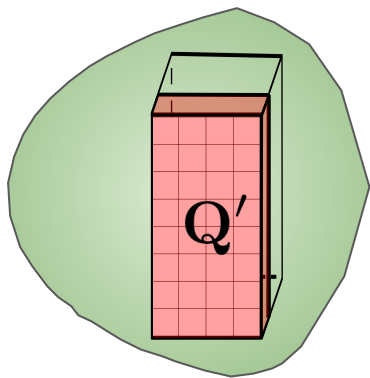
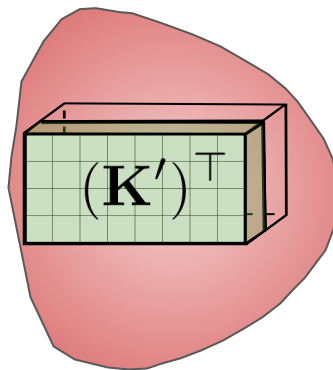


**softmax attention**



# Performers



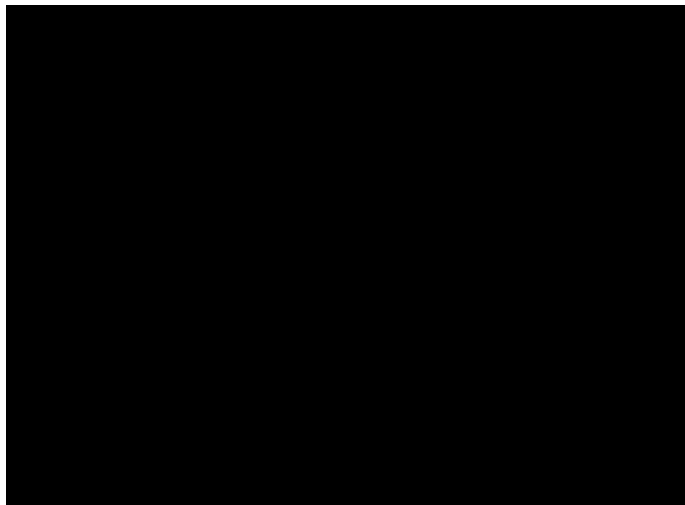
## Towards New Transformers' Revolution

Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, David Belanger, Lucy Colwell, Adrian Weller, Vikas Sindhvani

["Rethinking Attention with Performers"](#) [Google AI Blog Post](#) [code](#)

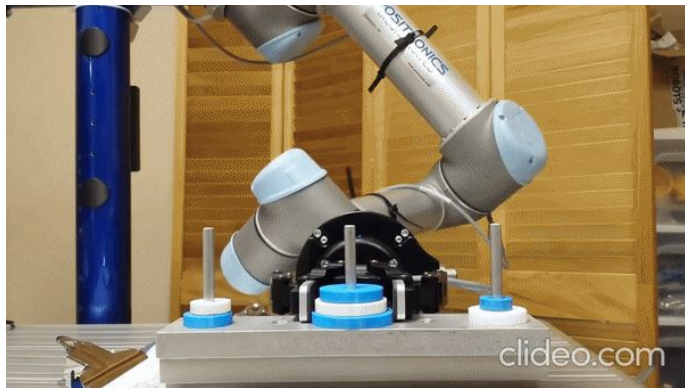
*to appear @ ICLR 2021, oral presentation,  
top 40 accepted papers*

# Why we need better Memorization & Attention in ML?



temporal/spatial attention

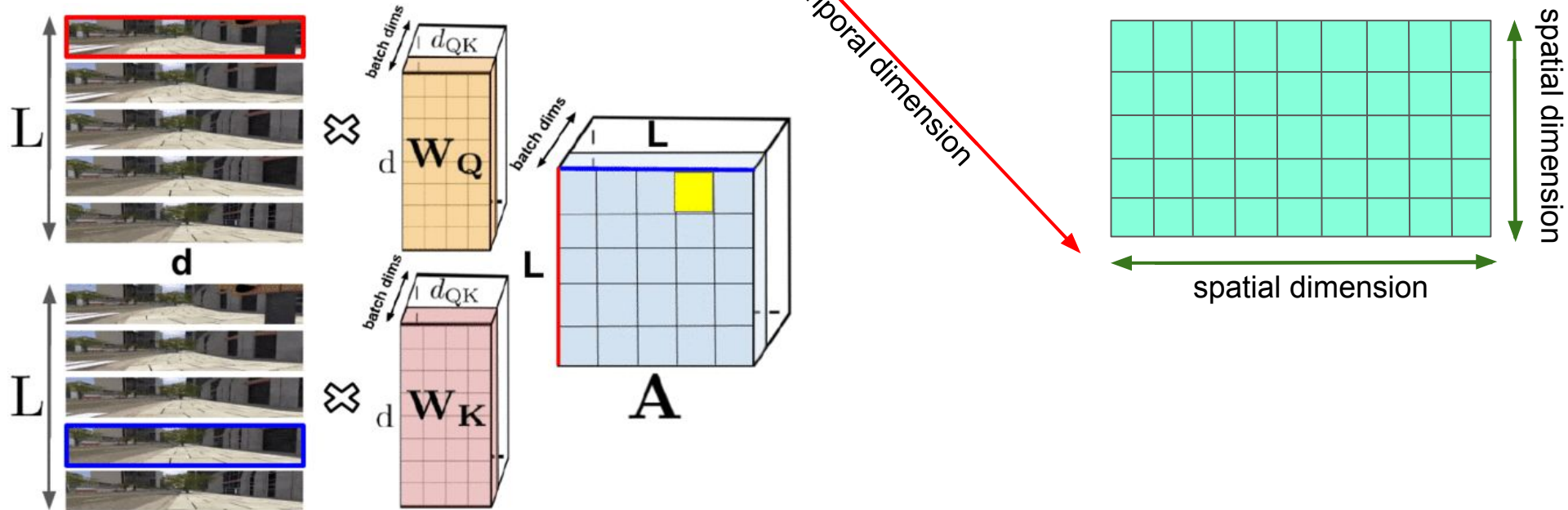
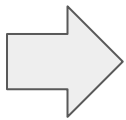
Fig. 1 DeepMind policy navigating simply by “sight”.

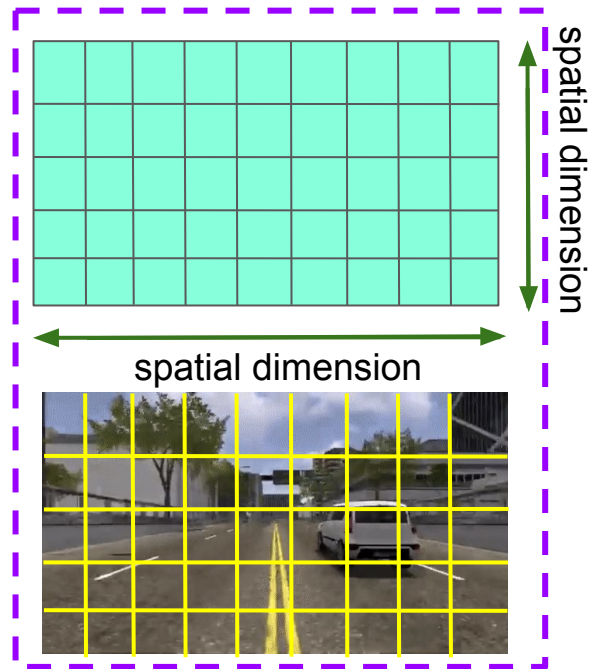
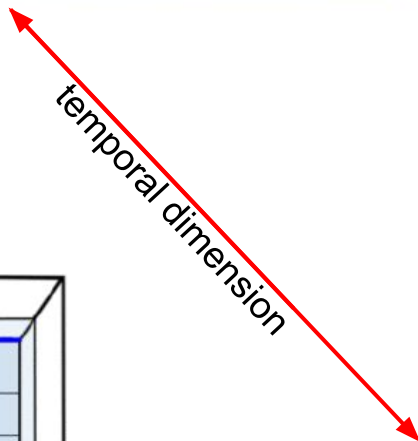
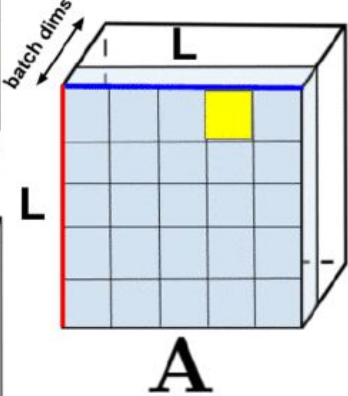
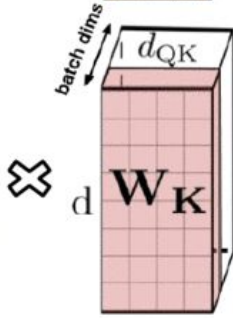
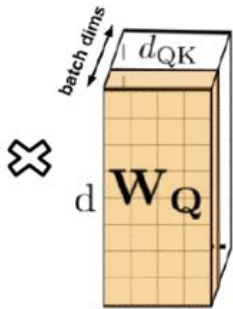
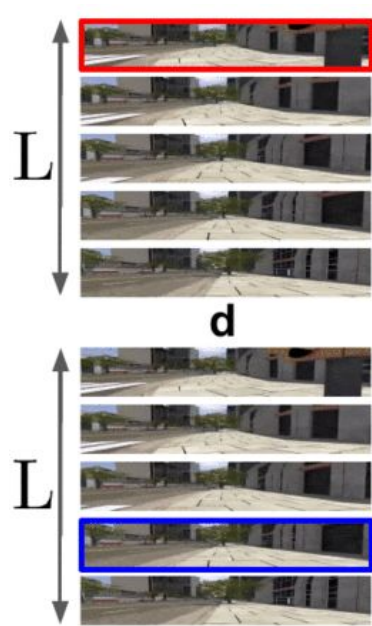
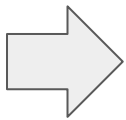


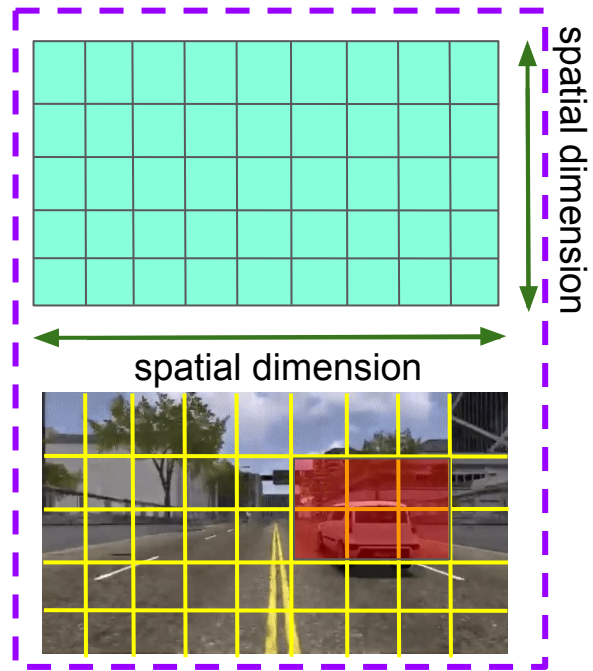
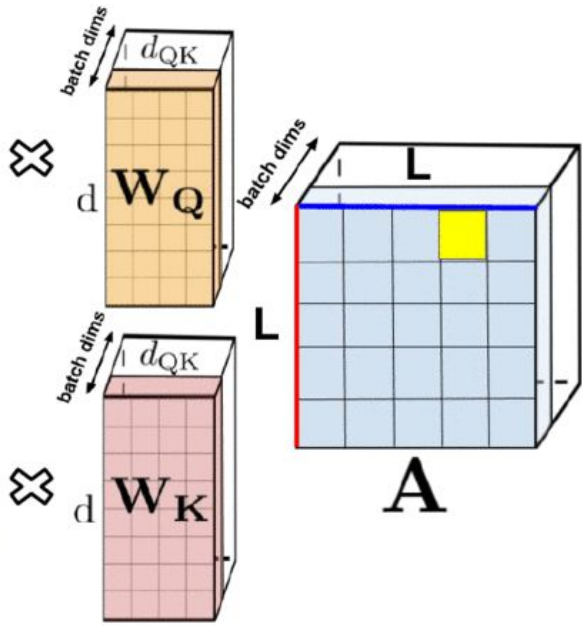
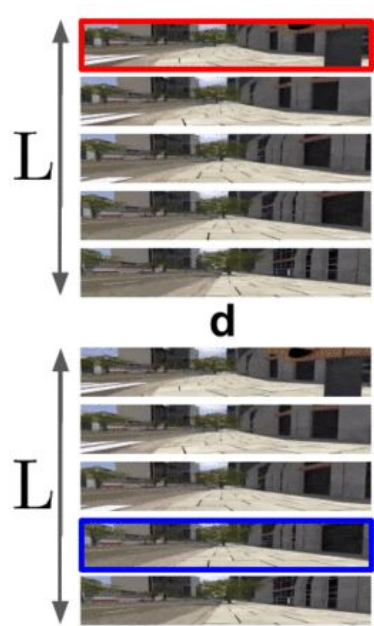
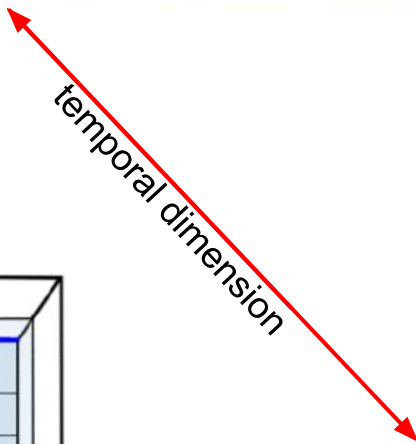
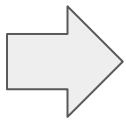
temporal attention

Fig. 2 Robotic arm “solving” Hanoi towers.

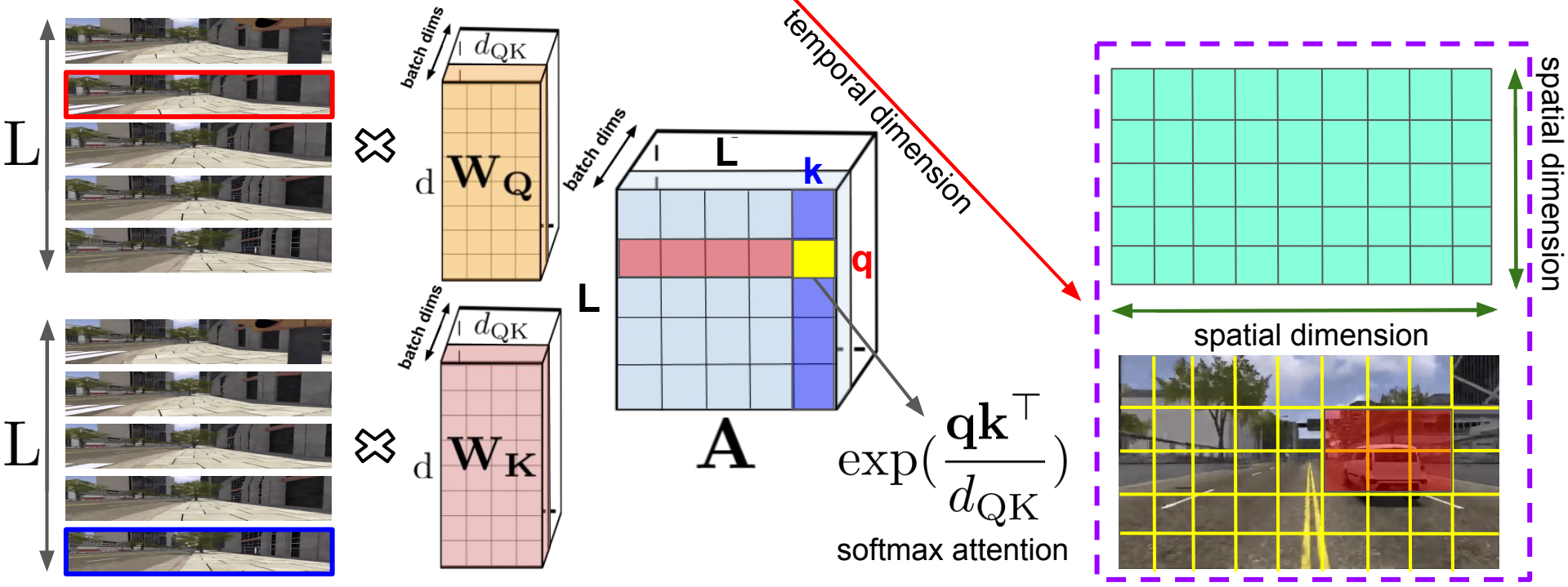
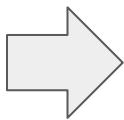
- *Lifelong Learning* requires going beyond purely-reactive Robotics tasks:  
  
*“Developmental Robotics: A Complex Dynamical System with Several Spatiotemporal scales.”*
- Memory is a key to AI and currently existing sequential recurrent architectures fail to memorize well.
- Lets learn how to attend to the world - is [Attention](#) all you need ? Different attention dimensions: **spatial & temporal**.
- Standard attention mechanisms are effectively parallelizable and avoid catastrophic forgetting, but are not scalable.  
  
