

#### Understanding Transformer from the Perspective of Associative Memory

Shu Zhong<sup>\*,†</sup>, Mingyu Xu<sup>\*</sup>, Tenglong Ao<sup>\*</sup>, Guang Shi<sup>†</sup>

ByteDance Seed

\*Equal contribution, <sup>†</sup>Corresponding authors



#### Part 1. Associative Memory and Its Capacity

#### Part 2. Update the Associative Memory

#### Part 3. Exploration of Next Generation Models

#### There is a contradiction between expressiveness and parallelism.

Q1: If you can **add** any **two** numbers in parallel each time, **how much** time do we need to calculate the sum of n elements ? The answer is.  $O(\log^{1} n)$ .

 $x_1 + x_2 + \dots + x_n$ 

Q2:If you can **add** any **two** numbers in parallel each time, **how much** time do we need to calculate the add of two matrix with size  $n \times n$ ? The answer is  $O(1) = O(\log^0 n)$ .

A + B

Time: Q1 > Q2 Operation in total: Q1 < Q2

#### There is a contradiction between expressiveness and parallelism.

Q1: If you can **add** any **two** numbers in parallel each time, **how much** time do we need to calculate the sum of *n* elements?  $x_1 + x_2 + \dots + x_n$  The answer is.  $O(\log^1 n)$ .

Q2:If you can add any two numbers in parallel each time, how much  
time do we need to calculate the add of two matrix with size 
$$n \times n$$
?  $A + B$   
The answer is  $O(1) = O(\log^9 n)$ .  
And,  $Or$ , Not gate  $2$   $O(\log^i n)$ .

There is a contradiction between expressiveness and parallelism.

Where is the popular model, such as Transformer?

Problem can solve without COT	Parallelism
Linear Programming	Р
Matrix Inverse	$NC^2$
Acyclic BCQ	LOGCFL
STCON	NL
USTCON	SL
S <sub>5</sub> , Python, Chess	$NC^1$
Parity Check	$TC^{0}$
Unary Language	AC <sup>0</sup>
Add Bool Matrix	NC <sup>0</sup>

The parallelism tradeoff: Limitations of log-precision transformers. TACL 2023. Chain of Thought Empowers Transformers to Solve Inherently Serial Problems. Arxiv 2024.



There is a contradiction between expressiveness and parallelism.

Where is the popular model, such as Transformer?

 $TC^0$  or  $AC^0$  (log or constant precision )

Problem can solve without COT	Parallelism
Linear Programming	Р
Matrix Inverse Acyclic BCQ	NC <sup>2</sup> LOGCFL
STCON	NL
S <sub>5</sub> , Python, Chess	SL NC <sup>1</sup>
Parity Check Unary Language	TC <sup>0</sup> AC <sup>0</sup>
Add Bool Matrix	$NC^{0}$

The parallelism tradeoff: Limitations of log-precision transformers. TACL 2023. Chain of Thought Empowers Transformers to Solve Inherently Serial Problems. Arxiv 2024.



One the most important reason that Transformer beat RNN:

## GPUs + High Parallelism



One the most important reason that Transformer beat RNN:

GPUs + High Parallelism

However, we know that there is a contradiction between expressiveness and parallelism.

Is there a model with more expressiveness than transformer but also can be parallelized in GPUs ?

#### DeltaNet: A model which beyond $TC^0$

#### 1980s and 1990s

variable weight wp. The output at time t is

$$p(t) = \sum_{i=1}^{n} w p_i(t) w_i(t).$$

The weights change over time according to the following equation: for  $i=1,\ldots,n$ ,

$$wp_{i}(t+1)=wp_{i}(t)+cp[z(t)-p(t-1)]x_{i}(t-1)$$

#### 2020s

From the perspective of fast weight programming (Irie et al., 2022a) and test-time training (Sun et al., 2024a) and regression (Wang et al., 2025), the hidden state **S** can be interpreted as a (fast) weight matrix, with the delta rule optimizing the online regression objective  $\mathcal{L}(\mathbf{S}_t) = \frac{1}{2} \|\mathbf{S}_t \mathbf{k}_t - \mathbf{v}_t\|^2$  via *test-time* stochastic gradient descent (SGD):

