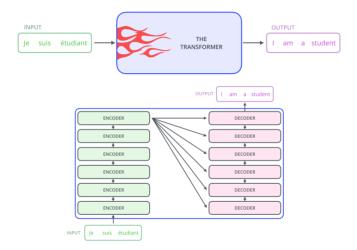
# Dissecting Attention: A Numerical Analyst's Perspective

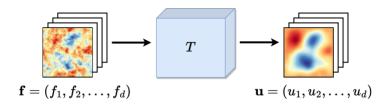
SCML Seminar KAUST, April 2024

### What is a Transformer?



The Transformer is a deep neural network architecture to solve the machine translation problem in Natural Language Processing. Source: Jay Alammar. *The Illustrated Transformers*.

#### Tensor2tensor



matrices to matrices of the same size.

A tensor 2 tensor DNN parametrizes the following discrete man for  $2I = (I^2(\Omega))^2$ 

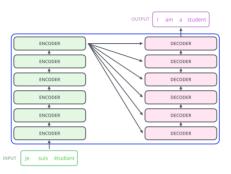
seg2seg or tensor2tensor in Neural Machine Translation maps various sized

• A tensor2tensor DNN parametrizes the following discrete map for  $\mathcal{H}=(L^2(\Omega))^d$ 

$$T_{\theta}: \mathbb{R}^{n \times d} \to \mathbb{R}^{n \times d}, \quad T: \mathcal{H} \to \mathcal{H}.$$

- Sentence in one language, embedded into high dimensional spaces, "translated" to another language's embedding after stacking multiple layers of the same module.
- Columns: numbers of latent/embedding dimension/channels (fixed in a given layer). Row: token embedding, patch embedding, or a DoF's embedding.
- The model can be trained on a lower "resolution" (n small) and evaluated at a higher "resolution" ( $n_{\text{eval}} \geq n$ ).

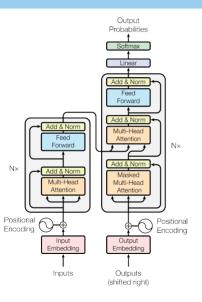
### What is a Transformer?



A Transformer consists of a sequence of encoder blocks with *identical* architectures, and decoder blocks with *identical* architectures. Source: Jay Alammar. *The Illustrated Transformers*.

- Like RNN in Neural Machine Translation, the Transformer block in each layer has the *same* architecture (and the number of parameters).
- Unlike CNN, after the initial embedding layer (from words to vector), the latent representations propagated in the hidden layers are of the *same* discretization size.

#### What is a Transformer?



#### Keywords:

- "Multi-head Attention".
- "Feedforward": fully connected MLP with shared weights at each position.
- "Add": skip-connection  $x \mapsto x + f(x)$ .
- "Norm": layer normalization (a learned diagonal columnscaling of the latent representation).
- "Positional Encoding": a hardcoded mapping to encode different positions in different latent dimensions.

### What is Multi-head Attention?

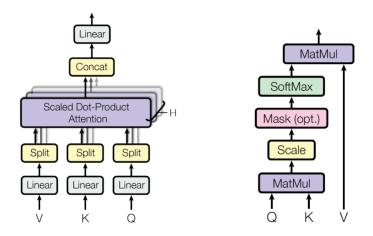
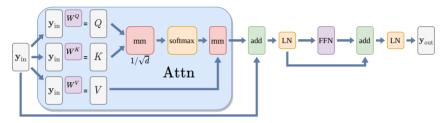


Image source: (Left) Multi-head Attention mapping. (Right) Scaled dot-product attention  $\operatorname{Softmax}(QK^{\top})V$ . Figure 2 in *Attention Is All You Need*.

# Single-head Self-Attention

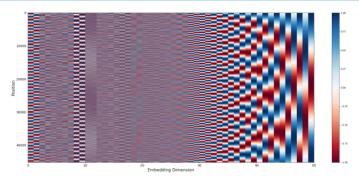


Self-attention mechanism in the classical Transformer in a single attention head.

