# LECTURE 6:
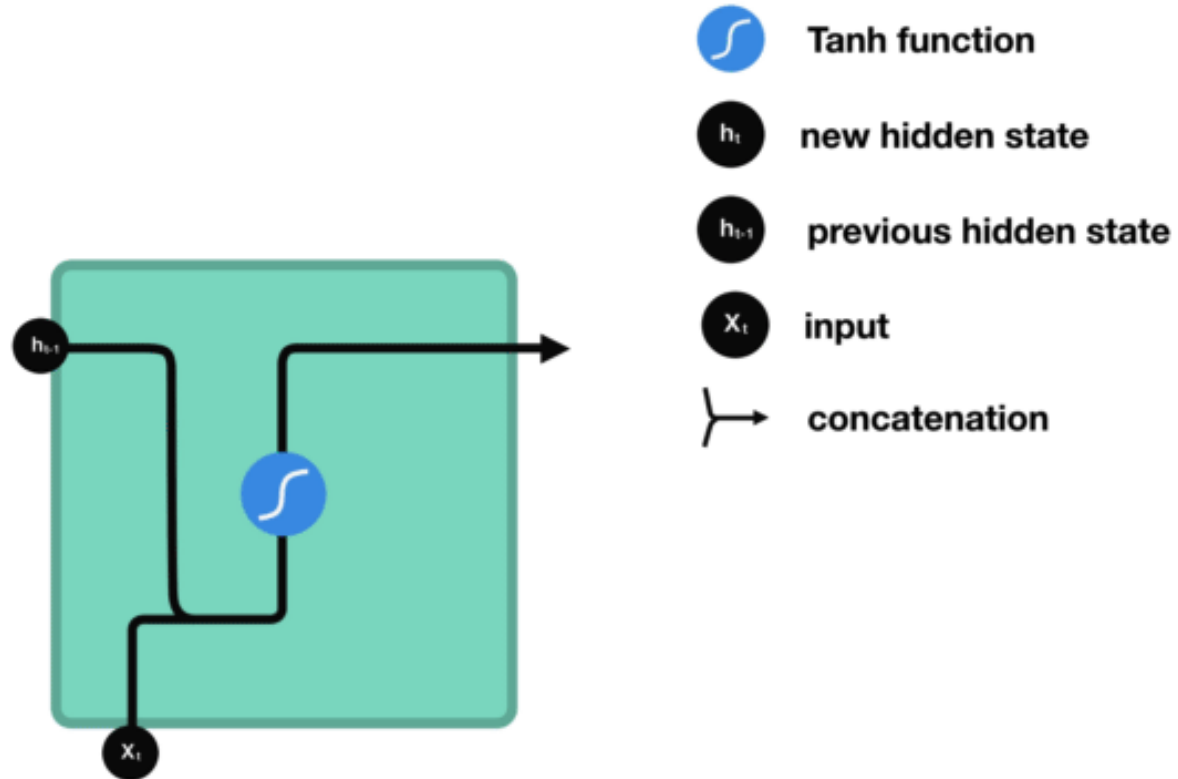
# ADVANCED RECURRENT NEURAL NETWORKS

University of Washington, Seattle

Fall 2024
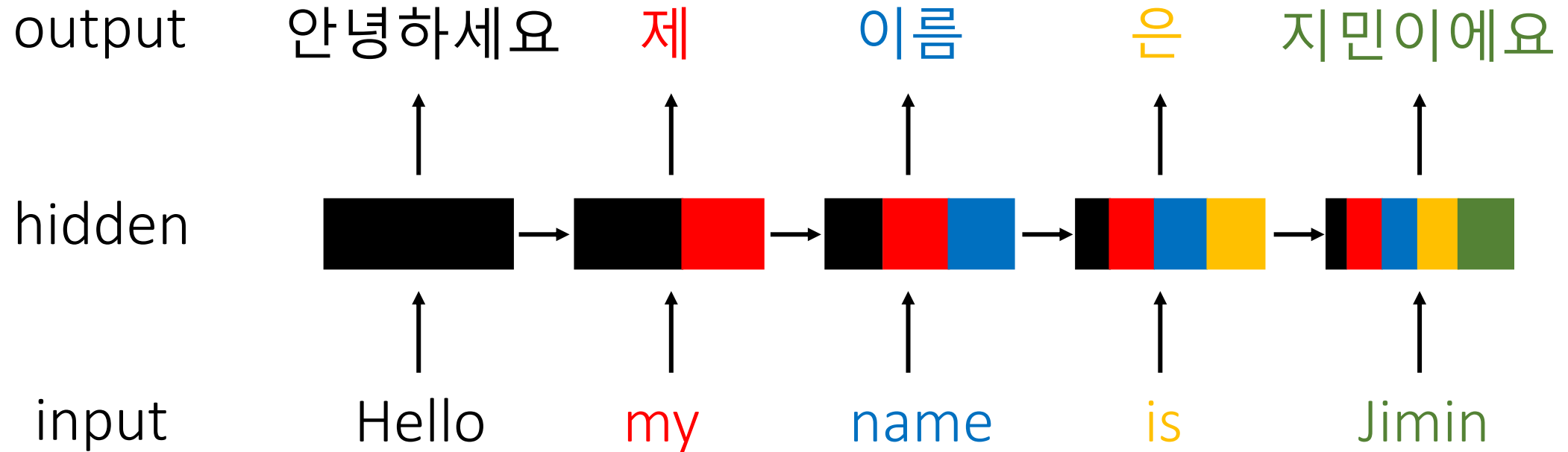
# Previously in EEP 596…

# Previously in EEP 596...

# OUTLINE

**Part 1: Gated RNNs**

- Need for Gated RNNs

- Long Short-Term Memory (LSTM)

- Gated Recurrent Unit (GRU)

**Part 2: Training Gated RNNs**

- Mini-batch Gradient in RNNs

- RNN extensions on LSTM/GRU

**Part 3: Encoder-Decoder RNNs**

- Many to many RNN Recap

- Encoder-Decoder Architecture

- Training Encoder-Decoder RNNs
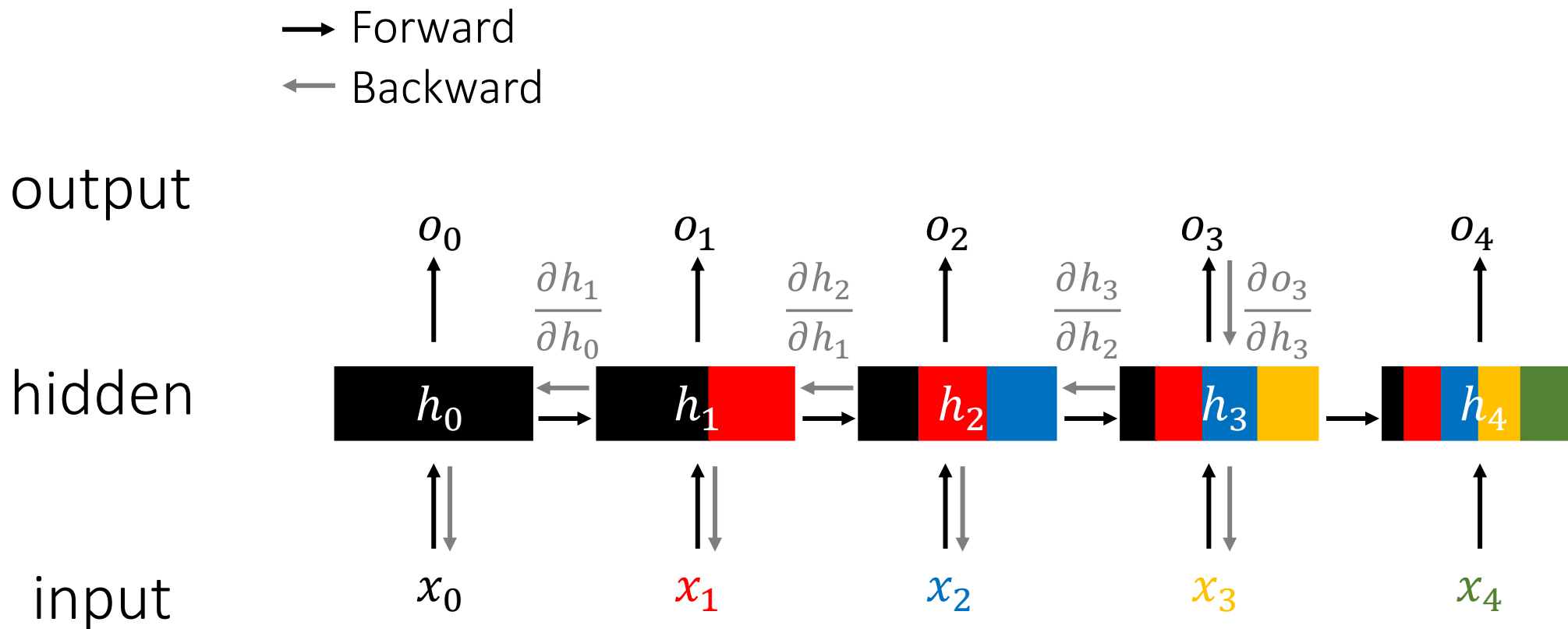
# GATED RNNs

