

HiPPO et S4 pour la modélisation de séquences avec dépendances à (très) longue portée

une unification des CNN et des RNN et une alternative aux Transformers

Séminaire TALia du 09/09/2022

1

arXiv > cs > arXiv:2008.07669

Computer Science > Machine Learning

[Submitted on 17 Aug 2020 (v1), last revised 23 Oct 2020 (this version, v2)]

HiPPO: Recurrent Memory with Optimal Polynomial Projections

Albert Gu, Tri Dao, Stefano Ermon, Atri Rudra, Christopher Re



2

arXiv > cs > arXiv:2110.13985

Computer Science > Machine Learning

[Submitted on 26 Oct 2021]

Combining Recurrent, Convolutional, and Continuous-time Models with Linear State-Space Layers

Albert Gu, Isys Johnson, Karan Goel, Khaled Saab, Tri Dao, Atri Rudra, Christopher Ré

3

arXiv > cs > arXiv:2111.00396

Computer Science > Machine Learning

[Submitted on 31 Oct 2021 (v1), last revised 5 Aug 2022 (this version, v3)]

Efficiently Modeling Long Sequences with Structured State Spaces

Albert Gu, Karan Goel, Christopher Ré



Stanford
MLSys
Seminar
#46

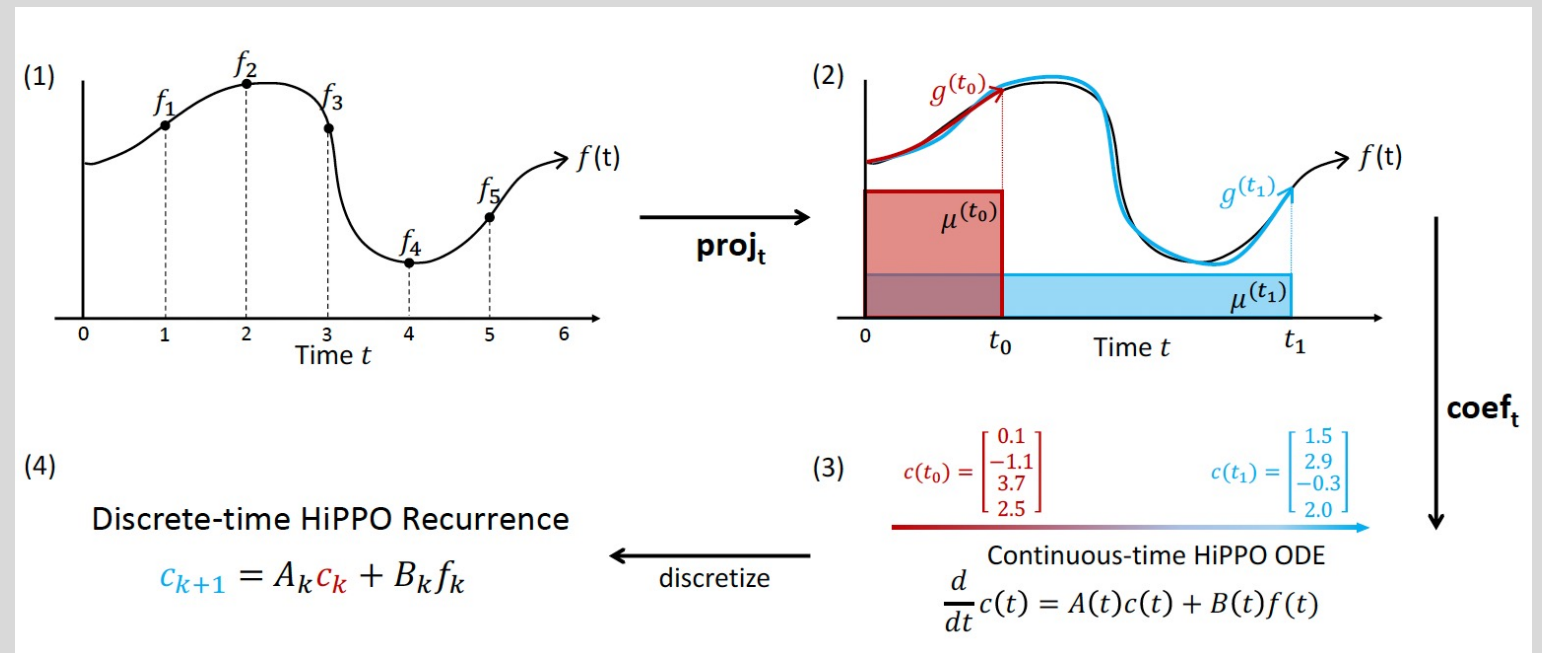
Efficiently Modeling Long Sequences with Structured State Spaces

Albert Gu



Plan

1. Le contexte historique
2. S4 relève le challenge des dépendances à longue portée
3. Les fondations théoriques de S4 avec HiPPO et SSM
4. Les 3 vues de S4 : CTM, RNN, CNN
5. Questions ouvertes pour le TAL

