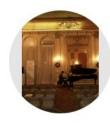
Rethinking Attention with Performers

Attention with linear time

허세민



Author

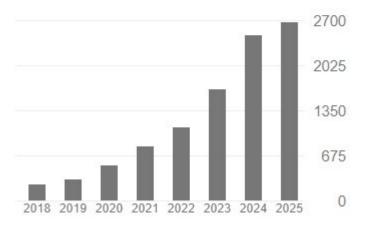


Krzysztof Choromanski

Google DeepMind Robotics & <u>Columbia University</u> columbia.edu의 이메일 확인됨 - <u>홈페이지</u>

robotics reinforcement learning efficient Transformers quasi Monte Carlo methods

인용	모두 보기	
	전체	2020년 이후
서지정보	10501	9606
h-index	40	37
i10-index	82	76



Rethinking attention with performers K Choromanski, V Likhosherstov, D Dohan, X Song, A Gane, T Sarlos, arXiv preprint arXiv:2009.14794	2615	2020
Rt-2: Vision-language-action models transfer web knowledge to robotic control B Zitkovich, T Yu, S Xu, P Xu, T Xiao, F Xia, J Wu, P Wohlhart, S Welker, Conference on Robot Learning, 2165-2183	1860	2023
Socratic models: Composing zero-shot multimodal reasoning with language A Zeng, M Attarian, B Ichter, K Choromanski, A Wong, S Welker, arXiv preprint arXiv:2204.00598	654	2022
Explaining how a deep neural network trained with end-to-end learning steers a car M Bojarski, P Yeres, A Choromanska, K Choromanski, B Firner, L Jackel, arXiv preprint arXiv:1704.07911	614	2017
Orthogonal random features FXX Yu, AT Suresh, KM Choromanski, DN Holtmann-Rice, S Kumar Advances in neural information processing systems 29	281	2016
End to end learning for self-driving cars. arXiv 2016 M Bojarski, D Del Testa, D Dworakowski, B Firner, B Flepp, P Goyal, arXiv preprint arXiv:1604.07316 103	268	2016
A theoretical and empirical comparison of gradient approximations in derivative-free optimization AS Berahas, L Cao, K Choromanski, K Scheinberg Foundations of Computational Mathematics 22 (2), 507-560	245	2022
Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities G Comanici, E Bieber, M Schaekermann, I Pasupat, N Sachdeva, I Dhillon, arXiv preprint arXiv:2507.06261	200	2025
Effective diversity in population based reinforcement learning J Parker-Holder, A Pacchiano, KM Choromanski, SJ Roberts Advances in Neural Information Processing Systems 33, 18050-18062	198	2020
Structured evolution with compact architectures for scalable policy optimization K Choromanski, M Rowland, V Sindhwani, R Turner, A Weller International Conference on Machine Learning, 970-978	170	2018
Es-maml: Simple hessian-free meta learning X Song, W Gao, Y Yang, K Choromanski, A Pacchiano, Y Tang arXiv preprint arXiv:1910.01215	151	2019

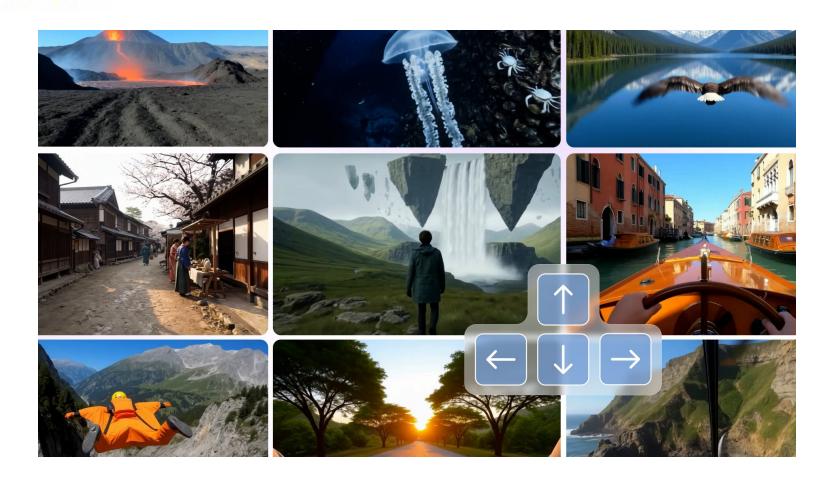
Contents

- 1. Introduction
- 2. Attention
- 3. Transformers are RNNs
- 4. Performers
- 5. Conclusion

Intro

Towards world Simulation

"DeepMind has ambitious plans to make massive generative models that simulate the world," Brooks wrote Monday morning. "I'm hiring for a new team with this mission."