$$\mathbf{S}_{t+1} = \mathbf{S}_t - \beta_t \nabla \mathcal{L}(\mathbf{S}_t) = \mathbf{S}_t - \beta_t (\mathbf{S}_t \boldsymbol{k}_t - \boldsymbol{v}_t) \boldsymbol{k}_t^{\mathsf{T}} = \mathbf{S}_t (\mathbf{I} - \beta_t \boldsymbol{k}_t \boldsymbol{k}_t^{\mathsf{T}}) + \beta_t \boldsymbol{v}_t \boldsymbol{k}_t^{\mathsf{T}}$$

$$WP_{t+1} = WP_t + cp \left(Z_t - WP_{t-1}w_{t-1}\right)x_{t-1}^{\mathsf{T}} \qquad S_t = S_{t-1} + \beta_t \left(v_t - S_{t-1}k_t\right)k_t^{\mathsf{T}}$$

#### https://web.cs.umass.edu/publication/docs/1980/UM-CS-1980-018.pdf

https://citeseerx.ist.psu.edu/document?repid=rep1&type=pdf&doi=2f0becffd2f44b198d28074d01722e4c7905dae2 https://www.cs.toronto.edu/~fritz/absps/fastweights.pdf *Parallelizing Linear Transformers with the Delta Rule over Sequence Length. NeurIPS 2024.* <sup>9</sup>



DeltaNet: A model which beyond  $TC^0$ 

1980s and 1990s

2020s

 $WP_{t+1} = WP_t + cp \left(Z_t - WP_{t-1}w_{t-1}\right)x_{t-1}^{\mathsf{T}} \qquad S_t = S_{t-1} + \beta_t \left(v_t - S_{t-1}k_t\right)k_t^{\mathsf{T}}$ 

Expect the internal output  $WP_{t-1}w_{t-1}/S_{t-1}k_t$  to be consistent with the feedback  $Z_t/v_t$ , and adjust the internal state WP/Sbased on the inconsistency  $(Z_t - WP_{t-1}w_{t-1})/(v_t - S_{t-1}k_t)$ with a learning rate  $cp/\beta_t$ .

Regarding this update rule, it can have expressive power beyond  $TC^0$ . Readers can refer to RWKV-7 "Goose" with Expressive Dynamic State Evolution. Arxiv 2025.



DeltaNet: A model beyond  $TC^0$ 

1980s and 1990s

2020s

 $WP_{t+1} = WP_t + cp \left(Z_t - WP_{t-1}w_{t-1}\right)x_{t-1}^{\mathsf{T}} \qquad S_t = S_{t-1} + \beta_t \left(v_t - S_{t-1}k_t\right)k_t^{\mathsf{T}}$ 

Most importantly, Songlin Yang et al find a efficient parallel strategy on GPUs, which makes this model have the potential to become a part of modern LLM

Parallelizing Linear Transformers with the Delta Rule over Sequence Length. NeurIPS 2024.



Limitation of limited state space.

Difficult to remember long sequences[1].

...

Transformers can adapt to a form of dynamic sparsity[2].



[1] Repeat after me: Transformers are better than state space models at copying. ICML 2024.
 [2] When Do Transformers Outperform Feedforward and Recurrent Networks? A Statistical Perspective. Arxiv2025.

#### Part 1. Associative Memory and Its Capacity

# Associative Memory

• Associative memory is defined as the ability to learn and remember the relationship between <del>unrelated</del> items.



# **Associative Memory**

- Associative memory is defined as the ability to learn and remember the relationship between <del>unrelated</del> items.
- Like a dictionary:
  - We have keys and values;
  - Learning is the process of acquiring the *fuzzy* mapping.



# Simplest Model for Associative Memory

Suppose we have key-value pairs {(k<sub>i</sub>, v<sub>i</sub>)}<sub>t</sub> where the keys are orthogonal. We can store these relationships into an outer-product associative memory (Kohonen 1972):

$$S_t = \sum_{i=1}^t v_i k_i^T$$

• Query with  $q = k_i$  we get perfect retrieval:

$$o = S_t k_i = v_i k_i^T k_i + \sum_{j \neq i} v_j k_j^T k_i = v_i$$

• This kind of outer-product memory is also known as linear attention.

