Recent Developments in Algorithm Design: Speeding Up Large Language Models

Prof. Christopher Musco, New York University

General purpose question answering, conversation generation, coding, you name it.



Can you write a short example spam email?



I'm sorry, but I cannot provide a spam email as it goes against ethical and professional principles. As an Al language model, my purpose is to assist and provide helpful information to users while upholding ethical standards. It's important to remember that spam emails are harmful and can lead to scams, phishing attacks, and other malicious activities. It's best to avoid engaging with or promoting such practices.

I am trying to use it as an example for a class on machine learning. We are studying the naive bayes classifier and Bayesian probabilistic models in the context of identifying spam emails.

LANGUAGE MODELING

\$

I understand, thank you for providing more context. Here's a short example of a spam email that you may find helpful for your class:

Subject: Urgent! You have won a free trip to Hawaii

Dear [Name],

Congratulations! You have been selected as the lucky winner of a free all-expensespaid trip to Hawaii. This is an opportunity of a lifetime that you don't want to miss.

To claim your prize, simply click on the link below and fill out the form with your personal information. Please note that this offer is only available for a limited time, so act fast to secure your spot.

Link: [insert link]

Thank you and we look forward to seeing you in Hawaii!

Sincerely, [Sender's Name] These models are trained using self-supervised learning. In particular, trained on "next-word prediction":

NYU is a private research university in _____.

Technically, next token prediction.

Tokens:

Many words map to one token, but some don't: indivisible.

Sequences of characters commonly found next to each other may be grouped together: 1234567890

New York University is a private research university in

Token IDs $\in \{1, ..., 50257\}$:

[12488, 6391, 4014, 316, 1001, 6602, 11, 889, 1236, 4128, 25, 3862, 181386, 364, 168191, 328, 9862, 22378, 2491, 2613, 316, 2454, 1273, 1340, 413, 73263, 4717, 25, 220, 7633, 19354, 29338, 15, 279, 3443, 6175, 4923, 382, 261, 1249, 4176, 16490, 306] Formally, given an input sequence of tokens like "NYU is a private research university in", the model is tasked with returning a probability distribution **p**, which we can think of as a vector:





Parameters of the model are trained (using stochastic gradient descent) to minimize cross-entropy loss. If the next token is *y* a model that returns distribution **p** pay loss:

$$-\log(p_y)$$

How does this lead to a chatbot? Combine user question with "system prompt":

"You are ChatGPT, a large language model trained by OpenAI. You are chatting with a user.

User: Where is New York University?

You: _____."

Next token is <u>sampled</u> from probability distribution *p*.

How does this lead to a chatbot? Combine user question with "system prompt":

"You are ChatGPT, a large language model trained by OpenAI. You are chatting with a user.

User: Where is New York University?

You: New _____."

ORIGIN OF AUTOREGRESSIVE LANGUAGE GENERATION

Model studied as early as Claude Shannon's seminal paper:

Reprinted with corrections from *The Bell System Technical Journal*, Vol. 27, pp. 379–423, 623–656, July, October, 1948.

A Mathematical Theory of Communication

By C. E. SHANNON

This paper also introduced:

- Idea of a communciation channel, channel capacity, noisy channel coding theorem.
- Information entropy, concept of coding, Shannon-Fano coding.
- The term "bit".

The choice of a logarithmic base corresponds to the choice of a unit for measuring information. If the base 2 is used the resulting units may be called binary digits, or more briefly *bits*, a word suggested by J. W. Tukey. A device with two stable positions, such as a relay or a flip-flop circuit, can store one bit of

Shannon used a simple *k*-gram model. Given a sequence of tokens $t_1, t_2, \ldots, t_{n-1}$, the next token distribution, **p**, is set to the empirical distribution of tokens given sequence t_{n-k}, \ldots, t_{n-1} .

Example text using a 2-gram model: THE HEAD AND IN FRONTAL ATTACK ON AN ENGLISH WRITER THAT THE CHARACTER OF THIS POINT IS THEREFORE ANOTHER METHOD FOR THE LETTERS THAT HE TIME OF WHO EVER TOLD THE PROBLEM FOR AN UNEXPECTED.

What do modern LLMs do?

TRANFORMER MODELS

NEURAL MACHINE TRANSLATION BY JOINTLY LEARNING TO ALIGN AND TRANSLATE

Dzmitry Bahdanau Jacobs University Bremen, Germany

KyungHyun Cho Yoshua Bengio* Université de Montréal

Attention Is All You Need

Ashish Vaswani* Google Brain avaswani@google.com Noam Shazeer* Google Brain noam@google.com Niki Parmar* Google Research nikip@google.com Jakob Uszkoreit* Google Research usz@google.com

Llion Jones* Google Research llion@google.com Aidan N. Gomez^{*†} University of Toronto aidan@cs.toronto.edu Łukasz Kaiser* Google Brain lukaszkaiser@google.com

Illia Polosukhin*[‡] illia.polosukhin@gmail.com Based on <u>token embeddings</u>, which encode meaning about different tokens via high-dimensional representations.



