

# Court Opinion Generation from Case Fact Description with Legal Basis

Quanzhi Li, Qiong Zhang

Alibaba Group, USA  
{quanzhi.li, qz.zhang}@alibaba-inc.com

## Abstract

In this study, we proposed an approach to automatically generating court view from the fact description of a legal case. This is a text-to-text natural language generation problem, and it can help the automatic legal document generation. Due to the specialty of the legal domain, our model exploits the charge and law article information in the generation process, instead of utilizing just the fact description text. The BERT model is used as the encoder and a Transformer architecture is used as decoder. To smoothly integrate these two parts together, we employ two separate optimizers for the two components during the training process. The experiments on two data sets of Chinese legal cases show that our approach outperforms other methods.

## Introduction

Previous studies in legal domain mainly focus on classifying legal documents, retrieving relevant cases, and predicting charge labels and law articles given the fact description of a case. An area that has not been explored much is the automatic generation of legal text, especially the text that have latent legal logic or knowledge, such as the court opinion section of a written judgment.

One of the most common legal document types is the document of written judgment. Almost for every case that is processed by the court or judges, there will be such a document. In the Chinese legal system, it is called “判决书” or “裁判文书”. In the judgment document, there is always a section called “court view” or “court opinion” (“本院认为” in Chinese), which explains the charge decision by the judges and usually also cites the applied law articles. This section is the core part of a legal judgment document. Another important section is the “fact description” (“案件事实”) part, which describes the facts of the case. Judges write the court view part based on the case facts, relevant law articles and charges. And following the court view part, there is a section called “court decision” or “judgment” (“判决结

果”), which is the court’s final decision based on the rationales described in the court view section. An example of a Chinese case judgment document is shown in Figure 1. This figure shows two parts of the judgment document, the fact description and the court opinion parts. The words in red indicates that it is the start of the corresponding part (fact description or court view). And the words in blue indicate the start of the final decision part.

This study focuses on generating the court view part, excluding the court final decision part. Automatic court opinion generation has many benefits and real applications. For example, it can reduce the workload of judges on writing court opinions, especially for simple cases but in large amount. It also benefits the individuals seeking legal advices. Nowadays, more and more people seek legal advice online. One common scenario is that the individual writes a paragraph of text that describes the facts, and wants to know the opinion of judges or legal professionals on that. Automatic court opinion generation is a good tool for this type of scenarios, especially for types of legal cases that are simple but with a large volume. Our algorithm is being deployed in several Chinese intelligent court systems and some legal consulting service platforms. These legal consulting service systems provide free online services for individuals, especially for the ones in the remote areas in China, where it is hard for low-income people to access in-person legal service, due to its expensiveness and the lack of legal professionals.

Court view generation is a text-to-text generation problem, but compared to other text generation problems, such as text summarization, it is more challenging, since it needs to deduce the latent information from the facts, and generate the opinion text that is both law and charge discriminative. In order to generate such type of court view, we proposed a court opinion generation model that has the following features: 1. It is based on the Transformer (Vaswani et al.,

### **Fact Description:**

经审理查明, 2018年7月6日晚, 被告人XXX携带刀具进入被害人刘某家中偷窃钱财, 在偷窃过程中被刘某及家人发现, 在逃跑过程中因和刘某撕打而持刀刺伤刘某, ...

After hearing, the court identified that in the evening of July 6<sup>th</sup>, 2018, the defendant XXX entered Liu YY's home with a knife, and tried to steal money and jewelries. He was seen by Liu and Liu's family members, and hurt Liu when he tried to escape. ...

### **Court Opinion:**

本院认为, 被告人XXX在盗窃过程中持有凶器, 在抗拒抓铺过程中致被害人轻伤, 其行为已经构成抢劫罪。依照中华人民共和国刑法第二百六十九条之规定, 判决如下: ...

The court holds that during the process of stealing the defendant XXX wound Liu YY with a knife and caused Liu minor injury. XXX's action has constituted the crime of robbery. According to the law article 269 of the Criminal Law of the People's Republic of China, the court decision is as follows: ...

Figure 1. A legal case example, showing the fact description and court view sections.

