

Hungry Hungry Hippos: Towards Language Modeling with State Space Models

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

設計された
A行列

状態空間モデル

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

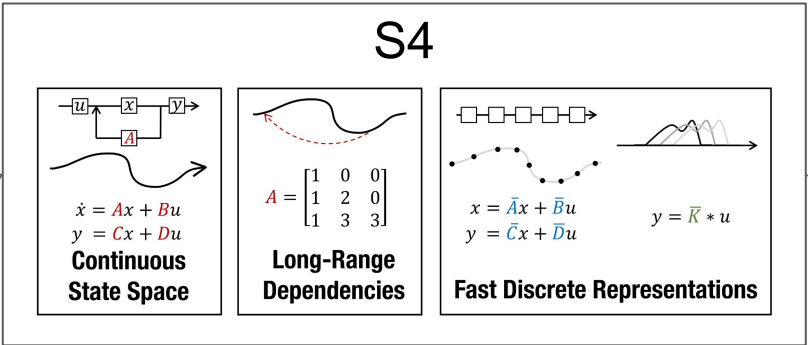
深層学習に導入

HiPPO

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

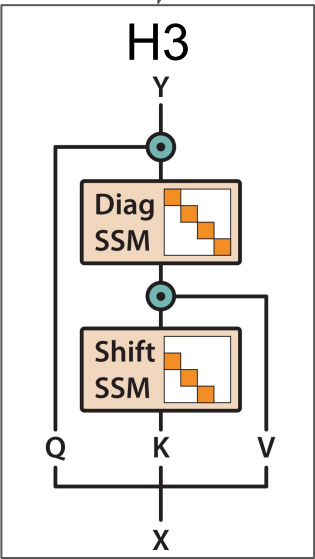
初期化に利用

系列を扱える深層学習モデル



言語に特化

CV, 音声, 時系列
ではS4が強い
⇔ 言語では
Transformerが強い

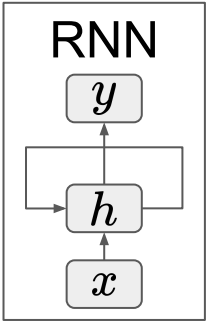
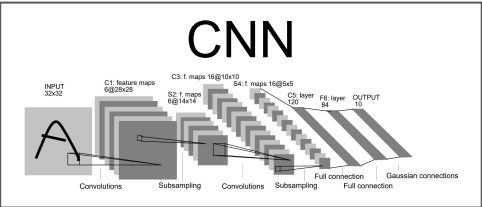


