

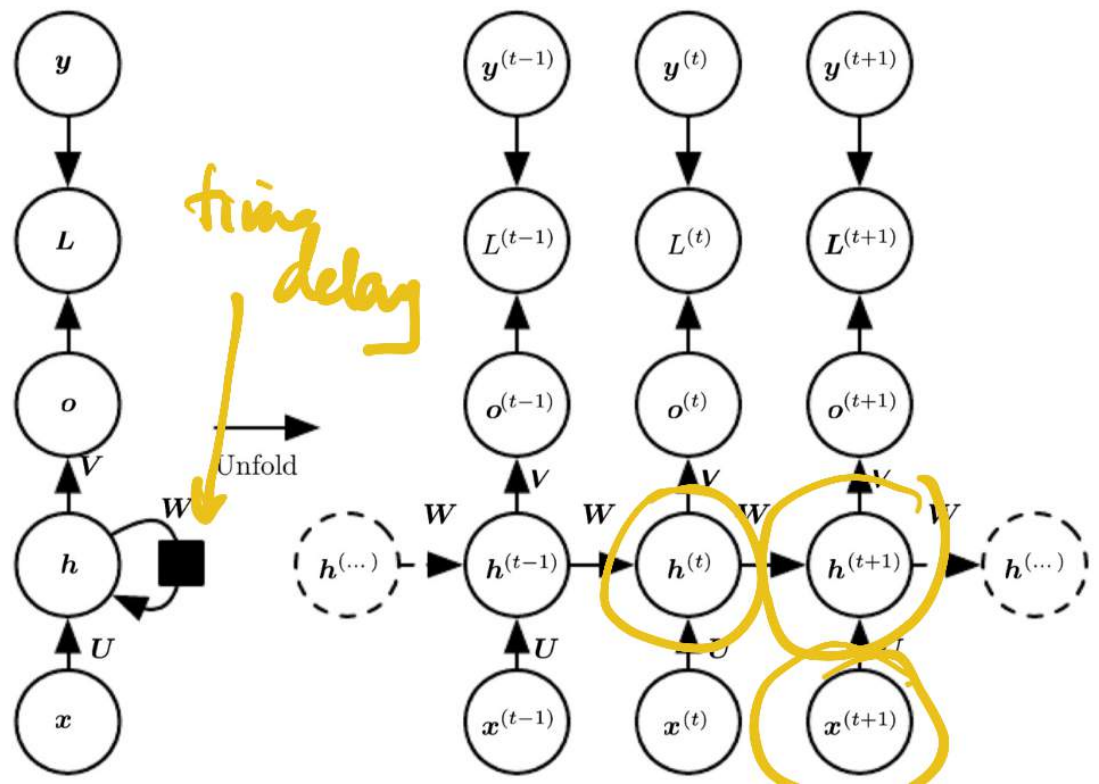
# Advanced Machine Learning, 3 Sep 2019

## Sequence input - RNNs

Speech processing, handwriting, bioinformatics

