Recent Efficiency Improvements to Transformers

David Woodruff

Carnegie Mellon University / Google

Outline

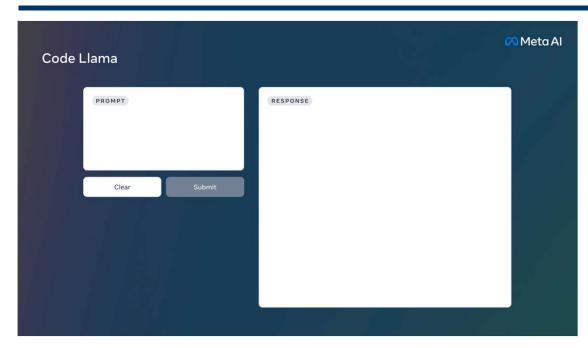
1. Background on Transformers and time complexity

2. HyperAttention

3. PolySketchFormer

4. Conclusions and Recent Work

Al Revolution



2023 IN REVIEW

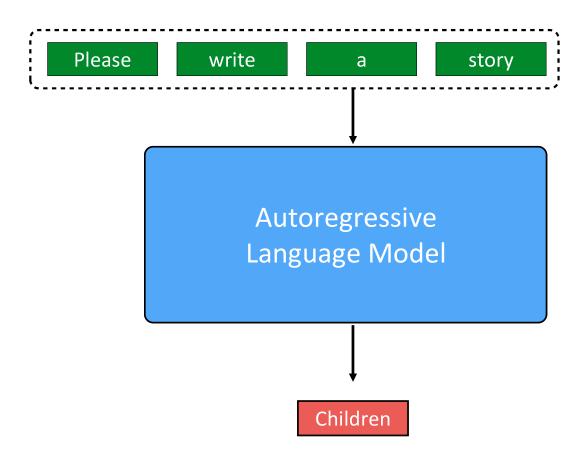
THE YEAR A.I. ATE THE INTERNET

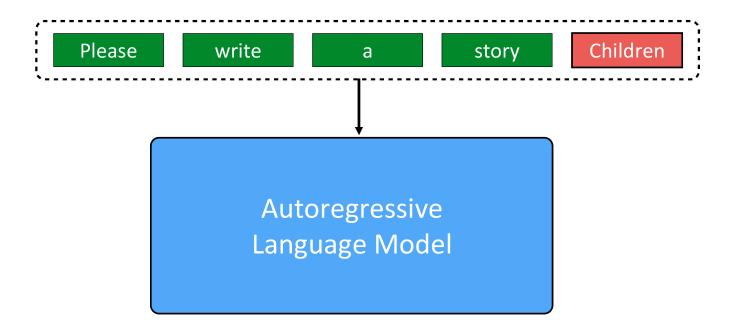
Call 2023 the year many of us learned to communicate, create, cheat, and collaborate with robots.

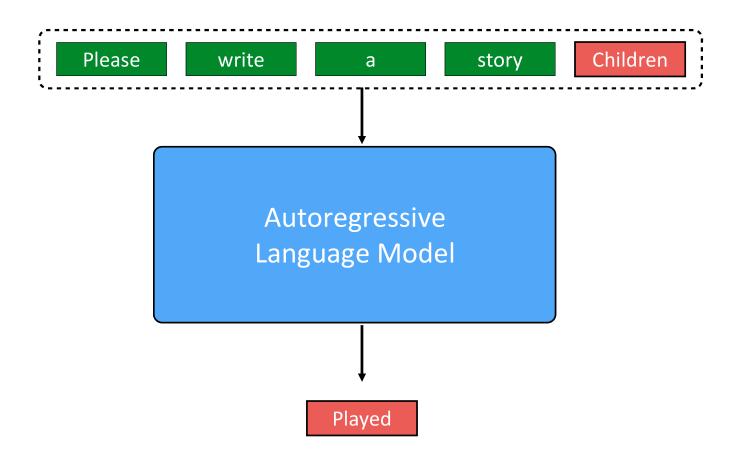
By Sue Halpern

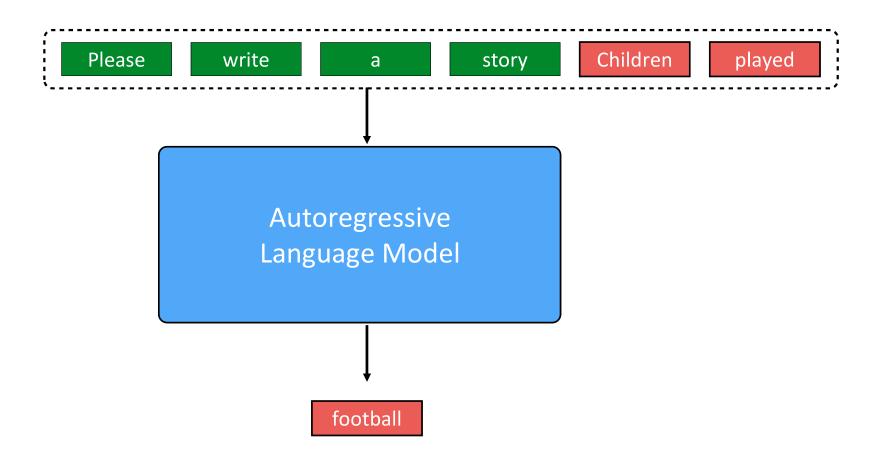
December 8, 2023

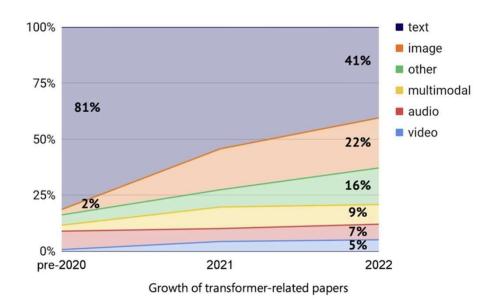




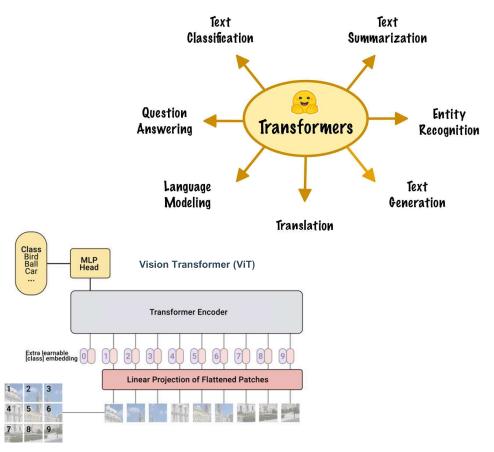








Transformer (Vaswani et al., 17')