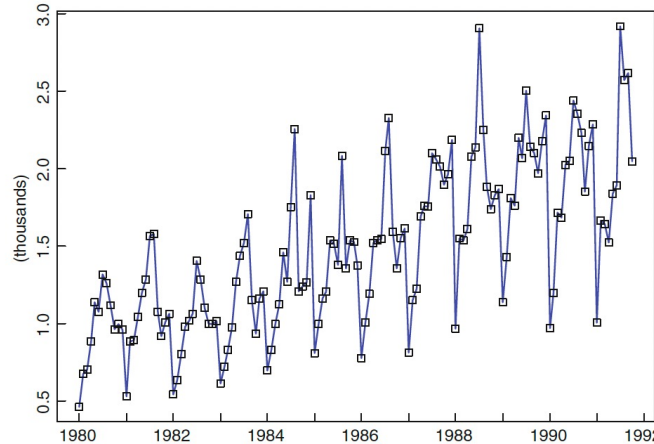
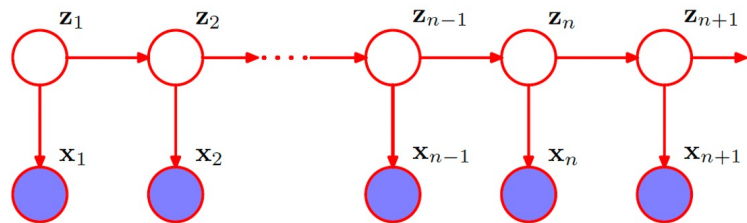


Les modèles de séquences – une longue histoire !

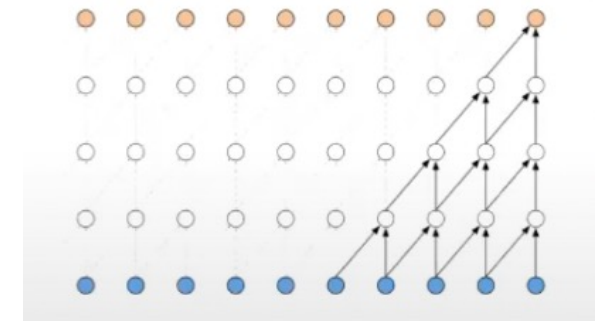
Modèles de séries temporelles ARMA

$$X_t = c + \varepsilon_t + \sum_{i=1}^p \varphi_i X_{t-i} + \sum_{i=1}^q \theta_i \varepsilon_{t-i}$$

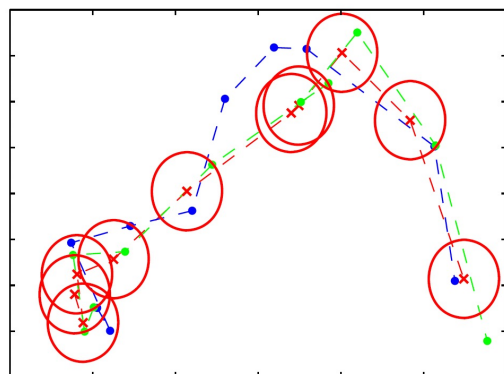
Hidden Markov Models (HMM)
Linear Dynamical Systems (LDS)



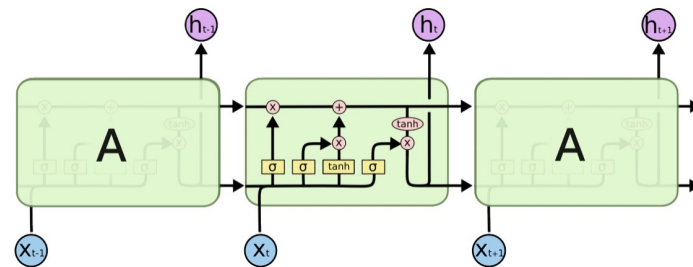
1D CNN



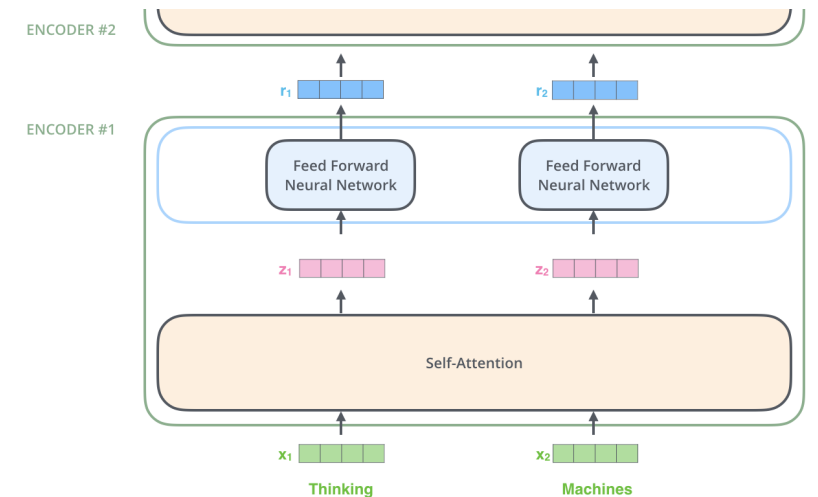
SSM, filtres de Kalman



RNN, LSTM, GRU, mécanisme d'attention



Transformers



S4 relève le challenge des dépendances à (très) longue portée

S4 versus Transformers les plus performants à l'entraînement

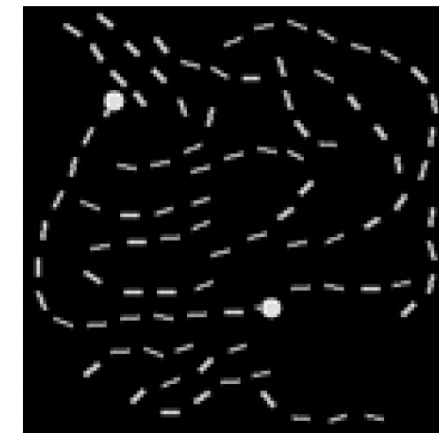
	LENGTH 1024		LENGTH 4096	
	Speed	Mem.	Speed	Mem.
Transformer	1×	1×	1×	1×
Performer	1.23×	0.43×	3.79×	0.086×
Linear Trans.	1.58×	0.37×	5.35×	0.067×
S4	1.58×	0.43×	5.19×	0.091×

