# *Deep Learning: Theory and Practice*

**Recurrent Neural Networks**

**30-04-2019**

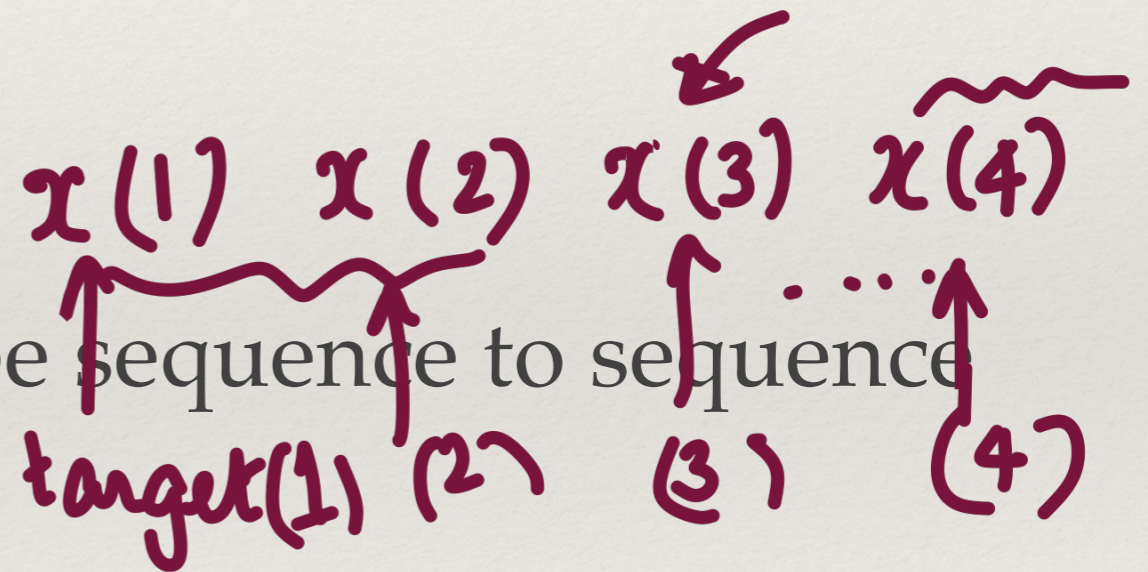# Introduction

- The standard DNN/CNN paradigms
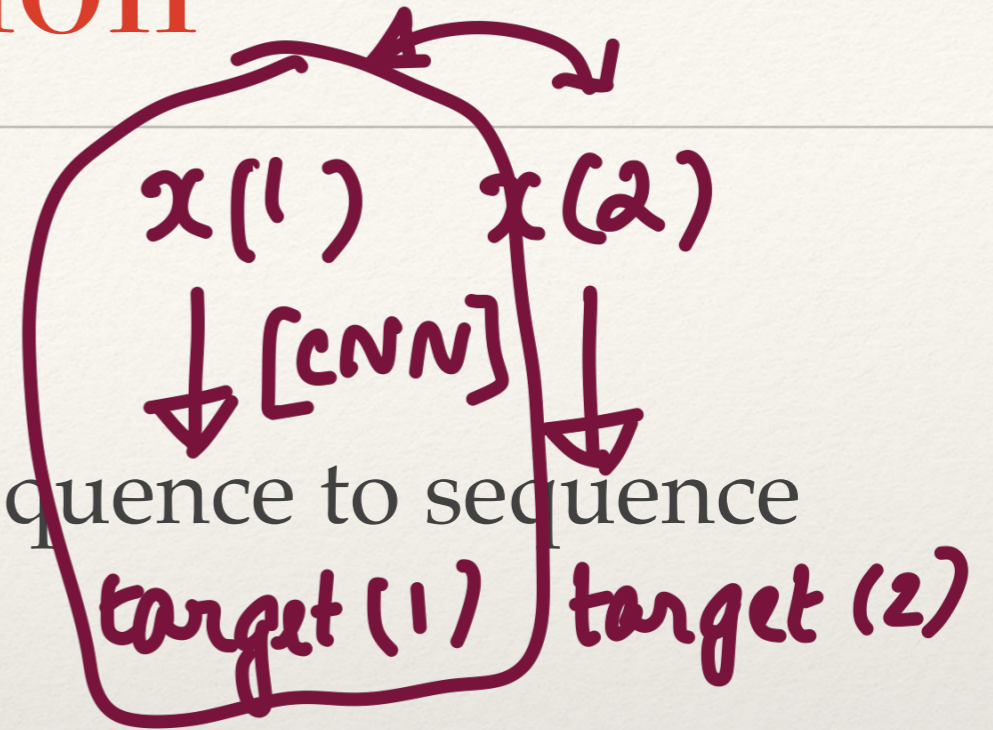  - (x,y) - ordered pair of data vectors/images (x) and target (t)

- 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.
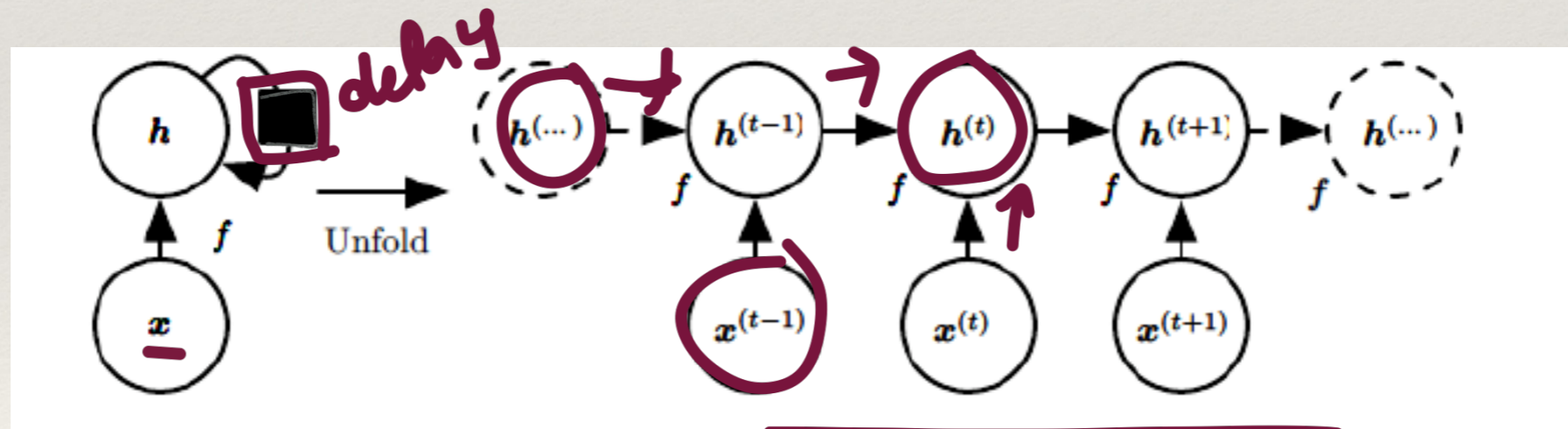
# Introduction to Recurrent Networks

**Recurrence**

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

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

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

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

Neural network

delay

Unfold

$h$

$x$

$f$

$h^{(\cdots)}$  $h^{(t-1)}$  $h^{(t)}$  $h^{(t+1)}$  $h^{(\cdots)}$

$f$   $f$   $f$   $f$

$x^{(t-1)}$  $x^{(t)}$  $x^{(t+1)}$

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

# Recurrent Networks



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

# Recurrent Networks

$X_1 = \{x_1^{(1)} \ldots \ldots x_1^{(T_1)}\}$ ← time

sequence index

$X_2 = \{x_2^{(1)} \ldots \ldots x_2^{(T_2)}\}$



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

$$L(\hat{y}(t), y(t))$$

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

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

# Back Propagation in RNNs

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

$$[x_1^{(1)}, x^{(2)}) \quad \ldots \quad x^{(30)}]$$
$$[x^{(1)}_2 \qquad x_2^{(17)}]$$ April, May

$$L\left(\{x^{(1)}, \ldots, x^{(\tau)}\}, \{y^{(1)}, \ldots, y^{(\tau)}\}\right)$$

$$= \sum_t L^{(t)}$$

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

## Model Parameters

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

## Gradient Descent

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

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

# Recurrent Networks



$(z)$

$y(t)$     $\underline{h}^{(t)} \rightarrow \underline{h}^{(t)} + \epsilon$

$i = y^{(t)}$

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

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

$(*)$

Backward Recurrence

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

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

# Back Propagation Through Time



Diagram labels and annotations:

