



# Introduction to Computer Vision

## Lecture 9 - Temporal Analysis I

Prof. He Wang

# Logistics

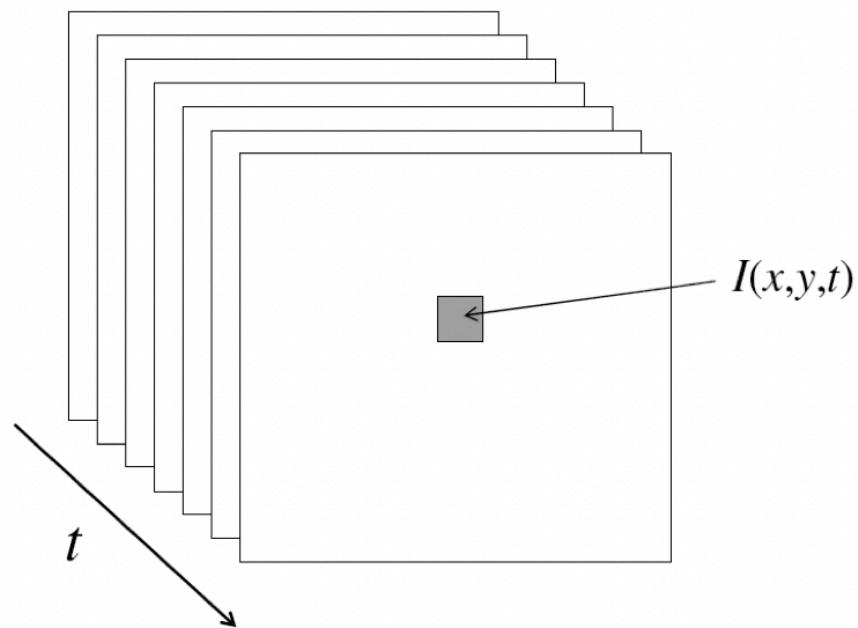
- Assignment 3 will release later this week.

# Motion and Optical Flow

Some slides are borrowed from Stanford CS131.

# From Single Image to Video

- A video is a sequence of frames captured over time
- Now our image data is a function of space ( $x, y$ ) and time ( $t$ )



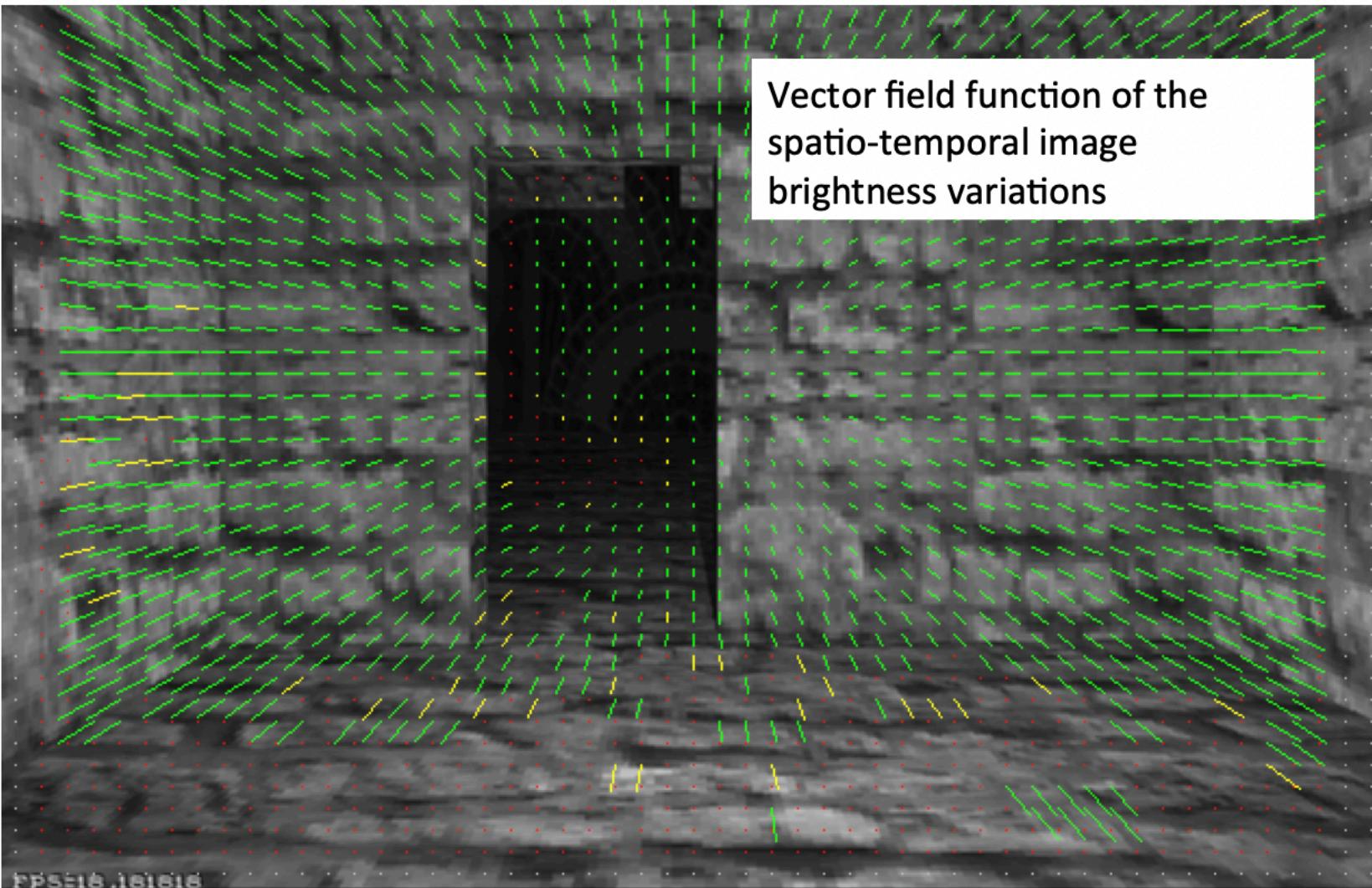
**Today let's focus on motions between two consecutive frames!**

# Optical Flow

- Definition: optical flow is the *apparent* motion of brightness patterns in the image
- Note: apparent motion can be caused by lighting changes without any actual motion
  - Think of a uniform rotating sphere under fixed lighting vs. a stationary sphere under moving illumination

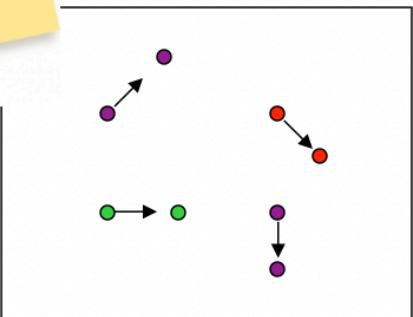
**GOAL:** Recover image motion at each pixel from optical flow

# Optical Flow

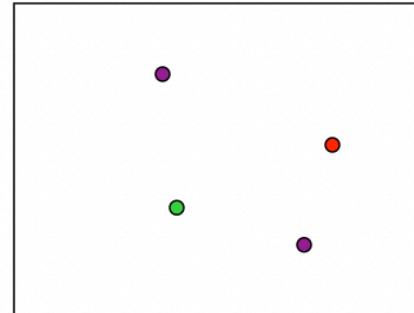


Picture courtesy of Selim Temizer - Learning and Intelligent Systems (LIS) Group, MIT

# Estimating Optical Flow



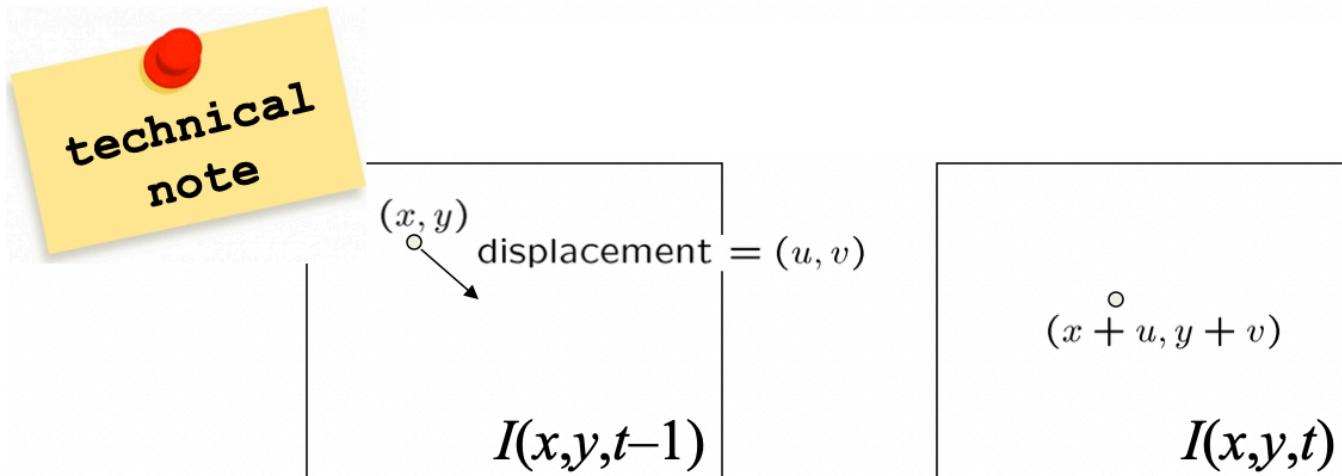
$I(x,y,t-1)$



$I(x,y,t)$

- Given two subsequent frames, estimate the apparent motion field  $u(x,y), v(x,y)$  between them
- Key assumptions
  - **Brightness constancy:** projection of the same point looks the same in every frame
  - **Small motion:** points do not move very far
  - **Spatial coherence:** points move like their neighbors

# The Brightness Constancy Constraint



- Brightness Constancy Equation:

$$I(x, y, t - 1) = I(x + u(x, y), y + v(x, y), t)$$

Linearizing the right side using Taylor expansion:

$$I(x + u, y + v, t) \approx I(x, y, t - 1) + I_x \cdot u(x, y) + I_y \cdot v(x, y) + I_t$$

Image derivative along x

$$I(x + u, y + v, t) - I(x, y, t - 1) = I_x \cdot u(x, y) + I_y \cdot v(x, y) + I_t$$

$$\text{Hence, } I_x \cdot u + I_y \cdot v + I_t \approx 0 \rightarrow \nabla I \cdot [u \ v]^T + I_t = 0$$

# The Brightness Constancy Constraint

Can we use this equation to recover image motion ( $u, v$ ) at each pixel?