S4 comme modèle de langue

Model	Params	Test ppl.	Tokens / sec
Transformer	247M	20.51	0.8K (1×)
GLU CNN	229M	37.2	-
AWD-QRNN	151M	33.0	-
LSTM + Hebb.	-	29.2	-
TrellisNet	180M	29.19	-
Dynamic Conv.	255M	25.0	-
TaLK Conv.	240M	23.3	-
S4	249M	21.28	48K (60×)

S4 pour la reconnaissance de la parole

	MFCC	RAW	0.5×
Transformer	90.75	✗	✗
Performer	80.85	30.77	30.68
ODE-RNN	65.9	✗	✗
NRDE	89.8	16.49	15.12
ExpRNN	82.13	11.6	10.8
LipschitzRNN	88.38	✗	✗
CKConv	95.3	71.66	65.96
WaveGAN-D	✗	96.25	✗
LSSL	93.58	✗	✗
S4	<u>93.96</u>	98.32	96.30

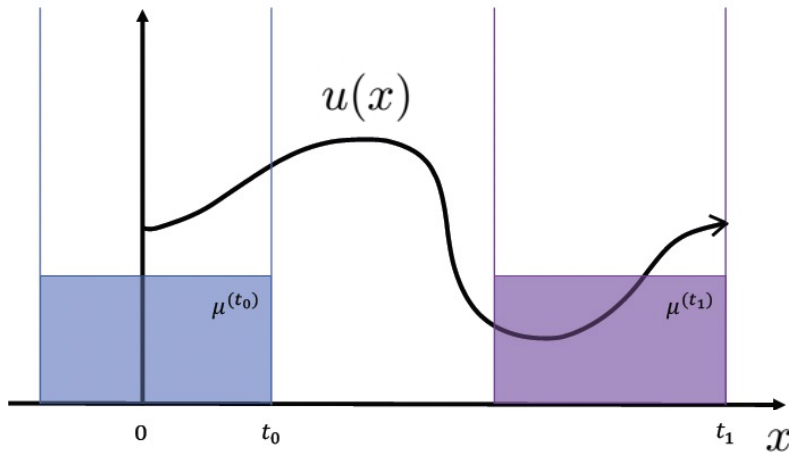


S4 sur les 6 tâches du benchmark Long Range Arena

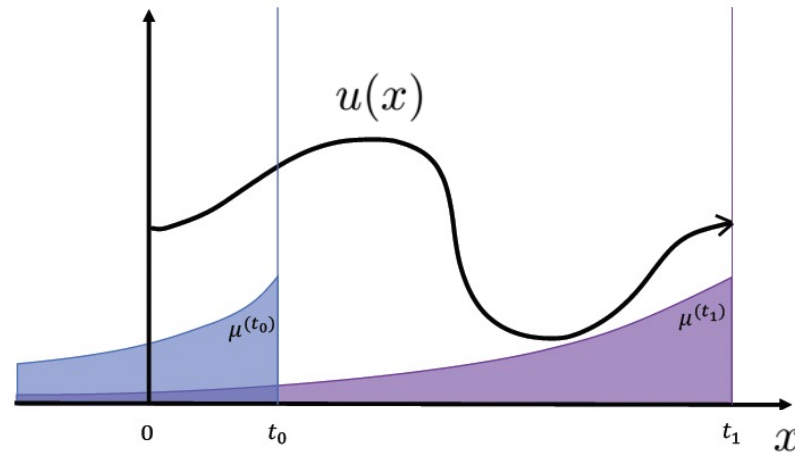
MODEL	LISTOPS	TEXT	RETRIEVAL	IMAGE	PATHFINDER	PATH-X	AVG
Transformer	36.37	64.27	57.46	42.44	71.40	✗	53.66
Reformer	<u>37.27</u>	56.10	53.40	38.07	68.50	✗	50.56
BigBird	36.05	64.02	59.29	40.83	74.87	✗	54.17
Linear Trans.	16.13	<u>65.90</u>	53.09	42.34	75.30	✗	50.46
Performer	18.01	65.40	53.82	42.77	77.05	✗	51.18
FNet	35.33	65.11	59.61	38.67	<u>77.80</u>	✗	54.42
Nyströmformer	37.15	65.52	<u>79.56</u>	41.58	70.94	✗	57.46
Luna-256	37.25	64.57	79.29	<u>47.38</u>	77.72	✗	<u>59.37</u>
S4	58.35	76.02	87.09	87.26	86.05	88.10	80.48

