# Deep Learning (1470)

## Randall Balestriero

**Class 11: Sequential Data and Language Modeling**

# Recap!

# Word Tokenizer

T h e y | w e n t | t o | t h e | g r o c e r y | s t o r e | a n d | b o u g h t | b r e a d

*Split on spaces → 9 tokens*

| "They" | "went" | "to" | "the" | "grocery" | "store" | "and" | "bought" | "bread" |
|--------|--------|------|-------|-----------|---------|-------|----------|---------|
| 0 | 0 | **1** | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | **1** | 0 | 0 | 0 |
| **1** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | **1** |
| 0 | 0 | 0 | **1** | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | **1** | 0 | 0 |
| 0 | **1** | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | **1** | 0 | 0 | 0 | 0 |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | **1** | ⋮ |

*vocab_size ≈ 50,000*

# Character Tokenizer

T | h | e | y | ␣ | w | e | n | t | ␣ | t | o | ␣ | t | h | e | ␣ | g | r | o    ...

*Split on each character → 47 tokens*

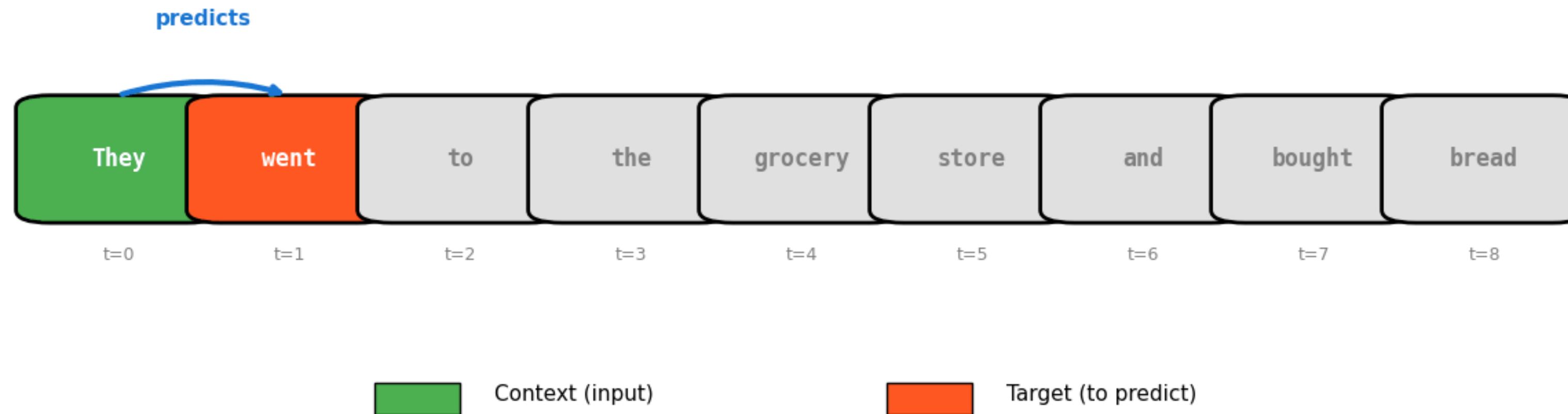**Vocabulary:**    '␣'  'a'  'b'  'c'  'd'  'e'  'g'  'h'  'n'  'o'  'r'  's'  't'  'u'  'w'  'y'    *(size = 16)*

0    1    2    3    4    5    6    7    8    9    10    11    12    13    14    15    16    17    18    19



