

Approximate Attention Abhinav Bhatele, Daniel Nichols



Announcements

- Interim report for the project is due on April 17 April 22
- Midterm is on April 10
- Office hours
 - Today I-2pm IRB 5237
- Wednesday 12-2pm Zoom (link on piazza)







- Determine how much tokens should "attend" to other tokens
- Consider a hashmap of our tokens







- Determine how much tokens should "attend" to other tokens
- Consider a hashmap of our tokens









- Determine how much tokens should "attend" to other tokens
- Consider a hashmap of our tokens









- Consider a hashmap of our tokens





- Determine how much tokens should "attend" to other tokens
- Consider a hashmap of our tokens







 $O = [q_i k_1, q_i k_2, \dots, q_i k_n]$



Addressing Quadratic Memory Constraints





Reducing Complexity in Attention





place	liko	home	
place	IIKC	nome	
			A/
	NM		

Reducing Complexity in Attention





place	like	home	
			X

Reducing Complexity with Sparsity







Reducing Complexity with Sparsity







Reducing Complexity with Locality





Abhinav Bhatele, Daniel Nichols (CMSC828G)

We can attend only to a neighborhood of tokens



Reducing Complexity with Structured Sparsity

We can try to find some structured sparse format that...





Abhinav Bhatele, Daniel Nichols (CMSC828G)

reduces compute and/or memory requirements...

and preserves accuracy.



Reducing Complexity with Dynamic Sparsity

We can figure out the sparsity at runtime

				~	-/	$-\lambda$	X.—		1
						11	11	1	9
					2			X	
				-/	1		-7	24	
					11		//		\mathbb{Z}
						\rightarrow	1		X
							$\langle \rangle$	1	1
								V)	6
				7					
					4				
				44-	H	X			
			1	1	11	77	\sim		
			1		Ŵ	7			
		17	-		47	7		XX	
	1	1		1	X.		1	\wedge	
				44	\rightarrow		44		
1	7		1	/			×		Ċ,
11	-					$\langle \rangle$		Y	1
/		1			1	6		7V	
		1			1		1	\sim	
	/			1	1		11		



Abhinav Bhatele, Daniel Nichols (CMSC828G)



Trade-off of better sparsity pattern vs more memory or compute



Reducing Complexity from Repeated Structure





How does this help us?

- We do not need to compute all of the attention values
 - Potentially sub-quadratic attention
- We do not need to store the whole KV-cache!
 - This can get quite large, becoming a memory bottleneck





Benchmarking Approximate Algorithms

- Accuracy/correctness is crucial to benchmark!
- How do we measure accuracy of approximate attention?
 - perplexity of model
 - performance on downstream benchmarks













K-V Cache Eviction













H2O: Approximate Attention

- How can we reduce the KV cache size without hurting accuracy?
 - 30B, 128 bs, 1024 sequence length needs ~180GB space for KV cache alone
- Not all attention values are needed, but how do we determine the most important?
- Across many popular LLMs attention is ~95% sparse



Attention Scores Follow A Power Law





A small subset of tokens are most influential for attention computation



Are the important tokens that important?





Only Retain Important Tokens

- Experimental results suggest keeping heavy-hitters and recent tokens in KV cache is
 - enough for high accuracy with big memory savings
- How do we determine heavy-hitters apriori?



Use Local Attention Scores to Approximate Global Attention

Use local attention steps during decoding phase to evict from KV cache









H2O Results: How does it do?

Across most of their experiments they can reduce KV cache to 20% size and get comparable results





Models "collapse" once an important token is evicted from K-V cache



Abhinav Bhatele, Daniel Nichols (CMSC828G)

It actually helps accuracy in a number of benchmarks



Combining with Other Methods

- Cache eviction policy can be combined with other types of sparse attention
- Get the benefits of other types of sparsity, while reducing KV cache size
- Generally improves all other methods



- Low rank approximations





- Low rank approximations
- Kernel methods





- Low rank approximations
- Kernel methods
- Hashing
 - "REFORMER: The Efficient Transformer" Kitaev et. al.





Use hashing to bucket similar tokens and compute dense attention within buckets



 $\mathbf{q}_1 \ \mathbf{q}_2 \ \mathbf{q}_3 \ \mathbf{q}_4 \ \mathbf{q}_5 \ \mathbf{q}_6$



(d) Chunked

q₅

- Low rank approximations
- Kernel methods
- Hashing
- Clustering
 - "Efficient Content-Based Sparse Attention with Routing Transformers" Roy et. al.





Abhinav Bhatele, Daniel Nichols (CMSC828G)



Use fast K-Means to cluster tokens and only do dense attention within clusters

(b) Strided attention

(c) Routing attention



- Low rank approximations
- Kernel methods
- Hashing
- Clustering
- Sinks and Memory Augmentation
 - "ETC: Encoding Long and Structured Inputs in Transformers" Ainslie et. al.



Abhinav Bhatele, Daniel Nichols (CMSC828G)



Full attention (unidirectional)



Quadratic complexity Local attention (bidirectional)



Linear complexity, but no connection between distant tokens



connection between distant tokens via global memory

- Low rank approximations
- Kernel methods
- Hashing
- Clustering
- Sinks and Memory Augmentation
- Recurrence
 - "Transformer-XL: Attentive Language Models Beyond a Fixed-Length Context" Dai et. al.





Abhinav Bhatele, Daniel Nichols (CMSC828G)

Process in chunks and carry some information between chunks

- Low rank approximations
- Kernel methods
- Hashing
- Clustering
- Sinks and Memory Augmentation
- Recurrence
- Hybrid
 - H2O







UNIVERSITY OF MARYLAND