インスピレーション

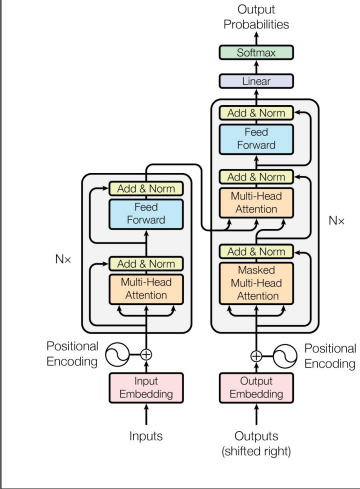
Linear Attention

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

variant



Transformer



論文の位置づけ

設計された
A行列

状態空間モデル

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

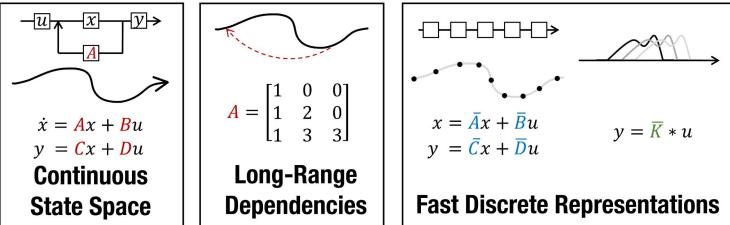
深層学習に導入

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

初期化に利用

系列を扱える深層学習モデル

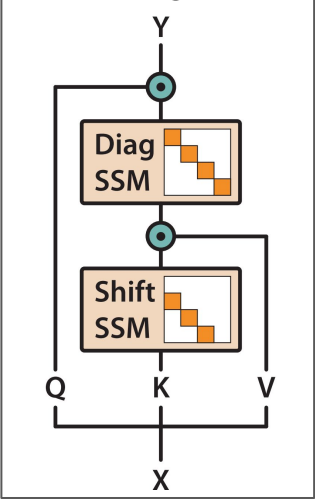
S4



言語に特化

CV, 音声, 時系列
ではS4が強い
⇔ 言語では
Transformerが強い

H3



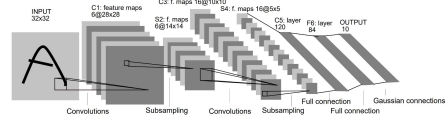
インスピレーション

Linear Attention

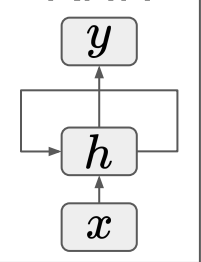
$$O_i = \frac{\phi(Q_i)^\top \sum_{j=1}^i \phi(K_j) V_j^\top}{\phi(Q_i)^\top \sum_{j=1}^i \phi(K_j)}$$

variant

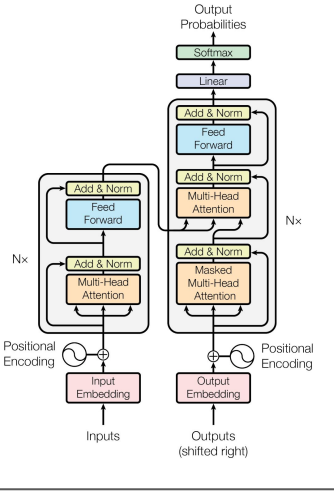
CNN



RNN



Transformer



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

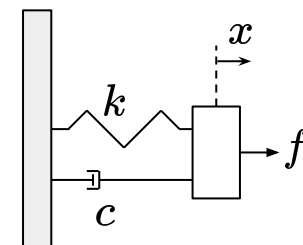
- 時系列解析や制御システムによく利用される

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

$$y = Cx + Du$$

例: 質点—バネ—ダンパー系

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



離散的な数値を扱えるようにSSMを離散化すると:

$$x_i = \bar{A}x_{i-1} + \bar{B}u_i \quad \bar{A} = (I - \Delta/2 \cdot A)^{-1}(I + \Delta/2 \cdot A)$$

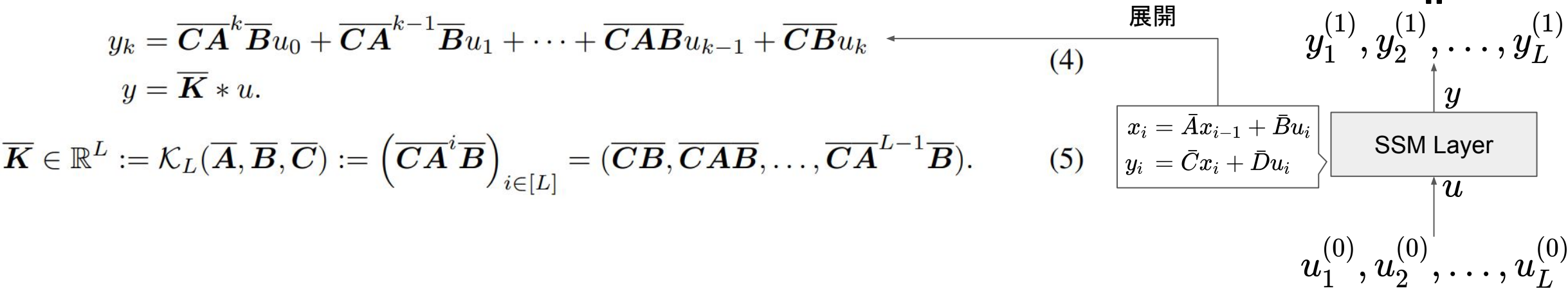
$$y_i = \bar{C}x_i + \bar{D}u_i \quad \bar{B} = (I - \Delta/2 \cdot A)^{-1}\Delta B \quad \bar{C} = C \quad \bar{D} = D$$