$dim = 16$

Average Sequence Length vs Vocabulary Size
for Different Tokenizers

# Unigram (N=1): Predict next token from 1 previous token

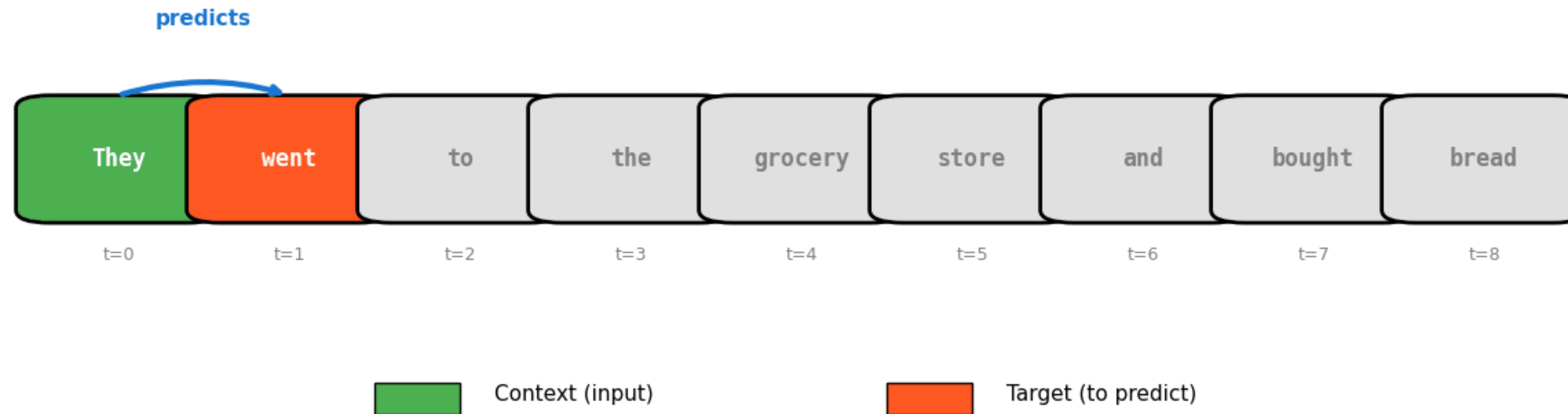**predicts**

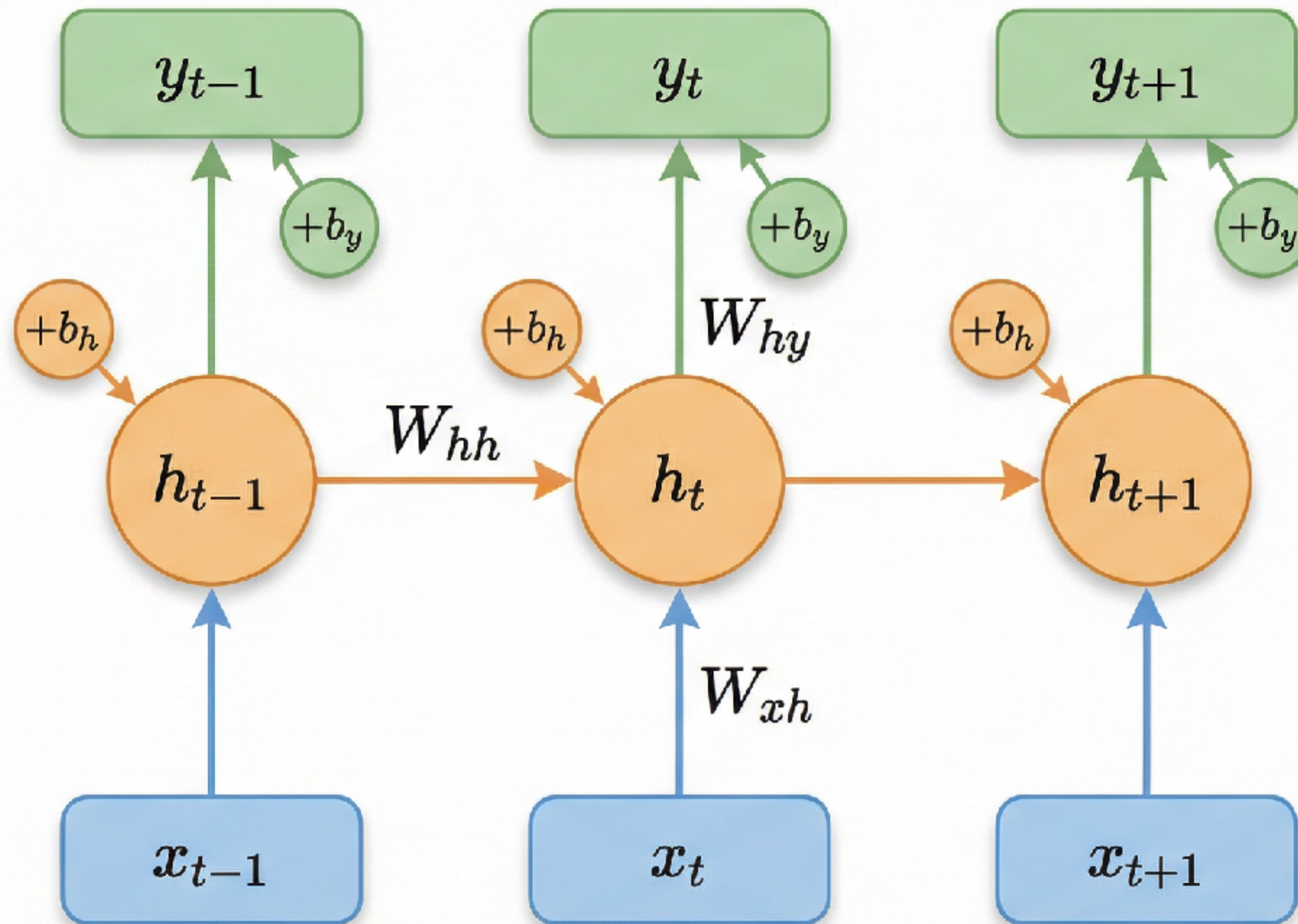| They | went | to | the | grocery | store | and | bought | bread |
|------|------|-----|-----|---------|-------|-----|--------|-------|
| t=0 | t=1 | t=2 | t=3 | t=4 | t=5 | t=6 | t=7 | t=8 |

■ Context (input)    ■ Target (to predict)

