# *Deep Learning: Theory and Practice*

**Recurrent Neural Networks**

**28-03-2019**

# Introduction

- ❖ The standard DNN/CNN paradigms

  - ❖ (x,y) - ordered pair of data vectors/images (x) and target (y)

- ❖ Moving to sequence data

  - ❖ (x(t),y(t)) where this could be sequence to sequence mapping task.

  - ❖ (x(t),y) where this could be a sequence to vector mapping task.

# Introduction

- Difference between CNNs/DNNs

  - $(x(t), y(t))$ where this could be sequence to sequence mapping task.

    - Input features / output targets are correlated in time.

    - Unlike standard models where each pair is independent.

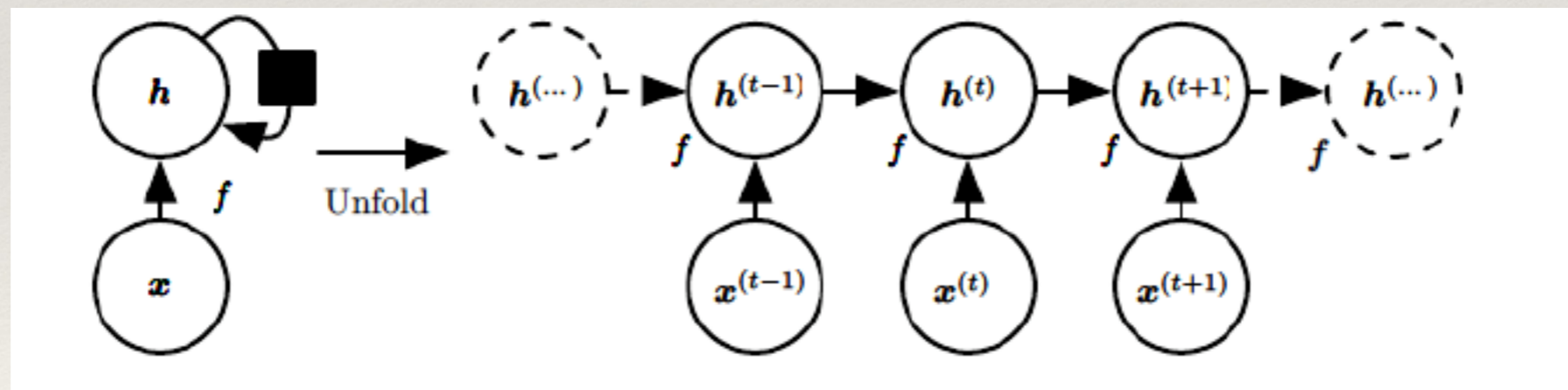    - Need to model dependencies in the sequence over time.

# Introduction to Recurrent Networks

$$s^{(t)} = f(s^{(t-1)}; \boldsymbol{\theta}),$$

$$
\begin{aligned}
s^{(3)} &= f(s^{(2)}; \boldsymbol{\theta}) \\
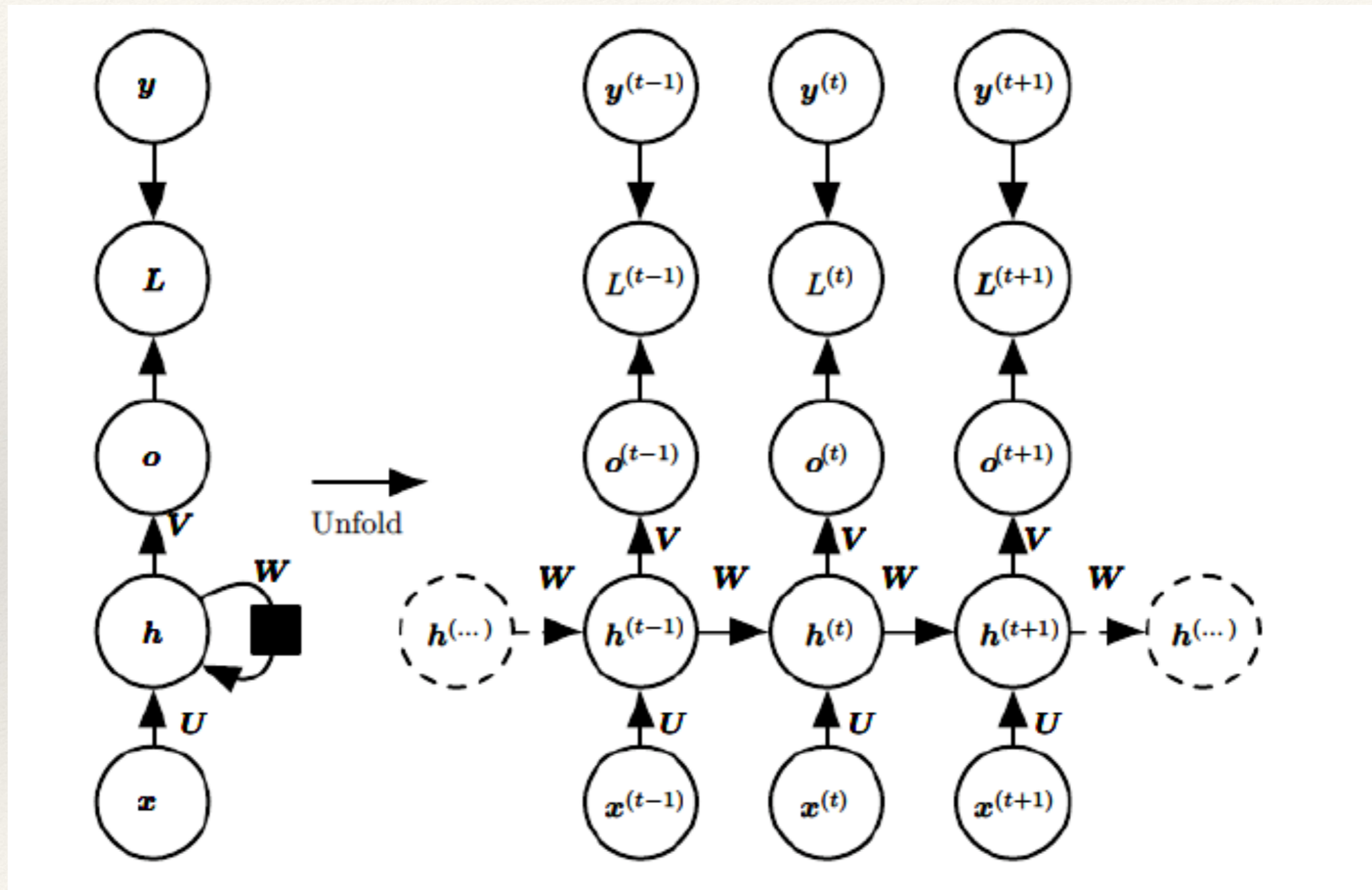&= f(f(s^{(1)}; \boldsymbol{\theta}); \boldsymbol{\theta})
\end{aligned}
$$

$$s^{(t)} = f(s^{(t-1)}, x^{(t)}; \boldsymbol{\theta}),$$

$$h^{(t)} = f(h^{(t-1)}, x^{(t)}; \boldsymbol{\theta}),$$



*"Deep Learning", Ian Goodfellow, Yoshua Bengio, Aaron Courville*

# Recurrent Networks



*"Deep Learning", Ian Goodfellow, Yoshua Bengio, Aaron Courville*

# Recurrent Networks



$$\begin{aligned}
\boldsymbol{a}^{(t)} &= \boldsymbol{b} + \boldsymbol{W}\boldsymbol{h}^{(t-1)} + \boldsymbol{U}\boldsymbol{x}^{(t)} \\
\boldsymbol{h}^{(t)} &= \tanh(\boldsymbol{a}^{(t)}) \\
\boldsymbol{o}^{(t)} &= \boldsymbol{c} + \boldsymbol{V}\boldsymbol{h}^{(t)} \\
\hat{\boldsymbol{y}}^{(t)} &= \mathrm{softmax}(\boldsymbol{o}^{(t)})
\end{aligned}$$

$$\begin{aligned}
&L\left(\{\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(\tau)}\}, \{\boldsymbol{y}^{(1)}, \ldots, \boldsymbol{y}^{(\tau)}\}\right) \\
&= \sum_t L^{(t)} \\
&= -\sum_t \log p_{\mathrm{model}}\left(y^{(t)} \mid \{\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(t)}\}\right)
\end{aligned}$$

*"Deep Learning", Ian Goodfellow, Yoshua Bengio, Aaron Courville*

# Back Propagation in RNNs

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

**Model Parameters**

$$U, V, W, b \text{ and } c$$

**Gradient Descent**

$$L\left(\{x^{(1)}, \ldots, x^{(\tau)}\}, \{y^{(1)}, \ldots, y^{(\tau)}\}\right)$$
$$= \sum_t L^{(t)}$$
$$= -\sum_t \log p_{\text{model}}\left(y^{(t)} \mid \{x^{(1)}, \ldots, x^{(t)}\}\right)$$

$$\frac{\partial L}{\partial L^{(t)}} = 1.$$

$$(\nabla_{o^{(t)}} L)_i = \frac{\partial L}{\partial o_i^{(t)}} = \frac{\partial L}{\partial L^{(t)}} \frac{\partial L^{(t)}}{\partial o_i^{(t)}} = \hat{y}_i^{(t)} - \mathbf{1}_{i,y^{(t)}}$$
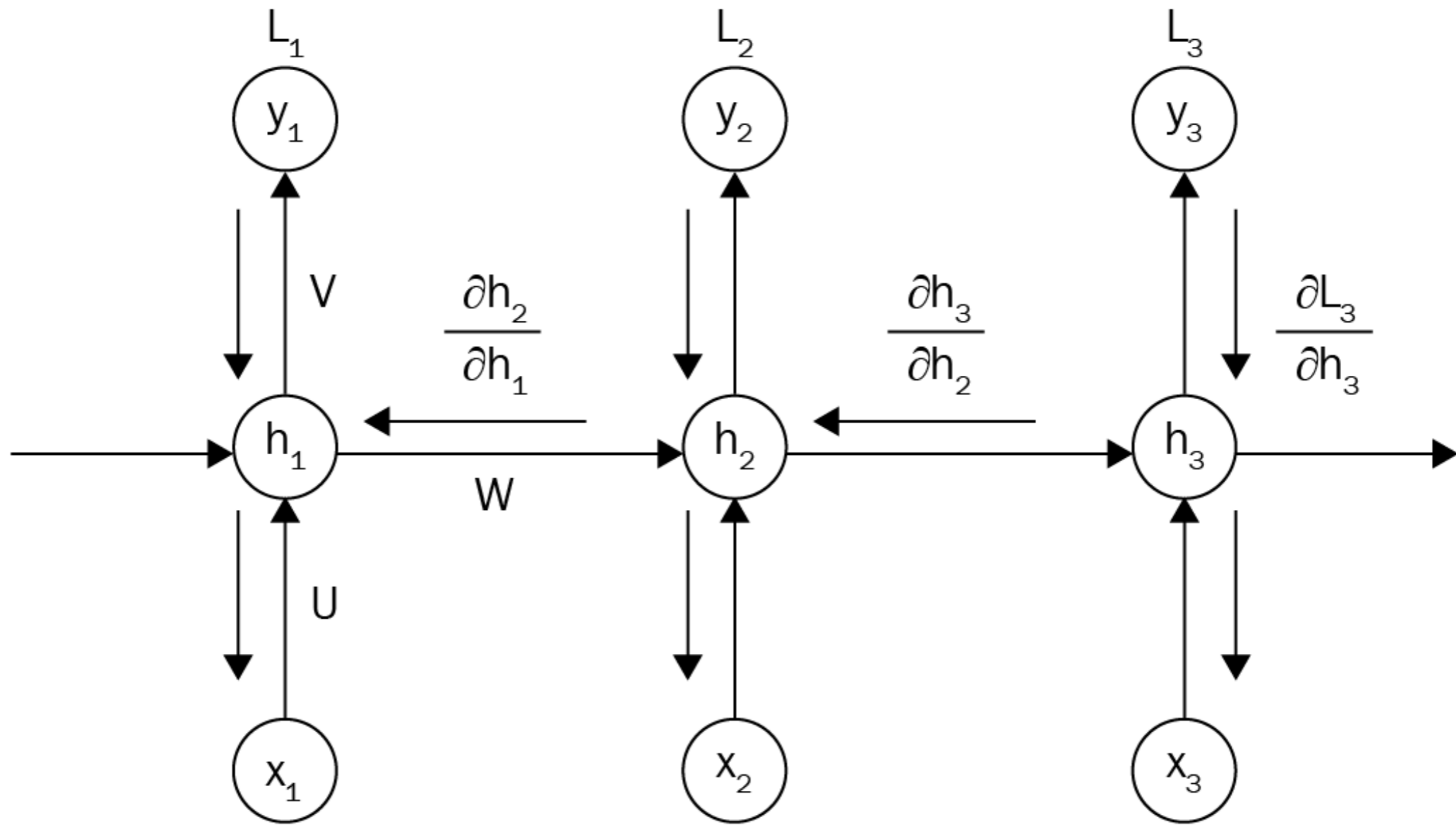
# Recurrent Networks



