

# Synthetic Bootstrapped Pretraining



*Zitong Yang*  
Stanford Statistics

# Collaborators



Aonan Zhang\*



Hong Liu



Tatsunori Hashimoto



Emmanuel Candès



Chong Wang

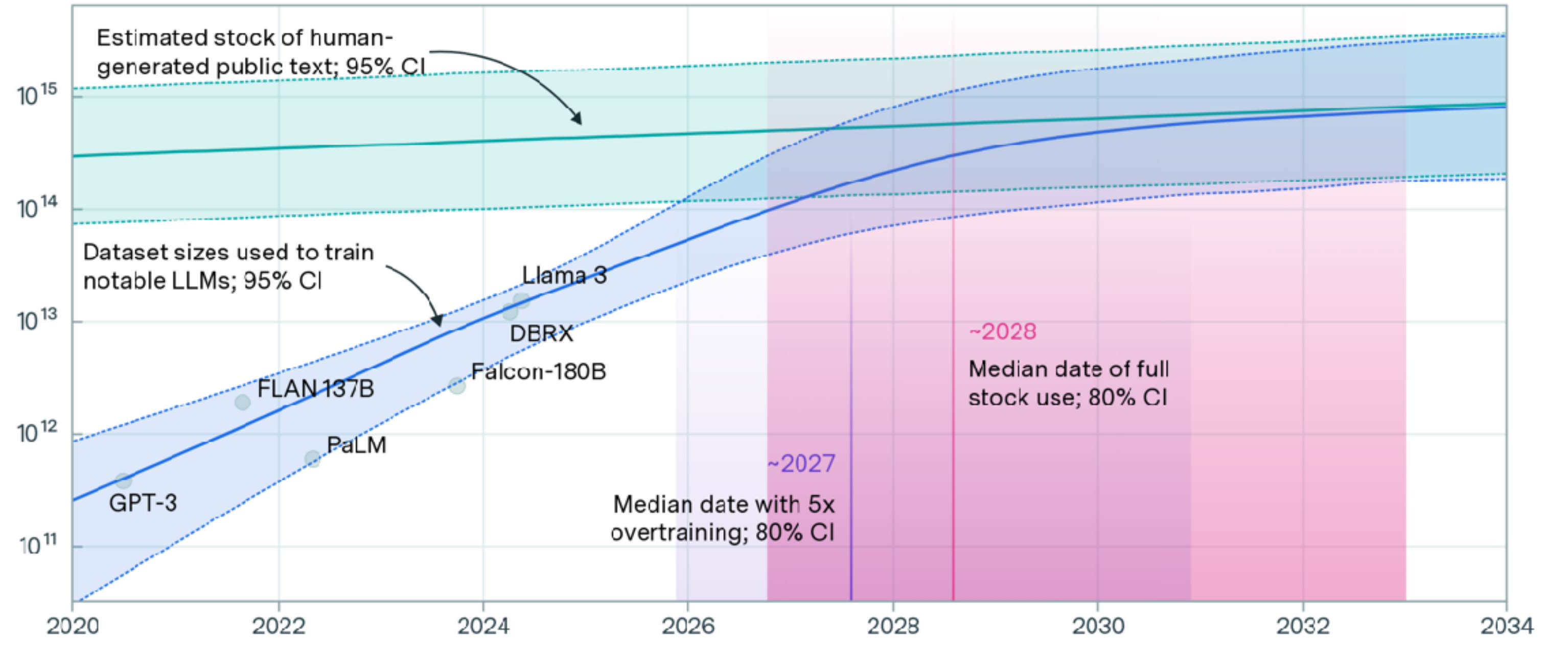
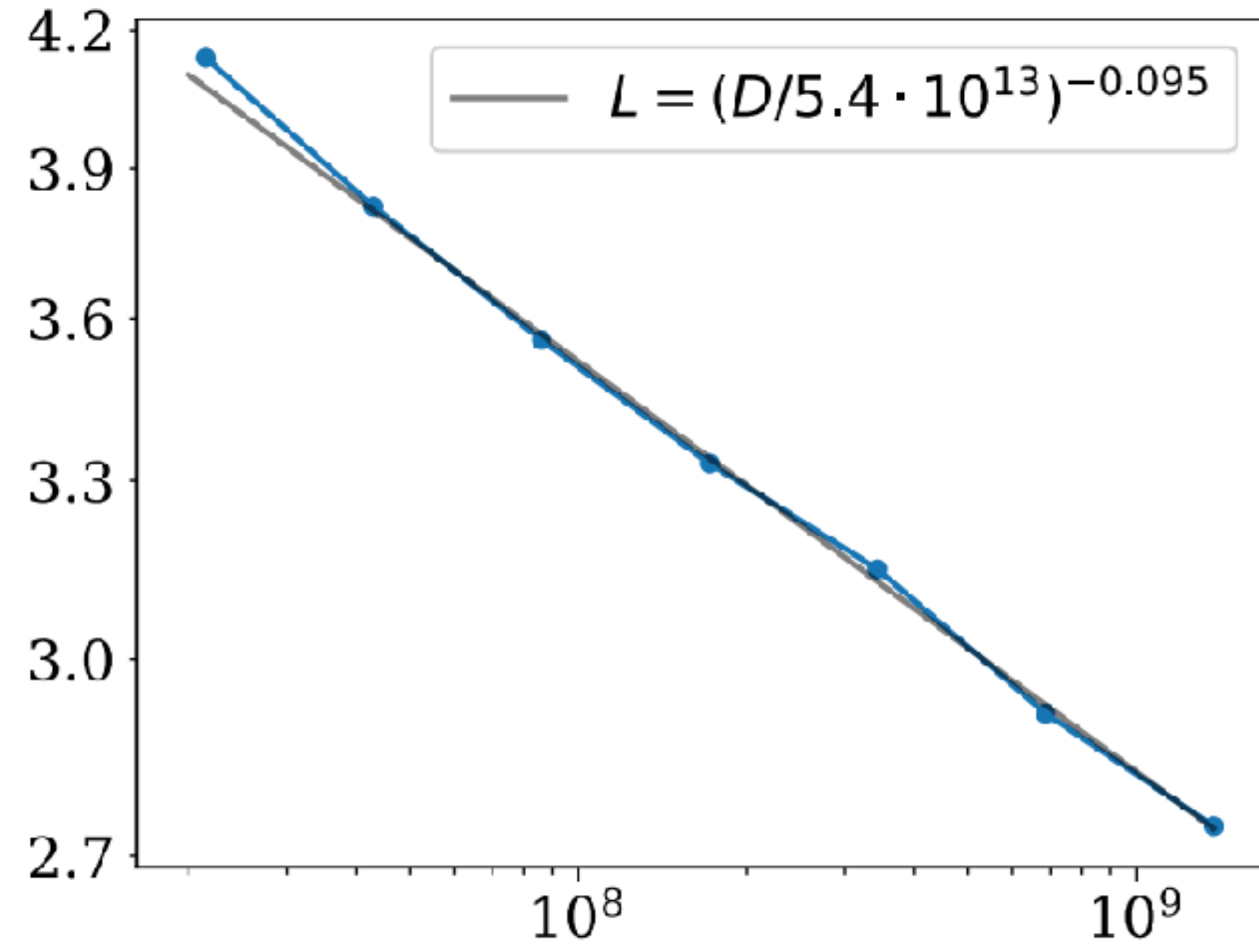


