

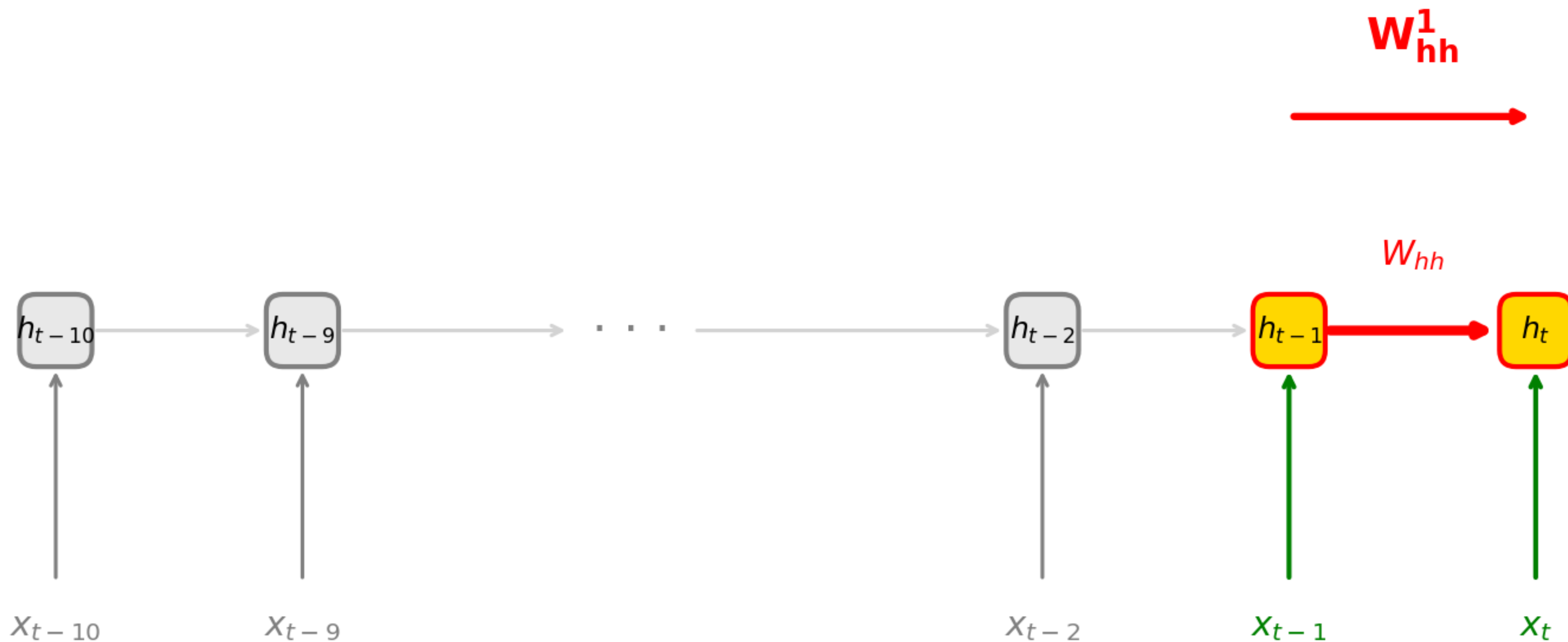
Deep Learning (1470)

Randall Balestriero

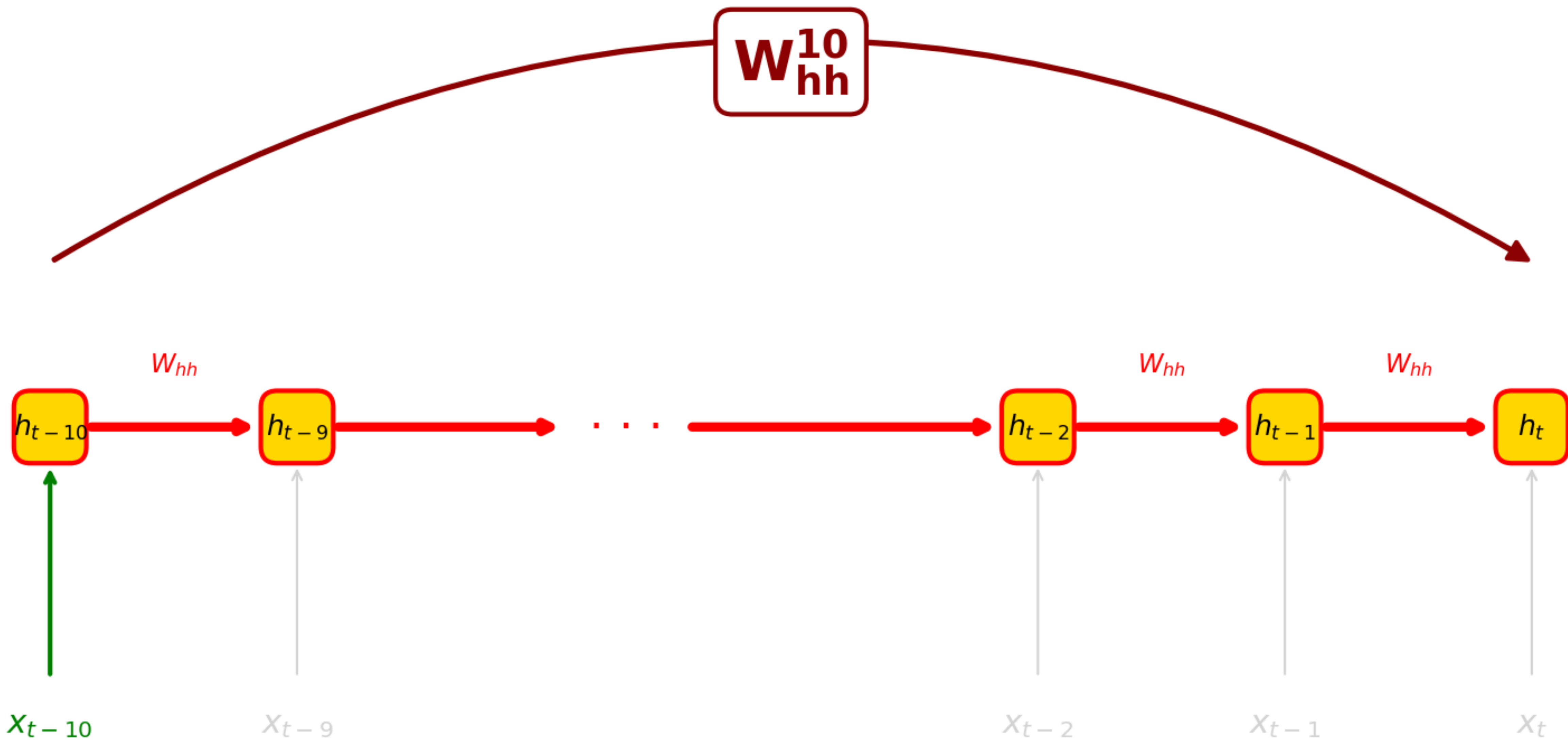
Class 12: seq2seq

Recap!

Information Flow from h_{t-1} to h_t
Contribution: $W_{hh}^1 \cdot h_{t-1}$



Information Flow from h_{t-10} to h_t
Contribution: $W_{hh}^{10} \cdot h_{t-10}$



Exponential Memory Loss: Example

$$\begin{bmatrix} 0.5 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 2 \end{bmatrix}$$

D

Exponential Memory Loss: Example

$$\begin{bmatrix} 0.5 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 2 \end{bmatrix} \longrightarrow \begin{bmatrix} 0.25 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 4 \end{bmatrix}$$

D **D²**

Exponential Memory Loss: Example

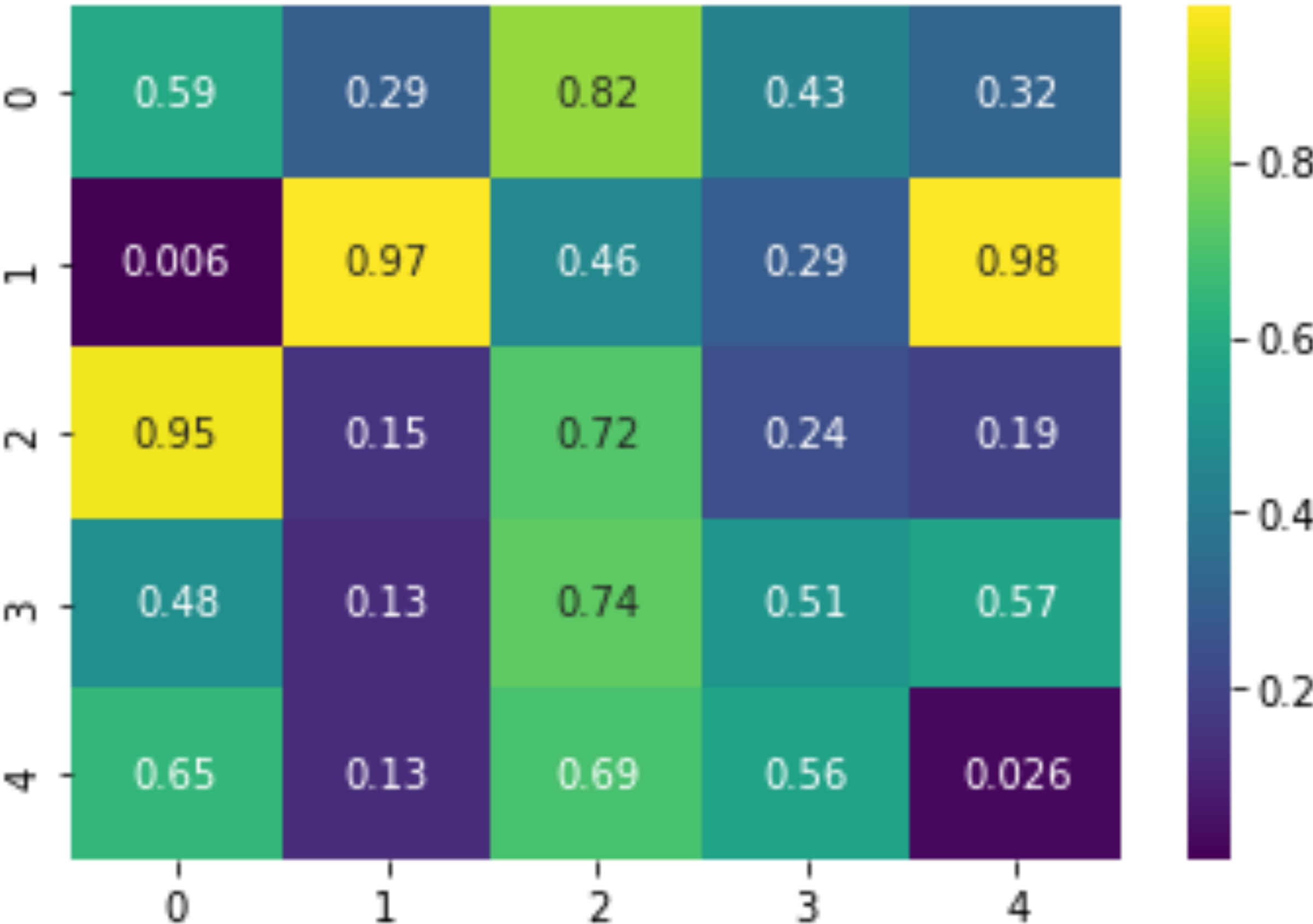
$$\begin{bmatrix} 0.5 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 2 \end{bmatrix} \longrightarrow \begin{bmatrix} 0.25 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 4 \end{bmatrix} \longrightarrow \begin{bmatrix} 0.125 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 8 \end{bmatrix}$$

$\mathbf{D} \qquad \qquad \mathbf{D}^2 \qquad \qquad \mathbf{D}^3$

Exponential Memory Loss: Example

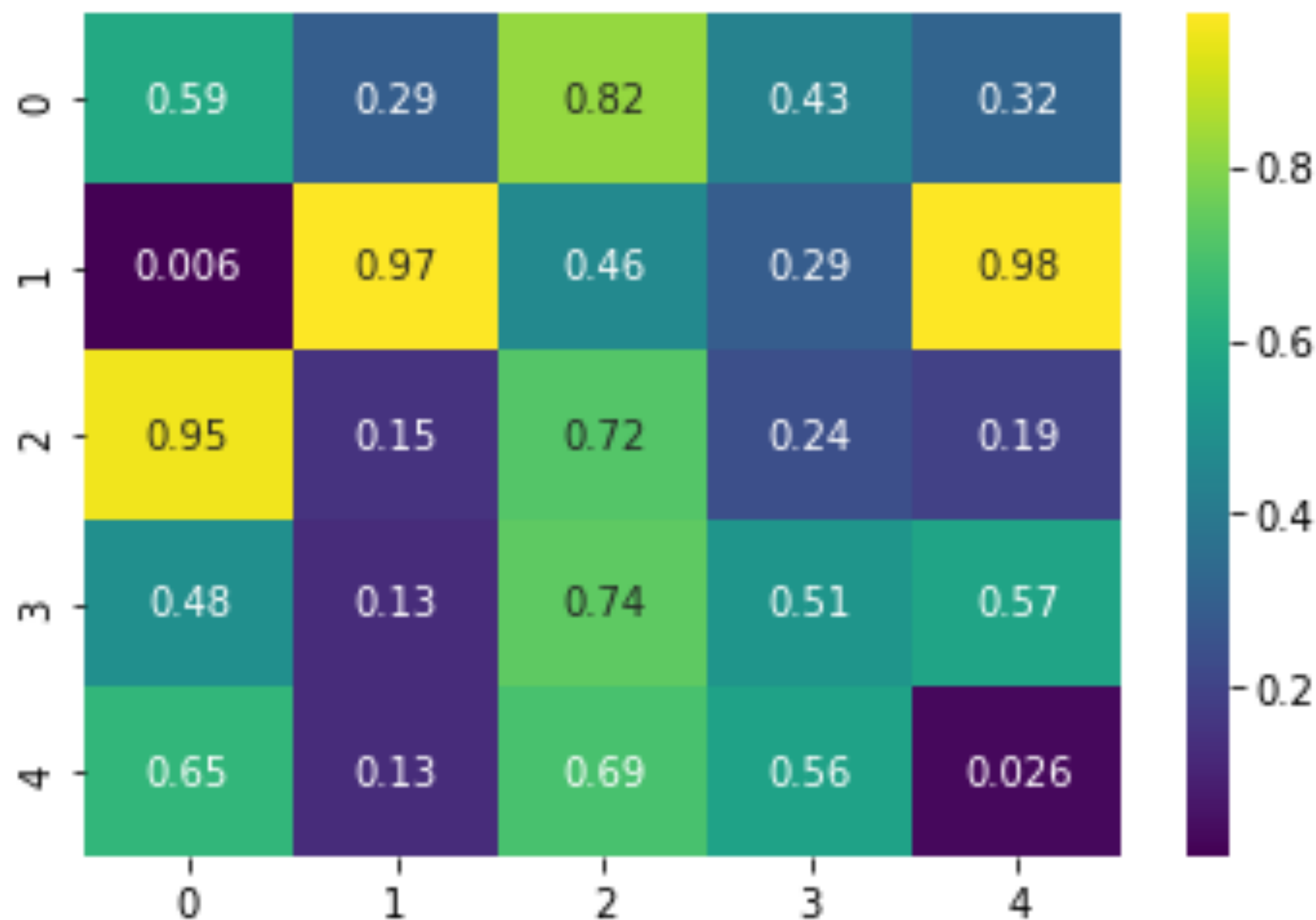
$$\begin{array}{ccccccc} \left[\begin{array}{ccc} 0.5 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 2 \end{array} \right] & \longrightarrow & \left[\begin{array}{ccc} 0.25 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 4 \end{array} \right] & \longrightarrow & \left[\begin{array}{ccc} 0.125 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 8 \end{array} \right] & \longrightarrow & \left[\begin{array}{ccc} 0.0625 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 16 \end{array} \right] \\ \mathbf{D} & & \mathbf{D}^2 & & \mathbf{D}^3 & & \mathbf{D}^4 \end{array}$$

Exponential Memory Loss: Example



What about now?