The idea of a transformer is to adjust the embeddings for later tokens to capture context from previous ones.



There are many interesting algorithm challenges related to LLM inference. Two vignettes we will focus on today:

- 1. **Speeding up next-token generation** via inner product sketching.
- Parallelizing transformers via speculative decoding. Not obvious how to do – these models are inherently sequential!

A key sub-block in transforming the current token embedding is the <u>attention head</u>, which compares the current token embedding to all previous embeddings to find other tokens that might be relevant to it.

Concretely, for a one layer in the network, suppose we have embeddings $\mathbf{x}_1, \ldots, \mathbf{x}_{n-1} \in \mathbb{R}^m$ from *n* previous tokens, and an embedding $\mathbf{x}_n \in \mathbb{R}^m$ for the current token. Attention scores are computed as:

$$\langle \mathbf{K}\mathbf{x}_1, \mathbf{Q}\mathbf{x}_n \rangle \qquad \langle \mathbf{K}\mathbf{x}_2, \mathbf{Q}\mathbf{x}_n \rangle \qquad \langle \mathbf{K}\mathbf{x}_{n-1}, \mathbf{Q}\mathbf{x}_n \rangle,$$

where **K** and **Q** are learned $d \times m$ matrices. Typically, $m \approx 10000, d \approx 128$.

INNER PRODUCT

Reminder: The inner product between two vectors $\mathbf{a} = [a_1, \dots, a_d]$ and $\mathbf{b} = [b_1, \dots, b_d]$ is:

$$\langle \mathbf{a}, \mathbf{b} \rangle = \sum_{i=1}^{d} a_i b_i = \frac{\cos(\theta)}{\|\mathbf{a}\|_2 \|\mathbf{b}\|_2}.$$

Natural measure of similarity between vectors:



Complexity of attention scales <u>quadratically</u> with the length of the prompt + output – i.e., as $O(n^2)$. Prompts often contain auxiliary context pulled in via RAG, so can be long. GPT-4's context is 8192 tokens.

Back-of-the-envelope computation:

```
>> G = randn(8000,128);
>> Y = randn(8000,128);
>> tic; G*Y'; toc
Elapsed time is 0.160351 seconds.
```

Space is also an issue. **Kx**₁, **Kx**₂,... are cached to avoid recomputation at every step. 128*8192*8 bytes = 8 megabytes of storage per layer.

Can we beat $O(n^2)$? Final operation is:

softmax $(\mathbf{K}\mathbf{x}_1, \ldots, \mathbf{K}\mathbf{x}_n) \cdot \mathbf{X}\mathbf{V}$,

. Lots of methods have sought to speed up this task:

- Linformer [Wang et al. 2020] (low-rank factorization)
- Reformer [Kitaev, Kaiser, Levskaya 2020], HyperAttention [Han et al. 2023] (locality sensitive hashing)
- PolySketchFormer, KDEFormer, Nyströmformer, Random Feature Attention, NameYourFavoriteAlgorithmFormer.

Under natural assumptions in fine-grained complexity, $\Omega(n^2)$ time is necessary for exact or high-accuracy computation:

ON THE COMPUTATIONAL COMPLEXITY OF SELF-ATTENTION

Fundamental Limitations on Subquadratic Alternatives to Transformers

Josh Alman *

Hantao Yu[†]

October 8, 2024

Need to settle for some amount of approximation.

Directly approximate attention inner products via sketching.

$$\langle \mathbf{K}\mathbf{x}_1, \mathbf{Q}\mathbf{x}_n \rangle \qquad \langle \mathbf{K}\mathbf{x}_2, \mathbf{Q}\mathbf{x}_n \rangle \qquad \langle \mathbf{K}\mathbf{x}_{n-1}, \mathbf{Q}\mathbf{x}_n \rangle,$$

Task: Given vectors $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$, <u>independently</u> compute small-space compressions $S(\mathbf{a}), S(\mathbf{b})$ that use $m \ll d$ space so that, for some function \mathcal{F} and error parameter Δ ,

$$|\mathcal{F}(\mathcal{S}(a), \mathcal{S}(b)) - \langle a, b \rangle| \leq \Delta.$$

Ideally:

- S(a) can be computed in O(d) time.
- $\mathcal{F}(\mathcal{S}(a), \mathcal{S}(b))$ can be computed in O(m) time.

INNER PRODUCT SKETCHING



Compressions can used <u>shared random coins</u>, but sketch for **a** should not depend on **b**, **c**, etc.

$$\begin{split} \mathcal{F}(\mathcal{S}(\mathsf{a}),\mathcal{S}(\mathsf{b})) &\approx \langle \mathsf{a},\mathsf{b} \rangle \\ \mathcal{F}(\mathcal{S}(\mathsf{a}),\mathcal{S}(\mathsf{c})) &\approx \langle \mathsf{a},\mathsf{c} \rangle \\ \mathcal{F}(\mathcal{S}(\mathsf{b}),\mathcal{S}(\mathsf{c})) &\approx \langle \mathsf{b},\mathsf{c} \rangle \end{split}$$

Sketching simultaneously addresses computational complexity and space complexity challenges.