Ruoming Pang

*\*Equal contribution*

<https://arxiv.org/pdf/2509.15248>

# Scaling wall



**Mindset:** use existing data more effectively

# Re-examining pretraining

## Pros

- ▶ Data scalability: crawling is more effective than data collection
- ▶ System scalability: sequence level, batch level, model level parallelizable
- ▶ Simplicity: maximum likelihood estimation from Sir Ronald A. Fisher

## Cons

- ▶ Data efficiency: amount of text processed by 13 y.o. is 100M tokens

**Belief:** advantage far outweigh limitation. We should preserve it as much as we can

# Where does the knowledge come from in pretraining?

## Thought experiment

- ▶ World with 5 tokens “A”, “B”, “C”, “D”, “E”
- ▶ Text documents with each token sampled u.a.r. [“BDECD...”, “ACEAC...”, ]
- ▶ Perform next token prediction with transformer LM: No meaningful learning signal

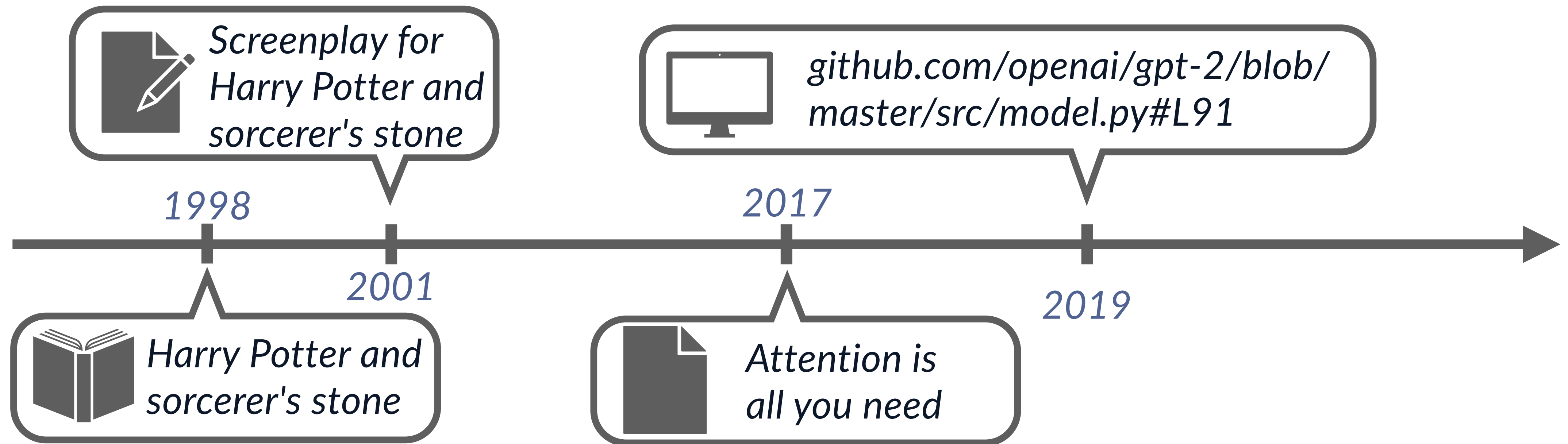
## Two views

- ▶ Statistical view: Natural language tokens has *statistical* correlations
- ▶ Computational view: Natural language has *compressible* patterns

Views aside: **structural correlation** enables learning

# Under-exploited sources of correlation

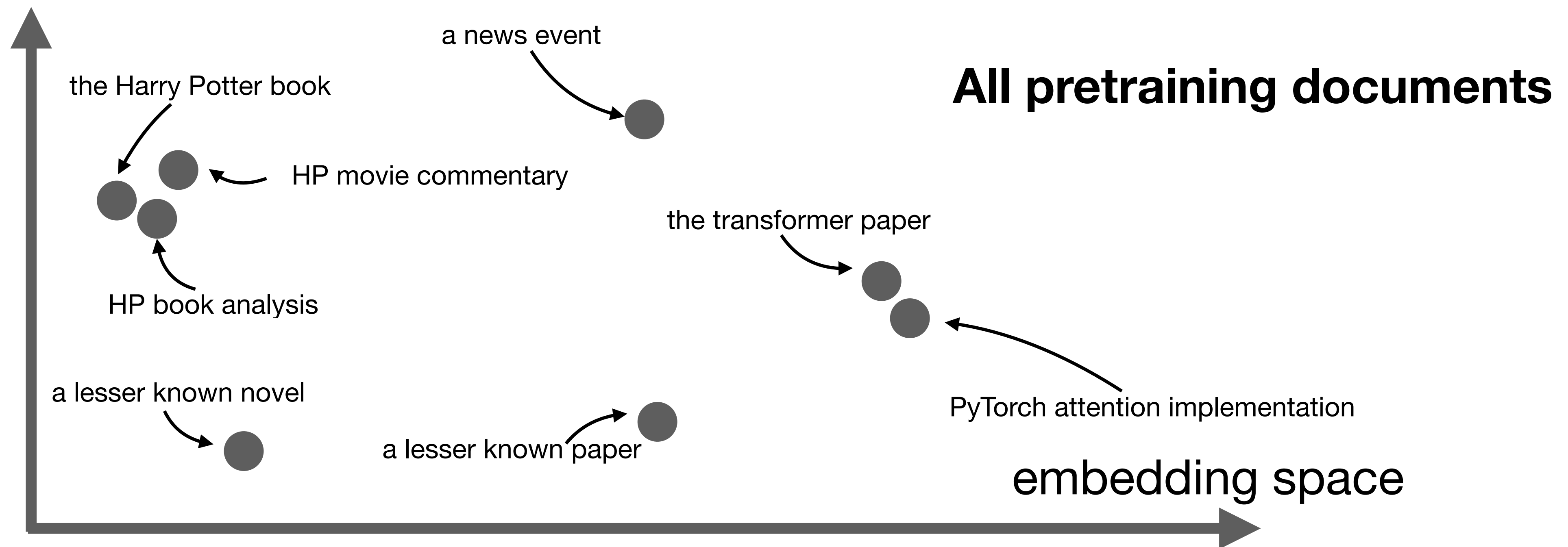
There exists *rich* correlations between documents



