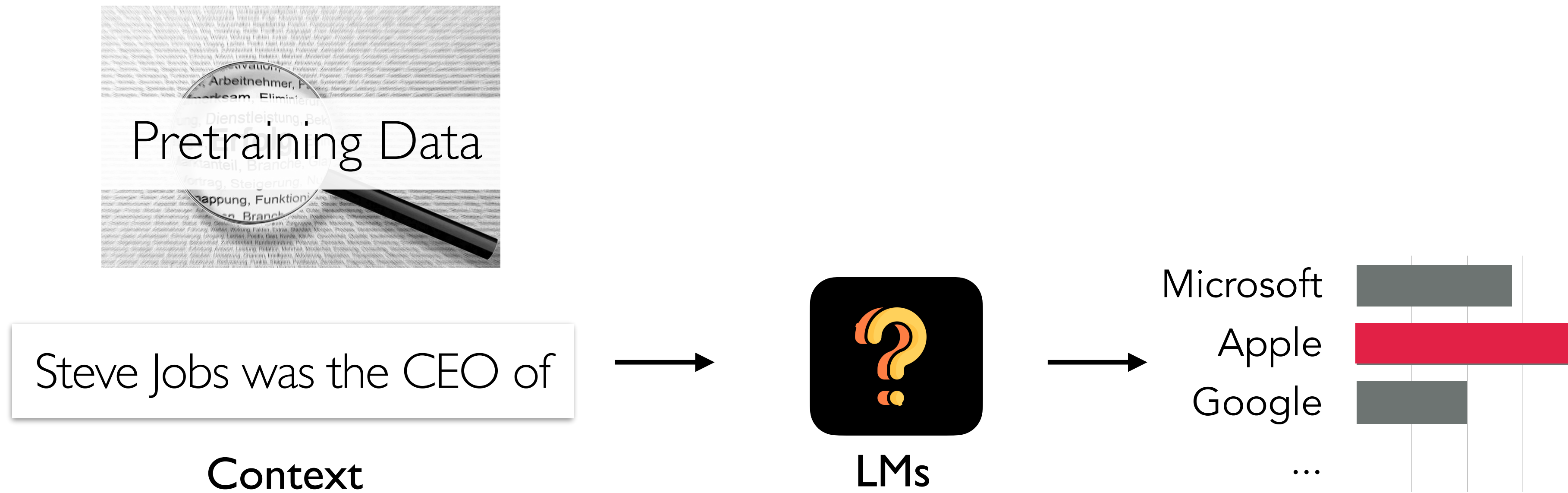# Beyond Monolithic Language Models

## Weijia Shi
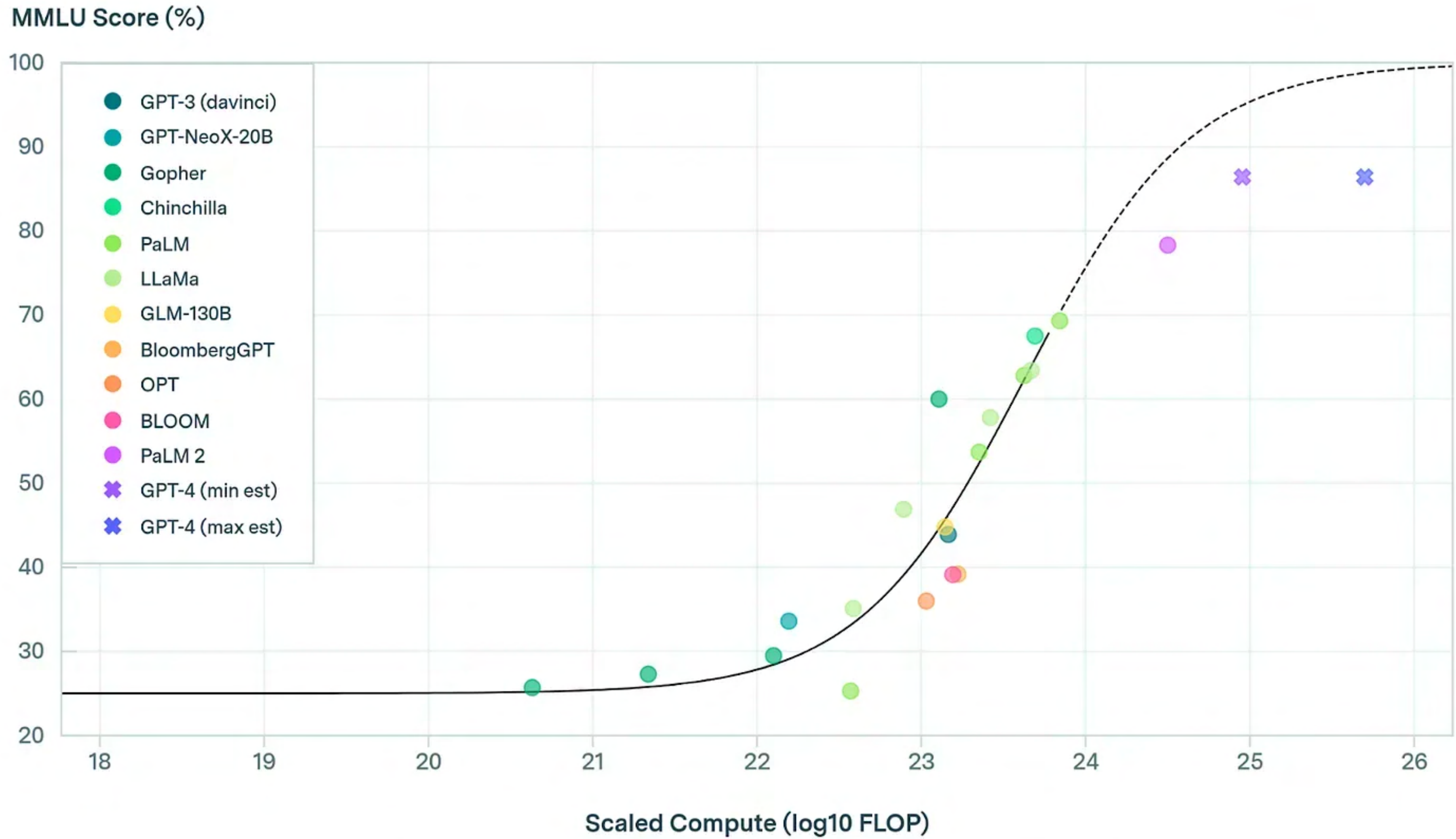
University of Washington

# Current Language Models (LMs): Monolithic

Pretraining Data

Steve Jobs was the CEO of

**Context**

**LMs**

Microsoft
Apple
Google
...

# Scaling

**?** *Scaling Up Monolithic LMs*

# Problems

Hallucinations

ARTIFICIAL INTELLIGENCE / TECH / POLICY

**Former Trump lawyer Michael Cohen accidentally cited fake court cases generated by AI** / Michael Cohen provided his lawyer with made-up court cases generated by Google Bard after mistaking the chatbot for a 'a super-charged search engine.'

# Problems

**Hallucinations**

**Hard to update new knowledge**

---

**ChatGPT** o1-mini >

> What is OpenAI o1?

◎ Thought for a few seconds ⌄

"OpenAI O1" does not correspond to any publicly recognized product, model, or initiative from OpenAI. It's possible that "O1" refers to one of the following:

1. **Typographical Error:** It might be a misspelling or shorthand for another OpenAI product, such as GPT-1 (the first generation of the Generative Pre-trained Transformer models) or a different version like GPT-3 or GPT-4.
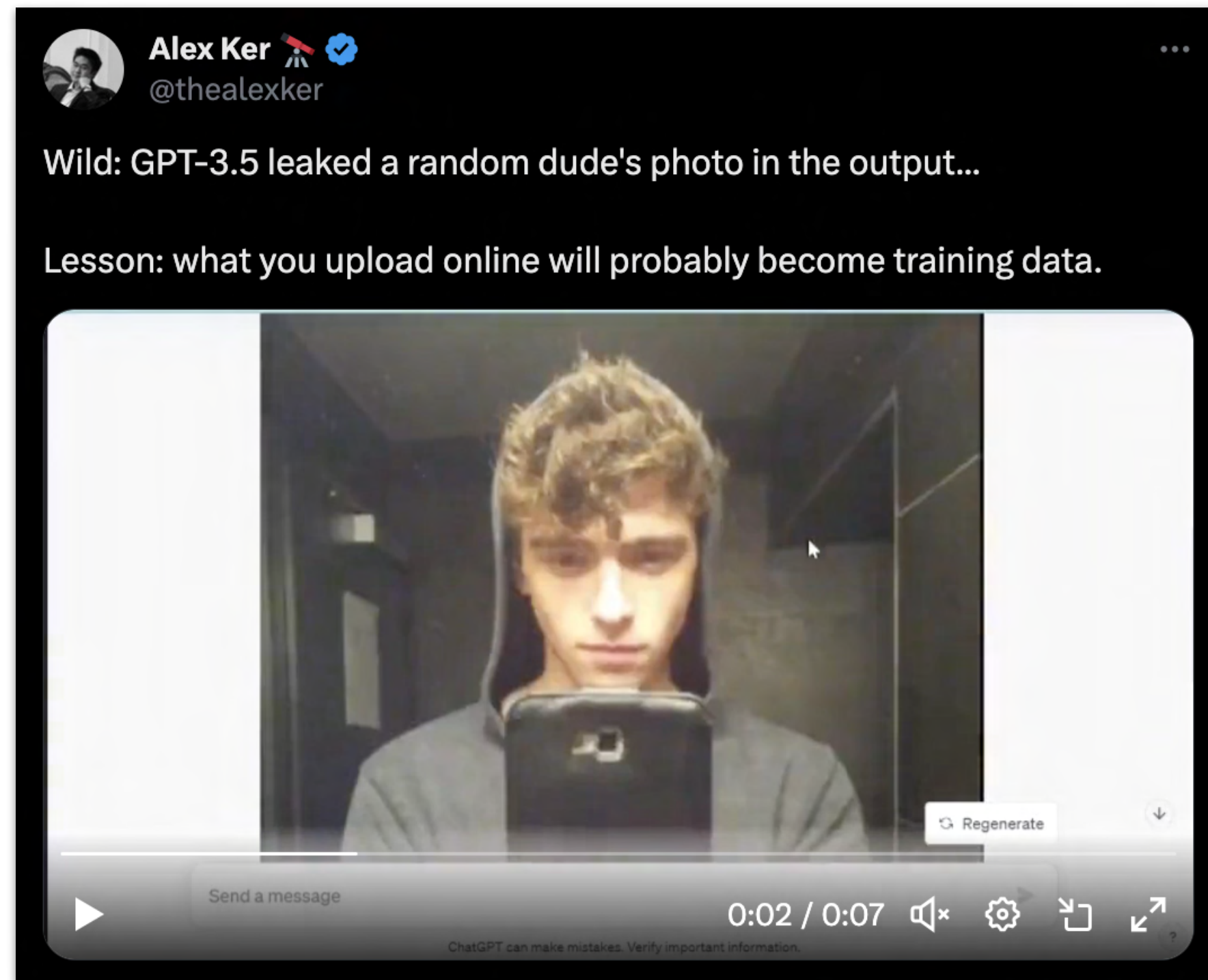
# Problems

Hard to update new knowledge

Copyright and privacy risks



Alex Ker
@thealexker

Wild: GPT-3.5 leaked a random dude's photo in the output...

Lesson: what you upload online will probably become training data.



**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.
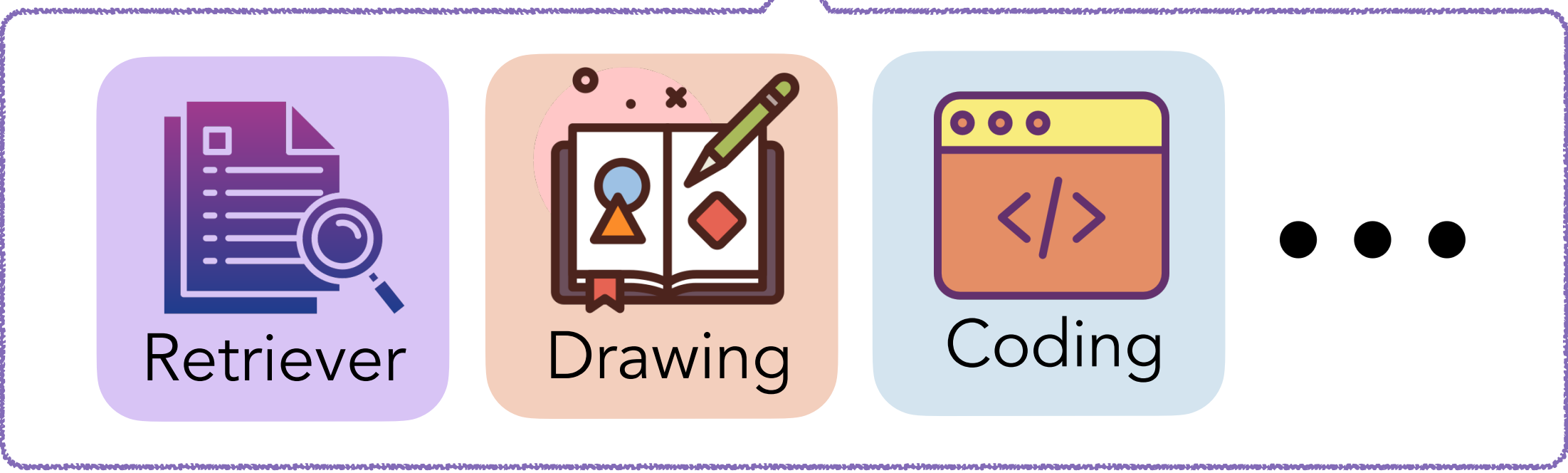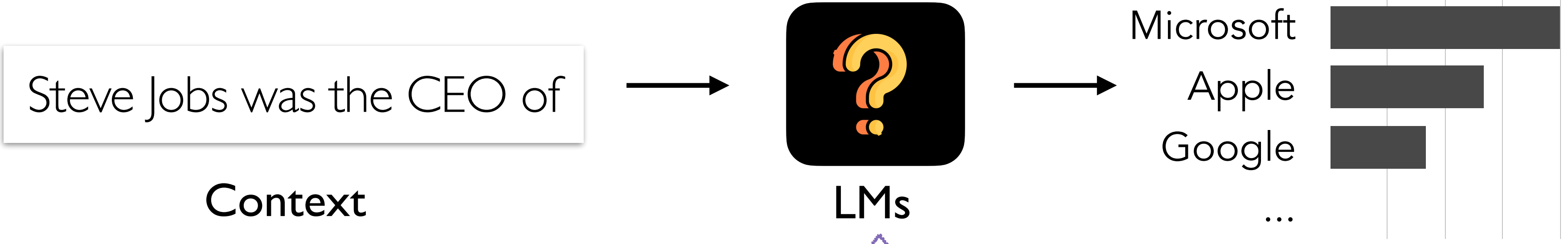
Dec. 27, 2023
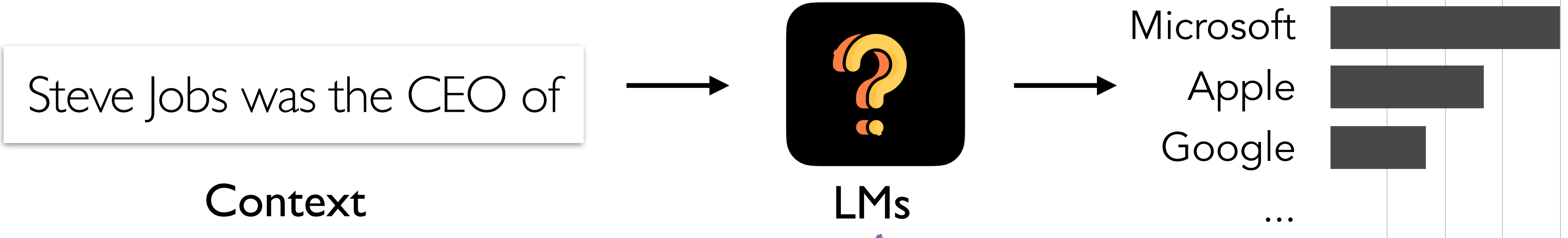
7

❌ *Scaling Up Monolithic LMs*

❌ *Scaling Up Monolithic LMs*

❓ *Alternative Paradigm*

# **Modularity**, *not Monoliths*

Steve Jobs was the CEO of

**Context**

LMs

Microsoft
Apple
Google
...

*Augmented Models*

Retriever   Drawing   Coding   • • •

Steve Jobs was the CEO of

**Context**

**LMs**

Microsoft
Apple
Google
...

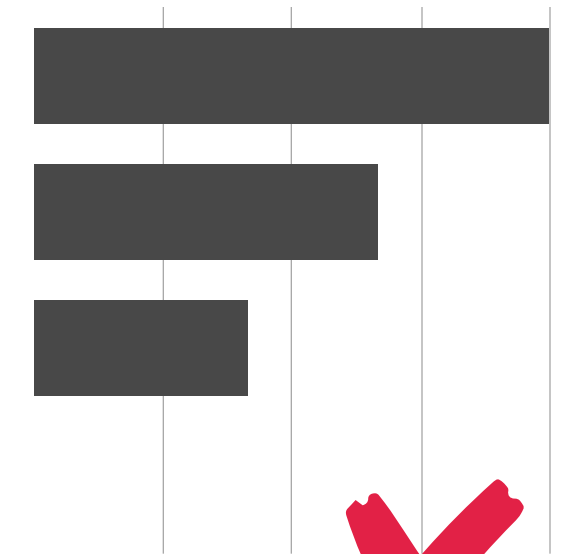Retriever   Drawing   Coding   • • •

*Augmented Models*

# Augmented Models

Context

Steve Jobs was the CEO of

LM

Microsoft

Apple

Google

...

# Augmented Models

Context

Steve Jobs was the CEO of

Datastore

Retriever

Jobs cofounded Apple in his parents' garage

LM

Microsoft
Apple
Google
...