# Full History: Predict next token from ALL previous tokens

**all context predicts next token**

| They | went | to | the | grocery | store | and | bought | bread |
|------|------|-----|-----|---------|-------|-----|--------|-------|
| t=0 | t=1 | t=2 | t=3 | t=4 | t=5 | t=6 | t=7 | t=8 |

■ Context (all previous)    ■ Target (to predict)

*Context size: 1*

# Unigram (N=1): Predict next token from 1 previous token

predicts

| They | went | to | the | grocery | store | and | bought | bread |
|------|------|-----|-----|---------|-------|-----|--------|-------|
| t=0 | t=1 | t=2 | t=3 | t=4 | t=5 | t=6 | t=7 | t=8 |

■ Context (input)    ■ Target (to predict)

# Full History: Predict next token from ALL previous tokens

all context predicts next token

| They | went | to | the | grocery | store | and | bought | bread |
|------|------|-----|-----|---------|-------|-----|--------|-------|
| t=0 | t=1 | t=2 | t=3 | t=4 | t=5 | t=6 | t=7 | t=8 |

■ Context (all previous)    ■ Target (to predict)

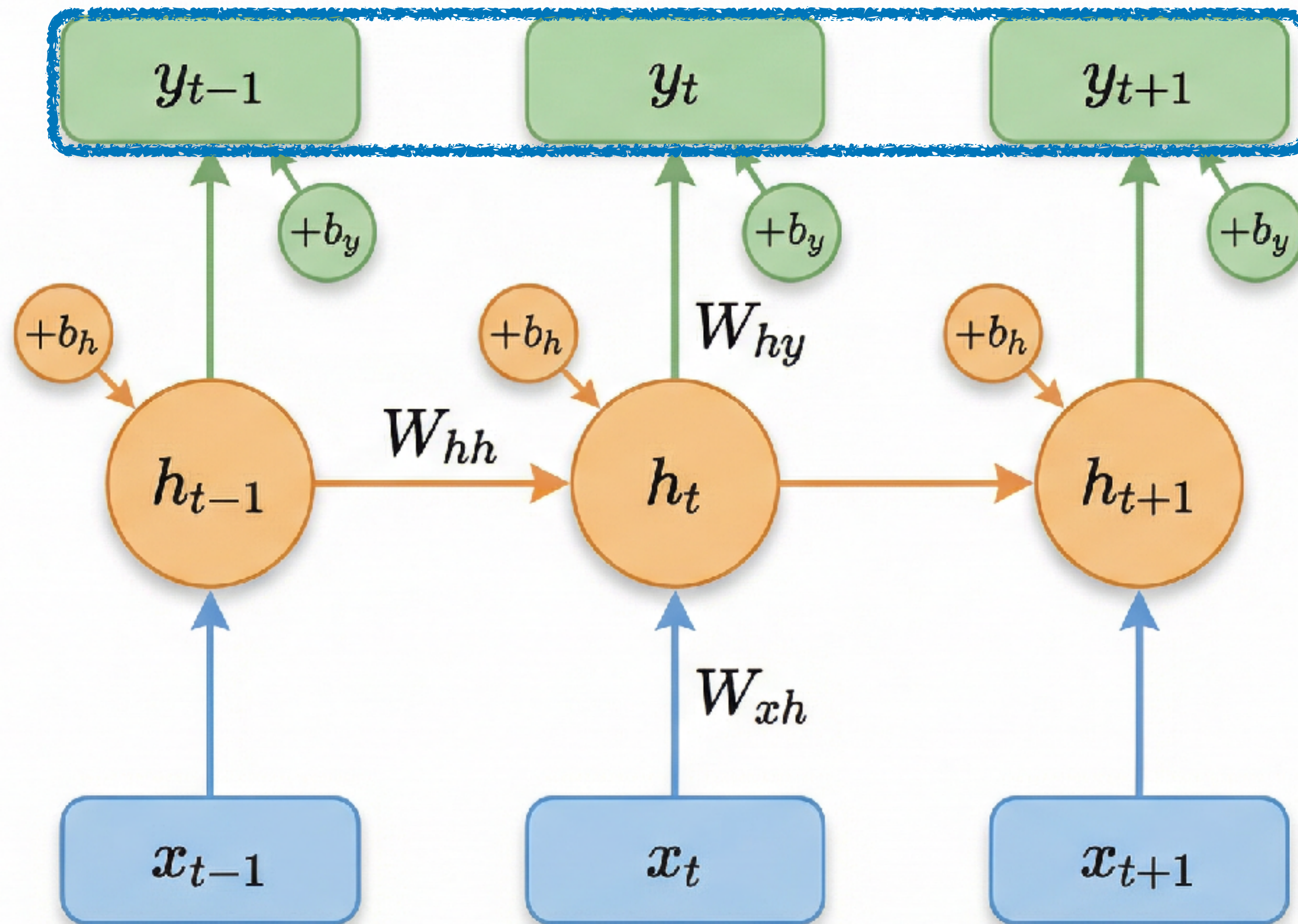*Context size: 1*

# Vanilla RNN Layer



$$h_t = \tanh(W_{xh} \cdot x_t + W_{hh} \cdot h_{t-1} + b_h)$$
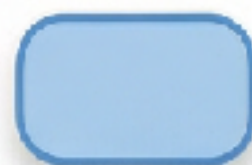$$y_t = W_{hy} \cdot h_t + b_y$$