SSMの深層学習への導入

離散化したSSMの行列をパラメタにして特徴量抽出に利用

- Self-Attentionと同様に系列の特徴量抽出に利用可能
- SSMの問題点: 計算が重たい

SSMを計算機で扱いやすいように畳み込み形式にすると:



S4の特徴

- SSMの初期値にHiPPO行列を使用
 - HiPPOとは: 記憶を保持できるように設計されたA行列
 - HiPPOを使うとうまく系列を表現できるようになる
- カーネルの \bar{A}^k の計算が重たい \Rightarrow 高速に計算可能にする
 - カーネルの計算が大変なのでモデルの訓練に時間がかかってしまうのをうまく緩和
 - $A = \Lambda - PQ^*$ で表せるようなクラスに限定して高速化
 - HiPPOも含有するように行列のクラスを決定
 - 導出は付録を参照.

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

論文の位置づけ

設計された
A行列

状態空間モデル

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

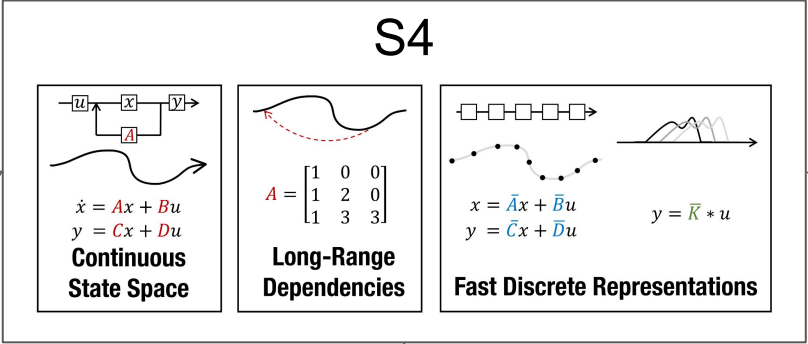
深層学習に導入

HiPPO

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

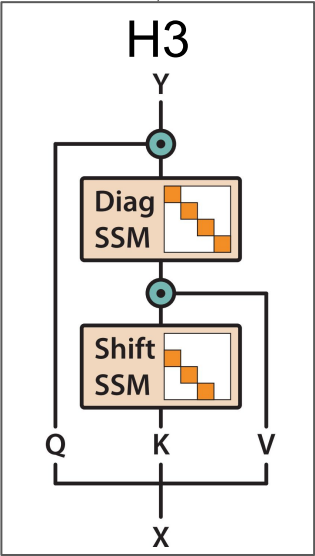
初期化に利用

系列を扱える深層学習モデル



言語に特化

CV, 音声, 時系列
ではS4が強い
⇔ 言語では
Transformerが強い

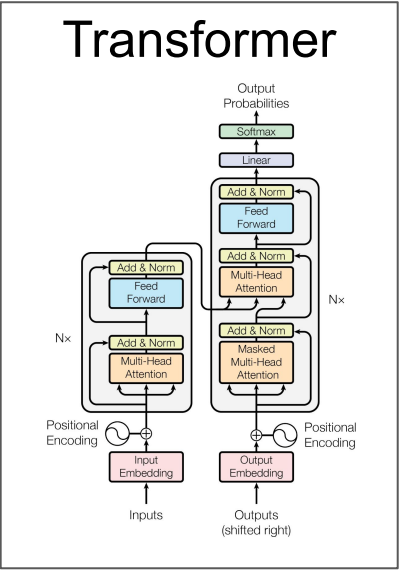
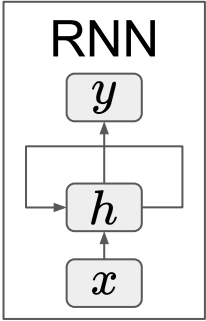
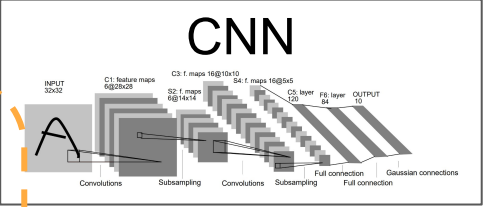