2017) architecture, but we use the BERT model (Devlin et al., 2019) as the fact description and law article encoder, to better exploit the grammar and semantic information of tokens learned by the pre-trained model. 2. Since the BERT encoder is pre-trained and the decoder is not, in order to smoothly integrate these two parts together for better generation performance, we employ two separate optimizers for these two components during the training process. 3. To make the court view more law and charge discriminative, we utilize the charge and law article information in the generation process; they are encoded and attended by the Transformer decoders. The above three points are also the main contributions of this study.

Automatic law article prediction has been studied by several studies, and there are already many studies on charge label prediction, so these two prediction tasks are not the focus of this study. In this study, we extract the charge label and law articles from the case documents. Jointly learning the court view generation and prediction of charge label and law articles will be one of our future studies.

## **Related Studies**

Previous studies on legal assistant system mainly focus on the fields of legal case retrieval (Chen et al. 2013; Raghav et al. 2016), case classification (Nallapati and Manning, 2008; Wang et al. 2018), legal case data set analysis (Locke and Zuccon, 2018; Du et al. 2019), charge label prediction and law article prediction (Liu and Liao, 2005; Liu and Hsieh, 2006; Liu et al. 2015; Luo et al. 2017; Yang et al. 2019b). They use either the traditional machine learning models or neural networks. Ye et al. (2018) use LSTM based sequence to sequence (seq2seq) model with attention (Sutskever et al. 2014; Vaswani et al. 2017) to generate court view from fact description.

Sequence to sequence learning has been used in a variety of language generation applications. It has attracted much attention in recent years due to the advance of deep learning. Our model also belongs to this widely used seq2seq paradigm (Sutskever et al. 2014). Rush et al. (2015) and Nallapati et al. (2016) were among the first to use the neural

encoder decoder structure in text summarization. See et al. (2017) enhance this model with a pointer generator network which allows it to copy words from the source text. Paulus et al. (2018) present a deep reinforced model for abstractive summarization which handles the coverage problem with an intra-attention mechanism. Celikyilmaz et al. (2018) propose an abstractive system where multiple agents represent the document together with a hierarchical attention mechanism for decoding. Narayan et al. (2018) propose an abstractive model suitable for single sentence summaries.

Pre-training has been widely used in natural language processing (NLP) tasks to learn better language representation, and several new pre-trained models have been published recently, such as BERT (Devlin et al. 2019), XLNet (Yang et al. 2019a), RoBERTa (Liu et al. 2019), ALBERT (Lan et al. 2019), ELMO (Peters, et al. 2018), etc. The pre-training on large amount of unlabeled data and fine-tuning with small scale labeled data are helpful for many tasks, and it is also used in the encoder part of our model in this work. Devlin et al. (2019) proposed BERT based on masked language modeling and next sentence prediction, and achieved state-of-the-art results on multiple NLP tasks. There are also some works on pre-training the encoder-decoder model for language generation (Rothe et al., 2019; Edunov et al., 2019; Liu and Lapata, 2019; Song et al., 2019). Liu and Lapata (2019) introduce a document-level encoder based on BERT and propose a fine-tuning schedule which adopts different optimizers for the encoder and the decoder. The main differences between our model and others are that we utilize the charge and law article information in the decoding process, in order to generate judge opinions that have more legal basis. Our model also uses the pre-trained BERT model as encoder and a non-pre-trained Transformer as decoder, and they use two separate optimizers.

## **Opinion Generation Model**

In this section, we first describe the high-level encoder-decoder structure, and then present the fact description encoder, law article encoding, and the Transformer decoder in details. Figure 2 shows the high-level model structure. The