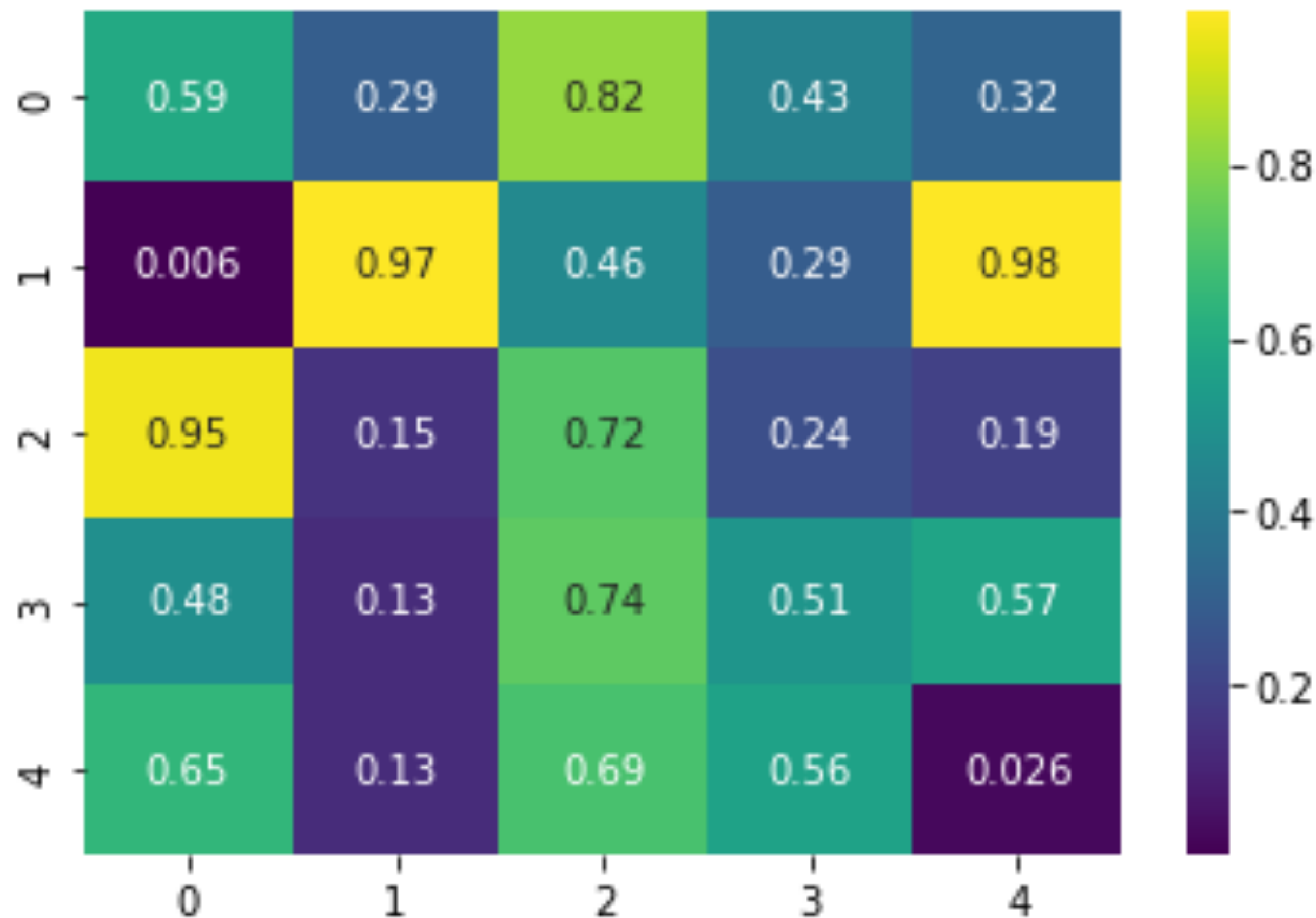
Exponential Memory Loss: Example



What about now?

$$W = U\Sigma V^{\top}$$

Exponential Memory Loss: Example

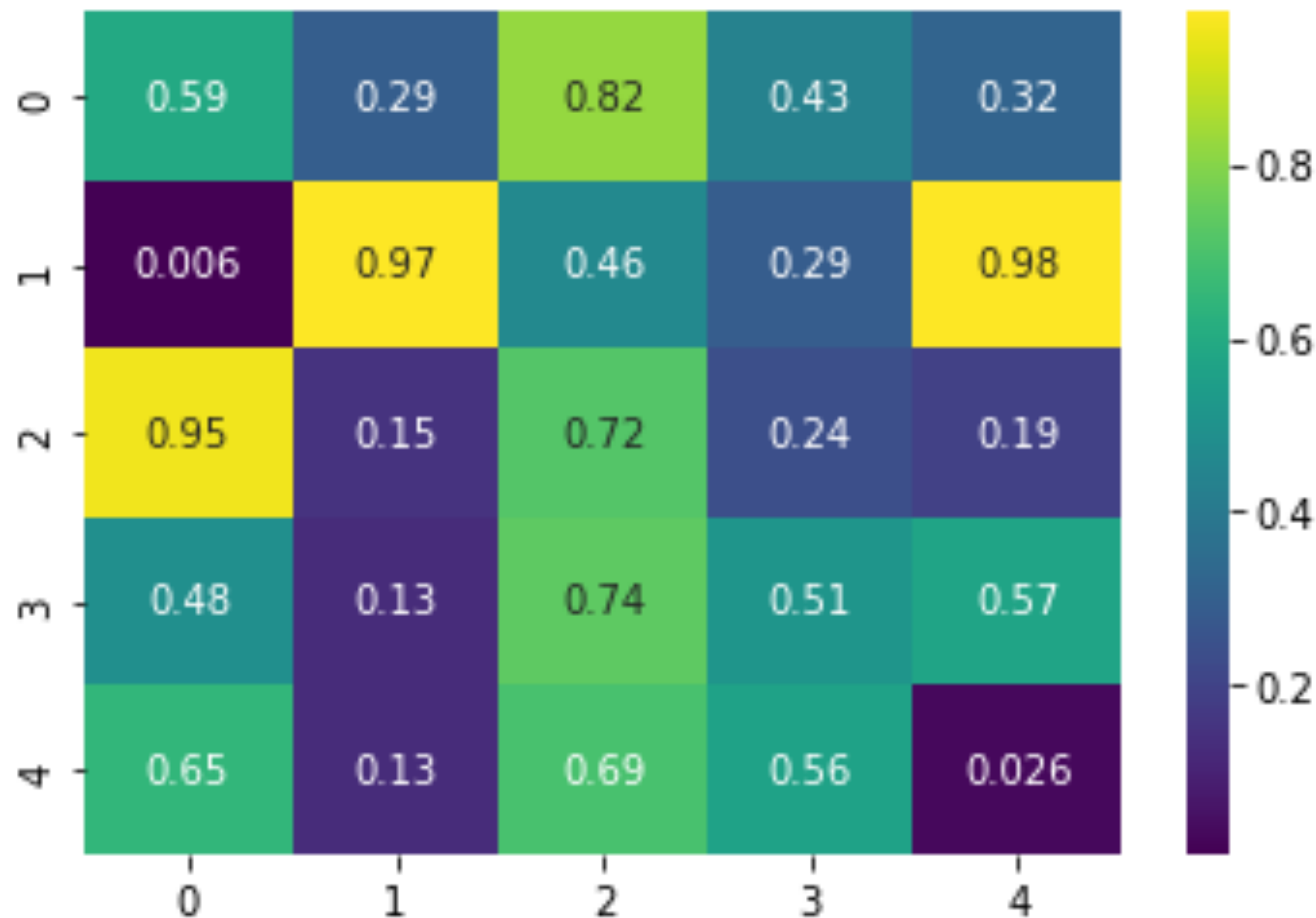


What about now?

$$W = U\Sigma V^{\top}$$

$$W^p = U\Sigma^p V^{\top}$$

Exponential Memory Loss: Example



What about now?

$$W = U\Sigma V^{\top}$$

$$W^p = U\Sigma^p V^{\top}$$

Can you come up with a condition to prevent explosion/collapse?

Beyond ORNNs/LSTMs/GRUs

- Instead of fixing the computation at time-step... we can?

Beyond ORNNs/LSTMs/GRUs

- Instead of fixing the computation at time-step... we can?

orthogonal like representation. We fix a pool size k , and then the update equations for this model are:

$$\begin{aligned}h_t &= \sigma(Ux_t + b) + Vh_{t-1} \\ y_t &= W_I h_t + W_P P_k(h_t)\end{aligned}\tag{6}$$

where if h is the kd dimensional vector $h = [h_1, \dots, h_{kd}]^T$, then $P(h)$ is the d dimensional vector defined by

$$P(h)_i = \sqrt{\sum_{j=k(i-1)+1}^{ki} h_j^2}$$

Recurrent Orthogonal Networks and Long-Memory Tasks

Mikael Henaff
New York University, Facebook AI Research

MBH305@NYU.EDU

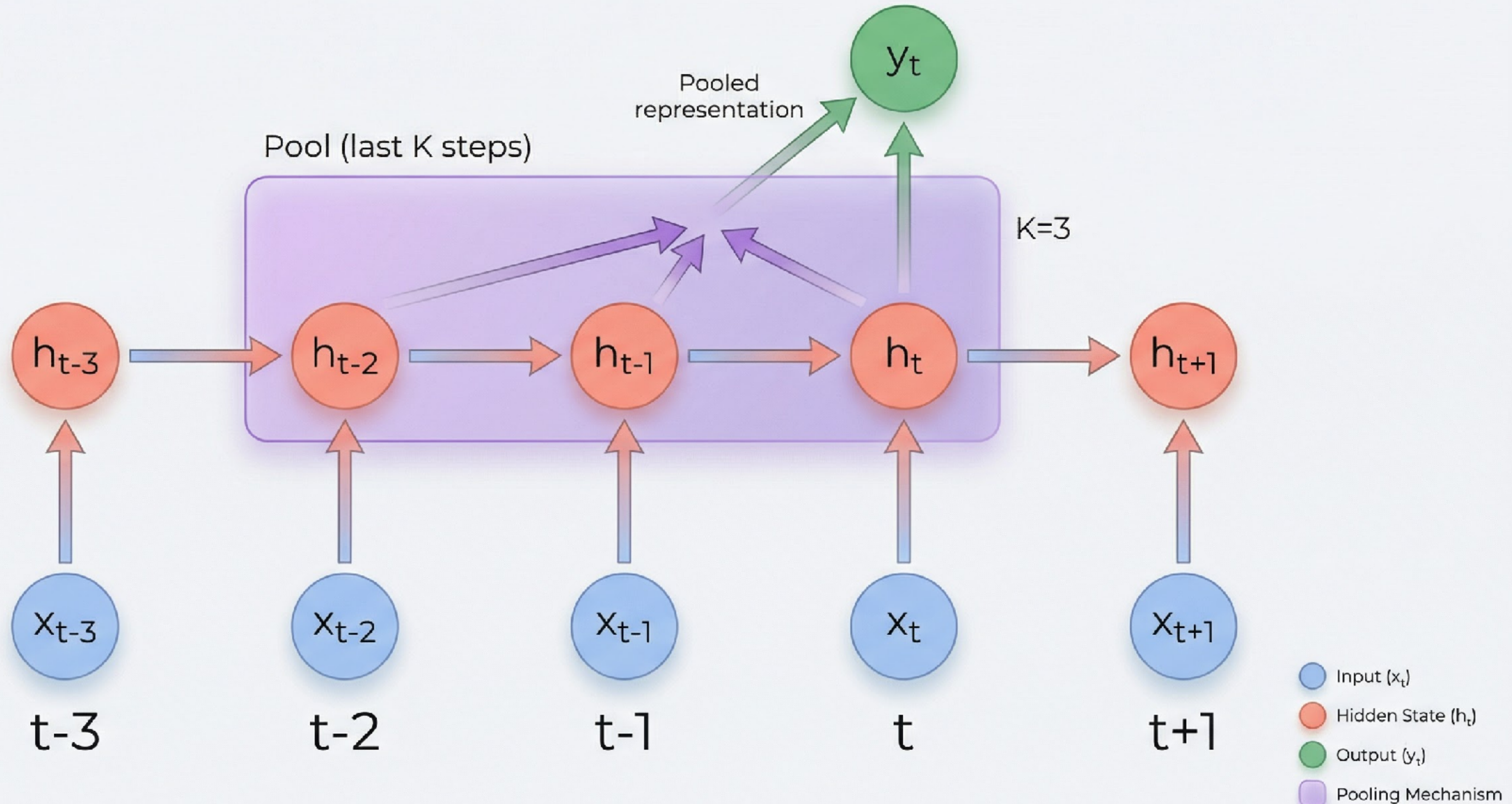
Arthur Szlam
Facebook AI Research

ASZLAM@FACEBOOK.COM

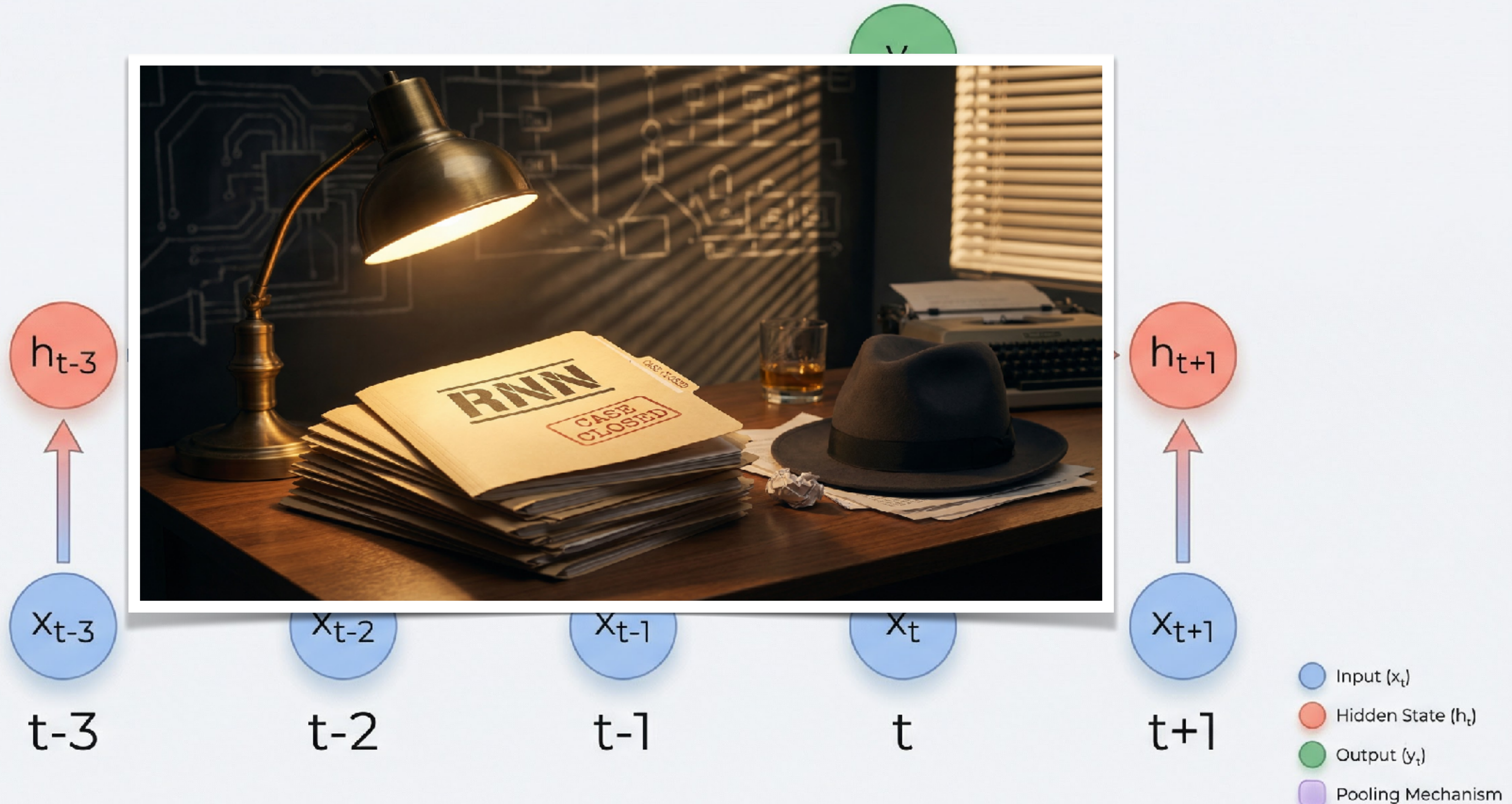
Yann LeCun
New York University, Facebook AI Research

YANN@CS.NYU.EDU

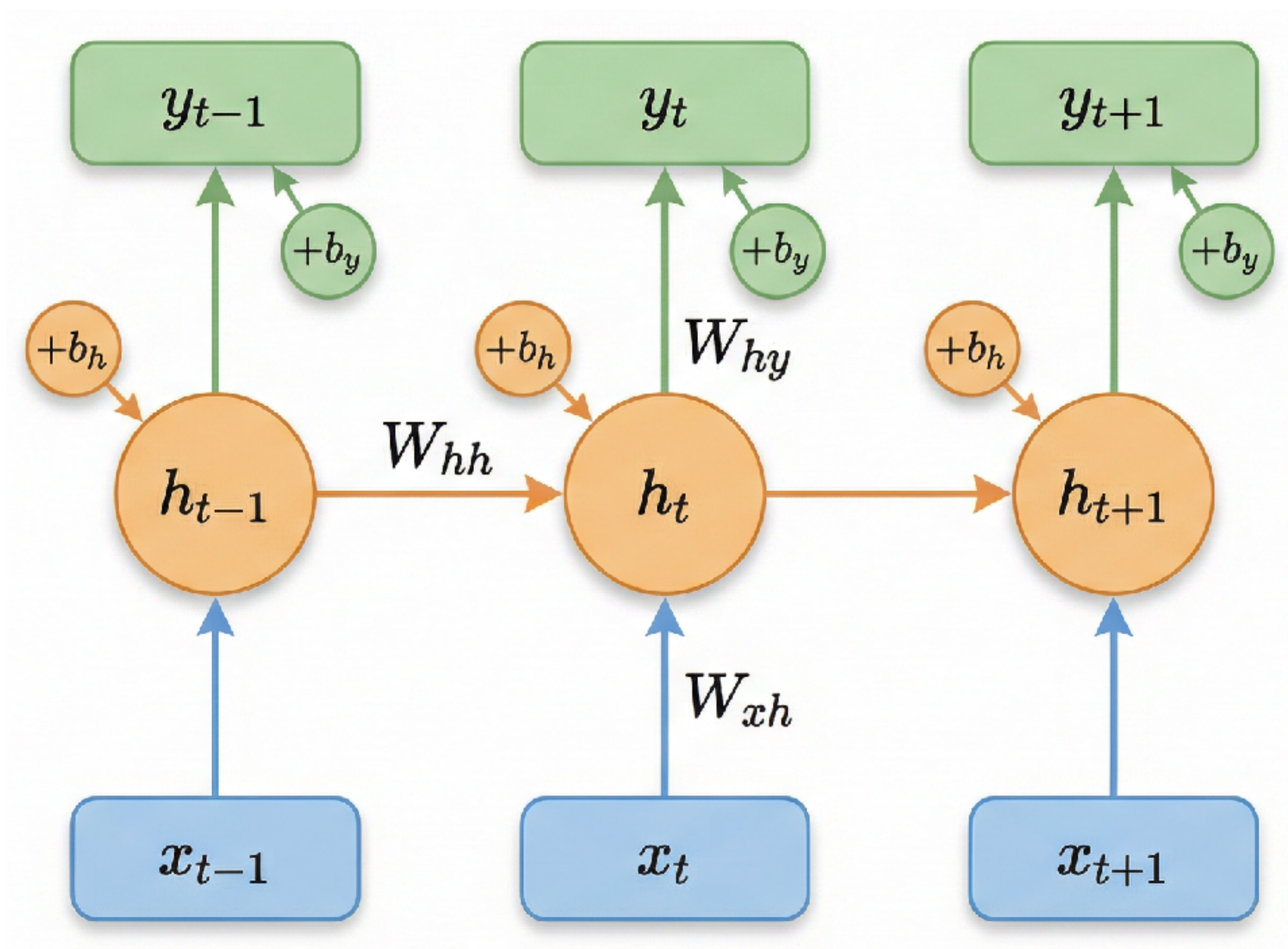
Recurrent Neural Network with Pooling Architecture