*“It uses more memory and more computation per real interaction...”* [DeepMind nav by sight]
- Lifelong Learning Robotics requires **long range contexts** with **no attention priors**.

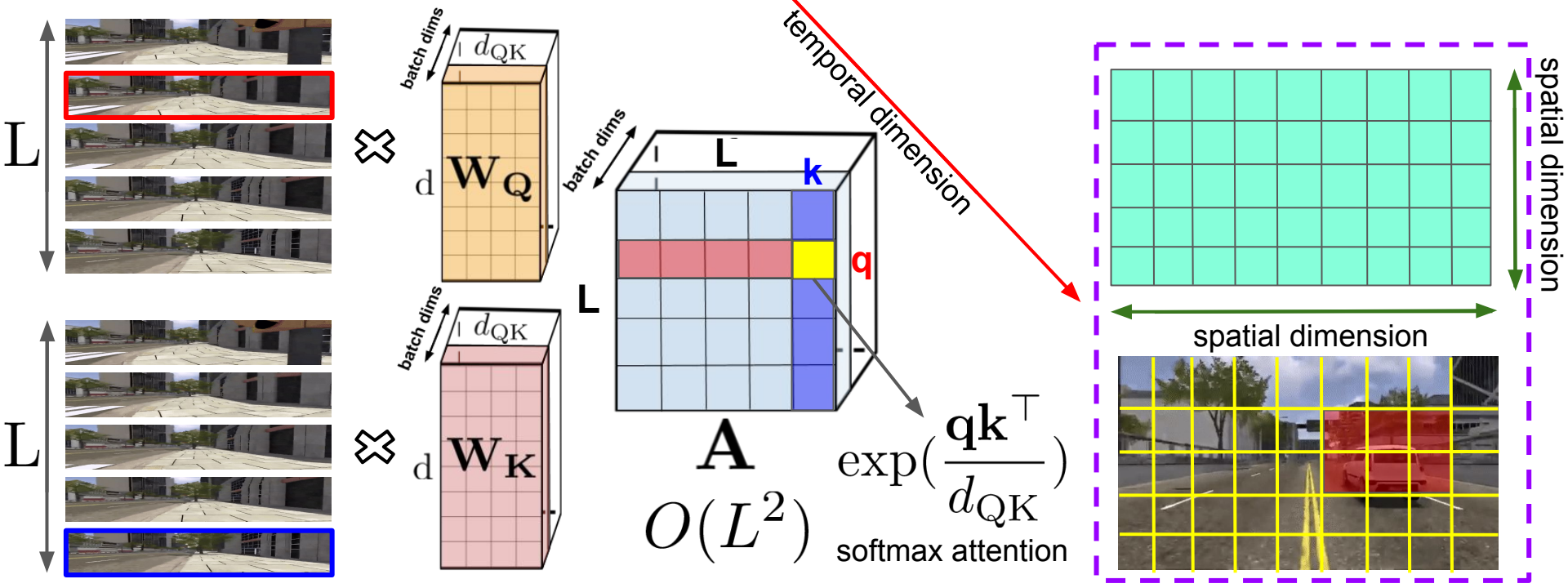
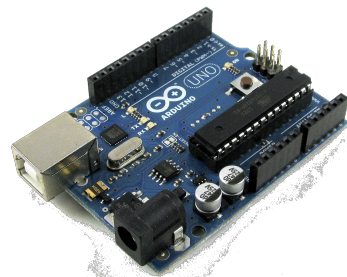
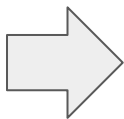




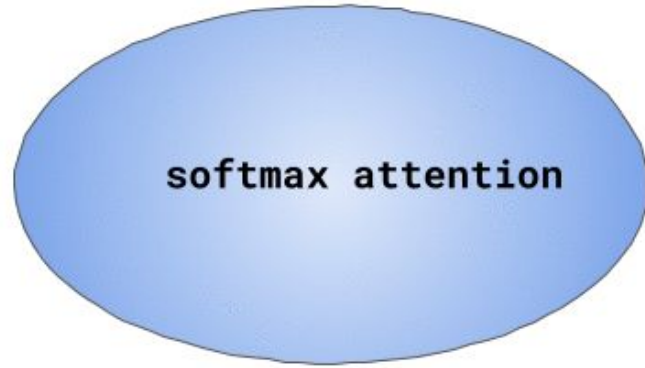












Performers = **Linear** Attention Modules Capable of Modeling  
Different Attention Kernels (Also **SOFTMAX** !)

# Performers performing on the Long Range Arena

## Long Range Arena Paper



| Model           | ListOps      | Text         | Retrieval    | Image        | Pathfinder   | Path-X | Avg          |
|-----------------|--------------|--------------|--------------|--------------|--------------|--------|--------------|
| Transformer     | 36.37        | 64.27        | 57.46        | 42.44        | 71.40        | FAIL   | <u>54.39</u> |
| Local Attention | 15.82        | 52.98        | 53.39        | 41.46        | 66.63        | FAIL   | 46.06        |
| Sparse Trans.   | 17.07        | 63.58        | <b>59.59</b> | <b>44.24</b> | 71.71        | FAIL   | 51.24        |
| Longformer      | 35.63        | 62.85        | 56.89        | 42.22        | 69.71        | FAIL   | 53.46        |
| Linformer       | 35.70        | 53.94        | 52.27        | 38.56        | <u>76.34</u> | FAIL   | 51.36        |
| Reformer        | <b>37.27</b> | 56.10        | 53.40        | 38.07        | 68.50        | FAIL   | 50.67        |
| Sinkhorn Trans. | 33.67        | 61.20        | 53.83        | 41.23        | 67.45        | FAIL   | 51.39        |
| Synthesizer     | <u>36.99</u> | 61.68        | 54.67        | 41.61        | 69.45        | FAIL   | 52.88        |
| BigBird         | 36.05        | 64.02        | <u>59.29</u> | 40.83        | 74.87        | FAIL   | <b>55.01</b> |
| Linear Trans.   | 16.13        | <b>65.90</b> | 53.09        | 42.34        | 75.30        | FAIL   | 50.55        |
| Performer       | 18.01        | <u>65.40</u> | 53.82        | <u>42.77</u> | <b>77.05</b> | FAIL   | <u>51.41</u> |
| Task Avg (Std)  | 29 (9.7)     | 61 (4.6)     | 55 (2.6)     | 41 (1.8)     | 72 (3.7)     | FAIL   | 52 (2.4)     |

Table 1: Experimental results on Long-Range Arena benchmark. Best model is in boldface and second best is underlined. All models do not learn anything on Path-X task, contrary to the Pathfinder task and this is denoted by FAIL. This shows that increasing the sequence length can cause seriously difficulties for model training. We leave Path-X on this benchmark for future challengers but do not include it on the Average score as it has no impact on relative performance.



| Model           | Steps per second  |                   |                   |                   | Peak Memory Usage (GB) |             |             |             |
|-----------------|-------------------|-------------------|-------------------|-------------------|------------------------|-------------|-------------|-------------|
|                 | 1K                | 2K                | 3K                | 4K                | 1K                     | 2K          | 3K          | 4K          |
| Transformer     | 8.1               | 4.9               | 2.3               | 1.4               | 0.85                   | 2.65        | 5.51        | 9.48        |
| Local Attention | 9.2 (1.1x)        | 8.4 (1.7x)        | 7.4 (3.2x)        | 7.4 (5.3x)        | 0.42                   | 0.76        | 1.06        | 1.37        |
| Linformer       | <u>9.3</u> (1.2x) | 9.1 (1.9x)        | 8.5 (3.7x)        | 7.7 (5.5x)        | <b>0.37</b>            | <b>0.55</b> | 0.99        | <b>0.99</b> |
| Reformer        | 4.4 (0.5x)        | 2.2 (0.4x)        | 1.5 (0.7x)        | 1.1 (0.8x)        | 0.48                   | 0.99        | 1.53        | 2.28        |
| Sinkhorn Trans  | 9.1 (1.1x)        | 7.9 (1.6x)        | 6.6 (2.9x)        | 5.3 (3.8x)        | 0.47                   | 0.83        | 1.13        | 1.48        |
| Synthesizer     | 8.7 (1.1x)        | 5.7 (1.2x)        | 6.6 (2.9x)        | 1.9 (1.4x)        | 0.65                   | 1.98        | 4.09        | 6.99        |
| BigBird         | 7.4 (0.9x)        | 3.9 (0.8x)        | 2.7 (1.2x)        | 1.5 (1.1x)        | 0.77                   | 1.49        | 2.18        | 2.88        |
| Linear Trans.   | 9.1 (1.1x)        | 9.3 (1.9x)        | <u>8.6</u> (3.7x) | 7.8 (5.6x)        | <b>0.37</b>            | <u>0.57</u> | <b>0.80</b> | <u>1.03</u> |
| Performer       | <b>9.5</b> (1.2x) | <b>9.4</b> (1.9x) | <b>8.7</b> (3.8x) | <b>8.0</b> (5.7x) | <b>0.37</b>            | 0.59        | <u>0.82</u> | 1.06        |

Table 2: Benchmark results of all Xformer models with a consistent batch size of 32 across all models. We report relative speed increase/decrease in comparison with the vanilla Transformer in brackets besides the steps per second. Memory usage refers to per device memory usage across each TPU device. Benchmarks are run on 4x4 TPU V3 Chips.

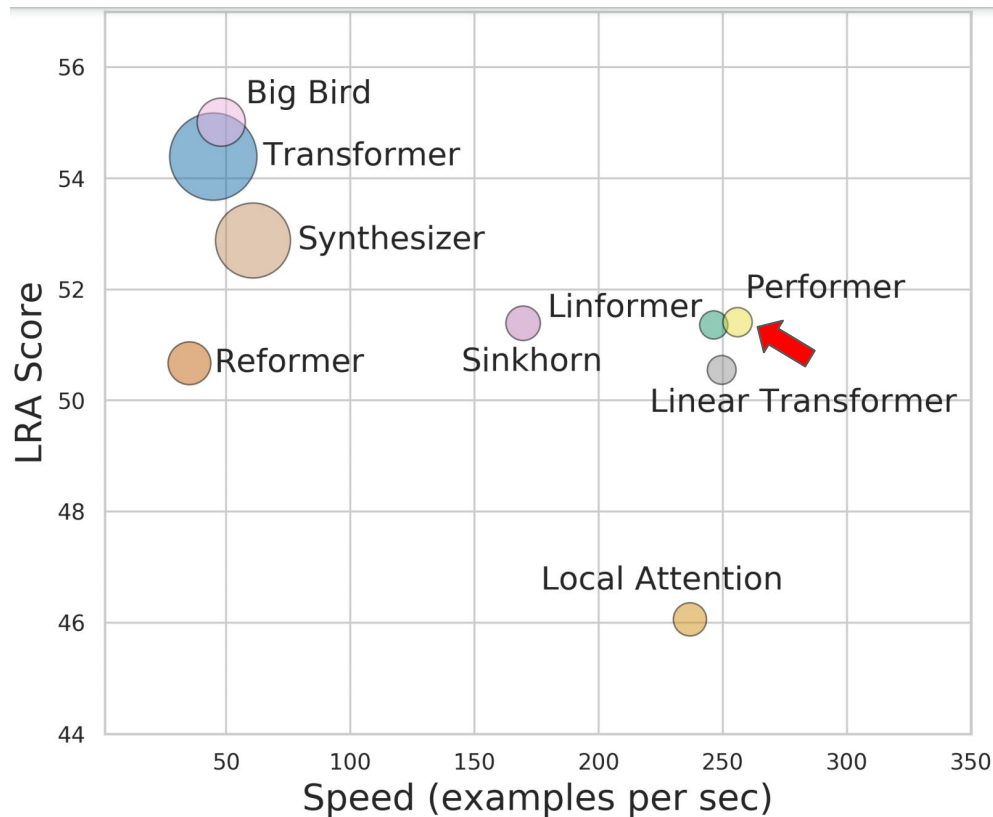


Figure 3: Performance ( $y$  axis), speed ( $x$  axis), and memory footprint (size of the circles) of different models.



# Performers in Vision

## Tokens-to-Token ViT: Training Vision Transformers from Scratch on ImageNet

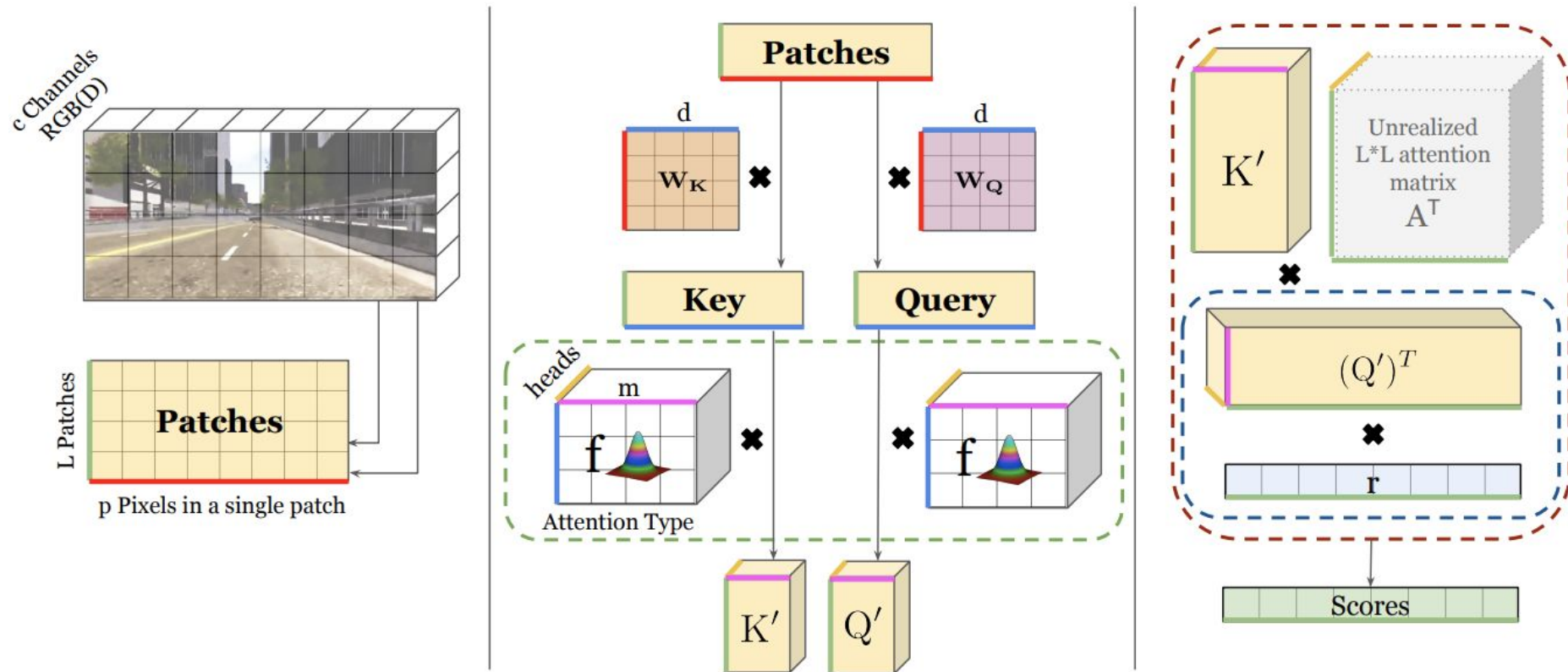
Yuan et al.

