# Topic Models
# & Recurrent Neural Networks

**Kuan-Yu Chen (陳冠宇)**

2018/04/12 @ TR-409, NTUST

# Semantic?

- In the context of NLP, long-span information is important!
    - Especially for language modeling

<div align="center">

一 枝 美麗 的 玫瑰 在 花盆

一 隻 可愛 的 小貓 在 花園

一 支 高貴 的 鋼筆 在 拍賣

</div>

- N-gram models are not an efficient strategy to capture the long-span information
    - From literal matching to semantic mapping

# Probabilistic Latent Semantic Analysis

- **Probabilistic Latent Semantic Analysis** also called
  - Probabilistic Latent Semantic Indexing (PLSI)
  - Aspect Model

- PLSA is a probabilistic counterpart of LSA
  - $P(d_j)$: the probability of selecting document $d_j$
  - $P(w_i|T_k)$: the probability of word $w_i$ condition on a latent factor/topic $T_k$
    - **Aspect!**
  - $P(T_k|d_j)$: the probability of a latent factor/topic $T_k$ generated by document $d_j$

# PLSA – 1

- The PLSA model is a latent variable model for co-occurrence data (i.e., each pair of word $w_i$ and document $d_j$) which associates an unobserved class variable (i.e., latent factor $T_k$)
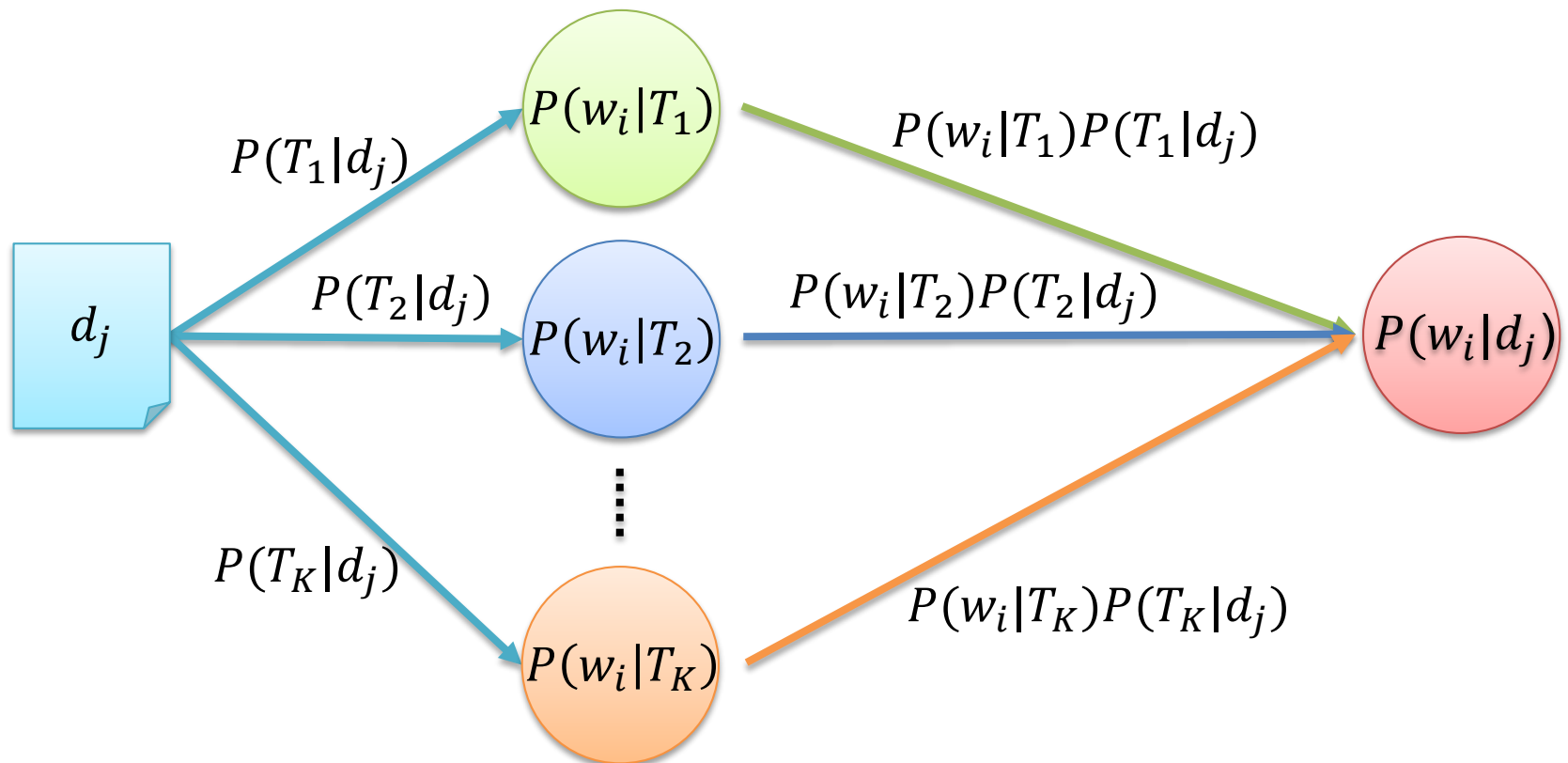
$$P(w_i, d_j) = P(d_j)P(w_i|d_j) = P(d_j) \sum_{k=1}^{K} P(w_i|T_k)P(T_k|d_j)$$