インスピレーション

Linear Attention

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

variant



背景: SSMとSelf-Attention

SSMとSelf-Attentionを比較してSSMに不足している部分を調査・改善する

- 性能面での比較:

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

} 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: $SSM(\text{ほぼ線形}) < Self-Attention(\text{二乗})$
- 実際の計算時間: $SSM < Self-Attention$

} 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

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既存のSSMはkey-valueのペアを保持できないと確認

- 対象タスク: 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>	<i>2</i>	20	10

- 評価: Attentionでできることが既存のSSMではできない

- 2層スタックしたモデルでタスクを解いた時の正解率で評価

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

- 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(特殊トークン)が発生してからのイベントを記憶できない
- トークン同士を比較できない

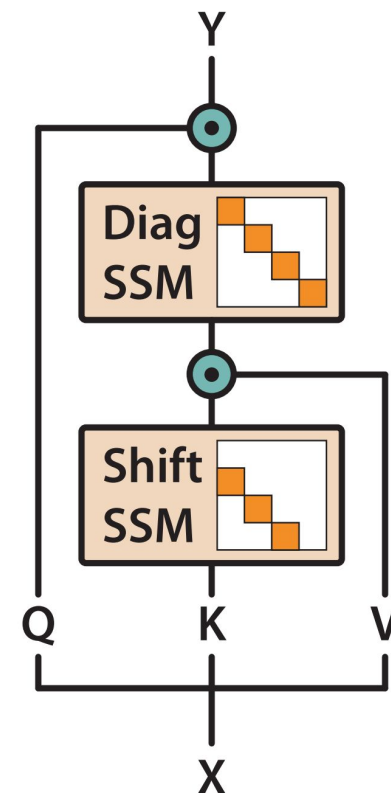
※ Attentionでは入力中のトークンをコピー可 & Attention scoreの計算でトークンの比較が可能

改善して提案するモデルは二つの能力がある

H3 layer

二つのSSMでトークン同士の比較と入力のコピーを実現

- $\text{SSM}_{\text{diag}} \cdot \text{SSM}_{\text{shift}}$: A行列をそれぞれ対角行列とシフト行列としたSSM
- $\text{SSM}_{\text{shift}}(\mathbf{K}) \odot \mathbf{V}$: シフトした系列と元の系列の要素積
 - $[0, x_1^{\text{key}}, \dots, x_{L-1}^{\text{key}}] \odot [x_1^{\text{val}}, x_2^{\text{val}}, \dots, x_L^{\text{val}}]$ のような項が作れる
 - ⇒ 他のトークンとのローカルな比較が可能に
 - ※ S4で扱う行列の範囲にシフト行列は含まれない
- $\mathbf{Q} \odot \text{SSM}_{\text{diag}}(\cdot)$: 対角行列のSSMで入力を保存 & 表現を作る
 - ⇒ 元の入力のコピーが実現できる
 - ⇒ 比較を全体でできるようになる



