

# Deletion-based Sentence Compression with Deep Reinforcement Learning

Yang Zhao

2018/2/9

# Overview

- Introduction
- Sentence compression with reinforcement learning
- Experiment
- Result

# Overview

- Introduction
- Sentence compression with reinforcement learning
- Experiment
- Result

# What is Sentence Compression

- Sentence compression aims to
  - use fewer words than the source sentence,
  - retain the most important information from the source sentence,
  - remain grammatical.

- Example:

S: A man suffered a serious head injury after a morning car crash today .

C: A man suffered a injury after a car crash .

# Applications of Sentence Compression

1. Compress text to be displayed on small screens like cellphone. *[Corston-Oliver, 2001]*
2. Generate subtitle for high-rate speech. *[Vandeghinste and Pan, 2004]*
3. Compress lengthy product titles for E-commerce. *[Wang et al., 2018]*

# Points of Sentence Compression

- Sentence compression aims to
  - use fewer words than the source sentence,
  - retain the most important information,
  - remain grammatical.

- 
- Length Constrains
    - Informativeness
      - Readability

S: A man suffered a serious head injury after a morning car crash today .

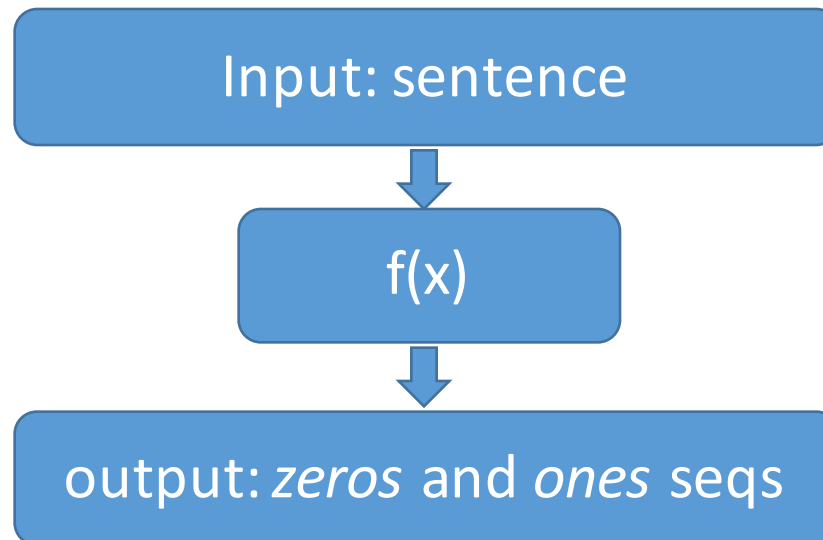
C: A man suffered a injury after a car crash .

# Problem formulation

S: A man suffered a serious head injury after a morning car crash today .

1 1 1 1 0 0 1 1 1 0 1 1 0

- Sequence labeling problem



# Previous works on sentence compression

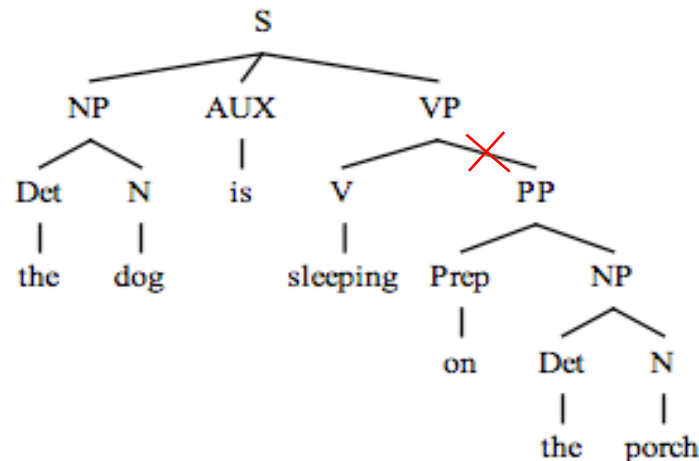
- Tree-pruning based approaches
  - Dependency tree pruning
- Machine-learning based approaches
  - CRF v.s. Recurrent neural network