| Models                        | Top1-Acc (%) | Params (M) | MACs (G) |
|-------------------------------|--------------|------------|----------|
| ResNet50 [15]                 | 76.2         | 25.5       | 4.3      |
| ResNet50*                     | 79.1         | 25.5       | 4.3      |
| <b>T2T-ViT-14</b>             | <b>80.6</b>  | 21.4       | 4.8      |
| <b>T2T-ViT<sub>t</sub>-14</b> | <b>80.7</b>  | 21.5       | 5.2      |
| ResNet101 [15]                | 77.4         | 44.6       | 7.9      |
| ResNet101*                    | 79.9         | 44.6       | 7.9      |
| <b>T2T-ViT-19</b>             | <b>81.2</b>  | 39.0       | 8.0      |
| <b>T2T-ViT<sub>t</sub>-19</b> | <b>81.4</b>  | 39.0       | 8.4      |
| ResNet152 [15]                | 78.3         | 60.2       | 11.6     |
| ResNet152*                    | 80.8         | 60.2       | 11.6     |
| <b>T2T-ViT-24</b>             | <b>81.8</b>  | 63.9       | 12.6     |
| <b>T2T-ViT<sub>t</sub>-24</b> | <b>82.2</b>  | 64.1       | 13.2     |

| Models                        | Tokens-to-Token module |       |            |          | T2T-ViT backbone |            |          | Model size |          |
|-------------------------------|------------------------|-------|------------|----------|------------------|------------|----------|------------|----------|
|                               | T2T Transformer        | Depth | Hidden dim | MLP size | Depth            | Hidden dim | MLP size | Params (M) | MACs (G) |
| ViT-S/16 [12]                 | -                      | -     | -          | -        | 8                | 786        | 2358     | 48.6       | 10.1     |
| ViT-B/16 [12]                 | -                      | -     | -          | -        | 12               | 786        | 3072     | 86.8       | 17.6     |
| ViT-L/16 [12]                 | -                      | -     | -          | -        | 24               | 1024       | 4096     | 304.3      | 63.6     |
| <b>T2T-ViT<sub>t</sub>-14</b> | Transformer            | 2     | 64         | 64       | 14               | 384        | 1152     | 21.5       | 5.2      |
| <b>T2T-ViT<sub>t</sub>-19</b> | Transformer            | 2     | 64         | 64       | 19               | 448        | 1344     | 39.0       | 8.4      |
| <b>T2T-ViT<sub>t</sub>-24</b> | Transformer            | 2     | 64         | 64       | 24               | 512        | 1536     | 64.1       | 13.2     |
| T2T-ViT-14                    | Performer              | 2     | 64         | 64       | 14               | 384        | 1152     | 21.4       | 4.8      |
| T2T-ViT-19                    | Performer              | 2     | 64         | 64       | 19               | 448        | 1344     | 39.0       | 8.0      |
| T2T-ViT-24                    | Performer              | 2     | 64         | 64       | 24               | 512        | 1536     | 63.9       | 12.6     |
| T2T-ViT-7                     | Performer              | 2     | 64         | 64       | 8                | 256        | 512      | 4.2        | 0.9      |
| T2T-ViT-10                    | Performer              | 2     | 64         | 64       | 10               | 256        | 512      | 5.6        | 1.2      |
| T2T-ViT-12                    | Performer              | 2     | 64         | 64       | 12               | 256        | 512      | 6.8        | 1.4      |

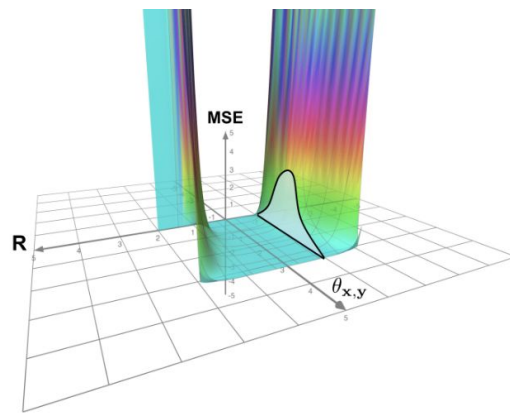
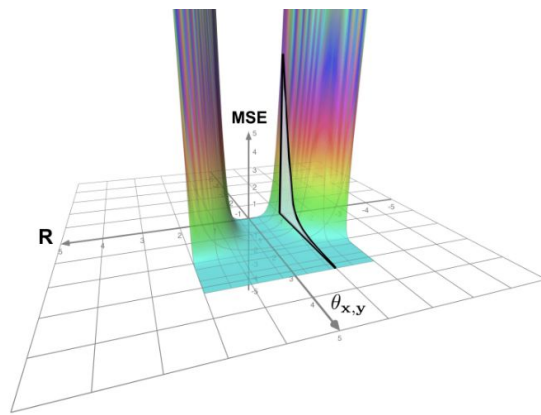
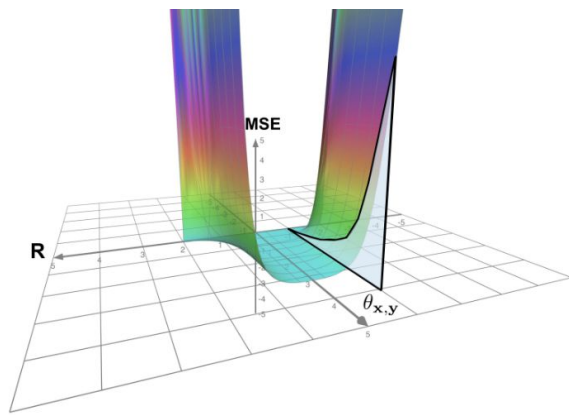
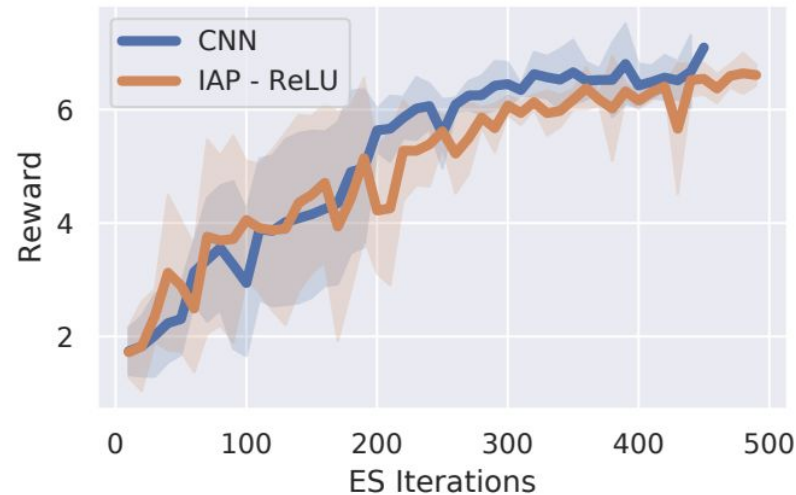
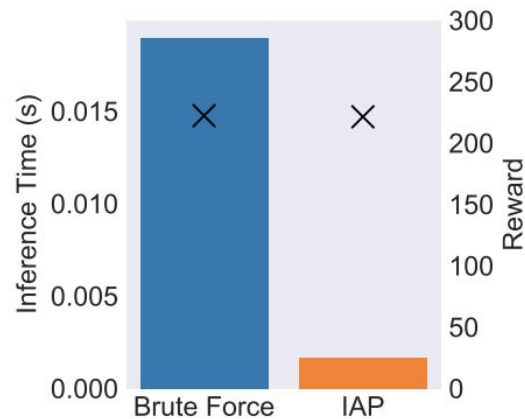
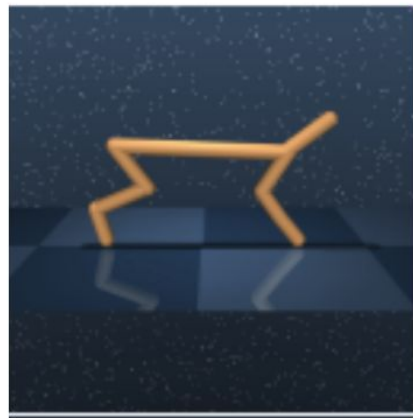
# Performers for Vision in Reinforcement Learning (IAP)

Unlocking Pixels for Reinforcement Learning via Implicit Attention Choromanski et al.



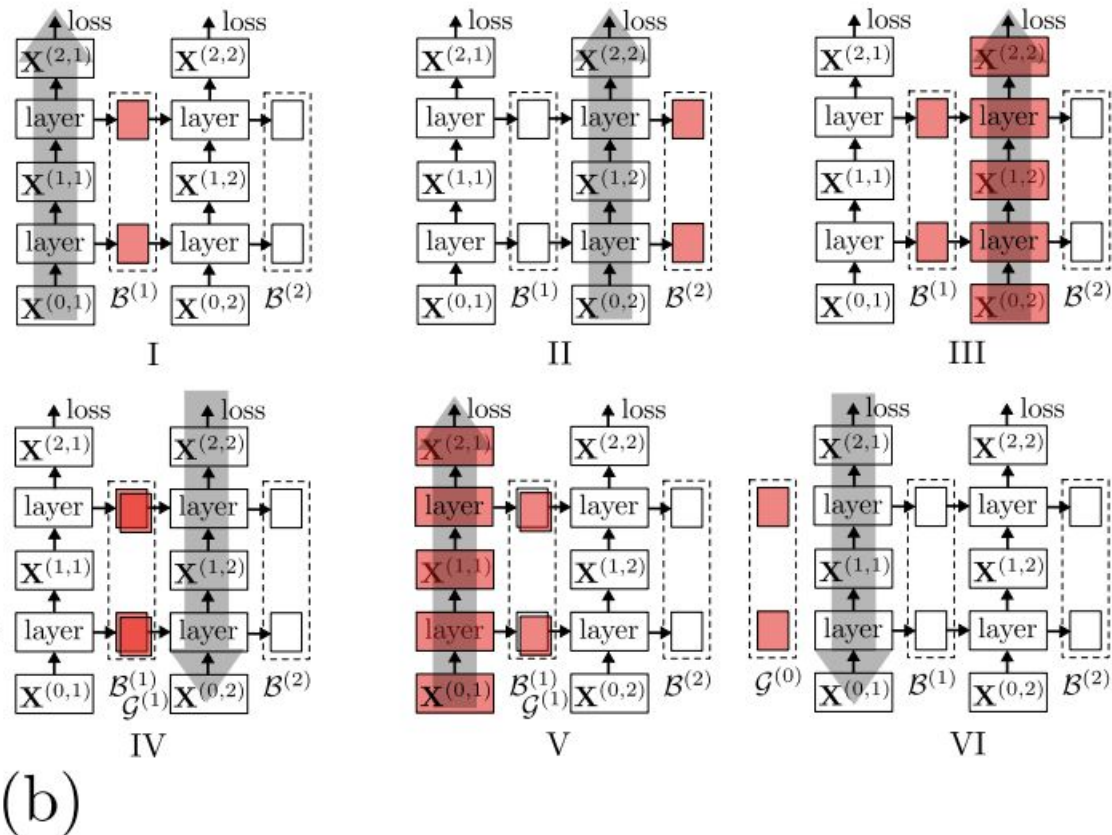
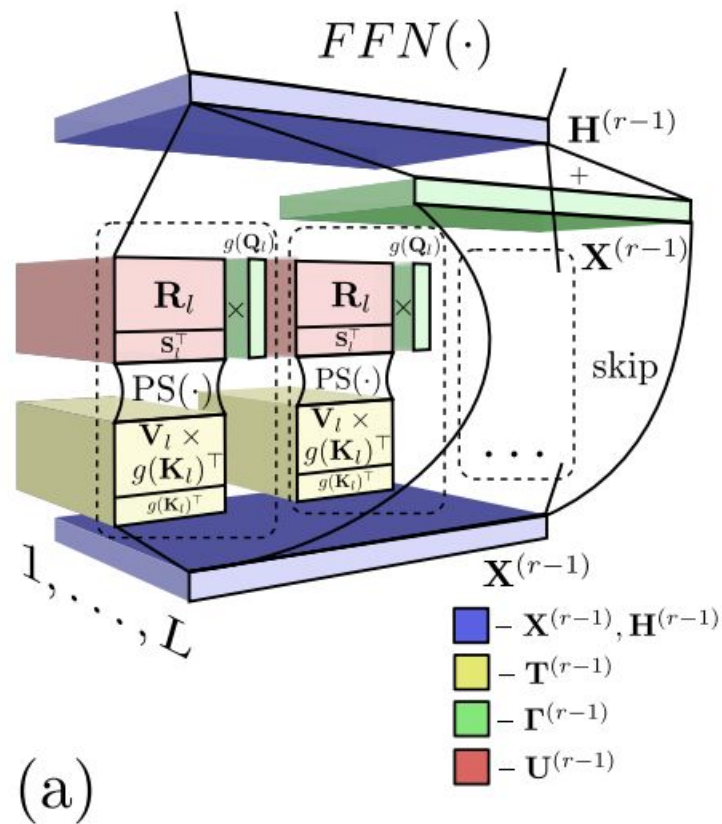
# Performers for Vision in Reinforcement Learning (IAP)

*Unlocking Pixels for Reinforcement Learning via Implicit Attention* Choromanski et al.

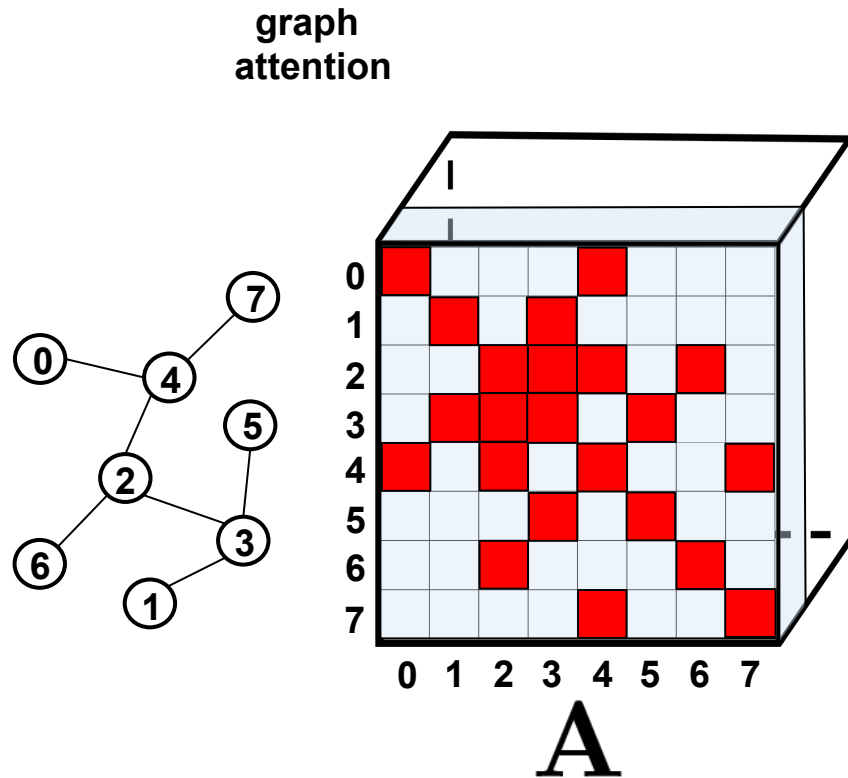
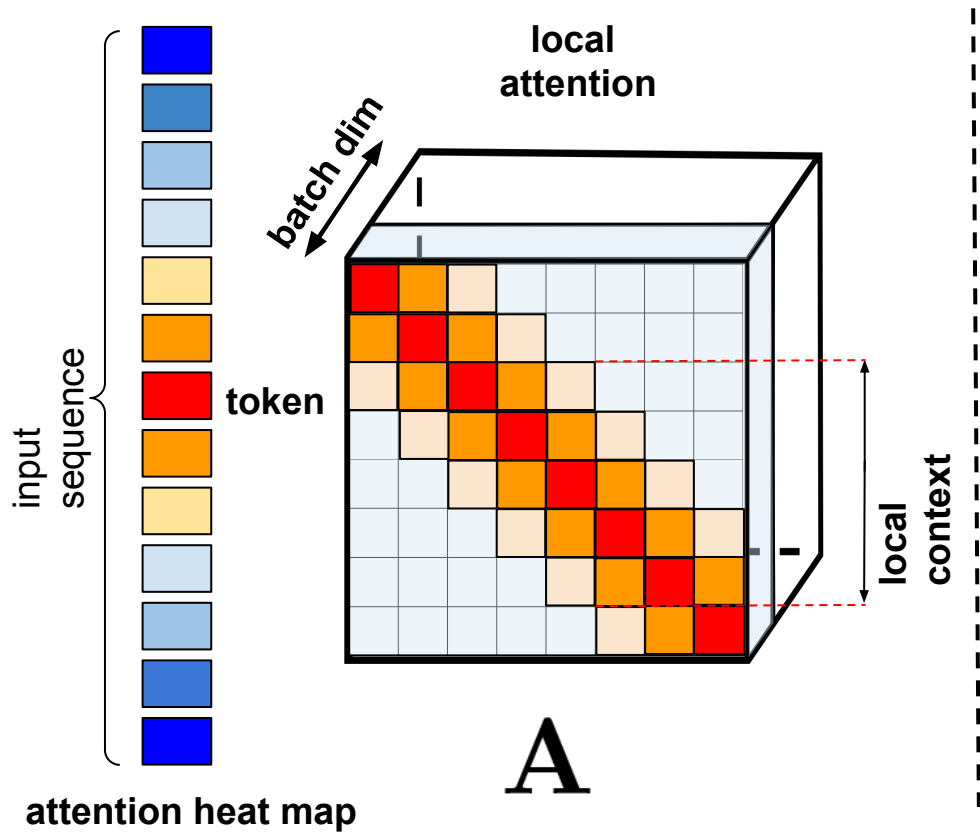


