# CS480/680
# Lecture 18: July 8, 2019
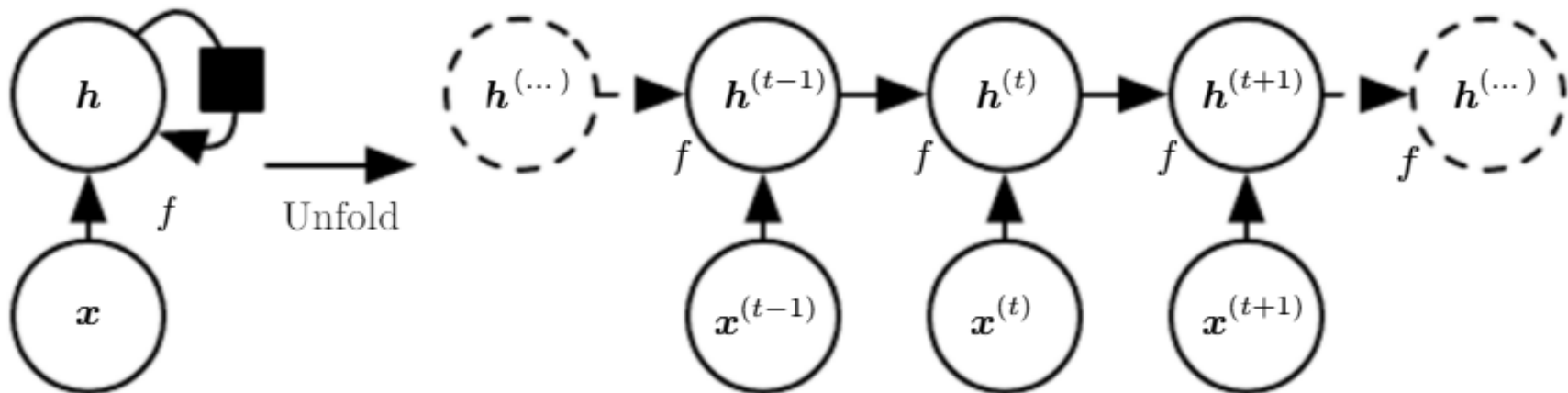
Recurrent and Recursive Neural Networks

[GBC] Chap. 10

# Variable length data

- Traditional feed forward neural networks can only handle fixed length data

- Variable length data (e.g., sequences, time-series, spatial data) leads to a variable # of parameters

- Solutions:
  - Recurrent neural networks
  - Recursive neural networks

# Recurrent Neural Network (RNN)

- In RNNs, outputs can be fed back to the network as inputs, creating a recurrent structure that can be unrolled to handle varying length data.