$$(\nabla_{\boldsymbol{o}^{(t)}} L)_i = \frac{\partial L}{\partial o_i^{(t)}} = \frac{\partial L}{\partial L^{(t)}} \frac{\partial L^{(t)}}{\partial o_i^{(t)}} = \hat{y}_i^{(t)} - \mathbf{1}_{i, y^{(t)}}$$

$$\nabla_{\boldsymbol{h}^{(\tau)}} L = \boldsymbol{V}^{\top} \nabla_{\boldsymbol{o}^{(\tau)}} L.$$

$$\nabla_{\boldsymbol{h}^{(t)}} L = \left( \frac{\partial \boldsymbol{h}^{(t+1)}}{\partial \boldsymbol{h}^{(t)}} \right)^{\top} (\nabla_{\boldsymbol{h}^{(t+1)}} L) + \left( \frac{\partial \boldsymbol{o}^{(t)}}{\partial \boldsymbol{h}^{(t)}} \right)^{\top} (\nabla_{\boldsymbol{o}^{(t)}} L)$$

$$= \boldsymbol{W}^{\top} (\nabla_{\boldsymbol{h}^{(t+1)}} L) \operatorname{diag} \left( 1 - \left( \boldsymbol{h}^{(t+1)} \right)^2 \right) + \boldsymbol{V}^{\top} (\nabla_{\boldsymbol{o}^{(t)}} L)$$
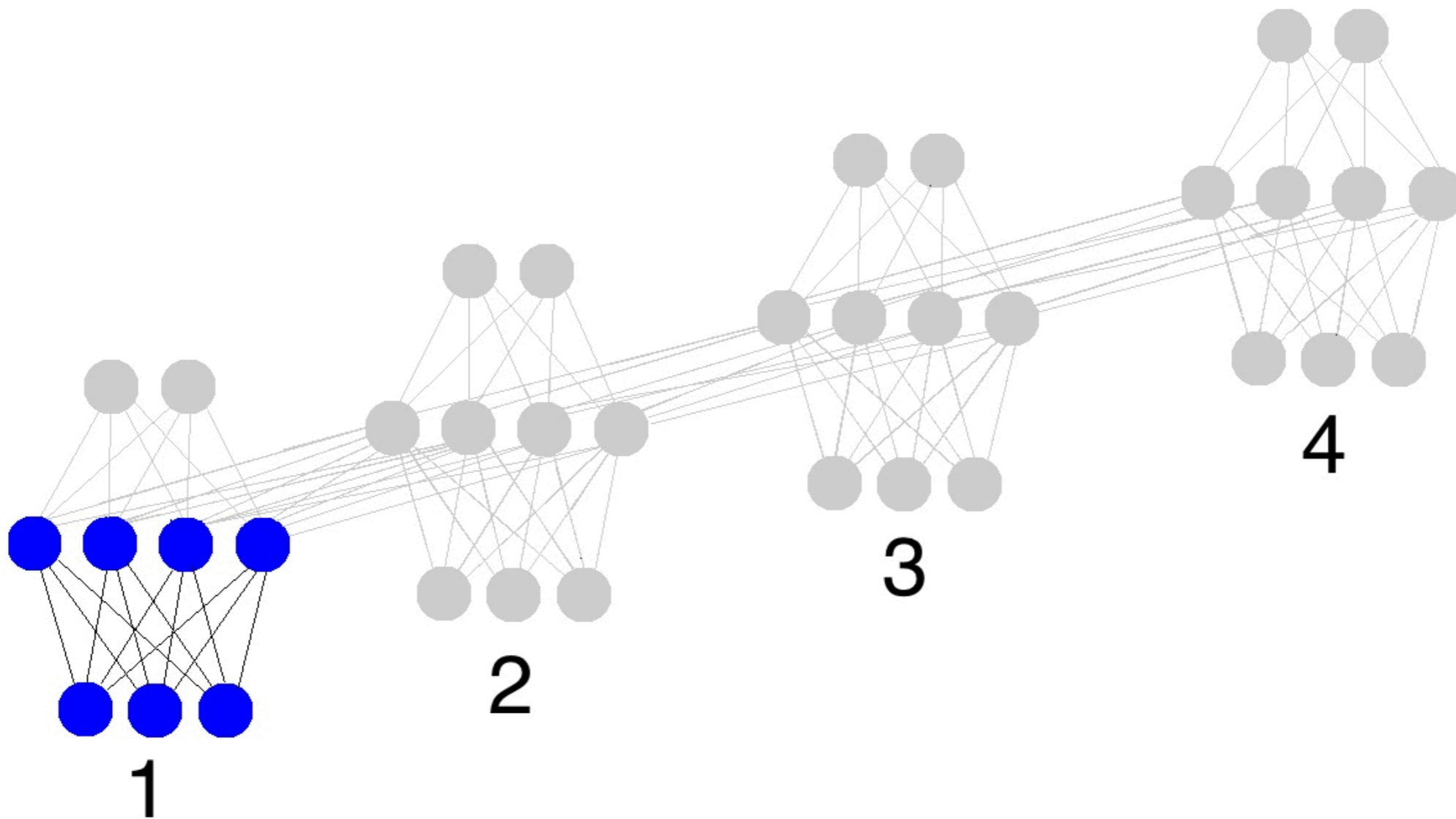
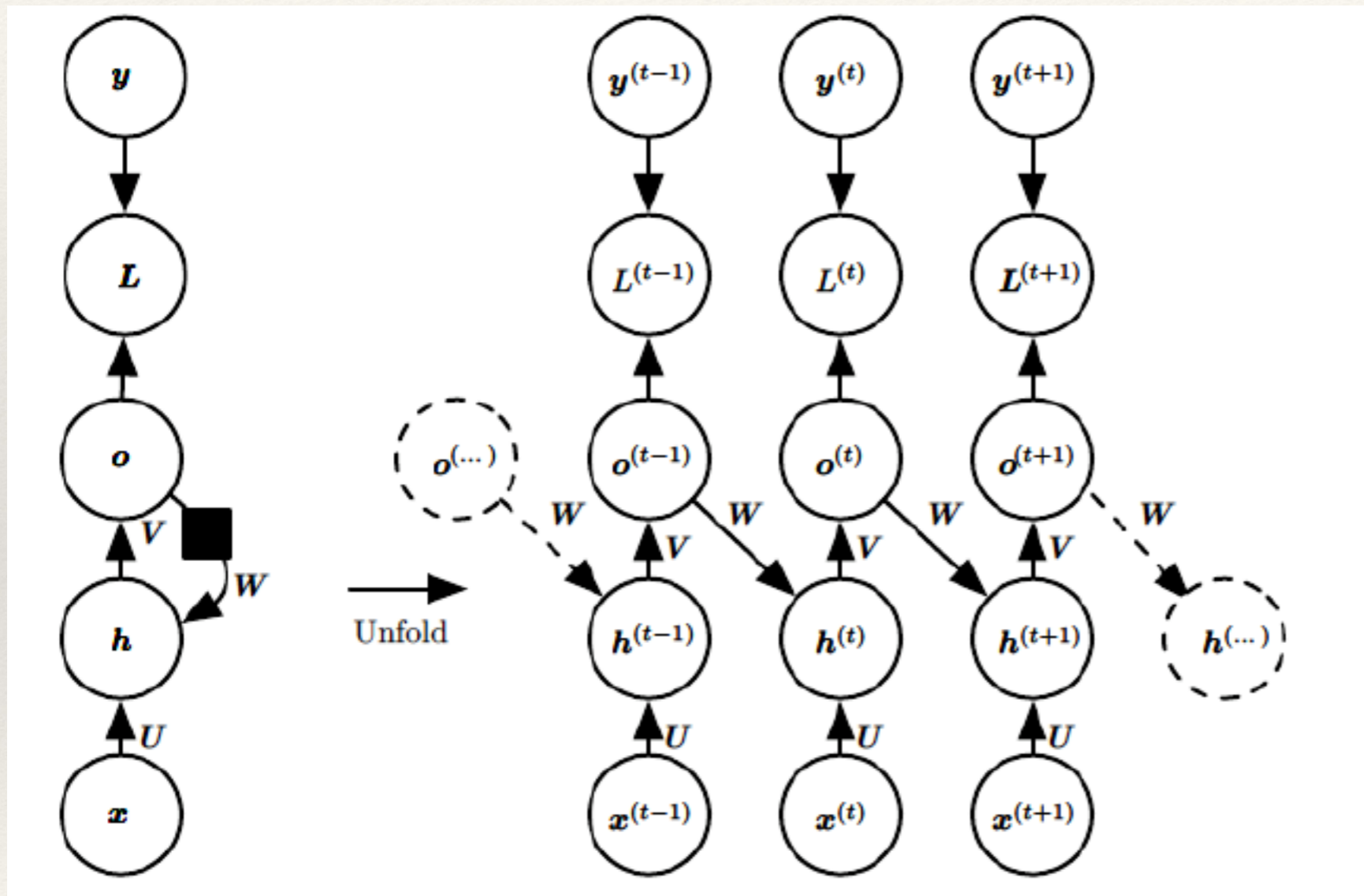# Back Propagation Through Time