**Technique:** take-advantage of inter-document correlation

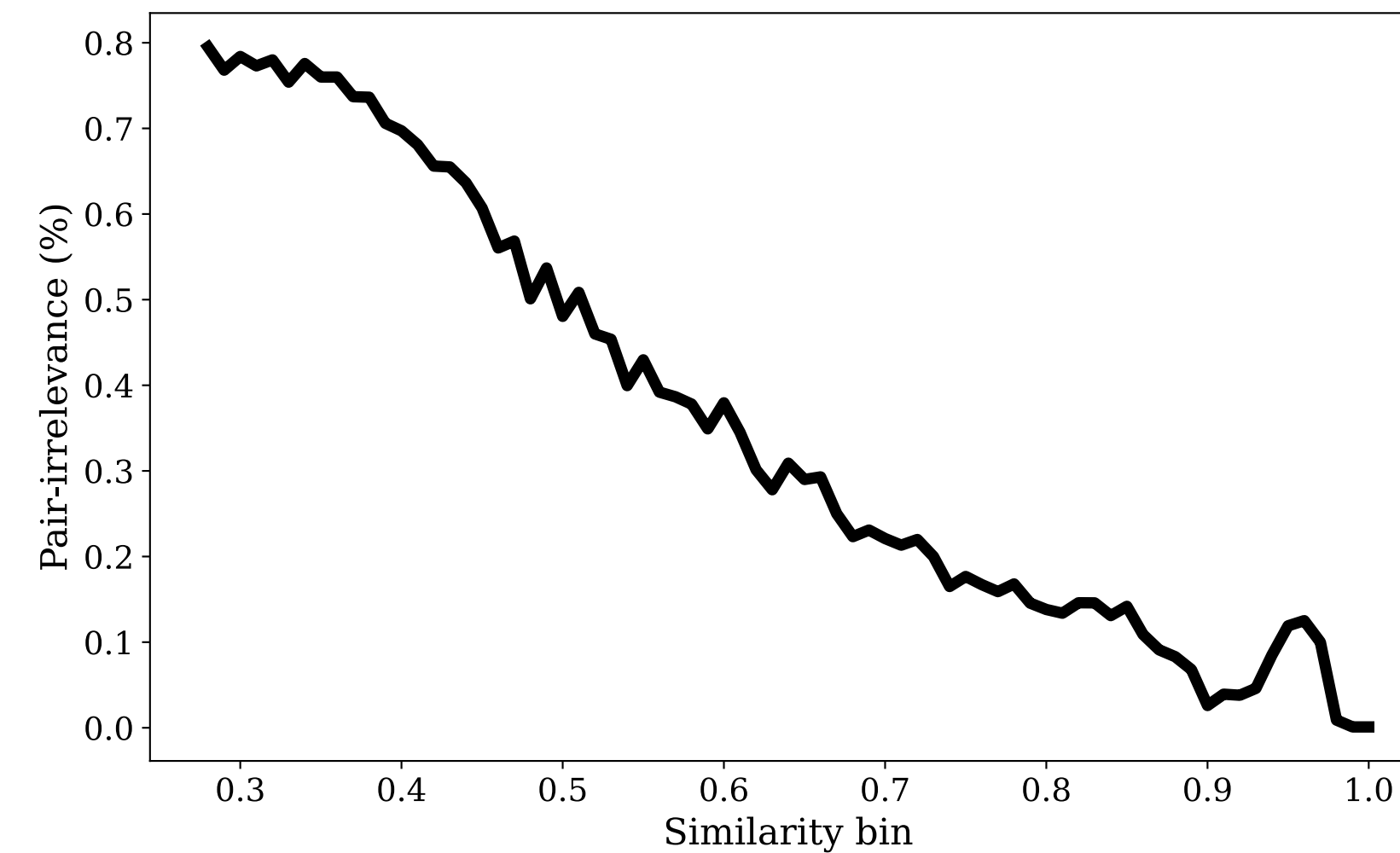
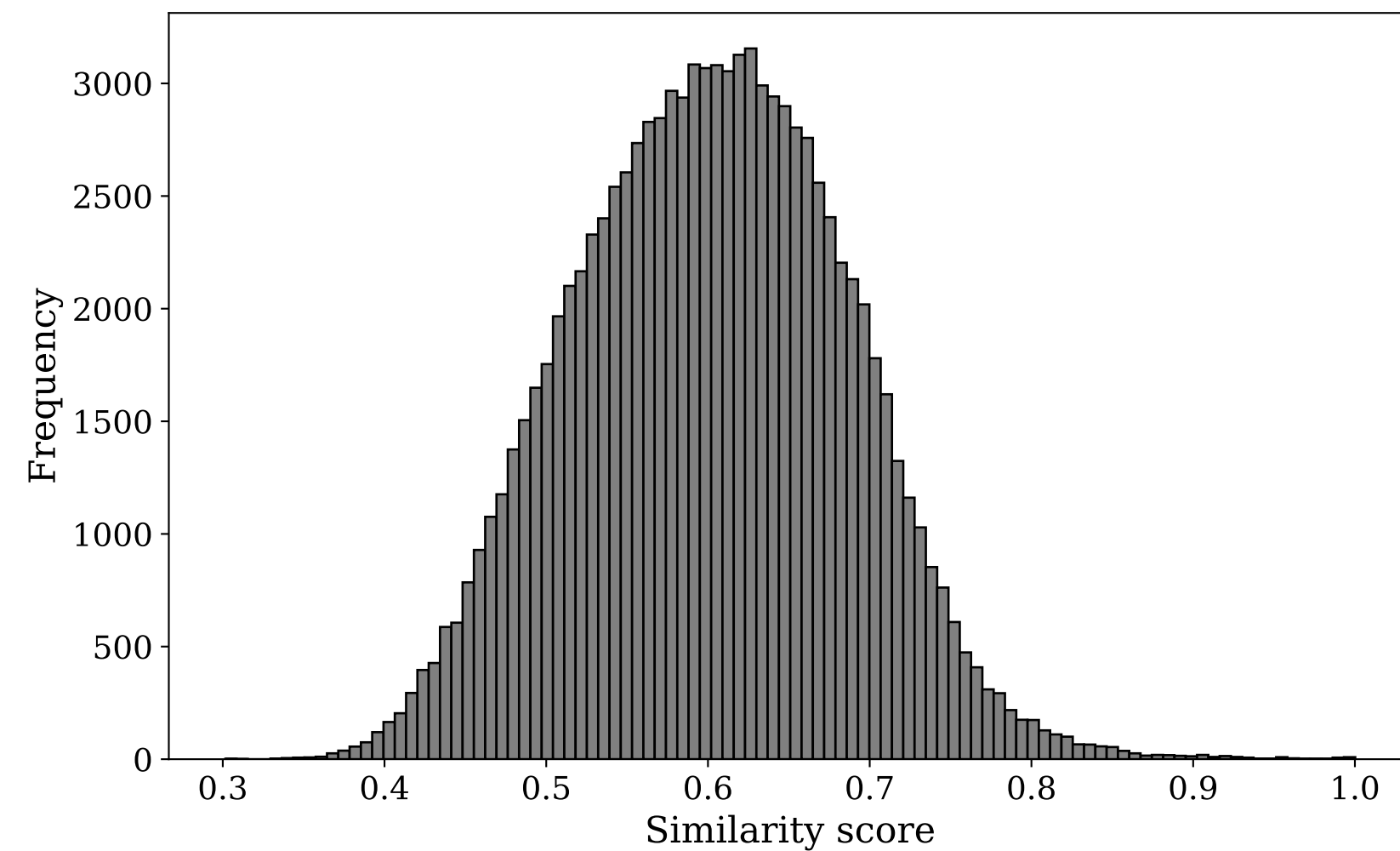
# Synthetic bootstrapped pretraining

1. **Nearest-neighbor pairing:** we used DCLM subset and Qwen-0.6B-Embedding



# Building a dataset of related documents

- ▶ Curate a dataset of document pairs when their similarity score is less than 0.75



- ▶ Similarity score reflects document relevance as judged by an LM

# Examples of related documents look like

doc1

## The Cultural Sites of Iran

With 196 countries and countless exciting destinations worldwide, there is so much to see in a very limited time. Even the most well-traveled person hardly gets to visit all and has to be selective. So, why should you consider visiting a country like Iran, especially when it comes to all those negative news and stereotypes surrounding it?