# Training

- Recurrent neural networks are trained by backpropagation on the unrolled network
  - E.g. backpropagation through time
- Weight sharing:
  - Combine gradients of shared weights into a single gradient
- Challenges:
  - Gradient vanishing (and explosion)
  - Long range memory
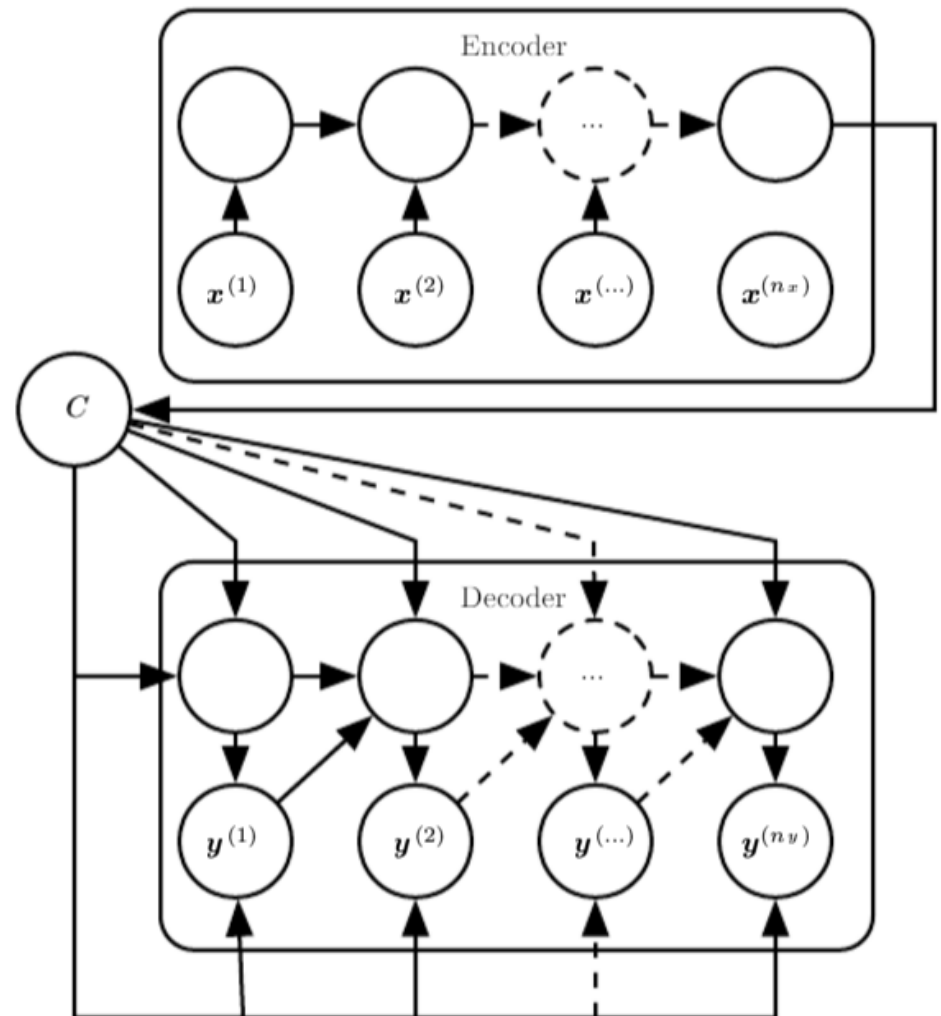  - Prediction drift

# RNN for belief monitoring

- HMM can be simulated and generalized by a RNN

# Bi-Directional RNN

- We can combine past and future evidence in separate chains

# Encoder-Decoder Model

- Also known as sequence2sequence
  - $x^{(i)}$: $i^{th}$ input
  - $y^{(i)}$: $i^{th}$ output
  - $c$: context (embedding)