# Previous works

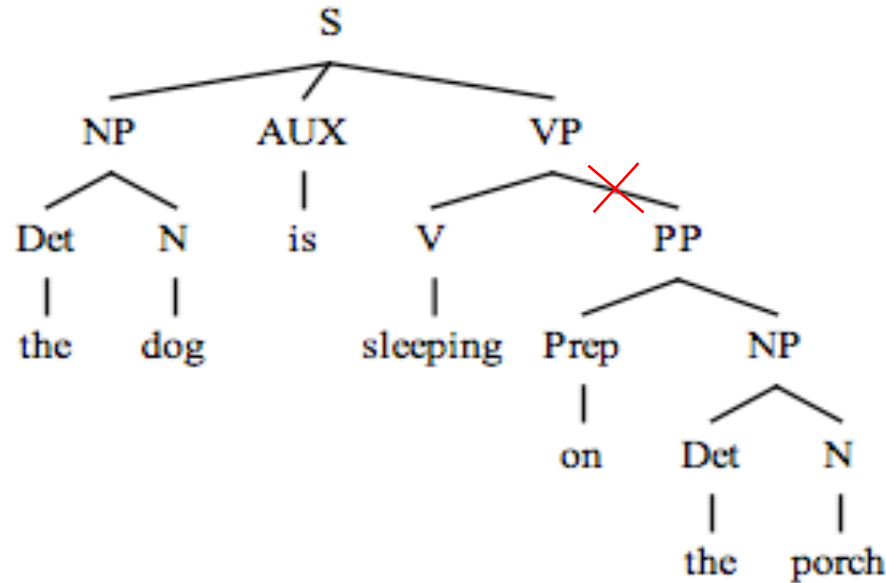
- Previous approaches for sentence compression are mainly rule-based or data-driven.

-Directly prune syntactic trees to generate compression  
(*Knight and Marcu, 2000; Kirkpatrick et al., 2011; Filippova and Altun, 2013; ...*)



*The dog is sleeping on the porch.*

# Problem with rule-based compression

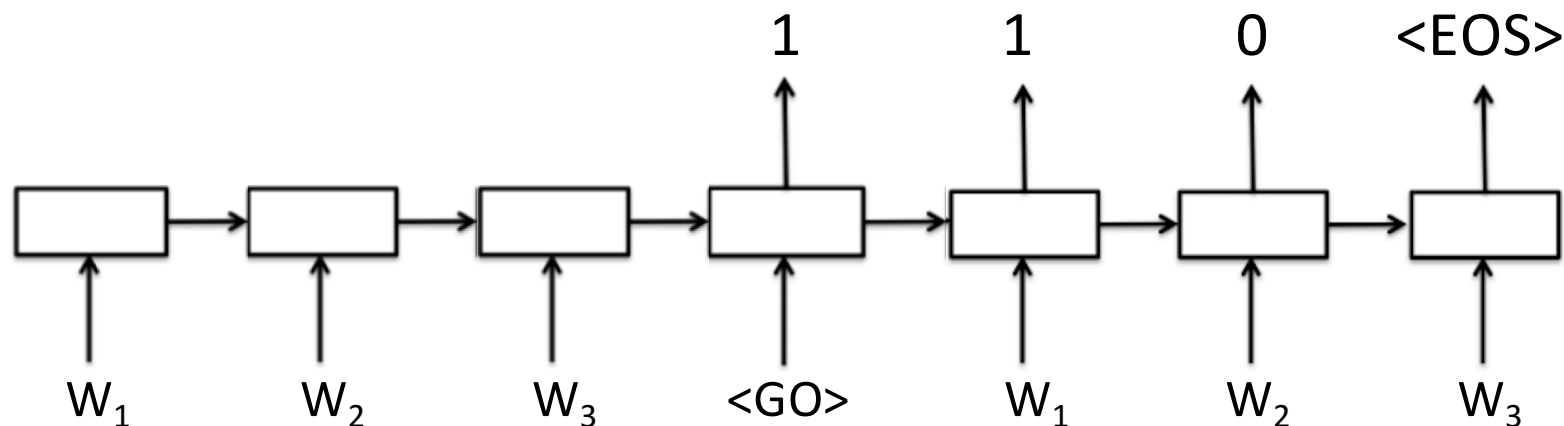


- If syntactic parse trees are incorrect, deleting sub-trees could yield wrong compression.

