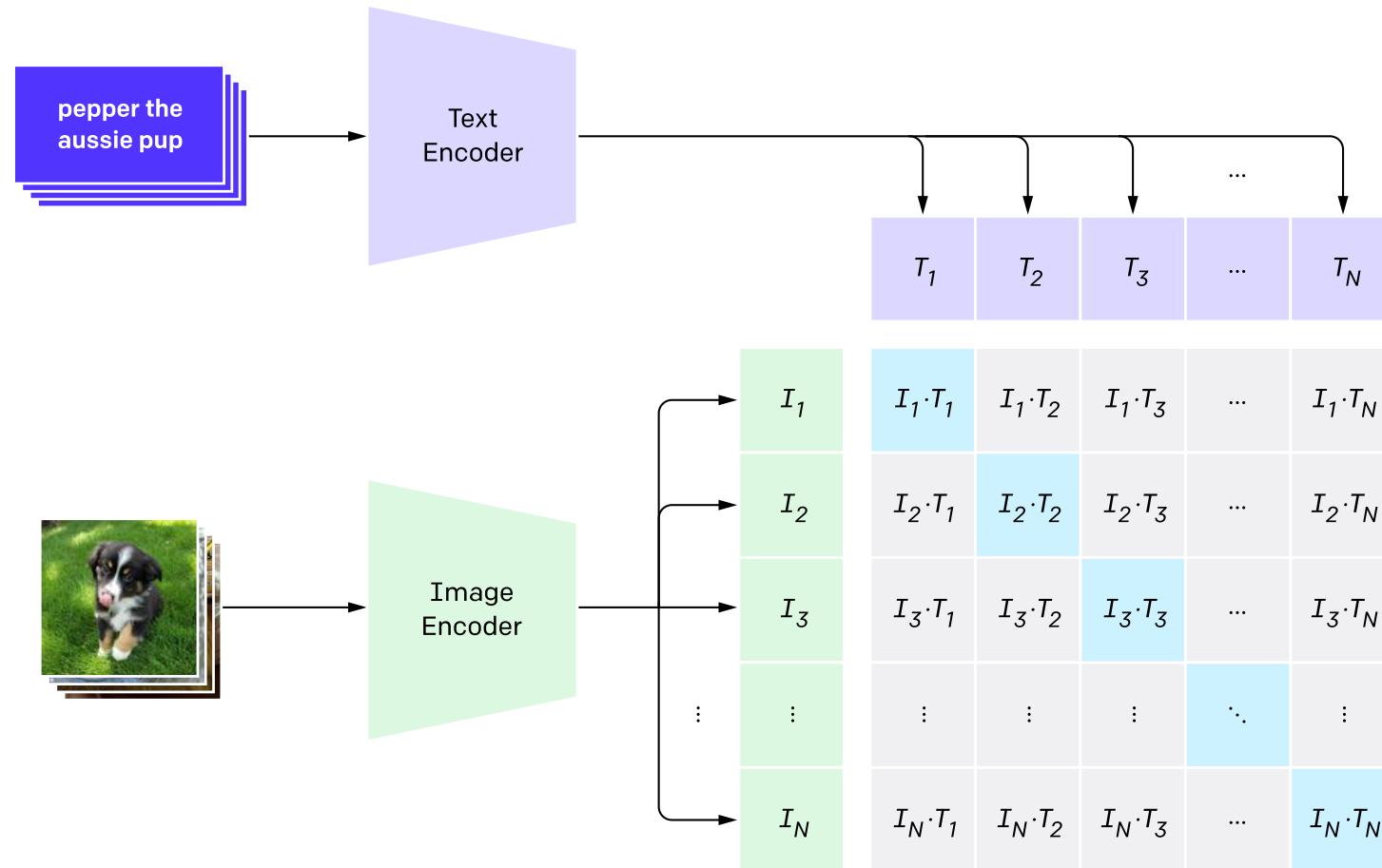
**lynx** (47.9%) Ranked 5 out of 200 labels



- × a photo of a **fox squirrel**.
- × a photo of a **mongoose**.
- × a photo of a **skunk**.
- × a photo of a **red fox**.
- ✓ a photo of a **lynx**.