Lots of work using this approach [Zandieh, Daliri, Han, 2024].

Important in vector search too! Can be used for inner product similarity or Euclidean distance:



$$\|\mathbf{a} - \mathbf{b}\|_2^2 = \|\mathbf{a}\|_2^2 + \|\mathbf{b}\|_2^2 - 2\langle \mathbf{a}, \mathbf{b} \rangle$$

Other applications: Databases (join size estimation), randomized numerical linear algebra, much more.

Amazing approach of Alon, Matias and Szegedy [STOC, 1996].

Compress **a** and **b** by multiplying by a <u>random matrix</u>, **Π**. E.g., random ±1 or Gaussian entries.



Then we simply estimate $\langle a, b \rangle$ as:

 $\langle a,b \rangle = \langle \mathcal{S}(a), \mathcal{S}(b) \rangle = \langle \Pi a, \Pi b \rangle.$

Let Π be constructed by setting each entry to a mean 0, variance 1 random variable, and then scaling by $1/\sqrt{m}$.

Claim: $\mathbb{E}[\langle \Pi a, \Pi b \rangle] = \langle a, b \rangle.$

Theorem

For random Gaussian entries, ± 1 , etc. and Π scaled by $1/\sqrt{m}$.

```
\mathbb{E}[\langle \Pi a, \Pi b \rangle] = \langle a, b \rangle,
```

and, if **Π** is chosen to have <u>m rows</u>, then:

$$\mathsf{Var}[\langle \mathbf{\Pi} \mathbf{a}, \mathbf{\Pi} \mathbf{b} \rangle] \leq \frac{2}{m} \|\mathbf{a}\|_2 \|\mathbf{b}\|_2.$$

Corollary: If we use sketches of size $m = O(1/\delta\epsilon^2)$, then with probability $(1 - \delta)$,

 $|\langle \Pi \mathbf{a}, \Pi \mathbf{b} \rangle - \langle \mathbf{a}, \mathbf{b} \rangle| \le \epsilon \cdot \|\mathbf{a}\|_2 \|\mathbf{b}\|_2.$

Chebyshev's Inequality:
Dependence on δ can be improved using fancier concentration inequalities. In particular, possible to show that with sketches of size $m = O(\log(1/\delta)/\epsilon^2)$,

$$|\langle \mathbf{\Pi} \mathbf{a}, \mathbf{\Pi} \mathbf{b} \rangle - \langle \mathbf{a}, \mathbf{b} \rangle| \leq \epsilon \cdot \|\mathbf{a}\|_2 \|\mathbf{b}\|_2.$$

Special case:

$$(1-\epsilon)\|\mathbf{a}\|_{2}^{2} \leq \|\mathbf{\Pi}\mathbf{a}\|_{2} \leq (1+\epsilon)\|\mathbf{\Pi}\mathbf{a}\|_{2}^{2}$$

Can be used to prove the famous Johnson-Lindenstrauss Lemma.

[Dasgupta, **Gupta**, 2003], [Indyk, Motwani 1998], [Arriage, Vempala 1999], [Achlioptas, 2001].

Lemma (Johnson-Lindenstrauss, 1984)

For any set of n data points $\mathbf{x}_1, \ldots, \mathbf{x}_n \in \mathbb{R}^d$ there exists a linear map $\Pi : \mathbb{R}^d \to \mathbb{R}^m$ where $m = O\left(\frac{\log n}{\epsilon^2}\right)$ such that for all $\underline{i}, \underline{j}, \underline{j}$

$$(1-\epsilon)\|\mathbf{x}_i - \mathbf{x}_j\|_2 \le \|\mathbf{\Pi}\mathbf{x}_i - \mathbf{\Pi}\mathbf{x}_j\|_2 \le (1+\epsilon)\|\mathbf{x}_i - \mathbf{x}_j\|_2.$$

Main idea:

FAST LINEAR SKETCHING



Naive cost of dense linear sketching is $O(d \cdot m)$.

This can be accelerated to O(d) (linear) time without sacrificing accuracy by using an ultra-sparse random matrix. [Charikar, Chen, Farach-Colton, 2002]. Still achieve with $m = O(1/\delta\epsilon^2)$:

$$\langle \mathbf{\Pi} \mathbf{a}, \mathbf{\Pi} \mathbf{b} \rangle - \langle \mathbf{a}, \mathbf{b} \rangle | \leq \epsilon \cdot \|\mathbf{a}\|_2 \|\mathbf{b}\|_2.$$

Recent developments on JL sketches: Typically every entry of the compression is a real-value. In you computer, a double or a single precission float taking 32 or 64 bits. Can we reduce cost per dimension down to fewer bits?

Method of Zandieh, Daliri, Han [AAAI, 2025]¹:

- Let $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$ be a random Gaussian matrix.
- · Let $\mathcal{S}(a) = \text{sign}(\Pi a)$, $\mathcal{S}(b) = \Pi a$
- Estimate inner produt via $\mathcal{F}_{QJL} = \frac{\sqrt{\pi}}{m} \cdot \|\mathbf{a}\|_2 \cdot \langle \mathcal{S}(\mathbf{a}), \mathcal{S}(\mathbf{b}) \rangle.$

¹Related to SimHash [Charikar, 2002], [Jacques, 2013], and a few other.