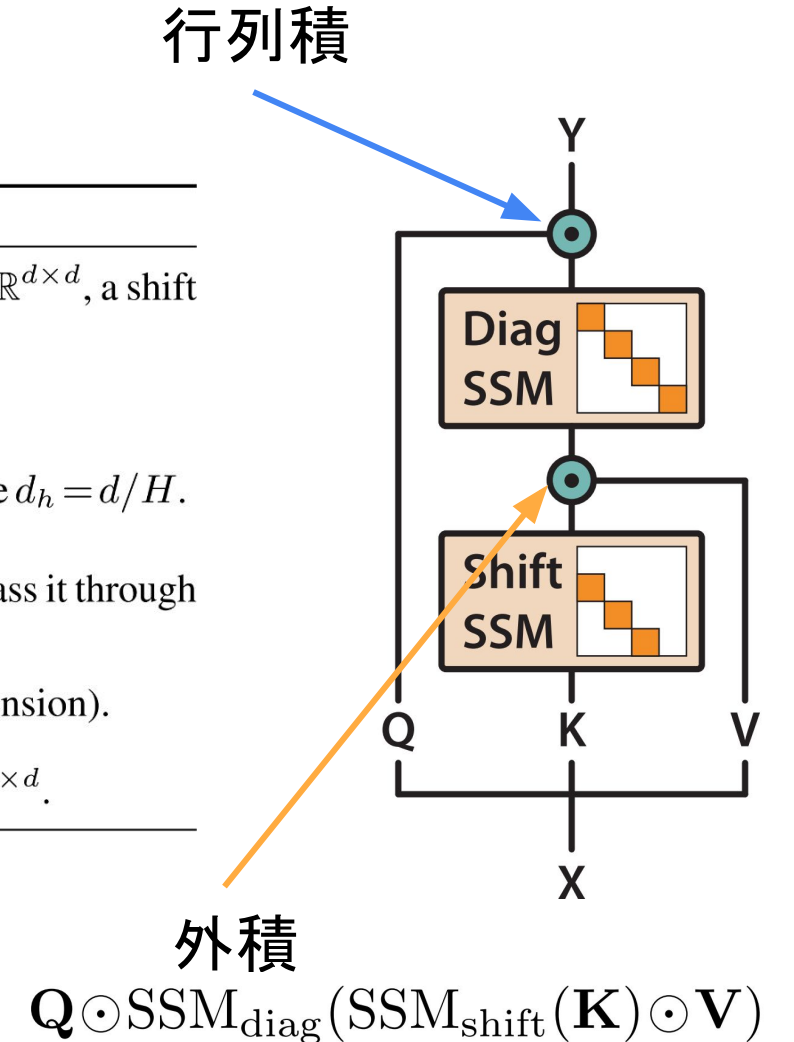
$$\mathbf{Q} \odot \text{SSM}_{\text{diag}}(\text{SSM}_{\text{shift}}(\mathbf{K}) \odot \mathbf{V})$$

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 SSM_{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: $\mathbf{K}\mathbf{V}^{(h)} = \text{SSM}_{\text{diag}}(\bar{\mathbf{K}}^{(h)}(\mathbf{V}^{(h)})^\top) \in \mathbb{R}^{N \times d_h \times d_h}$.
- 6: Batch-multiply by \mathbf{Q} : $\mathbf{O}^{(h)} = [\mathbf{Q}_1^{(h)} \mathbf{K}\mathbf{V}_1^{(h)}, \dots, \mathbf{Q}_N^{(h)} \mathbf{K}\mathbf{V}_N^{(h)}] \in \mathbb{R}^{N \times d_h}$ (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 $\mathbf{W}_O \in \mathbb{R}^{d \times d}$.

