

## Lecture 17

### • State space models Continued

→ How to support long-range dependancies?

$$\vec{h}_{t+1} = A \vec{h}_t + B \vec{x}_t$$

$$\vec{y}_t = C \cdot \vec{x}_t$$

} All information about the past can only be captured in  $\vec{h}_t$ . It has to fit.

Efficiency versus Efficacy tradeoff.

Larger  $\vec{h}_t \Rightarrow$  more memory/history tracking  
but also, more compute.

So want to make  $\vec{h}_t$  as large as possible, while still being able to run

---

S4

Efficiency via FFT - use of LTI systems.

What about long range dependancy?

Let  $A = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_N \end{bmatrix}$   $\vec{h}_t$  is  $N$  dimensional.

How is  $\lambda$  related to system stability?

$|\lambda| > 1 \rightarrow$  System explodes  $\rightarrow$  bad.

$|\lambda| < 1 \rightarrow$  history dies out fast.

$|\lambda| = 1 \rightarrow$  critical damping.

For  $\lambda \in \mathbb{R}$ ,  $|\lambda| = 1 \Rightarrow \lambda = \pm 1$ .

→ Only two choices!

→ Far apart from each other (So not GD friendly)

Solution: Complex eigenvalues!

$$\lambda = \exp(-\text{ReLU}(\lambda_{\mathbb{R}}) + j \lambda_{\mathbb{I}})$$

$$j = \sqrt{-1}$$

- Exact initialization matters, have to be mostly on or very close to unit circle with some inside unit circle
- Hippo-based initialization - like analytic pre-training.
- Freezing A with good initialization and only training the rest also seemed to work reasonably!

Table 4: (**Long Range Arena**) (*Top*) Original Transformer variants in LRA. Full results in Appendix D.2. (*Bottom*) Other models reported in the literature. *Please read Appendix D.5 before citing this table.*

MODEL	LISTOPS	TEXT	RETRIEVAL	IMAGE	PATHFINDER	PATH-X	AVG
Transformer	36.37	64.27	57.46	42.44	71.40	<b>X</b>	53.66
Reformer	<u>37.27</u>	56.10	53.40	38.07	68.50	<b>X</b>	50.56
BigBird	36.05	64.02	59.29	40.83	74.87	<b>X</b>	54.17
Linear Trans.	16.13	<u>65.90</u>	53.09	42.34	75.30	<b>X</b>	50.46
Performer	18.01	65.40	53.82	42.77	77.05	<b>X</b>	51.18
FNet	35.33	65.11	59.61	38.67	<u>77.80</u>	<b>X</b>	54.42
Nyströmformer	37.15	65.52	<u>79.56</u>	41.58	70.94	<b>X</b>	57.46
Luna-256	37.25	64.57	79.29	<u>47.38</u>	77.72	<b>X</b>	<u>59.37</u>
<b>S4</b>	<b>59.60</b>	<b>86.82</b>	<b>90.90</b>	<b>88.65</b>	<b>94.20</b>	<b>96.35</b>	<b>86.09</b>

Outperformed Transformer models on long-range benchmarks!

But didn't quite match/beat transformers on other tasks.

# Mamba / S6.

## State-space models

- Recurrent view  
(state-based)

"inference time"

- Input / Output view  
(convolution based)

"training time"

- Continuous-time view

"selectivity) gating!"

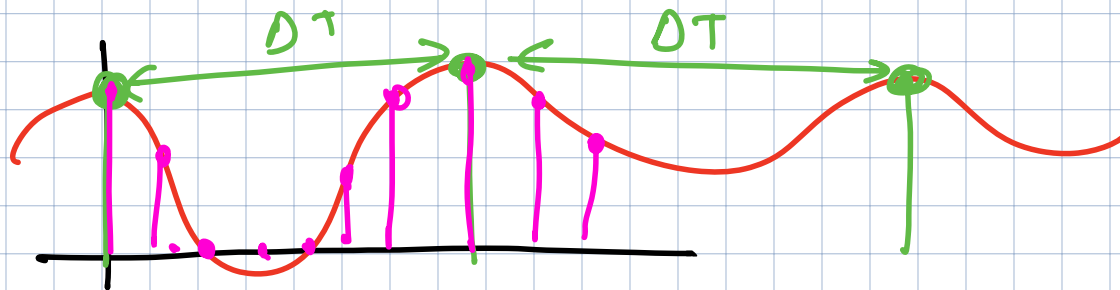
Discrete-time system:

$$\vec{h}_{t+1} = A \vec{h}_t + B \vec{u}_t$$
$$\vec{y}_{t+1} = C \cdot \vec{h}_t$$

Think of as a discretization of a continuous time system.  
(sampling)

$$\frac{d}{dt} \vec{h}(t) = A \vec{h}(t) + B \vec{u}(t) \quad \leftarrow \text{(Different } A, B)$$
$$\vec{y}(t) = C \cdot \vec{h}(t)$$

To convert from CT to DT, we need a sampling interval  $\Delta T$ .