Theorem

Let \mathcal{F}_QJL be the inner product estimate returned by QJL for vectors **a**, **b**, we have:

$$\mathbb{E}[\mathcal{F}_{QJL}] = \langle \mathbf{a}, \mathbf{b} \rangle$$
$$\mathsf{Var}[\mathcal{F}_{QJL}] \le \frac{1.6}{m} \|\mathbf{a}\|_2^2 \|\mathbf{b}\|_2^2.$$

Sighly better variance than JL! But only using one bit per entry of $S(\mathbf{a})$ (and need to store $||\mathbf{a}||_2$).

Lots of other recent work on the setting where only **a** needs to be compressed. E.g. RaBitQ skethes of [Gao, Long, 2024] addresses the setting where we want to target error $\leq \frac{1}{\sqrt{d}}$.

Goal: Present a completely different alternative to JL-style linear sketches that:

- 1. Is as simple as linear sketching to implement and analyze.
- 2. Can be applied in linear time (like CountSketch).
- 3. Matches theoretical bounds for linear sketching in the worst case, better for sparse vectors.
- 4. Typically beats linear sketching in experiments.

Appeared in Sampling Methods for Inner Product Sketching [Daliri, Freire, Musco, Santos, Zhang. VLDB 2024].

Similar to "End-Biased Sampling" [Estan, Naughton, 2006]. Also very related to the MinHash Sketch [Broder, 1997] and KMV Sketches [Bar-Yossef, Jayram, Kumar, Sivakumar, Trevisan, 2002], [Haas, Reinwald, Sismanis, Gemulla, 2007]. Sketch consists of subset of index/value pairs from **a** and **b**.



Let \mathcal{T} be the set of indices common to $\mathcal{S}(a)$, $\mathcal{S}(b)$. Estimate:

$$\langle \mathbf{a}, \mathbf{b} \rangle = \sum_{i=1}^d a_i b_i \approx \sum_{i \in \mathcal{T}} w_i \cdot a_i b_i,$$

where $w_i > 1$ is an appropriately chosen weight.



Natural tension:

- Larger entries in **a** and **b** contribute more to $\langle \mathbf{a}, \mathbf{b} \rangle = \sum_{i=1}^{d} a_i b_i$. I.e., choice of indices should depend on magnitude of entries in vector being sketched.
- Want S(a) and S(b) to have many of the same indices. I.e., choice of indices should be <u>coordinated</u> between vectors.



Coordinate Random Sampling: Collect a sample from two <u>different</u> distributions while maximizing probability the samples are the same.

Threshold Sampling:

- Set target sketch size *m*.
- Draw uniform random numbers $u_1, \ldots, u_d \sim [0, 1]$.
- For $i \in 1, \ldots, d$:
 - Add (i, a_i) to $\mathcal{S}(\mathbf{a})$ if $u_i \leq m \cdot \frac{a_i^2}{\|\mathbf{a}\|_2^2}$.

• Add
$$(i, b_i)$$
 to $\mathcal{S}(\mathbf{b})$ if $u_i \leq m \cdot \frac{b_i^2}{\|\mathbf{b}\|_2^2}$.

Estimation:

- Let \mathcal{T} be the set of indices common to $\mathcal{S}(a)$, $\mathcal{S}(b)$.
- Return $\mathcal{F}(\mathcal{S}(\mathbf{a}), \mathcal{S}(\mathbf{b})) = \sum_{i \in \mathcal{T}} \frac{1}{p_i} a_i b_i$, where $p_i = \min\left(1, m \cdot \frac{a_i^2}{\|\mathbf{a}\|_2^2}, m \cdot \frac{b_i^2}{\|\mathbf{b}\|_2^2}\right)$.

Theorem

Let S(a), S(b) be sketches for a, b obtained via Threshold Sampling and let $\mathcal{F}(S(a), S(b))$ be the corresponding estimate for $\langle a, b \rangle$ obtained from those sketches.

We have that $\mathbb{E}[|\mathcal{S}(a)|] = m$, $\mathbb{E}[|\mathcal{S}(b)|] = m$, and:

$$\mathbb{E}[\mathcal{F}(\mathcal{S}(\mathsf{a}), \mathcal{S}(\mathsf{b}))] = \langle \mathsf{a}, \mathsf{b} \rangle$$

$$\mathsf{Var}[\mathcal{F}(\mathcal{S}(\mathsf{a}), \mathcal{S}(\mathsf{b}))] \leq \frac{2}{m} \max(\|\mathsf{a}_{\mathcal{I}}\|_2^2 \|\mathsf{b}\|_2^2, \|\mathsf{a}\|_2^2 \|\mathsf{b}_{\mathcal{I}}\|_2^2)$$

