# Generating Text with Deep Reinforcement Learning, H. Guo et al, 2015.

## Lee Won-ho
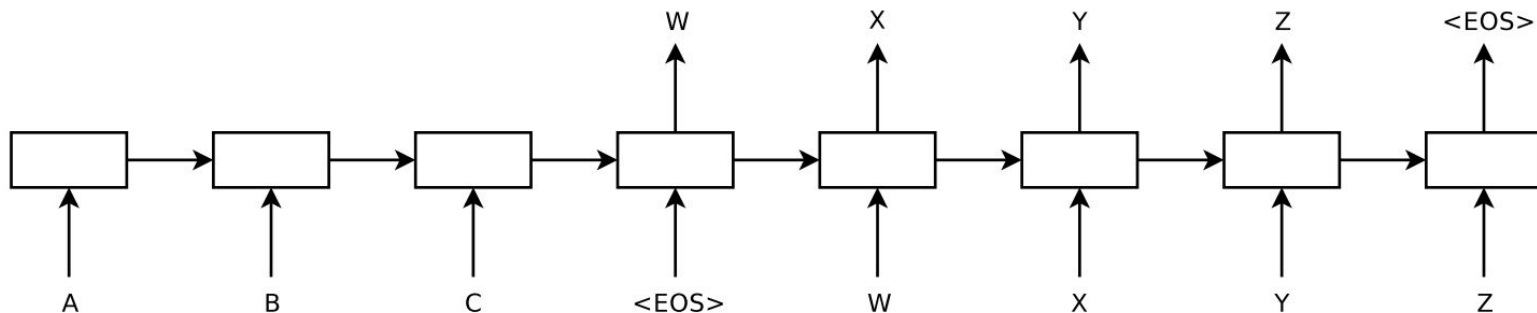gh9908@gmail.com

# Index

- Prerequisites

- Paper main concepts

- Experiment

# Prerequisites

- s2s learning

- Bidirectional LSTM

- Gradient Clipping
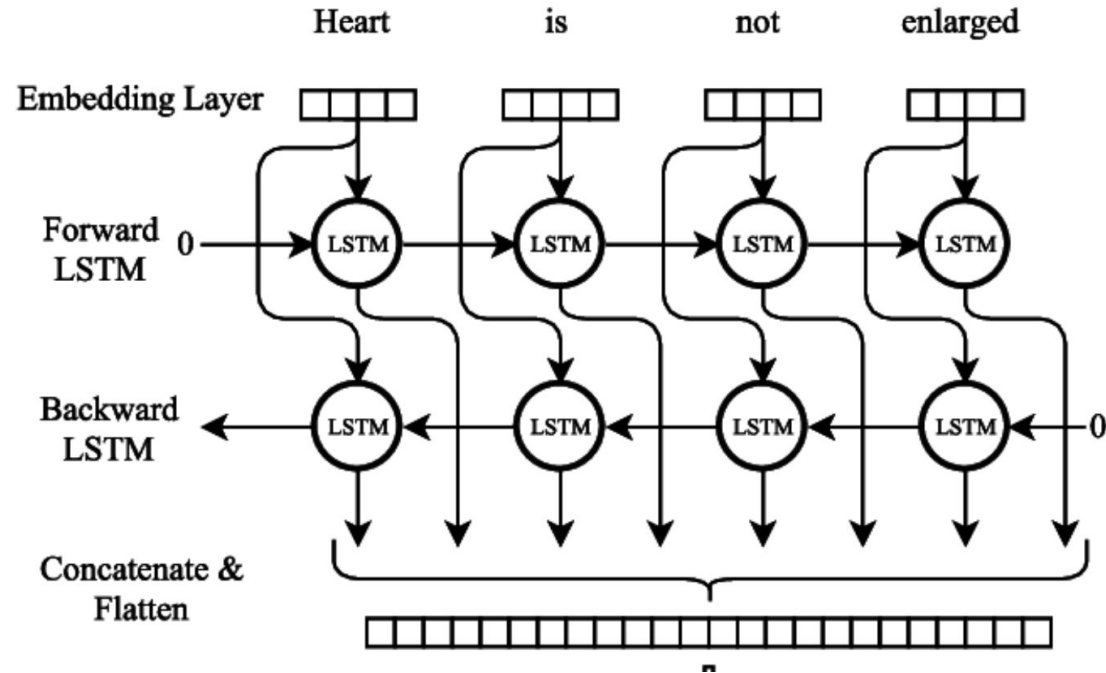
- DQN

- Replay Memory Sampling

# Sequence-to-Sequence Learning

- language model used in sequence prediction tasks

- one LSTM - encoder (read input sequence)