Titans

Titans: Learning to Memorize at Test Time

Ali Behrouz[†], Peilin Zhong[†], and Vahab Mirrokni[†]

Google Research
{alibehrouz, peilinz, mirrokni}@google.com

Prior Knowledge

- ✓ Test-Time Training (TTT)
 - 추론 시점(test time)에서 모델이 스스로 적응 할 수 있도록 하는 방법
 - 새로운 환경이나 분포 변화에 맞게 온라인 학습처럼 동작
 - Titan도 추론 시점에서 모델을 적응 시키는 전략을 활용함
- ✓ Linear Attention
 - 긴 시퀀스를 다룰 때, 기존 Attention의 $O(n^2)$ 복잡도를 줄이려는 접근
 - 시퀀스 길이가 길어져도 효율적으로 처리 가능
 - Titan은 효율적인 attention 구조를 필요로 하며, Linear Attention 아이디어가 그 기반이 됨

Attention

Attention

The	The
animal	animal
didn't	didn't
cross	cross
the	the
street	street
because	because
it	it
was	was
too	too
tired	tired

Attention
$$(Q, K, V) = \operatorname{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

Advantages

- 장거리 의존성을 효과적으로 포착할 수 있음
- 학습 과정에서 완전 병렬화가 가능함

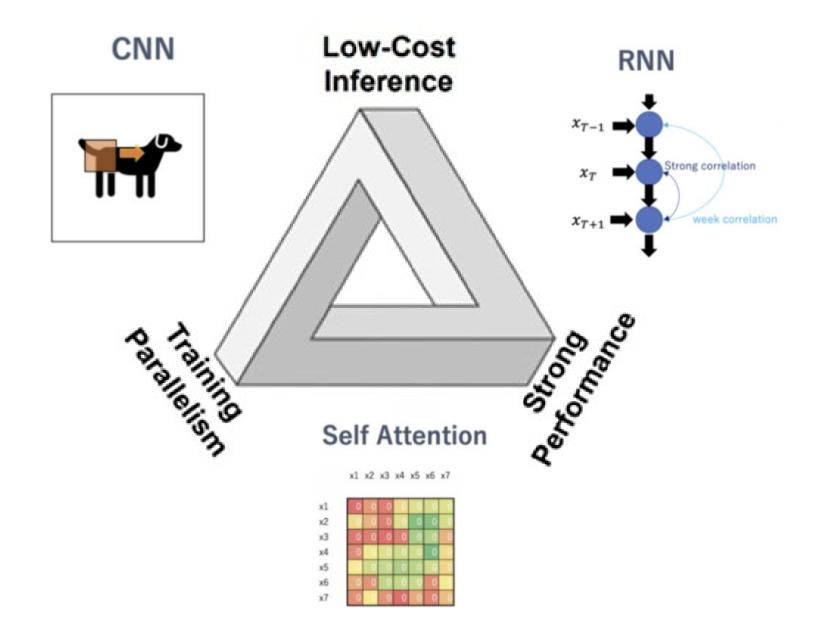
Limitations

- 시간과 메모리 복잡도가 제곱(Quadratic, O(n²))
- 매우 긴 시퀀스를 다루기 어려움

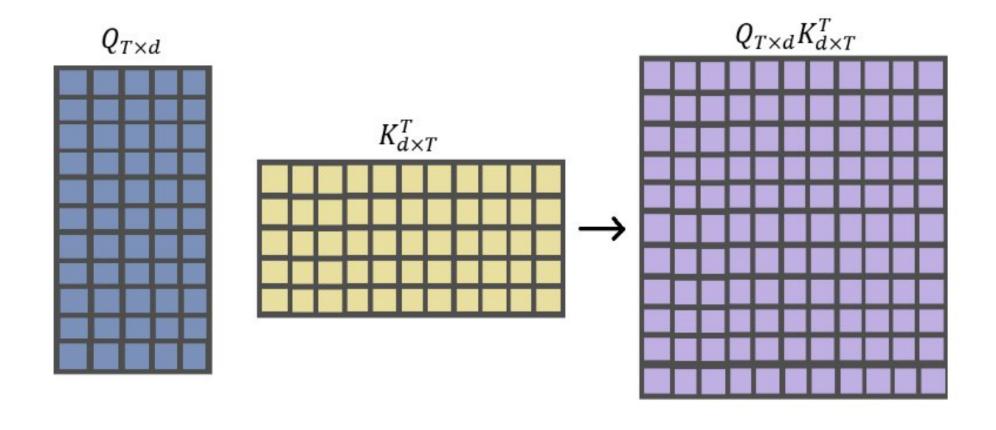
Alternatives

- FlashAttention
- Performer

Penrose Triangle



