

# CS 4248

## Natural Language Processing

**Professor NG Hwee Tou**  
**Department of Computer Science**  
**School of Computing**  
**National University of Singapore**  
**[nght@comp.nus.edu.sg](mailto:nght@comp.nus.edu.sg)**

# Materials

- NNM4NLP Chapter 17

# Fixed-Window Neural Language Model

- Input to multilayer perceptron: a sequence of  $k$  words
- Output: a probability distribution over the next word

# Fixed Window Neural Language Model

$$v(w) = \mathbf{E}_{[w]}$$

$$\mathbf{x} = [v(w_1); v(w_2); \dots; v(w_k)]$$

$$\mathbf{h} = g(\mathbf{x}\mathbf{W}^1 + \mathbf{b}^1)$$

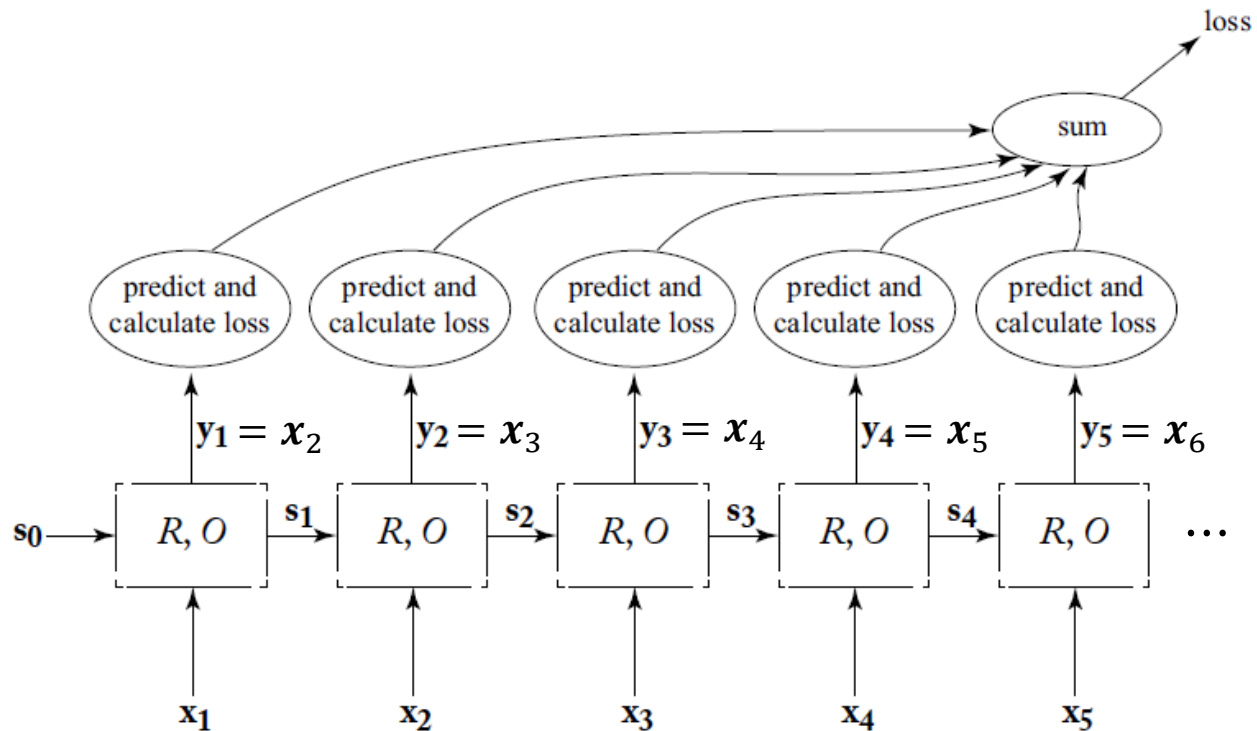
$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{h}\mathbf{W}^2 + \mathbf{b}^2)$$

$$w_i \in V \quad \mathbf{E} \in \mathbb{R}^{|V| \times d_w} \quad \mathbf{W}^1 \in \mathbb{R}^{k \cdot d_w \times d_{\text{hid}}} \quad \mathbf{b}^1 \in \mathbb{R}^{d_{\text{hid}}} \\ \mathbf{h} \in \mathbb{R}^{d_{\text{hid}}} \quad \mathbf{W}^2 \in \mathbb{R}^{d_{\text{hid}} \times |V|} \quad \mathbf{b}^2 \in \mathbb{R}^{|V|}$$