Need for Gated RNNs

Long Short-Term Memory (LSTM)

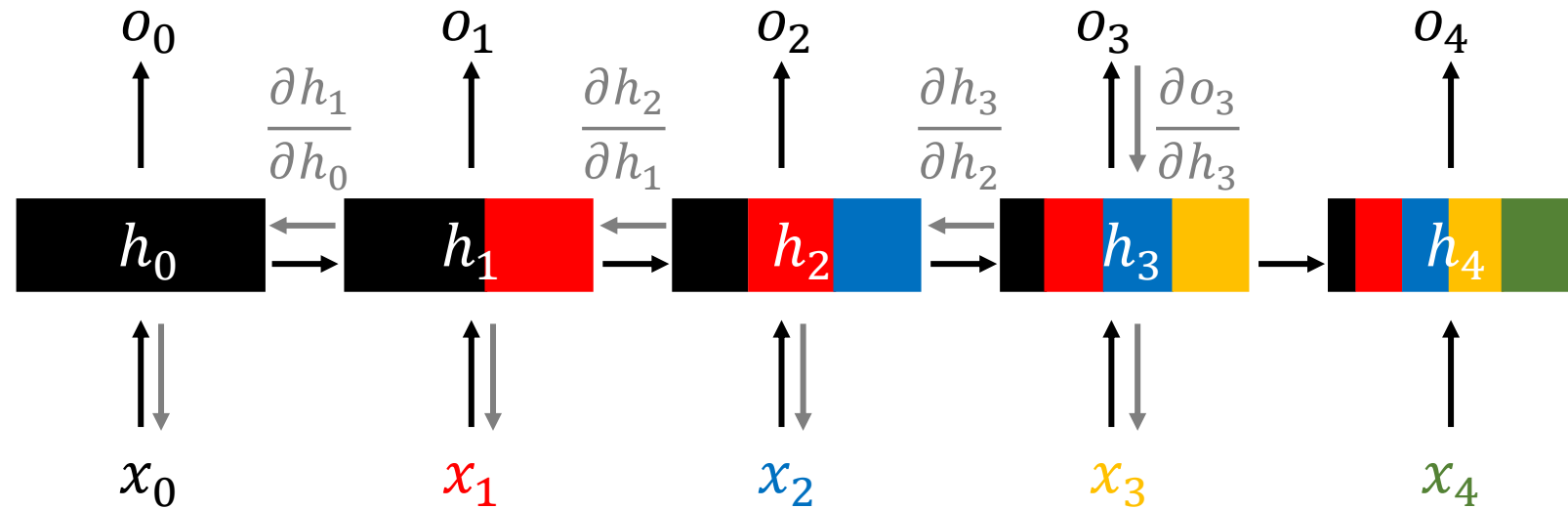Gated Recurrent Unit (GRU)

# Recap: Backpropagation in RNNs

# Recap: Backpropagation in RNNs

→ 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

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

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

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

input

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

Need for better RNN architecture capable of
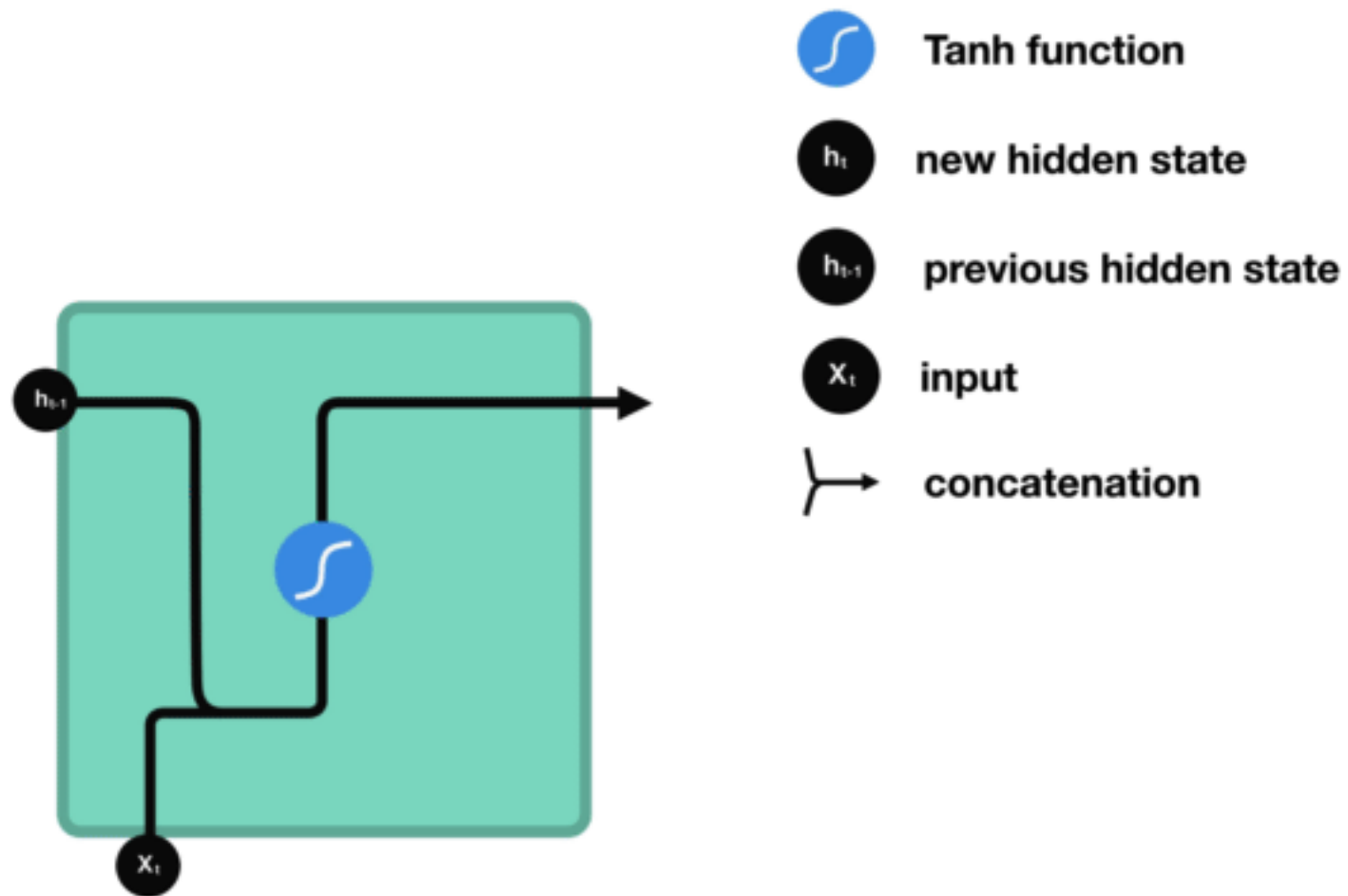processing longer sequence

# Gated RNNs



Vanilla RNN

LSTM

GRU

Neural Network Layer

Pointwise Operation

Vector Transfer

Concatenate

Copy

# Vanilla RNN



Vanilla RNN

# Vanilla RNN

# Vanilla RNN

Tanh function

$h_t$ new hidden state

$h_{t-1}$ previous hidden state

$X_t$ input

concatenation

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

# LSTM (Long Short-Term Memory)



LSTM

# 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
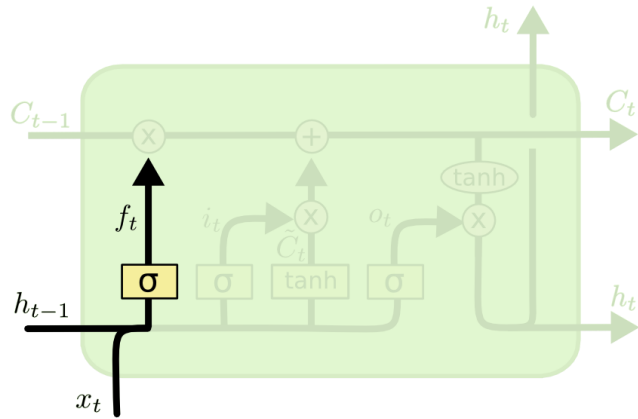
# LSTM: Detailed Architecture