- $oldsymbol{y}_{\mathsf{in}}, oldsymbol{y}_{\mathsf{out}} \in \mathbb{R}^{n imes d}$ , input/output embeddings; positional encodings added.
- Latent representations: query Q, key K, value V generated by 3 learnable matrices  $W^Q, W^K, W^V \in \mathbb{R}^{d \times d}$ :  $Q = yW^Q$ ,  $K = yW^K$ ,  $V = yW^V$ .
- The scaled dot-product attention:  $\operatorname{Attn}_s(\boldsymbol{y}) := \operatorname{Softmax}\left(d^{-1/2}QK^{\top}\right)V.$
- The full attention operator (add&norm, feedforward) is then

$$\operatorname{Attn}: \mathbb{R}^{n \times d} \to \mathbb{R}^{n \times d}, \quad \boldsymbol{z} = \boldsymbol{y} + \operatorname{Attn}_{s}(\boldsymbol{y}), \quad \boldsymbol{y} \mapsto \operatorname{Ln}\left(\boldsymbol{z} + g\left(\operatorname{Ln}(\boldsymbol{z})\right)\right)\right).$$

# Positional embedding



PE from Attention is All You Need.

- Positional embedding (PE):  $p \in \mathbb{R}^{n \times d}$  has the same dimension with the latent representation, and  $y \mapsto y + p$  for the y right after the input embedding.
- ullet M: maximum discretization size; c: channel index.

$$m{x}_{(i,c)} = \sin\left(rac{i}{M^{c/d}}
ight) ext{ if } c ext{ is even; } m{x}_{(i,2c+1)} = \cos\left(rac{i}{M^{(c-1)/d}}
ight) ext{ if } c ext{ is odd}$$

# Positional embeddings

Different interpretations of the latent representation (Query/Key/Value) which is an  $\mathbb{R}^{n\times d}$  matrix.

- Row: a high-dimensional embedding (vector representation) of a token.
- Column: certain discretization of "basis" or "frame" 2.

#### Open problems: Positional embeddings

- What role exactly does PE plays in attention?
- How PE shapes the topological structure of the latent representation space?
- How to design "nice" problem-oriented PE to achieve problem-specific attributes of traditional models?
- Is PE "≈" coordinates? ViT: DOSOVITSKIY et al. (2021); DeiT: TOUVRON et al. (2021); Swin: LIU et al. (2021).

<sup>&</sup>lt;sup>1</sup>L. Lu et al. (2021). "Learning nonlinear operators via DeepONet based on the universal approximation theorem of operators". In: *Nature machine intelligence*.

<sup>&</sup>lt;sup>2</sup>F. Bartolucci et al. (2023). "Representation Equivalent Neural Operators: a Framework for Alias-free Operator Learning". In: *Thirty-seventh Conference on Neural Information Processing Systems*.

# Single-layer Single-head Self-Attention

The full self-attention operator: n: discretization size, d: latent dimensions

Attn: 
$$\mathbb{R}^{n \times d} \to \mathbb{R}^{n \times d}$$
,  $\mathbf{y} \mapsto \mathbf{y}_{\text{out}}$ .

 $Attn(\cdot)$  consists the following operations:

- Adding PE (can be learnable):  $Pe: \mathbb{R}^{n \times d} \to \mathbb{R}^{n \times d}$ ,  $y \mapsto y + p(\theta)$
- Scaled dot-product attention:  $Q = yW^Q$ ,  $K = yW^K$ ,  $V = yW^V$

$$\left(\operatorname{Attn}_{s}(\boldsymbol{y})\right)_{i} = \sum_{j=1}^{n} \frac{\kappa(\boldsymbol{q}_{i}, \boldsymbol{k}_{j})}{\sum_{\ell=1}^{n} \kappa(\boldsymbol{q}_{i}, \boldsymbol{k}_{j'})} \boldsymbol{v}_{j}$$

where

$$\kappa(\cdot,\cdot):\mathbb{R}^d\times\mathbb{R}^d\to\mathbb{R}^+,\quad \kappa(\boldsymbol{q}_i,\boldsymbol{k}_j)=\exp(\boldsymbol{q}_i\cdot\boldsymbol{k}_j),$$
 notice  $\kappa(\boldsymbol{q}_i,\boldsymbol{k}_j)=\left(\boldsymbol{y}W^Q(\boldsymbol{y}W^K)^\top\right)_{ij}$ .