Here we're here to give you the reasons and to help you overcome your doubts and even encourage you to consider your next trip to Iran, this mysterious land as soon as you return to your home country!

Beautiful cities, friendly people, fabulous food, glorious architecture, Iran has delighted visitors for centuries with its World Heritage Sites, friendly towns and inspiring desert landscapes.

## Things to Do in Iran – Activities & Attractions

Iran is the land of four seasons, history and culture, souvenir and authenticity. This is not a tourism slogan, this is the reality inferred from the experience of visitors who have been impressed by Iran's beauties and amazing attractions.

Antiquity and richness of the Cultural Sites of Iran and civilization, the variety of natural and geographical attractions, four – season climate,  
...

History of Iran

doc2

Query Text: Home > FAQ Login / Register

Why should we spend our holiday in Iran?

Iran is a country, located in the Middle East, which can meet the various needs of tourists and satisfy their different tastes, due to its rich civilization, historical sites, geographic location, nature of the four seasons and diverse tourist attractions. Therefore, considering the high security and low cost of travel to the country, it is introduced as one of the major tourist destinations to spend holidays in.

Is Iran a safe travel destination?

One of the wrong assumptions about the country of Iran is in terms of its security. Despite its location in Asia and the Middle East, and neighboring countries like Iraq, Afghanistan and Pakistan, Iran is considered as one of the safest countries in the region. According to the international data, security in Iran is much more than a touristic country such as Turkey.

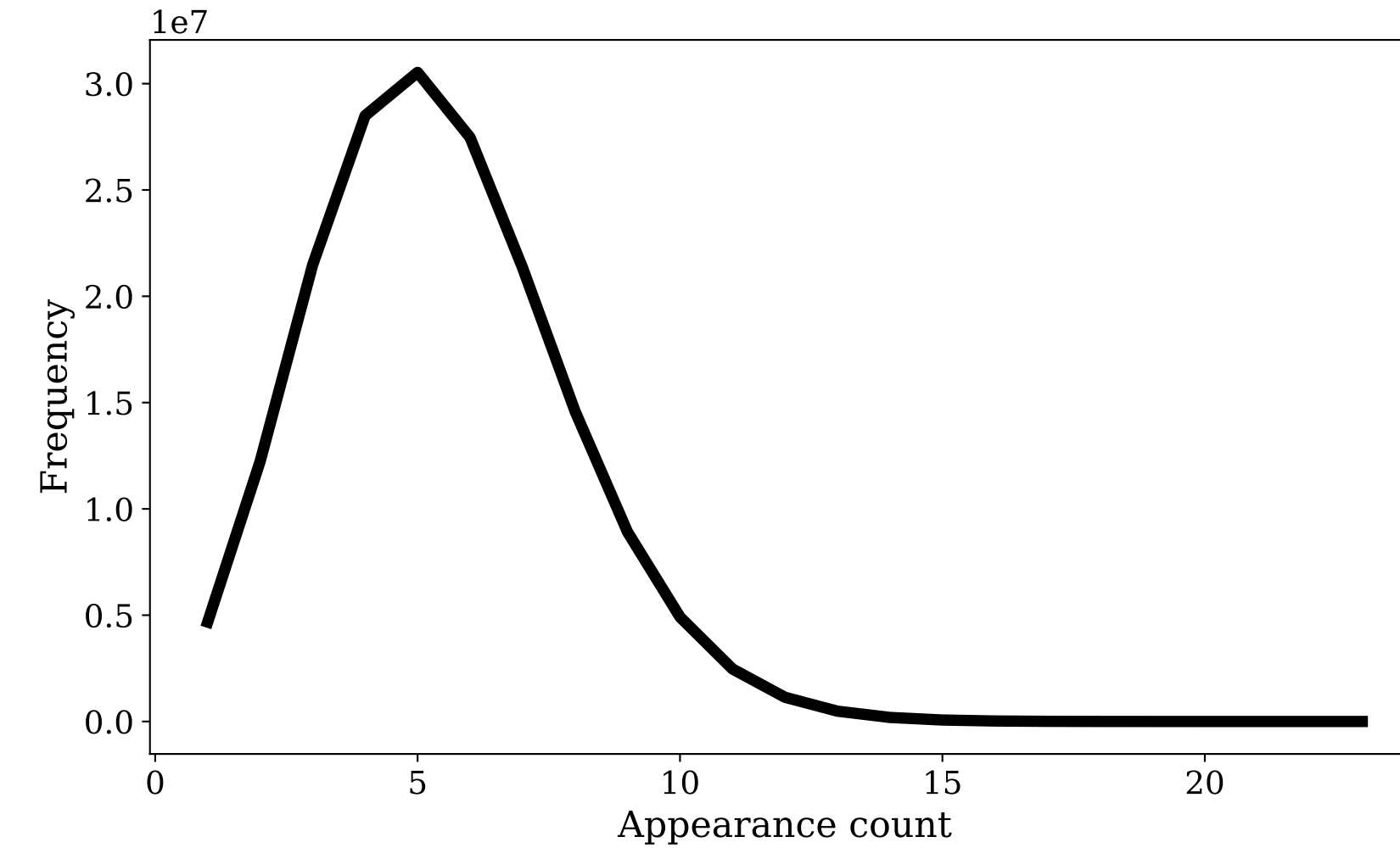
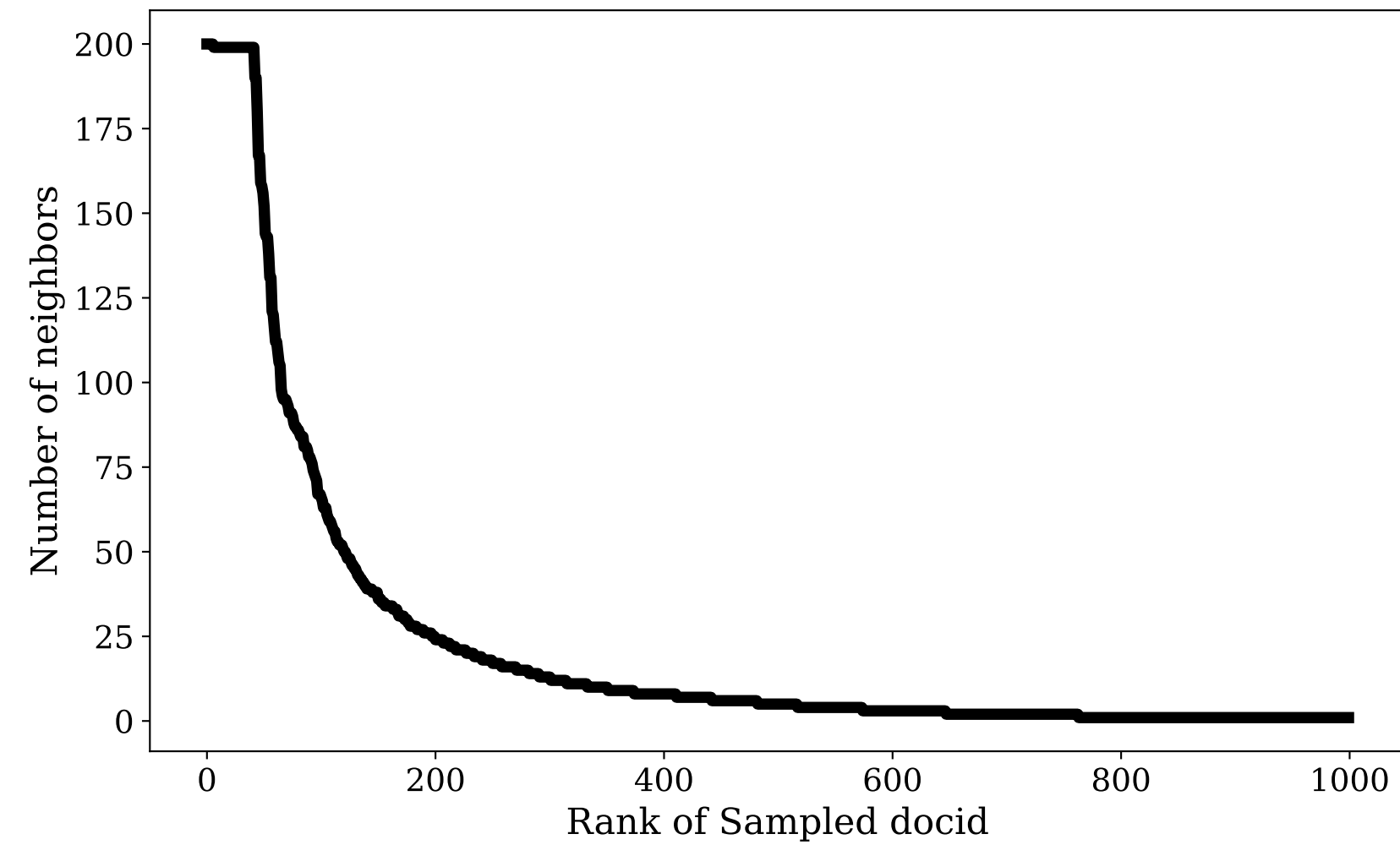
To confirm the statements made above, refer to websites like [www.travelriskmap.com](http://www.travelriskmap.com).