$L_1$, $L_2$, $L_3$

$y_1$, $y_2$, $y_3$ — $y_1$, $y_2$

initialization — $h_0$ — $h_1$

$V$

$\frac{\partial h_2}{\partial h_1}$, $\frac{\partial h_3}{\partial h_2}$, $\frac{\partial L_3}{\partial h}$

$h_1$, $h_2$, $h_3$ — $h_{z-1}$, $h(z)$

$W$

$U$

$x_1$, $x_2$, $x_3$ — $x_2$

forward pass $L(t)$

Backward pass

$$\left[ \frac{\partial L}{\partial h^{(\tau-1)}} \right]$$

$$L = \ldots\ldots + L^{(\tau-1)} + L^{(\tau)}$$

$$h_{\tau-1} + \epsilon$$

$$O_{\tau} = V h_{\tau} + c \quad ; \quad \hat{y}_{\tau} = \text{softmax}(O_{\tau})$$

$$L^{\tau} = CE[y_{\tau}, \hat{y}_{\tau}]$$

$$O_{\tau-1}$$

$$h_{\tau-1}$$

$$h_{(\tau-1)}$$

$$x^{(\tau-1)} \qquad x^{(\tau)}$$

$$\frac{\partial L}{\partial h_{\tau}} = V^{T}\left( \frac{\partial L}{\partial O_{\tau}} \right)$$

$$L^{(\tau)}$$

$$\frac{\partial L}{\partial h^{(\tau-1)}} = V^{T} \frac{\partial L}{\partial O_{\tau-1}}^{(t+1)} + W^{T} (\text{dia} v) \frac{\partial L}{\partial h_{\tau}}$$

$$\tanh$$

$$(L)$$

$$h_{t}$$

$$h^{(\tau)} = \tanh\left( W \cdot h^{(\tau-1)} + U x^{(\tau)} + b \right)$$

$$h_{t+1}$$

# Back Propagation Through Time

# Standard Recurrent Networks



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

# Other Recurrent Networks



**Teacher Forcing Networks**

# Recurrent Networks



Teacher
Forcing Networks
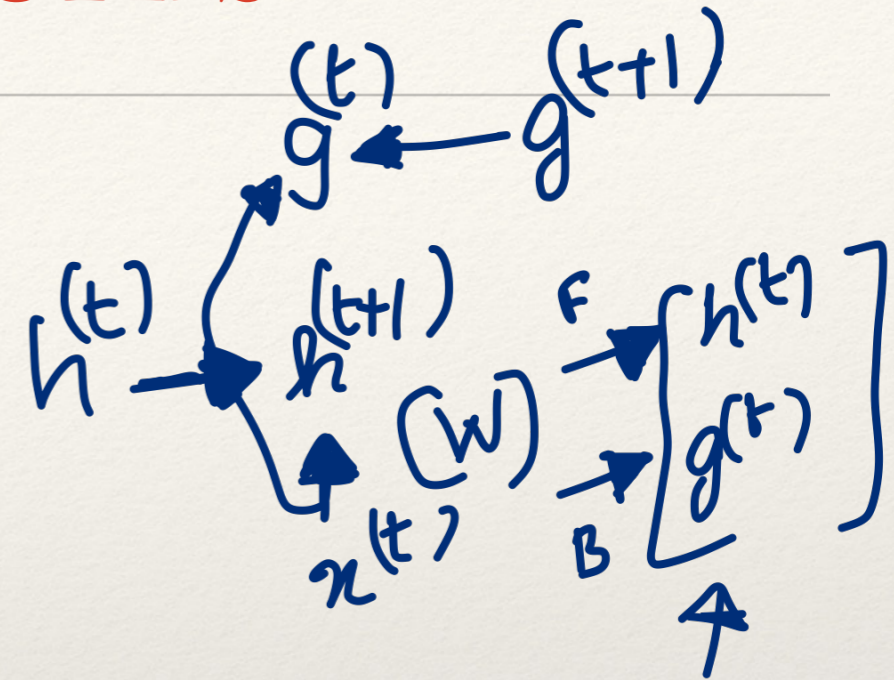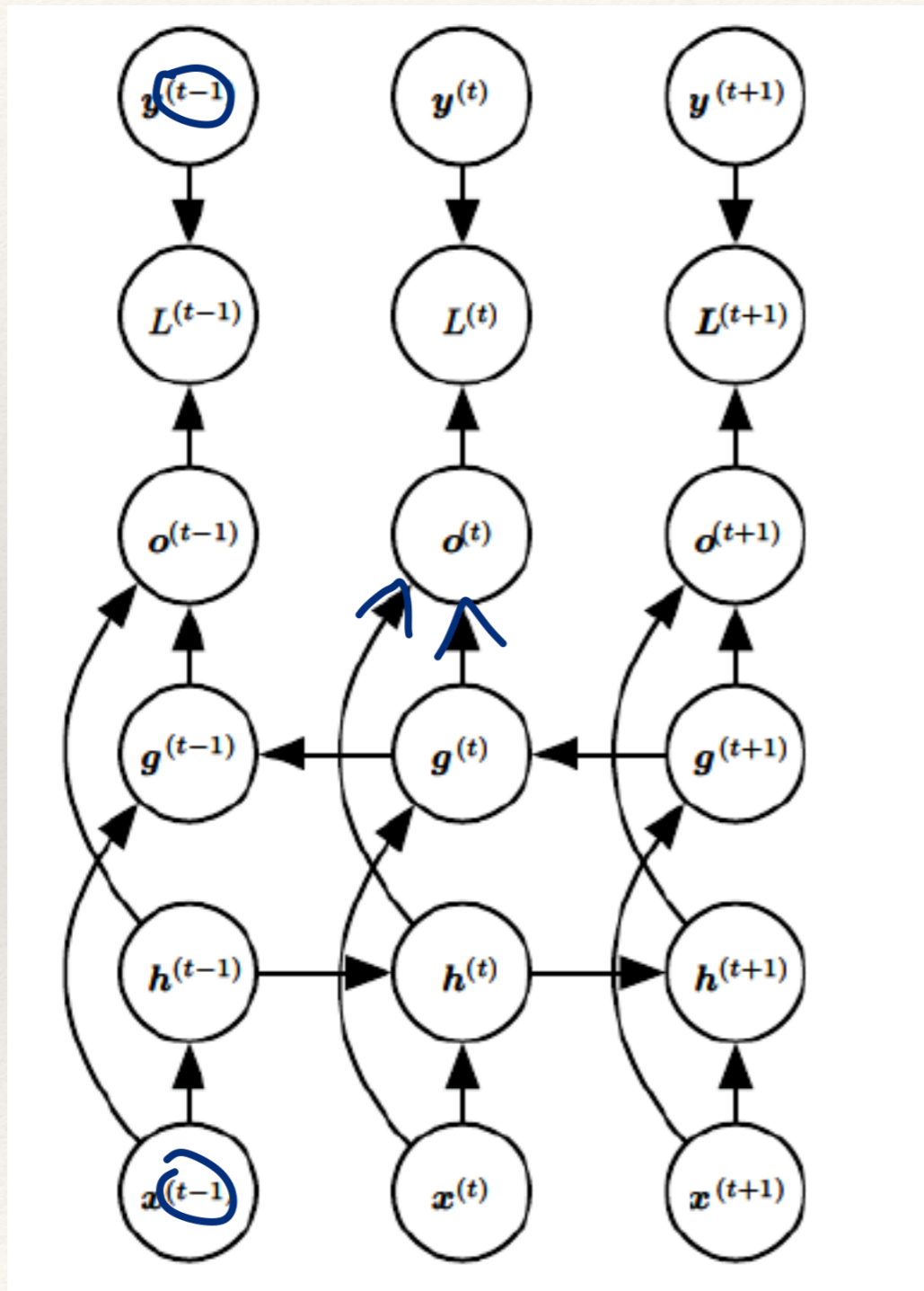
*"Deep Learning" Ian Goodfellow, Yoshua Bengio, Aaron Courville*

