

# Hungry Hungry Hippos: Towards Language Modeling with State Space Models

Daniel Y Fu, Tri Dao, Khaled Kamal Saab, Armin W Thomas, Atri Rudra, Christopher Re

第15回最先端NLP勉強会  
豊田工業大学  
牧野 晃平



<https://www.amazon.co.jp/Hasbro-Hungry-Hippos-%E4%B8%A6%E8%A1%8C%E8%BC%B8%E5%85%A5%E5%93%81/dp/B01M24DD0K>

H3

# Hungry Hungry Hippos: Towards Language Modeling with State Space Models

Daniel Y Fu, Tri Dao, Khaled Kamal Saab, Armin W Thomas, Atri Rudra, Christopher Re

第15回最先端NLP勉強会  
豊田工業大学  
牧野 晃平



<https://www.amazon.co.jp/Hasbro-Hungry-Hippos-%E4%B8%A6%E8%A1%8C%E8%BC%B8%E5%85%A5%E5%93%81/dp/B01M24DD0K>

# 論文の位置づけ

状態空間モデル  
[Brogan+74]

$$\dot{x} = Ax + Bu$$

$$y = Cx + Du$$

深層学習に導入

設計された  
A行列

HiPPO [Gu+20]

$$A_{nk} = - \begin{cases} (2n+1)^{1/2}(2k+1)^{1/2} & \text{if } n > k \\ n+1 & \text{if } n = k \\ 0 & \text{if } n < k \end{cases}$$

初期化に利用

## 系列を扱える深層学習モデル

S4 [Gu+22]

Continuous State Space

$$\dot{x} = Ax + Bu$$

$$y = Cx + Du$$

Long-Range Dependencies

$$A = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 2 & 0 \\ 1 & 3 & 3 \end{bmatrix}$$

Fast Discrete Representations

$$x = \bar{A}x + \bar{B}u$$

$$y = \bar{C}x + \bar{D}u$$

$$y = \bar{K} * u$$

言語に特化

計算速度: S4 > Att  
CV: S4  $\doteq$  Att  
音声: S4  $\doteq$  Att  
時系列: S4  $\doteq$  Att  
Att  
NLP: S4 < Att

H3 [Fu+23]

Diag SSM

Shift SSM

Q K V

X

インスピレーション

Linear Attention [Katharopoulos+20]

$$O_i = \frac{\phi(Q_i)^\top \sum_{j=1}^i \phi(K_j) V_j^\top}{\phi(Q_i)^\top \sum_{j=1}^i \phi(K_j)}$$

variant

CNN [LeCun+18]

INPUT 32x32

C1 feature maps 48x48x3

S2 16x16x3

S4 8x8x3

C3 16x16x3

S4 8x8x3

C4 16x16x3

S4 8x8x3

C5 16x16x3

S4 8x8x3

Full connection

Gaussian connections

RNN

y

h

x

Attention [Vaswani+19]

Positional Encoding

Input Embedding

Inputs

Positional Encoding

Output Embedding

Outputs (shifted right)

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Multi-Head Attention

Masked Multi-Head Attention

Nx

# 論文の位置づけ

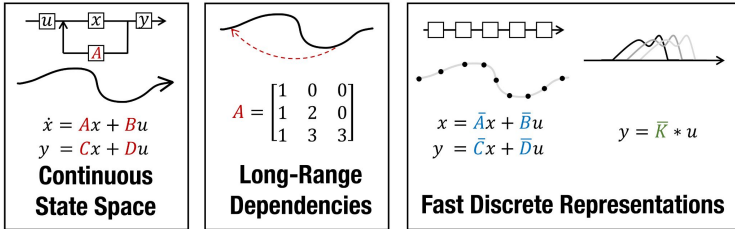
## 系列を扱える深層学習モデル

状態空間モデル  
[Brogan+74]

$$\begin{aligned}\dot{x} &= Ax + Bu \\ y &= Cx + Du\end{aligned}$$

深層学習に導入

S4 [Gu+22]



言語に特化

計算速度: S4 > Att

CV: S4 ≐

Att

音声: S4 ≐ Att

時系列: S4 ≐

Att

NLP: S4 < Att

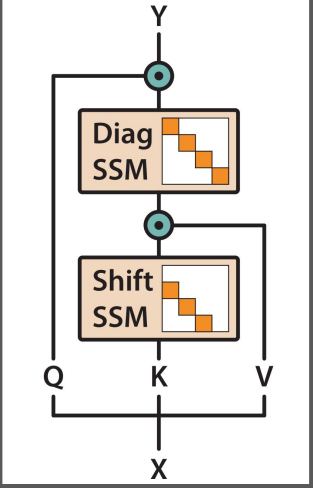
設計された  
A行列

初期化に利用

HiPPO [Gu+20]

$$A_{nk} = - \begin{cases} (2n+1)^{1/2}(2k+1)^{1/2} & \text{if } n > k \\ n+1 & \text{if } n = k \\ 0 & \text{if } n < k \end{cases}$$

H3 [Fu+23]



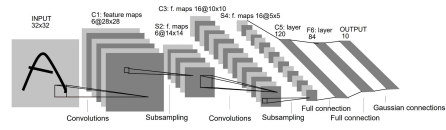
インスピレーション

variant

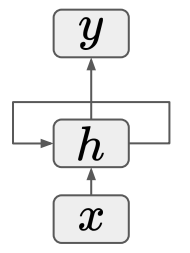
Linear Attention  
[Katharopoulos+20]

$$O_i = \frac{\phi(Q_i)^\top \sum_{j=1}^i \phi(K_j) V_j^\top}{\phi(Q_i)^\top \sum_{j=1}^i \phi(K_j)}$$

CNN [LeCun+18]

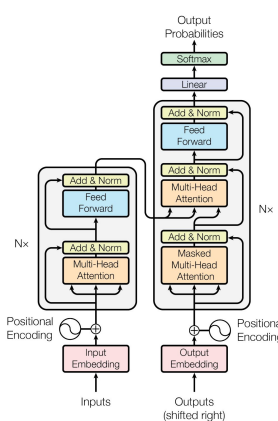


RNN



Attention

[Vaswani+19]