$$P(w_i|d_j) = \sum_{k=1}^{K} P(w_i, T_k|d_j) = \sum_{k=1}^{K} \frac{P(w_i, T_k, d_j)}{P(d_j)}$$

$$= \sum_{k=1}^{K} \frac{P(w_i, d_j|T_k)P(T_k)}{P(d_j)}$$

$$= \sum_{k=1}^{K} \frac{P(w_i|T_k)P(d_j|T_k)P(T_k)}{P(d_j)}$$

**Conditional Independence Assumption**
document and word are independent conditioned on the state of the associated latent variable

$$= \sum_{k=1}^{K} \frac{P(w_i|T_k)P(d_j, T_k)}{P(d_j)} = \sum_{k=1}^{K} P(w_i|T_k)P(T_k|d_j)$$

4

# PLSA − 2

- Thus, the modeling goal is to identify conditional probability mass functions $P(w_i|T_k)$ such that the document-specific word distributions $P(w_i|d_j)$ are as faithfully as possible approximated by convex combinations of these aspects

# PLSA – 3

- The training objective is defined to maximize the total log-likelihood of a given training collection

  - The model parameters are $P(d_j)$, $P(w_i|T_k)$, and $P(T_k|d_j)$

$$\mathcal{L} = \sum_{w_i \in V} \sum_{d_j \in \mathbf{D}} c(w_i, d_j) \log P(w_i, d_j)$$

$$= \sum_{w_i \in V} \sum_{d_j \in \mathbf{D}} c(w_i, d_j) \log \left( P(d_j) \sum_{k=1}^{K} P(w_i|T_k) P(T_k|d_j) \right)$$

# PLSA − 4

- By using the Expectation-Maximization algorithm
  - E-step

$$P(T_k|w_i, d_j) = \frac{P(w_i|T_k)P(T_k|d_j)}{\sum_{k=1}^{K} P(w_i|T_k)P(T_k|d_j)}$$

  - M-step

$$P(w_i|T_k) = \frac{\sum_{d_j \in \mathbf{D}} c(w_i, d_j)P(T_k|w_i, d_j)}{\sum_{i'=1}^{|V|} \sum_{d_j \in \mathbf{D}} c(w_{i'}, d_j)P(T_k|w_{i'}, d_j)}$$

$$P(T_k|d_j) = \frac{\sum_{i=1}^{|V|} c(w_i, d_j)P(T_k|w_i, d_j)}{\sum_{i'=1}^{|V|} c(w_{i'}, d_j)}$$

# PLSA – 5

- Consequently, for a given word sequence, $w_1, w_2, \ldots, w_T$, the joint probability in a language can be calculated by using PLSA

$$P(w_1, w_2, \ldots, w_T) = P(w_1) \prod_{t=2}^{T} P(w_t | w_1, w_2, \ldots, w_{t-1})$$

$$= P(w_1) \prod_{t=2}^{T} \left( \sum_{k=1}^{K} P(w_t | T_k) P(T_k | w_1, w_2, \ldots, w_{t-1}) \right)$$

  - Usually, we can combine the PLSA with the traditional n-gram models

    - Semantic matching and literal term matching

$$P(w_t | w_1, w_2, \ldots, w_{t-1}) = \alpha \cdot P(w_t | w_{t-n+1}, \ldots, w_{t-1}) +$$

$$(1 - \alpha) \cdot \sum_{k=1}^{K} P(w_t | T_k) P(T_k | w_1, w_2, \ldots, w_{t-1})$$

# PLSA − 6

- For a new history of words, $w_1, w_2, \ldots, w_{t-1} = H_1^{t-1}$, the **fold-in** strategy can be perform to obtain the topic distribution

  - The word distribution for each topic $P(w_i|T_k)$ is fixed

  - E-step

$$P(T_k|w_i, H_1^{t-1}) = \frac{P(w_i|T_k)P(T_k|H_1^{t-1})}{\sum_{k=1}^{K} P(w_i|T_k)P(T_k|H_1^{t-1})}$$

  - M-step

$$P(T_k|H_1^{t-1}) = \frac{\sum_{i=1}^{|V|} c(w_i, H_1^{t-1})P(T_k|w_i, H_1^{t-1})}{\sum_{i'=1}^{|V|} c(w_{i'}, H_1^{t-1})}$$

# Revisiting NNLM – 1

- The Neural Network Language Mode (NNLM) estimated a statistical ($n$-gram) language model for **predicting future words**

$$P(w_1, w_2, \ldots, w_T) \approx \prod_{t=1}^{T} P(w_t | w_{t-n+1}, \ldots, w_{t-1})$$

on

Dense Representation

Concatenate/Average

One-hot Representation

the

cat

sat

10

# Revisiting NNLM – 2.

$P(there\ are\ books\ on\ the\ table)$
$\approx \textcolor{red}{\boldsymbol{P(there)}}P(are|there)P(books|there\ are)P(on|are\ books)$
$P(the|books\ on)P(table|on\ the)$

there

**Output Layer**

a          zoo

**Hidden Layer**

**Concatenate**

**Embedding Layer**

**Input Layer**

\<s\>          \<s\>

# Revisiting NNLM – 2..

$P(there\ are\ books\ on\ the\ table)$
$\approx P(there)\textbf{\textit{P(are|there)}}P(books|there\ are)P(on|are\ books)$
$P(the|books\ on)P(table|on\ the)$