# Measuring Retrieval Error

• What if keys are not orthonormal?

$$o = S_t k_i = \underbrace{v_i (k_i^T k_i)}_{Signal \, \boldsymbol{v}_i c} + \underbrace{\sum_{j \neq i} v_j k_j^T k_i}_{j \neq i} \approx v_i$$

• Define inverse retrieval SNR (Signal-to-Noise Ratio):

$$SNR^{-1} = \mathbb{E}_{\boldsymbol{v}_j, \boldsymbol{k}_j} \left[ \frac{\|\boldsymbol{r}\|^2}{c^2 \|\boldsymbol{v}_i\|^2} \right]$$

• A larger value indicates a higher noise component, resulting in lower retrieval accuracy.

Noise **r** 

# Measuring Retrieval Error

• Considering keys and values are i.i.d. standard Gaussian vectors. We can quantitatively calculate the inverse SNR:

$$SNR_{Linear}^{-1} \approx \frac{N}{d}$$

- To attain a target SNR, the linear model width d must grow linearly with the sequence length N;
- This explains why linear attention usually has poor retrieval performance.

# Kernel Trick

• Introducing a kernel function  $\kappa(x, y) = \phi(x)^T \phi(y)$ .  $\phi(\cdot)$  maps keys into higher-dimensional space.

$$S_t = \sum_{i=1}^t v_i \phi(k_i)^T$$
$$o = S_t \phi(q) = \sum_{i=1}^t v_j \phi(k_i)^T \phi(q) = \sum_{i=1}^t v_i \kappa(k_i, q)$$

• Greater separability among the vectors leads to lower retrieval error.

# Kernel Trick

• Also, considering  $q = k_i$ , and assuming keys and values are i.i.d. standard Gaussian, we get:

$$SNR_{\kappa}^{-1} = N \frac{\mathbb{E}_{k_j} [\kappa^2(k_j, k_i)]}{\kappa^2(k_i, k_i)}$$

How much the kernel suppresses irrelevant features

How much the kernel amplifies the matched features

# Exp Kernel

• For  $\kappa(x, y) = \exp\left(\frac{x^T y}{\tau}\right)$  and  $\tau = \sqrt{d}$ , we get standard softmax attention without normalization:

$$o = S_t \phi(q) = \sum_{i=1}^t v_i \exp\left(\frac{k_i^T q}{\sqrt{d}}\right)$$

• Inverse SNR is:

$$SNR_{\exp}^{-1} = N \exp\left(-\frac{2(\tau-1)}{\tau^2}d\right)$$

• From d = O(N) in linear attention to  $d = O(\log^2 N)$ .

Understanding Transformer Reasoning Capabilities via Graph Algorithms. NeurIPS 2024.

# Exp Kernel

$$SNR_{\exp}^{-1} = N \exp\left(-\frac{2(\tau-1)}{\tau^2}d\right)$$

- Increase *d* improves retrieval (but in most cases *d* is enough).
- Reducing  $\tau$  but keeping it above 1 improves retrieval.
- Multihead is good for softmax attention, but not good for linear.



## ReLU Kernel

• FFN is associative memory with  $\kappa(x, y) = \text{ReLU}(x^T y)$ :

$$FFN(x) = \sum_{i=1}^{m} W_{V_i}^T \text{ReLU}(W_{K_i}x)$$

• Query is hidden state *x*, keys and values are learnable weights.

Inverse SNR is :  $SNR_{ReLU}^{-1} \approx \frac{N}{2d}$ We suspect a kernel with lower retrieval precision encourages a more polysemantic key-value memory, which is suitable for FFN.