Input Vector ($x$)     Hidden State ($h$)     Output Vector ($y$)
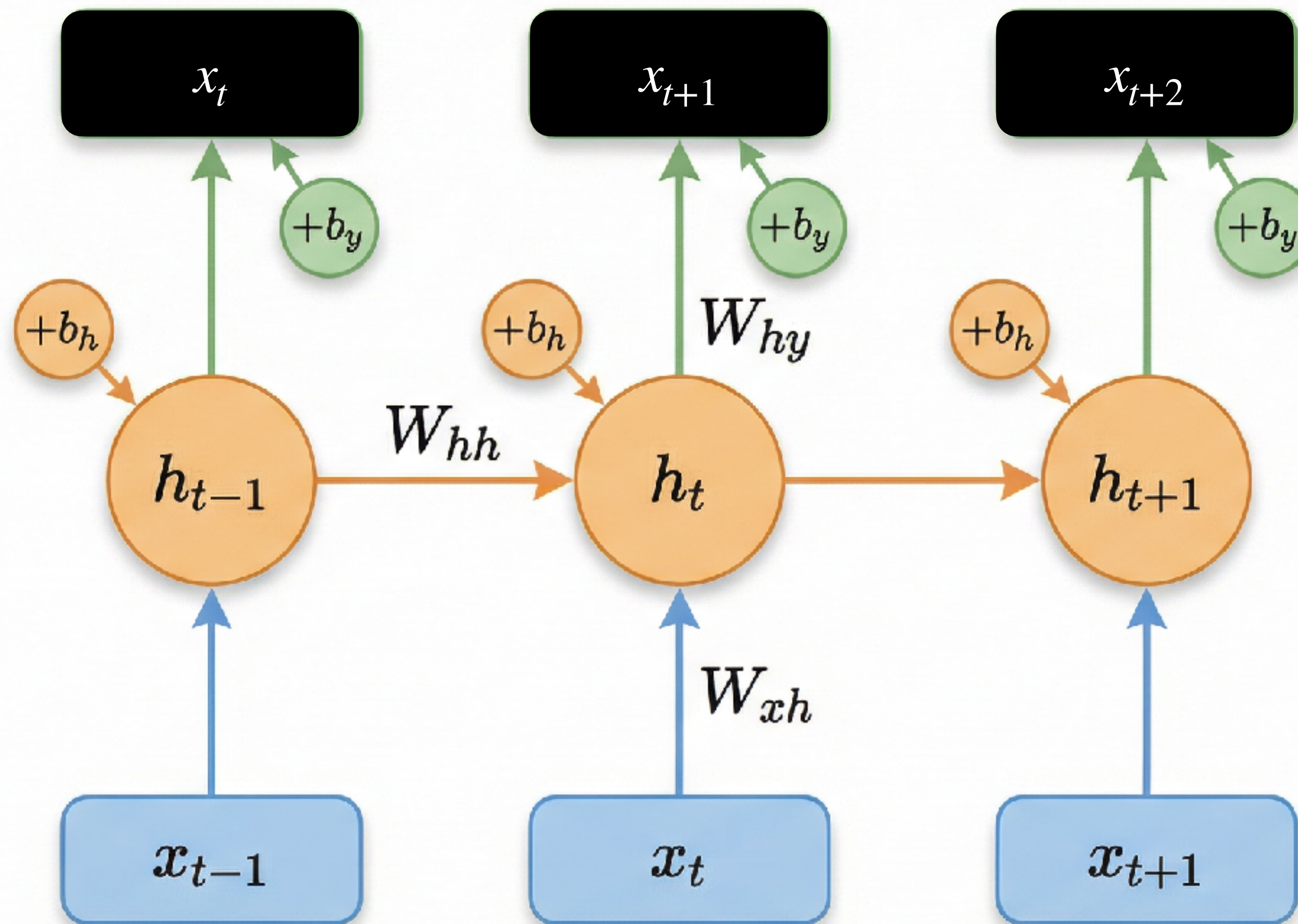
# Vanilla RNN Layer



**What are those?**

$$h_t = \tanh(W_{xh} \cdot x_t + W_{hh} \cdot h_{t-1} + b_h)$$

$$y_t = W_{hy} \cdot h_t + b_y$$

Input Vector ($x$)     Hidden State ($h$)     Output Vector ($y$)

