

Recent Efficiency Improvements to Transformers

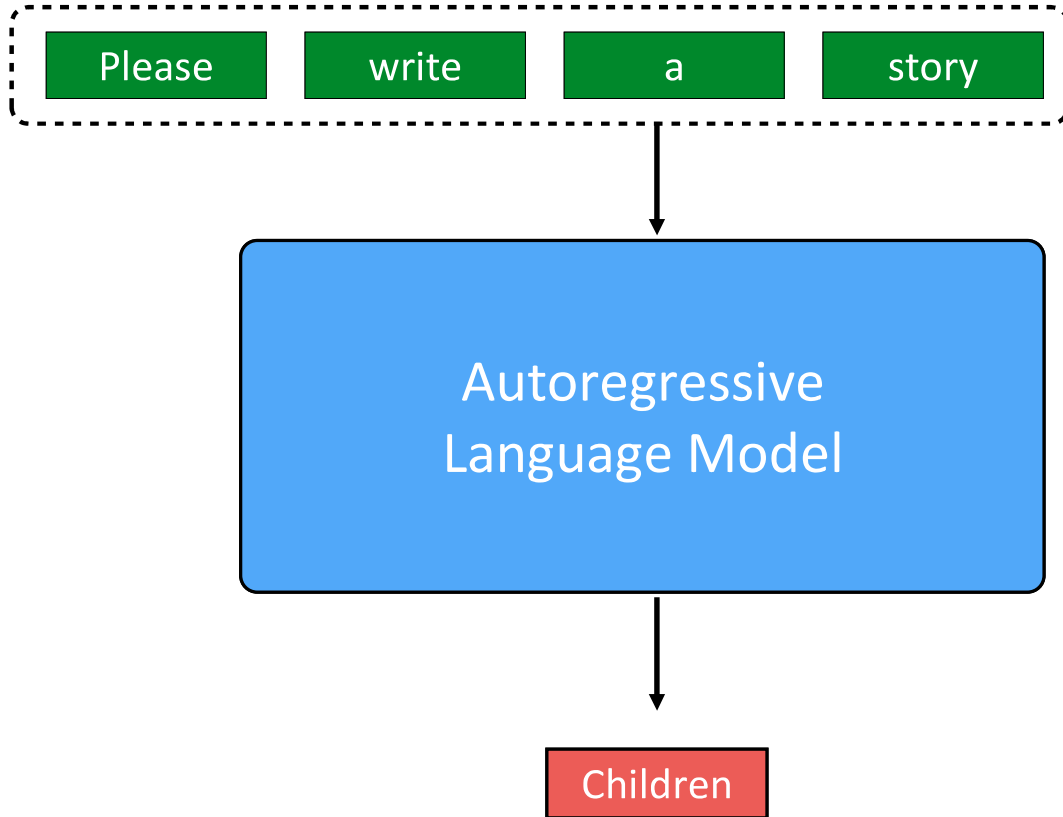
David Woodruff

Carnegie Mellon University / Google

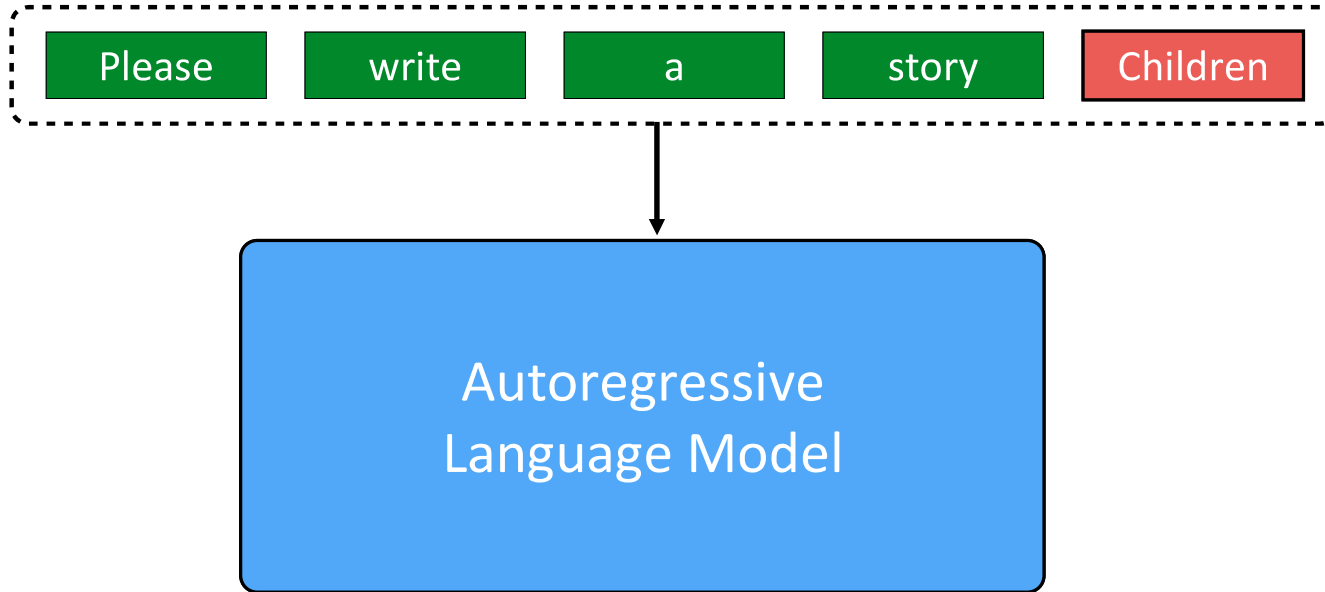
Outline

1. Background on Transformers and time complexity
2. HyperAttention
3. PolySketchFormer
4. Conclusions and Recent Work

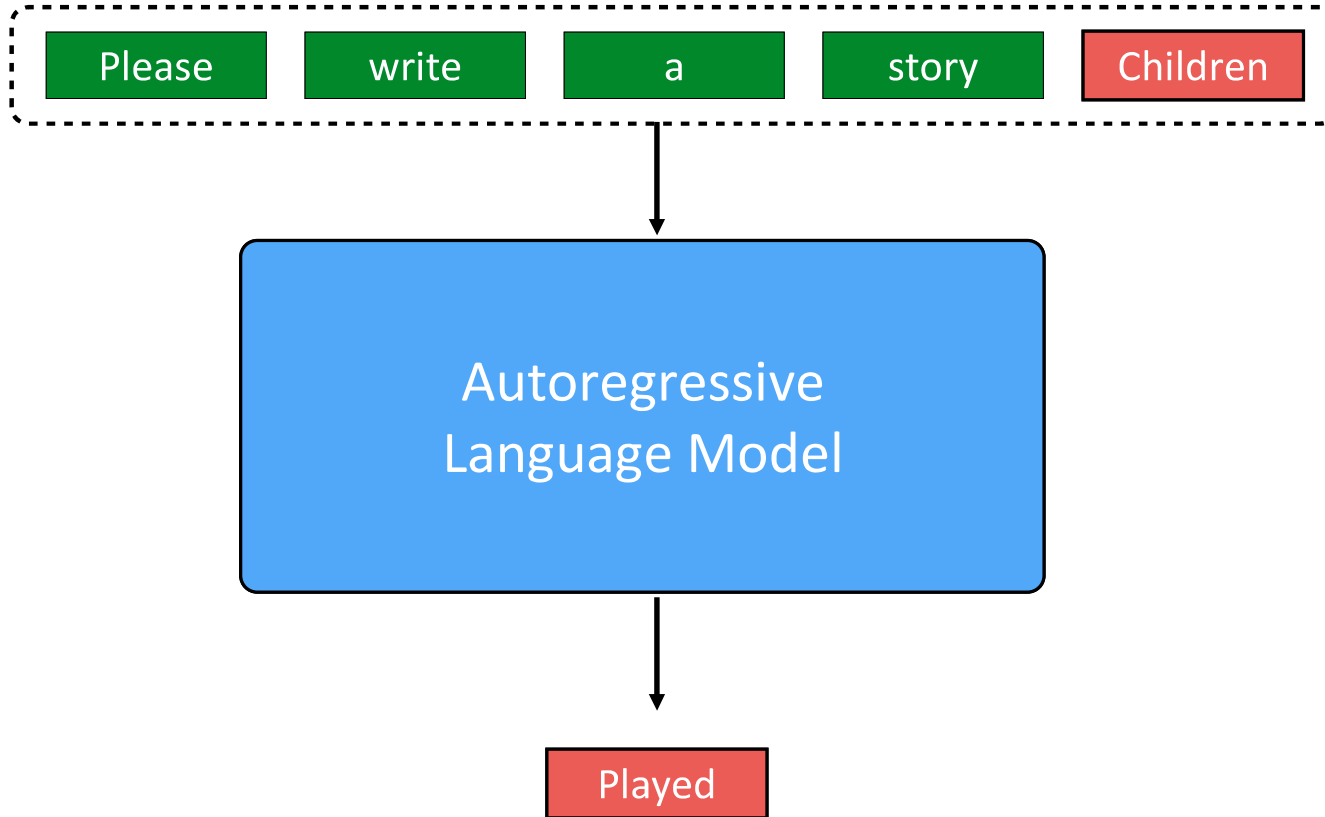
Autoregressive Models



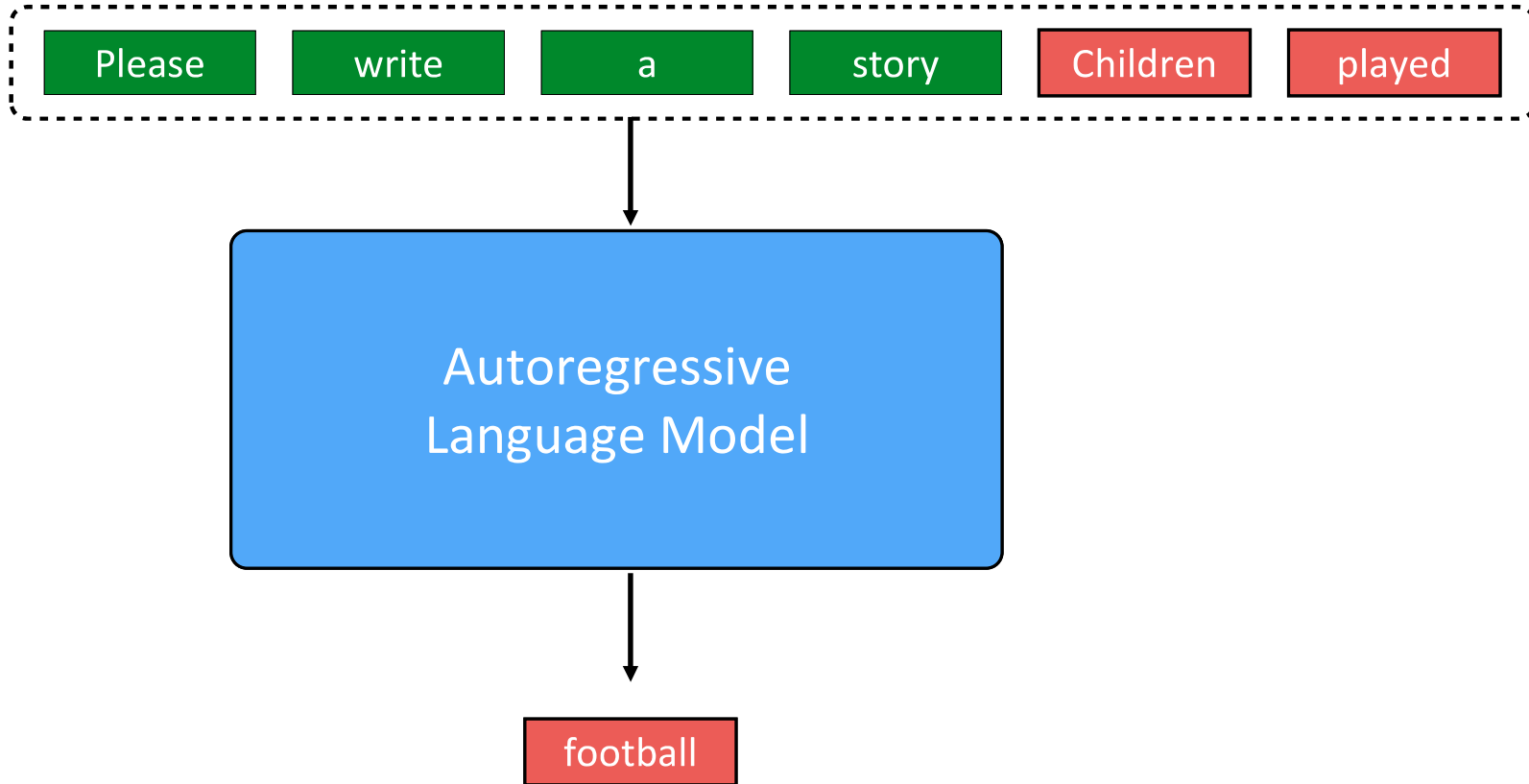
Autoregressive Models



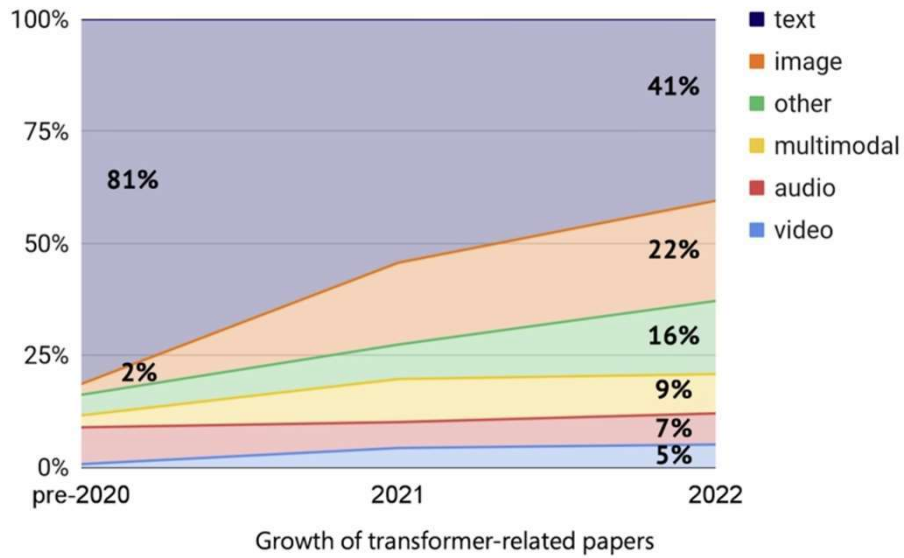
Autoregressive Models



