# Benefits of Rank in Attention Layers

Noah Amsel, Gilad Yehudai, Joan Bruna





## Motivation

- How do we get the most out of transformers?
- How do we set the hyperparameters?

What Can Transformers Learn In-Context? A Case Study of Simple Function Classes

Shivam Garg\* Stanford University shivamg@cs.stanford.edu

Dimitris Tsipras\* Stanford University tsipras@stanford.edu

Percy Liang Stanford University pliang@cs.stanford.edu

**Gregory Valiant** Stanford University valiant@stanford.edu





# Ancient History

- "Neural Machine Translation by Jointly Learning to Align and Translate" [Bahdanau, Cho & Bengio, 2016]
  - Given query  $\mathbf{q}$ , keys  $\{\mathbf{k}_i\}$ , and values  $\{\mathbf{v}_i\}$  in  $\mathbb{R}^r$
  - $\mathbf{s} = \mathbf{K}\mathbf{q}$  $s_i = \mathbf{k}_i \cdot \mathbf{q}$  $\alpha_j = \frac{\exp(s_j)}{\sum_k \exp(s_k)} \qquad \alpha = \operatorname{sm}(\mathbf{s})$ Attention coefficients:
  - Attention score:
  - Convex combo of values:
- $\sum \alpha_j \mathbf{v}_j$  $\mathbf{V}\alpha$

 $= V \operatorname{sm}(\mathbf{K}^{\mathsf{T}}\mathbf{q})$ 



## Singe-head Attention

# Given y and $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^d$ , let $\mathbf{q} = \mathbf{W}^Q \mathbf{y} \qquad \mathbf{K} = \mathbf{W}^K \mathbf{X} \qquad \mathbf{V} = \mathbf{W}^V \mathbf{X}$

where  $\mathbf{W}^Q$ ,  $\mathbf{W}^K$ ,  $\mathbf{W}^V$ ,  $\mathbf{W}^O \in \mathbb{R}^{r \times d}$ .

## $\left(\mathbf{W}^{O}\right)^{\mathsf{T}}\mathbf{W}^{V}\mathbf{X}\,\operatorname{sm}\left(\left(\mathbf{W}^{K}\mathbf{X}\right)^{\mathsf{T}}\mathbf{W}^{Q}\mathbf{y}\right)$

## Multi-head Attention

## Given **y** and $\mathbf{x}_1, \dots \mathbf{x}_N \in \mathbb{R}^d$

$$\sum_{h=1}^{H} \left( \mathbf{W}_{h}^{O} \right)^{\mathsf{T}} \mathbf{W}_{h}^{V} \mathbf{X}$$

## where $\mathbf{W}_{h}^{Q}$ , $\mathbf{W}_{h}^{K}$ , $\mathbf{W}_{h}^{V}$ , $\mathbf{W}_{h}^{O} \in \mathbb{R}^{r \times d}$ .

• Num parameters: 4*rHd* 

## $\operatorname{sm}\left(\left(\mathbf{W}_{h}^{K} \mathbf{X}\right)^{\mathsf{T}} \mathbf{W}_{h}^{Q} \mathbf{y}\right)$

## Multi-head Attention

# Given **y** and $\mathbf{x}_1, \dots \mathbf{x}_N \in \mathbb{R}^d$ $\sum_{h=1}^H \mathbf{W}_h \mathbf{X}_h$

where  $\mathbf{M}_h, \mathbf{W}_h \in \mathbb{R}^{d \times d}$  are rank-r.

• Num parameters: 4*rHd* 

 $\sum_{h=1}^{n} \mathbf{W}_{h} \mathbf{X} \operatorname{sm} \left( \mathbf{X}^{\mathsf{T}} \mathbf{M}_{h} \mathbf{y} \right)$ 

hido

## YearModel2017Attention is all you need [59]

num layers		's M	MLP depth			value rank	
den dim	ſ	MLP widt	:h	attn ran	k	num head	
d	L	w	D	r	$r_2$	H	
512	6	4d	2	64	r	d/r	
		Num params $= 4rHd = 4d^2$					
		Same as $H = 1$ , $r = d$					

lS 

48

Mean Squared Error 0.25 0.20 0.15 0.10 0.05 0.00



## r = d / H

# More Garg et al.





## Theory papers assume full-rank

- Practitioners use low r and H = d/r
- Theoreticians assume r = d and  $H \gg 1$ 
  - Is this a good proxy for real-world transformers?