# Recurrent Networks



**Multiple Input
Single Output**

# Recurrent Networks



**Single Input Multiple Output**

Autoregressive Models

# Recurrent Networks



**Bi-directional Networks**

$$\begin{bmatrix} W_1 & W_2 \end{bmatrix} \begin{bmatrix} h \\ g \end{bmatrix}$$

multiple input

single output

C A T        B(A)R

# Recurrent Networks



Encoder

$x^{(1)}$  $x^{(2)}$  $x^{(...)}$  $x^{(n_x)}$

$c$

sentence embedding

$c^1$ $c^2$ ..
$c^m$... $c^{n_y}$

Decoder

$y^{(1)}$  $y^{(2)}$  $y^{(...)}$  $y^{(n_y)}$

hindi

$x^{(t)}$
$\downarrow$
$y(t)$

$x^{(t)}$
$\downarrow$
$(y^m)$

**Sequence to Sequence Mapping Networks**

Encoder ← sequence to vector

Decoder ← vector to a new sequence

$$\textcircled{X} = \left\{ x^{(1)} \ldots \ldots \ x^{(T)} \right\} \rightarrow \left\{ \text{Eng Senten} \right\}$$

$$\textcircled{Y} = \left\{ y^{(1)} \ldots \ y^{(M)} \right\} \rightarrow \left\{ \text{Hindi Sentence} \right\}$$

$$X \rightarrow \boxed{\begin{array}{c} \text{Encoder} \\ \text{Decoder} \end{array}} \rightarrow Y$$