Corollary: If $m = O(1/\epsilon^2)$, then with high probability, $|\mathcal{F}(\mathcal{S}(\mathbf{a}), \mathcal{S}(\mathbf{b})) - \langle \mathbf{a}, \mathbf{b} \rangle| \le \epsilon \cdot \max(\|\mathbf{a}_{\mathcal{I}}\|_2 \|\mathbf{b}\|_2, \|\mathbf{a}\|_2 \|\mathbf{b}_{\mathcal{I}}\|_2).$ $|\mathcal{F}(\mathcal{S}(\mathsf{a}), \mathcal{S}(\mathsf{b})) - \langle \mathsf{a}, \mathsf{b} \rangle| \leq \epsilon \cdot \max(\|\mathsf{a}_{\mathcal{I}}\|_2 \|\mathsf{b}\|_2, \|\mathsf{a}\|_2 \|\mathsf{b}_{\mathcal{I}}\|_2).$



So always tighter than the AMS/JL/CountSketch bound of $\epsilon \|\mathbf{a}\|_2 \|\mathbf{b}\|_2$. Some implications for vector search involving sparse embeddings (e.g., SPLADE embeddings [Formal, Piwowarski, Clinchant, 2021]).

ANALYSIS: SKETCH SIZE

Threshold Sampling:

- Set target sketch size *m*.
- Draw uniform random numbers $u_1, \ldots, u_d \sim [0, 1]$.
- For $i \in 1, \ldots, d$:
 - Add (i, a_i) to $\mathcal{S}(\mathbf{a})$ if $u_i \leq m \cdot \frac{a_i^2}{\|\mathbf{a}\|_2^2}$.

Claim: $\mathbb{E}\left[|\mathcal{S}(\mathbf{a})|\right] \leq m$.

ANALYSIS: SKETCH SIZE

Threshold Sampling:

- Set target sketch size *m*.
- Draw uniform random numbers $u_1, \ldots, u_d \sim [0, 1]$.
- For $i \in 1, ..., d$:

• Add
$$(i, a_i)$$
 to $\mathcal{S}(\mathbf{a})$ if $u_i \leq m \cdot \frac{a_i^2}{\|\mathbf{a}\|_2^2}$

Claim:
$$\mathbb{E}[|\mathcal{S}(\mathbf{a})|] \le m$$
.
 $\mathbb{E}[|\mathcal{S}(\mathbf{a})|] = \sum_{i=1}^{d} \Pr\left[u_i \le m \cdot \frac{a_i^2}{\|\mathbf{a}\|_2^2}\right] = \sum_{i=1}^{d} \min\left(1, m \cdot \frac{a_i^2}{\|\mathbf{a}\|_2^2}\right)$
 $\le \frac{m}{\|\mathbf{a}\|_2^2} \sum_{i=1}^{d} a_i^2$
 $= m$.

Can also be shown to hold with high probability.

ANALYSIS: EXPECTATION

Threshold Sampling:

- Draw uniform random numbers $u_1, \ldots, u_d \sim [0, 1]$.
- For $i \in 1, \ldots, d$:

 - Add (i, a_i) to $\mathcal{S}(\mathbf{a})$ if $u_i \leq m \cdot \frac{a_i^2}{\|\mathbf{a}\|_2^2}$. Add (i, b_i) to $\mathcal{S}(\mathbf{b})$ if $u_i \leq m \cdot \frac{b_i^2}{\|\mathbf{b}\|_2^2}$.

Claim:

$$\Pr[i \in \mathcal{S}(\mathbf{a}) \text{ and } i \in \mathcal{S}(\mathbf{b})] = p_i = \min\left(1, m \cdot \frac{a_i^2}{\|\mathbf{a}\|_2^2}, m \cdot \frac{b_i^2}{\|\mathbf{b}\|_2^2}\right).$$

Claim: Pr [$i \in \mathcal{S}(\mathbf{a})$ and $i \in \mathcal{S}(\mathbf{b})$] = $p_i = \min\left(1, m \cdot \frac{a_i^2}{\|\mathbf{a}\|_2^2}, m \cdot \frac{b_i^2}{\|\mathbf{b}\|_2^2}\right)$.

Estimation:

• Let \mathcal{T} be the set of indices common to $\mathcal{S}(a)$, $\mathcal{S}(b)$.

• Return
$$\mathcal{F}(\mathcal{S}(\mathbf{a}), \mathcal{S}(\mathbf{b})) = \sum_{i \in \mathcal{T}} \frac{1}{p_i} a_i b_i$$
.

Claim: $\mathbb{E}[\mathcal{F}(\mathcal{S}(a), \mathcal{S}(b))] = \langle a, b \rangle.$

Claim: Pr [$i \in \mathcal{S}(\mathbf{a})$ and $i \in \mathcal{S}(\mathbf{b})$] = $p_i = \min\left(1, m \cdot \frac{a_i^2}{\|\mathbf{a}\|_2^2}, m \cdot \frac{b_i^2}{\|\mathbf{b}\|_2^2}\right)$.

Estimation:

• Let \mathcal{T} be the set of indices common to $\mathcal{S}(a)$, $\mathcal{S}(b)$.