Have feedback from earlier inputs

Hidden layer  
at time  $t+1$   
depends on  $h^{(t)}$

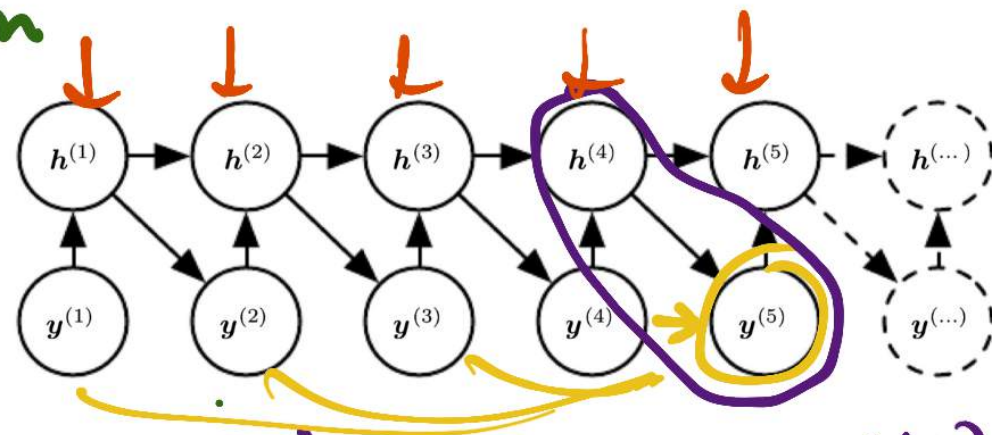
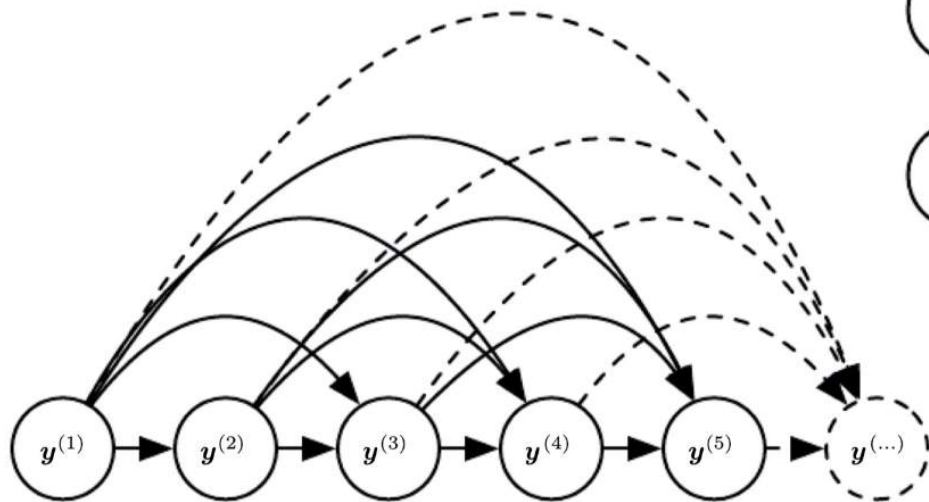