# Long-term Dependency Issues

# Vanishing/Exploding Gradients
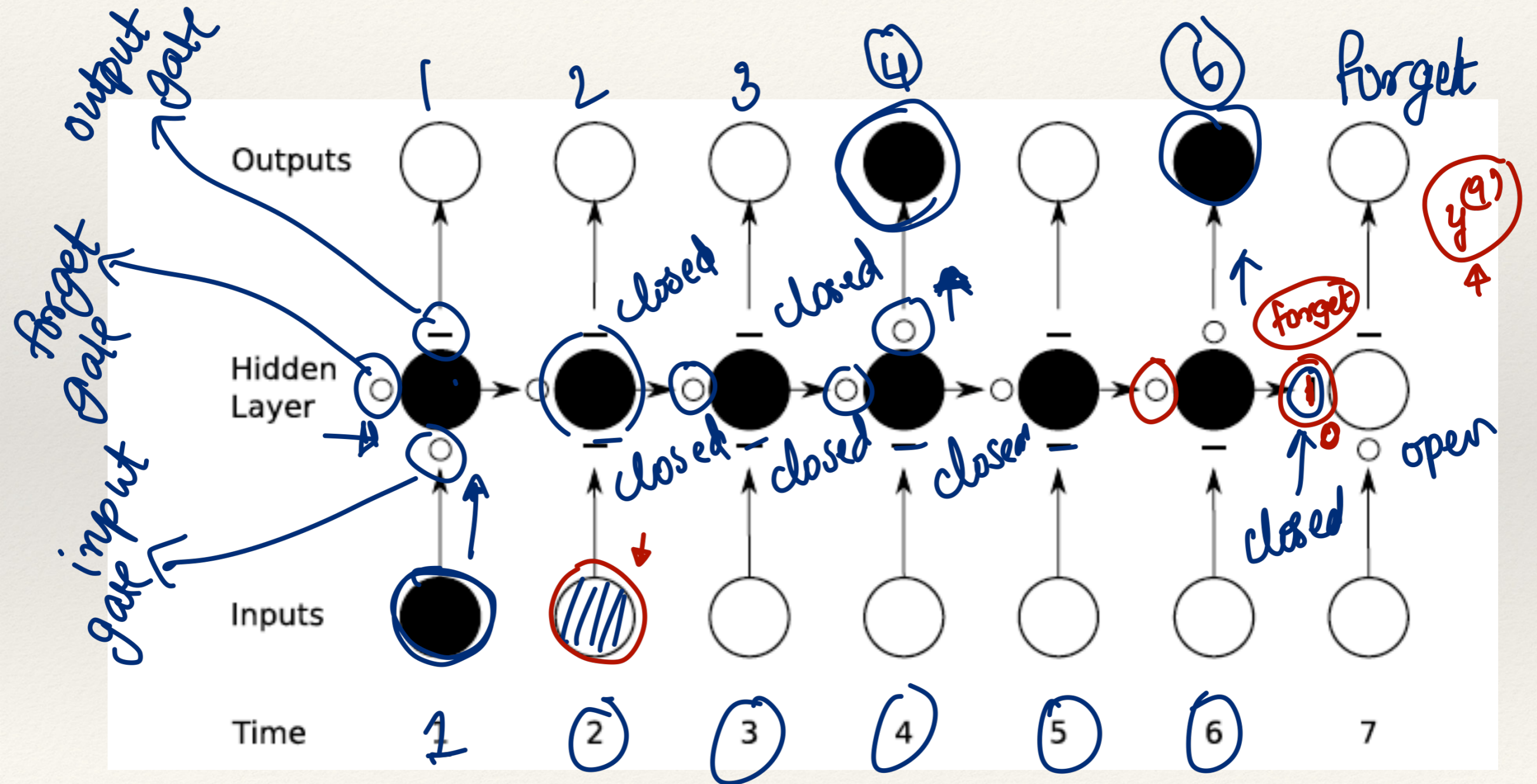


* Gradients either vanish or explode

* 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

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

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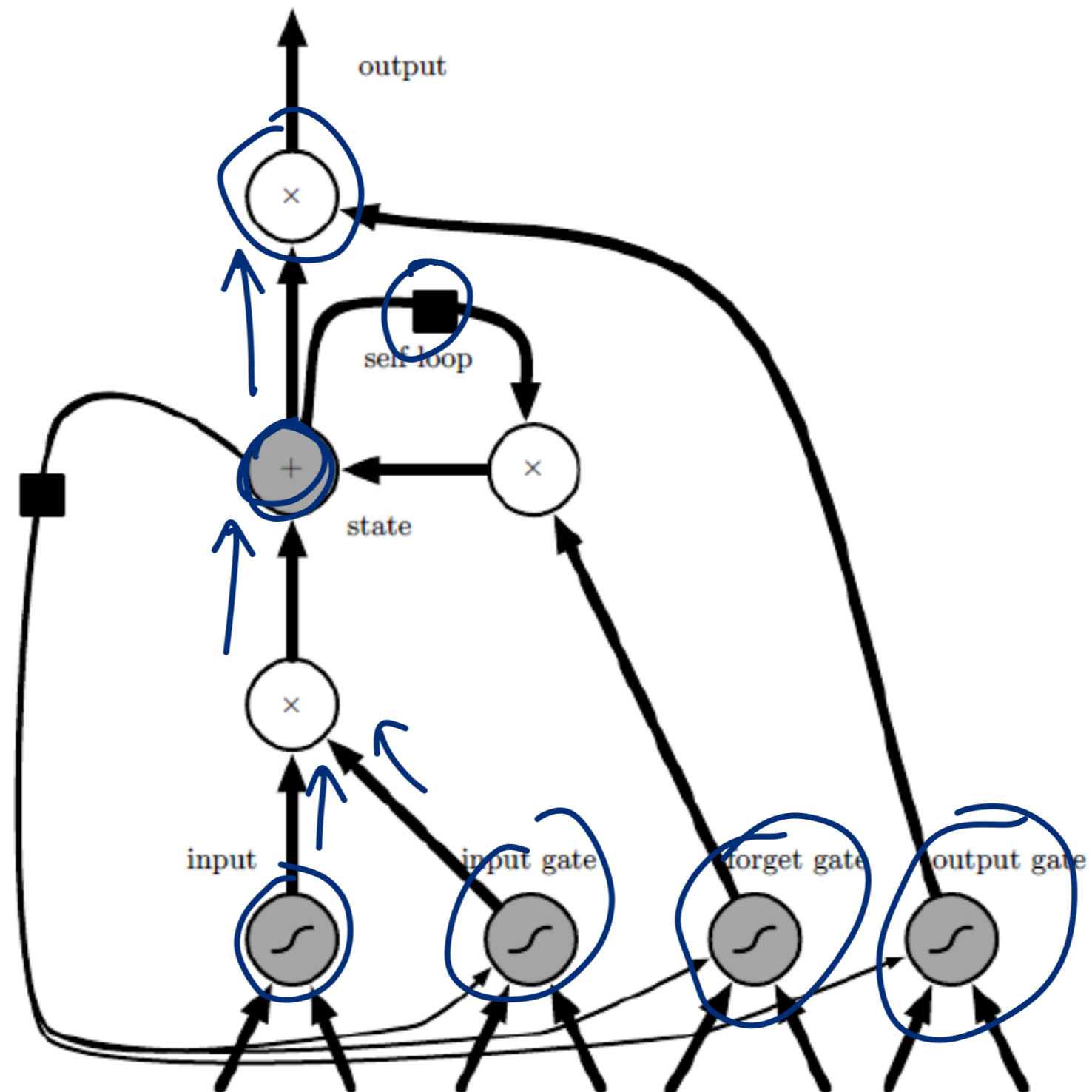
## Output Gate

$$a_\omega^t = \sum_{i=1}^{I} w_{i\omega}x_i^t + \sum_{h=1}^{H} w_{h\omega}b_h^{t-1} + \sum_{c=1}^{C} w_{c\omega}s_c^t$$
