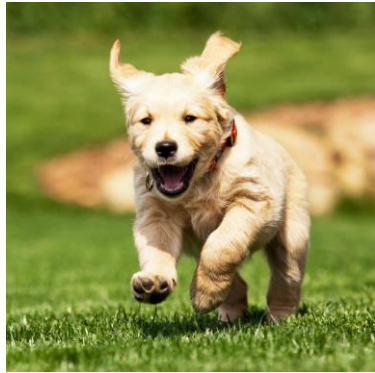# Recurrent Neural Networks

CS 6384 Computer Vision

Professor Yu Xiang

The University of Texas at Dallas
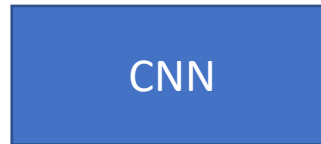
Some slides of this lecture are courtesy Stanford CS231n

# Single Images

- Convolutional neural networks



Image → CNN → High-level information
- Depth
- Object classes
- Object poses
- Etc.
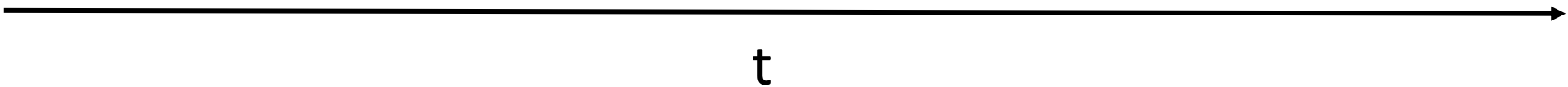
# Sequential Data

- Data depends on time
  - Video



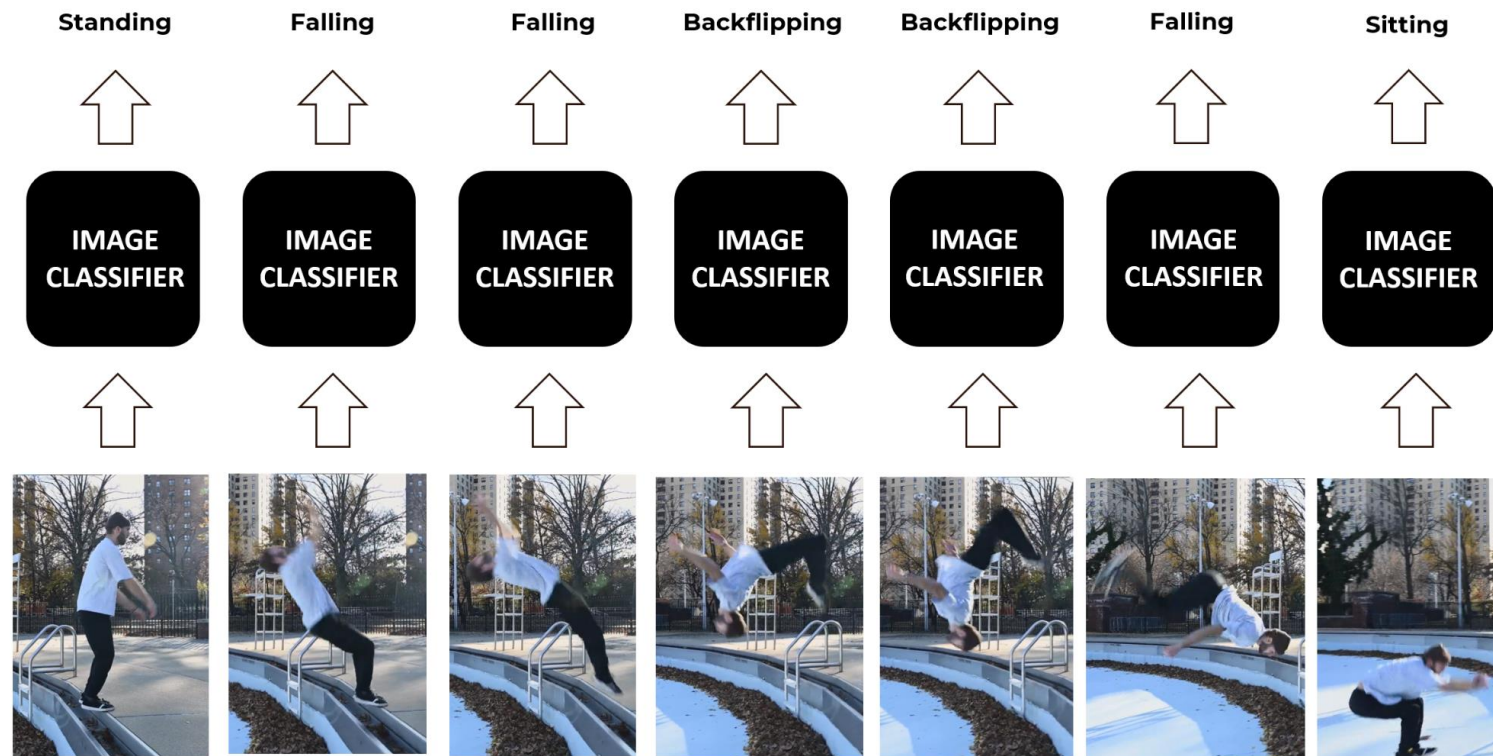t-1                t                t+1

  - Sentence

UT Dallas is a rising public research university in the heart of DFW.
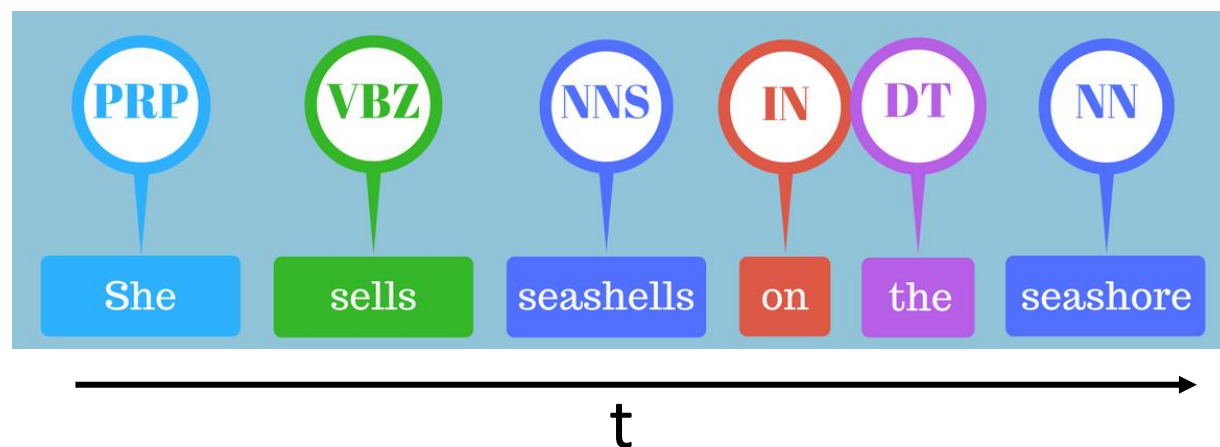
t