Complexity



Transformers are RNNs

Linear Attention

$$V_i' = \frac{\sum_{j=1}^{N} \sin(Q_i, K_j) V_j}{\sum_{j=1}^{N} \sin(Q_i, K_j)}.$$

Attention

$$V' = \operatorname{softmax}\left(\frac{QK^T}{\sqrt{D}}\right)V.$$

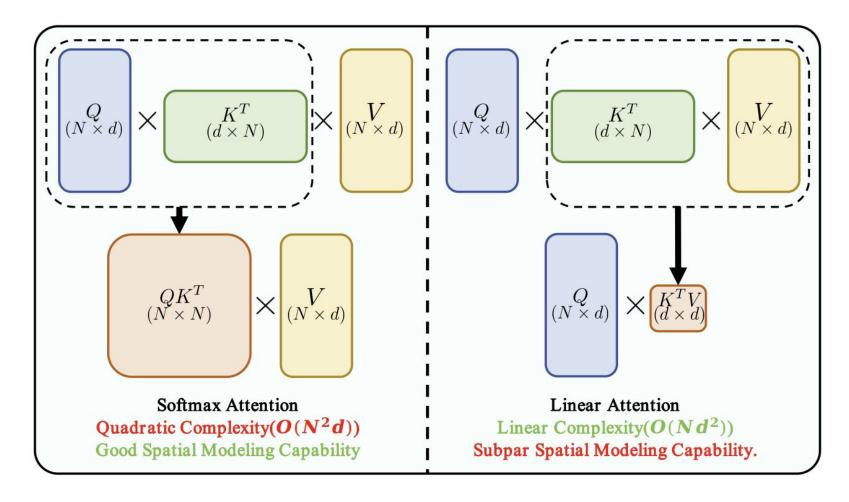
Linear Attention

$$V_{i}' = \frac{\sum_{j=1}^{N} \phi(Q_{i})^{T} \phi(K_{j}) V_{j}}{\sum_{j=1}^{N} \phi(Q_{i})^{T} \phi(K_{j})}$$

$$\downarrow$$

$$V_{i}' = \frac{\phi(Q_{i})^{T} \sum_{j=1}^{N} \phi(K_{j}) V_{j}^{T}}{\phi(Q_{i})^{T} \sum_{j=1}^{N} \phi(K_{j})}.$$

Linear Attention



$$\operatorname{softmax}\left(\frac{QK^{T}}{\sqrt{D}}\right)V. \qquad \qquad \phi\left(Q\right)\left(\phi\left(K\right)^{T}V\right).$$

Transformers are RNNs

$$s_{0} = 0,$$
 $s_{i} = \sum_{j=1}^{i} \phi(K_{j}) V_{j}^{T},$ $z_{0} = 0,$ $z_{i} = \sum_{j=1}^{i} \phi(K_{j}),$ $s_{i} = s_{i-1} + \phi(x_{i}W_{K})(x_{i}W_{V})^{T},$ $z_{i} = z_{i-1} + \phi(x_{i}W_{K}),$ $y_{i} = f_{l} \left(\frac{\phi(x_{i}W_{Q})^{T} s_{i}}{\phi(x_{i}W_{Q})^{T} z_{i}} + x_{i}\right).$

