# LECTURE 5:

# ADVANCED RECURRENT NEURAL NETWORKS

University of Washington, Seattle

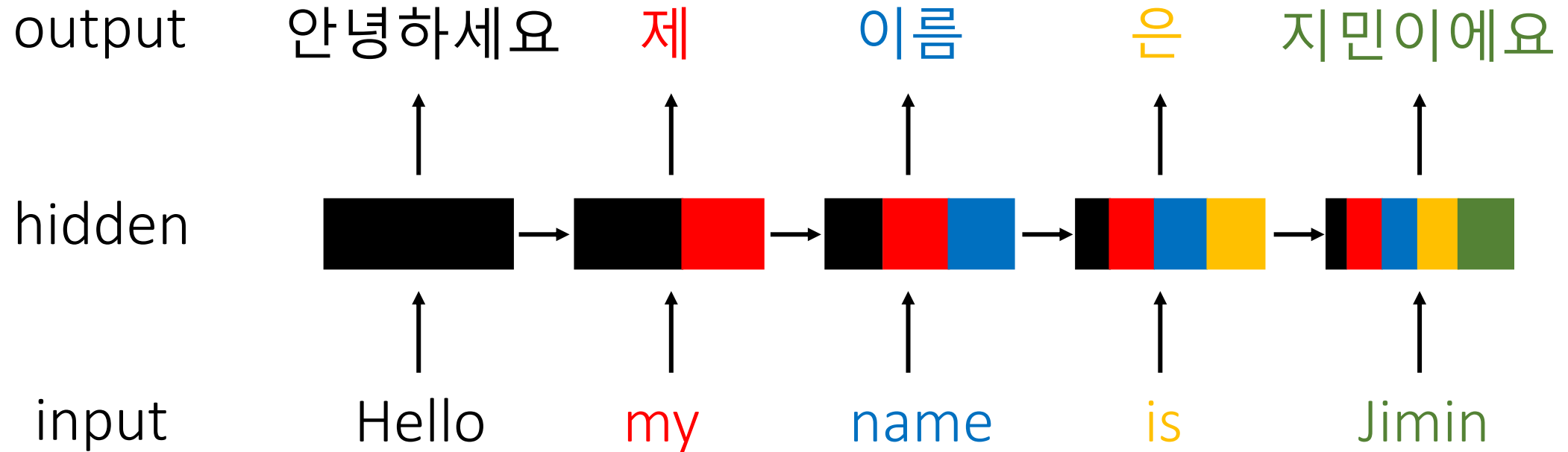Fall 2025

# Previously in EEP 596…



- **Tanh function**
- $h_t$   new hidden state
- $h_{t-1}$   previous hidden state
- $x_t$   input
- concatenation

# Previously in EEP 596…

# OUTLINE

**Part 1: Gated RNNs**

- Need for Gated RNNs

- Long Short-Term Memory (LSTM)

- Gated Recurrent Unit (GRU)

- RNN extensions on LSTM/GRU

**Part 2: Encoder-Decoder RNNs**

- Many to many RNN Recap

- Encoder-Decoder Architecture

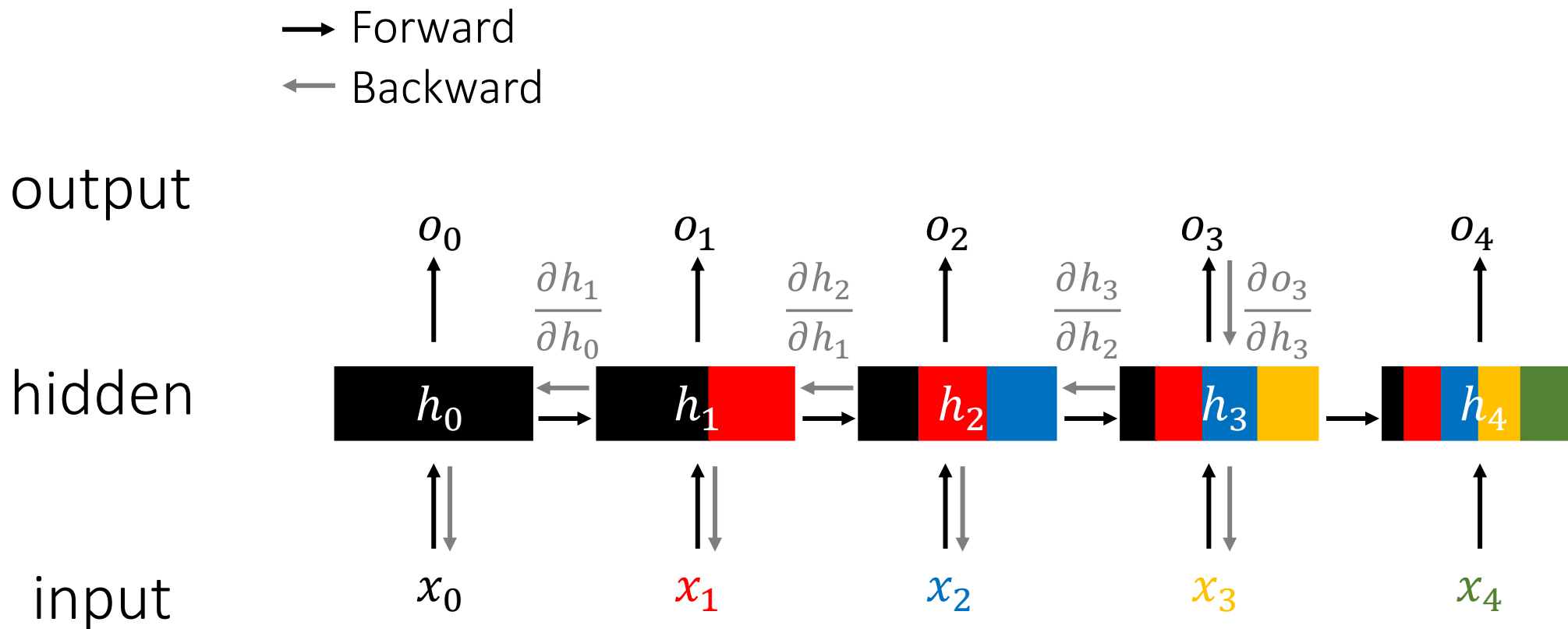- Training Encoder-Decoder RNNs

# GATED RNNs

Need for Gated RNNs

Long Short-Term Memory (LSTM)

Gated Recurrent Unit (GRU)
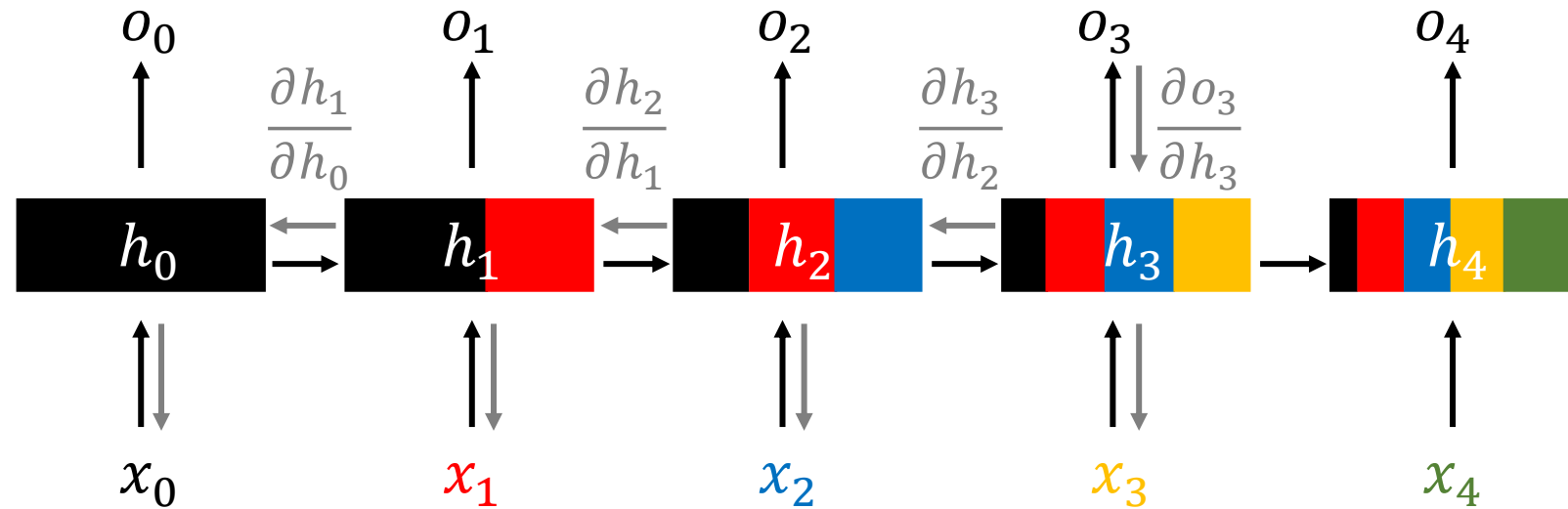
# Recap: Backpropagation in RNNs

# Recap: Backpropagation in RNNs



Backpropagation is performed backward in time

# Vanishing and Exploding Gradients

# Vanishing and Exploding Gradients



Forward

Backward

output

$o_0$  $o_1$  $o_2$  $o_3$  $o_4$

$\frac{\partial h_1}{\partial h_0}$  $\frac{\partial h_2}{\partial h_1}$  $\frac{\partial h_3}{\partial h_2}$  $\frac{\partial o_3}{\partial h_3}$

hidden

$h_0$  $h_1$  $h_2$  $h_3$  $h_4$  ...

input

$x_0$  $x_1$  $x_2$  $x_3$  $x_4$

Longer input sequence →
higher risk of Vanishing/Exploding Gradients!