HiPPO : mémorisation optimale et incrémentale de l'histoire

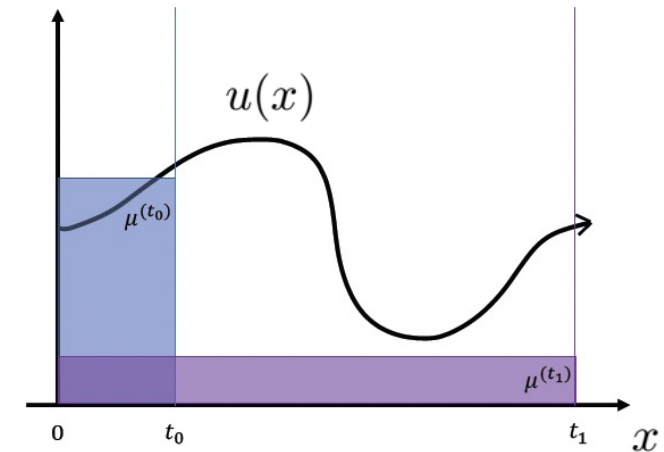
(1) choisir une pondération des évènements passés



Translated Legendre (LegT)



Translated Laguerre (LagT)



Scaled Legendre (LegS)

HiPPO : mémorisation optimale et incrémentale de l'histoire

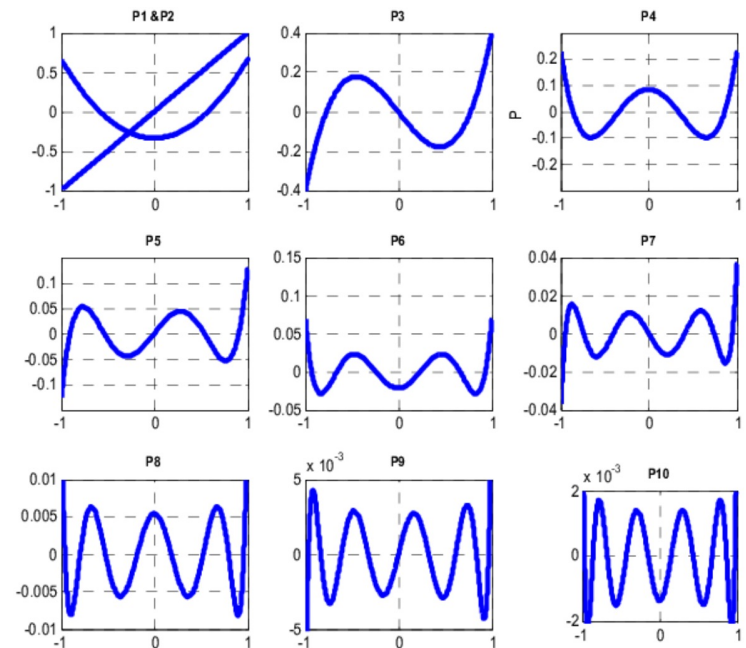
(2) choisir une base de fonctions $\{v_i^{(t)}\}_{i \leq N}$ pour approximer l'historique $u_{\leq t} := u(x)|_{x \leq t}$ par une combinaison

$$g^{(t)}(x) := \sum_{i=1}^N c_i(t) v_i^{(t)}(x)$$

qui minimise $\|u_{\leq t} - g^{(t)}\|_{L^2(\mu^{(t)})}$.

Pour des polynômes orthogonaux $\langle v_i^{(t)}, v_j^{(t)} \rangle_{\mu^{(t)}} = \delta_{ij}$ on a

$$c_i(t) = \langle v_i^{(t)}, u_{\leq t} \rangle_{\mu^{(t)}} = \int_0^t v_i^{(t)}(x) u_{\leq t}(x) d\mu^{(t)}(x)$$



HiPPO = **H**igher Order **P**olynomial **P**rojection **O**perators

HiPPO : mémorisation optimale et incrémentale de l'histoire

(3) écrire l'ED pour les coefficients $c(t) := (c_1(t), \dots, c_N(t)) \in \mathbb{R}^N$. On peut montrer que pour toute mesure $\mu^{(t)}$ on a une **ED linéaire** :

$$\frac{dc(t)}{dt} = \mathbf{A}(t)c(t) + \mathbf{B}(t)u(t), \quad \mathbf{A} \in \mathbb{R}^{N \times N}, \mathbf{B} \in \mathbb{R}^{1 \times N}$$

↑ mémoire instantanée compressée du signal $u(t)$ pour un budget N .
 ↑ signal scalaire

→ $c = \text{hippo}(u)$

non-trivial car il faut aussi dériver la mesure dans

$$c_i(t) = \int_0^t v_i^{(t)}(x) u_{\leq t}(x) d\mu^{(t)}(x)$$

Les matrices \mathbf{A} et \mathbf{B} sont **explicitement calculables** pour les $\mu^{(t)}$ classiques.

