

# Long-Context Language Models

Sam Pan

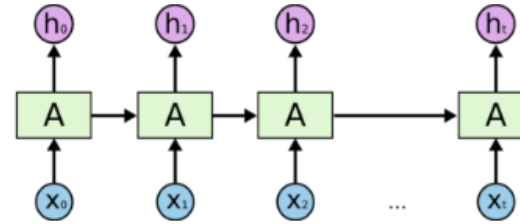
# LongNet: Scaling Transformers to 1B Tokens

**Jiayu Ding<sup>♣\*</sup> Shuming Ma<sup>♣\*</sup> Li Dong<sup>♣</sup> Xingxing Zhang<sup>♣</sup>**  
**Shaohan Huang<sup>♣</sup> Wenhui Wang<sup>♣</sup> Nanning Zheng<sup>♣†</sup> Furu Wei<sup>♣†</sup>**  
<sup>♣</sup>Microsoft Research  
<sup>♣</sup>Xi'an Jiaotong University

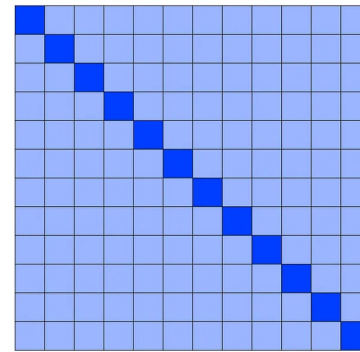
2023

# Self Attention

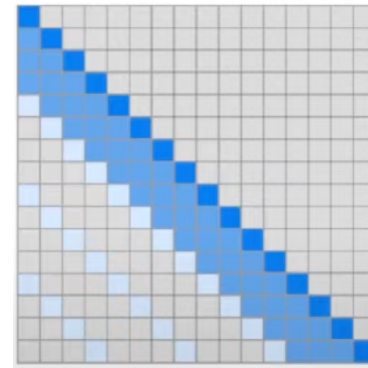
Recurrent



Vanilla Attention



Sparse Attention

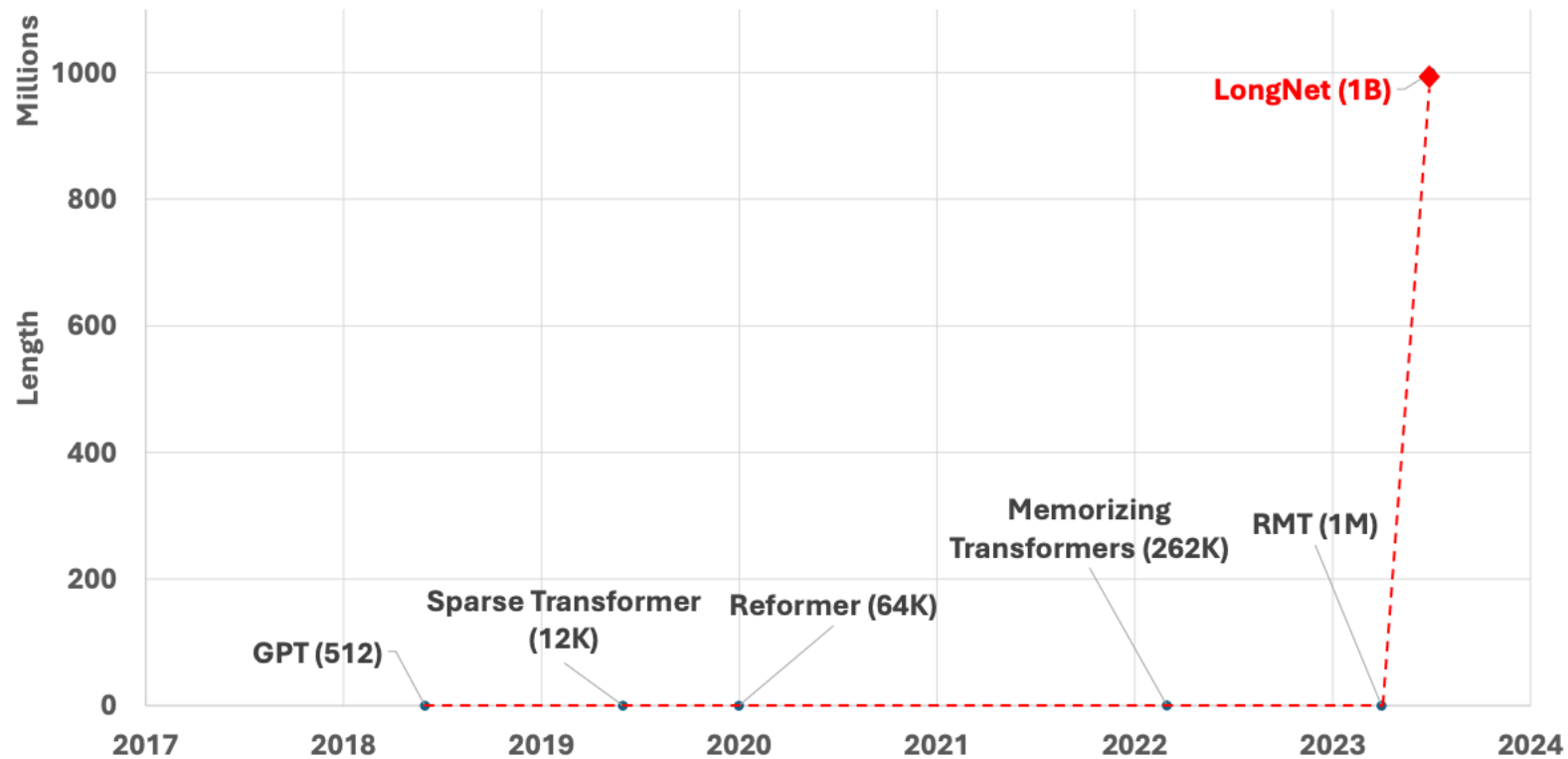


# Self Attention

Linear Complexity essential for mass scaling

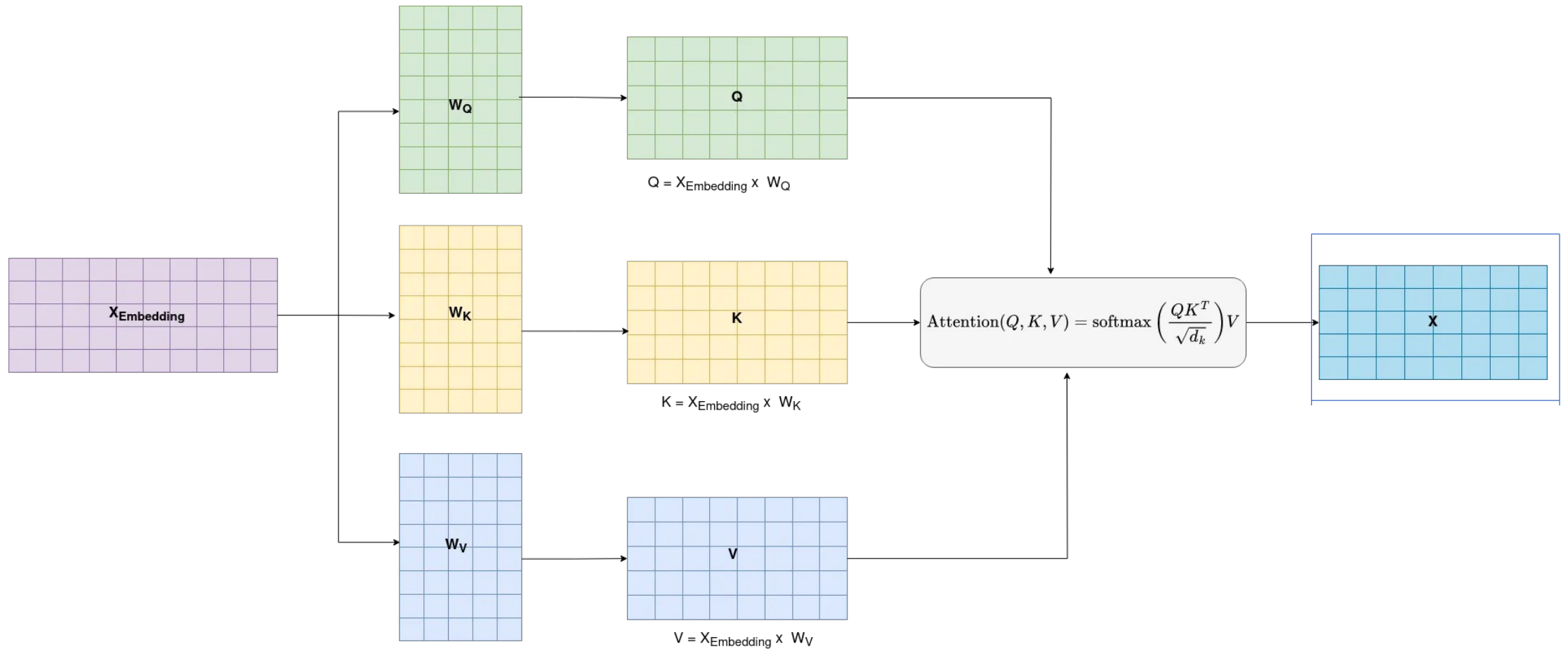
<b>Method</b>	<b>Computation Complexity</b>
Recurrent	$\mathcal{O}(Nd^2)$
Vanilla Attention	$\mathcal{O}(N^2d)$
Sparse Attention	$\mathcal{O}(N\sqrt{N}d)$
<b>Dilated Attention (This Work)</b>	$\mathcal{O}(Nd)$