# Vanishing and Exploding Gradients

Forward

Backward

output

$o_0$  $o_1$  $o_2$  $o_3$  $o_4$

$\frac{\partial h_1}{\partial h_0}$  $\frac{\partial h_2}{\partial h_1}$  $\frac{\partial h_3}{\partial h_2}$  $\frac{\partial o_3}{\partial h_3}$

hidden

$h_0$  $h_1$  $h_2$  $h_3$  $h_4$  $\cdots$
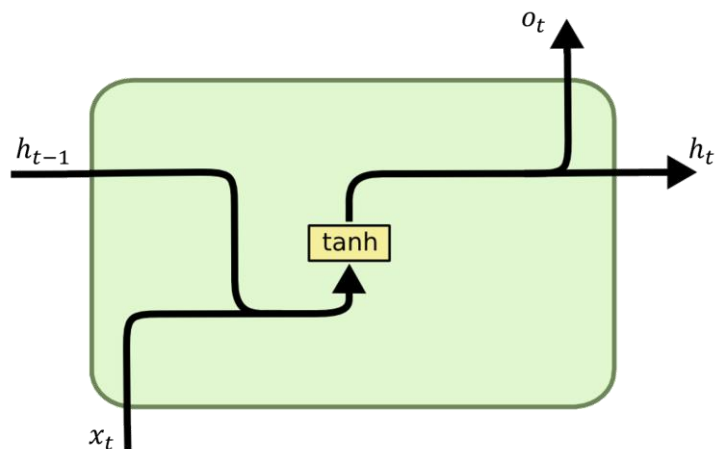
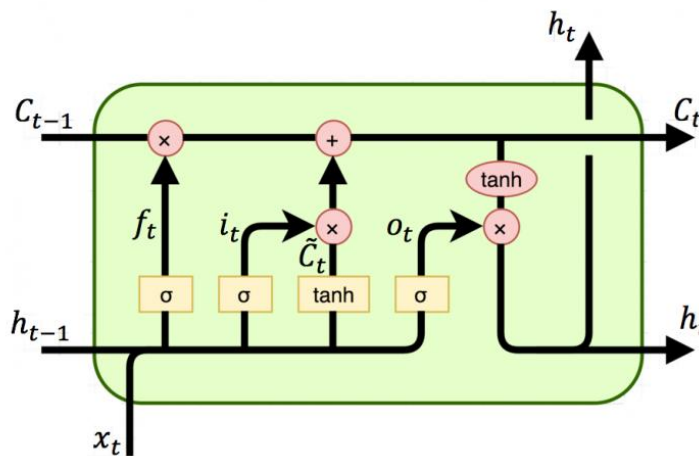input

$x_0$  $x_1$  $x_2$  $x_3$  $x_4$

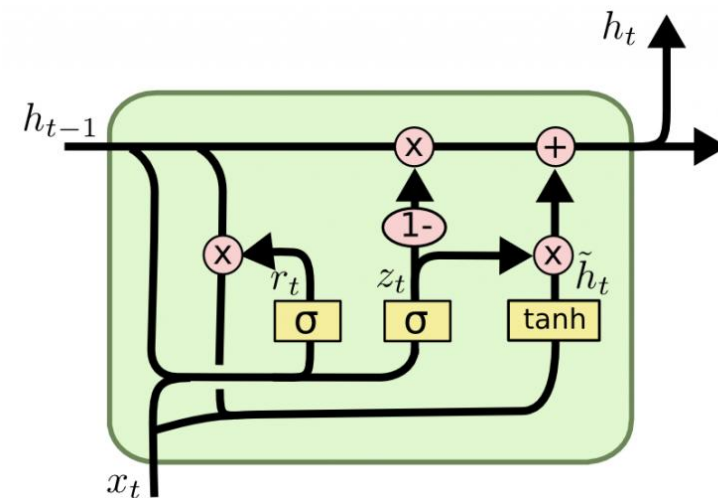Need for better RNN architecture capable of
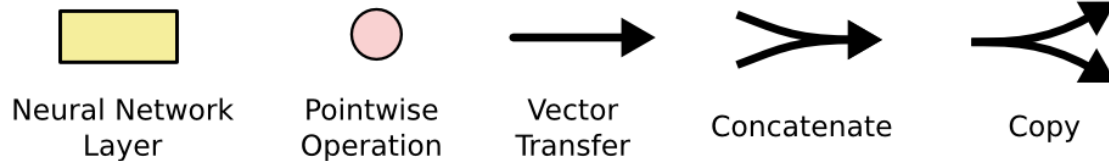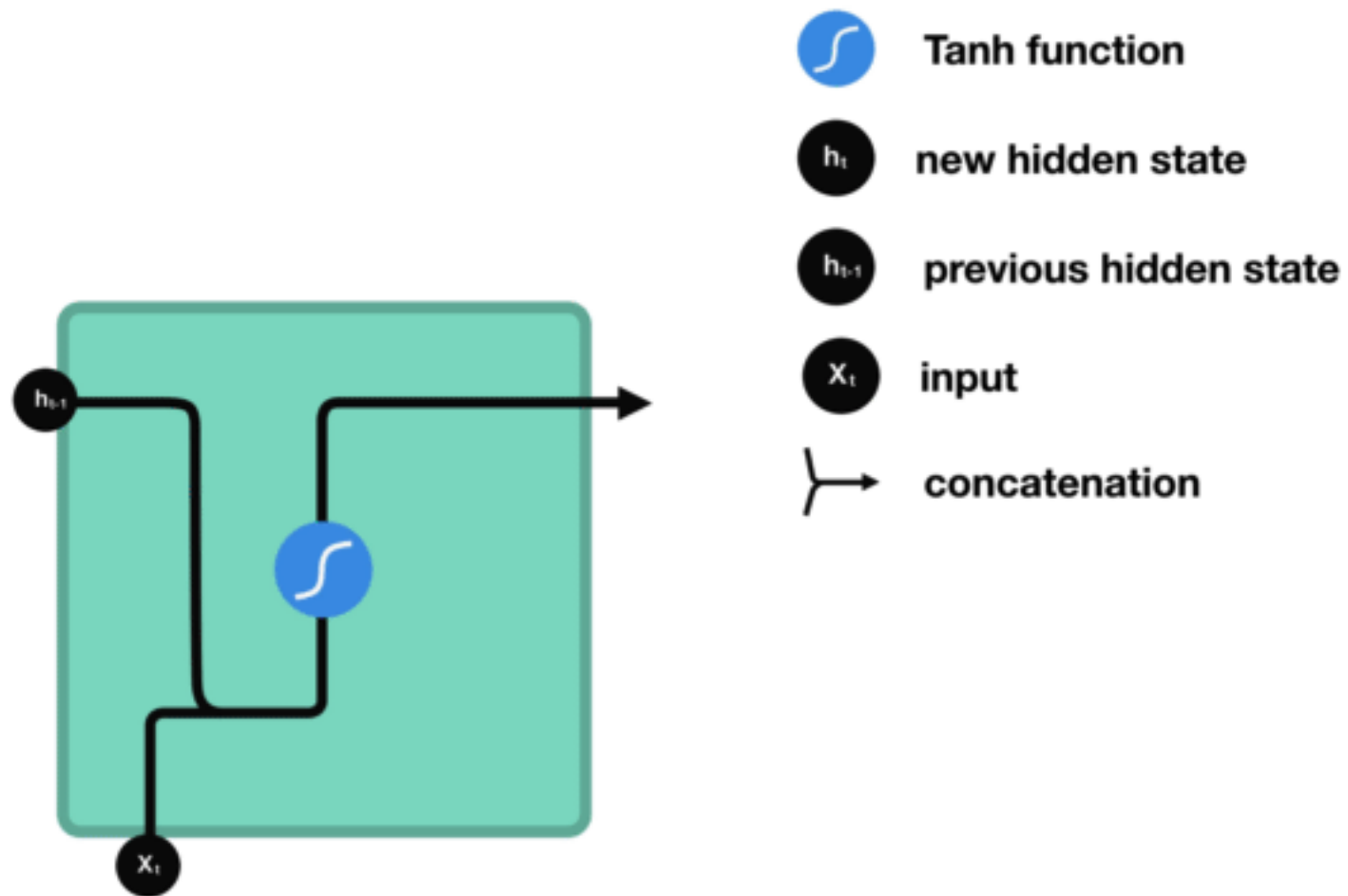processing longer sequence
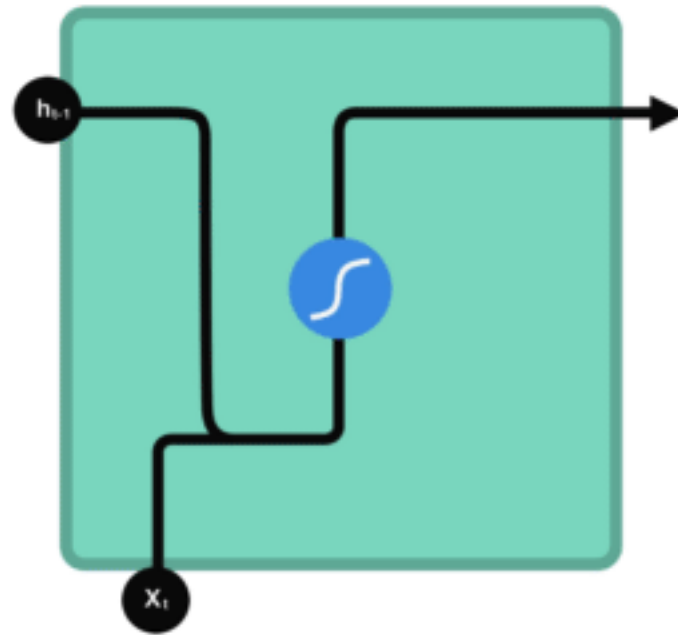
# Gated RNNs