Unfolded network

forward & backward propagation

Backward Propagation through time

Graphical interpretation



Dependence of  $y^{(t+n)}$   
on  $y^{(t)}$  is  
factorized via  $h^{(t)}$

$$y|x \rightarrow y|h \quad h=f(x)$$

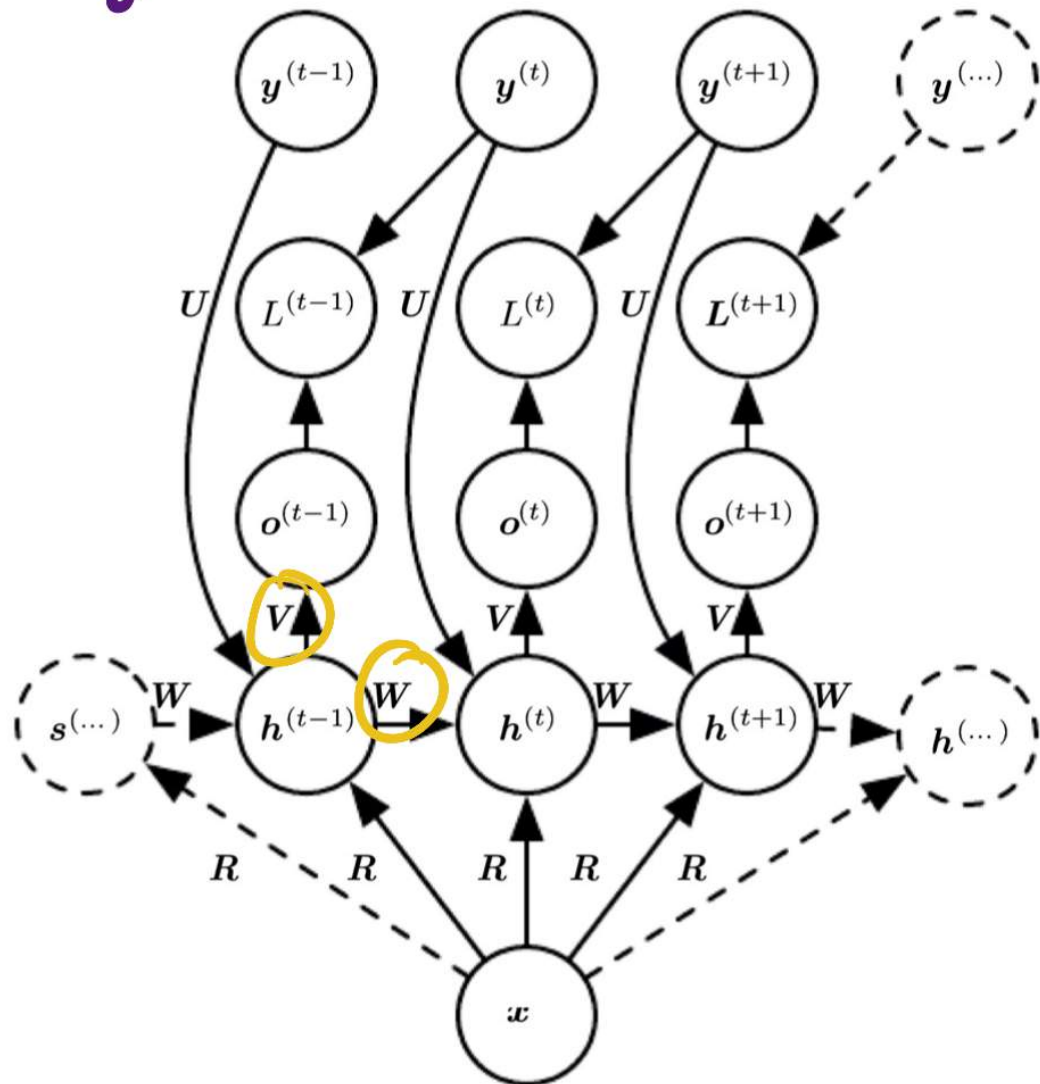
Input is generically called "Context"