# Back Propagation Through Time

# Standard Recurrent Networks



*"Deep Learning", Ian Goodfellow, Yoshua Bengio, Aaron Courville*

# Other Recurrent Networks



**Teacher Forcing Networks**

*"Deep Learning", Ian Goodfellow, Yoshua Bengio, Aaron Courville*
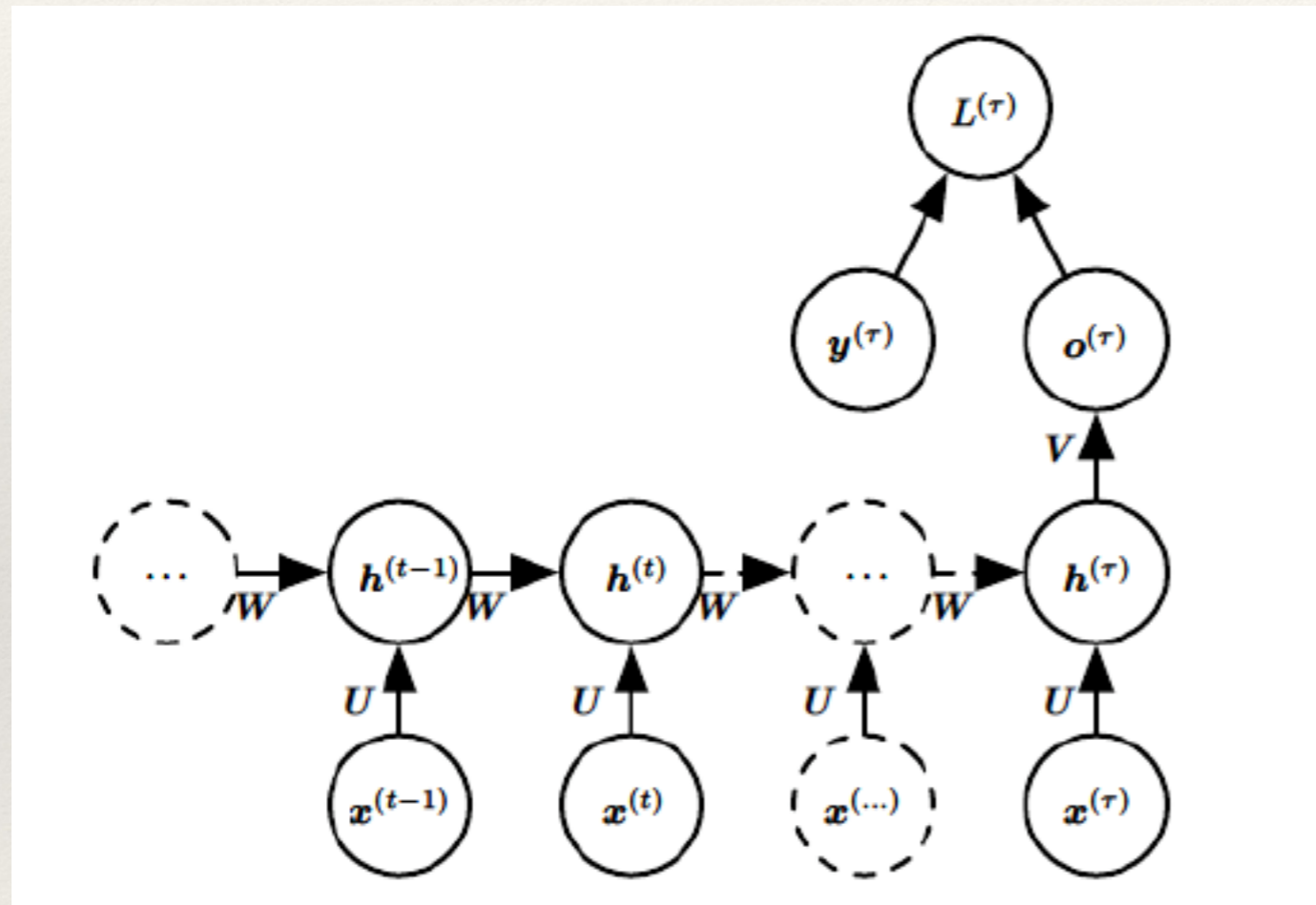
# Recurrent Networks
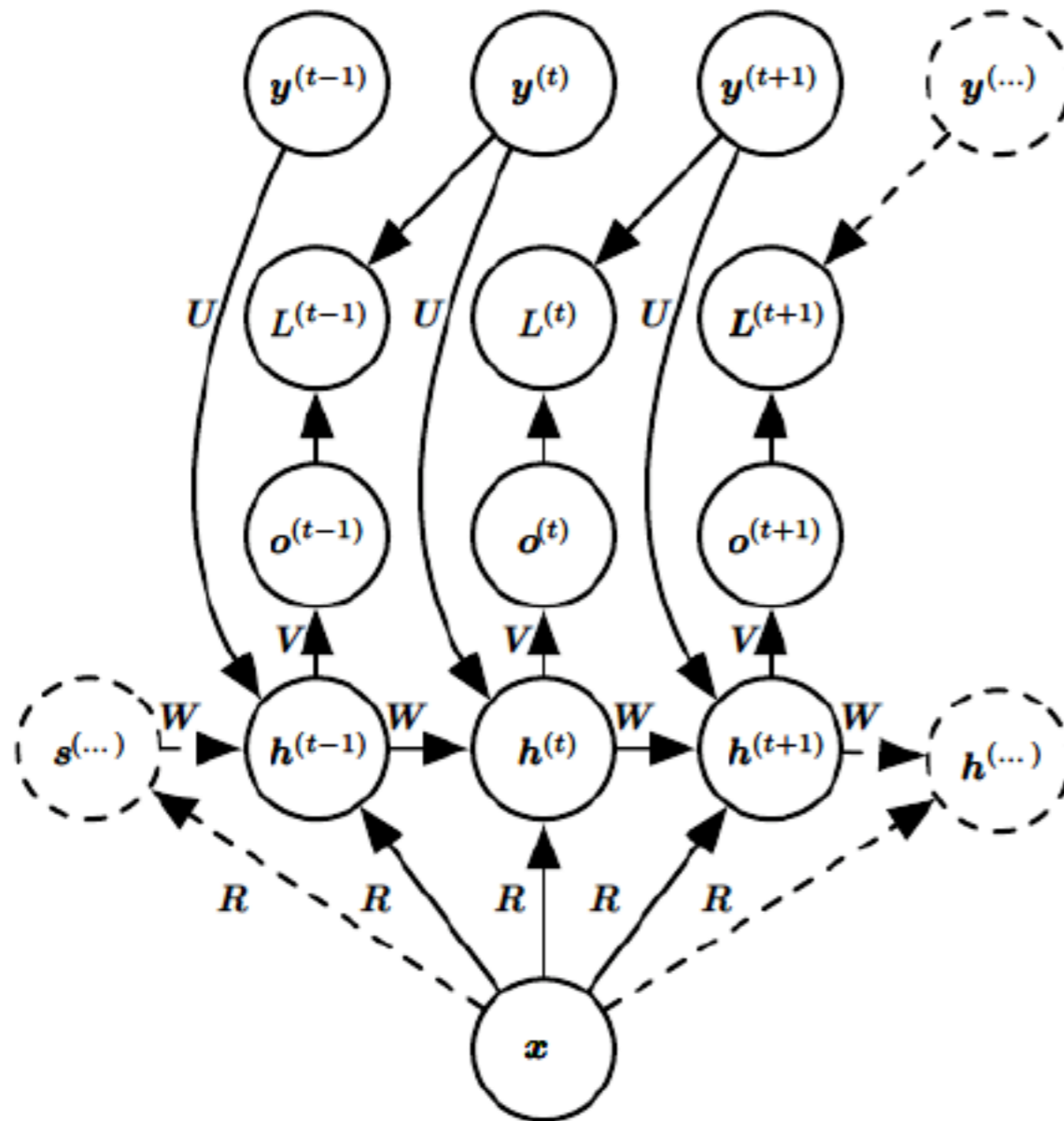