# Single-layer Single-head Self-Attention

The full self-attention operator:

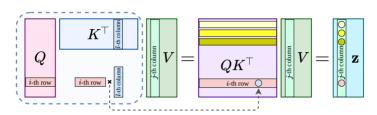
$$\operatorname{Attn}: \mathbb{R}^{n \times d} \to \mathbb{R}^{n \times d}, \quad \boldsymbol{z} = \boldsymbol{y} + \operatorname{Attn}_{\boldsymbol{s}}(\boldsymbol{y}), \quad \boldsymbol{y} \mapsto \operatorname{Ln}\left(\boldsymbol{z} + g\left(\operatorname{Ln}(\boldsymbol{z})\right)\right).$$

• Ln: layer normalization (LN), which has  $\gamma, oldsymbol{eta} \in \mathbb{R}^d$  learnable as follows

$$\operatorname{Ln}(oldsymbol{y}) := rac{oldsymbol{y} - oldsymbol{\mu}}{oldsymbol{\sigma}} \odot oldsymbol{\gamma} + oldsymbol{eta}, \ oldsymbol{\mu} := rac{1}{d} \sum_{j=1}^d oldsymbol{y}^j \in \mathbb{R}^n, \quad oldsymbol{\sigma}^2 := rac{1}{d} \sum_{j=1}^d (oldsymbol{y}^j - oldsymbol{\mu}) \odot (oldsymbol{y}^j - oldsymbol{\mu}) \odot (oldsymbol{y}^j - oldsymbol{\mu}) \in \mathbb{R}^n.$$

- Note all LN operations are done in an element-wise fashion. After LN is applied, each position in the discretization will roughly follow normal distribution.
- $g(\cdot)$ : a pointwise feedforward neural net with a "bottleneck" architecture in the base Transformer model.  $g(\cdot)$  has width d and is applied at each  $z_j$   $(j=1,\ldots,n)$ .

# Single-layer Single-head Self-Attention

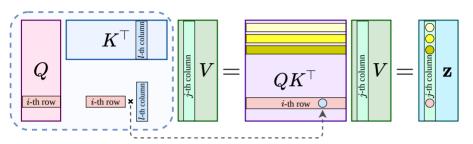


## (Still) Open problem

What exactly is the mechanics of the attention mechanism?

- Kernel interpretation, RKHS: TSAI et al. (2019), WRIGHT and GONZALEZ (2021), ZHANG et al. (2022).
- Fourier (change of basis): LI et al. (2021), NGUYEN, PHAM, et al. (2022).
- Low-rank or sparse: Y. XIONG et al. (2021), NGUYEN, SULIAFU, et al. (2021a), TAY et al. (2020), HAN et al. (2022).
- Random feature interpretation: CHOROMANSKI et al. (2021), PENG et al. (2021).
- Iterative "solver": YU et al. (2023)

## Scaled Dot-product Attention



$$(\mathbf{z}_i)_j = h \operatorname{Softmax}(QK^{\top})_{i\bullet} \mathbf{v}^j = hm_i^{-1} \exp\left(\mathbf{q}_i \cdot \mathbf{k}_1, \dots, \mathbf{q}_i \cdot \mathbf{k}_\ell, \dots, \mathbf{q}_i \cdot \mathbf{k}_n\right)^{\top} \cdot \mathbf{v}^j$$
$$= hm_i^{-1} \sum_{\ell=1}^n \exp(\mathbf{q}_i \cdot \mathbf{k}_\ell)(\mathbf{v}^j)_{\ell} \approx m^{-1}(x_i) \int_{\Omega} \kappa(x_i, \xi) v_j(\xi) d\xi,$$

The *i*-th row in the output computes approx. an integral transform with a non-symmetric normalized learnable low-rank "kernel" function  $\kappa(x,\xi)$ 

$$z(x) \approx \lambda v(x) + m^{-1}(x) \int_{\Omega} \kappa(x, \xi; \theta) v(\xi; \theta) d\xi, \quad \text{where } \mathbf{q}_i = q(x_i), \mathbf{k}_i = k(x_i), \mathbf{v}_i = v(x_i)$$

### Nonlocal Methods

Buades, Coll, and Morel<sup>3</sup> proposed that the denoising filter should depend on the signal! This is the earliest prototype of attention.

$$NLM[u](x) = \frac{1}{C(x)} \int_{\Omega} \exp\left(-\frac{1}{h^2} \int_{\Omega} G_{\alpha}(t) \left| u(x+t) - u(y+t) \right|^2 dt \right) u(y) dy,$$