# Architecture Symmetry

Component	Attention	$\operatorname{FFN}$	$[\boldsymbol{q}_{j,1}; \boldsymbol{q}_{j,2}; \ldots; \boldsymbol{q}_{j,n_h}] = \boldsymbol{q}_t,$
Kernel	Exp	ReLU / SiLU	$[oldsymbol{k}_{i,1};oldsymbol{k}_{i,2};\ldots;oldsymbol{k}_{i,n_h}]=oldsymbol{k}_i,$
Normalization	Yes	None	$[oldsymbol{v}_{i,1};oldsymbol{v}_{i,2};\ldots;oldsymbol{v}_{i,n_h}]=oldsymbol{v}_i,$
Multihead	Yes	None ———	$\rightarrow$ <u>t</u>
Sparsity	None	MoE	$oldsymbol{o}_{t,\cdot} = \sum oldsymbol{v}_{i,\cdot}  \kappa(oldsymbol{q}_{t,\cdot},oldsymbol{k}_{i,\cdot}),$
Gating	None	SwiGLU Gating	$\overline{i=1}$
			$oldsymbol{o}_t = oldsymbol{W}_O[oldsymbol{o}_{t,1};oldsymbol{o}_{t,2};\ldots;oldsymbol{o}_{t,n_h}].$

• We suspect FFN can benefit from multihead too [1].

# Architecture Symmetry

Component	Attention	FFN	
Kernel	Exp	ReLU / SiLU	
Normalization	Yes	None	
Multihead	Yes	None	
Sparsity —	None	MoE	
Gating	None	SwiGLU Gating	
$\mathrm{FFN}_{\mathrm{MoE}}(\boldsymbol{x}) = \sum_{k=1}^{H}$	$\int_{-1}^{\infty} g_e(\boldsymbol{x}) \left( \int_{-1}^{\infty} g_e(\boldsymbol{x}) \right) d\boldsymbol{x}$	$\sum egin{array}{c} oldsymbol{W}_{V_i}^{(e)}   ext{ReLU} \end{array}$	

• Integrating MoE into attention we get MoBA.

$$\operatorname{FFN}_{\operatorname{MoE}}(\boldsymbol{x}) = \sum_{e=1}^{E} g_e(\boldsymbol{x}) \left( \sum_{i \in \operatorname{Expert}(e)} \boldsymbol{W}_{V_i}^{(e)} \operatorname{ReLU}(\boldsymbol{W}_{K_i}^{(e)} \boldsymbol{x}) \right).$$
$$\boldsymbol{o}_{t,\operatorname{MoE}} = \sum_{e=1}^{E} g_e(\boldsymbol{q}_t, \{\boldsymbol{k}_i\}_{i \in \operatorname{Expert}(e)}) \left( \sum_{i \in \operatorname{Expert}(e)} \boldsymbol{v}_i \kappa(\boldsymbol{q}_t, \boldsymbol{k}_i) \right).$$

Moba: Mixture of block attention for long-context llms. Arxiv 2025

# Architecture Symmetry

Component	Attention	FFN	
Kernel	Exp	ReLU / SiLU	
Normalization	Yes	None	
Multihead	Yes	None	
Sparsity	None	MoE	
Gating	None	SwiGLU Gating	
		•	

$$g_i(x_{i:t}) = \prod_{j=i+1} \alpha_j(x_j)$$

$$egin{aligned} ext{FFN}_{ ext{SwiGLU}}(oldsymbol{x}) &= \sum_{i=1}^m oldsymbol{\mathcal{W}_{Vi}}_{oldsymbol{v}_i} oldsymbol{\mathcal{W}_{Gi}}^{ op} oldsymbol{x}) oldsymbol{\mathcal{S}wish}_{\kappa(\cdot,\cdot)} oldsymbol{\mathcal{W}_{Ki}}_{oldsymbol{k}_i} oldsymbol{x}, \ oldsymbol{k}_i oldsymbol{\psi}_{oldsymbol{t}_i}(oldsymbol{x}), \ oldsymbol{v}_{t, ext{gated}} &= \sum_{i=1}^t oldsymbol{v}_i \, g_i(oldsymbol{x}_{i:t}) \, \kappa(oldsymbol{q}_t, oldsymbol{k}_i), \end{aligned}$$

Forgetting transformer: Softmax attention with a forget gate. ICLR 2025

### Part 2. Update the Associative Memory

• Taking the linear model as an example,

• 
$$S_{t} = \underbrace{\sum_{i=1}^{t-1} v_{i} k_{i}^{\mathsf{T}}}_{S_{t-1}} + v_{t} k_{t}^{\mathsf{T}}$$
$$= \underbrace{S_{t-1}}_{S_{t-1}} - \underbrace{\left(-v_{t} k_{t}^{\mathsf{T}}\right)}_{\frac{\partial L_{t}}{\partial S_{t-1}}}$$