Transformers are RNNs

Advantages

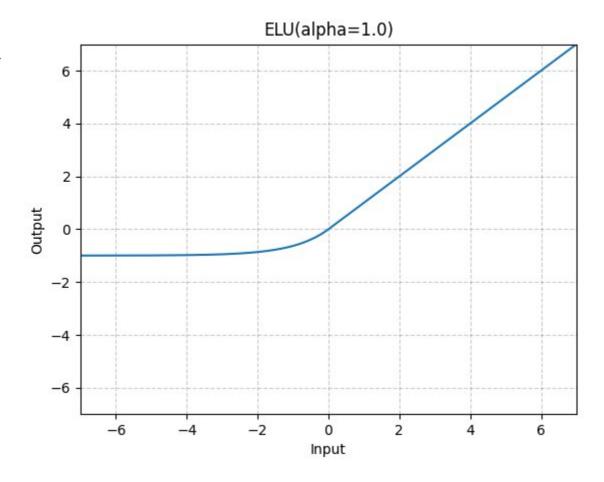
Performer

Motivation

- ✓ 수치적 불안정성
 - 근사된 Attention 값이 발산하거나 collapse 가능
- ✔ 근사 품질 저하
 - Softmax kernel을 정확히 근사하지 못함
- ✓ 분산 문제
 - 근사 분산이 매우 커져서 불안정

$$\phi\left(x\right) = \operatorname{elu}(x) + 1$$

$$V_{i}' = \frac{\phi(Q_{i})^{T} \sum_{j=1}^{N} \phi(K_{j}) V_{j}^{T}}{\phi(Q_{i})^{T} \sum_{j=1}^{N} \phi(K_{j})}.$$



Author

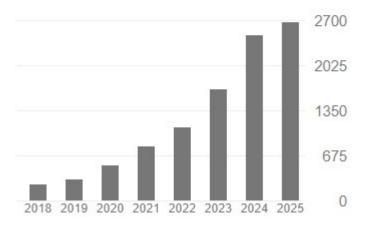


Krzysztof Choromanski

Google DeepMind Robotics & <u>Columbia University</u> columbia.edu의 이메일 확인됨 - <u>홈페이지</u>

robotics reinforcement learning efficient Transformers quasi Monte Carlo methods

인용		모두 보기	
	전체	2020년 이후	
서지정보	10501	9606	
h-index	40	37	
i10-index	82	76	



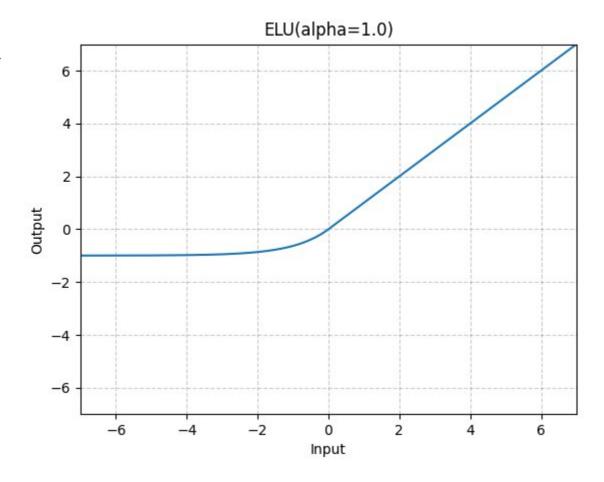
Rethinking attention with performers K Choromanski, V Likhosherstov, D Dohan, X Song, A Gane, T Sarlos, arXiv preprint arXiv:2009.14794	2615	2020
Rt-2: Vision-language-action models transfer web knowledge to robotic control B Zitkovich, T Yu, S Xu, P Xu, T Xiao, F Xia, J Wu, P Wohlhart, S Welker, Conference on Robot Learning, 2165-2183	1860	2023
Socratic models: Composing zero-shot multimodal reasoning with language A Zeng, M Attarian, B Ichter, K Choromanski, A Wong, S Welker, arXiv preprint arXiv:2204.00598	654	2022
Explaining how a deep neural network trained with end-to-end learning steers a car M Bojarski, P Yeres, A Choromanska, K Choromanski, B Firner, L Jackel, arXiv preprint arXiv:1704.07911	614	2017
Orthogonal random features FXX Yu, AT Suresh, KM Choromanski, DN Holtmann-Rice, S Kumar Advances in neural information processing systems 29	281	2016
End to end learning for self-driving cars. arXiv 2016 M Bojarski, D Del Testa, D Dworakowski, B Firner, B Flepp, P Goyal, arXiv preprint arXiv:1604.07316 103	268	2016
A theoretical and empirical comparison of gradient approximations in derivative-free optimization AS Berahas, L Cao, K Choromanski, K Scheinberg Foundations of Computational Mathematics 22 (2), 507-560	245	2022
Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities G Comanici, E Bieber, M Schaekermann, I Pasupat, N Sachdeva, I Dhillon, arXiv preprint arXiv:2507.06261	200	2025
Effective diversity in population based reinforcement learning J Parker-Holder, A Pacchiano, KM Choromanski, SJ Roberts Advances in Neural Information Processing Systems 33, 18050-18062	198	2020
Structured evolution with compact architectures for scalable policy optimization K Choromanski, M Rowland, V Sindhwani, R Turner, A Weller International Conference on Machine Learning, 970-978	170	2018
Es-maml: Simple hessian-free meta learning X Song, W Gao, Y Yang, K Choromanski, A Pacchiano, Y Tang arXiv preprint arXiv:1910.01215	151	2019