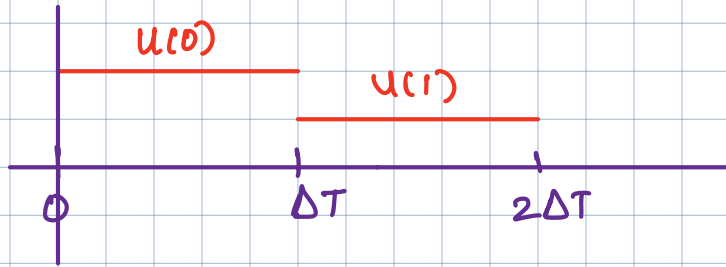
Scalar case:

$$\dot{x}(t) = a x(t) + b \cdot u(t).$$

Sampling time  $\Delta T$ .

Assume  $u(t)$  is piecewise constant for time  $\Delta T$ .

For small  $\Delta T$  this is reasonable.



Aside:

$a > 0 \Rightarrow$  explode

$a < 0 \Rightarrow$  decay to 0.

$$x(\Delta T) = \underbrace{e^{a \Delta T} \cdot x(0)}_{\text{exponential growth.}} + \underbrace{\frac{b}{a} (e^{a \Delta T} - 1)}_{\text{impact of input}} u(0).$$

Assume  $a < 0$

Then  $\Delta T$  close to 0  $\Rightarrow e^{a \Delta T} \approx 1 \Rightarrow$  <sup>Initial state?</sup> Remember history.

$\Delta T$  large  $\Rightarrow e^{a \Delta T}$  decaying exponential  $\Rightarrow$  forget past

$$\Delta T \text{ smallish} \Rightarrow e^{a\Delta T} - 1 \approx \underbrace{1 + a\Delta T}_{\text{Taylor expansion}} - 1 \approx a\Delta T.$$

$$\Rightarrow \frac{b}{a} (e^{a\Delta T} - 1) = b\Delta T.$$

$$\Delta T \text{ close to } 0 \Rightarrow b\Delta T \propto 0. \Rightarrow \begin{array}{l} \text{Effect of input?} \\ \text{input doesn't matter} \end{array}$$

$$\Delta T \text{ large} \Rightarrow b\Delta T \text{ large} \Rightarrow \text{input matters a lot.}$$

So by changing  $\Delta T$ , can change the relative importance of the 'fast'.

Idea: Choose  $\Delta T$  based on input (i.e. learn it).

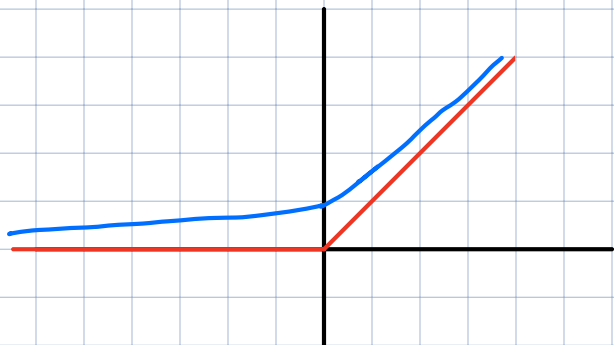
Try: Linear function of  $\vec{u}$ ?

$$\langle \vec{w}, \vec{u} \rangle$$

Try:  $\text{ReLU}(\langle \vec{w}, \vec{u} \rangle + b)$

Softplus:

$$\text{Softplus}(x) = \ln(1 + e^x)$$



if  $a=b=N=1$

$$g_t = \sigma(\text{Linear}(\vec{x}_t))$$

$$\vec{h}_t = g_t \vec{h}_{t-1} + (1-g_t) \vec{x}_t$$

So this gives us for each channel of input:

$$\vec{h}_{k+1} = A_k \vec{h}_k + B_k x_k$$

$$\vec{y}_{k+1} = C \vec{h}_k$$

$$A_k = \begin{bmatrix} e^{\lambda_1 \Delta T} & & \\ & e^{\lambda_2 \Delta T} & \\ & & \ddots \end{bmatrix}$$

But... have we lost anything?

$$B_k = \vec{b}_k = \vec{b} \Delta T.$$

$$A^k \rightarrow A_k A_{k-1} \dots A_1 A_0$$



Initialization:

$\lambda I$  blocks for  $A$ .

$$\lambda_n = -\frac{1}{2} + ni$$

$$\lambda_n = -\frac{1}{2}$$

$$\lambda_n = -(n+1)$$

## State-space models

- Recurrent view  
(state-based)

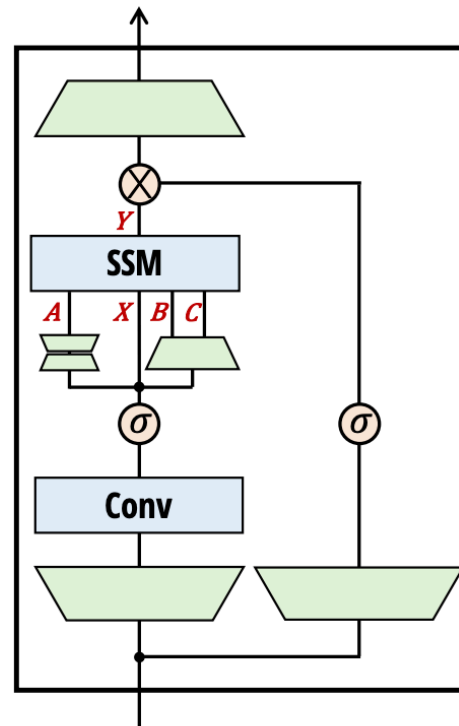
"inference time"

- Input / Output view  
(Convolution based)

"training time"

- Continuous-time view

"selectivity) gating"



**Sequential Mamba Block**