Feeding  
context to  
an RNN

Option 1

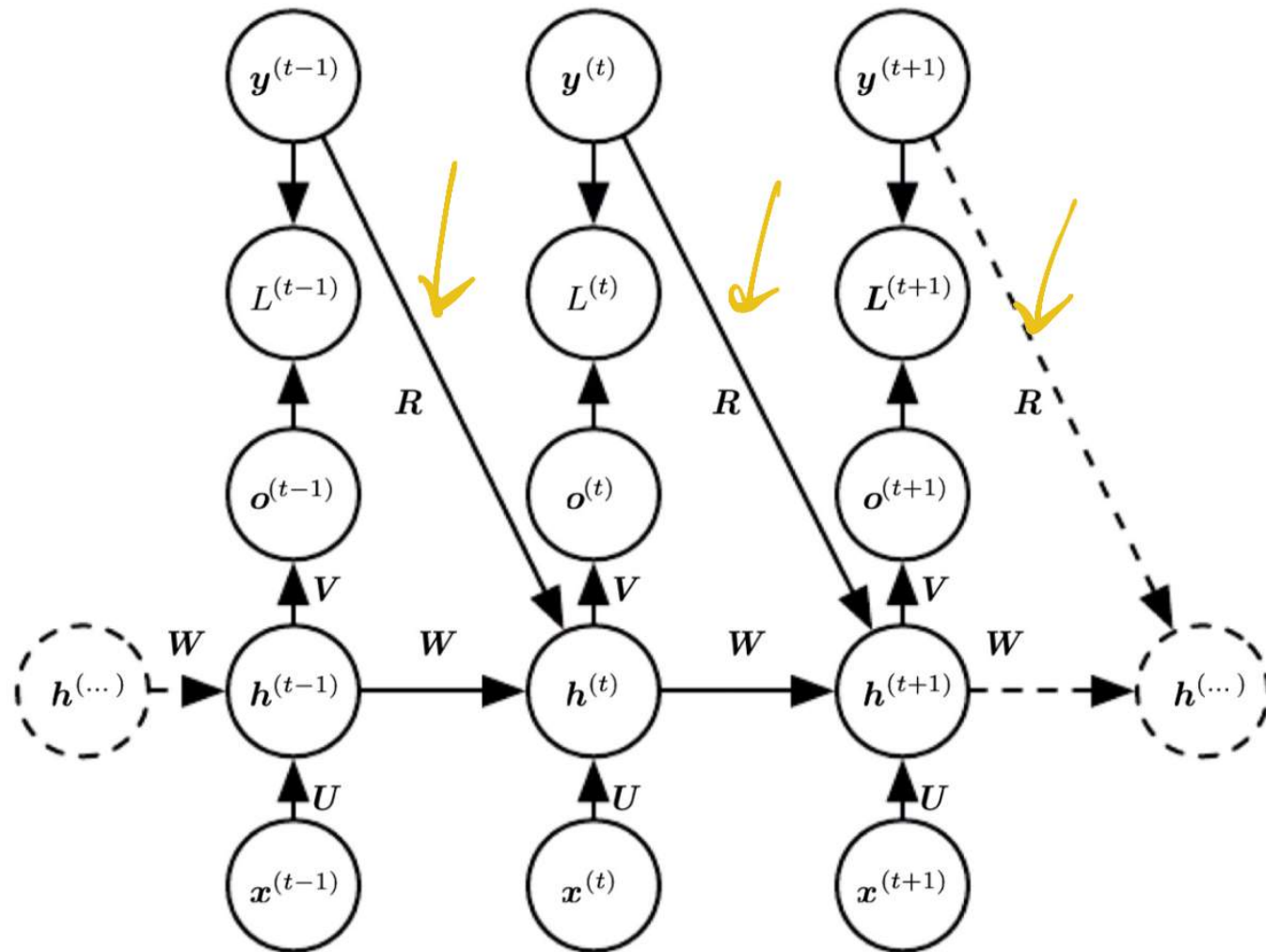
All nodes get

same  $x$



More "obvious" way  