**Cell state**
- Unique to LSTM
- Long term memory of the model

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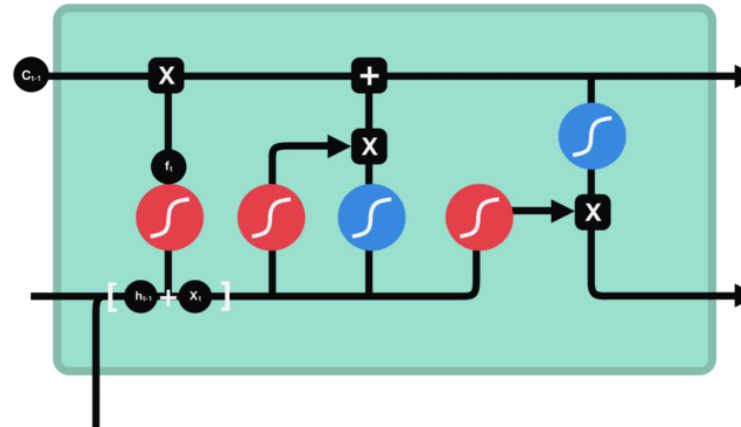
$$f_t = \sigma_g(W_f x_t + U_f h_{t-1} + b_f)$$

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

Forget gate layer
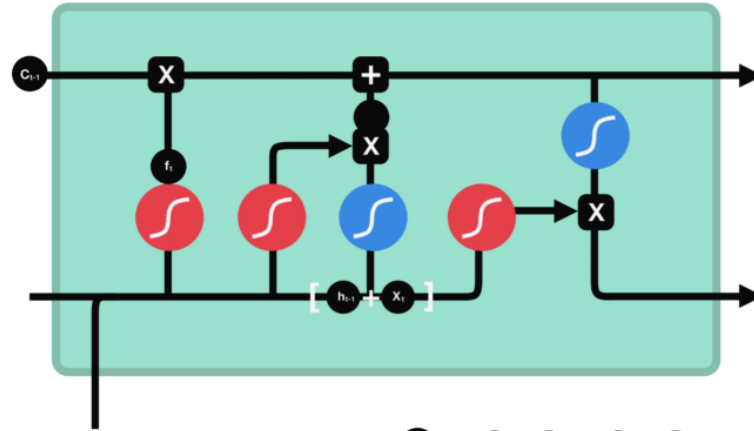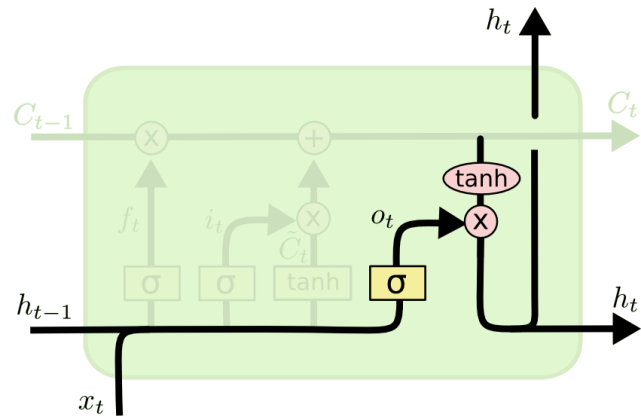
Input gate layer

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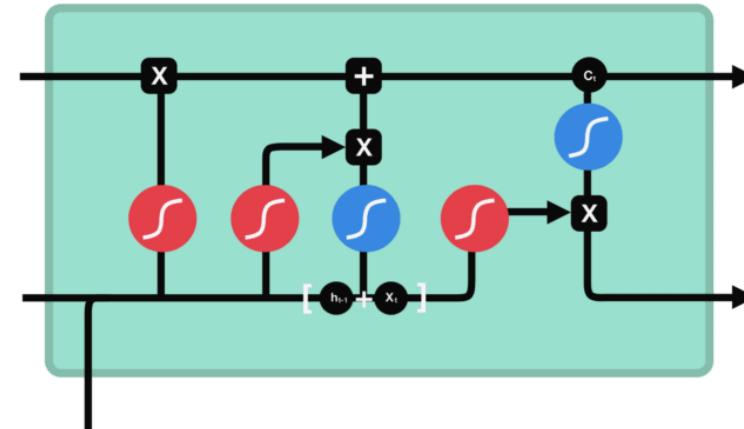
$$o_t = \sigma_g(W_o x_t + U_o h_{t-1} + b_o)$$



Output gate layer

# LSTM: Detailed Architecture

**Forget gate**

Decides what is relevant to keep from previous steps
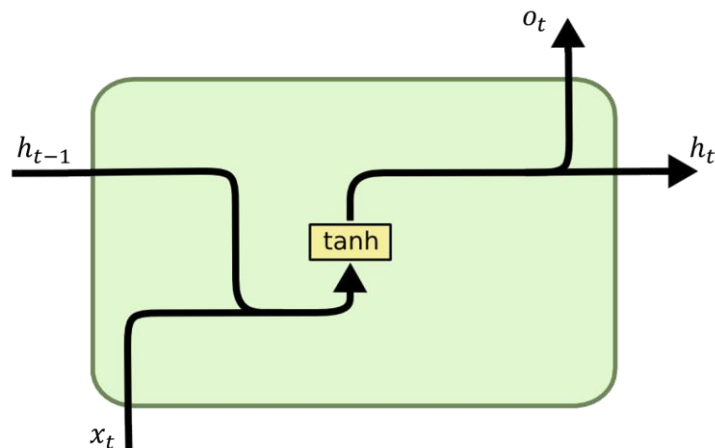
**Input gate**