# Revisiting NNLM – 2...

$P(there\ are\ books\ on\ the\ table)$
$\approx P(there)P(are|there)\textbf{\textcolor{red}{P(books|there\ are)}}P(on|are\ books)$
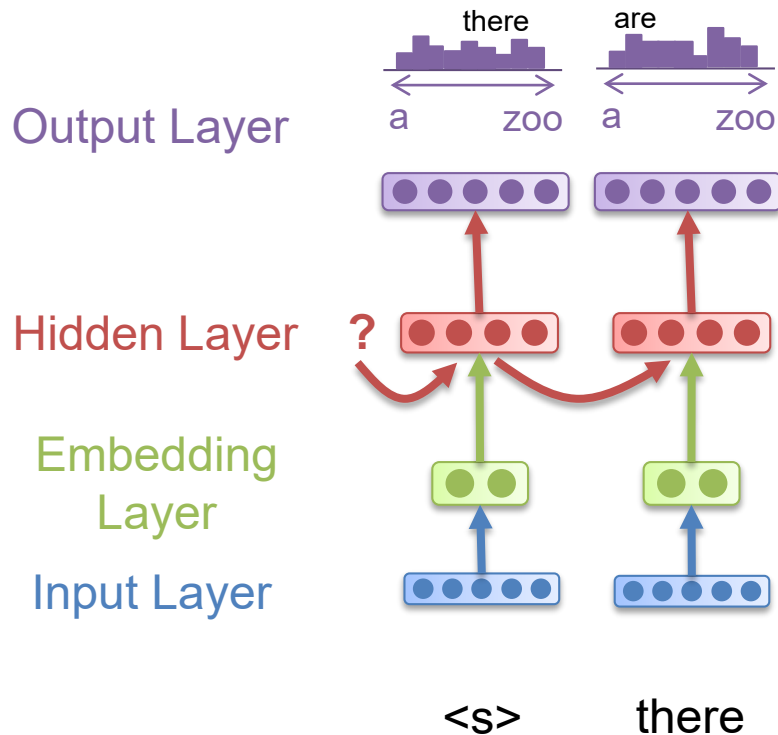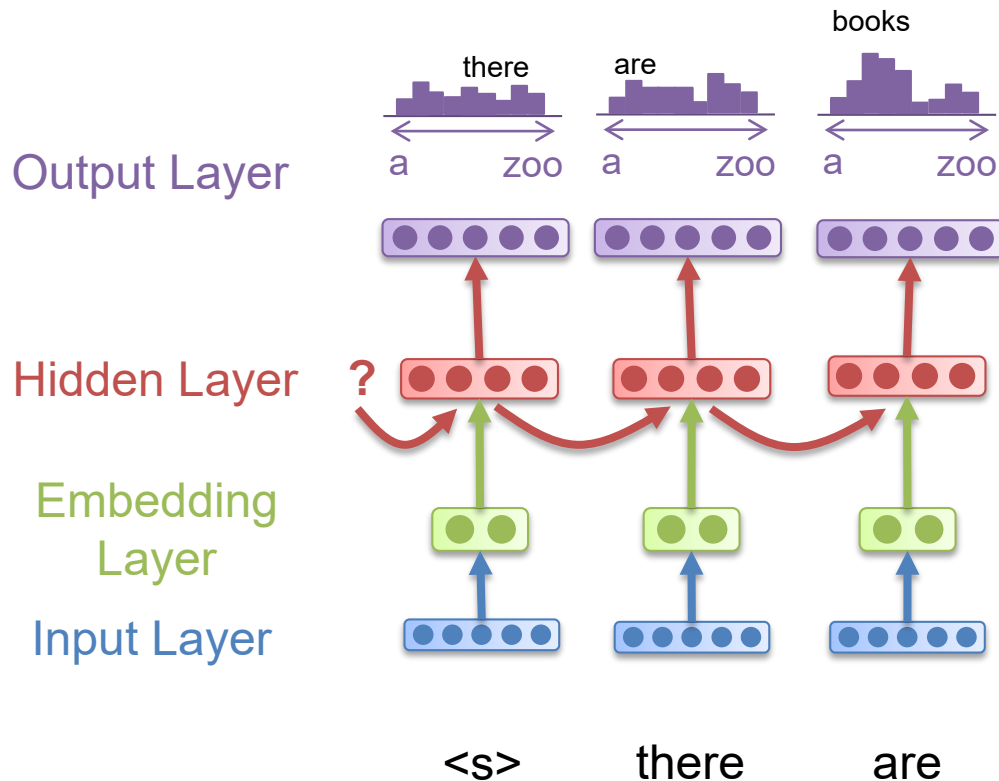$P(the|books\ on)P(table|on\ the)$

