

# Generating Classical Chinese Poems from Vernacular Chinese

Zhichao Yang<sup>1\*</sup>, Pengshan Cai<sup>1\*</sup>, Yansong Feng<sup>2</sup>, Fei Li<sup>1</sup>,  
Weijiang Feng<sup>3</sup>, Elena Suet-Ying Chiu<sup>1</sup>, Hong Yu<sup>1</sup>

<sup>1</sup> University of Massachusetts, MA, USA

{zhichaoyang, pengshancai}@umass.edu foxlf823@gmail.com  
chiu@llc.umass.edu hong\_yu@uml.edu

<sup>2</sup> Institute of Computer Science and Technology, Peking University, China  
fengyansong@pku.edu.cn

<sup>3</sup> College of Computer, National University of Defense Technology, China  
fengweijiang14@nudt.edu.cn

## Abstract

Classical Chinese poetry is a jewel in the treasure house of Chinese culture. Previous poem generation models only allow users to employ keywords to interfere the meaning of generated poems, leaving the dominion of generation to the model. In this paper, we propose a novel task of generating classical Chinese poems from vernacular, which allows users to have more control over the semantic of generated poems. We adapt the approach of unsupervised machine translation (UMT) to our task. We use segmentation-based padding and reinforcement learning to address under-translation and over-translation respectively. According to experiments, our approach significantly improve the perplexity and BLEU compared with typical UMT models. Furthermore, we explored guidelines on how to write the input vernacular to generate better poems. Human evaluation showed our approach can generate high-quality poems which are comparable to amateur poems.

## 1 Introduction

During thousands of years, millions of classical Chinese poems have been written. They contain ancient poets' emotions such as their appreciation for nature, desiring for freedom and concerns for their countries. Among various types of classical poetry, *quatrain poems* stand out. On the one hand, their aestheticism and terseness exhibit unique elegance. On the other hand, composing such poems is extremely challenging due to their phonological, tonal and structural restrictions.

Most previous models for generating classical Chinese poems (He et al., 2012; Zhang and Lapata, 2014) are based on limited keywords or characters at fixed positions (e.g., acrostic poems).

Since users could only interfere with the semantic of generated poems using a few input words, models control the procedure of poem generation. In this paper, we proposed a novel model for classical Chinese poem generation. As illustrated in Figure 1, our model generates a classical Chinese poem based on a vernacular Chinese paragraph. Our objective is not only to make the model generate aesthetic and terse poems, but also keep rich semantic of the original vernacular paragraph. Therefore, our model gives users more control power over the semantic of generated poems by carefully writing the vernacular paragraph.

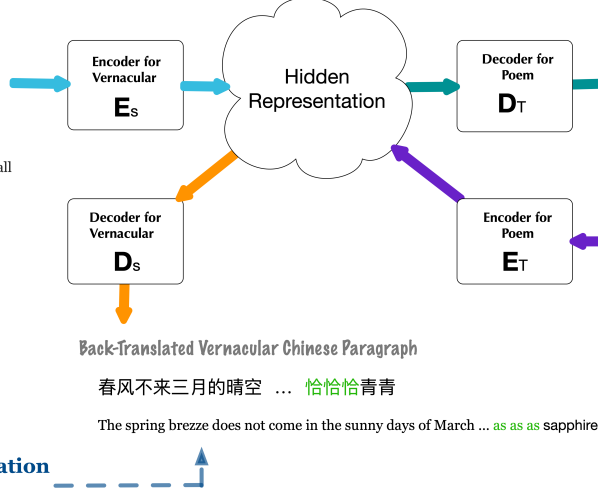
Although a great number of classical poems and vernacular paragraphs are easily available, there exist only limited human-annotated pairs of poems and their corresponding vernacular translations. Thus, it is unlikely to train such poem generation model using supervised approaches. Inspired by unsupervised machine translation (UMT) (Lample et al., 2018b), we treated our task as a translation problem, namely translating vernacular paragraphs to classical poems.

However, our work is not just a straight-forward application of UMT. In a training example for UMT, the length difference of source and target languages are usually not large, but this is not true in our task. Classical poems tend to be more concise and abstract, while vernacular text tends to be detailed and lengthy. Based on our observation on gold-standard annotations, vernacular paragraphs usually contain more than twice as many Chinese characters as their corresponding classical poems. Therefore, such discrepancy leads to two main problems during our preliminary experiments: (1) **Under-translation**: when summarizing vernacular paragraphs to poems, some vernacular sentences are not translated and ignored by our model. Take the last two vernacular sentences in Figure