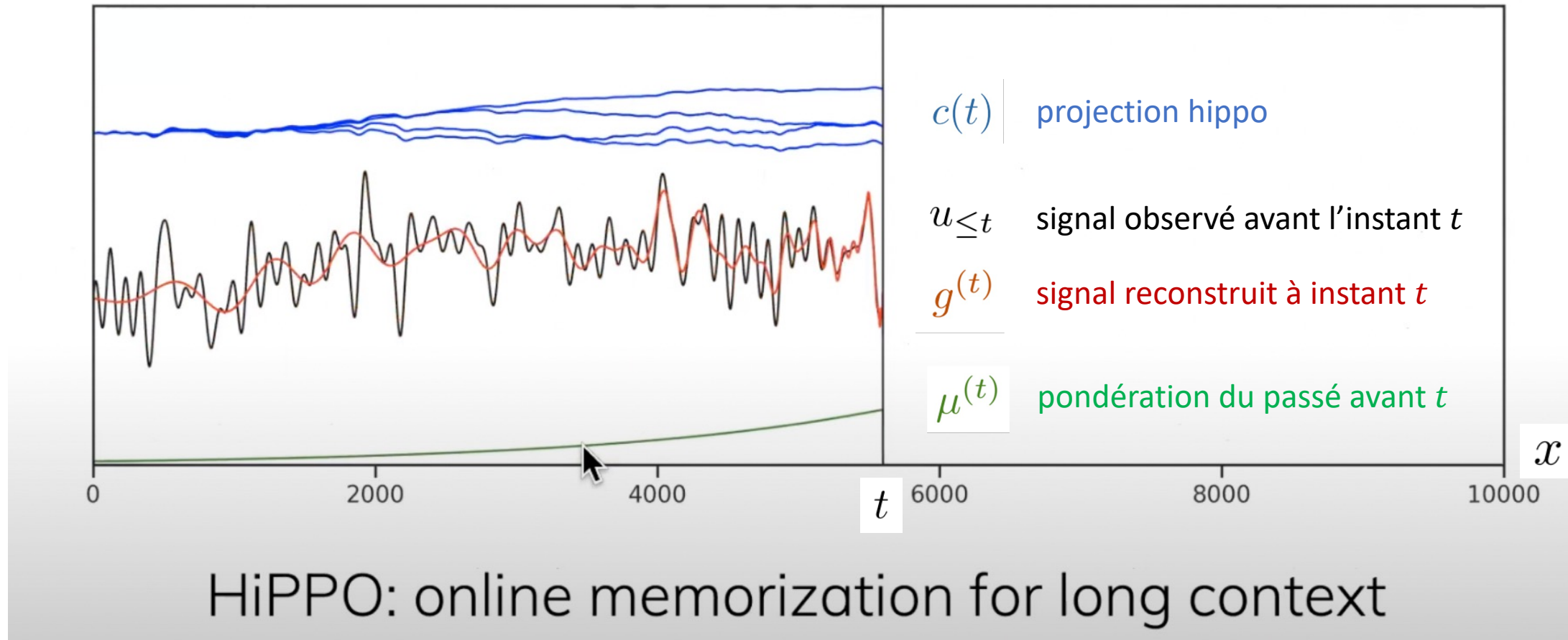
Pour LegS $\mu^{(t)} = \frac{1}{t} \mathbb{I}_{[0,t]}$

$$-t A_{nk} = \begin{cases} (2n+1)^{1/2} (2k+1)^{1/2} & \text{si } n > k \\ n+1 & \text{si } n = k \\ 0 & \text{si } n < k \end{cases}, \quad t B_n = (2n+1)^{\frac{1}{2}}$$