# Sub-Linear Memory: How To Make Performers Slim

*Likhoshesterov et al.*

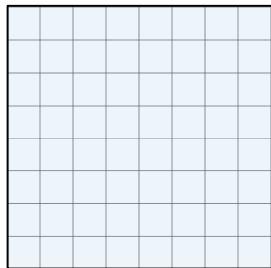


# What we Do NOT DO: Sparsifying Attention



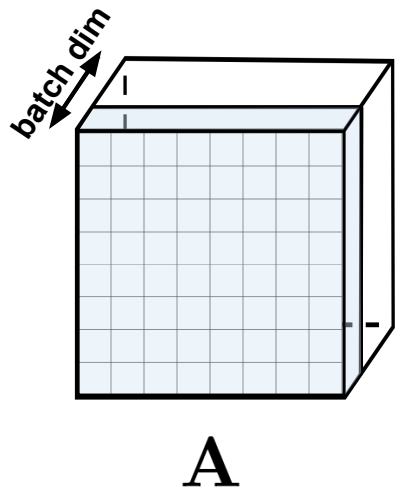


# Attention is Kernelizable - A Tale of Random Maps

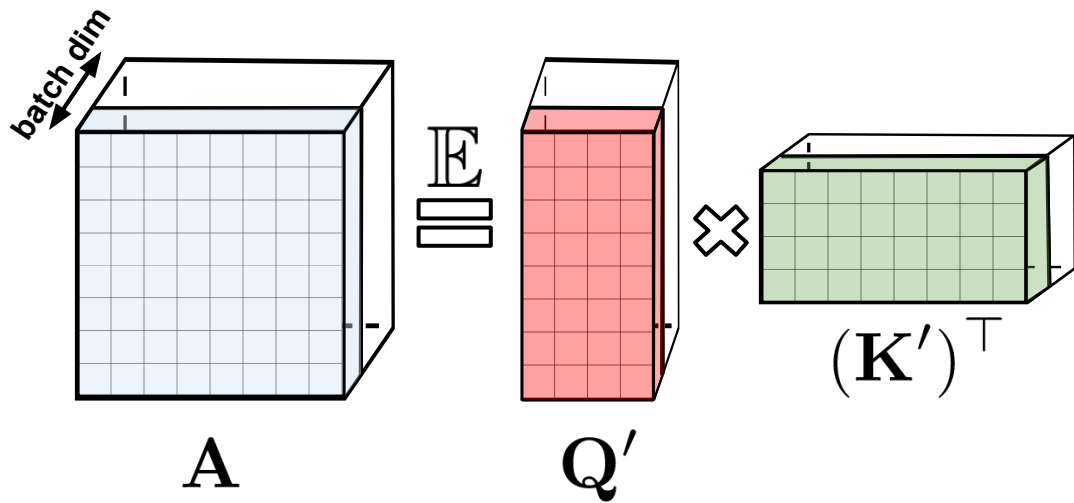


**A**

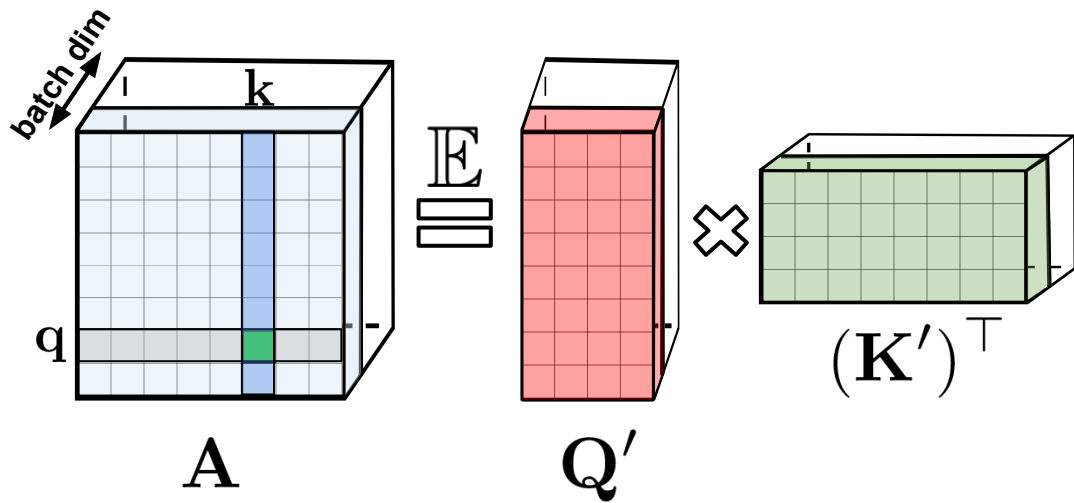
# Attention is Kernelizable - A Tale of Random Maps



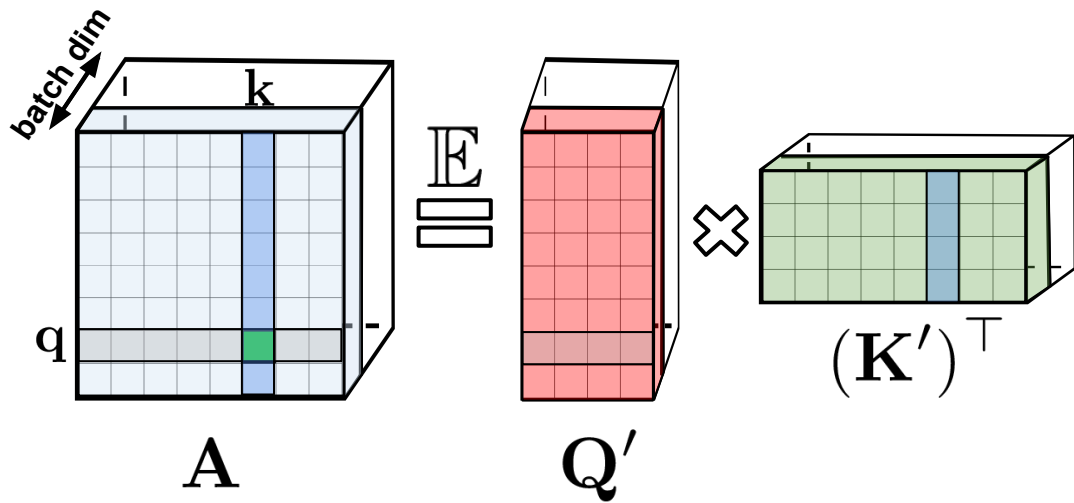
# Attention is Kernelizable - A Tale of Random Maps



# Attention is Kernelizable - A Tale of Random Maps

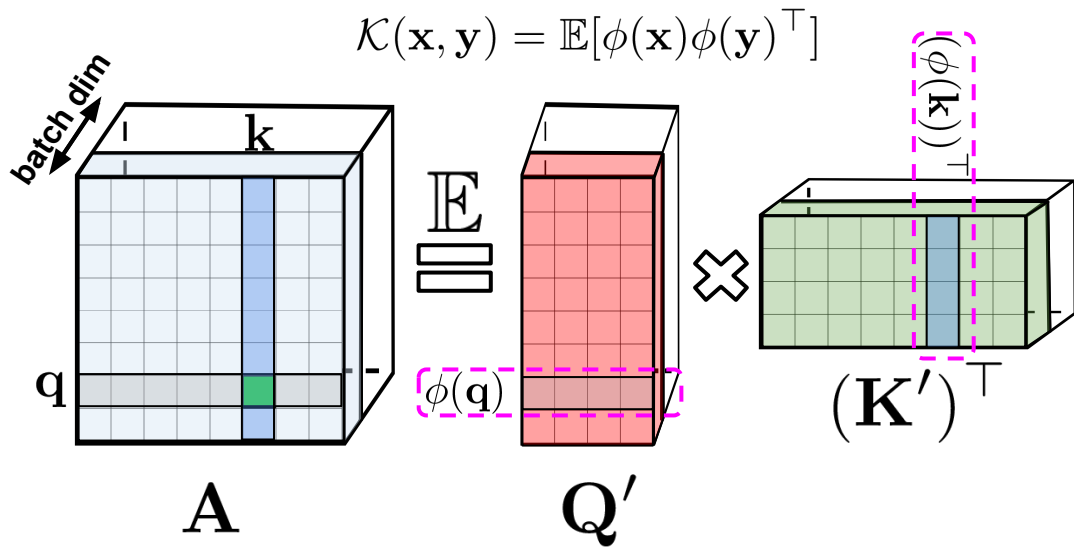


# Attention is Kernelizable - A Tale of Random Maps

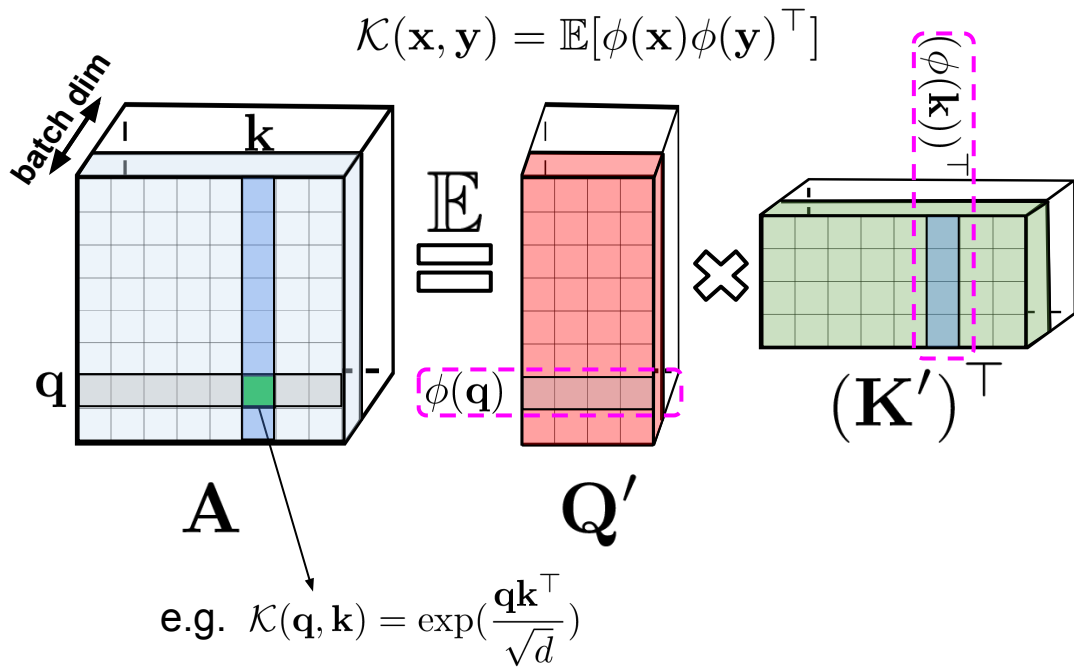




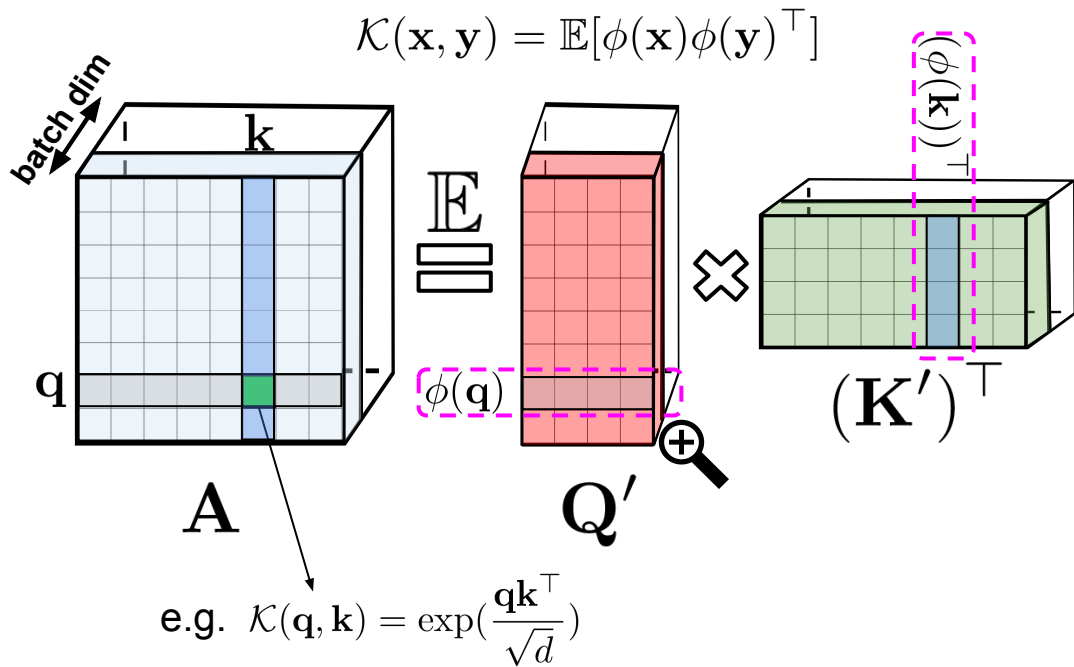
# Attention is Kernelizable - A Tale of Random Maps



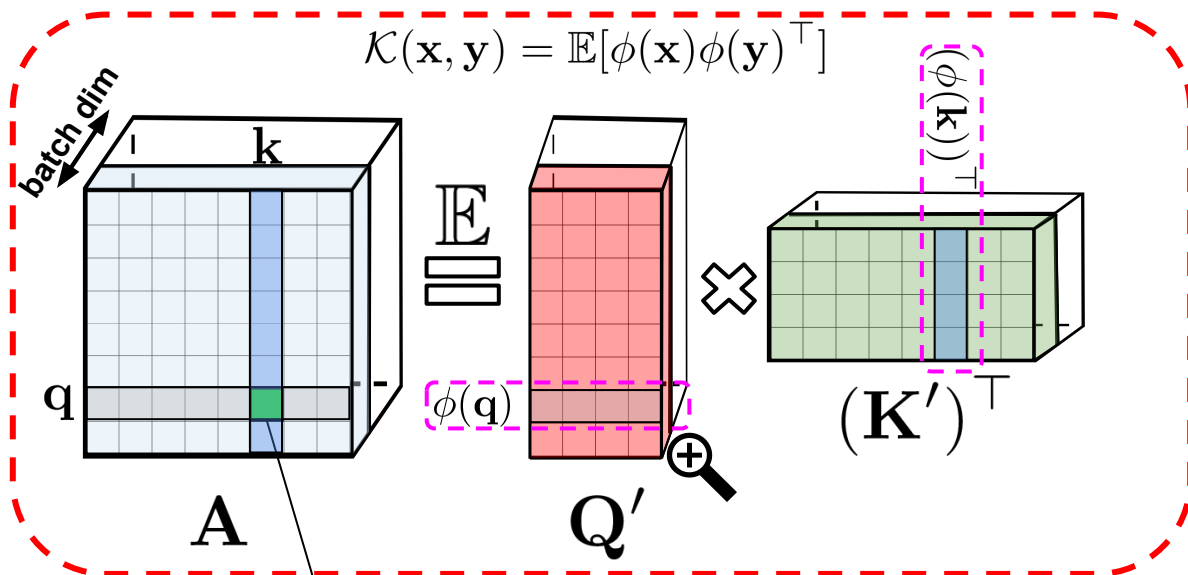
# Attention is Kernelizable - A Tale of Random Maps