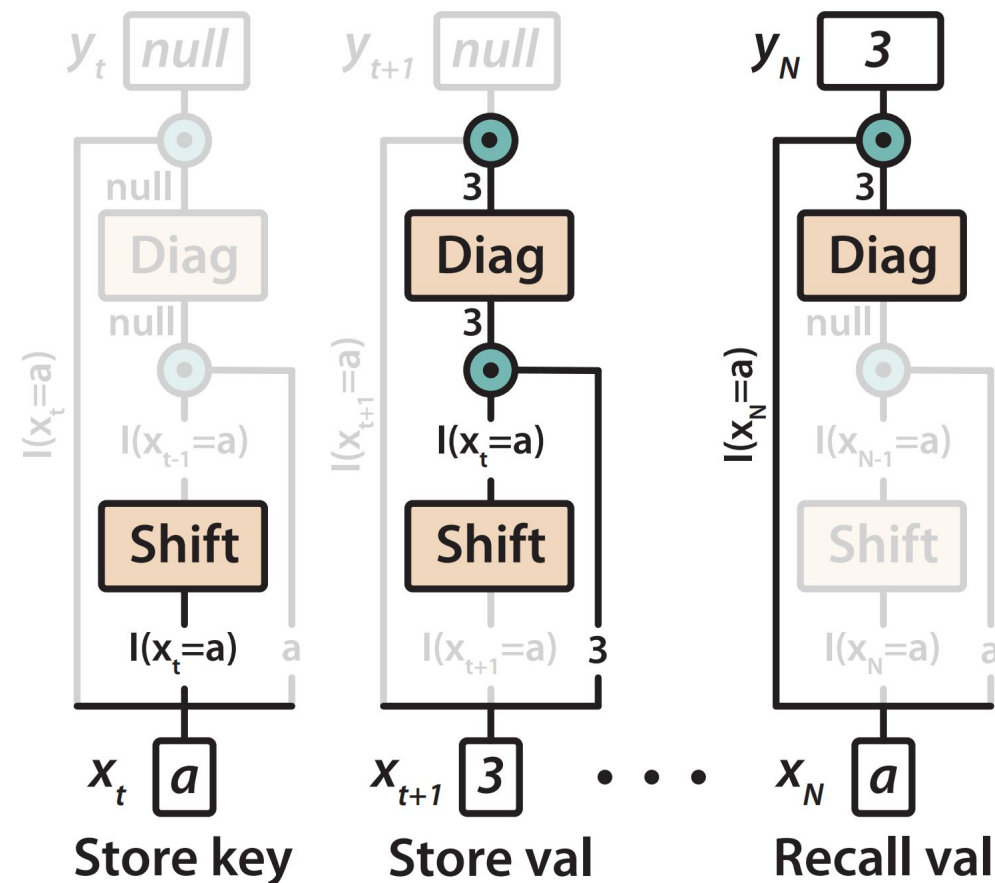


H3 layer: 動作例

- 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

設計された
A行列

状態空間モデル

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

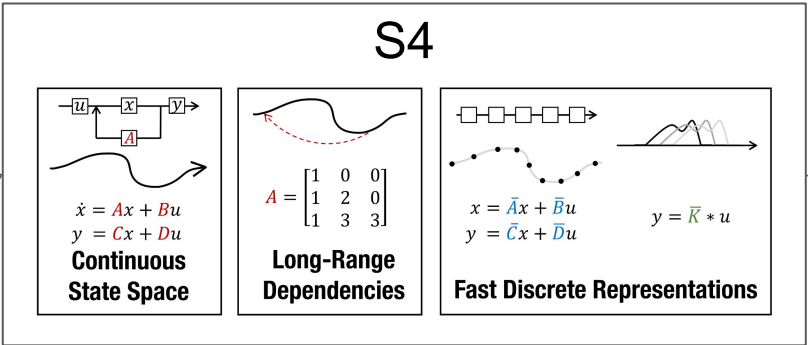
深層学習に導入

HiPPO

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

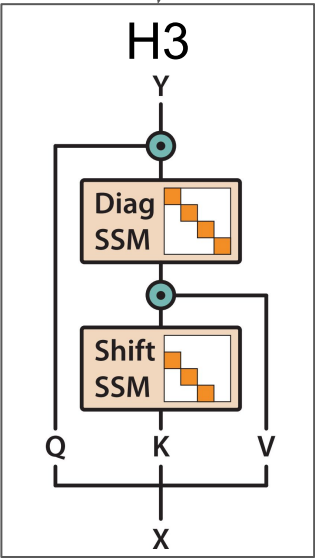
初期化に利用

系列を扱える深層学習モデル



言語に特化

CV, 音声, 時系列
ではS4が強い
⇔ 言語では
Transformerが強い

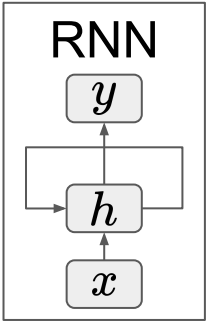
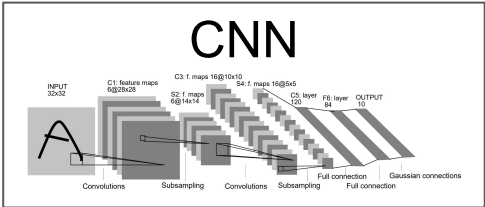