# Previous works

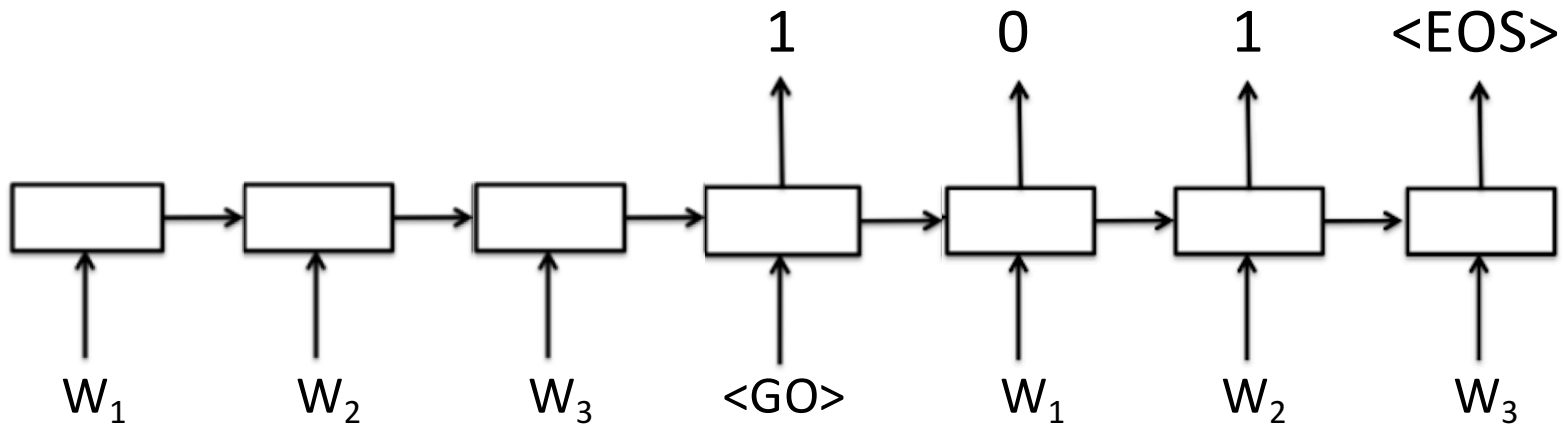
Input: The dog is sleeping on the porch

Output: 1 1 1 1 0 0 0



- Treat sentence compression as a sequence labeling problem (label “1”: kept; label “0”: removed ) (*Filippova et al., 2015; Klerke et al., 2016; ...* )

# Problem with data-driven compression



Problem: it is not able to consider the whole predicted compression as a whole

# Overview

- Introduction
- Sentence compression with reinforcement learning (RL)
- Experiment
- Results

# What the RL brings

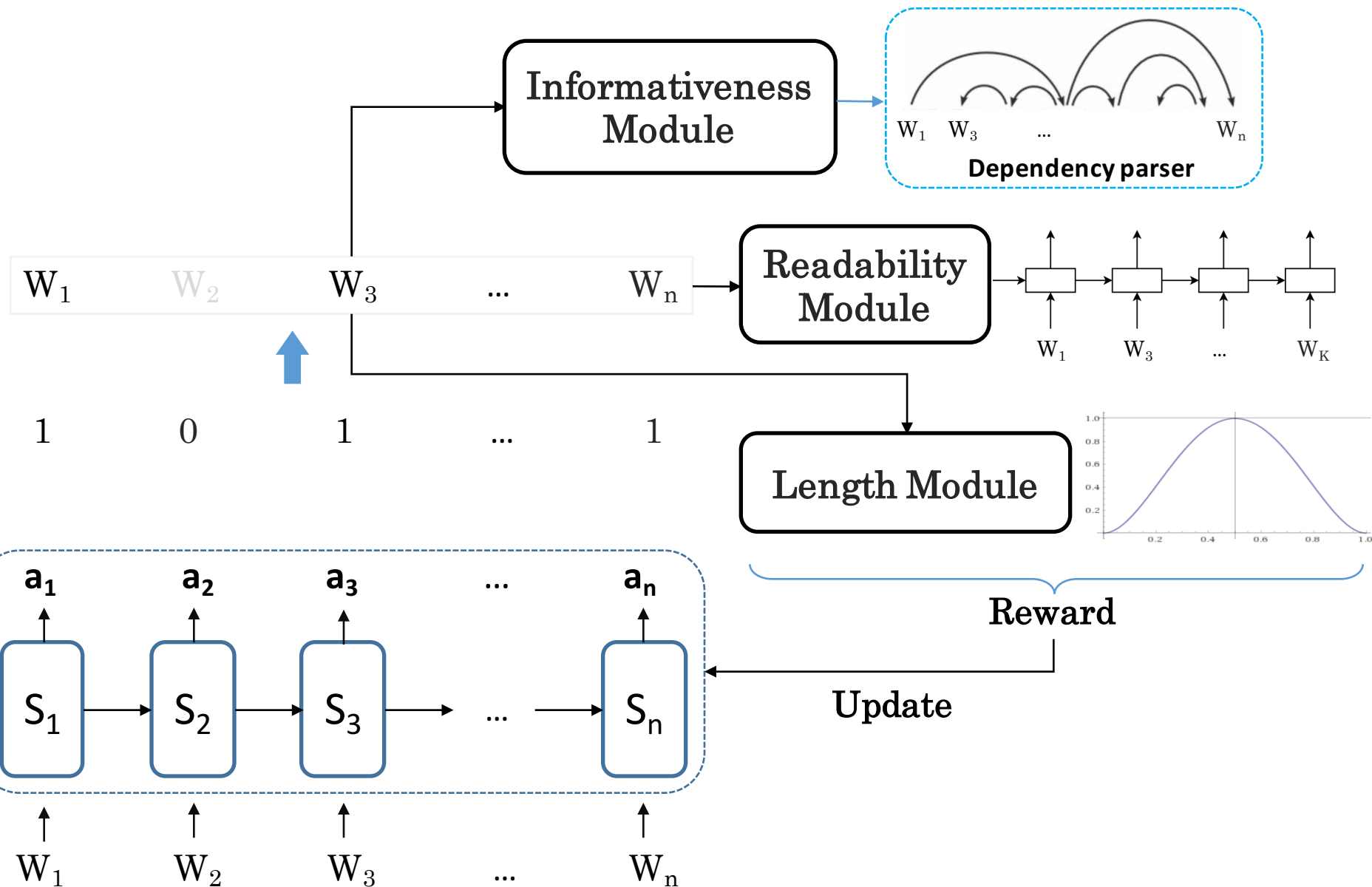
- Previous models are optimized to maximize the likelihood of training data, which may not well match the evaluation metrics that actually quantify the compression quality.
- Automatic evaluations in Sentence compression