What does "the rich civilization" mean, as mentioned about Iran?

According to documentation in some of the world history references,  
...

Travel guide to Iran

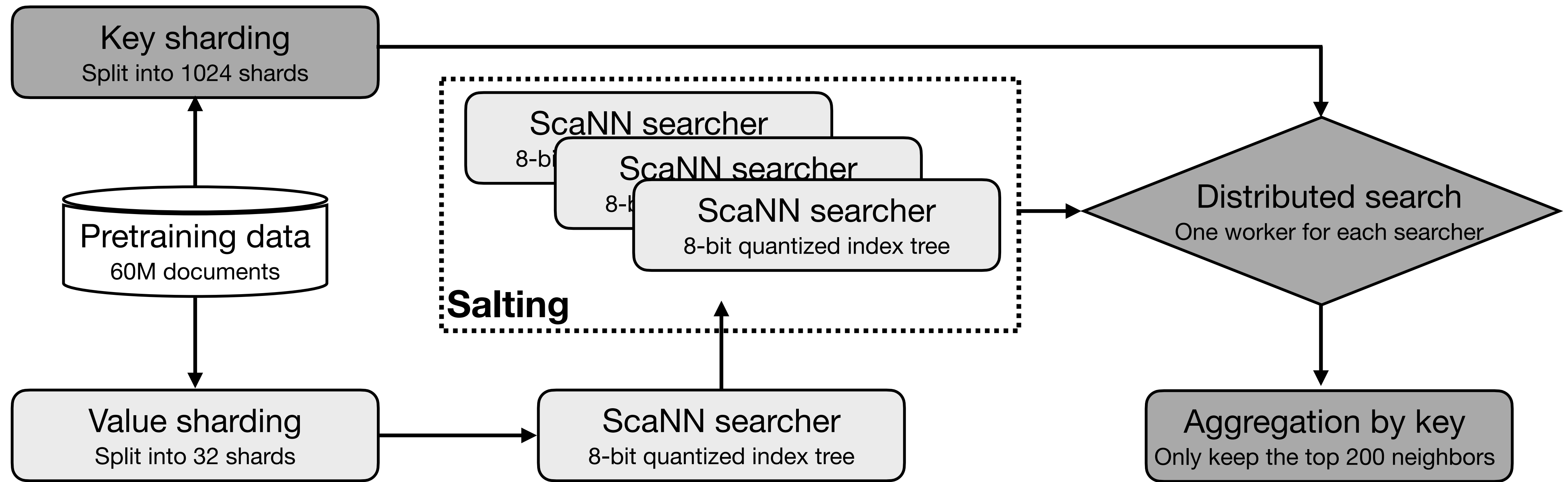
# Statistics for the nearest neighbor graph



- ▶ Documents sorted by *#neighbors* decays with a heavy tail
- ▶ Most of documents have 5 neighbors. Overall follows the shape of a Poisson distribution

# System design for nearest neighbor search (for experts)

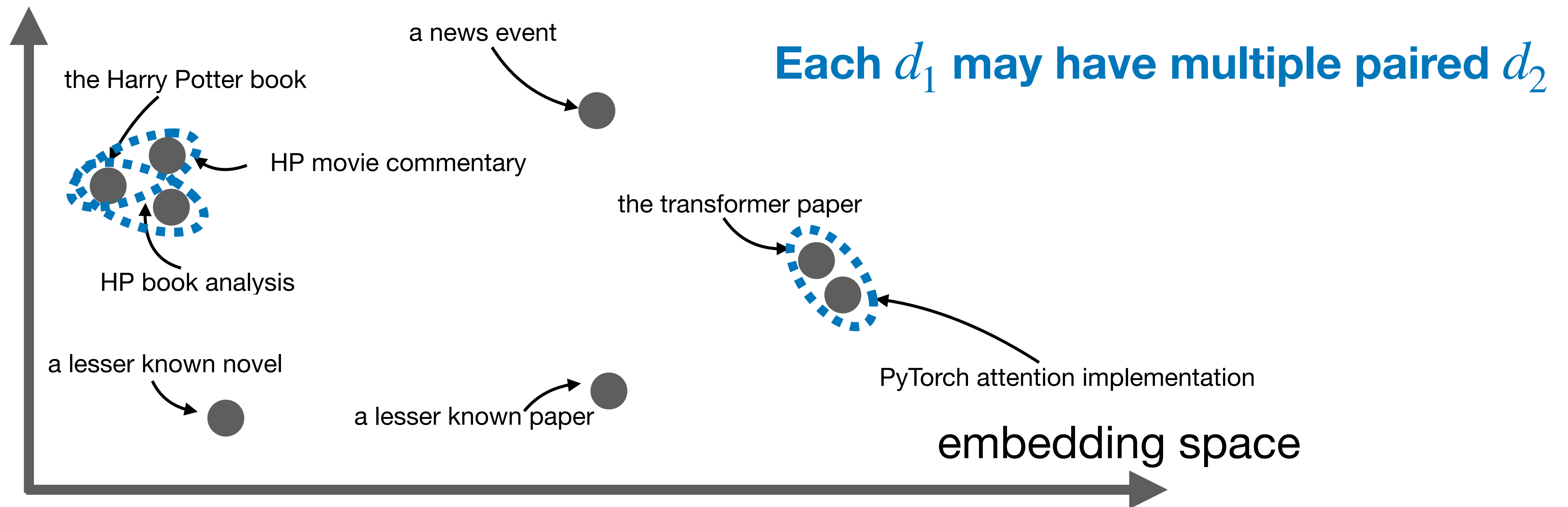
- ▶ Asymmetric key-value sharing; 8-bit quantization; Salted key for load balancing