# Attention is Kernelizable - A Tale of Random Maps



# Attention is Kernelizable - A Tale of Random Maps



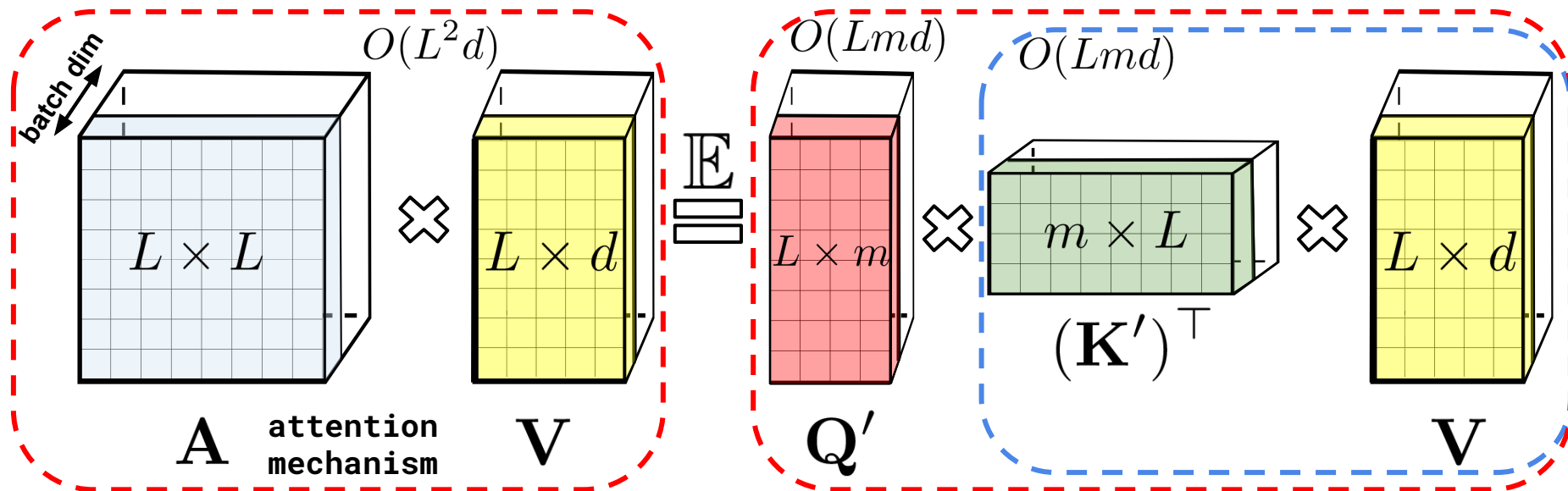
e.g.  $\mathcal{K}(\mathbf{q}, \mathbf{k}) = \exp(\frac{\mathbf{q}\mathbf{k}^\top}{\sqrt{d}})$

$$\left\{ \begin{array}{l} \phi(\mathbf{x}) = \frac{1}{\sqrt{m}} \underbrace{\exp(\frac{\|\mathbf{x}\|_2^2}{2\sqrt{d}})}_{\text{deterministic feature}} \cdot \underbrace{(\sin(\frac{\mathbf{x}}{d^{\frac{1}{4}}}\omega_1^\top), \dots, \sin(\frac{\mathbf{x}}{d^{\frac{1}{4}}}\omega_m^\top), \cos(\frac{\mathbf{x}}{d^{\frac{1}{4}}}\omega_1^\top), \dots, \cos(\frac{\mathbf{x}}{d^{\frac{1}{4}}}\omega_m^\top))}_{\text{random feature}} \\ \omega \sim \mathcal{N}(0, \mathbf{I}_d) \end{array} \right.$$

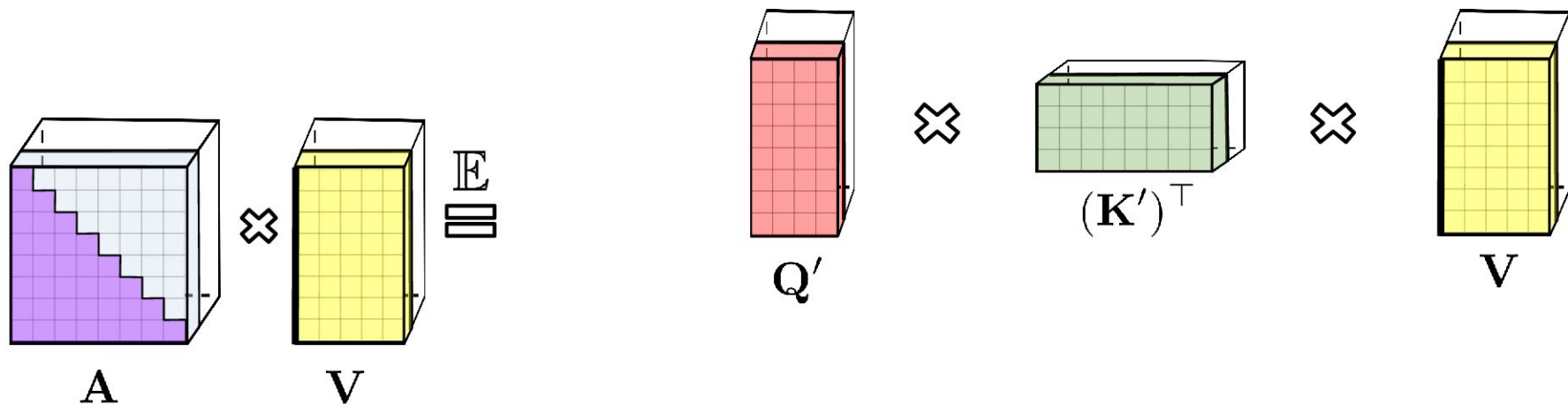
e.g. softmax kernel features

$\phi(\mathbf{q})$

# Associativity for Speedups and Space Compression

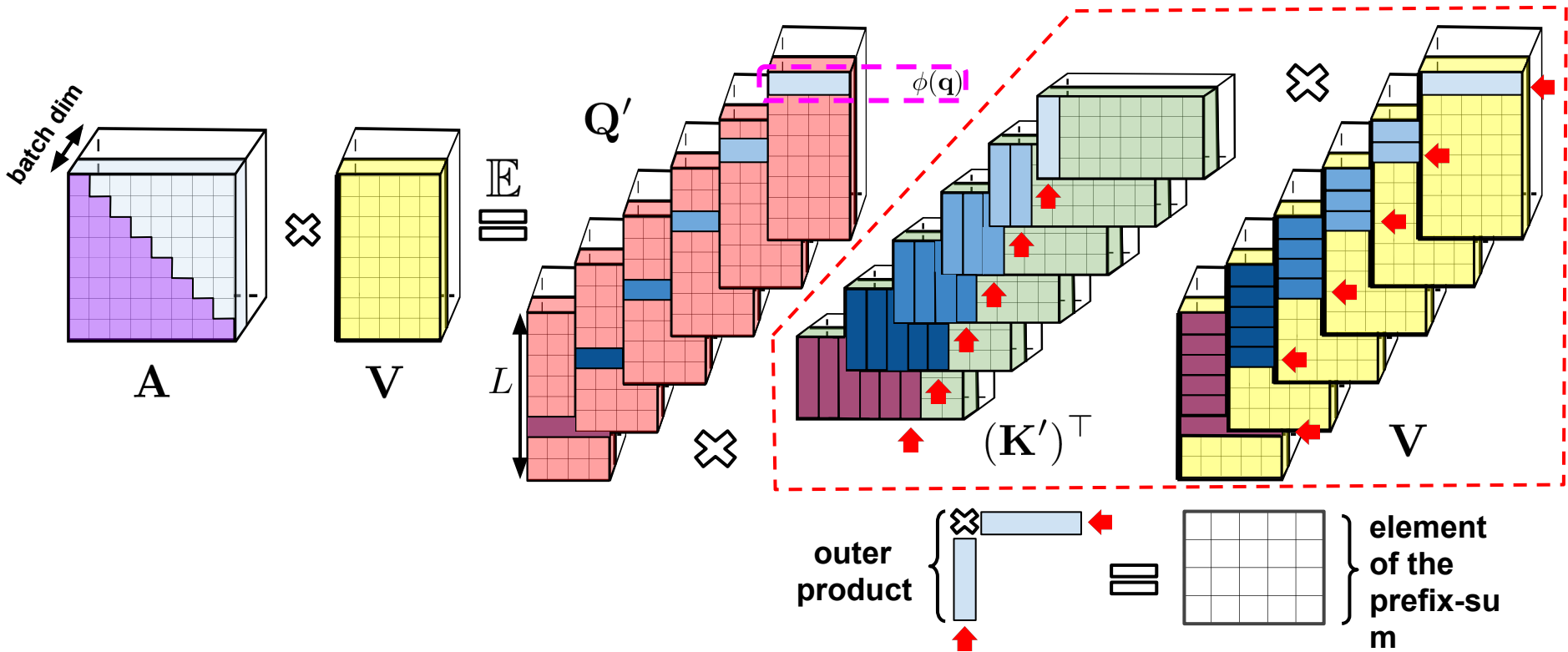


# Causal Transformers as Prefix-Sums Calculators





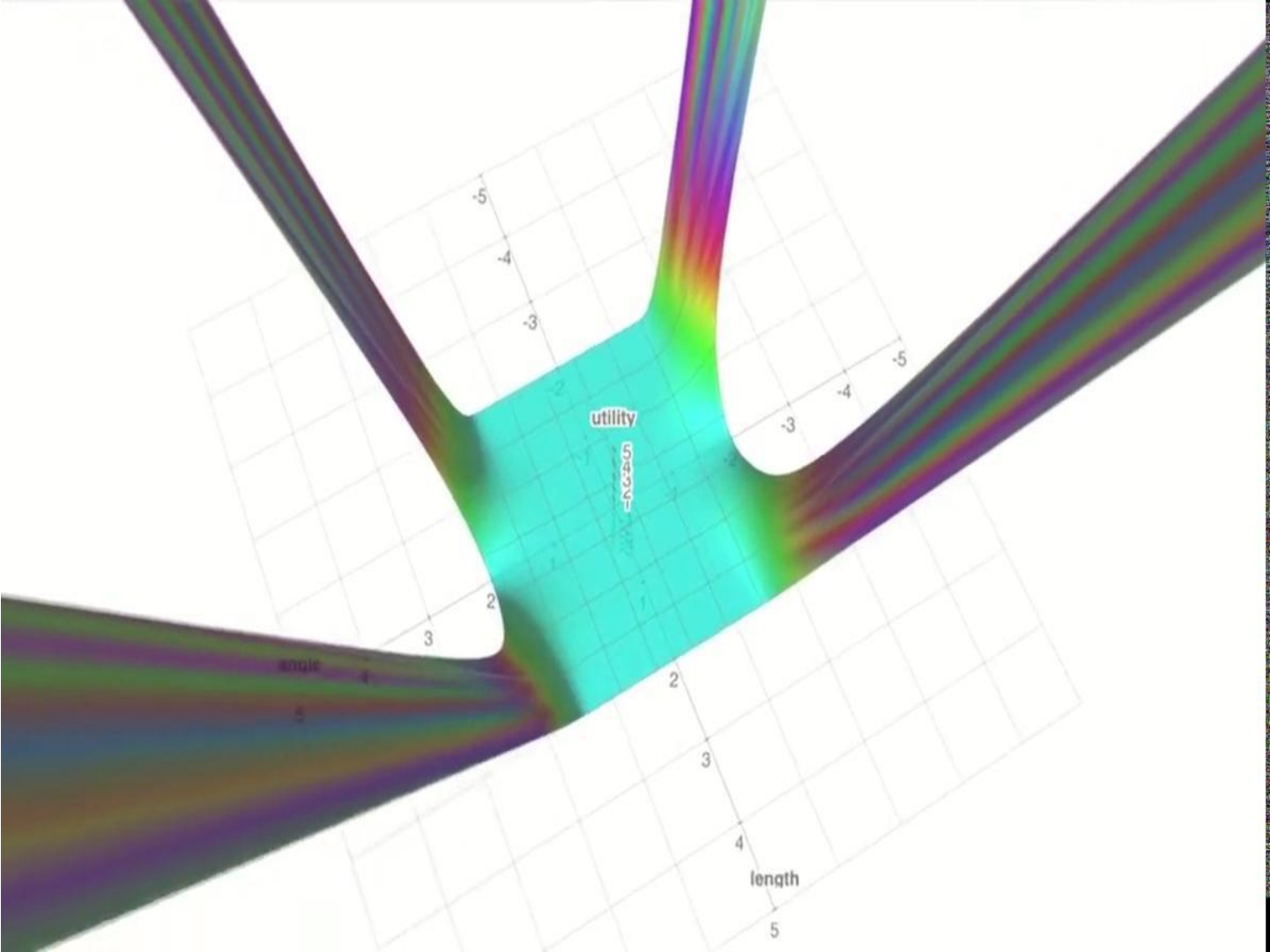
# Causal Transformers as Prefix-Sums Calculators



**Why this is not the end of the  
story ?**

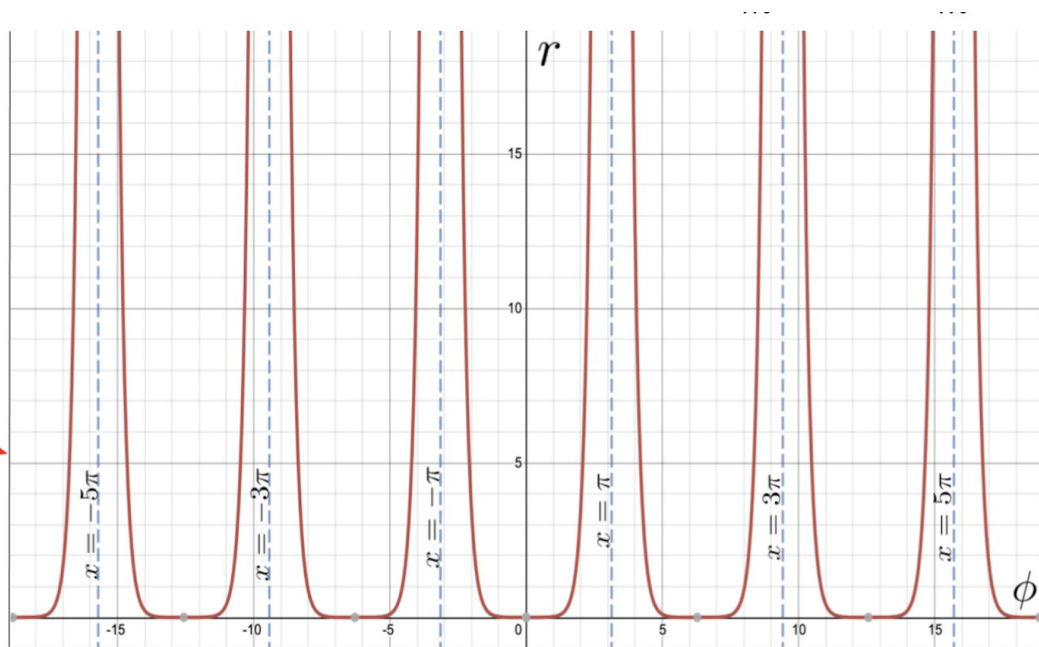
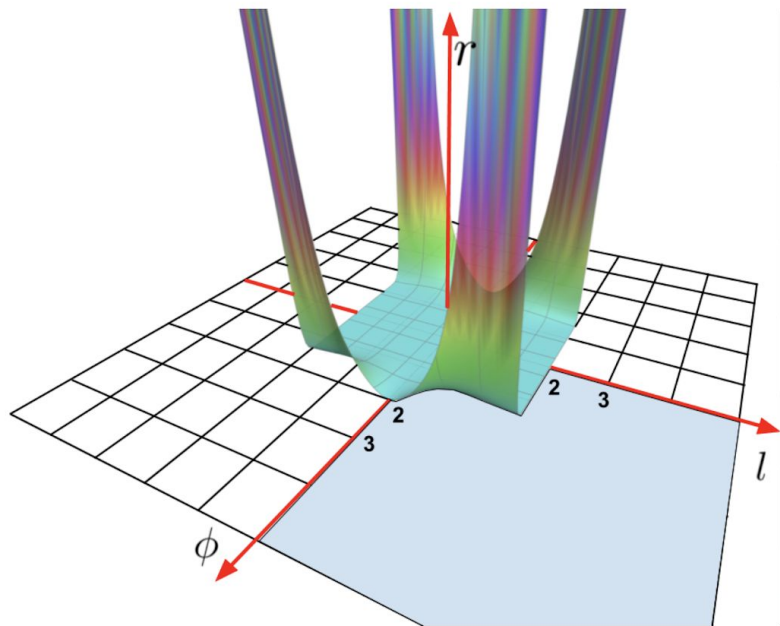
**FAVOR+**

$$\text{utility} = \frac{\text{MSE}(\text{SM}^{\text{trig}}(\mathbf{x}, \mathbf{y}))}{\text{MSE}(\text{SM}^+(\mathbf{x}, \mathbf{y}))}$$