# Vanilla RNN Layer


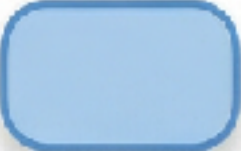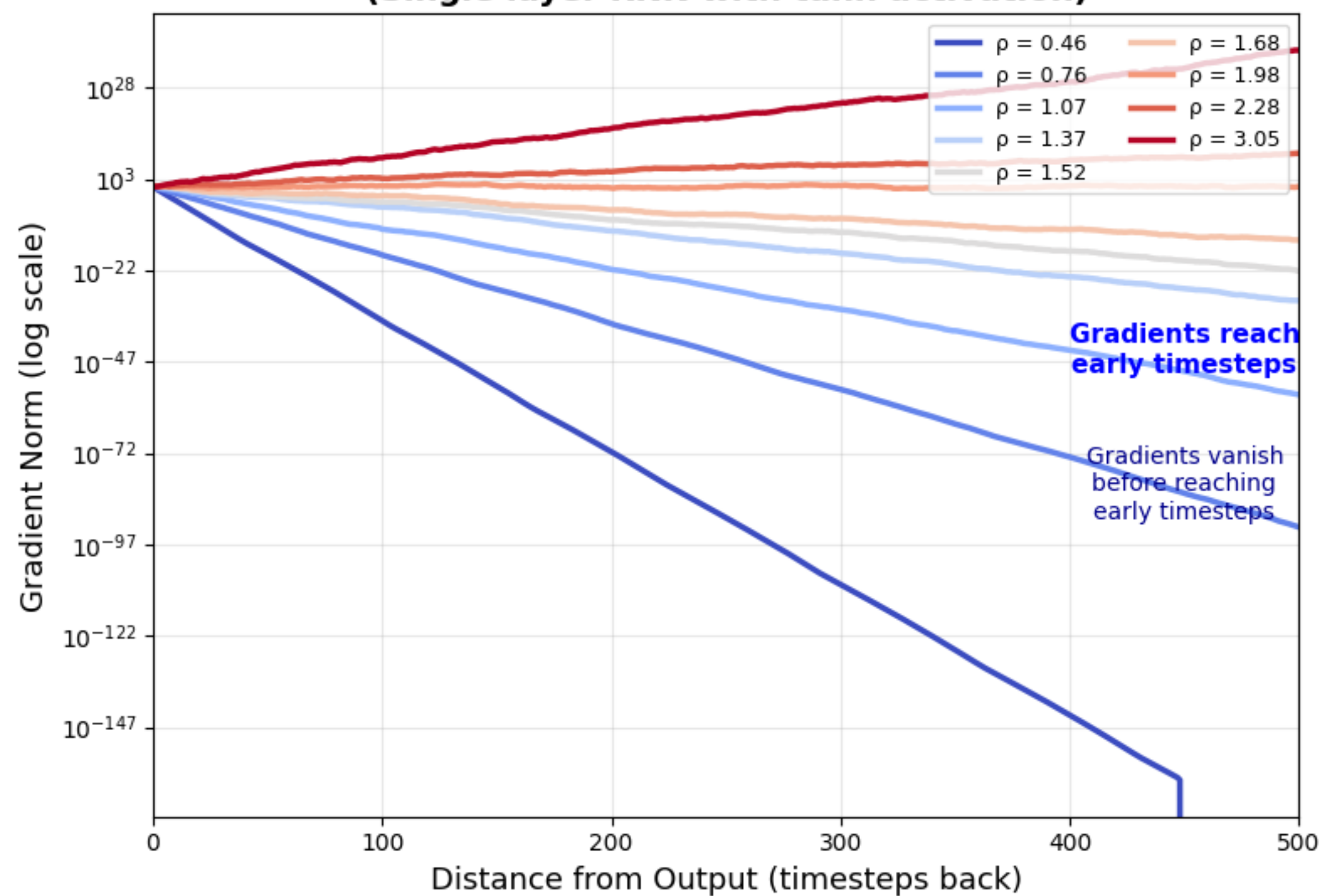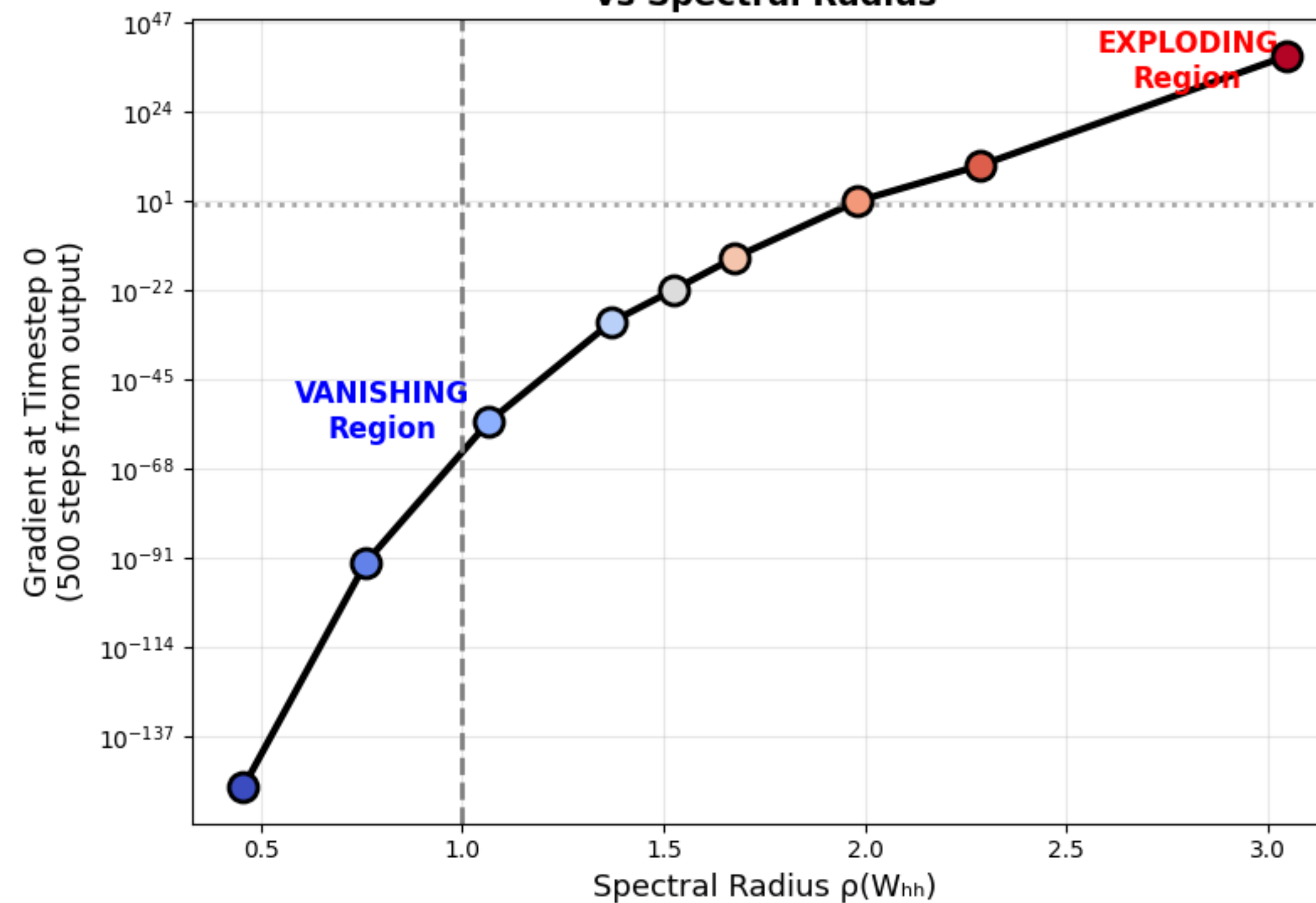
$$h_t = \tanh(W_{xh} \cdot x_t + W_{hh} \cdot h_{t-1} + b_h)$$