Decides what information is relevant to add from the current step

**Output Gate**
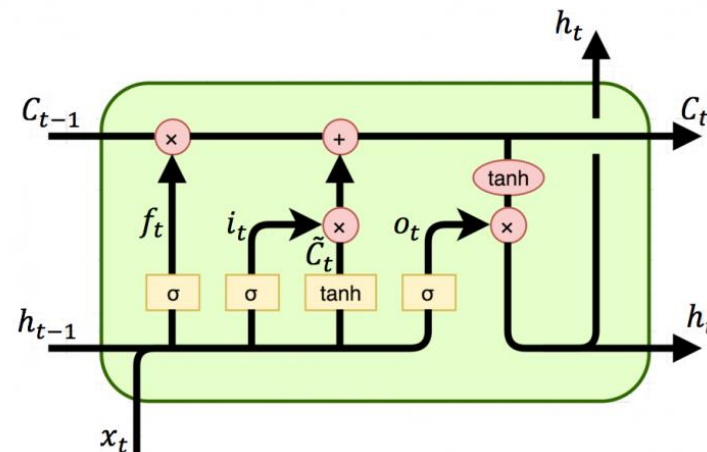
Determines what the next hidden state should be

# LSTM (Long Short-Term Memory)



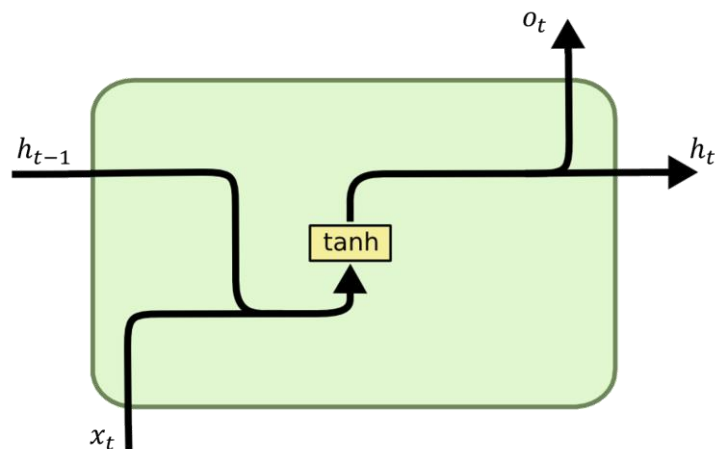Vanilla RNN

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

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

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



LSTM

$$\frac{\partial c_{t'}}{\partial c_t} = \prod_{k=1}^{t'-t} \sigma(v_{t+k}).$$
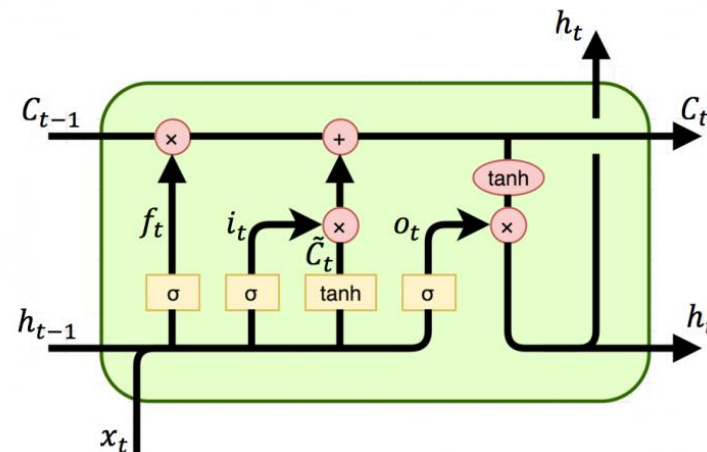
21

# LSTM (Long Short-Term Memory)



Vanilla RNN



LSTM

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

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

$$= \underbrace{w^{t'-t}}_{!!!} \prod_{k=1}^{t'-t} \sigma'(wh_{t'-k})$$
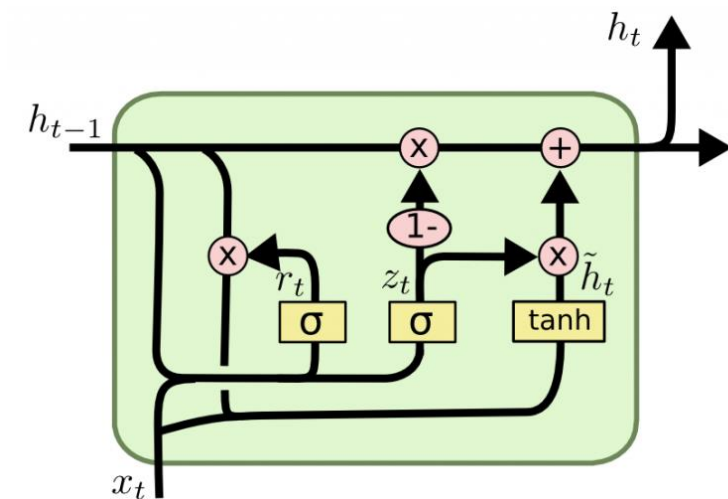
Gradient decays or grow exponentially
if $w \neq 1$

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

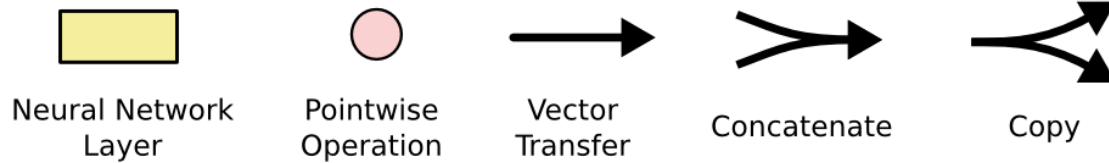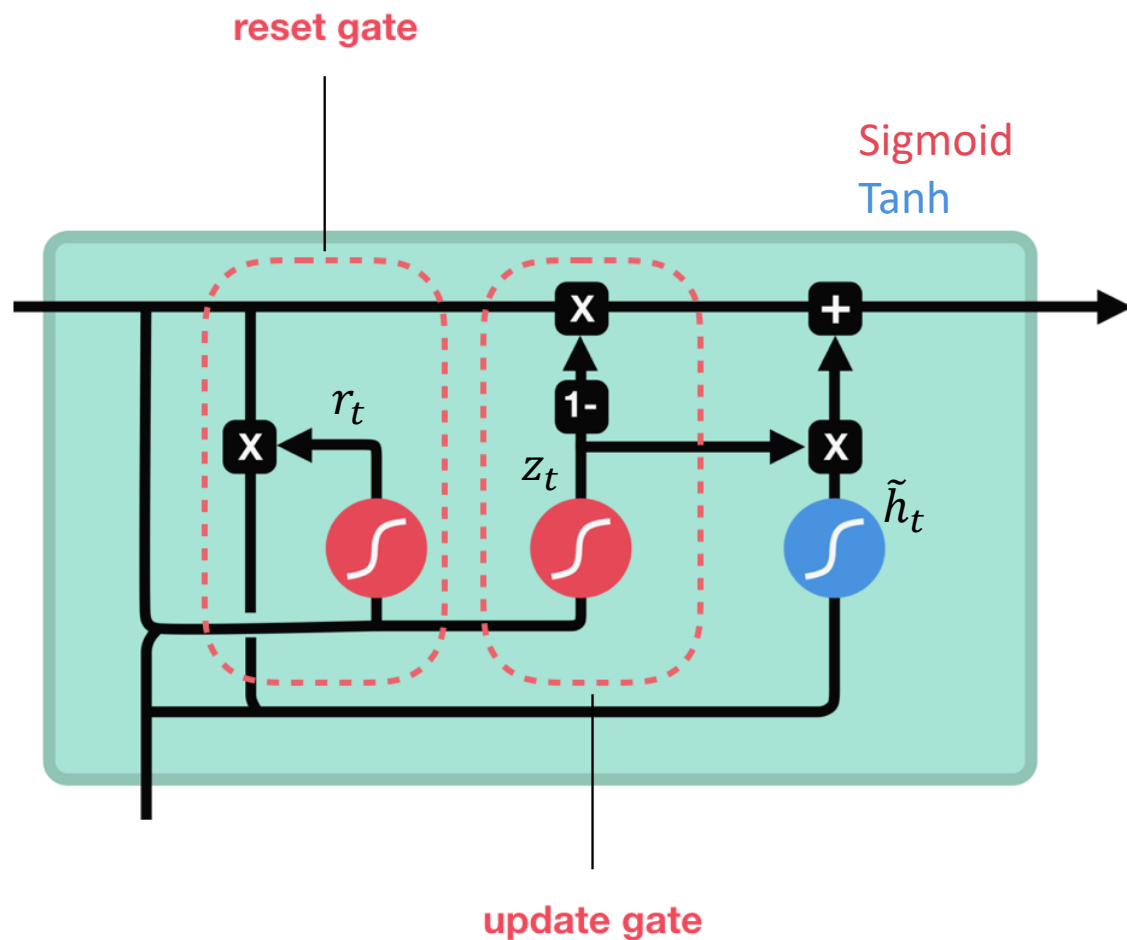No exponential decay or growth term

# Gated RNNs



GRU

Neural Network Layer    Pointwise Operation    Vector Transfer    Concatenate    Copy