\*Equal contribution

#### Vernacular Chinese Paragraph

东风不来，三月的柳絮不飞  
Without the east wind, the willow catkins  
in March do not flutter  
你的心如小小的寂寞的城  
Your heart is like the lonesome little town  
恰若青石的街道向晚  
Like its streets of cobblestones near nightfall  
跫音不响，三月的春帷不揭  
When footfalls are silent and the bed  
curtains of March not unveiled  
你的心是小小的窗扉紧掩  
Your heart is a little window tightly shut  
我达达的马蹄是美丽的错误  
My clattering hooves is a beautiful mistake  
我不是归人，是个过客  
I am not a homecoming man  
but a passing traveler



#### Generated Poem

春风不来三月晴，  
Shine in March but  
no spring breeze blows  
柳絮不飞霜叶轻。  
Light and small maple leaf  
with no catkin flows  
心如小窗寂无梦，  
My heart is like a small window,  
dreamless and quiet  
城如三月恰青青。  
As the city in March,  
as green as sapphire

#### Back-Translated Vernacular Chinese Paragraph

春风不来三月的晴空 ... 恰恰青青  
The spring breeze does not come in the sunny days of March ... as as as sapphire

Figure 1: An example of the training procedures of our model. Here we depict two procedures, namely back translation and language modeling. Back translation has two paths, namely  $E_S \rightarrow D_T \rightarrow E_T \rightarrow D_S$  and  $D_T \rightarrow E_S \rightarrow D_S \rightarrow E_T$ . Language modeling also has two paths, namely  $E_T \rightarrow D_T$  and  $E_S \rightarrow D_S$ . Figure 1 shows only the former one for each training procedure.

1 as examples, they are not covered in the generated poem. (2) **Over-translation**: when expanding poems to vernacular paragraphs, certain words are unnecessarily translated for multiple times. For example, the last sentence in the generated poem of Figure 1, *as green as sapphire*, is back-translated as *as green as as as sapphire*.

Inspired by the phrase segmentation schema in classical poems (Ye, 1984), we proposed the method of phrase-segmentation-based padding to handle with under-translation. By padding poems based on the phrase segmentation custom of classical poems, our model better aligns poems with their corresponding vernacular paragraphs and meanwhile lowers the risk of under-translation. Inspired by Paulus et al. (2018), we designed a reinforcement learning policy to penalize the model if it generates vernacular paragraphs with too many repeated words. Experiments show our method can effectively decrease the possibility of over-translation.

The contributions of our work are threefold:

- (1) We proposed a novel task for unsupervised Chinese poem generation from vernacular text.
- (2) We proposed using phrase-segmentation-based padding and reinforcement learning to address two important problems in this task, namely under-translation and over-translation.
- (3) Through extensive experiments, we proved

the effectiveness of our models and explored how to write the input vernacular to inspire better poems. Human evaluation shows our models are able to generate high quality poems, which are comparable to amateur poems.

## 2 Related Works

**Classical Chinese Poem Generation** Most previous works in classical Chinese poem generation focus on improving the semantic coherence of generated poems. Based on LSTM, Zhang and Lapata (2014) purposed generating poem lines incrementally by taking into account the history of what has been generated so far. Yan (2016) proposed a polishing generation schema, each poem line is generated incrementally and iteratively by refining each line one-by-one. Wang et al. (2016) and Yi et al. (2018) proposed models to keep the generated poems coherent and semantically consistent with the user’s intent. There are also researches that focus on other aspects of poem generation. (Yang et al. (2018) explored increasing the diversity of generated poems using an unsupervised approach. Xu et al. (2018) explored generating Chinese poems from images. While most previous works generate poems based on topic words, our work targets at a novel task: generating poems from vernacular Chinese paragraphs.

**Unsupervised Machine Translation** Compared

with supervised machine translation approaches (Cho et al., 2014; Bahdanau et al., 2015), unsupervised machine translation (Lample et al., 2018a,b) does not rely on human-labeled parallel corpora for training. This technique is proved to greatly improve the performance of low-resource languages translation systems. (e.g. English-Urdu translation). The unsupervised machine translation framework is also applied to various other tasks, e.g. image captioning (Feng et al., 2019), text style transfer (Zhang et al., 2018), speech to text translation (Bansal et al., 2017) and clinical text simplification (Weng et al., 2019). The UMT framework makes it possible to apply neural models to tasks where limited human labeled data is available. However, in previous tasks that adopt the UMT framework, the abstraction levels of source and target language are the same. This is not the case for our task.

**Under-Translation & Over-Translation** Both are troublesome problems for neural sequence-to-sequence models. Most previous related researches adopt the coverage mechanism (Tu et al., 2016; Mi et al., 2016; Sankaran et al., 2016). However, as far as we know, there were no successful attempt applying coverage mechanism to transformer-based models (Vaswani et al., 2017).