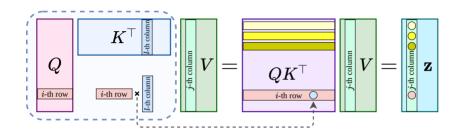
$$C(x) = \int_{\Omega} \exp\left(-\frac{1}{h^2} \int_{\Omega} G_{\alpha}(t) \left| u(x+t) - u(y+t) \right|^2 dt \right) dy.$$

Later this is generalized in papers by Osher and his postdoc using a variational perspective.<sup>4</sup>

<sup>&</sup>lt;sup>3</sup>A. Buades, B. Coll, and J.-M. Morel (2005). "A review of image denoising algorithms, with a new one". In: *Multiscale modeling & simulation*.

<sup>&</sup>lt;sup>4</sup>G. GILBOA and S. OSHER (2007). "Nonlocal linear image regularization and supervised segmentation". In: *Multiscale Modeling & Simulation*; G. GILBOA and S. OSHER (2009). "Nonlocal operators with applications to image processing". In: *Multiscale Modeling & Simulation*.

# Making Attention more efficient

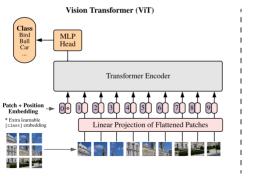


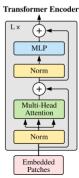
- The positive attention kernel  $\mathrm{Softmax}(QK^T)$  characterizes how each position's latent representation vector interact.
- This can be replaced by a simple Fourier kernel<sup>5</sup>.
- Computational cost of  $QK^T$  scales quadratically with respect to the number of positions n.

<sup>&</sup>lt;sup>5</sup> J. LEE-THORP, J. AINSLIE, I. ECKSTEIN, and S. ONTANON (2022). "FNet: Mixing Tokens with Fourier Transforms". In: Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies.

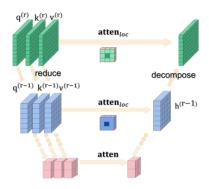
## Making Attention more efficient

- Low-rank: Nyströmformer (Y. XIONG et al. (2021)), Fast-Multipole method (NGUYEN, SULIAFU, et al. (2021b)).
- Sparsification: CHILD, GRAY, RADFORD, and SUTSKEVER (2019), locality-based feature maps in Reformer KITAEV, KAISER, and LEVSKAYA (2020),
- Linearization: WANG et al. (2020), RNN interpretation (KATHAROPOULOS, VYAS, PAPPAS, and FLEURET (2020)).
- Patchify: ViT (Dosovitskiy et al. (2021)).





## Hierarchically Nested Attention Neural Operator

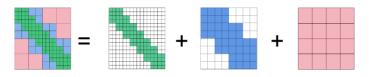


• Use local (windowed) attention to get the representation on each level:

$$\mathbf{atten}_{\mathsf{loc}}^{(m)}: \boldsymbol{v}_i^{(m)} = \sum_{j \in \mathcal{N}^{(m)}(i) \cup i} \mathcal{G}(\boldsymbol{q}_i^{(m)}, \boldsymbol{k}_j^{(m)}) \boldsymbol{v}_j^{(m)},$$

 $\mathcal{N}^{(m)}(i)$  contains m-th level neighbors of the i-th position in the discrete grid.

# Hierarchically Nested Attention Neural Operator



Aggregate the attention matrix (interaction) spanning multilevels:

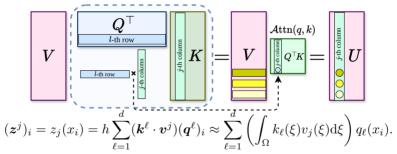
$$\left[\begin{array}{c} \vdots \\ \boldsymbol{h}_i^{(r)} \\ \vdots \end{array}\right] = \left(\sum_{m=1}^{r-1} (\mathbf{D}^{(r-1),\top} \cdots \mathbf{D}^{(m),\top} \mathbf{G}_{\mathsf{loc}}^{(m)} \mathbf{R}^{(m)} \cdots \mathbf{R}^{(r-1)}) + \mathbf{G}_{\mathsf{loc}}^{(r)} \right) \left[\begin{array}{c} \vdots \\ \boldsymbol{v}_i^{(r)} \\ \vdots \end{array}\right].$$

where the overall attention matrix on the finest level can be decomposed as

$$\mathbf{G}_h := \sum_{r=1}^{r-1} (\mathbf{D}^{(r-1), \top} \cdots \mathbf{D}^{(m), \top} \mathbf{G}_{\mathsf{loc}}^{(m)} \mathbf{R}^{(m)} \cdots \mathbf{R}^{(r-1)}) + \mathbf{G}_{\mathsf{loc}}^{(r)},$$