Recurrent Neural Network with Pooling Architecture



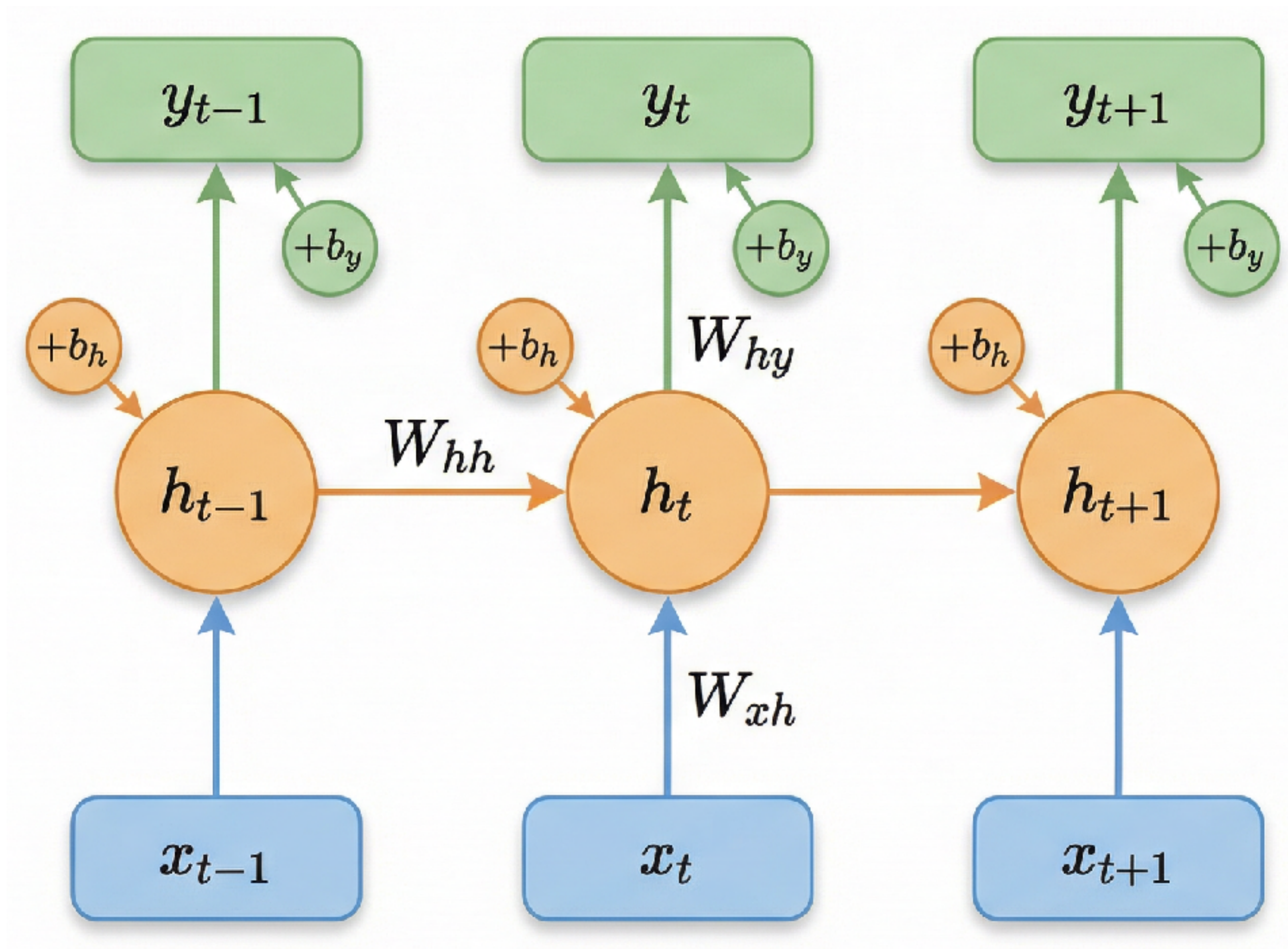
Back to Language Modeling



$$y_t = x_{t+1}$$

$$p(x_{t+1} | x_t, \dots, x_1) \approx f_{\theta}(x_t, h_{t-1})$$

Back to Language Modeling

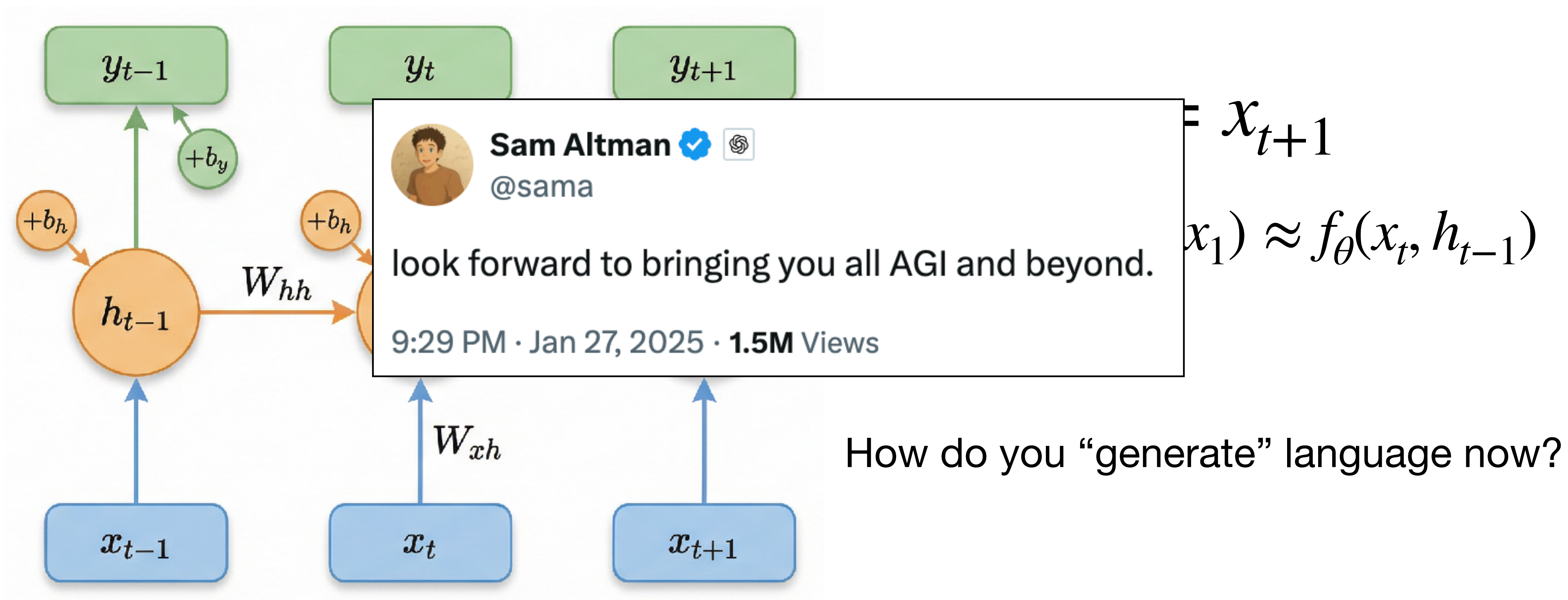


$$y_t = x_{t+1}$$

$$p(x_{t+1} | x_t, \dots, x_1) \approx f_{\theta}(x_t, h_{t-1})$$

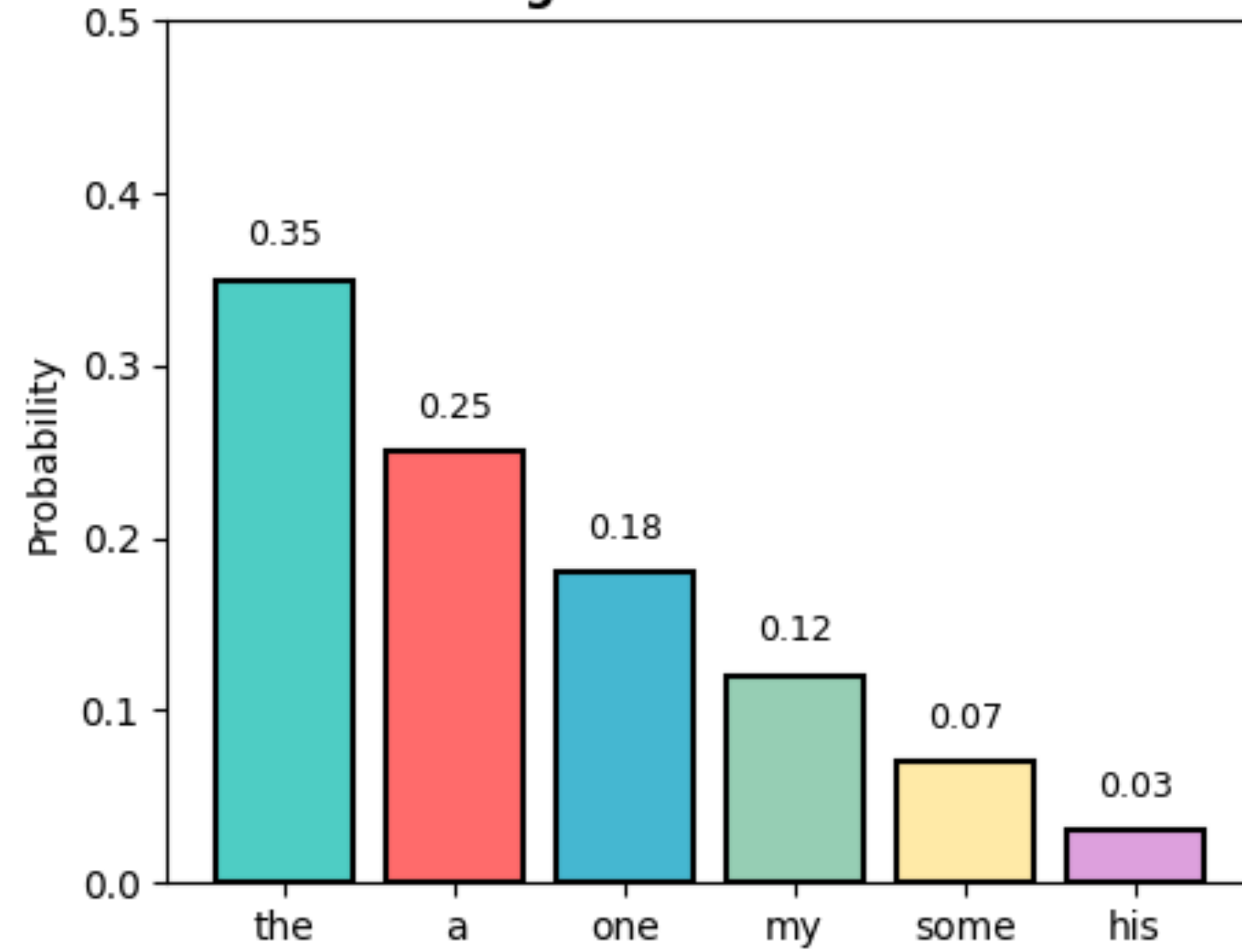
How do you “generate” language now?