# Fixed Window Neural LM

- Fixed window context: unable to use a larger context
- RNN language model
  - Able to use a context of **any length**
  - Model size does not increase with larger context

# RNN Language Model

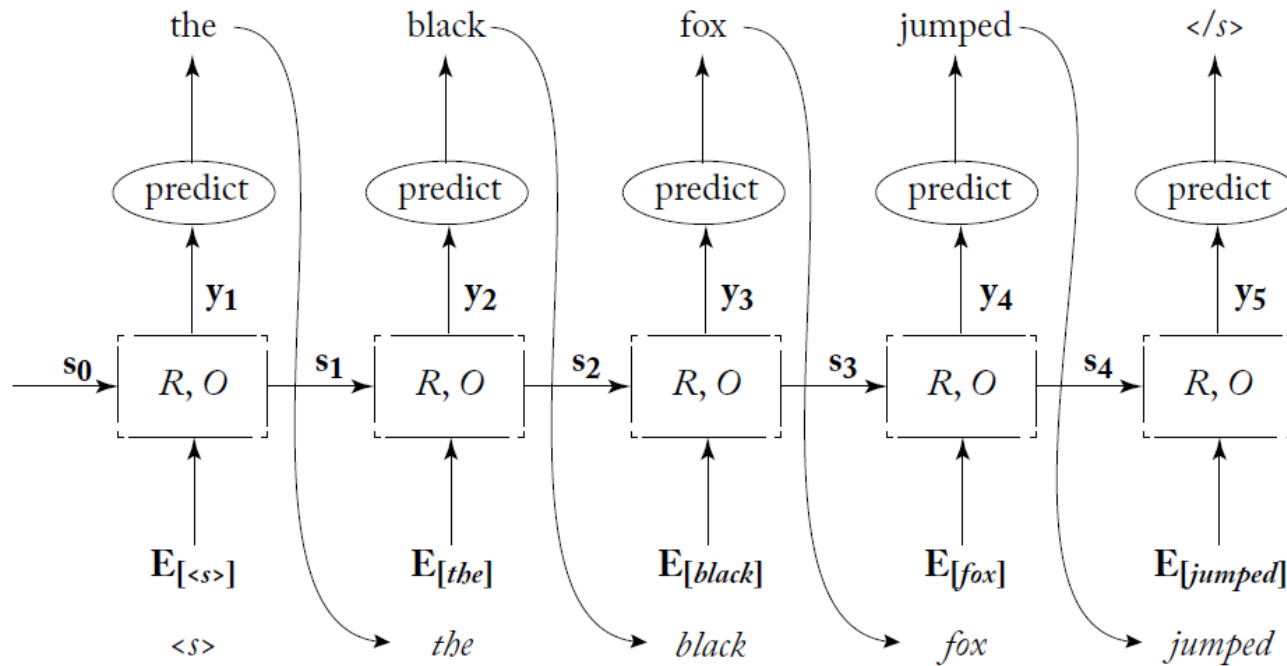


# RNN Language Model

- Categorical cross-entropy loss (negative log likelihood)
- $L = -\sum_t \log_2(\hat{y}_t)$
- RNN LMs are better language models (have lower perplexity) compared to n-gram LMs