**Teacher Forcing Networks**

# Recurrent Networks
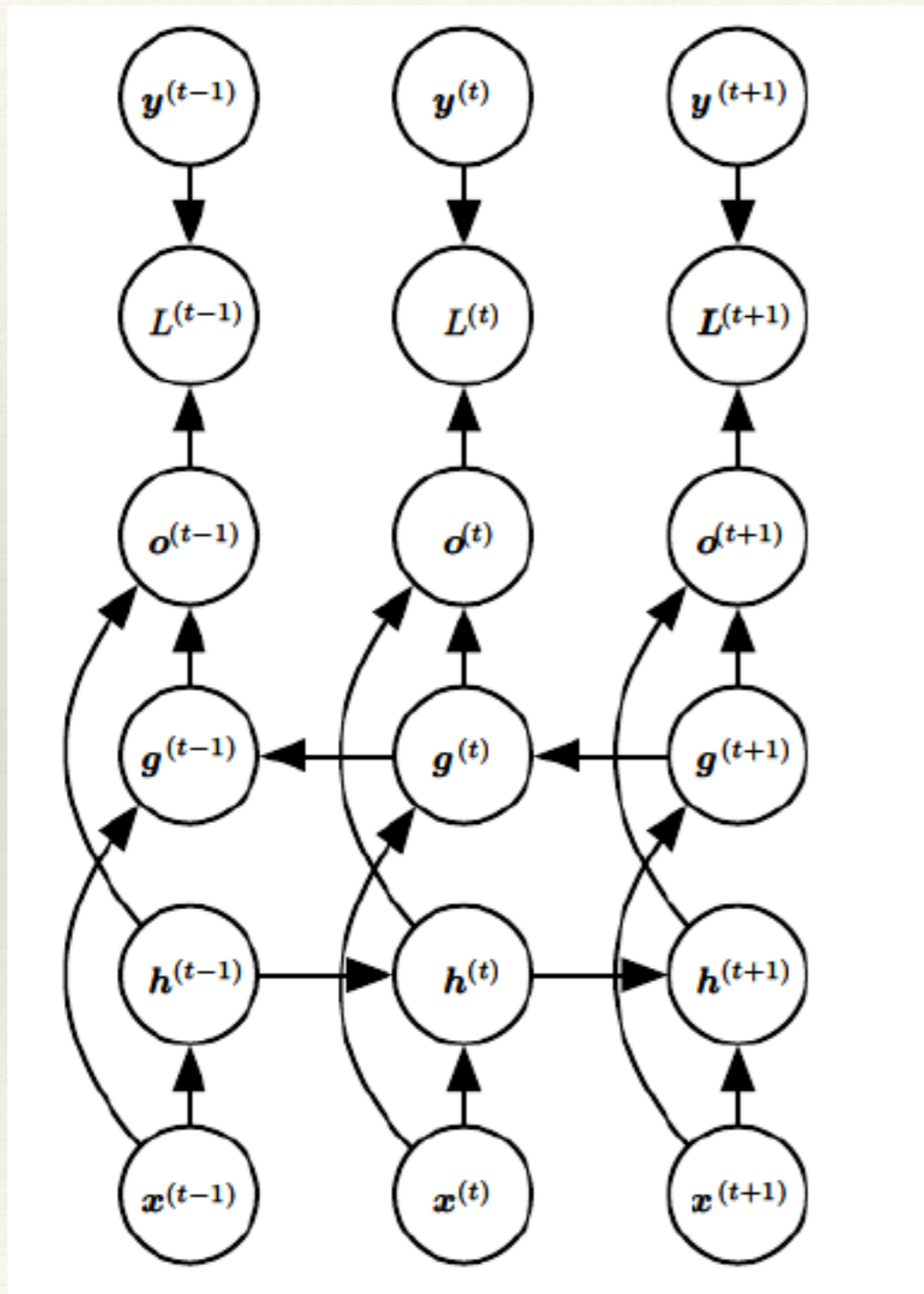


**Multiple Input
Single Output**

# Recurrent Networks
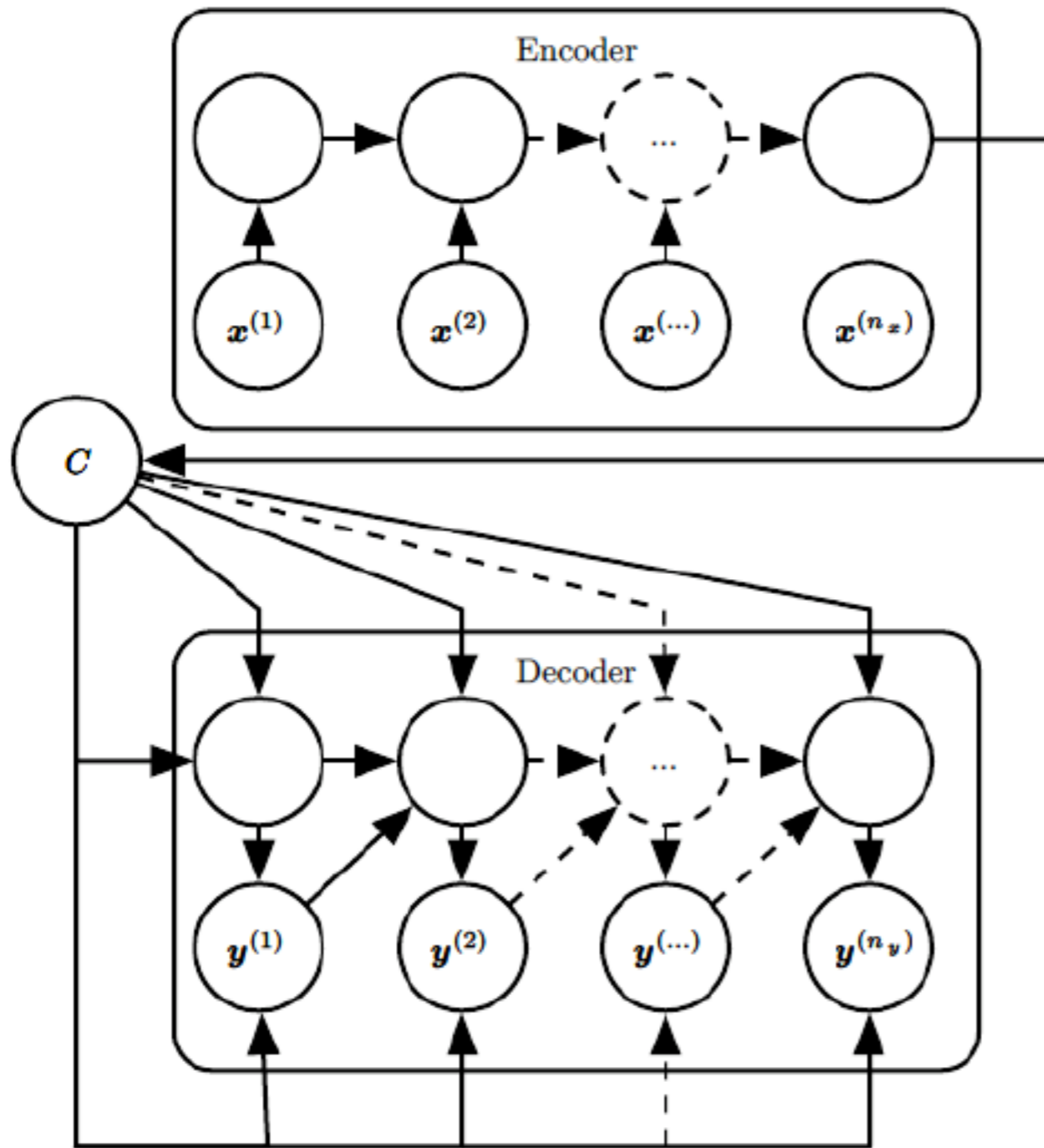


**Single Input Multiple Output**
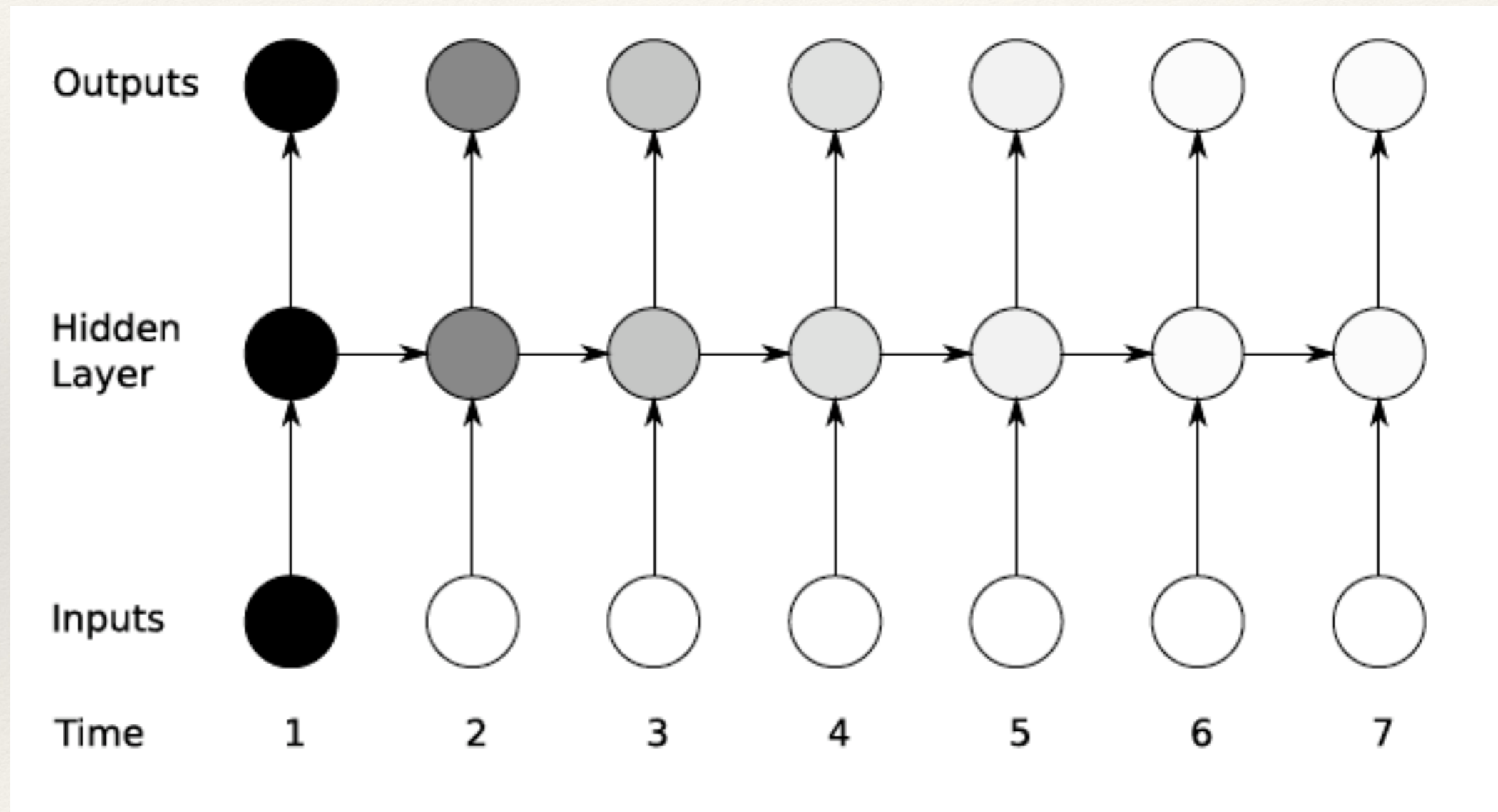
# Recurrent Networks



**Bi-directional Networks**

# Recurrent Networks



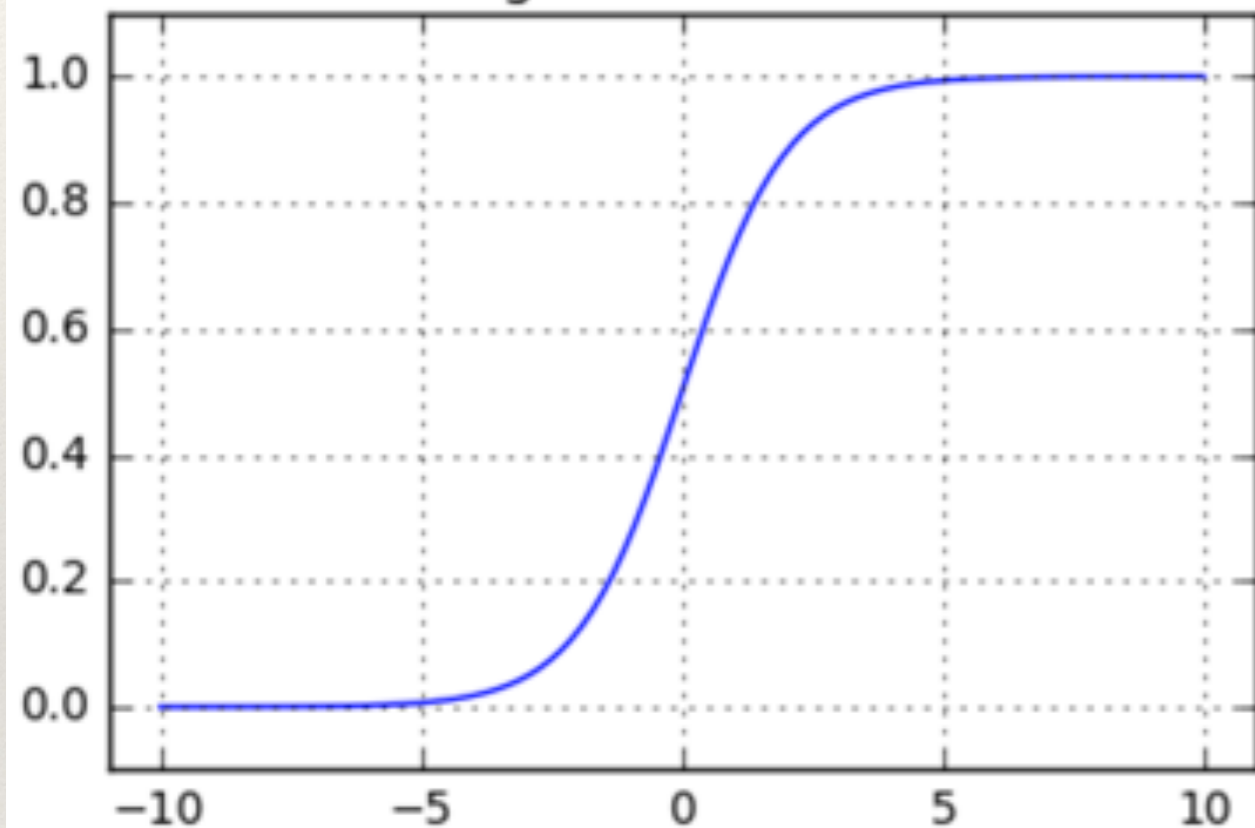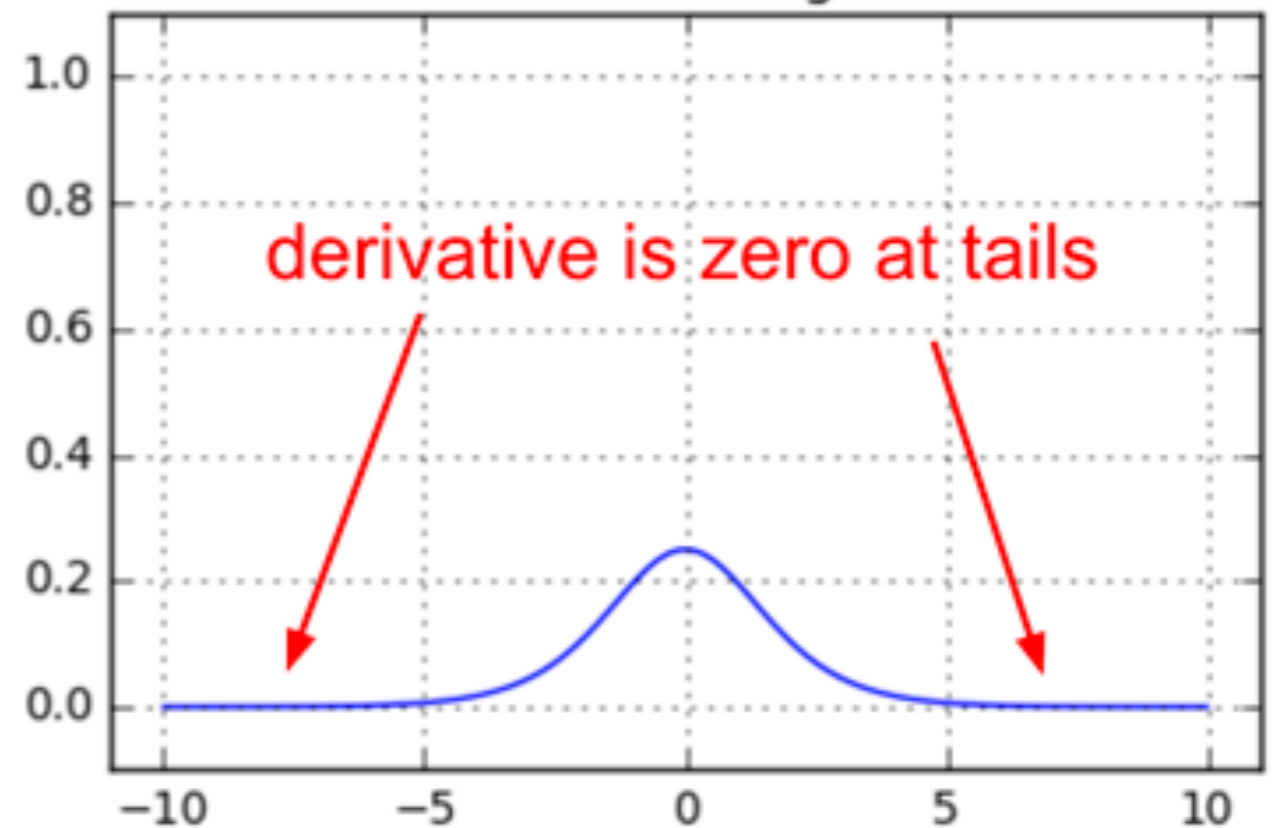**Sequence to Sequence Mapping Networks**