Back to Language Modeling

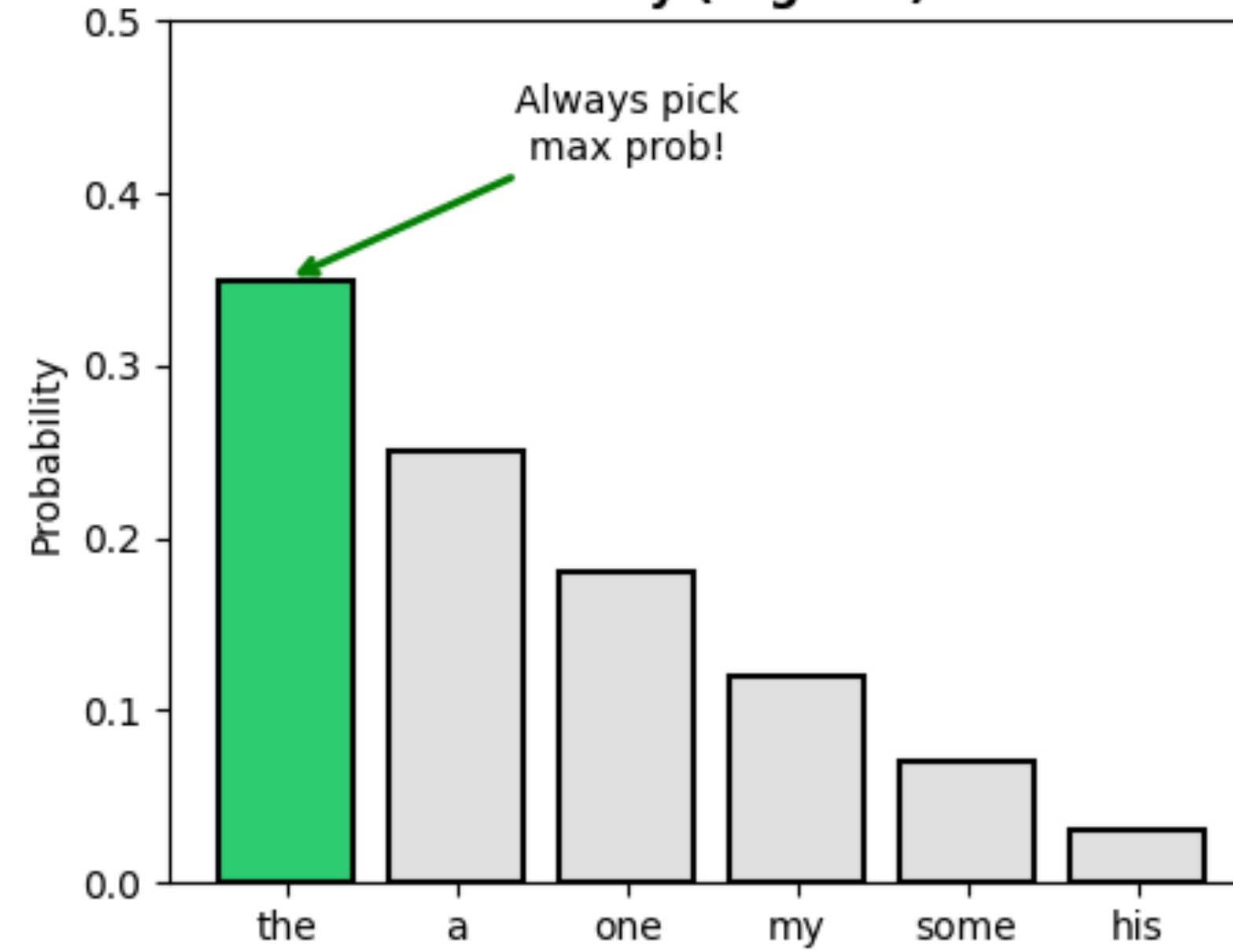


Sampling Strategies from Next-Token Distribution

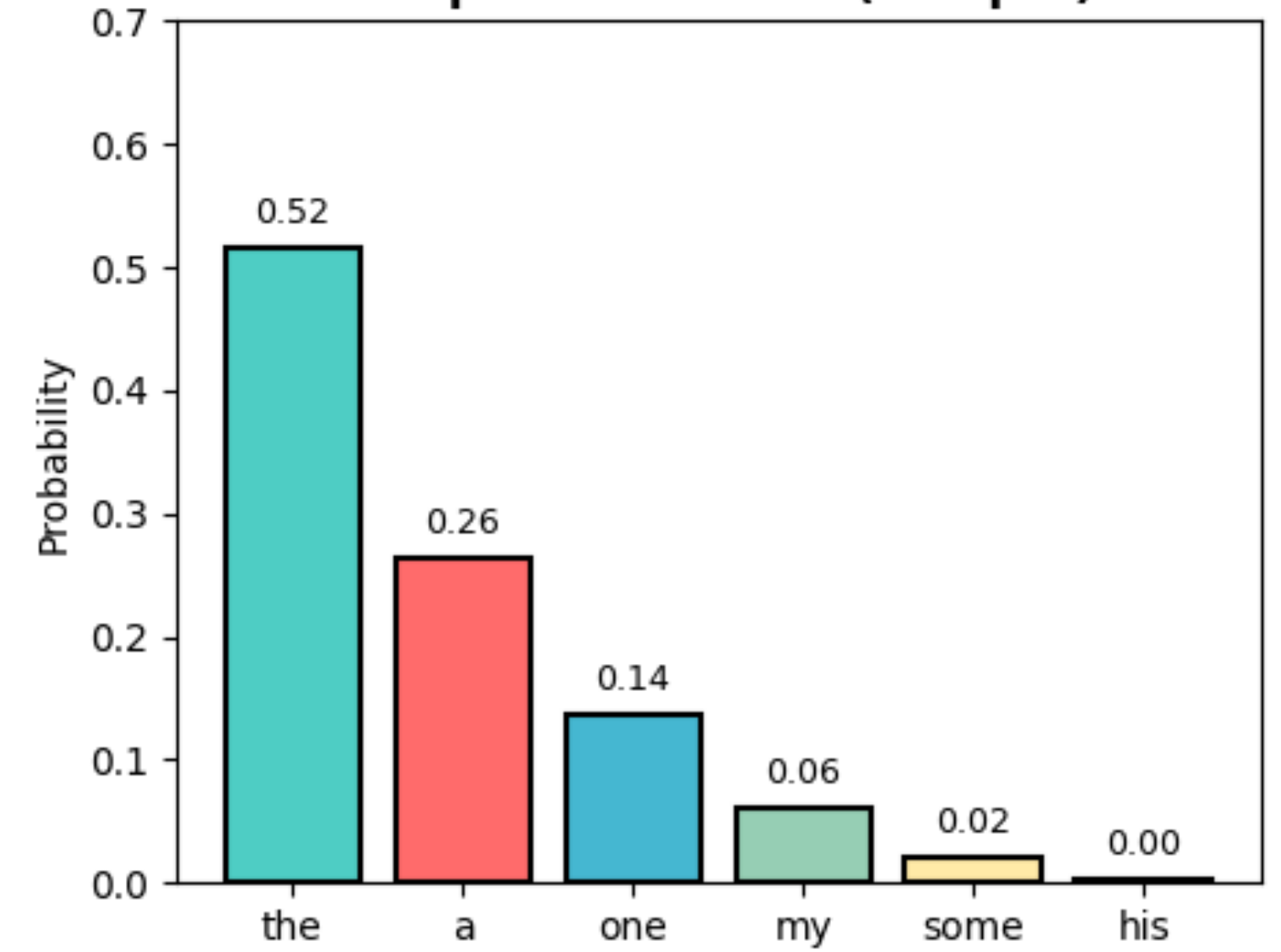
Original Distribution



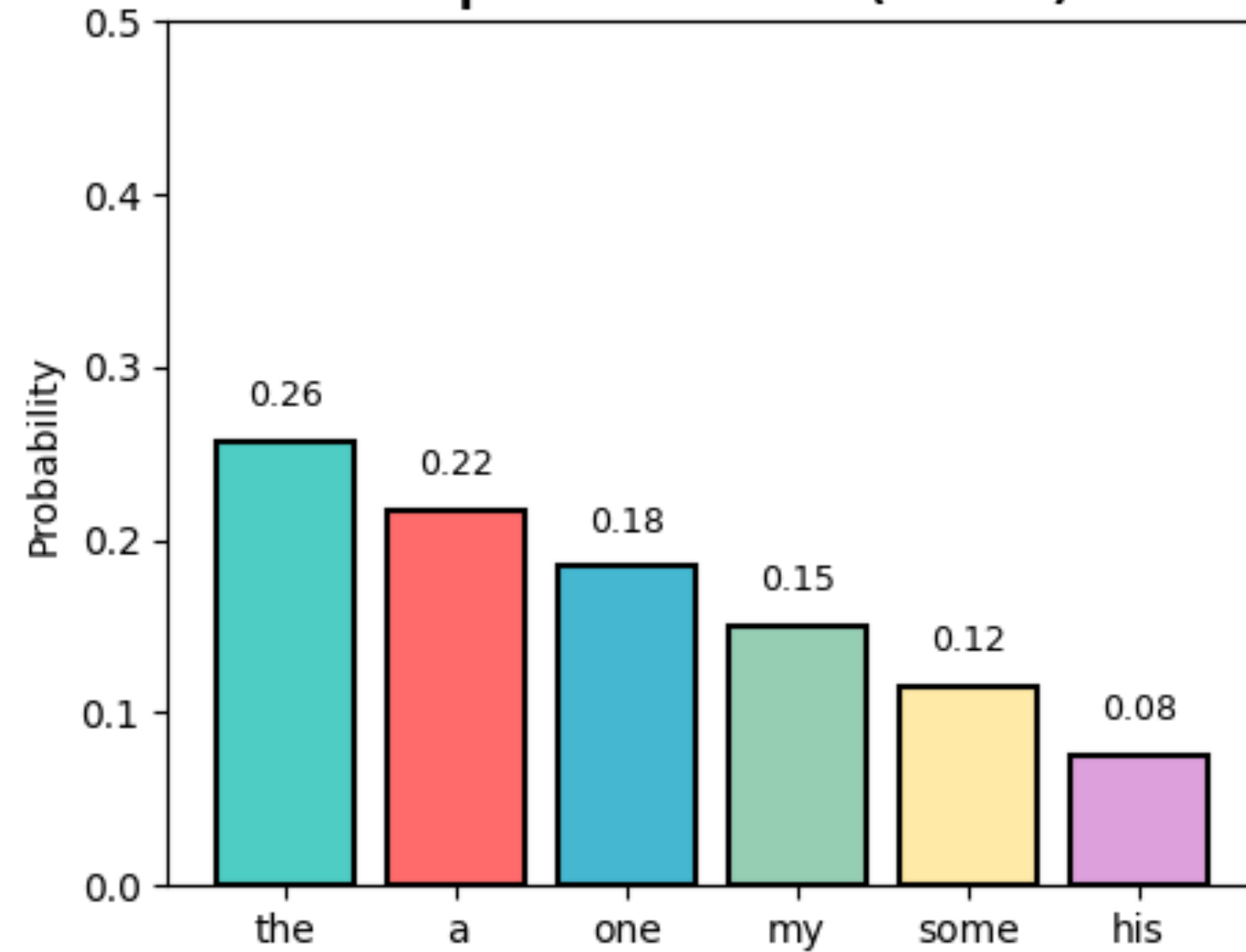
Greedy (argmax)



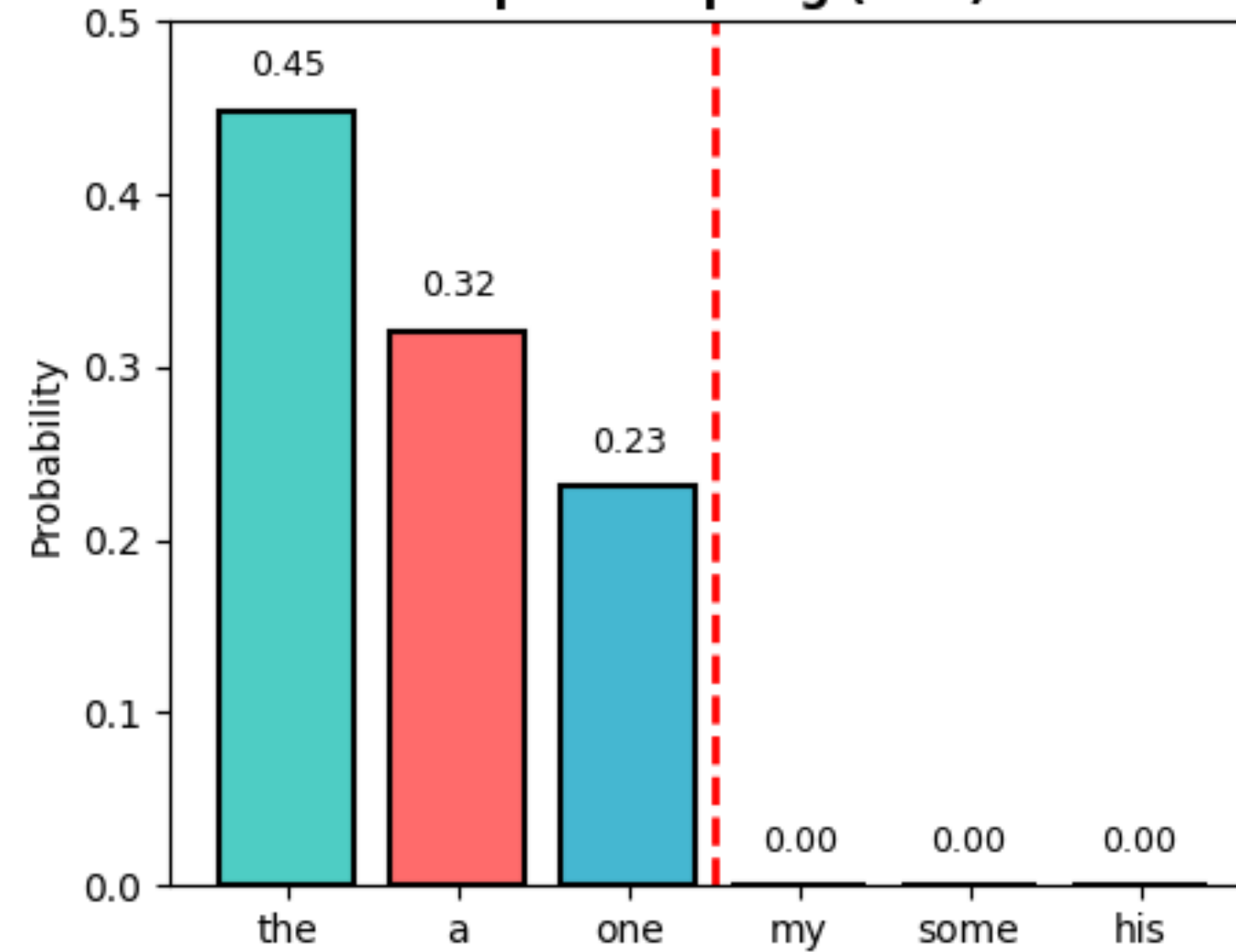
Temperature T=0.5 (sharper)



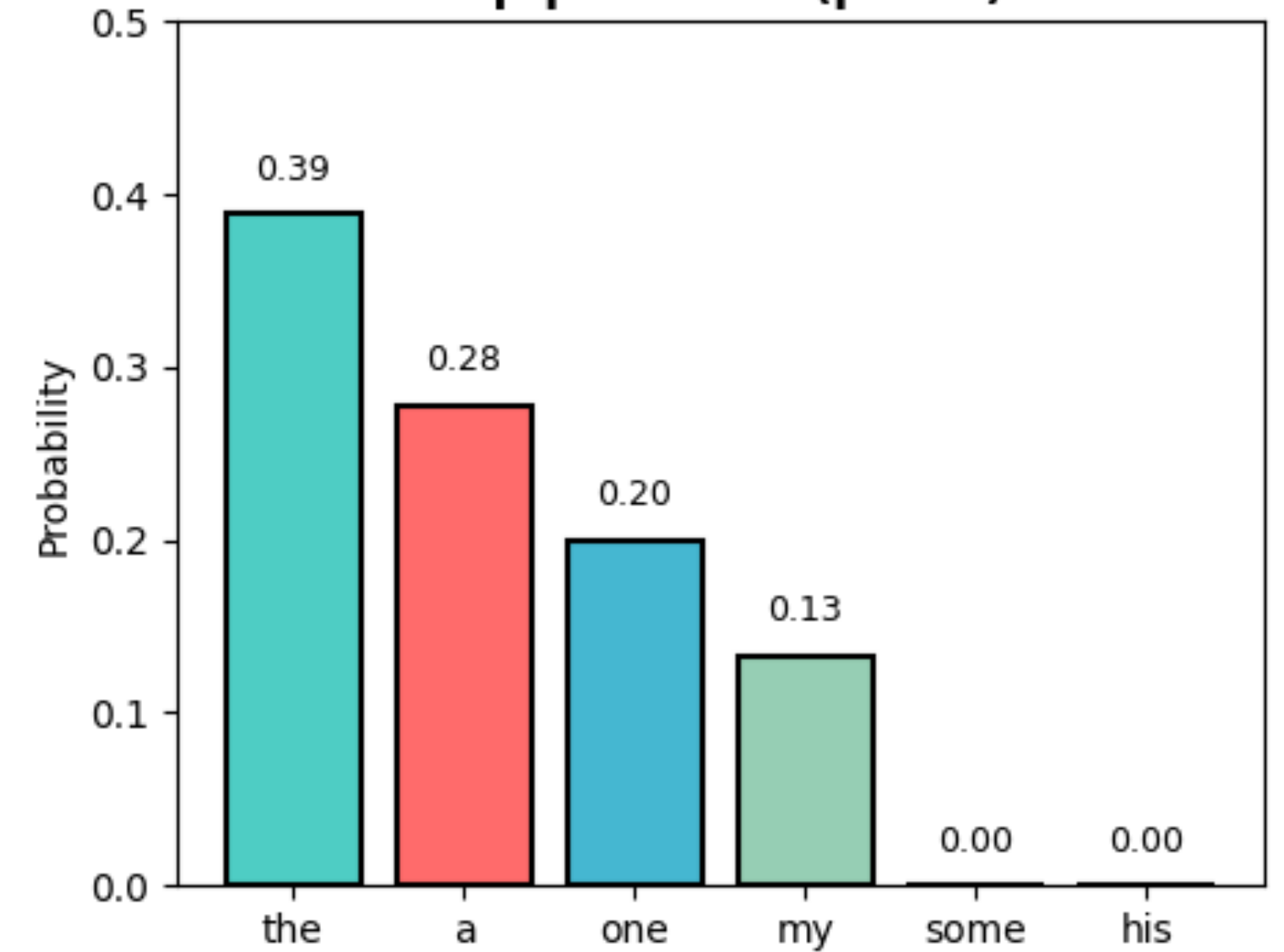
Temperature T=2.0 (flatter)



Top-k Sampling (k=3)

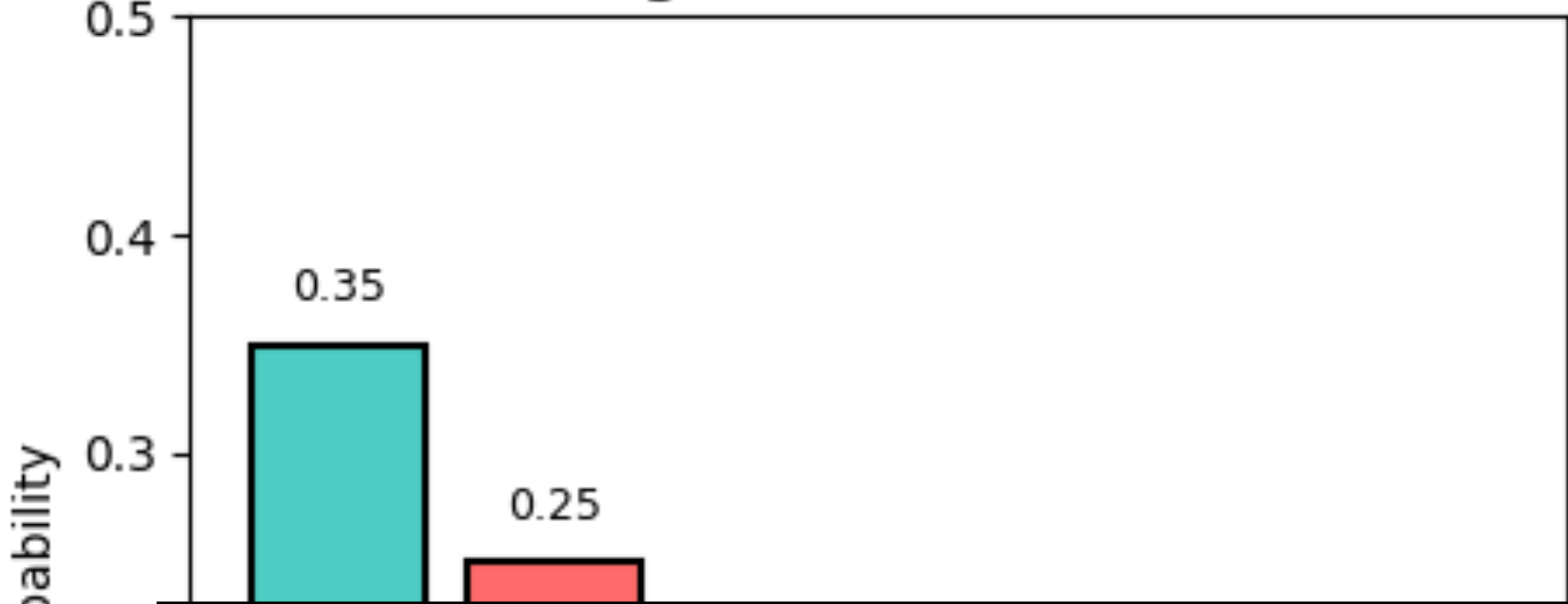


Top-p Nucleus (p=0.8)