• Taking the linear model as an example,

• 
$$S_{t} = \underbrace{\sum_{i=1}^{t-1} v_{i} k_{i}^{\mathsf{T}}}_{S_{t-1}} + v_{t} k_{t}^{\mathsf{T}}$$
$$= S_{t-1} - \underbrace{(-v_{t} k_{t}^{\mathsf{T}})}_{\frac{\partial L_{t}}{\partial S_{t-1}}} \qquad L_{t}(S_{t-1}) = - \langle S_{t-1} k_{t}, v_{t} \rangle$$

• Taking the linear model as an example,

• 
$$S_{t} = \underbrace{\sum_{i=1}^{t-1} v_{i} k_{i}^{\top} + v_{t} k_{t}^{\top}}_{S_{t-1}}$$
$$= \underbrace{S_{t-1} - \underbrace{\left(-v_{t} k_{t}^{\top}\right)}_{\frac{\partial L_{t}}{\partial S_{t-1}}}$$

$$L_t(S_{t-1}) = -\langle S_{t-1}k_t, v_t \rangle$$

The loss is unbounded.

**Solution 1**: regularize  $||S||_F$ 

• 
$$L_t(S_{t-1}) = -\langle S_{t-1}k_t, v_t \rangle + \frac{1}{2} ||\sqrt{1-\lambda_t} S_{t-1}||_F^2$$

**Solution 1**: regularize  $||S||_F$ 

• 
$$L_t(S_{t-1}) = -\langle S_{t-1}k_t, v_t \rangle + \frac{1}{2} ||\sqrt{1-\lambda_t} S_{t-1}||_F^2$$

• Update form:  $S_t = \lambda_t S_{t-1} + v_t k_t^{\top}$ 

This is gated linear model.

**Solution 1**: regularize  $||S||_F$ 

• 
$$L_t(S_{t-1}) = -\langle S_{t-1}k_t, v_t \rangle + \frac{1}{2} ||\sqrt{1-\lambda_t} S_{t-1}||_F^2$$

• Update form:  $S_t = \lambda_t S_{t-1} + v_t k_t^{\top}$ 

This is gated linear model.

• Gating implies a bias that older information is less important.

**Solution 2**: regularize ||*Sk*||

• 
$$L_t(S_{t-1}) = -\langle S_{t-1}k_t, v_t \rangle + \frac{1}{2} ||S_{t-1}k_t||^2$$
  
=  $\frac{1}{2} ||S_{t-1}k_t - v_t||^2 - \frac{1}{2} ||v_t||^2$ 

**Solution 2**: regularize ||*Sk*||

• 
$$L_t(S_{t-1}) = -\langle S_{t-1}k_t, v_t \rangle + \frac{1}{2} ||S_{t-1}k_t||^2$$
  
$$= \frac{1}{2} ||S_{t-1}k_t - v_t||^2 - \frac{1}{2} ||v_t||^2$$
$$:= \frac{1}{2} ||S_{t-1}k_t - v_t||^2$$

**Solution 2**: regularize ||*Sk*||

• 
$$L_t(S_{t-1}) = \frac{1}{2} ||S_{t-1}k_t - v_t||^2$$

• Update form:  $S_t = S_{t-1}(I - k_t k_t^{\mathsf{T}}) + v_t k_t^{\mathsf{T}}$ 

This is DeltaNet (Delta rule update).

**Solution 2**: regularize ||*Sk*||

• 
$$L_t(S_{t-1}) = \frac{1}{2} ||S_{t-1}k_t - v_t||^2$$

• Update form:  $S_t = S_{t-1}(I - k_t k_t^{\mathsf{T}}) + v_t k_t^{\mathsf{T}}$ 

This is DeltaNet (Delta rule update).

• How to understand the magic gate  $(I - k_t k_t^{\mathsf{T}})$ ?