# Trend of Transformer sequence lengths

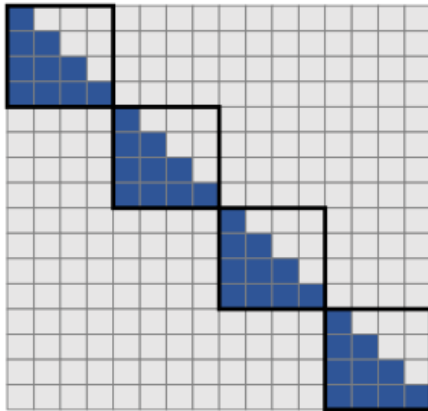


Sequence length, as the last atomic dimension of the neural network, is desirable to be unlimited

# Vanilla Attention

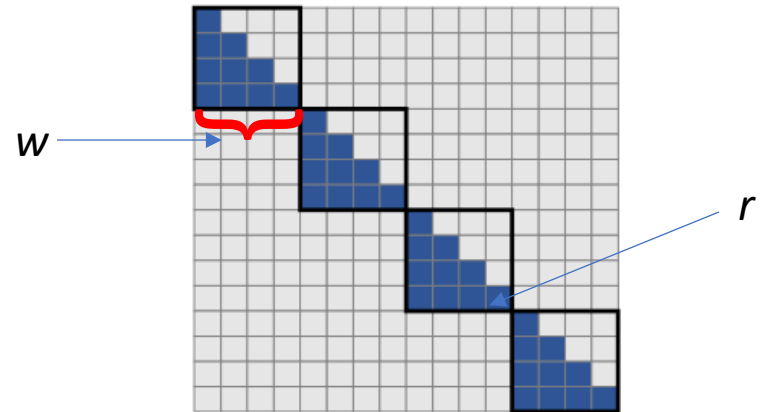


# LongNet: Dilated Attention



Segment Length: 4  
Dilated Rate: 1

# LongNet: Dilated Attention

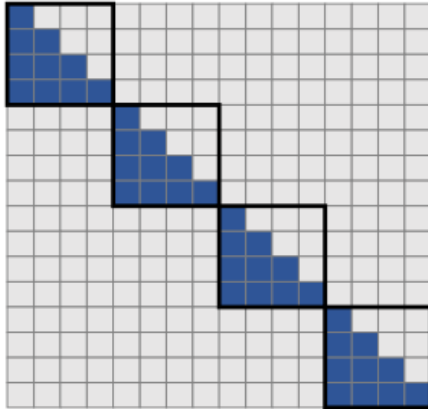


Segment Length: 4  
Dilated Rate: 1

Introduces new parameters  $w$  and  $r$   
 $w$  – segment length  
 $r$  – dilation rate



# LongNet: Dilated Attention



$w$  Segment Length: 4  
 $r$  Dilated Rate: 1

$$\tilde{Q}_i = [Q_{iw}, Q_{iw+r}, Q_{iw+2r}, \dots, Q_{(i+1)w-1}]$$

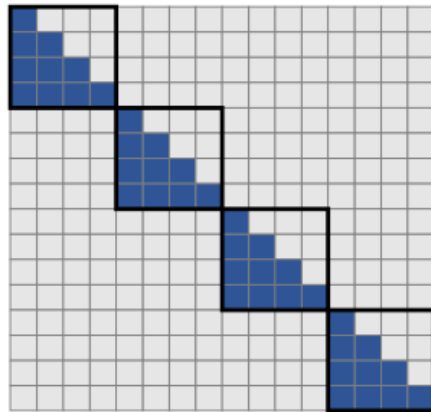
$$\tilde{K}_i = [K_{iw}, K_{iw+r}, K_{iw+2r}, \dots, K_{(i+1)w-1}]$$

$$\tilde{V}_i = [V_{iw}, V_{iw+r}, V_{iw+2r}, \dots, V_{(i+1)w-1}]$$

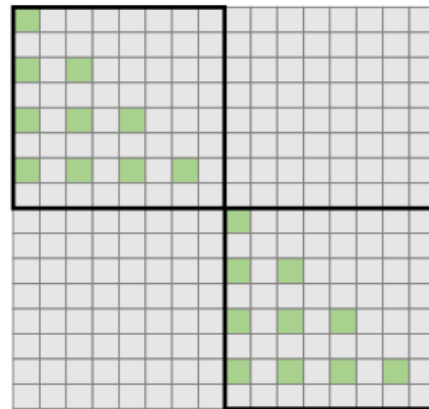
$$\tilde{O}_i = \text{softmax}(\tilde{Q}_i \tilde{K}_i^T) \tilde{V}_i$$