- Usage:
  - Machine translation
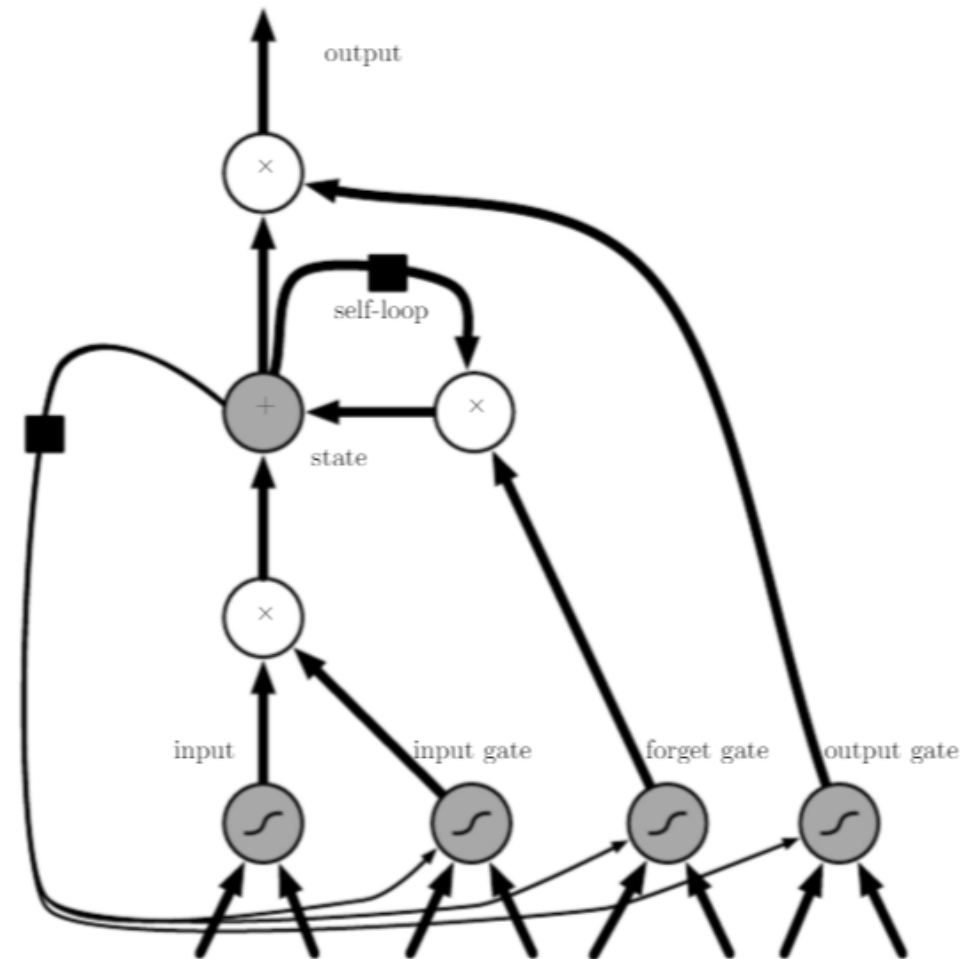  - Question answering
  - Dialog

# Machine Translation

- Cho, van Merrienboer, Gulcehre, Bahdanau, Bougares, Schwenk, Bengio (2014) Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation

| Source | Translation Model | RNN Encoder–Decoder |
|---|---|---|
| at the end of the | [a la fin de la] [f la fin des années] [être sup-primés à la fin de la] | [à la fin du] [à la fin des] [à la fin de la] |
| for the first time | [r © pour la premirère fois] [été donnés pour la première fois] [été commémorée pour la première fois] | [pour la première fois] [pour la première fois ,] [pour la première fois que] |
| in the United States and | [? aux ?tats-Unis et] [été ouvertes aux États-Unis et] [été constatées aux États-Unis et] | [aux Etats-Unis et] [des Etats-Unis et] [des États-Unis et] |
| , as well as | [?s , qu'] [?s , ainsi que] [?re aussi bien que] | [, ainsi qu'] [, ainsi que] [, ainsi que les] |
| one of the most | [?t ?l' un des plus] [?l' un des plus] [être retenue comme un de ses plus] | [l' un des] [le] [un des] |

# Long Short Term Memory (LSTM)

- Special gated structure to control memorization and forgetting in RNNs

- Mitigate gradient vanishing

- Facilitate long term memory

# Unrolled LSTM

- Picture

# LSTM cell in practice

- Adjustments:
  - Hidden state $h_t$ calledcell state $c_t$
  - Output $y_t$ called hidden state $h_t$