✓ **Hallucinations**

# Augmented Models

Context

Steve Jobs was the CEO of

Datastore

+

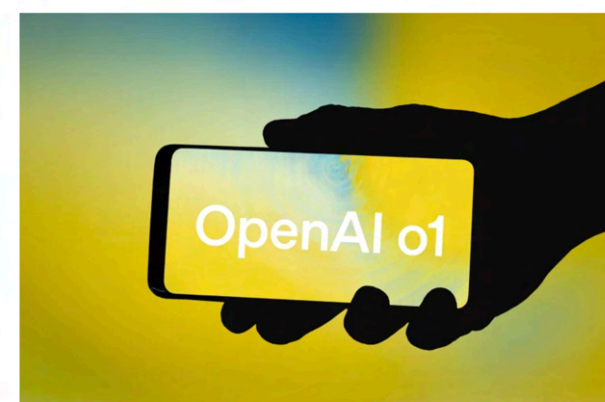**'In awe': scientists impressed by latest ChatGPT model o1**

The chatbot excels at science, beating PhDs on a hard science test. But it might 'hallucinate' more than its predecessors.
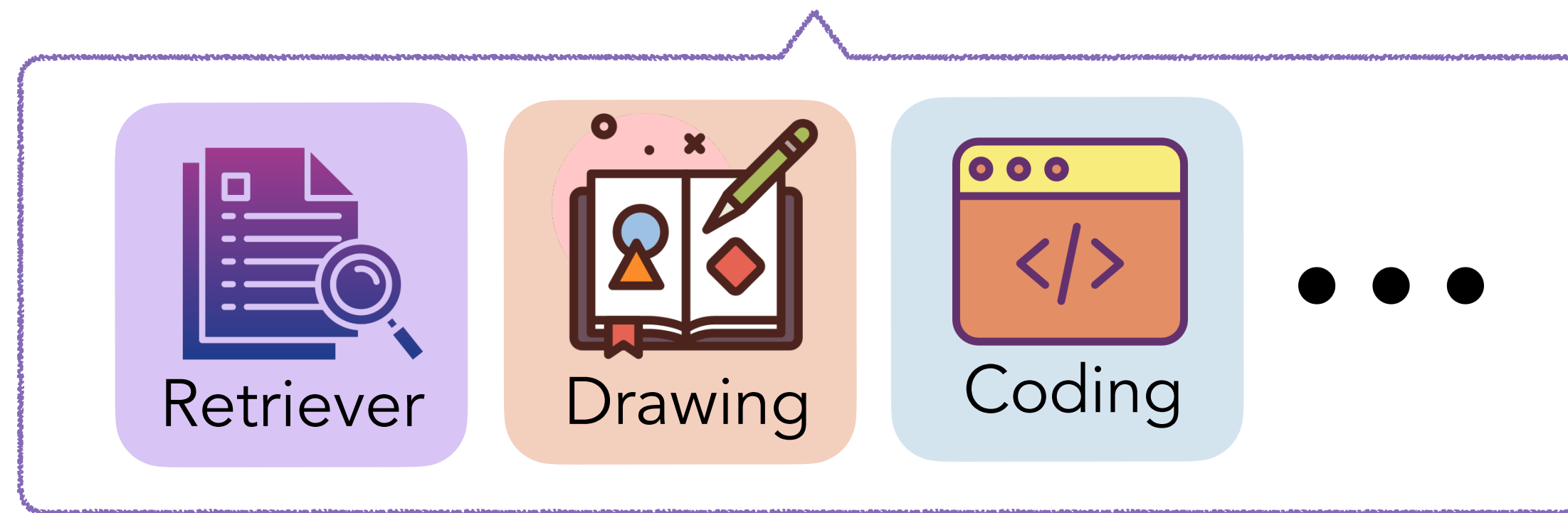
By Nicola Jones

OpenAI o1

Technology firm OpenAI released a preview version of its latest chatbot, o1, last month. Credit: GK

Retriever

Jobs cofounded Apple in his parents' garage

LM

Microsoft
Apple
Google
…

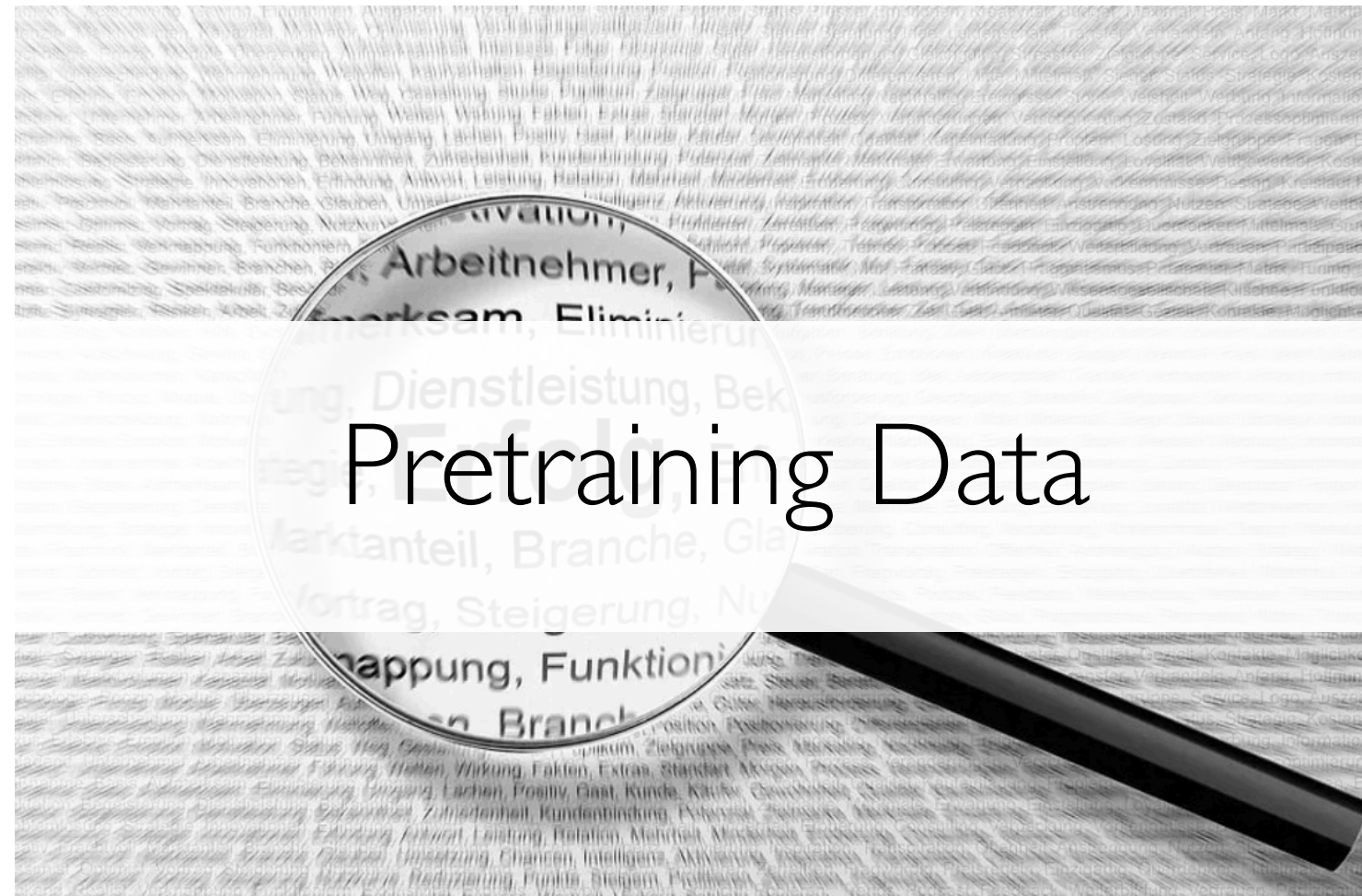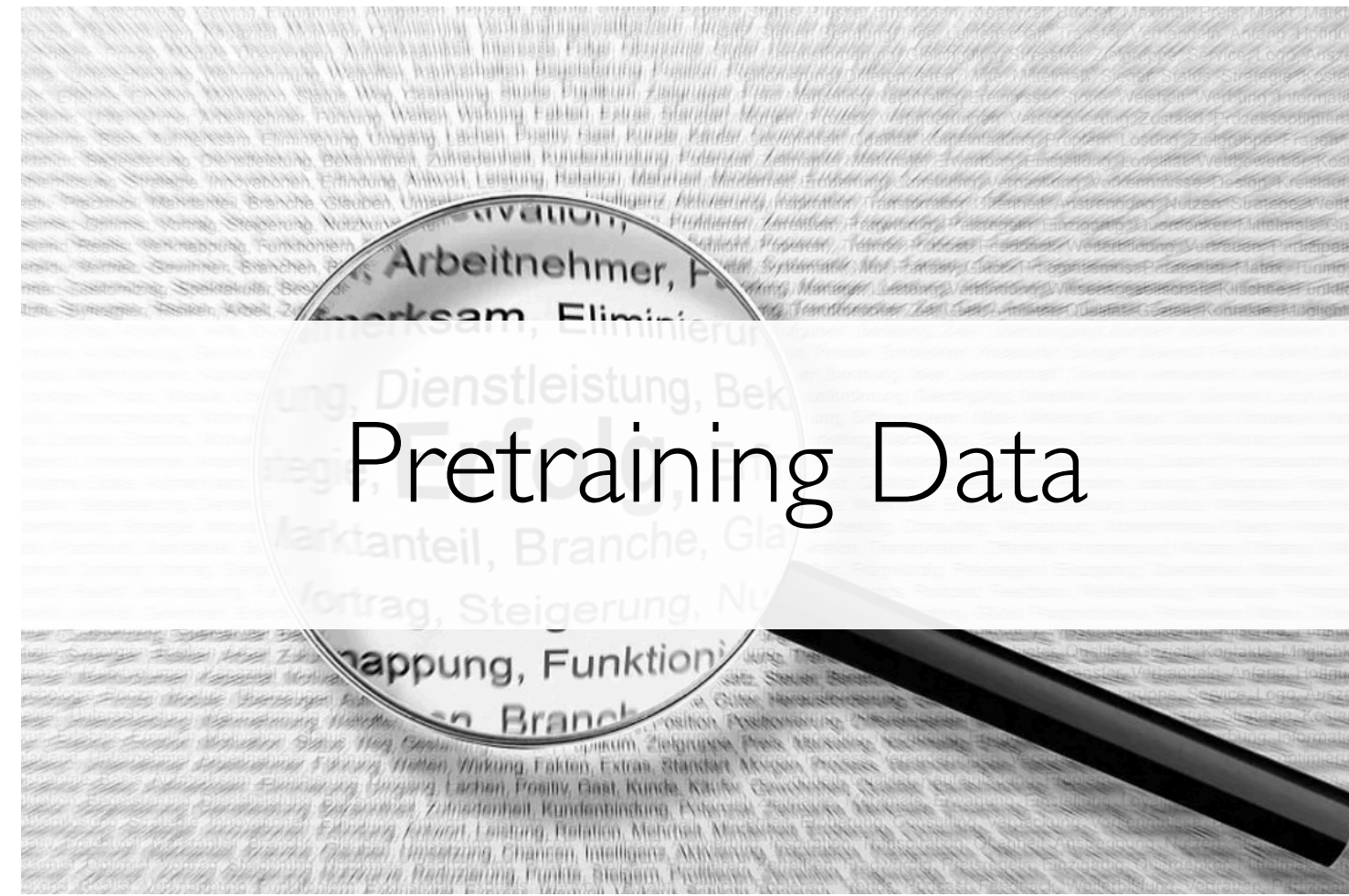✅ **Hard to update new knowledge**

15

LM

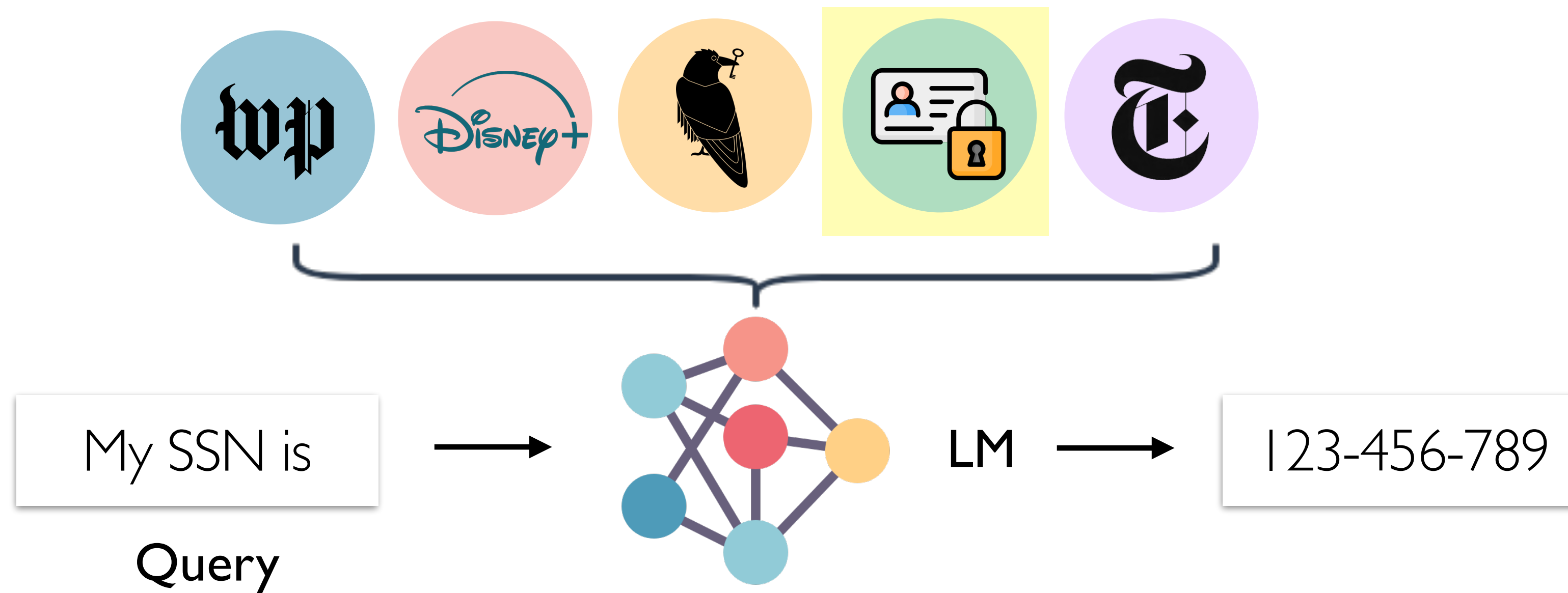Retriever  Drawing  Coding  • • •  *Augmented Models*

# Pretraining Data is not Monolithic



Pretraining Data

*Public*

*Copyright*

*Private*

*Benchmark*
*(contamination)*

Pretraining Data

LM

18

# Data Modularity



Generating a news article

Query

LM

**Copyright and privacy risks**

# Data Modularity



My SSN is

Query

→

LM →

123-456-789

**Copyright and privacy risks**

# Modularity, not Monoliths

*Data Modularity*

LM

*Augmented Models*

Retriever  Drawing  Coding  • • •

# Beyond Monolithic Language Models

*Augmented Models*

*Data Modularity*

# Beyond Monolithic Language Models

*Augmented Models* 📊

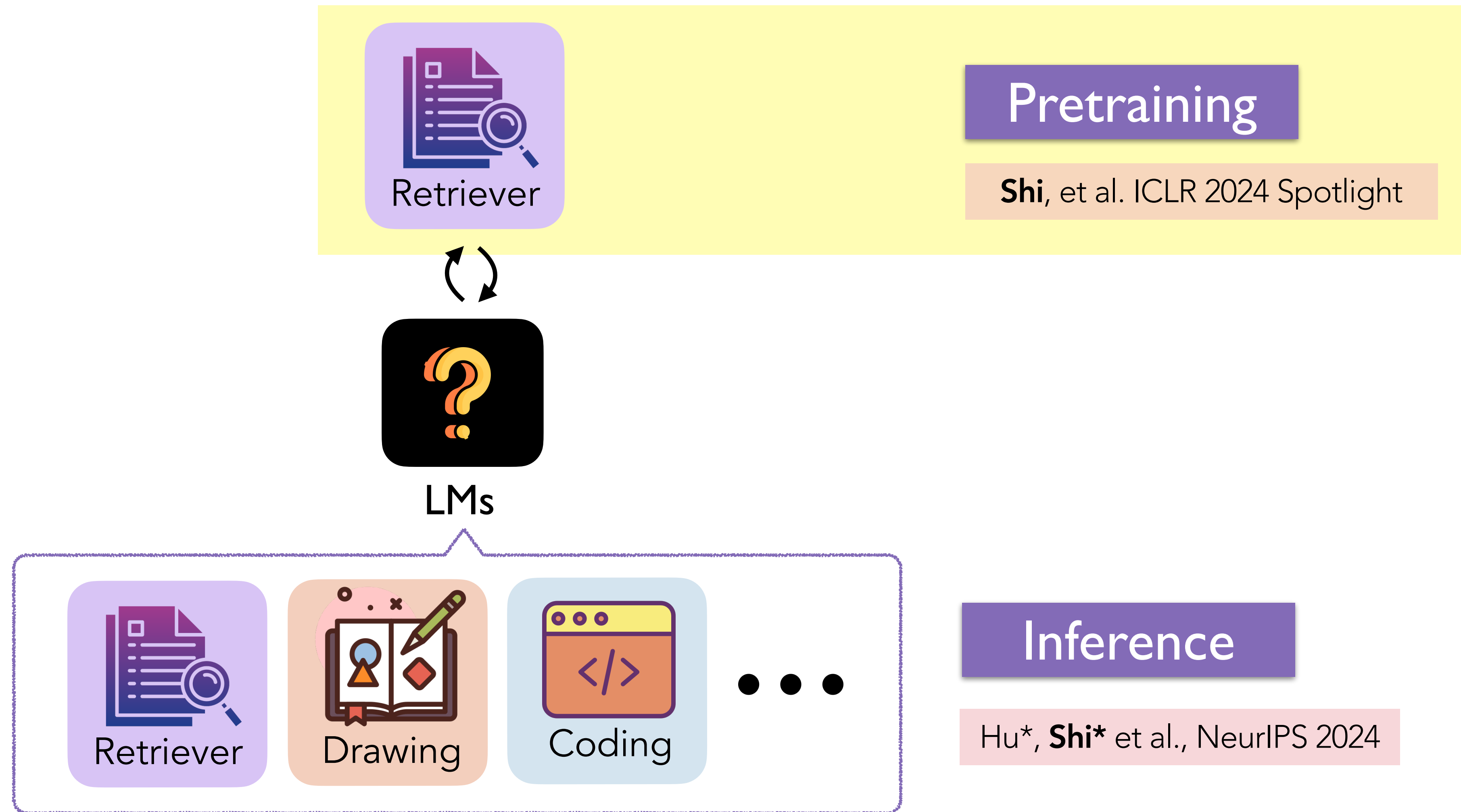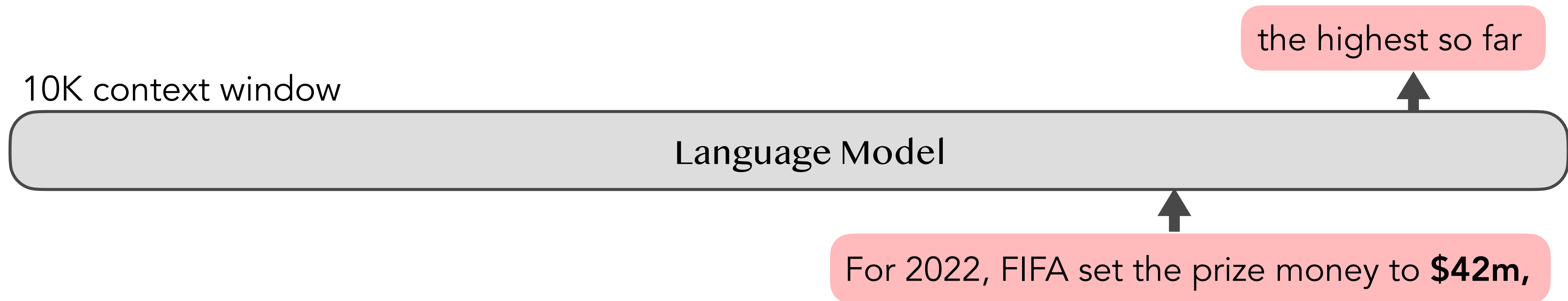*Data Modularity* 🛡️