$$y_t = W_{hy} \cdot h_t + b_y$$

Input Vector ($x$)    Hidden State ($h$)    Output Vector ($y$)

**Real Tanh RNN: Gradient Flow Through 500 Timesteps**
(Single layer RNN with tanh activation)

Gradient Norm (log scale)

| ρ = 0.46 | ρ = 1.68 |
| ρ = 0.76 | ρ = 1.98 |
| ρ = 1.07 | ρ = 2.28 |
| ρ = 1.37 | ρ = 3.05 |
| ρ = 1.52 | |

**Gradients reach early timesteps**

Gradients vanish before reaching early timesteps

Distance from Output (timesteps back)

**Gradient Reaching First Timestep**
vs Spectral Radius

Gradient at Timestep 0
(500 steps from output)

**EXPLODING Region**

**VANISHING Region**

Spectral Radius ρ(W_hh)

# Training Instability: One Gradient Step Destroys Careful Initialization

## Gradient Flow: Before vs After ONE Gradient Update



Legend:
- Before training (stable)
- After 1 step (lr=1e-03)
- After 1 step (lr=1e-02)
- After 1 step (lr=5e-02)

X-axis: Timesteps from Output

## Gradient Reaching First Timestep After ONE Update



Y-axis: Gradient at t=0 (500 steps back)
X-axis: Learning Rate

EXPLODING

STABLE

# Mitigating Long Context Issues in RNNs

Any suggestions?

# Mitigating Long Context Issues in RNNs

**Any suggestions?**

Adam EMA will be too slow

Any update can lead to collapse/explosion

Constrain the parameters singular values!

# Mitigating Long Context Issues in RNNs

**Efficient Orthogonal Parametrisation of Recurrent Neural Networks Using Householder Reflections**

Zakaria Mhammedi [1 2]   Andrew Hellicar [2]   Ashfaqur Rahman [2]   James Bailey [1]

**Unitary Evolution Recurrent Neural Networks**

Martin Arjovsky [*]                                    MARJOVSKY@DC.UBA.AR
Amar Shah [*]                                          AS793@CAM.AC.UK
Yoshua Bengio
Universidad de Buenos Aires, University of Cambridge,
Université de Montréal. Yoshua Bengio is a CIFAR Senior Fellow.