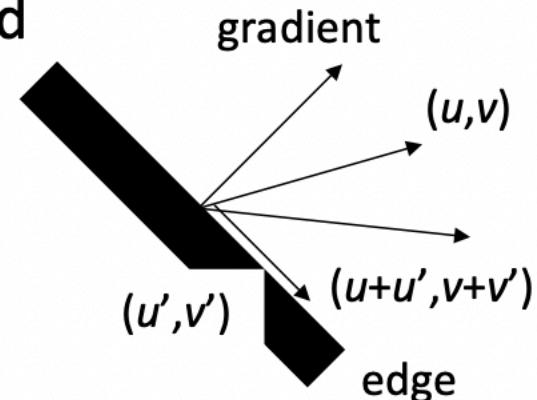
$$\nabla I \cdot [u \ v]^T + I_t = 0$$

- How many equations and unknowns per pixel?
  - One equation (this is a scalar equation!), two unknowns ( $u, v$ )

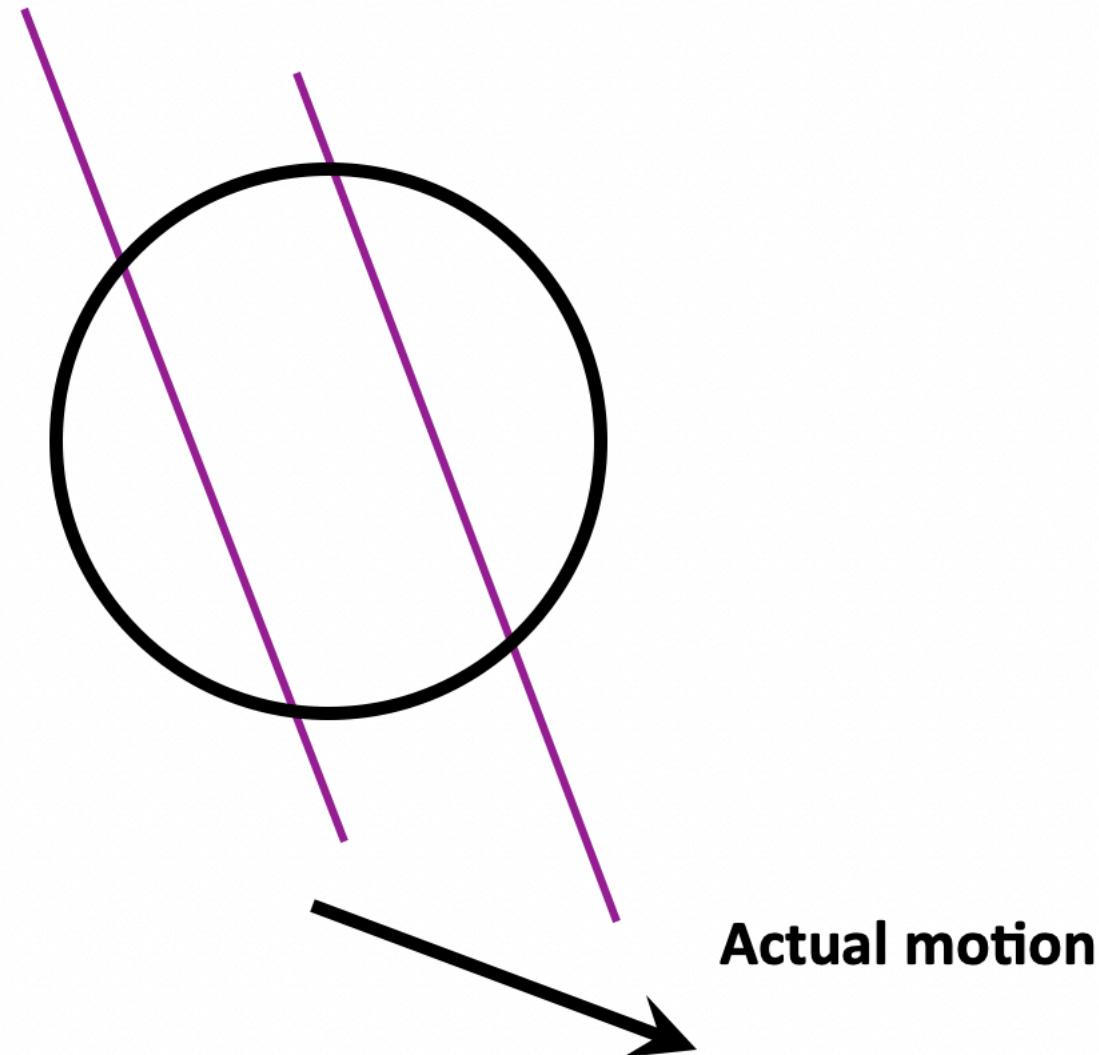
The component of the flow perpendicular to the gradient (i.e., parallel to the edge) cannot be measured

If  $(u, v)$  satisfies the equation,  
so does  $(u+u', v+v')$  if

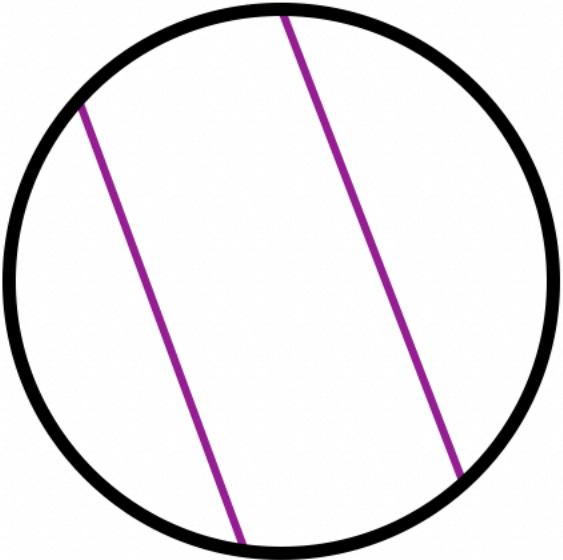
$$\nabla I \cdot [u' \ v']^T = 0$$



# The Aperture Problem



# The Aperture Problem



**Perceived motion**

# Barberpole Illusion



- [https://en.wikipedia.org/wiki/Barberpole\\_illusion](https://en.wikipedia.org/wiki/Barberpole_illusion)

# Solving the Ambiguity

- How to get more equations for a pixel?
- **Spatial coherence constraint:**
- Assume the pixel's neighbors have the same  $(u, v)$ 
  - If we use a 5x5 window, that gives us 25 equations per pixel

$$0 = I_t(\mathbf{p}_i) + \nabla I(\mathbf{p}_i) \cdot [u \ v]$$

$$\begin{bmatrix} I_x(\mathbf{p}_1) & I_y(\mathbf{p}_1) \\ I_x(\mathbf{p}_2) & I_y(\mathbf{p}_2) \\ \vdots & \vdots \\ I_x(\mathbf{p}_{25}) & I_y(\mathbf{p}_{25}) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} I_t(\mathbf{p}_1) \\ I_t(\mathbf{p}_2) \\ \vdots \\ I_t(\mathbf{p}_{25}) \end{bmatrix}$$

B. Lucas and T. Kanade. An iterative image registration technique with an application to stereo vision. In Proceedings of the International Joint Conference on Artificial Intelligence, pp. 674- 679, 1981.

# Lucas-Kanade Flow

- Overconstrained linear system:

$$\begin{bmatrix} I_x(p_1) & I_y(p_1) \\ I_x(p_2) & I_y(p_2) \\ \vdots & \vdots \\ I_x(p_{25}) & I_y(p_{25}) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} I_t(p_1) \\ I_t(p_2) \\ \vdots \\ I_t(p_{25}) \end{bmatrix}$$

$A \quad d = b$   
 $25 \times 2 \quad 2 \times 1 \quad 25 \times 1$

# Lucas-Kanade Flow

- Overconstrained linear system:

$$\begin{bmatrix} I_x(p_1) & I_y(p_1) \\ I_x(p_2) & I_y(p_2) \\ \vdots & \vdots \\ I_x(p_{25}) & I_y(p_{25}) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} I_t(p_1) \\ I_t(p_2) \\ \vdots \\ I_t(p_{25}) \end{bmatrix}$$

$A \quad d = b$   
 $25 \times 2 \quad 2 \times 1 \quad 25 \times 1$

Least squares solution for  $d$  given by  $(A^T A)^{-1} A^T b$

$$\begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} \sum I_x I_t \\ \sum I_y I_t \end{bmatrix}$$

$A^T A \qquad \qquad \qquad A^T b$

The summations are over all pixels in the  $K \times K$  window

# Condition for Solvability

– Optimal  $(u, v)$  satisfies Lucas-Kanade equation

$$\begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} \sum I_x I_t \\ \sum I_y I_t \end{bmatrix}$$

$A^T A$                                      $A^T b$

## When is This Solvable?

- $A^T A$  should be invertible
- $A^T A$  should not be too small due to noise
  - eigenvalues  $\lambda_1$  and  $\lambda_2$  of  $A^T A$  should not be too small
- $A^T A$  should be well-conditioned
  - $\lambda_1 / \lambda_2$  should not be too large ( $\lambda_1$  = larger eigenvalue)

Does this remind anything to you?

# Condition for Solvability

$M = A^T A$  is the *second moment matrix* !

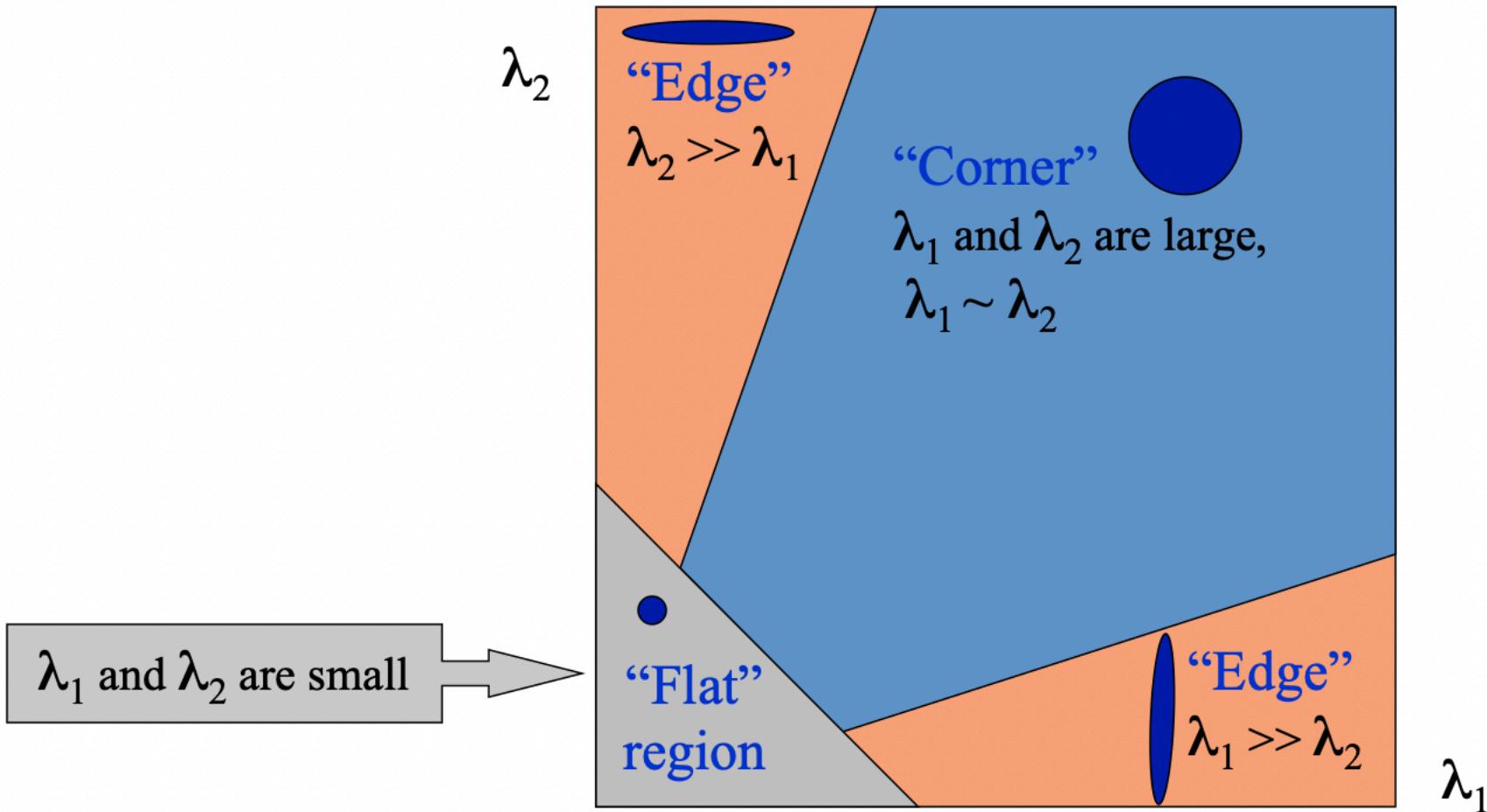
(Harris corner detector...)

$$A^T A = \begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix} = \sum \begin{bmatrix} I_x \\ I_y \end{bmatrix} [I_x \ I_y] = \sum \nabla I (\nabla I)^T$$

- Eigenvectors and eigenvalues of  $A^T A$  relate to edge direction and magnitude
  - The eigenvector associated with the larger eigenvalue points in the direction of fastest intensity change
  - The other eigenvector is orthogonal to it

# Interpreting the Eigenvalues

Classification of image points using eigenvalues of the second moment matrix:



# Edge



$$\sum \nabla I(\nabla I)^T$$

- gradients very large or very small
- large  $\lambda_1$ , small  $\lambda_2$

# Low-Texture Region



$$\sum \nabla I (\nabla I)^T$$

- gradients have small magnitude
- small  $\lambda_1$ , small  $\lambda_2$

# High Texture Region



$$\sum \nabla I (\nabla I)^T$$

- gradients are different, large magnitudes
- large  $\lambda_1$ , large  $\lambda_2$