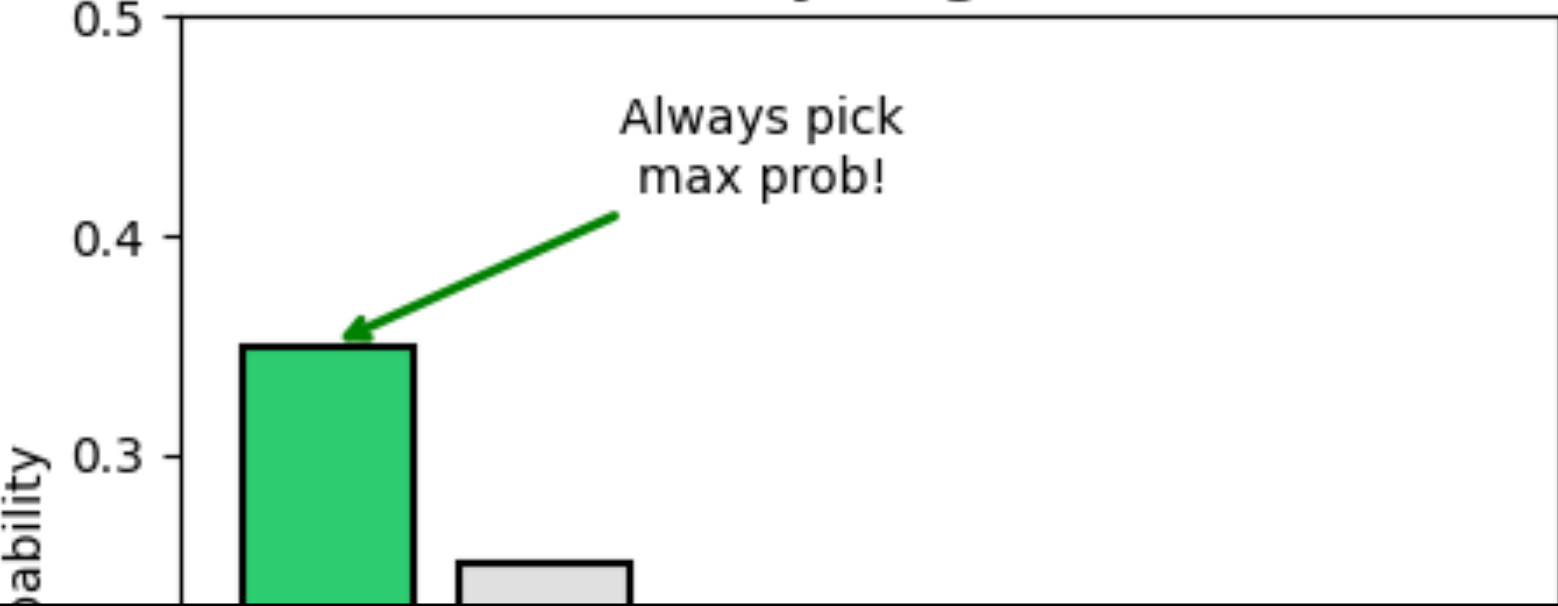


Sampling Strategies from Next-Token Distribution

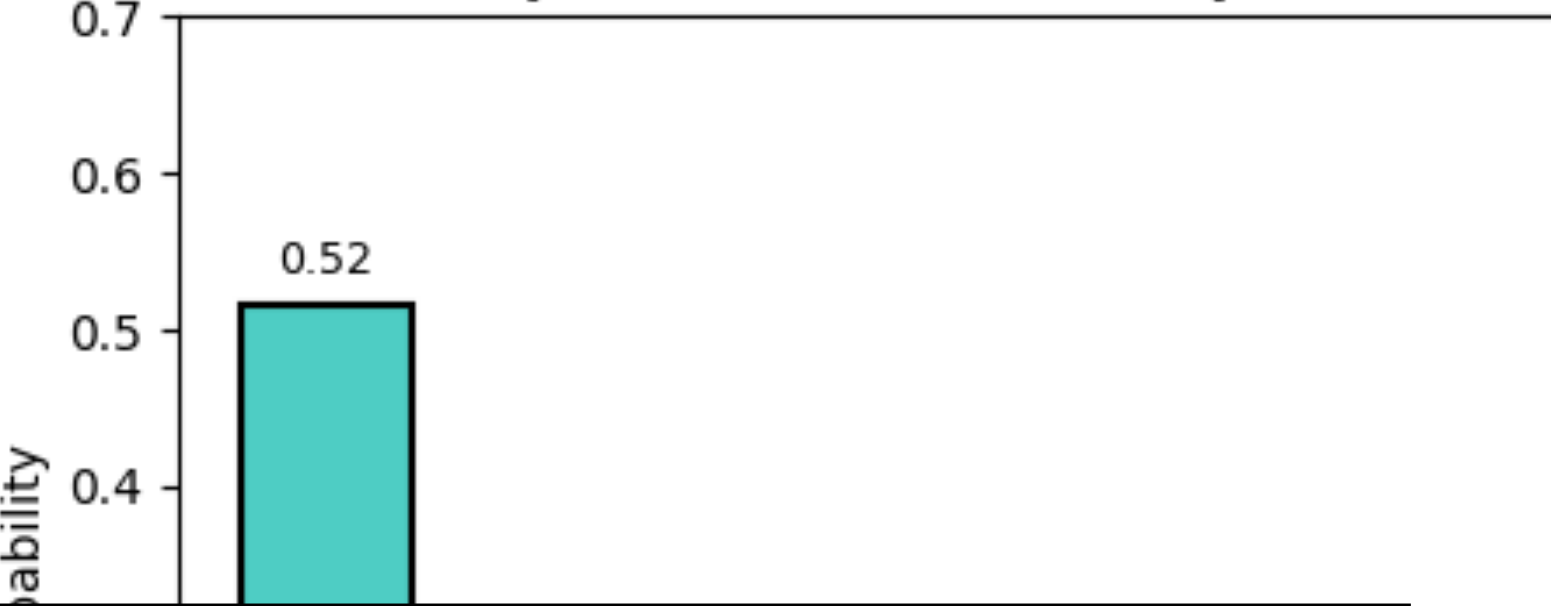
Original Distribution



Greedy (argmax)



Temperature T=0.5 (sharper)



do_sample

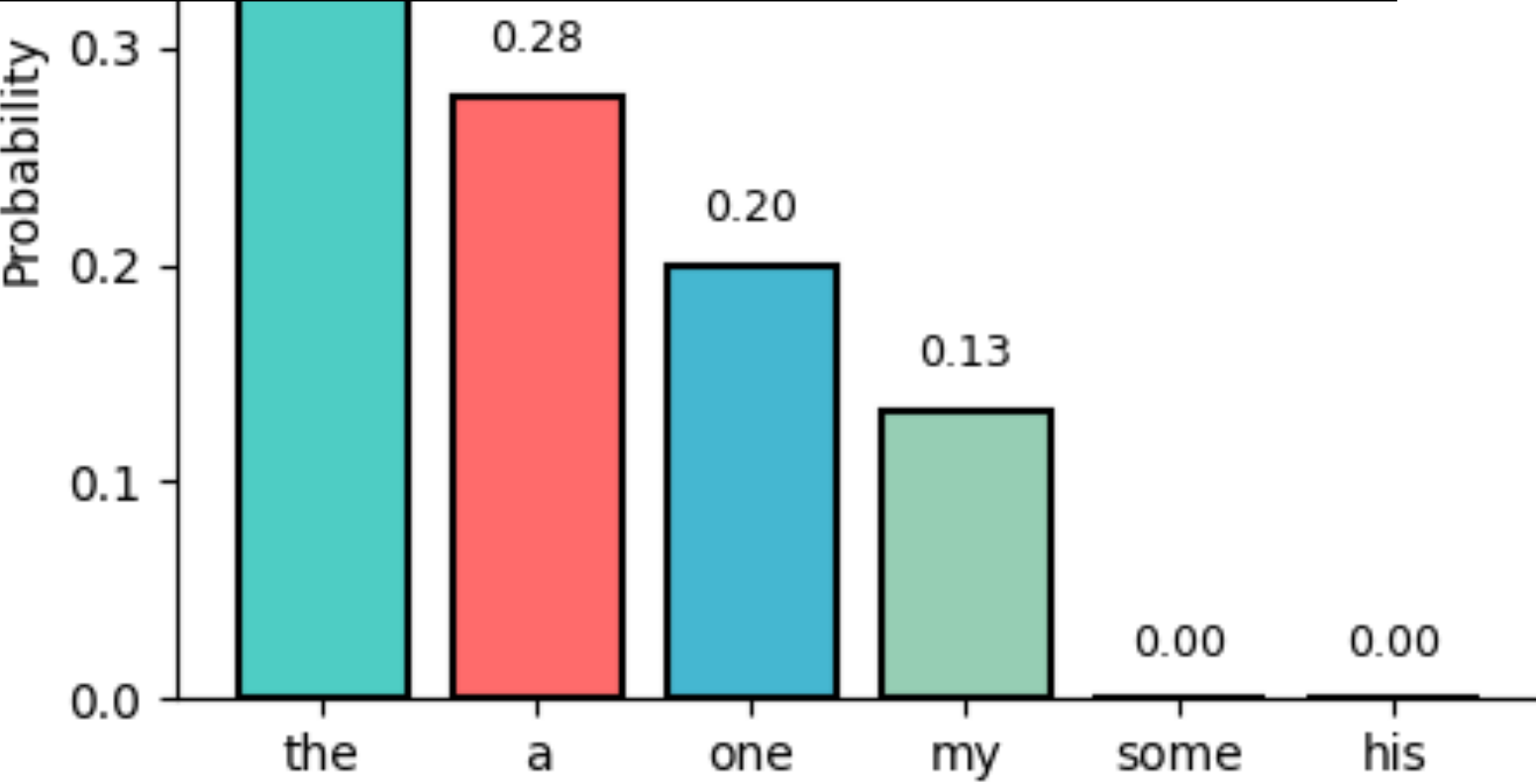
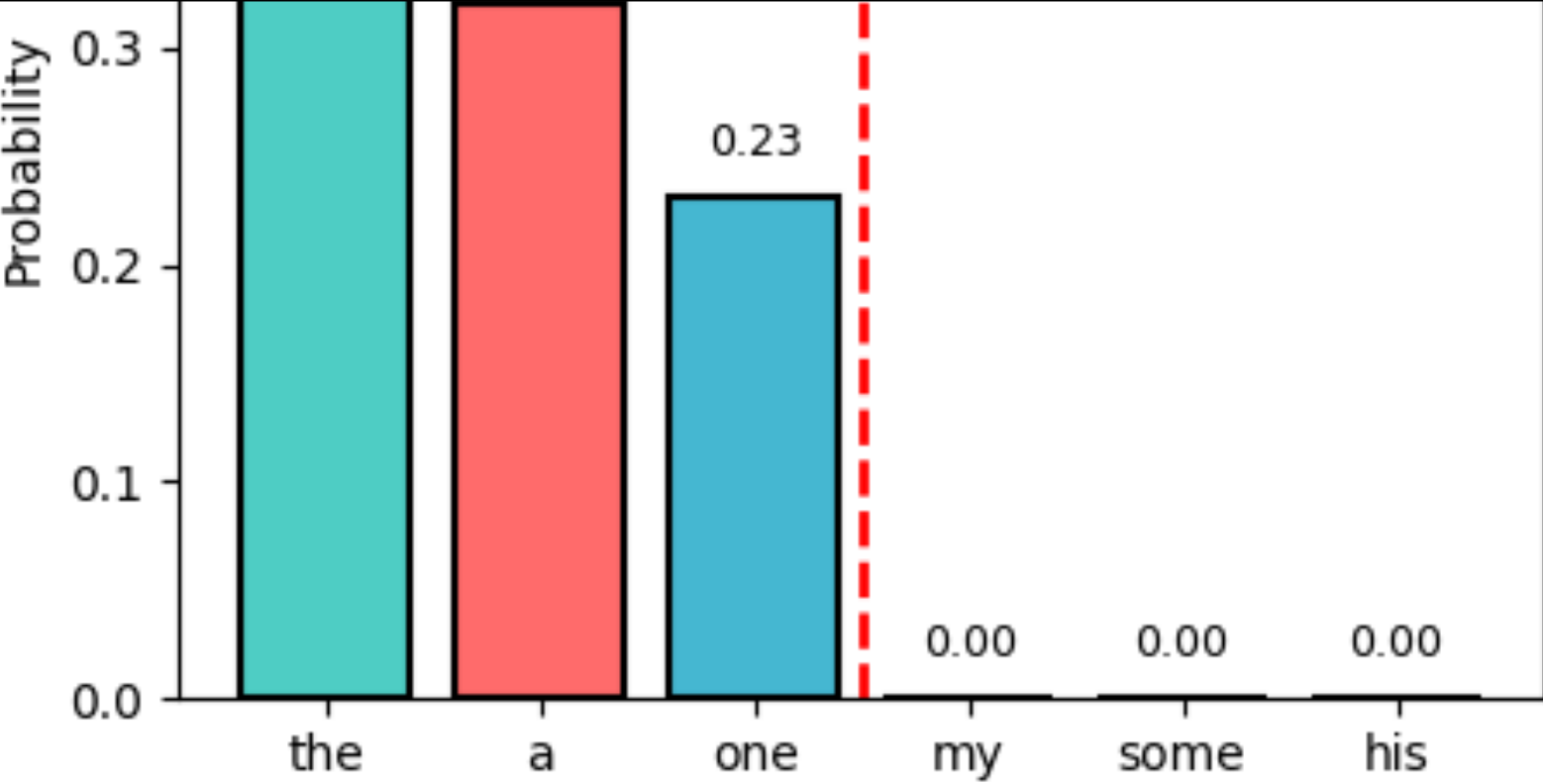
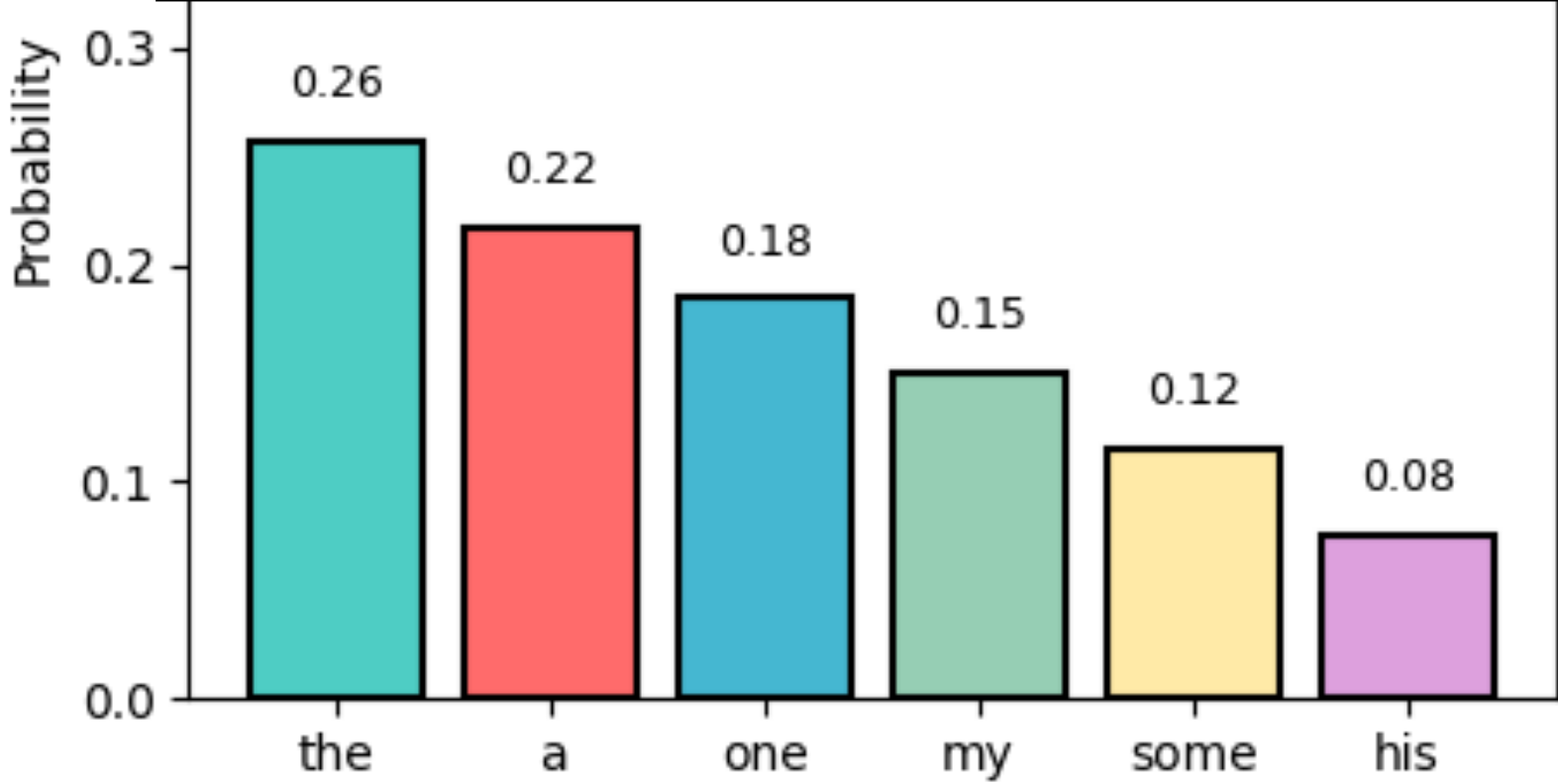
bool

Defines whether generation will sample the next token (True), or is greedy instead (False). Most use cases should set this flag to True. Check [this guide](#) for more information.

temperature

float

How unpredictable the next selected token will be. High values (>0.8) are good for creative tasks, low values (e.g. <0.4) for tasks that require “thinking”. Requires do_sample=True.



What about translation?