Gold	Aaron donald won the 2013 bronko nagurski trophy .
1	(ncsubj, win+ed : 3_VVD, donald: 2_NP1)
2	(dobj, win+ed : 3_VVD, trophy: 8_NN1)
3	(det, trophy: 8_NN1, the: 4_AT)
4	(ncmod, trophy: 8_NN1, 2013: 5_MC)
5	(ncmod, trophy: 8_NN1, bronko: 6_JJ)
6	(ncmod, trophy: 8_NN1, nagurski: 7_JJ)
7	(ncmod, donald: 2_NP1, Aaron: 1_NP1)
System	Aaron donald won .
1	(ncsubj, win+ed : 3_VVD, donald : 2_NP1)
2	(ncmod, donald: 2_NP1, Aaron : 1_NP1)

- $G$  = a set of grammatical relations in ground truth
- $S$  = a set of grammatical relations in system output
- $$F_1 = \frac{2|G \cap S|}{|G| + |S|}$$

*Grammatical relations yielded by dependency parser*

# Model Overview



# Key Modules of the Model

- Informativeness Module
- Readability Module
- Length Module

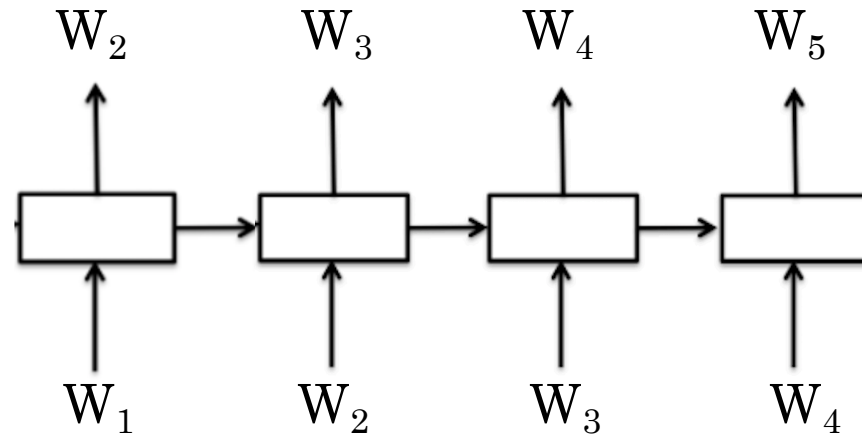


# Informativeness Module

- During training, when we apply dependency parser to both compressed sentence (S) and ground truth (G), two sets of grammatical relations would be yielded respectively
  - $G$  = a set of grammatical relations in ground truth
  - $S$  = a set of grammatical relations in system output
  - $F_1 = \frac{2|G \cap S|}{|G| + |S|}$