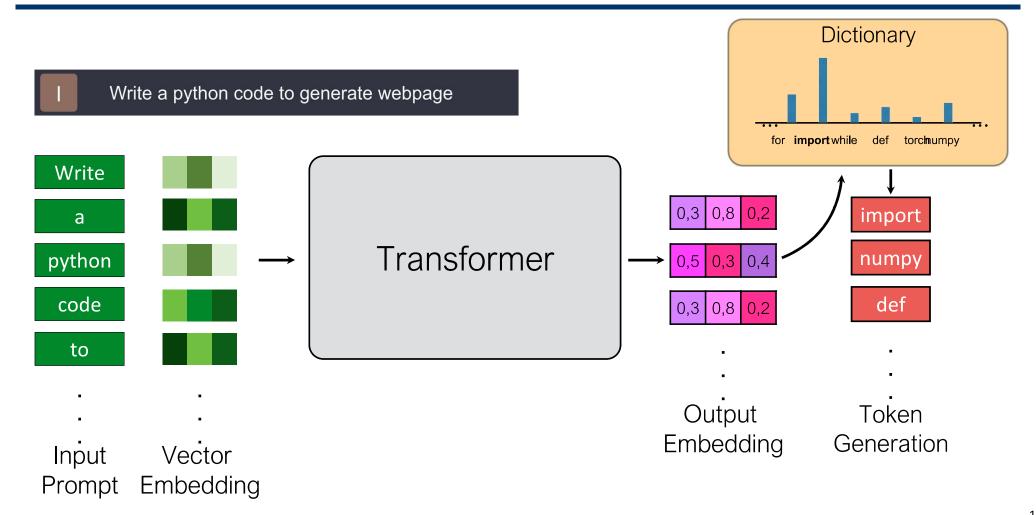


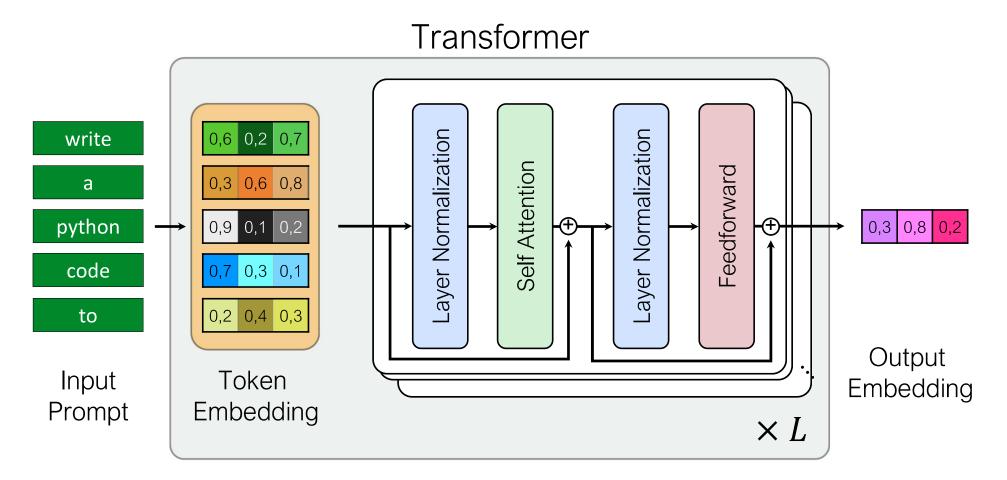
Write a python code to generate webpage Write a Transformer python code to Input Vector Prompt Embedding

Write a python code to generate webpage Write 0,3 0,8 0,2 a Transformer python 0,5 0,3 0,4 code 0,3 0,8 0,2 to Output **Embedding** Input Vector