インスピレーション

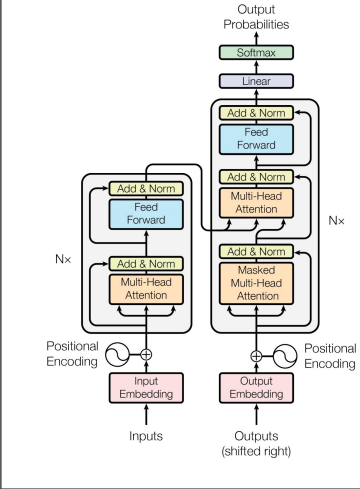
Linear Attention

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

variant



Transformer



H3 layer: Attentionとの比較

- Attentionの定式化:
 - 一般的なAttention [Vaswani]: $\text{Sim}(q, k) = e^{q^\top k}$
 - Linear Attention: $\text{Sim}(q, k) = \phi(q)^\top \phi(k)$
 - Linear AttentionのRNN標記:

$$z_i = \sum_{j=1}^i \phi(K_j) \in \mathbb{R}^d \quad S_i = \sum_{j=1}^i \phi(K_j) V_j^\top \in \mathbb{R}^{d \times d}$$
- 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 $S_{i+1} = S_i + \phi(K_{i+1}) V_{i+1}^\top$ $O_{i+1} = \phi(Q_{i+1})^\top S_{i+1}$
 $\Rightarrow \mathbf{A} = I \quad \mathbf{B} = I, u_i = \phi(K_i) V_i^\top, \mathbf{C}_i = \phi(Q_i)^\top$ としたときのSSM

$$\begin{aligned} x_i &= \mathbf{A}_i x_{i-1} + \mathbf{B}_i u_i, \\ y_i &= \mathbf{C}_i x_i + \mathbf{D}_i u_i. \end{aligned}$$
 - H3 layer: $\phi(K_i) = \text{SSM}_{\text{shift}}(K_i)$ としたときのLinear Attention

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

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

S4でHiPPOを含むように設計された行列: 対角行列+ベクトルの積 $A = \Lambda - PQ^*$

H3で使用している行列:

- 対角行列
 - シフト行列
- } HiPPOの一部 \rightarrow hungry HiPPO $\times 2 \rightarrow$ Hungry Hungry HiPPO \rightarrow H3

ただし... 筆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

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	100.0	100.0
Associative Recall	25.0	86.0	78.0	99.8	100.0

⇒ H3は入力のコピー&トークンの比較が可能

- 自然言語の言語モデリング
 - GPT-2 small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	19.6	24.9	24.0	19.8	20.6

⇒ H3でTransformerと同等・Hybridにすると性能向上に寄与

H3の評価: 事前学習言語モデル

事前学習となる言語モデルタスクで評価

- 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
- ⇒ 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
Hybrid H3-125M	8.8	20.9	23.7
GPT-2 medium (355M)	13.9*	17.0	21.8
Hybrid H3-355M	7.1	15.9	16.9
GPT-2 XL (1.5B)	12.4*	12.9	17.0
GPT-Neo-1.3B	6.2	13.1	13.3
Hybrid H3-1.3B	6.0	12.4	12.5
GPT-Neo-2.7B	5.7	11.7	11.5
Hybrid H3-2.7B	5.4	11.0	10.6

* GPT-2はPileを事前学習に使っていない

H3の評価: SuperGLUE

- 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の例はGPT-3の設定同様, プロンプトとして与える

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

3-shot

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

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

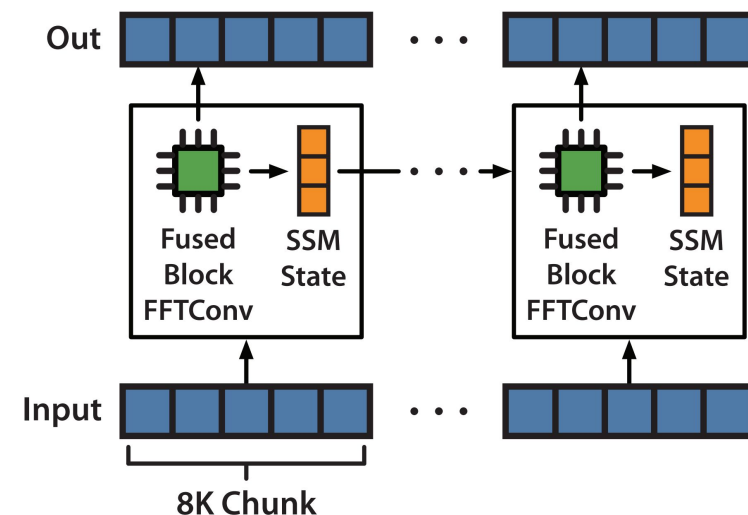
Table 7: Inference throughput on A100 80GB, 1.3B models. Batch size 64, prompt length 512, 1024, or 1536, and generating 128 tokens per sequence in the batch (i.e., 64×128 tokens in a batch). Hybrid H3 is up to $2.4\times$ faster than a Transformer of similar size in inference. The difference is larger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