# Long-term Dependency Issues

# Vanishing/Exploding Gradients
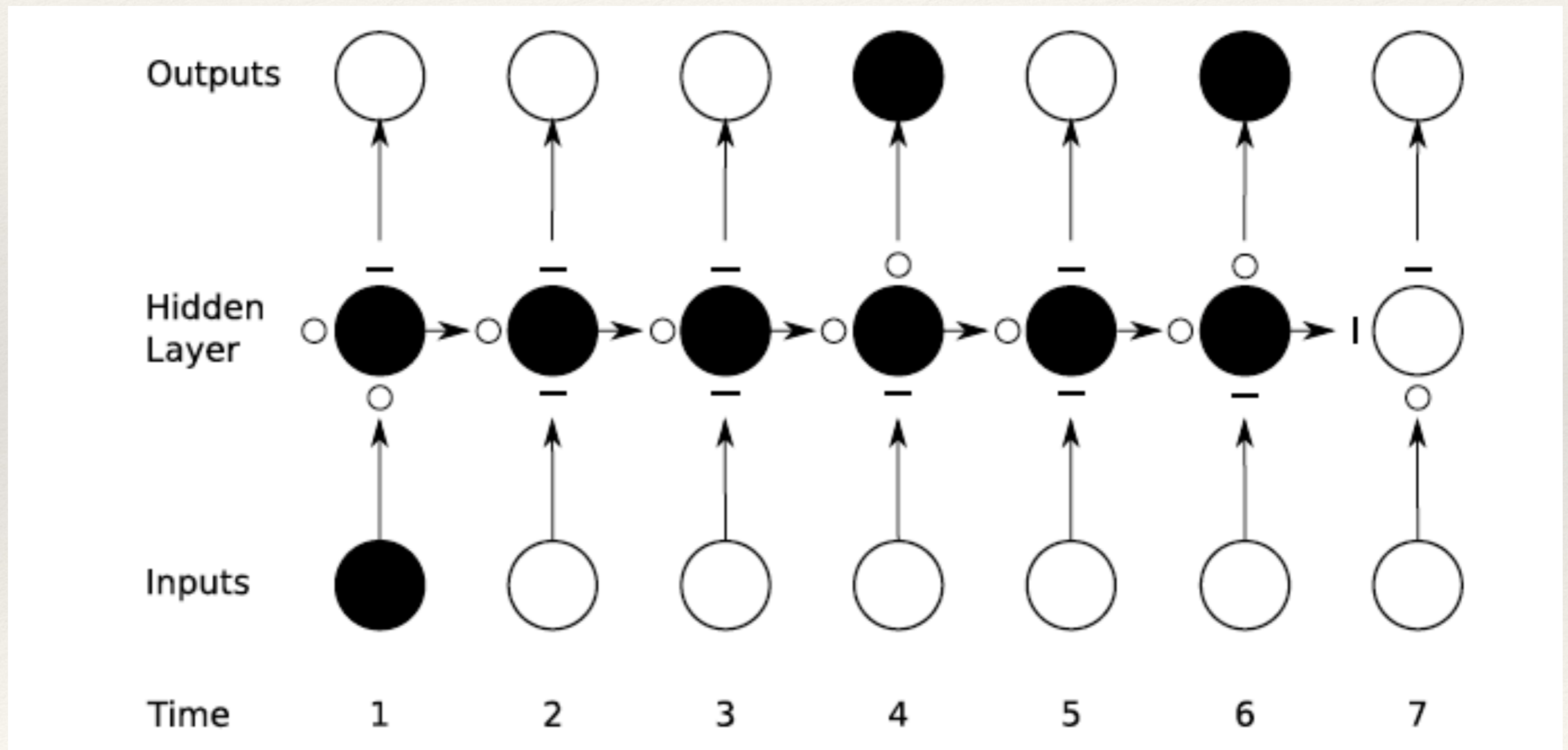


- ❖ Gradients either vanish or explode

  - ❖ Initial frames may not contribute to gradient computations or may contribute too much.

# Long-Short Term Memory

# LSTM Cell

**f - sigmoid function**
**g, h - tanh function**

## Input Gate

$$a_\iota^t = \sum_{i=1}^{I} w_{i\iota} x_i^t + \sum_{h=1}^{H} w_{h\iota} b_h^{t-1} + \sum_{c=1}^{C} w_{c\iota} s_c^{t-1}$$

$$b_\iota^t = f(a_\iota^t)$$

## Forget Gate

$$a_\phi^t = \sum_{i=1}^{I} w_{i\phi} x_i^t + \sum_{h=1}^{H} w_{h\phi} b_h^{t-1} + \sum_{c=1}^{C} w_{c\phi} s_c^{t-1}$$

$$b_\phi^t = f(a_\phi^t)$$

## Cell

$$a_c^t = \sum_{i=1}^{I} w_{ic} x_i^t + \sum_{h=1}^{H} w_{hc} b_h^{t-1}$$

$$s_c^t = b_\phi^t s_c^{t-1} + b_\iota^t g(a_c^t)$$

## Output Gate
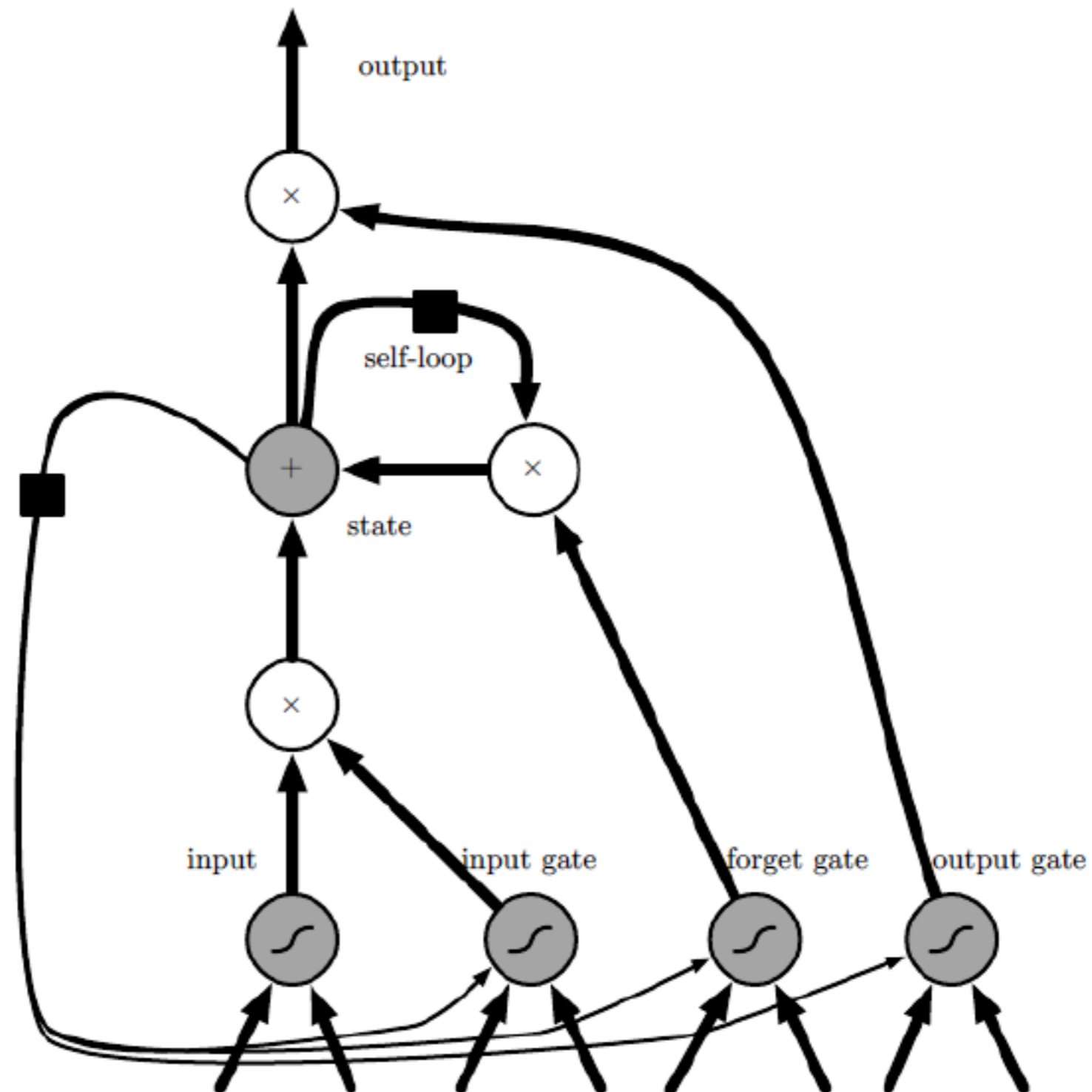
$$a_\omega^t = \sum_{i=1}^{I} w_{i\omega} x_i^t + \sum_{h=1}^{H} w_{h\omega} b_h^{t-1} + \sum_{c=1}^{C} w_{c\omega} s_c^t$$