Introduces new parameters  $w$  and  $r$

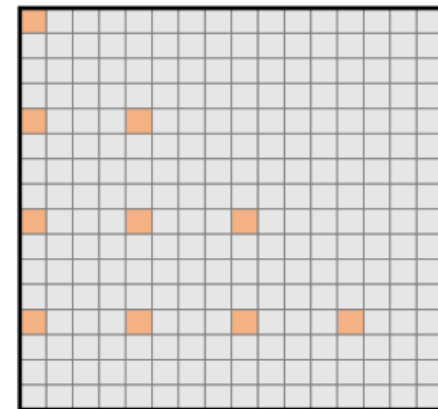
# LongNet: Dilated Attention



$w$  Segment Length: 4  
 $r$  Dilated Rate: 1



Segment Length: 8  
Dilated Rate: 2



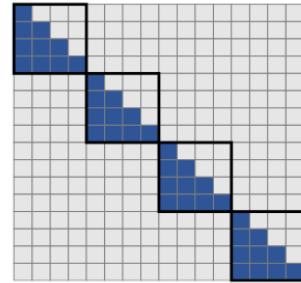
Segment Length: 16  
Dilated Rate: 4

Introduces new parameters  $w$  and  $r$

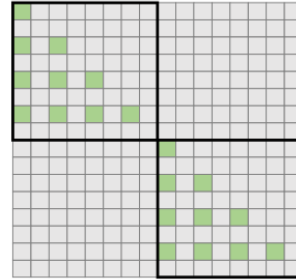
Geometric Sequence:  $w = \{w_0, w_1, w_2, \dots, N\} (w_i < w_{i+1} < N)$