### 3 Model

#### 3.1 Main Architecture

We transform our poem generation task as an unsupervised machine translation problem. As illustrated in Figure 1, based on the recently proposed UMT framework (Lample et al., 2018b), our model is composed of the following components:

- Encoder  $\mathbf{E}_s$  and decoder  $\mathbf{D}_s$  for vernacular paragraph processing
- Encoder  $\mathbf{E}_t$  and decoder  $\mathbf{D}_t$  for classical poem processing

where  $\mathbf{E}_s$  (or  $\mathbf{E}_t$ ) takes in a vernacular paragraph (or a classical poem) and converts it into a hidden representation, and  $\mathbf{D}_s$  (or  $\mathbf{D}_t$ ) takes in the hidden representation and converts it into a vernacular paragraph (or a poem). Our model relies on a vernacular texts corpus  $S$  and a poem corpus  $T$ . We denote  $S$  and  $T$  as instances in  $S$  and  $T$  respectively.

The training of our model relies on three procedures, namely *parameter initialization*, *language*

*modeling* and *back-translation*. We will give detailed introduction to each procedure.

**Parameter initialization** As both vernacular and classical poem use Chinese characters, we initialize the character embedding of both languages in one common space, the same character in two languages shares the same embedding. This initialization helps associate characters with their plausible translations in the other language.

**Language modeling** It helps the model generate texts that conform to a certain language. A well-trained language model is able to detect and correct minor lexical and syntactic errors. We train the language models for both vernacular and classical poem by minimizing the following loss:

$$\mathcal{L}^{lm} = \mathbb{E}_{S \in S} [-\log P(S|\mathbf{D}_s(\mathbf{E}_s(S_N)))] + \mathbb{E}_{T \in T} [-\log P(T|\mathbf{D}_t(\mathbf{E}_t(T_N)))] \quad (1)$$

where  $S_N$  (or  $T_N$ ) is generated by adding noise (drop, swap or blank a few words) in  $S$  (or  $T$ ).

**Back-translation** Based on a vernacular paragraph  $S$ , we generate a poem  $T_S$  using  $\mathbf{E}_s$  and  $\mathbf{D}_t$ , we then translate  $T_S$  back into a vernacular paragraph  $S_{T_S} = \mathbf{D}_s(\mathbf{E}_t(T_S))$ . Here,  $S$  could be used as gold standard for the back-translated paragraph  $S_{T_S}$ . In this way, we could turn the unsupervised translation into a supervised task by maximizing the similarity between  $S$  and  $S_{T_S}$ . The same also applies to using poem  $T$  as gold standard for its corresponding back-translation  $T_{S_T}$ . We define the following loss:

$$\mathcal{L}^{bt} = \mathbb{E}_{S \in S} [-\log P(S|\mathbf{D}_s(\mathbf{E}_t(T_S)))] + \mathbb{E}_{T \in T} [-\log P(T|\mathbf{D}_t(\mathbf{E}_s(S_T)))] \quad (2)$$

Note that  $\mathcal{L}^{bt}$  does not back propagate through the generation of  $T_S$  and  $S_T$  as we observe no improvement in doing so. When training the model, we minimize the composite loss:

$$\mathcal{L} = \alpha_1 \mathcal{L}^{lm} + \alpha_2 \mathcal{L}^{bt}, \quad (3)$$

where  $\alpha_1$  and  $\alpha_2$  are scaling factors.

#### 3.2 Addressing Under-Translation and Over-Translation

During our early experiments, we realize that the naive UMT framework is not readily applied to our task. Classical Chinese poems are featured for





|   | Training set | Validation set | Test set |
|---|--------------|----------------|----------|
| # Poems                                 | 163K         | 19K            | 487      |
| Average length of poems                 | 32.0         | 32.0           | 32.0     |
| # vernacular paragraphs                 | 337K         | 19K            | 487      |
| Average length of vernacular paragraphs | 71.8         | 76.8           | 73.3     |

Table 1: Statistics of our dataset

2017), we define the following loss function:

$$\mathcal{L}^{rl} = \mathbb{E}_{S \in \mathcal{S}} [(RR(S_{T_S}) - \tau) \log P(S | \mathbf{D}_s(\mathbf{E}_t(T_S)))], \quad (5)$$

where  $\tau$  is a manually set threshold. Intuitively, minimizing  $\mathcal{L}^{rl}$  is equivalent to maximizing the conditional likelihood of the sequence  $S$  given  $S_{T_S}$  if its *repetition ratio* is lower than the threshold  $\tau$ . Following (Wu et al., 2016), we revise the composite loss as:

$$\mathcal{L}' = \alpha_1 \mathcal{L}^{lm} + \alpha_2 \mathcal{L}^{bt} + \alpha_3 \mathcal{L}^{rl}, \quad (6)$$

where  $\alpha_1, \alpha_2, \alpha_3$  are scaling factors.