# Contrastive Language Image Pretraining (CLIP, 2021)

## 1. Contrastive pre-training



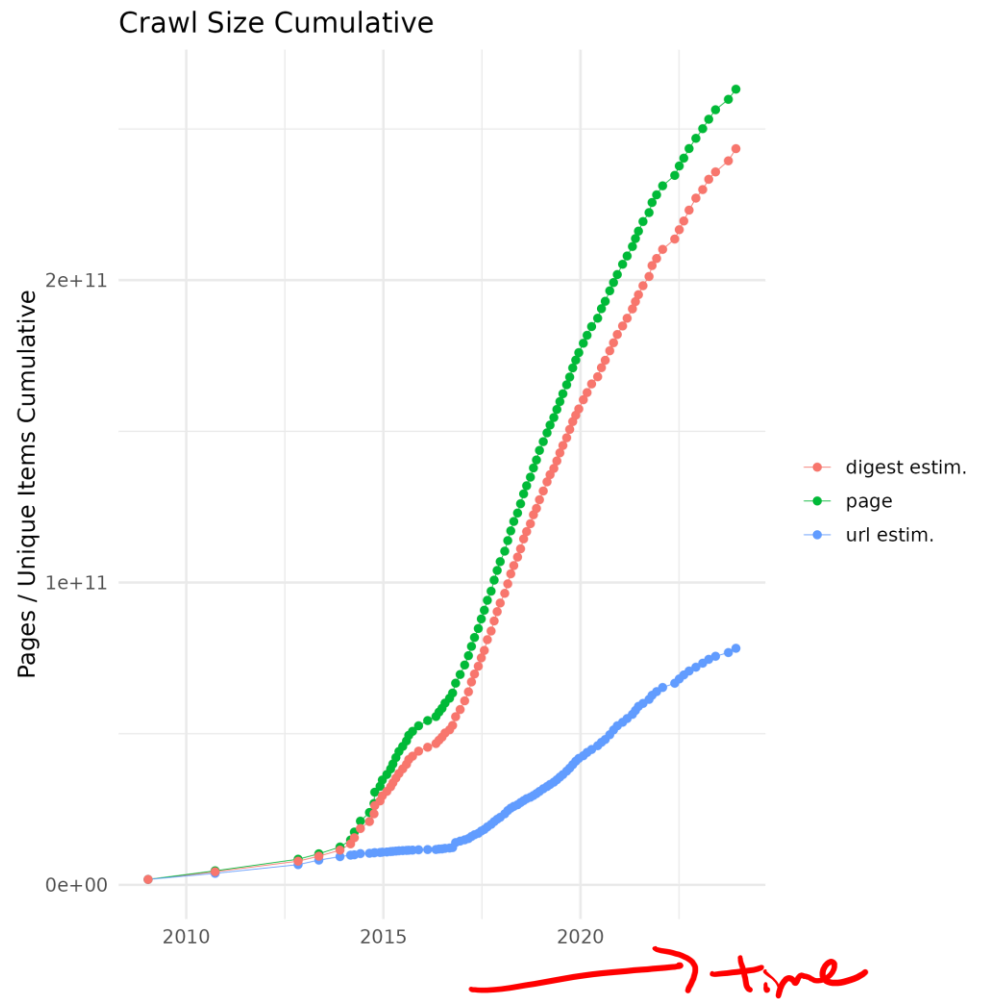
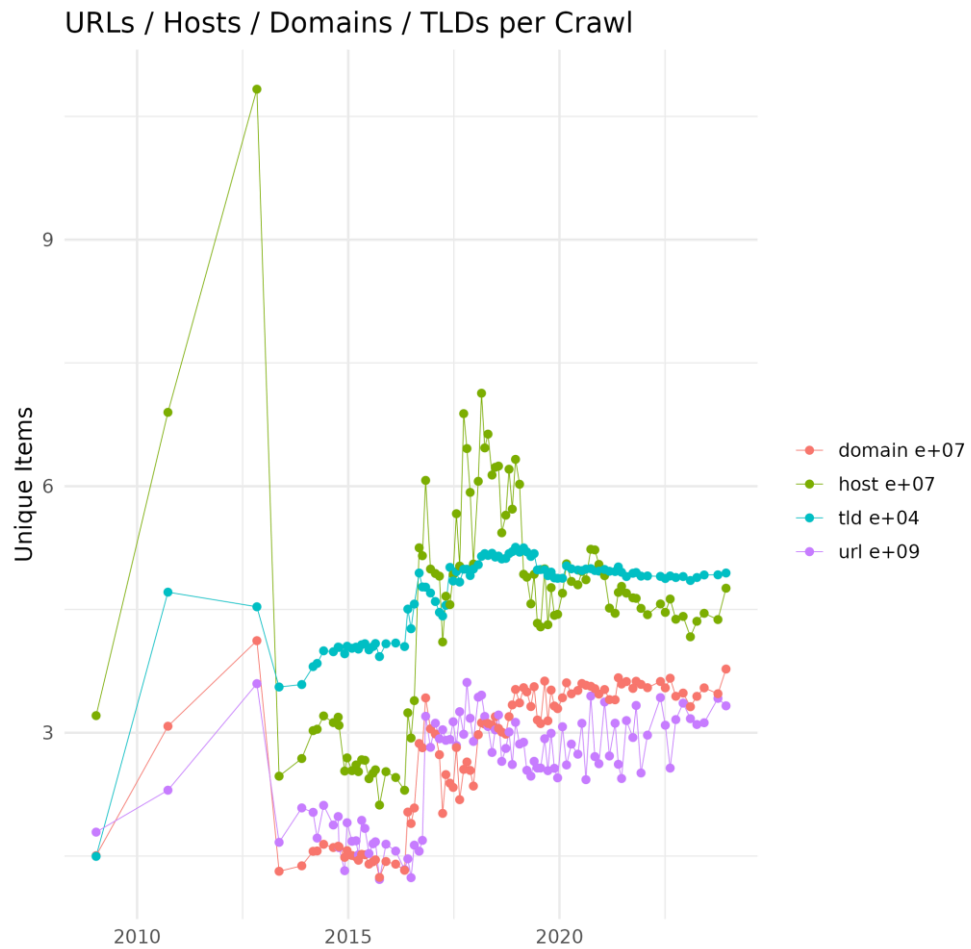
# A brief history of web text datasets