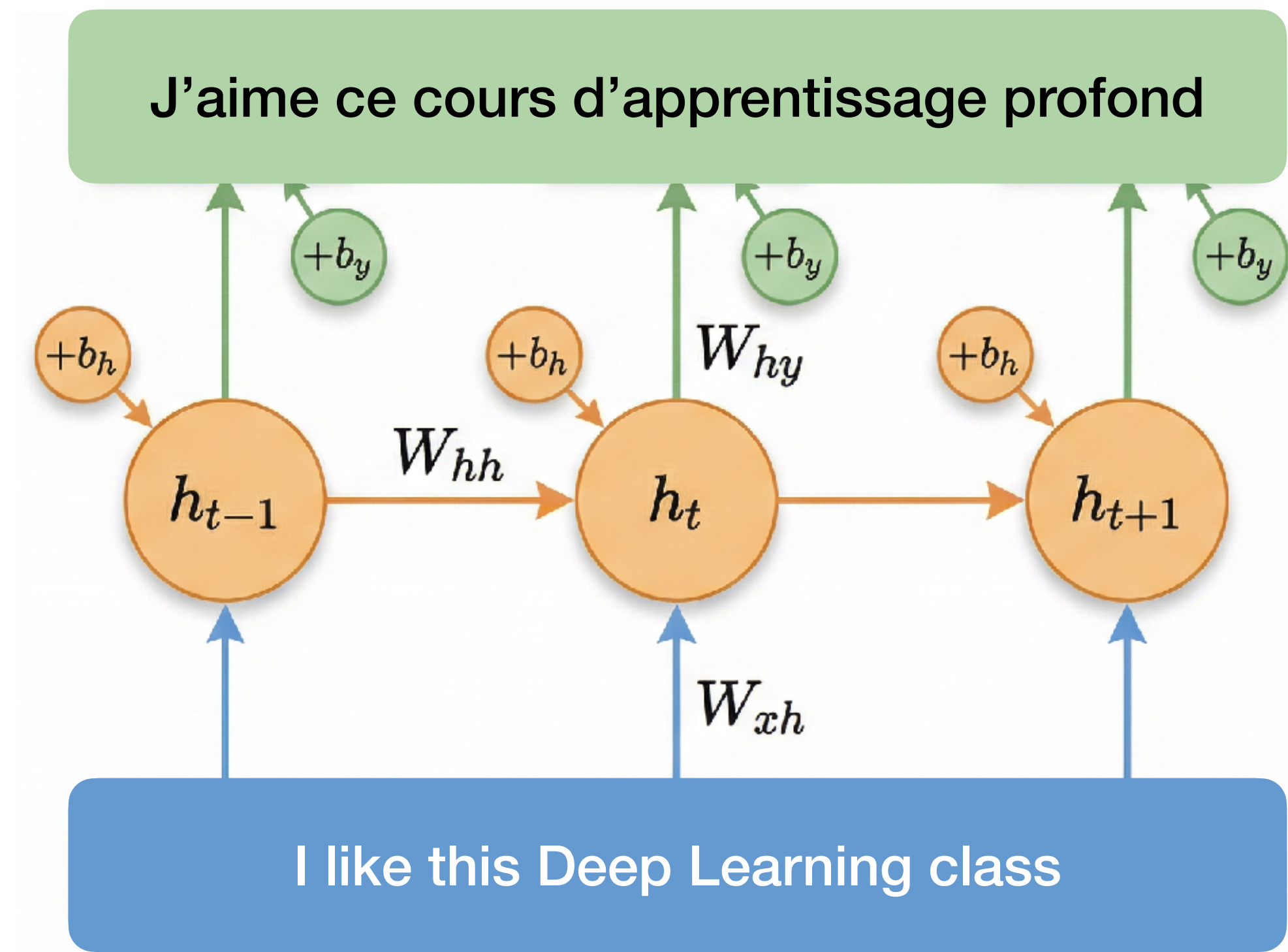
- I like this Deep Learning class <-> J'aime ce cours d'apprentissage profond

What about translation?

- I like this Deep Learning class \leftrightarrow J'aime ce cours d'apprentissage profond

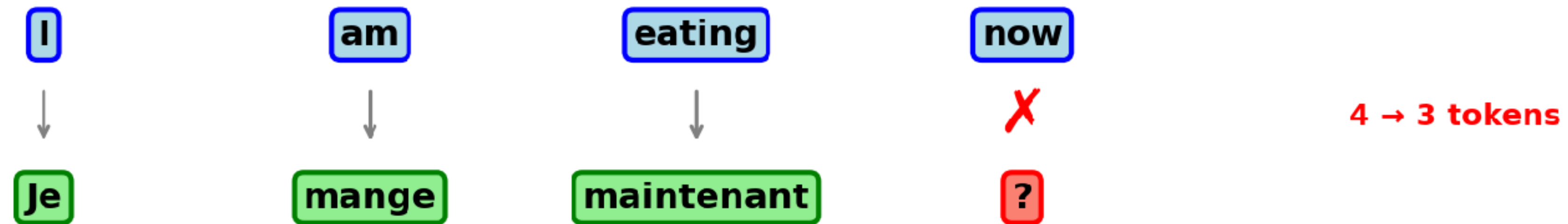
How can we do **Neural Machine Translation**?

What about translation?

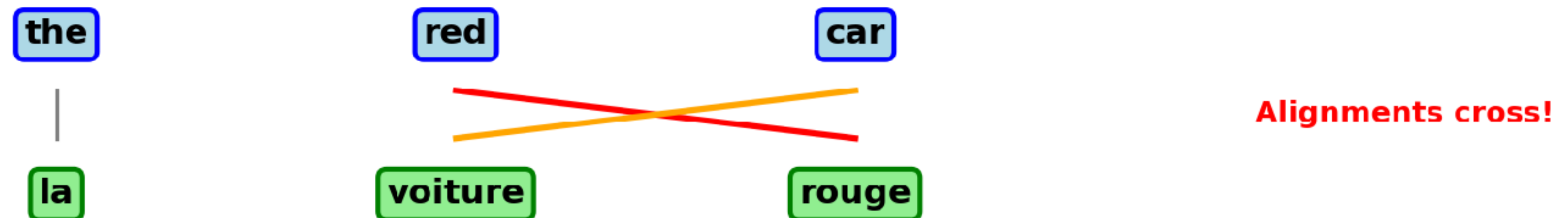


Why Any-to-Any Fails for Translation

① Length Mismatch



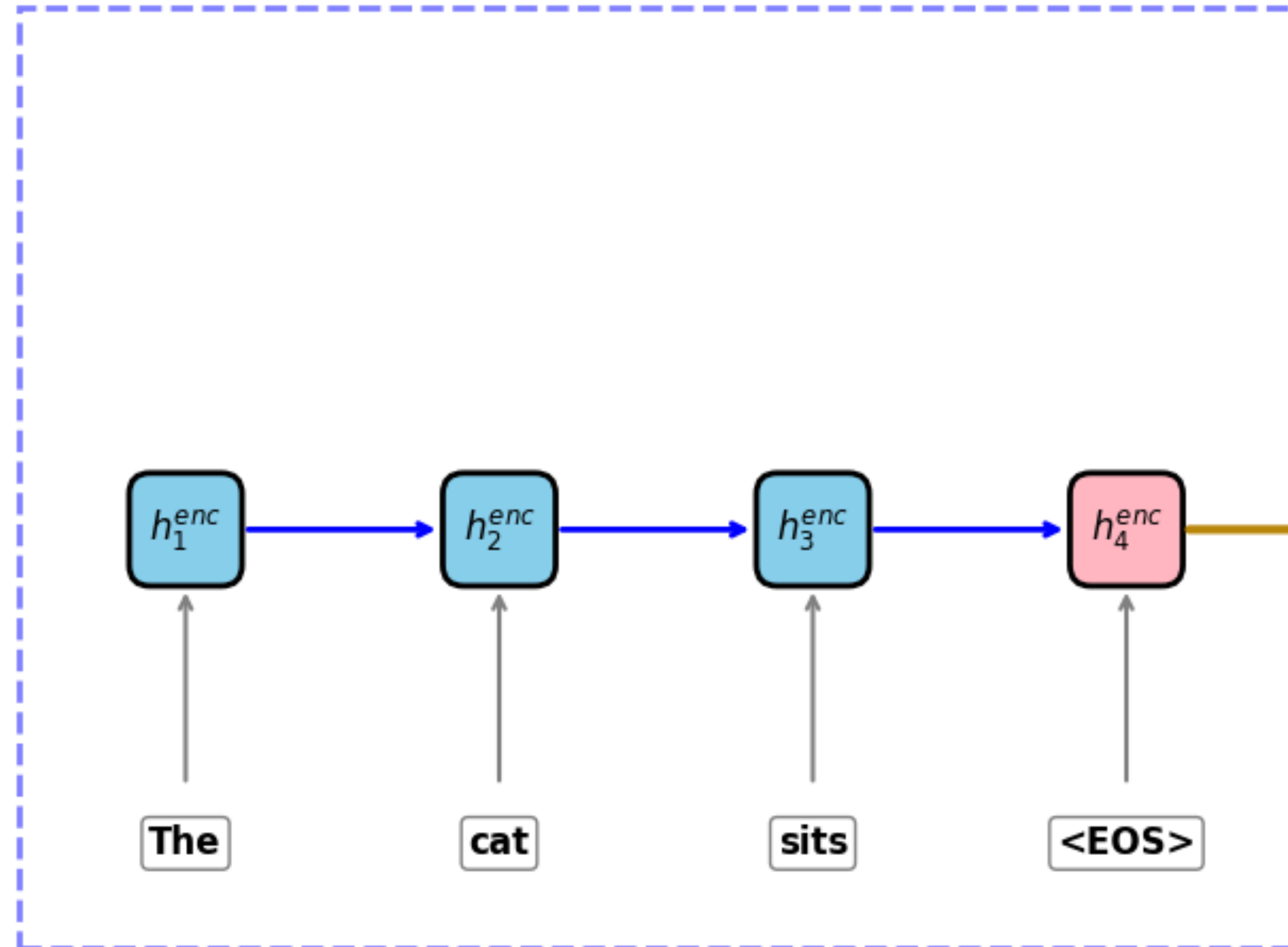
② Word Order Differs



→ Must see ALL source tokens before generating ANY target

Seq2Seq Translation with Autoregressive Decoding
Strict left-to-right generation: cannot peek at future tokens

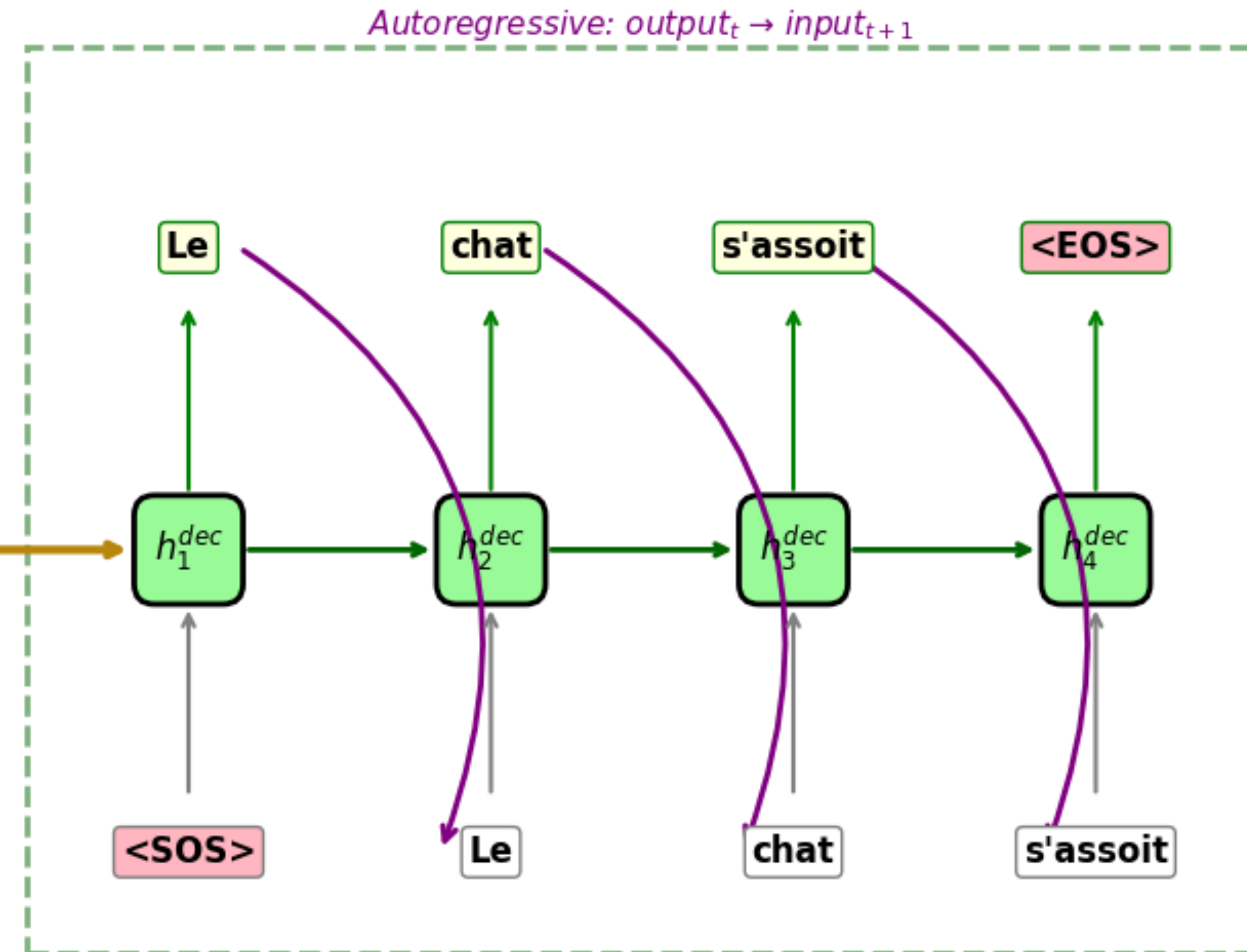
ENCODER



Source: English

Process ALL source first

DECODER



Target: French

THEN generate target sequentially

What about captioning?

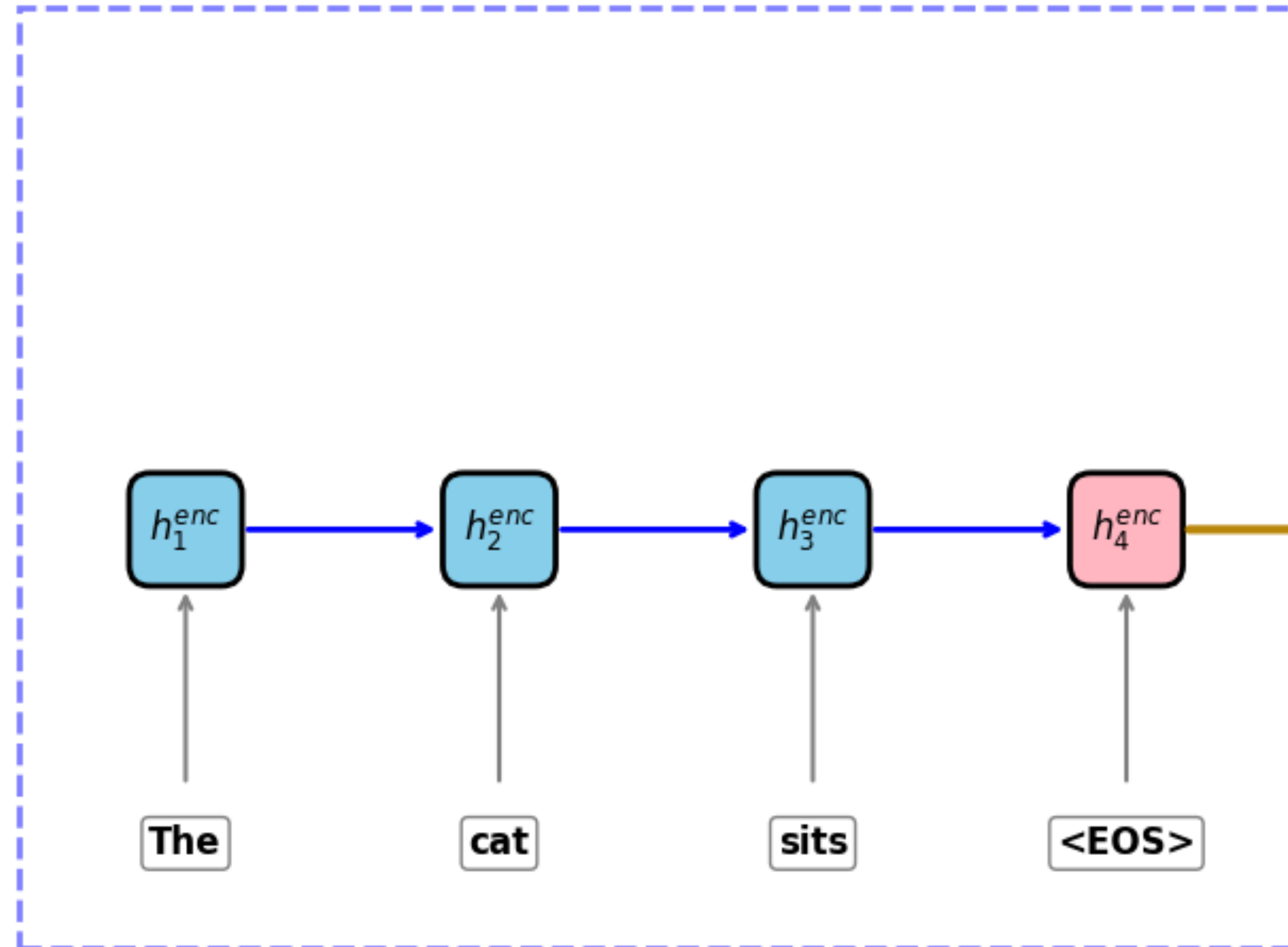
What about captioning?