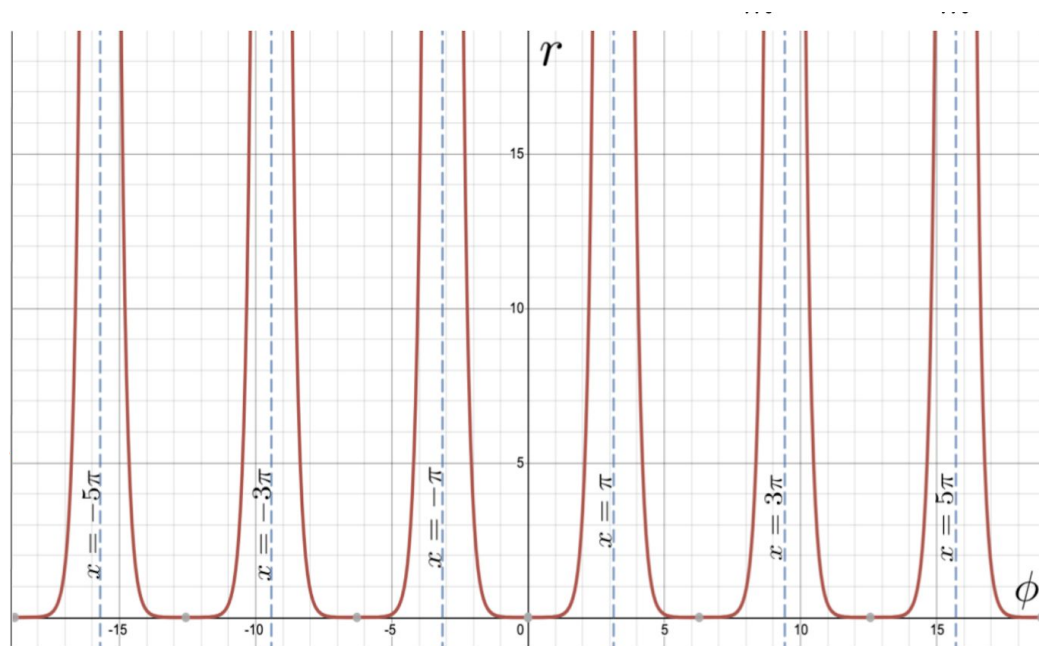
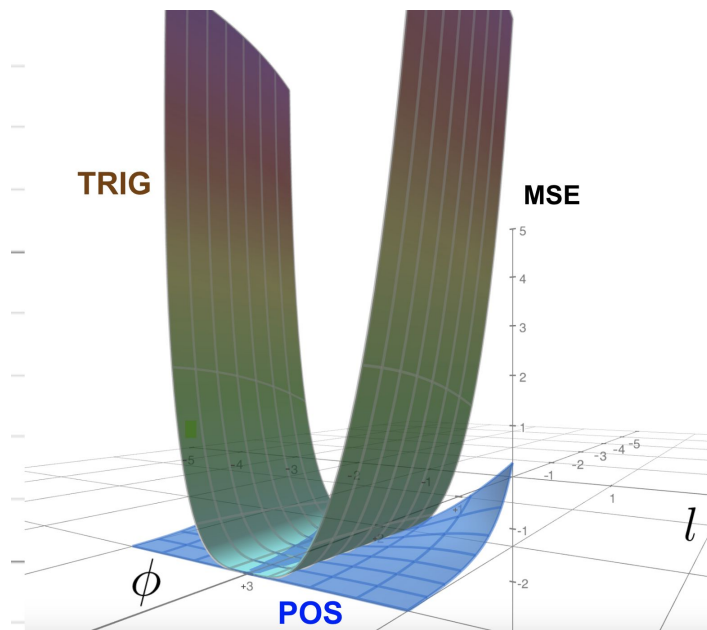
$$\text{length} = \|\mathbf{x}\| = \|\mathbf{y}\|$$

# Do we need to go beyond trigonometric features ?



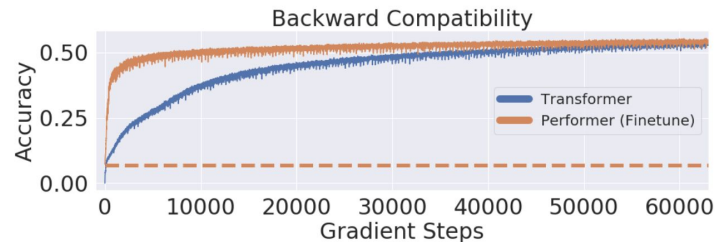
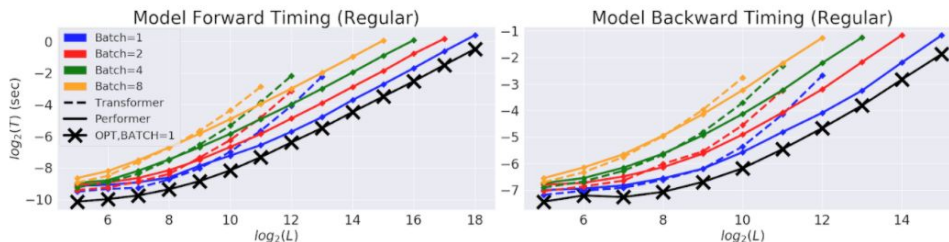
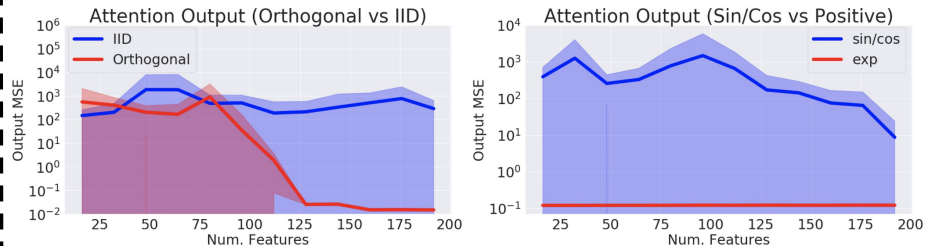
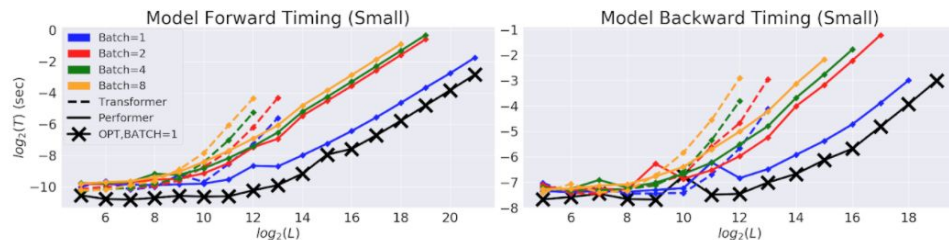
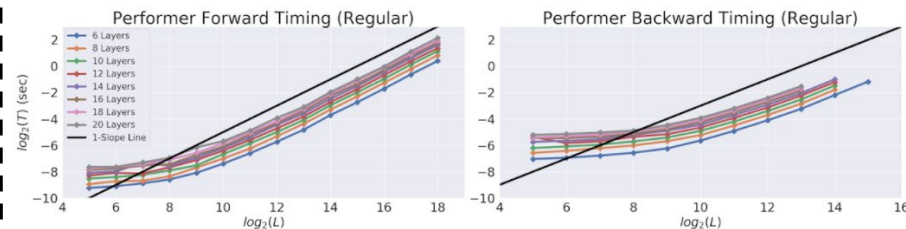
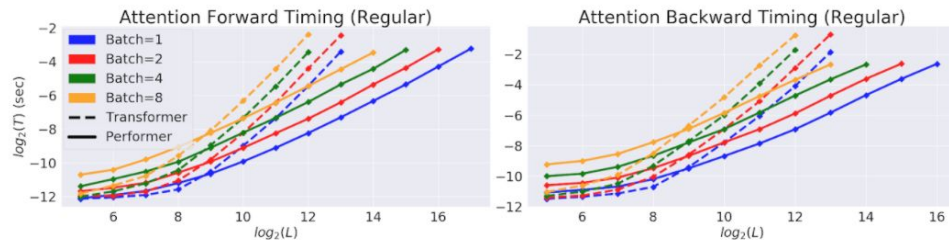
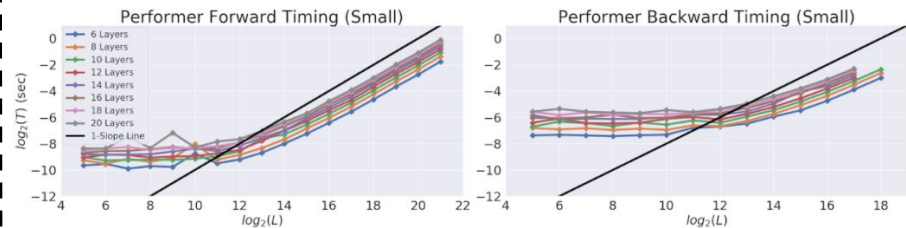
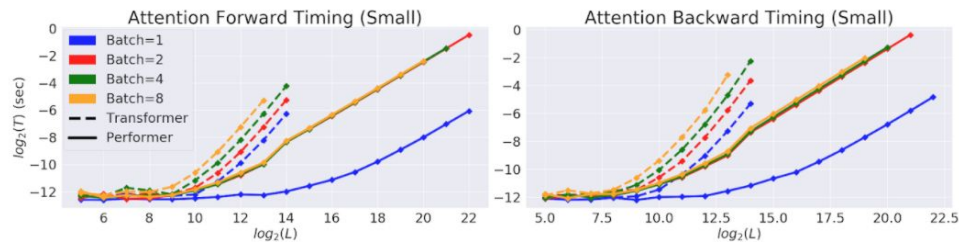
**Left:** Symmetrized (around origin) utility function  $r$  (defined as a ratio of the mean squared errors of estimators built on: trigonometric and positive random features) as a function of the angle  $\phi$  (in radians) between input feature vectors and their lengths  $l$ . Larger values indicate regions of  $(\phi, l)$ -space with better performance of positive random features. We see that for critical regions with  $\phi$  large enough (small enough softmax-kernel values) our method is arbitrarily more accurate than trigonometric random features. Plot presented for domain  $[-\pi, \pi] \times [-2, 2]$ . **Right:** The slice of function  $r$  for fixed  $l = 1$  and varying angle  $\phi$ .

# Do we need to go beyond trigonometric features ?



**Left:** Comparison of the mean squared errors (MSEs) of the estimators applying trigonometric random features (TRIG) and the one leveraging the mechanism of positive random features (POS) in the region of small softmax-kernel values.

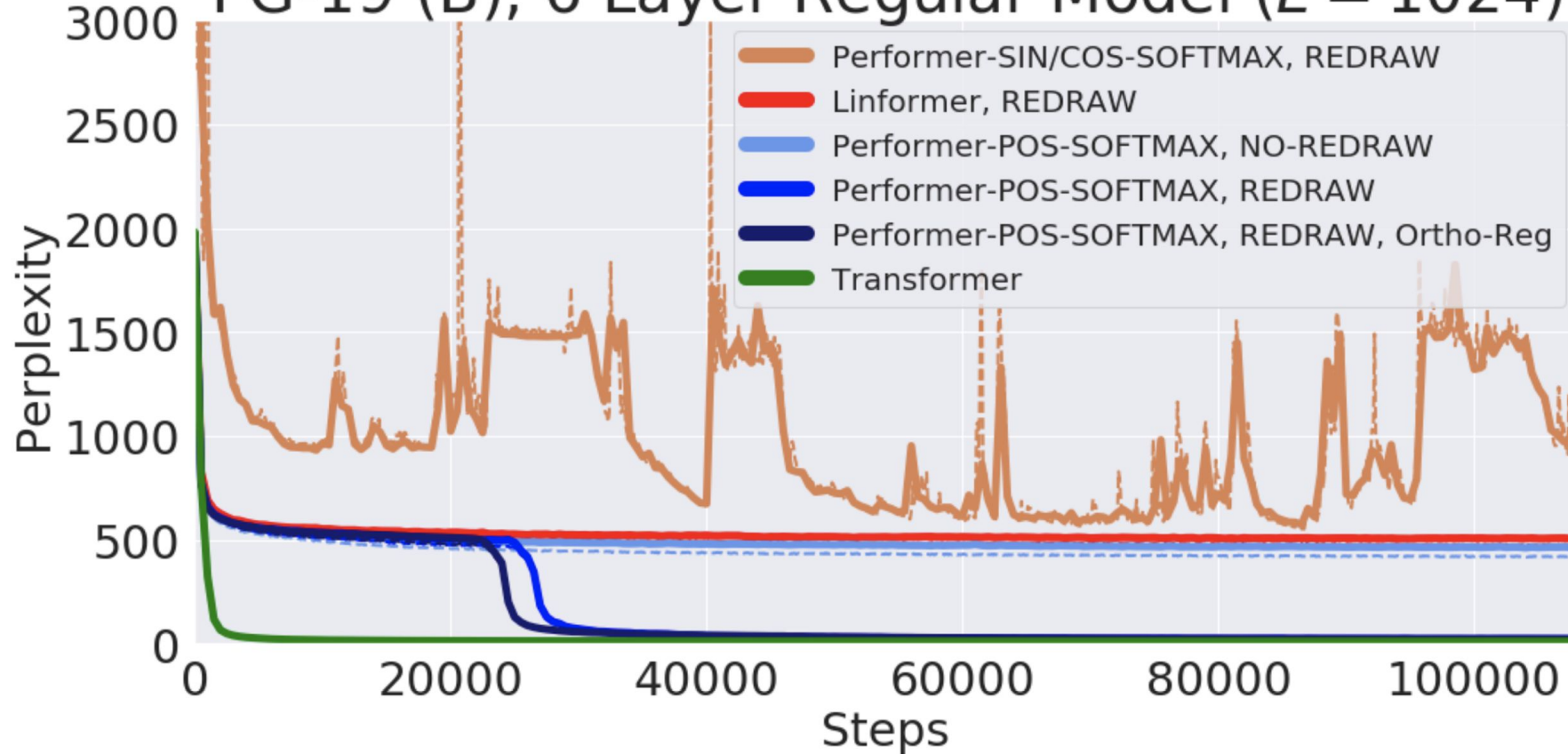
**Right:** The slice of function  $r$  for fixed  $l = 1$  and varying angle  $\phi$ .