# Measuring Motion

Image at frame  $t$

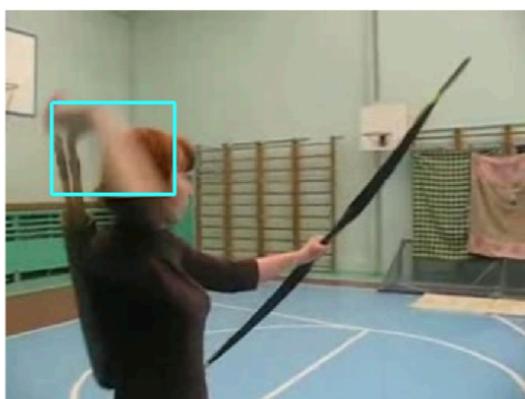
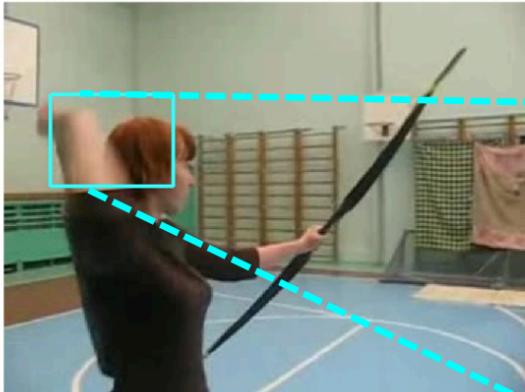
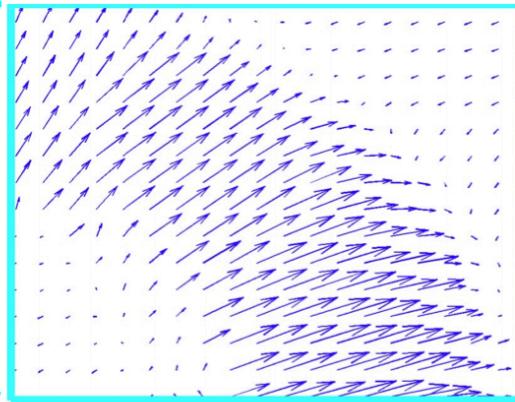


Image at frame  $t+1$

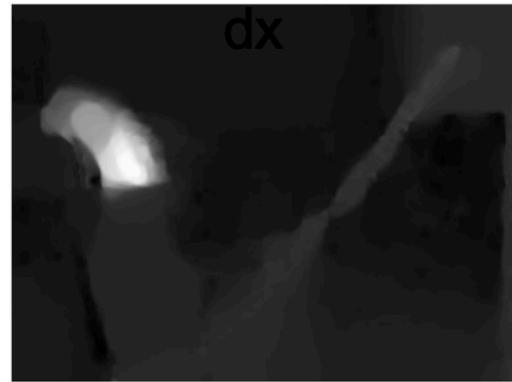
Optical flow gives a displacement field  $F$  between images  $I_t$  and  $I_{t+1}$



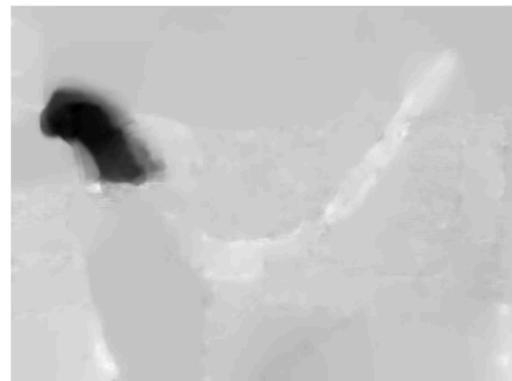
Tells where each pixel will move in the next frame:  
 $F(x, y) = (dx, dy)$   
 $I_{t+1}(x+dx, y+dy) = I_t(x, y)$

Optical Flow highlights  
**local motion**

Horizontal flow



$dx$



Vertical Flow  $dy$

Slide credit: Justin Johnson

# FlowNet: Learning Optical Flow with Convolutional Networks

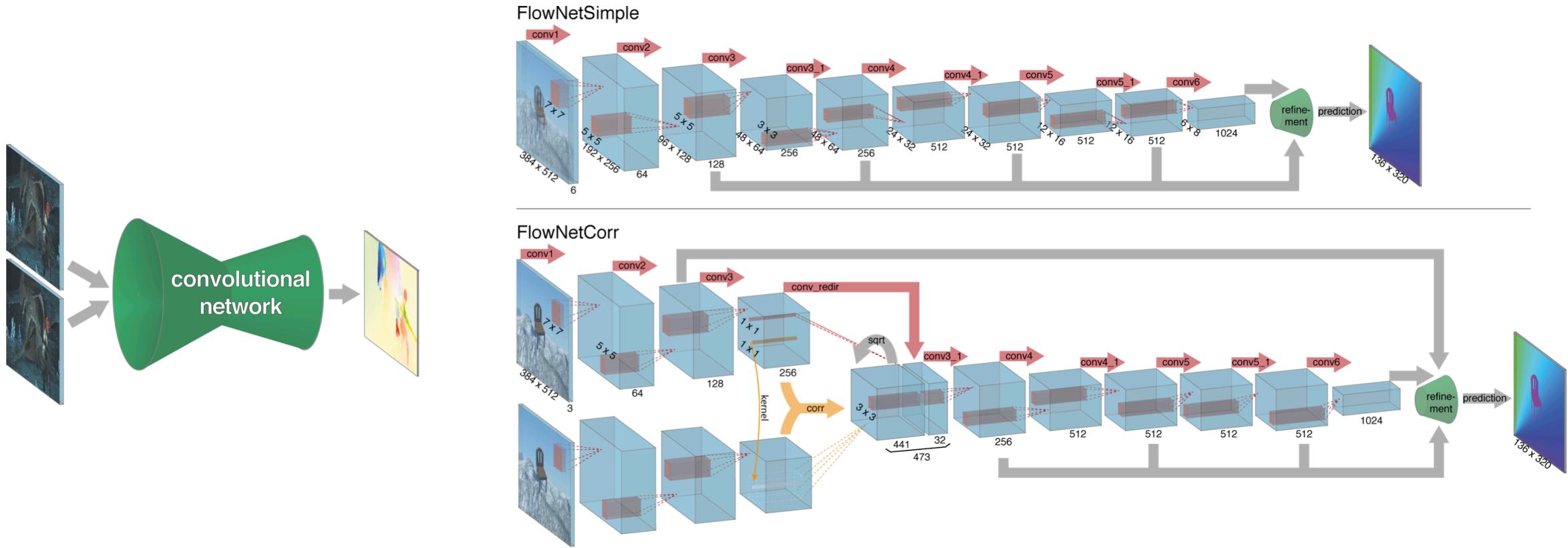


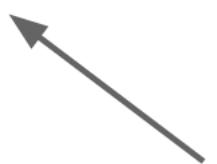
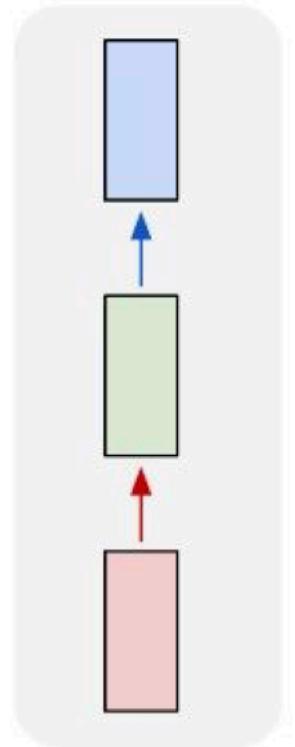
Figure 2. The two network architectures: FlowNetSimple (top) and FlowNetCorr (bottom).

# RNN

Some slides are borrowed from Stanford CS231N

# Single-Frame Neural Network

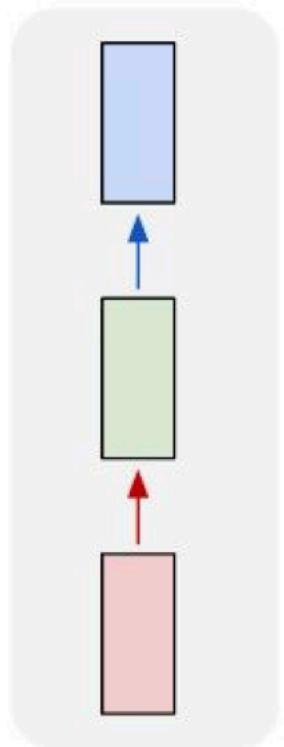
one to one



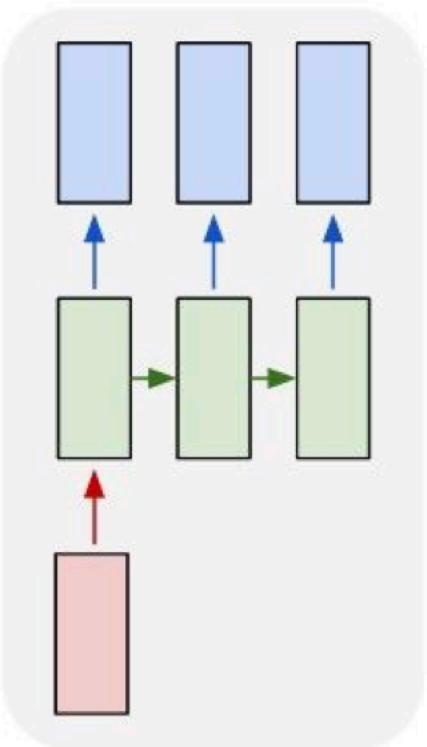
**Vanilla Neural Networks**

# Recurrent Neural Network: Process Sequential Data

one to one



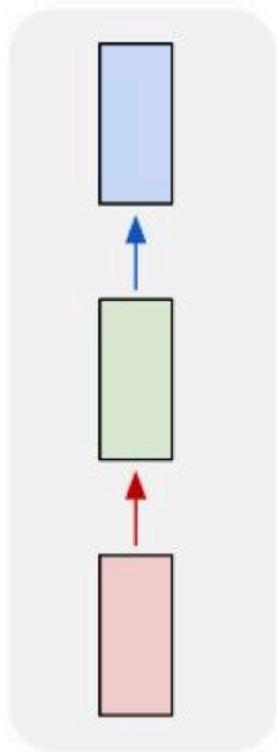
one to many



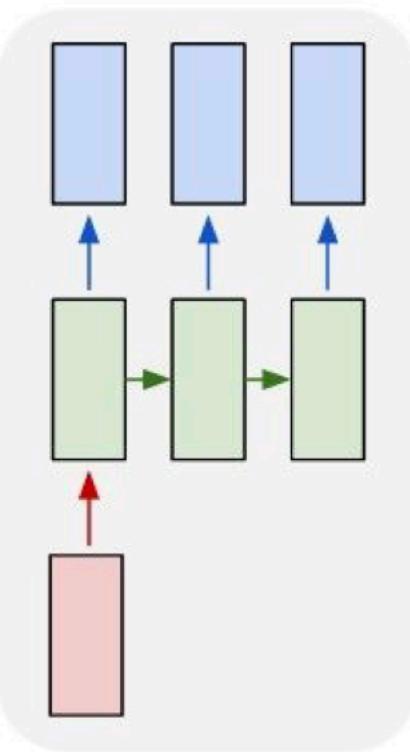
e.g. **Image Captioning**  
image -> sequence of words