$$b_\omega^t = f(a_\omega^t)$$
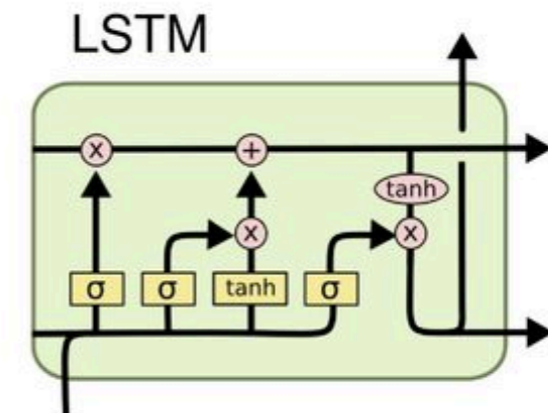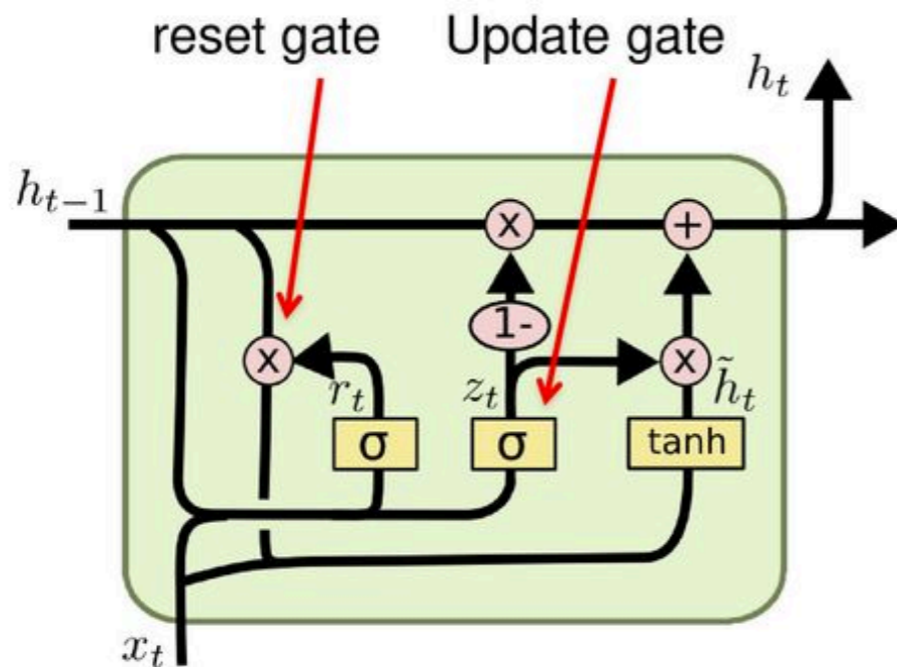
## LSTM output

$$b_c^t = b_\omega^t h(s_c^t)$$

# Long Short Term Memory Networks

# Gated Recurrent Units (GRU)

## GRU – gated recurrent unit

LSTM

(more compression)

reset gate    Update gate    $h_t$

$h_{t-1}$

$x_t$

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

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

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

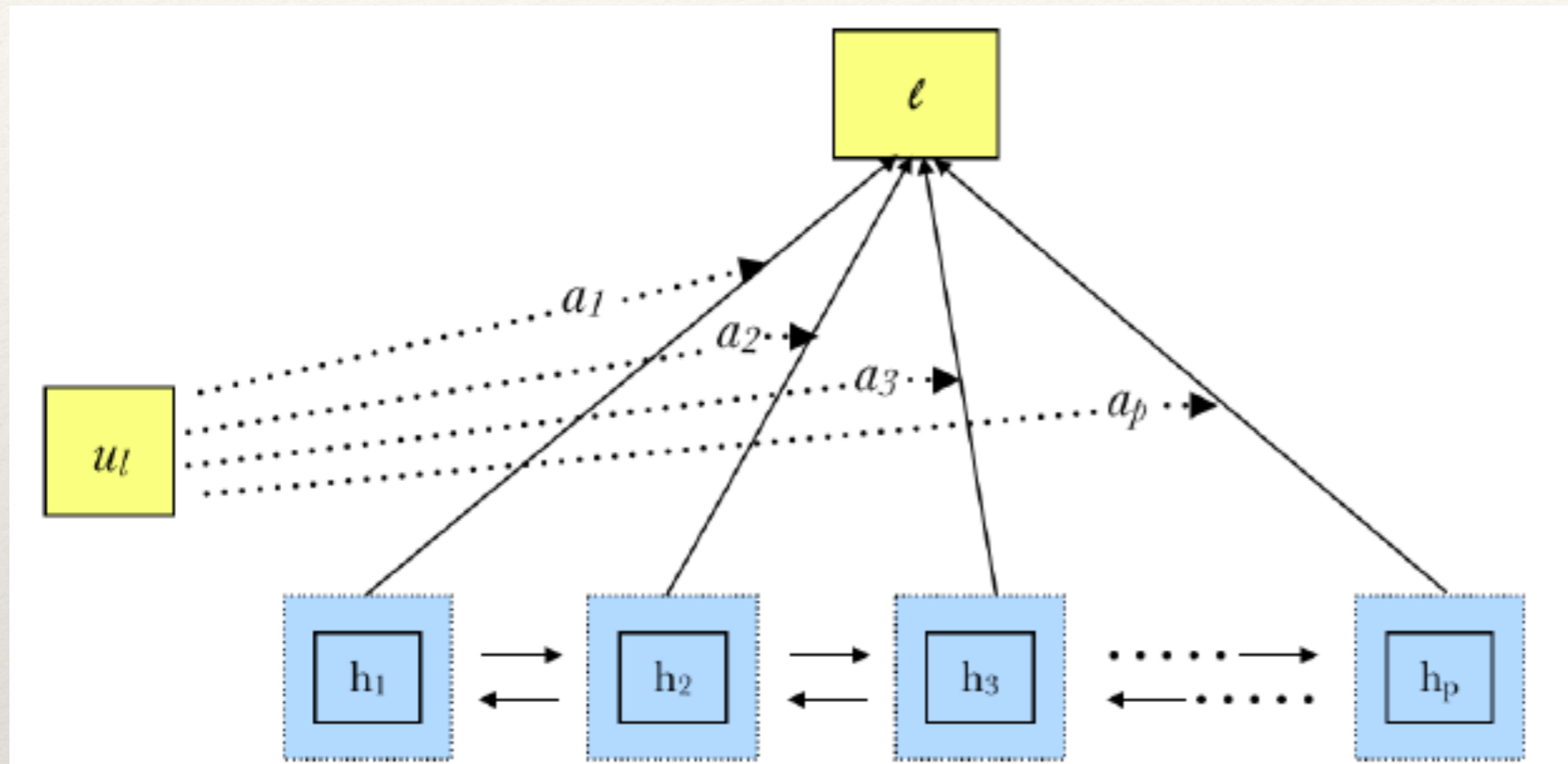$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

It combines the forget and input into a single update gate.
It also merges the cell state and hidden state. This is simpler
than LSTM. There are many other variants too.

X,*: element-wise multiply

# Attention in LSTM Networks



$$\mathbf{u}_t = tanh(\mathbf{W}_l \mathbf{h}_t + \mathbf{b}_l)$$

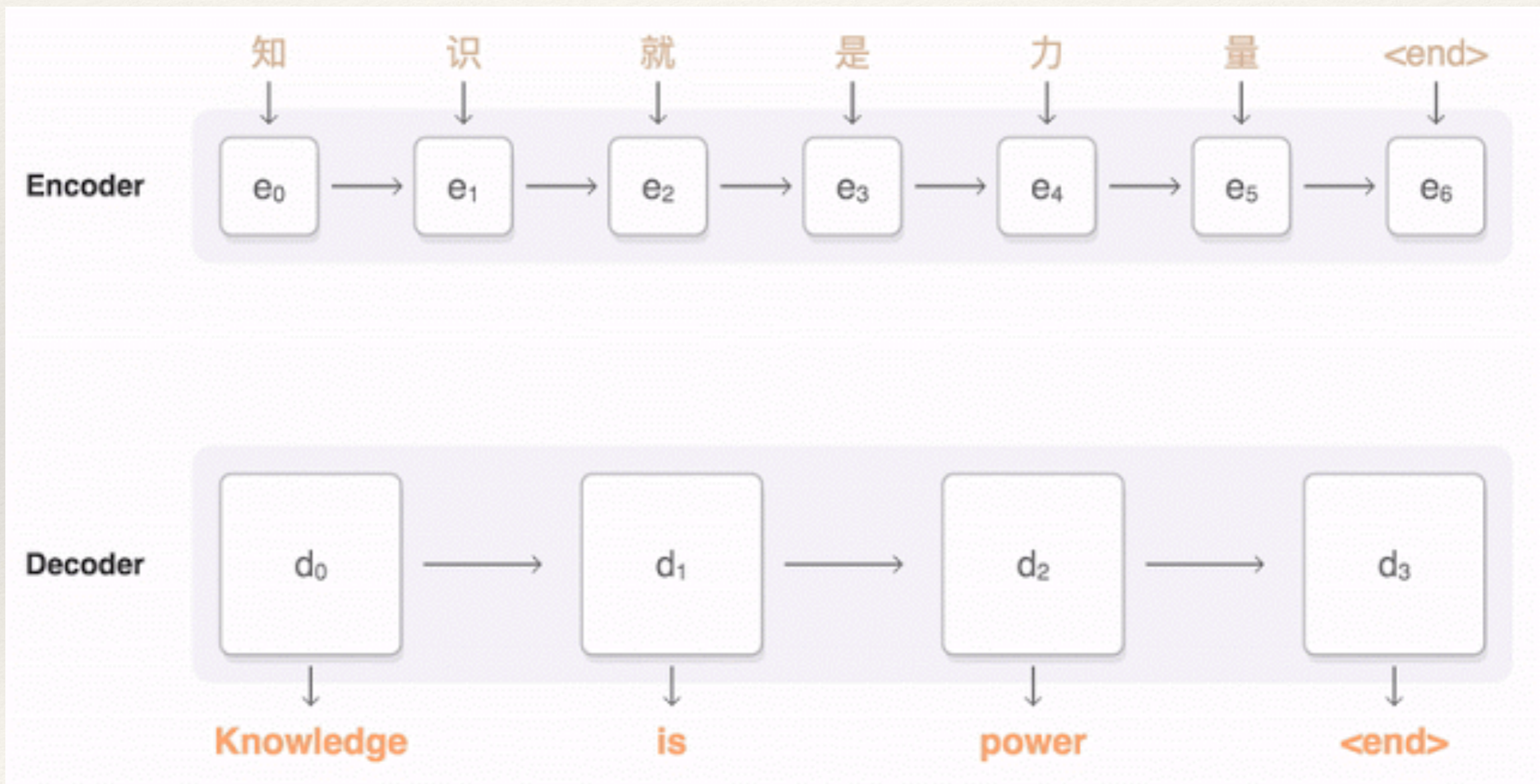$$a_t = \frac{exp(\mathbf{u}_t^T \mathbf{u}_l)}{\sum_t exp(\mathbf{u}_t^T \mathbf{u}_l)}$$

$$l = \sum_t a_t \mathbf{h}_t$$

- ❖ Attentions allows a mechanism to add relevance
  - ❖ Certain regions of the audio have more importance than the rest for the task at hand.