- Transformers (Vaswani et al., 2017)
  - WMT 2014 English-French dataset: 36M sentences
- BERT (Devlin et al., 2018)
  - BooksCorpus (800M words) + English Wiki (2,500M words, ~6M pages)
- GPT-2 (Radford et al., 2019)
  - WebText (8M docs): all outbound links from Reddit with 3+ karma
- C4 / T5 (Raffel et al., 2020)
  - Common Crawl web scrape

Data set	Size	GLUE	CNN4M	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ C4	745GB	83.28	<b>19.24</b>	80.88	71.36	<b>26.98</b>	<b>39.82</b>	<b>27.65</b>
C4, unfiltered	<u>6.1TB</u>	81.46	19.14	78.78	68.04	26.55	39.34	27.21
RealNews-like	35GB	<b>83.83</b>	<b>19.23</b>	80.39	72.38	<b>26.75</b>	<b>39.90</b>	<b>27.48</b>
WebText-like	17GB	<b>84.03</b>	<b>19.31</b>	<b>81.42</b>	71.40	<b>26.80</b>	<b>39.74</b>	<b>27.59</b>
Wikipedia	<u>16GB</u>	81.85	<b>19.31</b>	81.29	68.01	<b>26.94</b>	39.69	<b>27.67</b>
Wikipedia + TBC	20GB	83.65	<b>19.28</b>	<b>82.08</b>	<b>73.24</b>	<b>26.77</b>	39.63	<b>27.57</b>

Table 8: Performance resulting from pre-training on different data sets. The first four variants are based on our new C4 data set.

# The scale of the Common Crawl



# C4 (Colossal Clean Crawled Corpus)

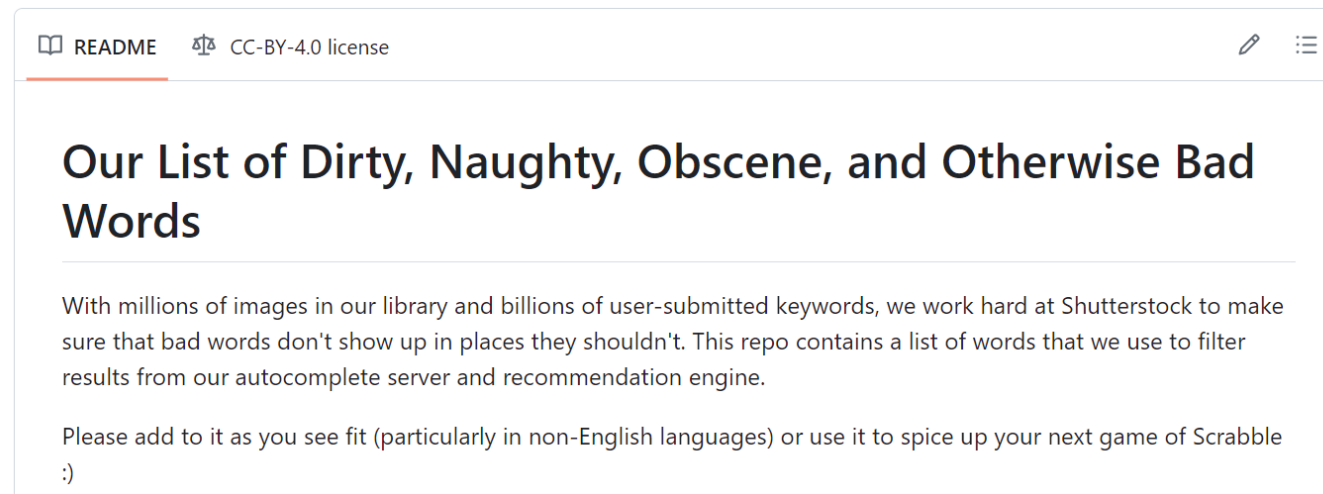
language, placeholder text, source code, etc.). To address these issues, we used the following heuristics for cleaning up Common Crawl's web extracted text:

- We only retained lines that ended in a terminal punctuation mark (i.e. a period, exclamation mark, question mark, or end quotation mark).
- We discarded any page with fewer than 5 sentences and only retained lines that contained at least 3 words.
- We removed any page that contained any word on the "List of Dirty, Naughty, Obscene or Otherwise Bad Words".<sup>6</sup>
- Many of the scraped pages contained warnings stating that Javascript should be enabled so we removed any line with the word Javascript.
- Some pages had placeholder "lorem ipsum" text; we removed any page where the phrase "lorem ipsum" appeared.
- Some pages inadvertently contained code. Since the curly bracket "{" appears in many programming languages (such as Javascript, widely used on the web) but not in natural text, we removed any pages that contained a curly bracket.
- To deduplicate the data set, we discarded all but one of any three-sentence span occurring more than once in the data set.

Additionally, since most of our downstream tasks are focused on English-language text, we used `langdetect`<sup>7</sup> to filter out any pages that were not classified as English with a probability of at least 0.99. Our heuristics are inspired by past work on using Common