## 4 Experiment

The objectives of our experiment are to explore the following questions: (1) How much do our models improve the generated poems? (Section 4.4) (2) What are characteristics of the input vernacular paragraph that lead to a good generated poem? (Section 4.5) (3) What are weaknesses of generated poems compared to human poems? (Section 4.6) To this end, we built a dataset as described in Section 4.1. Evaluation metrics and baselines are described in Section 4.2 and 4.3. For the implementation details of building the dataset and models, please refer to supplementary materials.<sup>1</sup>

### 4.1 Datasets

**Training and Validation Sets** We collected a corpus of poems and a corpus of vernacular literature from online resources. The poem corpus contains 163K quatrain poems from *Tang Poems* and *Song Poems*, the vernacular literature corpus contains 337K short paragraphs from 281 famous books, the corpus covers various literary forms including prose, fiction and essay. Note that our poem corpus and a vernacular corpus are not aligned. We further split the two corpora into a training set and a validation set.

<sup>1</sup>Our data and code is publicly available at <https://github.com/whaleloops/interpoetry>

**Test Set** From online resources, we collected 487 seven-character quatrain poems from *Tang Poems* and *Song Poems*, as well as their corresponding high quality vernacular translations. These poems could be used as gold standards for poems generated from their corresponding vernacular translations. Table 1 shows the statistics of our training, validation and test set.

### 4.2 Evaluation Metrics

**Perplexity** Perplexity reflects the probability a model generates a certain poem. Intuitively, a better model would yield higher probability (lower perplexity) on the gold poem.

**BLEU** As a standard evaluation metric for machine translation, BLEU (Papineni et al., 2001) measures the intersection of n-grams between the generated poem and the gold poem. A better generated poem usually achieves higher BLEU score, as it shares more n-gram with the gold poem.

**Human evaluation** While perplexity and BLEU are objective metrics that could be applied to large-volume test set, evaluating Chinese poems is after all a subjective task. We invited 30 human evaluators to join our human evaluation. The human evaluators were divided into two groups. The expert group contains 15 people who hold a bachelor degree in Chinese literature, and the amateur group contains 15 people who holds a bachelor degree in other fields. All 30 human evaluators are native Chinese speakers.

We ask evaluators to grade each generated poem from four perspectives: 1) *Fluency*: Is the generated poem grammatically and rhythmically well formed, 2) *Semantic coherence*: Is the generated poem itself semantic coherent and meaningful, 3) *Semantic preservability*: Does the generated poem preserve the semantic of the modern Chinese translation, 4) *Poeticness*: Does the generated poem display the characteristic of a poem and does the poem build good poetic image. The grading scale for each perspective is from 1 to 5.

| ID | Vernacular paragraph  | Generated poem  |
|----|---|---|
| 1  | 青山隐隐约约绿水千里迢迢，秋时已尽江南草木还未枯凋。二十四桥明月映照幽幽清夜，你这美人现在何处教人吹箫？  | 青山隐隐绿水光，千里秋时已尽藏。<br>江南草木还未枯，二十四桥幽夜香。  |
| 1  | Blue peak is blur and faint, green river meanders thousands miles, the southern grass has not dried up in the autumn. Twenty-four hole bridges and bright moons shine in the clear night. where do you beautifully teach people to flute? | Blue peak blurs and green river shines, thousands miles away, autumn has been hidden. the southern grass has not dried up in the autumn. Twenty-four hole bridges smells beautiful and quiet tonight. |
| 2  | 拂袖起舞于梦中徘徊，相思蔓上心扉。她眷恋梨花泪，静画红妆等谁归，空留伊人徐徐憔悴。   | 拂袖起舞梦徘徊，相思蔓上心扉开。<br>玉眷梨花泪痕静，画红等谁归去来。  |
| 2  | The sleeves danced in the dream, and the lovesickness was on the heart. She is in love with the tears of pears, and who is quietly wearing red makeup, only left alone to be languished slowly.   | The sleeves danced in the dream, the lovesickness appeared in the heart. Jade concerns tears of pears but the mark is still, wearing red makeup waiting for the one to come and go.                   |
| 3  | 窗外的麻雀在电线杆上多嘴，你说这一句很有夏天的感觉。手中的铅笔在纸上来来回回，我用几行字形容你是我的谁。  | 窗下麻姑灯火多，闲中说与万缘何。<br>夏频手把铅华纸，来往回头用几多。  |
| 3  | The sparrow outside the window is talking on the pole. You say this sentence makes you feel very summer. The pencil in my hand is writing back and forth on the paper. I only use a few lines to describe who you are to me.              | Under the window lie sparrow girls in this prosperous city, chit chatting about the destiny of the world. Summer hands over many drawing canvas, Looking back and forth, how many do you need?        |
| 4  | 雨天的屋瓦，浮漾湿湿的流光，灰而温柔，迎光则微明，背光则幽黯，对于视觉，是一种低沉的安慰。   | 雨余屋瓦浮漾湿，流光灰色暖相迎。<br>光则微明背则色，幽人黯黯对风清。  |
| 4  | The rainy days of the roof tiles are soaking wet and wet, gray and gentle, Facing the light, it is slightly bright, Against the light, it is pitch dark, For the concept of vision, it is a deep comfort.                                 | The excess rain makes roof tiles rippling, ambilight gray is warm and welcoming. Light is slightly bright, against is pure color. The person hides in dark but faces wind breeze.                     |
| 5  | 只要历史不阻断，时间不倒退，一切都会衰老。老就老了吧，安详地交给世界一副慈祥美。假饰天真是最残酷的自我糟践。  | 只要诸公不阻时，不倒退食一尘埃。<br>会衰老矣安分世，一副慈祥假此来。  |
| 5  | As long as history does not block, time does not go backwards, everything will age. It is fine to get old, and handing it to the world with kindness. Faking innocence is the cruelest self-destruction.                                  | As long as people do not block time, it will not go backwards and absorbs into a dust. People should stay chill and get old. Faking innocence is not the way to go.                                   |