Understand delta rule

• 
$$S_t = S_{t-1}(I - k_t k_t^{\mathsf{T}}) + v_t k_t^{\mathsf{T}} = \sum_{i=1}^{t-1} v_i k_i^{\mathsf{T}}(I - k_t k_t^{\mathsf{T}}) + v_t k_t^{\mathsf{T}}$$
  

$$= S_{t-1} + (v_t - \underbrace{\sum_{i=1}^{t-1} v_i(k_i^{\mathsf{T}} k_t)}_{\text{Info overlap with } v_t})k_t^{\mathsf{T}}$$

• Delta rule is a smarter gate for **erasing** overlapping historical info.

**Solution 3**: normalization (e.g.,  $\sum_{i=1}^{t} \kappa(q, k_i)$ )

• When *t* is sufficiently large:

$$L_t(S_{t-1}) \approx \frac{1}{t} \left( - \langle S_{t-1} k_t, v_t \rangle + \frac{1}{2} ||S_{t-1}||_F^2 \right)$$

• Besides regularizing  $||S||_F$ ,  $\frac{1}{t}$  also suppresses numerical explosion.

#### Part 3. Exploration of Next Generation Models

### What is the essential difference?

#### Solution :

• 1) Decay:  $S_{t+1} = S_t(\alpha_t I) + v_t k_t^{T}$ 

Data independent:RetNet,RWKV4Data dependent:Gated linear attention,Mamba2,RWKV6

• 2) Delete:

$$S_{t+1} = S_t (I - k_t k_t^{\mathsf{T}}) + v_t k_t^{\mathsf{T}}$$

Without decay: DeltaNet With decay: RWKV7, Gated DeltaNet

# Intuitively speaking:

The computation of  $\beta$  can be treated as a **prefix sum**, which allows for an efficient implementation. The computation of u seems to be inherently sequential, requiring an O(t) loop?



# Mathematically speaking:

The computation of  $\beta\,$  and u have different parallel complexity.

Three basic questions for determining Parallel Complexity:

1. What is the length of the critical path?  $O(\log^{k} n) \rightarrow \{TC^{k}, AC^{k}, NC^{k}\}$ 

2. How many fan-in one node can receive?  $2 \rightarrow NC$ unlimit  $\rightarrow \{AC, TC\}$ 

3.*Can we use a Threshold Gate out of OR/AND/NOT?* Yes  $\rightarrow$  TC No  $\rightarrow$  {NC, AC}



# Where is the popular model?



The higher the parallelism, the lower the expressiveness. We list representative tasks of common complexity next to them.

# Important tasks beyond $TC^0$

- *NC*<sup>1</sup>: Python, Chess, Entities...
- Beyond *NC*<sup>1</sup>:Graph connectivity ...



The illusion of state in state-space models. ICML 2024.

# When Transformer learn task beyond $TC^{0}$

 $S_5$ : Track the swap of 5 elements, which is beyond  $TC^0$ 

- Difficult to learn: context length 32 need 16 layers
- Difficult to generalize: accuracy 0 out of distribution



Implicit Language Models are RNNs: Balancing Parallelization and Expressivity. Arxiv 2025.

# DeltaFormer: Make Transformer Beyond *TC*<sup>0</sup>

• Rethink Transformer from delta rule:

Delta rule:  $S_t = S_{t-1}(I - k_t k_t^{\mathsf{T}}) + v_t k_t^{\mathsf{T}}$ 

+Kernel Trick:  $\kappa(k_i, k_j) = \phi(k_i)^T \phi(k_j)$  $S_t = S_{t-1}(I - \phi(k_t)\phi(k_t^T)) + v_t \phi(k_t^T)$ 

### Derivation of DeltaFormer

$$\kappa(k_i, k_j) = \phi(k_i)^T \phi(k_j)$$
$$S_t = S_{t-1}(I - \phi(k_t)\phi(k_t^\top)) + v_t \phi(k_t^\top)$$

By method of undetermined coefficients

General

form

$$u_t = v_t - \sum_{i=1}^{t-1} \kappa(k_t, k_i) u_i$$
$$o_t = \sum_{i=1}^t \kappa(q_t, k_i) u_i$$

$$u_t = \alpha_t v_t - \beta_t \sum_{i=1}^{t-1} \kappa_1(w_t, k_i) u_i$$
$$o_t = \sum_{i=1}^t \kappa_2(q_t, k_i) u_i$$