# GRU: Detailed Architecture



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

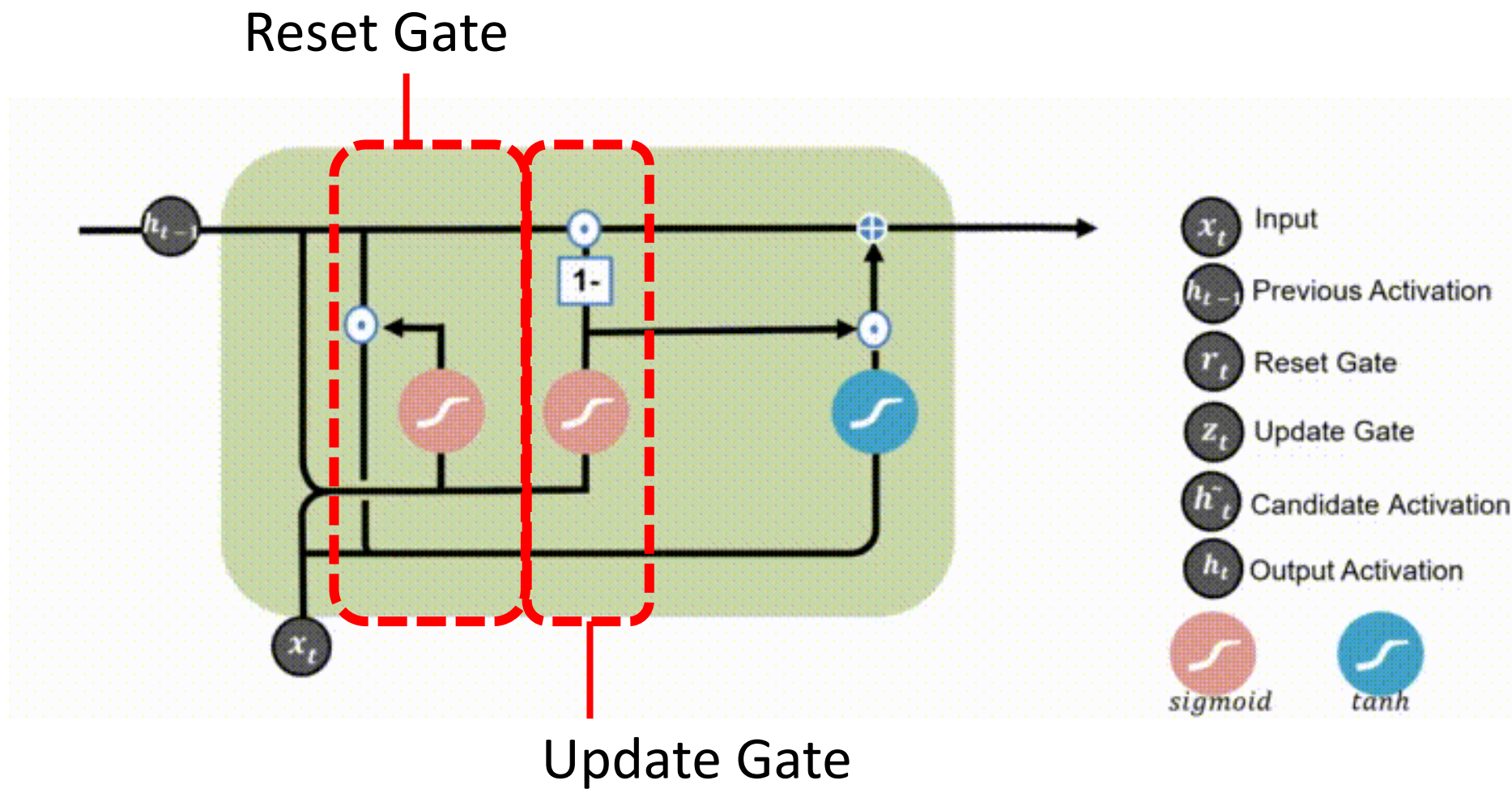$$r_t = \sigma \left( W_r \cdot [h_{t-1}, x_t] \right)$$

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

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

# 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

# LSTM: Detailed Architecture

**Update gate**

How much of the past information needs to be retained

**Reset gate**

How much of the past information to forget

# TRAINING GATED RNNs
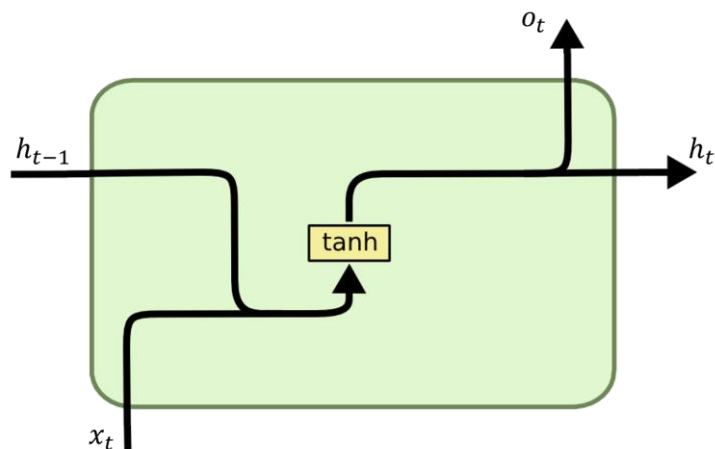
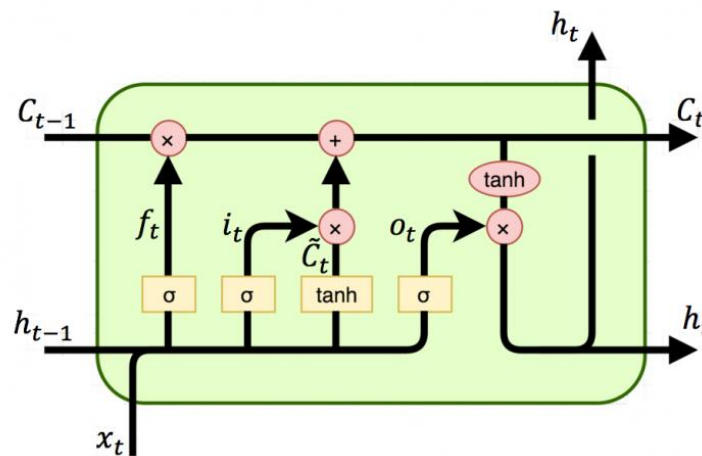Mini-batch Gradient in RNNs

RNN Extensions in LSTM/GRU

# Mini-batch Gradient in RNNs



# of features /timestep

Input 1

Input 2

Input 3

Input 4

batch Size

sequence Length

Mini-batch = set of sequences
Each timestep can be associated with an array

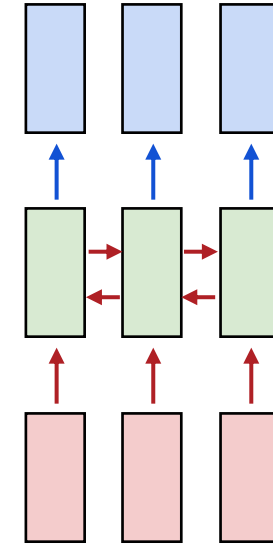# Gated RNNs



Vanilla RNN

LSTM

GRU

Inputs = $x_t$

Outputs = $f(h(t))$

# RNN Extensions in LSTM/GRU



Regular RNN

Deep RNN

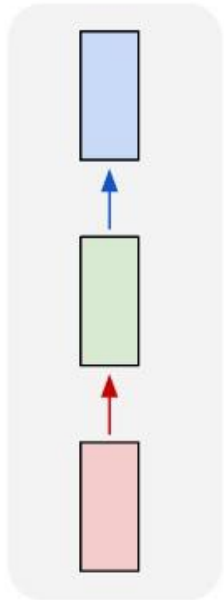Bi-directional RNN

# ENCODER-DECODER RNNs

Many-to-Many RNN Recap

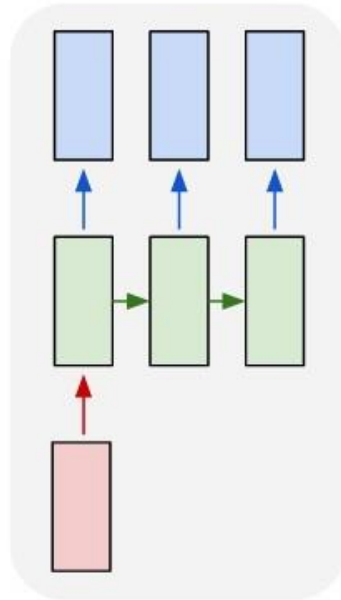Encoder-Decoder Architecture

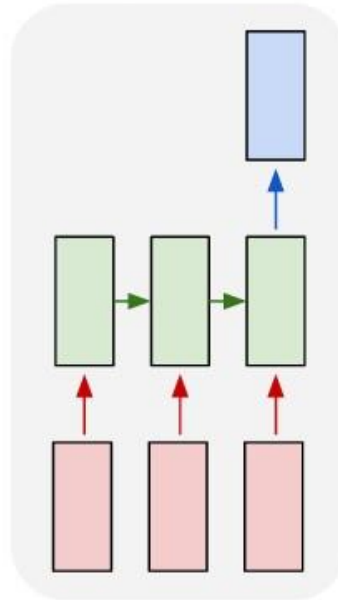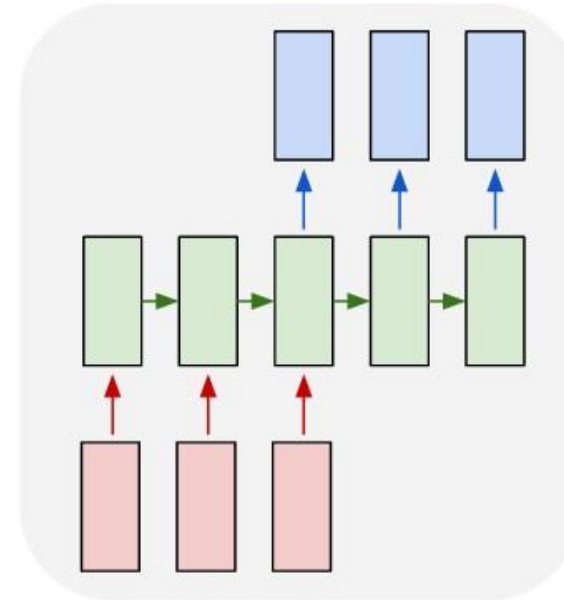Training Encoder-Decoder RNNs