# 状態空間モデル(State Space Model; SSM)

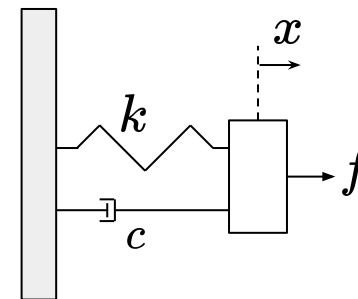
中間状態が必要な系を記述するためのモデル

- 時系列解析や制御システムによく利用される

$$\begin{aligned}\dot{x} &= Ax + Bu & u &: \text{入力} \\ y &= Cx + Du & x &: \text{中間状態} \\ & & y &: \text{出力}\end{aligned}$$

例: 質点—バネ—ダンパー系

$$\begin{aligned}\dot{\mathbf{x}} = \begin{bmatrix} \dot{x} \\ \ddot{x} \end{bmatrix} &= \begin{bmatrix} 0 & 1 \\ -c & -k \end{bmatrix} \begin{bmatrix} x \\ \dot{x} \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} f \\ y = x &= \begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} x \\ \dot{x} \end{bmatrix}\end{aligned}$$



離散的な系列を扱えるようにSSMを離散化すると:

$$\begin{aligned}x_i &= \bar{A}x_{i-1} + \bar{B}u_i & \bar{A} &= (I - \Delta/2 \cdot A)^{-1}(I + \Delta/2 \cdot A) \\ y_i &= \bar{C}x_i + \bar{D}u_i & \bar{B} &= (I - \Delta/2 \cdot A)^{-1}\Delta B \quad \bar{C} = C \quad \bar{D} = D\end{aligned}$$

## S4 [Gu+22]: SSMの深層学習への導入

行列をパラメタにして離散化したSSMをニューラルネットに導入

- S4の貢献: 深層学習で使うために訓練を早くした

- SSMを計算機で並列に扱うための畳み込み形式 [Gu+21]:

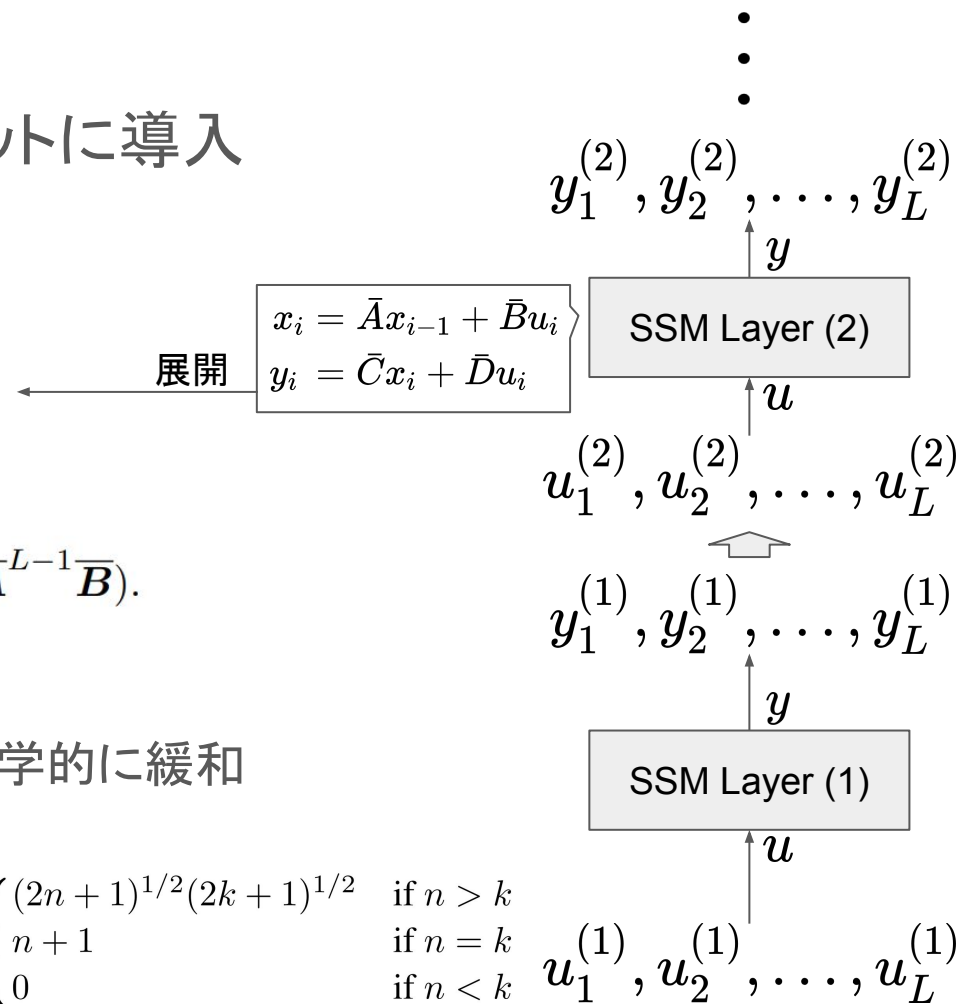
$$y_k = \overline{CA}^k \overline{B}u_0 + \overline{CA}^{k-1} \overline{B}u_1 + \cdots + \overline{CAB}u_{k-1} + \overline{CB}u_k$$
$$y = \overline{K} * u.$$

$$\overline{K} \in \mathbb{R}^L := \mathcal{K}_L(\overline{A}, \overline{B}, \overline{C}) := \left( \overline{CA}^i \overline{B} \right)_{i \in [L]} = (\overline{CB}, \overline{CAB}, \dots, \overline{CA}^{L-1} \overline{B}).$$

- カーネルの  $\overline{A}^k$  の計算が重たい  $\Rightarrow$  高速に計算可能にする
  - カーネルの計算で訓練に時間がかかってしまうのを数学的に緩和
- HiPPO [Gu+20] 行列をA行列の初期値にして性能向上
  - 記憶を保持できるように設計されたA行列

$\Rightarrow$  Attention weightの計算がない

$\Rightarrow$  系列長に対してほぼ線形の計算量で処理できる



$$A_{nk} = - \begin{cases} (2n+1)^{1/2}(2k+1)^{1/2} & \text{if } n > k \\ n+1 & \text{if } n = k \\ 0 & \text{if } n < k \end{cases}$$
$$u_1^{(1)}, u_2^{(1)}, \dots, u_L^{(1)}$$

# 論文の位置づけ

状態空間モデル  
[Brogan+74]

$$\dot{x} = Ax + Bu$$

$$y = Cx + Du$$

深層学習に導入

設計された  
A行列

HiPPO [Gu+20]

$$A_{nk} = - \begin{cases} (2n+1)^{1/2}(2k+1)^{1/2} & \text{if } n > k \\ n+1 & \text{if } n = k \\ 0 & \text{if } n < k \end{cases}$$