# Recurrent Neural Network: Process Sequential Data

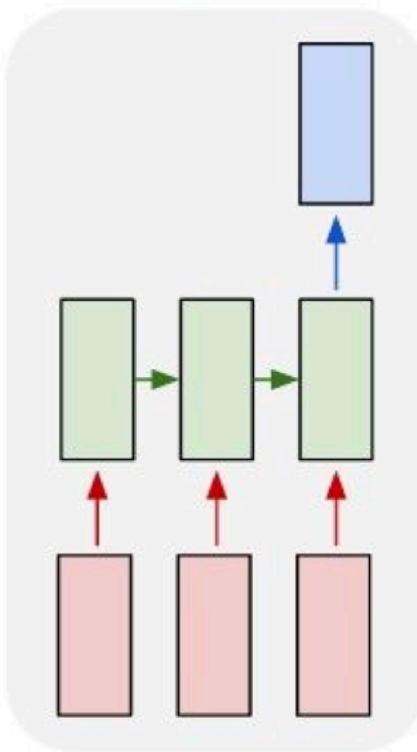
one to one



one to many



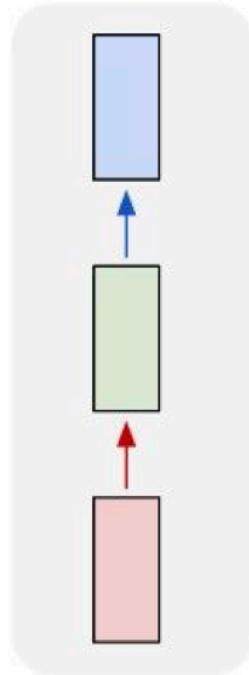
many to one



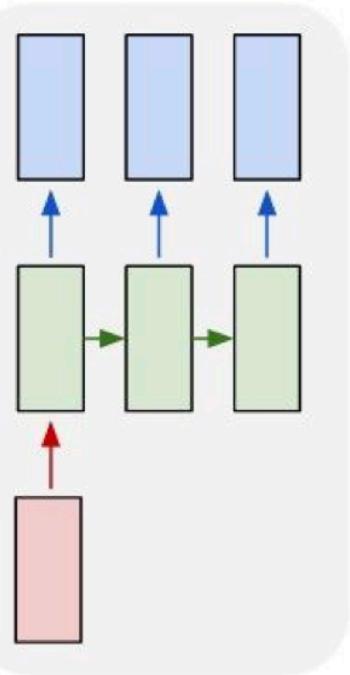
e.g. **action prediction**  
sequence of video frames -> action class

# Recurrent Neural Network: Process Sequential Data

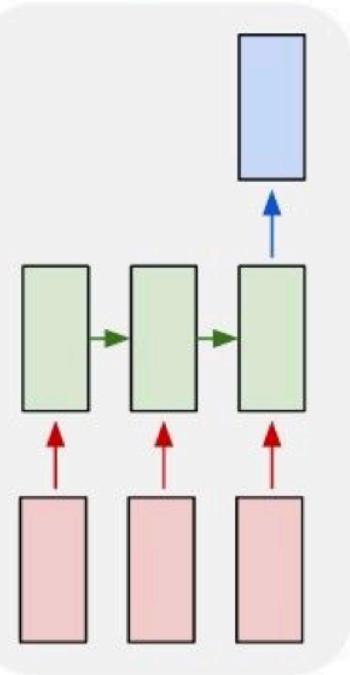
one to one



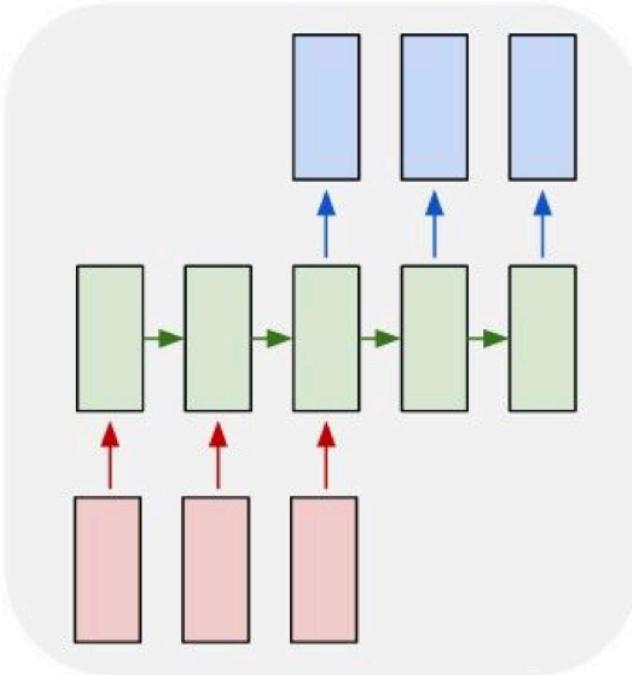
one to many



many to one

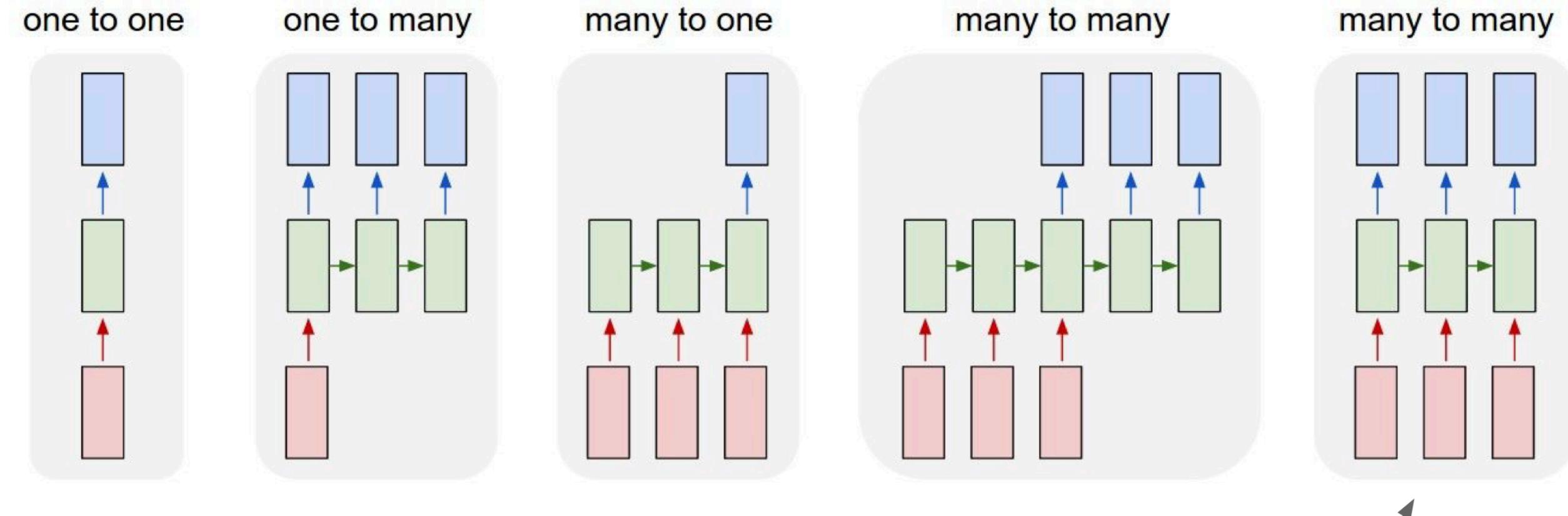


many to many



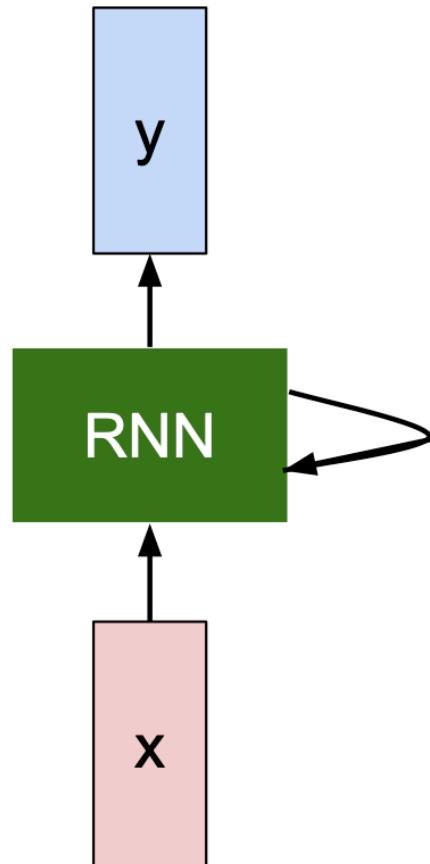
E.g. **Video Captioning**  
Sequence of video frames -> caption

