# Hierarchical Attention Networks for Document Classification
## Zichao Yang el al.

이종진

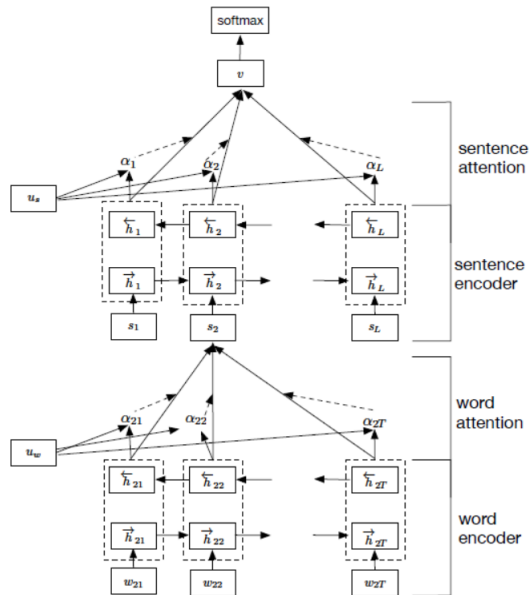**Seoul National University**

*ga0408@snu.ac.kr*

July 06, 2018

# Two characteristics

- ▶ Hierarchical document structure
  - → Word-level GRU and sentence-level GRU
- ▶ Informative words and sentences are different in a document
  - → Two attention mechanism.

# Model

# Word-level

- $w_{ijt}$ : word at time t in jth sentence of ith document
- Word encoder

$$x_{ijt} = W_e w_{ijt}$$
$$h_{ijt} = [\vec{GRU}(x_{ijt}), \overleftarrow{GRU}(x_{ijt})] \tag{1}$$

- Word Attention

$$u_{ijt} = tanh(W_w h_{ijt} + b_w)$$
$$\alpha_{ijt} = \frac{exp(u_{ijt}^T w_w)}{\sum exp(u_{ijt}^T w_w)} \tag{2}$$
$$s_{ij} = \sum \alpha_{ijt} h_{ijt}$$

- ▶ Sentence encoder

$$h_{ij} = [\vec{GRU}(s_{ij}), \overleftarrow{GRU}(s_{ij})] \tag{3}$$

- ▶ Sentence Attention

$$u_{ij} = tanh(W_w h_{ij} + b_w)$$
$$\alpha_{ijt} = \frac{exp(u_{ij}^T w_w)}{\sum exp(u_{ij}^T w_w)} \tag{4}$$
$$d_i = \sum \alpha_{ij} h_{ij}$$

- ▶ Classification

$$\hat{p}_i = softmax(W_c d_i + b_c) \tag{5}$$

- ▶ Loss : Negative log likelihood

$$L = -\sum p_i \log \hat{p}_i = softmax(W_c d_i + b_c) \tag{6}$$

# Experiments

▶ Yelp reviews / IMDB reviews / Yahoo answers / Amazon reviews

| Data set | classes | documents | average #s | max #s | average #w | max #w | vocabulary |
|----------|---------|-----------|------------|--------|------------|--------|------------|
| Yelp 2013 | 5 | 335,018 | 8.9 | 151 | 151.6 | 1184 | 211,245 |
| Yelp 2014 | 5 | 1,125,457 | 9.2 | 151 | 156.9 | 1199 | 476,191 |
| Yelp 2015 | 5 | 1,569,264 | 9.0 | 151 | 151.9 | 1199 | 612,636 |
| IMDB review | 10 | 348,415 | 14.0 | 148 | 325.6 | 2802 | 115,831 |
| Yahoo Answer | 10 | 1,450,000 | 6.4 | 515 | 108.4 | 4002 | 1,554,607 |
| Amazon review | 5 | 3,650,000 | 4.9 | 99 | 91.9 | 596 | 1,919,336 |

**Table 1:** Data statistics: #s denotes the number of sentences (average and maximum per document), #w denotes the number of words (average and maximum per document).
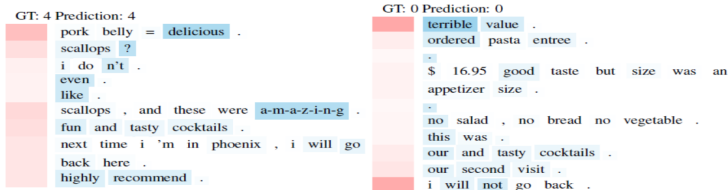
# Experiments

▶ Visualization of attention



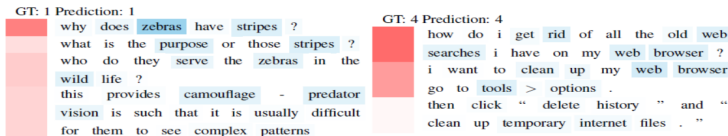**Figure 5:** Documents from Yelp 2013. Label 4 means star 5, label 0 means star 1.



**Figure 6:** Documents from Yahoo Answers. Label 1 denotes Science and Mathematics and label 4 denotes Computers and Internet.