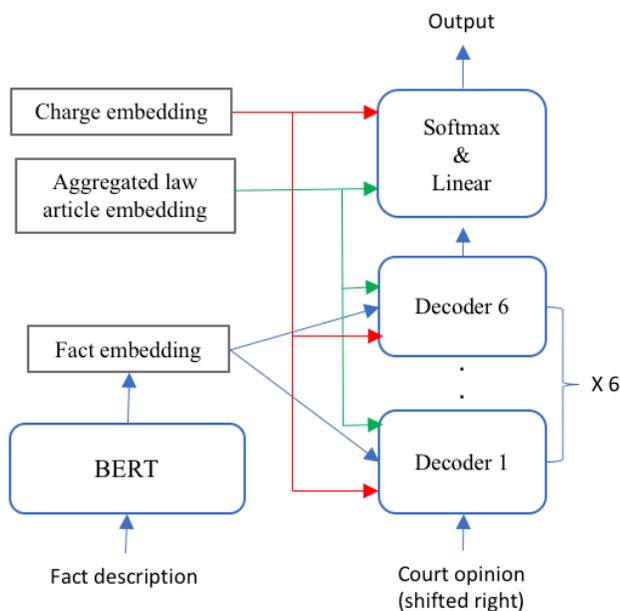


Figure 2. The high-level encoder-decoder structure.

fact description is encoded by a pre-trained BERT model, shown on the left side of the diagram. The charge label embedding is learned during the training process. Each case may have multiple law articles applied, and their embeddings are aggregated together to one vector representing all the articles for this case. The right side is the decoder part, which is based on the Transformer architecture, consisting of 6 layers of decoders. The fact embedding, charge embedding and the aggregated law article embedding are fed to each of the decoder layers, and used by the multi-head attention layer described later. The charge embedding and law article embedding are also used at the final softmax layer of the Transformer decoder, in order to add more context to help the decoder choose the correct token.

## Fact Description Encoder and Law Article Encoder

### Fact Description Encoding

We use BERT to encode fact description and law articles. BERT has been used to fine-tune various NLP tasks, but its application to text generation is not straightforward, since it is trained to predict single word and next sentence, not generating text sequence. This is why our model uses BERT to encode the fact description, but utilizes Transformer on the decoding side to generate text sequence. The BERT-Chinese model is fine-tuned by the fact descriptions from the training data, and the output from the last layer is fed to the decoder side.

### Law Article Encoding

In the Chinese legal system, each law article has a description defining the content of that article and in what situations

that article should be applied. When judges express their view on a case, their opinions are based on the facts and the applied law articles, and the languages they use are usually affected by the definition of the related law articles. Sometimes some sentences or phrases in the court view part are the same as the sentences or phrases in the articles. Therefore, we hypothesize that the textual content of law articles will help the court view generation.

For a law article, we also use the pre-trained BERT-Chinese model to generate its embedding by fine-tuning the model. Since a case may involve multiple law articles, we aggregate these law article embeddings, and generate a final representation and feed it into the Transformer decoder. For each article, we take the [CLS] token vector from the last layer to represent the article. And then for all the articles of this case, we take the average of their embeddings. The final embedding is fed into the decoder side for multi-head attention.

## Transformer Decoder

Figure 2 shows that there is a stack of six decoders in the decoding side. There are two multi-head attention layers in each decoder, one is a masked multi-head self-attention, and the other one is the decoder attended on the three types of contexts, i.e. the fact description embedding, charge label embedding and law article embedding. After each attention layer and the feedforward layer, there is a Add&Normalize layer.

### Multi-head Attention

In the multi-head attention structure, attention is computed not once but multiple times, in parallel and independently. The outputs are concatenated and linearly transformed. The following equation shows how one head of the scaled dot-product attention is computed:

$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (1)$$

where  $Q$  is a matrix that contains the query (vector representation of one word in the sequence),  $K$  are all the keys (vector representations of all the words in the sequence), and  $V$  are the values. For the first multi-head attention layer in each decoder,  $V$  consists of the same word sequence as  $Q$ . However, for the second attention module in each decoder, it considers the decoder sequence, the fact embedding, charge embedding and law article embedding, and therefore, here  $V$  is different from the sequence represented by  $Q$ .