# Revisiting NNLM – 2....
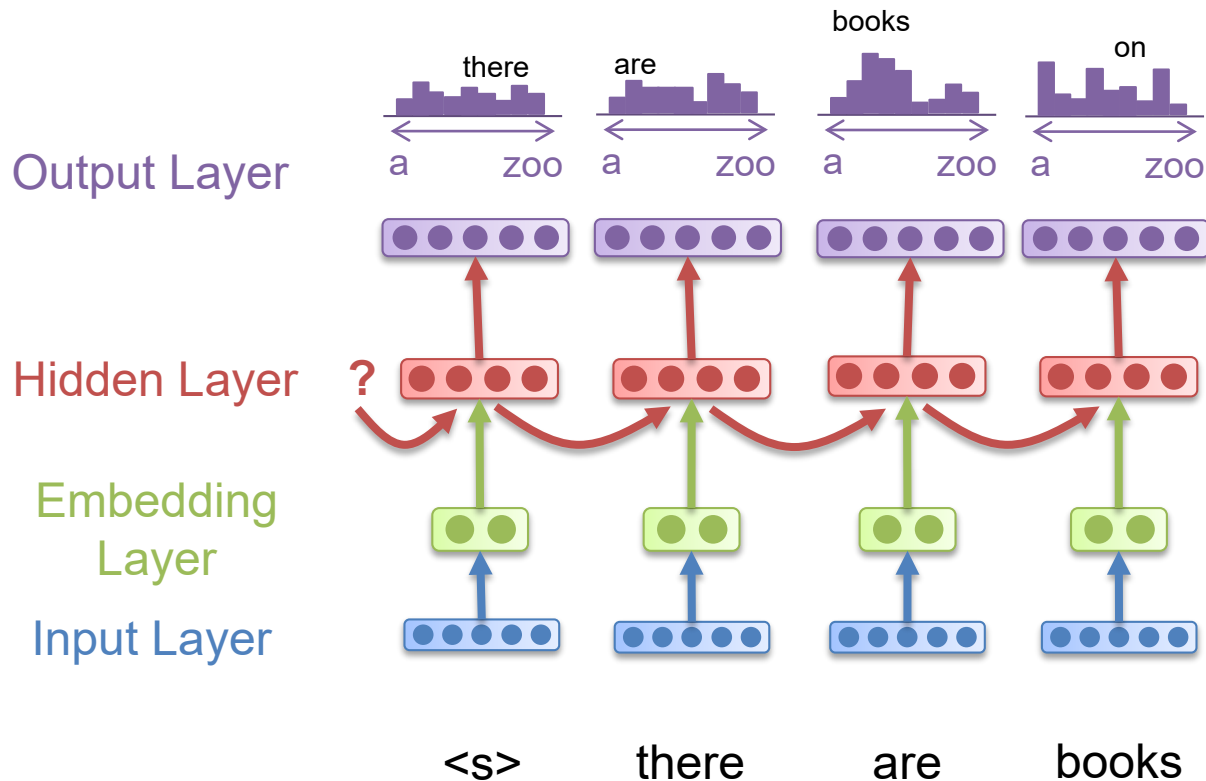
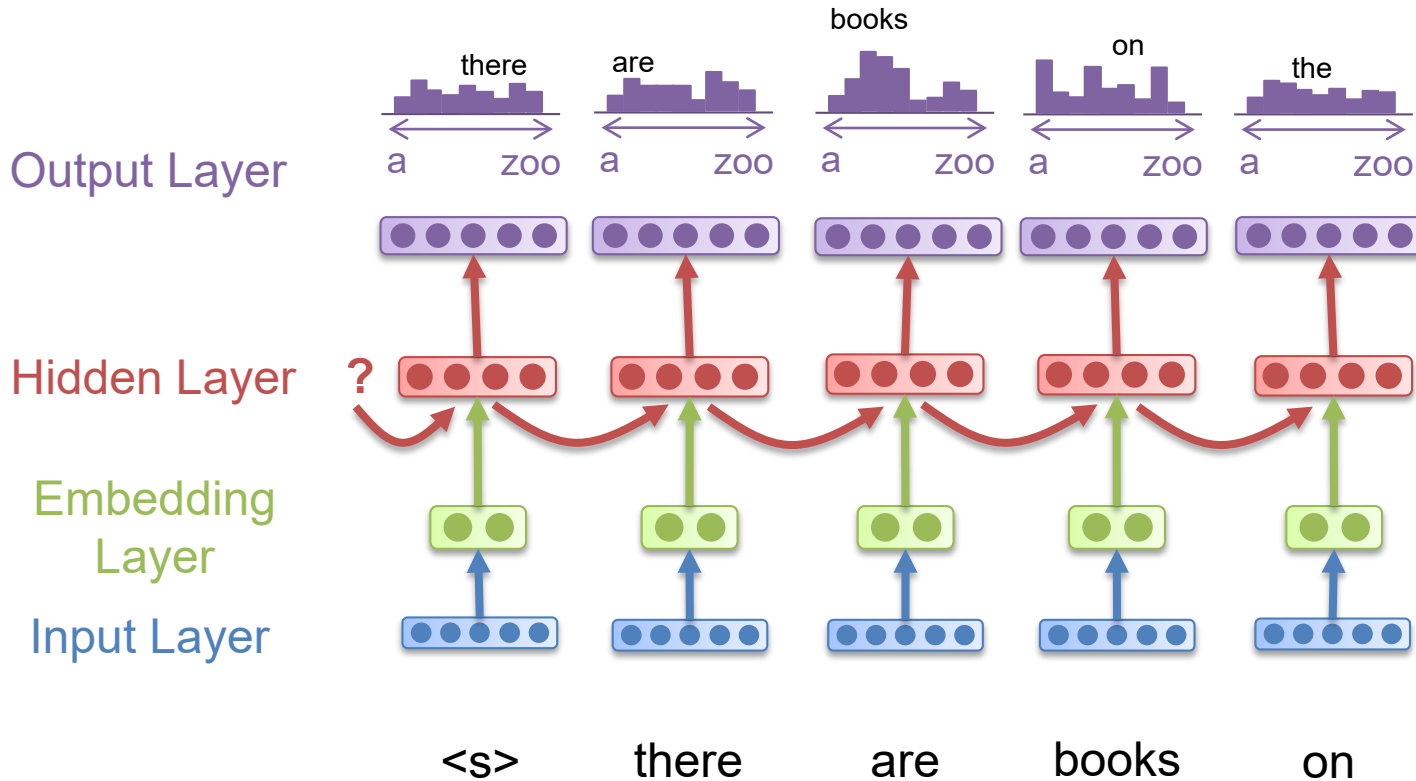$P(\text{there are books on the table})$