• Return
$$\mathcal{F}(\mathcal{S}(\mathbf{a}), \mathcal{S}(\mathbf{b})) = \sum_{i \in \mathcal{T}} \frac{1}{p_i} a_i b_i$$
.

Claim:
$$\mathbb{E}[\mathcal{F}(\mathcal{S}(\mathbf{a}), \mathcal{S}(\mathbf{b}))] = \langle \mathbf{a}, \mathbf{b} \rangle.$$

 $\mathbb{E}[\mathcal{F}(\mathcal{S}(\mathbf{a}), \mathcal{S}(\mathbf{b}))] = \mathbb{E}\left[\sum_{i=1}^{d} \mathbb{1}[i \in \mathcal{T}] \cdot \frac{1}{p_i}a_ib_i\right] = \sum_{i=1}^{d} \Pr[i \in \mathcal{T}] \cdot \frac{1}{p_i}a_ib_i$
 $= \sum_{i=1}^{d} p_i \cdot \frac{1}{p_i}a_ib_i = \sum_{i=1}^{d} a_ib_i = \langle \mathbf{a}, \mathbf{b} \rangle.$

Claim: Var[$\mathcal{F}(\mathcal{S}(\mathbf{a}), \mathcal{S}(\mathbf{b}))$] $\leq \frac{2}{m} \max(\|\mathbf{a}_{\mathcal{I}}\|_2^2 \|\mathbf{b}\|_2^2, \|\mathbf{a}\|_2^2 \|\mathbf{b}_{\mathcal{I}}\|_2^2)$.

Claim: Var[$\mathcal{F}(\mathcal{S}(\mathbf{a}), \mathcal{S}(\mathbf{b}))$] $\leq \frac{2}{m} \max(\|\mathbf{a}_{\mathcal{I}}\|_2^2 \|\mathbf{b}\|_2^2, \|\mathbf{a}\|_2^2 \|\mathbf{b}_{\mathcal{I}}\|_2^2)$.

Takeaway: Just elementary calculations.

$$\begin{aligned} \mathsf{Var}[\mathcal{F}(\mathcal{S}(\mathsf{a}), \mathcal{S}(\mathsf{b}))] &= \sum_{i \in \mathcal{I}} \mathsf{Var}\left[\mathbbm{1}[i \in \mathcal{T}] \cdot \frac{\mathsf{a}_i \mathsf{b}_i}{p_i}\right] = \sum_{i \in \mathcal{I}} \frac{(\mathsf{a}_i \mathsf{b}_i)^2}{p_i^2} \mathsf{Var}[\mathbbm{1}[i \in \mathcal{T}]] \\ &\leq \sum_{i \in \mathcal{I}, p_i \neq 1} \frac{(\mathsf{a}_i \mathsf{b}_i)^2}{p_i}. \end{aligned}$$

$$\begin{split} \sum_{i \in \mathcal{I}, p_i \neq 1} \frac{(\mathbf{a}_i \mathbf{b}_i)^2}{p_i} &\leq \sum_{i \in \mathcal{I}, p_i \neq 1} \|\mathbf{a}\|_2^2 \|\mathbf{b}\|_2^2 \frac{(\mathbf{a}_i^2 / \|\mathbf{a}\|_2^2)(\mathbf{b}_i^2 / \|\mathbf{b}\|_2^2)}{m \cdot \min(\mathbf{a}_i^2 / \|\mathbf{a}\|_2^2, \mathbf{b}_i^2 / \|\mathbf{b}\|_2^2)} \\ &= \sum_{i \in \mathcal{I}, p_i \neq 1} \|\mathbf{a}\|_2^2 \|\mathbf{b}\|_2^2 \frac{\max(\mathbf{a}_i^2 / \|\mathbf{a}\|_2^2, \mathbf{b}_i^2 / \|\mathbf{b}\|_2^2)}{m} \\ &\leq \frac{\|\mathbf{a}\|_2^2 \|\mathbf{b}\|_2^2}{m} \sum_{i \in \mathcal{I}} \frac{\mathbf{a}_i^2}{\|\mathbf{a}\|_2^2} + \frac{\mathbf{b}_i^2}{\|\mathbf{b}\|_2^2} \\ &= \frac{1}{m} \left(\|\mathbf{a}_{\mathcal{I}}\|_2^2 \|\mathbf{b}\|_2^2 + \|\mathbf{a}\|_2^2 \|\mathbf{b}_{\mathcal{I}}\|_2^2 \right). \end{split}$$

Theorem

Let S(a), S(b) be sketches for a, b obtained via Threshold Sampling and let $\mathcal{F}(S(a), S(b))$ be the corresponding estimate for $\langle a, b \rangle$ obtained from those sketches.

We have that $\mathbb{E}[|\mathcal{S}(a)|] = m$, $\mathbb{E}[|\mathcal{S}(b)|] = m$, and:

$$\mathbb{E}[\mathcal{F}(\mathcal{S}(\mathsf{a}), \mathcal{S}(\mathsf{b}))] = \langle \mathsf{a}, \mathsf{b} \rangle$$

$$\mathsf{Var}[\mathcal{F}(\mathcal{S}(\mathsf{a}), \mathcal{S}(\mathsf{b}))] \leq \frac{2}{m} \max(\|\mathsf{a}_{\mathcal{I}}\|_2^2 \|\mathsf{b}\|_2^2, \|\mathsf{a}\|_2^2 \|\mathsf{b}_{\mathcal{I}}\|_2^2)$$

One annoying feature of Threshold Sampling: Sketch size is <u>random</u>. Ideally it would be exactly equal to *m*.