## Related Work

# **Depth Separation for MLPs**

- Q: 2-layer NNs are universal approximators. So why are deep NNs better?
- A: Deep NNs are more parameter-efficient
  - Pay for some depth, save a <u>lot</u> of width
- Thm: There exists a function which
  - depth 3 MLP can represent with width  $d^{5}$
  - depth 2 MLP needs width  $e^{cd}$  to approximate
- [Eldan & Shamir '16] [Telgarsky '16] [Daniely '17] [Chatziafratis+ '19] ...

Construct a function showing that you can... MLP:

...pay a bit for depth, save a <u>lot</u> of width

**Multi-head attention layer:** ...pay a bit for rank, save a *lot* of heads

## **Separations in Expressive Power**

## Another rank separation result

- "Representational Strengths and Limitations of Transformers" [Sanford, Hsu, Telgarsky, 2023]
- Thm: A single layer of multi-head self-attention cannot compute Match3 unless  $rHp > \Omega(N)$ • Proved using communication complexity (p is precision)
- **Q:** Is low r a limitation, or just low rH?
- We extend hardness guarantee to...
  - Prohibitively large H
  - $p = \infty$
  - $\epsilon$ -approximation



## **Our Rank Separation**

target points  $\mathbf{X}_i$ source point y output  $f(\mathbf{x}_1, \dots, \mathbf{x}_N; \mathbf{y})$ 



## **Target Function: Nearest Neighbor**

- For  $\mathbf{y}, \mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{S}^{d-1}$ 
  - $f(\mathbf{x}_1, ..., \mathbf{x}_N; \mathbf{y}) := \text{argmin } \|\mathbf{x} \mathbf{y}\|_2$  $\mathbf{x} \in \{\mathbf{x}_{1},...,\mathbf{x}_{N}\}$ 
    - $= \operatorname{argmax} \mathbf{x}^{\mathsf{T}} \mathbf{y}$  $x \in \{x_1, ..., x_N\}$

    - $= \mathbf{X} hm(\mathbf{X}^{\mathsf{T}}\mathbf{y})$

 $= \mathbf{I}_d \mathbf{X} \operatorname{sm} \left( \mathbf{X}^{\mathsf{T}} \left( 10^{10} \cdot \mathbf{I}_d \right) \mathbf{y} \right)$ full ran

## Upper Bound

## For $\mathbf{y}, \mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{S}^{d-1}$ , the nearest neighbor function can be exactly represented by a single full-rank hardmax attention head.

"
$$r = d, F$$

H = 1 suffices"

## Lower Bound (Informal)

## If $H \lesssim (d/r)^{1/\epsilon}$ , then

 $\mathbb{E}_{\mathbf{x}_1 \perp \mathbf{x}_2, \mathbf{y} \sim \mathcal{U}(\mathbb{S}^{d-1})} \quad f(\mathbf{x}_1, \mathbf{x}_2; \mathbf{x}_2)$ 

$$\mathbf{y}) - \sum_{h=1}^{H} \operatorname{attn}_{h}(\mathbf{x}_{1}, \mathbf{x}_{2}; \mathbf{y}) \left\| \begin{array}{c} 2 \\ 2 \end{array} \geq \epsilon \right\|_{2}$$

## Generalizing Attention

Standard:



## where $\mathbf{W}_{h}^{Q}$ , $\mathbf{W}_{h}^{K}$ , $\mathbf{W}_{h}^{V}$ , $\mathbf{W}_{h}^{O} \in \mathbb{R}^{r \times d}$ .