- Update equations

Input gate: $i_t = \sigma(W^{(ii)}\bar{x}_t + W^{(hi)}h_{t-1})$
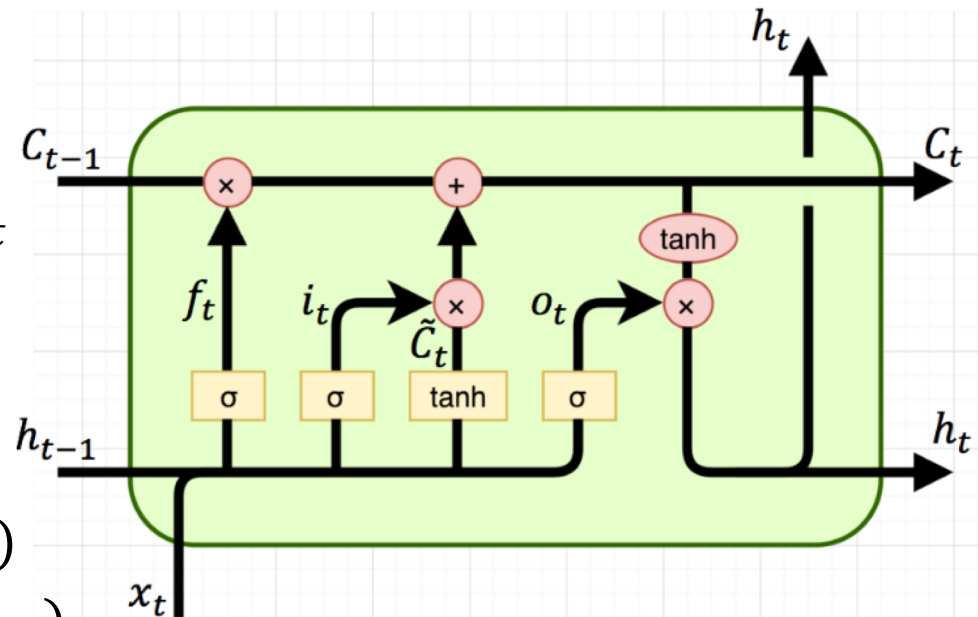
Forget gate: $f_t = \sigma(W^{(if)}\bar{x}_t + W^{(hf)}h_{t-1})$

Output gate: $o_t = \sigma(W^{(io)}\bar{x}_t + W^{(ho)}h_{t-1})$

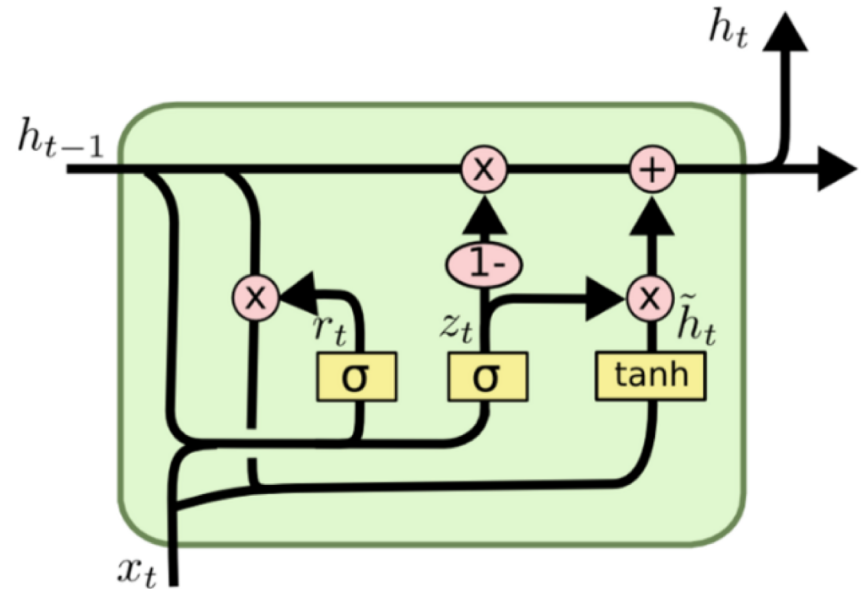Process input: $\tilde{c}_t = \tanh(W^{(i\tilde{c})}\bar{x}_t + W^{(h\tilde{c})}h_{t-1})$

Cell update: $c_t = f_t * c_{t-1} + i_t * \tilde{c}_t$

Output: $y_t = h_t = o_t * \tanh(c_t)$

CS480/680 Spring 2019 Pascal Poupart

# Gated Recurrent Unit (GRU)

- Simplified LSTM
  - No cell state
  - Two gates (instead of three)
  - Fewer weights