$r = \{1, r_1, r_2, \dots, r_k\}_k (1 < r_i < r_{i+1})$

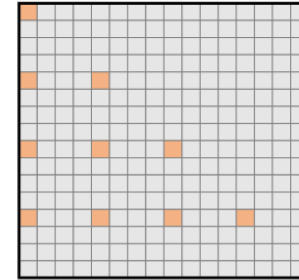
# LongNet: Dilated Attention



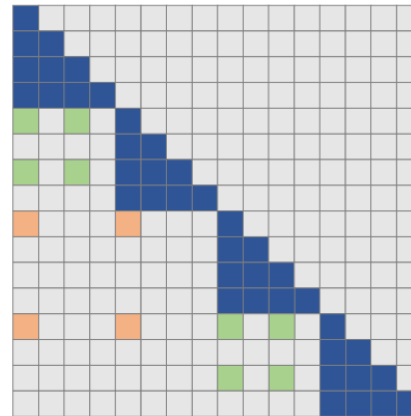
Segment Length: 4  
Dilated Rate: 1



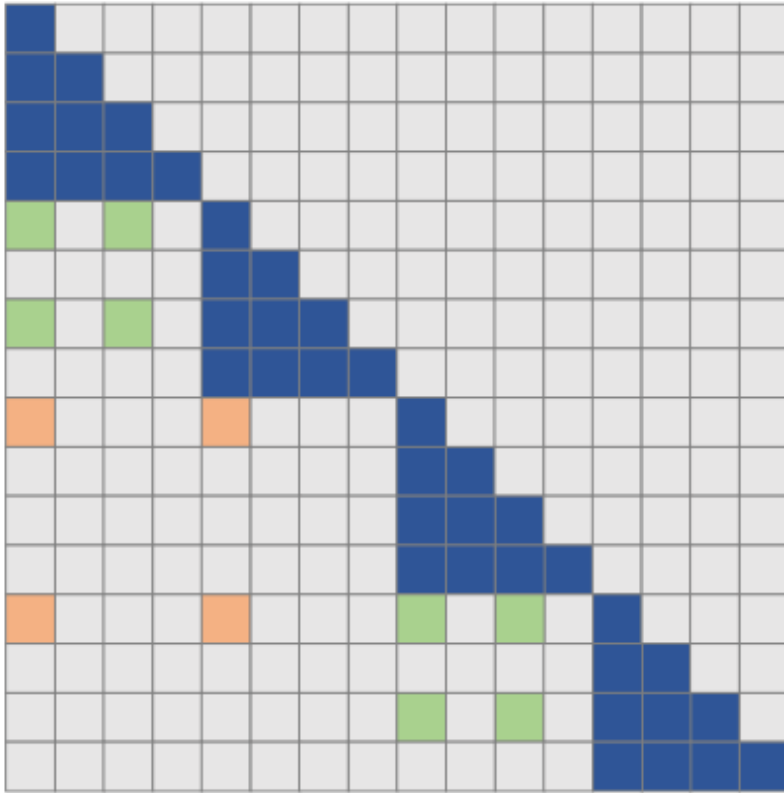
Segment Length: 8  
Dilated Rate: 2



Segment Length: 16  
Dilated Rate: 4



# Dilated Attention: Computational Complexity

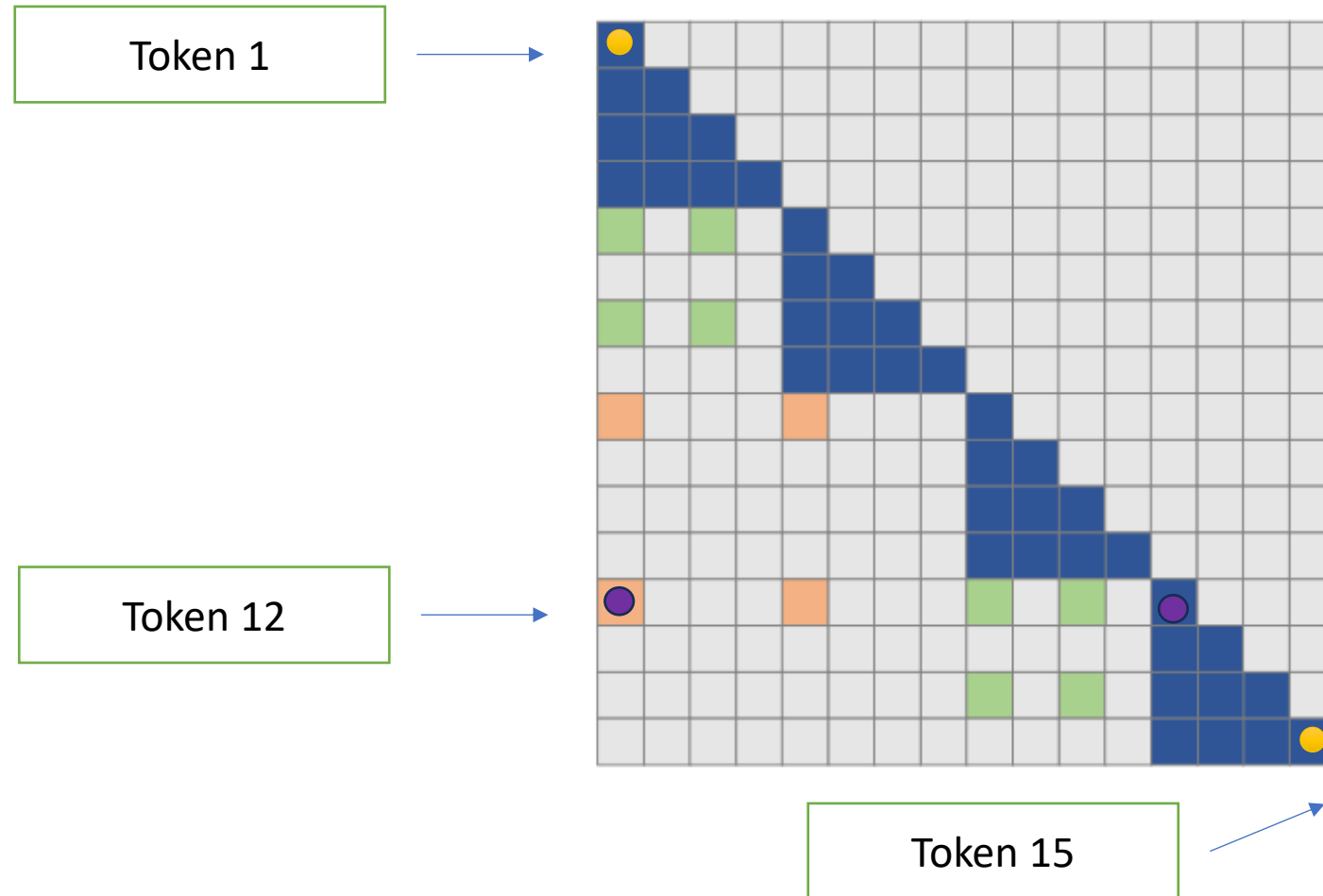


$$FLOPs = \frac{2N}{w} \left(\frac{w}{r}\right)^2 d = \frac{2Nwd}{r^2}$$

$$FLOPs = 2Nd \sum_{i=1}^k \frac{w_i}{r_i^2}$$

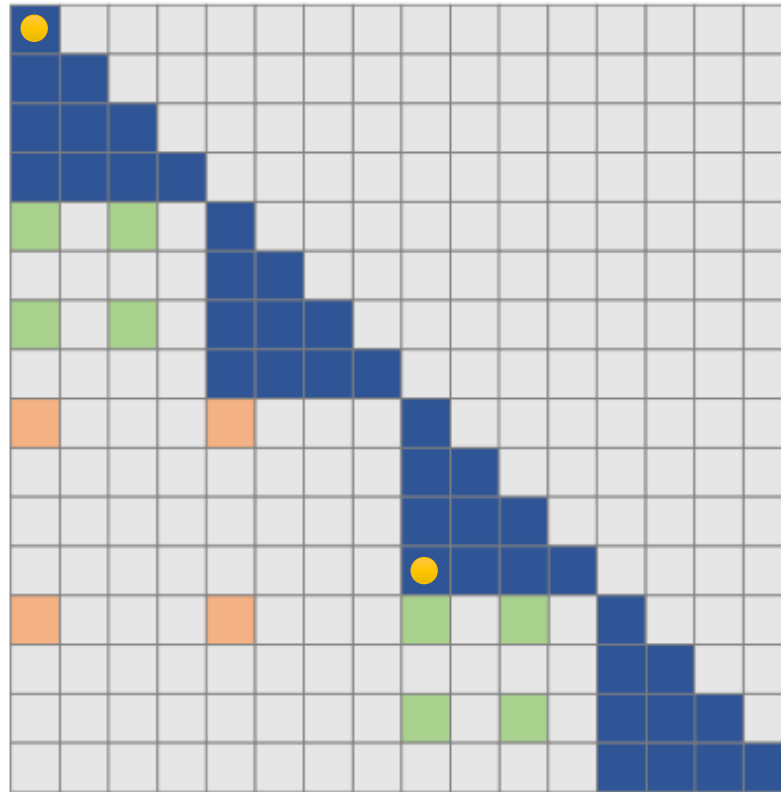
$$\mathcal{O}(Nd)$$

# Dilated Attention: Related Attention



# Dilated Attention: Related Attention

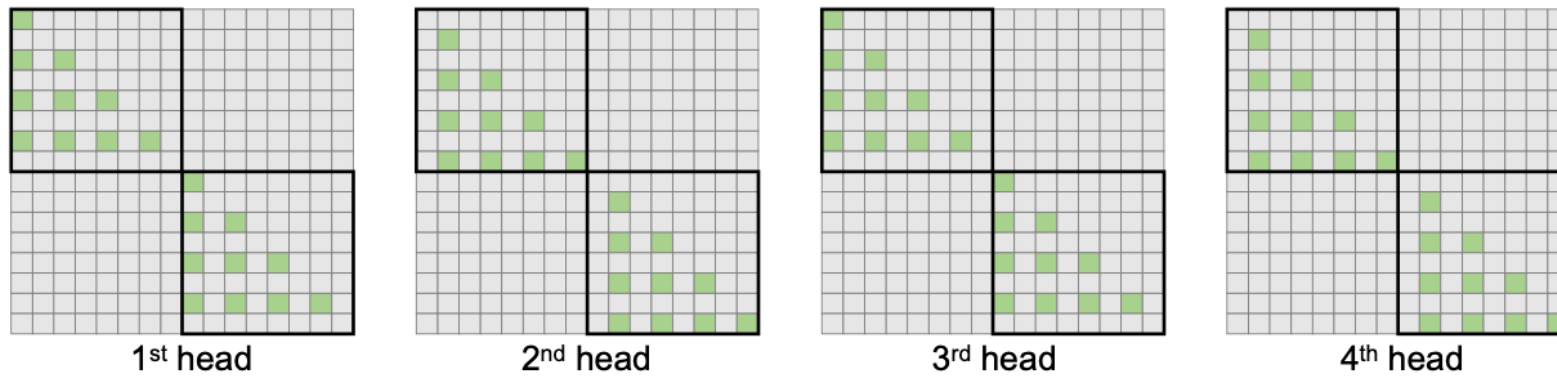
$$L = \mathcal{O}(\log N)$$



*Attention allocation decreases exponentially as the distance between tokens grows*