Prompt

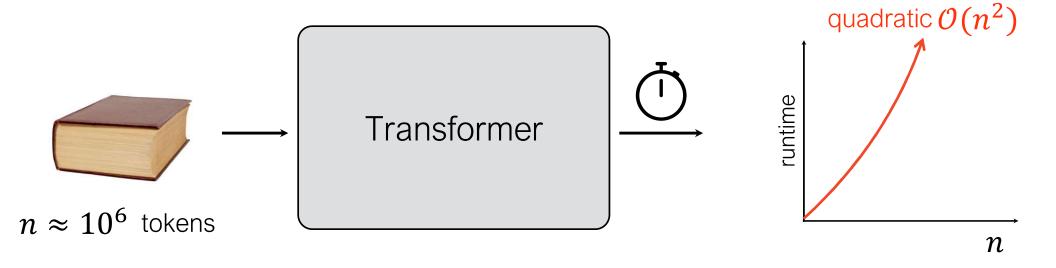
Embedding

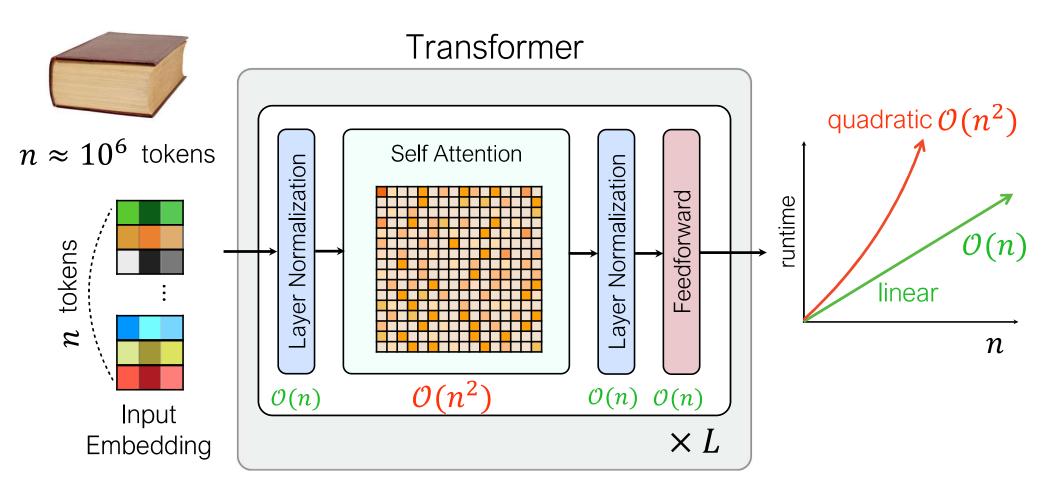


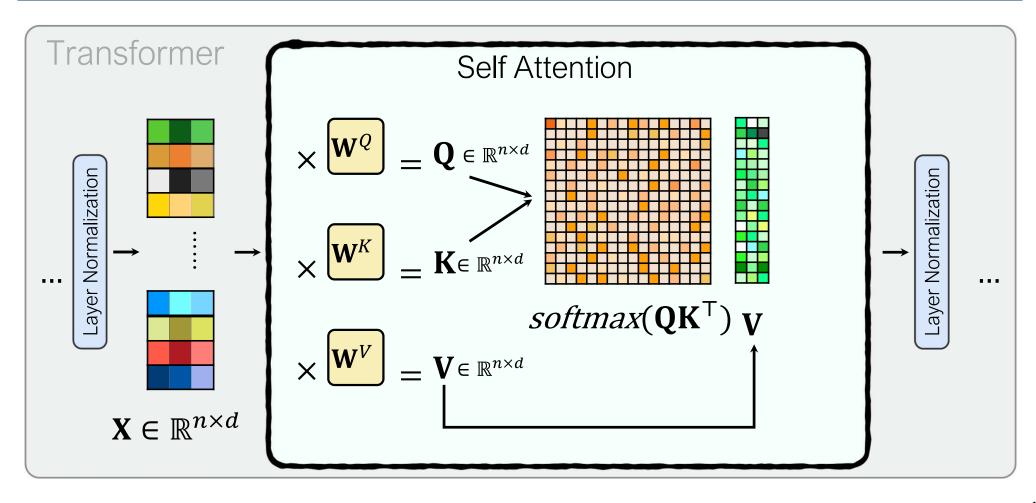


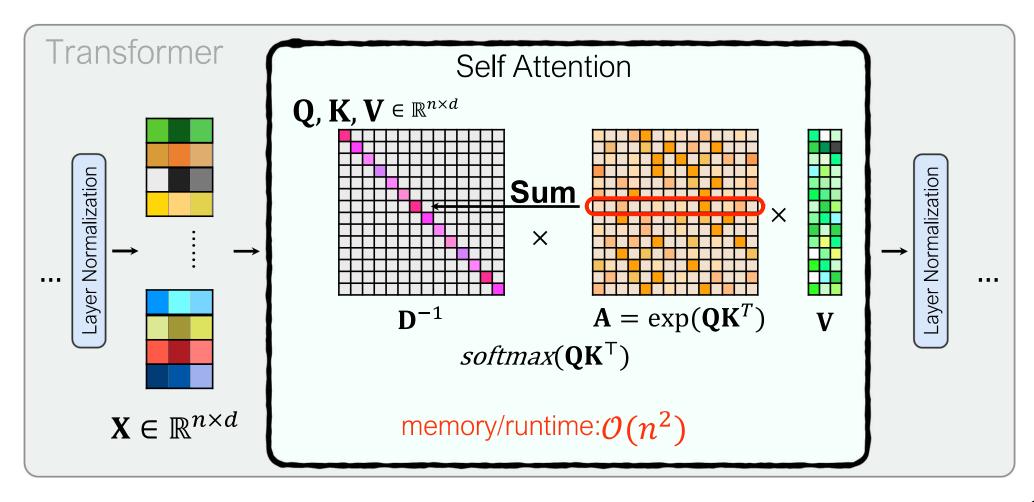
Self Attention write write а а python python X 0.1 0.1 0.1 0.4 0.3 code code to to Attention scores of previous tokens











Previous Work

Sparse Structure

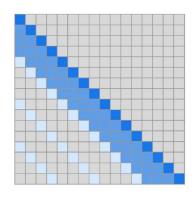
- ▶ Local Attention (Parmar et al., 18')
- ▶ Sparse Transformer (Child et al., 19')
- ▶ Longformer (Beltagy et al., 20')
- ▶ Reformer (Kitaev et al., 20')
- ▶ Sinkhorn Attention (Tay et al., 20')

Kernel Methods

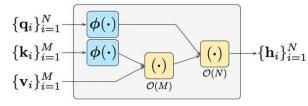
- ▶ Lambda network (Bello et al., 21')
- ▶ Performer (Choromanski et al., 21')
- Random Feature Attention (Peng et al., 21')
- ▶ Randomized Attention (Zheng et al., 22')

Low-rank Approximation

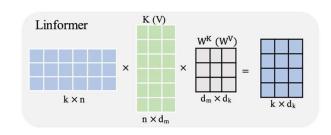
- ▶ Linformer (Wang et al., 20')
- Nystromformer (Xiong et al., 21')
- ▶ Nested Attention (Max et al., 21')



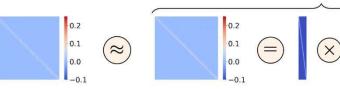
(b) Sparse Transformer (strided)



(b) Random feature attention.



softmax



Previous Work

Sparse Structure

- ▶ Local Attention (Parmar et al., 18')
- ▶ Sparse Transformer (Child et al., 19')
- ▶ Longformer (Beltagy et al., 20')
- ▶ Reformer (Kitaev et al., 20')
- ▶ Sinkhorn Attention (Tay et al., 20')

Kernel Methods

- ▶ Lambda network (Bello et al., 21')
- ▶ Performer (Choromanski et al., 21')
- ▶ Random Feature Attention (Peng et al., 21')
- ▶ Randomized Attention (Zheng et al., 22')

Low-rank Approximation

- ▶ Linformer (Wang et al., 20')
- ▶ Nystromformer (Xiong et al., 21')
- Nested Attention (Max et al., 21')

(Alman & Song 23') High quality (1/poly(n)) entrywise approximation of Att(Q, K, V) requires nearly quadratic time assuming SETH

Previous Works

Sparse Structure

- ▶ Local Attention (Parmar et al., 18')
- ▶ Sparse Transformer (Child et al., 19')
- ▶ Longformer (Beltagy et al., 20')
- ▶ Reformer (Kitaev et al., 20')
- ▶ Sinkhorn Attention (Tay et al., 20')

Kernel Methods

- ▶ Lambda network (Bello et al., 21')
- ▶ Performer (Choromanski et al., 21')
- ▶ Random Feature Attention (Peng et al., 21')
- ▶ Randomized Attention (Zheng et al., 22')

Low-rank Approximation

- ▶ Linformer (Wang et al., 20')
- ▶ Nystromformer (Xiong et al., 21')
- ▶ Nested Attention (Max et al., 21')

No End-to-End approximation (in some works)

 \triangleright Only approximate matrix $A = \exp(QK^T)$

Would like to:

- ightharpoonup Compute $Att \in \mathbb{R}^{n imes d}$ such that
- \parallel $Att Att(Q, K, V) \parallel_{op}$ is small

These methods do not support causal masking

(Alman & Song 23') High quality (1/poly(n)) entrywise approximation of Att(Q,K,V) is likely impossible in general

HyperAttention

Insu Han (Adobe), Rajesh Jayaram (Google), Amin Karbasi (Yale),

Vahab Mirrokni (Google), David Woodruff (CMU), Amir Zandieh

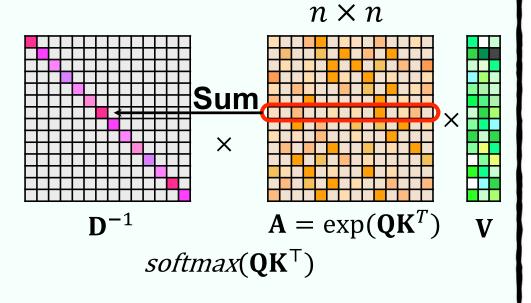
\mathbf{Q} , \mathbf{K} , $\mathbf{V} \in \mathbb{R}^{n \times d}$