## Galerkin-type Attention

While it makes sense to ask the kernel to be positive (similarity between rows), it
does not to ask the interaction between bases (columns) to be positive.



- Resembles a (learnable) Petrov-Galerkin projection.
- How latents interact is similar to the Channel Attention<sup>6</sup>.

<sup>&</sup>lt;sup>6</sup>S. Woo, J. Park, J.-Y. Lee, and I. S. Kweon (2018). "CBAM: Convolutional block attention module". In: *Proceedings of the European conference on computer vision (ECCV)*.

## Galerkin-type Attention

Consider *i*-th entry in the *j*-th column  $z^j$  of z, which is the inner product of the *i*-th row of Q and the *j*-th column of  $K^\top V$ :

$$(\boldsymbol{z}^j)_i = h \, \boldsymbol{q}_i^{\top} \cdot (K^{\top} V)_{\bullet i}$$

$$\boldsymbol{z}^{j} = h \left( \begin{array}{cccc} | & | & | & | \\ \boldsymbol{q}_{1} & \boldsymbol{q}_{2} & \cdots & \boldsymbol{q}_{n} \\ | & | & | & | \end{array} \right)^{\top} (K^{\top}V)_{\bullet j} = h \left( (K^{\top}V)_{\bullet j}^{\top} \begin{pmatrix} \boldsymbol{q}^{1} & \boldsymbol{\dots} \\ \boldsymbol{\dots} & \vdots & \boldsymbol{\dots} \\ \boldsymbol{q}^{d} & \boldsymbol{\dots} \end{pmatrix} \right)^{\top}$$

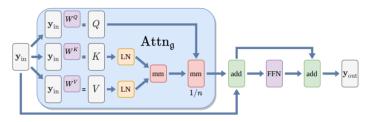
then

$$\boldsymbol{z}^j = h \sum^d \boldsymbol{q}^l (K^\top V)_{\ell j}, \quad \text{where} \ \ (K^\top V)_{\boldsymbol{\cdot} j} = \left(\boldsymbol{k}^1 \cdot \boldsymbol{v}^j, \boldsymbol{k}^2 \cdot \boldsymbol{v}^j, \cdots, \boldsymbol{k}^d \cdot \boldsymbol{v}^j\right)^\top.$$

Rewriting  $\langle v_i, k_l \rangle := (K^\top V)_{\ell i}$ 

$$z_j(x) := \sum_{l=0}^d \langle v_j, k_l \rangle \, q_l(x), \ \text{ for } j=1,\cdots,d, \ \text{ and } x \in \{x_i\}_{i=1}^n,$$

## Galerkin-projection inspired Attention



• Inspired by the Fourier transform to remove the softmax normalization to improve the computational efficiency: orthogonal  $\{q_j(\cdot)\}_{j=1}^d$ 

$$\min_{a_i} \left\| f - \sum_{i=1}^d a_i q_i(\cdot) \right\|_{\ell^2(\Omega)}^2, \quad \text{and} \quad z(x) := \sum_{\ell=1}^d \frac{(f,q_\ell)}{(q_\ell,q_\ell)} q_\ell(x),$$

• Inspired by the Gram matrix inverse normalization in the proof of the Ceá type lemma, and layer normalization modifications<sup>7</sup>.

<sup>&</sup>lt;sup>7</sup>R. XIONG et al. (2020). "On layer normalization in the transformer architecture". In: *International Conference on Machine Learning*. PMLR.