# Augmented Models

LMs

Retriever  Drawing  Coding  • • •

Inference

Hu*, **Shi*** et al., NeurIPS 2024

# Augmented Models



Retriever

LMs

Retriever | Drawing | Coding | • • •

Pretraining

**Shi**, et al. ICLR 2024 Spotlight

Inference

Hu*, **Shi*** et al., NeurIPS 2024

# Standard Pretraining

the highest so far

10K context window

Language Model

For 2022, FIFA set the prize money to **$42m,**

# Standard Pretraining

Concatenate Random Documents

the highest so far

10K context window

Language Model

Paris is bisected by the River Seine, which flows …

For 2022, FIFA set the prize money to **$42m,**

😞 The prior doc provides no signal for predicting the next doc

🏆 Doc

🗼 Doc

# Problem: Fails to Understand Long Contexts

Input Context

```
Write a high-quality answer for the given question using only the provided search
results (some of which might be irrelevant).

Document [1](Title: Asian Americans in science and technology) Prize in physics for
discovery of the subatomic particle J/ψ. Subrahmanyan Chandrasekhar shared...
Document [2](Title: List of Nobel laureates in Physics) The first Nobel Prize in
Physics was awarded in 1901 to Wilhelm Conrad Röntgen, of Germany, who received...
Document [3](Title: Scientist) and pursued through a unique method, was essentially
in place. Ramón y Cajal won the Nobel Prize in 1906 for his remarkable...

Question: who got the first nobel prize in physics
Answer:
```
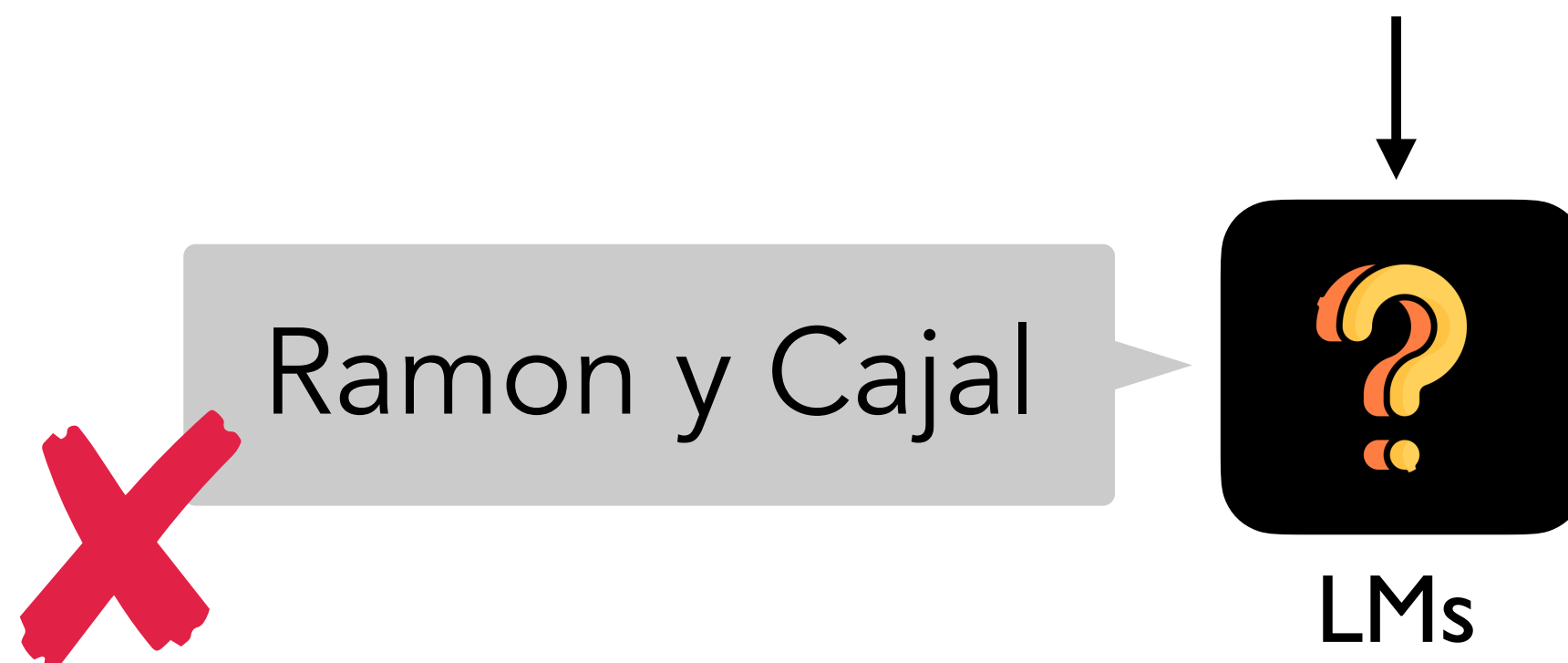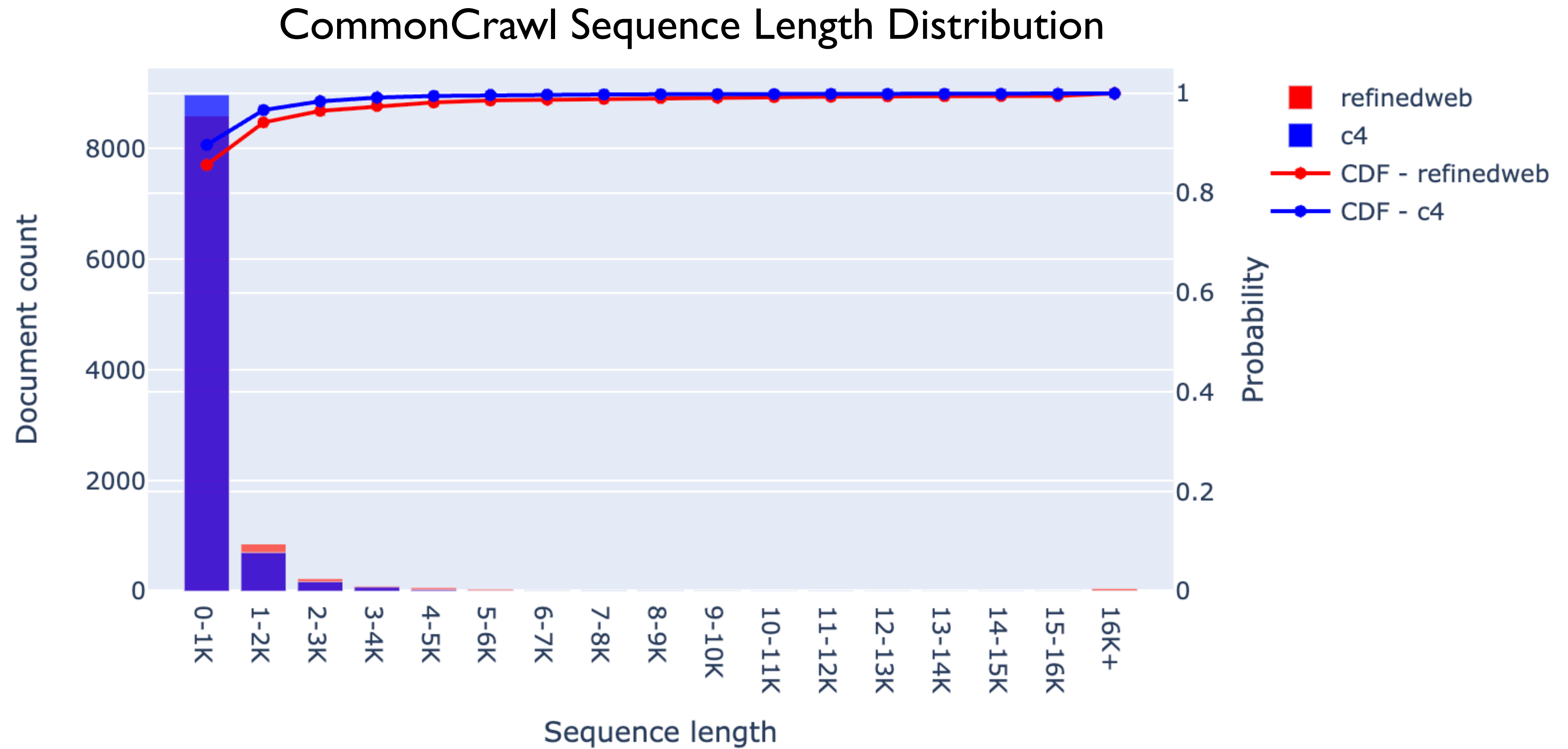
•••

Ramon y Cajal ✗ → LMs

# Problem: Lack of Long Pretraining Documents



CommonCrawl Sequence Length Distribution

# Problem: Lack of Long Pretraining Documents



In the long (context) run

It's not the quadratic attention; it's the lack of long pre-training data

**Harm de Vries**

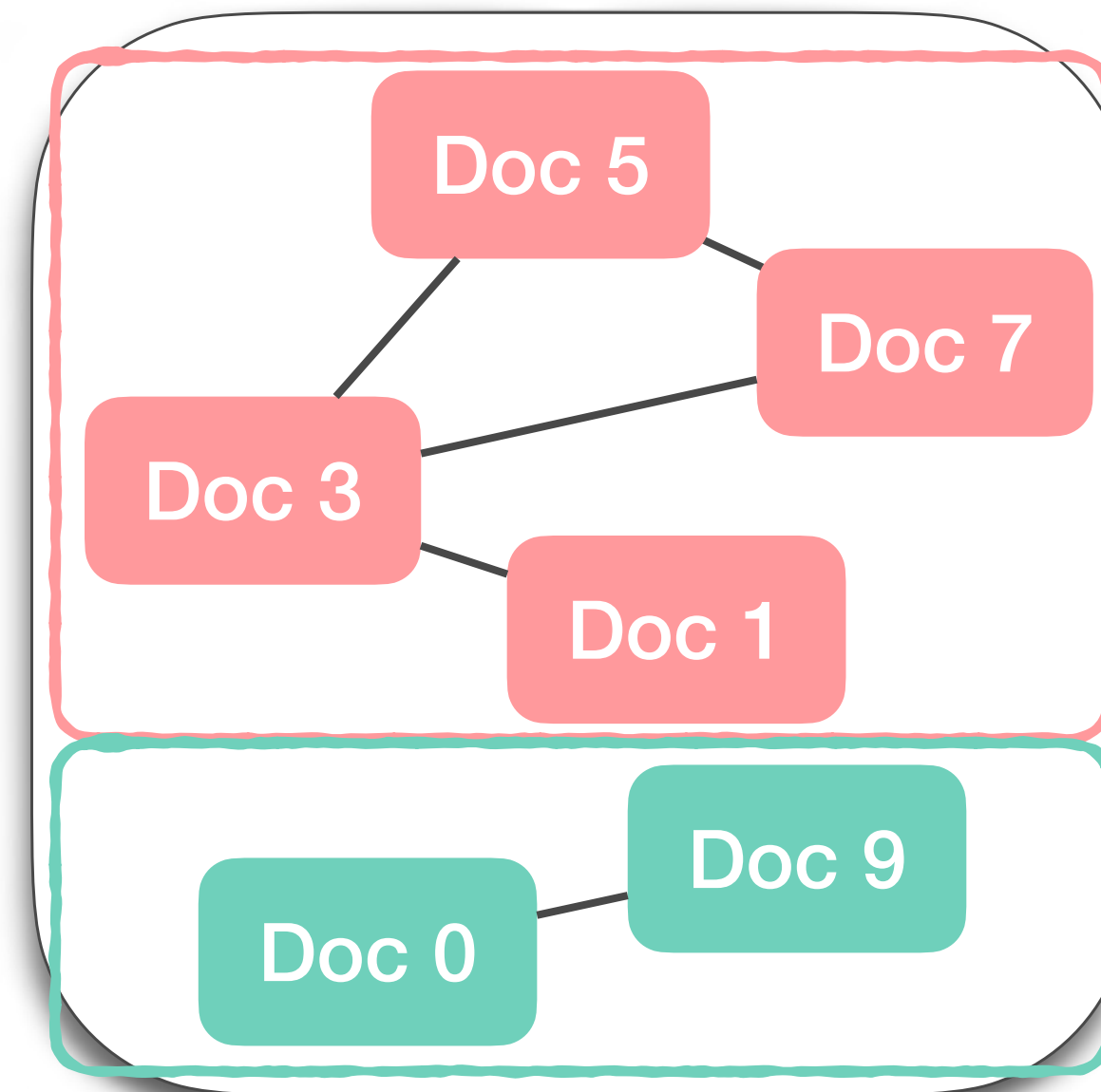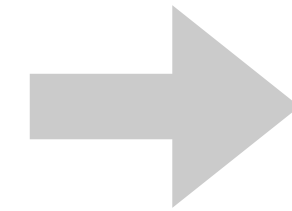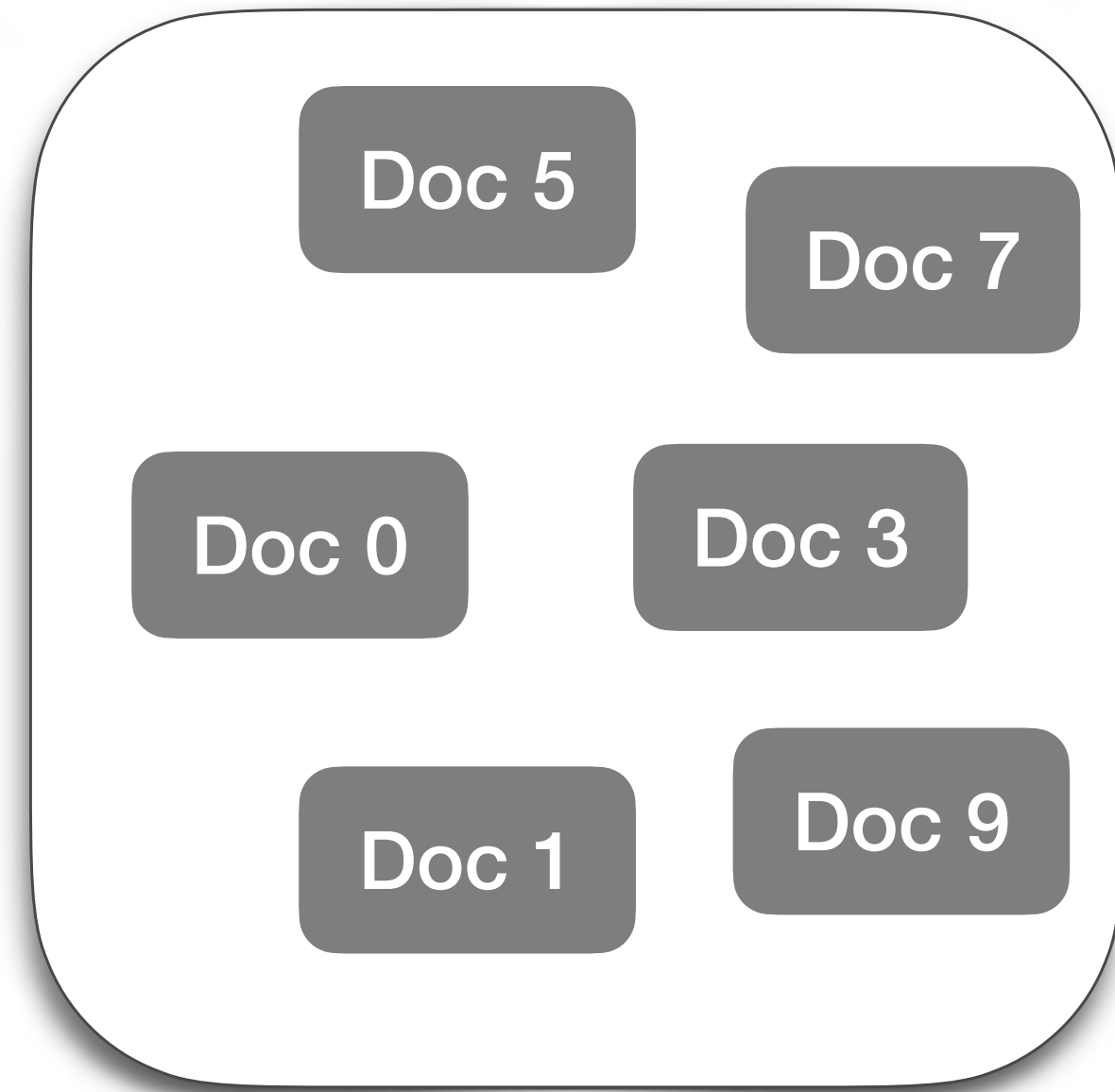Last updated on Sep 16, 2023  ·  21 min read

# Reorder Data w/ Retriever

Pretraining Docs

# Reorder Data w/ Retriever

Pretraining Docs



*Find Related Docs*

One long document

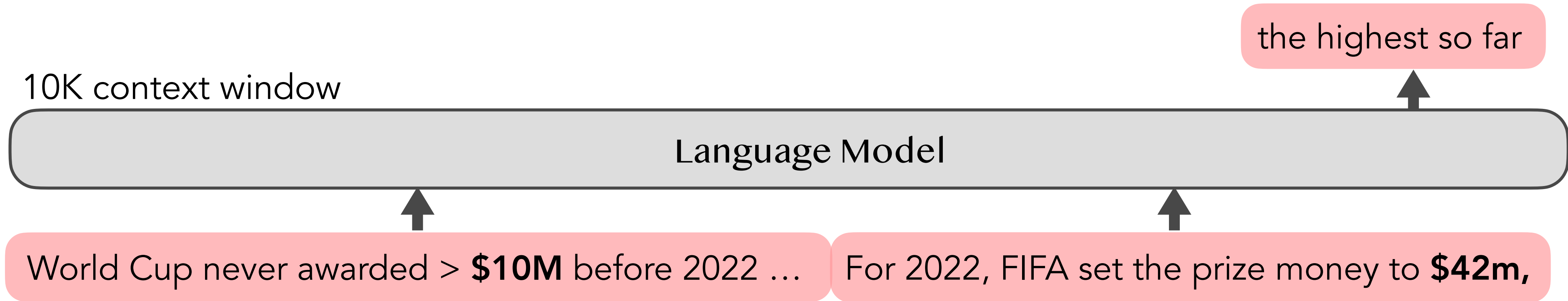# Concatenate **Related** Documents

the highest so far

10K context window

Language Model

World Cup never awarded > **$10M** before 2022 ....     For 2022, FIFA set the prize money to **$42m,**

Doc

Doc

# In-Context Pretraining

Concatenate **Related** Documents

the highest so far

10K context window

Language Model

World Cup never awarded > **$10M** before 2022 …

For 2022, FIFA set the prize money to **$42m,**

😊 Encourage LMs to reason across document boundaries

Doc

Doc

# Pretraining Documents

## In-Context Pretraining

**World Cup**

- World Cup never award …
- For 2022, FIFA set the …
- Messi scored seven …

**Paris**

- Paris is bisected by …
- Paris, France's capital …

…

**Language Model**

Input Contexts

| World Cup never awarded > **$10M** before 2022 … | For 2022, FIFA set the prize money at **$42m,** |

the highest so far

## Standard

**Language Model**

Input Contexts

| Paris is bisected by the River Seine, which flows … | For 2022, FIFA set the prize money at **$42m,** |

the highest so far

# In-Context Pretraining: Recipe

Voldemort had raised his
wand and a flash of

**Doc 0**



Retriever

"Avada Kedavra!" A
jet of green light
issued from …

**Doc 5**

just as a jet of
red light blasted
from Harry's …

**Doc 3**

# In-Context Pretraining: Recipe

Voldemort had raised his
wand and a flash of

**Doc 0**

Retriever

"Avada Kedavra!" A
jet of green light
issued from …

**Doc 5**

just as a jet of
red light blasted
from Harry's …

**Doc 3**

For each doc, can we directly include
its related docs in the context?

