

Multi-Head Attention Layer

- Input dimension • d
- Rank (query/key dimension) • r
- Number of heads • *H*

Standard Self-Attention Head:

 $\mathbf{X} = \begin{bmatrix} \mathbf{X}_1 & \cdots & \mathbf{X}_N \end{bmatrix}$ $\mathbf{Q} = \mathbf{W}_O \mathbf{X}, \quad \mathbf{K} = \mathbf{W}_K \mathbf{X}, \quad \mathbf{V} = \mathbf{W}_V \mathbf{X}$ $\mathbf{W}_{\boldsymbol{O}}^{\mathsf{T}} \mathbf{V}$ softmax $(\mathbf{K}^{\mathsf{T}} \mathbf{Q})$ $= \mathbf{X}^{\top} \left(\mathbf{W}_{K}^{\top} \mathbf{W}_{Q} \right) \mathbf{X}$

inputs in \mathbb{R}^d outputs in \mathbb{R}^d

Generalized Multi-Head Cross-Attention:

$$A_{H,r} (\mathbf{X}, \mathbf{y}) = \sum_{h=1}^{H} \mathbf{M}_{h} \mathbf{X} \ \phi_{h} (\mathbf{W}_{h} \mathbf{X}, \mathbf{y})$$

• $\mathbf{W}_{h} \text{ maps } \mathbb{R}^{d} \to \mathbb{R}^{r}$ (like $\mathbf{W}_{K} \text{ above}$)
• $\mathbf{M}_{h} \text{ maps } \mathbb{R}^{d} \to \mathbb{R}^{d}$ (like $\mathbf{W}_{O}^{\top} \mathbf{W}_{V} \text{ above}$)

- ϕ_h is any function that outputs a probability distribution over [N]
- "Cross-attention": query comes from \mathbf{y} , keys come from \mathbf{X}

Standard Hyperparameter Setting

Large <i>d</i>	Low rank $r \ll d$

Research Question: Is this always a wise choice?

	Attn Is All You Need (2017)	GPT-3 (2020)	Llama-2 (2023)
# layers	6	96	80
d	512	12,288	8,192
r	64	128	128
Н	d/r	d/r	d/r

Analogous Work: Depth Separation for MLPs

- Why do deep neural nets beat shallow ones?
- <u>Theorem</u>: Some functions on \mathbb{R}^d can be represented by a depth-3 MLP with poly(d) width, but cannot be approximated by depth-2 MLP unless width is exp(d)
- "Pay a bit for depth, save *many* parameters overall."

Multi-head attention with small query/ key dimension is ubiquitous, but 1-head attention with big queries/ keys can be more powerful.

Main Theorem

The "nearest neighbor" function on the *d*-sphere

$$f(\mathbf{x}_1, \dots, \mathbf{x}_N; \mathbf{y}) :=$$

can be represented by one full-rank attention head (r = d, H = 1),

$$f(\mathbf{x}_1, \ldots, \mathbf{x}_N; \mathbf{y}) =$$

but cannot be ε -approximated by low-rank attention (r < d) unless $H \gtrsim (d/r)^{1/\varepsilon}$:

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

Alternative Rank-Separation Theorem

A certain linear combination of biased nearest neighbor functions can be represented by full-rank biased attention with $H = d^2$, but not approximated by low-rank attention unless $H \gtrsim C^{d-r}$ / $\|V_h\|^2$

queries/keys/values in \mathbb{R}^r

H = d/r

Massive scaling up

Insignificant change

The Low-Rank Bottleneck in Attention Noah Amsel, Gilad Yehudai, Joan Bruna

- $\operatorname{arg\,min} \|\mathbf{x} \mathbf{y}\|_2$ $\mathbf{x} \in \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$

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

- $\mathbf{y} A_{H,r}(\mathbf{x}_1, \mathbf{x}_2; \mathbf{y}) \|^2 > \varepsilon$



Experiments

ask:	lr
	e
Nodel:	Li
Result:	Н
	n

Task:	Nea
Model:	Trar
Result:	1-la
	alw



Role of Depth

Conclusion

Try increasing r and H more than the standard setting, especially for non-textual data. Theorists shouldn't just assume full-rank attention.



n-context learning of linear functions in 20 dimensions, as in Garg et. al (2022): $(\mathbf{x}_1, \mathbf{w}^{\mathsf{T}} \mathbf{x}_1, \dots, \mathbf{x}_{40}, \mathbf{w}^{\mathsf{T}} \mathbf{x}_{40}, \mathbf{x}_{41}) \mapsto \mathbf{w}^{\mathsf{T}} \mathbf{x}_{41}$

inear embedding + 12-layer transformer with d = 48, H = d/r.

High-rank attention far outperforms low-rank attention with the same number of parameters.



earest neighbor function with d = 64, N = 16, iid Gaussian data Insformers with H = d/r, $d^{1.5}r$, or d^2/r .

ayer transformer succeeds iff using full rank attn (r = d). Full-rank is vays better, even vs. low-rank model with more parameters.

• Our main theorems concern 1-layer of attention.

• Is large rank still needed for multi-layer transformers?

• <u>Theorem</u>: We can approximate the nearest neighbor function on N = 2 points with 2 layers of rank-1 attention *if* we add extra dimensions for positional encodings.

• <u>Conjecture</u>: We cannot do this for general N.