## Model Training and Inference

In our model, the encoder is based on a pre-trained BERT model, and the decoder component is not pre-trained. It is obvious that there is a mismatch between these two components, because one is pre-trained and the other one needs to be completely trained. This may make the training process unstable, e.g. one component is underfitted and the other one

is overfitted. One way to handle this is to use two different optimizers for the two components. In our model implementation, two Adam optimizers are used. Each has its own learning rates and warmup steps. These values will be set so that the pre-trained BERT model should be fine-tuned with a smaller learning rate and decay, when the decoder becomes stable. This is to make sure that the BERT model to be fine-tuned with more accurate gradients. The learning rate update is illustrated by the following equation:

$$\text{learnRate} = \text{learnRate} * \min(s^{-0.5}, s * w^{-1.5}) \quad (2)$$

where  $s$  is the step and  $w$  is the warmup value.

For inference, we use beam search, whose size is set to 4, to find the best sequence. The generated word sequences will be ranked and the one with the largest value will be chosen.

## Experiments and Results

### Data Set Preparation and Evaluation Metrics

#### Data Sets

Our experiments of comparing different generation approaches were conducted on two data sets. Data set 1 was created by ourselves, and data set 2 was from (Ye et al., 2018). The data in both data sets are originally from the legal documents in China Judgements Online<sup>1</sup>. There is no duplicate between these two data sets. Data set 2 contains criminal cases, while data set 1 have many more cases and covers more sub-types. The basic statistics of the two data sets are presented in Table 1. The total number of unique law articles in these two data sets are 308, and there are 78 unique charges. On average, each case has about two law articles applied to. In this study, we only focus on the cases with one charge; the complicated cases are left for future study. The full list of charges in Chinese legal system is available online<sup>2</sup>.

#### Data Preparation

Similar to previous studies on Chinese legal cases, we extract the fact descriptions, court opinions, law articles and charge labels using regular expressions. These cases are written by judges, and the format are pretty standard. The paragraph started with “经审理查明” (“the court identified that”) is extracted as the fact description. The part between “本院认为” (“the court holds that”) and the decision “判决如下” (“the court decision is as follows”) are extracted as the rationales/opinions/views. Nearly all the samples in the two data sets match this pattern. In order to meet the text length limitation of pre-trained BERT, a threshold of 510 tokens is set up for both fact description and court opinion

parts. Facts and opinions longer than that are stripped, and they will be studied in the future. Charges and law articles are also easy to be extracted using regular expressions.

#### Special Tokens

We use “<DATE>” to replace dates in the data sets. In the Chinese legal system, numbers and entity names may affect the applied laws, charges and the final decision. For example, for a drug trafficking or bribery case, the amount of drug and money will definitely affect the charge, applied law articles and the number of years in prison. To handle this, we use a data-driven approach – we extracted all the numbers from the two data sets, and based on the analysis of these numbers, we split them into 20 buckets based on their magnitudes. These 20 buckets are represented by 20 symbols. Types of entities may also affect a case in terms of the charge, applied law articles and the final judgment. For example, in a bribery case, a government employee and a private company employee would get different charges. To handle this case, we converted the entities of private companies, state-owned companies and government agents to three different symbols. Government agents are easy to identify by using regular expression and rules. We use a lookup service, which provides lookup of most state-owned companies, to check if a company is state-owned. Individual names are converted to special symbol.

#### Law Article Data Set

As mentioned before, we need the definition of law articles to generate their embeddings. We processed 308 law articles. Their descriptions/definitions are collected from the website of the Chinese Highest Court<sup>3</sup>. One article example is shown in Figure 3. A charge label may connect to one or more articles, and one article can also connect to one or more charge labels.

#### Evaluation Metrics

To evaluate the effectiveness of our approach, we use a quantitative comparison between the true court opinion, the one extracted from the real case, and the automatically generated one. This approach has been used by many previous studies in text summarization tasks. We use both BLEU-4

Data set attribute	Data set 1	Data set 2
Training data	521,500	153,706
Validation data	25,000	9,152
Test data	30,000	9,123
Average number of tokens in fact description part	248.2	219.9
Average number of tokens in court opinion part	34.4	30.6

Table 1. Some statistics of the two data sets

<sup>1</sup> <http://wenshu.court.gov.cn>

<sup>2</sup> <https://www.zuiming.net/51.html>

<sup>3</sup> <http://www.dffy.com/faguixiazai/xingfa/200311/20031110213247.htm>  
[http://www.npc.gov.cn/wxzl/wxzl/2000-12/06/content\\_4470.htm](http://www.npc.gov.cn/wxzl/wxzl/2000-12/06/content_4470.htm)