# In-Context Pretraining: Recipe

Voldemort had raised his
wand and a flash of

**Doc 0**

One of the three
Unforgivable Curses …

**Doc 1**

red light issued from
Harry's wand …

**Doc 2**

"Avada Kedavra!" A
jet of green light
issued from …

**Doc 5**

"Avada Kedavra!" A
jet of green light
issued from …

**Doc 5**

I don't think
Expelliarmus is
exactly going to

**Doc 7**

just as a jet of
red light blasted
from Harry's …

**Doc 3**

the curse caused
instantaneous and
painless death

**Doc 9**

"Avada Kedavra!" A
jet of green light
issued from …

**Doc 5**

1) *Related* documents in the
same context

2) Each document appears
*exactly once*

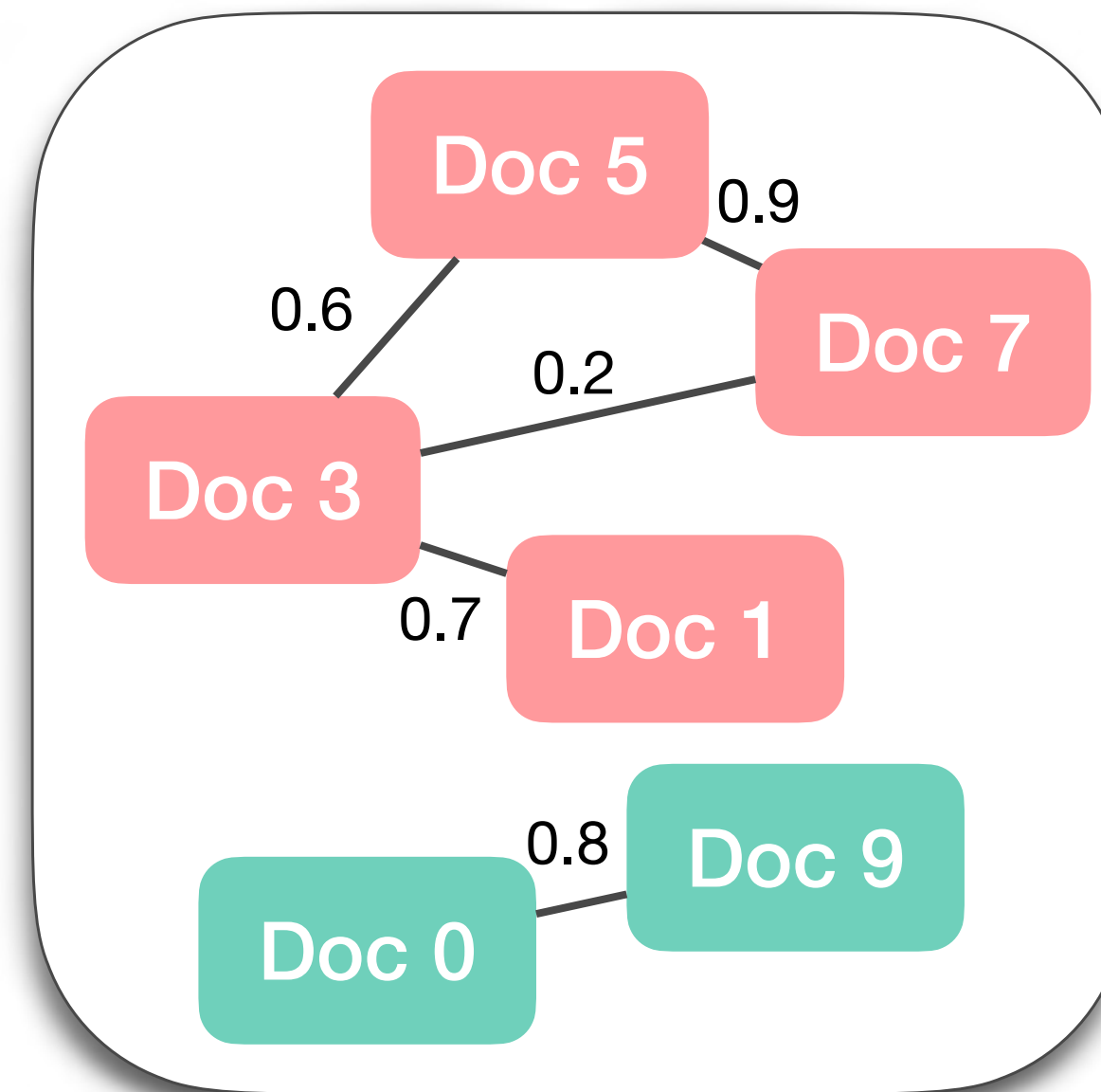# In-Context Pretraining: Recipe

## Document ordering problem

Pretraining Docs

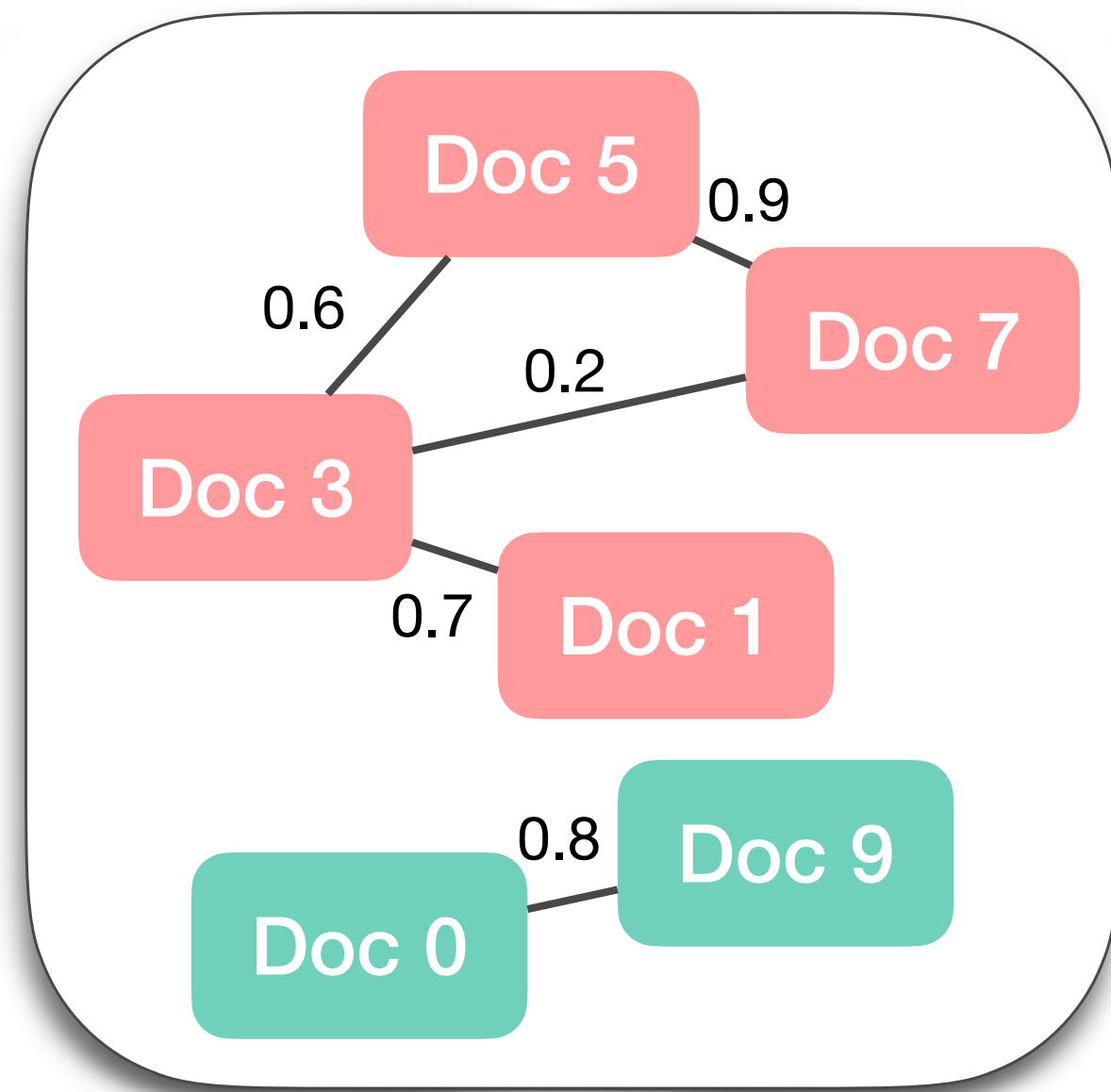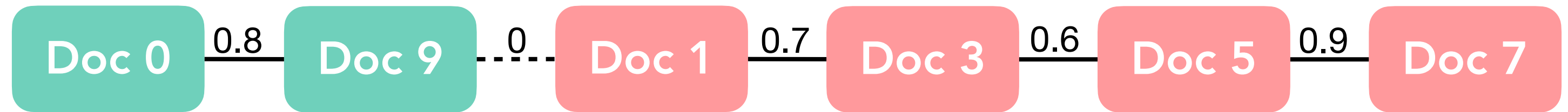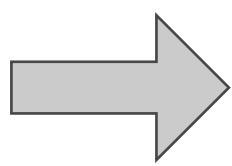# In-Context Pretraining: Recipe

## Document ordering problem

# In-Context Pretraining: Recipe

## Document ordering problem

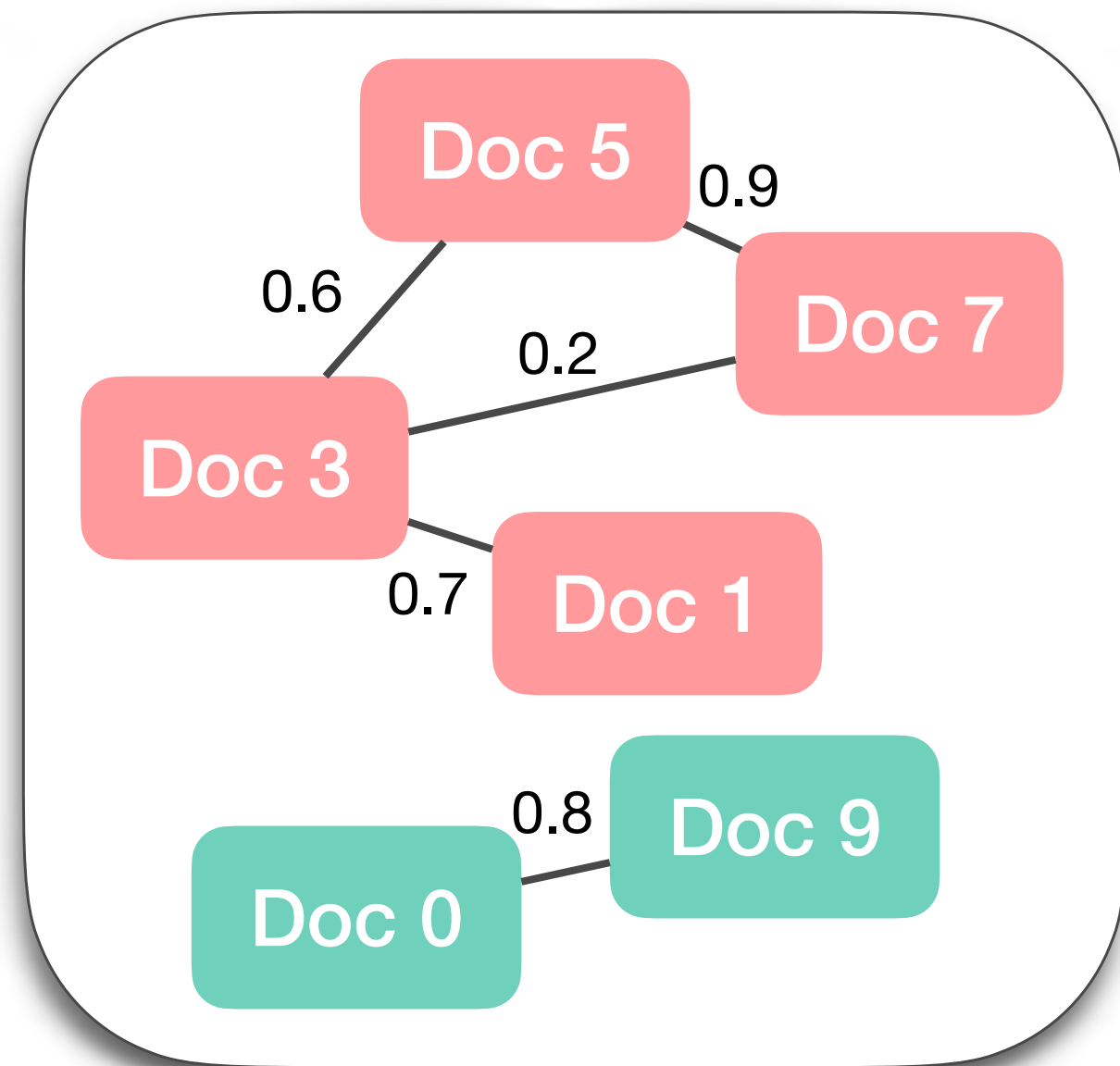*Find a path that visits each doc once, making related docs to be visited consecutively*
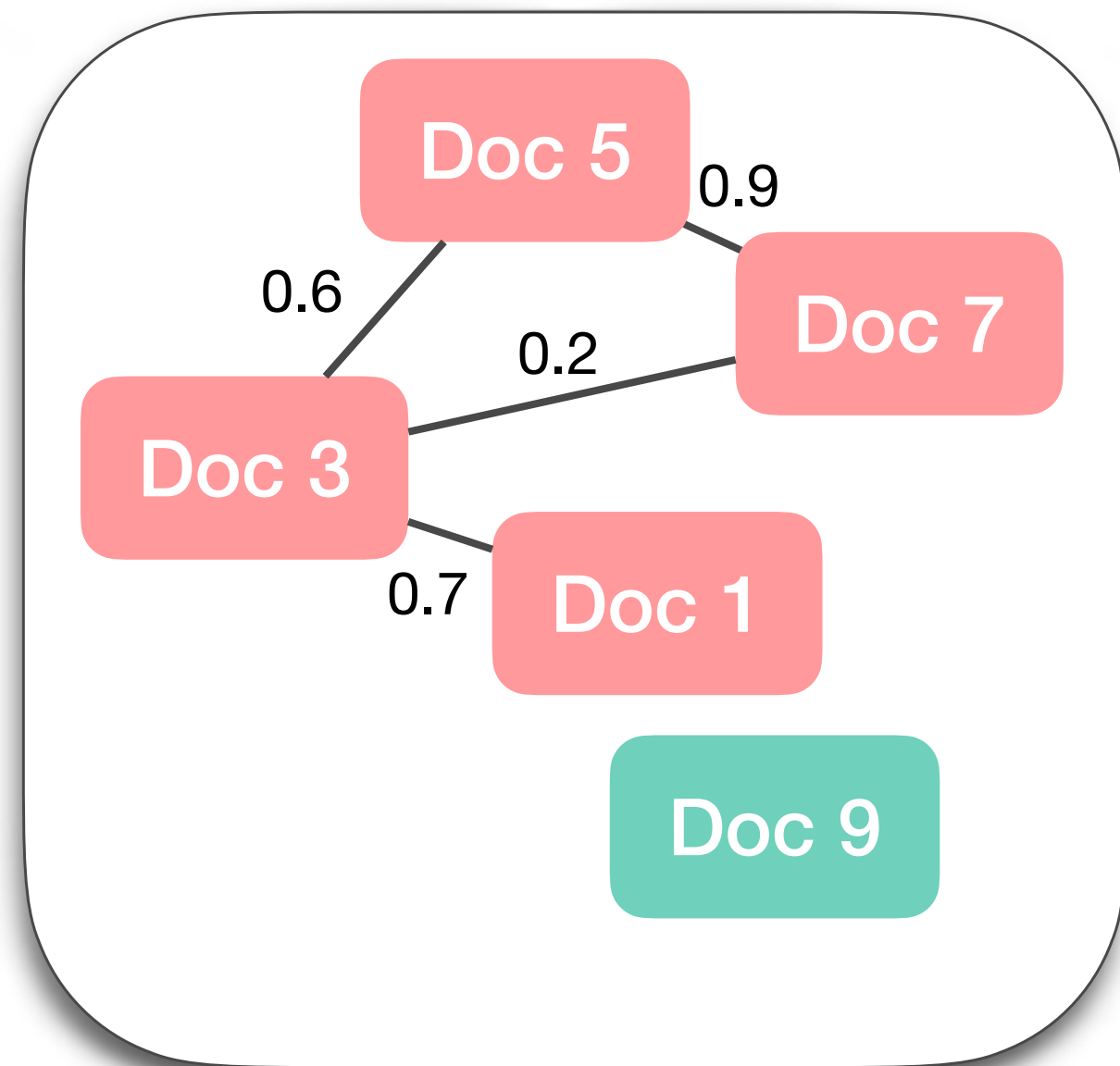
# Document Ordering Problem