$F_1$  as a reward is to feed into policy network during training

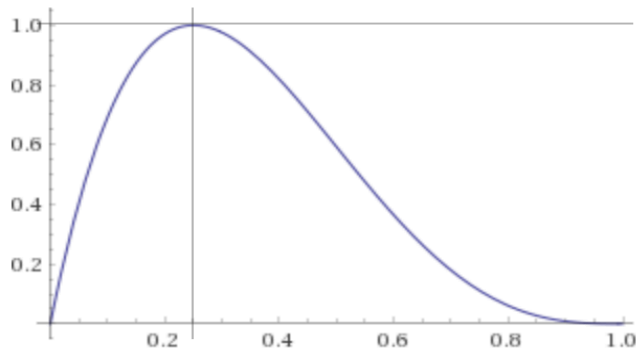
# Readability Module



$$R_{LM} = \exp \left( \frac{1}{|\hat{Y}|} \sum_{i=1}^{|\hat{Y}|} \log P_{LM}(\hat{y}_i | \hat{y}_{0:i-1}) \right)$$

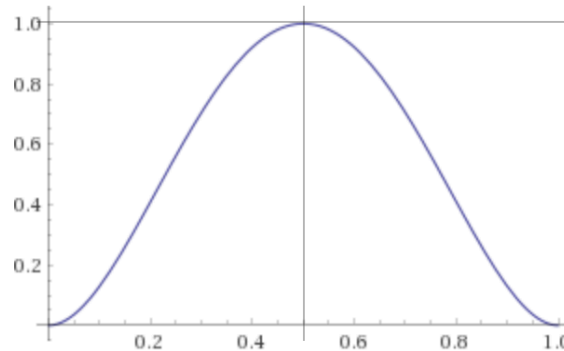
- It is the normalized sentence probability assigned by an LSTM language model trained on sentence compression corpus

# Length module - Compression Rate Reward Inspired by Beta Function



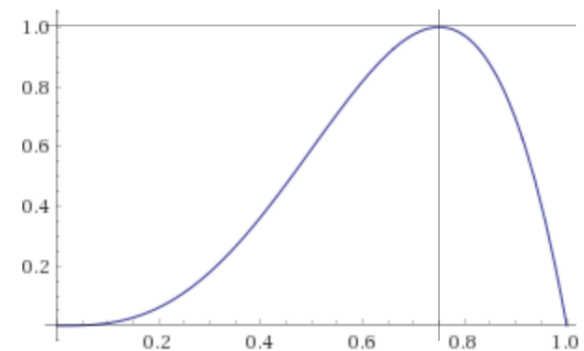
0.25

$$f(x) = \frac{2^8}{3^3} x^1 (1-x)^3$$



0.5

$$f(x) = 16x^2(1-x)^2$$



0.75

$$f(x) = \frac{2^8}{3^3} x^3 (1-x)^1$$

- Experimental results show that model can get to our specifying compression rate very quickly.

# Overview

- Introduction
- Sentence compression with reinforcement learning (RL)
- **Experiment**
- Result

# Dataset

GOOGLE News Dataset	
Domain	NEWS
size	10,000 sentence compression pairs
training/dev/test	8,000/1,000/1,000

# Training details

- Policy network is a two-layer bi-directional LSTM followed by a binary Softmax layer.
- Embedding size is 128.
- Hidden size is 128.
- Since the converge takes time, pre-training is used to speed up the training.

# Baselines

- Integer linear programming (ILP) (Clarke & Lapata, 2008)
- Long short-term memory network (LSTM) (Filippova et. al, 2015)
- ILP+LSTM (Wang et al., 2017)

# Preliminary Result

<b>Google News Dataset</b>	$F_1$	<b>Stanford-parser-<math>F_1</math></b>	<b>compression rate</b>
Original ILP (Clarke Lapata, 08)	56.0	-	0.50
LSTMs (Filippova et al., 15)	80.1	78.4	0.49
LSTMs+ILP(Wang et al., 17)	75.2	56.7	0.47
Our model	81.6	81.1	0.47

Table 1:  $F_1$  results based on ground-truth ( $F_1$ ) and grammatical relations ( $Stanf - F_1$ ).



# Next Step

- Conduct experiments on other datasets
- Replace the sequential language model with syntactic language model
- ...

Thank you!