**projUNN: efficient method for training deep networks with unitary matrices**

Bobak T. Kiani          Randall Balestriero          Yann LeCun
MIT                     Meta AI, FAIR                NYU & Meta AI, FAIR
bkiani@mit.edu          rbalestriero@fb.com          yann@fb.com
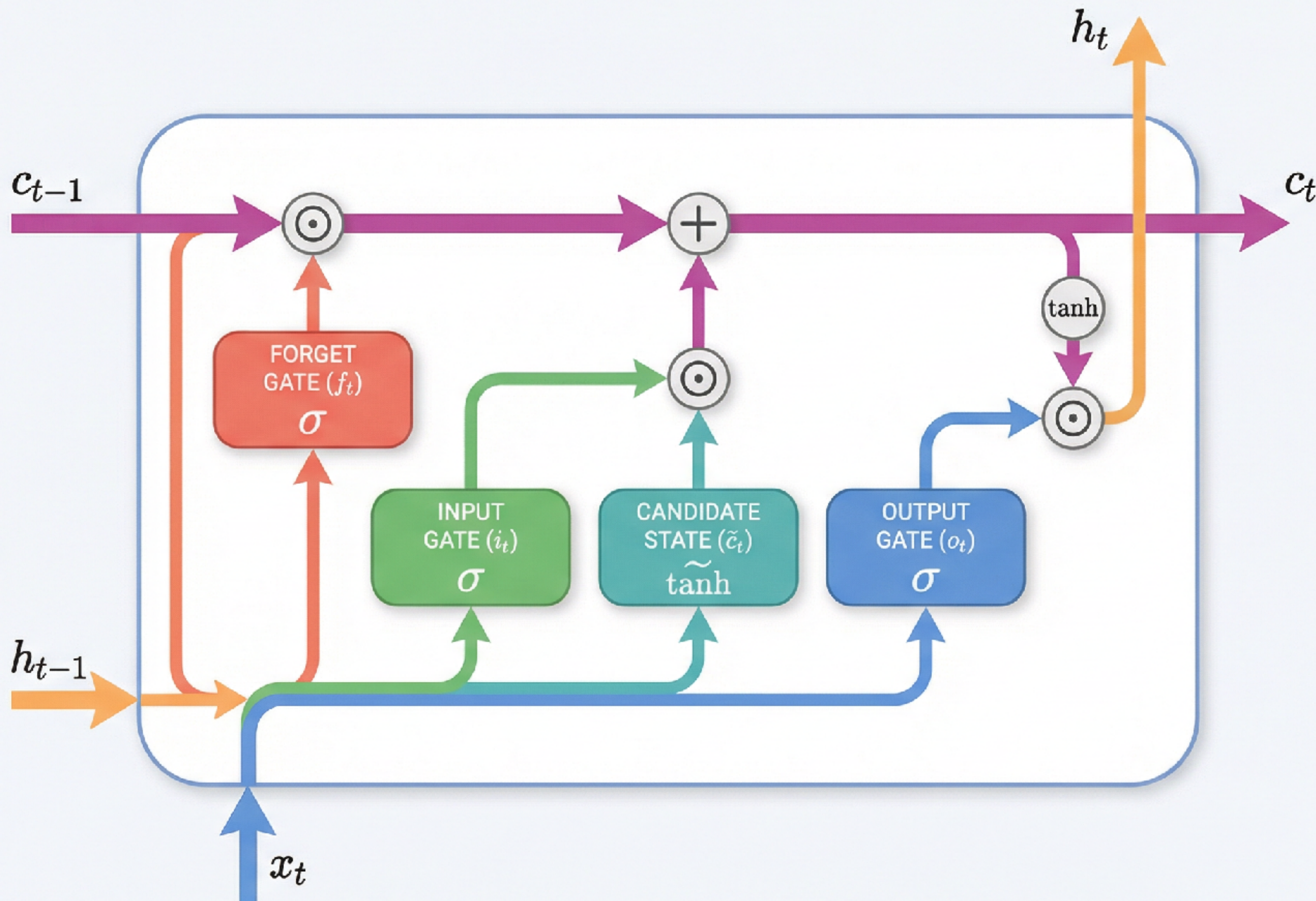
Seth Lloyd
MIT & Turing Inc.
slloyd@mit.edu

# Mitigating Long Context Issues in RNNs

**Any idea on how to enforce orthogonal/unitary W?**

# LSTM Cell



## Equations

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{c}_t = \tanh(W_c \cdot [h_{t-1}, x_t] + b_c)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$$