Vanilla RNN

LSTM

GRU

# Vanilla RNN



Vanilla RNN

# Vanilla RNN



Tanh function

$h_t$  new hidden state

$h_{t-1}$  previous hidden state

$x_t$  input

concatenation

# Vanilla RNN

**Tanh function**

$h_t$ **new hidden state**

$h_{t-1}$ **previous hidden state**

$x_t$ **input**

**concatenation**

$$a^{(t)} = b + Wh^{(t-1)} + Ux^{(t)}$$
$$h^{(t)} = \tanh(a^{(t)})$$
$$o^{(t)} = c + Vh^{(t)}$$
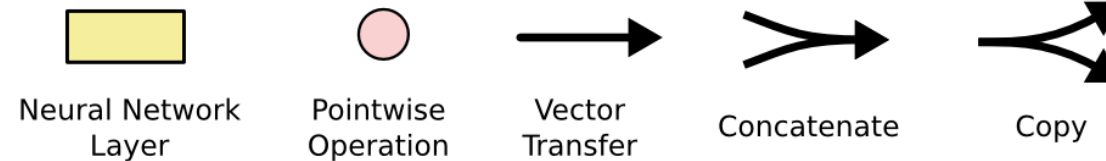$$\hat{y}^{(t)} = \mathrm{softmax}(o^{(t)})$$

# LSTM (Long Short-Term Memory)

Hochreiter, Sepp, and Jürgen Schmidhuber.
"Long short-term memory." *Neural computation* 9.8 (1997): 1735-1780.
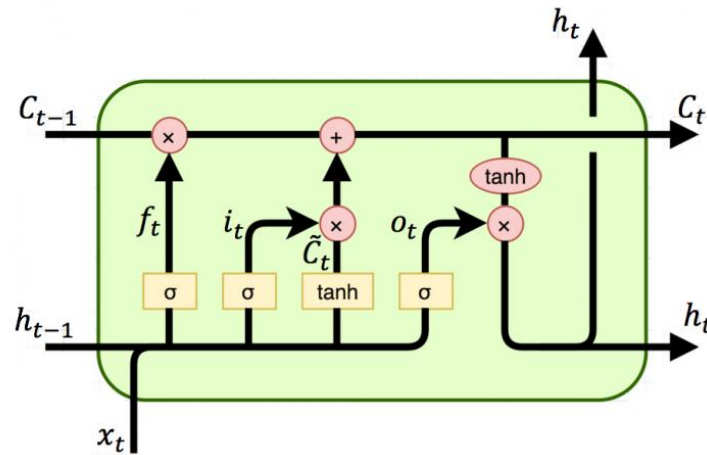


LSTM

| | | | | |
|---|---|---|---|---|
| Neural Network Layer | Pointwise Operation | Vector Transfer | Concatenate | Copy |

# LSTM (Long Short-Term Memory)



LSTM

$$f_t = \sigma_g(W_f x_t + U_f h_{t-1} + b_f)$$
$$i_t = \sigma_g(W_i x_t + U_i h_{t-1} + b_i)$$
$$o_t = \sigma_g(W_o x_t + U_o h_{t-1} + b_o)$$
$$c_t = f_t \circ c_{t-1} + i_t \circ \sigma_c(W_c x_t + U_c h_{t-1} + b_c)$$
$$h_t = o_t \circ \sigma_h(c_t)$$

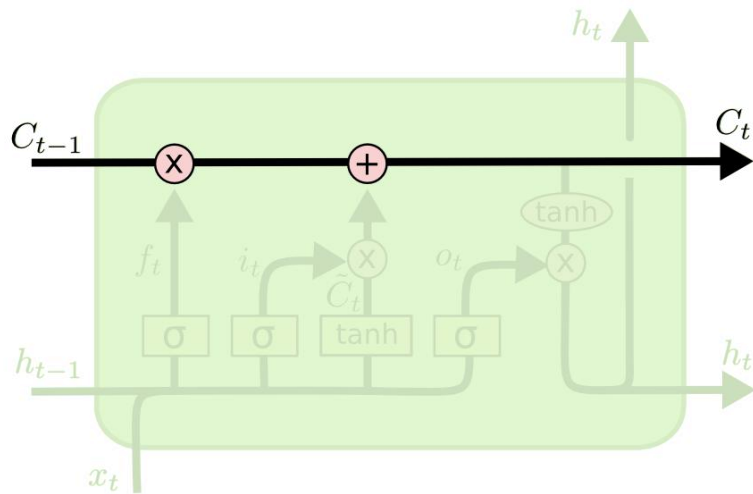Neural Network Layer    Pointwise Operation    Vector Transfer    Concatenate    Copy
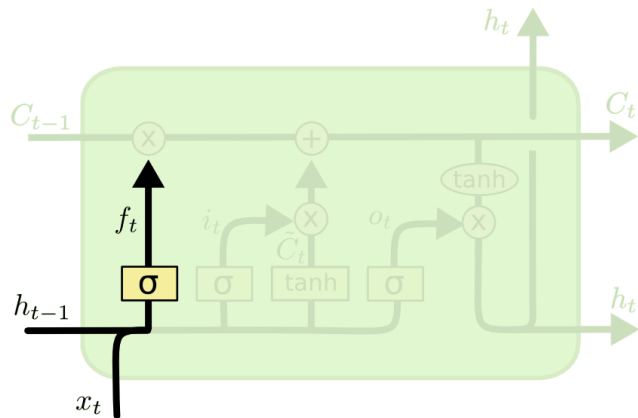
# LSTM: Detailed Architecture



## Cell state

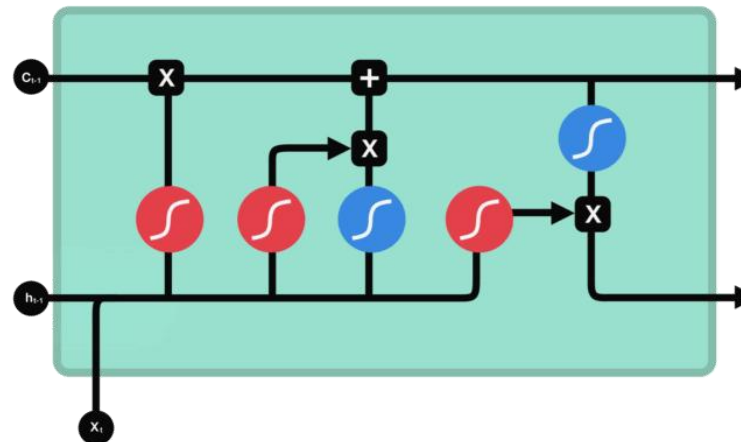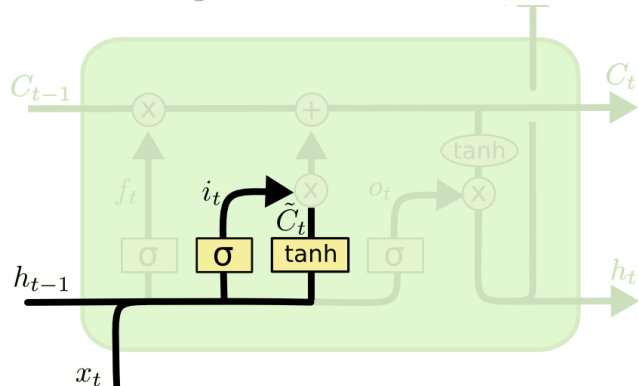- Unique to LSTM
- Long term memory of the model