Generalized:

 $\sum \mathbf{V}_{h} \mathbf{X} \boldsymbol{\phi}_{h} \left( \mathbf{K}_{h}^{\mathsf{T}} \mathbf{X}, \mathbf{y} \right)$ h=1

where  $\mathbf{V}_h \in \mathbb{R}^{d \times d}$ ,  $\mathbf{K}_h \in \mathbb{R}^{d \times r}$ ,  $\phi_h : \mathbb{R}^{r \times N} \times \mathbb{R}^d \to \Delta^N$ .

 $\sum_{h=1}^{N} \left( \mathbf{W}_{h}^{O} \right)^{\mathsf{T}} \mathbf{W}_{h}^{V} \mathbf{X} \operatorname{sm} \left( \left( \mathbf{W}_{h}^{K} \mathbf{X} \right)^{\mathsf{T}} \mathbf{W}_{h}^{Q} \mathbf{y} \right)$ 

## **Proof Sketch**

## Reduction to Scalar Function on $\mathbb{S}^{d-1} \times \mathbb{S}^{d-1}$



where  $g_h(\mathbf{x}, \mathbf{y}) = \tilde{\phi}_h(\mathbf{K}_h^{\mathsf{T}}\mathbf{x}, \mathbf{y})$ 

(follows by projecting onto  $\mathbf{x}_1 - \mathbf{x}_2$ )



## **Spherical Harmonics** An orthonormal basis for $L^2(\mathbb{S}^{d-1} \to \mathbb{R})$



# basis elements of degree  $\ell$  in on  $\mathbb{S}^{d-1}$  is  $N(d, \ell)$ 



## Spherical harmonic expansion of rank-1 functions

Spherical harmonic expansion:  $f = \sum$ 

• where 
$$\langle f, g \rangle = \mathbb{E}_{\|\mathbf{x}\|=1} \left[ f(\mathbf{x})g(\mathbf{x}) \right]$$

- Rank-1 function:  $f(\mathbf{x}) = \psi(\mathbf{x}^{\mathsf{T}}\mathbf{a})$  for  $\psi : \mathbb{R} \to \mathbb{R}$
- Theorem (Hecke-Funk):

 $\langle f, Y^i_{\ell} \rangle =$ 

• where  $P_{\ell}$  is the  $\ell^{th}$  Gegenbauer polynomial

$$\sum_{\ell=0}^{\infty} \sum_{i=1}^{N(d,\ell)} \langle f, Y_{\ell}^{i} \rangle \cdot Y_{\ell}^{i}$$

$$Y^i_{\ell}(\mathbf{a})\cdot \langle \psi, P_{\ell} \rangle$$



## Representing sgn $(\mathbf{x}^{\top}\mathbf{y})$ with spherical harmonics

- $\left\{Y_{\ell}^{i}\otimes Y_{\ell'}^{i'}\right\}_{\ell,\ell',i,i'}$  is orthonormal basis for  $L^{2}\left(\left(\mathbb{S}^{d-1}\times\mathbb{S}^{d-1}\right)\to\mathbb{R}\right)$
- By Hecke-Funk

$$\langle \operatorname{sgn}(\mathbf{x}^{\mathsf{T}}\mathbf{y}), Y_{\ell}^{i}(\mathbf{x})Y_{\ell'}^{i'}(\mathbf{y}) \rangle = \langle \operatorname{sgn} \rangle$$
  
 $\implies \operatorname{sgn}(\mathbf{x}^{\mathsf{T}}\mathbf{y}) = \sum_{\ell=0}^{\infty} \sum_{i=1}^{N(d,\ell)} \eta_{\ell}Y_{\ell}^{i}(\mathbf{x})Y_{\ell'}^{i}$ 

where  $\eta_{\ell} \sim \ell^{-1}$ 

 $\operatorname{gn}, P_{\ell} \rangle \cdot \langle Y_{\ell}^{i}, Y_{\ell'}^{i'} \rangle = \eta_{\ell} \cdot \begin{cases} 1 & \ell = \ell', i = i' \\ 0 & 0.W. \end{cases}$ 