$\approx P(\text{there})P(\text{are}|\text{there})P(\text{books}|\text{there are})\textcolor{red}{P(\text{on}|\text{are books})}$

$P(\text{the}|\text{books on})P(\text{table}|\text{on the})$

# Revisiting NNLM – 2.....

$P(there\ are\ books\ on\ the\ table)$
$\approx P(there)P(are|there)P(books|there\ are)P(on|are\ books)$
$\textbf{P(the|books\ on)}P(table|on\ the)$

# Revisiting NNLM – 2......

$P(there\ are\ books\ on\ the\ table)$
$\approx P(there)P(are|there)P(books|there\ are)P(on|are\ books)$
$P(the|books\ on)\textbf{\textit{P(table|on\ the)}}$

# From NNLM to RNNLM

- The hidden state can encapsulate the information of word usage (ordering)
  - Leverage the information!!

# RNNLM.

$P(there\ are\ books\ on\ the\ table)$
$\approx \textbf{\textit{P}}(\textbf{\textit{there}})P(are|there)P(books|there\ are)P(on|there\ are\ books)$
$P(the|there\ are\ books\ on)P(table|there\ are\ books\ on\ the)$

# RNNLM..

$P(there\ are\ books\ on\ the\ table)$
$\approx P(there)\textcolor{red}{\boldsymbol{P(are|there)}}P(books|there\ are)P(on|there\ are\ books)$
$P(the|there\ are\ books\ on)P(table|there\ are\ books\ on\ the)$

# RNNLM...

$P(there\ are\ books\ on\ the\ table)$
$\approx P(there)P(are|there)\textbf{\textit{P(books|there are)}}P(on|there\ are\ books)$
$P(the|there\ are\ books\ on)P(table|there\ are\ books\ on\ the)$

# RNNLM....

$P(there\ are\ books\ on\ the\ table)$
$\approx P(there)P(are|there)P(books|there\ are)\textcolor{red}{P(on|there\ are\ books)}$
$P(the|there\ are\ books\ on)P(table|there\ are\ books\ on\ the)$

# RNNLM.....

$P(there\ are\ books\ on\ the\ table)$
$\approx P(there)P(are|there)P(books|there\ are)P(on|there\ are\ books)$
$\textcolor{red}{P(the|there\ are\ books\ on)}P(table|there\ are\ books\ on\ the)$

# RNNLM……

$P(there\ are\ books\ on\ the\ table)$
$\approx P(there)P(are|there)P(books|there\ are)P(on|there\ are\ books)$
$P(the|there\ are\ books\ on)\textbf{\textit{P(table|there\ are\ books\ on\ the)}}$

# Recurrent Neural Network LM

- RNNLM has recently emerged as a promising modeling framework for several tasks
  - Both **word usage cues** and **long-span structural information** of word co-occurrence relationships can be take into account naturally

- The limitations of the feed-forward NNLM
  - Need to specify the context length
  - RNN can efficiently represent more complex patterns than shallow NNs

$\mathbf{w}_i$ : input layer

$w_i$

$\mathbf{s}_i$ : hidden layer

$\mathbf{y}_i$ : output layer

U

V

$\mathbf{s}_{i-1}$

# Compared with Topic Modeling

|  | RNNLM | Topic Models |
|---|---|---|
| Local and/or Long-span Information | Both bigram and long-span information | Long-span |
| Capture the Long-span Information | By the Hidden State | By EM Algorithm |
| The Combination Weight between Local and Long-span Information | Automatic Adaptation | Empirical Setting |
| The Importance of Each Word in the History | Automatic Learned | Equal Weight |
| Interpretability | No | Yes |

# Elman & Jordan Types

**Elman Type RNN**
**SimpleRNN**
**Vanilla RNN**

**Jordan Type RNN**

# RNN for LM

# RNN for Tagging

# RNN for Classification – 1

- Forward RNN

# RNN for Classification – 2



Joke

there    are    books    on    the    table

# RNN for Classification – 3

- Backward RNN

# RNN for Classification – 4

- Bi-directional RNN!!