# Removing “offensive” content

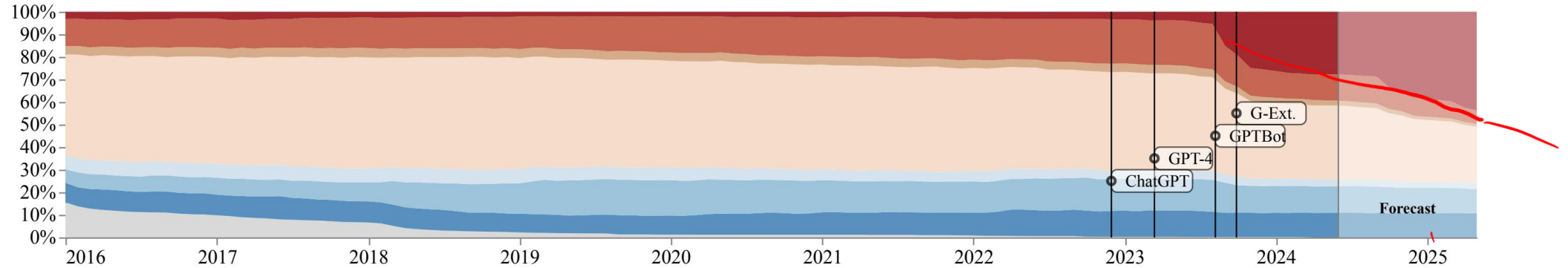
- In practice, list is mostly related to sexual/lewd content, not toxicity
- Majority of excluded docs relate to science, medicine, legal, etc.
- Disproportionate effect on mentions of sexual orientation



# Where is the data coming from?

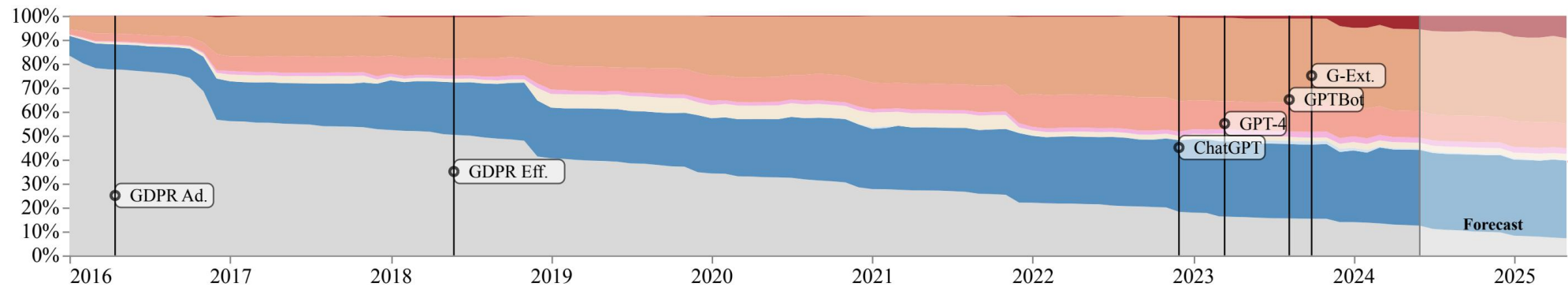
Website	Tokens (Billion)
<a href="https://en.wikipedia.org">en.wikipedia.org</a>	13.52
<a href="https://www.fanfiction.net">www.fanfiction.net</a>	11.68
<a href="https://www.nifty.org">www.nifty.org</a>	8.03
<a href="https://www.theguardian.com">www.theguardian.com</a>	5.44
<a href="https://edition.cnn.com">edition.cnn.com</a>	4.77
<a href="https://www.literotica.com">www.literotica.com</a>	4.54
<a href="https://stackoverflow.com">stackoverflow.com</a>	4.26
<a href="https://www.mumsnet.com">www.mumsnet.com</a>	3.95
<a href="https://tvtropes.org">tvtropes.org</a>	3.78
<a href="https://archiveofourown.org">archiveofourown.org</a>	3.77
<a href="https://doctrinepublishing.com">doctrinepublishing.com</a>	3.16
<a href="https://www.zerohedge.com">www.zerohedge.com</a>	3.15
<a href="https://www.scribd.com">www.scribd.com</a>	3.09
<a href="https://transcripts.cnn.com">transcripts.cnn.com</a>	3.06
<a href="https://www.theatlantic.com">www.theatlantic.com</a>	3.05
<a href="https://slashdot.org">slashdot.org</a>	2.98
<a href="https://chowhound.chow.com">chowhound.chow.com</a>	2.84
<a href="https://www.nbcnews.com">www.nbcnews.com</a>	2.70
<a href="https://www.reddit.com">www.reddit.com</a>	2.69
<a href="https://everything2.com">everything2.com</a>	2.69