# Multihead Attention

Different computation among different heads sparsify different parts of the query-key-value pairs

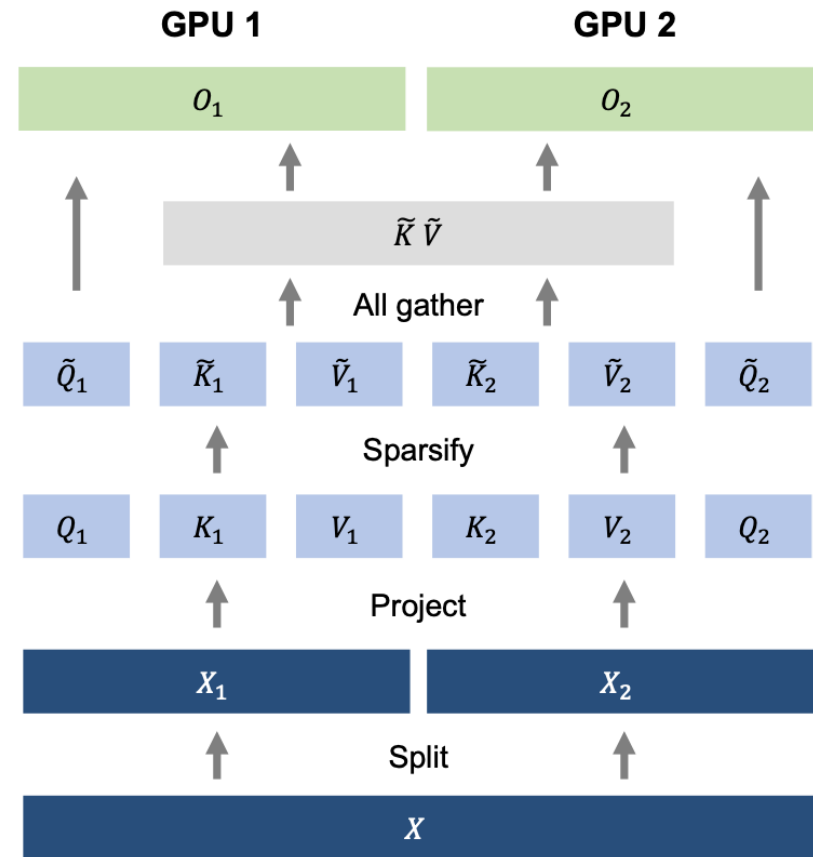


Segment Length: 8  
Dilated Rate: 2  
Heads: 4

# Distributed Training

Parallelize training by partitioning the sequence dimension

The computation and communication costs are nearly constant as the number of devices grows



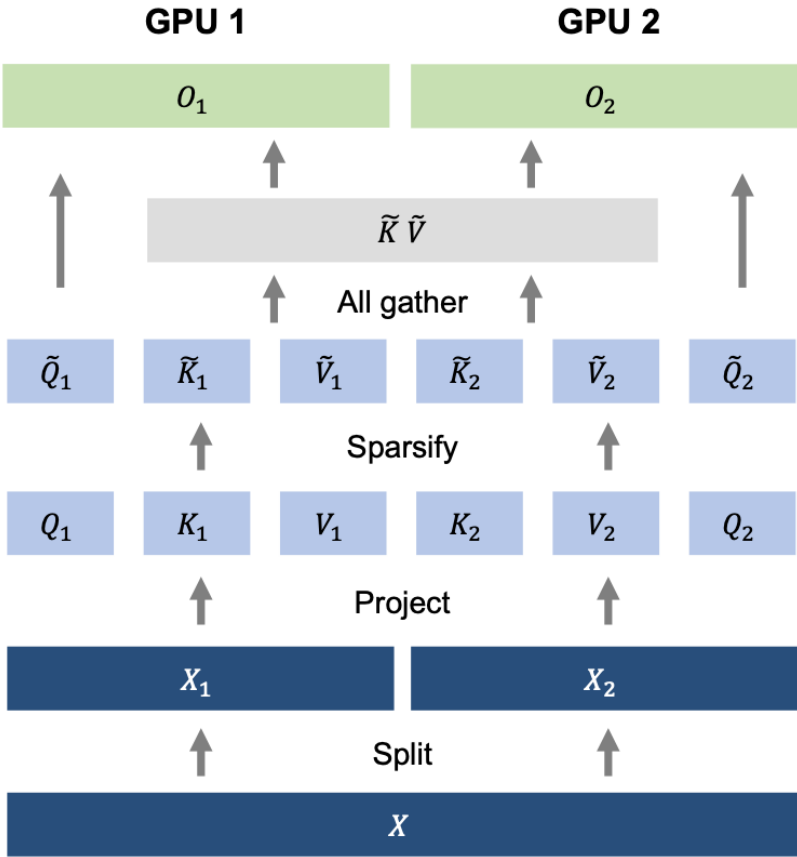


# Distributed Training

The concatenation of the outputs across different devices becomes the final attention output:

$$\tilde{O}_1 = \text{softmax}(\tilde{Q}_1 \tilde{K}^T) \tilde{V}, \quad \tilde{O}_2 = \text{softmax}(\tilde{Q}_2 \tilde{K}^T) \tilde{V}$$