# RNN Configurations

# Many-to-Many



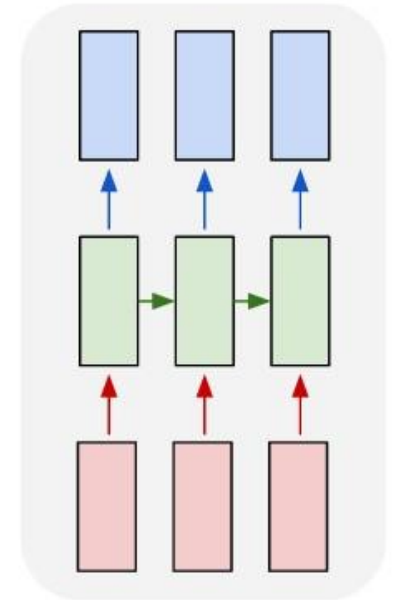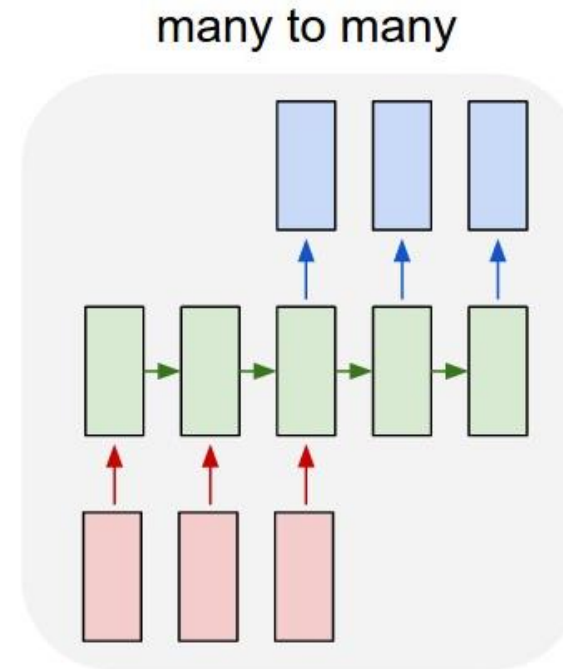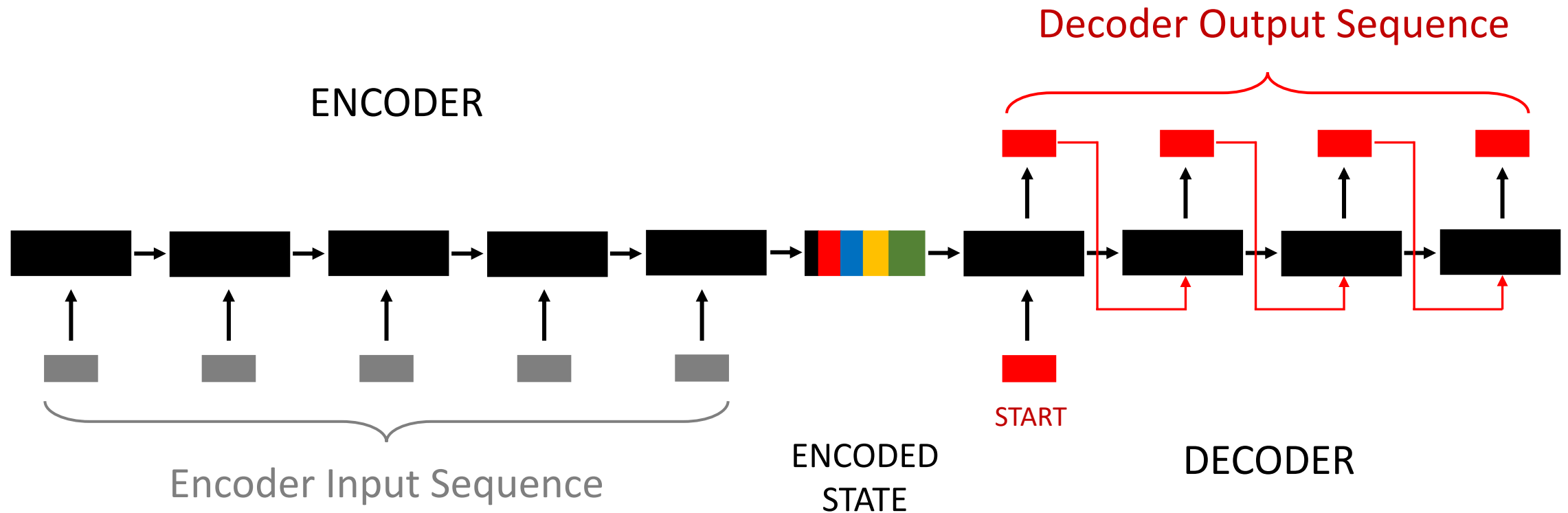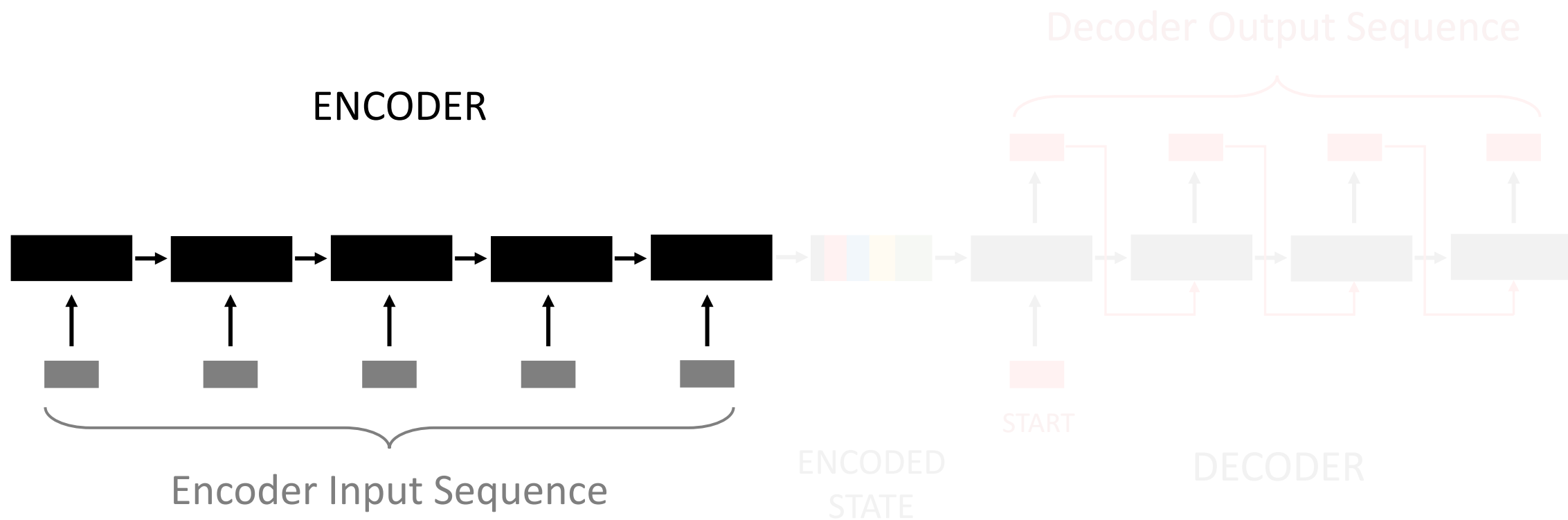many to many