Table 2: A few poems generated by our model from their corresponding vernacular paragraphs.

### 4.3 Baselines

We compare the performance of the following models: (1) *LSTM* (Hochreiter and Schmidhuber, 1997); (2) *Naive transformer* (Vaswani et al., 2017); (3) *Transformer + Anti OT* (RL loss); (4) *Transformer + Anti UT* (phrase segmentation-based padding); (5) *Transformer + Anti OT&UT*.

### 4.4 Reborn Poems: Generating Poems from Vernacular Translations

As illustrated in Table 2 (ID 1). Given the vernacular translation of each gold poem in test set, we generate five poems using our models. Intuitively, the more the generated poem resembles the gold poem, the better the model is. We report mean perplexity and BLEU scores in Table 3 (Where +Anti OT refers to adding the reinforcement loss to mitigate over-fitting and +Anti UT refers to adding phrase segmentation-based padding to mitigate under-translation), human evaluation results

in Table 4.<sup>2</sup>

According to experiment results, perplexity, BLEU scores and total scores in human evaluation are consistent with each other. We observe all BLEU scores are fairly low, we believe it is reasonable as there could be multiple ways to compose a poem given a vernacular paragraph. Among transformer-based models, both +Anti OT and +Anti UT outperforms the naive transformer, while Anti OT&UT shows the best performance, this demonstrates alleviating under-translation and over-translation both helps generate better poems. Specifically, +Anti UT shows bigger improvement than +Anti OT. According to human evaluation, among the four perspectives, our Anti OT&UT brought most score improvement in *Semantic preservability*, this proves our improvement on semantic preservability was most obvious to human evaluators. All transformer-

<sup>2</sup>We did not use LSTM in human evaluation since its performance is worse as shown in Table 3.

| Model       | Perplexity   | BLEU        | BLEU-1       | BLEU-2      | BLEU-3      | BLEU-4      |
|-------------|--------------|-------------|--------------|-------------|-------------|-------------|
| LSTM        | 118.27       | 3.81        | 39.16        | 6.93        | 1.58        | 0.49        |
| Transformer | 105.79       | 5.50        | 40.92        | 8.02        | 2.46        | 1.11        |
| +Anti OT    | 77.33        | 6.08        | 41.22        | 8.72        | 2.82        | 1.36        |
| +Anti UT    | 74.21        | 6.34        | 42.20        | <b>9.04</b> | <b>2.96</b> | 1.44        |
| +Anti OT&UT | <b>65.58</b> | <b>6.57</b> | <b>42.53</b> | 8.98        | <b>2.96</b> | <b>1.46</b> |

Table 3: Perplexity and BLEU scores of generating poems from vernacular translations. Since perplexity and BLEU scores on the test set fluctuates from epoch to epoch, we report the mean perplexity and BLEU scores over 5 consecutive epochs after convergence.

