

# Theta-CS：コーシー=シュワルツ不等式による前段スクリーニング

高槻高等学校 宮崎俊輔

## 1. 導入

- [2] TransformerのAttention機構は全トークン間の関係を内積計算で求めるため高精度だが、計算コストが膨大である。
- [1] Top- $\theta$  は各headに静的閾値を設けて不要なAttention Scoreを刈ることで、推論時の計算量を削減する手法である。さらに補正によって精度を維持しつつ訓練を要さない疎化を実現する。本研究では、コーシー=シュワルツ不等式を用いた安全な前段ふるいを導入し、閾値疎化の前に不要な内積計算を削減する。

## 2. 方法論

2.1,2.2は本研究の前提となるアーキテクチャである。

### 2.1 Attention機構<sup>[2]</sup>

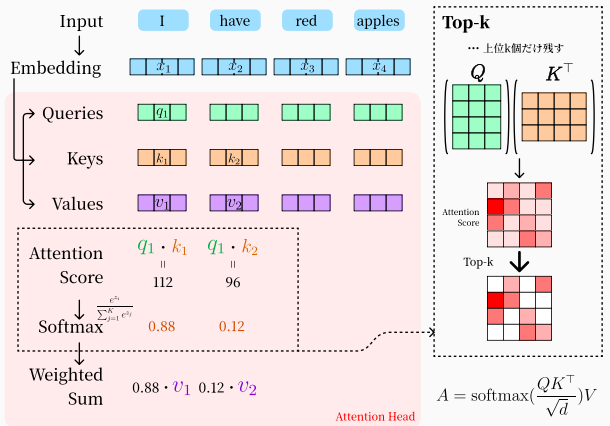


図1: Attention機構とTop-kの概略図

### 2.2 Top- $\theta$ <sup>[1]</sup>

… 事前計算されたk個の要素を保持する静的閾値 $\theta$ を使用することで、Top-kとほぼ同一の集合を保持する。

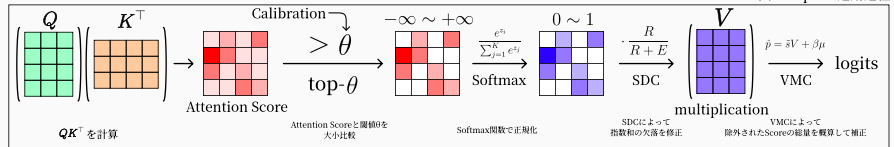


図2: Top- $\theta$ の適用過程

・SDC (Softmax Denominator Compensation)  
softmax関数の分母から省略された“指数の合計”を修正する

$$\text{softmax}(a)_i = \text{softmax}(\tilde{a})_i \cdot \frac{R}{R + E}$$

・VMC (V-Mean Compensation)  
除外されたScoreの総量を概算し、不足分を埋めて補正する

$$\mu = \frac{1}{n} \sum_{i=0}^{n-1} V_i, \quad \beta = 1 - \sum_{i=0}^{k-1} \tilde{s}_i \rightarrow \hat{p} = \tilde{s}V + \beta\mu$$

- Good: 訓練不要。短時間のしきい値校正だけで使える  
計算量とKVキャッシュの読み取りを削減可能
- Bad: 除外されたScoreの内積計算が無駄になる

Algorithm 1 Calibration( $\mathcal{C}, k, n$ ): 1-head threshold  
Require:  $\mathcal{C}$  (Calibration set of inputs)  
Require:  $k \in \mathbb{N}$  (elements to keep per attention row)  
Require:  $n \in \mathbb{R}$  (calibration offset in std\_devs)  
1:  $\Theta_k = \emptyset$  (empty sets of observed thresholds)  
2: for  $X \in \mathcal{C}$  do  
3: if  $X$ .prefill(X) then  
4:  $A = Q(X)K(X)^T$   
5:  $\mu = \text{mean}_i(\text{row}(A))$   
6: for  $r = k$  to  $n-1$  do  
7:  $\Theta_k \leftarrow \Theta_k \cup \{\text{quantile}_{1/2k}(A_{r,:})\}$   
8:  $A_k = \text{top}_k(A_k)$   
9: end for  
10:  $S = \text{row\_softmax}(A)$   
11:  $\text{chae}$  [generative decoding,  $X \in \mathbb{R}^d$ ]  
12:  $a = Q(X)K(X)^T(X)$   
13:  $n = \text{length}(a)$   
14:  $\Theta_k \leftarrow \Theta_k \cup \{\text{quantile}_{1/2k}(a)\}$   
15:  $\Theta_k = \text{top}_k(\Theta_k)$   
16:  $s = \text{softmax}(a)$   
17: end if  
18: end for  
19: return  $\Theta_k = \text{mean}(\Theta_k) + \alpha \cdot \text{std\_dev}(\Theta_k), \forall r$

図3: 閾値 $\theta$ を求める擬似コード

### 2.3 Theta-CS

… QK内積を計算する前にコーシー=シュワルツ不等式に従って明らかに小さい”あり得ない候補”を排除する

コーシー=シュワルツ不等式

$$|\langle x, y \rangle| \leq \|x\|_2 \|y\|_2$$

行列にして、

$$\|QK^T\|_2 \leq \|Q\|_2 \|K\|_2$$

ゆえに、

$$\|Q\|_2 \|K\|_2 \leq \theta \Rightarrow \|QK^T\|_2 \leq \theta$$

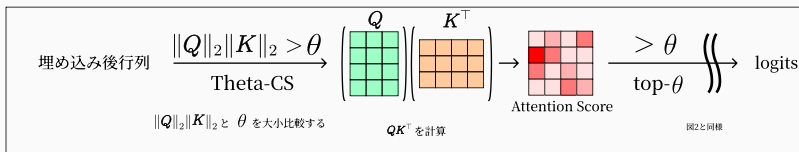


図4: Theta-CSの適用過程

- Good: 内積計算の前に候補を剪定することで計算量を削減  
実装が容易  
上限値に基づいてスクリーニングするため安全

ノルム

$$\vec{x} = (x_1, x_2, \dots, x_n)$$

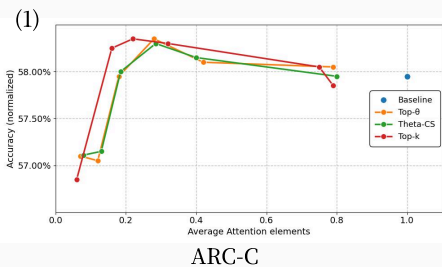
$$\|\vec{x}\|_p = \sqrt[p]{|x_1|^p + |x_2|^p + \dots + |x_n|^p}$$

Algorithm 2 Theta-CS PreFilter (pre-softmax, 1-head) with Top- $\theta$  thresholds  
Require:  $\{\theta_k\}_{k=1}^n$  (pre-softmax thresholds from Top- $\theta$  calibration)  
Require:  $\delta \in \mathbb{R}_{\geq 0}$  (safety margin)  
1:  $k\_norms \leftarrow$  empty vector (cache of  $\|k_i\|$  aligned with KV cache)  
2: for  $X \in \mathcal{C}$  or streaming inputs do  
3: if is\_prefill(X) then  
4:  $\{Q, K, V\} \leftarrow \text{project}(X)$  (shape:  $[n, d]$  each)  
5:  $k\_norms \leftarrow \|K\|_{\text{row}} \leftarrow 2$  (one-time compute for current prompt; append if continuing)  
6:  $n = \text{rows}(Q)$   
7: for  $r = 1$  to  $n$  do  
8:  $q \leftarrow Q_{r,:}$ ,  $\tau \leftarrow \theta - \|q\|$   
9:  $\theta \leftarrow \theta_r$ ,  $\tau \leftarrow \theta - \delta$   
10:  $m_r \leftarrow (q \cdot n - k\_norms \geq \tau)$  (Cauchy-Schwarz screen)  
11: if allfalse( $m_r$ ) then  
12:  $m_r[\text{arg max}(k\_norms)] \leftarrow \text{true}$   
13: end if  
14:  $K^{(\text{keep})} \leftarrow K[m_r,:]$ ,  $V^{(\text{keep})} \leftarrow V[m_r,:]$   
15:  $a^{(\text{keep})} \leftarrow K^{(\text{keep})} q^T$  (compute dot only on kept keys)  
16:  $a^{(\theta)} \leftarrow \text{threshold\_pre}(a^{(\text{keep})}, \theta)$  (Top- $\theta$  pre-softmax step)  
17:  $(\hat{s}, \text{SDC\_state}) \leftarrow \text{apply\_SDC\_and\_softmax}(a^{(\theta)}, \theta)$   
18:  $\hat{y}_r \leftarrow \text{apply\_VMC}(\hat{s}, V^{(\text{keep})}, \text{SDC\_state})$   
19: end for  
20: else  
21: {generative decoding: new token  $x \in \mathbb{R}^d$ }  
22:  $(q, k, v) \leftarrow \text{project}(x)$   
23: append  $k, v$  to KV cache; append  $\|k\|$  to  $k\_norms$   
24:  $q\_n \leftarrow \|q\|$ ,  $n \leftarrow \text{length}(k\_norms)$   
25:  $\theta \leftarrow \theta_n$ ,  $\tau \leftarrow \theta - \delta$   
26:  $m \leftarrow (q \cdot n - k\_norms \geq \tau)$   
27: if allfalse( $m$ ) then  
28:  $m[\text{arg max}(k\_norms)] \leftarrow \text{true}$   
29: end if  
30:  $K^{(\text{keep})} \leftarrow K[m,:]$ ,  $V^{(\text{keep})} \leftarrow V[m,:]$   
31:  $a^{(\text{keep})} \leftarrow K^{(\text{keep})} q^T$   
32:  $a^{(\theta)} \leftarrow \text{threshold\_pre}(a^{(\text{keep})}, \theta)$   
33:  $(\hat{s}, \text{SDC\_state}) \leftarrow \text{apply\_SDC\_and\_softmax}(a^{(\theta)}, \theta)$   
34:  $\hat{y} \leftarrow \text{apply\_VMC}(\hat{s}, V^{(\text{keep})}, \text{SDC\_state})$   
35: end if  
36: return  $\{\hat{y}_r\}$  or  $\hat{y}$  (same interface as Top- $\theta$  after pre-softmax)

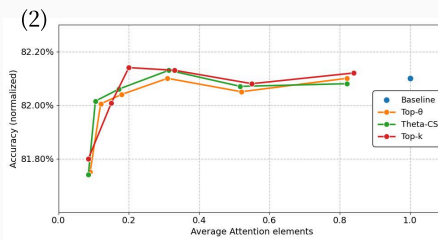
図5:  $\theta$ -CSの擬似コード

## 3. 評価実験

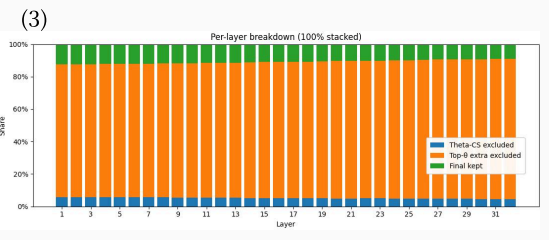
設定: モデル: LLaMA3-8B-Instruct  
ベンチマーク: ARC-C(prefill), HellaSwag(prefill), HumanEval(generative)



ARC-C



HellaSwag



Per-layer breakdown (HumanEval)

## 4. 考察

- (1)(2)から、  
・Theta-CS と Top- $\theta$  の精度がほとんど重なっていて安全性がある。

- (3)から、  
・先行除外率は5%程度であり、 $L^2 \times L^2$  のため高次元では上界が緩くなりすぎるためだろう。

## 5. 今後の展望

- ・今回はコーシー=シュワルツ不等式に則り  $L^2 \times L^2$  で上界を設定したが、下にあるようにコーシー=シュワルツ不等式を一般化したヘルダー不等式などまだまだ発展の余地があり、今後も数学的視点からの計算量削減に関する検証を進める。

ヘルダー不等式

$$1 \leq p, q \leq \infty, \quad 1/p + 1/q = 1$$

$$|\langle x, y \rangle| \leq \|x\|_p \|y\|_q$$

→

$p = 1, q = \infty$  で、

$$\|QK^T\|_1 \leq \|Q\|_1 \|K\|_\infty$$

ゆえに、

$$\|Q\|_1 \|K\|_\infty \leq \theta \Rightarrow \|QK^T\|_1 \leq \theta$$

- [1] "Top-Theta Attention" (Benzo Andrei et al., 2025)  
[2] "Attention Is All You Need" (Ashish Vaswani et al., 2017)  
[3] https://github.com/ashishvaswani/transformer  
[4] https://mathlandscap.com/holder/  
[5] https://mathlandscap.com/cauchy-schwarz/  
[6] https://mathlandscap.com/math-640/