# Sequential Data Labeling

- Video frame labeling



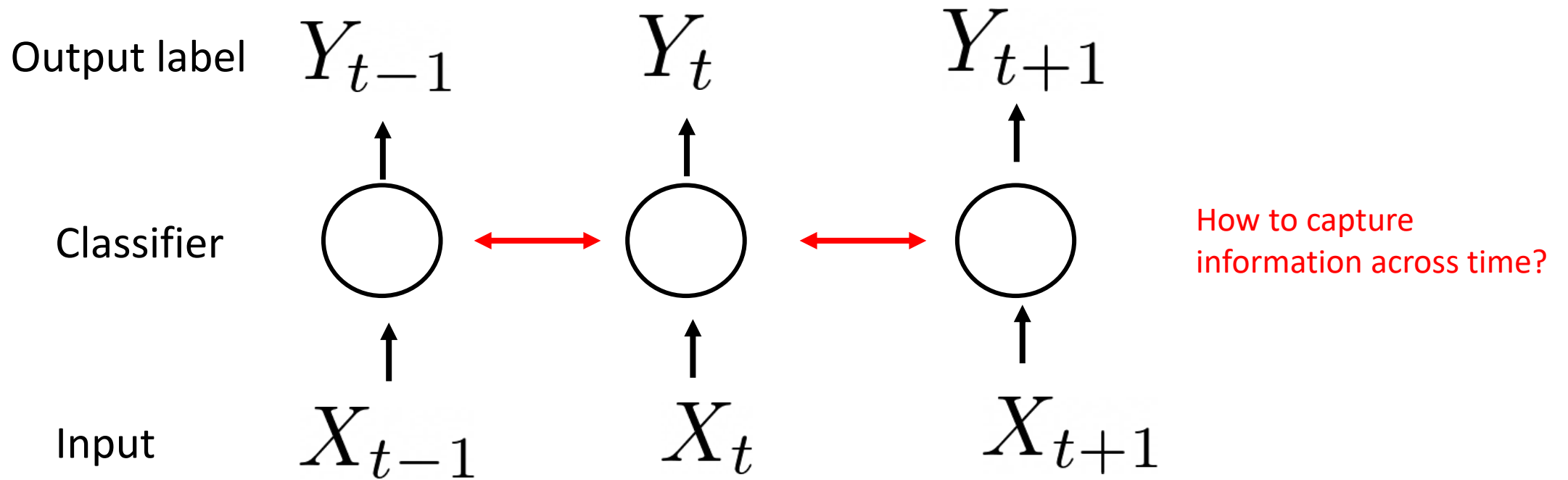https://bleedai.com/human-activity-recognition-using-tensorflow-cnn-lstm/

# Sequential Data Labeling

- Part-of-speech tagging (grammatical tagging)



| Tag | Meaning | English Examples |
|------|---------|------------------|
| ADJ | adjective | *new, good, high, special, big, local* |
| ADP | adposition | *on, of, at, with, by, into, under* |
| ADV | adverb | *really, already, still, early, now* |
| CONJ | conjunction | *and, or, but, if, while, although* |
| DET | determiner, article | *the, a, some, most, every, no, which* |
| NOUN | noun | *year, home, costs, time, Africa* |
| NUM | numeral | *twenty-four, fourth, 1991, 14:24* |
| PRT | particle | *at, on, out, over per, that, up, with* |
| PRON | pronoun | *he, their, her, its, my, I, us* |
| VERB | verb | *is, say, told, given, playing, would* |
| . | punctuation marks | *. , ; !* |
| X | other | *ersatz, esprit, dunno, gr8, univeristy* |

# Sequential Data Labeling

Output label $Y_{t-1}$ $Y_t$ $Y_{t+1}$

Classifier

Input $X_{t-1}$ $X_t$ $X_{t+1}$

How to capture information across time?

# Recurrent Neural Networks

Output label

$$Y_{t-1} \qquad Y_t \qquad Y_{t+1}$$

RNN $\qquad$ RNN $\qquad$ RNN

Internal state
(memory)

$$\mathbf{h}$$

Input

$$X_{t-1} \qquad X_t \qquad X_{t+1}$$

# Hidden State Update



$$\mathbf{h}_t = f_W(\mathbf{h}_{t-1}, \mathbf{x}_t)$$

Updating function with parameters W

Hidden state at time t

Hidden state at time t-1

Input at time t

# Using the Hidden State

RNN $\mathbf{h}$

$Y$

$X$
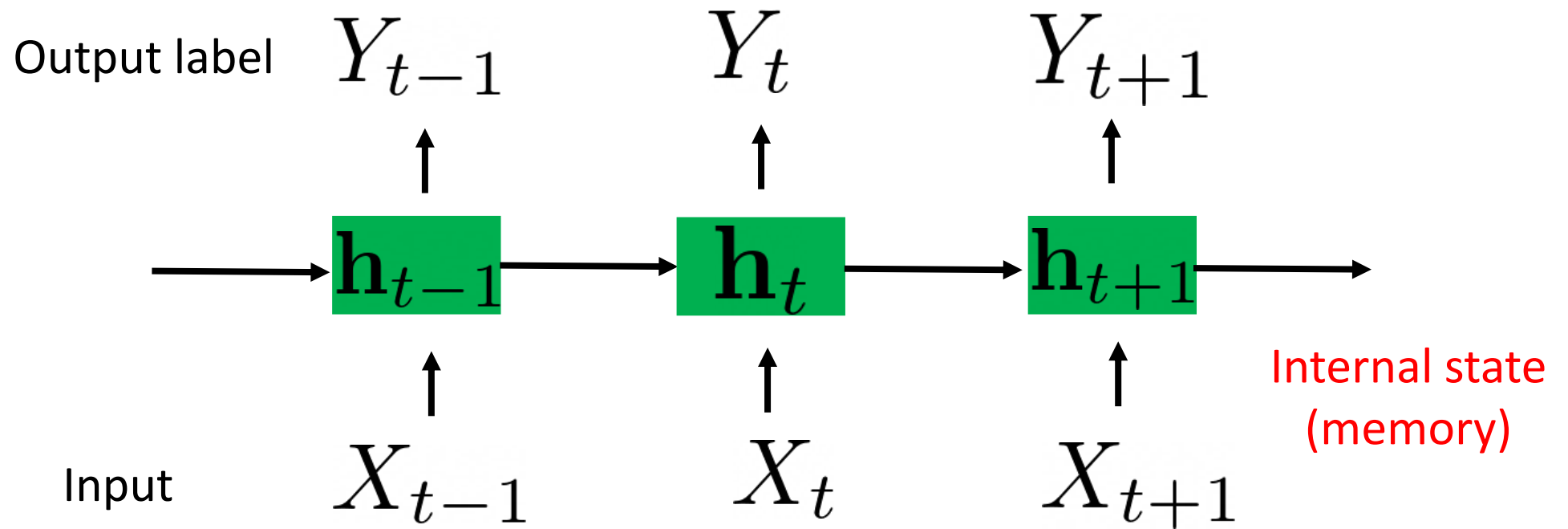
$$\mathbf{h}_t = f_W(\mathbf{h}_{t-1}, \mathbf{x}_t)$$