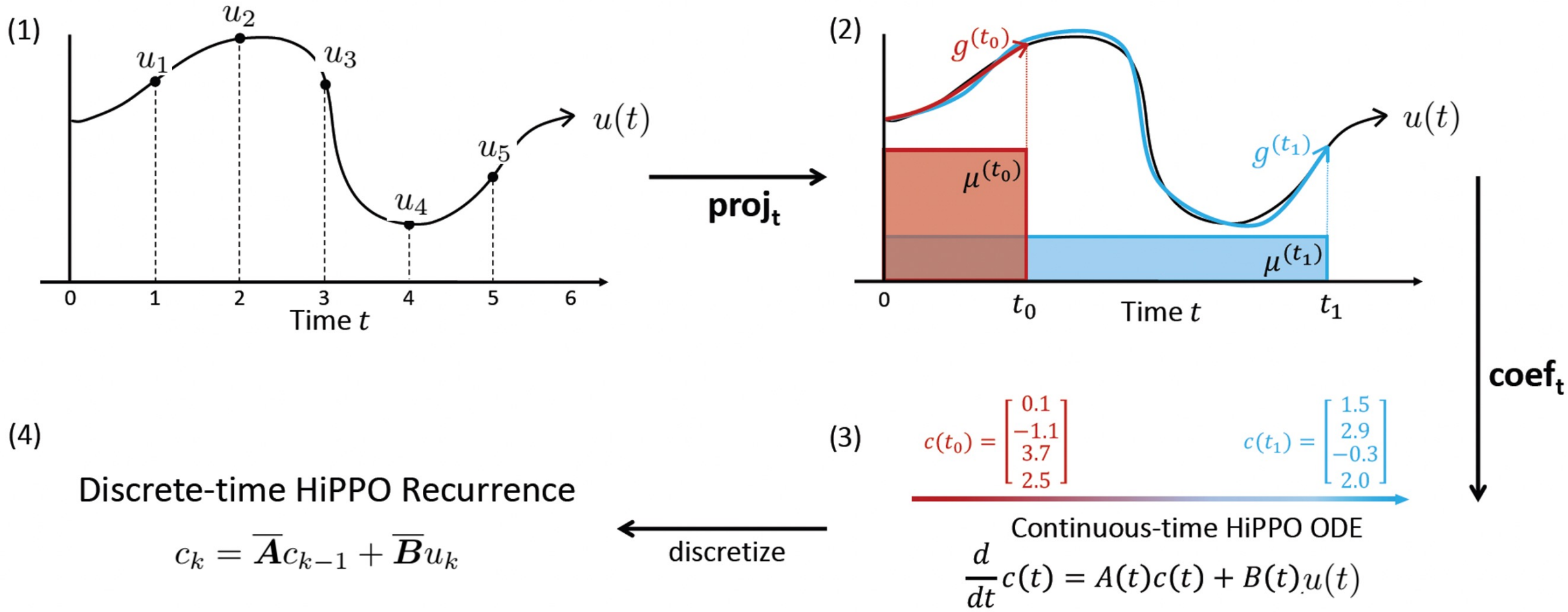
↑ **Matrice HiPPO**

Pour LegS on a **équivariance sous les changements d'échelle** : $\gamma : t \rightarrow \alpha t$ $\text{hippo}(u \circ \gamma) = \text{hippo}(u) \circ \gamma$

HiPPO : reconstruction à la volée du signal à partir de l'état



HiPPO : une séquence comme discrétisation d'un signal



State Space Models (SSM)

Inconvénient de hippo : on passe d'un signal $u(t)$ 1-D à un état $c(t)$ N -D :

Solution : SSM extraire un signal $y(t)$ 1-D de l'état avec une projection affine de l'état $c(t)$:

SSM cas continu (ED)

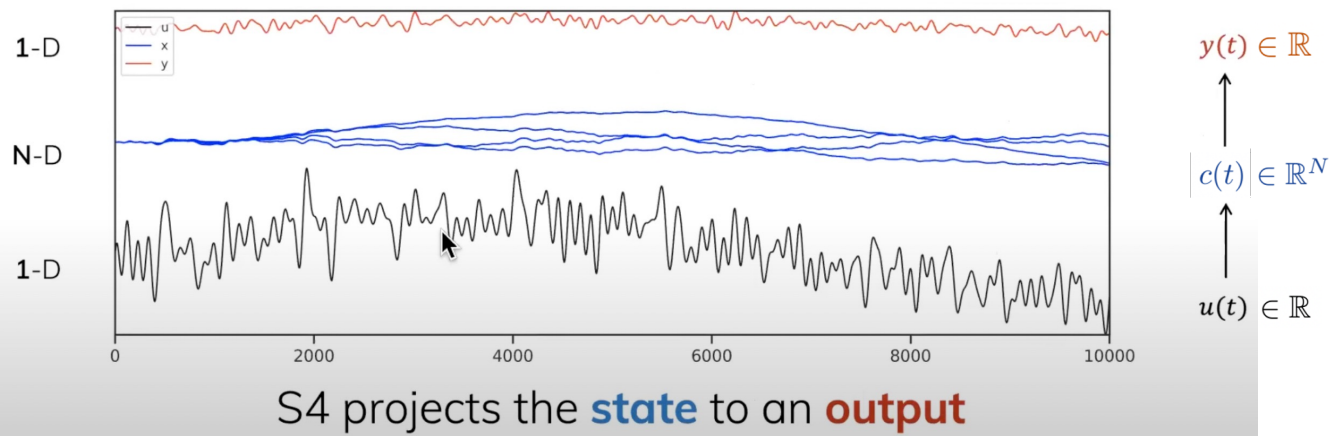
$$\begin{aligned} c'(t) &= \mathbf{A}c(t) + \mathbf{B}u(t) \\ y(t) &= \mathbf{C}c(t) + \mathbf{D}u(t) \end{aligned}$$

Récurrance cas discret

$$\begin{aligned} c_k &= \bar{\mathbf{A}}c_{k-1} + \bar{\mathbf{B}}u_k \\ y_k &= \bar{\mathbf{C}}c_k + \bar{\mathbf{D}}u_k \end{aligned}$$

Motivation : (Théorème) pour toute mesure $\mu^{(t)}$ raisonnable, l'état $c(t)$ vérifie une équation $c'(t) = \mathbf{A}c(t) + \mathbf{B}u(t)$

$$y(t) = \mathbf{C}c(t) + \mathbf{D}u(t)$$



Idée : apprendre implicitement la mesure $\mu^{(t)}$ via \mathbf{A} et \mathbf{B} plutôt que la prescrire.

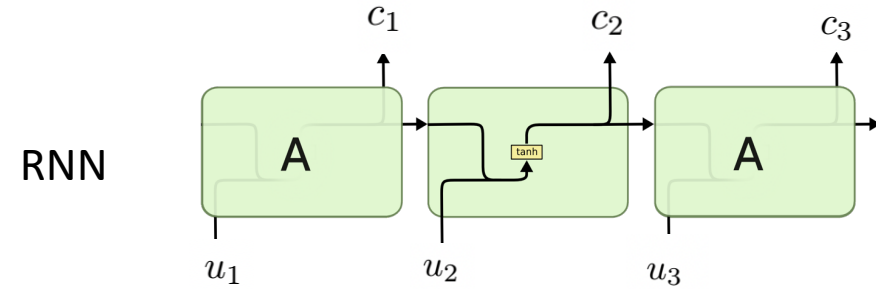
Difficulté : en pratique il est computationnellement impossible d'apprendre \mathbf{A} , en tous cas naïvement.