# Encoder – Decoder Networks with Attention

# Attention Models

# Attention - Speech Example

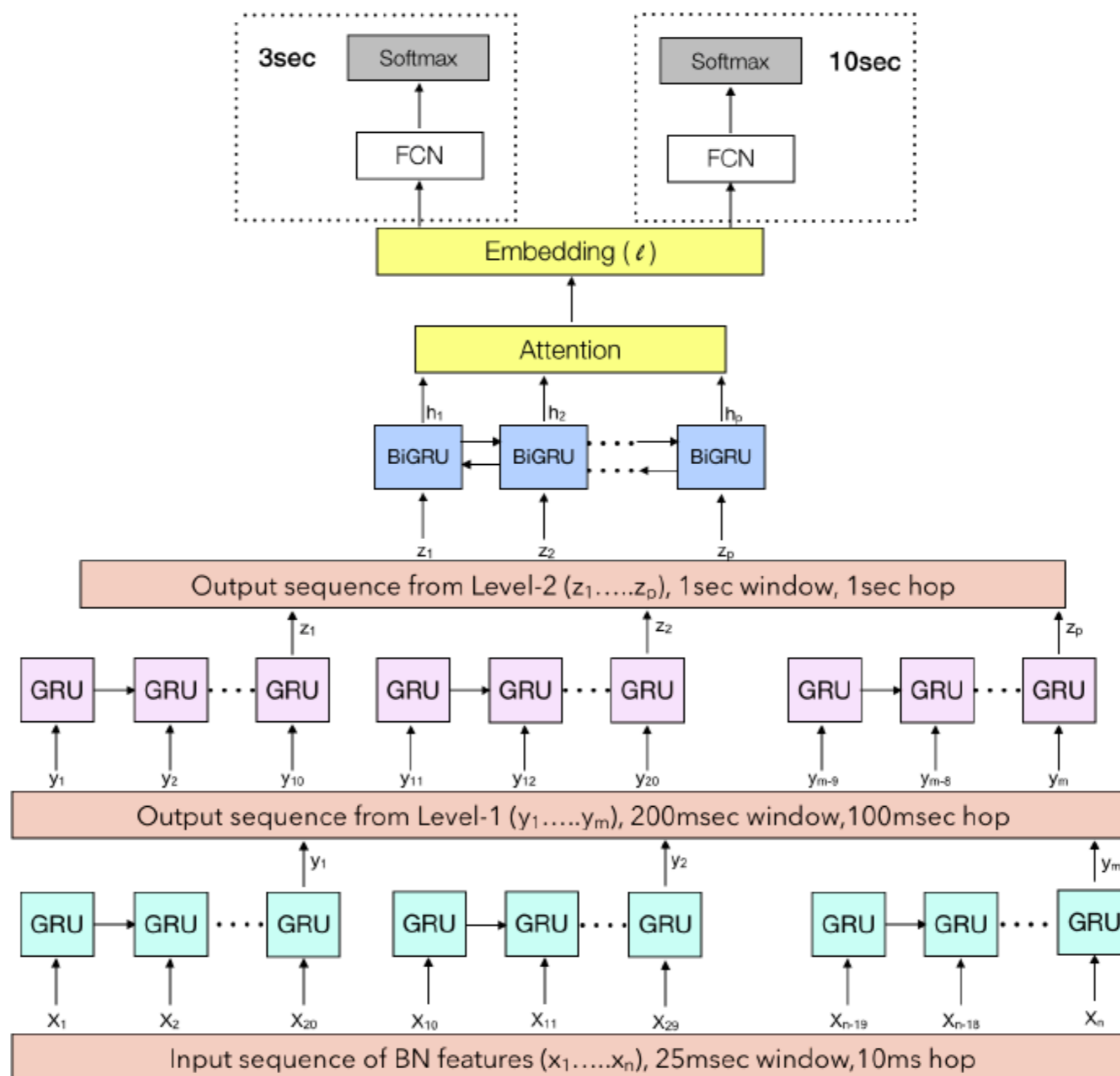From our lab [part of ICASSP 2019 paper].

# Language Recognition Evaluation

Table 1: LRE17 training set : target languages, language clusters and total number of hours.

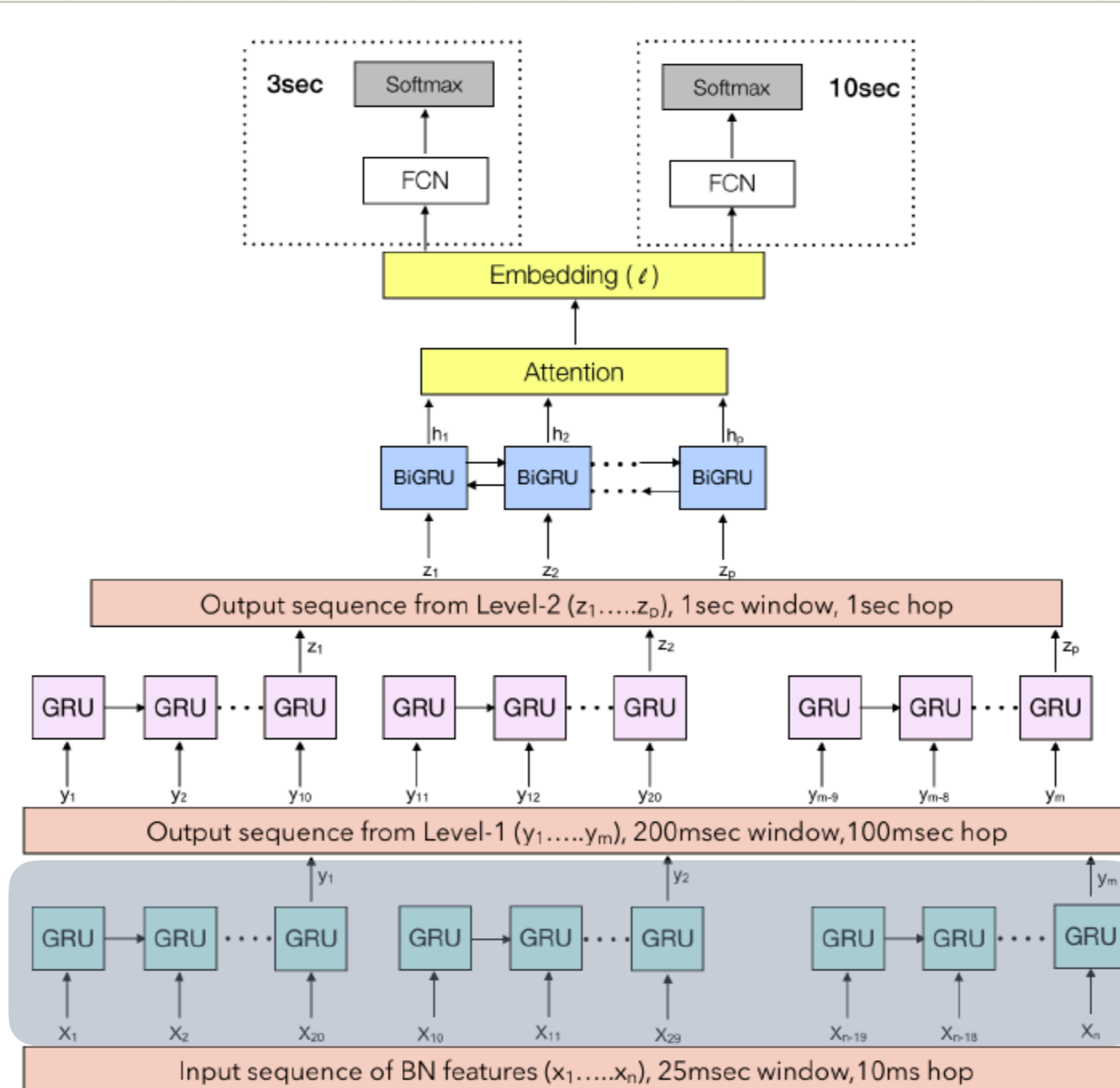| Cluster | Target Languages | Hours |
|---|---|---|
| Arabic | Egyptian Arabic (ara-arz) | 190.9 |
| | Iraqi Arabic (ara-acm) | 130.8 |
| | Levantine Arabic (ara-apc) | 440.7 |
| | Maghrebi Arabic (ara-ary) | 81.8 |
| Chinese | Mandarin (zho-cmn) | 379.4 |
| | Min Nan (zho-nan) | 13.3 |
| English | British English (eng-gbr) | 4.8 |
| | General American English (eng-usg) | 327.7 |
| Slavic | Polish (qsl-pol) | 59.3 |
| | Russian (qsl-rus) | 69.5 |
| Iberian | Caribbean Spanish (spa-car) | 166.3 |
| | European Spanish (spa-eur) | 24.7 |
| | Latin American Continental Spanish (spa-lac) | 175.9 |
| | Brazilian Portuguese (por-brz) | 4.1 |