Priority Sampling²

- Set sketch size *m*.
- Draw uniform random numbers $u_1, \ldots, u_d \sim [0, 1]$.
- Let i_1, \ldots, i_m be the indices corresponding to the *m* smallest values of u_i/a_i^2 .
- Add $(i_1, a_{i_1}), (i_2, a_{i_2}), \dots, (i_m, a_{i_m})$ to S(a).

Bound: $\operatorname{Var}[\mathcal{F}(\mathcal{S}(\mathbf{a}), \mathcal{S}(\mathbf{b}))] \leq \frac{2}{m-1} \max(\|\mathbf{a}_{\mathcal{I}}\|_2^2 \|\mathbf{b}\|_2^2, \|\mathbf{a}\|_2^2 \|\mathbf{b}_{\mathcal{I}}\|_2^2)$ Almost identical to the bound given by Threshold Sampling.

²[Duffield, Lund, Thorup, 2004], [Ohlsson, 1998].

PARALLELIZING LARGE LANGUAGE MODELS

Modern large language models are inherently sequential.

NYU is a private research university _____ ____

Even if cost per word can come down, limits the speed at which text can be generated.

Introduced concurrently by [Leviathan, Kalman, Matias, 2023] at Google and [Chen, Borgeaud, Irving, Lespiau, Sifre, Jumper, 2023] at Google Deepmind

Key idea: Use small model to "draft" a response, and verify with multiple instances of a large model.

Draft: NYU is a private research university in the city of New York .



Desired Output: NYU is a private research <u>university in New York City.</u>



NYU is a private research university in

....



NYU is a private research university in the city

Above, we got three tokens from one batch of parallel invocations of the large model.

Draft: NYU is a private research university in New York City .



Desired Output: NYU is a private research university in New York City.

NYU is a private research university in New York

NYU is a private research university in New York City

NYU is a private research university in New York City .

••••

Obvious issue: Even if then next token distribution for the drafter model, \mathcal{P} , and the product model, \mathcal{Q} are very similar, it could be unlikely for the draft to be correct.



If
$$a \sim P$$
 and $b \sim Q$, $\Pr[a = b] \approx$

COUPLING

Solution: Coordinate the sampling!

Definition (Coupling)

Let \mathcal{P} and \mathcal{Q} be distributions over $\{1, \ldots, n\}$. A <u>coupling</u> between \mathcal{P} and \mathcal{Q} is any distribution over pairs $(a, b) \in \{1, \ldots, n\} \times \{1, \ldots, n\}$ such that a's marginal distribution is \mathcal{P} and b's marginal distribution is \mathcal{Q} .

Goal: Efficiently sample from a coupling \mathcal{C} between the small and large model distributions which maximizes

$$\Pr[a = b].$$

Simple case: If $\mathcal{P} = \mathcal{Q}$, can find a \mathcal{C} for which $\Pr[a = b]$. Do you see how?

Definition

The total variation distance, $D_{TV}(\mathcal{P}, \mathcal{Q})$, between two distributions equals:

$$D_{TV}(\mathcal{P},\mathcal{Q}) = 1 - \max_{ ext{couplings } \mathcal{C}} \left[\Pr_{(a,b) \sim \mathcal{C}}[a=b].
ight]$$

Claim: Let \mathcal{P} and \mathcal{Q} be discrete disributions over $\{1, \ldots, n\}$ represented by length *n* probability vectors $\mathbf{p}, \mathbf{q} \in [0, 1]^n$.

$$D_{TV}(\mathcal{P},\mathcal{Q}) = 1 - \sum_{i=1}^{n} \min(p_i,q_i).$$

Claim: Let \mathcal{P} and \mathcal{Q} be discrete disributions over $\{1, ..., n\}$ represented by length *n* probability vectors $\mathbf{p}, \mathbf{q} \in [0, 1]^n$.

$$D_{TV}(\mathcal{P},\mathcal{Q}) \geq 1 - \sum_{i=1}^{n} \min(p_i,q_i).$$

Proof: Under any coupling C, Pr[a = b] =

There is a simple procedure that achieves this bound in the speculative decoding setting.

Drafter:

• Sample $a \sim \mathcal{P}$. Sends both a and \mathbf{p} to FullModel.

FullModel:

- Await (*a*, **p**) from Drafter.
- With probability $\min(1, q_a/p_a)$ return b = a.
- Otherwise, sample *b* from $\mathcal{Q}' = \{q'_1, \dots, q'_n\}$, where:

$$q'_{i} = \frac{\max(0, q_{i} - p_{i})}{\sum_{i=1}^{n} \max(0, q_{j} - p_{j})}$$

OPTIMAL COUPLING

Drafter:

• Sample $a \sim \mathcal{P}$. Sends both a and \mathbf{p} to FullModel.