SSM : les vues récurrentes et convolutionnelles

Représentation récurrente

$$\begin{aligned} \text{SSM} \quad c_k &= \overline{A}c_{k-1} + \overline{B}u_k \\ y_k &= \overline{C}c_k + \overline{D}u_k \end{aligned}$$

Mapping seq2seq **linéaire** $(u_k)_{l \in \mathbb{N}} \rightarrow (c_k)_{l \in \mathbb{N}}$



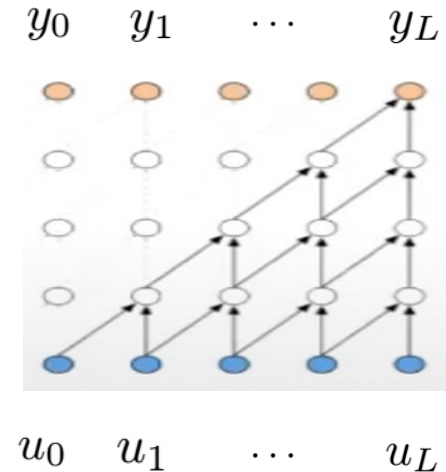
Mapping seq2seq **non-linéaire** $(u_k)_{l \in \mathbb{N}} \rightarrow (c_k)_{l \in \mathbb{N}}$

Efficace en prédiction, coût constant mais **non parallélisable**

Représentation convolutionnelle (on déroule le SSM)

$$\begin{aligned} c_0 &= \overline{B}u_0 & c_1 &= \overline{A}\overline{B}u_0 + \overline{B}u_1 & c_2 &= \overline{A}^2\overline{B}u_0 + \overline{A}\overline{B}u_1 + \overline{B}u_2 & \dots \\ y_0 &= \overline{C}\overline{B}u_0 & y_1 &= \overline{C}\overline{A}\overline{B}u_0 + \overline{C}\overline{B}u_1 & y_2 &= \overline{C}\overline{A}^2\overline{B}u_0 + \overline{C}\overline{A}\overline{B}u_1 + \overline{C}\overline{B}u_2 & \dots \\ & & & & y_k &= \overline{C}\overline{A}^k\overline{B}u_0 + \overline{C}\overline{A}^{k-1}\overline{B}u_1 + \dots + \overline{C}\overline{A}\overline{B}u_{k-1} + \overline{C}\overline{B}u_k \end{aligned}$$

$$y = \overline{K} * u \quad \overline{K} := (\overline{C}\overline{B}, \overline{C}\overline{A}\overline{B}, \dots, \overline{C}\overline{A}^{L-1}\overline{B}) \in \mathbb{R}^L \quad \text{un très gros filtre CNN}$$



Efficace à l'entraînement (futur connu) mais avec **un coût qui croit avec la taille du contexte**

SSM pour \mathbf{A} quelconque est inapprenable en pratique!

La précision des prédictions du modèle SSM dépend énormément de \mathbf{A}

Un SSM avec une initialisation aléatoire de $\mathbf{A} \in \mathbb{R}^{N \times N}$ est très **peu performant** :
 Sur un benchmark classique (*sequential MNIST*) : précision misérable de **60%**

Mais avec \mathbf{A} initialisé avec la **matrice HiPPO** : précision SOTA de **98%** !

Le coût d'entraînement d'un modèle SSM est prohibitif si l'on procède naïvement

Chaque multiplication par \mathbf{A} coute $O(N^2)$. Pour une séquence de longueur L on a :

- Un coût de calcul $O(LN^2)$ pour calculer le filtre $\overline{\mathbf{K}} := (\overline{\mathbf{CB}}, \overline{\mathbf{CAB}}, \dots, \overline{\mathbf{CA}^{L-1}\mathbf{B}}) \in \mathbb{R}^L$
- Une empreinte mémoire de $O(LN)$ pour stocker $(c_k)_{k \in [L]}$.

S4 = SSM + HiPPO + algo pour matrices \mathbf{A} de type NPLR

S4 = Structured State Space Sequence Model

Théorème 1 : toutes les matrices HIPPO issues des mesures classiques $\mu^{(t)}$ sont **NPLR** (**N**ormal **P**lus **L**ow **R**ank).

$$\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\dagger + \mathbf{P}\mathbf{Q}^\top, \quad \mathbf{\Lambda} = \text{diag}[\lambda_1, \dots, \lambda_N] \in \mathbb{C}^N$$

$$\mathbf{U} \in \mathbb{C}^{N \times N}, \quad \mathbf{U}\mathbf{U}^\dagger = \mathbf{U}^\dagger\mathbf{U} = \mathbb{I}$$

$$\mathbf{P}, \mathbf{Q} \in \mathbb{R}^{N \times r}, \quad r = 1, 2$$

Théorème 2 (S4 récurrent) : un pas de récurrence coûte $O(N)$au lieu de $O(N^2)$ naïvement !

Théorème 3 (S4 convolution) : le filtre $\overline{\mathbf{K}}$ demande $\tilde{O}(N + L)$ opérations et un espace mémoire de $O(N + L)$.