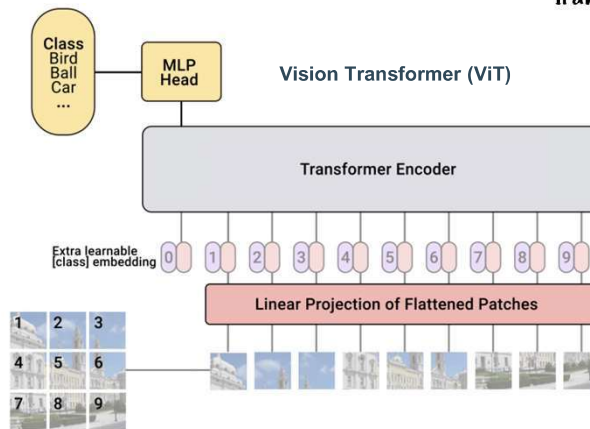
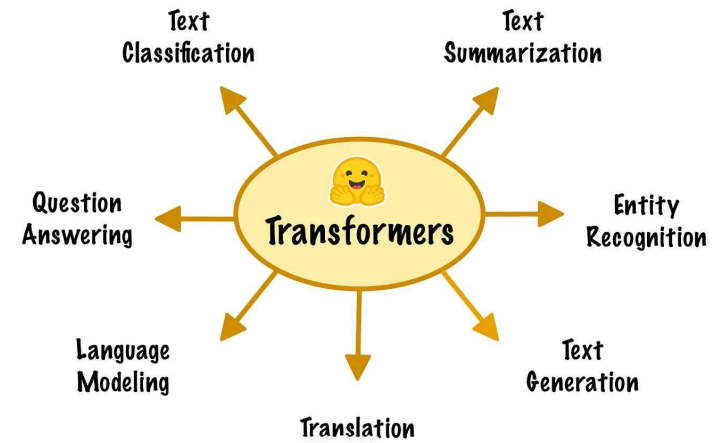
Autoregressive Models



Transformers

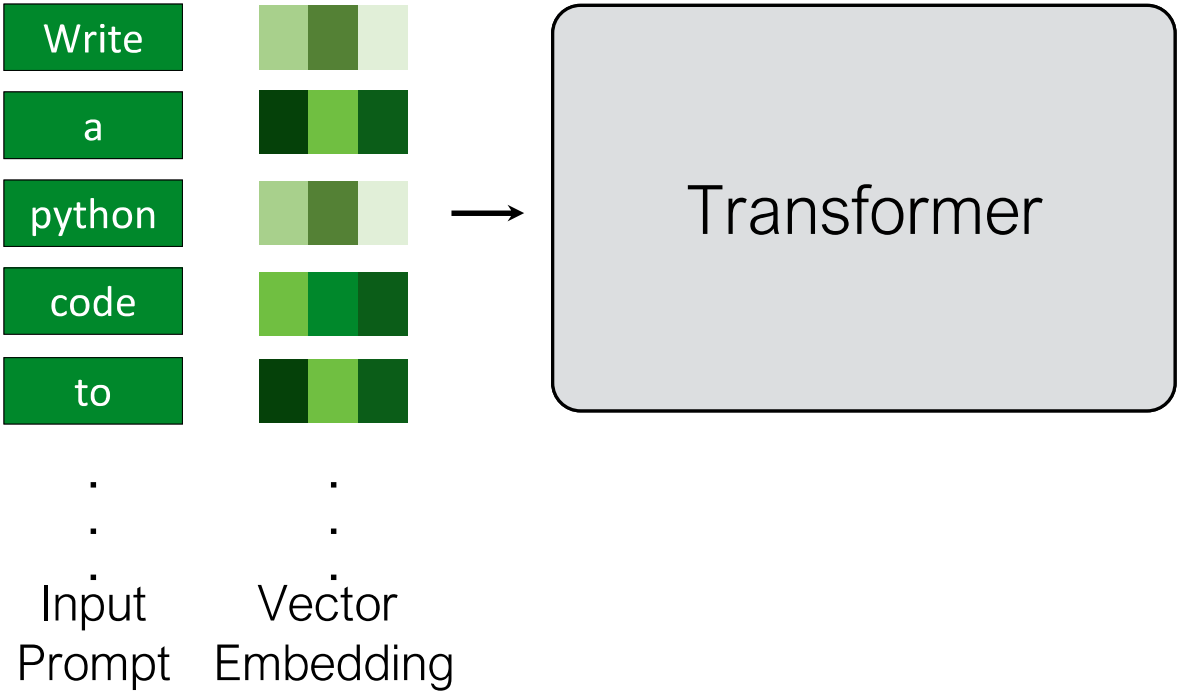


Transformer (Vaswani et al., 17')



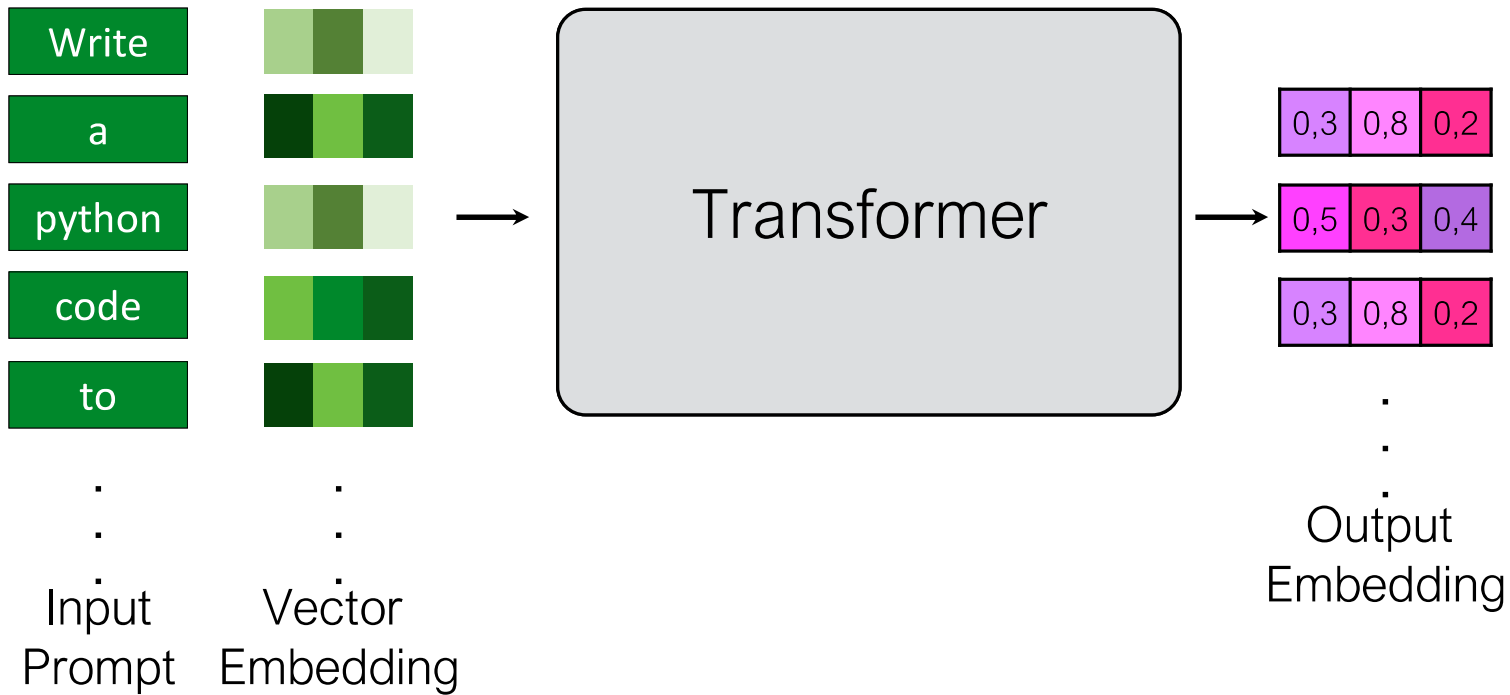
Transformers

I Write a python code to generate webpage



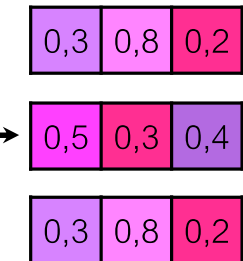
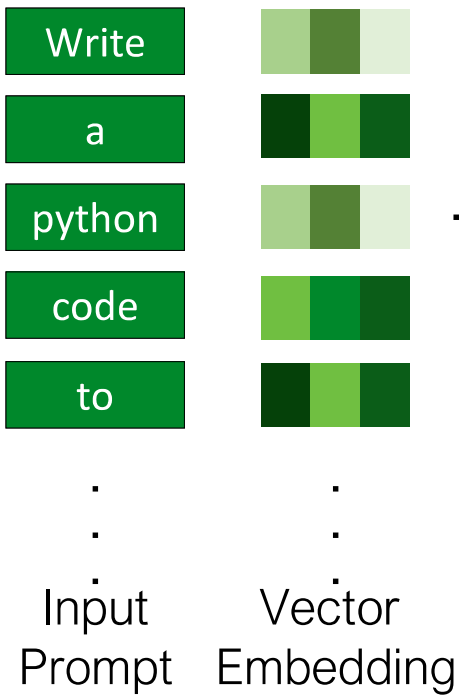
Transformers

I Write a python code to generate webpage

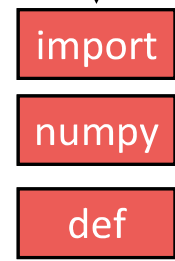
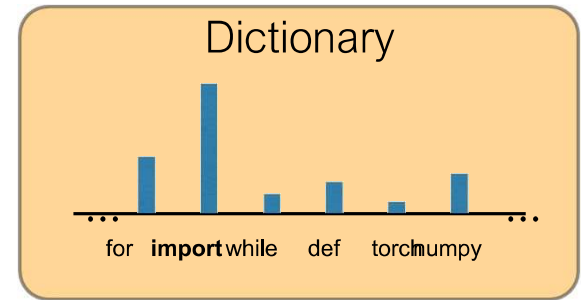