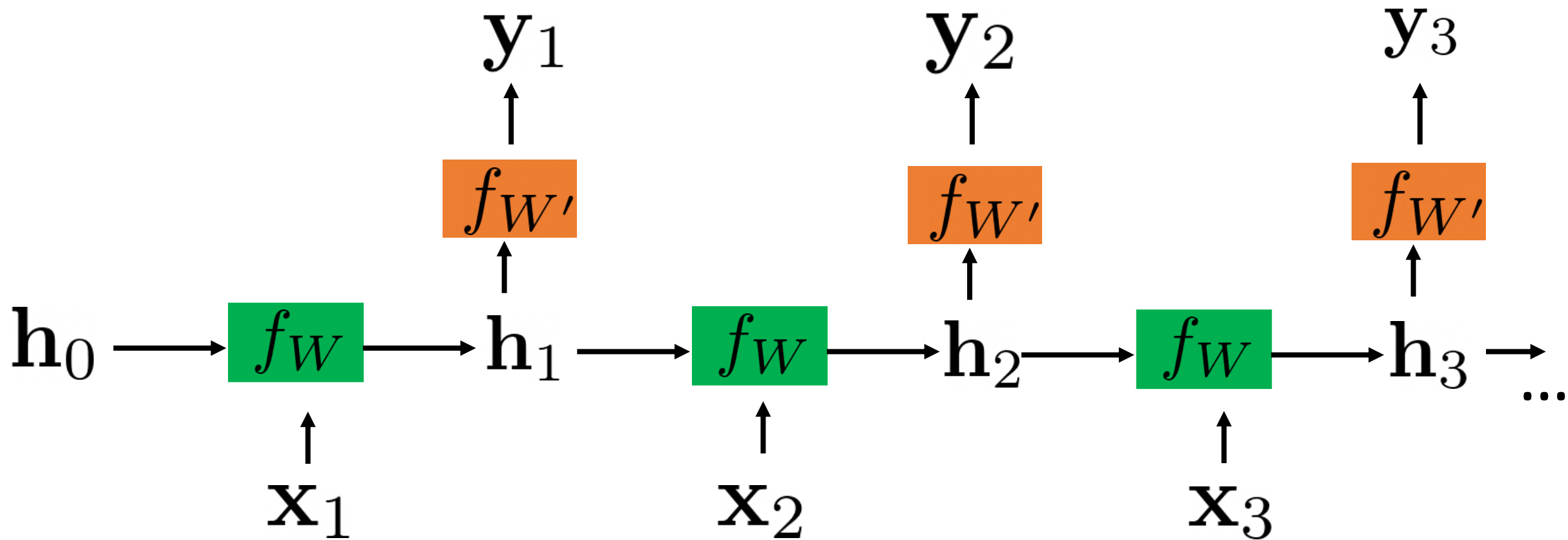$$\mathbf{y}_t = f_{W'}(\mathbf{h}_t)$$

# Recurrent Neural Networks

Output label $\quad Y_{t-1} \qquad\qquad Y_t \qquad\qquad Y_{t+1}$



Internal state (memory)

Input $\qquad X_{t-1} \qquad\qquad X_t \qquad\qquad X_{t+1}$

# Vanilla RNN

### Hidden state updating rule

$$\mathbf{h}_t = \tanh(W_{hh}\mathbf{h}_{t-1} + W_{xh}\mathbf{x}_t)$$

RNN

$$m \times 1 \qquad m \times m \quad m \times 1 \quad m \times n \quad n \times 1$$

**tanh**   tanh(x)

$$\frac{e^{2x}-1}{e^{2x}+1}$$

$$\mathbf{y}_t = W_{hy}\mathbf{h}_t$$
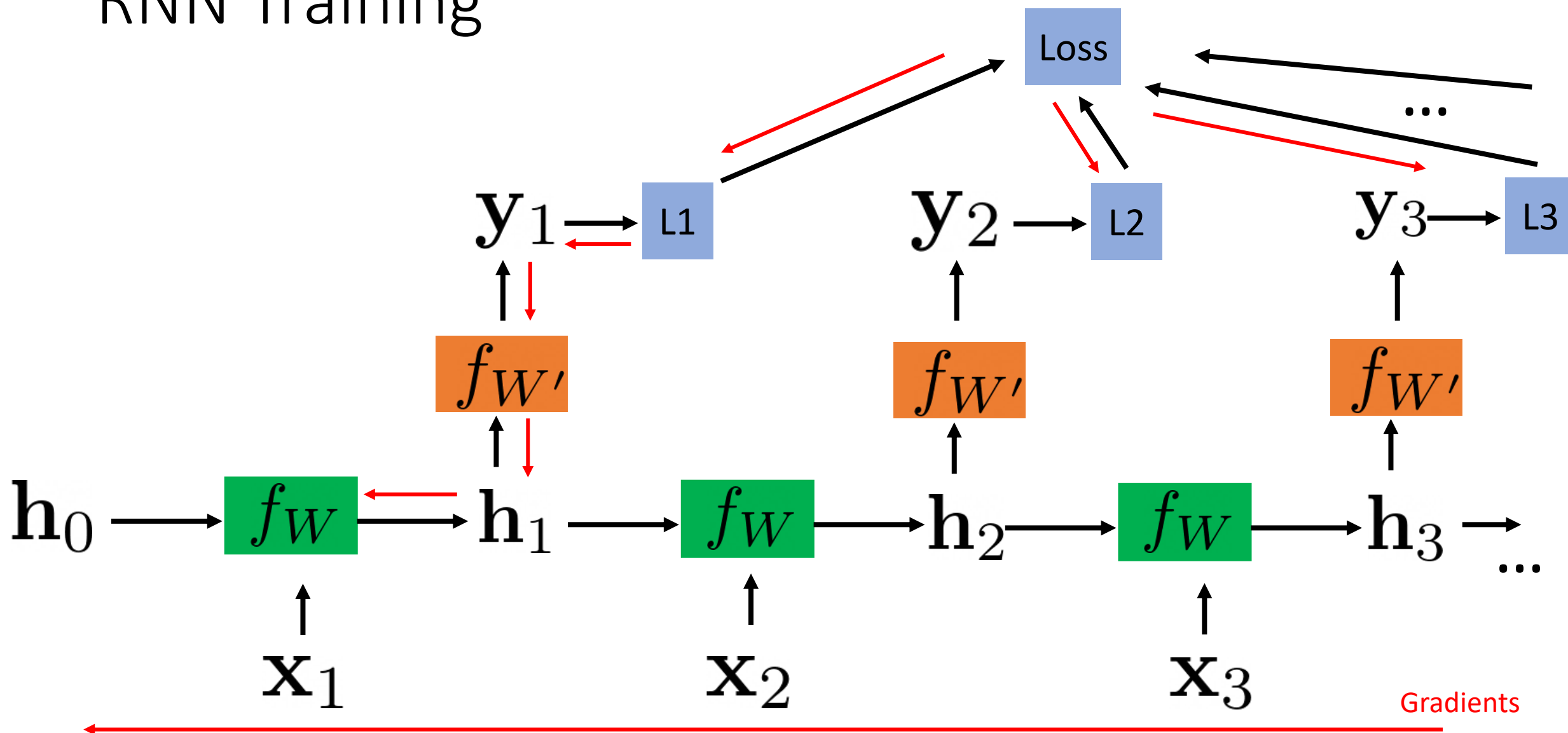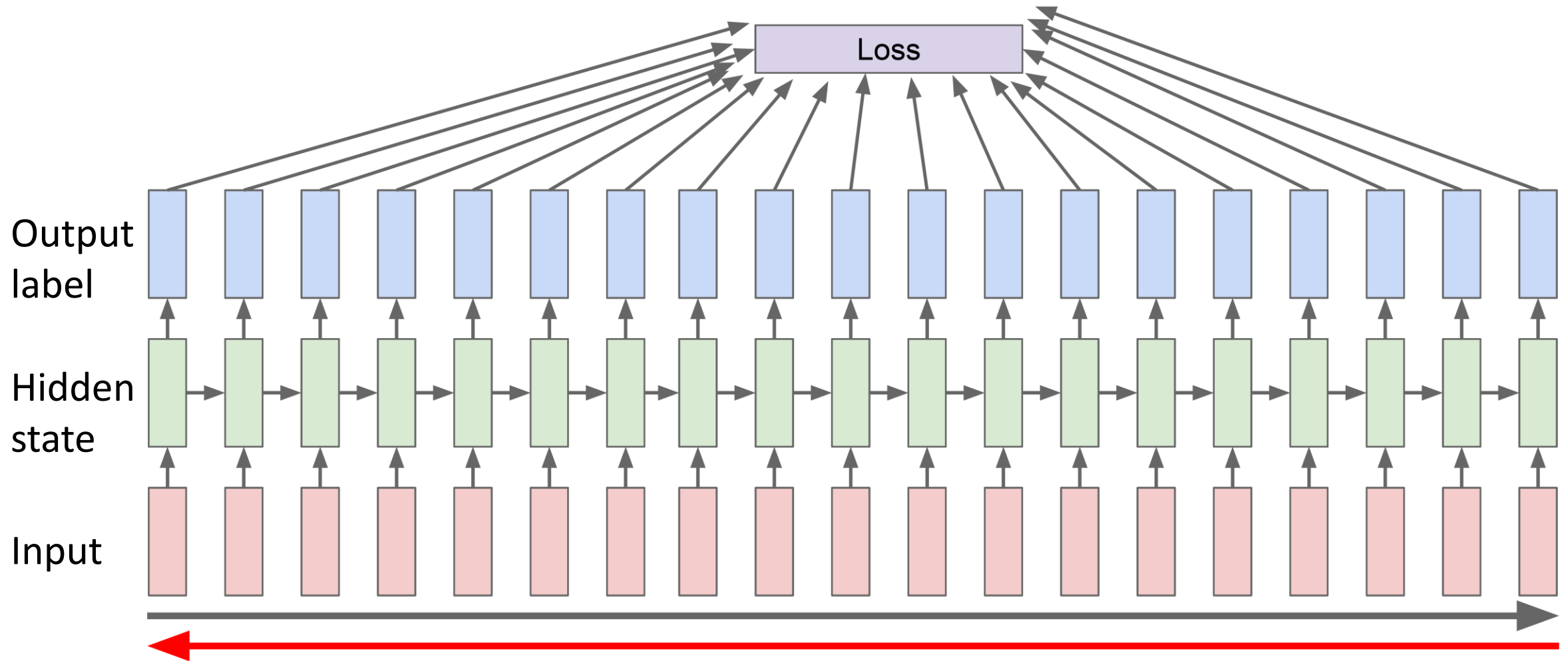
$$l \times 1 \qquad m \times 1$$