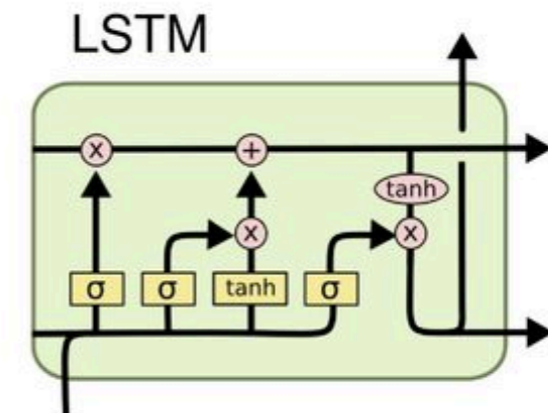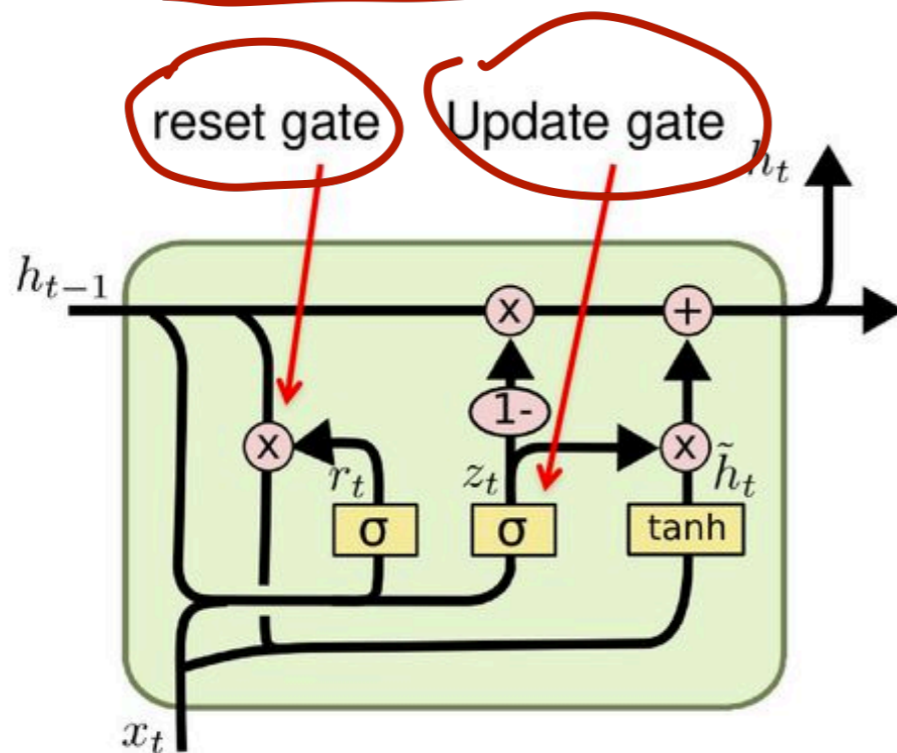
$$b_\omega^t = f(a_\omega^t)$$

## LSTM output

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

# Long Short Term Memory Networks

# Gated Recurrent Units (GRU)



LSTM

## GRU – gated recurrent unit

(more compression)

reset gate    Update gate

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

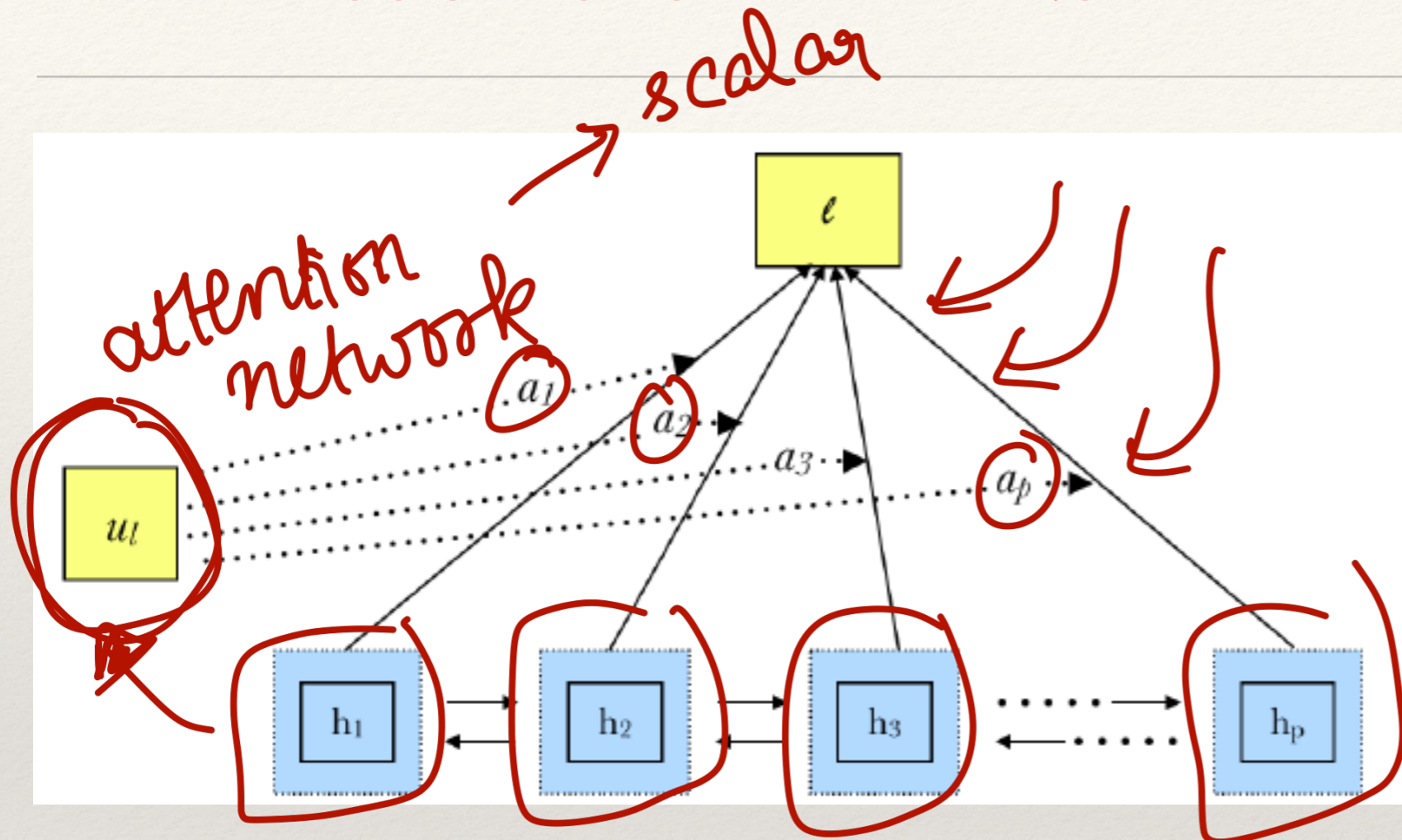$$r_t = \sigma\left(W_r \cdot [h_{t-1}, x_t]\right)$$