Transformers

I Write a python code to generate webpage

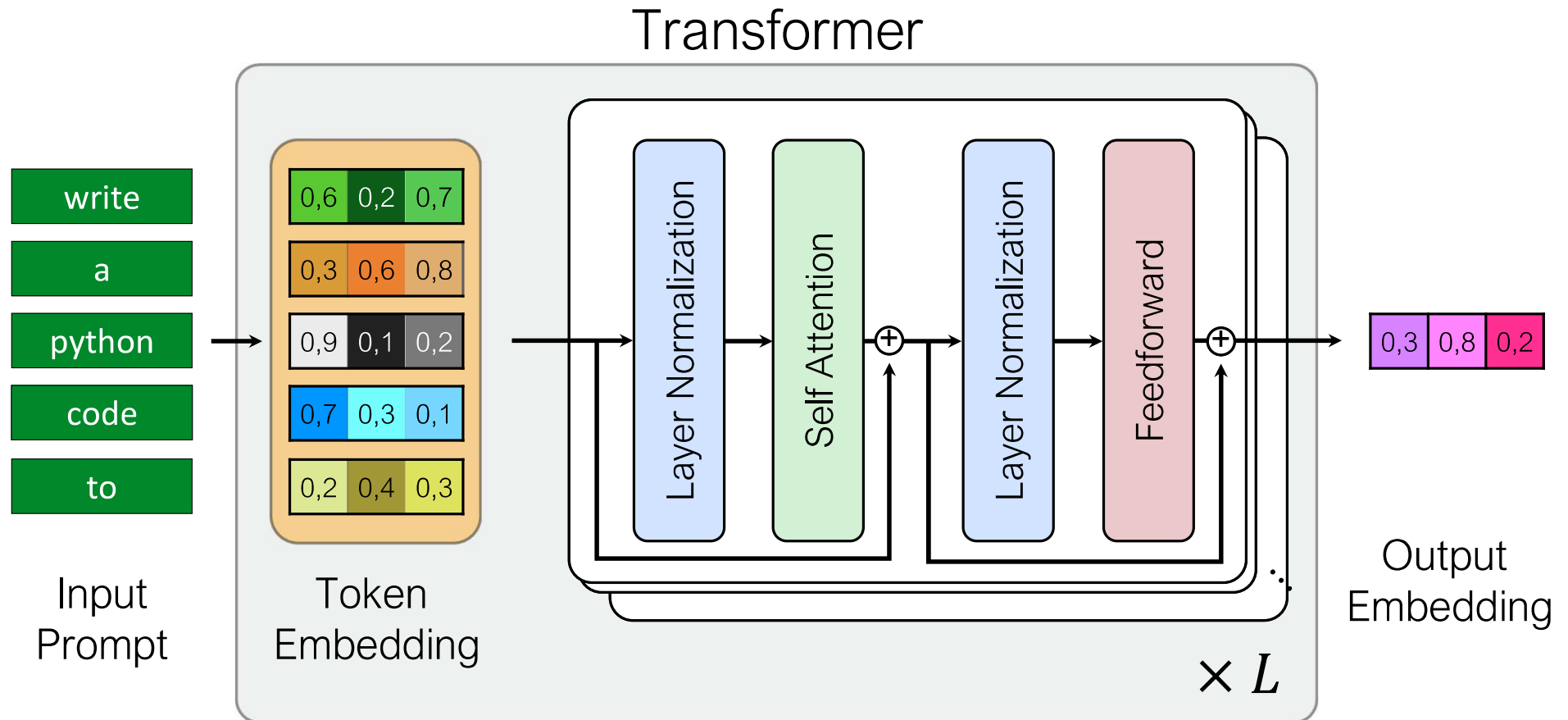


Output Embedding

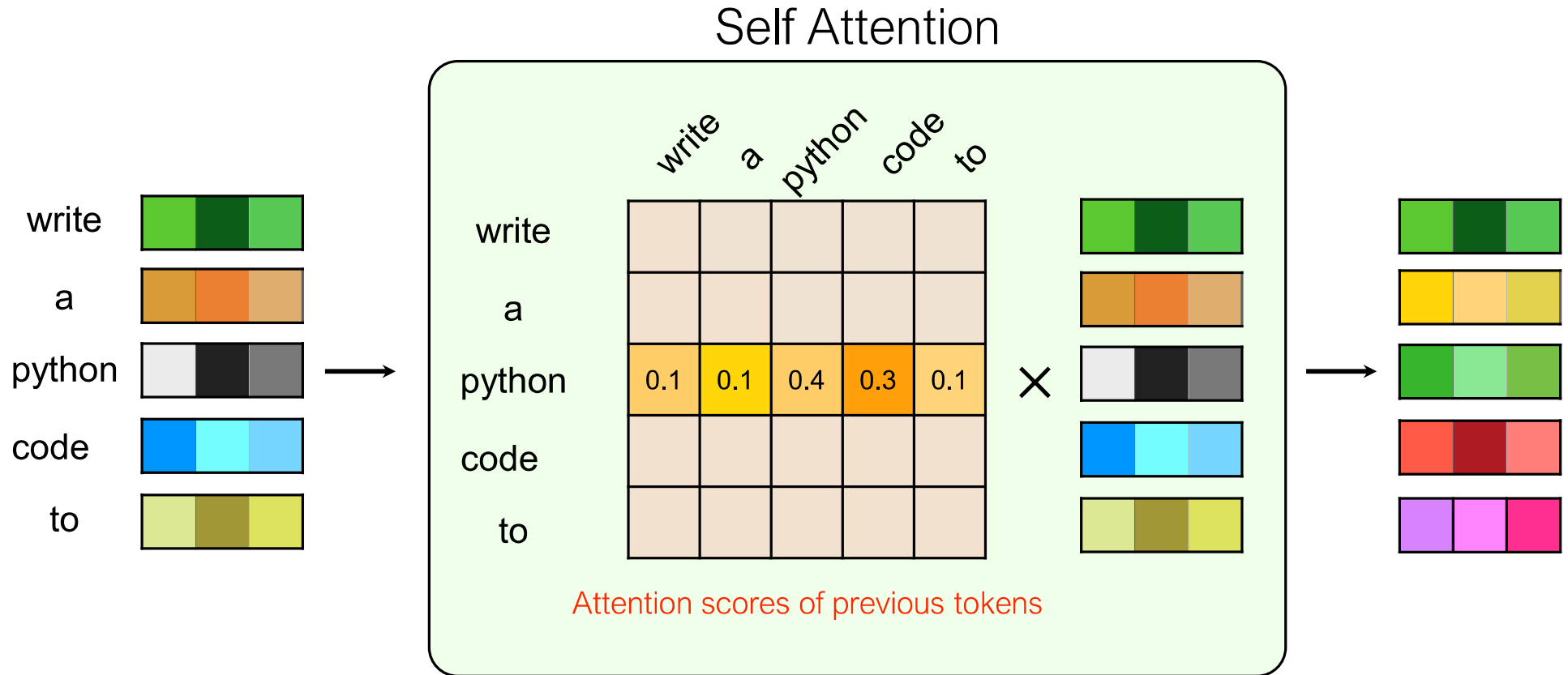


Token Generation

Transformers



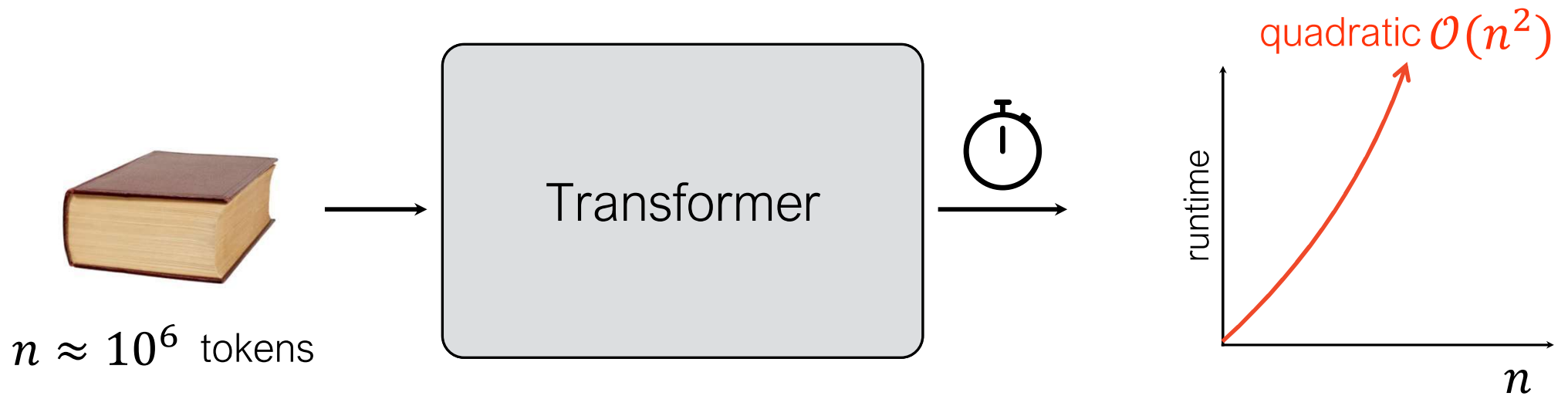
Transformers



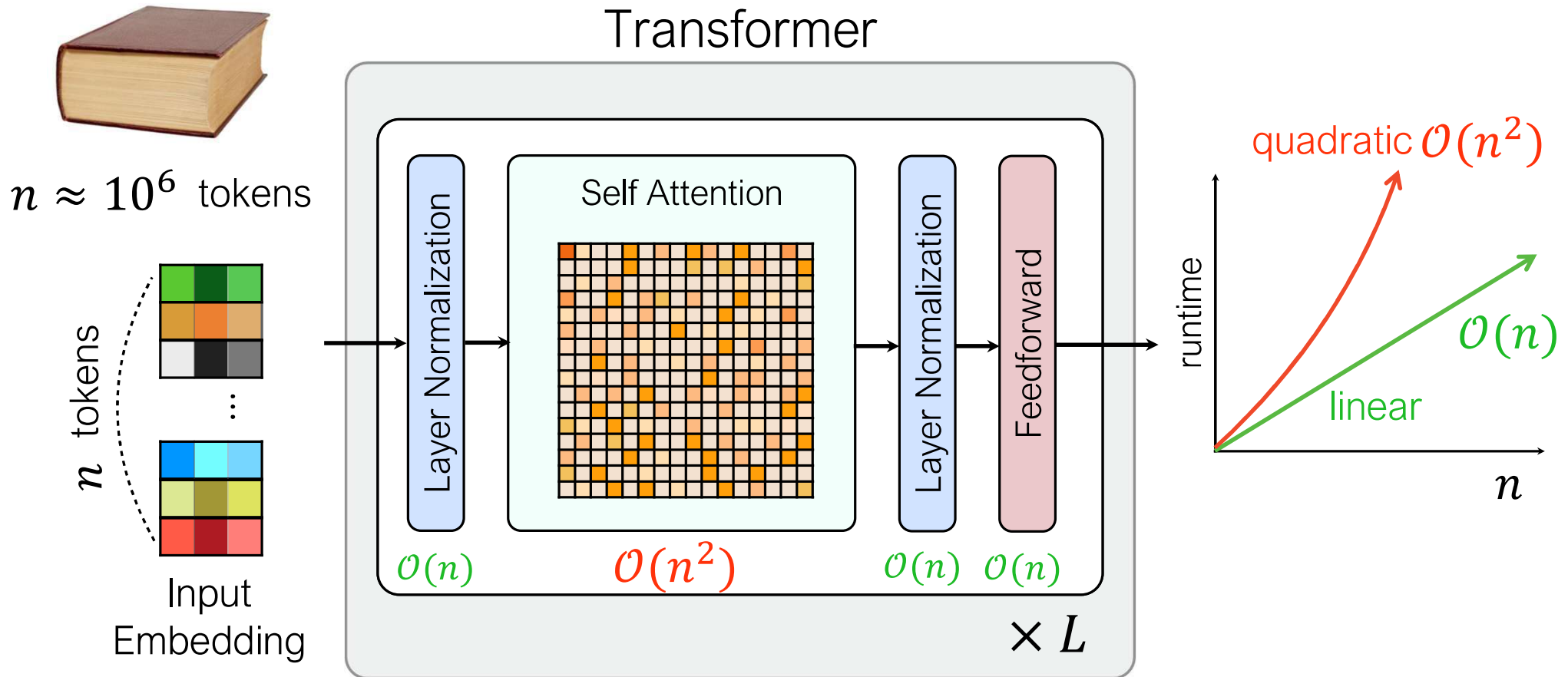
Algorithmic Acceleration



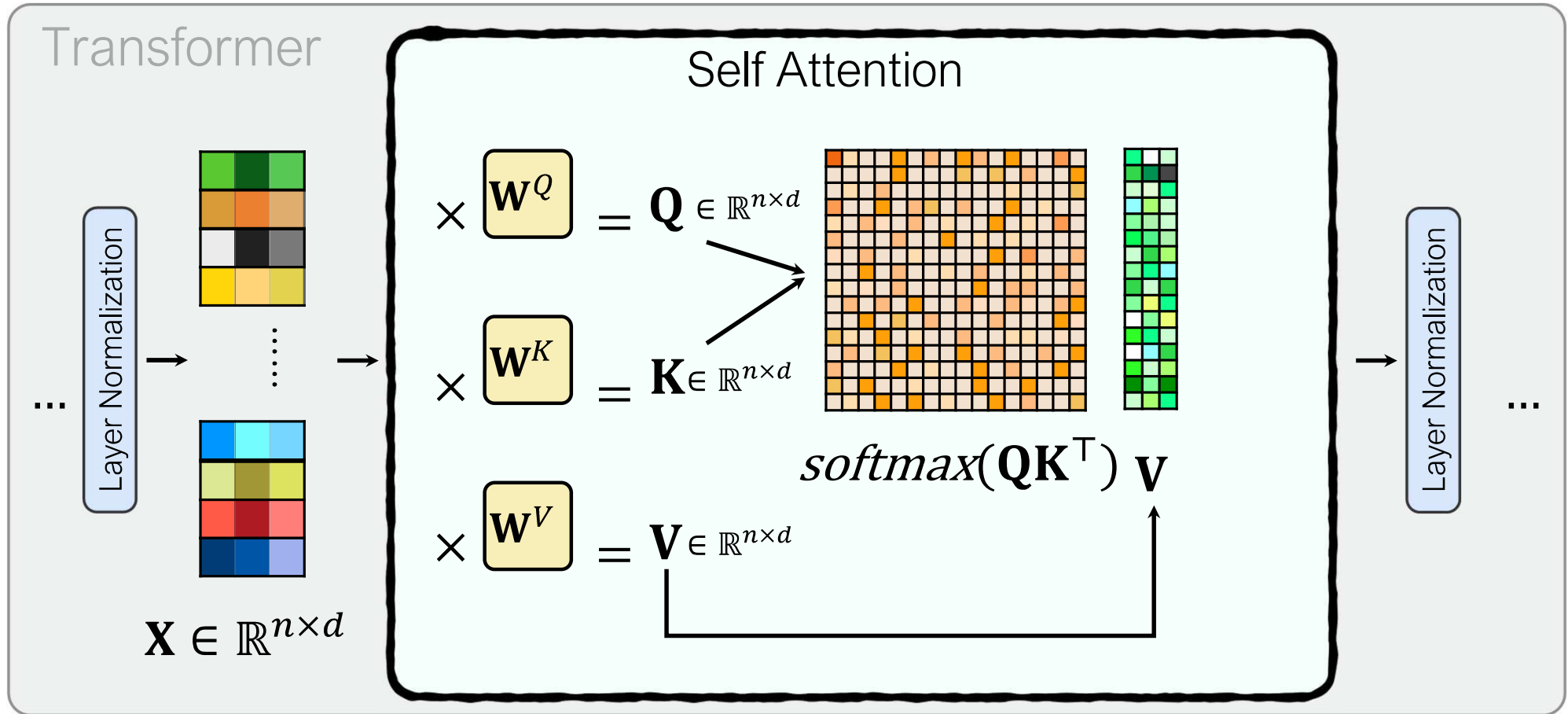
Algorithmic Acceleration



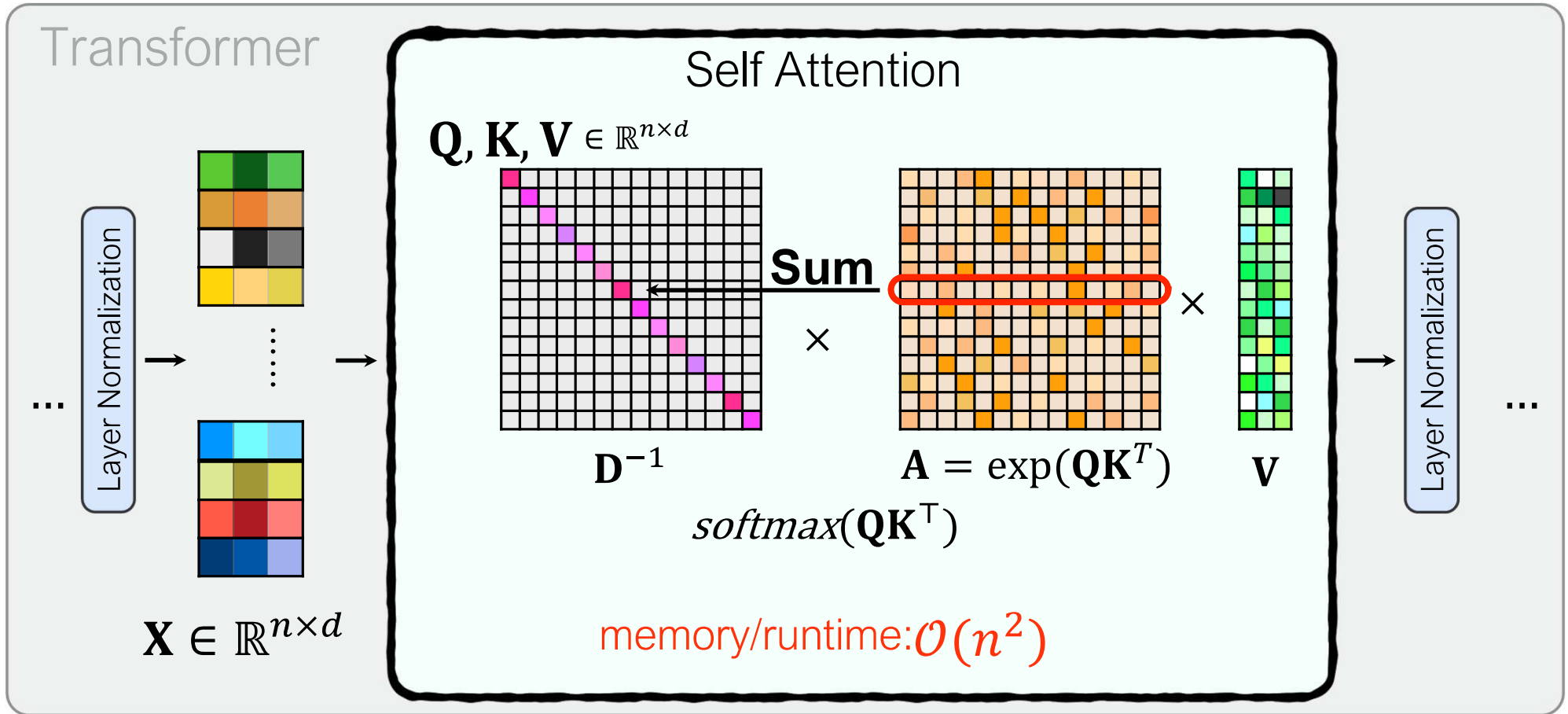
Algorithmic Acceleration



Algorithmic Acceleration



Algorithmic Acceleration



Previous Work

- **Sparse Structure**

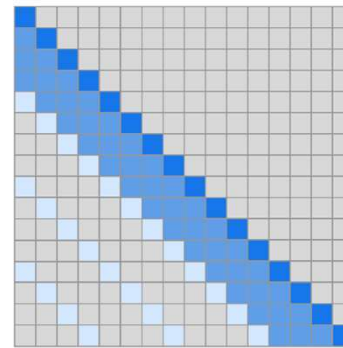
- ▶ Local Attention (Parmar et al., 18')
- ▶ Sparse Transformer (Child et al., 19')
- ▶ Longformer (Beltagy et al., 20')
- ▶ Reformer (Kitaev et al., 20')
- ▶ Sinkhorn Attention (Tay et al., 20')

- **Kernel Methods**

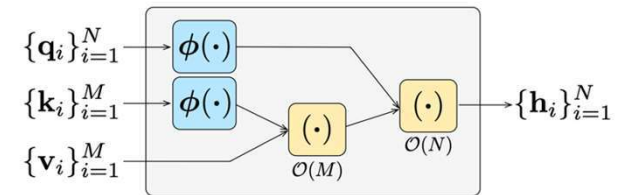
- ▶ Lambda network (Bello et al., 21')
- ▶ Performer (Choromanski et al., 21')