| Model       | Fluency     | Semantic coherence | Semantic preservability | Poeticness  | Total        |
|-------------|-------------|--------------------|-------------------------|-------------|--------------|
| Transformer | 2.63        | 2.54               | 2.12                    | 2.46        | 9.75         |
| +Anti OT    | 2.80        | 2.75               | 2.44                    | 2.71        | 10.70        |
| +Anti UT    | 2.82        | 2.82               | 2.86                    | 2.85        | 11.35        |
| +Anti OT&UT | <b>3.21</b> | <b>3.27</b>        | <b>3.27</b>             | <b>3.28</b> | <b>13.13</b> |

Table 4: Human evaluation results of generating poems from vernacular translations. We report the mean scores for each evaluation metric and total scores of four metrics.

based models outperform LSTM. Note that the average length of the vernacular translation is over 70 characters, comparing with transformer-based models, LSTM may only keep the information in the beginning and end of the vernacular. We anticipated some score inconsistency between expert group and amateur group. However, after analyzing human evaluation results, we did not observed big divergence between two groups.

#### 4.5 Interpoetry: Generating Poems from Various Literature Forms

Chinese literature is not only featured for classical poems, but also various other literature forms. *Song lyric*(宋词), or *ci* also gained tremendous popularity in its palmy days, standing out in classical Chinese literature. *Modern prose*, *modern poems* and *pop song lyrics* have won extensive praise among Chinese people in modern days. The goal of this experiment is to transfer texts of other literature forms into quatrain poems. We expect the generated poems to not only keep the semantic of the original text, but also demonstrate terseness, rhythm and other characteristics of ancient poems. Specifically, we chose 20 famous fragments from four types of Chinese literature (5 fragments for each of modern prose, modern poems, pop song lyrics and Song lyrics). We try to As no ground truth is available, we resorted to human evaluation with the same grading standard in Section 4.4.

Comparing the scores of different literature forms, we observe Song lyric achieves higher scores than the other three forms of modern literature. It is not surprising as both Song lyric and

quatrain poems are written in classical Chinese, while the other three literature forms are all in vernacular.

Comparing the scores within the same literature form, we observe the scores of poems generated from different paragraphs tends to vary. After carefully studying the generated poems as well as their scores, we have the following observation:

1) In classical Chinese poems, poetic images (意象) were widely used to express emotions and to build artistic conception. A certain poetic image usually has some fixed implications. For example, *autumn* is usually used to imply sadness and loneliness. However, with the change of time, poetic images and their implications have also changed. According to our observation, if a vernacular paragraph contains more poetic images used in classical literature, its generated poem usually achieves higher score. As illustrated in Table 2, both paragraph 2 and 3 are generated from pop song lyrics, paragraph 2 uses many poetic images from classical literature (e.g. pear flowers, makeup), while paragraph 3 uses modern poetic images (e.g. sparrows on the utility pole). Obviously, compared with poem 2, sentences in poem 3 seems more confusing, as the poetic images in modern times may not fit well into the language model of classical poems.

2) We also observed that poems generated from descriptive paragraphs achieve higher scores than from logical or philosophical paragraphs. For example, in Table 2, both paragraph 4 (more descriptive) and paragraph 5 (more philosophical) were

| Literature form | Fluency     | Semantic coherence | Semantic preservability | Poeticness  | Total       |
|-----------------|-------------|--------------------|-------------------------|-------------|-------------|
| Prose           | 2.52        | 2.30               | 2.30                    | 2.32        | 9.44        |
| Modern poem     | 2.37        | 2.34               | 2.01                    | 2.16        | 8.88        |
| Pop song lyric  | 2.40        | 2.31               | 2.24                    | 2.42        | 9.37        |
| Song lyric      | <b>2.62</b> | <b>2.54</b>        | <b>2.26</b>             | <b>2.49</b> | <b>9.91</b> |

Table 5: Human evaluation results for generating poems from various literature forms. We show the results obtained from our best model (Transformer+Anti OT&UT).

selected from famous modern prose. However, compared with poem 4, poem 5 seems semantically more confusing. We offer two explanations to the above phenomenon: **i.** Limited by the 28-character restriction, it is hard for quatrain poems to cover complex logical or philosophical explanation. **ii.** As vernacular paragraphs are more detailed and lengthy, some information in a vernacular paragraph may be lost when it is summarized into a classical poem. While losing some information may not change the general meaning of a descriptive paragraph, it could make a big difference in a logical or philosophical paragraph.

#### 4.6 Human Discrimination Test

We manually select 25 generated poems from vernacular Chinese translations and pair each one with its corresponding human written poem. We then present the 25 pairs to human evaluators and ask them to differentiate which poem is generated by human poet.<sup>3</sup>

As demonstrated in Table 6, although the general meanings in human poems and generated poems seem to be the same, the wordings they employ are quite different. This explains the low BLEU scores in Section 4.3. According to the test results in Table 7, human evaluators only achieved 65.8% in mean accuracy. This indicates the best generated poems are somewhat comparable to poems written by amateur poets.