# “Consent in crisis” (2024)



## Robots.txt Restrictions

- Full restrictions
- Pattern-based restrictions
- Disallow private directories
- Other restrictions
- Crawl delay specified
- Sitemap provided
- No restrictions or sitemap
- No Robots.txt



## ToS Restrictions

- No Crawling & AI
- No Crawling
- No AI
- Non-Commercial Use
- Non-Compete
- No Re-Distribution
- Conditional Use
- Unrestricted Use
- No Terms Pages

# *The Times Sues OpenAI and Microsoft Over A.I. Use of Copyrighted Work*

Millions of articles from The New York Times were used to train chatbots that now compete with it, the lawsuit said.

## WARNER MUSIC GROUP STRIKES 'LANDMARK' DEAL WITH SUNO; SETTLES COPYRIGHT LAWSUIT AGAINST AI MUSIC GENERATOR

2.3K  
SHARES



in



NOVEMBER 25, 2025

BY MURRAY STASSEN

## AI firm Anthropic agrees to pay authors \$1.5bn to settle piracy lawsuit

5 September 2025

Share Save

Model: Playground v2.5  
Prompt: "Mario"



Model: Playground v2.5  
Prompt: "Videogame, Plumber"



Model: DALL·E 3  
Prompt: "Videogame, Plumber"



Model: Playground v2.5  
Prompt: "Batman"



Model: Playground v2.5  
Prompt: "Gotham, Superhero"



Model: DALL·E 3  
Prompt: "Gotham, Superhero"

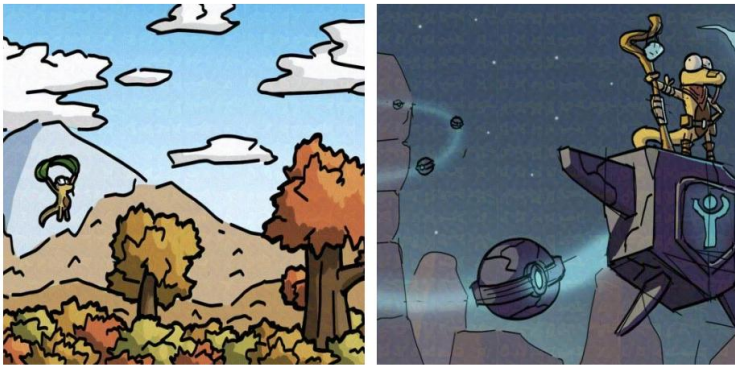


Original artwork by @nulevoy



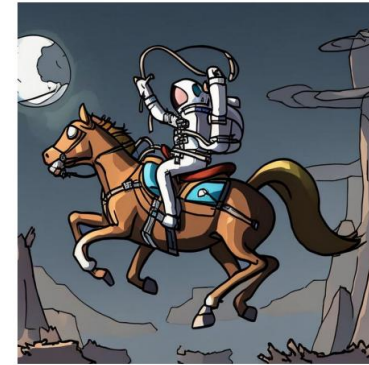
... so artists only release protected art

Protected artwork

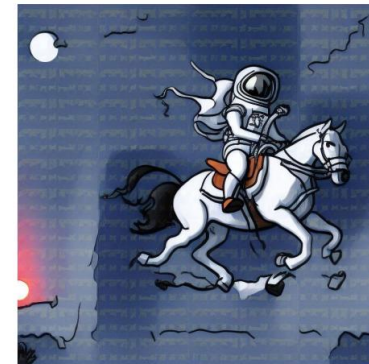


An astronaut riding a horse

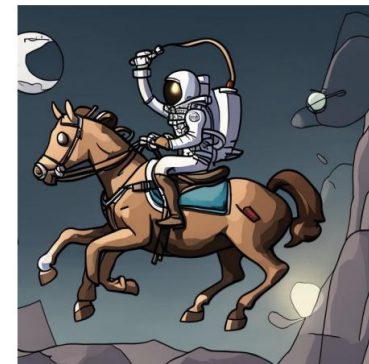
Public art can be used to mimic artists...