# Recurrent Neural Network: Process Sequential Data

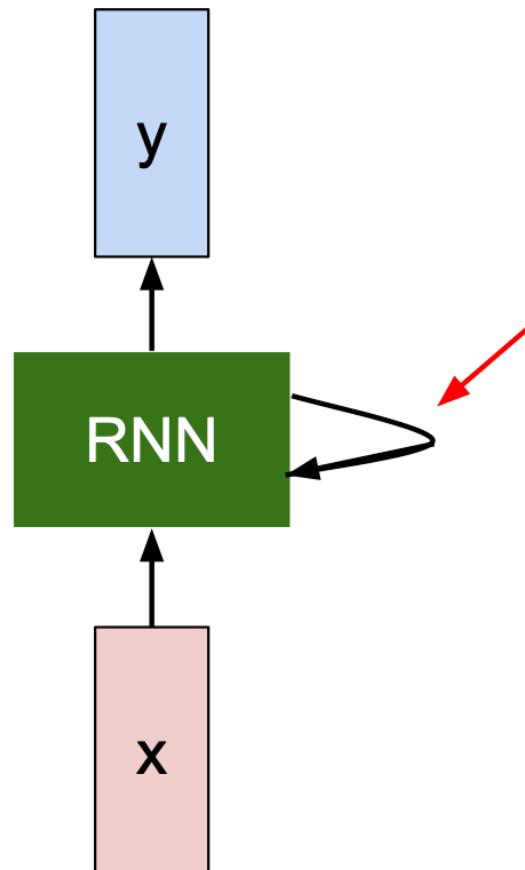


e.g. **Video classification on frame level**

# Recurrent Neural Network

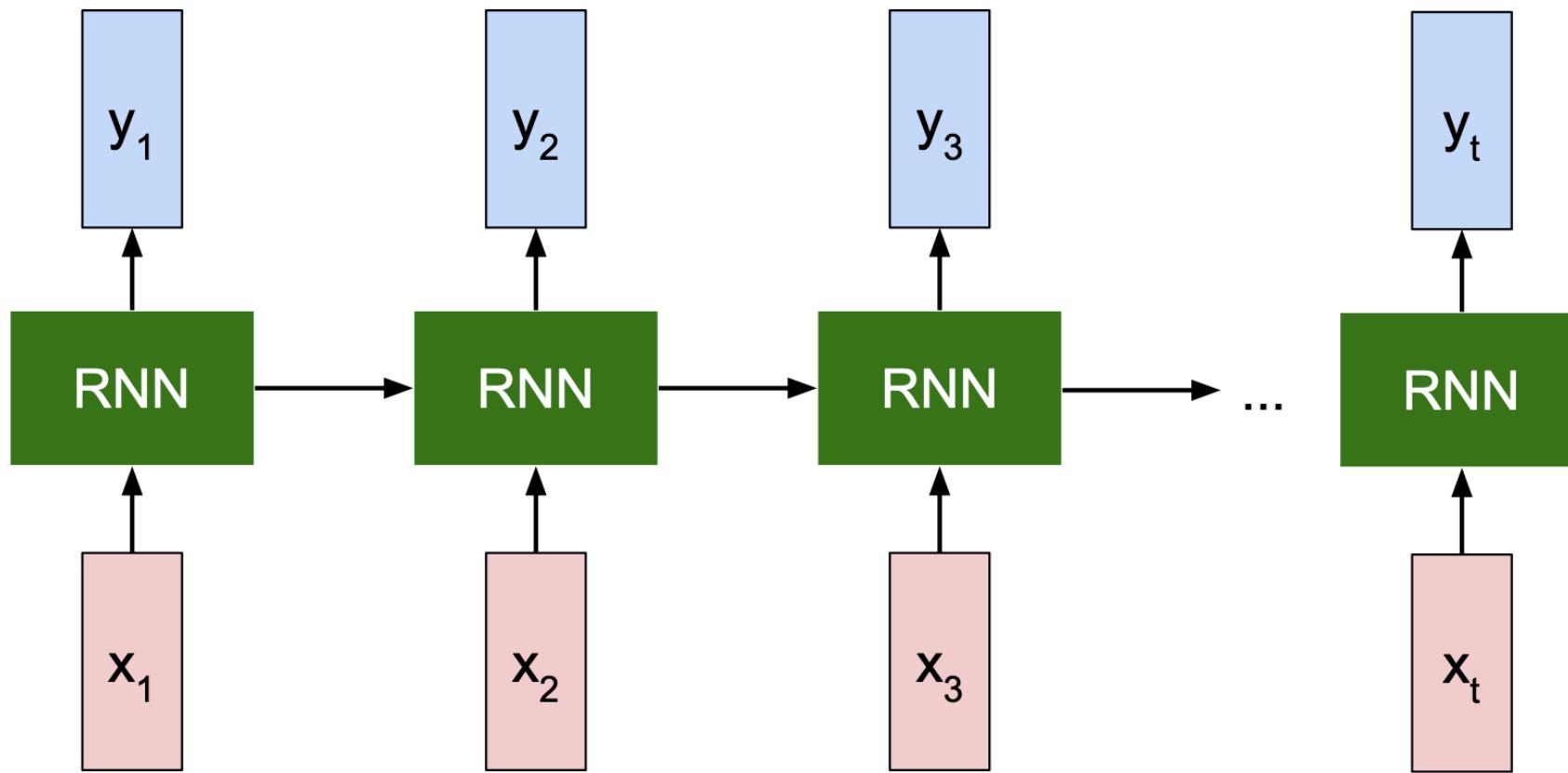


# Recurrent Neural Network



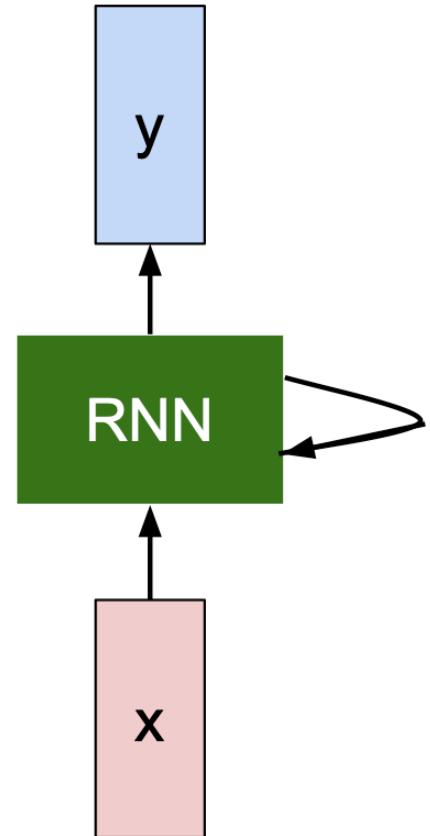
Key idea: RNNs have an “internal state” that is updated as a sequence is processed

# Unrolled RNN



# RNN Hidden State Update

We can process a sequence of vectors  $\mathbf{x}$  by applying a **recurrence formula** at every time step:

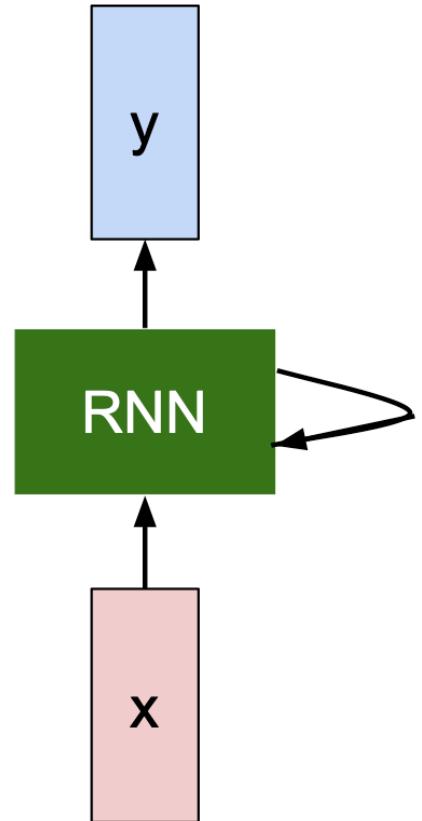


# RNN Output

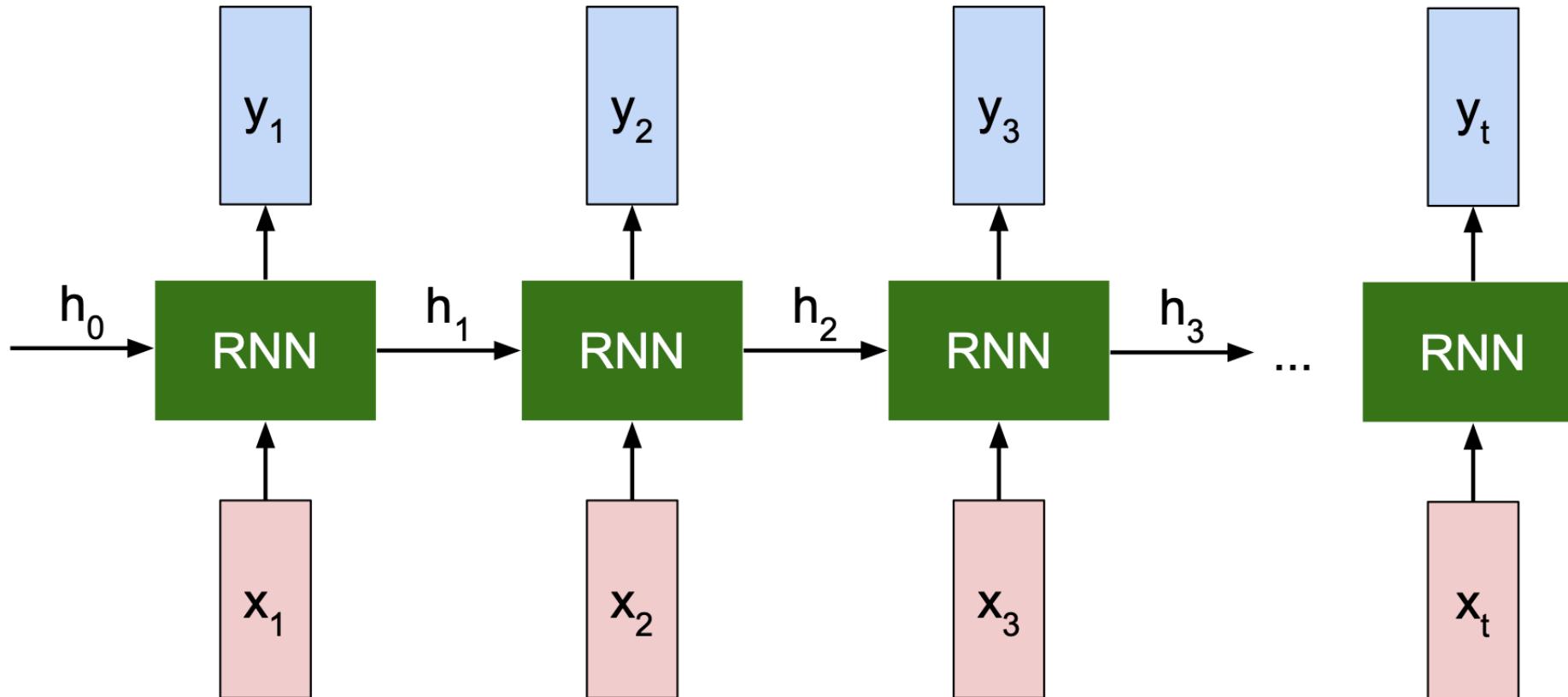
We can process a sequence of vectors  $x$  by applying a **recurrence formula** at every time step:

$$y_t = f_{W_{hy}}(h_t)$$

output                          new state  
another function  
with parameters  $W_o$



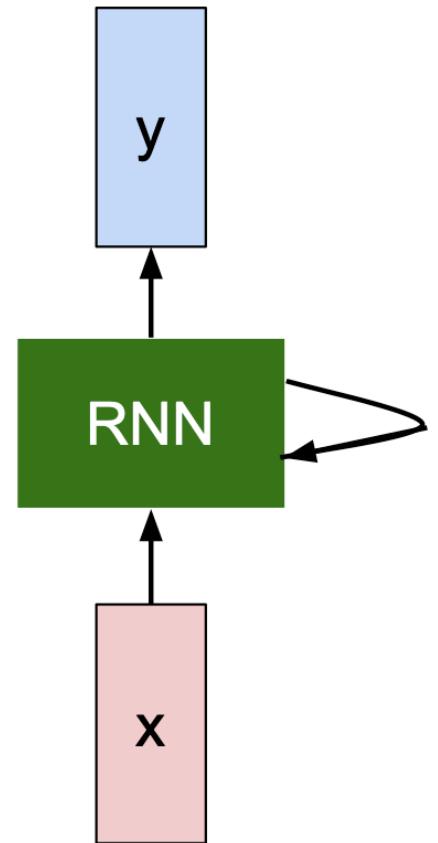
# Recurrent Neural Network



# Recurrent Neural Network

We can process a sequence of vectors  $\mathbf{x}$  by applying a **recurrence formula** at every time step:

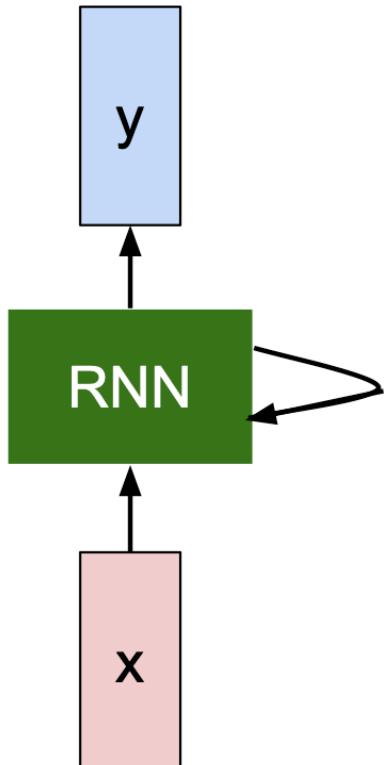
$$h_t = f_W(h_{t-1}, x_t)$$



Notice: the same function and the same set of parameters are used at every time step.

# Vanilla RNN

The state consists of a single “*hidden*” vector  $\mathbf{h}$ :