is one input per  
time instant

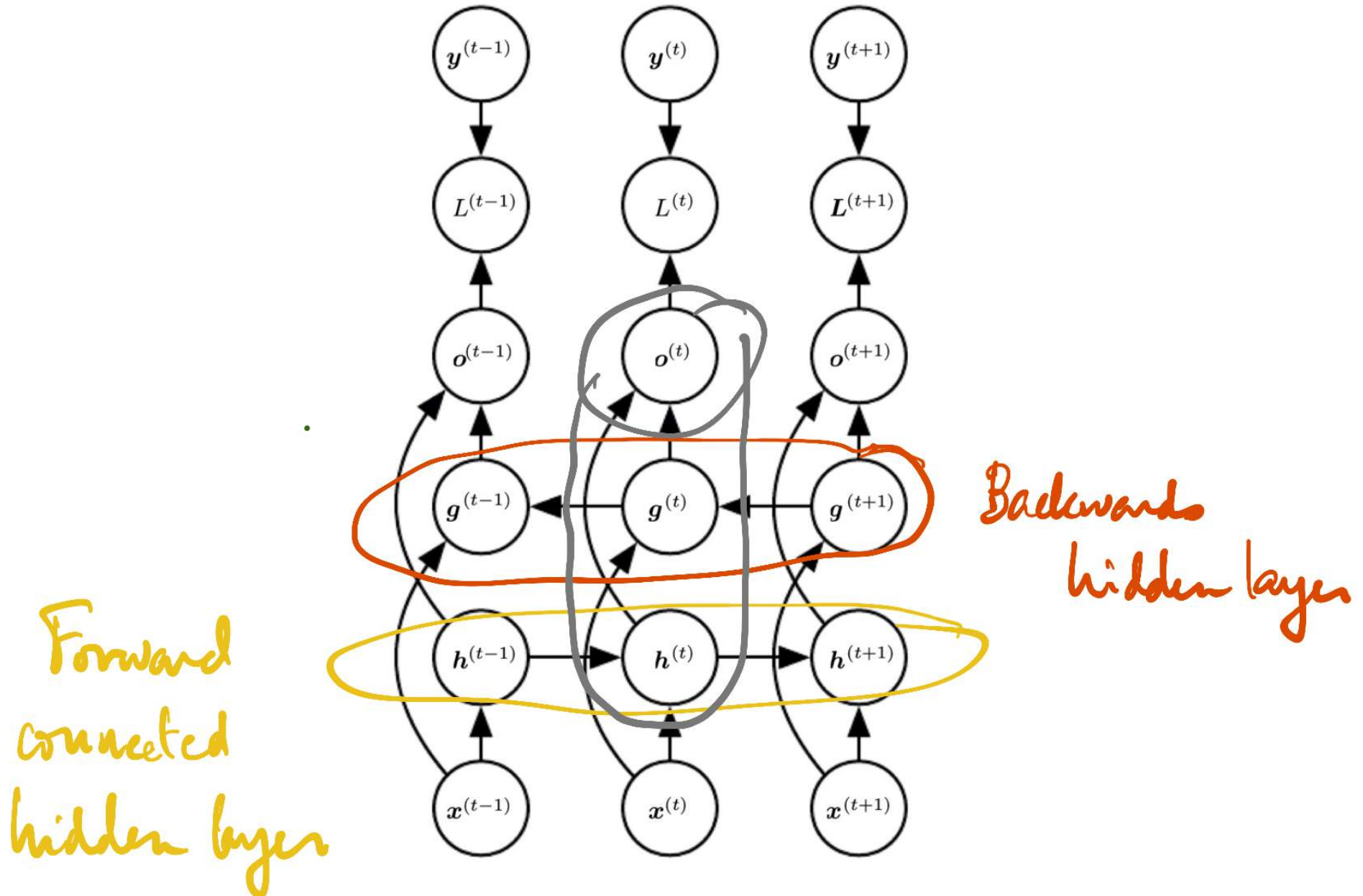
Explicit influence of  
 $y^{(t)}$  on  $y^{(t+1)}$