joke

there    are    books    on    the    table

# RNN for Classification – 5

- Multi-Layers RNN

# Long Short-Term Memory (LSTM)
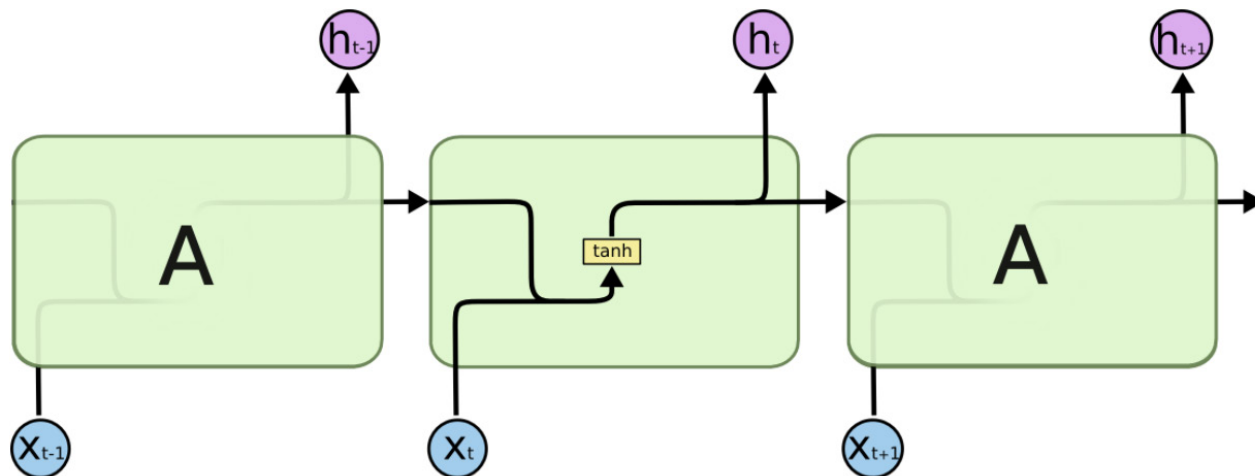
- Learning to Forget!
  - RNN
    - The classic model

  - LSTM
    - Learning to forget
    - Capture longer information
    - Very slow in practice

# Vanilla RNN

- RNN is hard to capture long-term dependencies

$$h_t = \tanh(W\,[h_{t-1}, x_t] + b)$$



| | | | | |
|---|---|---|---|---|
| Neural Network Layer | Pointwise Operation | Vector Transfer | Concatenate | Copy |

35

# LSTM.

- The key to LSTMs is the **cell state**
  - The horizontal line running through the top of the diagram
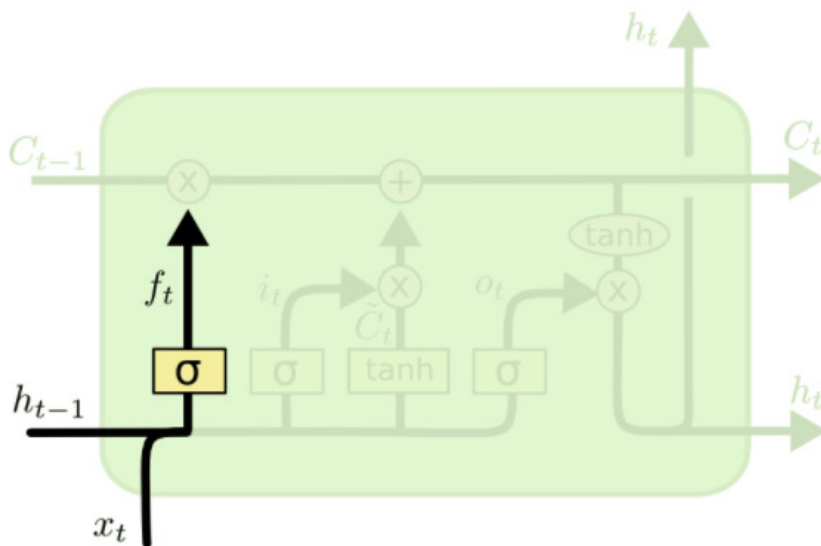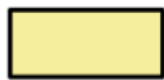  - It's very easy for information to just flow along it **unchanged**



36

# LSTM..

- The **forget gate** is to decide what information we're going to throw away from the cell state
  - $f_t = 1$: completely keep the information
  - $f_t = 0$: completely get rid of the information



$$f_t = \sigma\left(W_f \cdot [h_{t-1}, x_t] + b_f\right)$$

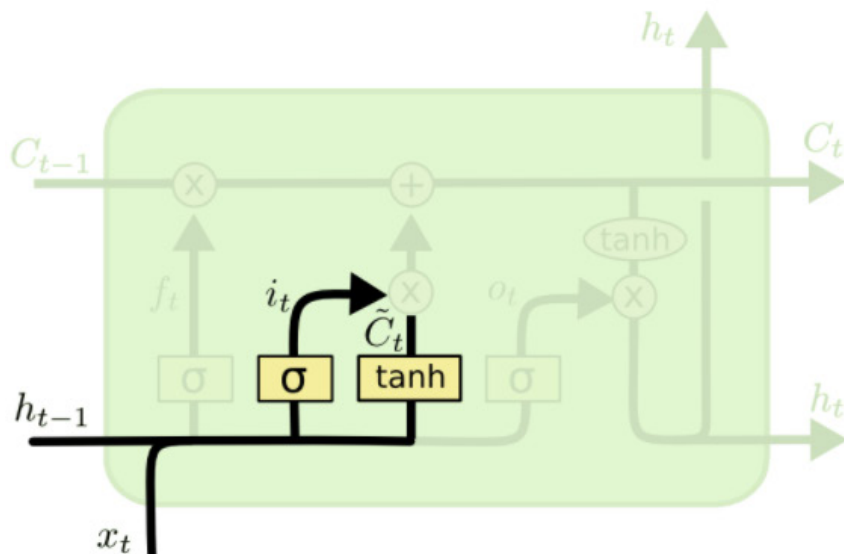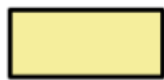| Neural Network Layer | Pointwise Operation | Vector Transfer | Concatenate | Copy |