S4 : intégration dans une architecture de RN

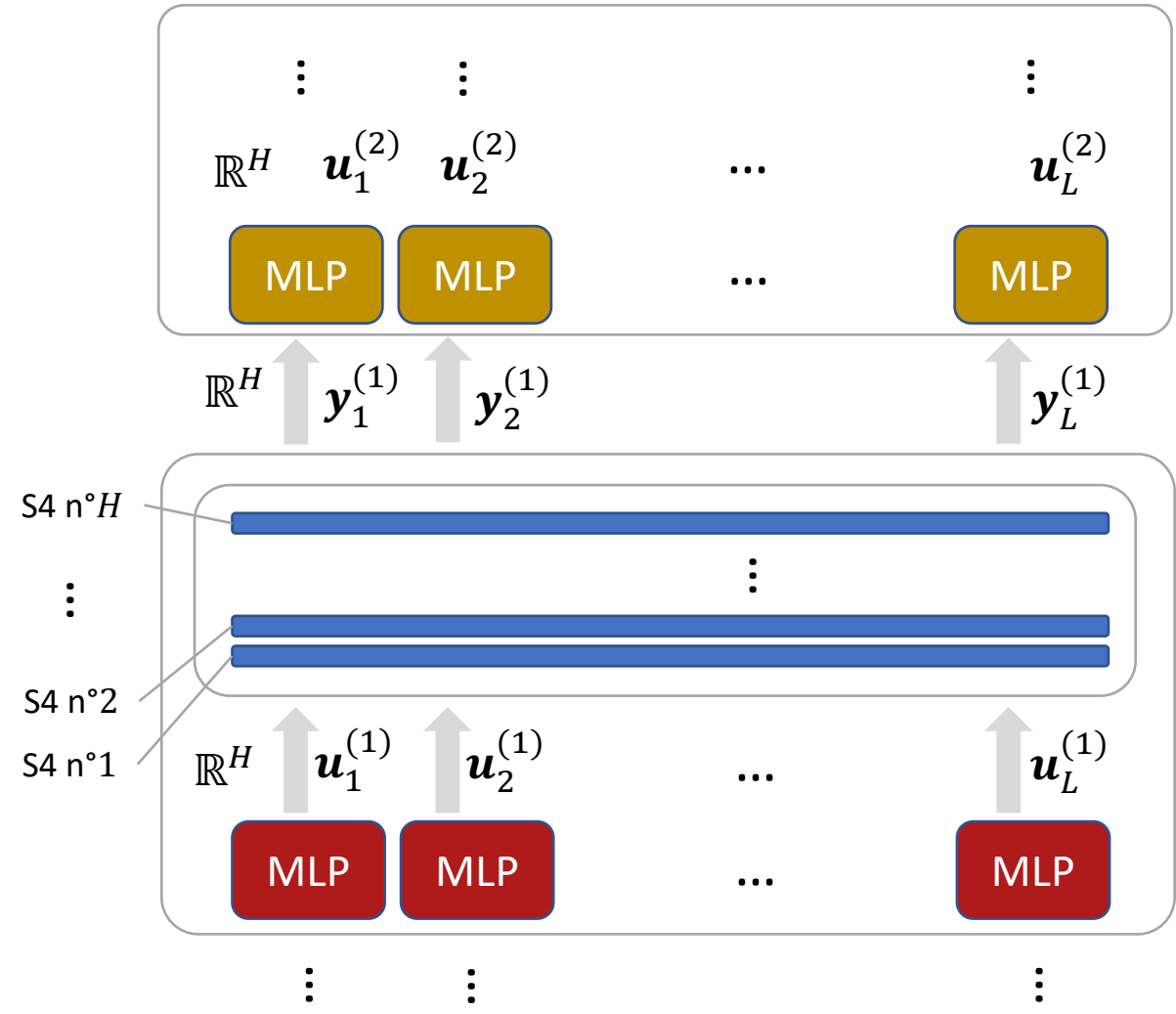
Les paramètres de S4 :

$$A = U\Lambda U^\dagger + PQ^\top = U[\Lambda - \underbrace{(U^\dagger P)}_{\tilde{P}} \underbrace{(U^\dagger Q)^\dagger}_{\tilde{Q}^\dagger}]U^\dagger$$

$$\Lambda = \text{diag}[\lambda_1, \dots, \lambda_N] \in \mathbb{C}^N$$

$$B \in \mathbb{C}^N, \quad C \in \mathbb{C}^N, \quad \tilde{P} \in \mathbb{C}^N, \quad \tilde{Q} \in \mathbb{C}^N. \quad \left. \vphantom{B, C, \tilde{P}, \tilde{Q}} \right\} 5N \text{ params}$$

Si on veut un mapping entre des séquences de vecteurs $\in \mathbb{R}^H$ plutôt que des séquences de nombres on instancie simplement H modèles S4 indépendants !



S4 : à retenir

1. Un modèle de séquence **générique** (= ne demande aucun biais inductif, aucun ajustement d'architecture ou aucun pré-traitement) applicable à tous types des séquences (textes, images, audio, séries temporelles, ADN,...)
2. S4 : SSM + HiPPO + algo optimisé tirant parti de la structure NPLR des matrices A « utiles »
Les matrices « utiles » sont précisément celles pour lesquelles le calcul est faisable !
3. Une approche du problème des LRD **solidement motivée théoriquement** : « *No dirty tricks!* »
4. Repose sur la paramétrisation judicieuse NPLR des matrices A d'un SSM suggérées par HiPPO
5. Très **performant en pratique** grâce à :
 - un entraînement efficace avec la vue CNN parallélisable $O(L + N)$
 - une prédiction rapide avec la vue RNN indépendante de la longueur L de la séquence
 - la robustesse vis-à-vis des fréquences d'échantillonnages (invariance d'échelle)
6. Une alternative rafraichissante aux Transformers !

Questions et discussion : S4 pour le TAL ?