SSMの計算の高速化: 概要

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つの方法でメモリに載る系列長の計算を高速化
 - 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 []
 - 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

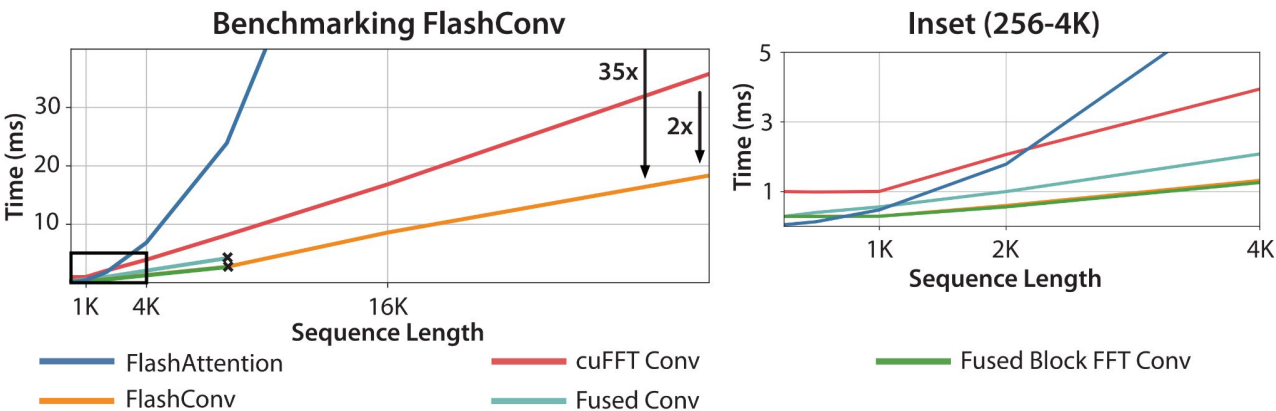


高速化の評価

計算速度を既存のSSM S4と比較
(FlashConvはSSMなら適用可)

- 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

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×

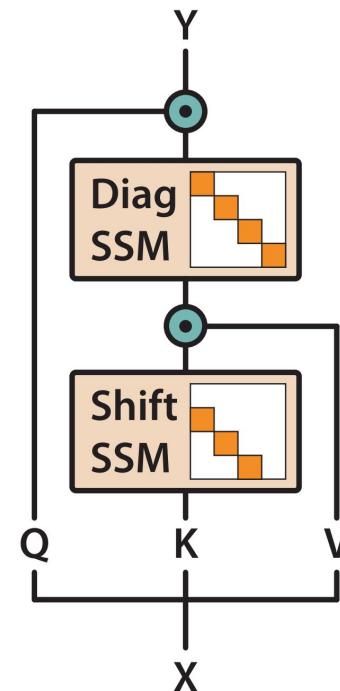


まとめ

- SSMがTransformerに言語面だけ劣っている ⇒ 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 ⇒ 更に大きく
- TransformerとH3の組み合わせがよかった ⇒ いい組合せを探す



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

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) $A = PQ^*$ で表現
 - A : 対角行列 $P, Q \in \mathbb{C}^{N \times 1}$ ($PQ^* \in \mathbb{C}^{N \times N}$)
 - 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(cf. HiPPO)
2. カーネルの母関数がべき級数(離散フーリエ変換)だと思って計算
- 3.

付録:SSMの計算の高速化—Fused Block FFTConv

付録:SSMの計算の高速化—State Passing

付録: いろいろなタスクでの評価