# RNN Computation Graph



The same set of weights for different time steps $f_W$ $f_{W'}$

# RNN Training

# Backpropagation through Time



Output label

Hidden state
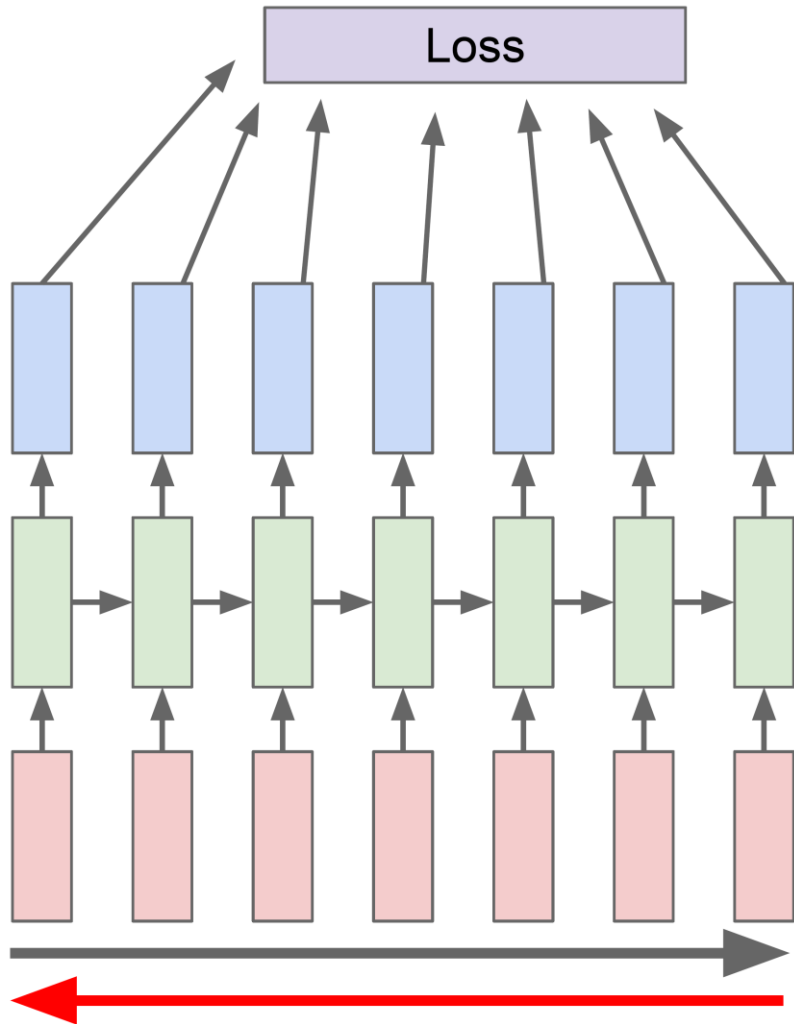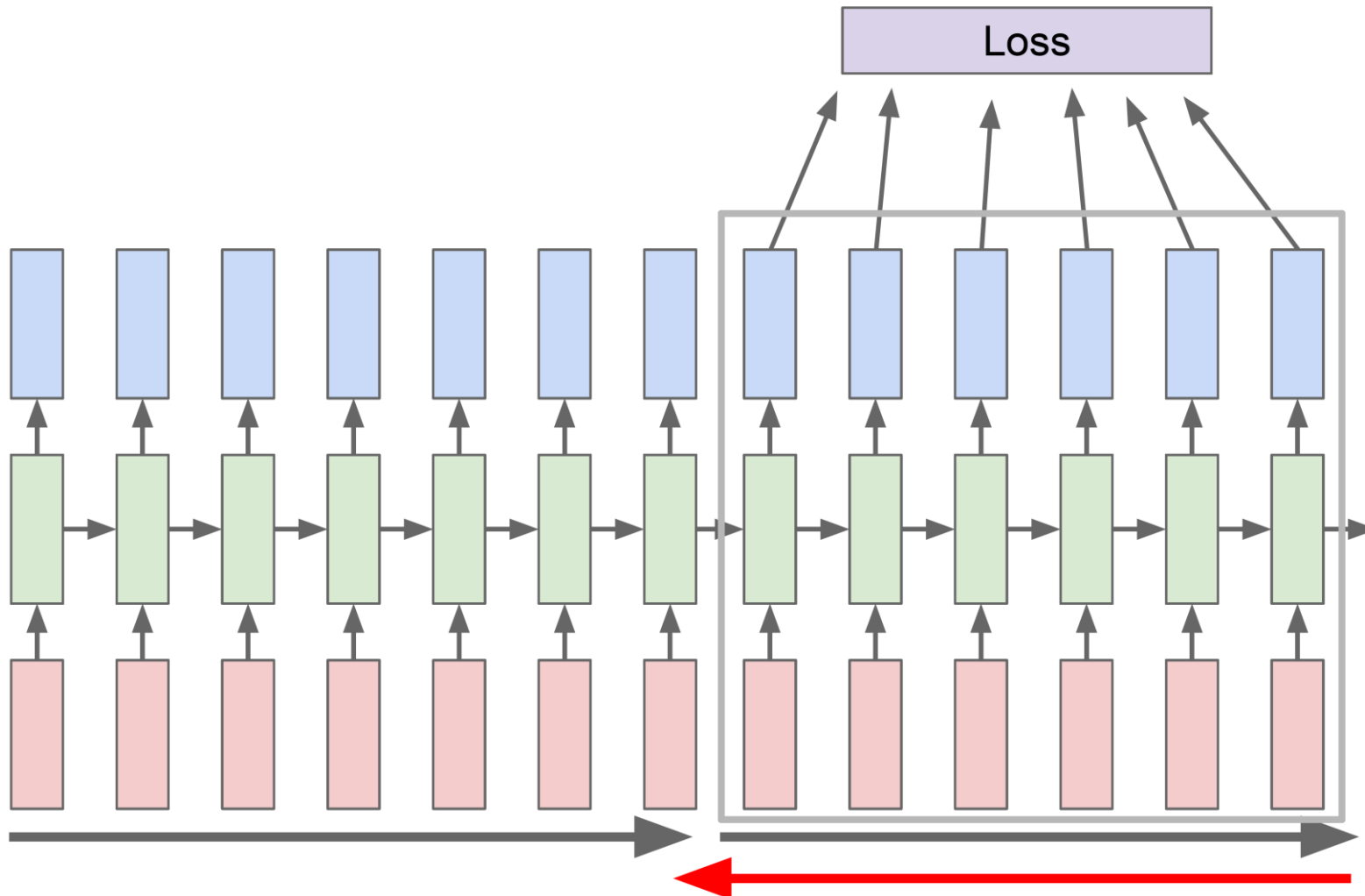
Input