$$\tilde{O} = [\tilde{O}_1, \tilde{O}_2]$$

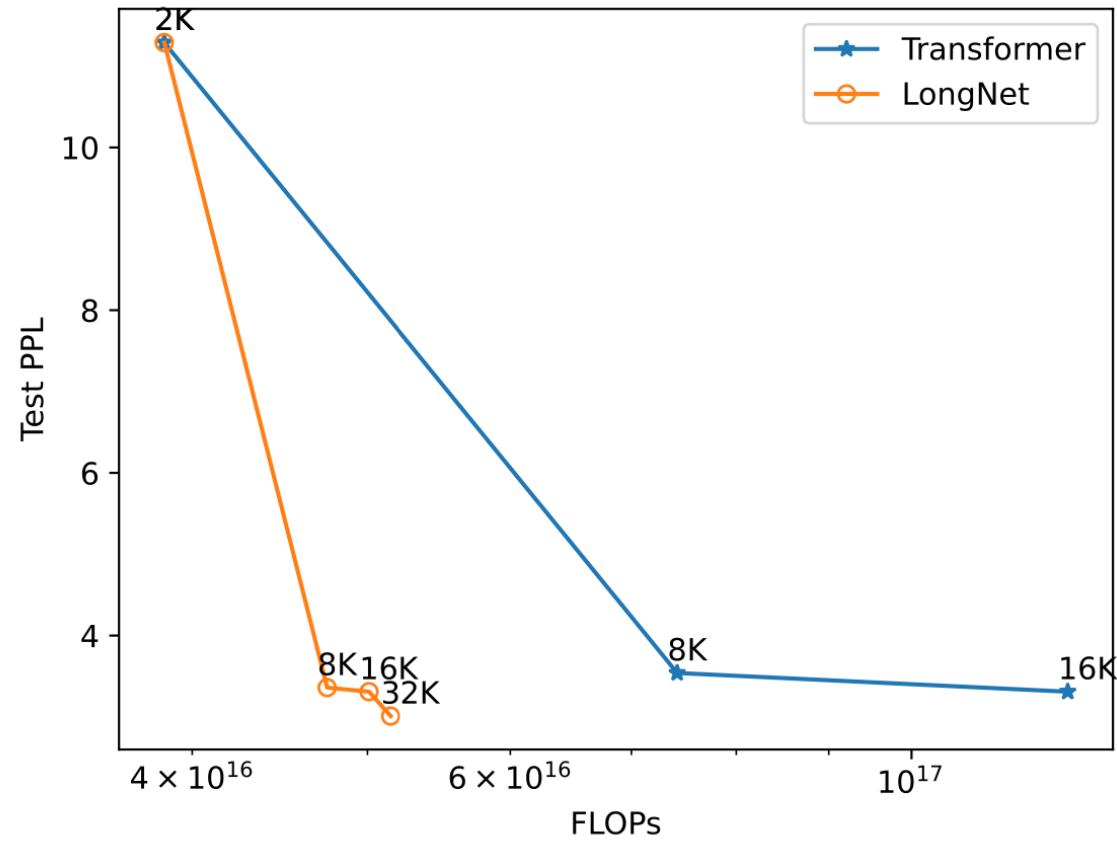


# Scaling up to 1B tokens

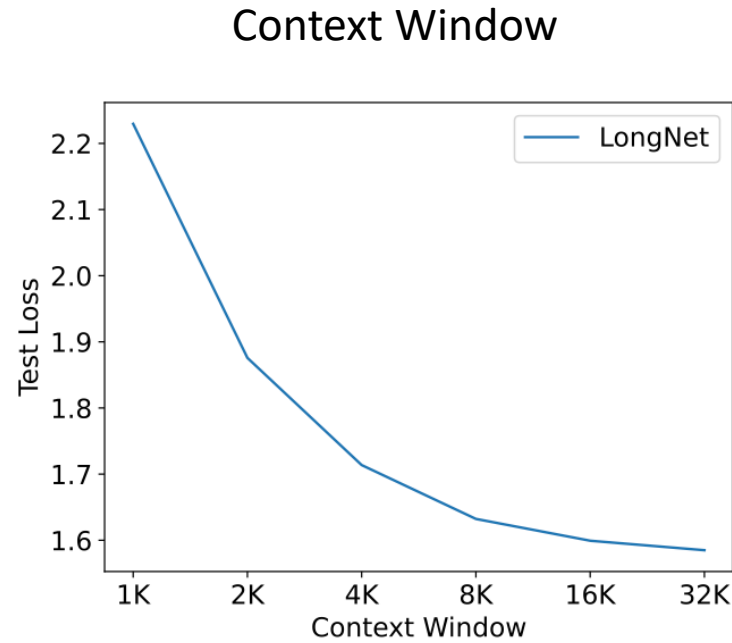
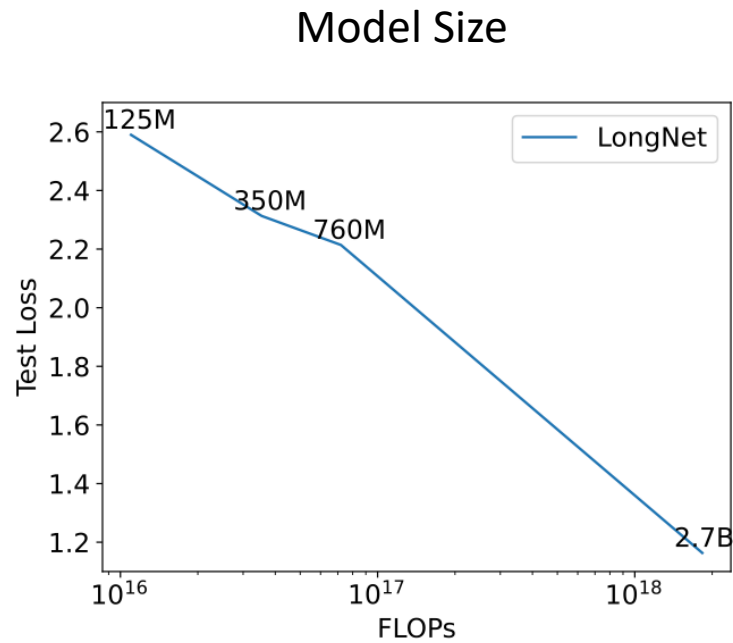
Increasing the sequence length during training generally leads to a better language model

Model	Length	Batch	Github		
			2K	8K	32K
Transformer [VSP <sup>+</sup> 17]	2K	256	4.24	5.07	11.29
Sparse Transformer [CGRS19]	8K	64	4.39	3.35	8.79
<b>LONGNET (ours)</b>			4.23	3.24	3.36
Sparse Transformer [CGRS19]	16K	32	4.85	3.73	19.77
<b>LONGNET (ours)</b>			4.27	3.26	3.31
Sparse Transformer [CGRS19]	32K	16	5.15	4.00	3.64
<b>LONGNET (ours)</b>			4.37	3.33	3.01

# Results



# Scaling up Model Size



Dense Transformer is not a prerequisite for scaling the language models

# Memorizing Transformers

Yuhuai Wu, Markus N. Rabe, DeLesley Hutchins, Christian Szegedy

`{yuhuai, mrabe, delesley, szegedy}@google.com`

2022

# Growing Knowledge Base

Theorem database in mathematics

Codebase in program synthesis

**Lemma 3.3** Let  $\varphi \in L^2(\mathbb{K})$  be a refinable function with refinement mask  $m_0(\xi)$  such that condition (3.5) is satisfied. Then,  $P_j = \sum_{k \in \mathbb{N}_0} |\langle f, \varphi_{j,k} \rangle|^2 < \infty$ , for any function  $f \in L^2(\mathbb{K})$  and

$$(i) \quad \lim_{j \rightarrow \infty} P_j = \|f\|^2; \quad (ii) \quad \lim_{j \rightarrow -\infty} P_j = 0$$

where  $\varphi_{j,k}(x) = q^{j/2} \varphi(p^{-j}x - u(k))$ ,  $j \in \mathbb{Z}, k \in \mathbb{N}_0$ .