- ▶ 155M CPU hours over 60M documents with heavy memory and utilization optimization

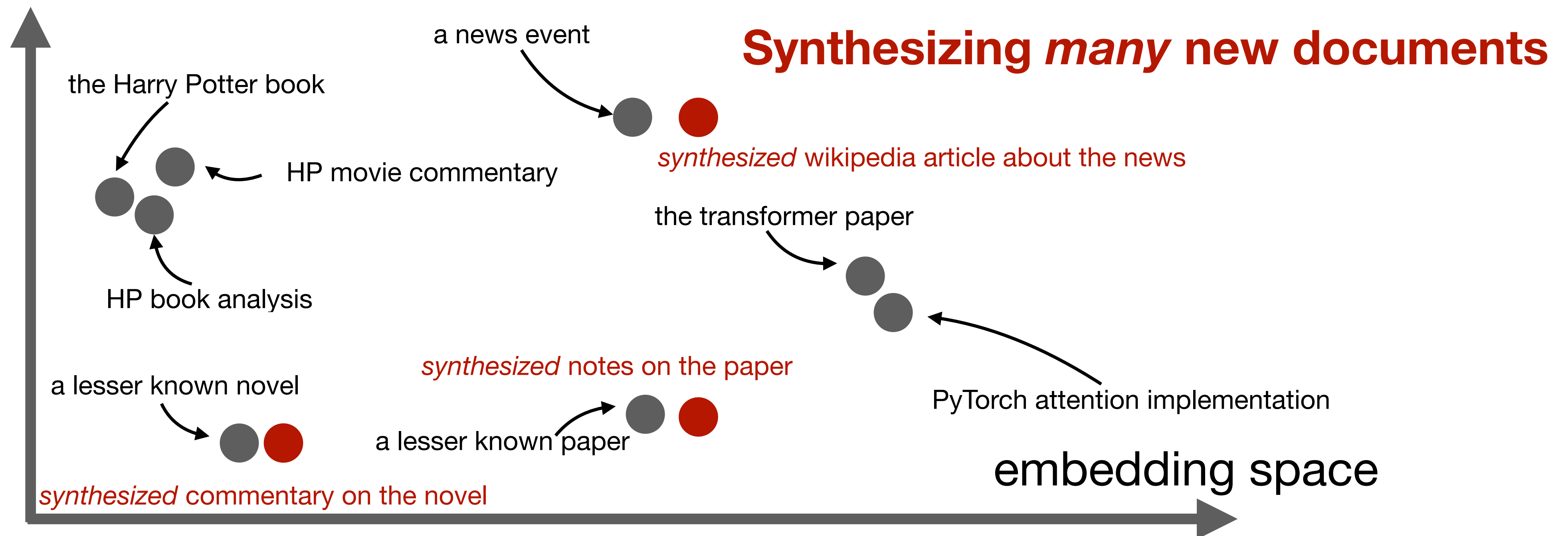
# Synthetic bootstrapped pretraining

1. **Nearest-neighbor pairing:** we used DCLM subset and Qwen-0.6B-Embedding
2. **Synthesizer tuning:** SFT-like objective  $p_{\theta}(d_1 | d_2)$  initialized at pretrained checkpoint



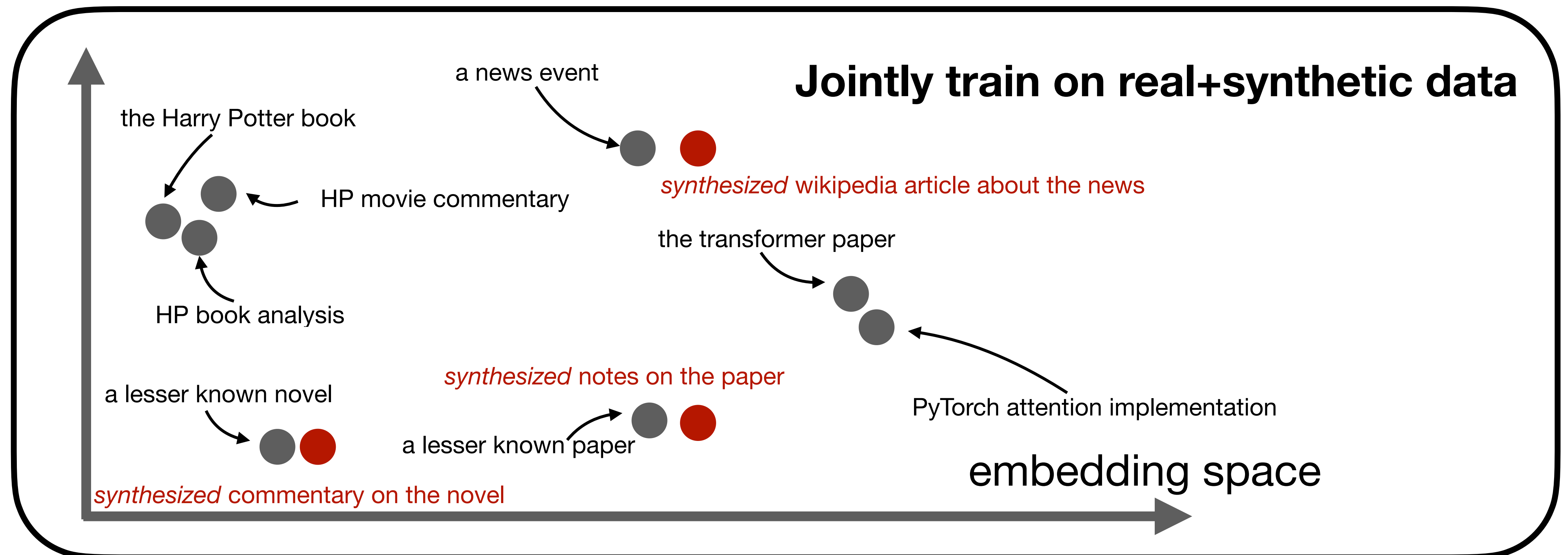
# Synthetic bootstrapped pretraining

1. **Nearest-neighbor pairing:** we used DCLM subset and Qwen-0.6B-Embedding
2. **Synthesizer tuning:** SFT-like objective  $p_{\theta}(d_1 | d_2)$  initialized at pretrained checkpoint
3. **Synthesis at scale:** Temperature=1 allows each document to have varied synthesis