What is the problem in this training paradigm?
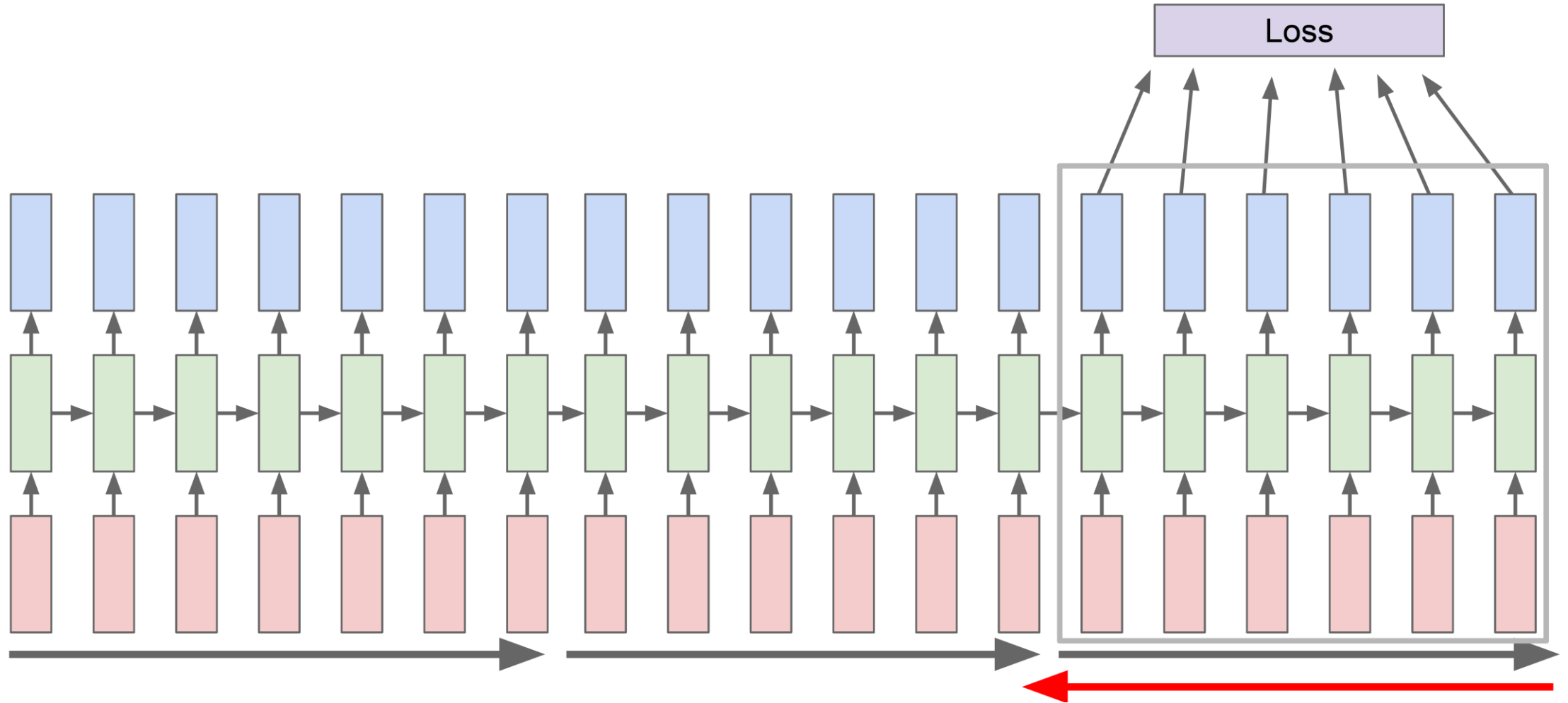
# Truncated Backpropagation through Time



Run forward and backward through chunks of the sequence instead of whole sequence
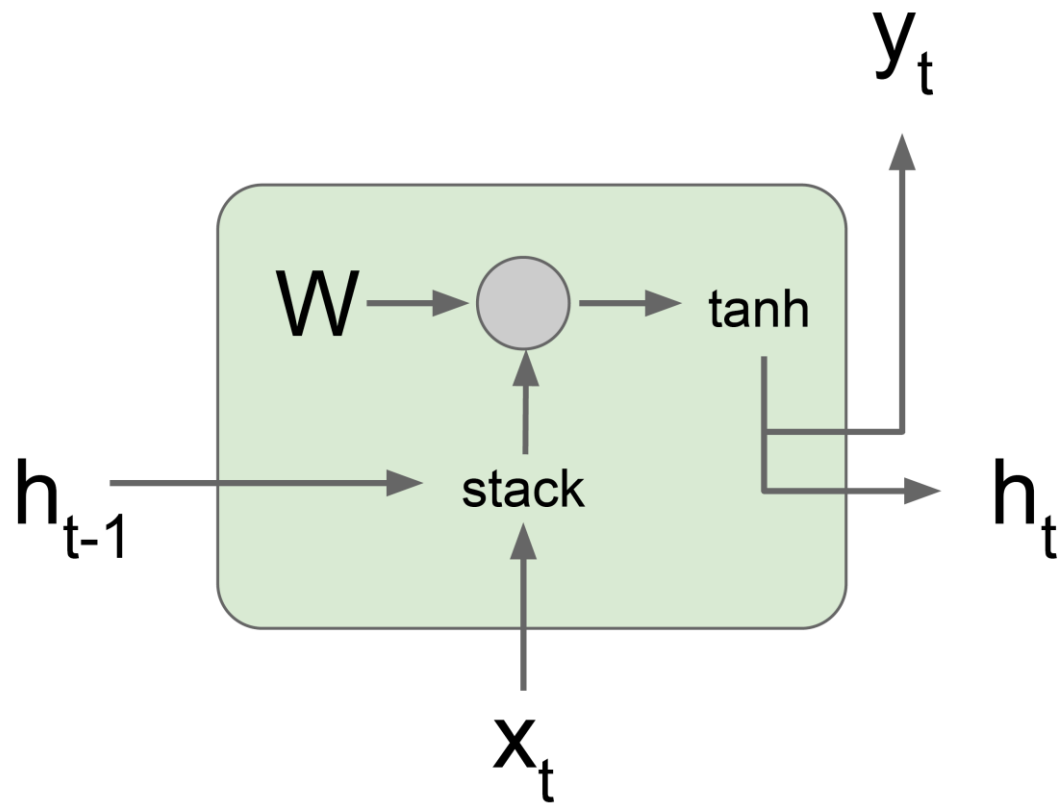
# Truncated Backpropagation through Time



Carry hidden states forward in time forever, but only backpropagate for some smaller number of steps

# Truncated Backpropagation through Time

# Vanilla RNN Gradient Flow



$$\mathbf{h}_t = \tanh(W_{hh}\mathbf{h}_{t-1} + W_{xh}\mathbf{x}_t)$$

$$= \tanh\left(\begin{pmatrix} W_{hh} & W_{hx} \end{pmatrix}\begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right)$$

$$= \tanh\left(W\begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right)$$

Yu Xiang

# Vanilla RNN Gradient Flow

Backpropagation from $h_t$ to $h_{t-1}$ multiplies by W (actually $W_{hh}^T$)



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

$$= \tanh\left(\begin{pmatrix} W_{hh} & W_{hx} \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right)$$