Motivation

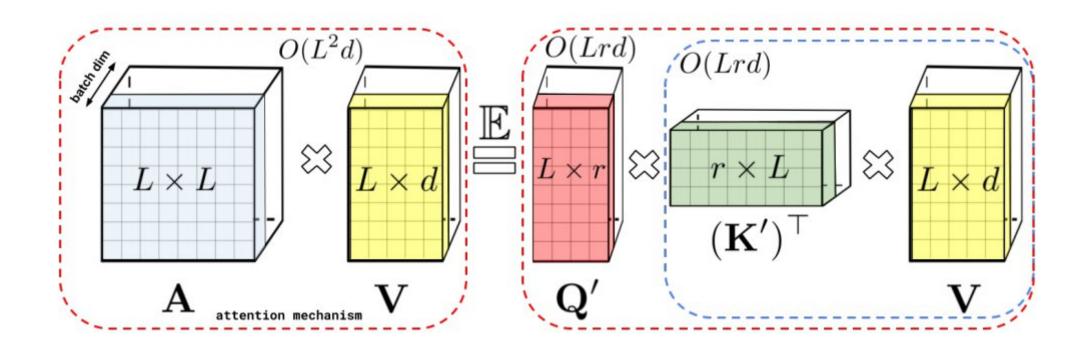
- ✓ 수치적 불안정성
 - 근사된 Attention 값이 발산하거나 collapse 가능
- ✔ 근사 품질 저하
 - Softmax kernel을 정확히 근사하지 못함
- ✓ 분산 문제
 - 근사 분산이 매우 커져서 불안정

$$\phi\left(x\right) = \operatorname{elu}(x) + 1$$

$$V_{i}' = \frac{\phi(Q_{i})^{T} \sum_{j=1}^{N} \phi(K_{j}) V_{j}^{T}}{\phi(Q_{i})^{T} \sum_{j=1}^{N} \phi(K_{j})}.$$



FAVOR+ (FA)



FAVOR+ (R+)

Lemma 1 (Positive Random Features (PRFs) for Softmax). For $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, $\mathbf{z} = \mathbf{x} + \mathbf{y}$ we have:

$$SM(\mathbf{x}, \mathbf{y}) = \mathbb{E}_{\omega \sim \mathcal{N}(0, \mathbf{I}_d)} \left[\exp \left(\omega^\top \mathbf{x} - \frac{\|\mathbf{x}\|^2}{2} \right) \exp \left(\omega^\top \mathbf{y} - \frac{\|\mathbf{y}\|^2}{2} \right) \right] = \Lambda \mathbb{E}_{\omega \sim \mathcal{N}(0, \mathbf{I}_d)} \cosh(\omega^\top \mathbf{z}), (7)$$

where $\Lambda = \exp(-\frac{\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2}{2})$ and \cosh is hyperbolic cosine. Consequently, softmax-kernel admits a positive random feature map unbiased approximation with $h(\mathbf{x}) = \exp(-\frac{\|\mathbf{x}\|^2}{2})$, l = 1, $f_1 = \exp$ and $\mathcal{D} = \mathcal{N}(0, \mathbf{I}_d)$ or: $h(\mathbf{x}) = \frac{1}{\sqrt{2}} \exp(-\frac{\|\mathbf{x}\|^2}{2})$, l = 2, $f_1(u) = \exp(u)$, $f_2(u) = \exp(-u)$ and the same \mathcal{D} (the latter for further variance reduction). We call related estimators: \widehat{SM}_m^+ and $\widehat{SM}_m^{\text{hyp+}}$.