Protections prevent naive mimicry



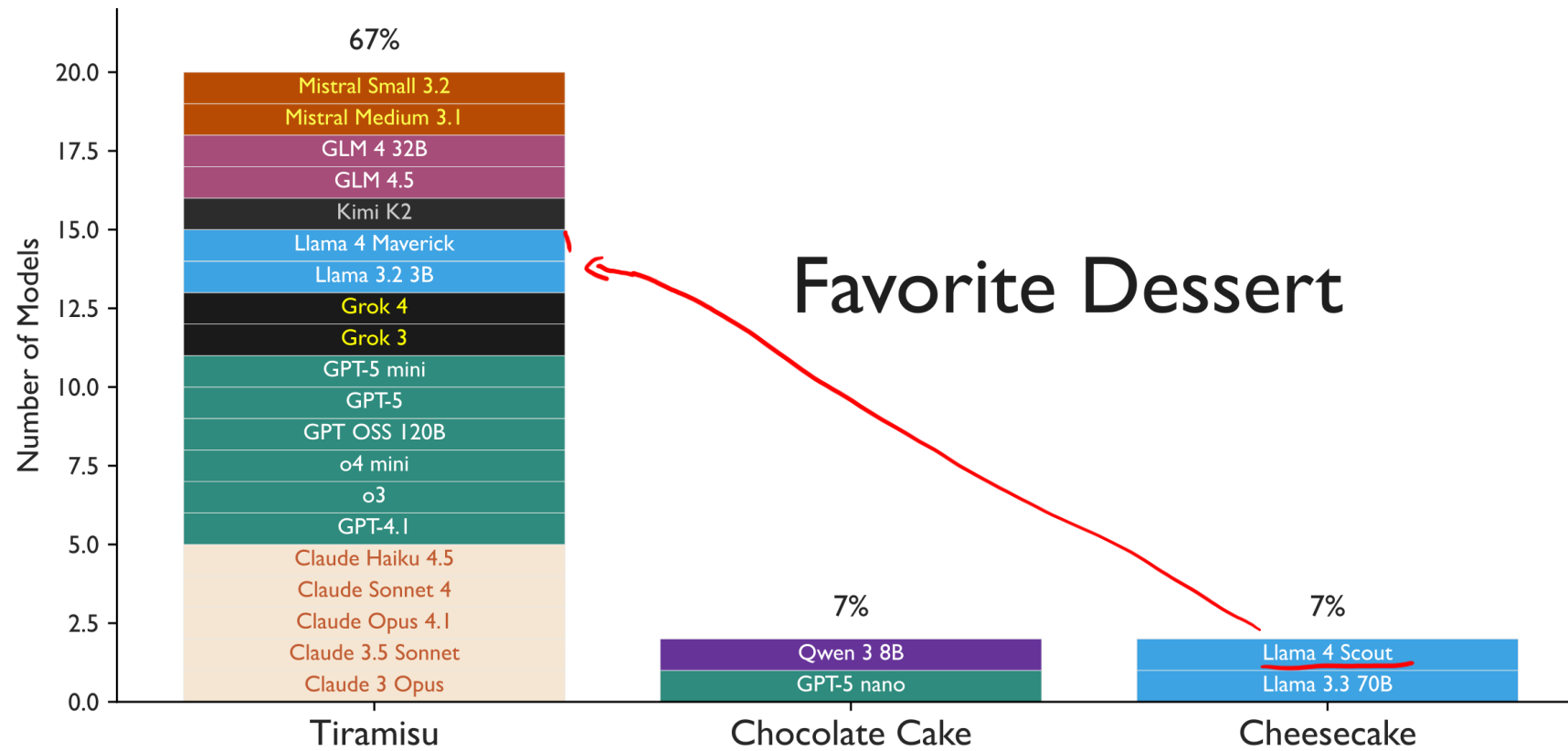
Protections fail against robust mimicry (ours)



# Algorithmic monoculture

The rise of algorithms used to shape societal choices has been accompanied by concerns over *monoculture*—the notion that choices and preferences will become homogeneous in the face of algorithmic curation. One of many canonical articulations of this concern was expressed in *The New York Times* by Farhad Manjoo (1), who wrote: “Despite the barrage of choice, more of us are enjoying more of the same songs, movies and TV shows.” Because of algorithmic curation, trained on collective social feedback (2), our choices are converging.