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

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

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.
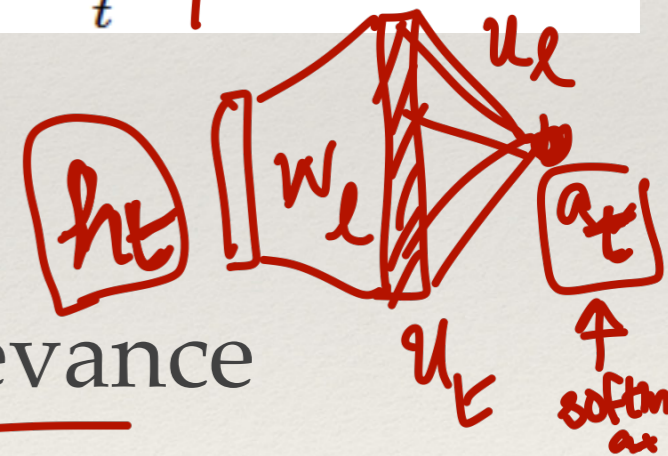
X,*: element-wise multiply

# Attention in LSTM Networks
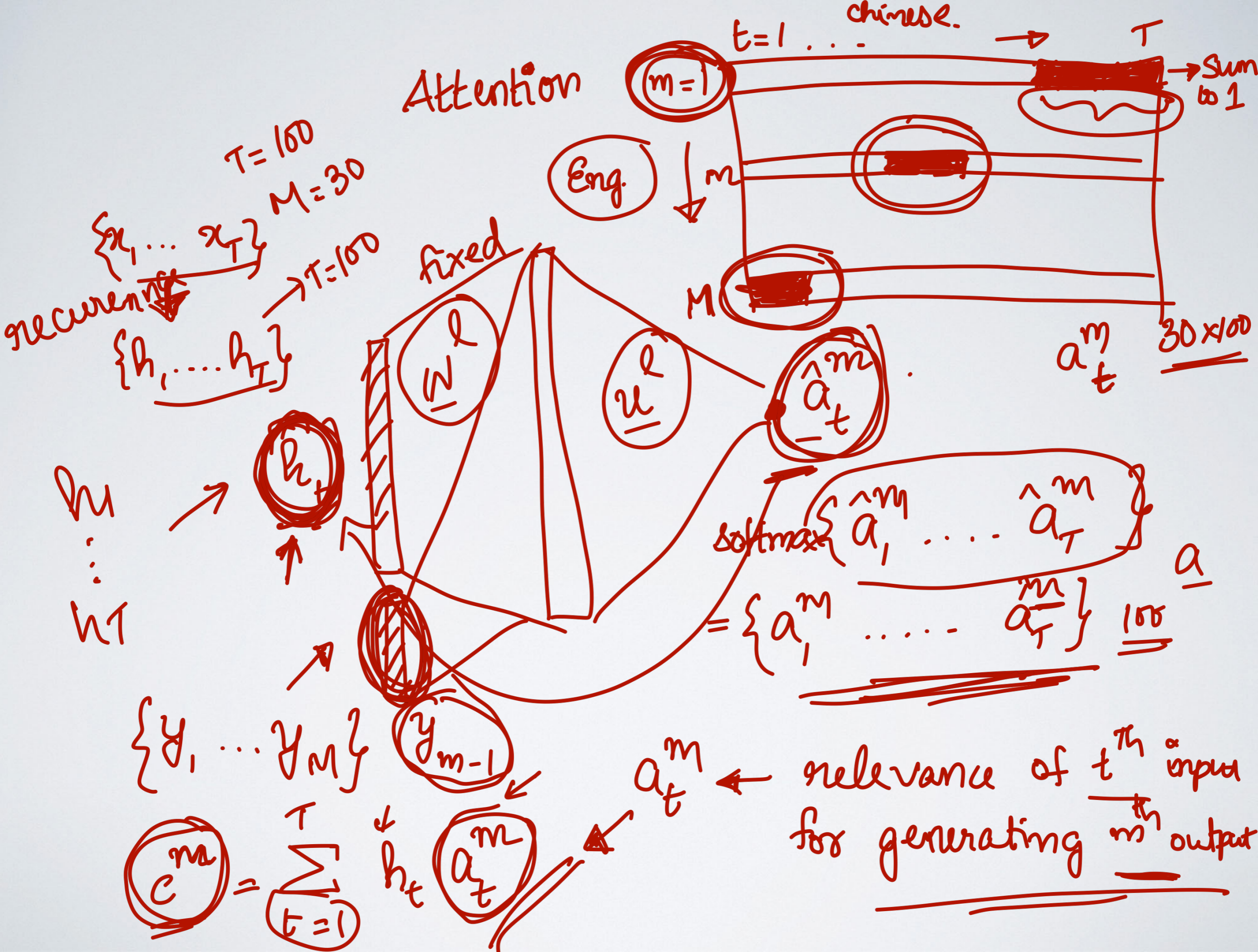


$$\mathbf{u}_t = tanh(\mathbf{W}_l \mathbf{h}_t + \mathbf{b}_l)$$

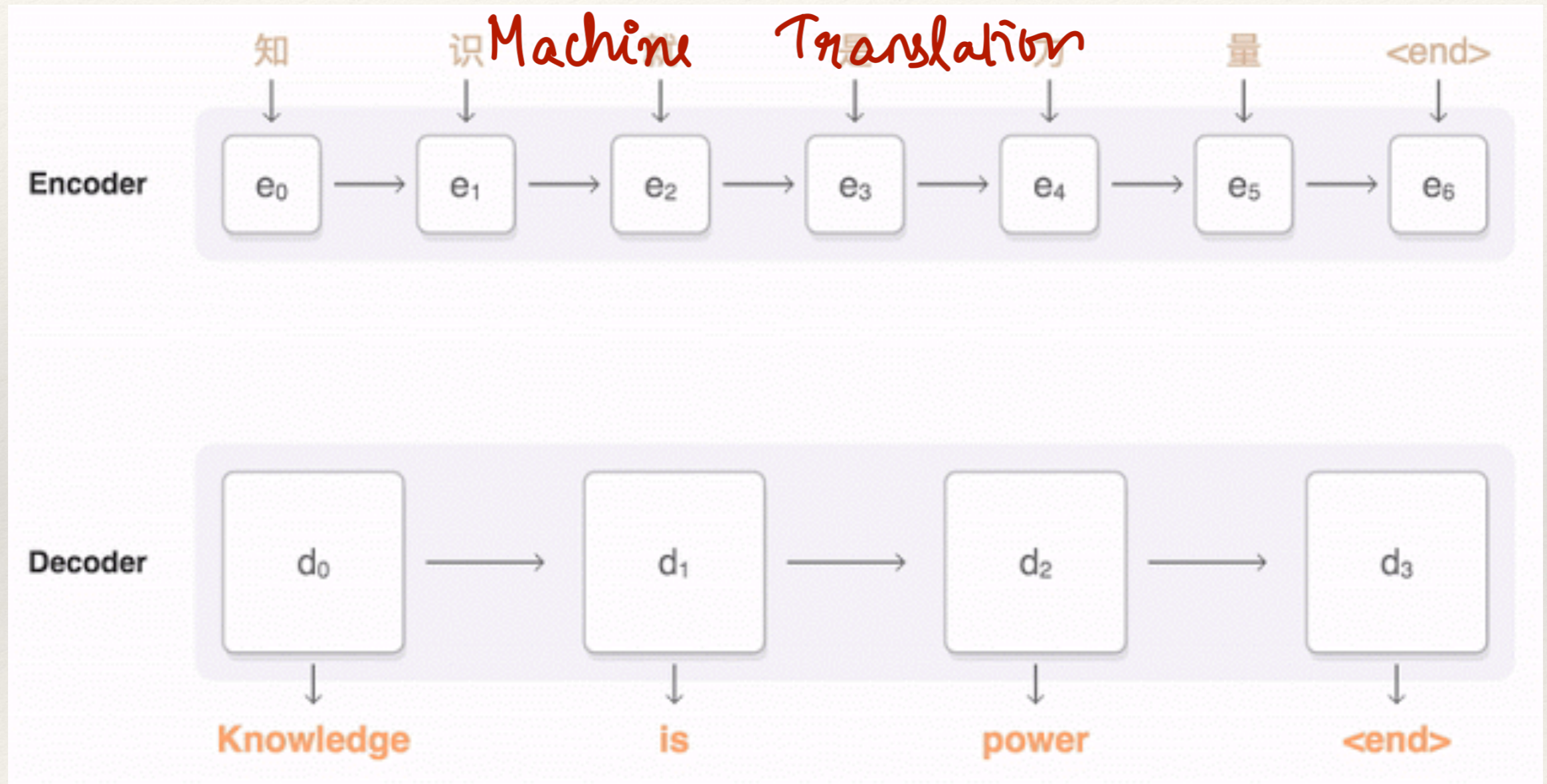$$a_t = \frac{exp(\mathbf{u}_t^T \mathbf{u}_l)}{\sum_t exp(\mathbf{u}_t^T \mathbf{u}_l)}$$

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

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

# Attention

$T=100$
$M=30$

$\{x_1 \ldots x_T\}$

recurrent $\rightarrow T=100$

$\{h_1, \ldots, h_T\}$

fixed

Eng.

$m=1$    $t=1 \ldots$   chinese. $\rightarrow$   $T$

$\rightarrow$ Sum to 1

$\downarrow m$

$M$

$a^m_t$   $30 \times 100$

$W^l$   $u^l$   $\hat{a}^m_t$ .