- ▶ Random Feature Attention (Peng et al., 21')
- ▶ Randomized Attention (Zheng et al., 22')

- **Low-rank Approximation**

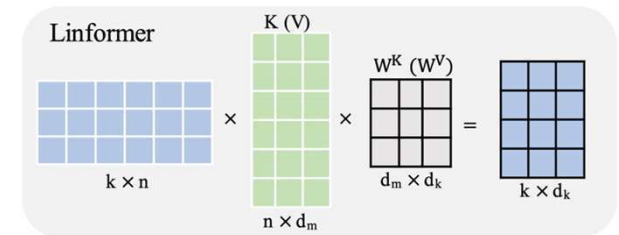
- ▶ Linformer (Wang et al., 20')
- ▶ Nystromformer (Xiong et al., 21')
- ▶ Nested Attention (Max et al., 21')



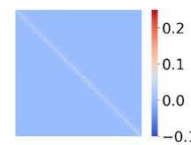
(b) Sparse Transformer (strided)



(b) Random feature attention.

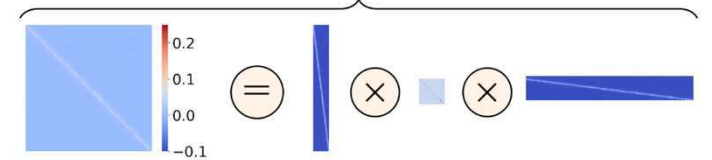


softmax



\approx

Nyström approximation



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)

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

Would like to:

- Compute $\tilde{Att} \in \mathbb{R}^{n \times d}$ such that
- $\| \tilde{Att} - Att(Q, K, V) \|_{op}$ 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

Algorithmic Acceleration

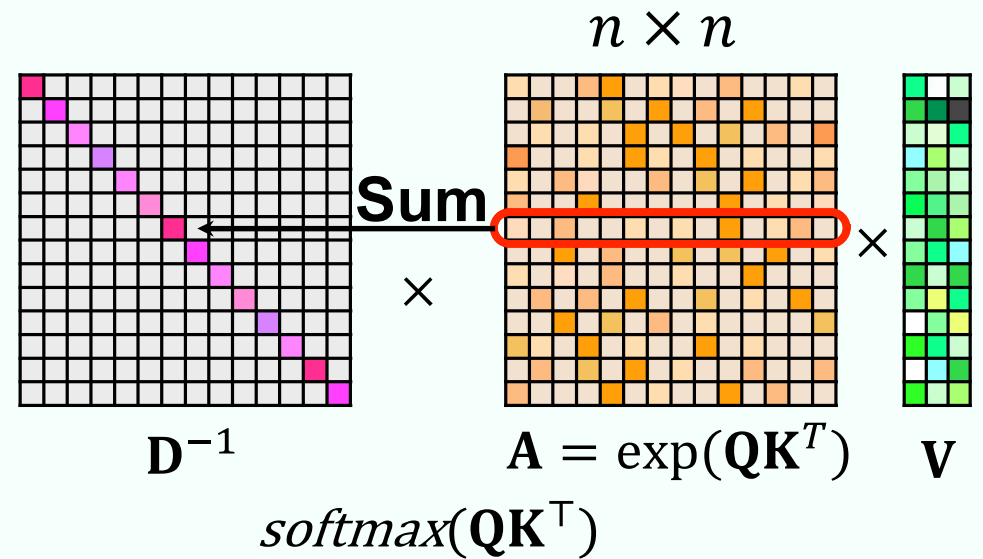
$\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{n \times d}$