# LSTM: Detailed Architecture
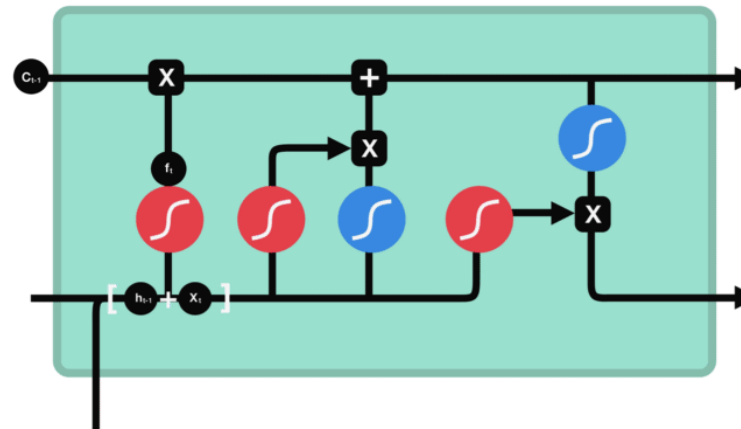
$$f_t = \sigma_g(W_f x_t + U_f h_{t-1} + b_f)$$



$$i_t = \sigma_g(W_i x_t + U_i h_{t-1} + b_i)$$



Forget gate layer
(information to forget)

Input gate layer
(information to keep)

# LSTM: Detailed Architecture

$$c_t = f_t \circ c_{t-1} + i_t \circ \sigma_c(W_c x_t + U_c h_{t-1} + b_c)$$



Update cell state
(Forget + Add new info)

$$o_t = \sigma_g(W_o x_t + U_o h_{t-1} + b_o)$$

$$h_t = o_t \circ \sigma_h(c_t)$$



Output gate layer
(Form new hidden state)

**Forget gate** $f_t$
Decides what is relevant to keep from previous steps

**Input gate** $i_t$
Decides what information is relevant to add from the current step
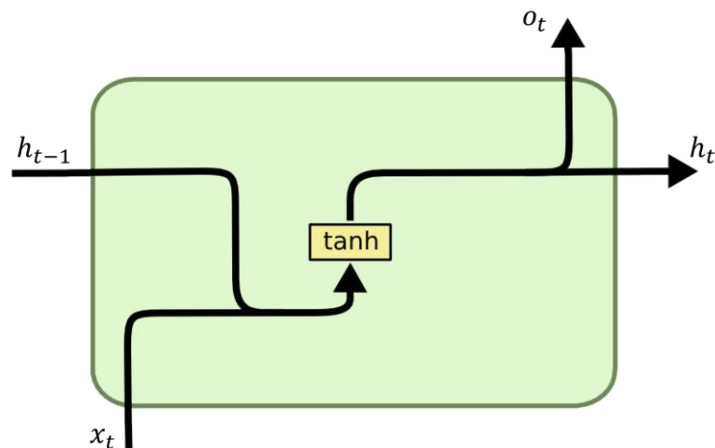
**Output Gate** $o_t$
Combined with cell state to determines next hidden state

LSTM divides **original hidden state** into
1. **long-term memory (cell state)** and 2. **context (LSTM hidden state)**
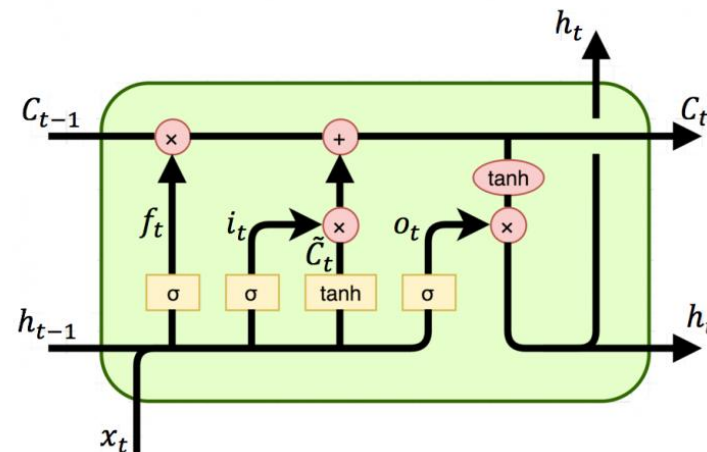
# LSTM (Long Short-Term Memory)



Vanilla RNN

$$h_t = \sigma(wh_{t-1}).$$
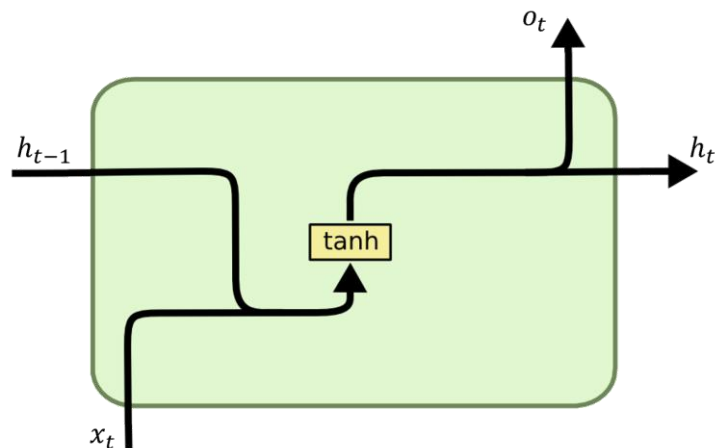
$$\frac{\partial h_{t'}}{\partial h_t} = \prod_{k=1}^{t'-t} w\sigma'(wh_{t'-k})$$

$$= \underbrace{w^{t'-t}}_{!!!} \prod_{k=1}^{t'-t} \sigma'(wh_{t'-k})$$

LSTM

$$\frac{\partial c_{t'}}{\partial c_t} = \prod_{k=1}^{t'-t} \sigma(v_{t+k}).$$
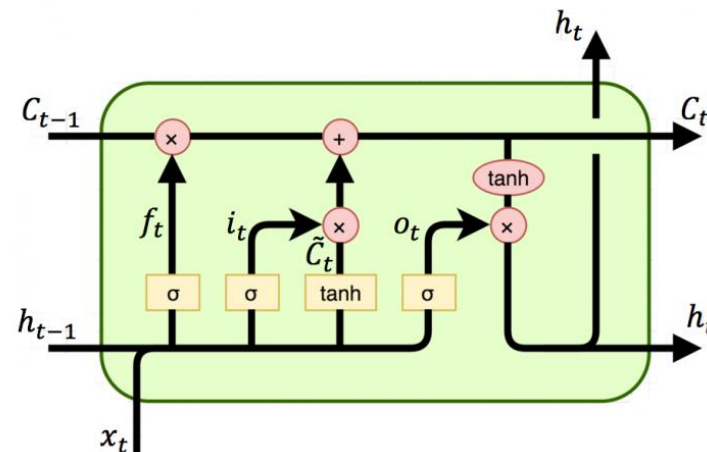
# LSTM (Long Short-Term Memory)



Vanilla RNN

$$h_t = \sigma(wh_{t-1}).$$

$$\frac{\partial h_{t'}}{\partial h_t} = \prod_{k=1}^{t'-t} w\sigma'(wh_{t'-k})$$

$$= \underbrace{w^{t'-t}}_{!!!} \prod_{k=1}^{t'-t} \sigma'(wh_{t'-k})$$

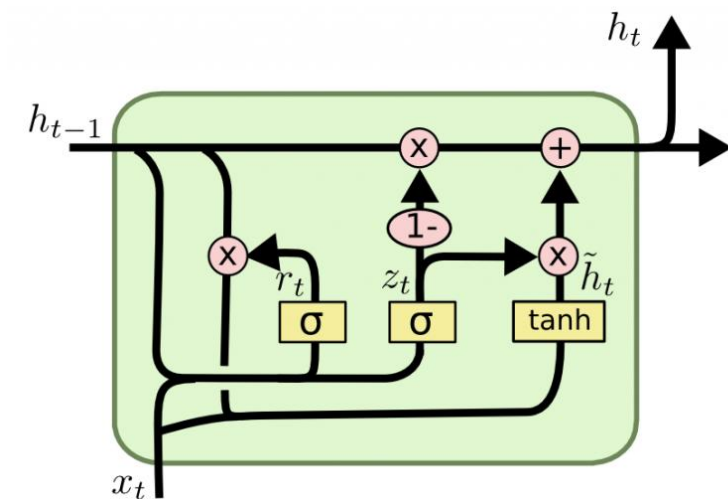Gradient decays or grow exponentially
if $w \neq 1$