## Expansion of the head functions

- $\phi_h(\mathbf{K}_h^{\mathsf{T}}\mathbf{x},\mathbf{y})$  only cares about an *r*-dimensional projection of  $\mathbf{x}$ 
  - $\Rightarrow$  ortho to many spherical harmonic

• Lemma:  $\phi_h(\mathbf{K}_h^{\mathsf{T}}\mathbf{x}, \mathbf{y})$  is orthogonal to  $Y_{\mathscr{C}} \otimes Y_{\mathscr{C}}$  for all  $Y_{\mathscr{C}}$  in  $\mathscr{C}^{\mathsf{th}}$  harmonic

- out of total dimension  $N(d, \ell)$
- All H heads are spanned by  $\leq H \cdot M(d, \ell)$  harmonics

onics, e.g. 
$$\phi_h\left(\begin{bmatrix}x_1\\x_2\end{bmatrix}, \mathbf{y}\right) \perp x_3^5$$

except for a subspace of dimension  $M(d, \ell) \leq \binom{r+\ell}{\ell}$ 

 $\mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \mathcal{U}(\mathbb{S}^{d-1})} \left( \operatorname{sgn} \left( \mathbf{x}^{\mathsf{T}} \mathbf{y} \right) - \sum_{h=1}^{H} g_{h}(\mathbf{x}, \mathbf{y}) \right)^{\mathsf{T}}$ 

 $= \left\| \sum_{\ell=0}^{\infty} \left[ \sum_{i=1}^{N(d,\ell)} \eta_{\ell} Y_{\ell}^{i} \otimes Y_{\ell}^{i} - \sum_{h=1}^{H} \sum_{i=1}^{M(d,\ell)} ?_{h,\ell,i} Y_{\ell}^{i} \otimes Y_{\ell}^{i} \right] \right\|^{-1}$ 

 $\geq \left\| \sum_{\ell=1}^{\infty} \sum_{\ell=1}^{N(d,\ell)} \eta_{\ell} \cdot Y_{\ell}^{i} \otimes Y_{\ell}^{i} \right\|^{2} = \sum_{\ell=1}^{\infty} \eta_{\ell}^{2} \left( N(d,\ell) - H \cdot M(d,\ell) \right)$  $\mathcal{L}=0 \ i=H \cdot M(d, \ell)$ 

# Combining

ℓ=0

## Combining

• Very roughly,

$$\sum_{\ell=0}^{\infty} \eta_{\ell}^2 \left( N(d,\ell) - H \cdot M(d,\ell) \right) \gtrsim \sum_{\ell=0}^{\infty} \ell$$

So unless  $H > d^p / r^p$  we are left with



 $\sum \ell^{-2} \left( d^{\ell} - Hr^{\ell} \right)$ 

h 
$$\sum_{\ell > p} l^{-2} = p^{-1}$$
 error

## [End Proof Sketch]

## Lower Bound

constants c, c', C and C' such that if either of the following sets of assumptions hold:

think of  $H \lesssim (d/r)^{1/\epsilon}$  $H \leq C \cdot 2^{d - (r+1)\log_2(2d/r)}$ .  $H \leq \frac{1}{2} \left( \frac{1}{2e} \cdot \frac{d}{r + C'/\epsilon} \right)^{C'/\epsilon} .$ 

1. High-accuracy regime:  $r \leq d - 3$ ,  $\epsilon \leq \frac{c}{d+1}$ , and 2. High-dimensional regime:  $d \ge 5$ ,  $\epsilon \ge \frac{c'}{d-2e^2 \cdot r}$  and