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: $\mathcal{O}(n^2)$

Algorithmic Acceleration

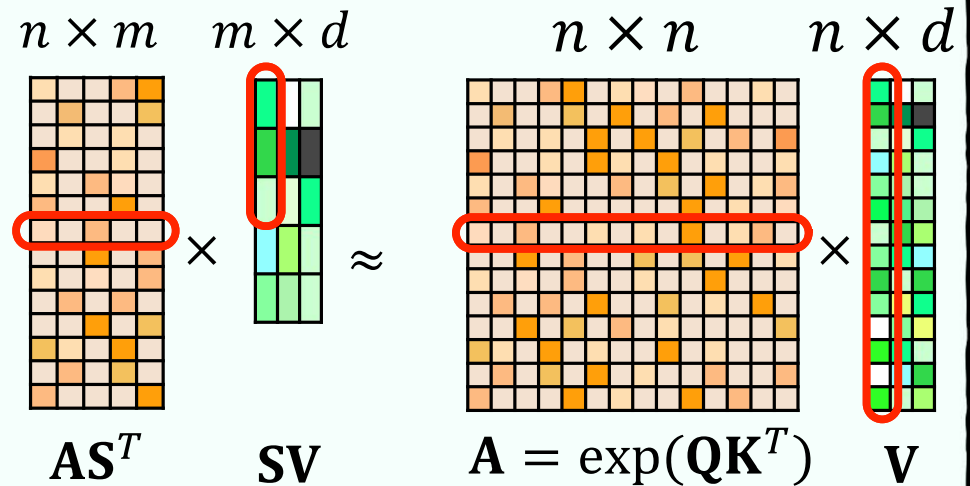
$Q, K, V \in \mathbb{R}^{n \times d}$

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 $\propto \|v_i\|_2^2 \rightarrow m \approx \text{srank}(\text{softmax}(QK^T)) \cdot d$



Algorithmic Acceleration

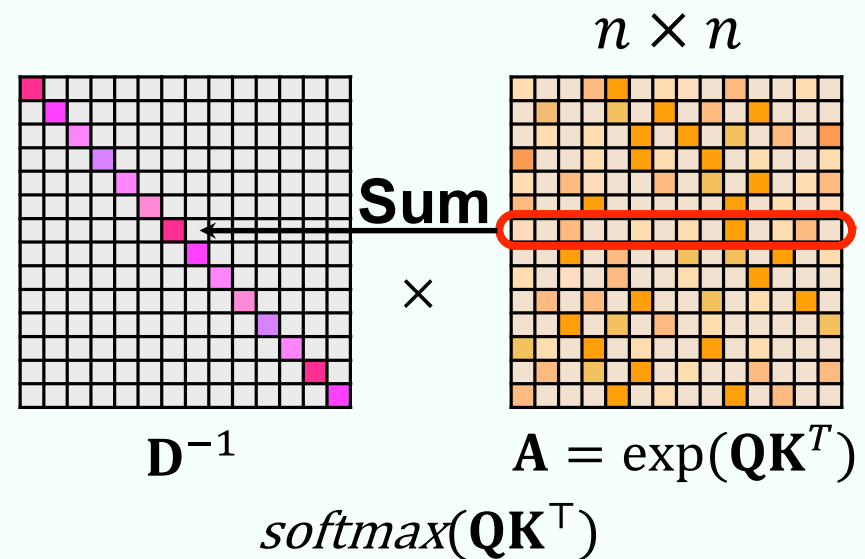
$\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{n \times d}$

Self Attention

1. Approximate

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

2. Compute a row sampling sketch $S \in \mathbb{R}^{m \times n}$ where row i is sampled with probability $\propto \|v_i\|_2^2 \rightarrow m \approx \text{srnk}(\text{softmax}(\mathbf{Q}\mathbf{K}^T)) \cdot d$

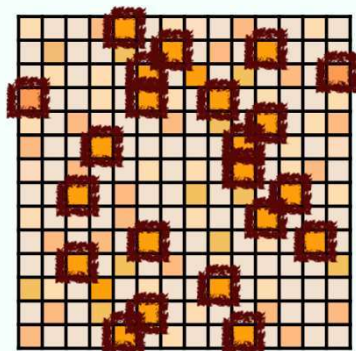


Algorithmic Acceleration

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

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

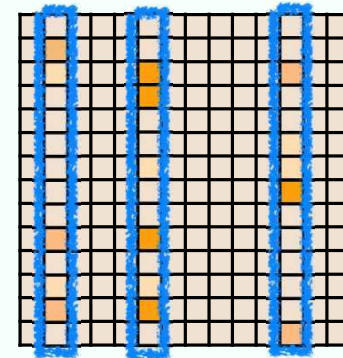
more important



less important

+

Estimate 'Light' elements of \mathbf{A} via uniform column sampling

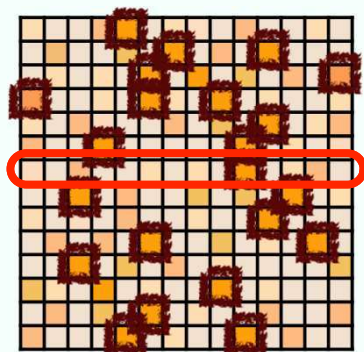


Algorithmic Acceleration

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

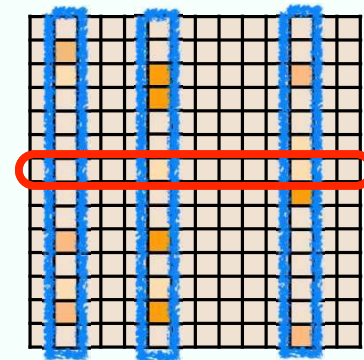
Find 'Heavy' elements of $\mathbf{A} = \exp(\mathbf{QK}^T)$

more important



less important

Estimate 'Light' elements of \mathbf{A} via uniform column sampling



+

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

Algorithmic Acceleration

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

$$\| \mathit{softmax}(\mathbf{QK}^\top)\mathbf{V} - \mathit{Att} \|_{op} \leq \varepsilon \| \mathit{softmax}(\mathbf{QK}^\top) \|_{op} \| \mathbf{V} \|_{op}$$

- Column norm bound non-trivial – allows for entries as large as $\frac{1}{n^{2-o(1)}}$ in $\mathit{softmax}(\mathbf{QK}^\top)$
- 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

Algorithmic Acceleration

n	1	1	1	1	1	1	1	1	1	1	1
n	1	1	1	1	1	1	1	n	1	1	1
n	1	1	1	1	1	1	1	1	1	1	n
n	1	1	1	1	1	1	1	1	1	1	1
n	1	1	1	1	1	1	1	1	1	1	1
n	1	1	1	1	1	1	n	1	1	1	1
n	1	1	1	1	1	1	1	1	1	1	1
n	1	1	1	1	1	1	1	1	1	1	1
n	1	1	1	1	1	1	1	1	1	1	1
n	1	n	1	1	1	1	1	1	1	1	1
n	1	1	1	1	1	1	1	n	1	1	1
n	1	1	1	n	1	1	1	1	1	1	1

$$A = \exp(QK^T)$$

2n											
	3n										
		3n									
			2n								
				2n							
					3n						
						2n					
							2n				
								2n			
									3n		
										3n	
											3n

D

Algorithmic Acceleration

1/2										
1/3					1/3					
1/3									1/3	
1/2										
1/2										
1/3					1/3					
1/2										
1/2										
1/2										
1/3	1/3									
1/3					1/3					
1/3			1/3							

1
0
0
0
0
0
0
0
0
0
0
0
0

1/2
1/3
1/3
1/2
1/2
1/3
1/2
1/2
1/2
1/3
1/3
1/3

× =

$softmax(QK^T) = D^{-1}A \quad V$

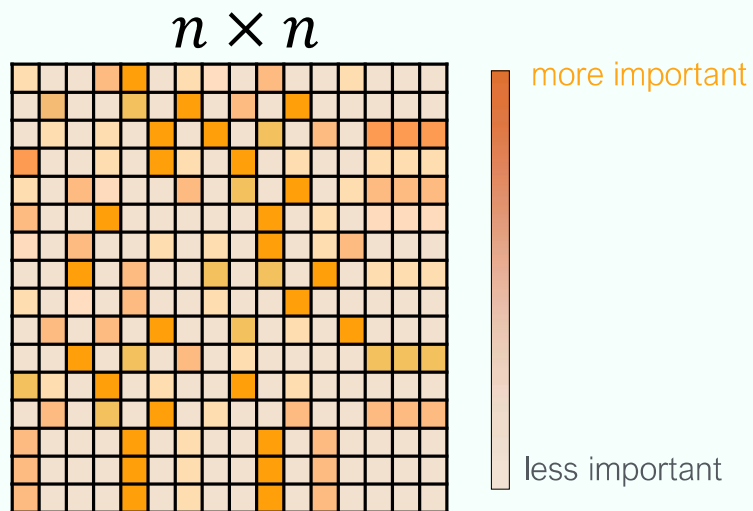
- $\| softmax(QK^T) \|_{op}^2 \approx n$
 - $\| V \|_{op}^2 = 1$
- $\| softmax(QK^T)V - Attn \|_{op}^2 \leq n/10$

Hardness: $n^{2-o(1)}$

Bounded column norms in $softmax(QK^T)$ avoids this hard instance!

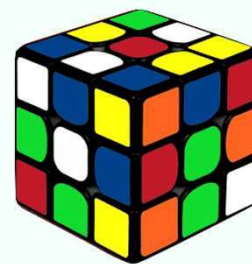
Algorithmic Acceleration

Finding Heavy contributions in practice

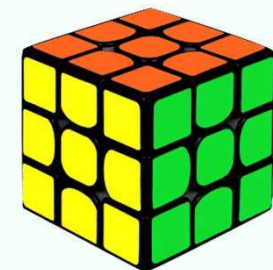
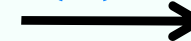


$$\mathbf{A} = \exp(\mathbf{Q}\mathbf{K}^T)$$

memory/runtime: $\mathcal{O}(n^2)$



$\mathcal{O}(n)$ time

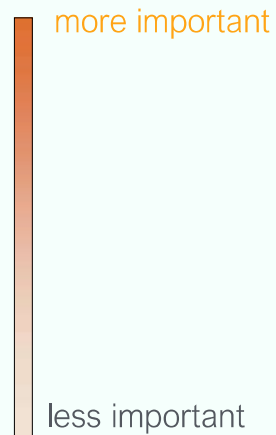
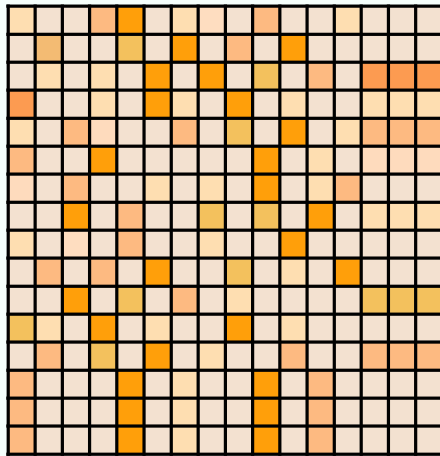


A GPU-friendly algorithm to compute heavy entries and minimize I/O

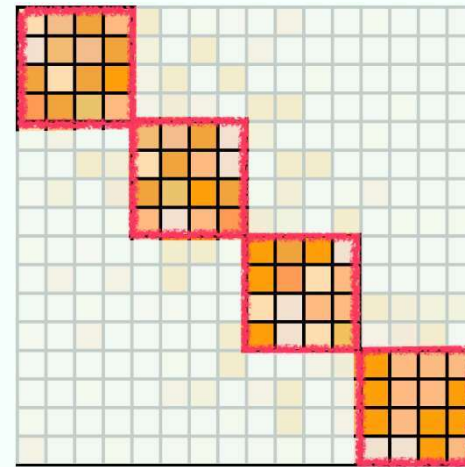
Algorithmic Acceleration

Finding Heavy contributions in practice

$n \times n$



\approx



$$\mathbf{A} = \exp(\mathbf{Q}\mathbf{K}^T)$$

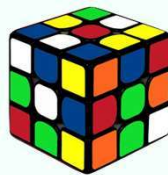
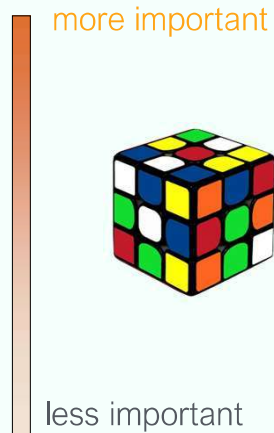
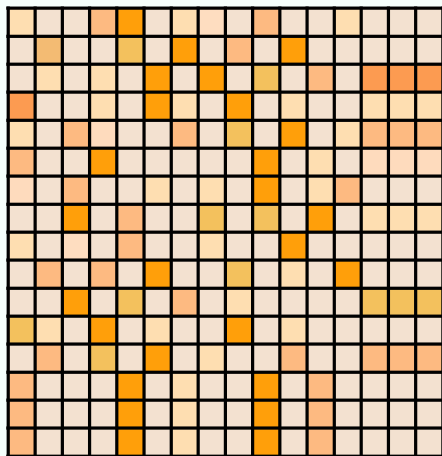
memory/runtime: $\mathcal{O}(n^2)$