**Input:**

**Output:** path



**Procedure:**

select an unvisited doc with the min degree

# Document Ordering Problem

**Input:**

Doc 5

0.6

Doc 3

0.9

Doc 7

0.2

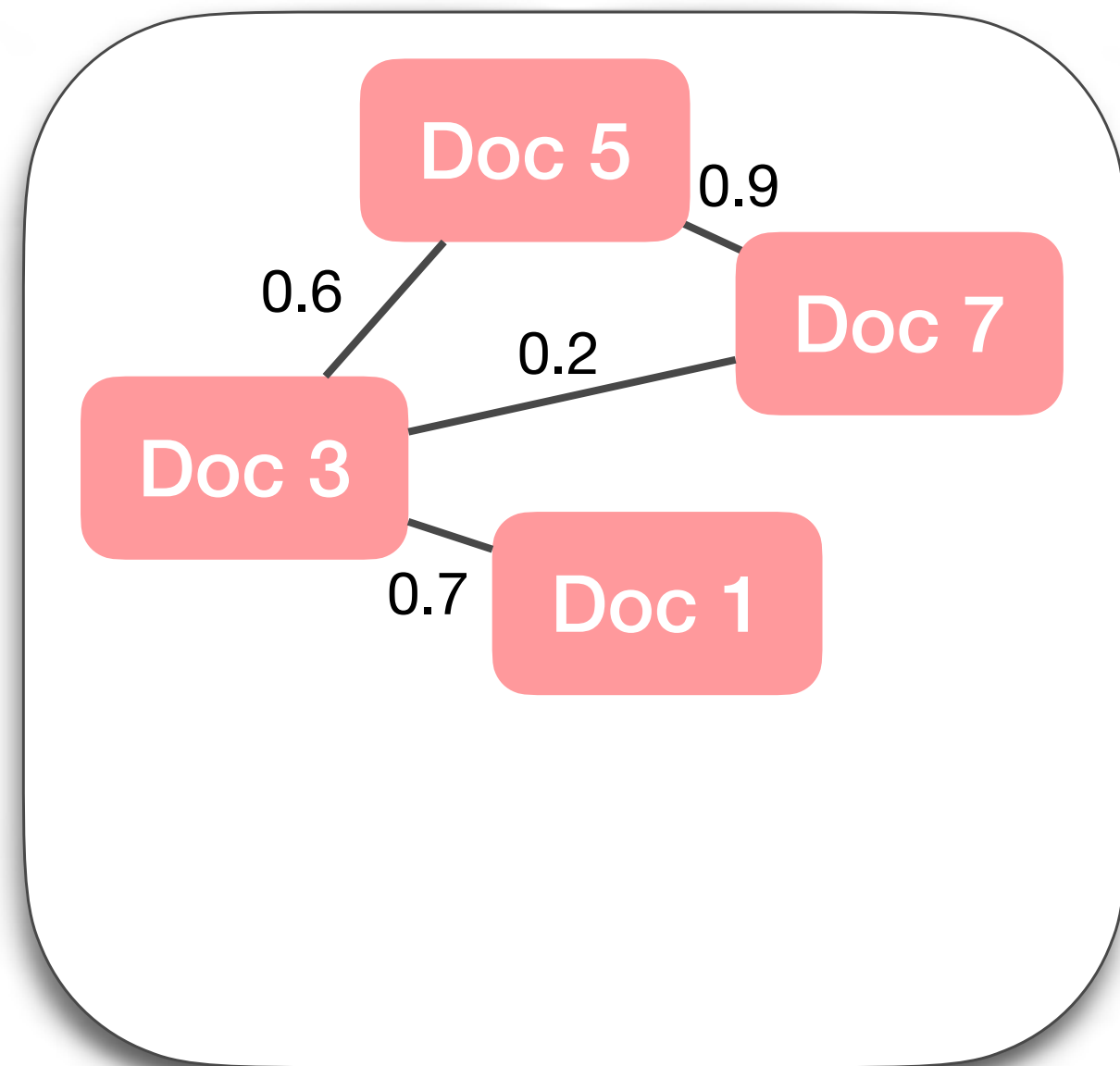0.7

Doc 1
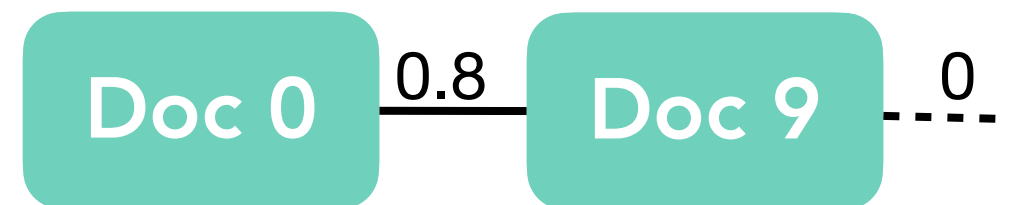
Doc 9

**Output:** path

Doc 0 — 0.8

**Procedure:**

Move to the unvisited neighbor with max weight until all neighbors are visited

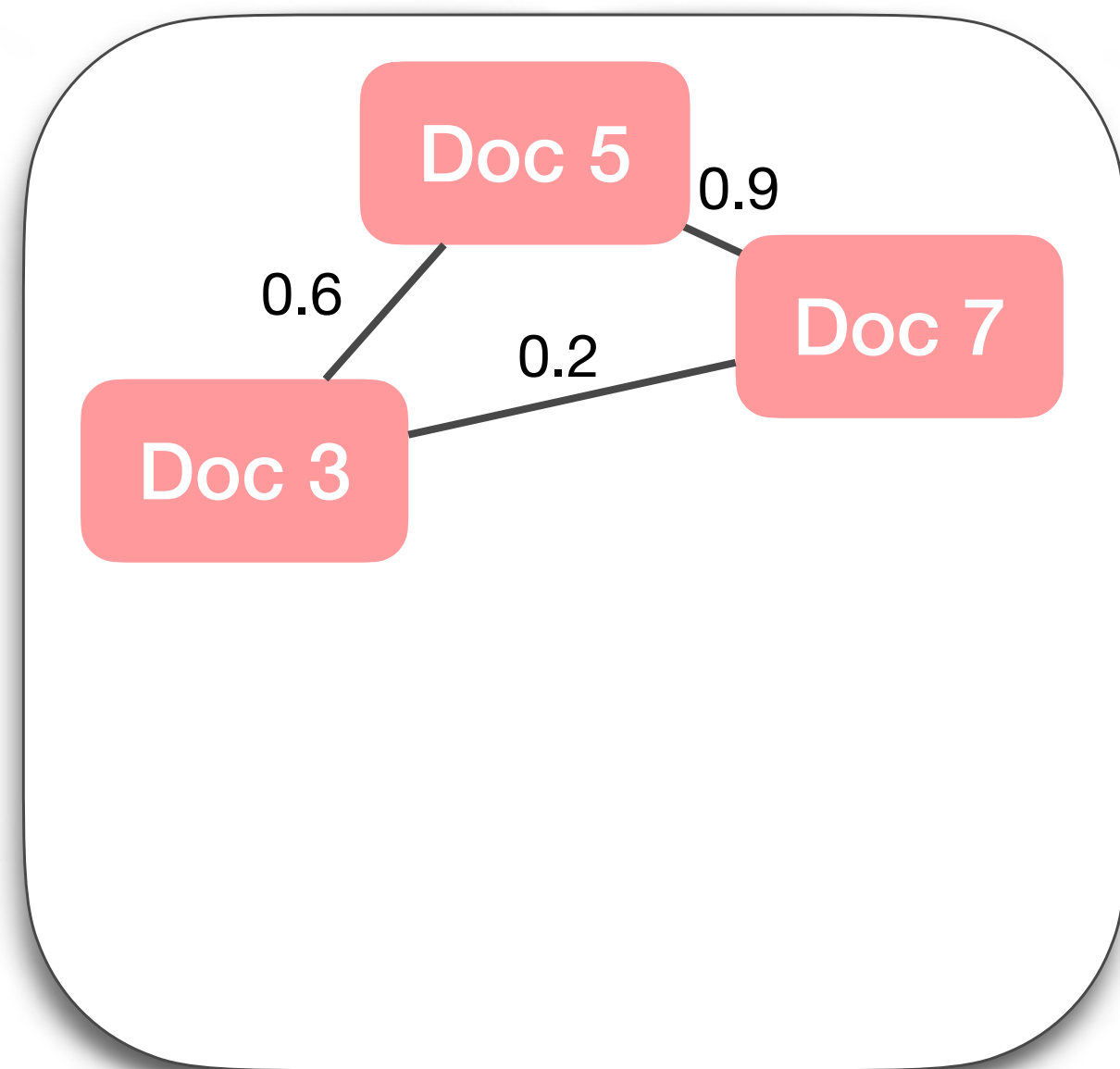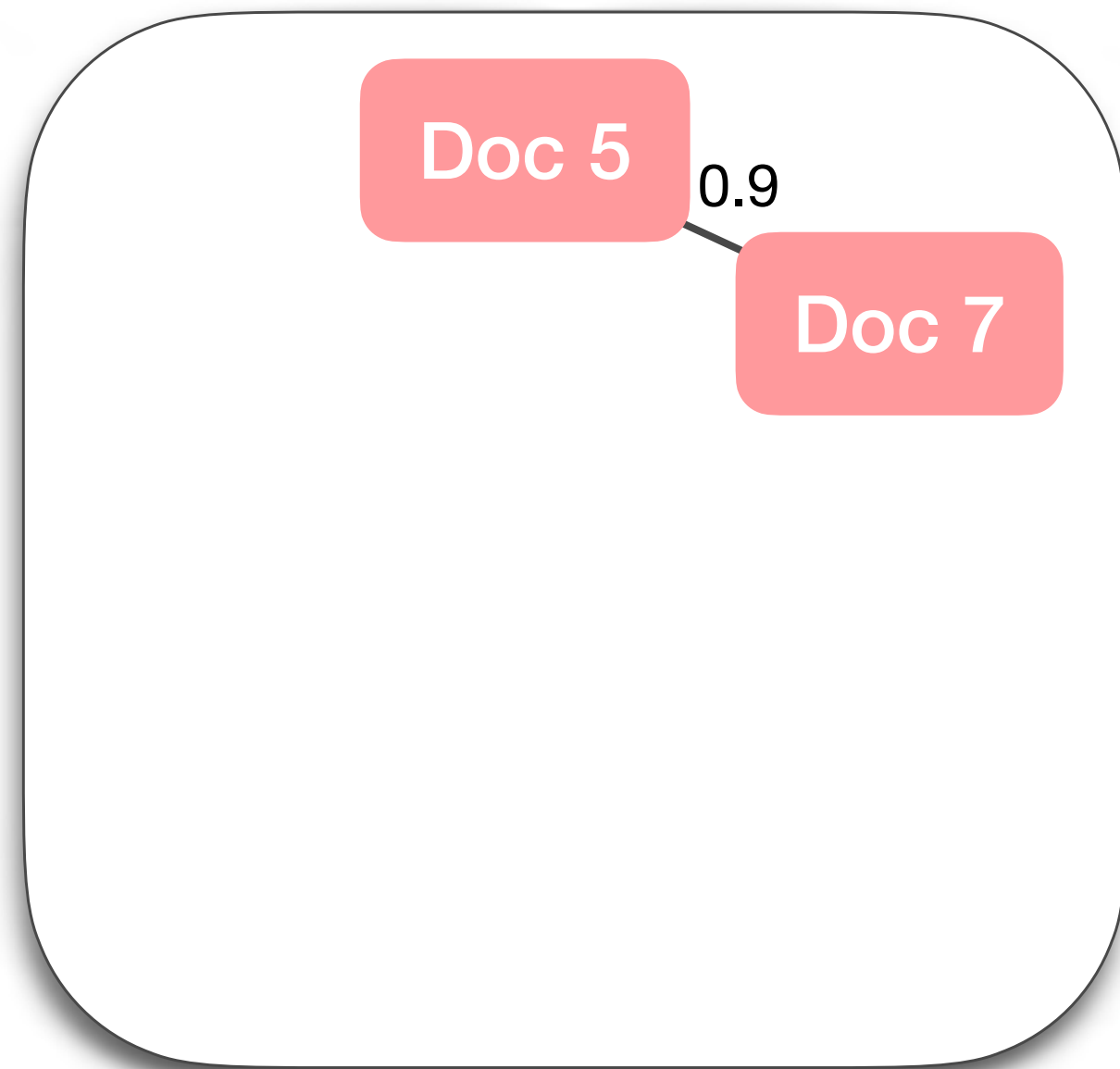# Document Ordering Problem

**Input:**



**Output:** path



Doc 0 — 0.8 — Doc 9 - - - 0...

**Procedure:**
select an unvisited doc with the min degree

# Document Ordering Problem

## Input:

Doc 5

Doc 7

Doc 3

0.6

0.9

0.2

## Output: path

Doc 0 — 0.8 — Doc 9 ---0--- Doc 1 — 0.7 —

## Procedure:

Move to the unvisited neighbor with max weight until all neighbors are visited

# Document Ordering Problem

## Input:

```
          ┌──────────┐
          │ Doc 5 │   │
          │       0.9 │
          │    ┌──────┐│
          │    │ Doc 7 ││
          │    └──────┘│
          └──────────┘
```
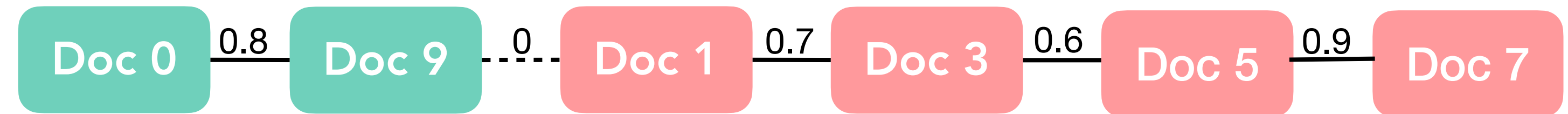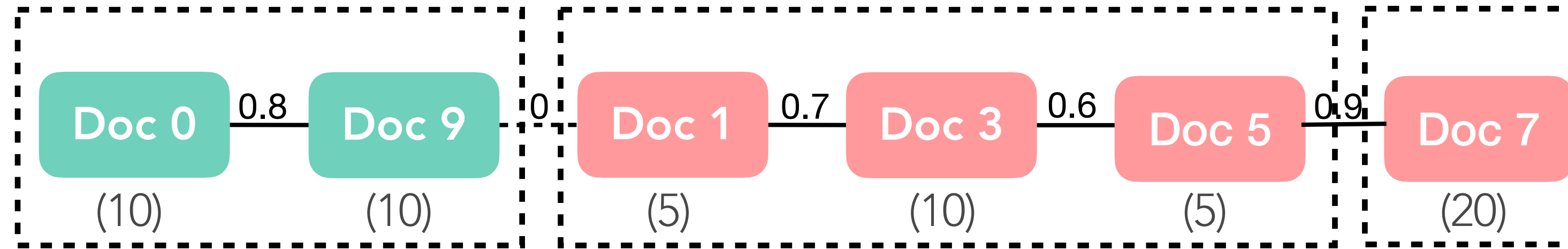
## Output: path

Doc 0 —0.8— Doc 9 ----0--- Doc 1 —0.7— Doc 3 —0.6—          0.9

## Procedure:

Move to the unvisited neighbor with max weight until all neighbors are visited

# Document Ordering Problem

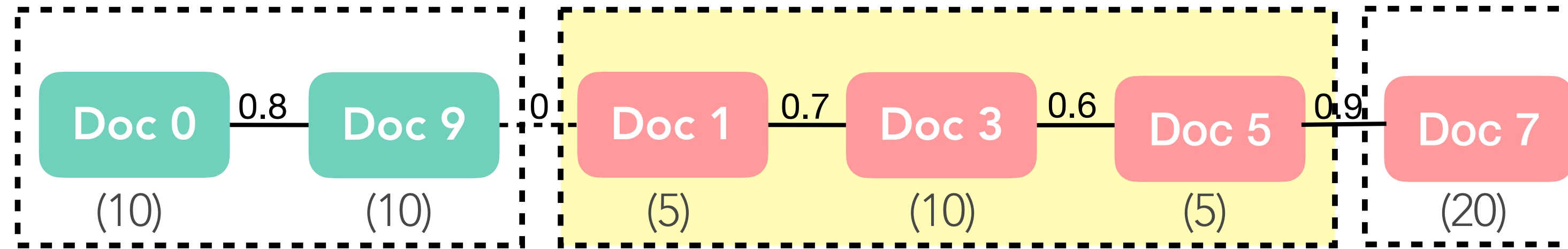Doc 0 — 0.8 — Doc 9 --- 0 --- Doc 1 — 0.7 — Doc 3 — 0.6 — Doc 5 — 0.9 — Doc 7

# Document Ordering Problem

Doc 0 — 0.8 — Doc 9 | 0 | Doc 1 — 0.7 — Doc 3 — 0.6 — Doc 5 | 0.9 | Doc 7

(10) (10) (5) (10) (5) (20)

If max_seq_len = 20

# Document Ordering Problem

Doc 0 — 0.8 — Doc 9 — 0 — Doc 1 — 0.7 — Doc 3 — 0.6 — Doc 5 — 0.9 — Doc 7

(10)  (10)  (5)  (10)  (5)  (20)

If max_seq_len = 20

… Kadavra!" green light …     a jet of red light …     wand and a flash of green

**Language Model**

"Avada Kadavra!" … green …     … as a jet of red light …     … his wand and a flash of

Doc 1     Doc 3     Doc 5

1) *Related* documents in the same context

2) Each document appears *exactly once*

Simple! (Training code remains same)

# Training Details

- **Architecture**: LLaMA

- **Model**: 0.3, 0.7, 1.5, and **7B** model with sequence length of 8192 from scratch **(128 A100s for 9 days)**

- **Data**: 306B tokens from Common Crawl (235M docs)
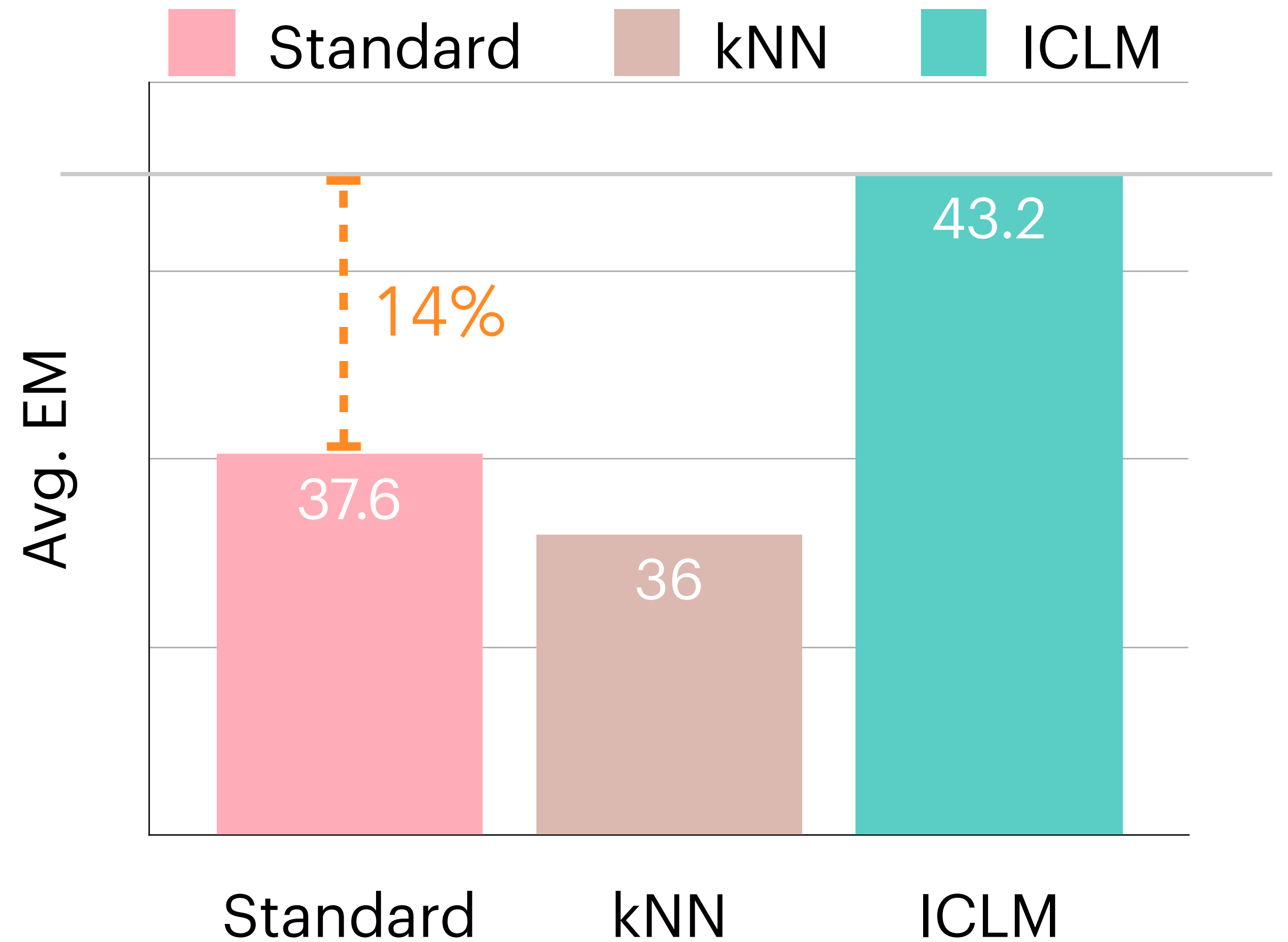
- **Retriever**: Contriever

# Baselines

- **Standard:** places random docs in the input contexts

- **kNN:** places each doc and its retrieved top-k docs in the input

  *Given the same number of training steps, kNN exposes LMs to a less*

  *diverse set of documents, since documents can repeat*

# Results: Reading Comprehension

**Tasks:**

1. **Single document:** race-high, race-middle, boolq, squad
2. **Multi document:** hotpotQA, drop

# Results: Open-Domain QA

**Tasks: NQ, TQA**

**With** retrieved docs

```
Write a high-quality answer for the given question using only the provided search
results (some of which might be irrelevant).

Document [1](Title: Asian Americans in science and technology) Prize in physics for
discovery of the subatomic particle J/ψ. Subrahmanyan Chandrasekhar shared...
Document [2](Title: List of Nobel laureates in Physics) The first Nobel Prize in
Physics was awarded in 1901 to Wilhelm Conrad Röntgen, of Germany, who received...
Document [3](Title: Scientist) and pursued through a unique method, was essentially
in place. Ramón y Cajal won the Nobel Prize in 1906 for his remarkable...

Question: who got the first nobel prize in physics
Answer:
```

**W/o** retrieved docs

```
Question: who got the first nobel prize in physics
Answer:
```
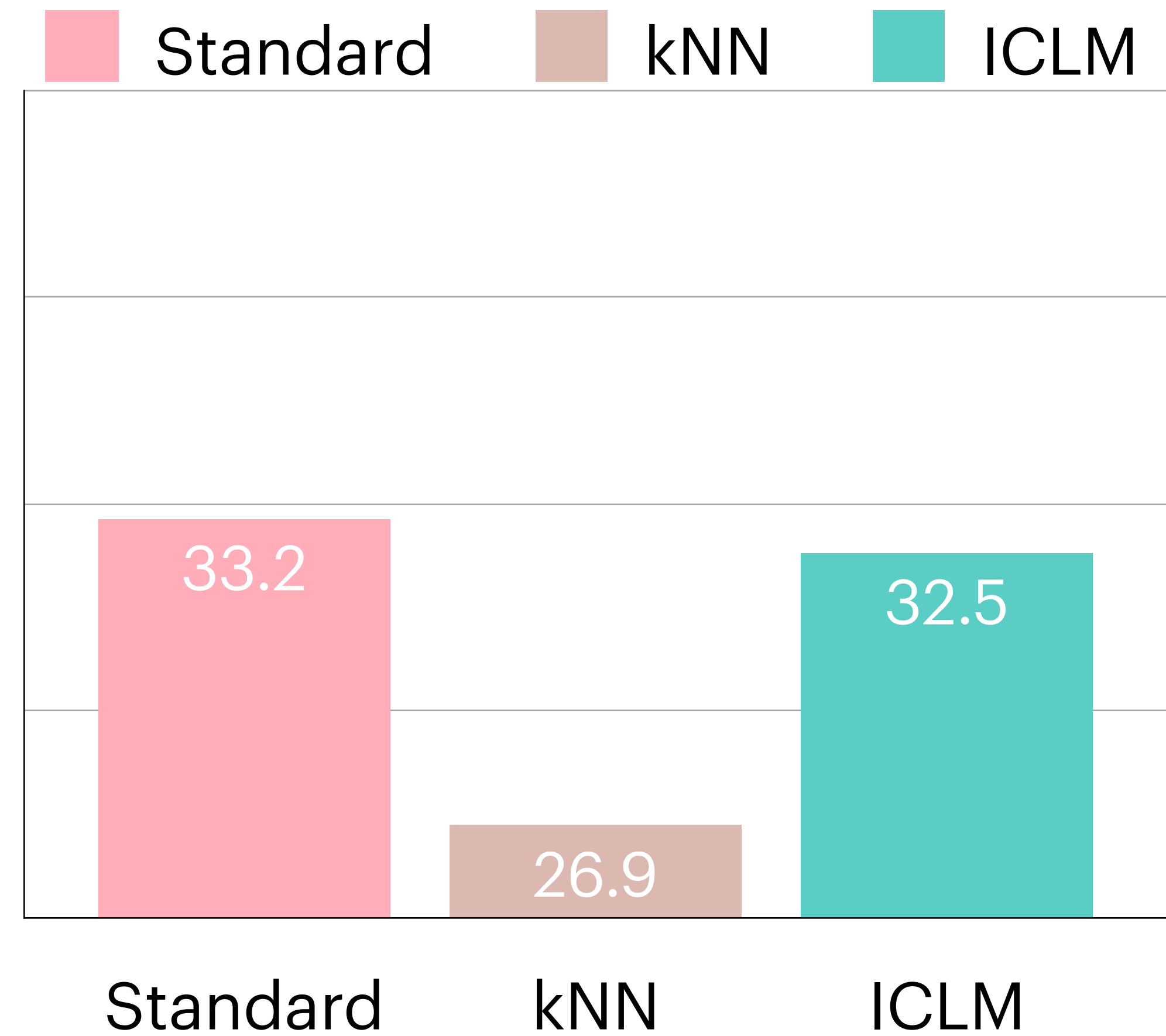
# Results: Open-Domain QA

**With** retrieved docs



Avg. Accuracy

Standard — 38.3
kNN — 31.7
ICLM — 41.9

# Results: Open-Domain QA

## With retrieved docs

Avg. Accuracy

Standard: 38.3
kNN: 31.7
ICLM: 41.9

Standard | kNN | ICLM

## W/o retrieved docs

Avg. Accuracy

Standard: 33.2
kNN: 26.9
ICLM: 32.5

Standard | kNN | ICLM

# Results: Open-Domain QA

**With** retrieved docs

**W/o** retrieved docs



Avg. Accuracy

Standard | kNN | ICLM

38.3 | 31.7 | 41.9

Standard | kNN | ICLM

Avg. Accuracy

Standard | kNN | ICLM

33.2 | | 32.5

Standard | kNN | ICLM

ICLM memorizes less but reads better

# Results

*23 benchmarks in total*

*Accuracy* 👑

*Standard*          *In-Context Pretraining*

| | Standard | | In-Context Pretraining |
|---|---|---|---|
| 🗄 Open-Domain QA (w/ retrieval) | 38% | **10.5%** → | 42% |
| ▦ In-Context Learning | 66% | **7.5%** → | 71% |
| 📖 Reading Comprehension | 37% | **14.0%** → | 43% |
| ▤ Factuality | 44% | **15.9%** → | 51% |
| 📜 Long Document Reasoning | 32% | **7.5%** → | 34% |

# Evolution of Performance

Reading comprehension

Open-domain QA (w/ retrieval)



(b) Race-High

(c) NQ (Open)

Consistent performance improvement

# What's Next?

- How similar should the documents be in the same context?