LSTM

$$\frac{\partial c_{t'}}{\partial c_t} = \prod_{k=1}^{t'-t} \sigma(v_{t+k}).$$

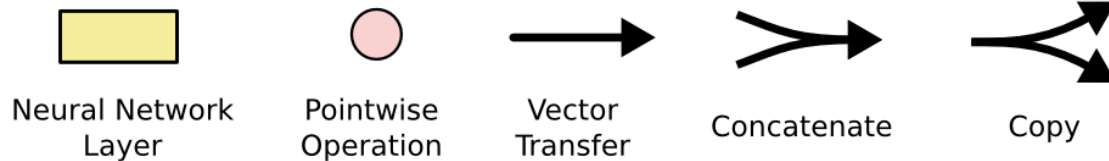No exponential decay or growth term

# Gated RNNs

Chung, Junyoung, et al. "Empirical evaluation of gated recurrent neural networks on sequence modeling." *arXiv preprint arXiv:1412.3555* (2014).
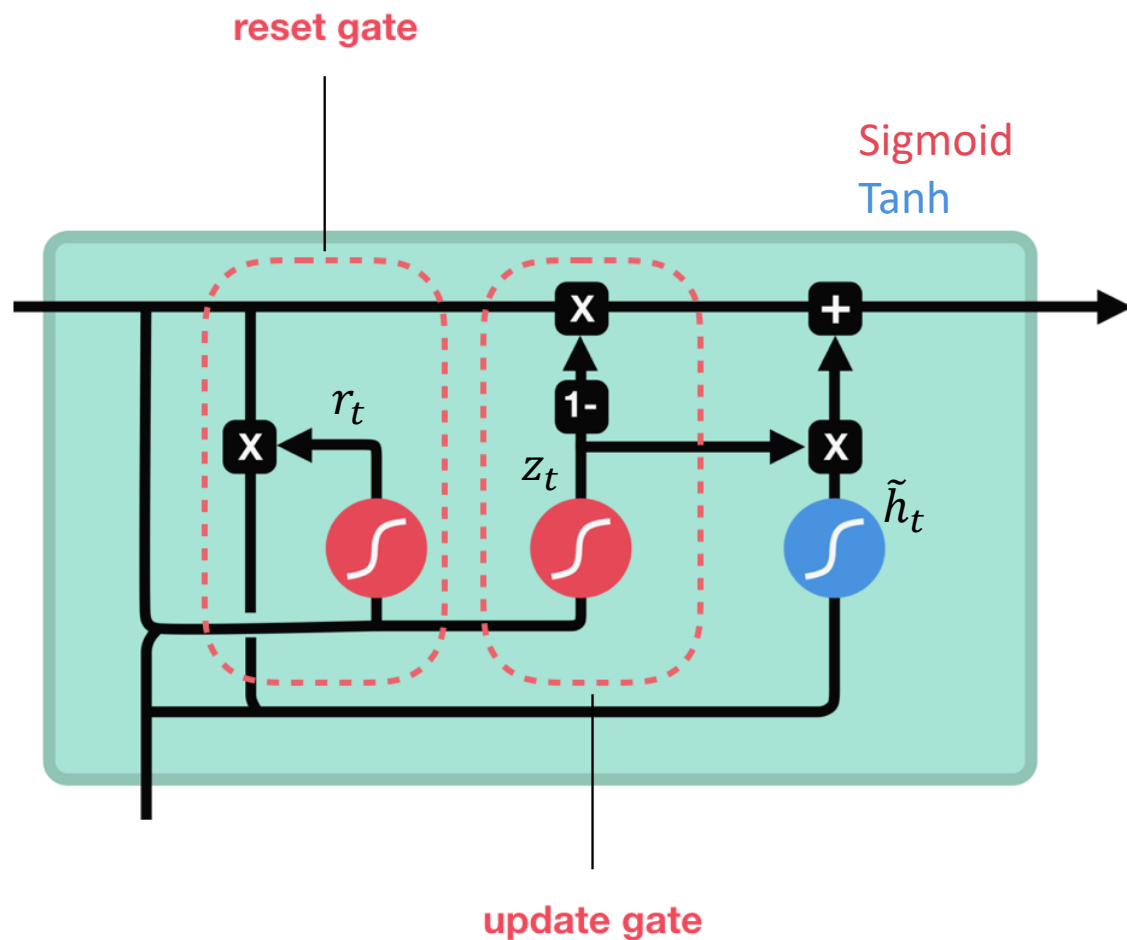
GRU

Neural Network Layer | Pointwise Operation | Vector Transfer | Concatenate | Copy

# GRU: Detailed Architecture



$$z_t = \sigma\left(W_z \cdot [h_{t-1}, x_t]\right) \qquad \in [0,1]$$
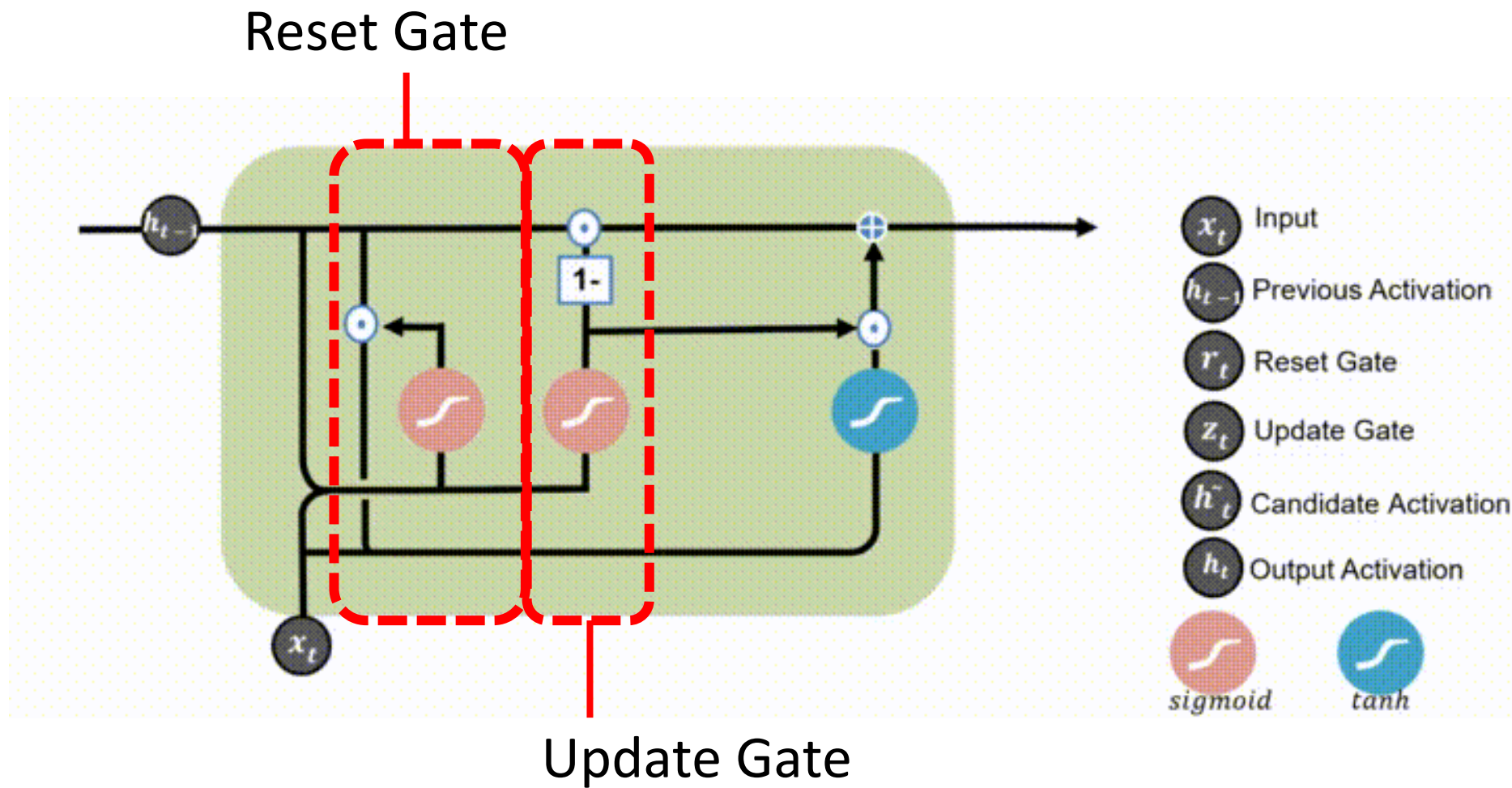
$$r_t = \sigma\left(W_r \cdot [h_{t-1}, x_t]\right) \qquad \in [0,1]$$

$$\tilde{h}_t = \tanh\left(W \cdot [r_t * h_{t-1}, x_t]\right) \quad \in [-1,1]$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