- Update equations

Reset gate: $r_t = \sigma(W^{(ir)}\bar{x}_t + W^{(hr)}h_{t-1})$

Update gate: $z_t = \sigma(W^{(iz)}\bar{x}_t + W^{(hz)}h_{t-1})$

Process input: $\tilde{h}_t = \tanh\left(W^{(i\tilde{h})}\bar{x}_t + r_t * \left(W^{(h\tilde{h})}h_{t-1}\right)\right)$

Hidden state update: $h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$
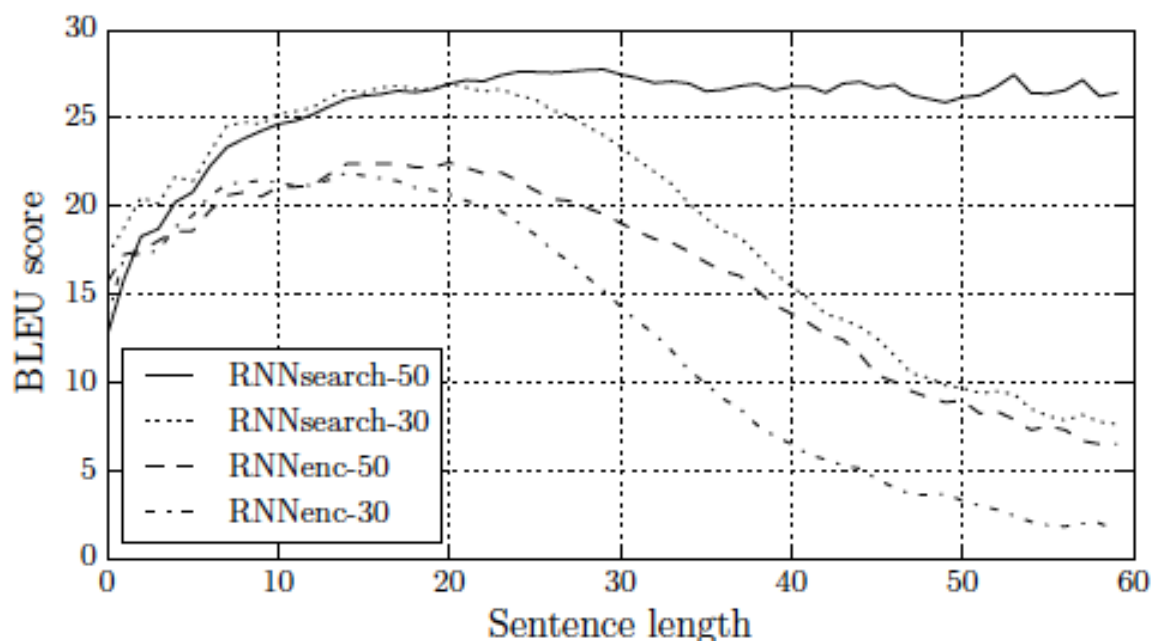
Output: $y_t = h_t$

# Attention

- Mechanism for alignment in machine translation, image captioning, etc.

- Attention in machine translation: align each output word with relevant input words by computing a softmax of the inputs

  - Context vector $c_i$: weighted sum of input encodings $h_j$

    $$c_i = \sum_j a_{ij} h_j$$

  - Where $a_{ij}$ is an alignment weight between input encoding $h_j$ and output encoding $s_i$

    $$a_{ij} = \frac{\exp(alignment(s_{i-1}, h_j))}{\sum_{j'} \exp(alignment(s_{i-1}, h_{j'}))} \text{ (softmax)}$$