FAVOR+ (R+)

F.1 PROOF OF LEMMA 1

Proof. We first deduce that for any $a, b \in \mathbb{R}^d$

$$SM(\mathbf{x}, \mathbf{y}) = \exp(\mathbf{x}^{\top} \mathbf{y}) = \exp(-\|\mathbf{x}\|^2/2) \cdot \exp(\|\mathbf{x} + \mathbf{y}\|^2/2) \cdot \exp(-\|\mathbf{y}\|^2/2)$$

Next, let $w \in \mathbb{R}^d$. We use the fact that

$$(2\pi)^{-d/2} \int \exp(-\|\boldsymbol{w} - \boldsymbol{c}\|_2^2/2) d\boldsymbol{w} = 1$$

for any $c \in \mathbb{R}^d$ and derive:

$$\exp(\|\boldsymbol{x} + \boldsymbol{y}\|^{2}/2) = (2\pi)^{-d/2} \exp(\|\boldsymbol{x} + \boldsymbol{y}\|^{2}/2) \int \exp(-\|\boldsymbol{w} - (\boldsymbol{x} + \boldsymbol{y})\|^{2}/2) d\boldsymbol{w}$$

$$= (2\pi)^{-d/2} \int \exp(-\|\boldsymbol{w}\|^{2}/2 + \boldsymbol{w}^{\top}(\boldsymbol{x} + \boldsymbol{y}) - \|\boldsymbol{x} + \boldsymbol{y}\|^{2}/2 + \|\boldsymbol{x} + \boldsymbol{y}\|^{2}/2) d\boldsymbol{w}$$

$$= (2\pi)^{-d/2} \int \exp(-\|\boldsymbol{w}\|^{2}/2 + \boldsymbol{w}^{\top}(\boldsymbol{x} + \boldsymbol{y})) d\boldsymbol{w}$$

$$= (2\pi)^{-d/2} \int \exp(-\|\boldsymbol{w}\|^{2}/2) \cdot \exp(\boldsymbol{w}^{\top}\boldsymbol{x}) \cdot \exp(\boldsymbol{w}^{\top}\boldsymbol{y}) d\boldsymbol{w}$$

$$= \mathbb{E}_{\omega \sim \mathcal{N}(\mathbf{0}_{d}, \mathbf{I}_{d})} [\exp(\omega^{\top}\boldsymbol{x}) \cdot \exp(\omega^{\top}\boldsymbol{y})].$$

That completes the proof of the first part of the lemma. An identity involving hyperbolic cosine function is implied by the fact that for every $\mathbf{u} \in \mathbb{R}^d$ and $\omega \sim \mathcal{N}(0, \mathbf{I}_d)$ the following is true:

$$\mathbb{E}[\exp(\omega^{\top}\mathbf{u})] = \sum_{i=0}^{\infty} \frac{\mathbb{E}[(\omega^{\top}\mathbf{u})^{2i}]}{(2i)!} = \frac{1}{2} \sum_{i=0}^{\infty} \frac{\mathbb{E}[(\omega^{\top}\mathbf{u})^{2i}] + \mathbb{E}[(-\omega^{\top}\mathbf{u})^{2i}]}{(2i)!}.$$
 (12)

The cancellation of the odd moments $\mathbb{E}[(\omega^{\top}\mathbf{u})^{2i+1}]$ follows directly from the fact that ω is taken from the isotropic distribution (i.e. distribution with pdf function constant on each sphere). That completes the proof.

FAVOR+ (O)

Theorem 1 (regularized versus softmax-kernel). Assume that the L_{∞} -norm of the attention matrix for the softmax-kernel satisfies: $\|\mathbf{A}\|_{\infty} \leq C$ for some constant $C \geq 1$. Denote by \mathbf{A}^{reg} the corresponding attention matrix for the regularized softmax-kernel. The following holds:

$$\inf_{i,j} \frac{\mathbf{A}^{\text{reg}}(i,j)}{\mathbf{A}(i,j)} \ge 1 - \frac{2}{d^{\frac{1}{3}}} + o\left(\frac{1}{d^{\frac{1}{3}}}\right), \text{ and } \sup_{i,j} \frac{\mathbf{A}^{\text{reg}}(i,j)}{\mathbf{A}(i,j)} \le 1.$$
 (9)

Furthermore, the latter holds for $d \geq 2$ even if the L_{∞} -norm condition is not satisfied, i.e. the regularized softmax-kernel is a universal lower bound for the softmax-kernel.