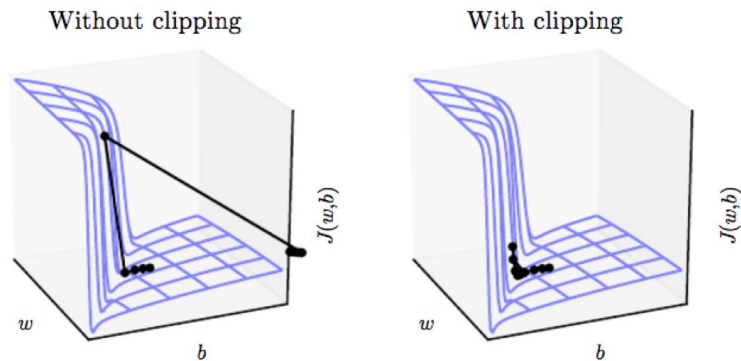- another LSTM - decoder (extract the output sequence)

# Bidirectional LSTM

- sequence processing model that consists of two LSTMs
- one LSTM - forward direction , the other LSTM - backwards direction
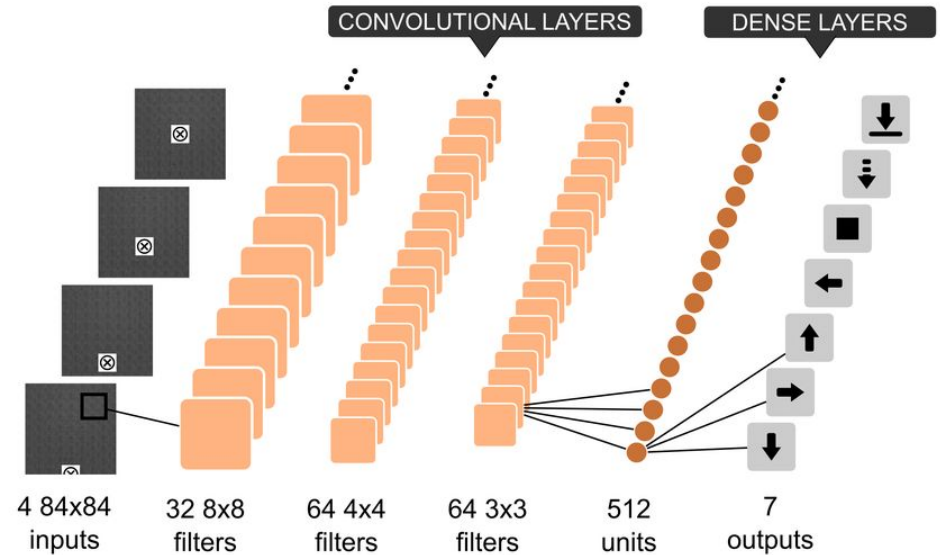- model can know whole information of the sequences

# Gradient Clipping

- two issues with properly training RNN
  - Vanishing Gradient
  - Exploding Gradient

- clip the norm **|| g ||** of the gradient **g** before a parameter update



Without clipping

With clipping

$$\text{if } \|\mathbf{g}\| > v \text{ then } \mathbf{g} \leftarrow \frac{\mathbf{g}^v}{\|\mathbf{g}\|}$$

On the difficulty of training Recurrent Neural Networks, R.Pascanu et al,2013

# DQN

- approximates a state-value function in a Q-Learning framework with a neural network

# Replay Memory Sampling

- store the agent's experiences at each time-step in a data-set pooled over many episodes into a replay memory and then, sample the memory

- Experience Replay
  - randomly sample in memory pool for a minibatch
- Prioritized Experience Replay
  - measured by the magnitude of TD error
  - stochastic prioritization

Prioritized Experience Replay, T.Schaul et al, 2016

# Abstract

- The **aim** here is to enable the **decoder to first tackle easier portions** of the sequences, and then turn to cope with difficult parts

- encoder-decoder is employed to automatically create features to represent the **internal states** of and formulate a list of potential **actions** for the DQN

- The DQN learns to make decision on which action will be selected from the list to **modify the current decoded sequence**.