Randall winning gold at the 100m (sub10s)

Seq2Seq Translation with Autoregressive Decoding
Strict left-to-right generation: cannot peek at future tokens

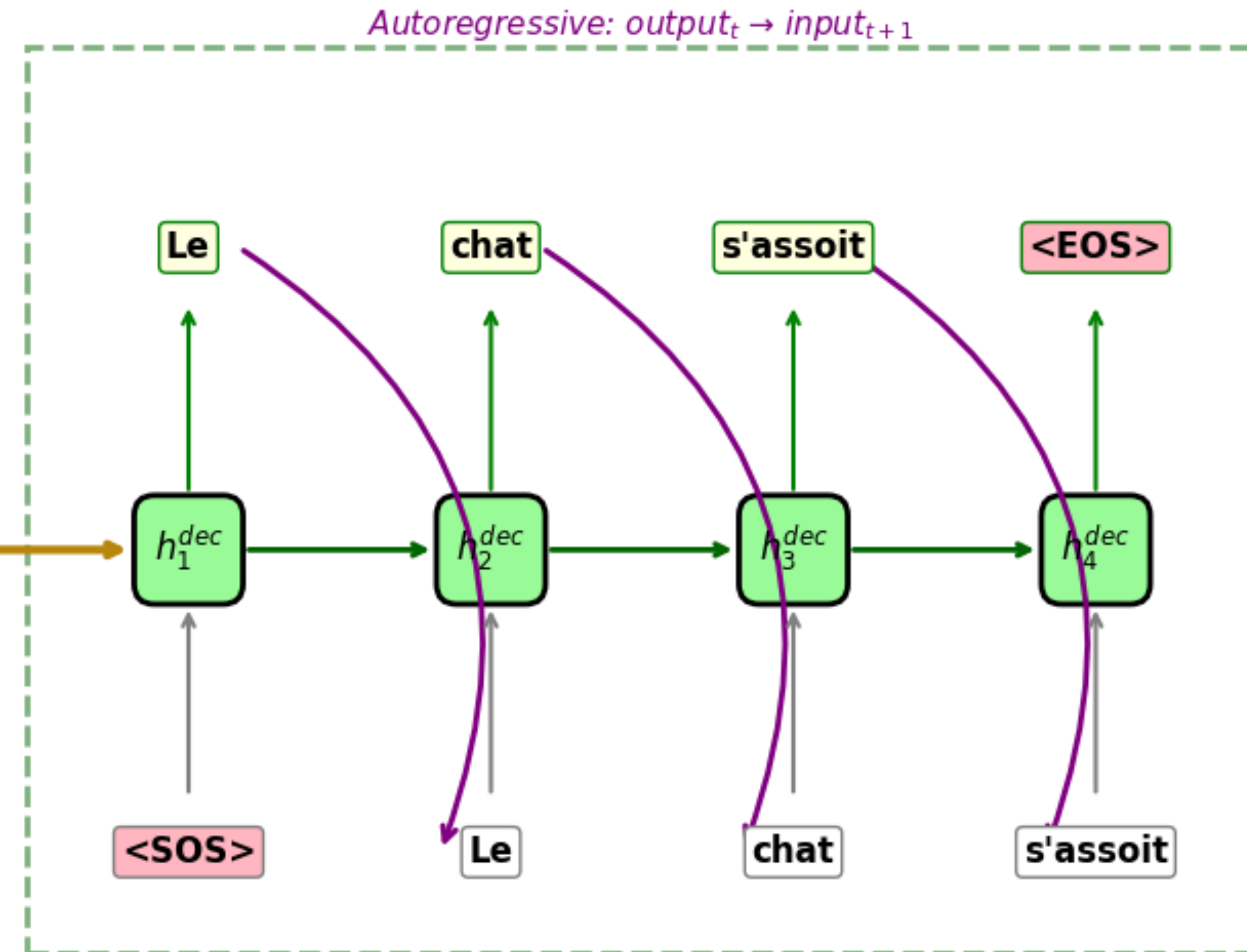
ENCODER



Source: English

Process ALL source first

DECODER



Target: French

THEN generate target sequentially

Seq2Seq Translation with Autoregressive Decoding

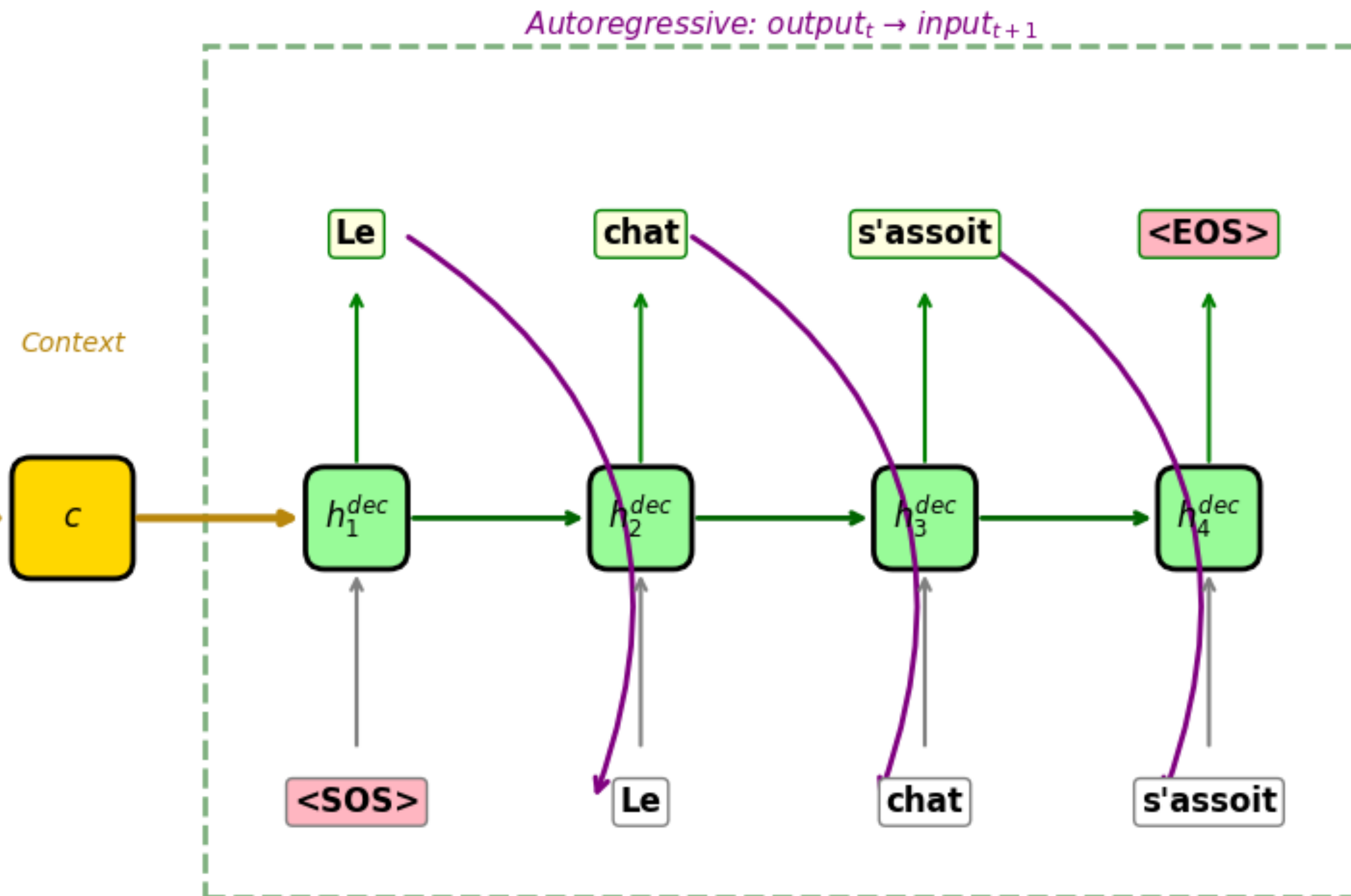
Strict left-to-right generation: cannot peek at future tokens

ENCODER



Source: Vision

DECODER



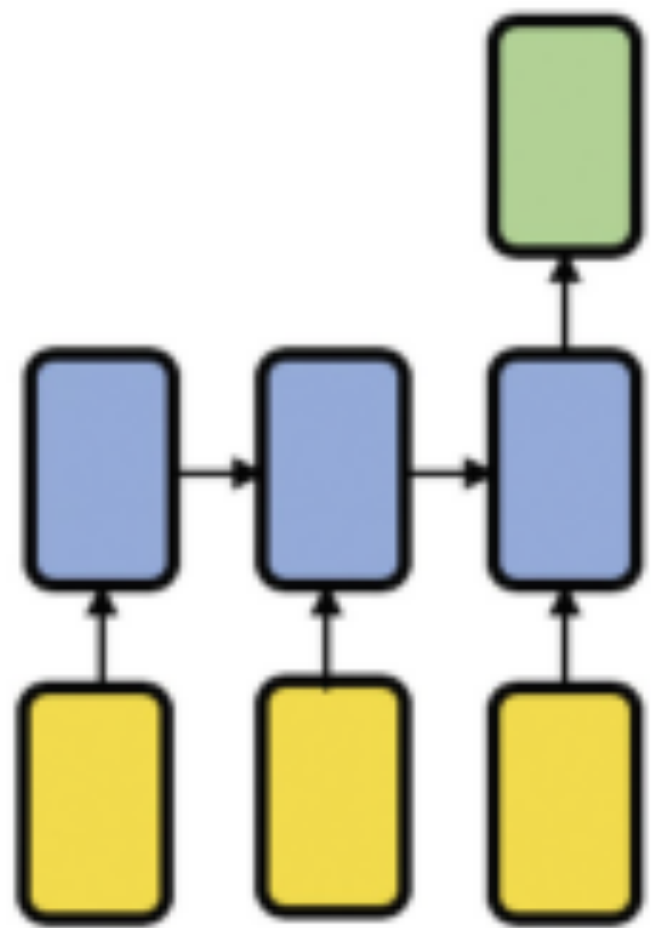
Target: French

Process ALL source first

THEN generate target sequentially

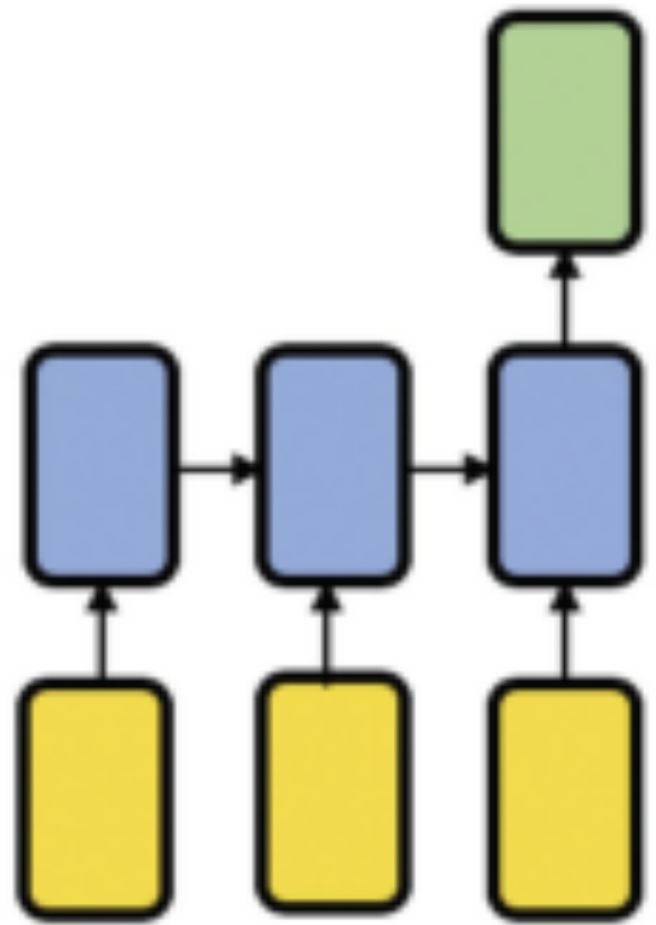
The Sequence Modeling Zoo

The Sequence Modeling Zoo

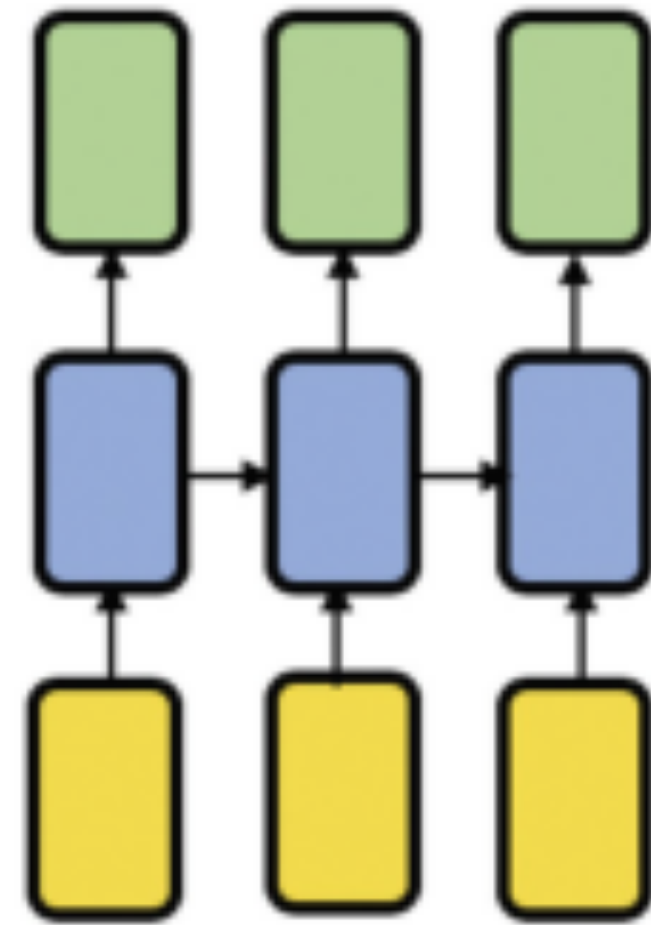


Many-to-One
(e.g., sequence
classification)

The Sequence Modeling Zoo

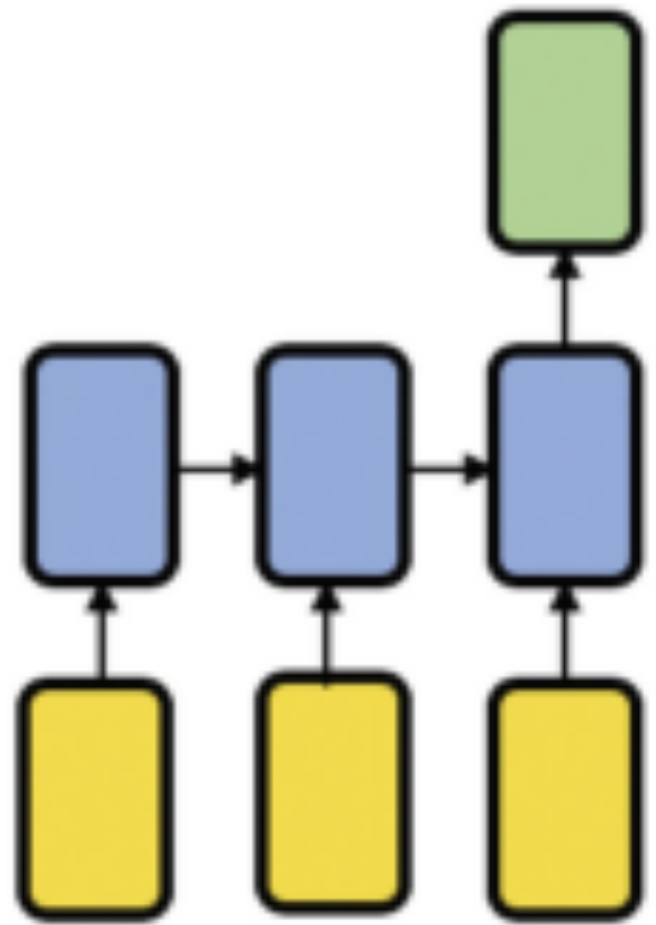


Many-to-One
(e.g., sequence
classification)

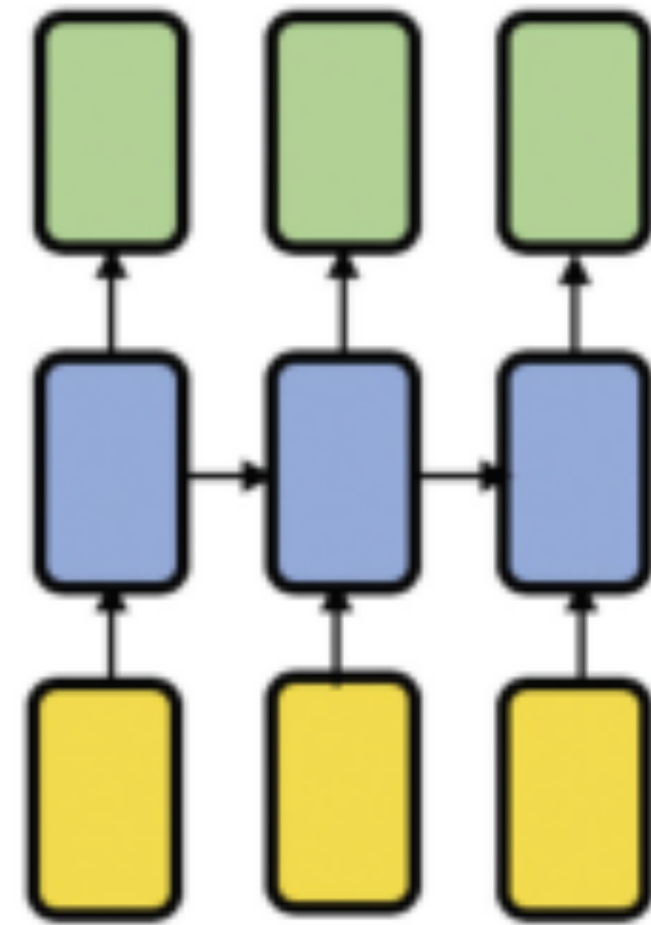


Many-to-Many
(e.g., sequential
prediction)