# Information Flow in GRU



Reset Gate

Update Gate

Input $x_t$

Previous Activation $h_{t-1}$

Reset Gate $r_t$

Update Gate $z_t$

Candidate Activation $\tilde{h}_t$

Output Activation $h_t$

sigmoid    tanh

# Information Flow in GRU

**Reset gate** $r_t$

How much of the past information $h_{t-1}$ should be retained with respect to new input $x_t$ to form new $\tilde{h}_t$ candidate

**Update gate** $z_t$

$(1 - z_t)$ How much of the past information $h_{t-1}$ should be discarded

$(z_t)$ How much of new information $\tilde{h}_t$ should make into final $h_t$

**GRU** keeps original RNN input-output structure $(x_t, h_t)$ by letting $z_t$ to handle both data **retention** and **attrition**

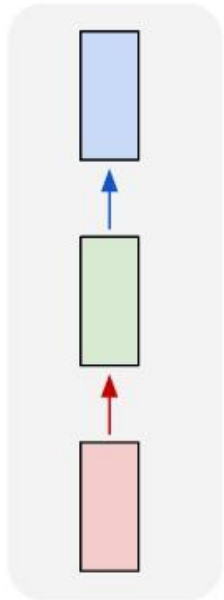# ENCODER-DECODER RNNs

Many-to-Many RNN Recap

Encoder-Decoder Architecture
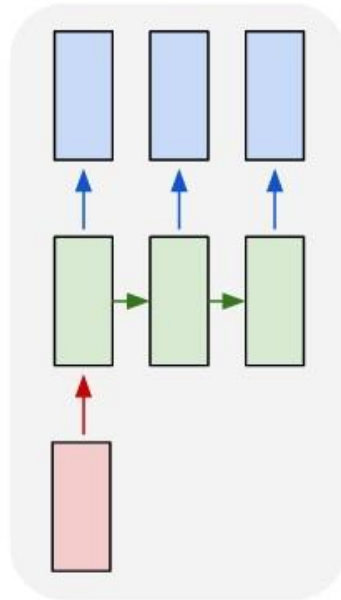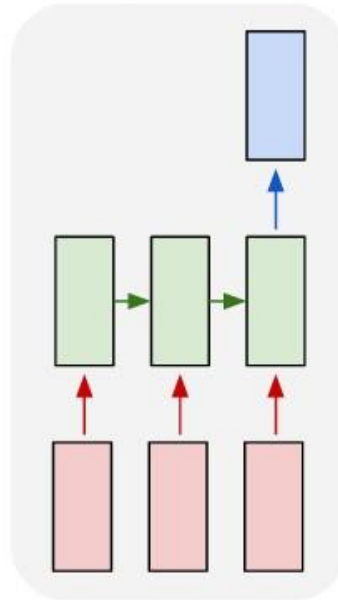
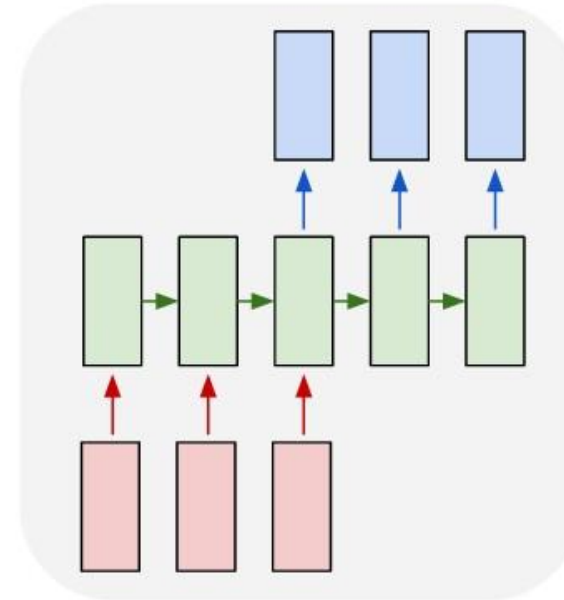Training Encoder-Decoder RNNs
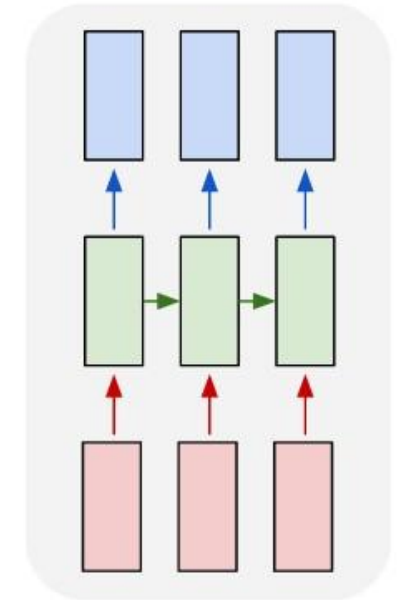
# RNN Configurations
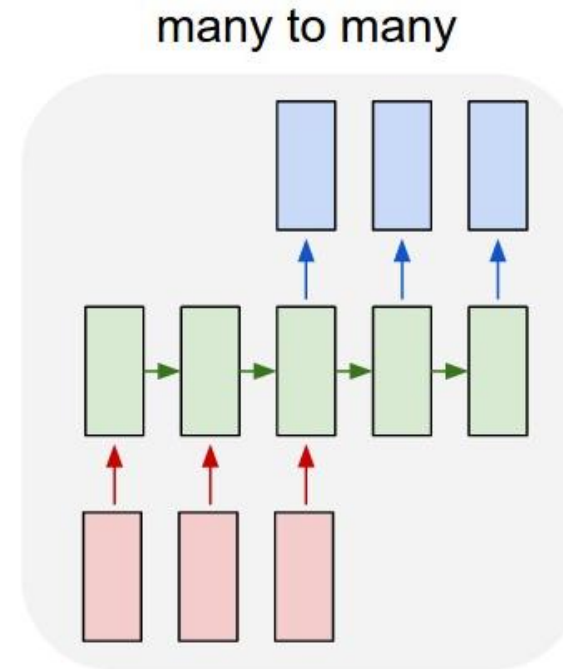