Sometimes we need to see the future

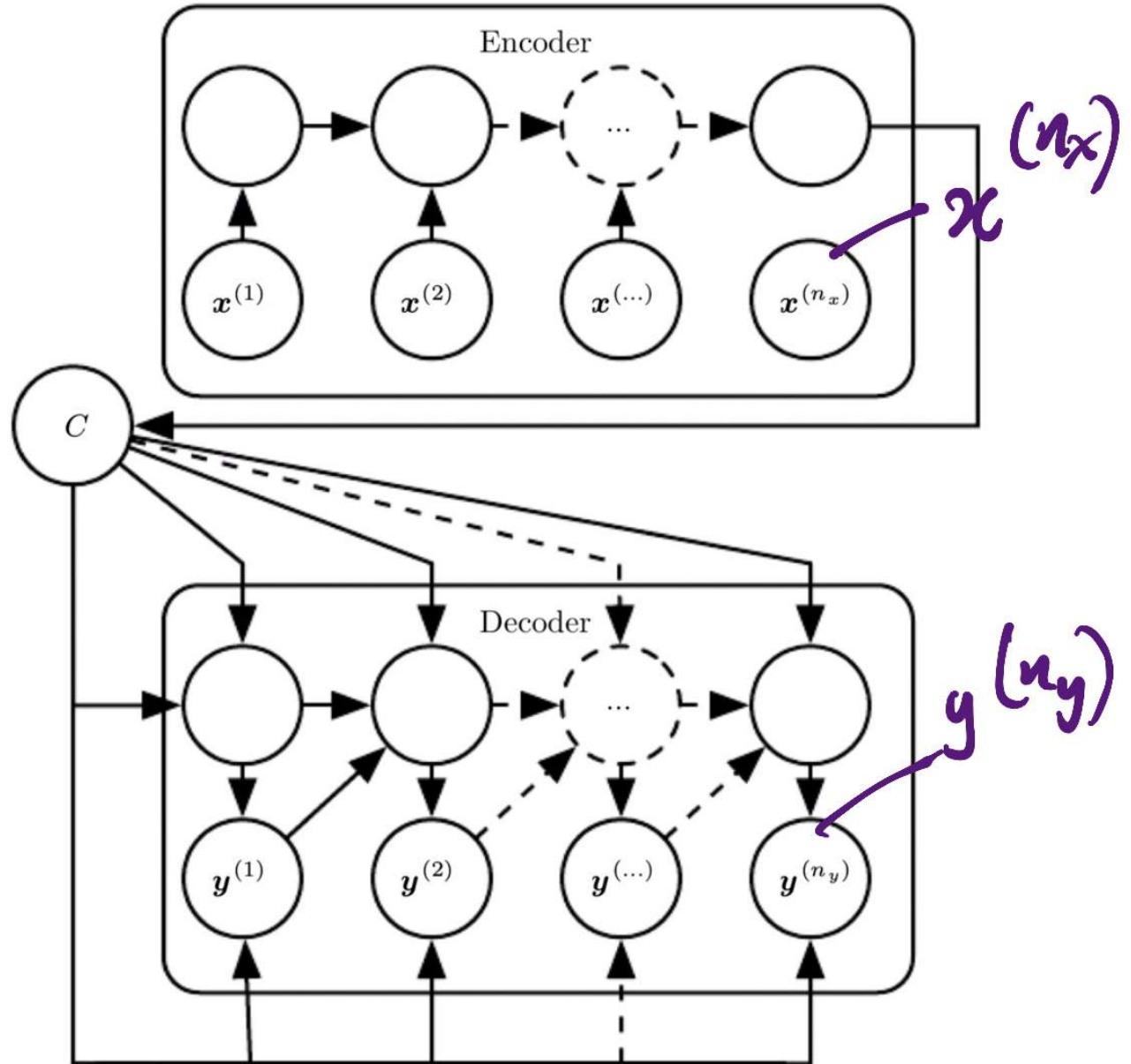
## Bidirectional RNNs



What if output sequence length is different from input sequence length?

Separate input from output

Encoder-Decoder



The problem of "long term" memory

Gradients vanish (or explode)