Self Attention

1. Approximate

$$D_{i,i} = \sum_{j \in [n]} A_{i,j} = \sum_{j \in [n]} \exp(\langle q_i, k_j \rangle)$$

2. Approximate matrix product $A \cdot V$



memory/runtime: $O(n^2)$

$\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{n \times d}$

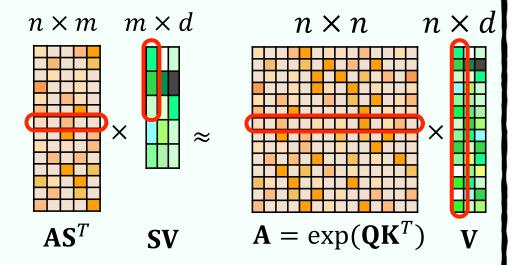
Self Attention

1. Approximate

$$D_{i,i} = \sum_{j \in [n]} A_{i,j} = \sum_{j \in [n]} \exp(\langle q_i, k_j \rangle)$$

2. Compute a row sampling sketch $S \in \mathbb{R}^{m \times n}$ where row i is sampled with probability $\propto ||v_i||_2^2 \to m$

$$\approx srank(softmax(QK^{\mathsf{T}})) \cdot d$$



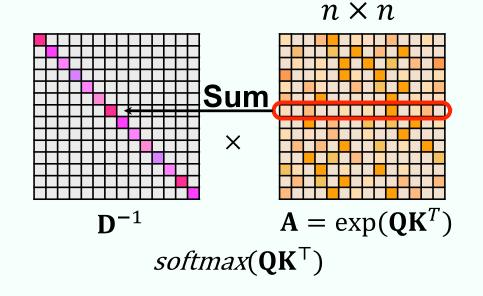
$\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{n \times d}$

Self Attention

1. Approximate

$$D_{i,i} \approx \sum_{j \in [n]} A_{i,j} = \sum_{j \in [n]} \exp(\langle q_i, k_j \rangle)$$

2. Compute a row sampling sketch $S \in \mathbb{R}^{m \times n}$ where row i is sampled with probability $\ll ||v_i||_2^2 \to m \approx srank(softmax(QK^\top)) \cdot d$

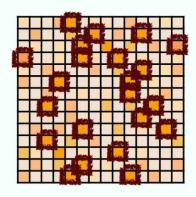




Find 'Heavy' elements of $\mathbf{A} = \exp(\mathbf{Q}\mathbf{K}^T)$

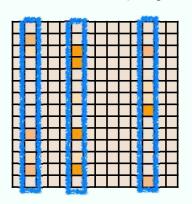
more important

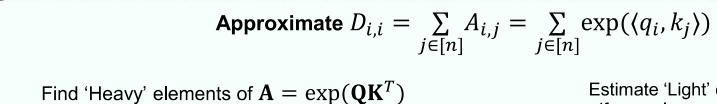
less important



+

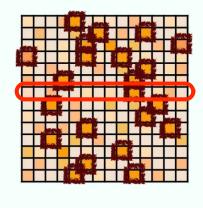
Estimate 'Light' elements of **A** via uniform column sampling





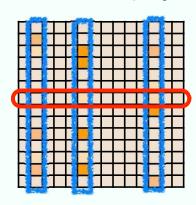
more important

less important



+

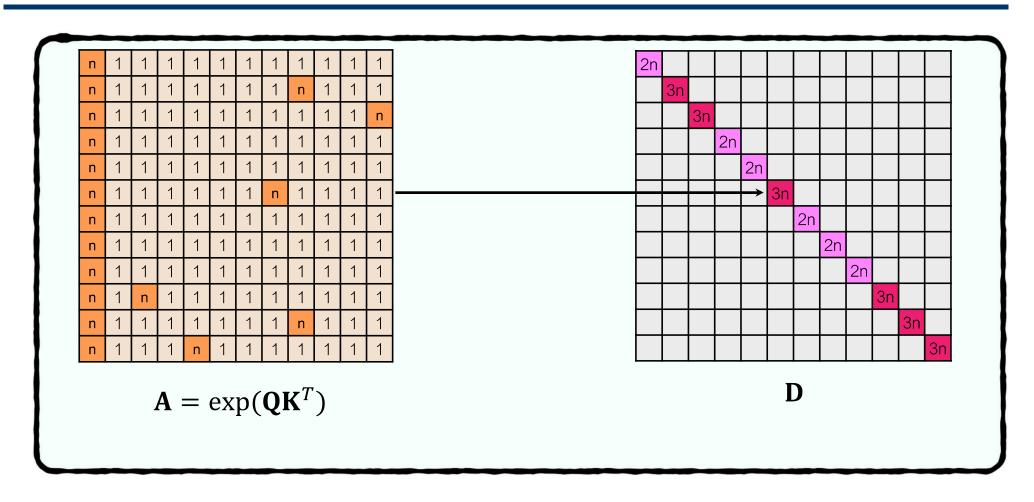
Estimate 'Light' elements of **A** via uniform column sampling