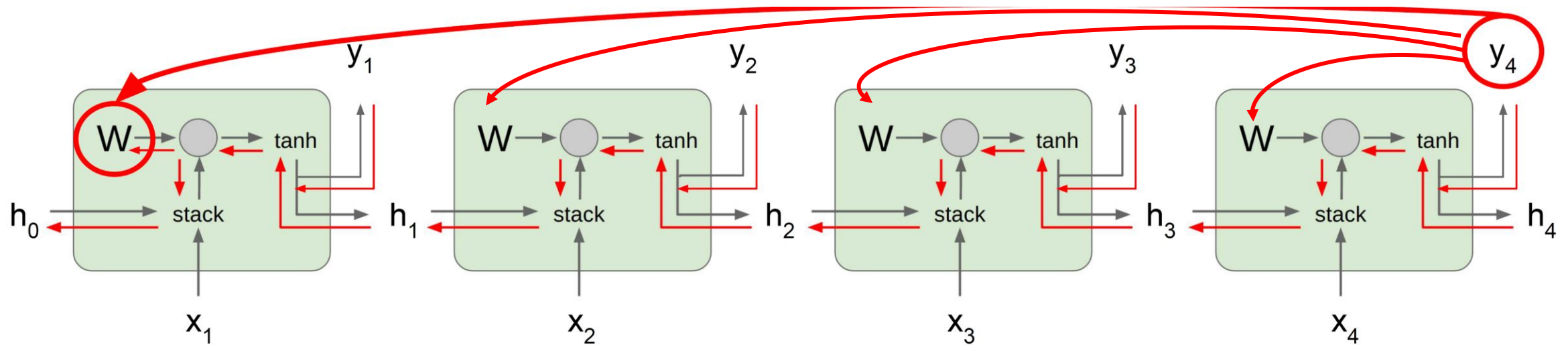$$= \tanh\left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right)$$

$$\frac{\partial h_t}{\partial h_{t-1}} = tanh'(W_{hh} h_{t-1} + W_{xh} x_t) W_{hh}$$
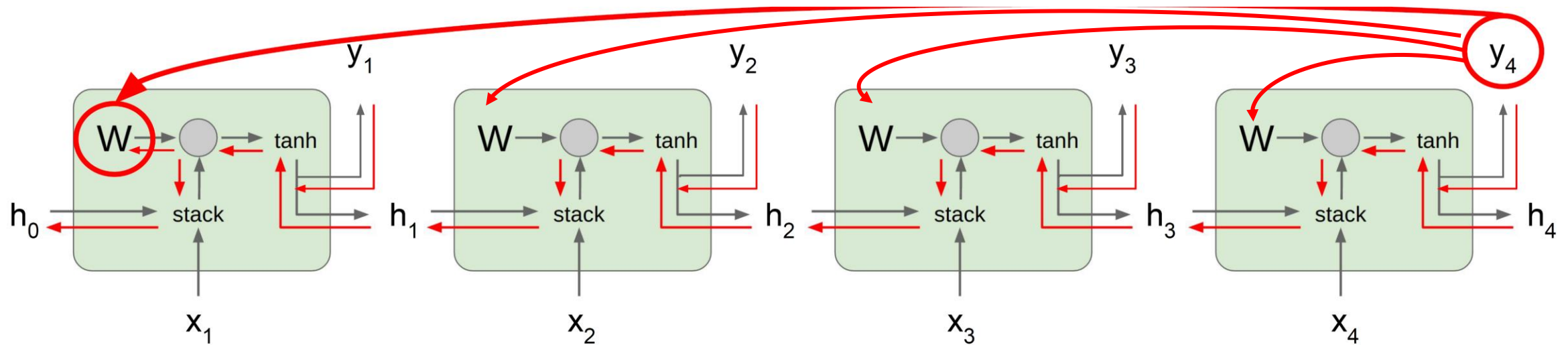
# Vanilla RNN Gradient Flow



$$\frac{\partial L}{\partial W} = \sum_{t=1}^{T} \frac{\partial L_t}{\partial W}$$

$$\frac{\partial L_T}{\partial W} = \frac{\partial L_T}{\partial h_T} \frac{\partial h_t}{\partial h_{t-1}} \cdots \frac{\partial h_1}{\partial W} = \frac{\partial L_T}{\partial h_T} \left( \prod_{t=2}^{T} \frac{\partial h_t}{\partial h_{t-1}} \right) \frac{\partial h_1}{\partial W}$$

# Vanilla RNN Gradient Flow



$$\frac{\partial L_T}{\partial W} = \frac{\partial L_T}{\partial h_T}\left(\prod_{t=2}^{T} \frac{\partial h_t}{\partial h_{t-1}}\right)\frac{\partial h_1}{\partial W}$$

- Vanishing gradients $\qquad \|\frac{\partial h_t}{\partial h_{t-1}}\|_2 < 1$

- Exploding gradients $\qquad \|\frac{\partial h_t}{\partial h_{t-1}}\|_2 > 1$

https://en.wikipedia.org/wiki/Matrix_norm

# Vanilla RNN Gradient Flow

- Exploding gradients  $\|\dfrac{\partial h_t}{\partial h_{t-1}}\|_2 > 1$

  - Gradient clipping

  ```
  grad_norm = np.sum(grad * grad)
  if grad_norm > threshold:
      grad *= (threshold / grad_norm)
  ```

- Vanishing gradients  $\|\dfrac{\partial h_t}{\partial h_{t-1}}\|_2 < 1$

  - Change RNN architecture
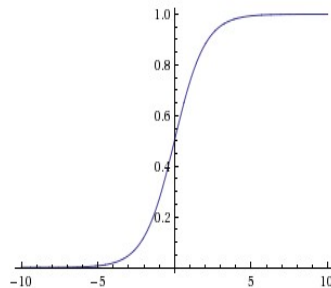
# Long Short Term Memory (LSTM)

Vanilla RNN

$$h_t = \tanh\left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right)$$

LSTM

Input gate

forget gate

output gate

gate gate

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

Cell $\quad c_t = f \odot c_{t-1} + i \odot g$

Hidden state $\quad h_t = o \odot \tanh(c_t)$

**Sigmoid**

$$\sigma(x) = 1/(1 + e^{-x})$$

Store Cell and hidden states

# Long Short Term Memory (LSTM)



$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$
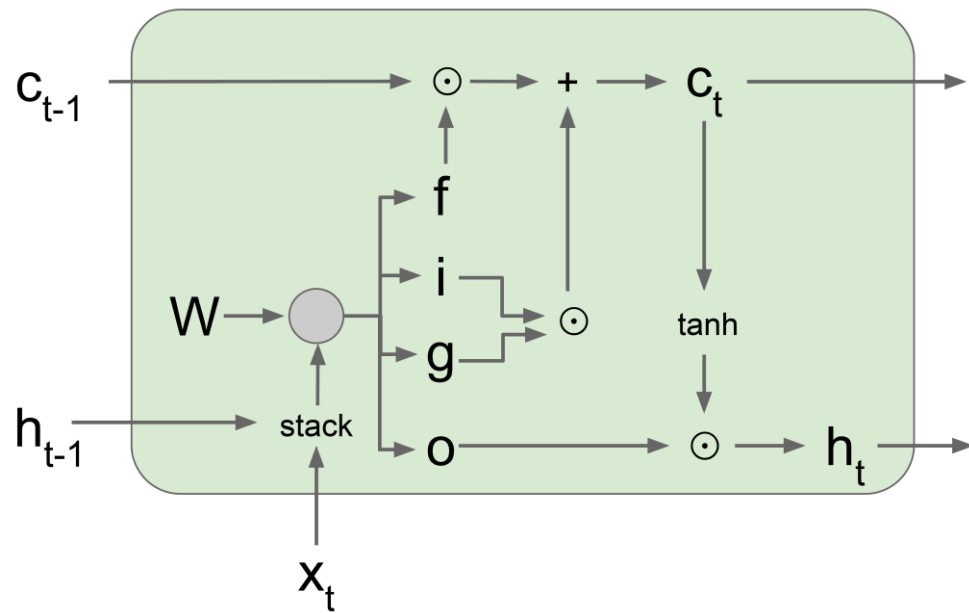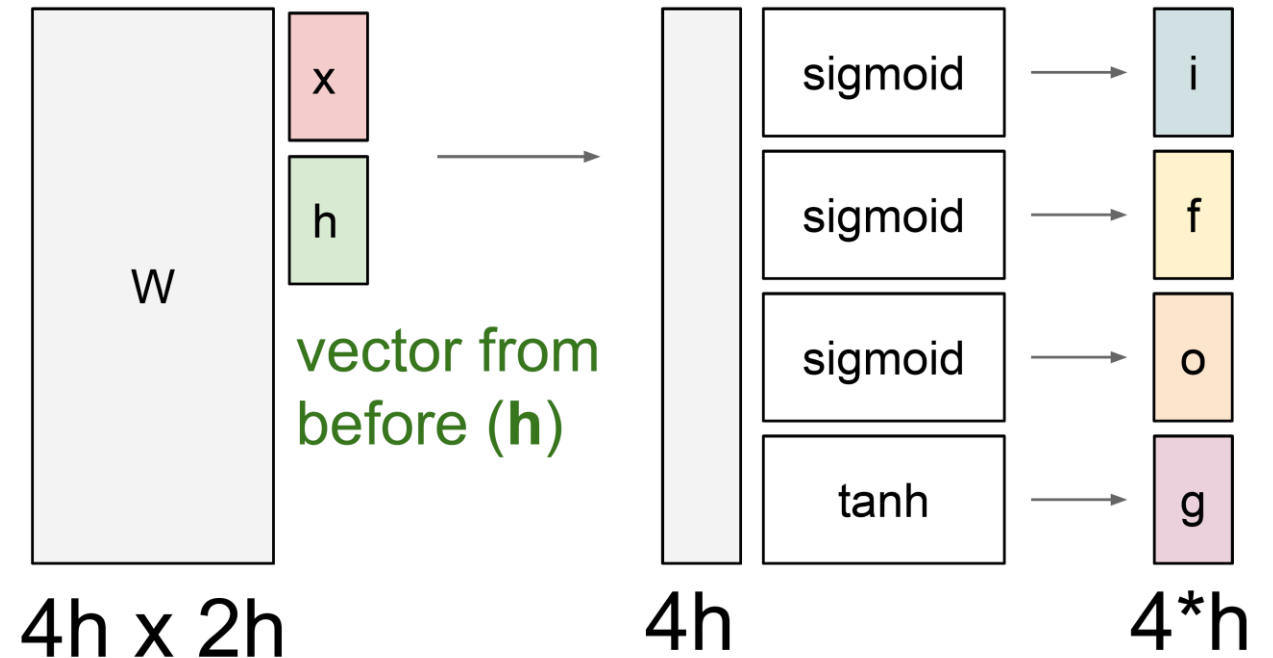
$$c_t = f \odot c_{t-1} + i \odot g$$

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

vector from before (**h**)

4h x 2h

4h

4*h

- **g**: Gate gate, how much to write to cell
- **i**: Input gate, whether to write to cell
- **f**: Forget gate, whether to erase cell
- **o**: Output gate, how much to reveal cell
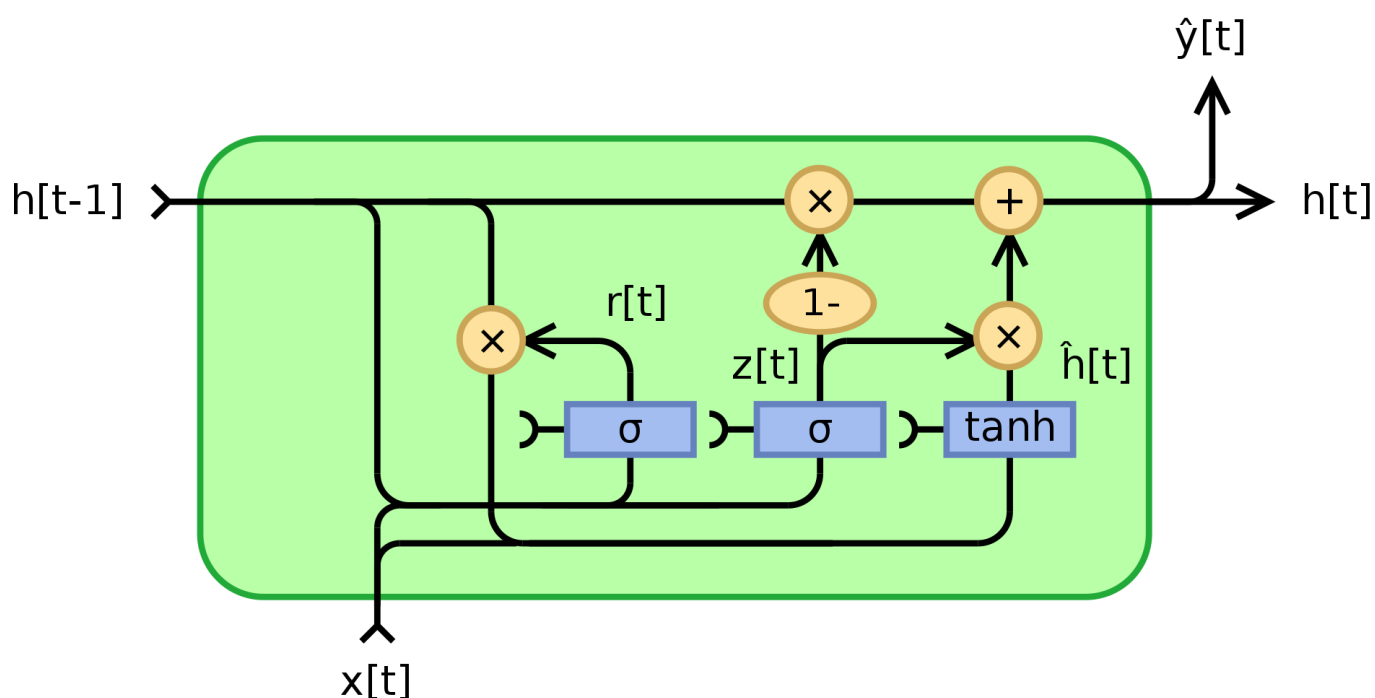
# Long Short Term Memory (LSTM)

- Make the RNN easier to preserve information over many steps
  - E.g., f = 1 and i = 0

  - This is difficult for vanilla RNN

- LSTM does not guarantee that there is no vanishing or exploding gradient

- It provides an easier way to learn long-distance dependencies

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

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

# Gated Recurrent Unit (GRU)



$$z_t = \sigma_g(W_z x_t + U_z h_{t-1} + b_z)$$

$$r_t = \sigma_g(W_r x_t + U_r h_{t-1} + b_r)$$

$$\hat{h}_t = \phi_h(W_h x_t + U_h(r_t \odot h_{t-1}) + b_h)$$

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

- $x_t$: input vector
- $h_t$: output vector
- $\hat{h}_t$: candidate activation vector
- $z_t$: update gate vector
- $r_t$: reset gate vector
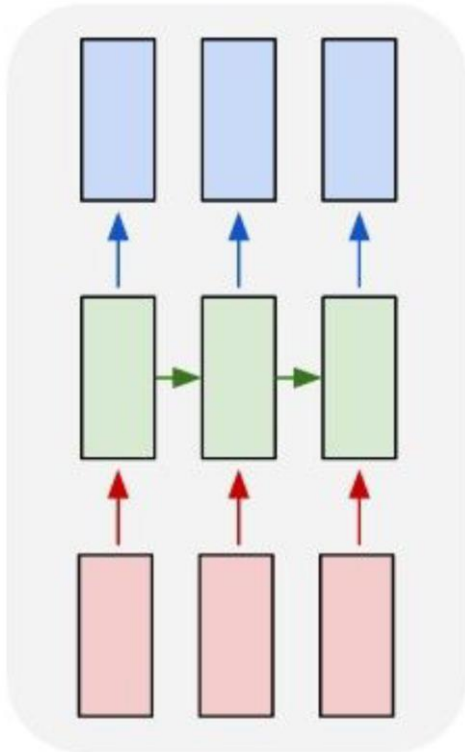- $W$, $U$ and $b$: parameter matrices and vector