- compared to a left-to-right greedy beam search LSTM decoder, the proposed method **outperformed the baseline** when decoding unseen sentences

# Algorithm

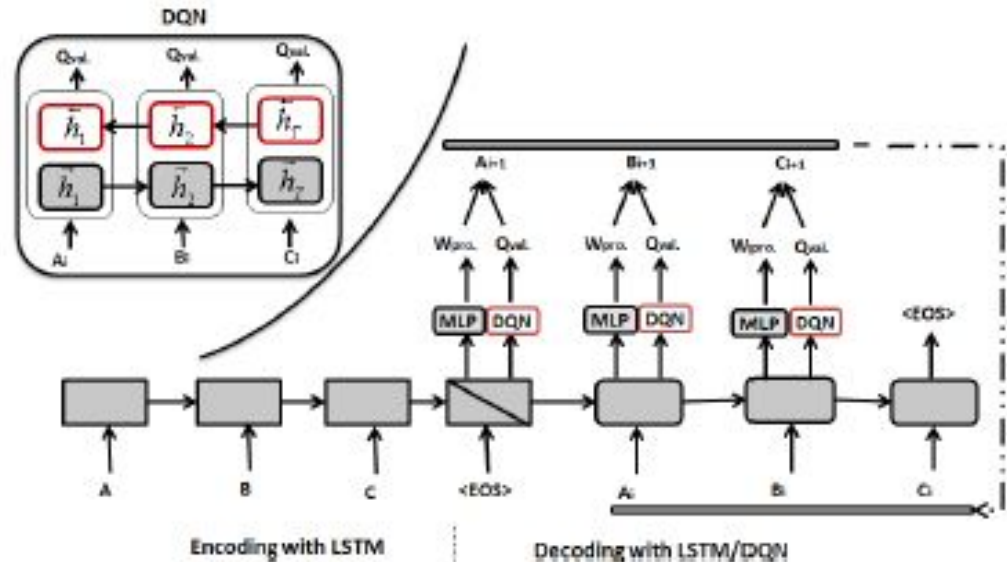**Algorithm 1** Generating Text with Deep Q-Network