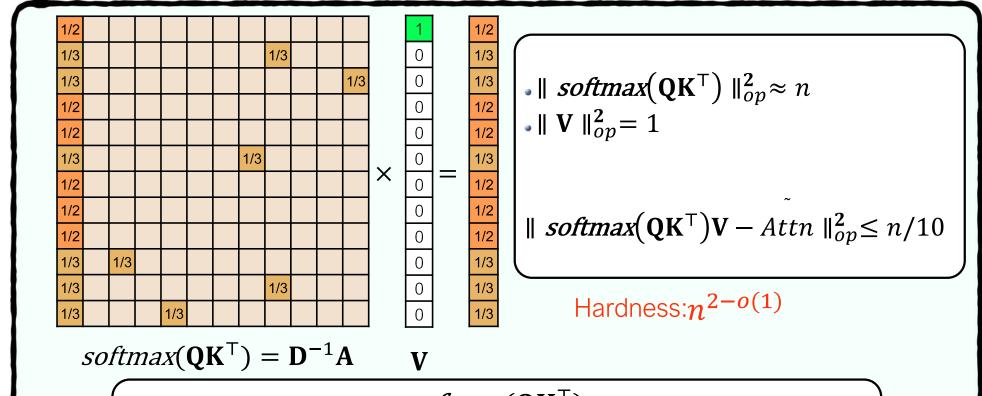


 $D_{i,i}$ =contribution of heavy elements + contribution of light elements

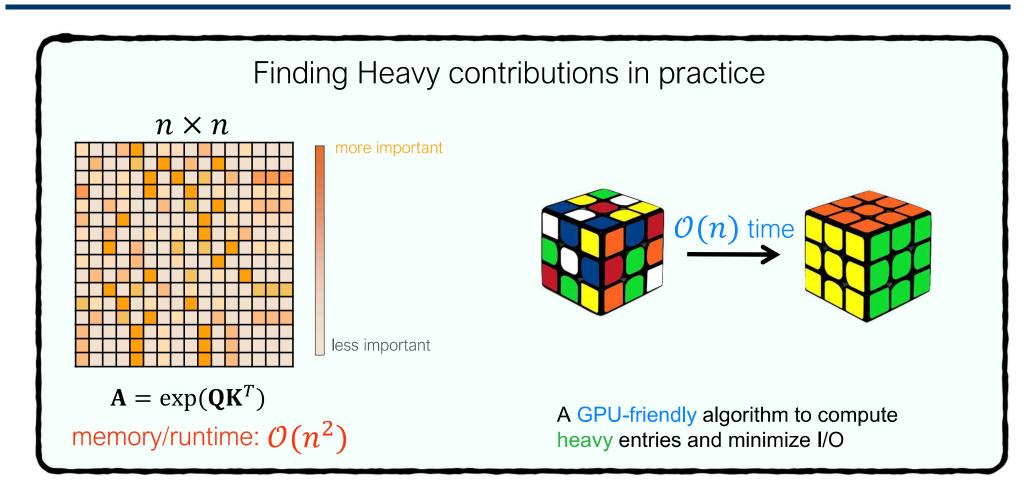
Theorem (informal). If the maximum squared column norm in $softmax(\mathbf{Q}\mathbf{K}^{\top})$ is $\frac{1}{n^{1-o(1)}}$ and the ratio of max and min row sums in $A = exp(\mathbf{Q}\mathbf{K}^{\top})$ after removing heavy elements is $n^{o(1)}$, then Att can be computed in $O(dn^{1+o(1)})$ time with: $\|softmax(\mathbf{Q}\mathbf{K}^{\top})\mathbf{V} - Att\|_{op} \le \varepsilon \|softmax(\mathbf{Q}\mathbf{K}^{\top})\|_{op} \|\mathbf{V}\|_{op}$

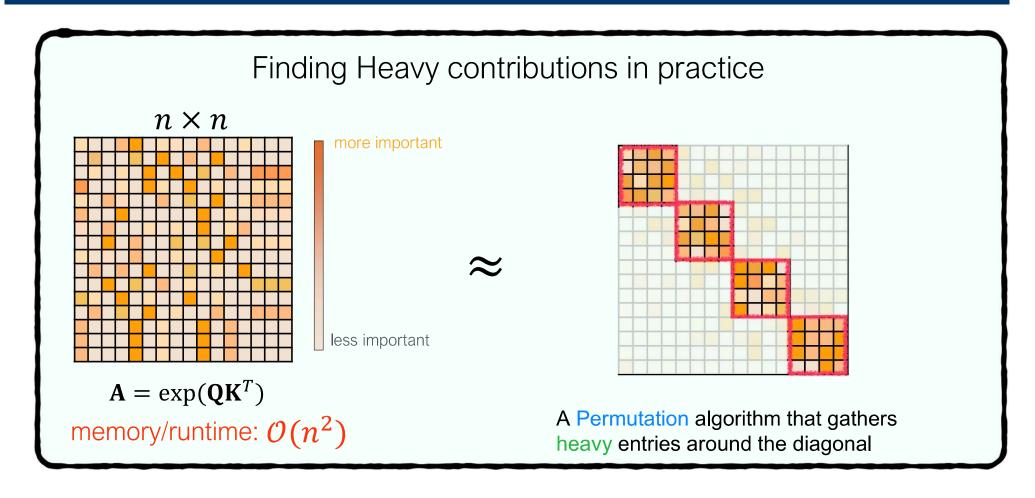
- Column norm bound non-trivial allows for entries as large as $\frac{1}{n^{\frac{1}{2}-o(1)}}$ in $softmax(QK^T)$
- Estimating the contribution of light elements is non-trivial
- Tested assumption of squared column norms in first attention layer of T2T-ViT on ImageNet
- For chatglm2-6b-32k and LongBeach, only the lexicographically first few columns had large norm