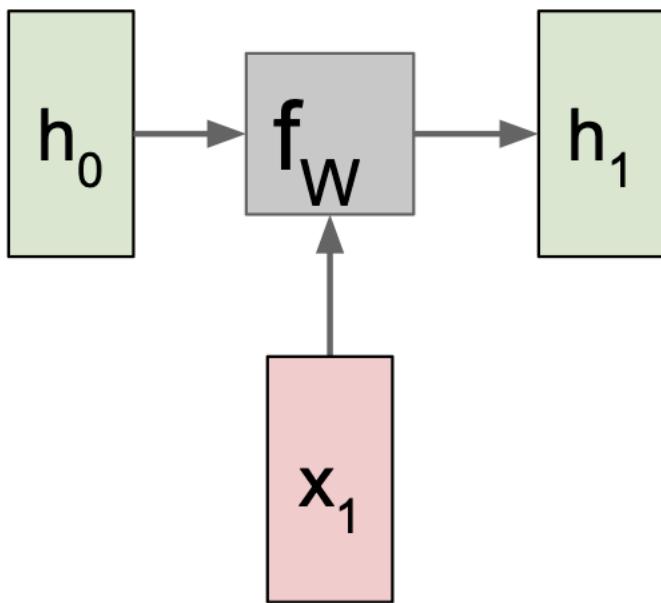
$$h_t = f_W(h_{t-1}, x_t)$$

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

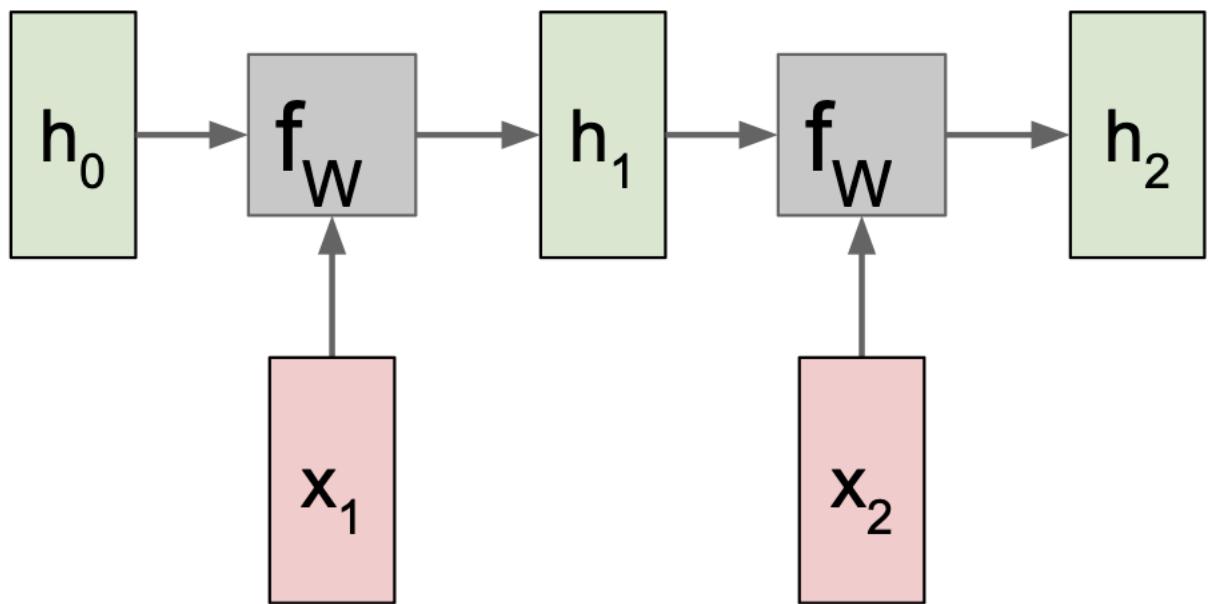
$$y_t = W_{hy}h_t$$

Sometimes called a “Vanilla RNN” or an  
“Elman RNN” after Prof. Jeffrey Elman

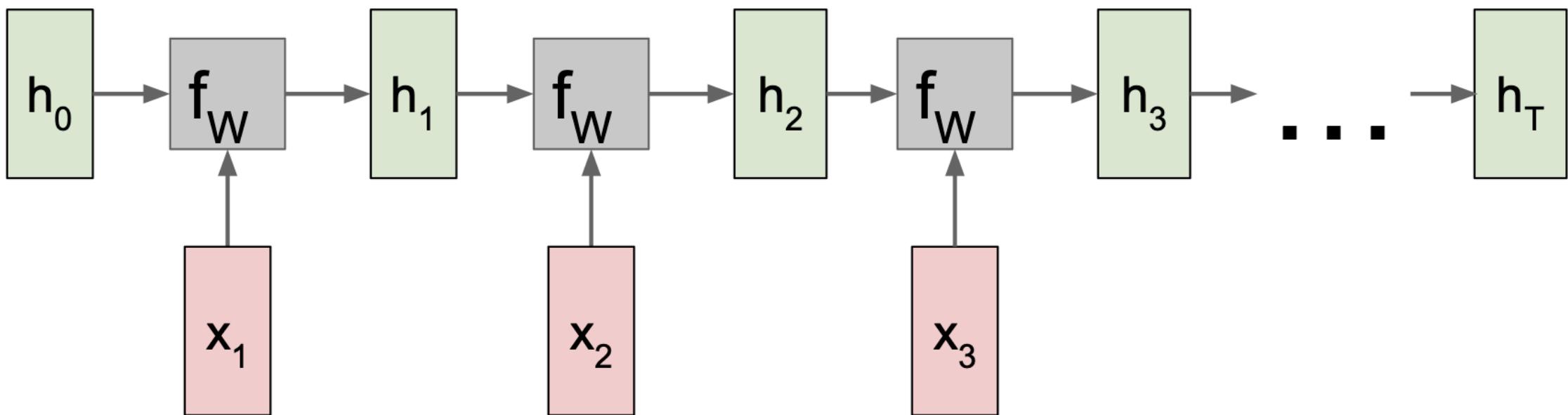
# RNN: Computational Graph



# RNN: Computational Graph

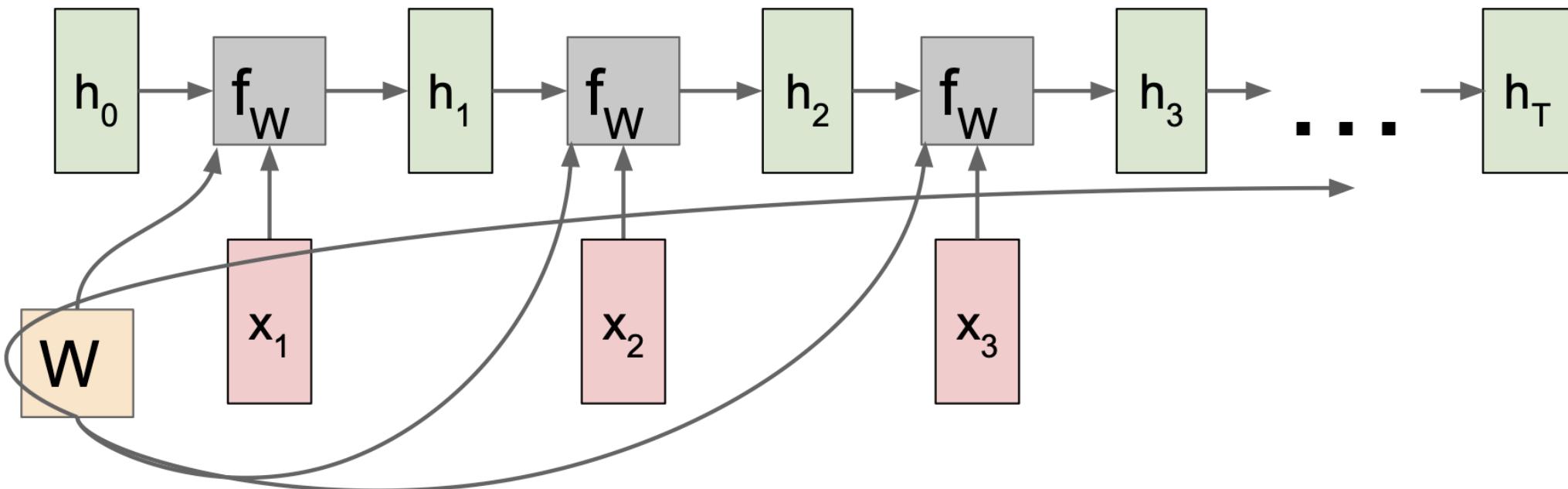


# RNN: Computational Graph



# RNN: Computational Graph

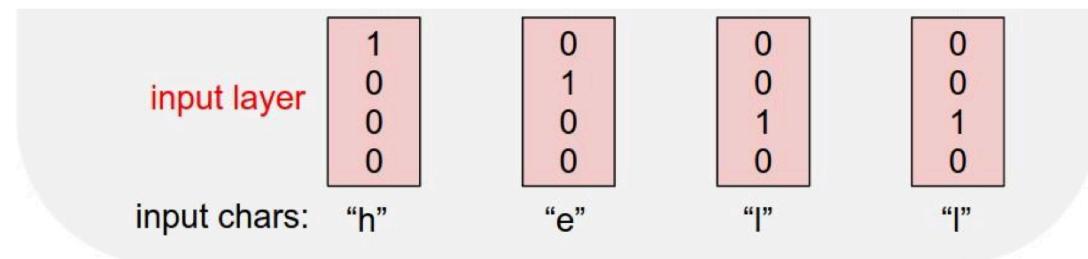
Re-use the same weight matrix at every time-step



# Character-Level Language Model

Vocabulary:  
[h,e,l,o]

Example training  
sequence:  
“hello”

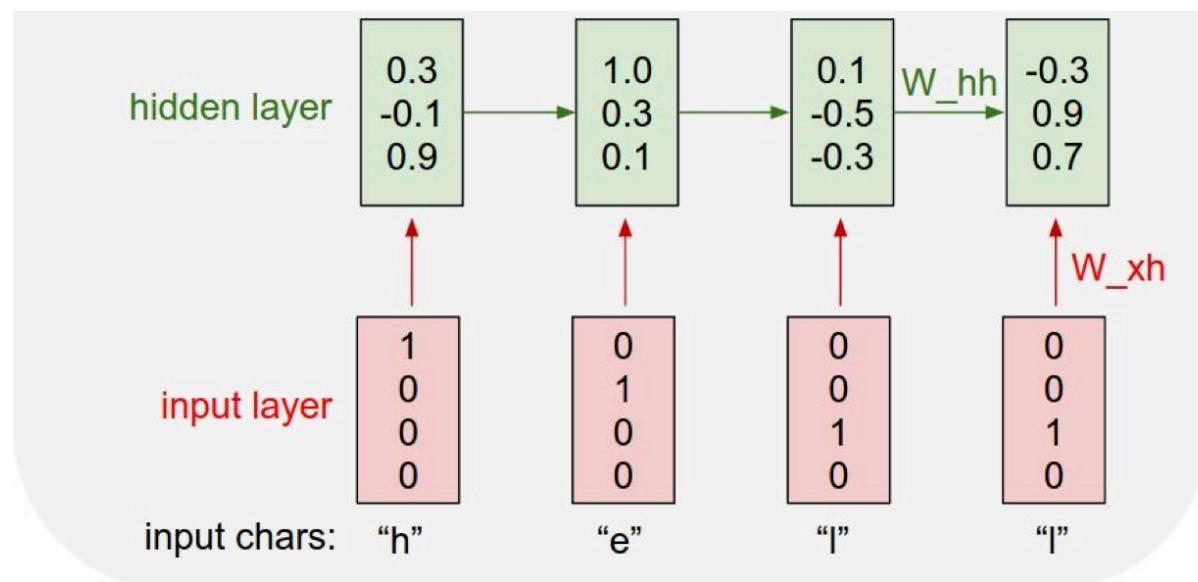


# Character-Level Language Model

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

Vocabulary:  
[h,e,l,o]

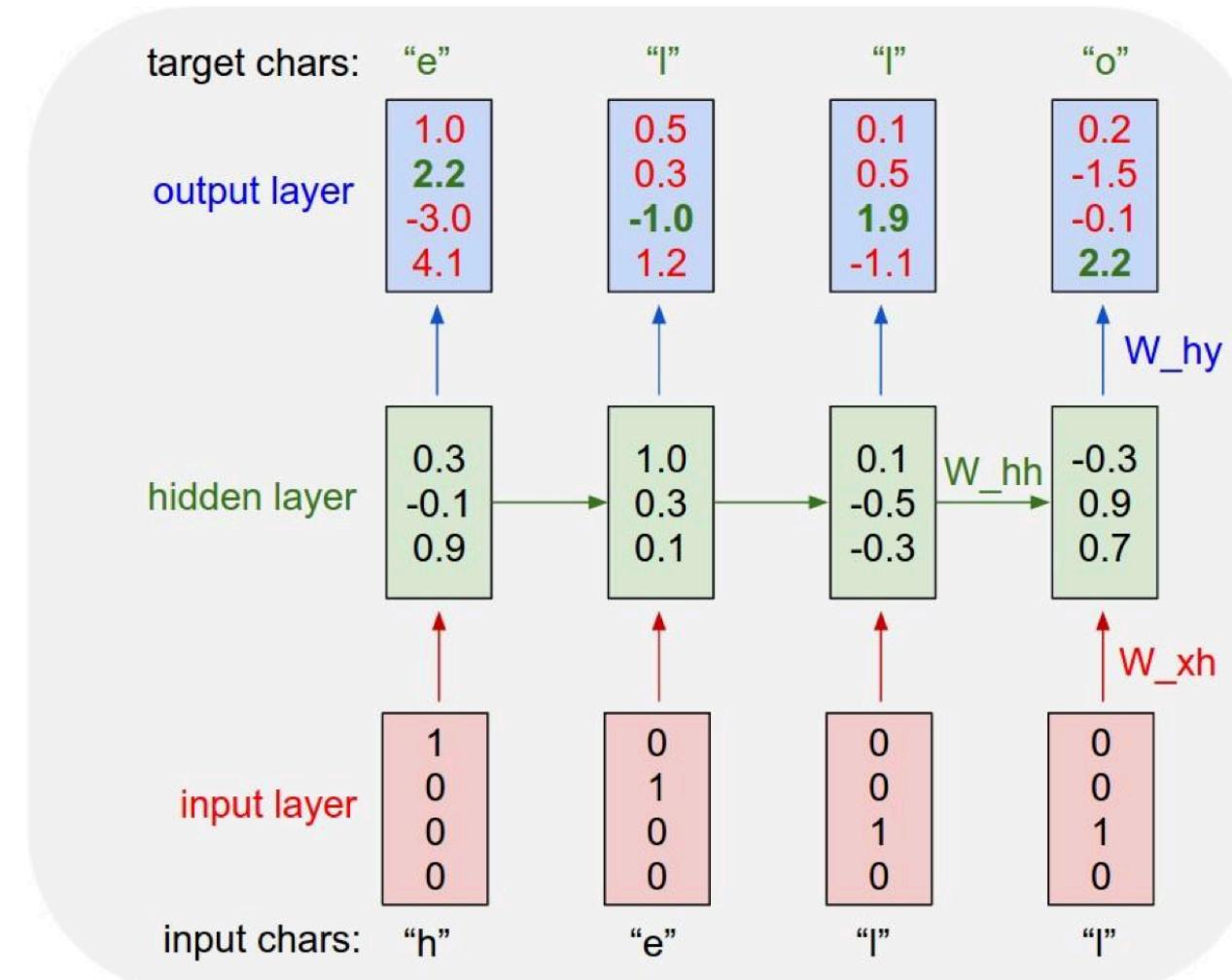
Example training  
sequence:  
“hello”