Hard to model strong dependency on distant past

Increasing nonlinearity as time progresses

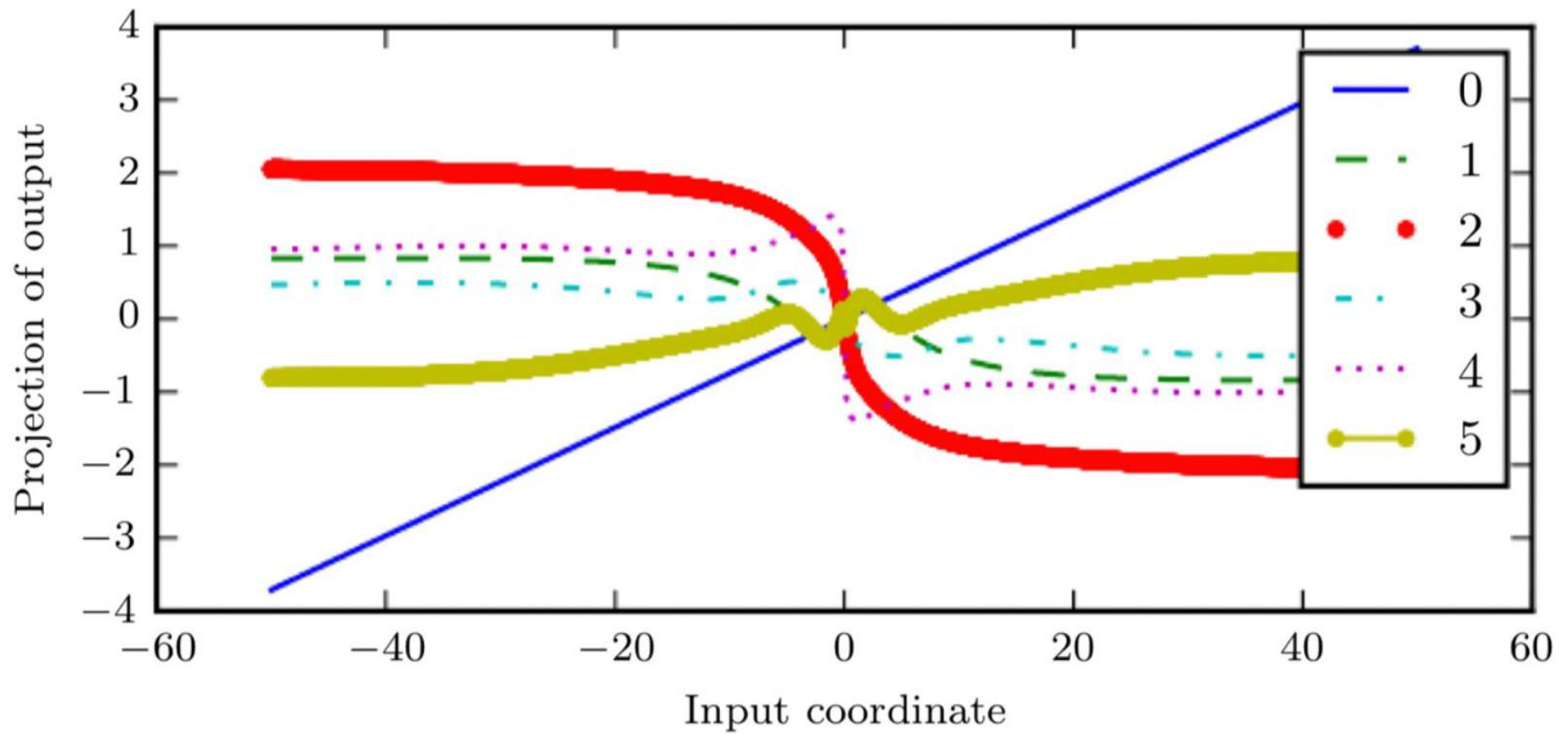
$$h^{(t)} = W^T h^{(t-1)}$$

decompose as  $Q \Lambda Q^T$

orthogonal

eigen decomposition

eigenvalue diagonal matrix



Fix tanh activation & explore outputs  
as function of input after many layers



$$h^{(t+1)} = Q \Lambda Q^T h^{(t)}$$

Compose across  $n$  layers

$$h^{(n+1)} = \underbrace{(Q \Lambda Q^T)}_{\mathbb{I}} \underbrace{(Q \Lambda Q^T)}_{\mathbb{I}} \dots (Q \Lambda Q^T) h^{(1)}$$

$$= Q \Lambda^n Q^T h^{(1)}$$

↳ eigenvalues away from 1  
will vanish

How to remember "long term"?

① Vary  $W$  across time. How?

② Skip connections

Connect layer  $t$  to  $t+d$  directly

③ "Leaky" connection

Make the feedback linear

## Moving average

$$h^{(t)} = \alpha h^{(t-1)} + (1-\alpha) i^{(t)}$$

Vary  $\alpha$ :

$\alpha \approx 1$  : More influence of past

$\alpha \rightarrow 0$  : More influence of present

## New architecture using "gates"

1. Is the input still valid?

2. Is the past of  $h$  still valid?

3. Should I continue to emit output?

These decisions should be time varying



# LSTM

Long short term memory