https://en.wikipedia.org/wiki/Gated_recurrent_unit

# GRUs vs. LSTMs

- Both have a forget gate

- GRU has fewer parameters, no output gate

- GRUs have similar performance compared to LSTMs, have shown better performance on certain datasets
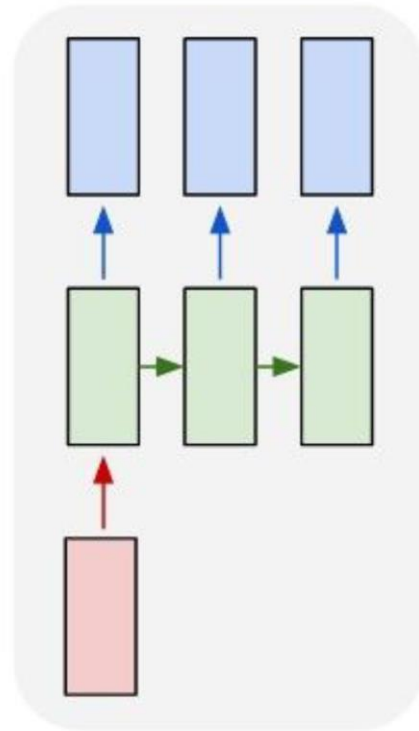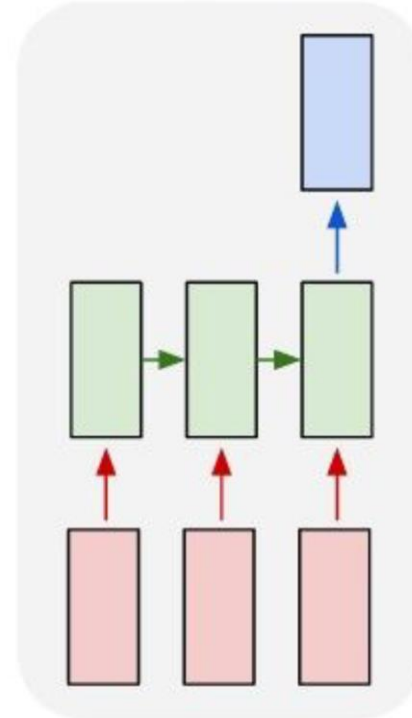
# Recurrent Neural Networks



**many to many**

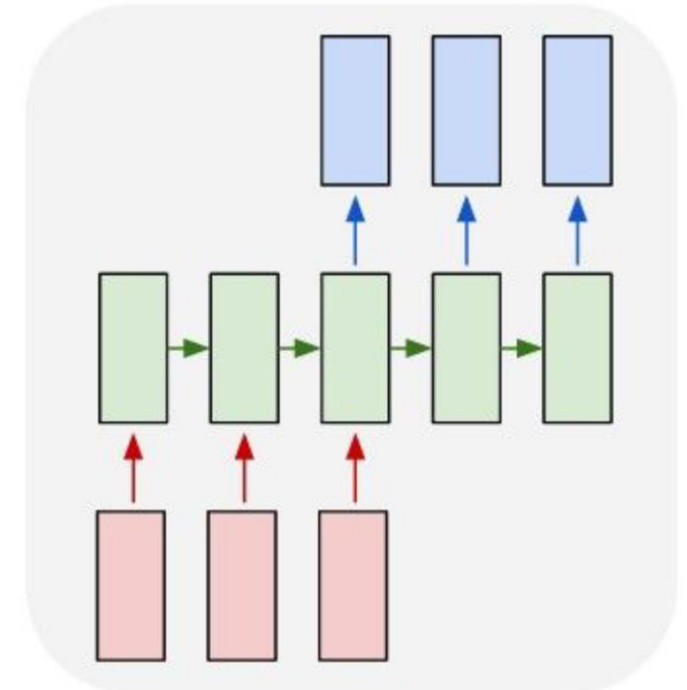E.g., action recognition on video frames

**one to many**

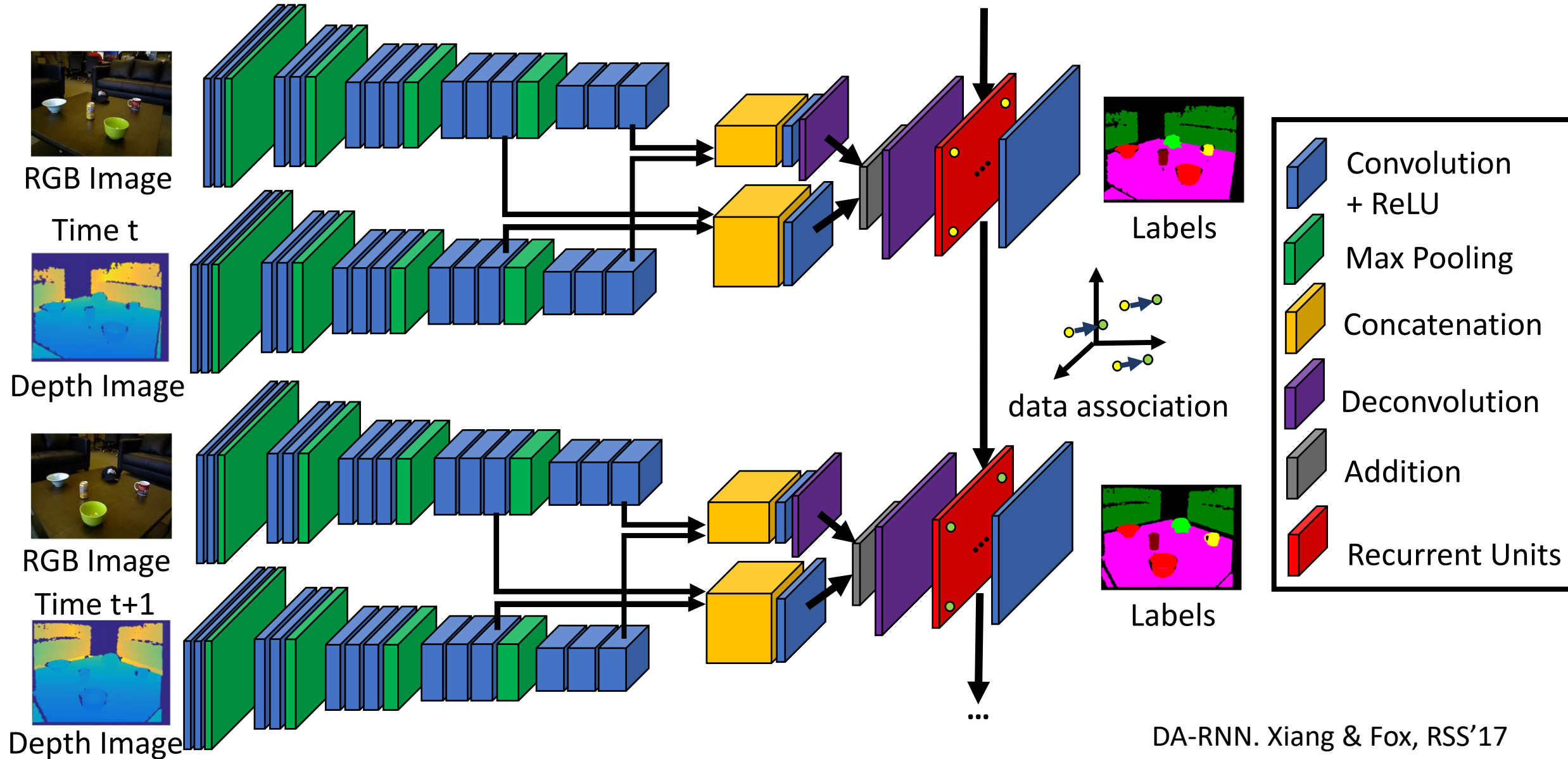E.g., image captioning, image -> sequences of words

**many to one**

E.g., action prediction, sequences of frames -> action class

**many to many**

E.g., Video Captioning Sequence of video frames -> caption

# Recurrent Units on CNN Features



RGB Image
Time t

Depth Image

RGB Image
Time t+1

Depth Image

data association

Labels

Labels

Convolution + ReLU

Max Pooling

Concatenation

Deconvolution

Addition

Recurrent Units

DA-RNN. Xiang & Fox, RSS'17

# Summary

- RNNs can be used for sequential data to capture dependencies in time

- LSTMs and GRUs are better then vanilla RNNs

- It is difficult to capture long-term dependencies in RNNs

- Use transformers  (next lecture)

# Further Reading

- Stanford CS231n, lecture 10, Recurrent Neural Networks
  http://cs231n.stanford.edu/

- Long Short Term Memory
  https://www.researchgate.net/publication/13853244_Long_Short-term_Memory

- Gated Recurrent Units https://arxiv.org/pdf/1412.3555.pdf