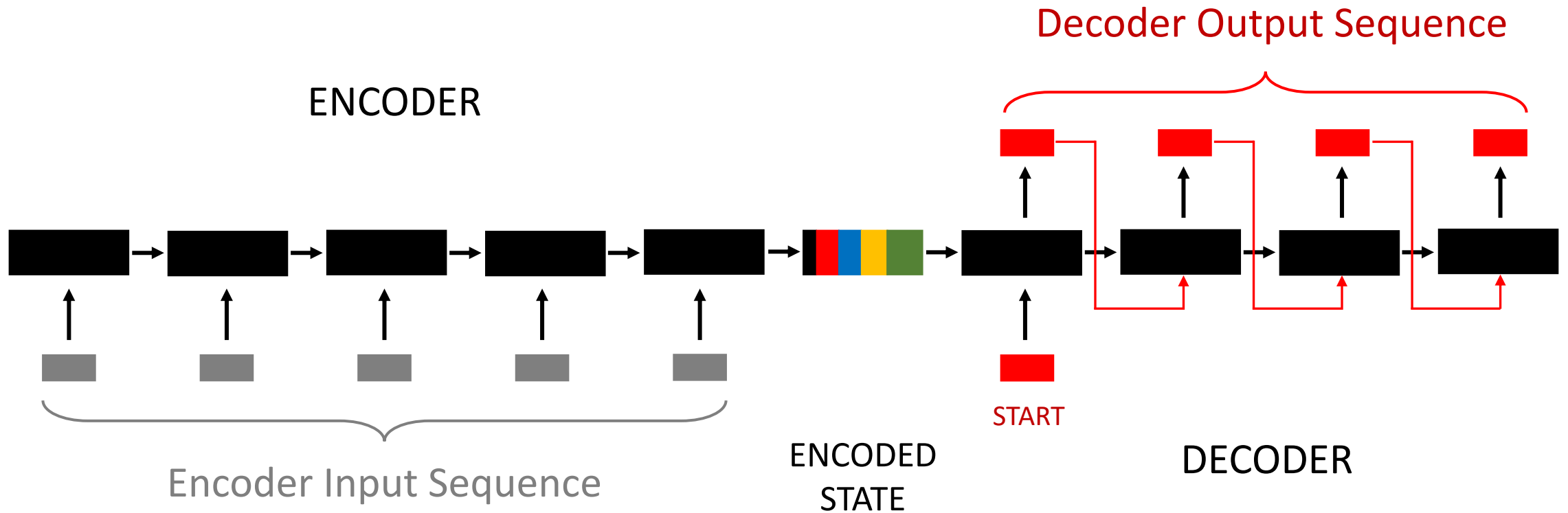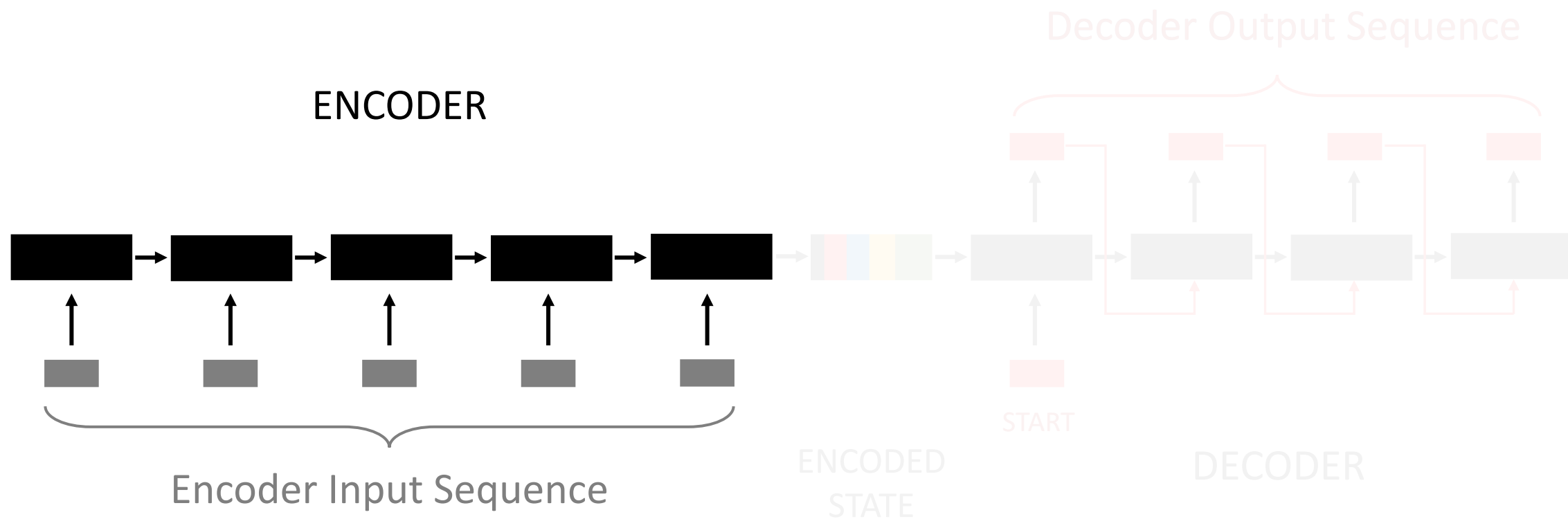one to one    one to many    many to one    many to many    many to many
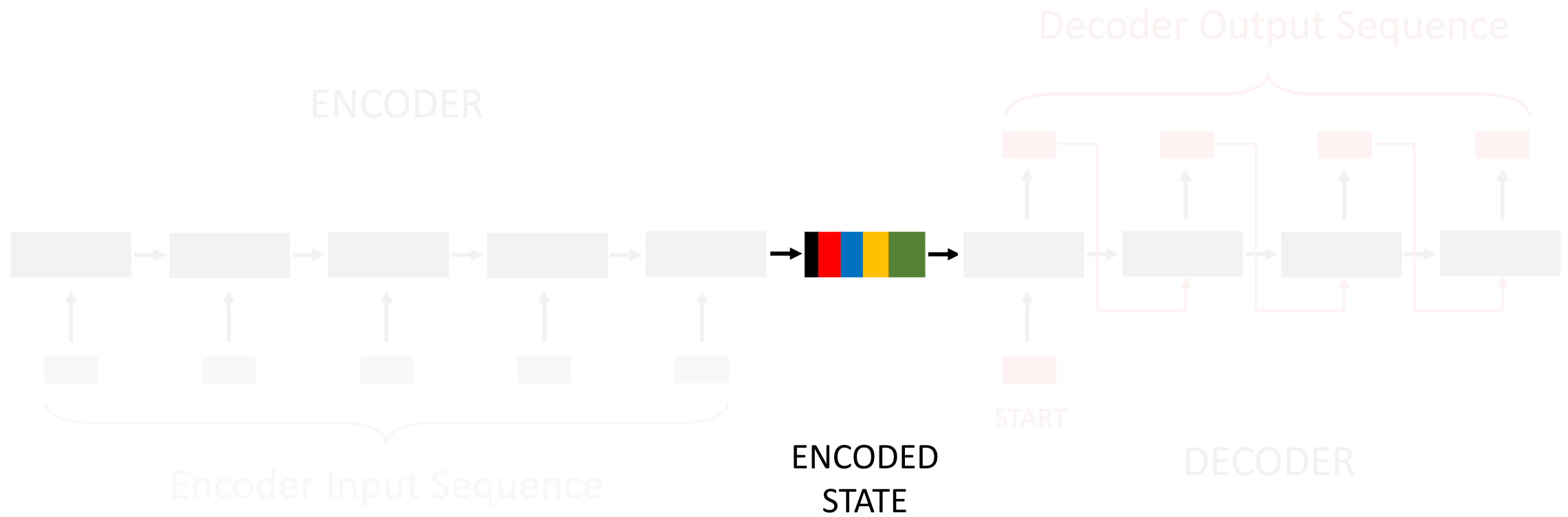
# Many-to-Many