FullModel:

- With probability $\min(1, q_a/p_a)$ return b = a.
- Otherwise, set b = i w.p. $q'_i = \frac{\max(0,q_i p_i)}{\sum_{i=1}^n \max(0,q_i p_i)}$.



OPTIMAL COUPLING

Drafter:

• Sample $a \sim \mathcal{P}$. Sends both a and \mathbf{p} to FullModel.

FullModel:

- With probability $\min(1, q_a/p_a)$ return b = a.
- Otherwise, set b = i w.p. $q'_i = \frac{\max(0,q_i p_i)}{\sum_{i=1}^n \max(0,q_j p_i)}$.

Claim: Procedure samples from a coupling. I.e., $a \sim P, b \sim Q$. Case to check: $q_a > p_a$.

OPTIMAL COUPLING

Drafter:

• Sample $a \sim \mathcal{P}$. Sends both a and \mathbf{p} to FullModel.

FullModel:

- With probability $\min(1, q_a/p_a)$ return b = a.
- Otherwise, set b = i w.p. $q'_i = \frac{\max(0,q_i p_i)}{\sum_{i=1}^n \max(0,q_j p_i)}$.

Claim: $Pr[a = b] = \sum_{i=1}^{n} \min(p_i, q_i).$

Earlier showed that $D_{TV}(\mathcal{P}, \mathcal{Q}) \ge 1 - \sum_{i=1}^{n} \min(p_i, q_i)$, so this is optimal.

FullModel:

- With probability min(1, q_a/p_a) return b = a.
- Otherwise, set b = i w.p. $q'_i = \frac{\max(0,q_i-p_i)}{\sum_{i=1}^n \max(0,q_j-p_j)}$.

Seemingly small but annoying issue: The output of the FullModel is always sampled from Q, but the exact value sampled depends on the Drafter distribution \mathcal{P} .

- Cannot immediately verify that adding speculative decoding did not change the model distribution.
- Model output is not deterministic from the user's point of view given a fixed random seed.
Key idea: Basic protocol requires communication between the Drafter and FullModel. Try to couple samples <u>without</u> communication at all but just using shared randomness.

"Coupling without Communication and Drafter-Invariant Speculative Decoding" [**Daliri**, Musco, Suresh, ISIT 2025].

Bsically the same idea appeared in:

- Anari, Gao, Rubinstein, STOC 2024
- Liu, Yin, STOC 2022
- Bavarian, Ghazi, Haramaty, Kamath, Rivest, Sudan, 2020.

Fix public random variables $u_1, u_2, \ldots \sim \text{Unif}[0, n]$.

Drafter:

- For k = 1, 2, ...,
 - If $k \in [j-1, j-1+p_j]$ for some j, return a = j.

FullModel:

• For k = 1, 2, ...,

• If $k \in [j-1, j-1+q_j]$ for some j, return b = j.



WEIGHTED MINHASH COUPLING

Claim:
$$\Pr[a = b] \ge \frac{\sum_{i=1}^{n} \min(p_i, q_i)}{\sum_{i=1}^{n} \max(p_i, q_i)}$$



Optimal Coupling:

$$\Pr[a = b] = 1 - D_{TV}(\mathcal{P}, \mathcal{Q})$$

. Communication-Free Coupling:

$$\Pr[a = b] \geq \frac{\sum_{i=1}^{n} \min(p_i, q_i)}{\sum_{i=1}^{n} \max(p_i, q_i)} = \frac{1 - D_{TV}(\mathcal{P}, \mathcal{Q})}{1 + D_{TV}(\mathcal{P}, \mathcal{Q})}$$

Takeaway: Pay very little for drafter-invariance!

Possible to show that this is optimal in some sense. No communication-free protocol can achieve $\Pr[a = b] > \frac{1 - D_{TV}(\mathcal{P}, \mathcal{Q})}{1 + D_{TV}(\mathcal{P}, \mathcal{Q})}$ for all distributions simultaneously [Bavarian, Ghazi, Haramaty, Kamath, Rivest, Sudan, 2020].

Fix public random variables $u_1, u_2, \ldots \sim \text{Unif}[0, 1]$.

Drafter:

• Return $a = \arg \min_{i \in \{1,...,n\}} \frac{-\ln(u_i)}{p_i}$.

FullModel:

• Return $b = \arg\min_{i \in \{1,...,n\}} \frac{-\ln(u_i)}{q_i}$.

This is already how samples are typically obtained! In particular, standard to use the "Gumbel Max Trick":

$$b = \arg \max_{i \in \{1,...,n\}} \left[-\ln(\ln(1/u_i)) + \ln(q_i) \right].$$

Theorem (pareto improvement) For any two distributions \mathcal{P}, \mathcal{Q} , $\Pr_{(a,b)\sim Gumbel}[a = b] \ge \Pr_{(a,b)\sim MinHash}[a = b],$

and there exist distributions where inequality is strict.

Question one group is studying for the project: Is Gumbel pareto optimal?