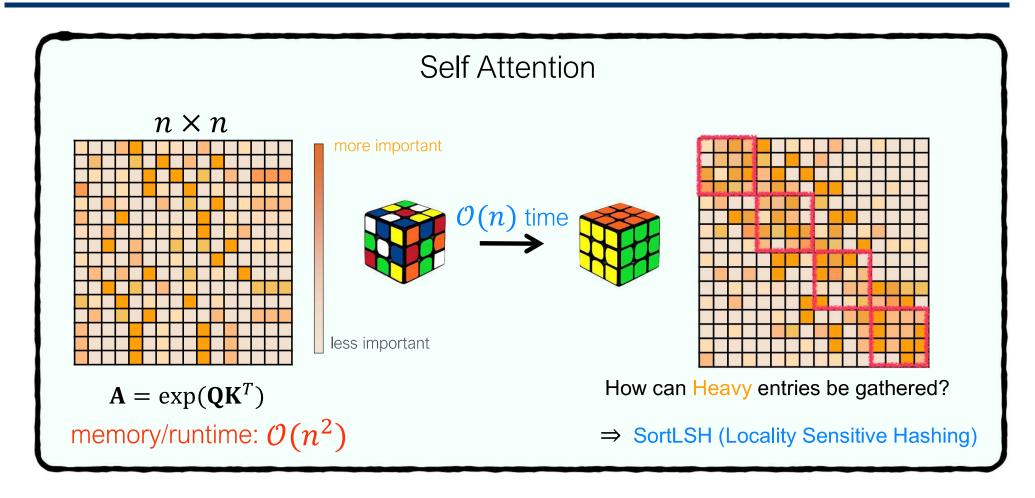


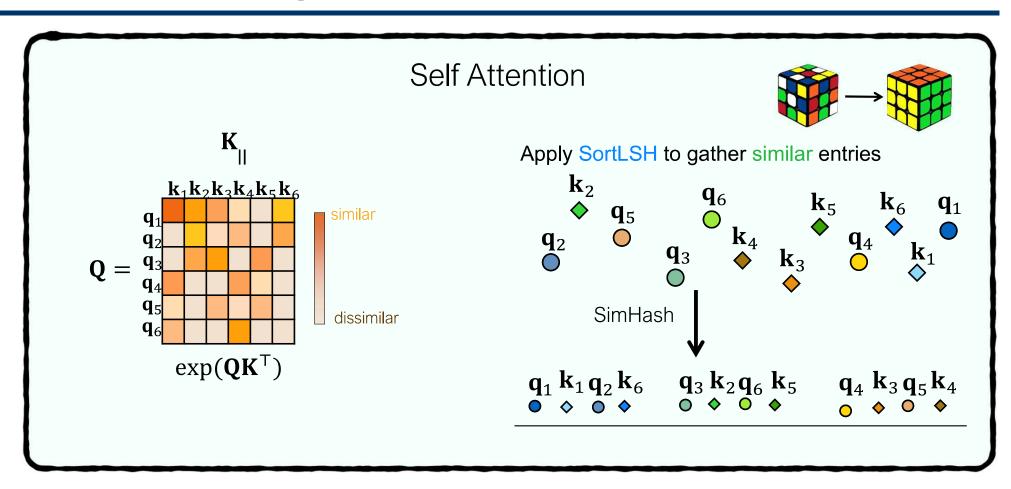


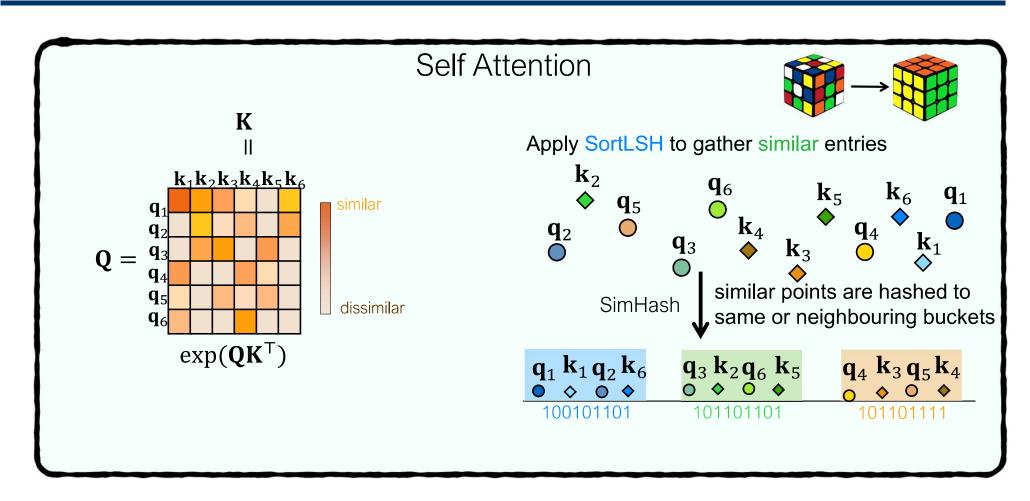
Bounded column norms in $softmax(\mathbf{QK}^{\mathsf{T}})$ avoids this hard instance!