We interviewed evaluators who achieved higher than 80% accuracy on their differentiation strategies. Most interviewed evaluators state they realize the sentences in a human written poem are usually well organized to highlight a theme or to build a poetic image, while the correlation between sentences in a generated poem does not seem strong. As demonstrated in Table 6, the last two sentences in both human poems (marked as red) echo each other well, while the sentences in machine-

<sup>3</sup>We did not require the expert group’s participation as many of them have known the gold poems already. Thus using their judgments would be unfair.

|         |   |
|---------|---|
| Human   | 黄沙磧里客行迷，四望云天直下低。<br>Within yellow sand moraine guest travels lost,<br>looking around found sky and clouds low.<br>为言地尽天还尽，行到安西更向西。<br>It is said that earth and sky ends here,<br>however I need to travel more west than anxi.                                       |
| Machine | 异乡客子黄沙迷，雁路迷寒云向低。<br>Guest in yellow sand gets lost,<br>geese are lost because clouds gets low.<br>只道山川到此尽，安西还要更向西。<br>It is said that mountains end here,<br>however anxi is even more west.  |
| Human   | 绝域从军计惘然，东南幽恨满词笺。<br>It’s hard to pay for the ambition of the military field,<br>the anxiety of situation in southeast is all over poems.<br>一箫一剑平生志，负尽狂名十五年。<br>A flute and sword is all I care about in my life,<br>15 years have failed the reputation of "madman". |
| Machine | 从军疆场志难酬，令人怅望东南州。<br>It’s hard to fulfill my ambition on the military field,<br>the situation in the southeast states are troublesome.<br>形容仗剑敌平戎，情怀注满赋雪愁。<br>I would like to use my sword to conquer my enemy,<br>yet my feelings are full of worry like the snow.    |

Table 6: Examples of generated poems and their corresponding gold poems used in human discrimination test.

generated poems seem more independent. This gives us hints on the weakness of generated poems: While neural models may generate poems that resemble human poems lexically and syntactically, it’s still hard for them to compete with human beings in building up good structures.

## 5 Discussion

**Addressing Under-Translation** In this part, we wish to explore the effect of different phrase segmentation schemas on our phrase segmentation-based padding. According to Ye (1984), most seven-character quatrain poems adopt the 2-2-3 segmentation schema. As shown in examples in Figure 3, we compare our phrase segmentation-based padding (2-2-3 schema) to two less common schemas (i.e. 2-3-2 and 3-2-2 schema) we report our experiment results in Table 8.



| Accuracy | Value |
|----------|-------|
| Min      | 52.0  |
| Max      | 84.0  |
| Mean     | 65.8  |

Table 7: The performance of human discrimination test.

青山隐隐水迢迢

- 2-2-3: 青山 <p><p> 隐隐 <p><p> 水迢迢 <p><p><p>
- 2-3-2: 青山 <p><p> 隐隐水 <p><p><p> 迢迢 <p><p>
- 3-2-2: 青山隐 <p><p><p> 隐水 <p><p> 迢迢 <p><p>

Figure 3: Examples of different padding schemas.

The results show our 2-2-3 segmentation-schema greatly outperforms 2-3-2 and 3-2-2 schema in both perplexity and BLEU scores. Note that the BLEU scores of 2-3-2 and 3-2-2 schema remains almost the same as our naive baseline (Without padding). According to the observation, we have the following conclusions: 1) Although padding better aligns the vernacular paragraph to the poem, it may not improve the quality of the generated poem. 2) The padding tokens should be placed according to the phrase segmentation schema of the poem as it preserves the semantic within the scope of each phrase.

**Addressing Over-Translation** To explore the effect of our reinforcement learning policy on alleviating over-translation, we calculate the *repetition ratio* of vernacular paragraphs generated from classical poems in our validation set. We found *naive transformer* achieves 40.8% in repetition ratio, while our *+Anti OT* achieves 34.9%. Given the repetition ratio of vernacular paragraphs (written by human beings) in our validation set is 30.1%, the experiment results demonstrated our RL loss effectively alleviate over-translation, which in turn leads to better generated poems.

## 6 Conclusion

In this paper, we proposed a novel task of generating classical Chinese poems from vernacular paragraphs. We adapted the unsupervised machine translation model to our task and meanwhile proposed two novel approaches to address the under-translation and over-translation problems. Experiments show that our task can give users more controllability in generating poems. In addition, our approaches are very effective to solve the prob-

| Padding schema | Perplexity   | BLEU        |
|----------------|--------------|-------------|
| 2-2-3          | <b>74.21</b> | <b>6.34</b> |
| 2-3-2          | 83.12        | 5.49        |
| 3-2-2          | 85.66        | 5.75        |

Table 8: Perplexity and BLEU scores of different padding schemas.

lems when the UMT model is directly used in this task. In the future, we plan to explore: (1) Applying the UMT model in the tasks where the abstraction levels of source and target languages are different (e.g., unsupervised automatic summarization); (2) Improving the quality of generated poems via better structure organization approaches.

## References

- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *ICLR*.
- Sameer Bansal, Herman Kamper, Adam Lopez, and Sharon Goldwater. 2017. Towards speech-to-text translation without speech recognition. In *EACL*.
- Pi-Chuan Chang, Michel Galley, and Christopher D. Manning. 2008. Optimizing chinese word segmentation for machine translation performance. In *WMT@ACL*.
- Kyunghyun Cho, Bart van Merriënboer, Çağlar Gülçehre, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using rnn encoder-decoder for statistical machine translation. In *EMNLP*.
- Yang Feng, Lin Ma, Wei Liu, and Jiebo Luo. 2019. Unsupervised image captioning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Jing He, Ming Zhou, and Long Jiang. 2012. Generating chinese classical poems with statistical machine translation models. In *AAAI*.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural Computation*, 9:1735–1780.
- Guillaume Lample, Alexis Conneau, Ludovic Denoyer, and Marc’Aurelio Ranzato. 2018a. Unsupervised machine translation using monolingual corpora only. In *ICLR*.
- Guillaume Lample, Myle Ott, Alexis Conneau, Ludovic Denoyer, and Marc’Aurelio Ranzato. 2018b. Phrase-based & neural unsupervised machine translation. In *EMNLP*.
- Haitao Mi, Baskaran Sankaran, Zhiguo Wang, and Abe Ittycheriah. 2016. Coverage embedding models for neural machine translation. In *EMNLP*.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2001. Bleu: a method for automatic evaluation of machine translation. In *ACL*.
- Romain Paulus, Caiming Xiong, and Richard Socher. 2018. A deep reinforced model for abstractive summarization. In *ICLR*.
- Steven J. Rennie, Etienne Marcheret, Youssef Mroueh, Jarret Ross, and Vaibhava Goel. 2017. Self-critical sequence training for image captioning. *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1179–1195.
- Baskaran Sankaran, Haitao Mi, Yaser Al-Onaizan, and Abe Ittycheriah. 2016. Temporal attention model for neural machine translation. *CoRR*, abs/1608.02927.
- Zhaopeng Tu, Zhengdong Lu, Yang P. Liu, Xiaohua Liu, and Hang Li. 2016. Modeling coverage for neural machine translation. In *ACL*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *NIPS*.
- Zhe Wang, Wei He, Hua Wu, Haiyang Wu, Wei Li, Haifeng Wang, and Enhong Chen. 2016. Chinese poetry generation with planning based neural network. In *COLING*.
- Wei-Hung Weng, Yu-An Chung, and Peter Szolovits. 2019. [Unsupervised clinical language translation](#). In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD ’19*, pages 3121–3131, New York, NY, USA. ACM.
- Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Lukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Gregory S. Corrado, Macduff Hughes, and Jeffrey Dean. 2016. Google’s neural machine translation system: Bridging the gap between human and machine translation. *CoRR*, abs/1609.08144.
- Linli Xu, Liang Jiang, Chuan Qin, Zhe Wang, and Dongfang Du. 2018. How images inspire poems: Generating classical chinese poetry from images with memory networks. In *AAAI*.
- Rui Yan. 2016. i, poet: Automatic poetry composition through recurrent neural networks with iterative polishing schema. In *IJCAI*.
- Cheng Yang, Maosong Sun, Xiaoyuan Yi, and Wenhao Li. 2018. Stylistic chinese poetry generation via unsupervised style disentanglement. In *EMNLP*.
- Jiaying Ye. 1984. *Poem Criticism with Jialin*. Zhonghua shu ju, Beijing, China.
- Xiaoyuan Yi, Maosong Sun, Ruoyu Li, and Zonghan Yang. 2018. Chinese poetry generation with a working memory model. In *IJCAI*.
- Xingxing Zhang and Mirella Lapata. 2014. Chinese poetry generation with recurrent neural networks. In *EMNLP*.
- Zhirui Zhang, Shuo Ren, Shujie Liu, Jianyong Wang, Peng Chen, Mu Li, Ming Zhou, and Enhong Chen. 2018. Style transfer as unsupervised machine translation. *CoRR*, abs/1808.07894.