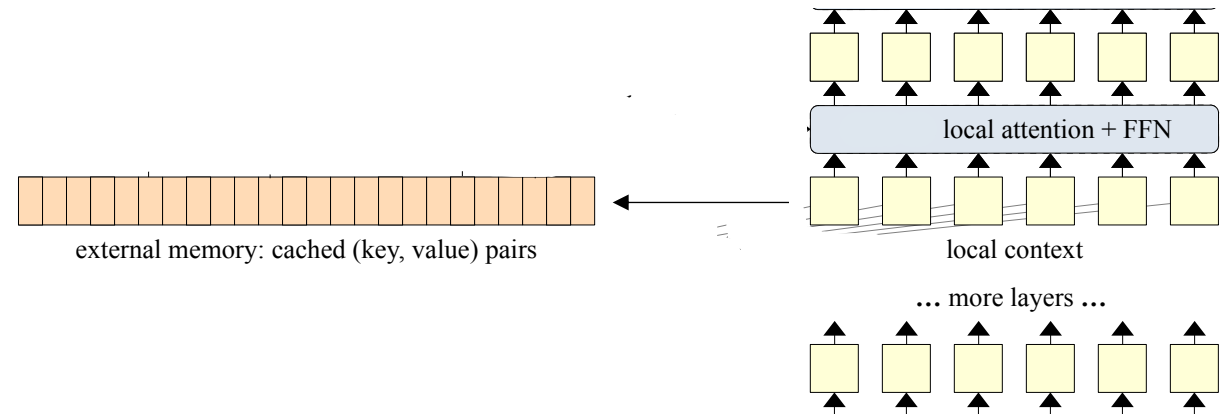
**Proof:** Implementation of Plancherel and Parseval formulae yields

$$\begin{aligned} & \sum_{k \in \mathbb{N}_0} |\langle f, \varphi_{j,k} \rangle|^2 \\ &= q^j \sum_{k \in \mathbb{N}_0} \left| \int_{\mathbb{K}} \hat{f}(\xi) \overline{\hat{\varphi}(p^j \xi)} \chi_k(p^j \xi) d\xi \right|^2 \\ &= q^j \sum_{k \in \mathbb{N}_0} \left| \int_{B^{-j}} \left[ \sum_{n \in \mathbb{N}_0} \hat{f}(\xi + p^{-j}u(n)) \overline{\hat{\varphi}(p^j(\xi + p^{-j}u(n)))} \right] \chi_k(p^j \xi) d\xi \right|^2 \\ &= \int_{B^{-j}} \left| \sum_{n \in \mathbb{N}_0} \hat{f}(\xi + p^{-j}u(n)) \overline{\hat{\varphi}(p^j(\xi + p^{-j}u(n)))} \right|^2 d\xi \\ &= \|F_j\|^2, \end{aligned} \tag{3.8}$$



# Memorizing Transformers

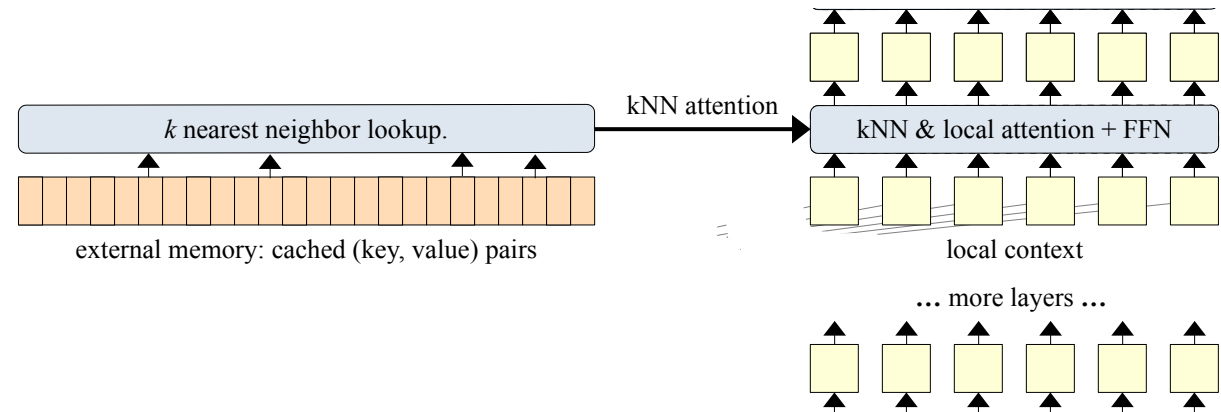
Maintain a memory previously generated keys and values