Then, for any choice of H rank-r generalized attention heads  $\phi_h : \mathbb{R}^{r \times 2} \to \Delta^1, V_h \in \mathbb{R}^{d \times d}, K_h \in \mathbb{R}^{d \times d}$  $\mathbb{R}^{d \times r}$  the error of approximating the nearest neighbor function is bounded as follows

$$\underset{\substack{\boldsymbol{x}_{1}, \boldsymbol{x}_{2} \sim \mathcal{D}_{2}(\mathbb{S}^{d-1})\\\boldsymbol{y} \sim \text{Unif}(\mathbb{S}^{d-1})}}{\mathbb{E}} \left\| f(\boldsymbol{X}; \boldsymbol{y}) - \sum_{h=1}^{H} \boldsymbol{V}_{h} \boldsymbol{X} \phi_{h} \left( \boldsymbol{K}_{h}^{\top} \boldsymbol{X}, \boldsymbol{y} \right) \right\|_{2}^{2} \geq \epsilon ,$$

**Theorem 2** (Low-Rank Approximation Lower Bounds, Equivariant Case). There exist universal

error of approximation is bounded as follows:

$$\mathbb{E}_{\substack{\boldsymbol{x}_1, \boldsymbol{x}_2 \sim \mathcal{D}_2(\mathbb{S}^{d-1}) \\ \boldsymbol{y} \sim \mathcal{N}\left(0, \frac{1}{\sqrt{d}}I\right)}} \left[ \left\| T(\boldsymbol{x}_1, \boldsymbol{x}_2, \boldsymbol{y}) - f^*(\boldsymbol{x}_1, \boldsymbol{x}_2, \boldsymbol{y}) \right\|^2 \right] > \frac{1}{40}$$

## Alternative Lower Bound

There exists a target function  $f^*$  such that unless  $H \cdot \max \|\mathbf{V}_h\| \lesssim c^{d-r}$ , the

### • Differences: bias in upper bound, dependence on $\|\mathbf{V}_h\|$ , proof techniques

## What about more layers?

- Low-rank attention layers are weaker, even for large  ${\cal H}$
- **Q:** Is a low-rank *transformer* weaker? ... idk
- Construction: *modified* rank-1 transformer works for N = 2
- Conjecture: low-rank transformer fails for large N, unlike full-rank one

## Modified Attention Construction: L = 2, N = 2

- Majority vote of many such heads is correct with high probability
- To tally the votes, need extra "index" and "scratchpad" dimensions

• Input: 
$$\begin{bmatrix} \mathbf{x}_1 \\ 1 \\ 0 \end{bmatrix}$$
,  $\begin{bmatrix} \mathbf{x}_2 \\ -1 \\ 0 \end{bmatrix}$ ,  $\begin{bmatrix} \mathbf{y} \\ 0 \\ 0 \end{bmatrix}$ 

- Attn layer 2: look up the target  $\mathbf{x}_1$  or  $\mathbf{x}_2$  whose sign matches the tally

Random rank-1 head hm  $(\mathbf{X}^{\mathsf{T}}\mathbf{q}\mathbf{q}^{\mathsf{T}}\mathbf{y})$  guesses correctly w.p.  $\sim \frac{1}{2} + \frac{1}{\sqrt{d}}$ 

• Attn layer 1: each head votes. Sum the votes in index dimension. Save in y's scratchpad dimension



## Experiments

## **Experimental Setup**

- Off-the-shelf multilayer transformers from PyTorch (but  $H = d^c/r$ , RMSNorm)
- "Farthest neighbor" with self-attention  $\mathbf{x}_1, \dots, \mathbf{x}_N \sim (\mathbb{S}^{d-1})$
- No positional encodings (yes biases)
- Fix d = 64, N = 16
- Best of 5 runs





Rank



# What does full-rank attention learn? • Expected: $\mathbf{I}_d \mathbf{X} \operatorname{sm} \left( \mathbf{X}^{\mathsf{T}} \left( -10^{10} \cdot \mathbf{I}_d \right) \mathbf{y} \right)$







Rank

# Role of M



Rank

## Conclusion

## Low-rank attention is fundamentally weaker than full-rank attention, even for $H \gg d/r$



# **Open / in progress**

- Can we prove hardness for L > 1? Is our conjecture true?
- Other tradeoffs in transformer hyperparameters besides r and H
- Are there hard problems that look more like text (not isotropic?)

• Is the lower bound tight? (pretty much: just use random rank-1 heads)