# Synthetic bootstrapped pretraining

1. **Nearest-neighbor pairing:** we used DCLM subset and Qwen-0.6B-Embedding
2. **Synthesizer tuning:** SFT-like objective  $p_{\theta}(d_1 | d_2)$  initialized at pretrained checkpoint
3. **Synthesis at scale:** Temperature=1 allows each document to have varied synthesis



# Visualization of synthesized document

## Real document

HOME > Chowhound > California >

### Coffeehouses in SD

I am always looking for good coffeehouses in San Diego. I am looking for coffeehouses with (obviously) excellent coffeebeans, baristas who know their business (good microfoam, good crema, they should know the different ratios for steamed milk, foamed milk, espresso for the cappuccino, latte etc., etc), good snacks (paninis, cakes etc), good (cold) flavored coffeedrinks and good ambience. So far I haven't found many who are good but would like to hear more suggestion.

Caffe Calabria: By far the best coffeehouse in SD. Excellent coffee and baristas (best cappuccino in SD,), great paninis. Only downside is that they close very early.

Café Lestat: Good example what happens with good coffee and average baristas. Lestat is also using coffeebeans from Caffé Calabria but it is sad to see that they don't know how to make excellent espresso, cappuccino.

E Street Café: Average espresso but strange flavored ...

## Synthesis I

Coffee = Love

On a recent trip to San Diego, one thing I wanted to find was coffee. I was worried it might be difficult with the huge decrease in Asian flavors, but I was pleasantly surprised. Many of the Mexican, Italian, and Cuban flavors that are pretty common in California are now in the mix. While I did find a few highly recommended spots, the majority of the locations I checked out were nothing to write home about. At the least, they were not being sold out, which I consider a plus.

One major factor that really needs to be taken into account is the type of machine being used. There are two types of espresso machines: the cheap ones that spit out lots of pressure and don't have any control over the shot, and the expensive ones that create amazing shots and utilize the types of superfine grinds you usually find in Italy. The latter types are much harder to find and cost significantly more.

Another factor to consider is the type of beans being used. If they are done correctly, great espresso can be made from any type of bean, but it's a rare experience to find a place that does ...

## Synthesis II

Discuss as:

Coffeehouses: San Diego. The best?

Café Lestat, by far, is a stand-out. (photo by Thomas Moriarty)

If the coffee scene in San Diego sounds like the one in New York, there's a reason. The upper echelons of NYC's drink scene are the same ones who have taken over the majority of the coffee houses in San Diego. This is probably partly due to the fact that both cities have large student populations, both are home to large coffee drinker bases, both have a large number of espresso shops, and both offer top-notch, hand-crafted cappuccinos. But if there is one city in America that deserves to have its name above that of New York, it's San Diego.

There are just under 100 coffee shops in San Diego, with almost half of them located on University Ave. alone. So finding the perfect coffee shop is crucial. We spent a whole day just roaming around the area, hunting for the best.

In terms of the coffee itself, it's hard to beat Café Lestat. The baristas are amazing and their methods are pristine ...

# Experiment design

## Data, model, and evaluation

- ▶ Data: de-duplicated DCLM dataset
- ▶ Model: Llama 3 architecture with additional QK-norm
- ▶ Evaluation: 6 QA accuracies and 3 perplexity evaluation

## Compute-matched comparison

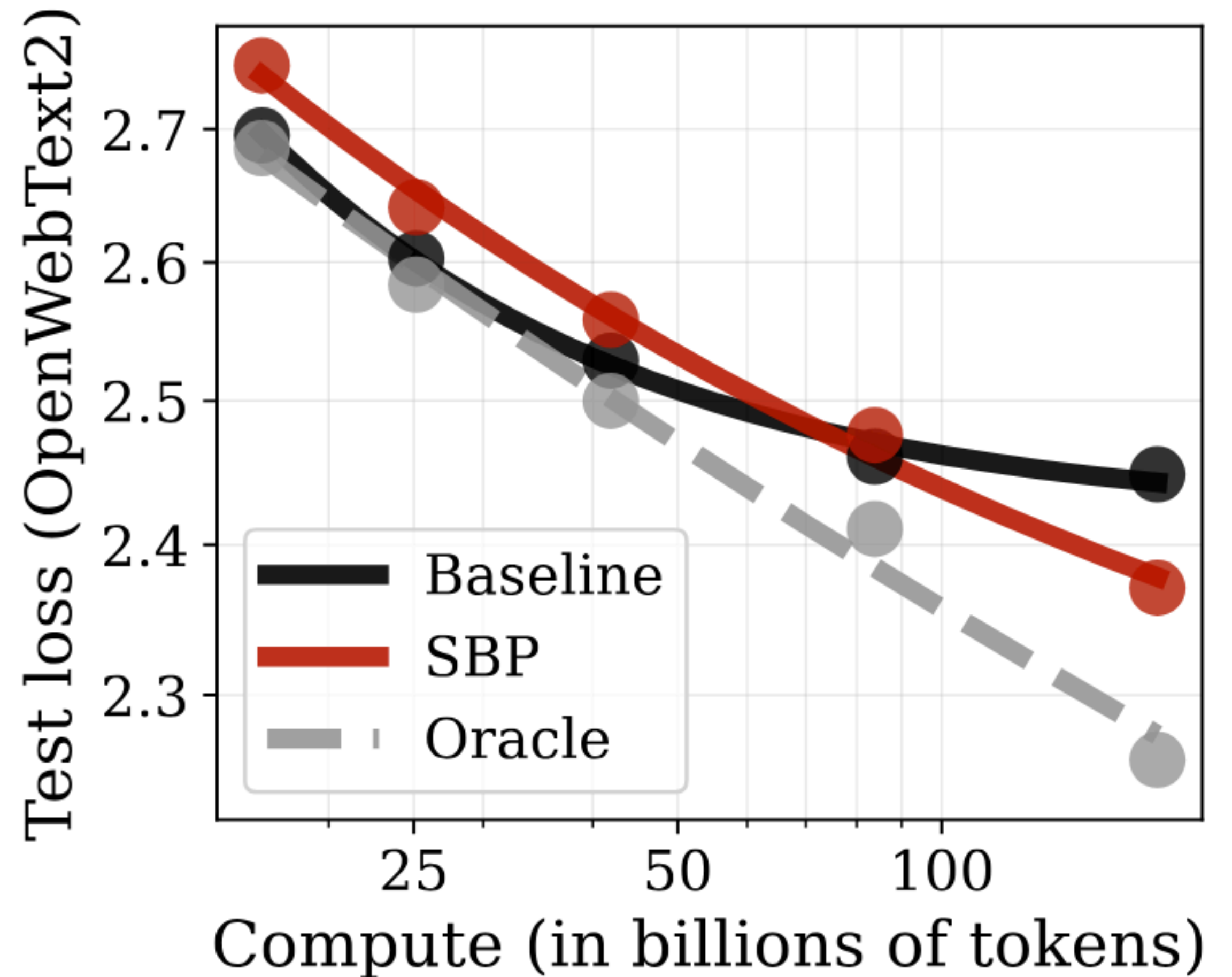
- ▶ Baseline: 20 times repetition
- ▶ SBP: Same compute, same data source
- ▶ Oracle: 20x additional data, no repetition during pretraining