# A preliminary result on the Galerkin-type Attention

## Theorem (Approximation capacity of a single layer of Galerkin attention <sup>8</sup>)

 $\mathbb{Q}_h \subset \mathcal{Q}$  and  $\mathbb{V}_h \subset \mathcal{V}$  are the current approximation space, suppose there exists a continuous key-to-value map that is bounded below on the discrete approximation space, i.e., the functional norm of  $v \mapsto b(q,v)$  is bounded below for any q, then for  $g_\theta$  consists a Galerkin attention composed with a channel reduction map

$$\min_{\theta} \|f - g_{\theta}(\mathbf{y})\| \leq \underbrace{c^{-1}}_{\|b(q,\cdot)\|_{\mathbb{V}_h'} \geq c} \underbrace{\min_{q \in \mathbb{Q}_h} \max_{v \in \mathbb{V}_h} \frac{|b(\Pi f - q, v)|}{\|v\|}}_{\text{(Error of the Petrov-Galerkin projection)}} + \underbrace{\|f - \Pi f\|}_{\text{(Consistency)}}.$$

- Intepretation: for a "query" (a function in a Hilbert space), to deliver the best approximator in "value" (trial space), the "key" space (test space) has to be big enough so that for every value there is a key to unlock it.
- discrete Ladyzhenskaya-Babuška-Brezzi inf-sup condition: why Transformers have capacity to generalize so well with respect to the length of the sequence.

<sup>&</sup>lt;sup>8</sup> C. (2021). "Choose a Transformer: Fourier or Galerkin". In: Advances in Neural Information Processing Systems (NeurIPS)

# Sketch of the proof

• For the continuous bilinear form  $b(\cdot,\cdot): \mathcal{Q} \times \mathcal{V} \to \mathbb{R}$ ,  $b(q,\cdot): v \mapsto b(q,v)$  is bounded below on  $\mathbb{V}_h \subset \mathcal{V}$  for any  $q \in \mathbb{Q}_h$ :

$$c||q||_{\mathcal{H}} \le \sup_{v \in \mathbb{V}_h} \frac{|b(q,v)|}{||v||_{\mathcal{V}}}.$$

The Riesz map by  $b(\cdot,\cdot)$  from the value space to the key space is injective (or key-to-value is surjective). We can verify c is independent of the sequence length for the scaled dot-product attention (without softmax).

• Consider an incoming function f's projection in  $\mathbb{Q}_h$ :  $f_h$  (query). By the inf-sup condition above,

$$||f_h - g_{\theta}(\mathbf{y})||_{\mathcal{H}} \le c^{-1} \sup_{\mathbf{v} \in \mathbb{V}_h} \frac{|b(f_h - g_{\theta}(\mathbf{y}), \mathbf{v})|}{||\mathbf{v}||_{\mathcal{V}}}.$$

The rest is just to show

$$\min_{\theta} \max_{v \in \mathbb{V}_h} \frac{|b(f_h - g_{\theta}(\mathbf{y}), v)|}{\|v\|_{\mathcal{V}}} \le \min_{q \in \mathbb{Q}_h} \max_{v \in \mathbb{V}_h} \frac{|b(f_h - q, v)|}{\|v\|_{\mathcal{V}}}.$$

# Sketch of the proof

Solving the min-max problem

$$\min_{q \in \mathbb{Q}_h} \max_{v \in \mathbb{V}_h} \frac{\left| \langle \Phi f_h, v \rangle_h - b(q, v) \right|}{\|v\|_{\mathcal{H}}}$$

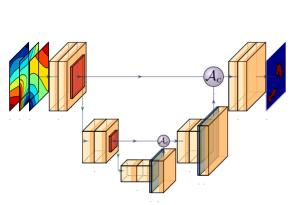
is equivalent to solving the following saddle point problem:

$$\begin{cases} \langle w, v \rangle_h + b(p, v) = \langle \Phi f_h, v \rangle_h, & \forall v \in \mathbb{V}_h, \\ b(q, w) = 0, & \forall q \in \mathbb{Q}_h. \end{cases}$$

• Lastly, the scaled dot-product attention has capacity to represent this best approximator's vector representation  $\boldsymbol{p}$ : let  $\Lambda = \operatorname{blkdiag}\left\{(BM^{-1}B^{\top})^{-1},0\right\}$ , B and M be the Gram (mass) matrices associated with  $b(\cdot,\cdot)$  and  $\langle\cdot,\cdot\rangle_h$ 

$$\begin{split} \widetilde{Q} &:= \mathbf{y}\widetilde{W}^Q \leftarrow \mathbf{y}W^QU, \\ \widetilde{K} &:= \mathbf{y}\widetilde{W}^K \leftarrow \mathbf{y}W^QU\Lambda, \\ \widetilde{V} &:= \mathbf{y}\widetilde{W}^V \leftarrow \mathbf{y}W^VM^{-1}, \\ \text{and } \boldsymbol{p} &= \widetilde{Q}(\widetilde{K}^T\widetilde{V})\boldsymbol{\zeta}. \end{split}$$