$h_t$

$hu$
$\vdots$
$hT$

softmax $\{\hat{a}^m_1 \ldots \hat{a}^m_T\}$

$= \{a^m_1 \ldots a^m_T\}$   $\underline{100}$   $\underline{a}$

$y_{m-1}$

$\{y_1 \ldots y_M\}$

$c^m = \sum_{t=1}^{T} h_t \quad a^m_t$

$a^m_t \leftarrow$ relevance of $t^{th}$ input for generating $m^{th}$ output

# Encoder – Decoder Networks with Attention

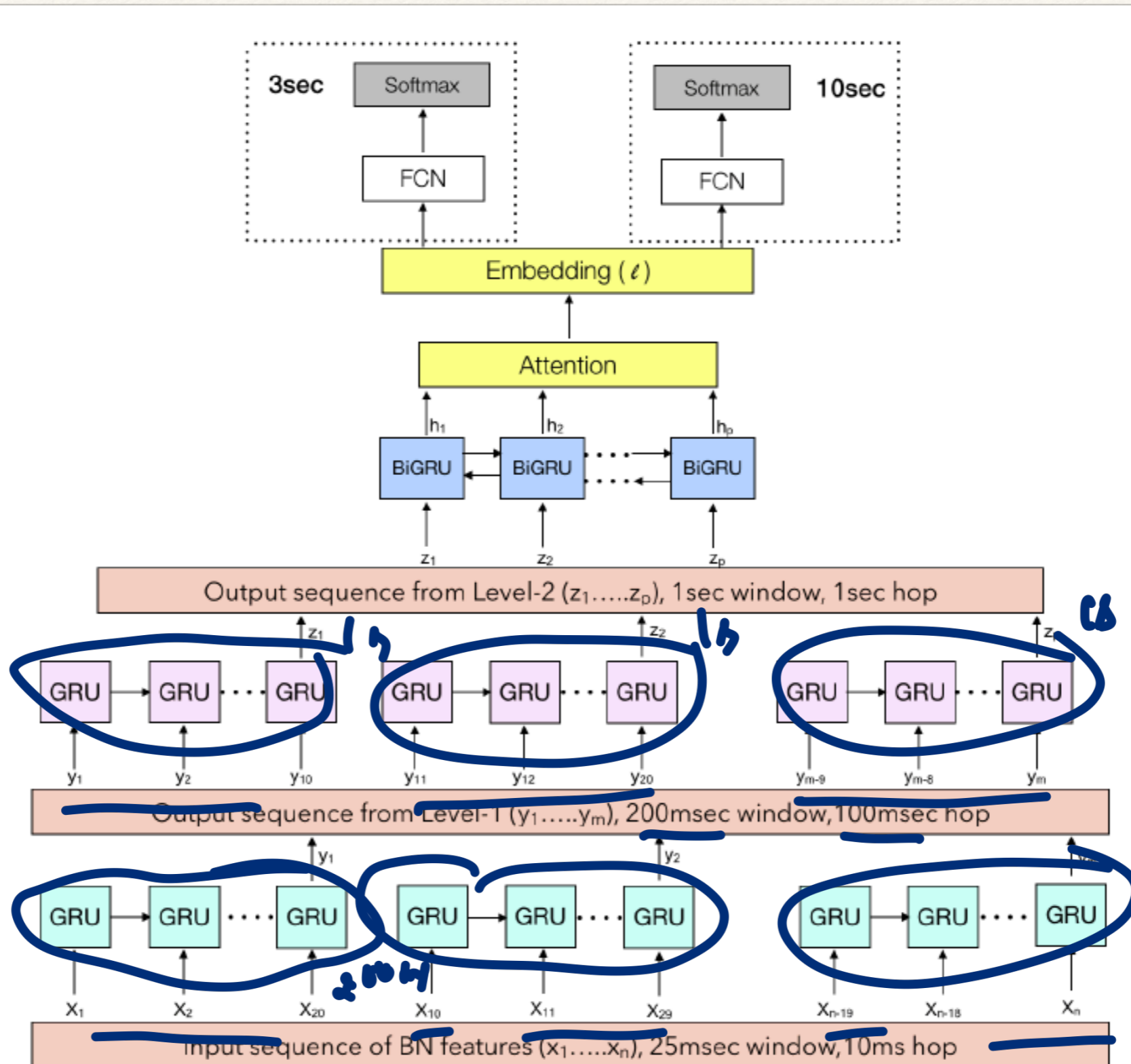# Attention Models

# Attention - Speech Example

From our lab [part of ICASSP 2019 paper].
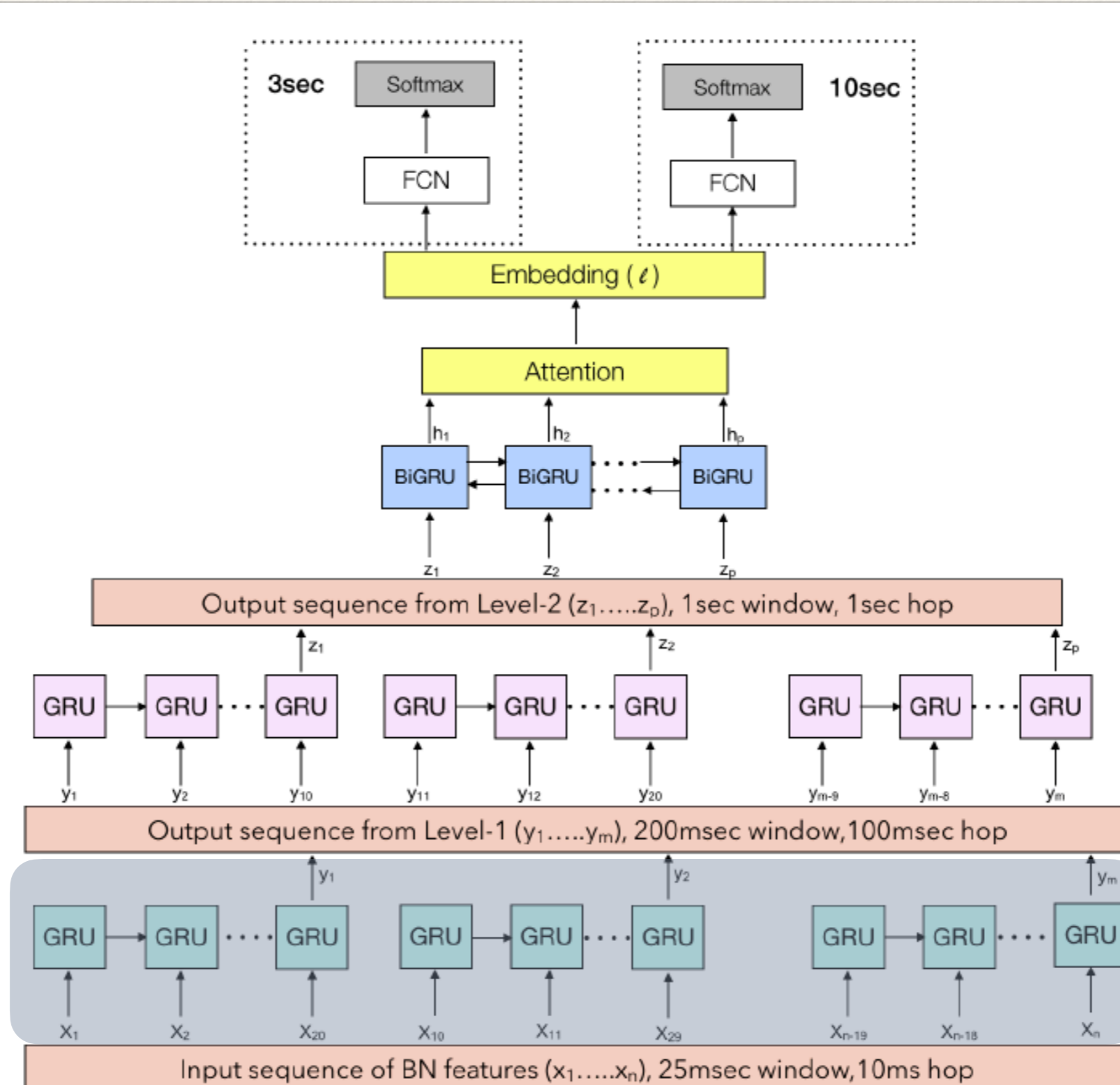
# Language Recognition Evaluation

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

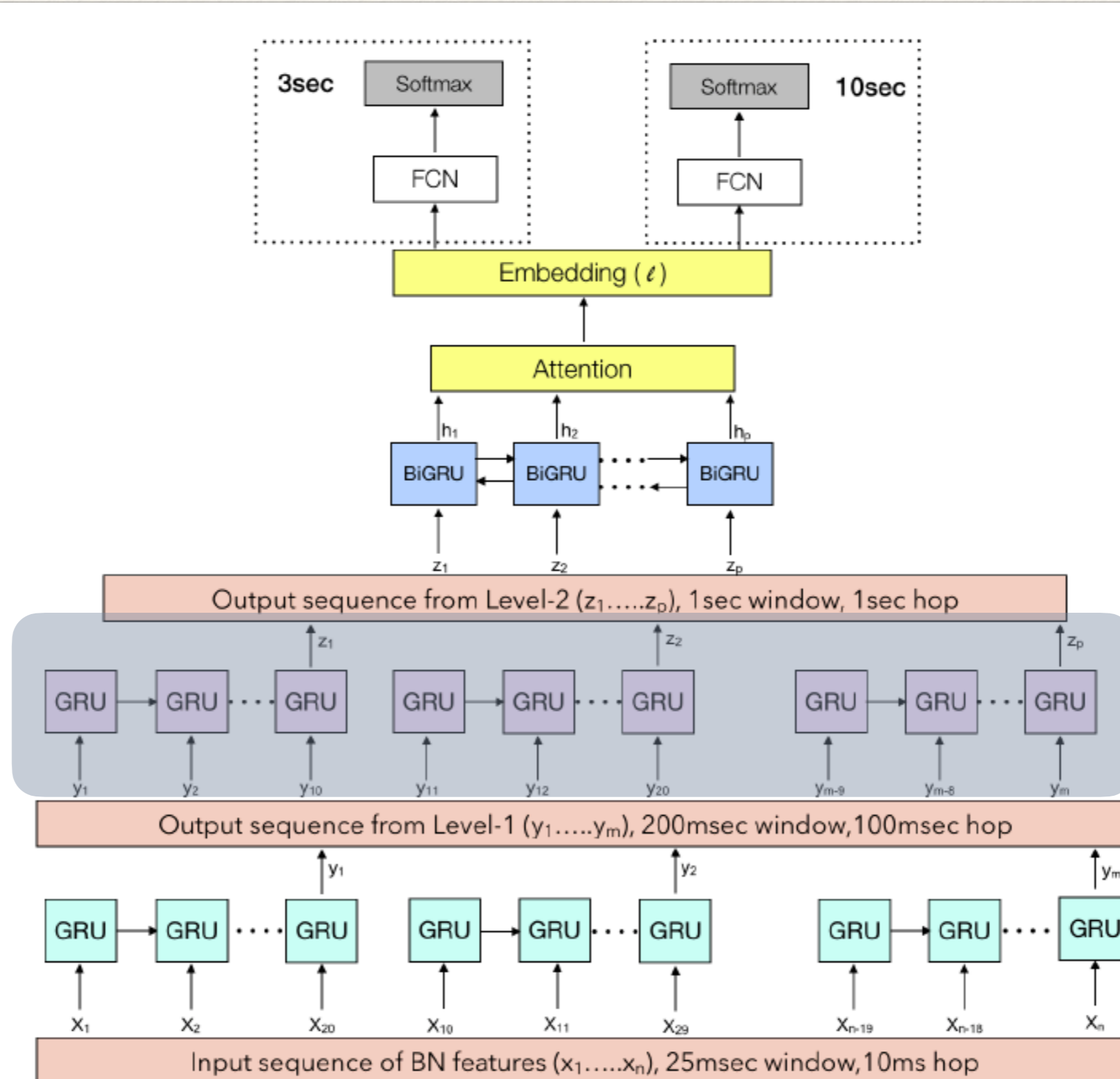| Cluster | Target Languages | Hours |
|---|---|---|
| Arabic | Egyptian Arabic (ara-arz) | 190.9 |
| | Iraqi Arabic (ara-acm) | 130.8 |
| | Levantine Arabic (ara-apc) | 440.7 |
| | Maghrebi Arabic (ara-ary) | 81.8 |
| Chinese | Mandarin (zho-cmn) | 379.4 |
| | Min Nan (zho-nan) | 13.3 |
| English | 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 |