# RNN Generators

## Sequence generation





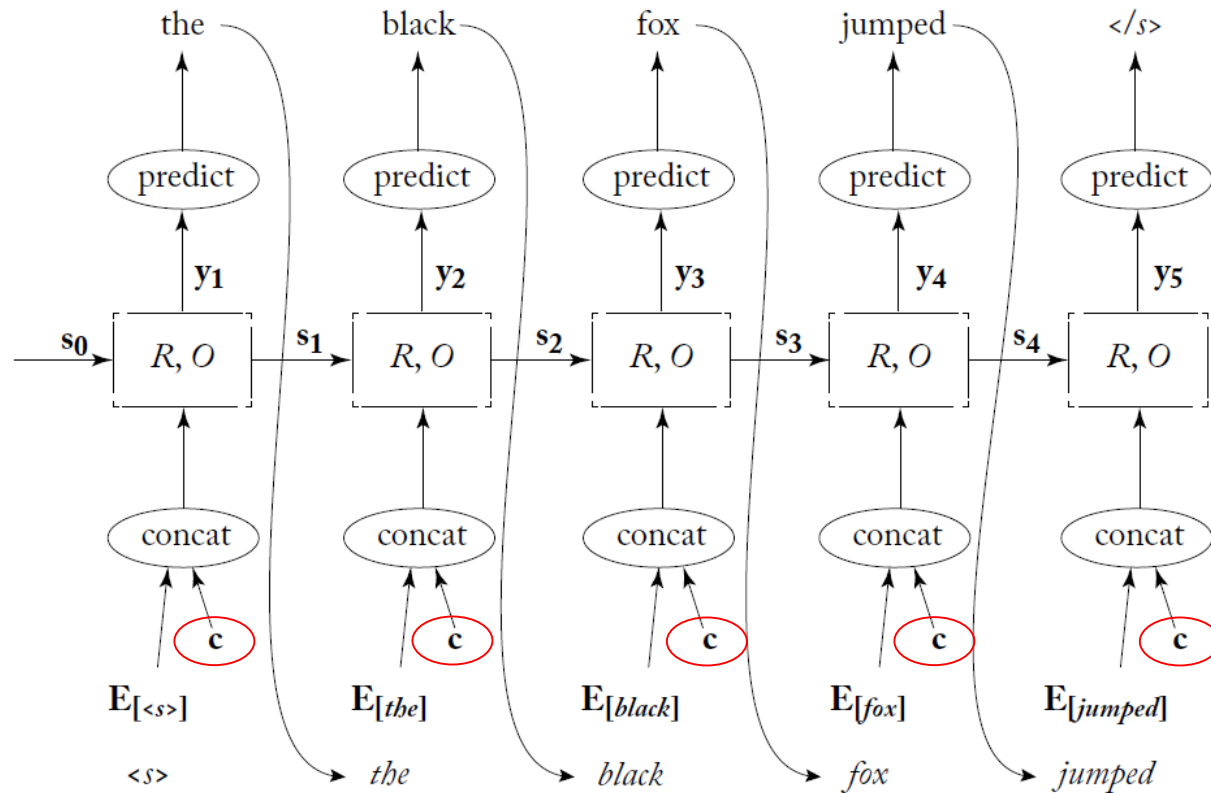
# RNN Generators

$$\mathbf{s}_j = R(\mathbf{s}_{j-1}, \hat{\mathbf{t}}_j)$$

$$\mathbf{y}_j = O(\mathbf{s}_j)$$

$$p(\hat{\mathbf{t}}_{j+1} | \hat{\mathbf{t}}_{1:j}) = \text{softmax}(\text{MLP}(\mathbf{y}_j))$$