# Memorizing Transformers

Maintain a memory previously generated keys and values

Approximate attention into memory via kNN for scalability



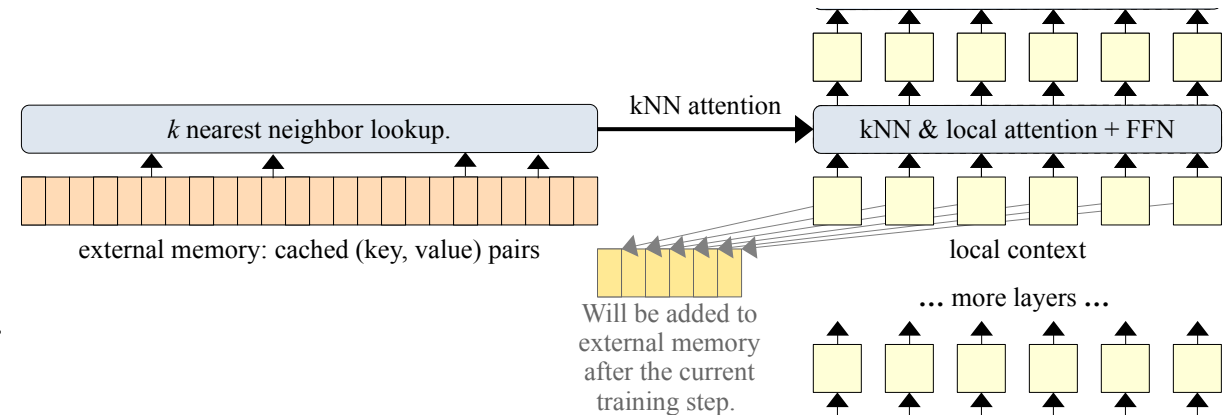


# Memorizing Transformers

Maintain a memory previously generated keys and values

Approximate attention into memory via kNN for scalability

Fast maximum inner product search algorithm –  $O(\log N)$



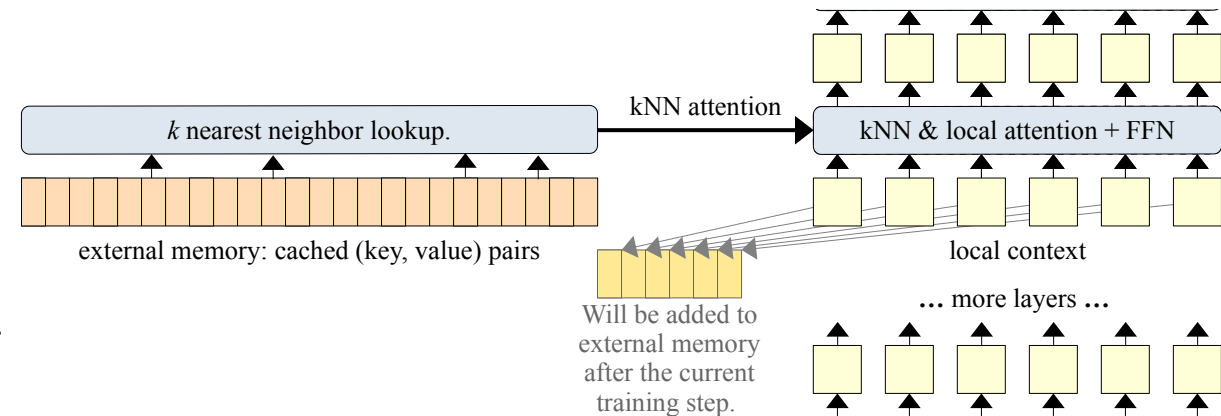
# Memorizing Transformers

Maintain a memory previously generated keys and values

Approximate attention into memory via kNN for scalability

Fast maximum inner product search algorithm –  $O(\log N)$

Do not backpropagate into memory



# Memory update

- Split sequence into subsequences in consecutive order
- Store keys and values in memory after subsequence is read
- If memory is not empty, the current subsequence can attend to the memory

# Attending to both context and memory

- Originally only attend to context's key and value
- Perform a top-k attention using fast maximum inner search into memory
  - Retrieve the top-k keys and value and perform attention on them
- Allows for large scale memory without computational constraint
- This is approximate k-NN search

# Attending to both context and memory

$$\mathbf{V}_a = \mathbf{V}_m \odot g + \mathbf{V}_c \odot (1 - g)$$

$\mathbf{V}_a$  - Next token combined result of attention

$\mathbf{V}_m$  - Attention from external memory

$\mathbf{V}_c$  - Attention from context

$g$  – learnable parameter between 0 and 1

# Results

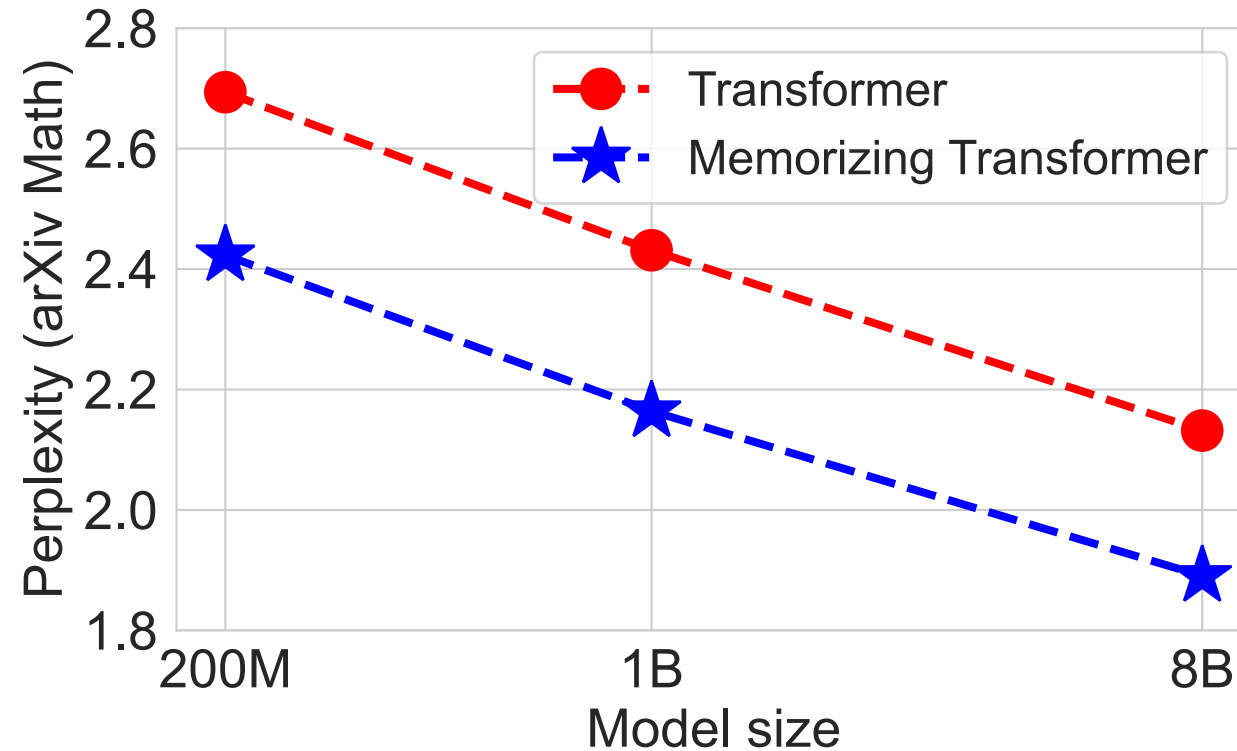
Context	Memory	XL cache	arXiv	PG19	C4(4K+)	GitHub	Isabelle
512	None	None	3.29	13.71	17.20	3.05	3.09
2048	None	None	2.69	12.37	14.81	2.22	2.39
512	None	512	2.67	12.34	15.38	2.26	2.46
2048	None	2048	2.42	11.88	14.03	2.10	2.16
512	1536	None	2.61	12.50	14.97	2.20	2.33
512	8192	None	2.49	12.29	14.42	2.09	2.19
512	8192	512	2.37	11.93	14.04	2.03	2.08
512	65K	512	2.31	11.62	14.04	1.87	2.06
2048	8192	2048	2.33	11.84	13.80	1.98	2.06
2048	65K	2048	<b>2.26</b>	<b>11.37</b>	<b>13.64</b>	<b>1.80</b>	<b>1.99</b>

Adding external memory results in substantial gains across datasets and architectures

Increasing the size of the memory increases the benefit of the memory

# Scaling Model

- Smaller Memorizing Transformer with just 8k tokens improves perplexity
- 8K Memory attained results comparable to 5-8x larger model



# Finetuning transformer to use memory

Finetune 20K steps to obtain 85% of the benefits brought by memory