many to many

# Encoder-Decoder Architecture

Decoder Output Sequence

ENCODER



Encoder Input Sequence

ENCODED
STATE

START

DECODER

ENCODER

ENCODER INPUT SEQUENCE

ENCODED STATE

DECODER

Decoder Output Sequence

START

Encoder Input Sequence

# Encoded State

ENCODER

Decoder Output Sequence

ENCODED
STATE

Encoder Input Sequence

START

DECODER

# Decoder



Decoder Output Sequence

ENCODER

Encoder Input Sequence

ENCODED STATE

START

DECODER

# Encoder-Decoder Architecture
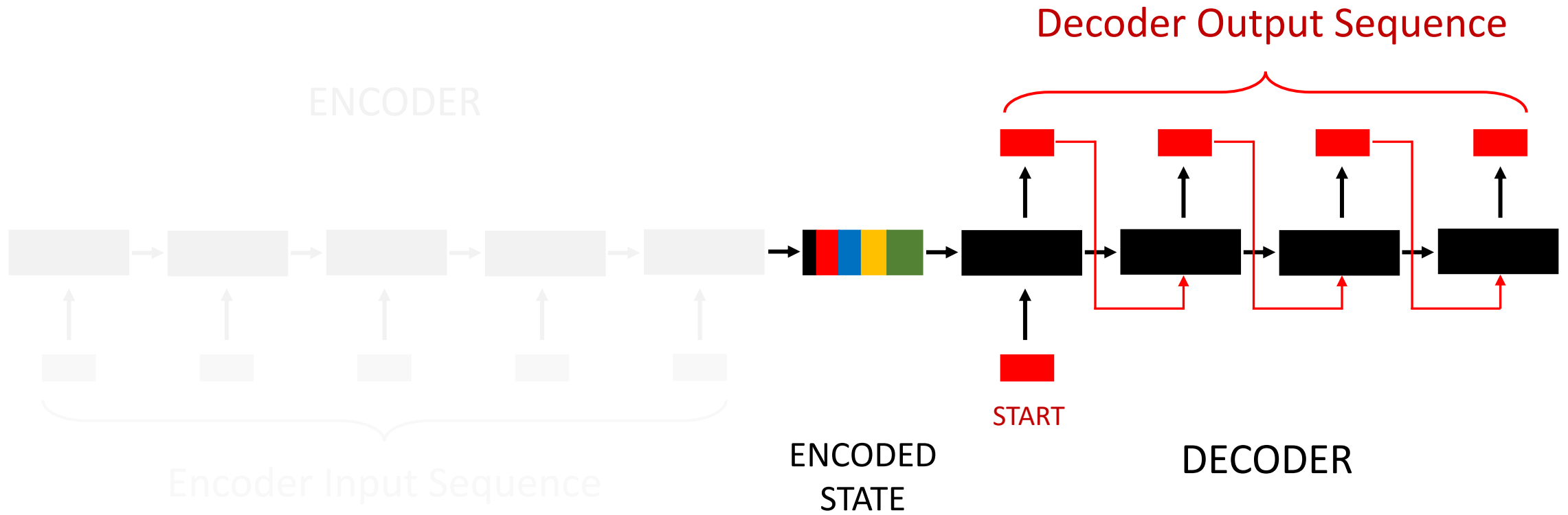
**Decoder Output Sequence**

**ENCODER**



ENCODED
STATE

START

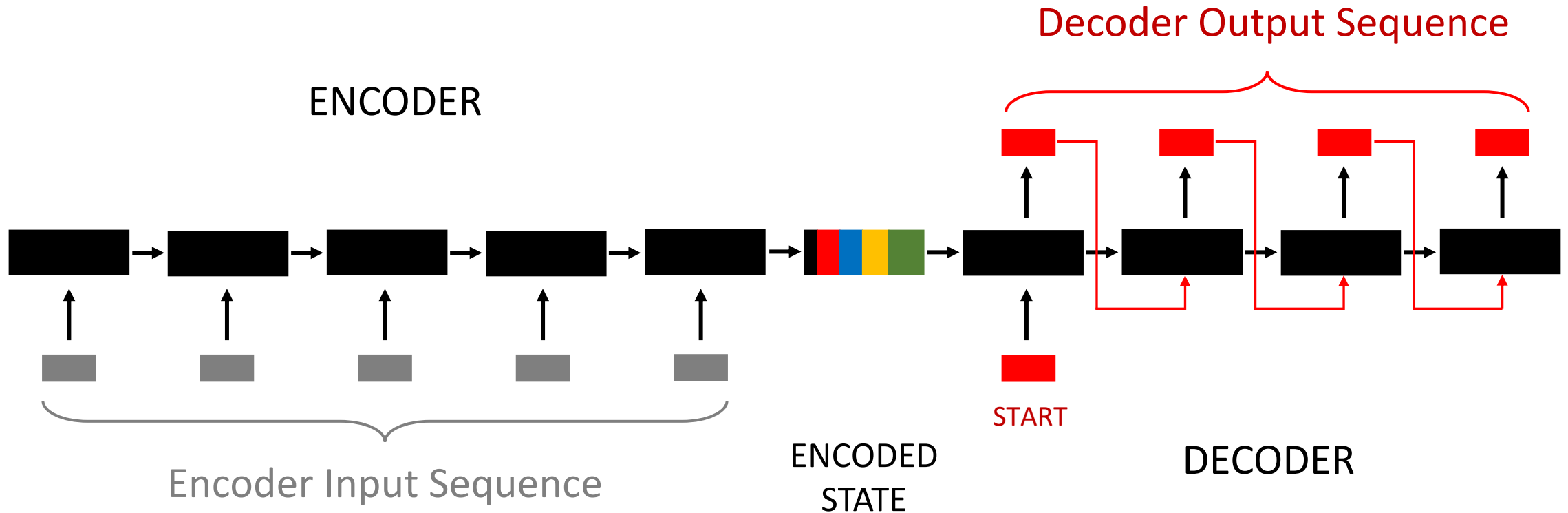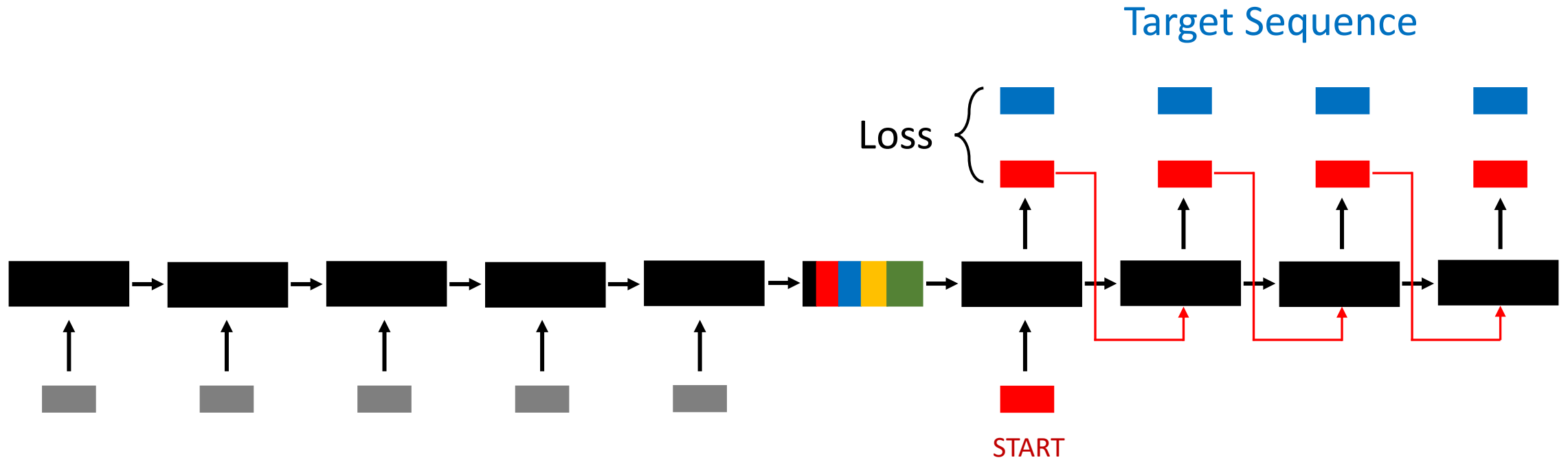DECODER

Encoder Input Sequence

Input sequence length to Encoder (Tx) can be different from the output sequence length of Decoder (Ty)
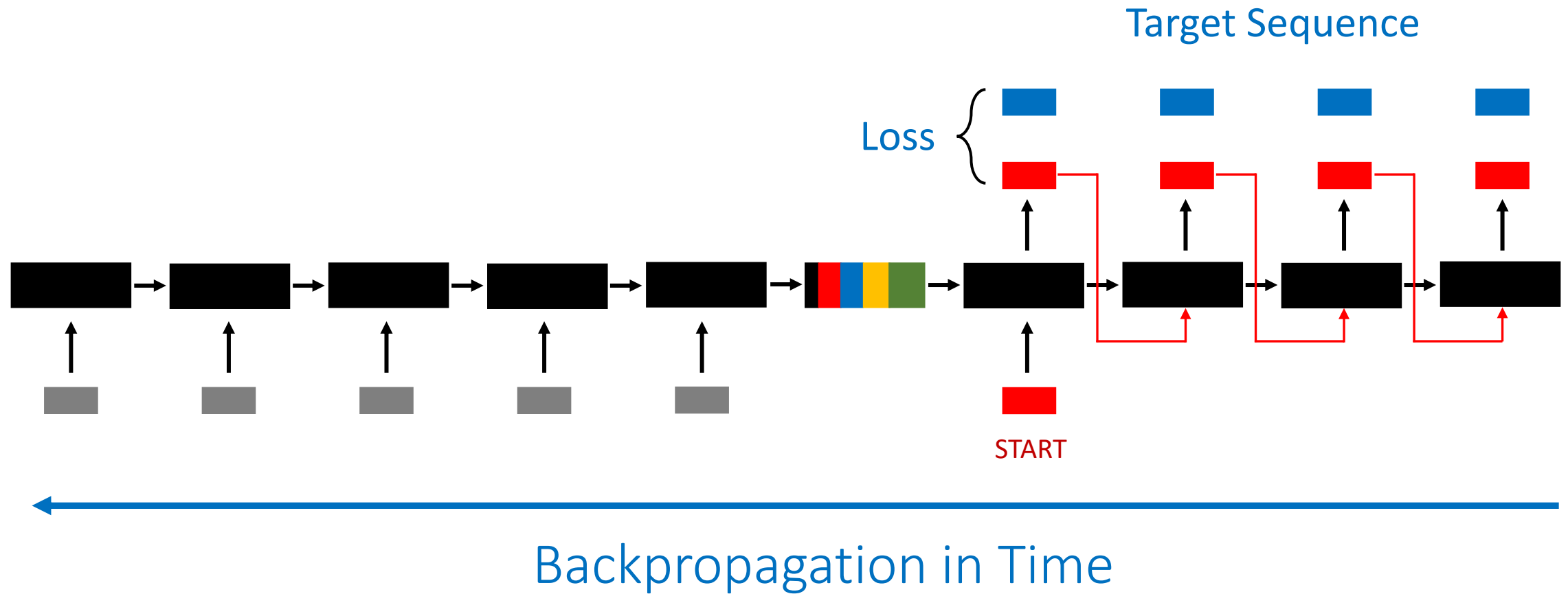
# TRAINING ENCODER-DECODER

# Training Encoder-Decoder

# Training Encoder-Decoder

# Next episode in EEP 596…

Attention and Transformer