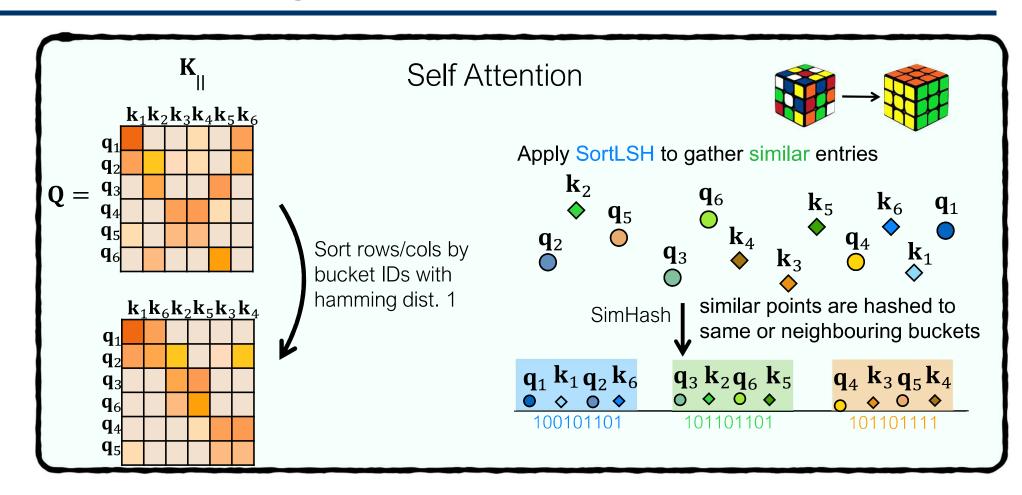


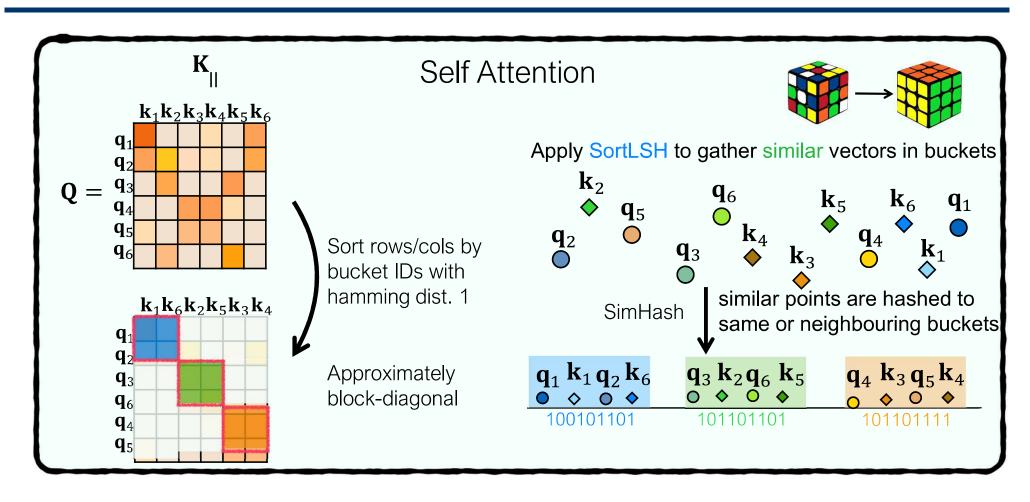


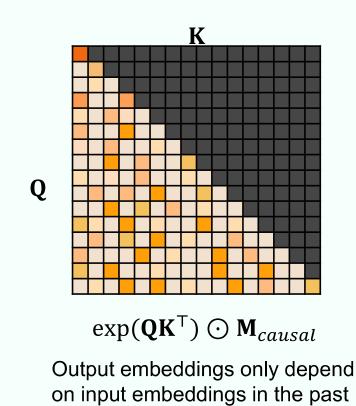




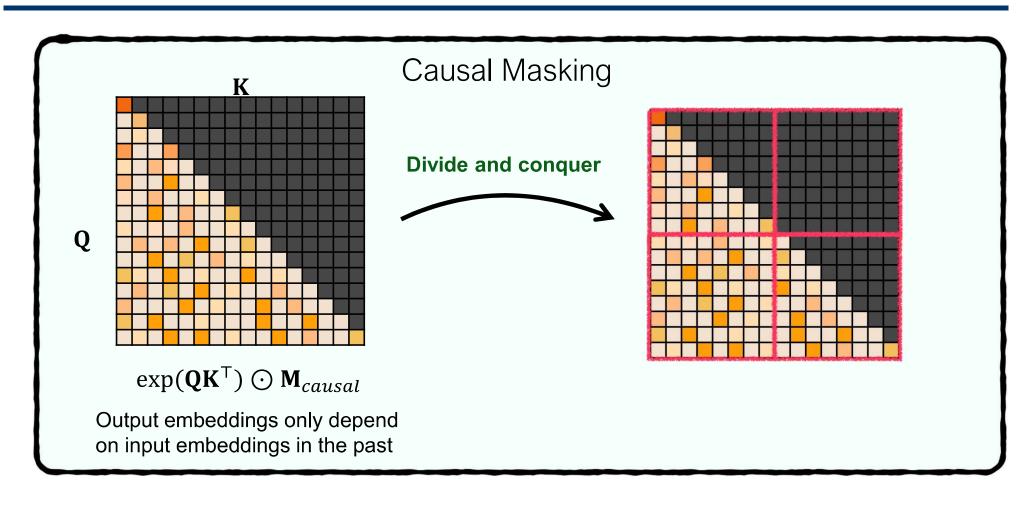


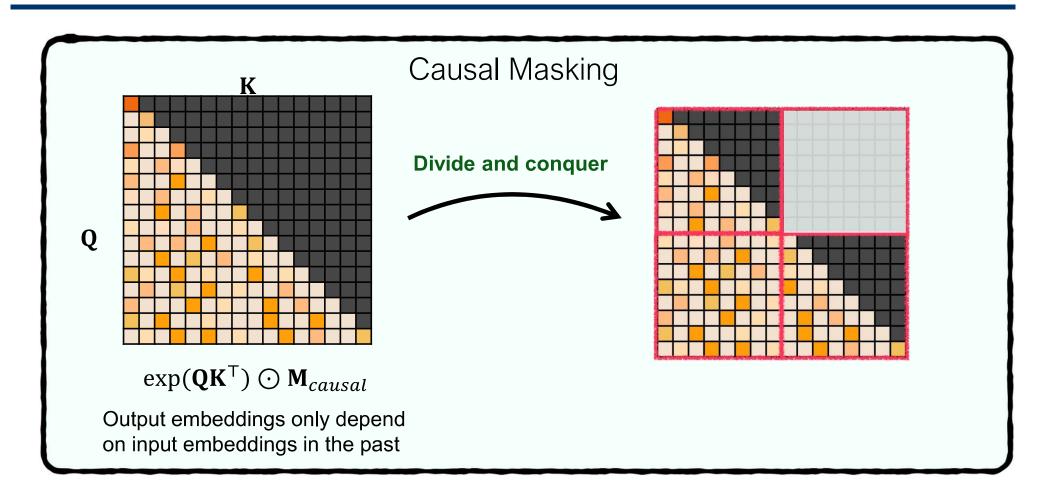






Causal Masking









HyperAttention: Long-context Attention in Near-Linear Time

Insu Han Yale University insu.han@yale.edu

> Vahab Mirrokni Google Research mirrokni@google.com

Rajesh Jayaram Google Research rkjayaram@google.com

> David P. Woodruff CMU, Google Research dwoodruf@cs.cmu.edu

Amin Karbasi Yale University, Google Research amin.karbasi@yale.edu

> Amir Zandieh Independent Researcher amir.zed512@gmail.com

Dialog:

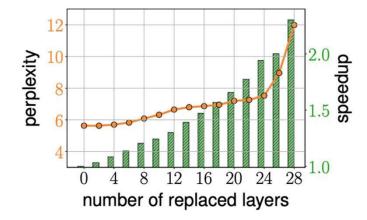
Marisol: it's so sweet he had been waiting

Jackie: we don't know yet when we'll get married but you are all invited ofc

Carlita: PLEASE don't pick June, I'll be in Canada then Eunica: I hate weddings but I'll make an exception

Marisol: can't wait!

LongBench datasets with n = 32768



PolySketchFormer

Praneeth Kacham, Vahab Mirrokni, Peilin Zhong (Google Research)

Generalizations of Softmax Attention

- Let $sim(q,k) \ge 0$ be an arbitrary function that measures similarity between the query q and key k
- Attention mechanism w.r.t sim is

$$o_j = \sum_{i \le j} \frac{sim(q_j, k_i)}{\sum_{i' \le j} sim(q_j, k_{i'})} v_i$$

• Softmax:
$$sim(q, k) \doteq \exp(\langle q, k \rangle)$$

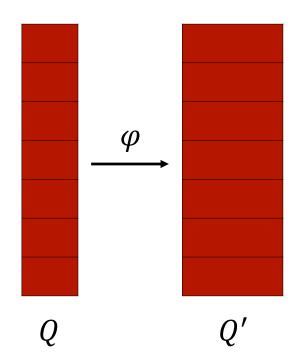
$$\sum_{i \leq j} \frac{\exp(\langle q_j, k_i \rangle)}{\sum_{i' \leq i} \exp(\langle q_i, k_{i'} \rangle)} v_i$$

Kernel View of Attention

- Suppose φ is such that $sim(q, k) = \langle \varphi(q), \varphi(k) \rangle$
- If $Q' = \varphi(Q)$ and $K' = \varphi(K)$, output is

$$D^{-1} \cdot LT(Q' \cdot (K')^{\mathsf{T}}) \cdot V$$

- Here LT is the lower triangular part for the causal setting
- Why write this way?
 - Linear time algorithm for computing $LT(A \cdot B^{\mathsf{T}}) \cdot C$
 - ullet Runtime depends on output dimension of $arphi(\cdot)$
- What about $oldsymbol{arphi}$ for softmax?
 - No finite dimensional feature maps



Previous Work

- ullet Performer (Choromanski et al.,) uses a finite-dimensional map $oldsymbol{arphi}$ to approximate exponential
 - ullet Vectors with larger norms require $oldsymbol{arphi}$ with larger dimension
- Other works consider arbitrary φ instead of first defining $sim(\cdot,\cdot)$
 - $\varphi(x) \doteq elu(x) + 1$ (Katharopoulos et al. '20), $\varphi(x) \doteq relu(x)$
 - Model quality is worse compared to softmax
- Is softmax necessary? Do other functions with similar properties work?
- Consider $sim(q, k) = \langle q, k \rangle^p$ where $p \ge 2$ is an even integer
 - Always ≥ 0
 - Increases as $\langle q, k \rangle$ goes up

7

Feature map for Polynomials

- A finite dimensional φ such that $\langle \varphi(q), \varphi(k) \rangle = \langle q, k \rangle^p$?
 - $\varphi: x \mapsto x^{\otimes p}$
 - If $x \in \mathbb{R}^h$, then $x^{\otimes p} \in \mathbb{R}^{h^p}$
 - $\bullet (x^{\otimes p})_{(i_1,i_2,\dots,i_p)} = x_{i_1} \cdot x_{i_2} \cdot \dots \cdot x_{i_p}$

$$\langle q^{\otimes p}, k^{\otimes p} \rangle = \langle q, k \rangle^p$$

Linear Attention using Polynomials

- Given $Q, K, V \in \mathbb{R}^{n \times h}$
 - ullet Compute $Q^{igotimes p}$ and $K^{igotimes p}$
 - $LT(Q^{\otimes p} \cdot (K^{\otimes p})^{\mathsf{T}}) \cdot V$ in $O(nh^{p+1})$ time
- Typically, h = 64, 128, 256
 - Too expensive even for p=4
- Use sketching to approximate!