37

# LSTM...

- The **input gate** is to decide which value we will update
  - A candidate vector, which contains the new information, will also be create



$$i_t = \sigma \left( W_i \cdot [h_{t-1}, x_t] \; + \; b_i \right)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] \; + \; b_C)$$

| Neural Network Layer | Pointwise Operation | Vector Transfer | Concatenate | Copy |

38
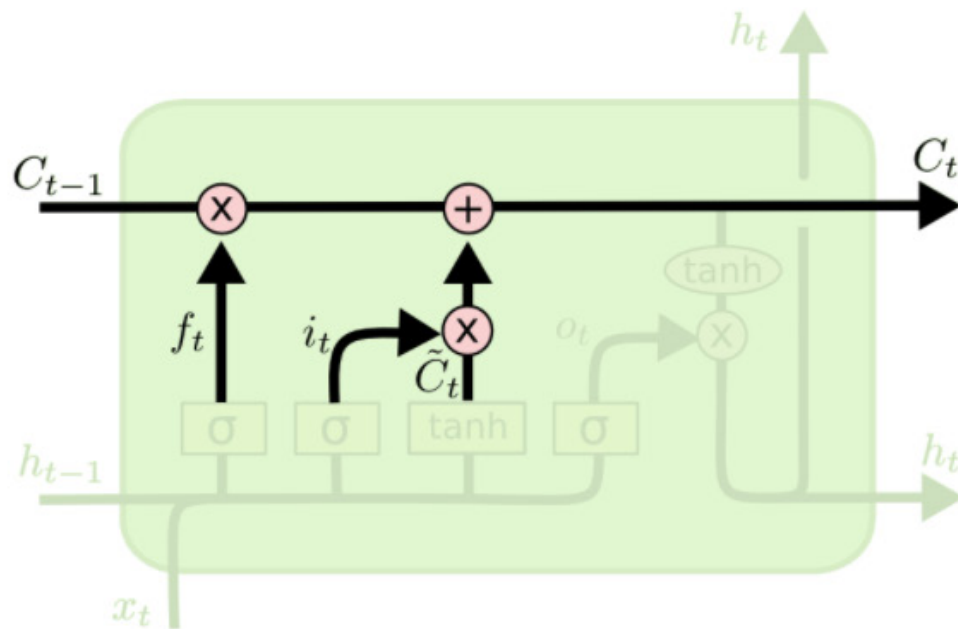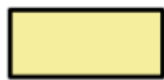
# LSTM....

- Update the cell state!



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

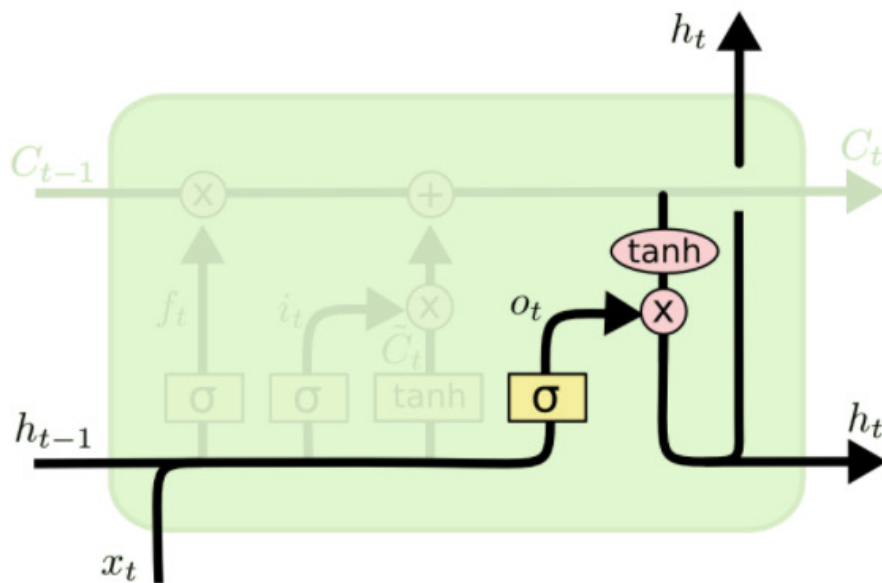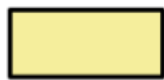| Neural Network Layer | Pointwise Operation | Vector Transfer | Concatenate | Copy |
|---|---|---|---|---|