# Result: 40% improvement attained by the oracle

Benchmark	200B-scale			1T-scale		
	Baseline	SBP	Oracle	Baseline	SBP	Oracle
<i>Perplexity on held-out data ↓</i>						
OpenWebText2	5.74	-0.53	-1.02	4.51	-0.02	-0.12
LAMBADA	6.87	-0.85	-1.86	4.33	-0.03	-0.22
Five-shot MMLU	3.83	-0.36	-0.51	3.17	-0.06	-0.05
<i>QA accuracy ↑</i>						
ARC-Challenge (0-shot)	35.32	+1.28	+2.82	42.66	+1.62	+3.84
ARC-Easy (0-shot)	68.94	+2.65	+4.29	75.63	+0.42	+2.11
SciQ (0-shot)	90.50	+1.00	+2.40	93.20	+0.80	+0.50
Winogrande (0-shot)	60.14	+1.90	+5.53	65.19	+1.42	+2.92
TriviaQA (1-shot)	22.51	+3.36	+7.37	36.07	+0.25	+0.59
WebQS (1-shot)	8.56	+3.74	+10.83	19.34	+0.54	+0.44
<b>Average QA accuracy</b>	<b>47.66</b>	<b>+2.32</b>	<b>+5.54</b>	<b>55.35</b>	<b>+0.84</b>	<b>+1.73</b>

# Training dynamics

- ▶ At beginning, baseline and oracle performs similarly. SBP is worse than them because it uses synthetic data
- ▶ Later on, baseline and oracle begins to diverge but SBP still follows a linear trend
- ▶ Towards the end, baseline plateaus but SBP continues to decrease



# Synthetic data quality

	<b>Repetition</b> ↓	<b>Duplicate@1M</b> ↓	<b>Non-factual</b> ↓	<b>Pair-irrelevance</b> ↓	<b>Pair-copying</b> ↓
<b>200B-scale</b>	4.3%	0.8%	15.1%	25.6%	0.1%
<b>1T-scale</b>	3.9%	0.8%	8.7%	7.8%	0.9%
<b>Real data</b>	1.8%	0.7%	1.8%	n.a.	n.a.

- ▶ Better data quality with larger scale
- ▶ Synthesized data is not mere repetition

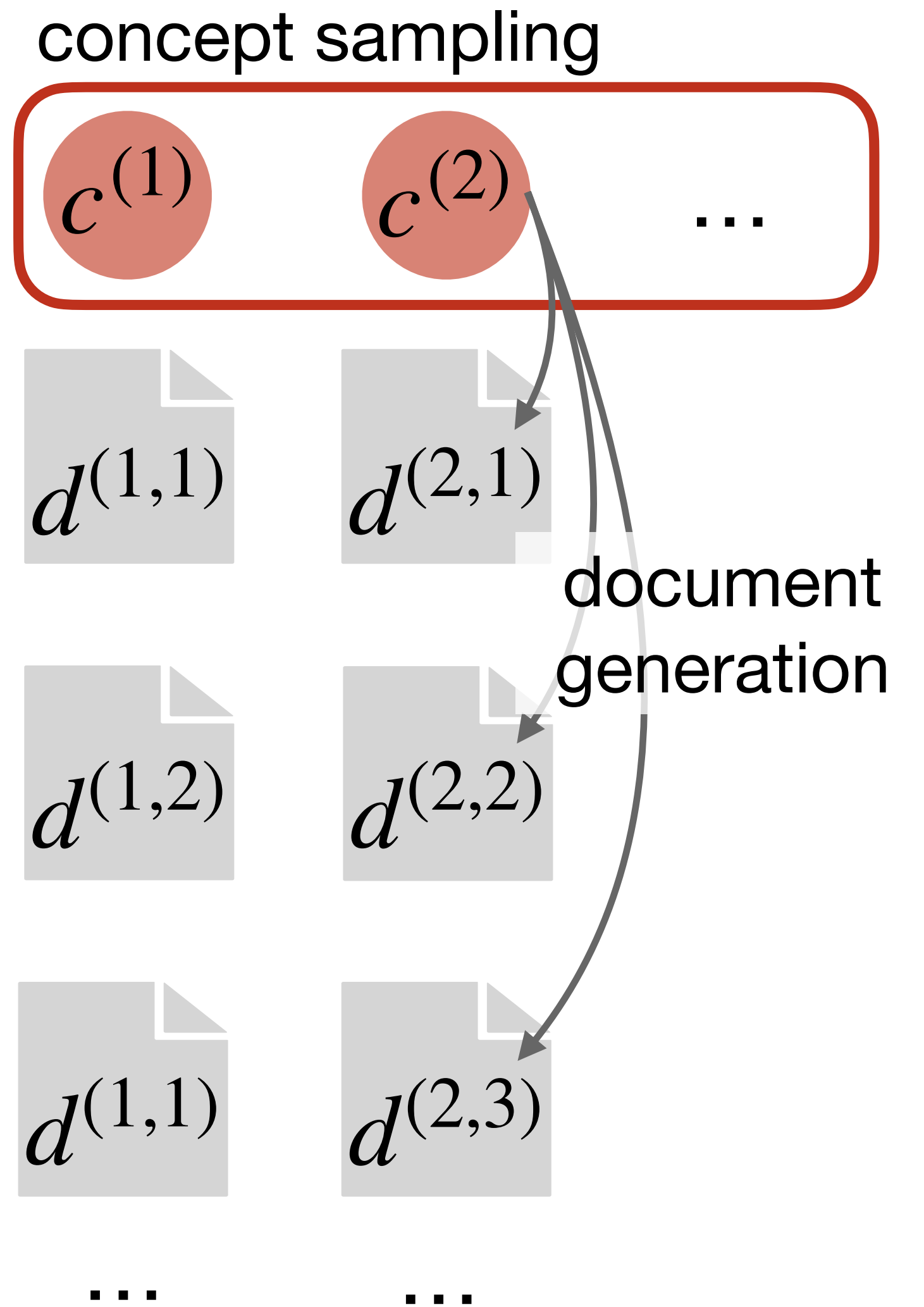
# Bayesian interpretation

- ▶ Model structural correlation as probabilistic dependence
- ▶ Pretraining learns the marginal

$$P(d) = \int P(d | c)P(c)dc$$

- ▶ Synthesizer-tuning learns the posterior

$$P(d_2 | d_1) = \int P(d_2 | c)P(c | d_1)dc$$



Transformer's ignorance of an explicit parametric structure is its greatest blessing

# Searching for weaker forms of self-supervision...

**unsupervised learning**

**supervision**