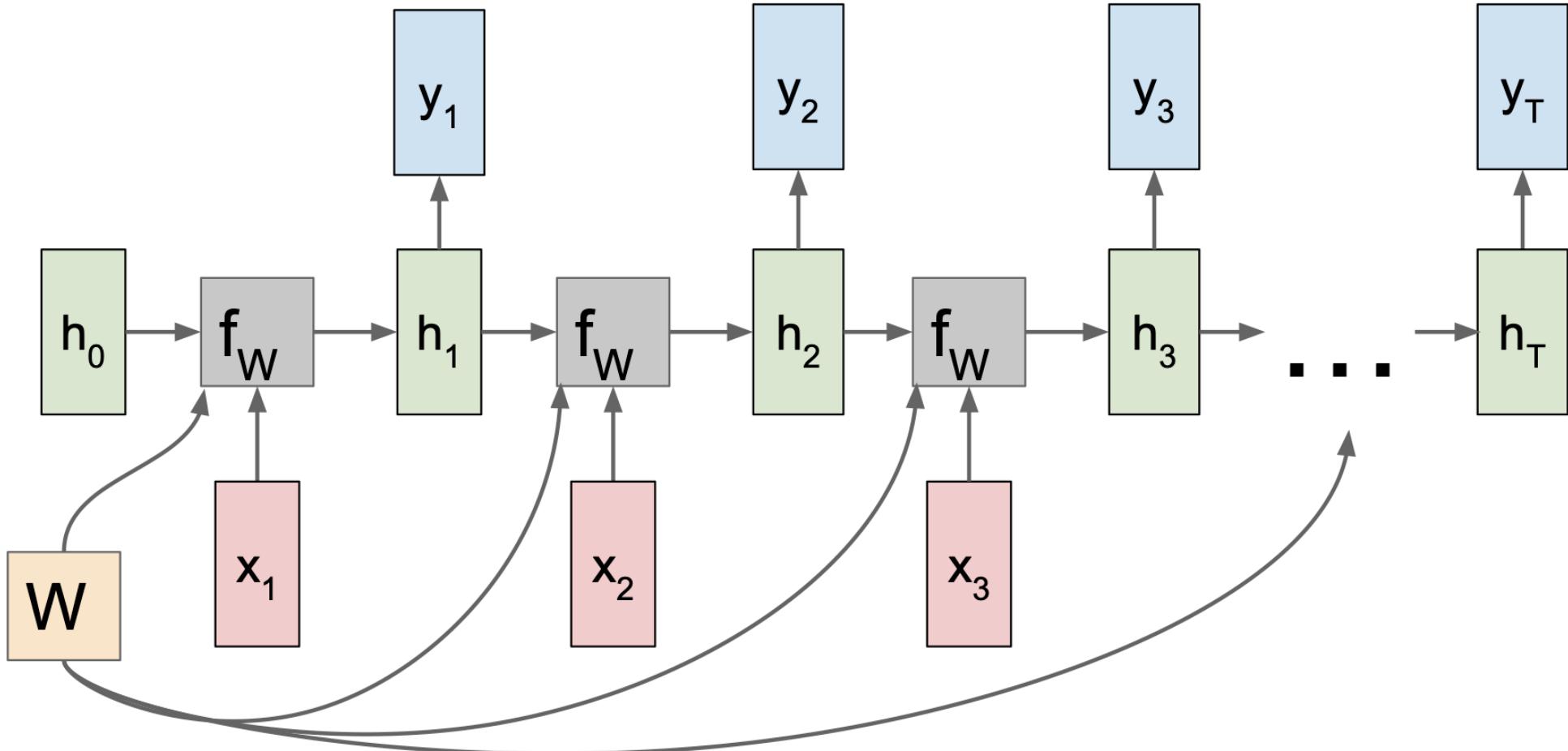
# Character-Level Language Model

Vocabulary:  
[h,e,l,o]

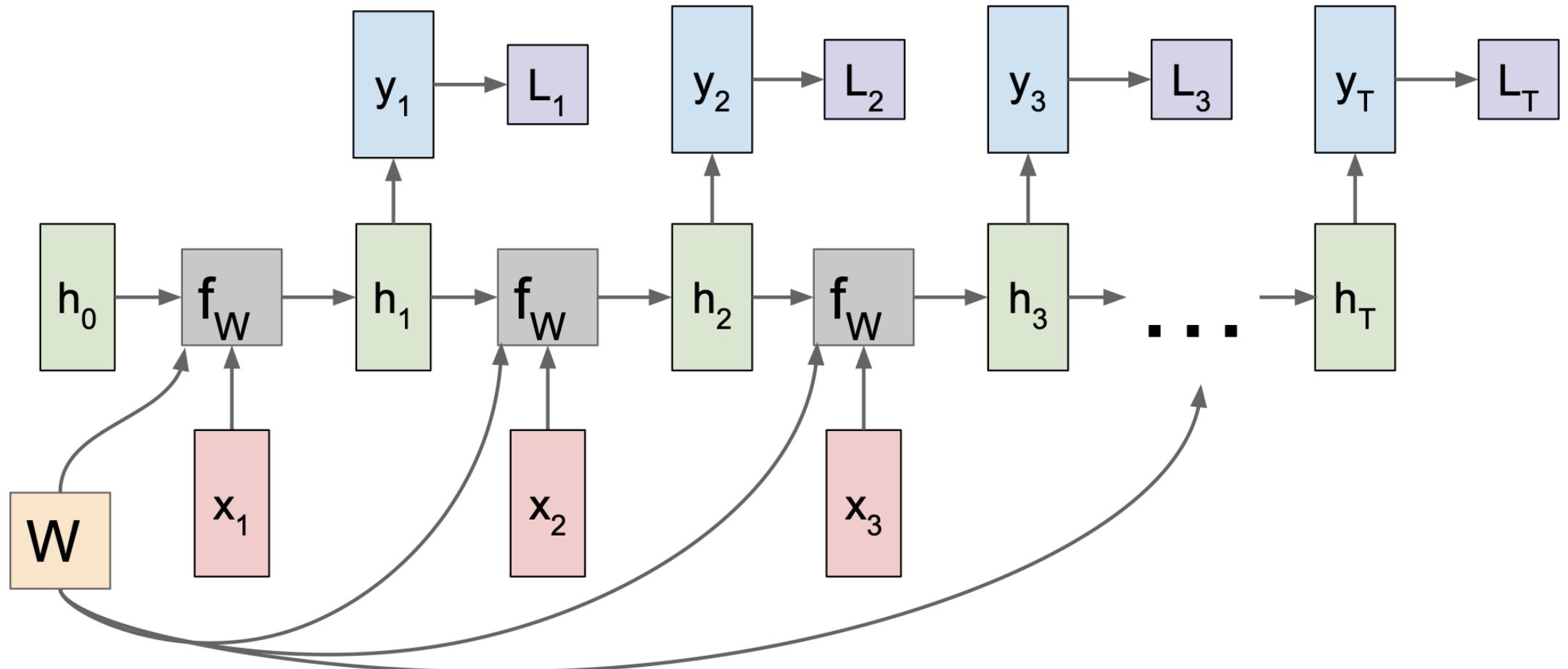
Example training  
sequence:  
“hello”