1: Initialize replay memory $D$; initialize EnLSTM, DeLSTM, and DQN with random weights
2: **Pretraining Encoder-Decoder LSTMs**
3: **for** epoch = 1,M **do**
4:     randomize given training set with sequence pairs $< X, Y >$.
5:     **for** each sequence pair $EnSen^k \in X$ and $TaSen^k \in Y$ **do**
6:         Encode $EnSen^k$ with EnLSTM, and then predict the next token (e.g., word) in $TaSen^k$ with DeLSTM.
7:     **end for**
8: **end for**
9: **Training Q-value function**
10: **for** epoch = 1,U **do**
11:     **for** each sequence pair $EnSen^k \in X$ and $TaSen^k \in Y$ (with length $l$) **do**
12:         feed $EnSen^k$ into pretrained encoder-decoder LSTMs; obtain the decoded sequence $DeSen_0^k$
13:         **for** iteration i = 1, 2l **do**
14:             **if** random() $< \epsilon$ **then**
15:                 select a random action $a_t$ (e.g., word $w$) at time step $t$ of $DeSen_t^k$ (selection biases to incorrect decoded tokens)
16:             **else**
17:                 compute $Q(s_t, a)$ for all actions using DQN; select $a_t = argmaxQ(S_t, a)$, resulting in a new token $w$ for the $t$-th token in $DeSen_t^k$
18:             **end if**
19:             replace $DeSen_t^k$ with $w$, resulting $DeSen_{t+1}^k$
20:             compute the similarity of $DeSen_{t+1}^k$ and $TaSen^k$, resulting reward score $r_t$
21:             store transition tuple $[s_t, a_t, r_t, s_{t+1}]$ in replay memory $D$; $s_t = [EnSen_t^k, DeSen_t^k]$.
22:             random sample of transition $[s_t, a_t, r_t, s_{t+1}]$ in $D$
23:             **if** $r_t > \sigma$ (preset BLEU score threshold) **then**
24:                 $q_t = r_t$; current sequence decoding successfully complete.
25:             **else**
26:                 $q_t = r_t + \lambda max_a, Q(s', a'; \theta_{t-1})$
27:             **end if**
28:             perform gradient descent step on only the DQN network $(q_t - Q(s, a; \theta_t))^2$
29:         **end for**
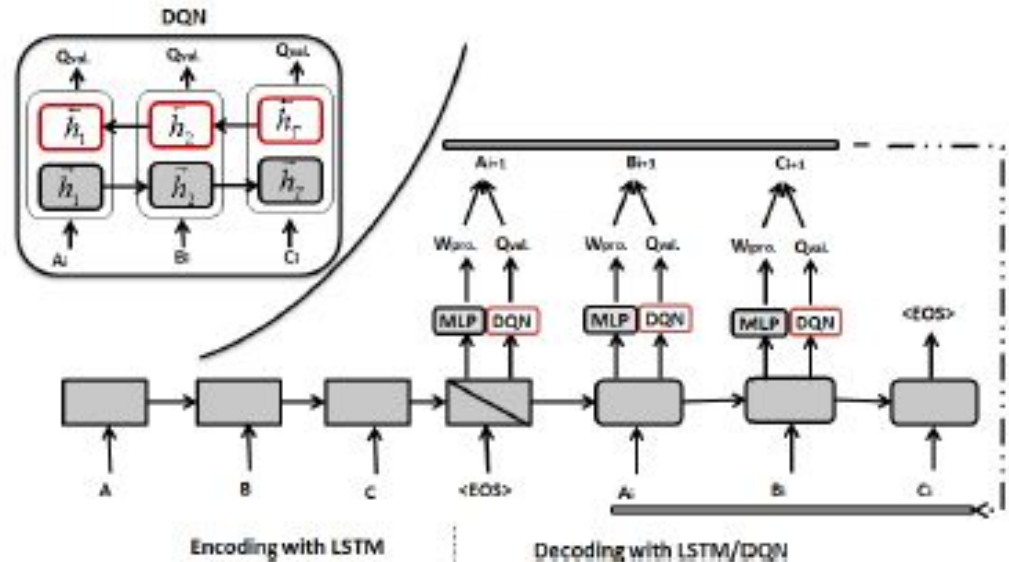30:     **end for**
31: **end for**

# StateGF (State Generation Function)

- The encoder-decoder LSTM network is depicted as gray-filled rectangles in Figure 1. For descriptive purpose, we named this State Generation Function (denoted as StateGF) under the context of DQN
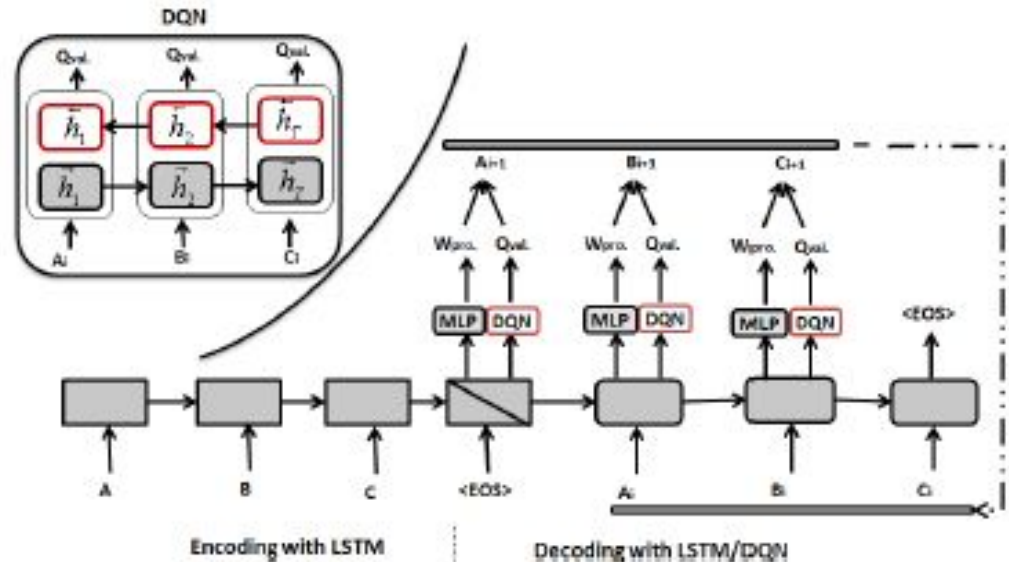
# StateGF (State Generation Function)