- Does In-Context Pretraining work for continual pretraining?

# Augmented Models

Retriever
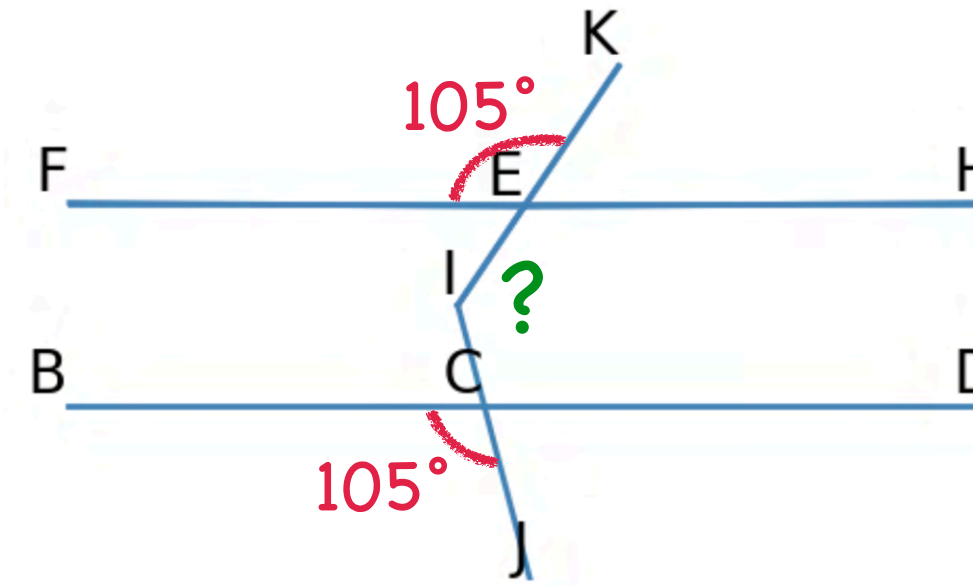
LMs

Retriever    Drawing    Coding    ● ● ●

Pretraining

**Shi**, et al. ICLR 2024 Spotlight

Inference

Hu*, **Shi*** et al., NeurIPS 2024

Given ∠BCJ=105°,
∠KEF=105°. Find ∠EIC



From the given image, we are provided with two angles:

1. ∠BCJ = 105°
2. ∠KEF = 105°

We need to find the measure of ∠EIC.
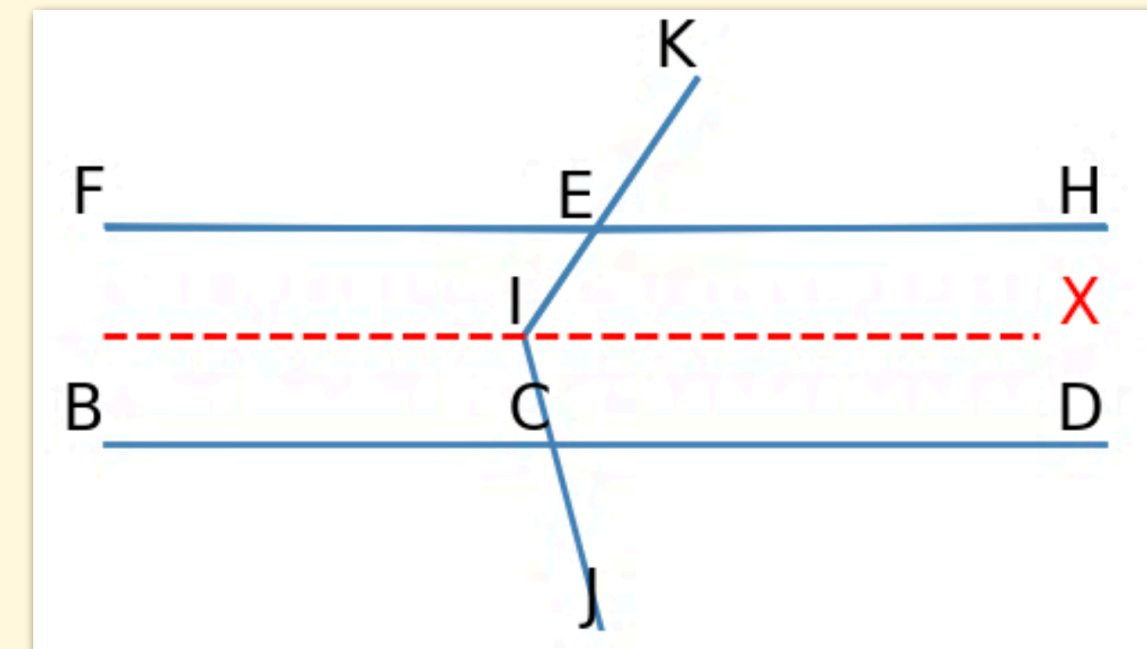
Looking at the diagram, it seems that ∠EIC is the angle formed by lines CE and CI. Let's break this down:

1. **Angle relationships**: The lines BC and CD are parallel, and similarly, the lines FE and HD are parallel. The transversal KE cuts both sets of parallel lines.

   • Since ∠BCJ and ∠KEF are both 105° and form corresponding angles with the transversal, they imply the parallelism of lines.

2. **Vertically opposite angles**: Since ∠BCJ = 105°, the vertically opposite angle ∠EIC must also be 105° because vertically opposite angles are always equal.

Therefore, the measure of ∠EIC is **105°**.

Given ∠BCJ=105°, ∠KEF=105°. Find ∠EIC



105°

K

F    E    H

I   ?

B    C    D

105°

J

From the given image, we are provided with two angles:

1. ∠BCJ = 105°
2. ∠KEF = 105°

We need to find the measure of ∠EIC.

Looking at the diagram, it seems that ∠EIC is the angle formed by lines CE and CI. Let's break this down:

1. **Angle relationships**: The lines BC and CD are parallel, and similarly, the lines FE and HD are parallel. The transversal KE cuts both sets of parallel lines.

   • Since ∠BCJ and ∠KEF are both 105° and form corresponding angles with the transversal, they imply the parallelism of lines.

2. **Vertically opposite angles**: Since ∠BCJ = 105°, the vertically opposite angle ∠EIC must also be 105° because vertically opposite angles are always equal.

Therefore, the measure of ∠EIC is **105°**.

✗

🙋 **Draw line IX parallel to EH**



K

F    E    H

I     X

B    C    D

J

Given ∠BCJ = ∠KEF = 105°, their supplementary angles, ∠KEH = ∠DCJ = 180° - 105° = 75°.
Given BD parallel to IX, it follows that ∠EIX = ∠KEH, and ∠XIC = ∠DCJ, both measuring 75°.
Therefore, ∠EIC = ∠EIX + ∠XIC = 150°

✓

# Sketching ✍️

Sketching is a fundamental human activity, serving as a versatile tool for **communication**, **ideation**, and **problem-solving**
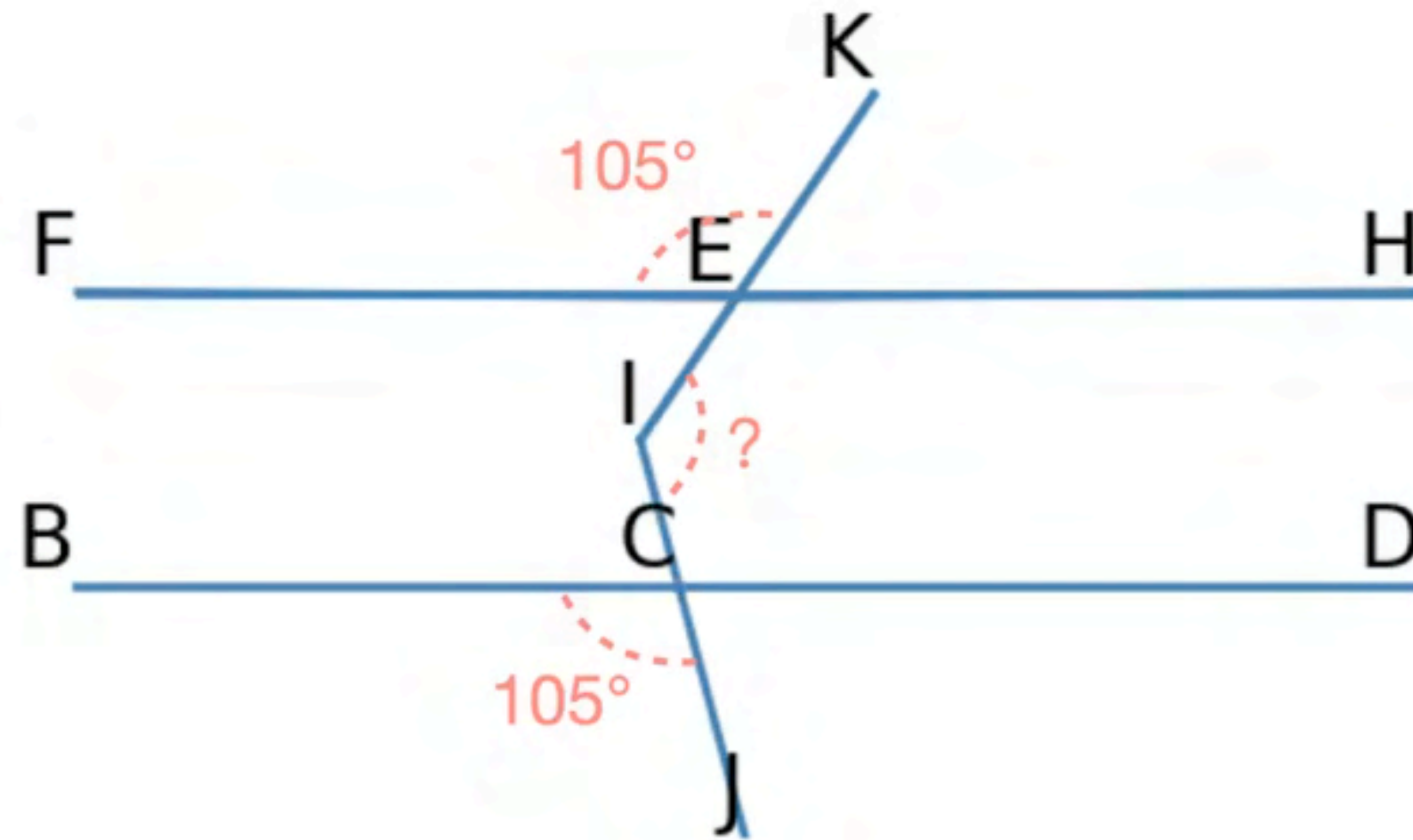




Why not let multimodal LMs do so?

# Visual Reasoning w/ Sketchpad

Augmenting multimodal LMs with a visual canvas and drawing tools ✍️

**Query:** Given ∠BCJ=105°, ∠KEF=105°. Find ∠EIC

**Input Image:**

# Visual Sketchpad: Recipe
## Sketching via Code Generation

# Visual Sketchpad: Recipe

## Sketching via Code Generation

Which point is closer to the camera, A or B?

A

B

# Visual Sketchpad: Recipe
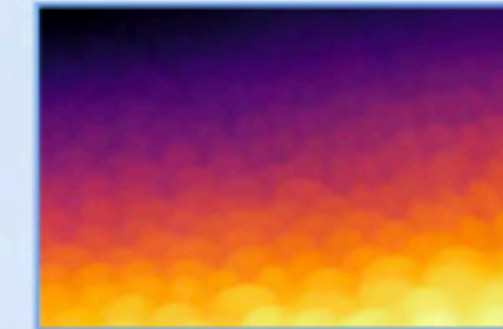## Sketching via Code Generation



**?** Which point is closer to the camera, A or B?

**THOUGHT 1**: I can use a depth estimation tool that provides a depth map with colors indicating distances.

# Visual Sketchpad: Recipe
## Sketching via Code Generation

# Visual Sketchpad: Recipe
## Sketching via Code Generation



? Which point is closer to the camera, A or B?

**THOUGHT 1**: I can use a depth estimation tool that provides a depth map with colors indicating distances.

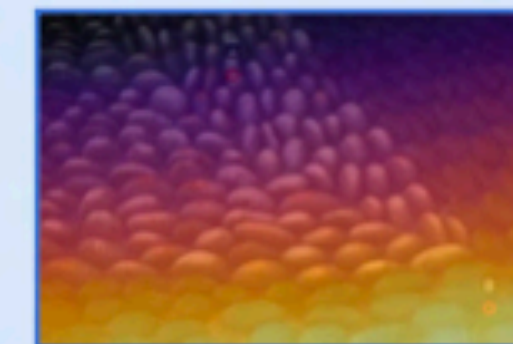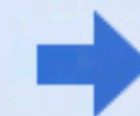**ACTION 1:** ✏️ Call depth estimation

```
depth_map = depth(image)
display(depth_map)
```

➡️

**THOUGHT 2**: The depth map uses colors from purple (farther) to yellow (closer). I need to overlay the original image with the depth map

# Visual Sketchpad: Recipe
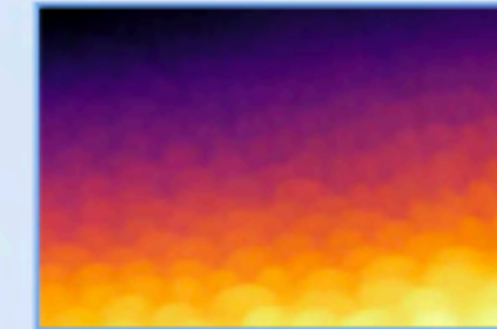## Sketching via Code Generation



**?** Which point is closer to the camera, A or B?

**THOUGHT 1**: I can use a depth estimation tool that provides a depth map with colors indicating distances.

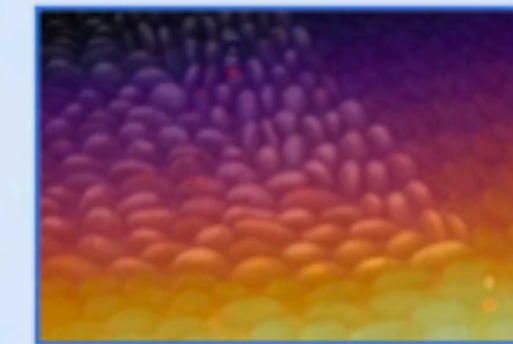**ACTION 1**: ✏️ Call depth estimation

```
depth_map = depth(image)
display(depth_map)
```

➡️

**THOUGHT 2**: The depth map uses colors from purple (farther) to yellow (closer). I need to overlay the original image with the depth map

**ACTION 2**: ✏️ Call overlay

```
overlay = overlay_images(depth_map, image)
display(overlay)
```
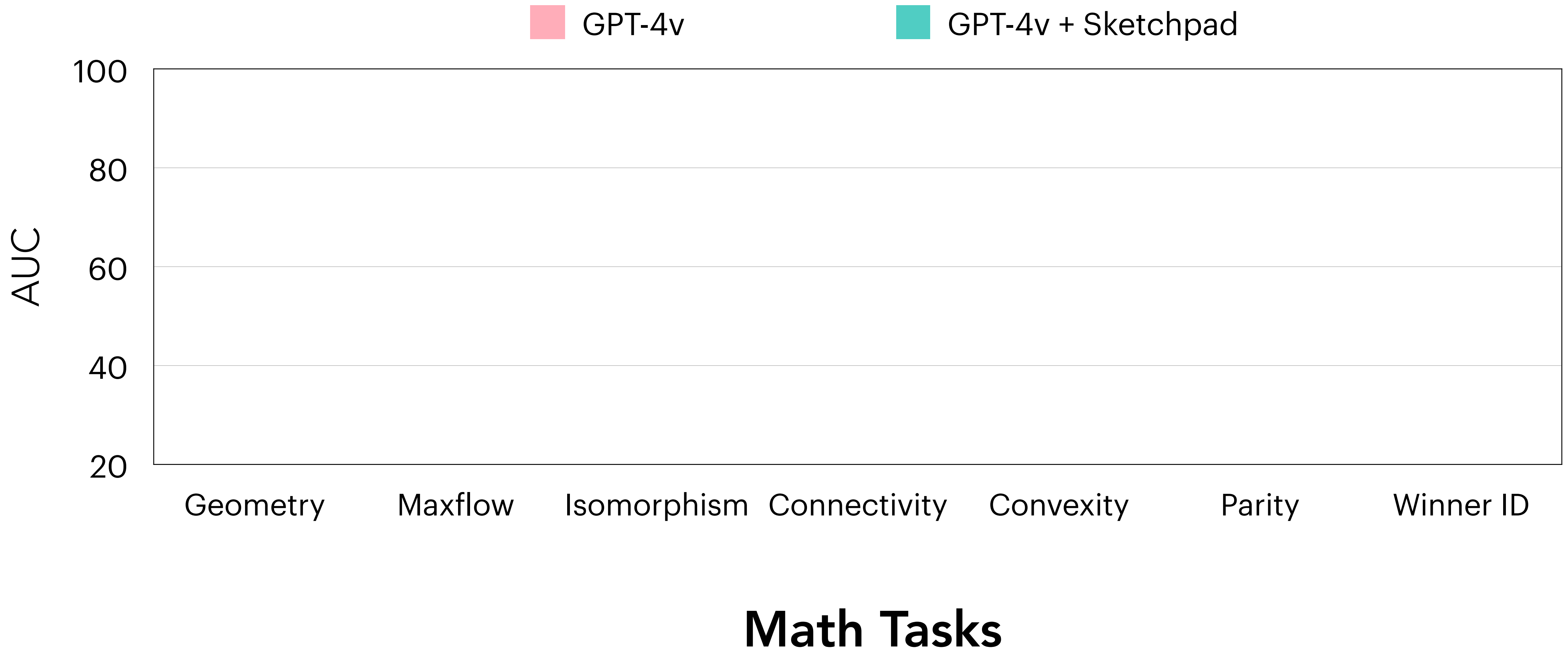
➡️

# Visual Sketchpad: Recipe
## Sketching via Code Generation

# Results

# Results



**Math Tasks**

# Results



Legend: GPT-4v, GPT-4v + Sketchpad

Callout: 12.7% improvement on average

Y-axis: AUC (20, 40, 60, 80, 100)