# Cultural homogenization



# Cultural homogenization

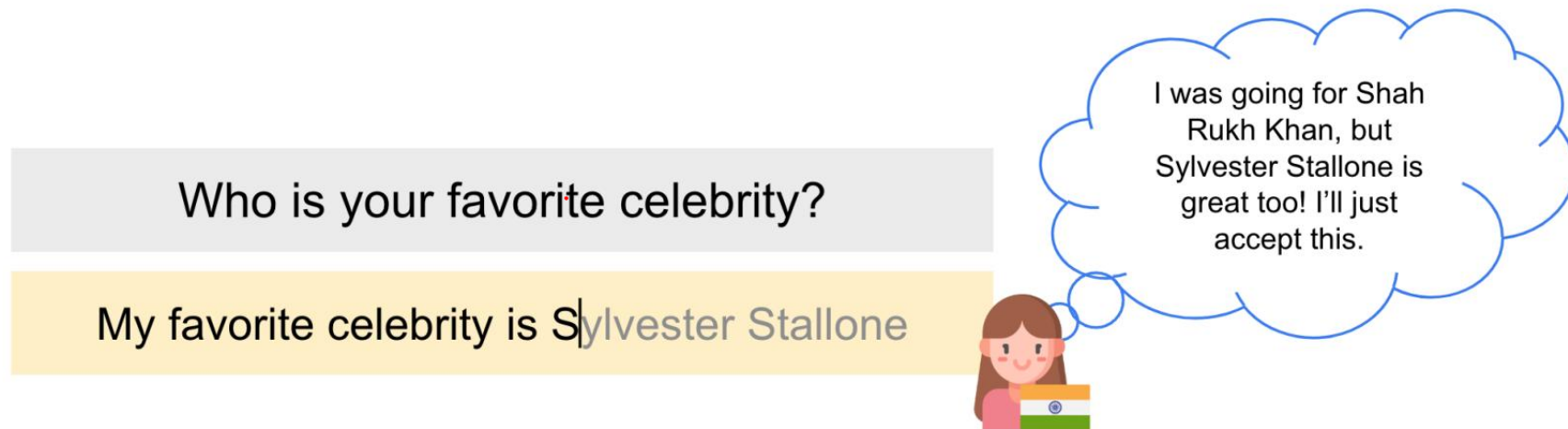


Fig. 1. A representation of the potential cultural homogenization from Western-centric AI models

# Takeaways

- The story of ML has been about increasing generality and homogenization: learning algorithms → architectures → models
- We are seeing unprecedented growth in the adoption and investment in AI.
- How will AI transform society?

# What's next?

- Machine learning
  - Graphical models
  - Interactive learning
  - Deep learning + deep learning theory
  - Advanced ML
  - Statistics, linear algebra, optimization, probability, analysis
- Domains
  - Reinforcement learning, robotics
  - Natural language processing
  - Computer vision
  - Computational biology, neuroscience
- Research
  - CSE599 seminars
  - CSE390R
  - Capstones
  - Join research labs: <https://www.cs.washington.edu/findingresearch>

# Announcements

- Final exam
  - Monday 12/8, 10:30am to 12:20pm
    - Come 5+ minutes before, we will start the exam at 10:30am sharp
  - Two locations this time
    - G20 (usual classroom) for 446 students
    - G01 for 546 students
  - Bring your Husky ID
  - You may bring a cheatsheet (same rules, 1 piece of paper, both sides)
  - All material tested except for guest lecture

# Thank you!

On behalf of Sewoong and the TAs, thanks for joining us this quarter!

Please give us both your feedback (by Dec 7):

Pang Wei: <https://uw.iasystem.org/survey/314304>

Sewoong: <https://uw.iasystem.org/survey/314303>

We will award 1% extra credit to all students who complete it.

446 students: Please give feedback to your section leader as well.