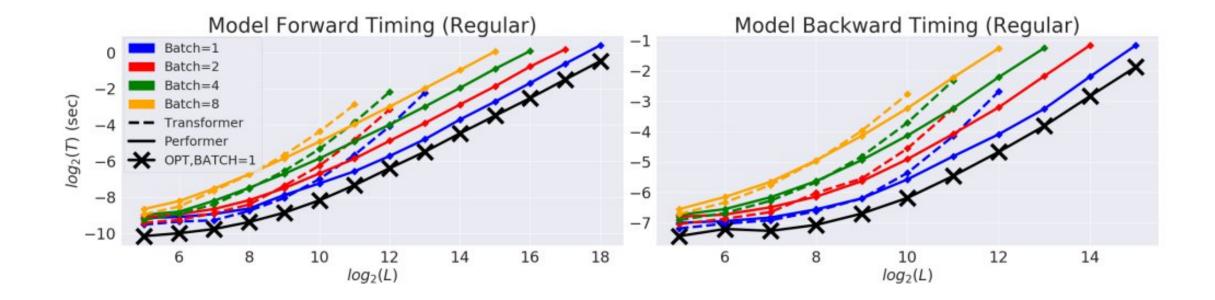
Theorem 2. If $\widehat{SM}_m^{\text{ort}+}(\mathbf{x}, \mathbf{y})$ stands for the modification of $\widehat{SM}_m^+(\mathbf{x}, \mathbf{y})$ with orthogonal random features (and thus for $m \leq d$), then the following holds for any d > 0:

$$MSE(\widehat{SM}_{m}^{\text{ort+}}(\mathbf{x}, \mathbf{y})) \le MSE(\widehat{SM}_{m}^{+}(\mathbf{x}, \mathbf{y})) - \frac{2(m-1)}{m(d+2)} \left(SM(\mathbf{x}, \mathbf{y}) - \exp\left(-\frac{\|\mathbf{x}\|^{2} + \|\mathbf{y}\|^{2}}{2}\right) \right)^{2}.$$
(10)

Furthermore, completely analogous result holds for the regularized softmax-kernel SMREG.

Experiments

FAVOR+ (FA)



FAVOR+ (OR+)

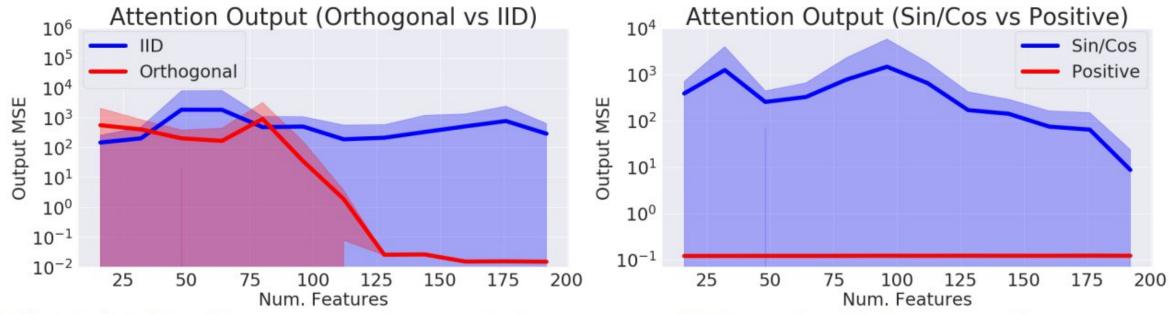


Figure 4: MSE of the approximation output when comparing Orthogonal vs IID features and trigonometric $\sin/\cos vs$ positive features. We took L=4096, d=16, and varied the number of random samples m. Standard deviations shown across 15 samples of appropriately normalized random matrix input data.

Conclusion & Broader Impact

Conclusion

- ✓ 핵심 아이디어
 - FAVOR+
- ✓ 핵심 기여
 - Attention을 $O(n^2) \rightarrow O(n)$ 으로 근사
 - 무편향, 안정성, 수렴성을 수학적으로 보장
 - 다양한 실험에서 Transformer와 동등하거나 더 나은 성능
- ✓ 의의와 영향
 - **대규모 시퀀스 처리 가능 →** 긴 문장, 유전체 데이터, 영상 등 확장성
 - 자원 효율적 → 메모리/시간 절약, 친환경적 AI 연구 기여 가능
 - 연구적 영향 → 이후 Linear Attention 계열 연구(Linformer, Nyströmformer, FlashAttention 등)의 기초를 마련
 - 한계 인식: approximation error 존재, 특정 task에서는 성능 저하 가능

감사합니다