  - Alignment example: $alignment(s_{i-1}, h_j) = s_{i-1}^T h_j$

# Attention

- Picture

# Machine Translation with Bidirectional RNNs, LSTM units and attention
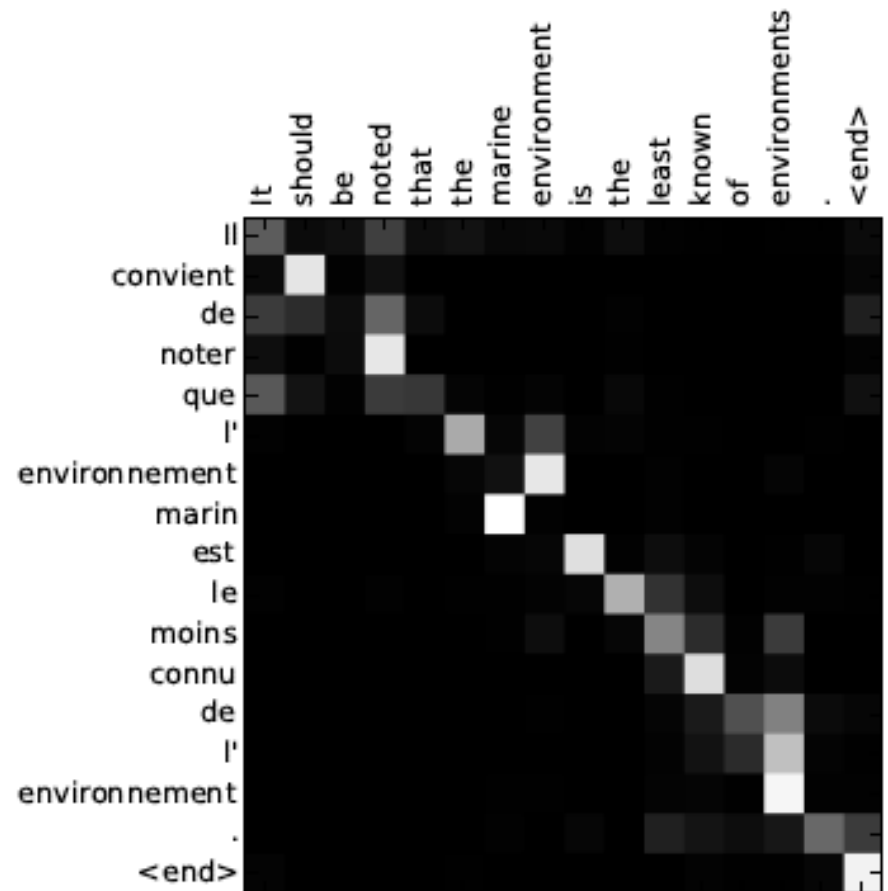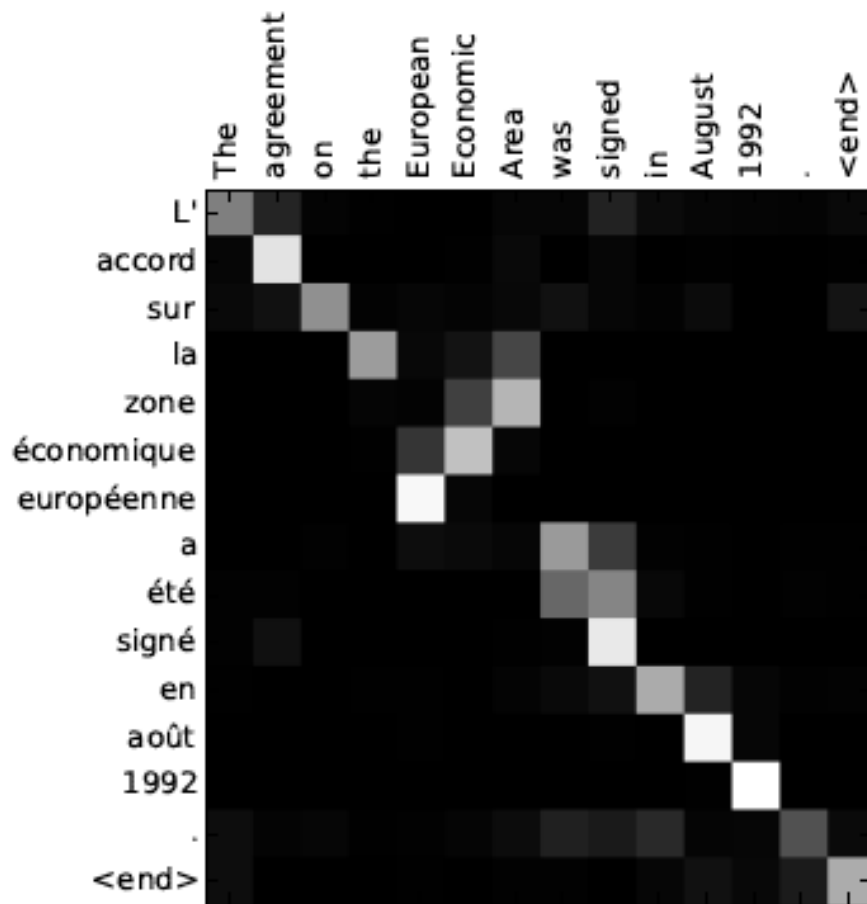
- Bahdanau, Cho, Bengio (ICLR-2015)