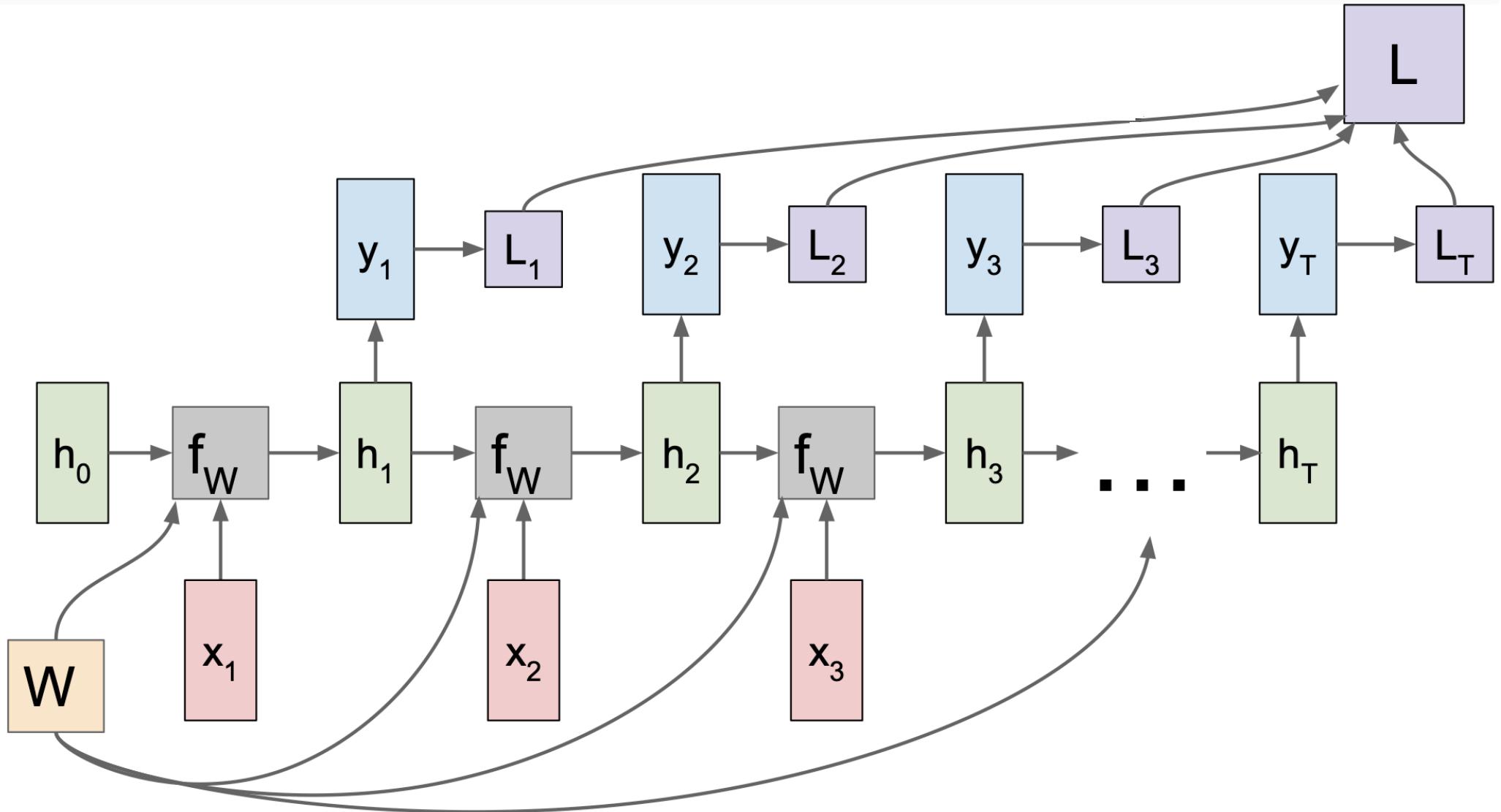
# Many-to-Many RNN: Computational Graph



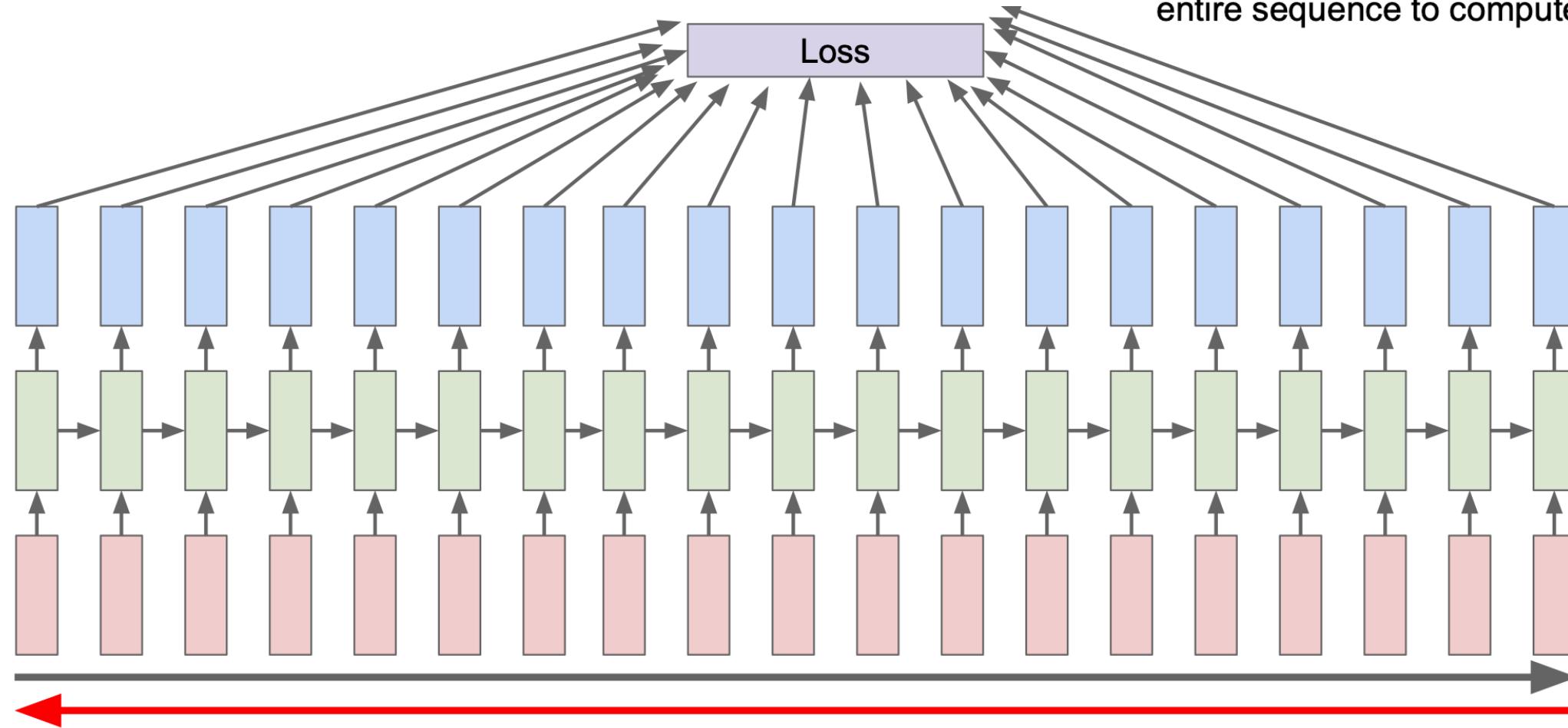
# Many-to-Many RNN: Computational Graph



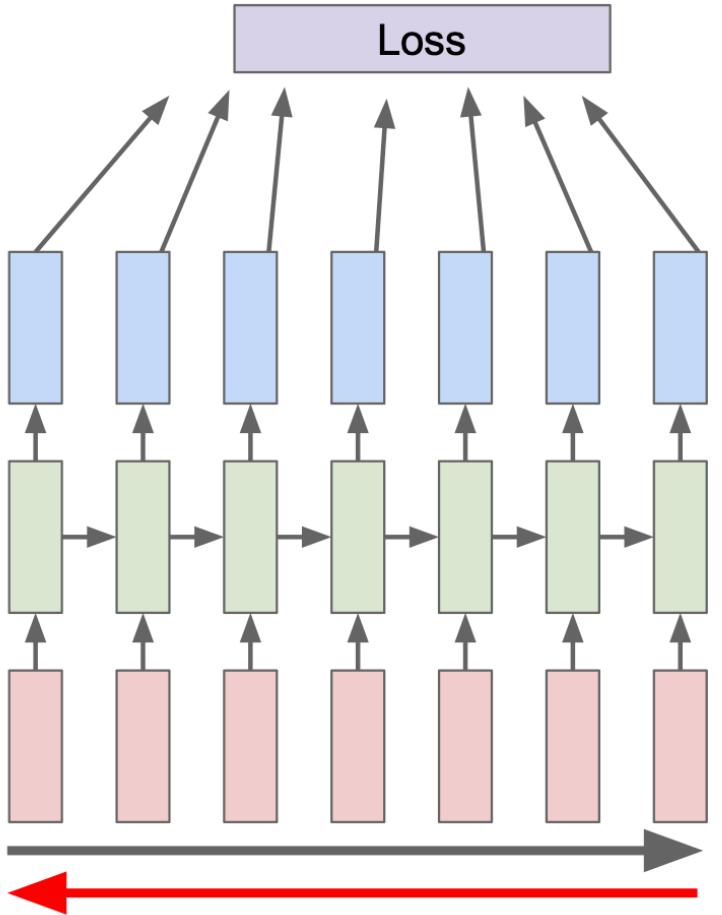
# Many-to-Many RNN: Computational Graph



# Backpropagation through Time

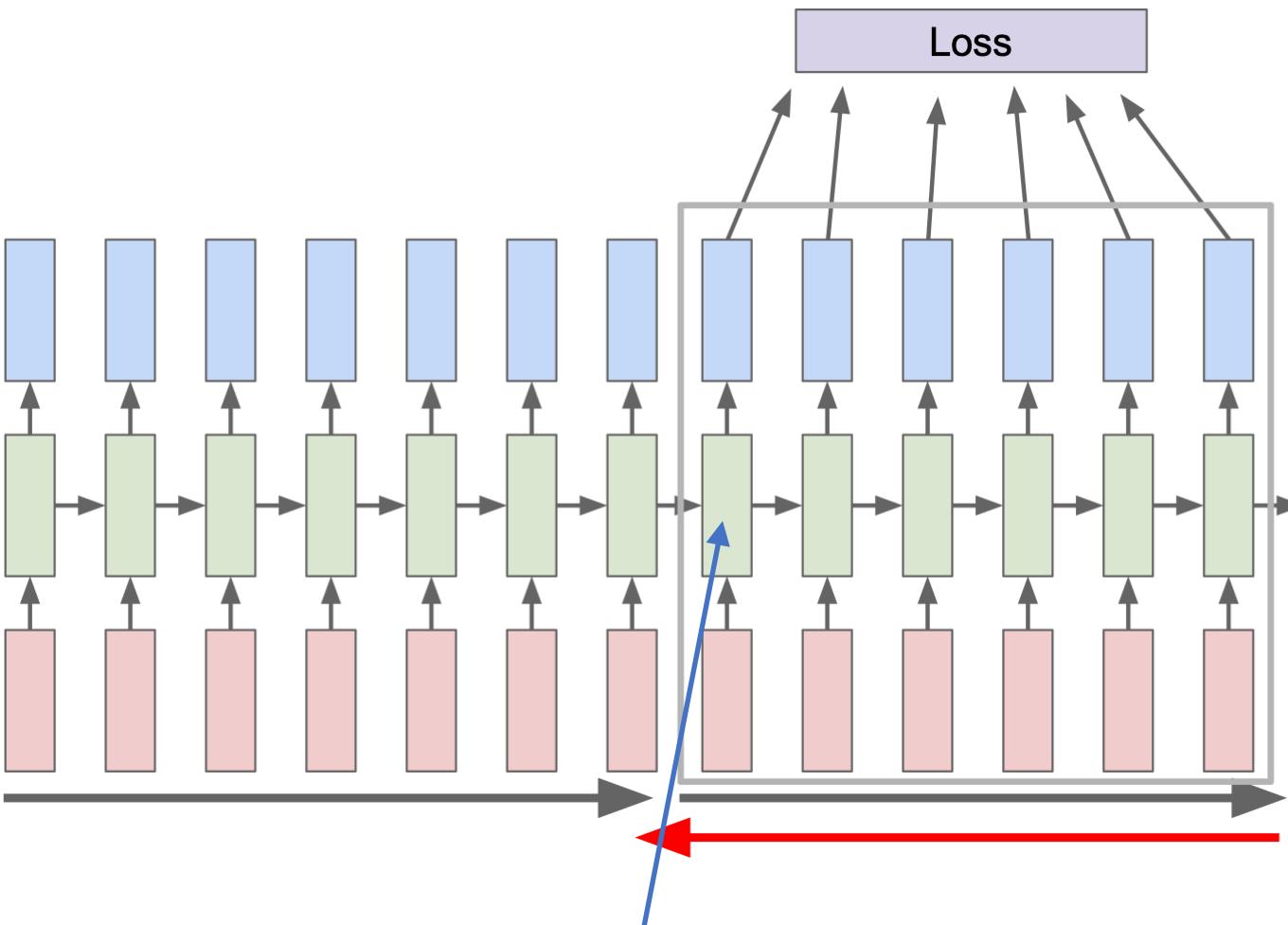


# Truncated Backpropagation through Time



Run forward and backward  
through chunks of the  
sequence instead of whole  
sequence

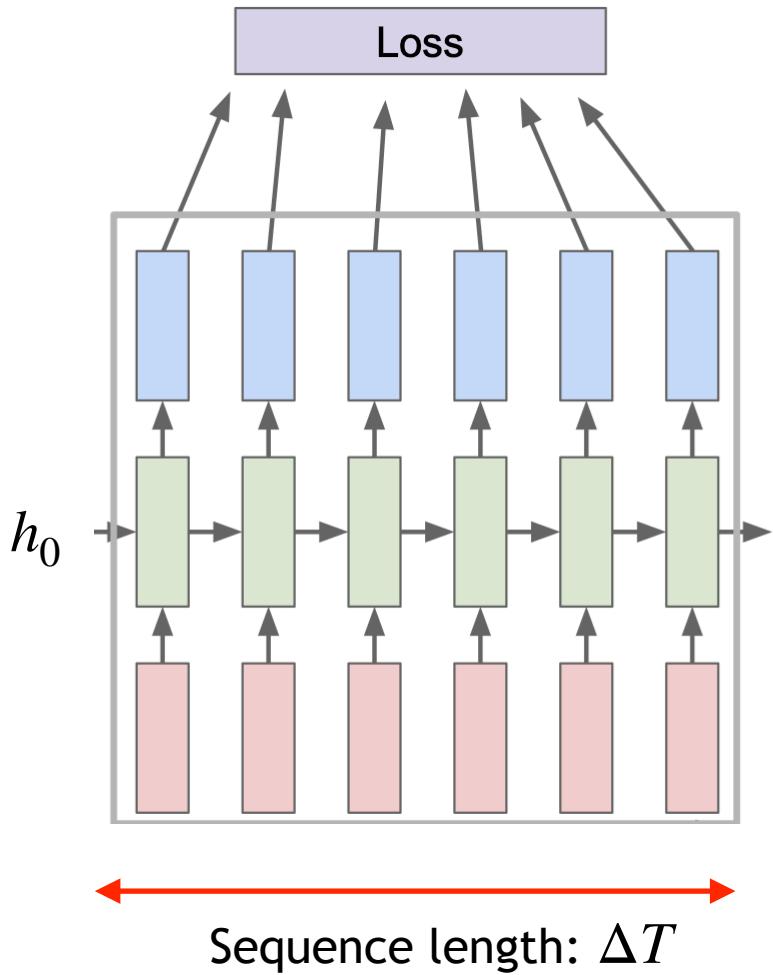
# Truncated Backpropagation through Time



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

$h_t$  or in practice can use  $h_0$  if  $t$  is so very large such that we even don't want to compute  $h_t$  during training

# Truncated Backpropagation through Time



If use  $h_0$ , then training only enforces the temporal relationship inside the window, whose length is **sequence length**.

Long-term relationship longer than  $\Delta T$  will probably not be learned.



# Introduction to Computer Vision

Next Week: Labor Day Break  
Next lecture: Lecture 10,  
Temporal Analysis II