- encode the sequence using one LSTM (denoted as EnLSTM), reading into the tokens one timestep at a time

- this encode process results in a fixed dimensional vector representation for the whole sentence $h_N^{en}$



Encoding with LSTM | Decoding with LSTM/DQN

# StateGF (State Generation Function)

- the resulted vector is used as the initial state of another LSTM (denoted as DeLSTM) for decoding to generate the target sequence
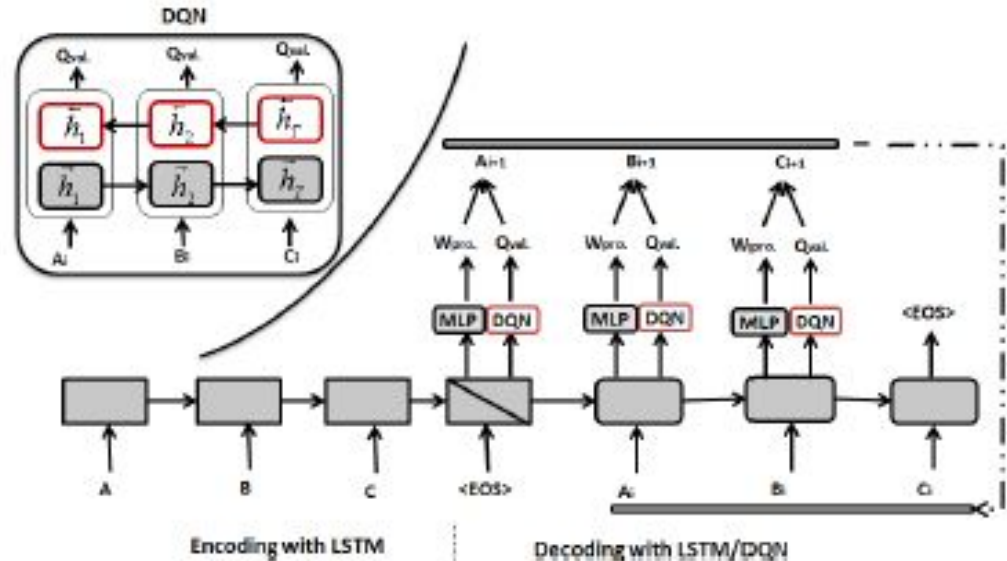
# StateGF (State Generation Function)

- the DeLSTM creates a sequence of hidden state for each time step

- each of these hidden vectors is fed into a Softmax function to produce a distribution over the C possible classes (e.g., words in a vocabulary or dictionary)

$$P(W_{pro}^t = c | EnSen, \vartheta) = \frac{exp(w_c^T h_t^{de})}{\sum_{c=1}^{C} exp(w_c^T h_t^{de})}$$

$$1/|S| \sum_{(X,Y) \in S} log p(Y|X)$$

$$\hat{Y} = \underset{Y}{\mathrm{argmax}} \, p(Y|X)$$



Encoding with LSTM · Decoding with LSTM/DQN

# BLEU Score for DQN Reward

- Reward is calculated based on the closeness between the target sentence and the decoded output sentence after the DQN takes an action