RNNsearch: with attention
RNNenc: no attention

- Bleu: BiLingual Evaluation Understudy
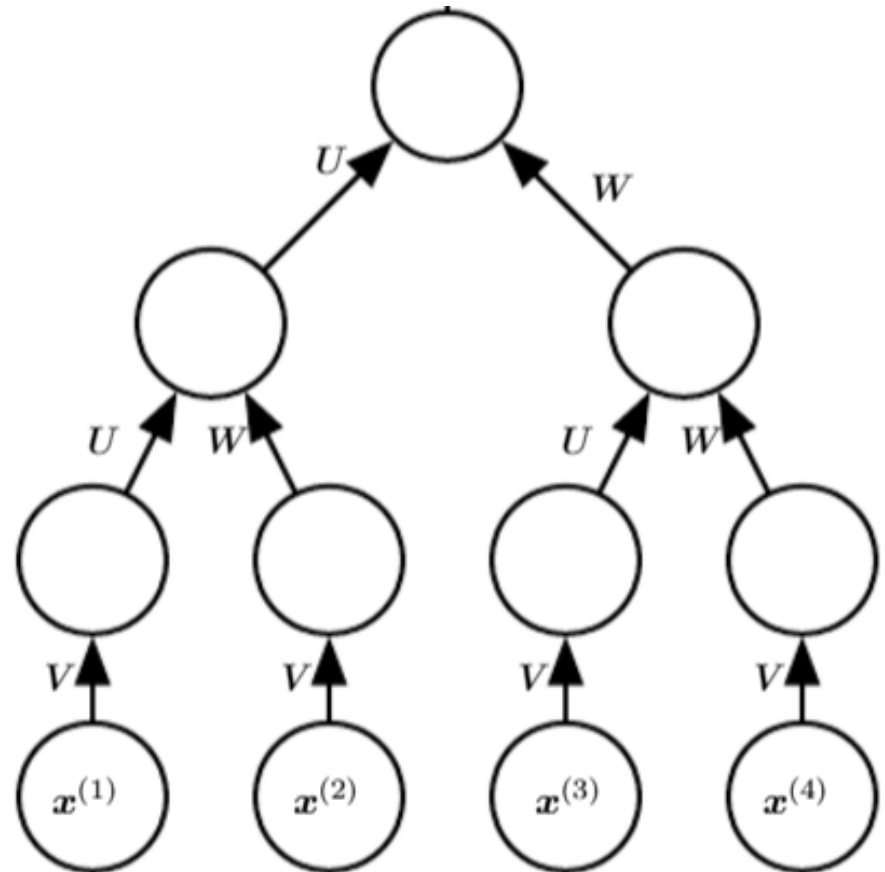  - Percentage of translated words that appear in ground truth

# Alignment example

- Bahdanau, Cho, Bengio (ICLR-2015)

# Recursive Neural network

- Recursive neural networks generalize recurrent neural networks from chains to trees.

- Weight sharing allows trees of different sizes to fit variable length data.

- What structure should the tree follow?

# Example: Semantic Parsing

- Use a parse tree or dependency graph as the structure of the recursive neural network
- Example: