

---

# NEWS HEADLINE-CONTENT GENERATOR: BUILDING SEQ2SEQ ARCHITECTURE USING NUMPY FROM SCRATCH

---

A PREPRINT

**Hao-Chien, Hung**

Institute of Atomic and Molecular Sciences, Academia Sinica  
Taipei, 10617, Taiwan  
hchungdelta@gmail.com

February 20, 2019

## ABSTRACT

The whole machine learning model in this paper is written using Numpy, which is available on Github (see footnote). The basic LSTM seq2seq architecture is adopted in this model, the self-attention mechanism and the layer normalization are applied to enhance the performance and the training efficiency. Furthermore, parallel computing is implemented to speed up training. In this paper, the main task is to test whether the model can be used for news headline-content generation. Only 306k news headline-content data are used for training, however, even with relatively small data, the results are reasonably good and indicate the promising future for further developing.

1

**Keywords** Natural Language Processing · Deep Learning · Attention Mechanism · Layer Normalization

## 1 Introduction

Recently, there is an unprecedented progress in natural language processing (NLP), which attributes to the rapidly developing neural network architecture and the computing power increases in the order of magnitude just within these few decades. The dominant learning model in NLP is the encoder-decoder neural network, which is known as seq2seq model, has a multifaceted capability to perform a variety of tasks, such as machine translation, speech recognition, and text generation.

Nowadays, most of the seq2seq models are using long-short-term memory(LSTM)[1] or gated recurrent unit (GRU)[2] as the basic building block, which is less likely to suffer from loss of information compared to the traditional recurrent neural network(RNN). Subsequently, the introduction of attention mechanisms [3, 4] further boost the performance of seq2seq model. More recently, a new attention mechanism called self-attention have been introduced by Ashish Vaswani *et al.* [5], which a set of key-value pairs are adopted. In this paper, a different version of self-attention is implemented, which will be discussed later.

To train the LSTM seq2seq model is computationally expensive, and it potentially suffers from internal covariate shift as the neural network become deeper, which further slow down the training. Fortunately, layer normalization [6] has been introduced to efficiently solve this problem. Layer normalization normalizes the activities of the neurons and hence mitigates the risk of vanishing gradient problem during training. Compare to batch normalization [7], layer normalization doesn't require a large amount of batch during training and doesn't need to deal with different input sequence length. Therefore, layer normalization is generally considered more suitable to be combined with LSTM-based model rather than batch normalization.

In this paper, the main purpose is to test whether the seq2seq model can be used as a news headline-content generator. That is, input the headline, the model can generate the content only dependent on the headline we input. The paper is

---

<sup>1</sup>Code is available on Github [https://github.com/hchungdelta/Simple\\_NN\\_API/tree/master/NN\\_v.2.1\\_news\\_generator](https://github.com/hchungdelta/Simple_NN_API/tree/master/NN_v.2.1_news_generator)

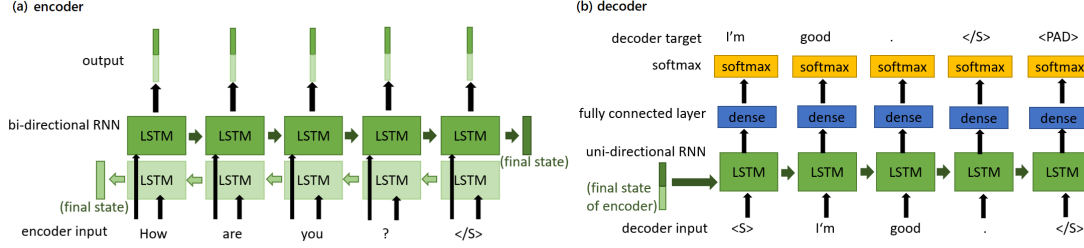


Figure 1: The schematic diagrams of seq2seq model. (a) The encoder, where the Bi-RNN mechanism is used. (b) The decoder, where the uni-RNN mechanism is used.  $\langle S \rangle$ ,  $\langle /S \rangle$ , and  $\langle PAD \rangle$  represent the start token, end token, and padding token, respectively.

organized as follows. Section 2 presents the architecture and the basic algorithm of this model. In section 3, the results and discussions are presented. Summary and conclusion are included in section 4.

## 2 Neural Network Architecture

### 2.1 LSTM

In this paper, the LSTM mechanism [1] is adopted as the basic building blocks of my seq2seq model. The LSTM mechanism is shown below:

$$\text{forget gate} : hf_t = \sigma(W_{fh}h_{t-1} + W_{fx}x_t + b_f) \quad (1)$$

$$\text{input gate} : hi_t = \sigma(W_{ih}h_{t-1} + W_{ix}x_t + b_i) \quad (2)$$

$$\text{output gate} : ho_t = \sigma(W_{oh}h_{t-1} + W_{ox}x_t + b_o) \quad (3)$$

$$\text{candidate output} : hc_t = \tanh(W_{oc}h_{t-1} + W_{cx}x_t + b_c) \quad (4)$$

$$\text{hidden state} : c_t = hf_t * c_{t-1} + hi_t * hc_t \quad (5)$$

$$\text{hidden output} h = ho_t * \tanh(c_t) \quad (6)$$

A number of LSTM form the neural network sequence. In seq2seq model, two sequences are needed, which are known as the encoder and the decoder.

### 2.2 Seq2Seq Model

For the encoder, the Bi-directional RNN [8] is implemented to process the inputs. For the decoder, since future inputs aren't available, uni-directional RNN is used. A simple schematic diagram is depicted in Fig. 1.

### 2.3 Layer Normalization

Layer normalization [6] is implemented to improve the efficiency of training. In the early stage, batch normalization [7] is also adopted to train the model. However, the results indicate that it potentially makes the model less stable, since the lengths of input data are usually different, which make batch normalization less suitable for RNN. Therefore, only layer normalization is used in this paper. The mechanism of layer normalization is shown in the following:

$$\text{mean} : \mu_d = \frac{1}{d} \sum_{i=0}^d x_i \quad (7)$$

$$\text{variance} : \sigma_d^2 = \frac{1}{d} \sum_{i=0}^d (x_i - \mu_d)^2 \quad (8)$$

$$\text{normalization} : \hat{x} = \frac{x - \mu_d}{\sqrt{\sigma_d^2 + \epsilon}} \quad (9)$$

$$\text{scale/shift} : LN(x) = \gamma \hat{x} + \beta \quad (10)$$

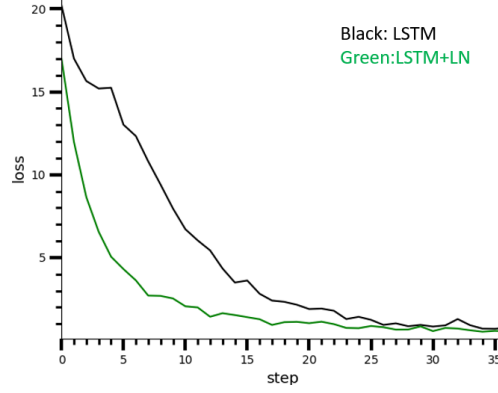


Figure 2: The comparison of the training efficiency of pure LSTM model and LSTM model with layer normalization.

$d$  represents the depth of the input.  $i$  in  $x_i$  represents the  $i - th$  scalar value in this input.  $\gamma$  and  $\beta$  are both vectors to re-scale the input. As the weights and the biases in LSTM, the parameters of these two vectors are shared among the RNN sequence. The backpropagation is similar to those used in batch normalization [7]. As the following expresses:

$$\frac{\partial L}{\partial \gamma} = \sum_{i=0}^d \frac{\partial L}{\partial LN(x)} \cdot \hat{x} \quad (11)$$

$$\frac{\partial L}{\partial \beta} = \sum_{i=0}^d \frac{\partial L}{\partial LN(x)} \quad (12)$$

$$\frac{\partial L}{\partial \hat{x}} = \frac{\partial L}{\partial LN(x)} \cdot \gamma \quad (13)$$

$$\frac{\partial L}{\partial \sigma_d^2} = \sum_{i=0}^d \frac{\partial L}{\partial \hat{x}} \cdot (x - \mu_d) \cdot \frac{-1}{2} (\sigma_d^2 + \epsilon)^{\frac{-3}{2}} \quad (14)$$

$$\frac{\partial L}{\partial \mu_d} = \sum_{i=0}^d \frac{\partial L}{\partial \hat{x}} \cdot \frac{-1}{\sqrt{\sigma_d^2 + \epsilon}} + \frac{\partial L}{\partial \sigma_d^2} \cdot \frac{\sum_{i=0}^d -2(x - \mu_d)}{d} \quad (15)$$

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial \hat{x}} \cdot \frac{1}{\sqrt{\sigma_d^2 + \epsilon}} + \frac{\partial L}{\partial \sigma_d^2} \cdot \frac{2(x - \mu_d)}{d} + \frac{\partial L}{\partial \mu_d} \cdot \frac{1}{d} \quad (16)$$

After applying layer normalization, the LSTM mechanism now should be rewritten as:

$$hf_t = \sigma(LN_{fh}(W_{fh}h_{t-1}) + LN_{fx}(W_{fx}x_t) + b_f) \quad (17)$$

$$hi_t = \sigma(LN_{ih}(W_{ih}h_{t-1}) + LN_{ix}(W_{ix}x_t) + b_i) \quad (18)$$

$$ho_t = \sigma(LN_{oh}(W_{oh}h_{t-1}) + LN_{ox}(W_{ox}x_t) + b_o) \quad (19)$$

$$hc_t = \tanh(LN_{ch}(W_{oc}h_{t-1}) + LN_{cx}(W_{cx}x_t) + b_c) \quad (20)$$

$$c_t = hf_t * c_{t-1} + hi_t * hc_t \quad (21)$$

$$h = ho_t * \tanh(LN_o(c_t)) \quad (22)$$

The comparison of training efficiency of the model with and without layer normalization is displayed in Fig. 2.

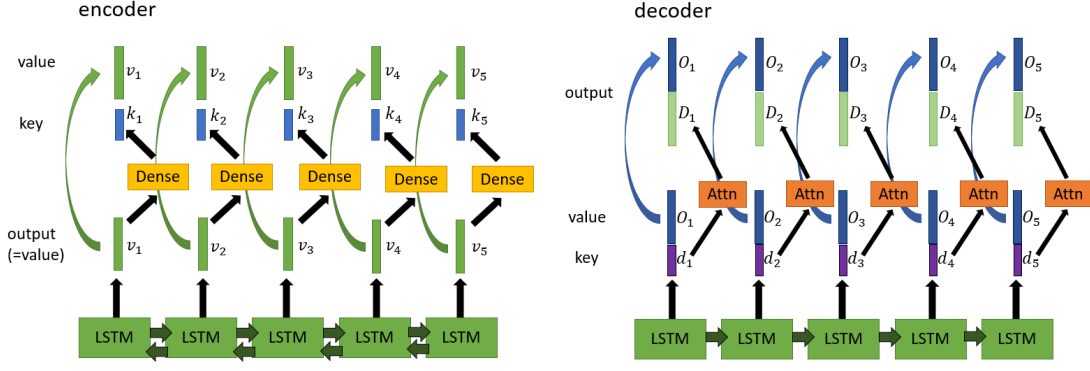


Figure 3: The schematic diagram of the attention mechanism implemented in this paper

## 2.4 Attention mechanism

Dot-product attention mechanism [5] is adopted in this paper, although a little change in algorithm is made. As illustrated in 3. The encoder output and the decoder output are separated into key parts ( $k_i$  for encoder key,  $d_i$  for decoder key) and value parts ( $v_i$  for the encoder value,  $O_i$  for the decoder value). Afterward, the original decoder value and the attention output are concatenated as the final decoder output.

The algorithm is as follows:

$$score_{ij} = s_{ij} = d_i \cdot k_j \quad (23)$$

$$alpha_{ij} = \alpha_{ij} = \frac{e^{s_{ij}}}{\sum_{x=0}^{x_{max}} e^{s_{ix}}} \quad (24)$$

$$attention\ output_i = D_i = \sum_{x=0}^{x_{max}} \alpha_{ix} v_x \quad (25)$$

$$output_i = concat(O_i, D_i) \quad (26)$$

For the backpropagation:

$$\frac{\partial L_i}{\partial v_j} = \frac{\partial L_i}{\partial D_i} \frac{\partial D_i}{\partial v_j} = \frac{\partial L_i}{\partial D_i} \alpha_{ij} \quad (27)$$

$$\frac{\partial L_i}{\partial d_i} = \frac{\partial L_i}{\partial D_i} \frac{\partial D_i}{\partial d_i} = \frac{\partial L_i}{\partial D_i} \sum_{x=0}^{x_{max}} \frac{\partial}{\partial d_i} \alpha_{ix} v_x = \frac{\partial L_i}{\partial D_i} \left[ \sum_{x=0}^{x_{max}} \alpha_{ix} k_x v_x - \sum_{x=0}^{x_{max}} \sum_{y=0}^{x_{max}} \alpha_{ix} \alpha_{iy} k_y v_x \right] \quad (28)$$

$$\frac{\partial L_i}{\partial k_j} = \frac{\partial L_i}{\partial D_i} \frac{\partial D_i}{\partial k_j} = \frac{\partial L_i}{\partial D_i} \frac{\partial}{\partial k_j} \sum_{x=0}^{x_{max}} \alpha_{ix} k_x v_x = \frac{\partial L_i}{\partial D_i} [\delta_{ij} d_i v_j - \alpha_{ix} \alpha_{ij} d_i v_x] \quad (29)$$

In this model, since the encoder training data is much shorter than the decoder, only a single-headed attention mechanism is implemented, rather than multi-headed attention mechanism. Fig. 4 compares the efficiency of the model with and without attention mechanism.

## 2.5 Orthogonal initialization

Rather than randomly generate the weight matrices as initial states, in this paper, the weight matrices are generated using the concept called orthogonal initialization. It can potentially help the learning model recognize the difference between each input.

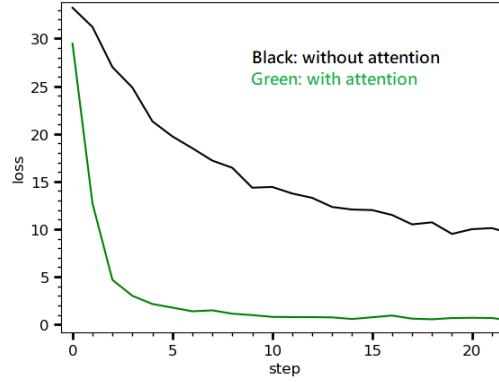


Figure 4: The comparison of the training efficiency of the model with and without attention mechanism.

## 2.6 Parallel Computing

Parallel computing is implemented by mpi4py. In current development, only embarrassing parallel computing is available in this code, since it is effortless to apply to all the learning models without further adjustments. All the nodes share the same weights ( $W$ ) and biases ( $b$ ) throughout the training while exposing to different training data. For instance, if four nodes are used for training and the batch is set to be 32, training data no.1 to no.32 will be delivered to node one, no.33 to no.64 will be delivered to node two, and so on. Hence at each step, the total amount of 128 training data are used for training. Afterward, the master node will sum up all the gradients of weights ( $dW$ ) and biases ( $db$ ) to optimize the model. As the following equations show.

$$dW_{sum} = \sum_{node=1}^{node_{max}} dW_{node} \quad (30)$$

$$db_{sum} = \sum_{node=1}^{node_{max}} db_{node} \quad (31)$$

The summed gradients of weights and biases will be sent to the optimizer for updating the original weights and biases (See Sec. 2.7).

The schematic diagram (Fig. 5 (a)) visualizes the notion of the embarrassing parallel computing in this code. If one wants to take an average of the gradients to avoid overshooting, just simply divided the learning rate by the amount of the nodes. Fig. 5 (b) depicts the loss-step curves for the same learning model while using a different amount of nodes during training. The plot clearly indicates the parallel computing is useful and efficient. Until now, there is no dramatic deterioration in accuracy found when using parallel computing.

## 2.7 Optimizer

Stochastic Gradient Descent (SGD), Momentum, and Adaptive Moment Estimation (Adam) optimizer [9] are available in this code. In this paper, only Adam optimizer is used for training since the efficiency of Adam can outperform the other two approaches nearly in the order of magnitude. The learning rate ( $lr$ ) is set to be 0.0015.  $\epsilon$ ,  $\beta_1$ , and  $\beta_2$  are equal to  $1e^{-8}$ , 0.9, and 0.999, respectively. The following shows how the Adam optimizer works,  $t$  represents the step during training (here only take the gradients of the weights for example, the same goes for biases).

$$dW_{t+1} = dW_t - lr \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon} \quad (32)$$

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}, \quad \hat{v}_t = \frac{v_t}{1 - \beta_2^t} \quad (33)$$

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) dW_t \quad (34)$$

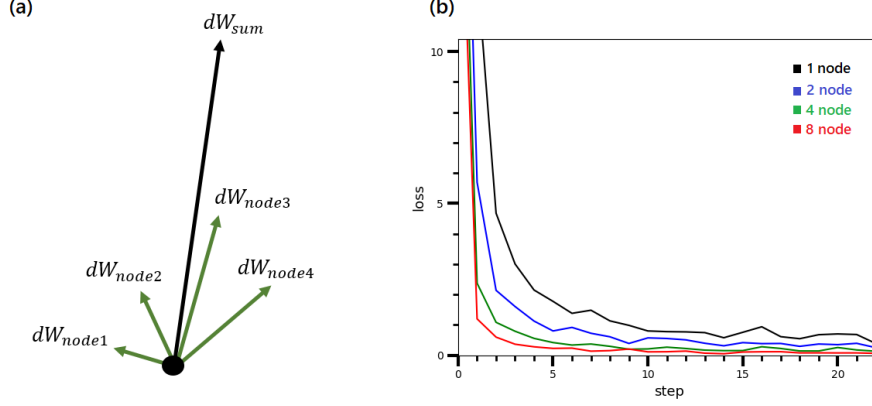


Figure 5: (a) The schematic diagram of embarrassing parallel computing implemented in this model. (b) The comparison of the loss-step curves of the same learning model trained with different amount of nodes.

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) dW_t^2 \quad (35)$$

The gradient clipping [10] is also implemented to help model converge more smoothly. In this paper, the threshold is set to be 1.5. The equation is shown below:

$$dW_{ij} = \begin{cases} threshold, & \text{if } dW_{ij} > threshold \\ dW_{ij}, & \text{otherwise} \\ -threshold, & \text{if } dW_{ij} < -threshold \end{cases} \quad (36)$$

## 2.8 Word2vec

The word vector has been proved to be more efficient, and serve as a better representation for the words compared to the simple one-hot vector. First, jieba is applied to separate the Chinese sentence into vocabularies, and the digits are separated one by one (for instance, 2010 will be expressed as 2, 0, 1, and 0). Afterward, word2vec(using genism) is implemented to generate the word vector. To avoid including too many vocabularies in the dictionary, the threshold is set to be 50 (the lower limit of the number of occurrences), below which the vocabularies are labeled as "unknown". The depth of the word vector is 300 units. A few hundred MB news data was provided to train the word vector, as it turns out that nearly 69k vocabularies are over the threshold and be included in the dictionary.

## 2.9 Model

The neural network architecture used in this paper is illustrated in Fig. 6. The hidden units in LSTM are set to be 800 (for biRNN it is 1600 in total). The attention depth is 300 (depth of  $k_i$  and  $d_i$ ). The input length (headline) and output length (content) are 24 and 90, respectively. The batch size is 32 during training. With only 1 CPU, it cost nearly a month to train the model. The results and discussions are detailed in the next section.

## 3 Results and Discussions

Unlike machine translation, there is no universal standard to judge the performance of news content-generation machine learning model. The evaluation of the performance of this model must be based on human judgment. In this section, two tests are carried out. The first test is to create a set of headlines, and then feed into the model to generate content, the evaluation is based on the performance of these generated results. The second test is to use recent news headlines as input (these data are not in training data) and check the consistency of our results and the real content of these news. Two criteria are considered to judge the performance of the model, the first is the grammar of the sentence, whether the sentence is grammatically correct or not. The second criterion is the consistency of the content and the headline, to see whether the model can generate relative content based on the input headline.

For the first test, the results are shown in Fig. 7. I create a few headlines(all of them are sarcastic) and feed them into the model. One can see that the result is quite reasonable. The model can understand the input and tries to generate

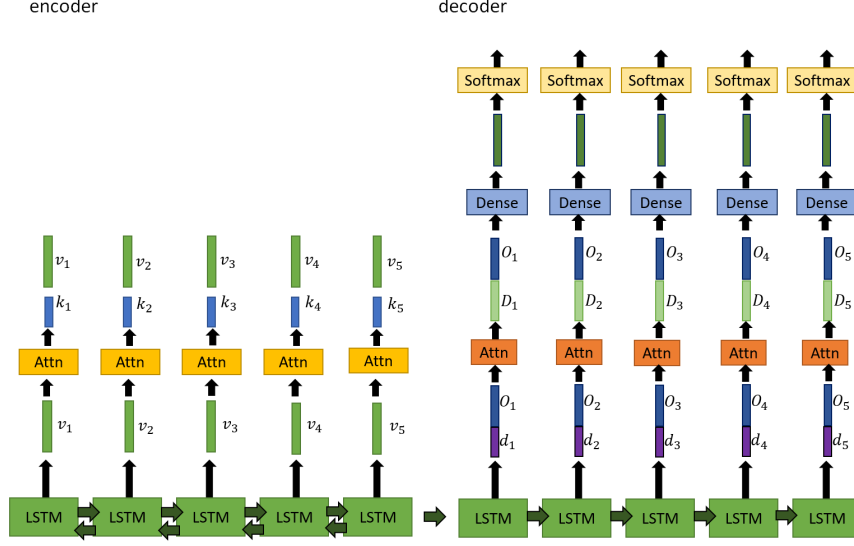


Figure 6: The neural network architecture used in this paper.

content based on it. Of course, the model absolutely needs to be further improved, too many unknown tokens appears, and the sentence is not so grammatically correct. Those problems will be discussed later.

For the second test, the results are shown in Fig. 8 and Fig. 9. Although the generated content is not consistent with real content (as expected), these generated content are still reasonable and interesting. If we modify these content, such as replace the unknown tokens with certain relative words and correct the grammar errors, this model is possible to reach the same level of performance as the journalists.

To check whether the attention mechanism works, the attention weight matrix  $\alpha$  is visualized in Fig. 10. Obviously, the relative words have larger weights. (桃園:桃園/機場, 機場:機場, 水上:水上, 樂園:樂園). However, since the amount of training data is inadequate, some words in the content might pay attention to a non-relative term.

To enhance the performance, we need to improve the model from two perspectives. From algorithm, and from training data. From the algorithm perspective, one important issue is the unknown vocabulary, of course, includes all the vocabularies in the training data is inefficient. Perhaps the best solution is to adopt a word and character-hybrid learning model [11]. From the training data perspective, there are some problems need to be solved to enhance the performance of the model. (1) date & location (2) different publishers (3) new information (4) the draft of the news.

- (1) date & location: The date when the news published is also crucial. For instance, the president will change over a certain period. The same goes to the location, if the mayor is mentioned in the headline, it would be ambiguous that who we are referred to. Since we don't provide the information such as in which city we live in, it is highly possible that the model will generate content that is less relevant to our expectation. To improve consistency, the date & location are needed to be embedded in the model.
- (2) different publishers: What is intriguing is that, even with the same headline, different news publishers can generate totally different content depending on their political tendency to a certain party. The persona embedding mechanism can be applied to solve this issue.
- (3) new information: To understand a new statement is also an important issue. For example, if the news used as training data are published before 2010, it is unlikely that the model can recognize the term "machine learning".
- (4) the draft of the news: Perhaps to include the draft of the news, such as interview and some other information, can drastically boost the performance. Since the draft is highly relative to the content, and it is the raw material used to construct the news.

## 4 Conclusion

In this paper, a news headline-content generation seq2seq model is introduced. Although the training data are relatively small (about tens of MB) compare to the state-of-the-art seq2seq model, the results indicate the output content is quite reasonable. It guarantees news headline-content to be a fruitful field for further investigation and research. Although some challenges still remain as mentioned in 3, if we adopted a better machine learning algorithm and much more data for training, the performance, and the effectiveness are expectable. The whole machine learning algorithm is written in Numpy, it is available on Github.

[https://github.com/hchungdelta/Simple\\_NN\\_API/tree/master/NN\\_v.2.1\\_news\\_generator](https://github.com/hchungdelta/Simple_NN_API/tree/master/NN_v.2.1_news_generator)

## References

- [1] Sepp Hochreiter and Jürgen Schmidhuber. Long-short term memory. pages 1735–1780. *Neural Computation* 9(8), 1997.
- [2] Kyunghyun Cho, Bart van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. On the properties of neural machine translation: Encoder–decoder approaches. In *eprint arXiv:1409.1259*, 2014.
- [3] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. In *eprint arXiv:1409.0473*, 2014.
- [4] Minh-Thang Luong, Hieu Pham, and Christopher D. Manning. Effective approaches to attention-based neural machine translation. In *eprint arXiv:1508.04025*, 2015.
- [5] Ashish Vaswani, Noam Shazeer, Niki Parmar, Llion Jones Jakob Uszkoreit, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *eprint arXiv:1706.03762v5*, 2017.
- [6] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization. In *eprint arXiv:1607.06450*, 2016.
- [7] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *eprint arXiv:1502.03167*, 2015.
- [8] Mike Schuster and Kuldip K. Paliwal. Bidirectional recurrent neural networks. volume 45. *IEEE TRANSACTION ON SIGNAL PROCESSING* 11, 1997.
- [9] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *eprint arXiv:1412.6980*, 2014.
- [10] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural networks. In *eprint arXiv:1211.5063*, 2012.
- [11] Minh-Thang Luong and Christopher D. Manning. Achieving open vocabulary neural machine translation with hybrid word-character models. In *eprint arXiv:1604.00788*, 2016.



```

headline (original): 韓國瑜:復興禿頭風潮推廣高雄特色禿頭產業
headline (input): 韓國瑜:復興禿頭風潮推廣高雄特色禿頭產業</S>
generated context (output):
    國民黨高雄市長參選人韓國瑜今天(20日)一同政見, 以前可以說是年輕人的年輕人,
    可以打李佳芬的就要找出他的選票, 是他推動的「特色產業」, 將以他的政見特色產業, 也創造出在高雄特色音樂實力。
    </S>
=====
headline (original): 桃園機場擬建立水上樂園
headline (input): 桃園機場擬建立水上樂園</S>
generated context (output):
    桃園機場<UNK>, 桃園市政府今天啟動「第二<UNK>水上樂園站」, 6月29日在高雄巨蛋體驗<UNK><UNK>主題,
    屆時也會建立公園站急<UNK>體驗水上樂園。</S>
=====
headline (original): 學校申請不上恐導致學生無法入學
headline (input): 學校申請不上恐導致學生無法入學</S>
generated context (output):
    為解決學生申請入學管道, <UNK><UNK>教育部長, 今天在<UNK>4<UNK>於學生宿舍申請不說,
    已違反規定, 且不只要<UNK>學生, 不能參加學校考試, 讓學生入學也能向教育部申請。</S>
=====
headline (original): 如何有效減肥又省錢?專家:別吃就好
headline (input): 如何有效減肥又省錢?專家:別吃就好</S>
generated context (output):
    近日, 不少人認為身材的<UNK>受<UNK>為影響孩子的疾病之一, 近日有1<UNK>人整理出不到1000<UNK>的機會,
    如今人又再一次吃就好, 專家表示, 只要治療的機會, 大部分的人都會<UNK>「<UNK>」,
    因此幾乎都有可能被人放的<UNK>, 因此<UNK>為「身材」的角度說明</S>
=====
headline (original): 高雄民眾:走在路上被石油噴到很不爽
headline (input): 高雄民眾:走在路上被石油噴到很不爽</S>
generated context (output):
    你知道高雄什麼呢?高雄有民眾不滿, 在網路上<UNK>跑去看影片, 赫然發現,
    民眾的民族<UNK>竟然是<UNK>的, 而且一條<UNK>, 被網友質疑是襲擊, <UNK>了<UNK>,
    不過被網友看了很不爽, 真是「<UNK>的很好嗎」。</S>
=====

```

Figure 7: Some samples to demonstrate the performance of the model.

<p>headline (input): 7年增14.5萬&lt;UNK&gt;桃園「人口淨遷入」&lt;UNK&gt;全台&lt;/S&gt; generated context (output): 桃園市人口去年疊整&lt;UNK&gt;&lt;UNK&gt;高達14.5萬元, 桃園市29%、桃園市、新莊、新莊、新莊與桃園, 桃園市也是許多人口的共同績效。 桃園市也是新莊地區, 但因為桃園市去年多發生7年起, &lt;UNK&gt;成長1倍, &lt;UNK&gt;及時需1年的長人口增加, &lt;UNK&gt;成長率達15.5萬元, &lt;UNK&gt;&lt;UNK&gt;兆7.</p> <p>headline (input): 羅志祥無預警「全面退出」&lt;UNK&gt;粉絲驚呆爆內幕&lt;/S&gt; generated context (output): 香港本土&lt;UNK&gt;「&lt;UNK&gt;」羅志祥今年2月突然宣布來台發展自家地退出, 雖然沒有退出&lt;UNK&gt;, 不過昨(3日)晚赴&lt;UNK&gt;突然退出, 羅志祥此舉也突然退出演藝圈。&lt;/S&gt;</p> <p>headline (input): 疑為趕搭公車&lt;UNK&gt;7旬兄妹闖紅燈過馬路遭&lt;UNK&gt;&lt;/S&gt; generated context (output): 苗栗一名70歲的黃姓男子, 23日早上跟妻子&lt;UNK&gt;&lt;UNK&gt;, 行經彰化&lt;UNK&gt;路口停等紅燈時, 發現路邊沒有停下來, 要他加速逃逸, &lt;UNK&gt;&lt;UNK&gt;發現, 立刻通報警方, 對他攔下, 發現他酒測值0.41、10、15歲, &lt;UNK&gt;坦承是因為自己闖紅燈, 10月與一名&lt;UNK&gt;, 結果在公車撞上&lt;/S&gt;</p> <p>headline (input): 「自殺大樓」一躍而下&lt;UNK&gt;女大生沾滿死者血... 淒厲慘叫昏倒&lt;/S&gt; generated context (output): 彰化一名23歲的女學生, 晚間8點多在中山高&lt;UNK&gt;大樓內發現大樓墜樓, 1樓後2樓房間&lt;UNK&gt;, 就在大樓內, 臉、頸部、四肢甚至毛頭部, 一度緊張急送急, 由於2人都是在&lt;UNK&gt;, 紅&lt;UNK&gt;女子難以行走, 警方研判, 嫌犯疑似在大樓頂樓工作, &lt;UNK&gt;&lt;UNK&gt;, 進而在樓梯間徘徊, 不知不覺時就發住院, 大樓</p> <p>headline (input): 燈會竟出現「18禁」畫面? 童好奇問爸:這在幹嘛&lt;/S&gt; generated context (output): 台灣&lt;UNK&gt;&lt;UNK&gt;縣的台灣18日前發展出「18&lt;UNK&gt;&lt;UNK&gt;」, 近日在台灣有不少外國觀光客來台, 前往&lt;UNK&gt;玩「&lt;UNK&gt;」畫面, 期間更拍下《&lt;UNK&gt;的&lt;UNK&gt;》, 讓不少網友驚呼「&lt;UNK&gt;&lt;UNK&gt;」, 但台灣人也說, 台灣燈會是大&lt;UNK&gt;18禁&lt;UNK&gt;的觀光, 景點若是「&lt;UNK&gt;」, 如此有趣的成果也有, 讓人不禁</p>	<p>headline (original): 7年增14.5萬 桃園「人口淨遷入」冠居全台 real context: 人口遷徙, 攸關各地區域發展情形, 內政部最近公布一項統計, 顯示7年來桃園增加近14萬5千人口, 高居全國第一, 對比距離不遠的首都台北, 人口出現負成長, 有從台北搬到桃園的人直說, 當地吸引他們的關鍵就在房價。</p> <p>headline (original): 羅志祥無預警「全面退出」 粉絲驚呆爆內幕 real context: 藝人羅志祥昨(18日)突然在Instagram發限時動態, 無預警宣布「全面退出」經營12年的潮牌「STAGE」, 震驚粉絲們。</p> <p>headline (original): 疑為趕搭公車7旬兄妹闖紅燈過馬路遭撞亡 real context: 基隆(18日)發生死亡車禍, 一對70多歲兄妹, 疑似為了趕搭公車橫越馬路, 92歲老翁來不及反應撞上, 兄妹兩人陸續宣告不治, 警方依「過失傷害致死」罪嫌將老翁移送偵辦。</p> <p>headline (original): 「自殺大樓」一躍而下 女大生沾滿死者血... 淒厲慘叫昏倒 real context: 一名資深房仲日前分享親身經歷的靈異事件, 高雄三民區有一棟業界著名的「自殺大樓」, 10多年前有一名男子跳樓, 當時他目睹有一名女騎士遭波及, 全身沾滿跳樓男的血, 當場崩潰慘叫、隨即昏倒; 房東則癱坐在管理室旁的地板, 喃喃念著「我的房子2年死了3個人」。</p> <p>headline (original): 燈會竟出現「18禁」畫面? 童好奇問爸:這在幹嘛 real context: 元宵節腳步近, 北中南燈會都已開跑, 但今年同樣也是褒貶不一! 像是今年2019台灣燈會在屏東大鵬灣登場, 試營運以來就大受好評, 其中無人機空中秀及結合水域的主燈秀彩排, 更是讓人驚豔。不過對比高雄燈會日前才因蓮花燈挨酸, 現在南投竟遭爆出有「18禁」的畫面, 讓網友直呼感覺「怪怪的」。</p>
--	---

Figure 8: (part one) Use recent news headlines (not in training data) as input, and compare to the real content. Right: results, left: recent news data.

<p>headline (input): 洗手台就是&lt;UNK&gt;！清潔工躲廁所淒涼用餐&lt;UNK&gt;惹網友不捨&lt;/S&gt; generated context (output): 不少&lt;UNK&gt;在浴室廁所，蹲在地板上貼著「&lt;UNK&gt;」，但其實不少人都會擔心，假期&lt;UNK&gt;小孩不想廁所，廁所也不多；但其實不少人都會貼出，有人在廁所貼上「&lt;UNK&gt;&lt;UNK&gt;手套」，但近日卻有網友發現，一個「&lt;UNK&gt;」的貼&lt;UNK&gt;，竟是清潔公司的「&lt;UNK&gt;」，讓不少網友看了直呼心疼，紛紛跪求</p> <p>headline (input): 女大生頭痛昏迷4天&lt;UNK&gt;醒來已當媽驚喊：孩子不是我的！&lt;/S&gt; generated context (output): 遇到這種誇張的事情，讓自己相當無奈！一名女大生發文，日前在&lt;UNK&gt;路上，遇到一名女大生，對方詢問「要帶小孩不是我的孩子不是妳嗎？」讓她無奈表示「孩子不是我的！」，讓她覺得相當尷尬，忍不住在臉書發文「小孩不是我的寶貝不是我的！」&lt;/S&gt;</p> <p>headline (input): 韓國瑜猛轟中央&lt;UNK&gt;呂秀蓮：半年後就原形畢露&lt;/S&gt; generated context (output): 年底大選將近，前副總統呂秀蓮日前在臉書宣布，「縣市長選舉的10個民進黨國家是徵收&lt;UNK&gt;」，但其中&lt;UNK&gt;的前副秘書長&lt;UNK&gt;，在&lt;UNK&gt;綠的&lt;UNK&gt;，於昨(19日)召開第2000名「&lt;UNK&gt;&lt;UNK&gt;」時，蔡英文表示「他們都很難執政的是很好的」。&lt;/S&gt;</p> <p>headline (input): 天氣一週「3變」！明稍&lt;UNK&gt;四又回暖&lt;UNK&gt;週末再轉濕涼&lt;/S&gt; generated context (output): 持續在太平洋高壓勢力的台灣北部近海，加上&lt;UNK&gt;&lt;UNK&gt;逐漸影響，天氣一週後，各地高溫再上升3度以上，因此現在到週四(3日、3日)都有3個高濕站，不過週末假期又是多雲到晴的天氣，不過在下午2&lt;UNK&gt;後雷雨，局部短暫陣雨後，外出記得攜帶雨具備用，&lt;UNK&gt;&lt;UNK&gt;2到3、28度；由於20日、</p>	<p>headline (original):洗手台就是飯桌 ！ 清潔工躲廁所淒涼用餐惹網友不捨 real context: 繁榮社會背後的悲哀？近日有網友在香港拍下了這張照片，只見公廁內有一名清潔工，放了張塑膠椅就坐在洗手台前吃著飯菜，淒涼的畫面讓人不捨！也讓這名網友忍不住感嘆「香港怎麼會變成這樣」。</p> <p>headline (original): 女大生頭痛昏迷4天 醒來已當媽驚喊：孩子不是我的！ real context: 英國曼徹斯特發生一宗離奇事件，18歲女大生艾波妮(Ebony Stevenson)日前因為頭痛和多次癲癇發作，送醫後整整昏迷4天，豈料當她甦醒後，竟已成為人母！當護理師把女寶寶送到她懷中時，艾波妮驚恐拒絕，傻眼表示「她不可能是我的。」</p> <p>headline (original): 韓國瑜猛轟中央 呂秀蓮：半年後就原形畢露 real context: 高雄市長韓國瑜日前表示，「台灣鬼混20多年，完全沒有在進步」，此話一出又再度引發爭議，然而前副總統呂秀蓮今(18日)在接受廣播節目專訪時批「什麼叫鬼混？」她質問，為什麼以前大家不認識韓國瑜，現在才爆紅？韓過去也當過立委好幾年，是不是都在鬼混呢？</p> <p>headline (original): 天氣一週「3變」！明稍涼週四又回暖週末再轉濕涼 real context: 天氣一週多變！今(12日)起東北季風減弱，天氣明顯回暖，各地回到穩定的天氣型態，北部白天回溫幅度大，今高溫來到26度，中南部高溫30度，但要留意早晚溫差大。而據中央氣象局預報，明(13日)東北季風稍增強，北部及東北部稍轉涼，其他地區仍是多雲到晴，北部高溫約20至22度，下降3到6度；而週五晚間又會有另一波東北風南下，溫度將略降。</p>
<p>headline (input): 車前突現人影！&lt;UNK&gt;女子暗夜過馬路嚇壞駕駛&lt;/S&gt; generated context (output): &lt;UNK&gt;過馬路，被人發現不只見他，竟然直接在大馬路上，一名就看到人騎單車的男子，突然被對&lt;UNK&gt;，說他車子突然&lt;UNK&gt;，突然來到路邊，差點被撞上，趕緊報警，警方到場，男子終於被救出，而他卻阿嬤也傳，警方調查，目前&lt;/S&gt;</p> <p>headline (input): 快訊／台中女老師失蹤三日&lt;UNK&gt;警公布&lt;UNK&gt;協尋&lt;/S&gt; generated context (output): 最新消息！國立&lt;UNK&gt;一名李姓老師，2日在&lt;UNK&gt;海域發現31歲的印尼&lt;UNK&gt;，失蹤後，家屬急忙出面，警方公布監視器畫面發現，有2&lt;UNK&gt;少年在巷口，失蹤後，警方根據報案&lt;UNK&gt;協尋，同時也找到死者，警方也聯繫中關人，當天也聯繫不上，也感謝台中市警察局，也提醒女幼兒園沒有&lt;UNK&gt;</p>	<p>headline (original): 車前突現人影！女子暗夜過馬路嚇壞駕駛 real context: 行人夜間違規穿越馬路真的很危險，桃園市有一名駕駛凌晨開車在路上，突然一名女子就這麼從道路的中線冒了出來，還好駕駛閃得快，才沒有撞到女子，警方表示這樣的行為，已經違反交通管理處罰條例，將要開罰。</p> <p>headline (original): 快訊／台中女老師失蹤三日 警公布特徵齊協尋 real context: 台中沙鹿北勢國中一名女老師，15日清晨騎機車出門，疑似要到學校上課途中失聯，親友苦尋不到向警方報案，警方調閱監視器，今天(18日)派出空拍機找尋沒進展，晚間公布老師穿著特徵，希望民眾協助提供線索給清水分局。</p>

Figure 9: (part two) Use recent news headlines (not in training data) as input, and compare to the real content. Right: results, left: recent news data.

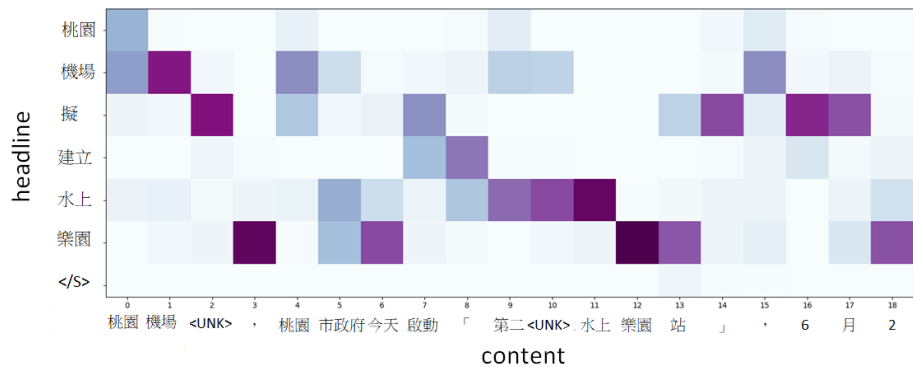


Figure 10: Visualizing the attention weight matrix alpha. The deeper the color is, the heavier the weight is put on that term.