# LSTM.....

- The **output gate** is a filter to select what information the model is going to output



$$o_t = \sigma\left(W_o\left[h_{t-1}, x_t\right] + b_o\right)$$

$$h_t = o_t * \tanh\left(C_t\right)$$

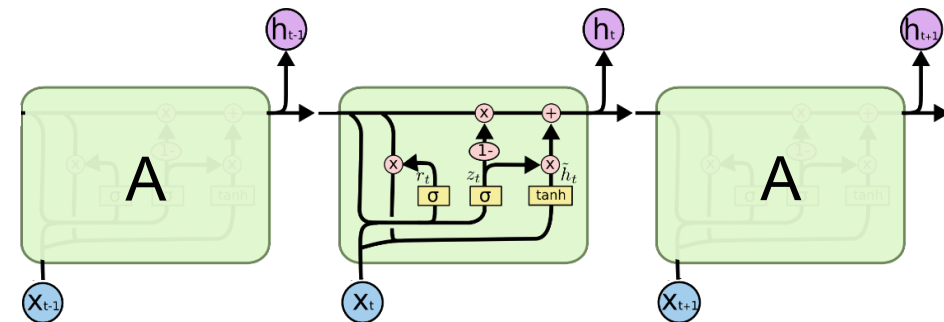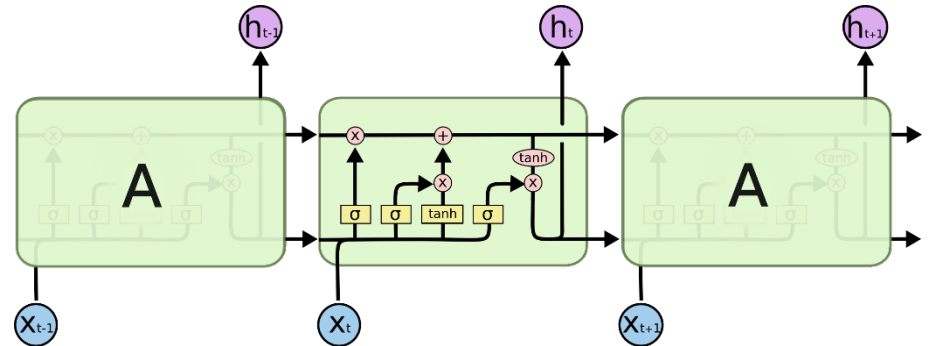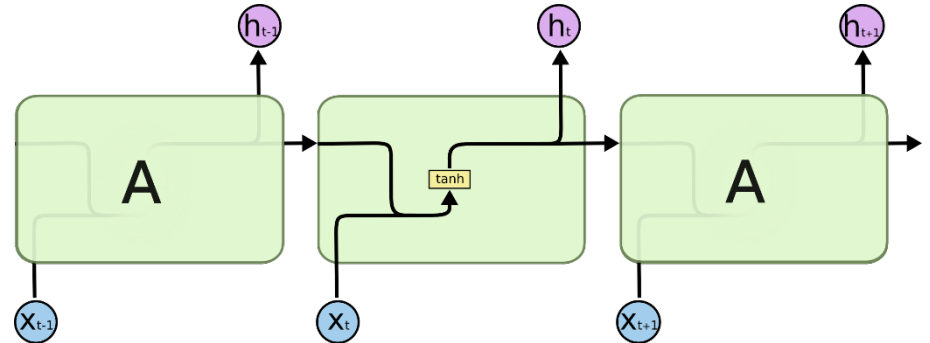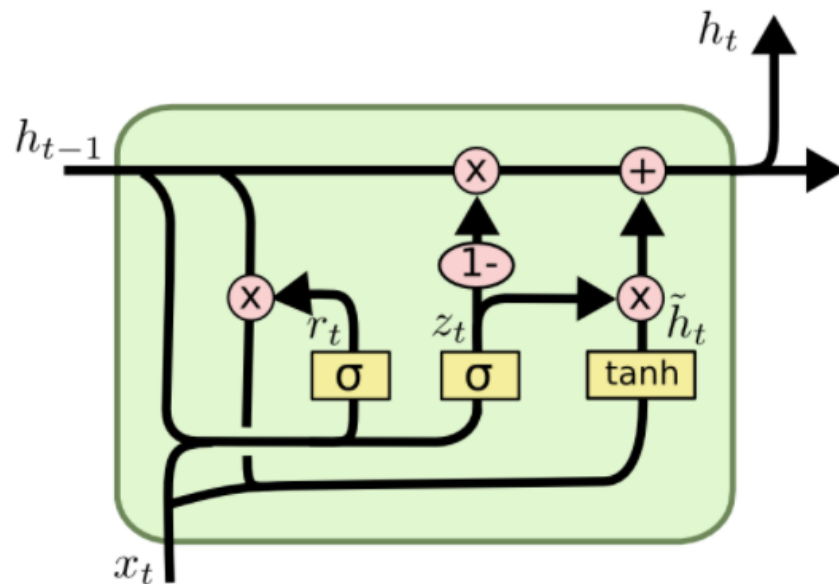| Neural Network Layer | Pointwise Operation | Vector Transfer | Concatenate | Copy |

# Gated Recurrent Unit (GRU)

- RNN
  - The classic model

- LSTM
  - Learning to forget
  - Capture longer information
  - Very slow in practice

- GRU
  - A balanced choice

# GRU

- GRU combines the forget and input gates into a **update gate**
- It also merges the cell state and hidden state, and a **reset gate** is used to control the previous information
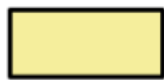
$$z_t = \sigma\left(W_z \cdot [h_{t-1}, x_t]\right) \textbf{ Update Gate}$$

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

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

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

| Neural Network Layer | Pointwise Operation | Vector Transfer | Concatenate | Copy |
|---|---|---|---|---|

42

# Questions?



**kychen@mail.ntust.edu.tw**