# Encoder-Decoder Architecture

# Encoder

ENCODER

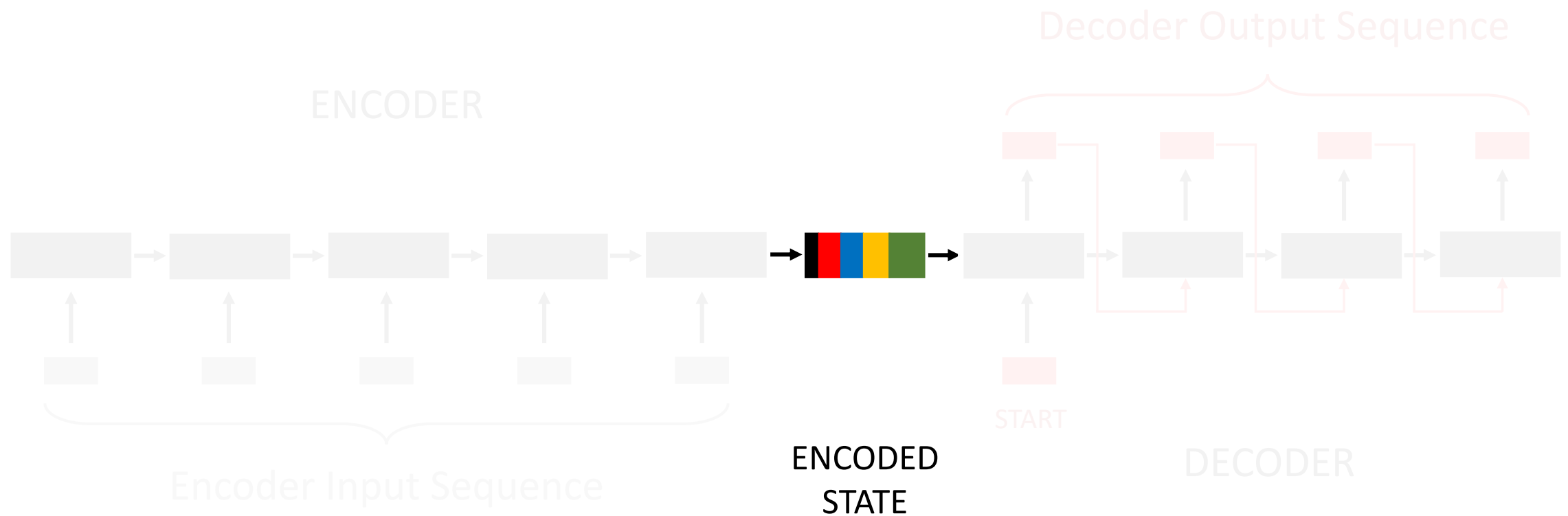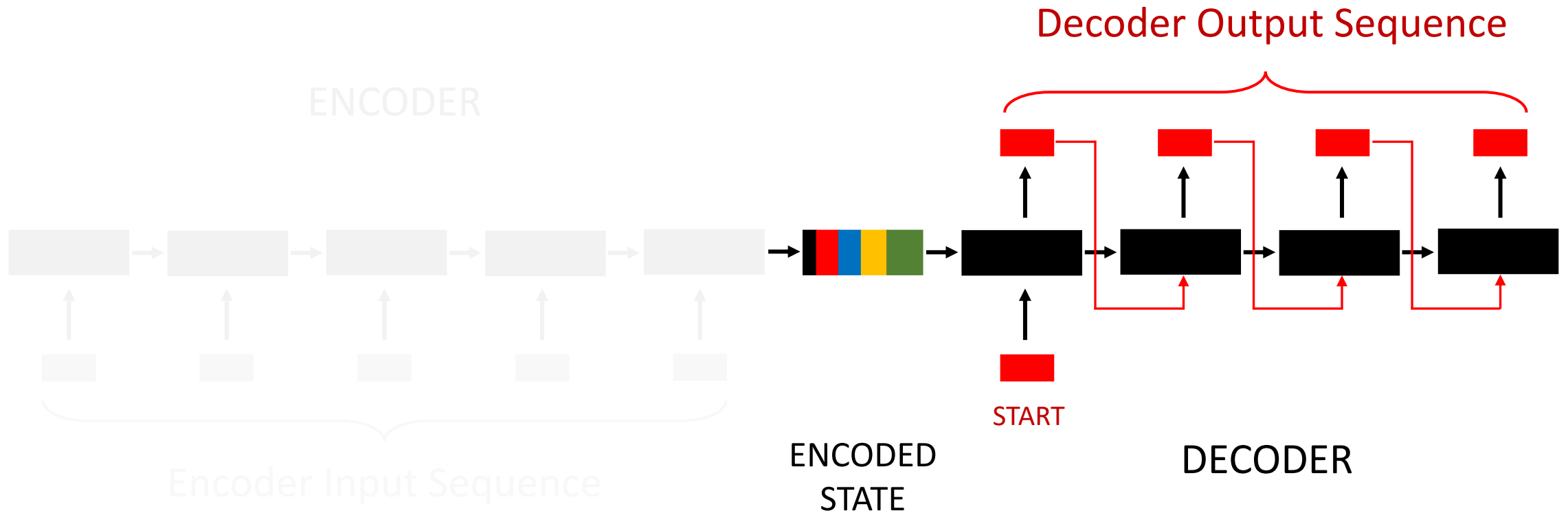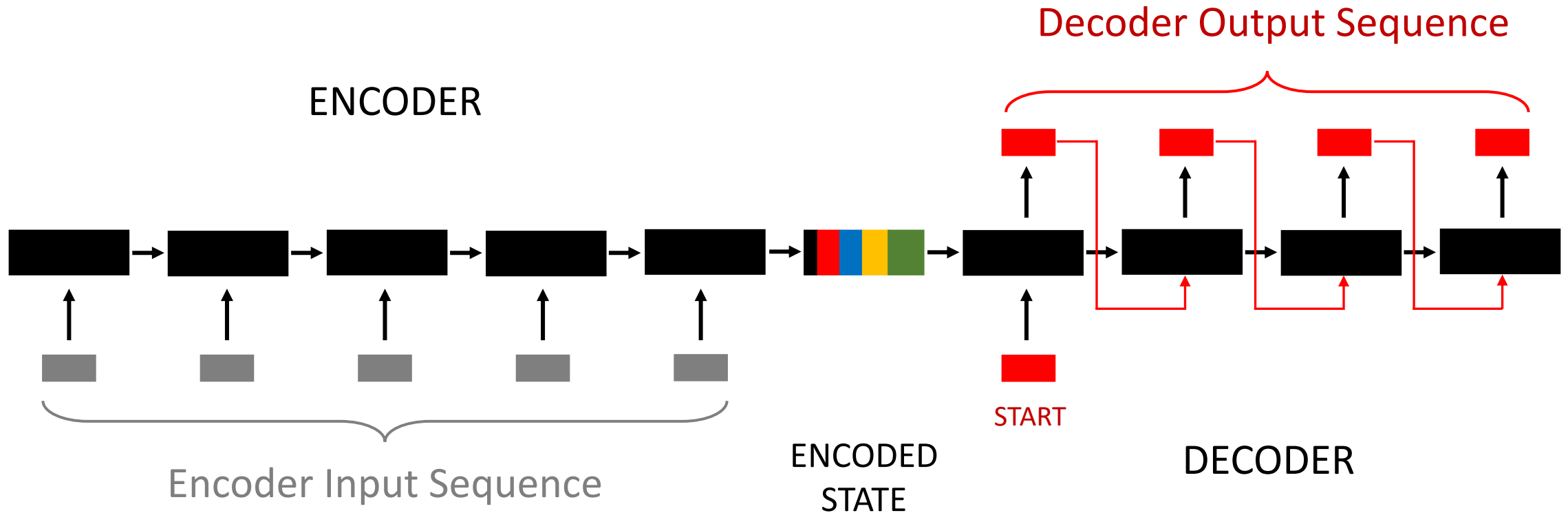Decoder Output Sequence