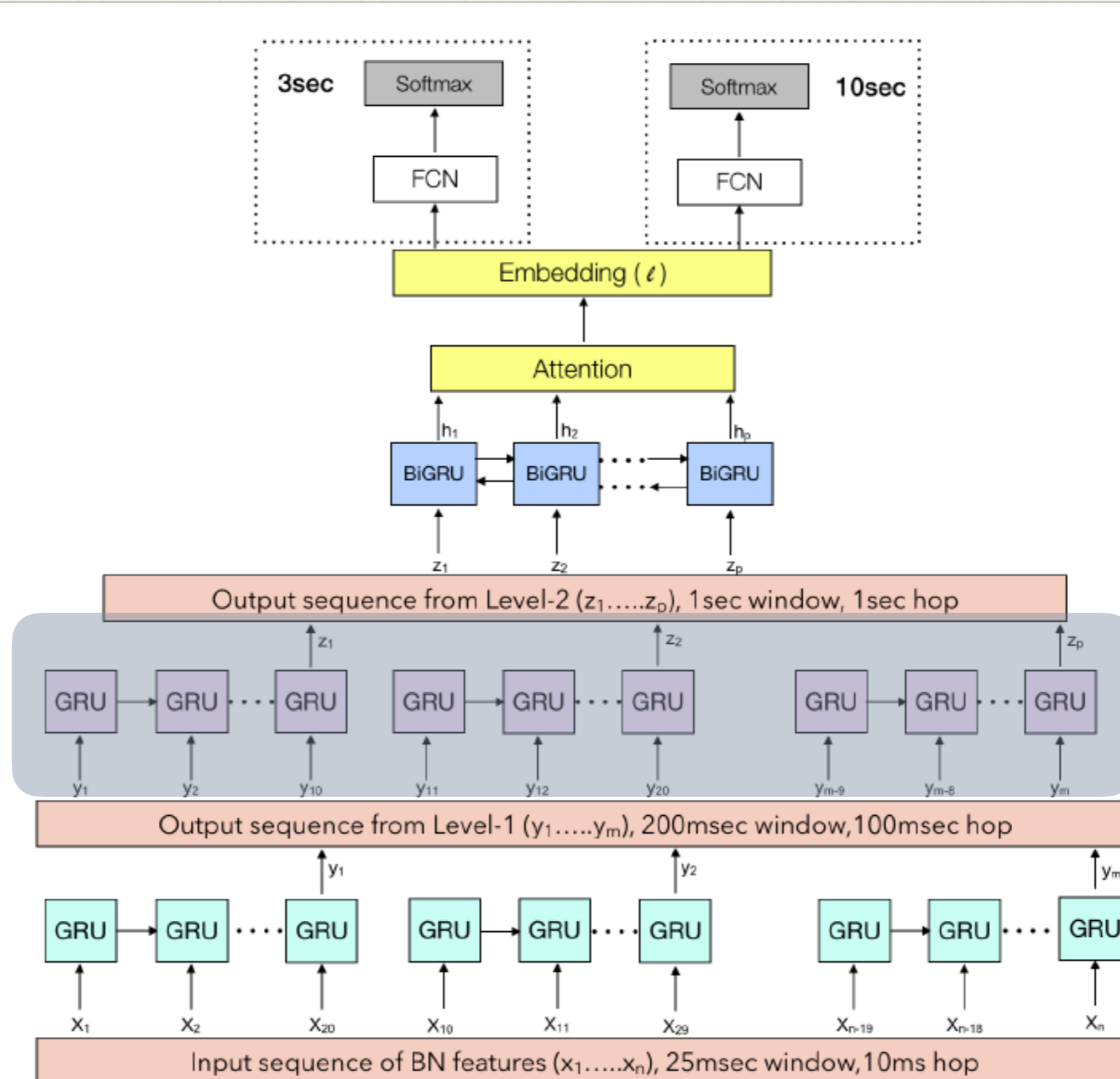# End-to-end model using GRUs and Attention

# Proposed End-to-End Language Recognition Model

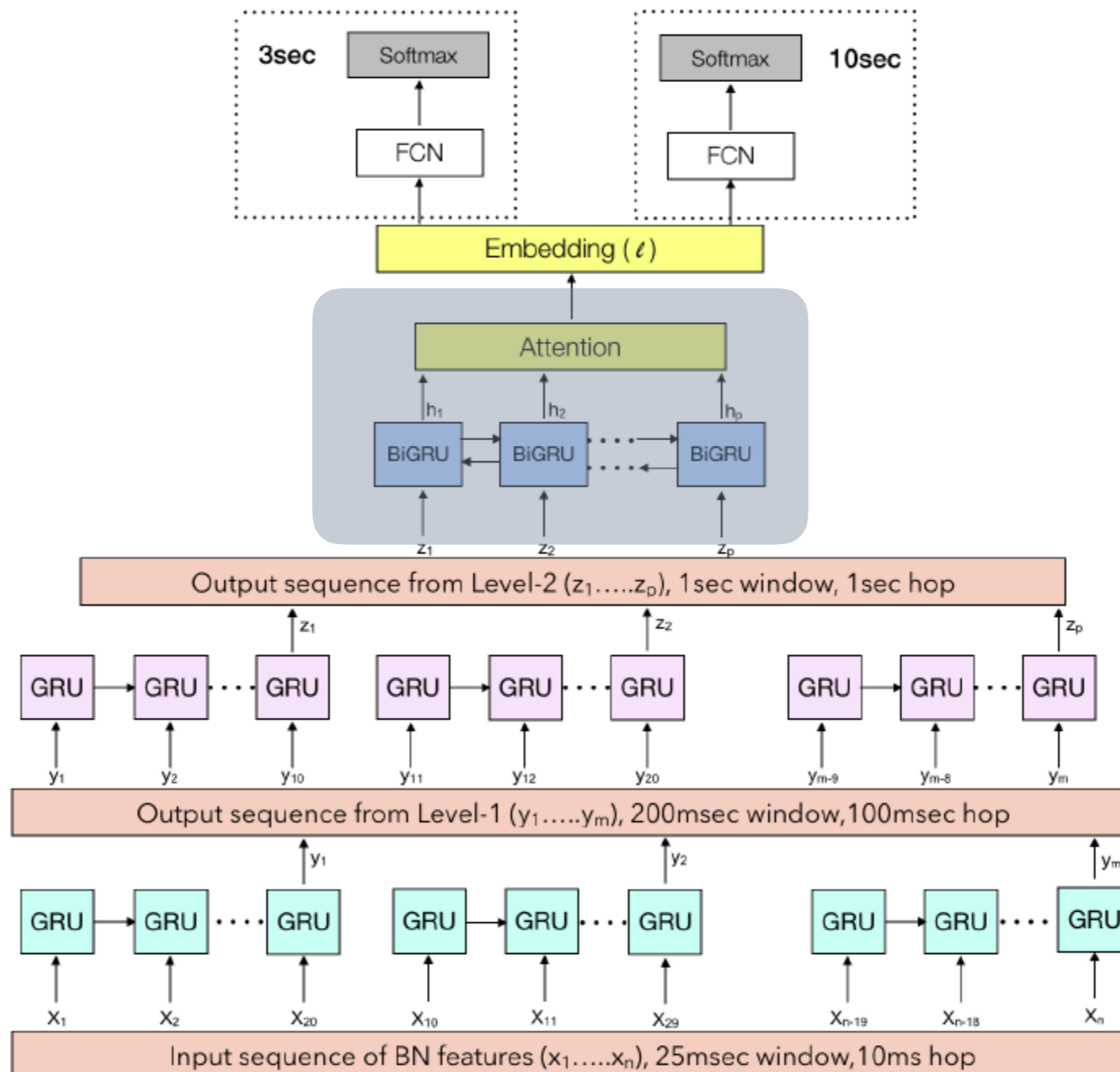# Proposed End-to-End Language Recognition Model

# Proposed End-to-End Language Recognition Model

# Language Recognition Evaluation

- State-of-art models use the input sequence directly.

- We proposed the attention model - Attention weighs the importance of each short-term segment feature for the task.

**Attention Weight**
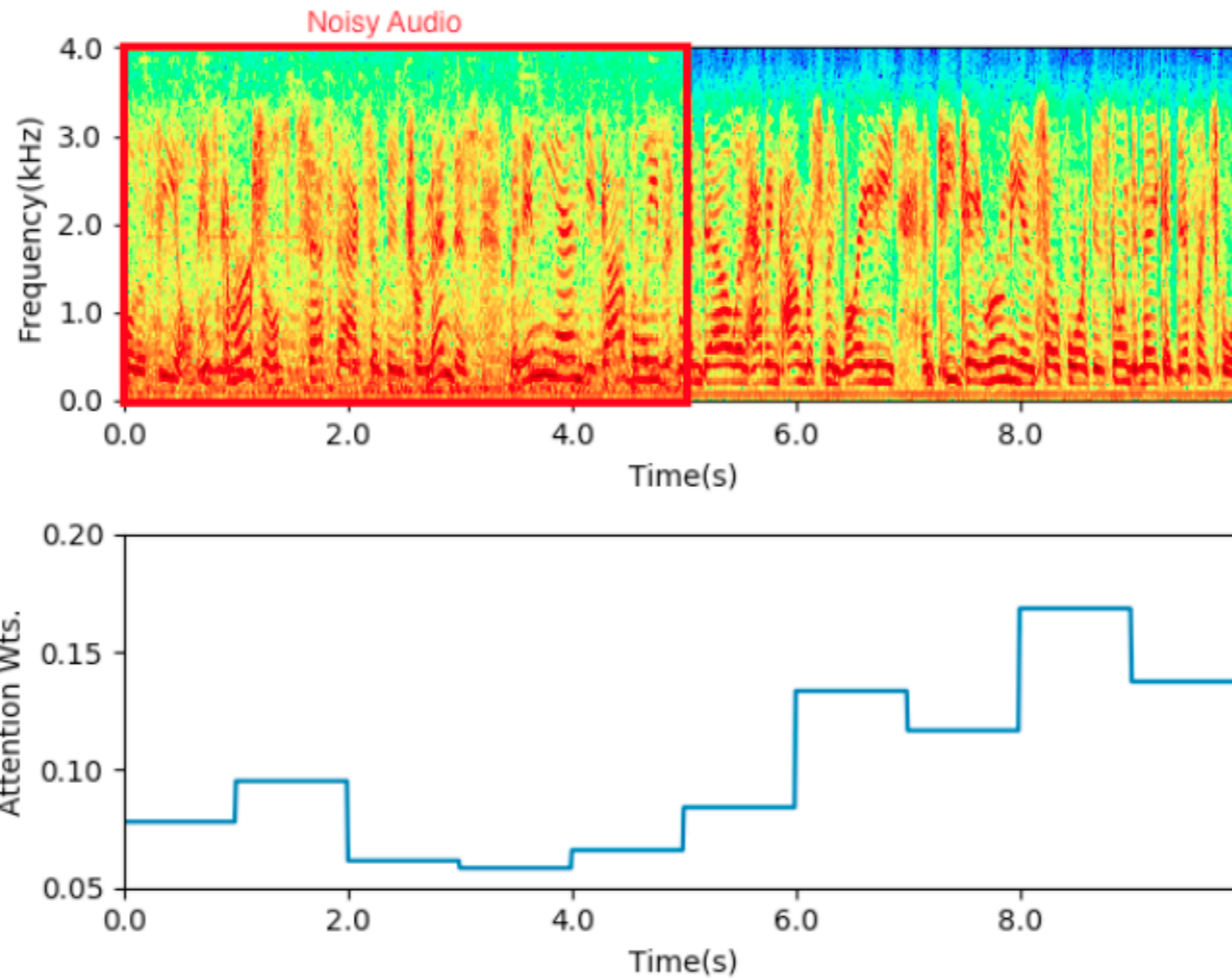
0-3s : O...One muscle at all, it was terrible

3s-4s : .... ah .... ah ....

4s - 9s : I couldn't scream, I couldn't shout, I couldn't even move my arms up, or my legs

9s -11s : I was trying me hardest, I was really really panicking.

Bharat Padi, et al. "End-to-end language recognition using hierarchical gated recurrent networks", under review 2018.

# Language Recognition Evaluation

# Language Recognition Evaluation

**Table 3**. Approximate computational time in seconds for ten 30sec eval files using a single CPU. Machine Specification: 32 CPU, 8 core, 2 thread Intel x86_64 machine with 16 GB Nvidia Quadro P5000 GPU cards.

|     | ivec. [19] | LSTM [16] | HGRU |
|-----|------------|-----------|------|
| CPU | 12         | 51        | 8    |
| GPU | 12         | 11.5      | 1.5  |

**Table 4**. LID accuracy in % for additional experiments with multiple speakers speaking the same language and the experiments without any SAD information.

| Cond. | i-vec. [19] | HGRU |
|-------|-------------|------|
| Multi-Speaker | 60.6 | **67.7** |
| Without SAD information | 49.7 | **52.7** |