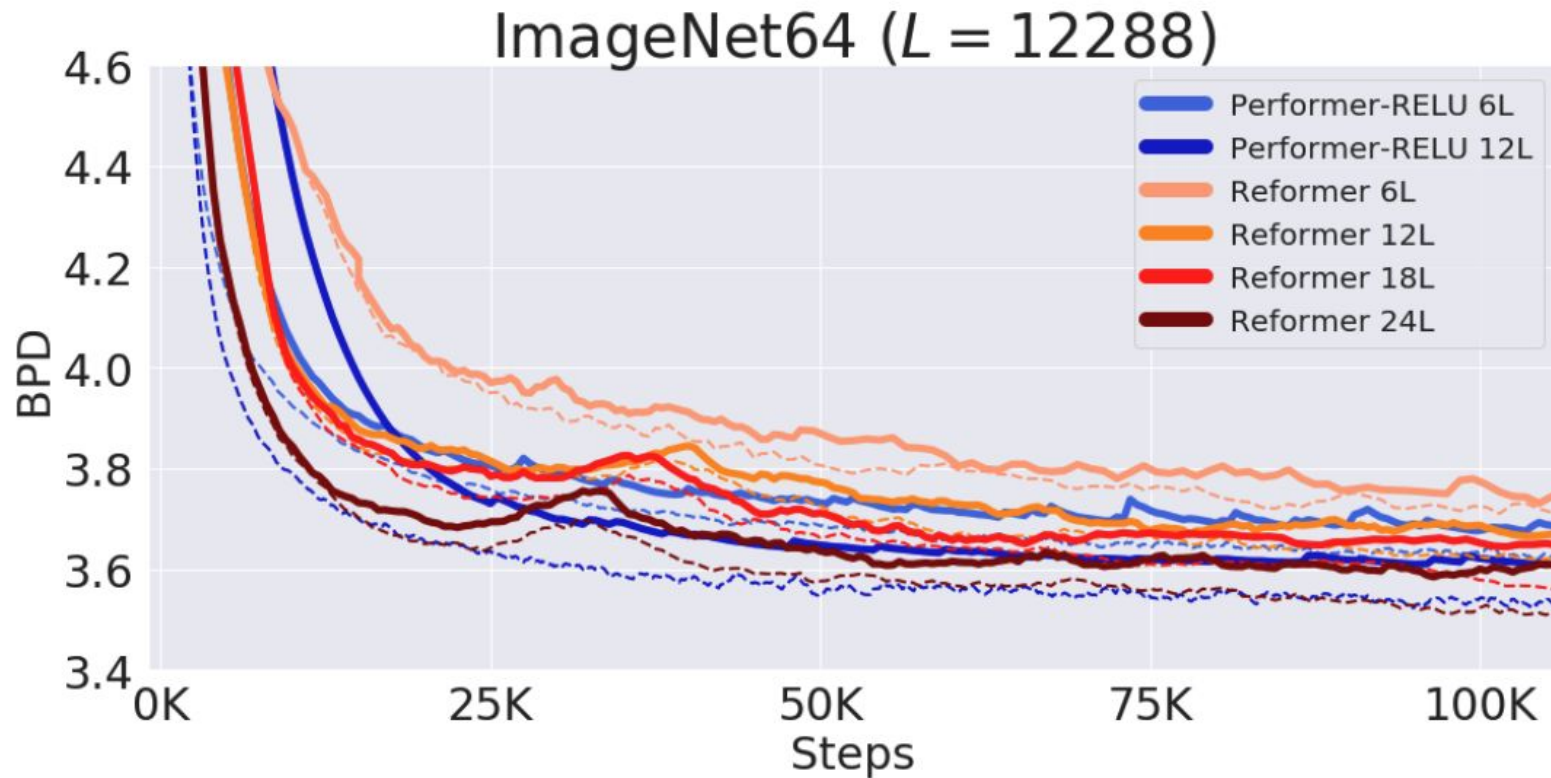


# Positive vs trigonometric random features in practice

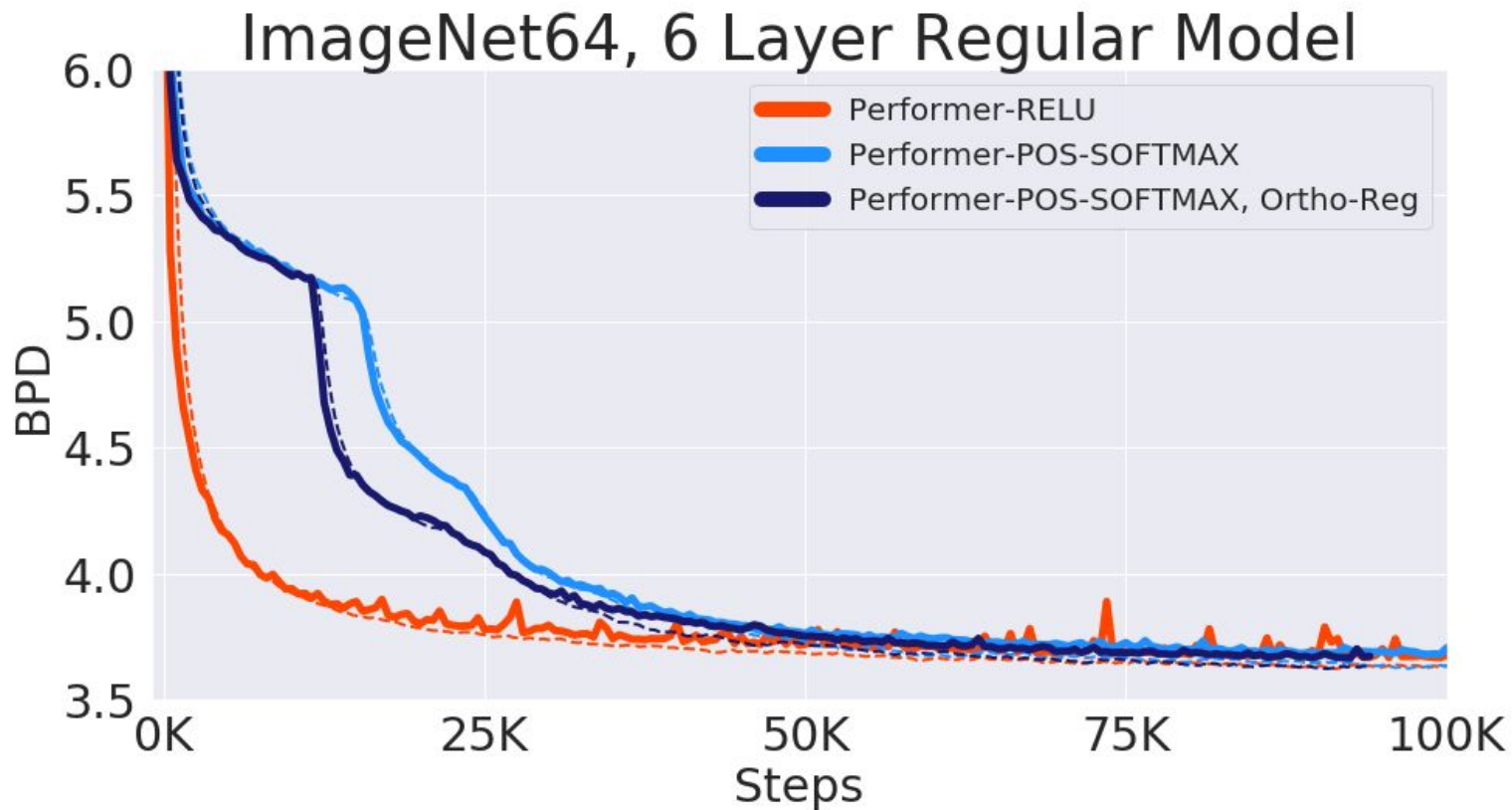
PG-19 (B), 6 Layer Regular Model ( $L = 1024$ )



# Performers on ImageNet64 - pixel predictions models



# Performers on ImageNet64 - Approx. Softmax vs ReLU

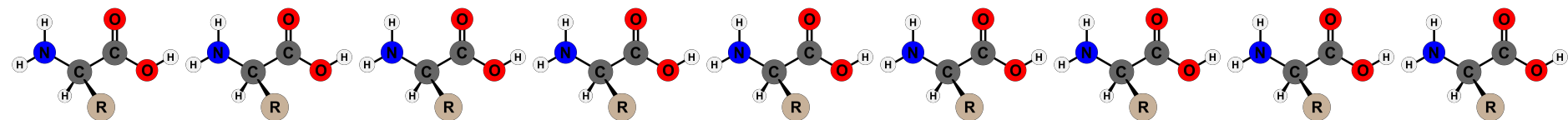




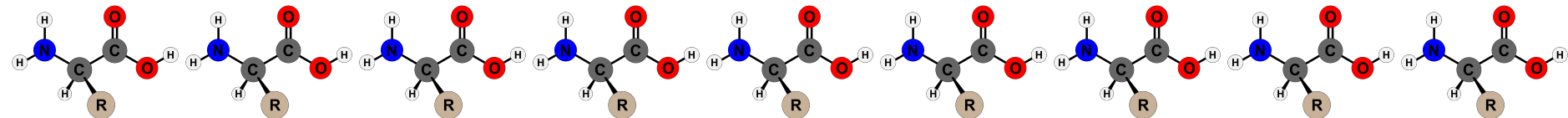
# Performers for Protein Design



# Performers for Protein Design

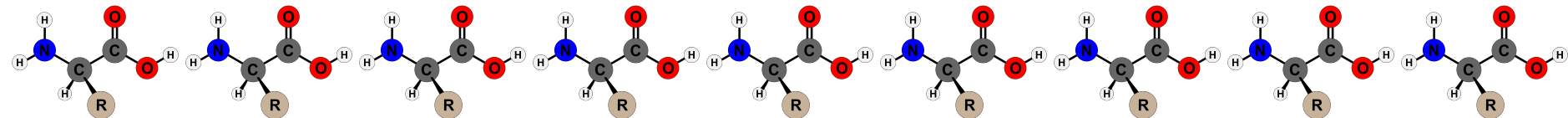




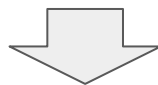
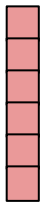


$\mathbf{X}_i$

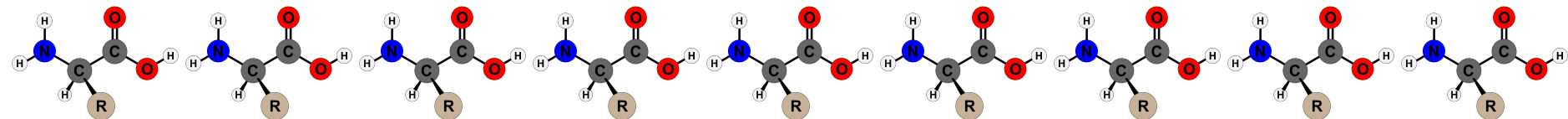
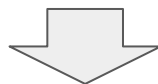
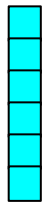
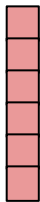




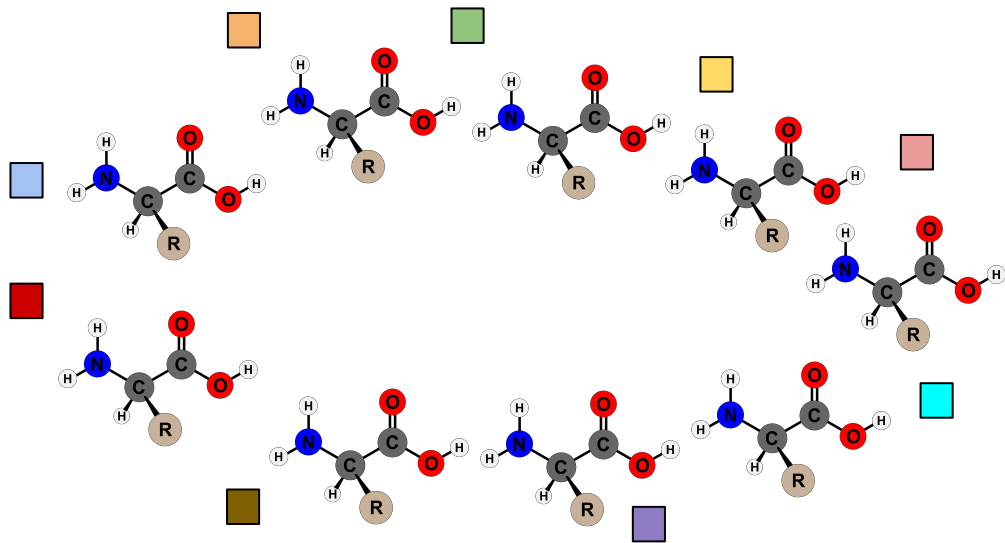
$\mathbf{X}_i$

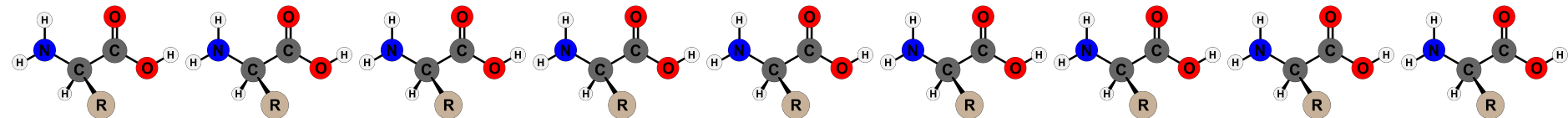
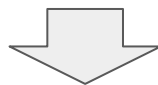
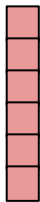


**tertiary structure**

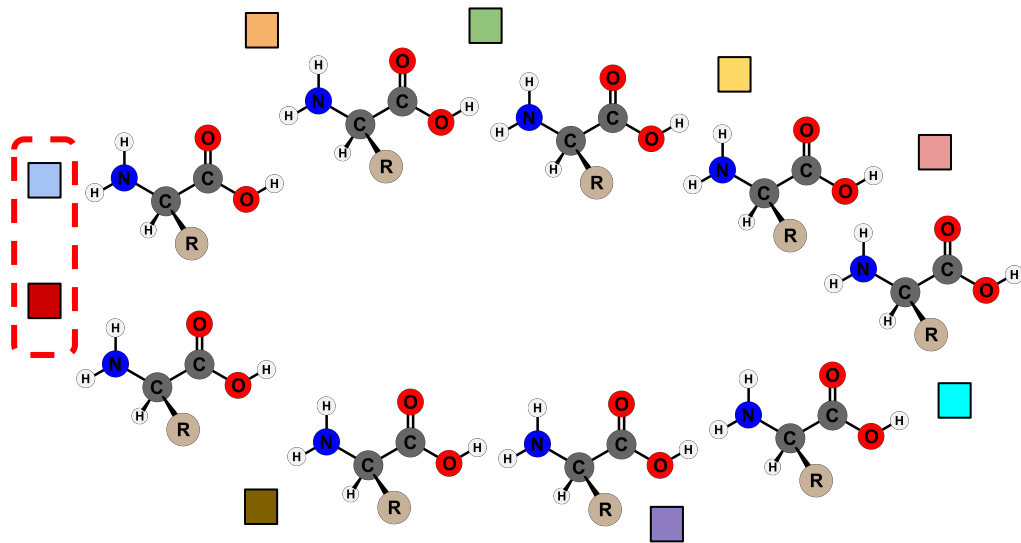

 $X_i$ 


tertiary structure




 $X_i$ 


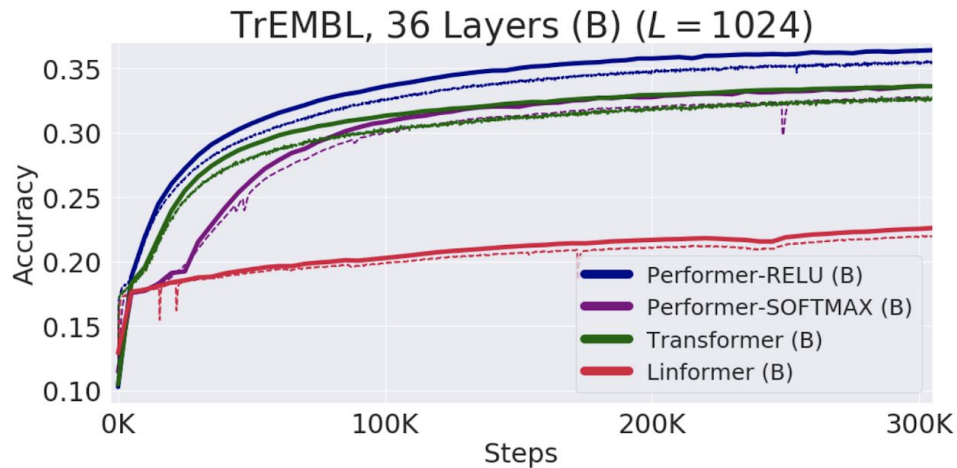
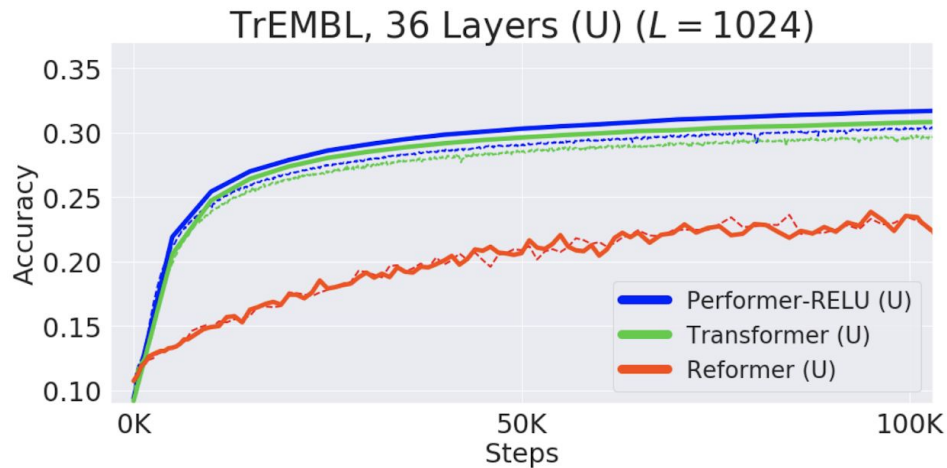
tertiary structure