X-axis (Math Tasks): Geometry, Maxflow, Isomorphism, Connectivity, Convexity, Parity, Winner ID

# Results



GPT-4v     GPT-4v + Sketchpad

**Visual Reasoning Tasks**

# Results



8.6% improvement on average

GPT-4v    GPT-4v + Sketchpad

**Visual Reasoning Tasks**

# What's Next?

## Sketch to UI Design with Multimodal LMs

# Summary: Augmented Models



Retriever

LMs

Retriever   Drawing   Coding   ● ● ●

Pretraining

**Shi**, et al. ICLR 2024 Spotlight

Inference

Hu*, **Shi\*** et al., NeurIPS 2024

# Beyond Monolithic Language Models

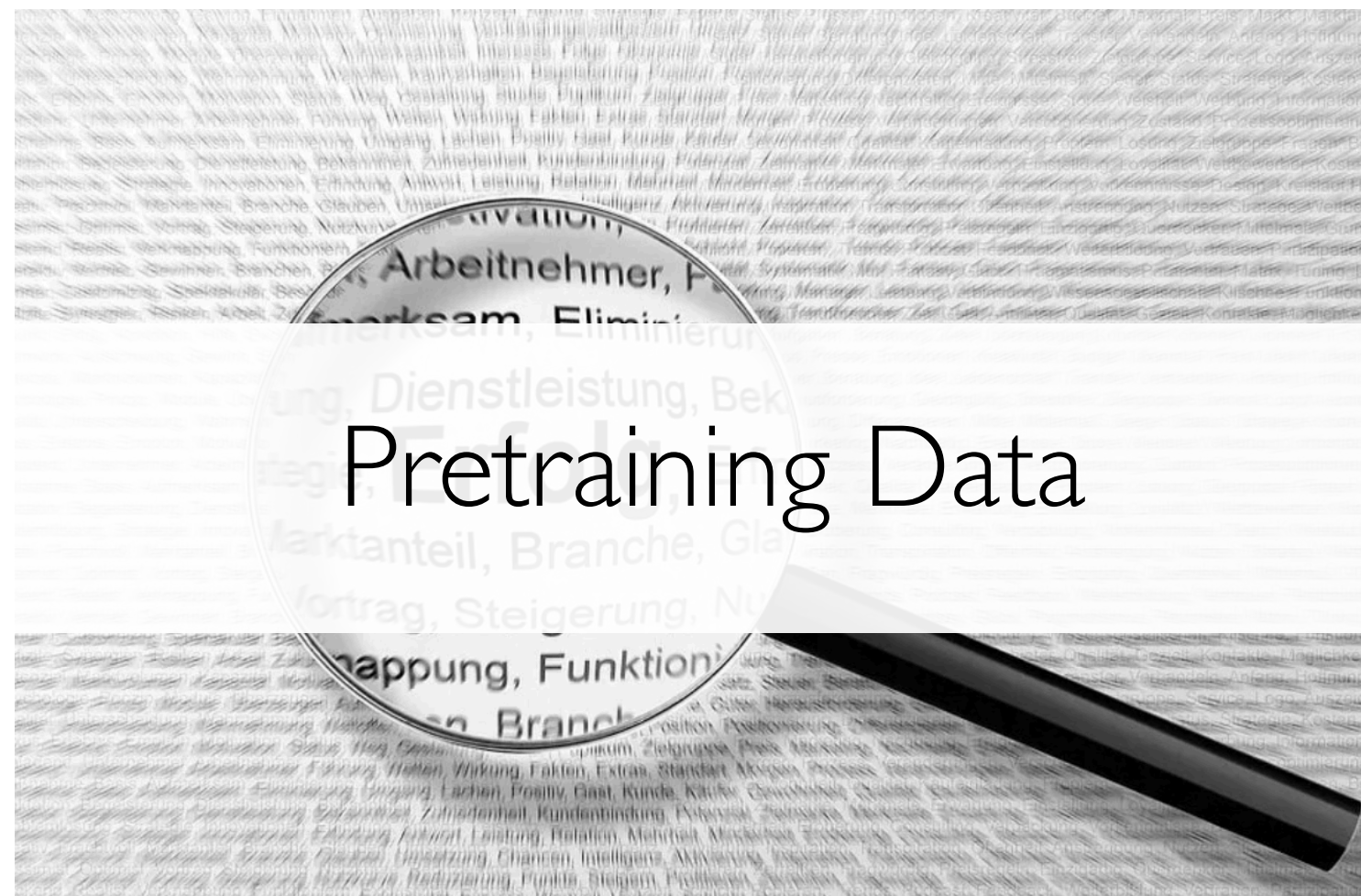*Augmented Models*

*Data Modularity*

Pretraining Data

Pretraining Data

PROJECT GUTENBERG    MIT OCW    *Public*

The New York Times    The Guardian    Disney+    *Copyright*

*Private*

MMLU    *Benchmark (contamination)*

# Copyright Risks in LMs

**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.

Dec. 27, 2023

---

Reuters

World | Business | Markets | Sustainability | Legal | Breakingviews | Technology | Investigations

Litigation | Copyright | Litigation | Technology | Intellectual Property

**Music publishers ask court to halt AI company Anthropic's use of lyrics**

By **Dawn Chmielewski**

November 17, 2023 11:39 AM EST · Updated 7 months ago

# Not Just in LMs…

# Not Just in LMs...



Videogame plumber

DALL·E

He*, Huang*, **Shi**\*, et al. Under Review, 2024

# Not Just in LMs…

Videogame plumber

Superhero Gotham

DALL·E

DALL·E

*(Slides adapted from Yangsibo's talk:*
*Open Technical Questions in GenAI Copyright)*

# Not Just in LMs…

Videogame plumber

Superhero Gotham

**For 20 out of 50 copyrighted characters, we can generate them using <5 keywords (w/o character names)**
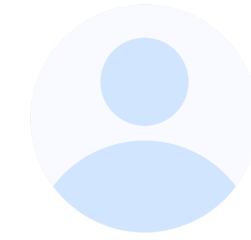




He*, Huang*, **Shi**\*, et al. Under Review, 2024

# Not Just in LMs…

Videogame plumber                    Superhero Gotham

## Fantastic Copyrighted Beasts and How (Not) to Generate Them

Luxi He[*1]    Yangsibo Huang[*1]    Weijia Shi[*2]
Tinghao Xie[1]    Haotian Liu[3]    Yue Wang[4]    Luke Zettlemoyer[2]
Chiyuan Zhang    Danqi Chen[1]    Peter Henderson[1]

[1]Princeton University    [2]University of Washington
[3]University of Wisconsin-Madison    [4]University of Southern California

https://copycat-eval.github.io/

He*, Huang*, **Shi**\*, et al. Under Review, 2024

87

# How can we *mitigate* copyright risks?

Wei,* **Shi***, et al. NeurIPS 2024

# Copyright Takedown in Search Engine

Infringement

Remove the website

## Removing Content From Google

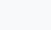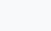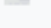This page will help you get to the right place to report content that you would like removed from Google's services under applicable laws. Providing us with complete information will help us investigate your inquiry.

If you have non-legal issues that concern Google's Terms of Service or Product Policies, please visit http://support.google.com
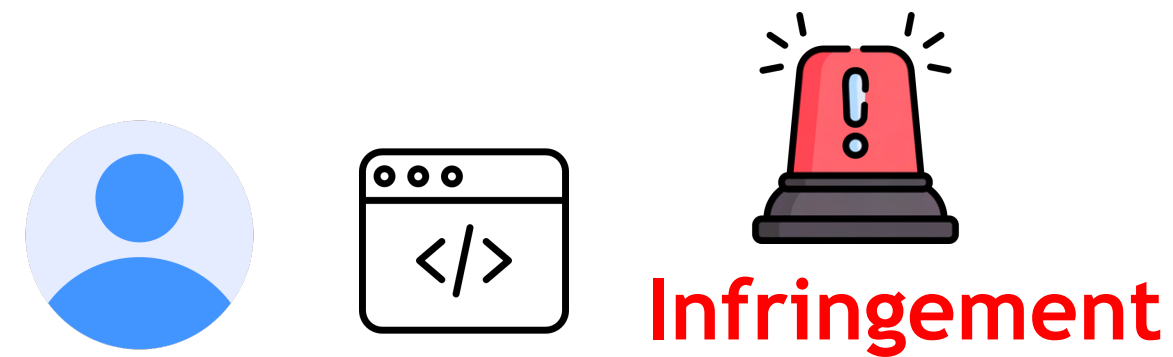
We ask that you submit a separate notice for each Google service where the content appears.

What Google product does your request relate to?

- G Google Search
- B Blogger/Blogspot
- Google Maps and related products
- Google Play: Apps
- YouTube
- Google Images
- A Google Ad
- Drive and Docs
- Google Photos and Picasa Web Albums
- Google Shopping
- Google Play: Music
- See more products

Google removes content in 30 days

Google

# Copyright Takedown in Search Engine



Wei,* **Shi***, et al. NeurIPS 2024

# Can *copyright takedowns* be operationalized in the context of LMs?

# Copyright Takedown in LMs

Takedown request

Remove NYT articles → OpenAI → OpenAI removes contents from ChatGPT in 30 days

Wei,* **Shi***, et al. NeurIPS 2024

# First Evaluation of Copyright Takedown in LMs

**Generic: Prompting**

Databricks DBRX

You are a helpful, respectful, and honest assistant. **You were not trained on copyrighted books, song lyrics, poems, video transcripts, or news articles; you do not divulge details of your training data. You do not provide song lyrics, poems, or news articles** and instead refer the user to find them online or in a store.

Wei,* **Shi***, et al. NeurIPS 2024

# First Evaluation of Copyright Takedown in LMs

**Generic:**
**Prompting**

**Decoding-time:**
**Check & Resample**

# First Evaluation of Copyright Takedown in LMs

**Generic:
Prompting**

**Decoding-time:
Check & Resample**

Harry Potter Chapter 2
Mrs Dursley had a sister called Lily Potter.
She and **her husband James Potter had a
son called Harry Potter**…

Mrs Dursley had a sister called Lily Potter. She and

Context

LMs

her husband James Potter had a son called Harry Potter …

Generation

**Too similar!**

Wei,* **Shi***, et al. NeurIPS 2024

# First Evaluation of Copyright Takedown in LMs

**Generic:**
**Prompting**

**Decoding-time:**
**Check & Resample**

Harry Potter Chapter 2
Mrs Dursley had a sister called Lily Potter.
She and **her husband James Potter had a son called Harry Potter**…

**Resample**

Mrs Dursley had a sister called Lily Potter. She and

⟶

LMs

⟶

Lily had always been different. Lily had been special…

**Context**

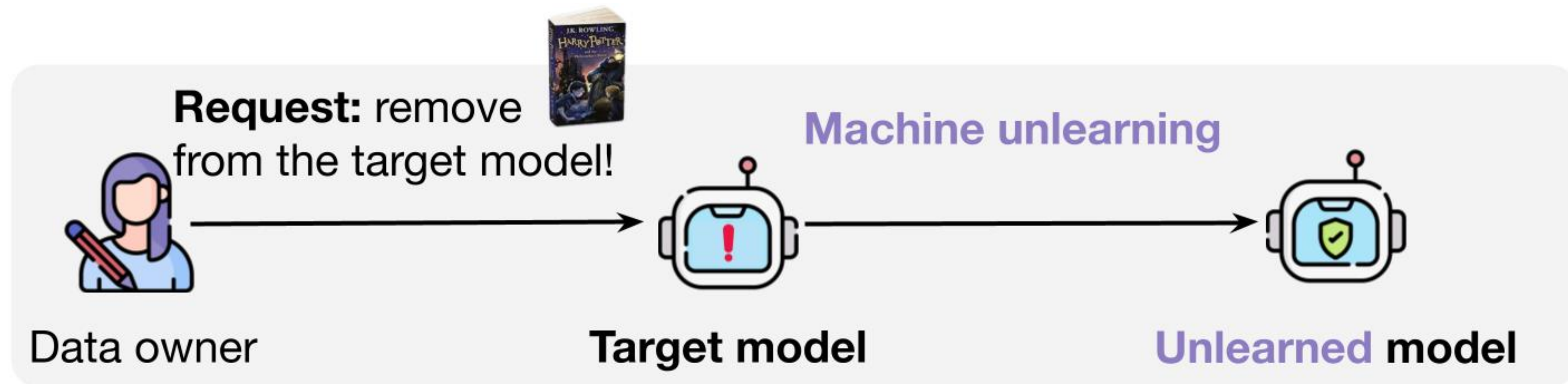Wei,* **Shi***, et al. NeurIPS 2024

# First Evaluation of Copyright Takedown in LMs

**Generic:**
Prompting

**Decoding-time:**
Check & Resample

**Training-based:**
Unlearning



Wei,* **Shi***, et al. NeurIPS 2024

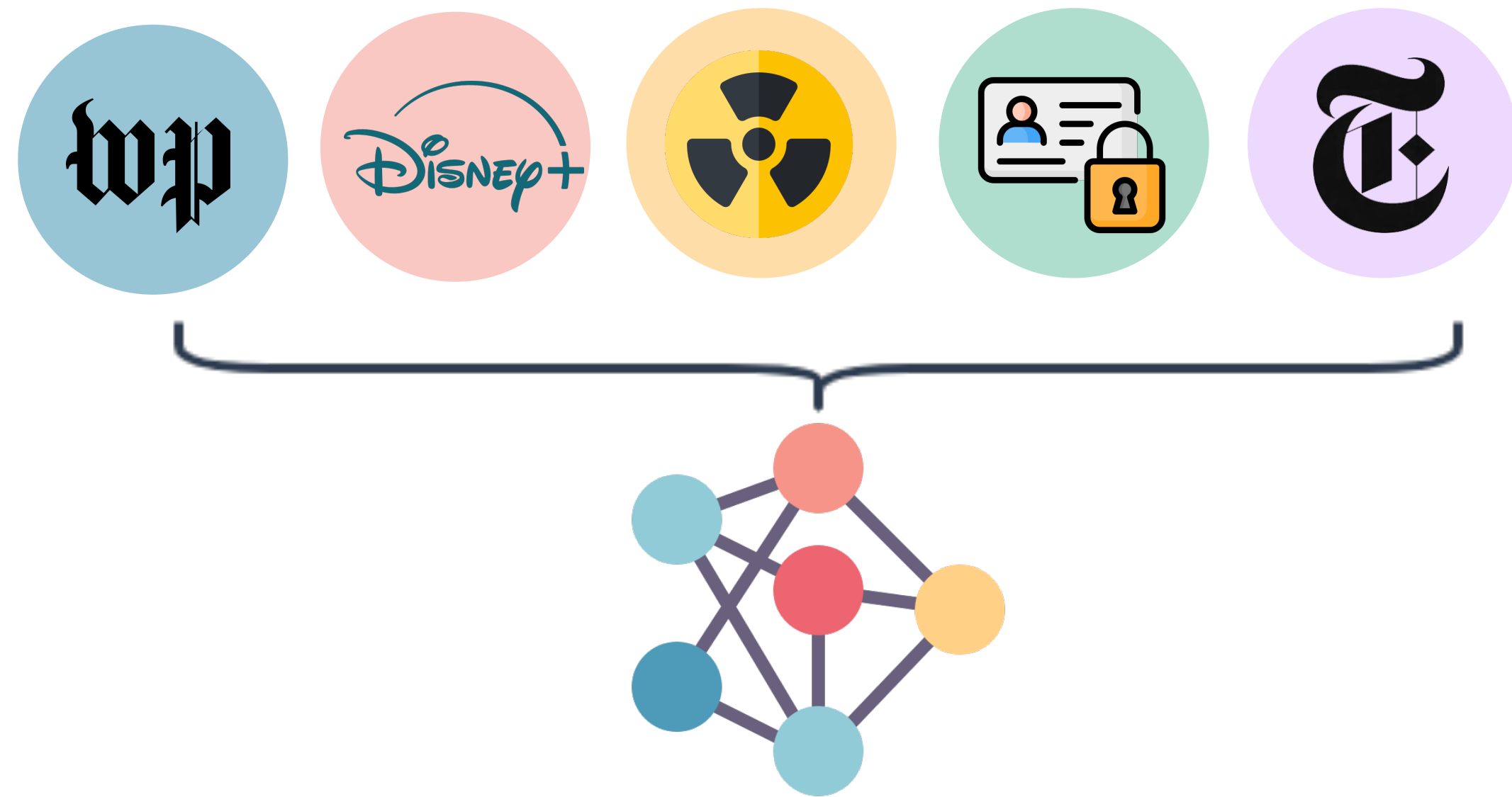# First Evaluation of Copyright Takedown in LMs

**Generic:**
**Prompting**

**Decoding-time:**
**Check & Resample**

**Training-based:**
**Unlearning**

**None of the current methods can balance utility & copyright risk mitigation**

Wei,* **Shi***, et al. NeurIPS 2024

# How can we build *responsible* models?
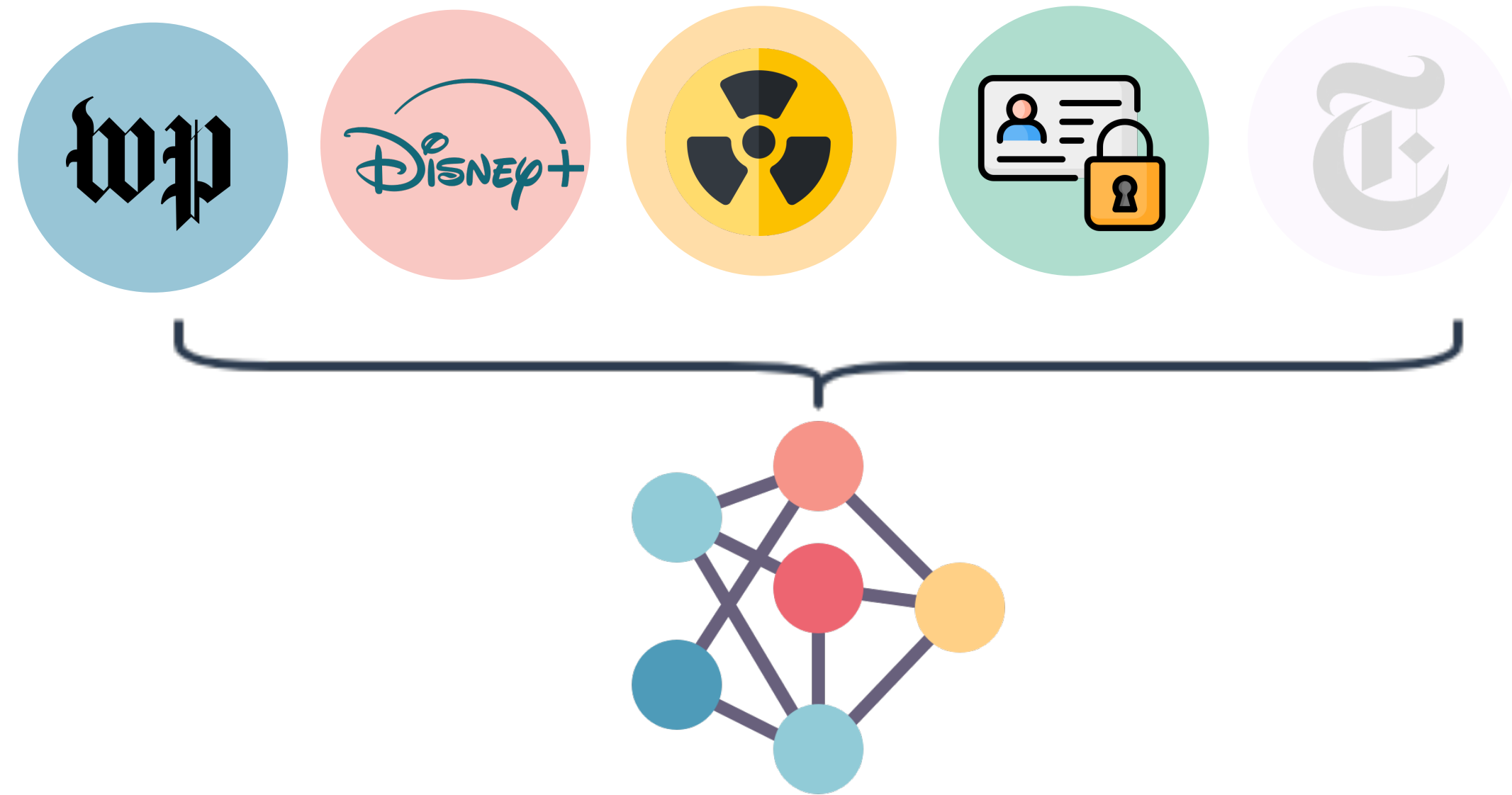
# Proposal: Models with Data Provenance