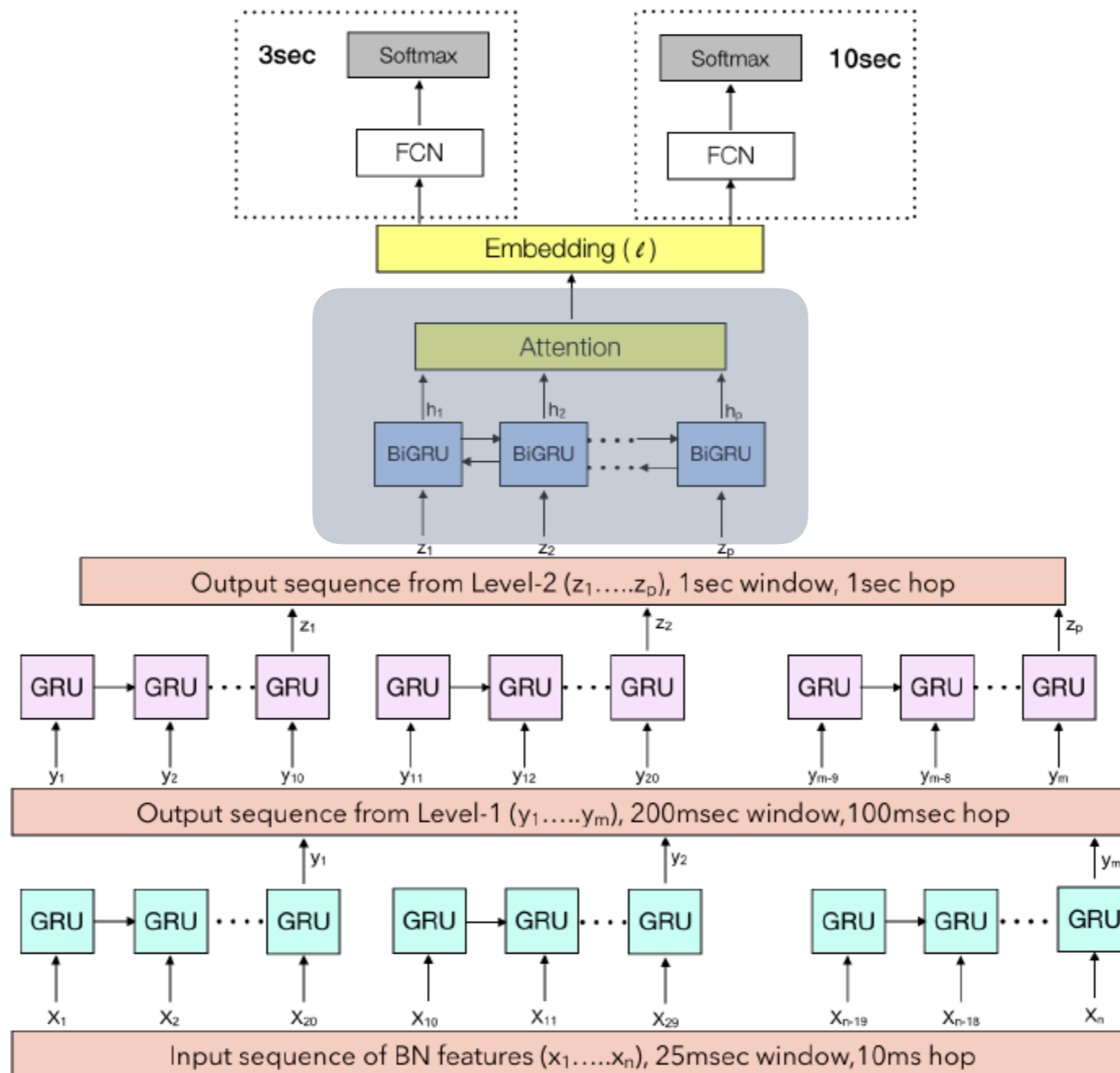# End-to-end model using GRUs and Attention

# 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 th importance of each short-term segment feature for the task.
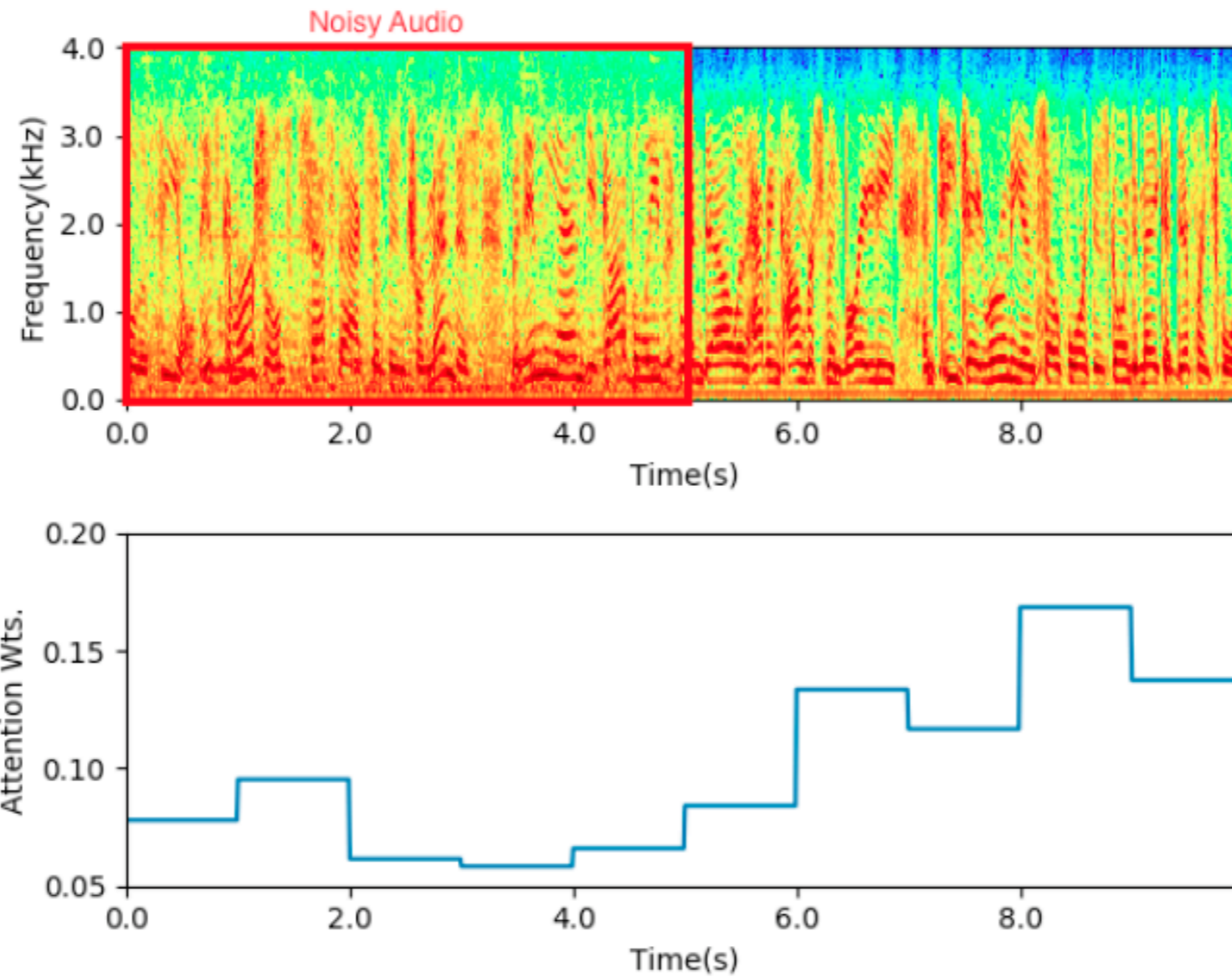
**Attention Weight**

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** |