# Conditioned Generation



# Conditioned Generation

$$\mathbf{s}_j = R(\mathbf{s}_{j-1}, [\hat{\mathbf{t}}_j; \mathbf{c}])$$

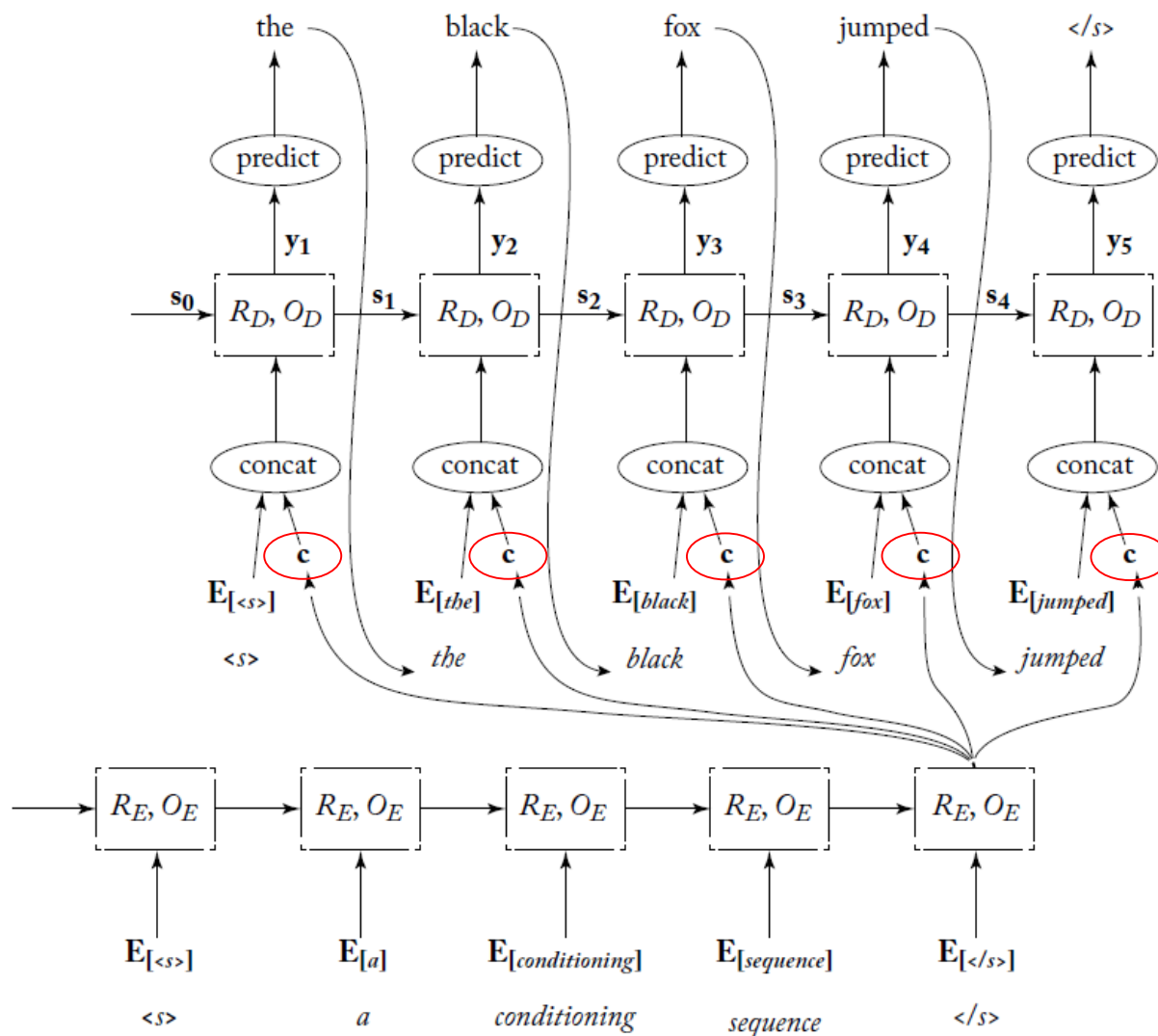
$$\mathbf{y}_j = O(\mathbf{s}_j)$$

$$p(\hat{t}_{j+1} | \hat{t}_{1:j}, \mathbf{c}) = \text{softmax}(\text{MLP}(\mathbf{y}_j))$$

# Sequence to Sequence Model

- Aka encoder-decoder model
- Application: neural machine translation, grammatical error correction
  - Translate a source sentence  $x_{1:n}$  into a target sentence  $t_{1:m}$
- Encoder:  $c = \text{RNN}_{\text{enc}}(x_{1:n})$
- Decoder: a conditioned generator RNN
- The whole model (encoder and decoder) is trained end-to-end

# Sequence to Sequence Model



# Seq2Seq Model with Attention

$$\mathbf{s}_j = R_{\text{dec}}(\mathbf{s}_{j-1}, [\hat{\mathbf{t}}_j; \mathbf{c}^j])$$

$$\mathbf{y}_j = O_{\text{dec}}(\mathbf{s}_j)$$

$$p(\hat{t}_{j+1} | \hat{t}_{1:j}, \mathbf{x}_{1:n}) = \text{softmax}(\text{MLP}^{\text{out}}(\mathbf{y}_j))$$

# Seq2Seq Model with Attention

$$\mathbf{c}^j = \sum_{i=1}^n \alpha_{[i]}^j \cdot \mathbf{c}_i \quad \leftarrow \text{attend}$$

$$\mathbf{c}_{1:n} = \text{biRNN}_{\text{enc}}^*(\mathbf{x}_{1:n})$$

$$\bar{\alpha}_{[i]}^j = \mathbf{v} \cdot \tanh([\mathbf{s}_{j-1}; \mathbf{c}_i] \mathbf{U} + \mathbf{b}) \quad \leftarrow \text{MLP}^{\text{att}}$$

$$\alpha^j = \text{softmax}(\bar{\alpha}_{[1]}^j, \dots, \bar{\alpha}_{[n]}^j)$$

# Seq2Seq Model with Attention

- Attention weights  $\alpha_{[i]}^j$  reveal which parts  $i$  of the source sentence the decoder finds relevant at output step  $j$ .
- Improved interpretability