### Efficient chunk-wise implementation

• Recurrent: 
$$u_t = v_t - \sum_{i=1}^{t-1} \kappa(k_t, k_i) u_i$$

• Parallelize:  $U = V - AU \implies U = (I + A)^{-1}V$ 

• Chunk-wise: 
$$U_C = V_c - A_c U_c - A_{c,p} U_p$$
  
 $\Box > U_C = (I + A_c)^{-1} (V_c - A_{c,p} U_p)$ 

# Efficient chunk-wise implementation

#### **Current Implementation**

- Recurrent (Serial): O(T) loop
- Parallelize(Torch trsv): Too much I/0.
- Chunk-wise(dfpa): O(T/C) loop with  $O(T^2D + TCD + TC^2)$  flops

#### Performance

Expected performance results on H800:

• [ B, H, T, D ] = 2, 32, 8192, 128

Forward:

serial	time	279.908 ms
Torch trsv	time	102.156 ms
dfpa	time	12.705 ms

Backward:

Torch	trsv	backward	time	275.659 ms
dfpa		backward	time	25.713 ms

Reference time consumption in training

∂ ms
ms
ms
ms

# Make Transformer beyond $TC^0$

• More expressive:

#### Can track the exchange of nelements with $d = O(\log n)$ . In the Lemma 2 of rwkv7, they tracking 5 elements used 5 dimensions

• Native compression:

If we read out and rewrite KU cache every O(n) step, the actually KU cache is O(nlogn).

#### Theorem 1 (State Exchange)

**Assumption 1:** There exist n state points on a d-dimensional unit sphere, and the absolute value of the inner product of any two distinct state points is less than or equal to  $\epsilon(d, n)$ , which means:

$$\|m{x}_1,m{x}_2,\dots,m{x}_n\in\mathbb{R}^d\quad s.t.\quad \|m{x}_i\|_2=1(orall i),\quad \max_{i
eq i}|m{x}_i^{ op}m{x}_j|\leq\epsilon(d,n)<rac{1}{8}.$$

Assumption 2: There is a function f satisfies:

$$\forall \boldsymbol{x} \in \{-1, 0, 1, 2\}, \forall \tilde{\boldsymbol{x}} \in U(\boldsymbol{x}, 4\epsilon(d, n)): \quad f(\tilde{\boldsymbol{x}}) = \boldsymbol{x}.$$

Consider initializing n key-value pairs as  $\{(k_1, v_1), \ldots, (k_n, v_n)\}$ . The keys  $\{k_1, \ldots, k_n\}$  lie on a d-dimensional unit sphere and satisfies Assumption 1, which means:

 $\forall i, j \in \{1, \dots, n\}, i \neq j: \quad \|\boldsymbol{k}_i\|_2 = 1, \quad |\boldsymbol{k}_i^\top \boldsymbol{k}_j| \leq \epsilon(n, d).$ 

Define an attention mechanism as follows:

$$oldsymbol{v}_t = oldsymbol{v}_t - \sum_{i=1}^{t-1} f(oldsymbol{k}_i^{ op} oldsymbol{k}_t) oldsymbol{u}_i, \quad oldsymbol{o}_t = \sum_{i=1}^t f(oldsymbol{q}_t^{ op} oldsymbol{k}_i) oldsymbol{u}_i,$$

where  $f(\cdot)$  satisfies Assumption 2 and it is noted that  $\forall i \in \{1, ..., n\}$ , since  $f(\mathbf{k}_i^\top \mathbf{k}_i) = 0$ , we have  $\mathbf{u}_i = \mathbf{v}_i$ .