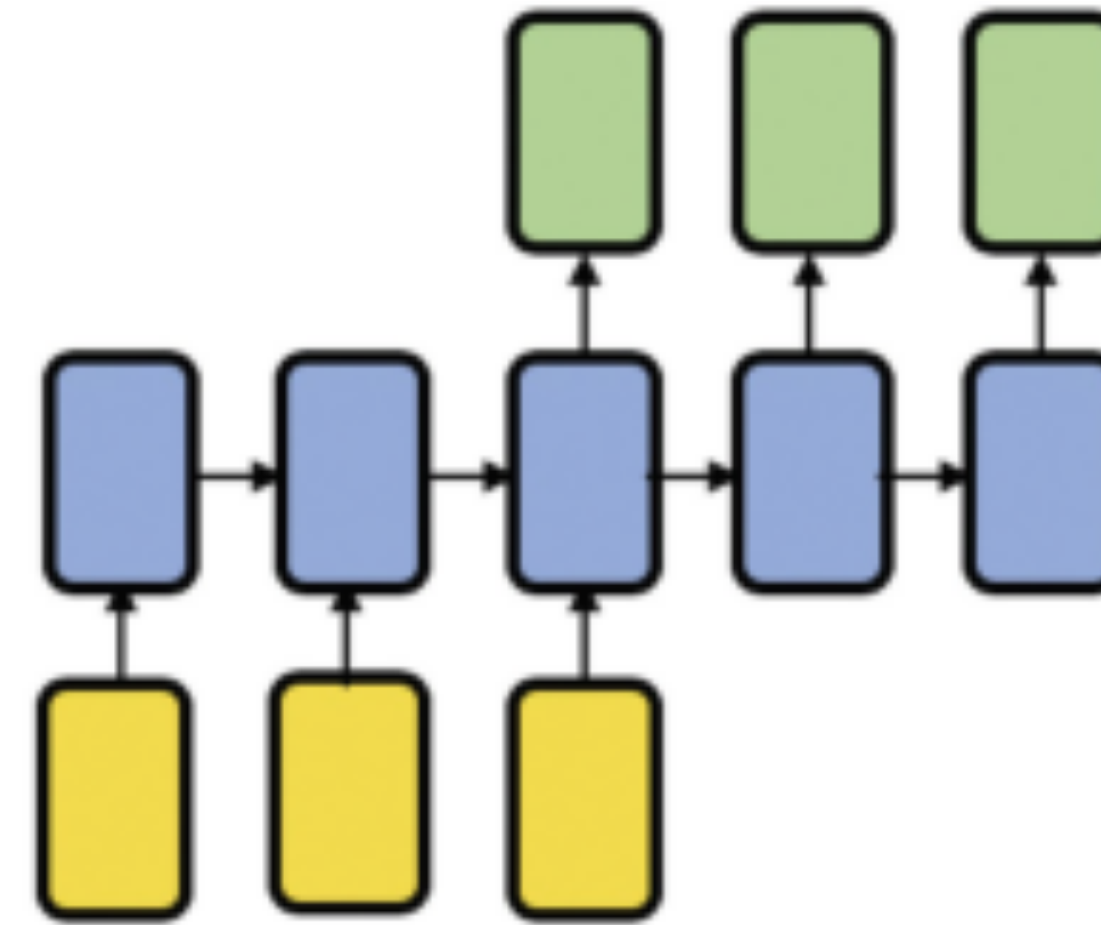
The Sequence Modeling Zoo



Many-to-One
(e.g., sequence
classification)

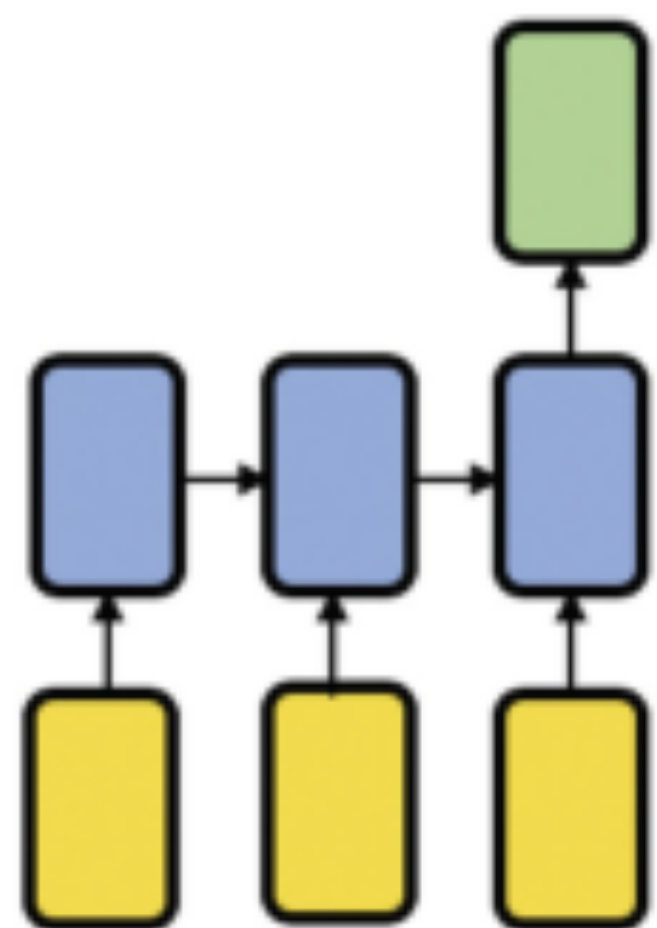


Many-to-Many
(e.g., sequential
prediction)

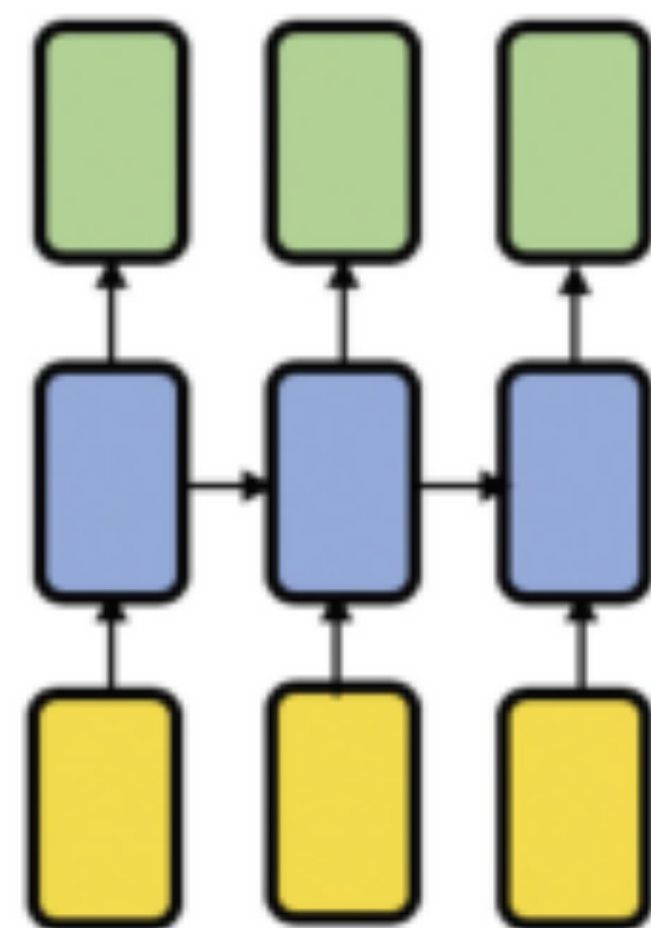


Many-to-Many
(e.g., seq2seq)

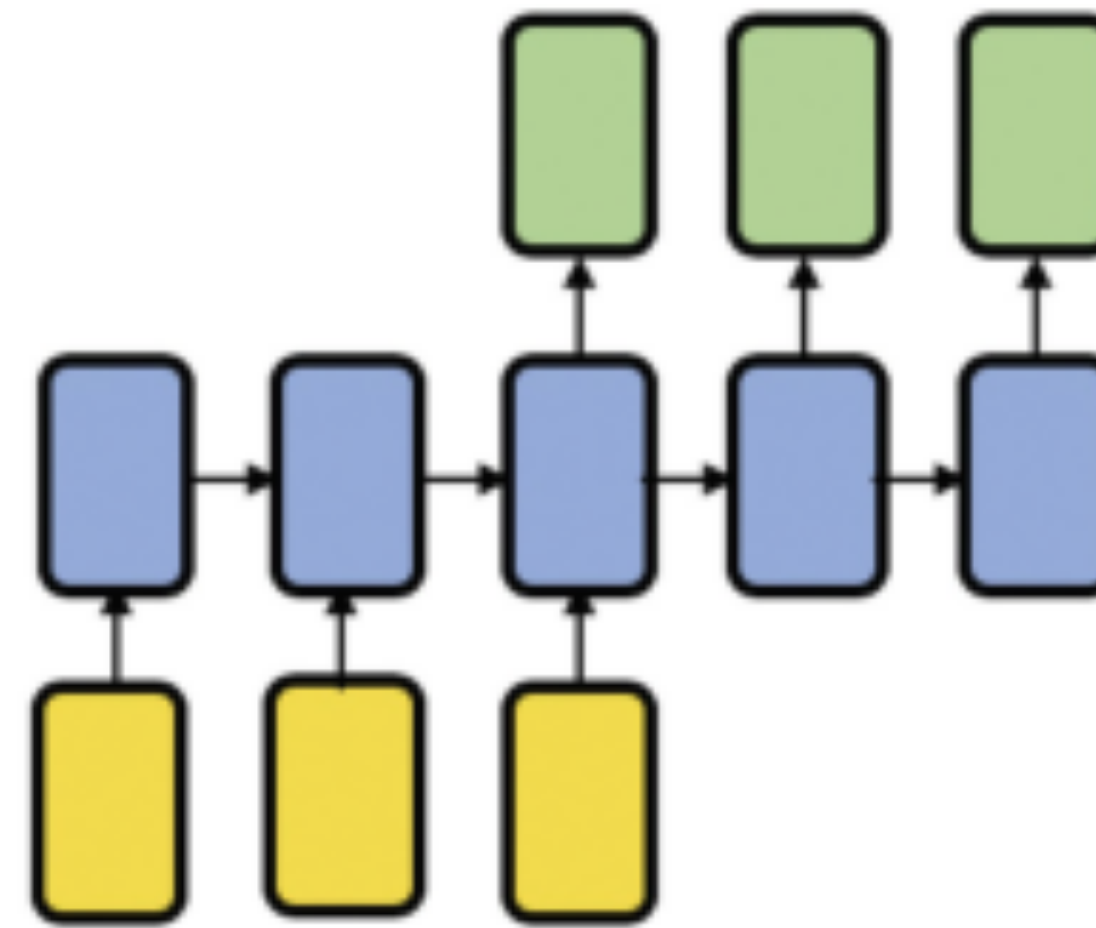
The Sequence Modeling Zoo



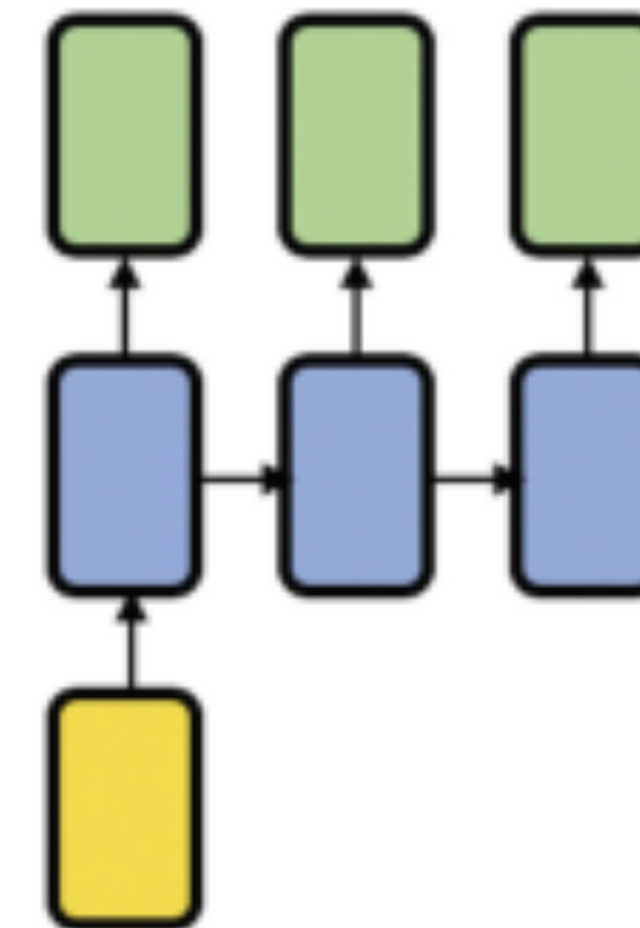
Many-to-One
(e.g., sequence
classification)



Many-to-Many
(e.g., sequential
prediction)

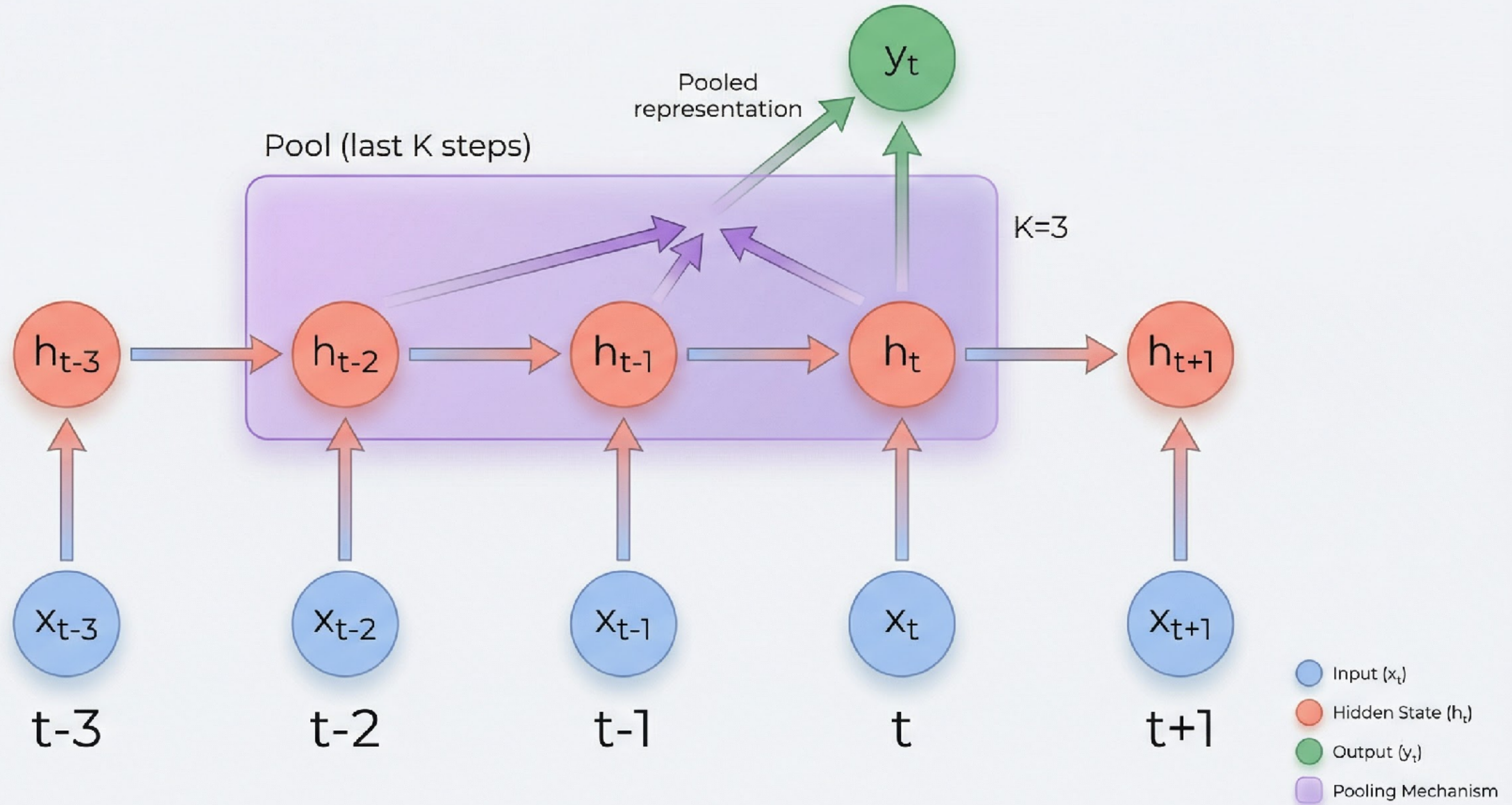


Many-to-Many
(e.g., seq2seq)



One-to-Many
(e.g., image
to text)

Recurrent Neural Network with Pooling Architecture



See you Friday!