# What we have already done...



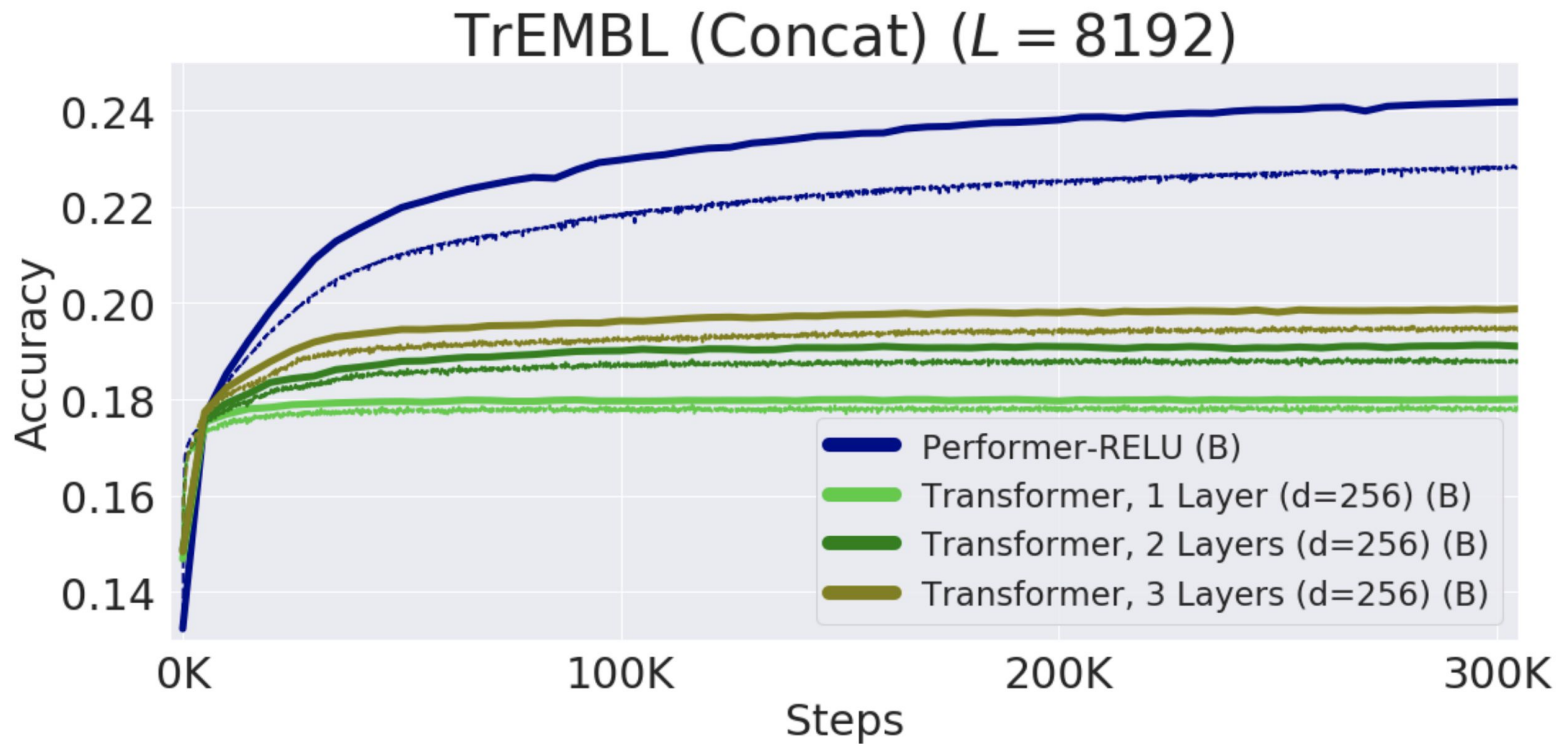
# Performers on moderate-size biological sequences



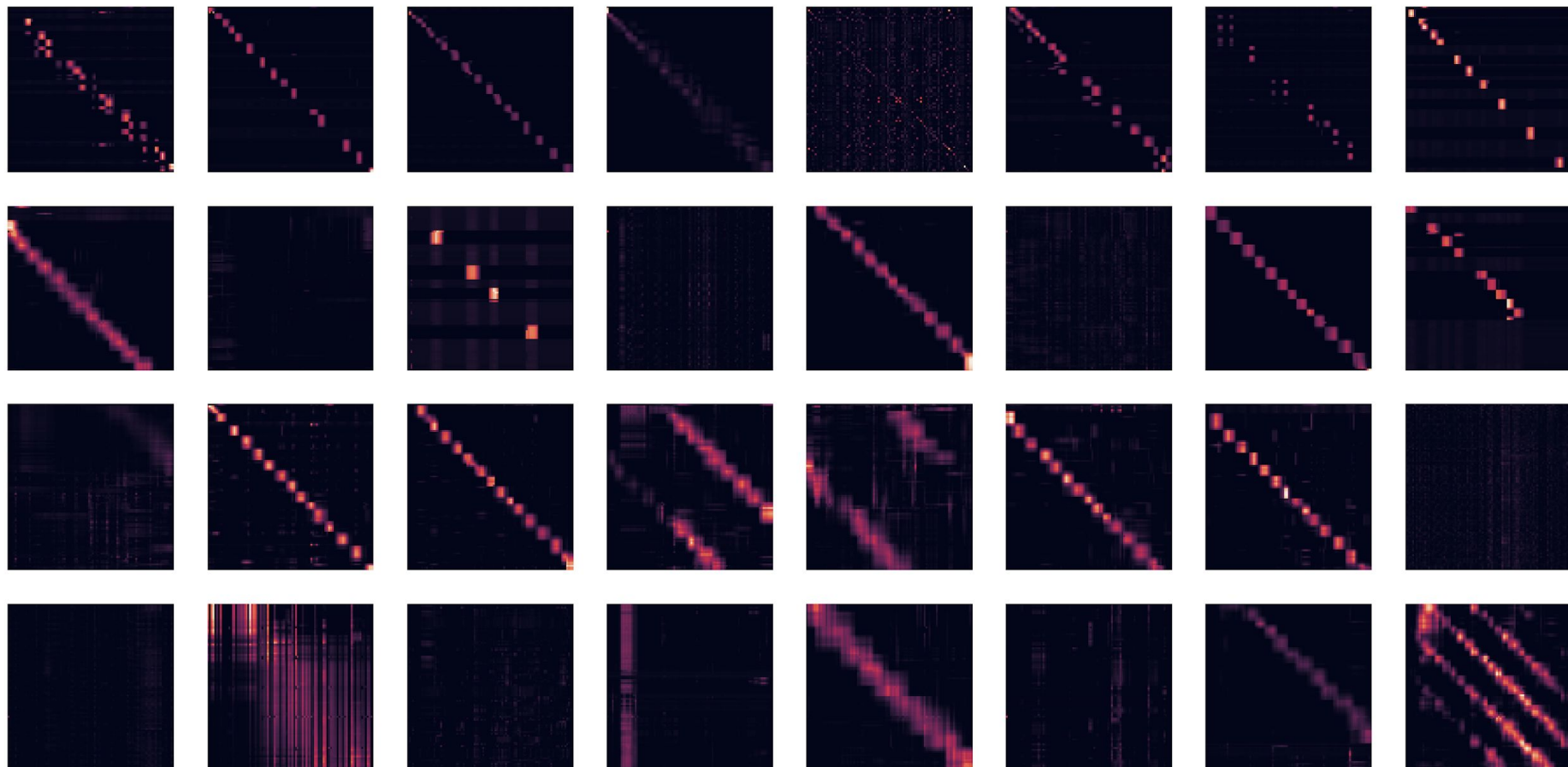
Train = Dashed, Validation = Solid, Unidirectional = (U), Bidirectional = (B). For TrEMBL, we used the exact same model parameters ( $n_{heads}, n_{layers}, d_{ff}, d$ ) = (8, 36, 1024, 512) from (Madani et al., 2020) for all runs. For fairness, all TrEMBL experiments used 16x16 TPU-v2's. Batch sizes were maximized for each separate run given the compute constraints. Hyperparameters & extended results including dataset statistics, out of distribution evaluations, and visualizations will appear soon in the extended version of the paper.

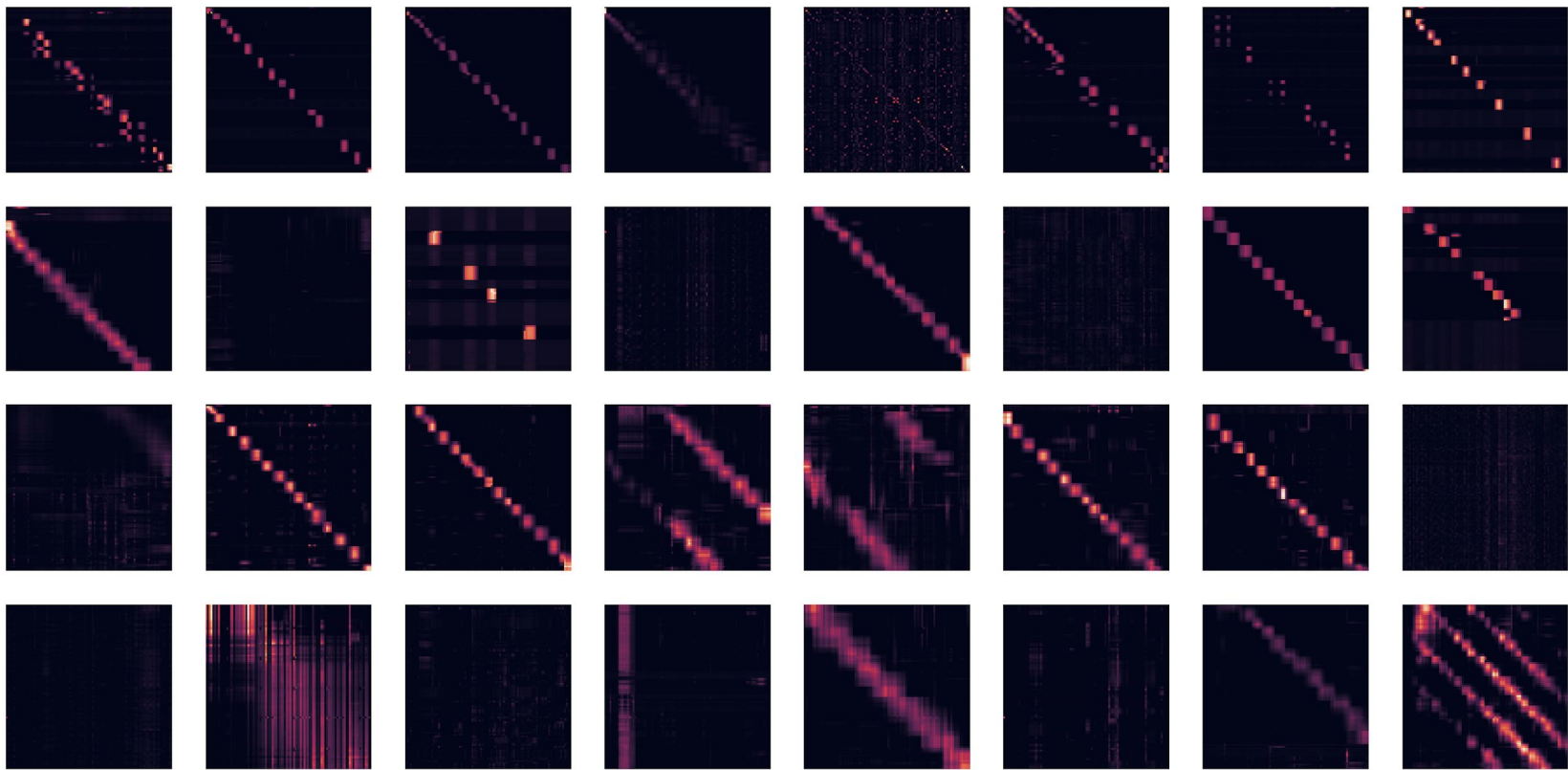


# Performers on long biological sequences: towards modeling complexes of proteins - proof of concept

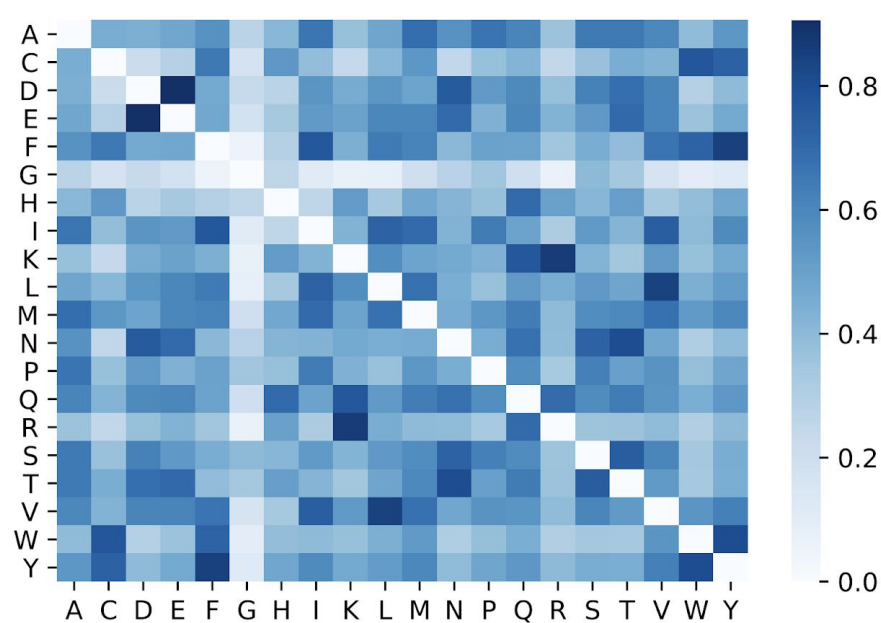
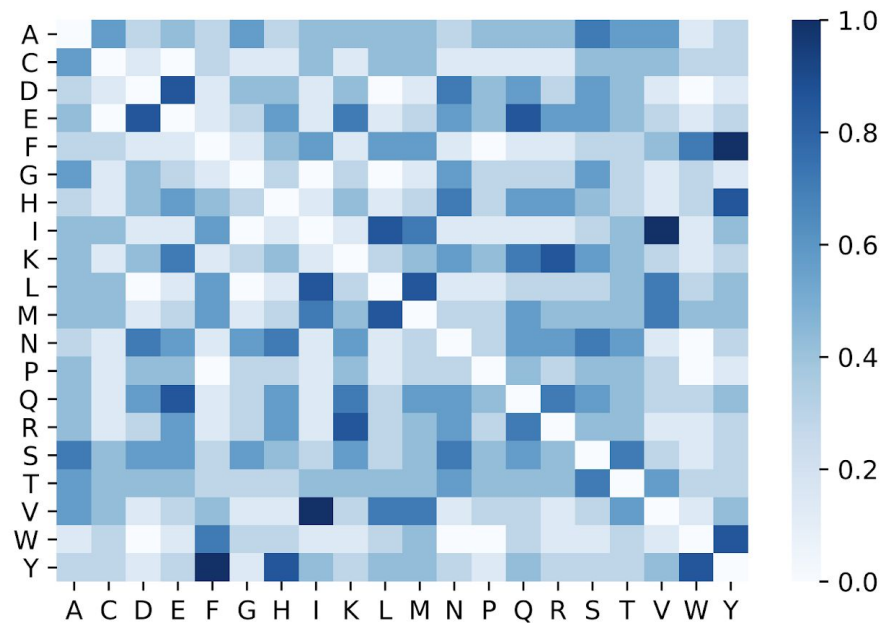


# What do Performers attend to ?



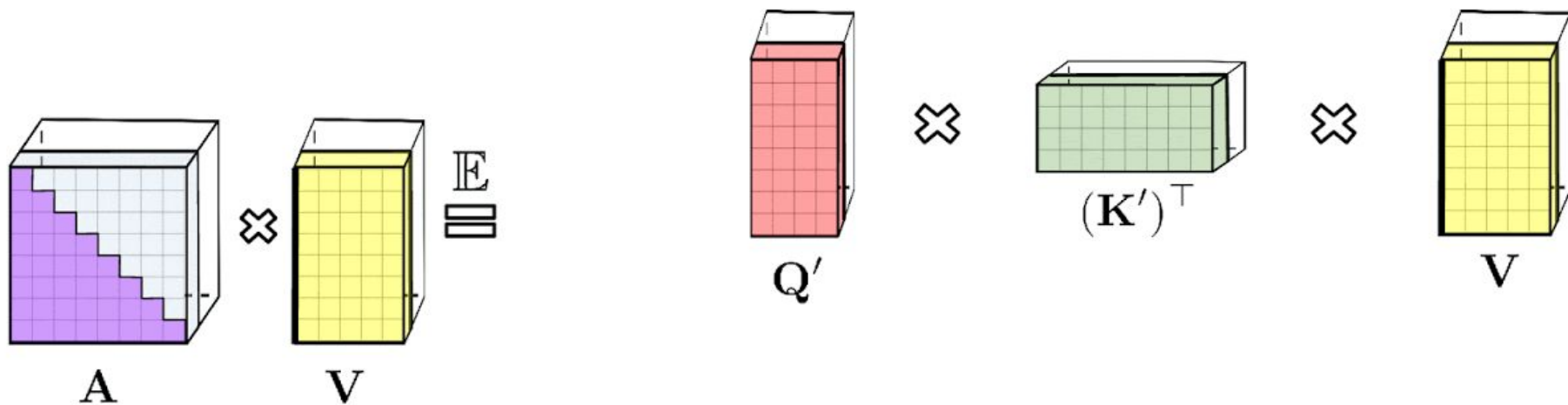


We show the attention matrices for the first 4 layers and all 8 heads (each row is a layer, each column is head index, each cell contains the attention matrix across the entire BPT1\_BOVIN protein sequence). Note that many heads show a diagonal pattern, where each node attends to its neighbors, and some heads show a vertical pattern, where each head attends to the same fixed positions.



Amino acid similarity matrix estimated from attention matrices aggregated across a small subset of sequences, as described in Vig et al. (Vig et al., 2020). The sub-figures correspond respectively to: (1) the normalized BLOSUM matrix, (2) the amino acid similarity estimated via a trained Performer model. Note that the Performer recognizes highly similar amino acid pairs such as (D, E) and (F, Y).

# Thank you for the Attention !



**Fig.** Linearized softmax causal attention as a prefix-sum computation engine.