At the current step t, t > n, the value corresponding to  $\mathbf{k}_i$  is denoted by  $\tilde{\mathbf{v}}_i, i \in \{1, \ldots, n\}$ . Note that, after t - 1 - n exchanges,  $\tilde{\mathbf{v}}_i$  is not necessarily equal to the initially assigned  $\mathbf{v}_i$ .  $\forall 1 \le t_2 < t_1 \le n$ , to exchange the stored values  $\tilde{\mathbf{v}}_{t_1}$  and  $\tilde{\mathbf{v}}_{t_2}$  corresponding to  $\mathbf{k}_{t_1}$  and  $\mathbf{k}_{t_2}$ , it suffices to construct:

$$\boldsymbol{k}_t = \boldsymbol{k}_{t_1} - \boldsymbol{k}_{t_2}, \quad \boldsymbol{v}_t = 0$$

When retrieving the values:

• Query  $\boldsymbol{q}_t = \boldsymbol{k}_{t_1}$ , then  $\boldsymbol{o}_t = \tilde{\boldsymbol{v}}_{t_2}$ ;

• Query 
$$\boldsymbol{q}_t = \boldsymbol{k}_{t_2}$$
, then  $\boldsymbol{o}_t = \tilde{\boldsymbol{v}}_{t_1}$ ;

• Query  $q_t = k_{t_3}, 1 \le t_3 \le n, t_3 \ne t_1, t_2, then o_t = \tilde{v}_{t_3}.$ 

This implies the exchange of values corresponding to  $k_{t_1}$  and  $k_{t_2}$  is completed.

- DeltaFormer with any kernel surpass Transformer in tracking.
- Accuracy: 8 layers Transformer < 0.50,

1 layers DeltaFormer can reach 1.00.



- Nonlinear kernel has more capacity than linear.
- Track n elements with dim d < n:</li>
   linear → drop much , nonlinear > 0.95.



**Figure 5** Comparison of DeltaFormer with  $\kappa(x, y) = x^{\top}y$  and  $\kappa(x, y) = \lfloor x^{\top}y \rfloor$ .

- GQA like method enhance nonlinear kernel with the same KU cache
- Track 5 elements with a kv head with dim 3, but we have more "query" w.



**Figure 8** Comparison of DeltaFormer with different  $\kappa_1$  and  $\kappa_2$  under different query heads numbers.

- Curriculum learning: when accuracy reach 0.99, context length × 2
- With curriculum learning: RoPE can NOT reach 1.00, NoPE can reach 1.00.
- Without curriculum learning: RoPE can learn, NoPE can NOT learn.



**Figure 6** Comparison of DeltaFormer using different learning strategy and position embedding. Each use  $\kappa_1(x, y) = \lfloor x^\top y \rfloor$ . "32 to 256" means that the initial training length is 32, which means the number of swaps is 32. When the accuracy reaches 0.99, the training length will be doubled until it reaches 256. Each color gradient in the image represents a doubling of the training length. And "256" means that the model is trained on a training length of 256 from the beginning. The y-axis reflects the accuracy at the current training length.

### Learn from data – Reachability of DAG

• The reachability of directed acyclic graphs: Transformer get 0.8, DeltaFormer can reach 1.00.



**Figure 7** Comparison of Transformer and DeltaFormer using different similarity functions  $\kappa_1(\cdot)$  for performing swapping tasks. For  $\kappa_2(\cdot)$ , we use the softmax function to maintain consistency with Transformer. Pay attention to the scale of the y-axis.

# Learn from data – Dyck grammar

• Example: "((()())(()))"



Context length = 64, next token prediction, compress Dyck grammar



Better compression ratio

 $u_t = v_t + u_{t-1}$  Can record whether there are more left parentheses than right parentheses at present.

# Learn from data – From Dyck to real code

• 14B total parameters model, 500B training token



Total loss: DeltaFormer lead baseline by 0.005 Code loss: DeltaFormer lead baseline by 0.03

This implies that the model's ability to compress more complex knowledge exceeds that of the Transformer.

# Learn from data – Induction heads

• Better learn induction heads: "... AB... A" predict "B"



The calculation of u can to some extent replace the function of copy heads, thereby promoting the formation of induction heads

## DeltaFormer – Future work

- **Optimization**. -- Stability training and scaling ...
- Fine grained design. --GQA / Different kernel / Head
   Nums / Head Dims / PE / Sparse...
- Implementation. -- Fast..
- **Practical**. -- Expressivity/ KU cache/ Flops Trade off..
- More expressive but also practical model.



# Thanks for listening!

# Q & A