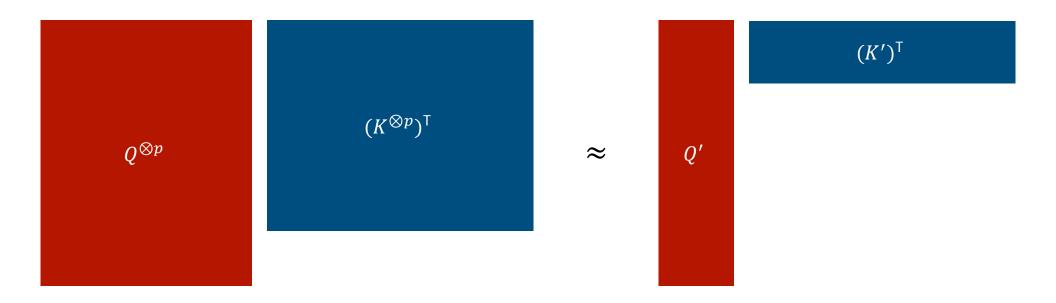
Sketching for Approximate Matrix Multiplication

Want to compute

$$LT(Q^{\otimes p} \cdot (K^{\otimes p})^{\mathsf{T}}) \cdot V$$

- ullet $Q^{igotimes p}$ and $K^{igotimes p}$ can have a large number of columns
- Can we compute matrices Q' and K' such that $Q^{\bigotimes p} \cdot (K^{\bigotimes p})^{\mathsf{T}} \approx Q' \cdot (K')^{\mathsf{T}}$?
 - Ahle et al. '20 give a fast sketch called TensorSketch
 - Can approximate using $LT(Q' \cdot (K')^T) \cdot V$

Sketching for Approximate Matrix Multiplication



Matrix Sketching

- ullet Never have to compute the matrices $Q^{igotimes p}$, $K^{igotimes p}$ and just use Q' and K'
- Can simply compute $LT(Q' \cdot (K')^{\mathsf{T}}) \cdot V$ in linear time
- Does this work?
 - Model training fails to converge
- Non-negativity
 - $Q' \cdot (K')^{\mathsf{T}}$ can have negative entries, whereas entries of $Q^{\bigotimes p} \cdot (K^{\bigotimes p})^{\mathsf{T}}$ are ≥ 0

Solving Issue of Negative Entries

- Consider $Q^{\prime\prime}=(Q^\prime)^{\bigotimes 2}$ and $K^{\prime\prime}=(K^\prime)^{\bigotimes 2}$
 - Q', K' are sketches for degree p/2
- All entries of $Q''\cdot (K'')^{\mathsf{T}}$ are non-negative! They are of the form $\langle q',k'\rangle^2\geq 0$
- Show that if Q' and K' have an approximate matrix product property for degree p/2, then Q'' and K'' have a similar guarantee for degree p
 - $\| Q'' \cdot (K'')^{\mathsf{T}} Q^{\bigotimes p} (K^{\bigotimes p})^{\mathsf{T}} \|_F$ is small
- Compute $LT(Q'' \cdot (K'')^T) \cdot V$
- The model converges!

Other Optimizations

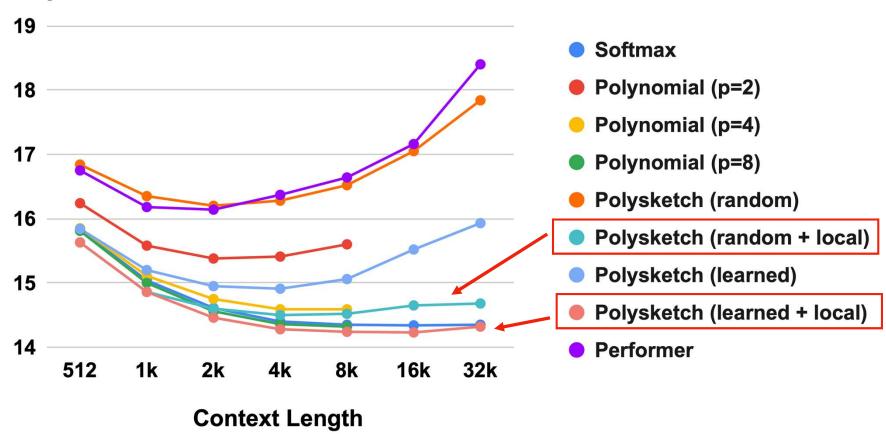
• TensorSketch is a random sketch – instead, treat the sketch as learnable parameters

• When computing $LT(A \cdot B^T) \cdot C$, use block multiplication and cumulative sums

Compute diagonal blocks exactly as such blocks are sensitive to approximation

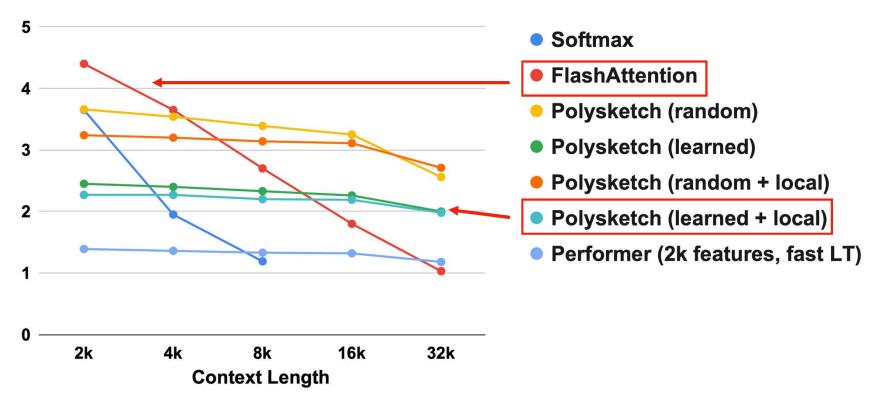
Model Perplexities

Perplexities on Wiki-40B



Training Latencies

Train steps/sec of different mechanisms



Conclusions and Future Work

In practice,

- 1. FlashAttention is an optimized implementation of full softmax attention and is used heavily
- 2. HyperAttention on pretrained models seems to reduce quality too much. Fine-tuning can increase quality but it depends on the hash bucket sizes
- 3. PolySketchFormer has not been tried on very long contexts. Hyperattention allows more state to be maintained for long contexts by increasing the number of heavy entries stored

Recently [Kannan, Bhattacharyya, Kacham, W] use tools from randomized linear algebra to show for a large class of loss functions, there is a small subset of keys so that any heavy attention score involves a key from that subset. Still being tested!