A **Permutation** algorithm that gathers **heavy** entries around the diagonal

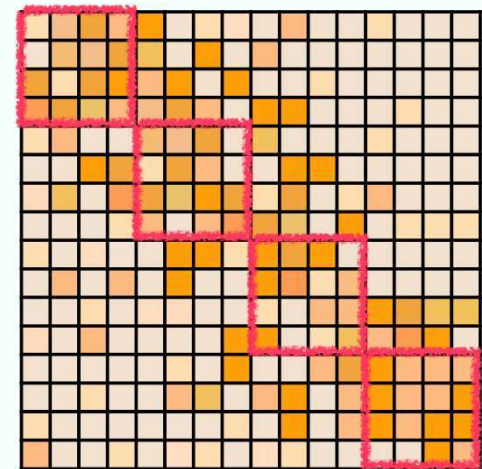
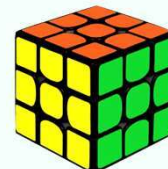
Algorithmic Acceleration

Self Attention

$n \times n$



$O(n)$ time



$$\mathbf{A} = \exp(\mathbf{Q}\mathbf{K}^T)$$

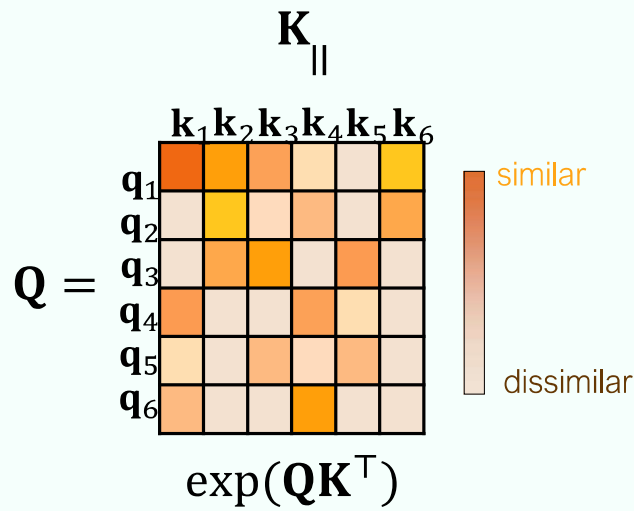
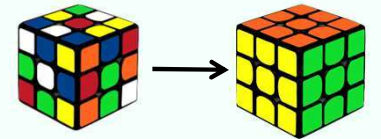
memory/runtime: $O(n^2)$

How can **Heavy** entries be gathered?

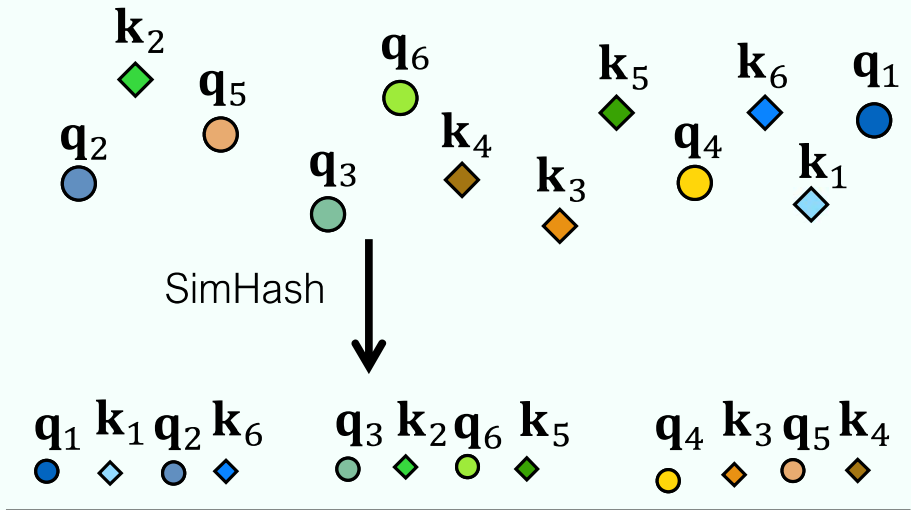
\Rightarrow SortLSH (Locality Sensitive Hashing)

Algorithmic Acceleration

Self Attention

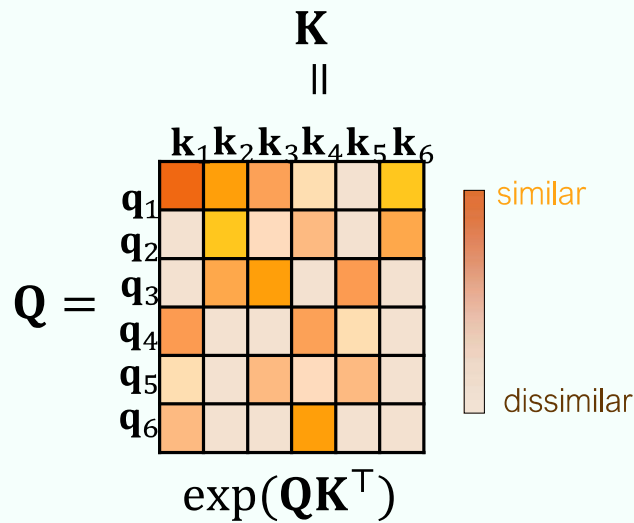
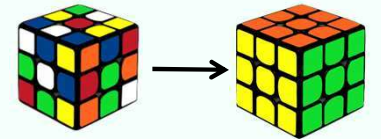


Apply **SortLSH** to gather **similar** entries

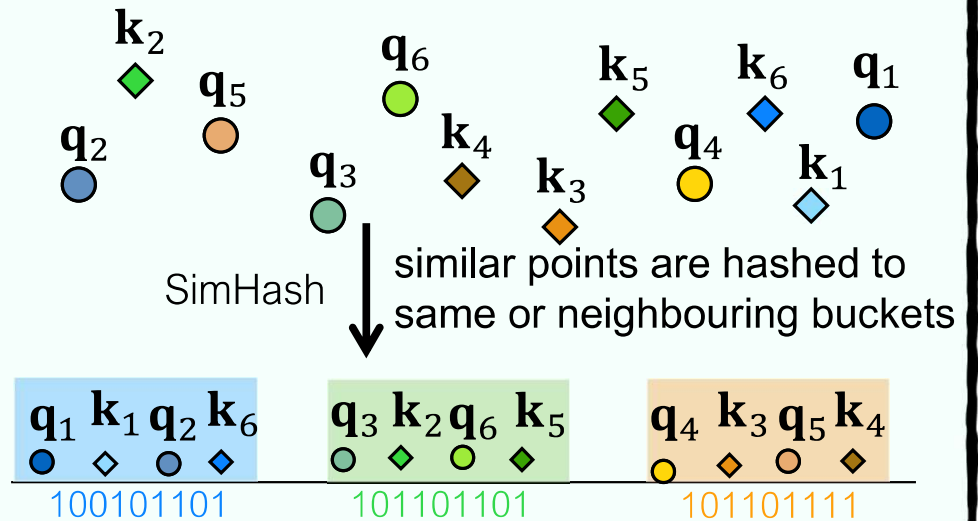


Algorithmic Acceleration

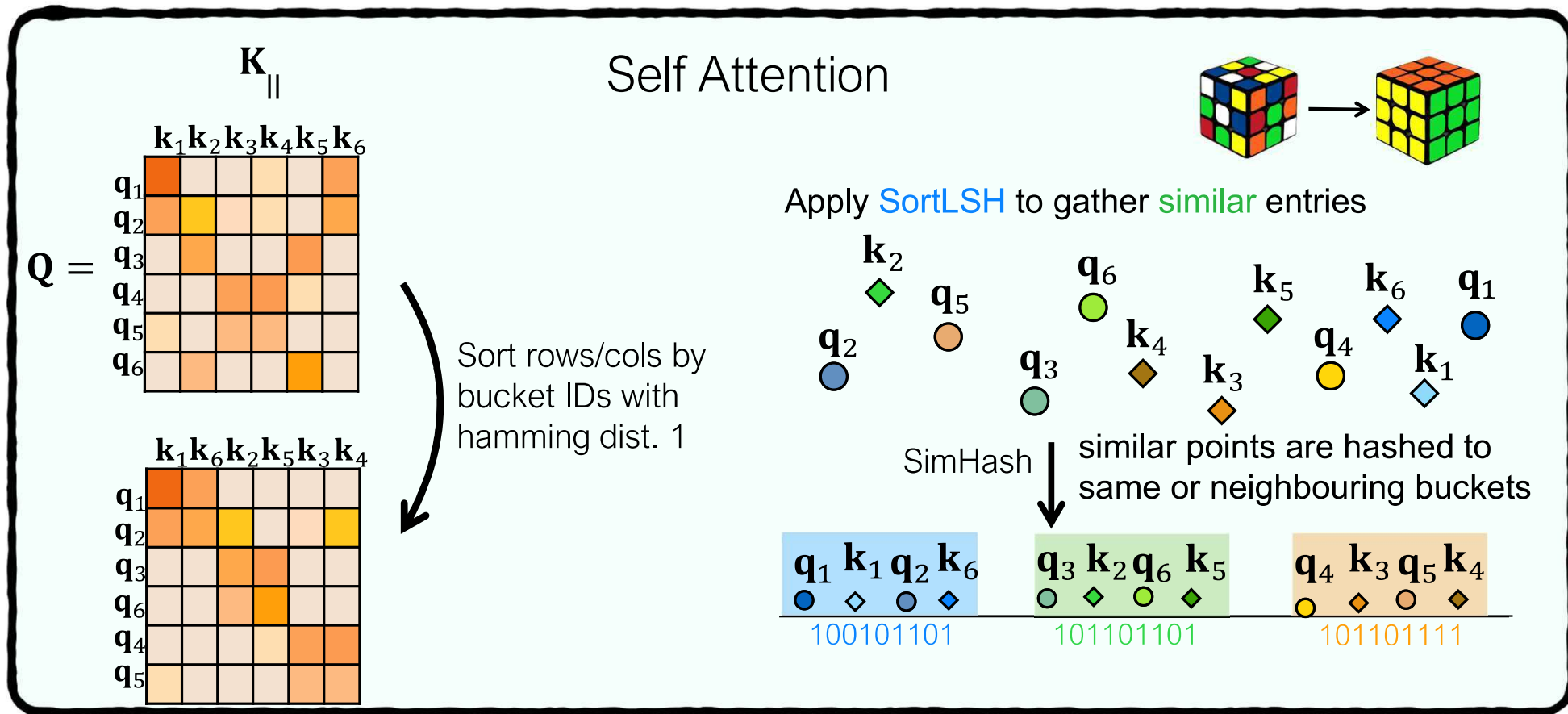
Self Attention



Apply **SortLSH** to gather **similar** entries

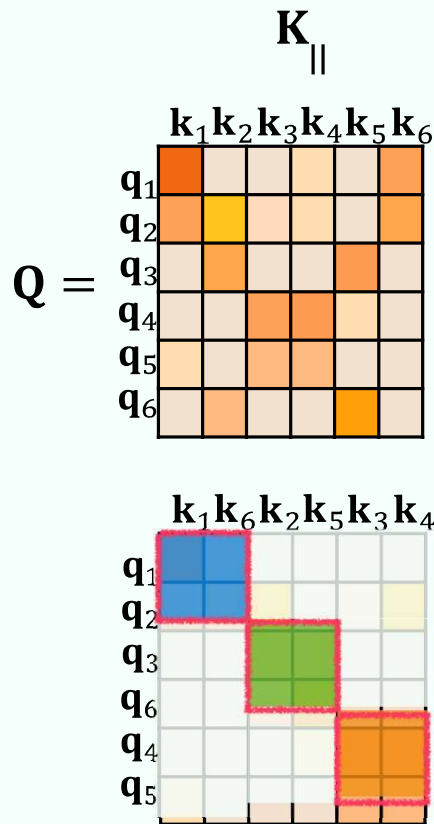
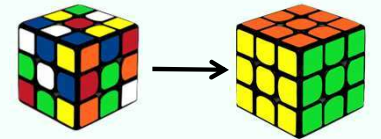