- measure the score difference between the current iteration and the previous iteration

# Empirical Observations on Model Design

- Separating State Generation Function from DQN

- Pre-training the State Generation Function

- Updating with Replay Memory Sampling

- Importance of Supervised Softmax Signal

- Simultaneously Updating with Both Softmax and Q-value

# Separating StateGF from DQN

- Using DQN to approximate the Q-value function equals to train a network against moving targets because here the network's targets depend on the network itself

- Suppose, for a given input feed, the StateGF would generate a different output sequence each time for the DQN

- the DQN network has to also deal with a moving state function involving text with ***very high dimensionality***

# Pre-training the StateGF

- two empirical techniques are employed to ensure that we have a deterministic network for generating states for DQN

    - pre-training StateGF

    - iteratively DeSen is used by the DeLSTM as next input

# Updating with Replay Memory Sampling

- Our studies also indicate that, performing updates to the Q-value function using transitions from the current training sentence causes the network to stronglty **overfit the current input sentence**.

- Transition Tuple

$$[(EnSen_i, DeSen_i), \widehat{y_t^i}, r_i, ([EnSen_i, DeSen_{i+1}])$$

# Importance of Supervised Softmax Signal

- the whole network, including the LSTMs and DQN, only receive the error signals from the Q-value predictions

- **without the supervised signal** the DQN was very **difficult to learn.**

  - moving StateGF and a moving Q-value target function

  - word probability list for each output of the DeLSTM is changing and unreliable

# Simultaneously Updating with Both Softmax and Q-value Error

- If during training the DQN, we not only update the DQN as discussed previously, but also update the state generation functions.

- We found that the network could be **easily bias to the state generation functions** since the Softmax error signal is very strong and more reliable

- we could bias towards the learning of DQN, but this would introduce one more **tricky parameter for tuning**, then we have an **indeterministic** state generation function again

# Dataset

- randomly select 12,000 sentences, with max length of 30, from the Billion Word Corpus

- for train  10,000 sentences

- for validation 1,000 sentences

- for test 1,000 sentence from training set, 1,000 unseen sentences

# Training and Testing Detail

- For each sentence with length of $l$ , we allow DQN to edit the sentence with $2l$ iterations

- Since the maximal length of a sentence in our experiment is 30, the DQN has at most 31 output nodes.

- the DQN can choose one of the 30 top words, each corresponding to a time step at the DeLSTM, as its action, or take the 31st action which indicates not modification is needed
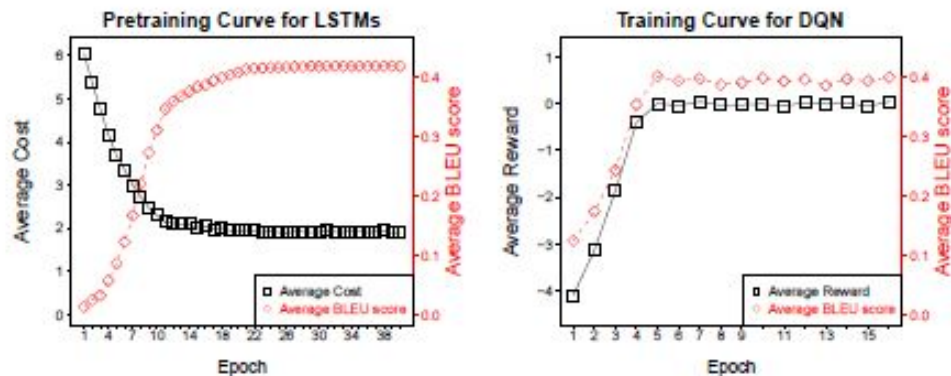
# Experimental Results



Figure 2: The evolution of cost for training the StateGF and reward for training the DQN.

| Testing systems | LSTM decoder | DQN |
|---|---|---|
| Average SmoothedBLEU on sentences IN the training set | 0.425 | 0.494 |
| Average SmoothedBLEU on sentences NOT in the training set | 0.107 | 0.228 |

Table 1: Experimental results for decoding the seen and unseen sentences in testing.

# Experimental Results



**Effect of DQN Exploration**