初期化に利用

## 系列を扱える深層学習モデル

S4 [Gu+22]

Continuous State Space

$$\dot{x} = Ax + Bu$$

$$y = Cx + Du$$

Long-Range Dependencies

$$A = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 2 & 0 \\ 1 & 3 & 3 \end{bmatrix}$$

Fast Discrete Representations

$$x = \bar{A}x + \bar{B}u$$

$$y = \bar{C}x + \bar{D}u$$

$$y = \bar{K} * u$$

言語に特化

計算速度: S4 > Att  
CV: S4  $\doteq$  Att  
音声: S4  $\doteq$  Att  
時系列: S4  $\doteq$  Att  
NLP: S4 < Att

H3 [Fu+23]

Diag SSM

Shift SSM

Q K V

X

インスピレーション

Linear Attention [Katharopoulos+20]

$$O_i = \frac{\phi(Q_i)^\top \sum_{j=1}^i \phi(K_j) V_j^\top}{\phi(Q_i)^\top \sum_{j=1}^i \phi(K_j)}$$

variant

CNN [LeCun+18]

INPUT 32x32

C1 feature maps 48x48x3

S2 16x16x3

S4 16x16x3

C3 16x16x3

S4 16x16x3

C5 16x16x3

S4 16x16x3

Full connection

Gaussian connections

RNN

y

h

x

Attention [Vaswani+19]

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

Output Embedding

Outputs (shifted right)

# 背景: SSMとAttention

SSMとAttentionを比較してSSMに不足している部分を調査・改善したい

- 性能面での比較:

- CV・音声・時系列の性能:  $SSM \div Attention$
- NLPの性能:  $SSM < Attention$

} 既存のSSMにできないことがあるはず

⇒ 人工言語を使ったタスクで性能調査をして, 問題点の対策をする

- 計算速度面での比較:

- 系列長に対する計算量:  $SSM(\text{ほぼ線形}) \ll Attention(\text{二乗})$
- 実際の計算時間:  $SSM < Attention$

} 計算機的な問題があるはず

⇒ 計算機に寄り添った計算方法にする



# 事前調査：人工言語を利用したSSMの性能調査

二種類の人工的なタスクでSSMができないことを見つける

- タスク: In-context learningの能力を調査するためのタスクを使用

Task	Input	Output	Sequence Length	Vocab Size
Induction Head	<i>a b c d e f g h i ... x y z</i>	<i>f</i>	30	20
Associative Recall	<i>a 2 c 4 b 3 d 1 a</i>	2	20	10

[Olson+22]・[Ba+16]  
を参考に作成

- 評価: 2層のモデルでタスクを解いた時の正解率で評価

- Attentionでできることが既存のSSMではできない

Task	Random	S4D	Gated State Spaces	Attention
Induction Head	5.0	35.6	6.8	<b>100.0</b>
Associative Recall	25.0	86.0	78.0	<b>100.0</b>

- この調査でわかったこと: SSMで何ができないか

- あるイベント後のトークンを記憶出来ない      ∵ 入力をそのままコピーする機構がない
- トークン同士を比較できない      ∵ トークン同士を比べる計算をしていない

※ Attentionでは入力中のトークンをコピー可 & Attention weightの計算でトークンの比較が可能

# 提案手法: H3 layer

二つのSSMで入力のコピーとトークン同士の比較を実現

- $\text{SSM}_{\text{diag}} \cdot \text{SSM}_{\text{shift}}$ : A行列をそれぞれ対角行列とシフト行列としたSSM

※ S4で扱う行列の範囲にシフト行列は含まれない

- $\text{SSM}_{\text{shift}}(\mathbf{K}) \odot \mathbf{V}$ : シフトした系列と元の系列の要素積

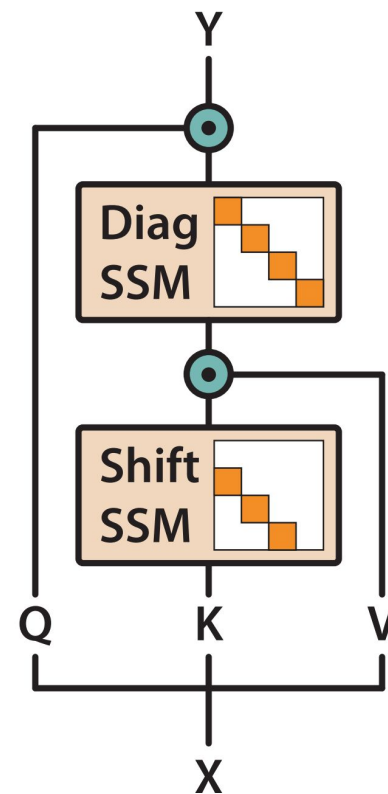
- 過去の特徴を集約した表現が作れる: 直感的には以下のような項が作れる

$$\left[ 0, x_1^{\text{key}}, \dots, \sum_{k=1}^l x_k^{\text{key}}, \dots, \sum_{k=1}^L x_L^{\text{key}} \right] \odot [x_1^{\text{val}}, x_2^{\text{val}}, \dots, x_l^{\text{val}}, \dots, x_L^{\text{val}}]$$

⇒ Attention weightのような気持ちで過去のトークンと比較

- $\mathbf{Q} \odot \text{SSM}_{\text{diag}}(\cdot)$ : 過去をまとめた表現を現在の表現と比較

⇒ 入力をコピー出来つつトークンの比較を実現



$$\mathbf{Q} \odot \text{SSM}_{\text{diag}}(\text{SSM}_{\text{shift}}(\mathbf{K}) \odot \mathbf{V})$$

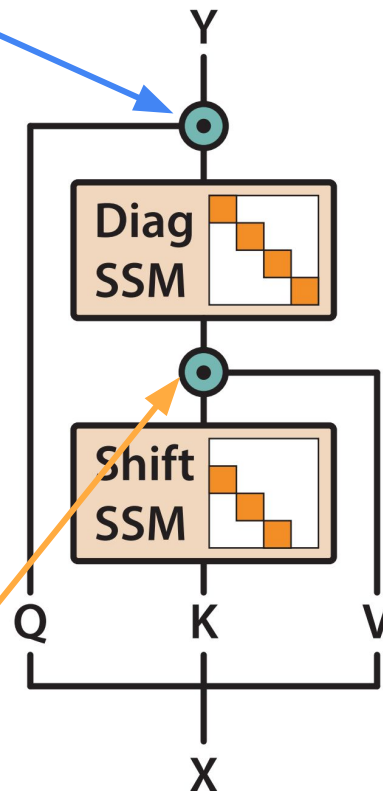
# 提案手法: 正確なH3 layer

## Algorithm 1 H3 Layer

**Require:** Input sequence  $u \in \mathbb{R}^{N \times d}$  from the previous layer, weight matrices  $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V, \mathbf{W}_O \in \mathbb{R}^{d \times d}$ , a shift SSM  $\text{SSM}_{\text{shift}}$ , a diagonal SSM  $\text{SSM}_{\text{diag}}$ , head dimension  $d_h$ .

- 1: Compute  $\mathbf{Q} = u\mathbf{W}_Q, \mathbf{K} = u\mathbf{W}_K, \mathbf{V} = u\mathbf{W}_V \in \mathbb{R}^{N \times d}$ .
- 2: Pass  $\mathbf{K}$  through the shift SSM:  $\bar{\mathbf{K}} = \text{SSM}_{\text{shift}}(\mathbf{K}) \in \mathbb{R}^{N \times d}$ .
- 3: Split  $\mathbf{Q}, \bar{\mathbf{K}}, \mathbf{V}$  into  $H$  “heads”  $(\mathbf{Q}^{(h)}, \bar{\mathbf{K}}^{(h)}, \mathbf{V}^{(h)})$  for  $h = 1, \dots, H$ , each a sequence of  $N$  vectors of size  $d_h = d/H$ .
- 4: **for**  $1 \leq h \leq H$  **do**
- 5:   Take the batched outer product  $\bar{\mathbf{K}}^{(h)}(\mathbf{V}^{(h)})^\top \in \mathbb{R}^{N \times d_h \times d_h}$  (batched in the  $N$ -dimension) and pass it through a diagonal SSM:  $\mathbf{KV}^{(h)} = \text{SSM}_{\text{diag}}(\bar{\mathbf{K}}^{(h)}(\mathbf{V}^{(h)})^\top) \in \mathbb{R}^{N \times d_h \times d_h}$ .
- 6:   Batch-multiply by  $\mathbf{Q}$ :  $\mathbf{O}^{(h)} = [\mathbf{Q}_1^{(h)} \mathbf{KV}_1^{(h)}, \dots, \mathbf{Q}_N^{(h)} \mathbf{KV}_N^{(h)}] \in \mathbb{R}^{N \times d_h}$  (batched in the  $N$ -dimension).
- 7: **end for**
- 8: Concatenate the output  $\mathbf{O}^{(h)}$  of each head, and multiply by the output projection matrix  $\mathbf{W}_O \in \mathbb{R}^{d \times d}$ .

行列積



外積

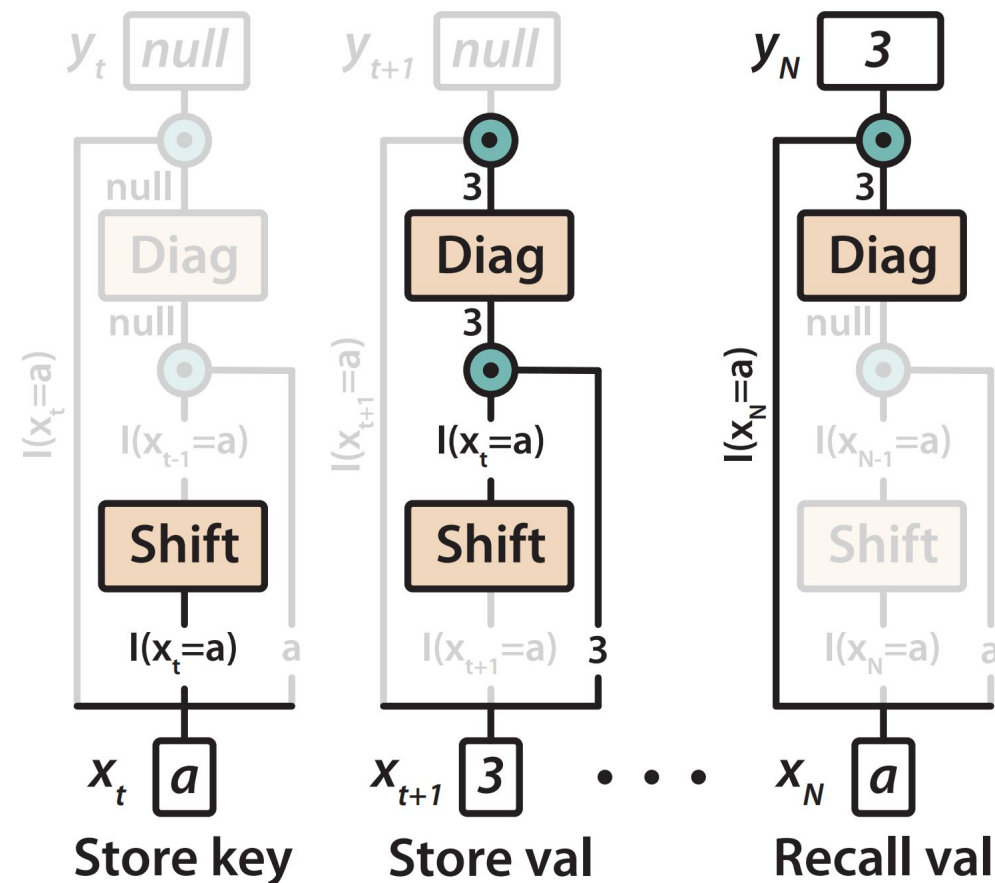
$$\mathbf{Q} \odot \text{SSM}_{\text{diag}}(\text{SSM}_{\text{shift}}(\mathbf{K}) \odot \mathbf{V})$$

## 提案手法: H3 layerの動作例

- Store key: シフトした系列を作成
- Store val: シフトされた系列と元の系列を比較

⇒  $a$  と  $3$  の積が計算できて, ペアとして扱える

- Recall val: 入力中から  $a$  に対応するトークンを参照する



...  $a3$  ...  $a$ ?

# 論文の位置づけ

状態空間モデル  
[Brogan+74]

$$\dot{x} = Ax + Bu$$
$$y = Cx + Du$$

深層学習に導入

設計された  
A行列

HiPPO [Gu+20]

$$A_{nk} = - \begin{cases} (2n+1)^{1/2}(2k+1)^{1/2} & \text{if } n > k \\ n+1 & \text{if } n = k \\ 0 & \text{if } n < k \end{cases}$$

初期化に利用

## 系列を扱える深層学習モデル

S4 [Gu+22]

$\dot{x} = Ax + Bu$   
 $y = Cx + Du$   
Continuous State Space

$A = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 2 & 0 \\ 1 & 3 & 3 \end{bmatrix}$   
Long-Range Dependencies

$x = \bar{A}x + \bar{B}u$   
 $y = \bar{C}x + \bar{D}u$   
Fast Discrete Representations

言語に特化

H3 [Fu+23]

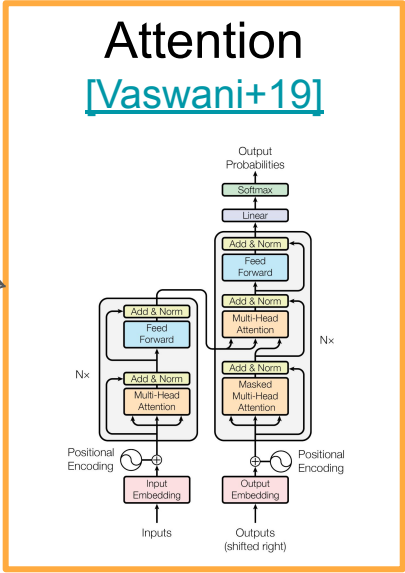
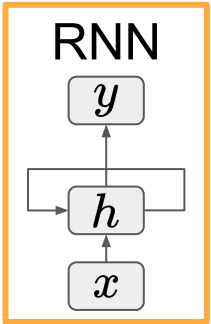
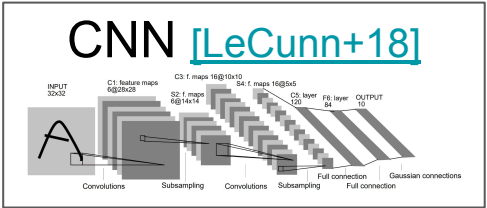
計算速度: S4 > Att  
CV: S4  $\doteq$  Att  
音声: S4  $\doteq$  Att  
時系列: S4  $\doteq$  Att  
Att  
NLP: S4 < Att

インスピレーション

Linear Attention [Katharopoulos+20]

$$O_i = \frac{\phi(Q_i)^\top \sum_{j=1}^i \phi(K_j) V_j^\top}{\phi(Q_i)^\top \sum_{j=1}^i \phi(K_j)}$$

variant



# 提案手法: H3 layerとAttentionの比較

- Attentionの定式化:

- 一般的なAttention [\[Vaswani+19\]](#):  $\text{Sim}(q, k) = e^{q^\top k}$

- Linear Attention:  $\text{Sim}(q, k) = \phi(q)^\top \phi(k)$

- Linear AttentionのRNN標記:

$$z_i = \sum_{j=1}^i \phi(K_j) \in \mathbb{R}^d \quad S_i = \sum_{j=1}^i \phi(K_j) V_j^\top \in \mathbb{R}^{d \times d}$$

$$O_i = \frac{\sum_{j=1}^i \text{Sim}(Q_i, K_j) V_j}{\sum_{j=1}^i \text{Sim}(Q_i, K_j)} \in \mathbb{R}^d$$

$$O_i = \frac{\phi(Q_i)^\top \sum_{j=1}^i \phi(K_j) V_j^\top}{\phi(Q_i)^\top \sum_{j=1}^i \phi(K_j)}$$

$$O_i = \frac{\phi(Q_i)^\top S_i}{d_i}$$

- 線形時変システム: パラメタを時間に合わせて用意するSSM→

- Linear Attention:  $d_i = 1$ すると  $S_{i+1} = S_i + \phi(K_{i+1}) V_{i+1}^\top$   $O_{i+1} = \phi(Q_{i+1})^\top S_{i+1}$

⇒  $\mathbf{A} = I$   $\mathbf{B} = I$ ,  $u_i = \phi(K_i) V_i^\top$ ,  $\mathbf{C}_i = \phi(Q_i)^\top$  としたときのSSM

$$x_i = \mathbf{A}_i x_{i-1} + \mathbf{B}_i u_i,$$

$$y_i = \mathbf{C}_i x_i + \mathbf{D}_i u_i.$$

$$\begin{aligned} x_{i+1} &= \mathbf{A} x_i + \mathbf{B} \phi(K_i) V_i^\top \\ y_{i+1} &= \mathbf{C} x_i, \end{aligned}$$

- H3 layer:  $\phi(K_i) = \text{SSM}_{\text{shift}}(K_i)$ としたときのLinear Attention

⇒ Linear Attention と似たような特徴を持てる!

# 「Hungry」HiPPO の意味

HiPPO: 記憶を保持できるように設計されたA行列

$$A_{nk} = - \begin{cases} (2n+1)^{1/2}(2k+1)^{1/2} & \text{if } n > k \\ n+1 & \text{if } n = k \\ 0 & \text{if } n < k \end{cases}$$

H3で使用している行列:

- 対角行列
  - シフト行列
- } HiPPOの一部 → hungry HiPPO × 2 → Hungry Hungry HiPPO → H3

ただし... 著者本人はブレインストーミングで出てきた遊びみたいな名前だっている

- [https://www.youtube.com/watch?v=x\\_Z9fzYCIB0&t=1191s](https://www.youtube.com/watch?v=x_Z9fzYCIB0&t=1191s)

# H3の評価: 言語モデル

Transformer・既存のSSMと比較して評価する

- 人工言語を利用したタスク(再掲)

Task	Random	S4D	Gated State Spaces	H3	Attention
Induction Head	5.0	35.6	6.8	<b>100.0</b>	<b>100.0</b>
Associative Recall	25.0	86.0	78.0	99.8	<b>100.0</b>

⇒ H3は入力のコピー&トークンの比較が可能

- 自然言語 (OpenWebText) の言語モデリング
  - GPT-2 small [\[Radford+19\]](#) と同等のパラメタ(12層)で評価
  - Hybrid: Transformer層を2・8層目 (2・2+L/2層目) に挿入したモデル

H3	H3 Hybrid (2 Attn)	S4D	GSS	GSS Hybrid (2 Attn)	Transformer
21.0	<b>19.6</b>	24.9	24.0	19.8	20.6

⇒ H3でTransformerと同等・Hybridにすると性能向上



# H3の評価: 事前学習モデル

## 言語モデルで事前学習して評価

- 設定:
  - 同一サイズのHybridモデルのH3を用意してPileコーパス [\[Gao+20\]](#) で事前学習
  - 他のコーパスについてはZero-shotで評価のみ行った場合のスコア
  - ベースラインは事前学習済みのチェックポイントを使用

- 結果:
    - GPT-Neoを超える性能が得られた
- ⇒ SSMの事前学習がうまくできている
- Zero-shotの評価についても優れている

Model	Pile	Zero-shot	
		OpenWebText	WikiText103
GPT-2 small (125M)	19.0*	22.6	29.9
GPT-Neo-125M	9.4	22.6	26.3
<b>Hybrid H3-125M</b>	<b>8.8</b>	<b>20.9</b>	<b>23.7</b>
GPT-2 medium (355M)	13.9*	17.0	21.8
<b>Hybrid H3-355M</b>	<b>7.1</b>	<b>15.9</b>	<b>16.9</b>
GPT-2 XL (1.5B)	12.4*	12.9	17.0
GPT-Neo-1.3B	6.2	13.1	13.3
<b>Hybrid H3-1.3B</b>	<b>6.0</b>	<b>12.4</b>	<b>12.5</b>
GPT-Neo-2.7B	5.7	11.7	11.5
<b>Hybrid H3-2.7B</b>	<b>5.4</b>	<b>11.0</b>	<b>10.6</b>

\* GPT-2はPileを事前学習に使っていない

# H3の評価: SuperGLUE [\[Wang+19\]](#)

- 設定:

- 選択肢のうち尤度が高いものを選択して出力とする
- Few-shotはGPT-3 [\[Brown+20\]](#)の設定同様, プロンプトで与える

- 結果:

- Zero/Few-shot の場合で共に過半数のタスクで既存モデルを上回った

## Zero-shot

Model	WSC	WIC	RTE	CB	MultiRC	ReCoRD	BoolQ	COPA	Average
OPT-125M	<b>39.4</b>	<u>52.0</u>	<u>48.7</u>	<u>37.4</u>	<u>58.9</u>	<u>44.9</u>	<u>59.6</u>	<u>60.0</u>	<u>50.1</u>
GPT-Neo-125M	<u>36.5</u>	<b>53.6</b>	<u>53.1</u>	<u>41.1</u>	<b>59.9</b>	<u>39.6</u>	<b>62.2</b>	<u>60.0</u>	<u>50.8</u>
<b>Hybrid H3-125M</b>	<b>39.4</b>	<u>51.4</u>	<b>59.2</b>	<b>48.2</b>	<u>51.4</u>	<b>55.0</b>	<u>59.6</u>	<b>67.0</b>	<b>53.9</b>
GPT-2 medium (355M)	<u>50.0</u>	<b>52.0</b>	<u>51.3</u>	<u>28.6</u>	<b>59.5</b>	<u>53.3</u>	<u>61.0</u>	<u>65.0</u>	<u>52.6</u>
OPT-350M	<b>53.5</b>	<u>50.8</u>	<u>53.4</u>	<u>35.7</u>	<u>58.9</u>	<u>51.4</u>	<u>60.9</u>	<u>60.0</u>	<u>53.1</u>
<b>Hybrid H3-355M</b>	<u>37.5</u>	<u>51.7</u>	<b>55.2</b>	<b>41.1</b>	<b>59.5</b>	<b>62.3</b>	<b>61.5</b>	<b>69.0</b>	<b>54.7</b>
OPT-1.3B	<u>36.5</u>	<u>49.5</u>	<b>53.4</b>	<b>39.3</b>	<b>58.3</b>	<u>61.8</u>	<u>55.0</u>	<u>69.0</u>	<u>52.9</u>
GPT-Neo-1.3B	<u>41.3</u>	<u>50.0</u>	<u>52.3</u>	<u>33.9</u>	<u>57.9</u>	<u>55.5</u>	<u>59.9</u>	<u>66.0</u>	<u>52.1</u>
<b>Hybrid H3-1.3B</b>	<b>52.9</b>	<b>50.3</b>	<b>53.4</b>	<u>33.9</u>	<u>58.2</u>	<b>67.8</b>	<b>61.7</b>	<b>74.0</b>	<b>56.5</b>
OPT-2.7B	<b>51.0</b>	<u>50.8</u>	<u>50.5</u>	<u>41.1</u>	<u>57.4</u>	<u>65.9</u>	<u>60.9</u>	<u>66.0</u>	<u>55.5</u>
GPT-Neo-2.7B	<u>37.5</u>	<u>50.0</u>	<u>52.3</u>	<b>50.0</b>	<b>59.1</b>	<u>60.0</u>	<b>61.1</b>	<u>67.0</u>	<u>54.6</u>
<b>Hybrid H3-2.7B</b>	<u>36.5</u>	<b>51.3</b>	<b>57.0</b>	<u>37.5</u>	<u>58.7</u>	<b>71.3</b>	<b>61.1</b>	<b>81.0</b>	<b>56.8</b>

## 3-shot

Model	WSC	WIC	RTE	CB	MultiRC	ReCoRD	BoolQ	COPA	Average
OPT-125M	<u>36.5</u>	<b>50.2</b>	<u>47.3</u>	<u>44.6</u>	<b>57.9</b>	<u>44.9</u>	<u>41.9</u>	<u>60.0</u>	<u>47.9</u>
GPT-Neo-125M	<u>38.5</u>	<u>50.0</u>	<u>53.1</u>	<u>17.9</u>	<u>56.3</u>	<u>39.6</u>	<b>62.1</b>	<u>60.0</u>	<u>47.2</u>
<b>Hybrid H3-125M</b>	<b>43.3</b>	<u>49.1</u>	<b>58.1</b>	<b>51.8</b>	<u>48.9</u>	<b>55.0</b>	<u>56.1</u>	<b>67.0</b>	<b>53.7</b>
GPT-2 medium (355M)	<u>36.5</u>	<b>50.5</b>	<u>48.0</u>	<u>8.9</u>	<u>43.5</u>	<u>53.3</u>	<u>58.8</u>	<u>65.0</u>	<u>45.6</u>
OPT-350M	<u>37.5</u>	<u>50.0</u>	<u>45.8</u>	<b>44.6</b>	<u>49.8</u>	<u>51.4</u>	<b>61.7</b>	<u>60.0</u>	<u>50.1</u>
<b>Hybrid H3-355M</b>	<b>42.3</b>	<u>47.5</u>	<b>50.5</b>	<u>28.6</u>	<b>59.7</b>	<b>62.3</b>	<u>60.5</u>	<b>69.0</b>	<b>52.6</b>
OPT-1.3B	<b>44.2</b>	<b>51.1</b>	<u>53.4</u>	<u>16.1</u>	<b>59.9</b>	<u>62.1</u>	<u>38.3</u>	<u>70.0</u>	<u>49.4</u>
GPT-Neo-1.3B	<u>35.6</u>	<u>50.6</u>	<u>47.3</u>	<b>32.1</b>	<b>59.9</b>	<u>55.7</u>	<b>61.2</b>	<u>67.0</u>	<u>51.2</u>
<b>Hybrid H3-1.3B</b>	<u>36.5</u>	<u>49.2</u>	<b>55.2</b>	<u>23.2</u>	<u>59.3</u>	<b>67.6</b>	<u>56.9</u>	<b>76.0</b>	<b>53.0</b>
OPT-2.7B	<u>44.2</u>	<u>50.5</u>	<b>53.4</b>	<u>17.9</u>	<u>59.2</u>	<u>66.0</u>	<b>62.0</b>	<u>71.0</u>	<u>53.0</u>
GPT-Neo-2.7B	<b>49.0</b>	<b>51.9</b>	<u>51.6</u>	<u>21.4</u>	<u>57.0</u>	<u>60.0</u>	<u>56.0</u>	<u>68.0</u>	<u>51.9</u>
<b>Hybrid H3-2.7B</b>	<u>36.5</u>	<u>45.6</u>	<u>47.3</u>	<b>46.4</b>	<b>59.4</b>	<b>71.1</b>	<u>60.6</u>	<b>77.0</b>	<b>55.5</b>

# H3の評価: スループット

単位時間あたりに出力可能なトークン数(スループット)を計測

- 目的: 再帰的なH3モデルとAuto RegressiveなTransformerと比べてどれだけ早く生成が可能かを確認
- 結果: スループットは約2.4倍  $\Rightarrow$  より高速に言語モデルが可能

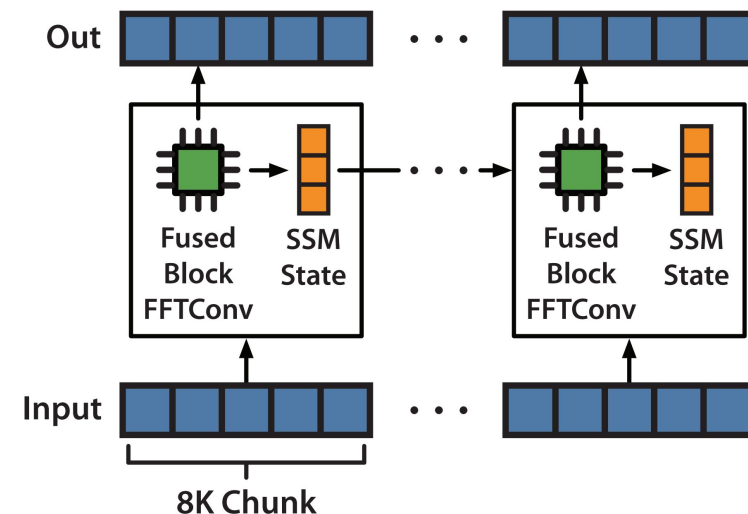
Table 7: Inference throughput on A100 80GB, 1.3B models. Batch size 64, prompt length 512, 1024, or 1536, and generating 128 tokens per sequence in the batch (i.e.,  $64 \times 128$  tokens in a batch). Hybrid H3 is up to  $2.4\times$  faster than a Transformer of similar size in inference. The difference is larger for longer sequences.

Tokens/s	Prompt length 512	Prompt length 1024	Prompt length 1536
Transformer-1.3B	1340	770	520
Hybrid H3-1.3B	1980	1580	1240

# FlashConv: SSMの計算の高速化

SSMの計算をGPUに寄り添った計算になるように工夫 → SSM全体に利用可

- Fused Block FFTConv: 2つの方法で計算を高速化
  - Kernel Fusion: 計算順を工夫してメモリへのread/writeを減らす
    - AttentionでいうFlashAttention [\[Dao+22\]](#)
  - Block FFT: 計算途中に出現するブロック対角行列をTensorコアで計算
- State Passing: 系列がメモリに乗らない場合に分けて計算
  - RNNのような再帰的なモデルなので、計算済みの状態を渡せば任意のサイズのチャンクに分割可能  
⇒ メモリに載るようにチャンクに分割すれば長い系列も扱える



# FlashConvの評価

## 計算速度を既存のS4で比較

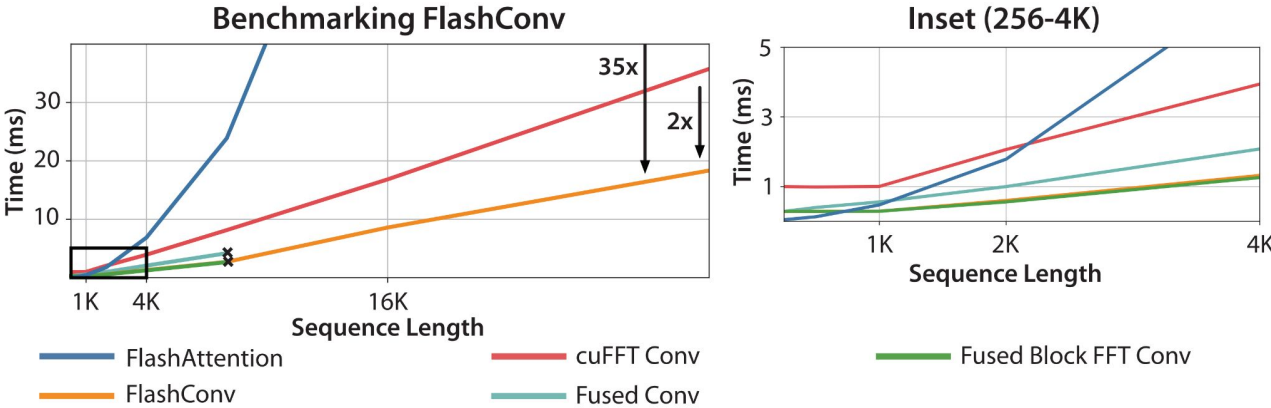
- Long-Range Arena ベンチマークでの速度評価
  - 既存のSoTAのS4を2倍の速度で動作させられるようになった
- Forward+Backwardの実行時間
  - 系列長に対してほぼ線形にしか実行時間は増加しない
  - CUDA比で2倍, Attention比で35倍の速さで実行可能

LRAの実行速度比較

Models	Speedup
Transformer	1×
FlashAttention (Dao et al., 2022b)	2.4×
Block-sparse FlashAttention (Dao et al., 2022b)	2.8×
S4 (Gu et al., 2022c)	2.9×
S4 with FLASHCONV	5.8×

LRAの性能評価

Model	ListOps	Text	Retrieval	Image	Pathfinder	Path-X	Avg
S4D (Gu et al., 2022b)	<b>58.3</b>	87.3	90.7	<b>87.5</b>	<b>93.6</b>	<b>92.3</b>	<b>85.0</b>
H3	57.5	<b>88.2</b>	<b>91.0</b>	87.3	93.0	91.8	84.8

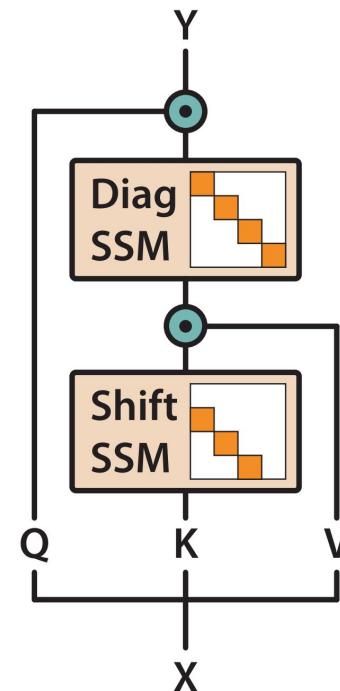


# まとめ

- SSMがTransformerに言語面だけ劣っている  $\Rightarrow$  原因の調査とその改善をした
  - 調査方法: 人工言語を使ったタスクを解いて不足した能力を割り出す
  - 調査結果: 入力をコピーして出力する機構と入力トークン同士を比較する能力が不十分
  - 改善方法: 二種類の行列を使ったSSMでLinear Attentionを真似る
  - 評価結果: 人工言語を使ったタスクだけでなく自然言語処理タスクにおいてもTransformerと同等以上の性能を示した
- GPUに寄り添った計算方法でSSMの計算コストを下げた

## 今後の課題

- 1.3BモデルまではSSMのパラメタをうまく使えた  $\Rightarrow$  更に大きく
- TransformerとH3の組み合わせがよかった  $\Rightarrow$  いい組合せを探す



$$\mathbf{Q} \odot \text{SSM}_{\text{diag}}(\text{SSM}_{\text{shift}}(\mathbf{K}) \odot \mathbf{V})$$

# 参考文献

[Brogan+74] Brogan. 1974. Modern control theory.

[[Fu+23](#)] Fu et al. 2023. Hungry Hungry Hippos: Towards Language Modeling with State Space Models. In Proc. of ICLR.

[[Gu+22](#)] Gu et al. 2022. Efficiently Modeling Long Sequences with Structured State Spaces. In Proc. of ICLR.

[[Gu+20](#)] Gu et al. 2020. HiPPO: Recurrent Memory with Optimal Polynomial Projections. In Proc. of NeurIPS.

[[Vaswani+19](#)] Vaswani et al. 2019. Attention is All you Need. In Proc. of NeurIPS.

[[Katharopoulos+20](#)] Katharopoulos et al. 2020. Transformers are RNNs: Fast Autoregressive Transformers with Linear Attention. In Proc. of ICML.

[[Gu+21](#)] Gu et al. 2021. Combining Recurrent, Convolutional, and Continuous-time Models with Linear State Space Layers. In Proc. of ICLR.

[[Olson+22](#)] Olson et al. 2022. In-context Learning and Induction Heads. Transformer Circuits Thread.

[[Ba+16](#)] Ba et al. 2016. Using Fast Weights to Attend to the Recent Past. In Proc. of NIPS.

[[Radford+19](#)] Radford et al. 2019. Language Models are Unsupervised Multitask Learners. In OpenAI blog.

[[Gao+20](#)] Gao et al. 2020. The Pile: An 800GB Dataset of Diverse Text for Language Modeling. In ArXiv.

[[Wang+19](#)] Wang et al. 2019. SuperGLUE: A Stickier Benchmark for General-Purpose Language Understanding Systems. In Proc. of NeurIPS.

[[Brown+20](#)] Brown et al. 2020. Language Models are Few-Shot Learners. In Proc. of NeurIPS.

[[Dao+22](#)] Dao et al. 2022. FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness. In Proc. of NeurIPS.

# 付録



# 付録:SSMの畳み込みカーネルの計算効率化

行列のクラスを計算しやすいものに限定して高速な計算を計算を可能にする

1.  $A$  を正規行列+低rank行列 (Normal Plus Low Rank; NPLR)  $A = \Lambda + PQ^*$ で表現
  - $\Lambda$ : 対角行列  $P, Q \in \mathbb{C}^{N \times L}$  ( $PQ^* \in \mathbb{C}^{N \times N}$ )
  - 行列のクラスを表現力が十分かつ計算しやすい形に限定する
2. カーネルの母関数をべき級数(離散フーリエ変換, DFT)  $\sum_{j=0}^{L-1} \bar{\mathbf{K}}_j \zeta^j$  計算
  - Woodbury identity (Sherman-Morrison-Woodburyの公式)によって累乗(L回の行列積)が逆行列を一つ計算すればよくなる
3. 逆フーリエ変換(iFFT)して計算したいカーネルを獲得
  - カーネルを対角行列に限定した場合, 高速に計算が可能な問題に帰着 (Cauchy Kernel)

# 付録:SSMの計算の高速化—Fused Block FFTConv

- Kernel Fusion: 計算順を工夫してメモリへのread/writeを減らす
  - FFT畳み込み(FFTConv)の通常の実装のボトルネックは畳み込みカーネルのread/write  
⇒ 畳み込みカーネルを計算してしまっ、メモリに載せてしまう  
= AttentionでいうFlashAttention [\[Dao+22\]](#)
- Block FFT: 計算途中に出現するブロック対角行列をTensorコアで計算
  - ブロック対角行列をブロックで分けて、小さい行列で計算するようにして、計算量を削減
  - Tensorコアには $16 \times 16$ 行列の積を高速に計算可能 ⇒ 小さい行列の積を計算するときの特