Algorithmic Acceleration



Algorithmic Acceleration

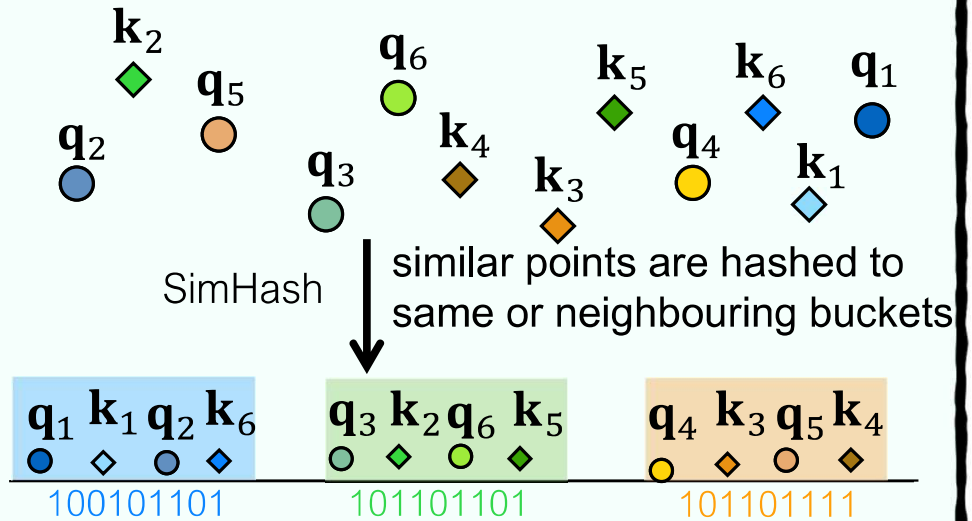
Self Attention



Sort rows/cols by bucket IDs with hamming dist. 1

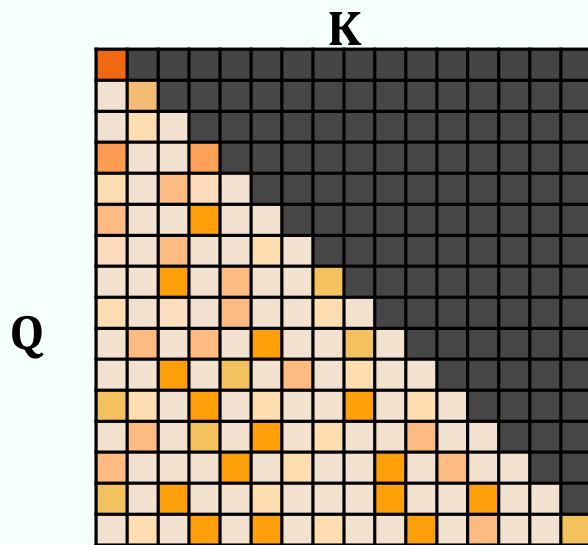
Approximately block-diagonal

Apply **SortLSH** to gather **similar** vectors in buckets



Algorithmic Acceleration

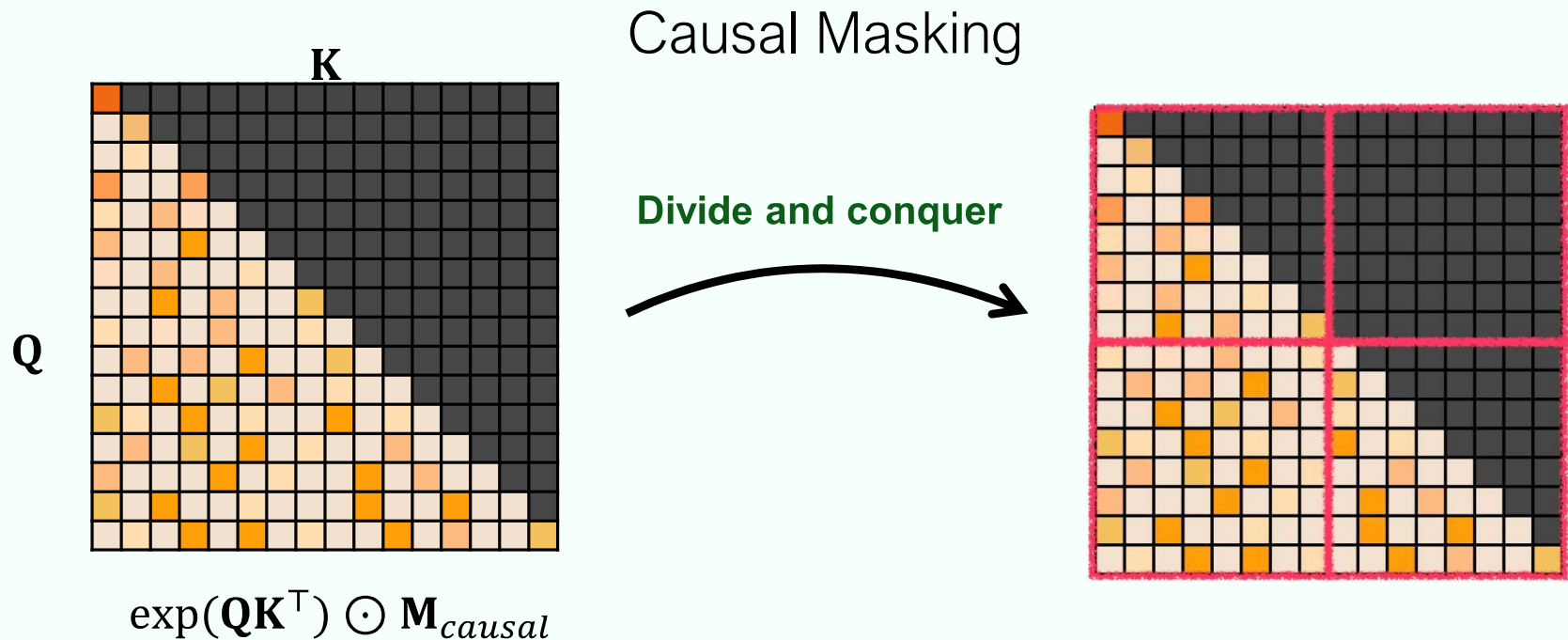
Causal Masking



$$\exp(\mathbf{QK}^T) \odot \mathbf{M}_{causal}$$

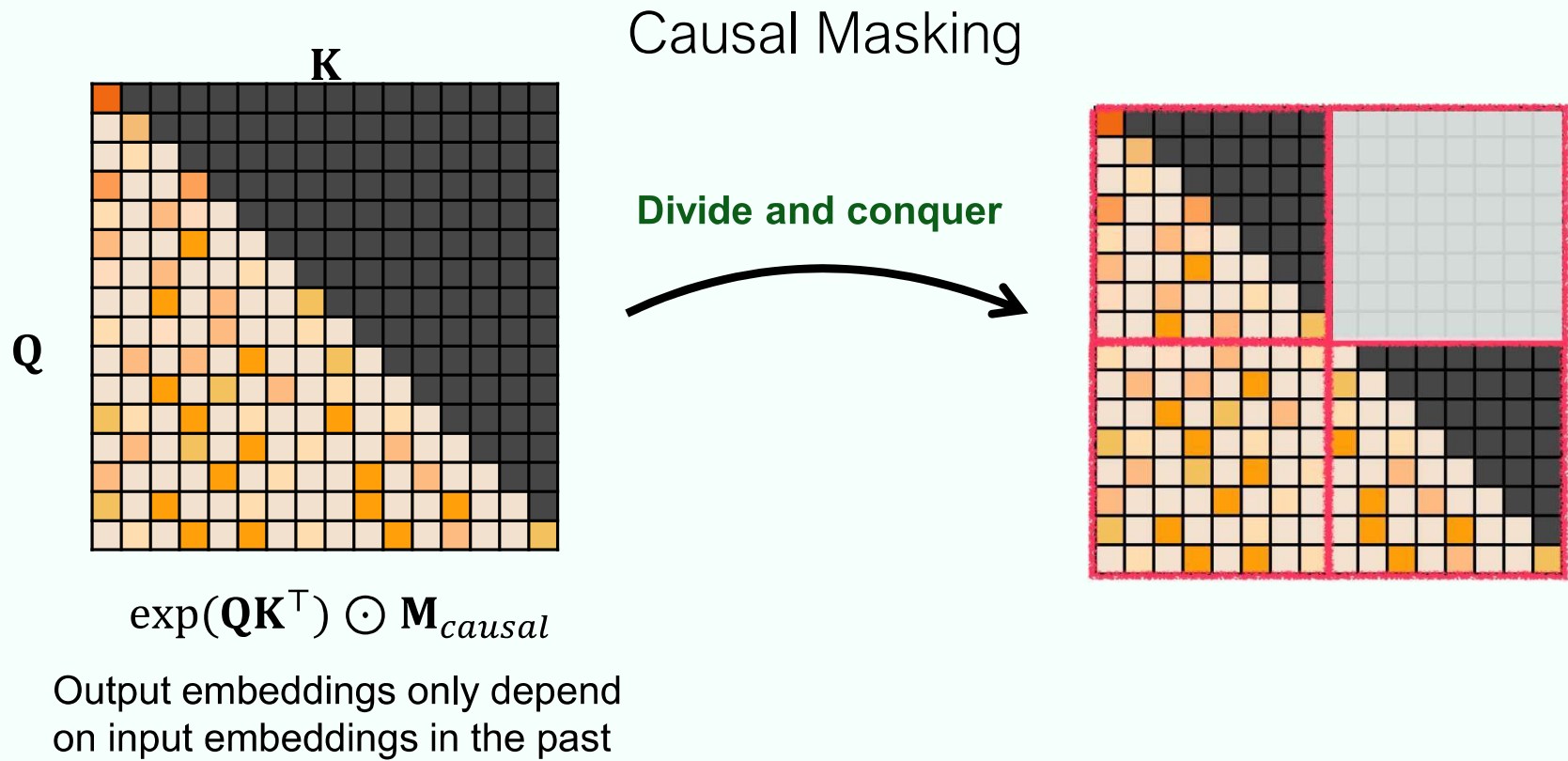
Output embeddings only depend
on input embeddings in the past

Algorithmic Acceleration



Output embeddings only depend on input embeddings in the past

Algorithmic Acceleration



Algorithmic Acceleration



+

HyperAttention: Long-context Attention in Near-Linear Time

Insu Han
Yale University
insu.han@yale.edu

Rajesh Jayaram
Google Research
rkjayaram@google.com

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

Vahab Mirrokni
Google Research
mirrokni@google.com

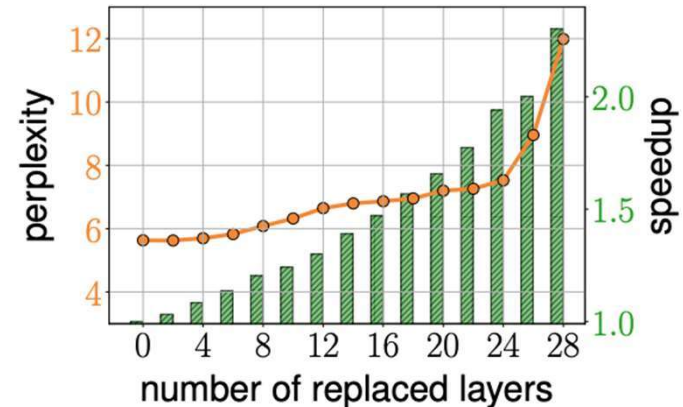
David P. Woodruff
CMU, Google Research
dwoodruf@cs.cmu.edu

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

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

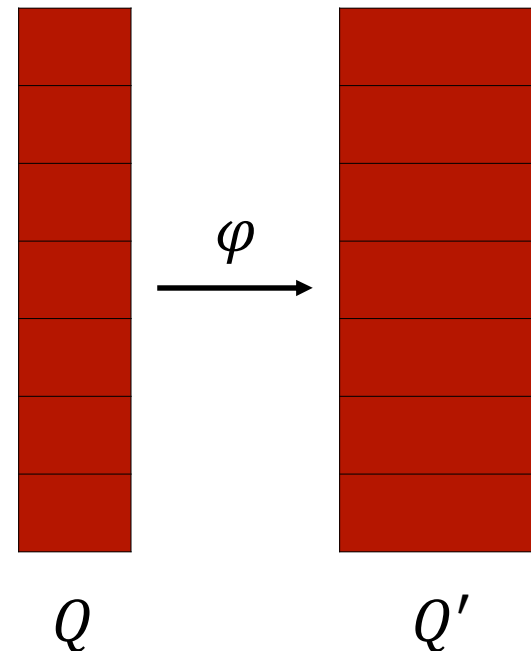
- Softmax: $sim(q, k) \doteq \exp(\langle q, k \rangle)$