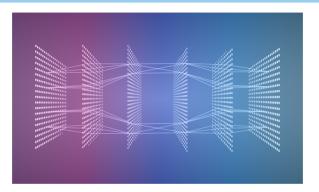
## Making linear Attention more efficient



- U-Net meta-architecture 9.
- Input: the concatenation of discretizations of  $\phi$  and  $\nabla \phi$ .
- Output: the approximation to the index map  $\mathcal{I}^D$ .
- □: 3 × 3 convolution + ReLU;
- ■: normalization;
- ■: interpolation;
- cross attention from the coarse grid to the fine grid;
- input and output discretized functions.

<sup>&</sup>lt;sup>9</sup>O. Ronneberger, P. Fischer, and T. Brox (2015). "U-net: Convolutional networks for biomedical image segmentation". In: *International Conference on Medical image computing and computer-assisted intervention*. Springer

# Representational capcity



- The embedding layer in Transformer constructs an ultra-high-dimensional vector representation of each token in a sentence or each patch in an image.
- In the encoder layer, this (latent) representation interacts with itself nonlinearly to get a "better" representation.
- This interaction can be position-wise (row-wise), or channel-wise (column-wise).

Image source: Transformers: What They Are and Why They Matter, Mehreen Saeed.

# Representational capcity

#### Open problem about representational capacity

How to prove the universal representation (approximation) theorem for Transformer when the number of layers increase?

- What is representational power of (stacked) attention layer exactly<sup>10</sup>?
- Random feature model<sup>11</sup>: each channel (column) of the latent representation is similar to an RF-RR model

$$m{f} \mapsto \Phi(m{f};m{ heta}) = rac{1}{d} \sum_{j=1}^d lpha_j(m{ heta}) g(m{f};m{ heta})$$

where

$$oldsymbol{ heta} = rgmin rac{1}{N} \sum_{i=1}^N \|oldsymbol{u}_i - \Phi(oldsymbol{f}_i; heta)\|_V^2 + ext{regularizations}$$

<sup>&</sup>lt;sup>10</sup>C. YuN et al. (2020). "Are Transformers universal approximators of sequence-to-sequence functions?" In: International Conference on Learning Representations.

<sup>&</sup>lt;sup>11</sup>A. Rahimi and B. Recht (2008). "Weighted sums of random kitchen sinks: Replacing minimization with randomization in learning". In: *Advances in neural information processing systems*.

# Representational capcity and positional embedding

- PE plays an important role in shaping the representational capacity.
- Learnable PE<sup>12</sup>, rotational-invariant PE<sup>13</sup>, etc.

### Theorem (Universal representater theorem (informal simplified version)<sup>14</sup>)

Given fixed n and d, the function class of Transformers  $\{u(\boldsymbol{y}): u(\boldsymbol{y}) = g(\boldsymbol{y} + \boldsymbol{x}), \text{ where } g := g_{\ell} \circ \cdots \circ g_1\}$  with the absolute fixed  $PE \ \boldsymbol{x}$  is a universal approximator for continuous functions that map a compact domain in  $\mathbb{R}^{n \times d}$  to  $\mathbb{R}^{n \times d}$ .

#### Open problem: representational capacity

Can the theoretical results on the approximation capacity of Transformer with different PEs be reflected in specially designed experiments?

<sup>&</sup>lt;sup>12</sup>J. Gehring et al. (2017). "Convolutional sequence to sequence learning". In: *International conference on machine learning*. PMLR.

<sup>&</sup>lt;sup>13</sup> J. Su et al. (2021). "Roformer: Enhanced transformer with rotary position embedding". In: arXiv preprint arXiv:2104.09864.

<sup>&</sup>lt;sup>14</sup>S. Luo et al. (2022). "Your Transformer May Not be as Powerful as You Expect". In: *Advances in Neural Information Processing Systems*