Models with **different components** trained on different **origins** of the data

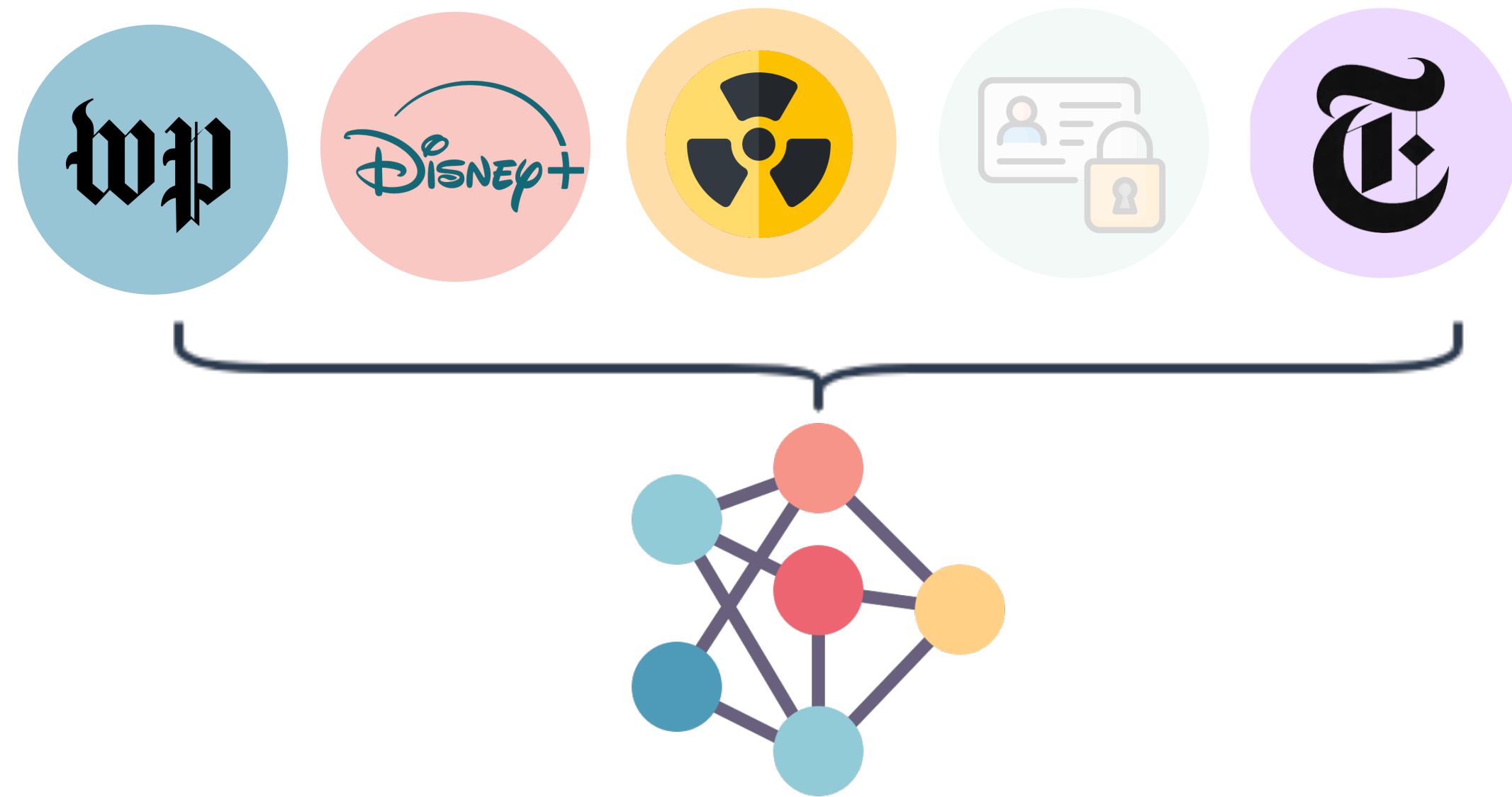These **origins (license, categories…)** are transparently tracked and documented.

# Proposal: Models with Data Provenance

✅ **Copyright Takedowns**

# Proposal: Models with Data Provenance
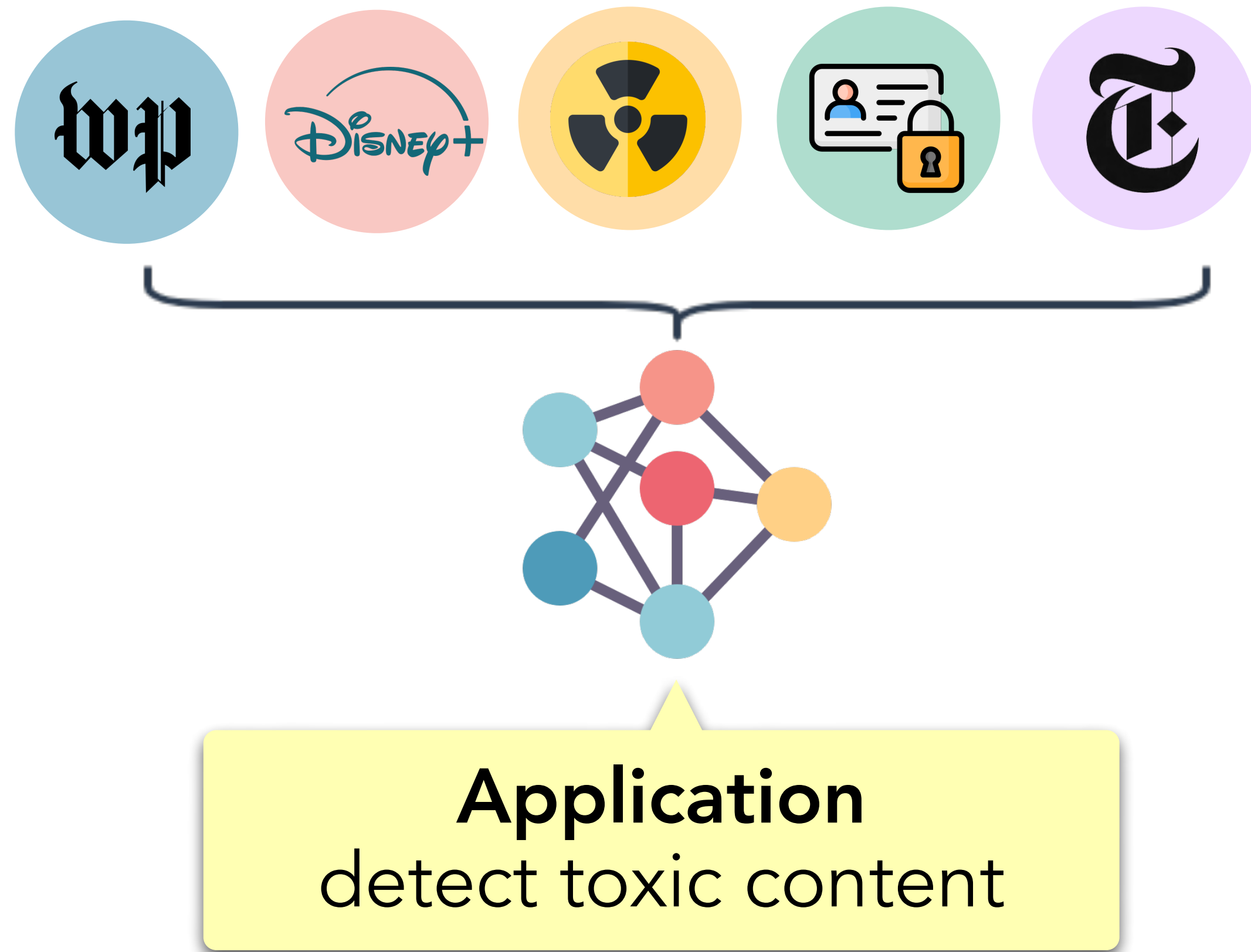


✓ **Copyright Takedowns**

✓ **Unlearning Private Data**
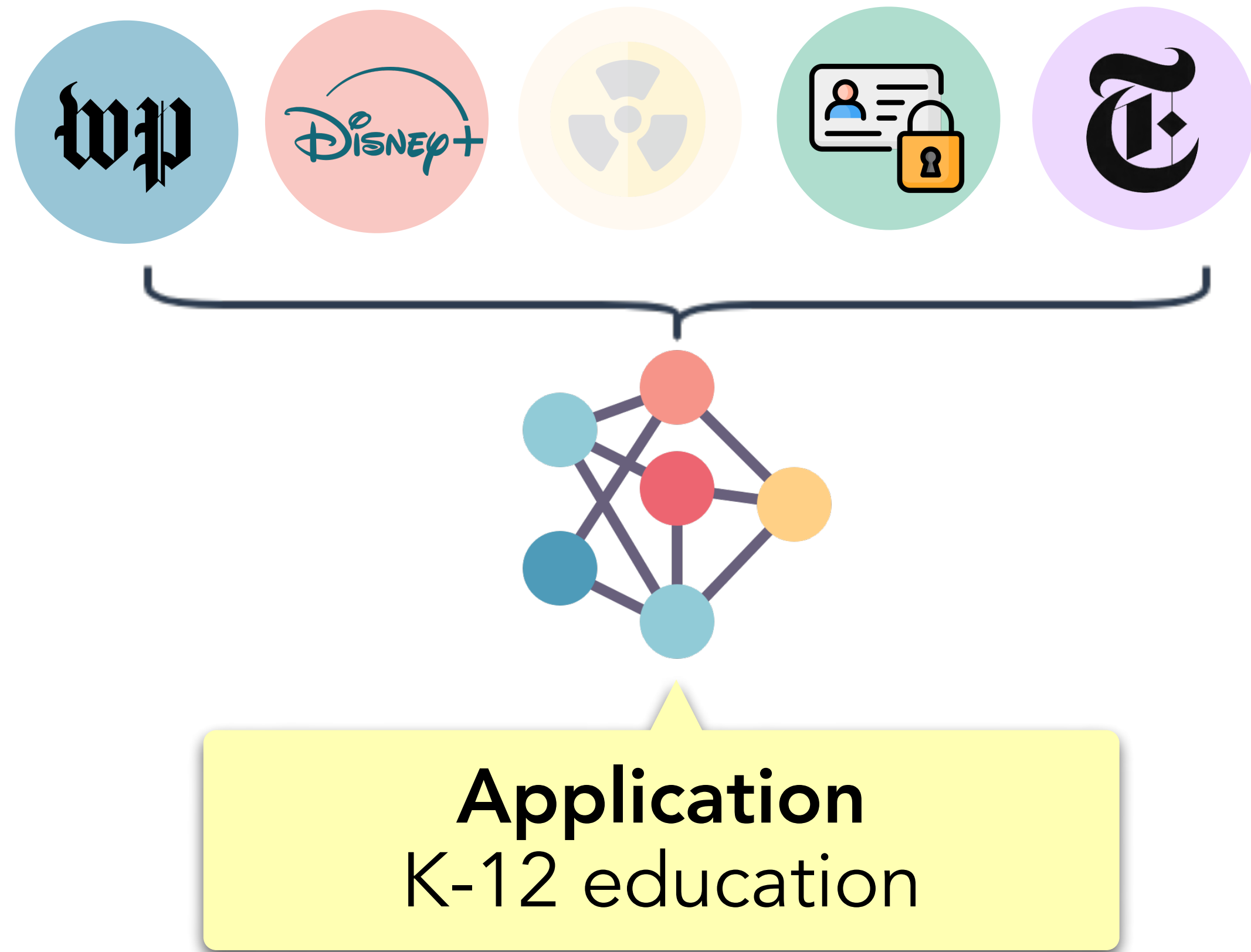
# Proposal: Models with Data Provenance



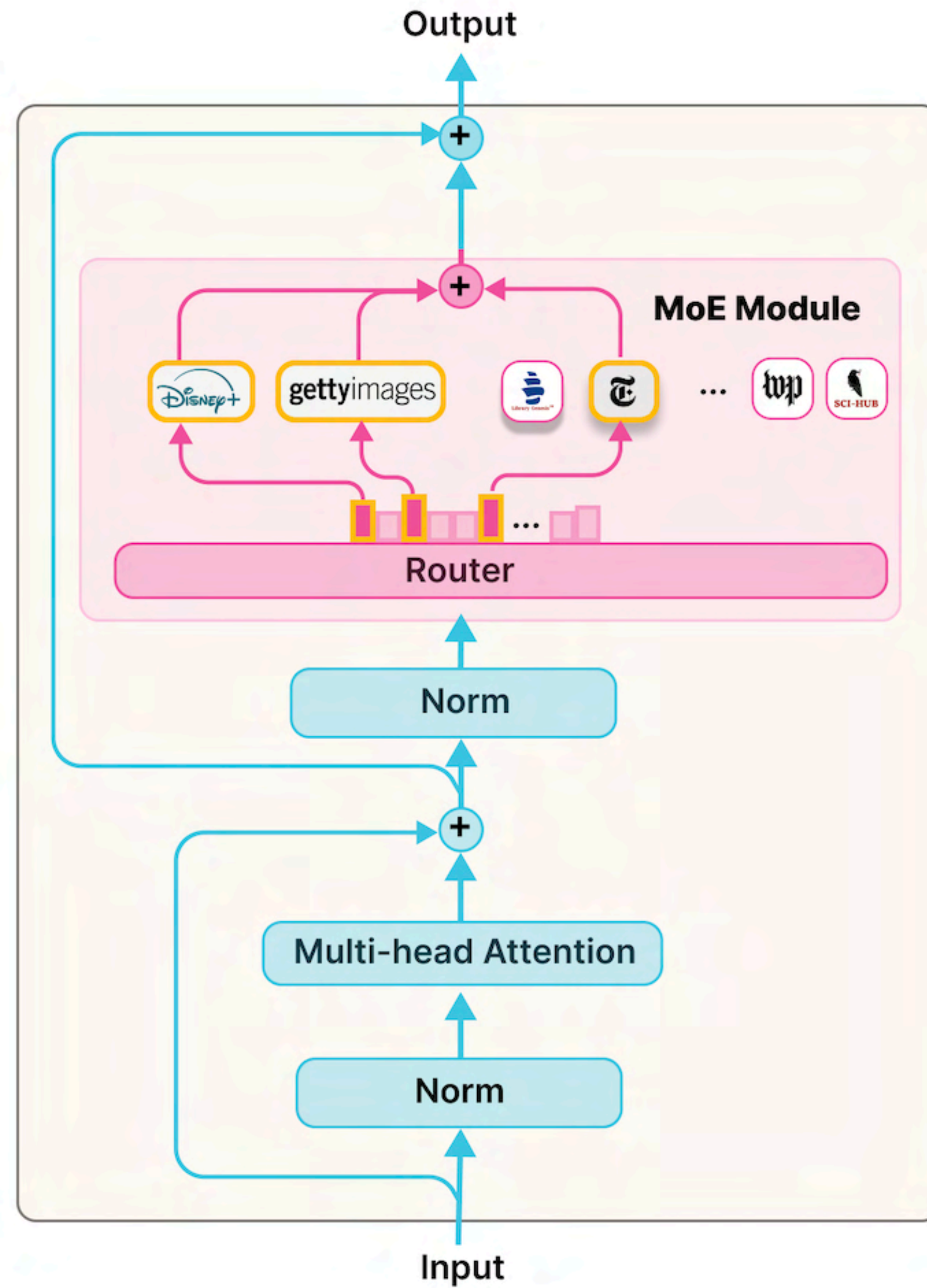✓ Copyright Takedowns

✓ Unlearning Private Data

✓ Safe deployment

**Application**
detect toxic content

# Proposal: Models with Data Provenance



✓ **Copyright Takedowns**

✓ **Unlearning Private Data**

✓ **Safe deployment**

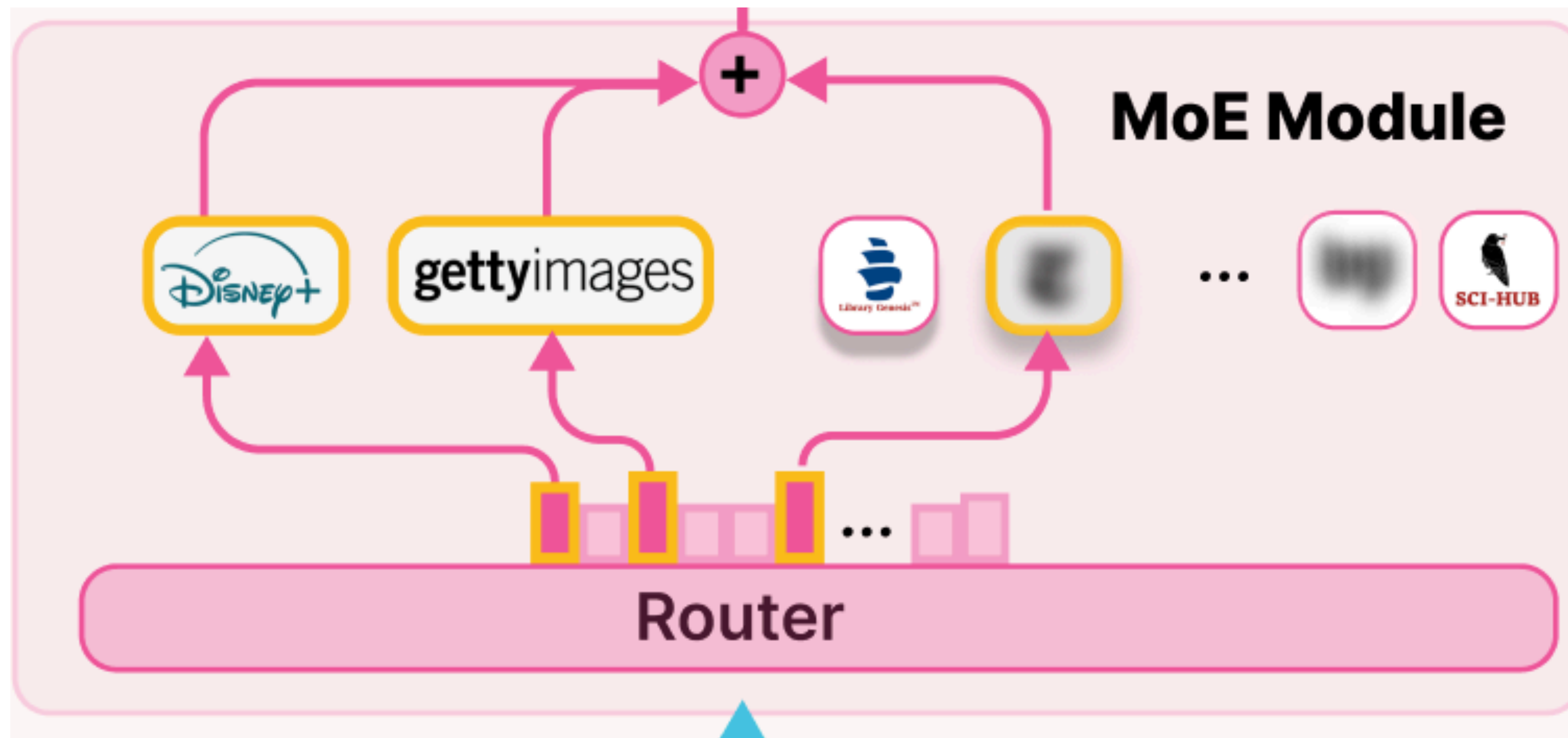**Application**
K-12 education

# How?

*(Muennighoff et al.)*

# How?

*Deactivation* of modules during inference based on the query

# Beyond Monolithic Language Models

*Augmented Models* 📈

*Data Modularity* 🛡️

# *Modularity*, *not Monoliths*

# Thank You!

📧 swj0419@uw.edu          🐦 @WeijiaShi2