Kernel View of Attention

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

$$D^{-1} \cdot \text{LT}(Q' \cdot (K')^\top) \cdot V$$

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



Previous Work

- Performer (Choromanski et al.,) uses a finite-dimensional map φ to approximate exponential
 - Vectors with larger norms require φ 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 $p \geq 2$ is an even integer
 - Always ≥ 0
 - Increases as $\langle q, k \rangle$ goes up

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 $x^{\otimes p} \in \mathbb{R}^{h^p}$
 - $(x^{\otimes p})_{(i_1, i_2, \dots, i_p)} = x_{i_1} \cdot x_{i_2} \cdot \dots \cdot x_{i_p}$

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

Linear Attention using Polynomials

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

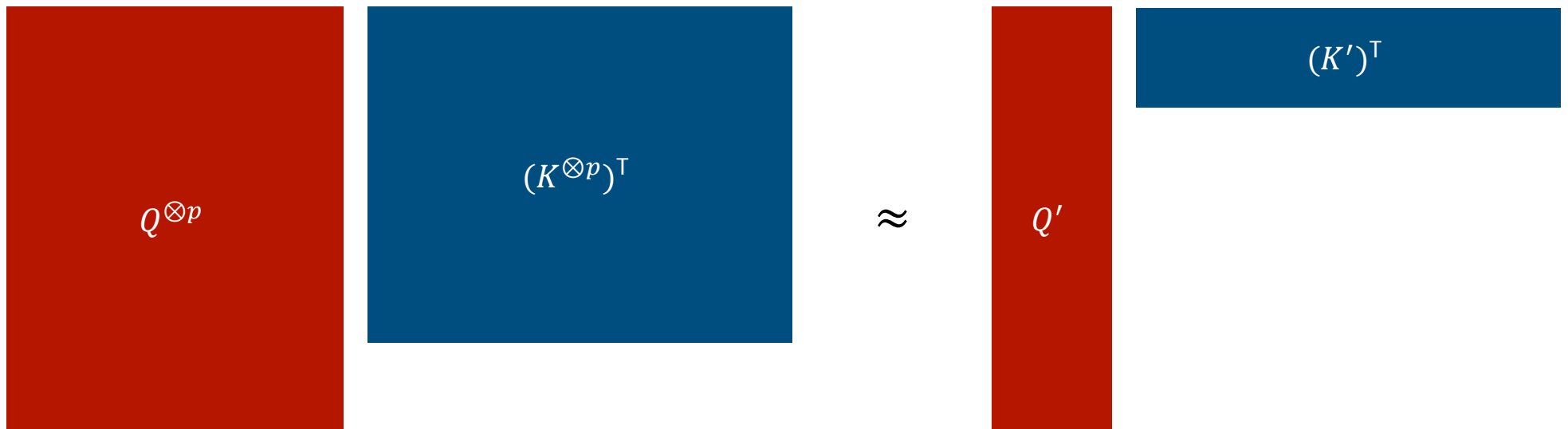
Sketching for Approximate Matrix Multiplication

- Want to compute

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

- $Q^{\otimes p}$ and $K^{\otimes p}$ can have a large number of columns
- Can we compute matrices Q' and K' such that $Q^{\otimes p} \cdot (K^{\otimes p})^T \approx Q' \cdot (K')^T$?
 - 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

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

Solving Issue of Negative Entries

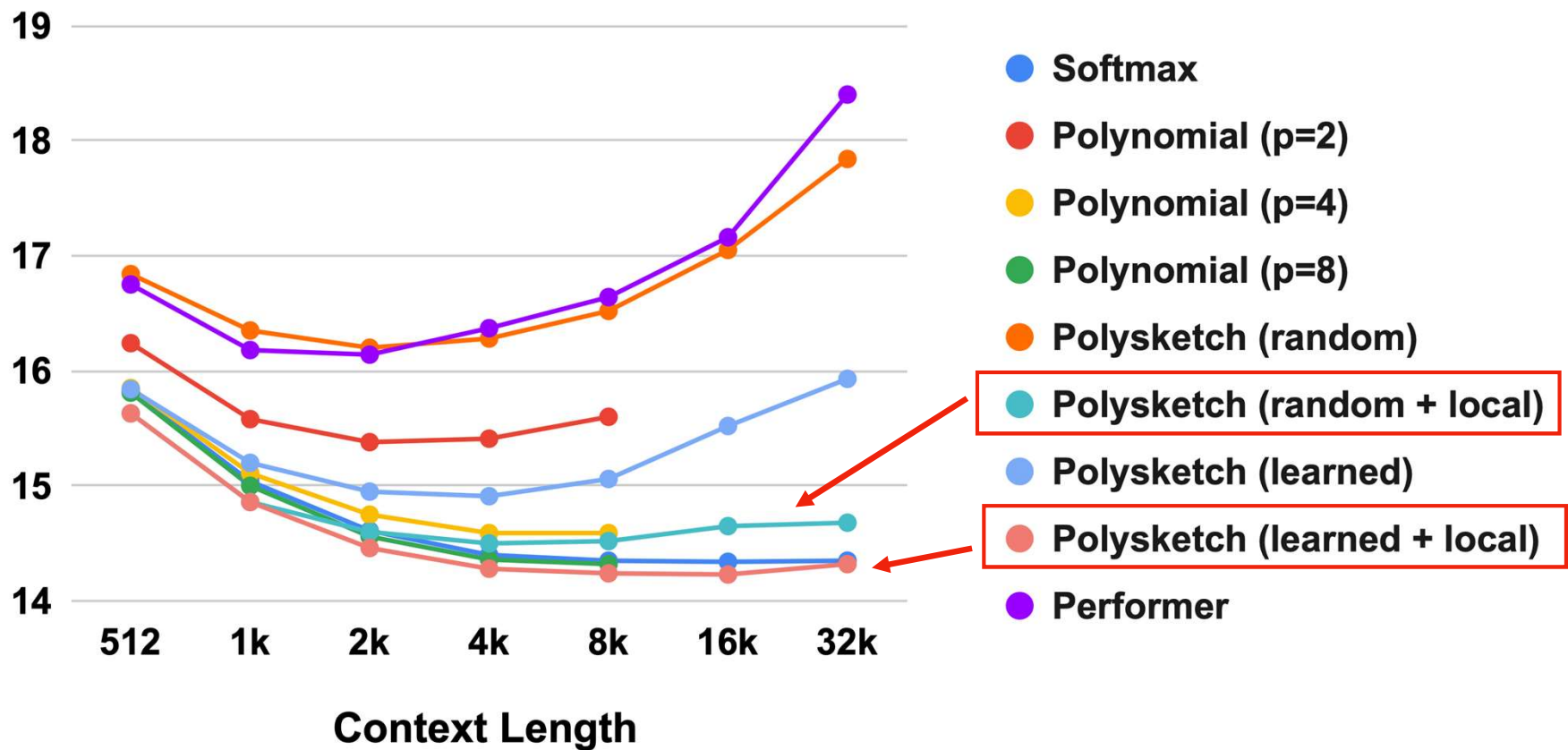
- Consider $Q'' = (Q')^{\otimes 2}$ and $K'' = (K')^{\otimes 2}$
 - Q', K' are sketches for degree $p/2$
- All entries of $Q'' \cdot (K'')^{\top}$ are non-negative! They are of the form $\langle q', k' \rangle^2 \geq 0$
- Show that if Q' and K' have an approximate matrix product property for degree $p/2$, then Q'' and K'' have a similar guarantee for degree p
 - $\| Q'' \cdot (K'')^{\top} - Q^{\otimes p} (K^{\otimes p})^{\top} \|_F$ is small
- Compute $LT(Q'' \cdot (K'')^{\top}) \cdot V$
- 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

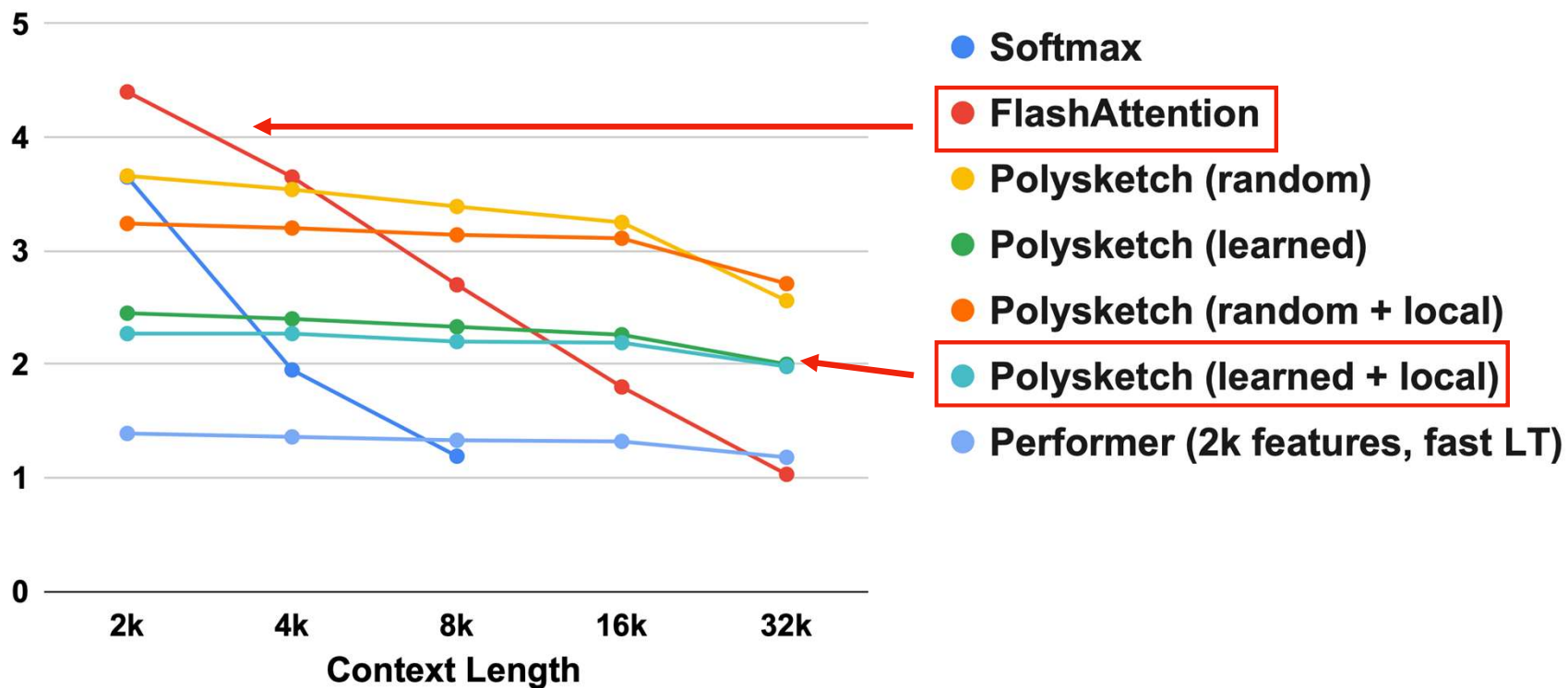
Model Perplexities

Perplexities on Wiki-40B



Training Latencies

Train steps/sec of different mechanisms



Conclusions and Future Work

- In practice,
 - FlashAttention is an optimized implementation of softmax attention and used heavily
 - HyperAttention on pretrained models may reduce quality too much. Fine-tuning increases quality but depends on the hash bucket sizes – still being tested!
 - Expect PolySketchFormer to do worse than softmax on very long contexts, but still may be useful for shorter contexts
- [Kannan, Bhattacharyya, Kacham, W] use tools from randomized linear algebra to show for a class of sym functions, there is a small subset of keys so that any heavy attention score involves a key from that subset. Still being tested!
- Open questions
 - Can we achieve linear time for Hyperattention with weaker assumptions?
 - Can we design a non-negative tensorsketch for PolySketchFormer without squaring the embedding dimension?
 - Are there other natural assumptions that allow us to break the quadratic time worst case hardness?