START

ENCODED
STATE

DECODER

Encoder Input Sequence

# Encoded State

# Decoder

# Encoder-Decoder Architecture

**Decoder Output Sequence**
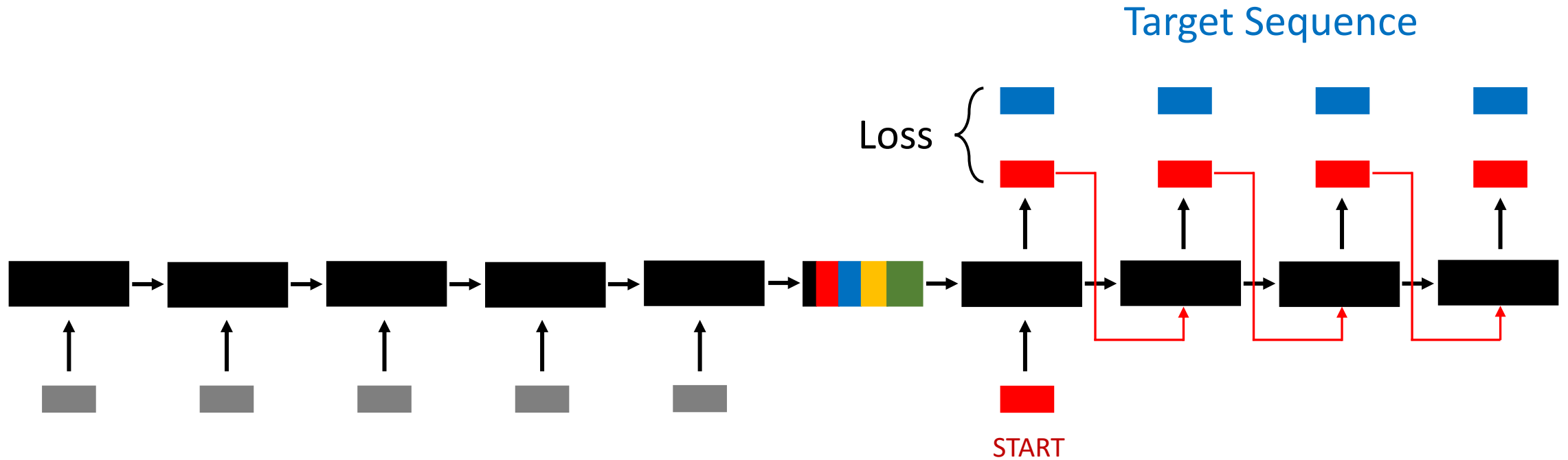
ENCODER

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

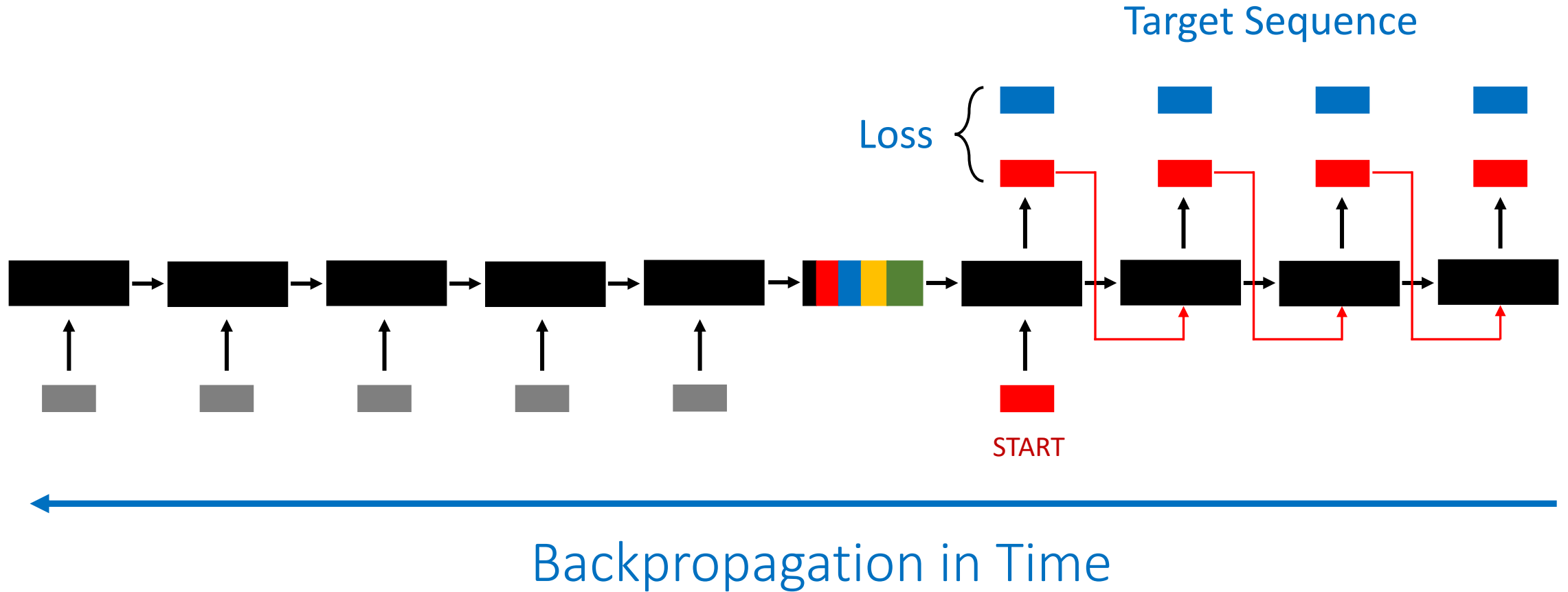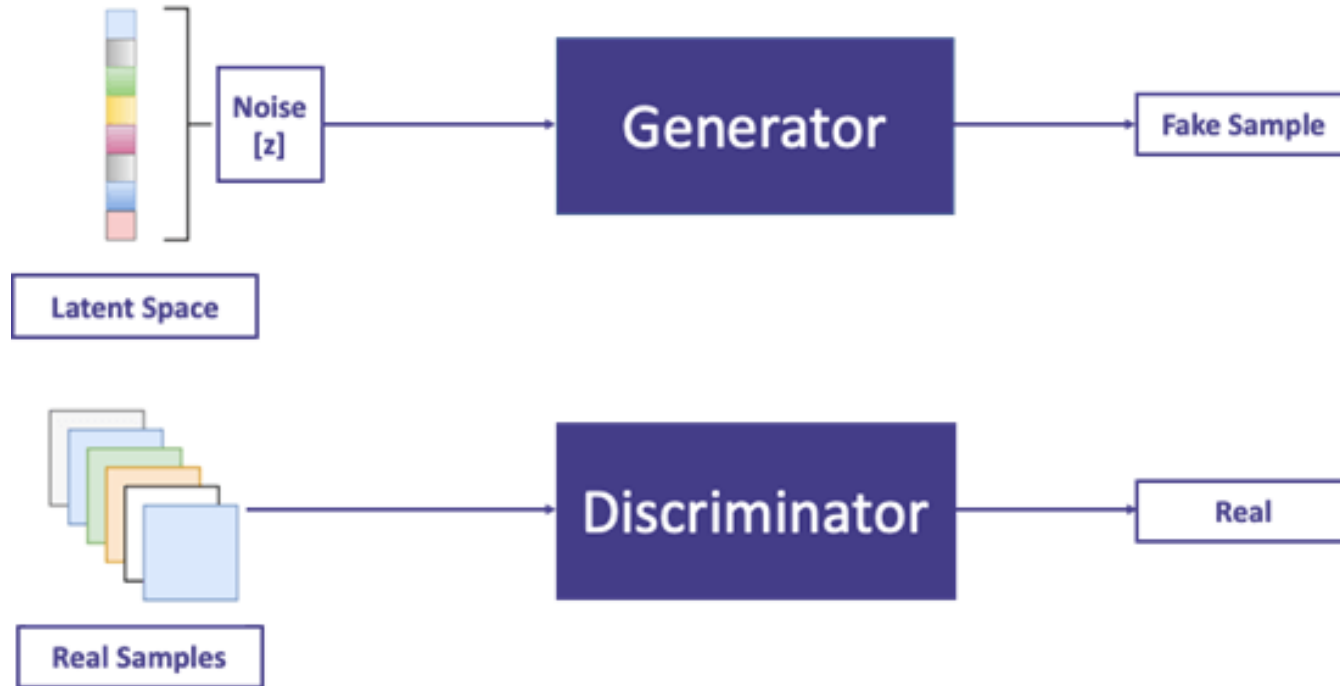# Training Encoder-Decoder



Target Sequence

Loss

START

# Training Encoder-Decoder



Target Sequence

Loss

START

Backpropagation in Time

# Next episode in EEP 596…



**Generative Adversarial Networks**