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t \odot \tanh(c_t)$$

## Legend

FORGET GATE ($f_t$) — Decides what to discard from cell state

INPUT GATE ($i_t$) — Decides what new info to store

CANDIDATE STATE ($\tilde{c}_t$) — Proposes new values

OUTPUT GATE ($o_t$) — Decides what to output

$c_{t-1}/c_t$ (Cell State) — Memory Highway

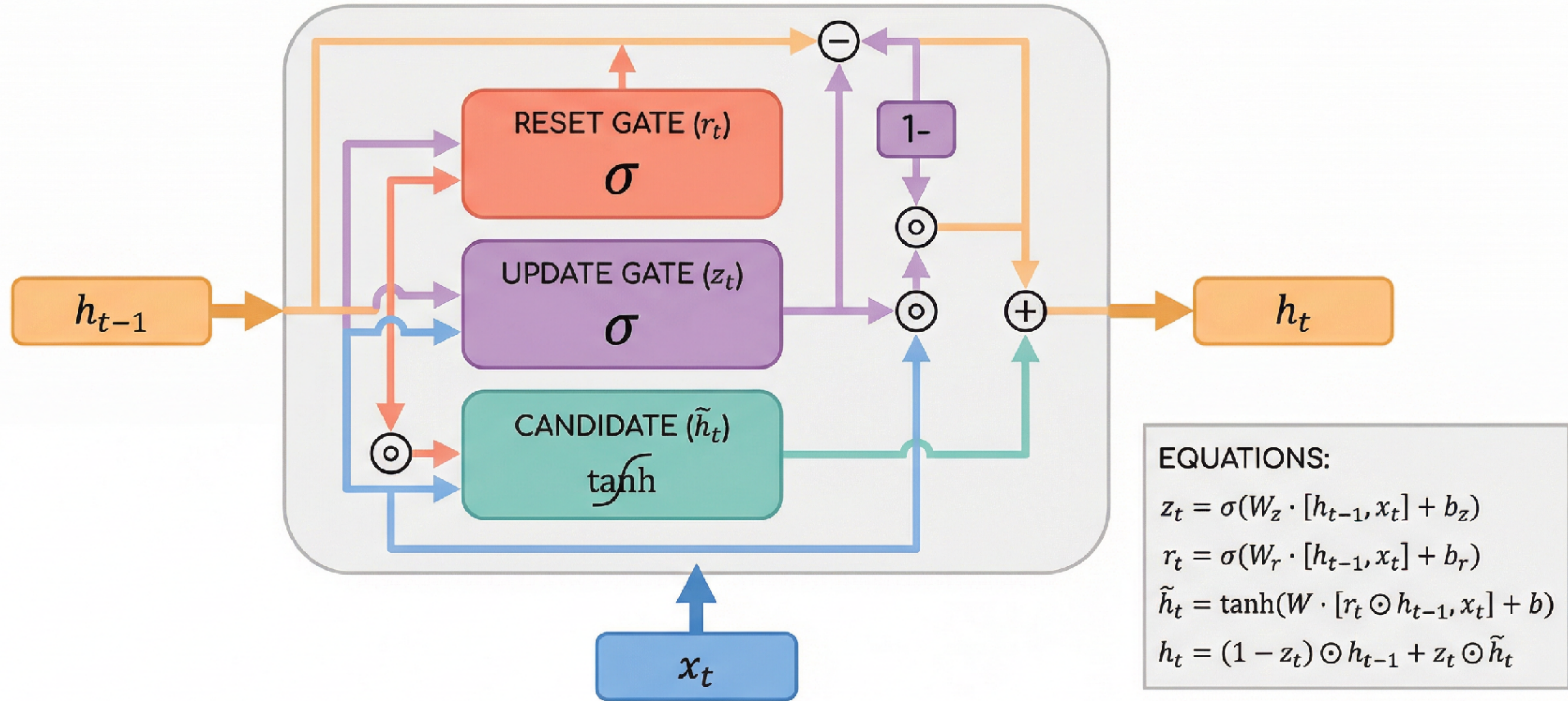$\sigma$ — Sigmoid

$\tanh$ — Hyperbolic Tangent

$\odot$ — Element-wise Multiplication

$+$ — Addition

# GRU Cell



RESET GATE ($r_t$)
$\sigma$

UPDATE GATE ($z_t$)
$\sigma$

CANDIDATE ($\tilde{h}_t$)
tanh

1-

$h_{t-1}$

$x_t$

$h_t$

EQUATIONS:

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

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

$$\tilde{h}_t = \tanh(W \cdot [r_t \odot h_{t-1}, x_t] + b)$$

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

LEGEND:
- Reset gate: controls how much past to use for candidate - Update gate: controls interpolation between old and new
"GRU: 2 gates, no separate cell state"

# Try it!

## RNN #

```
class torch.nn.RNN(input_size, hidden_size, num_layers=1, nonlinearity='tanh',
bias=True, batch_first=False, dropout=0.0, bidirectional=False, device=None,
dtype=None)
```
[source]

## LSTM

```
class torch.nn.LSTM(input_size, hidden_size, num_layers=1, bias=True, batch_first=False,
dropout=0.0, bidirectional=False, proj_size=0, device=None, dtype=None) #
```
[source]

## GRU

```
class torch.nn.GRU(input_size, hidden_size, num_layers=1, bias=True, batch_first=False,
dropout=0.0, bidirectional=False, device=None, dtype=None)
```
[source]

# Questions?