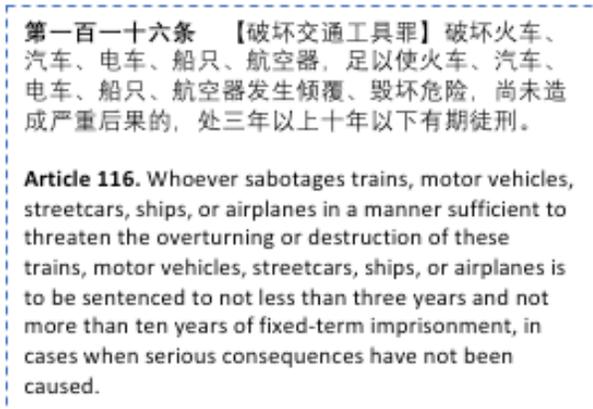


Figure 3. One example of Chinese law articles.

score (Papineni et al. 2002) and three ROUGE scores (Lin and Hovy, 2003) as the evaluation metrics.

ROUGE works by comparing an automatically produced text against the reference text. It basically measures the overlap of N-grams between the system and reference texts. However, the original ROUGE measure does not tell you much as a metric. To get a good quantitative value, in the context of ROUGE, we compute precision and recall using the overlap, and then report the F1-measure of ROUGE. In this study, we use ROUGE-1, which is based on unigram overlap, ROUGE-2, which is based on bigram overlap, and also ROUGE-L, which measures the longest common subsequence of reference and generated texts, to compute the F1 scores.

## Compared Methods and Experiment Settings

### Compared Methods

To evaluate our proposed approach, we compared it to the following approaches, including the basic approaches and the state-of-the-art ones.

- **Random & Random<sup>+charge&law</sup>**: this method is to randomly select court opinions from the training set. The Random<sup>+charge&law</sup> version will select court view from the training cases that have the same charge label and law articles as the test case. We want to see how much improvement we can get by limiting the case pool to the ones with same charge and law articles. Random method is to give the low bound performance of the court opinion generation task.
- **BM25 & BM25<sup>+charge&law</sup>**: this method uses BM25 score (Robertson and Zaragoza, 2009) to retrieve the case whose fact description has the highest BM25 score, and use its court view as the result. The rationale behind this method is that similar fact descriptions may have the similar court views. Similar to The Random<sup>+charge&law</sup>, The BM25<sup>+charge&law</sup> will retrieve the court

view from the training cases that have the same charge label and law articles as the test case.

- **Seq2seq-LSTM-attention**: this approach is based on the seq2seq model with attention, and both the encoder and decoder use a bidirectional-LSTM model (Sutskever et al., 2014; Vaswani et al., 2017; Ye et al., 2018). Attention mechanism can catch the important input information for the current output sequence.
- **Seq2seq-LSTM-attention<sup>+charge&law</sup>**: this method uses the Seq2seq-LSTM-attention method described above, and also it includes the charge label and law article information in the decoder side.
- **DCA**: Celikyilmaz et al. (2018) use multiple encoders to represent the document together with a hierarchical attention mechanism for decoding. Their proposed Deep Communicating Agents (DCA) model is trained end-to-end with reinforcement learning.

### Experimental Settings and Training

Like previous studies, we use the validation data set to tune our model and hyper-parameters. The maximal length is 510 for both the fact descriptions and court views, although most court views in our data sets have less than 60 words. Model performance was checked on the validation set after every 300 batches training. Training process was terminated if the model performance is not improved for successive 10 times.

**Encoder.** We use the BERT-Chinese pre-trained model as the encoder for encoding fact description. This model uses 12 encoder layers, and the embedding size is 768 for the input token, the position embedding and the text segmentation embedding. The multi-head attention has 12 heads, drop out is 0.1, L2 decay rate is 0.01, and activation function is Gaussian Error Linear Unit (GELU). The aggregated law article embedding is generated by the approach described before.

**Decoder.** The decoder side has six layers, as illustrated in Figure 2. The input is shifted one token position to the right, utilizing a teacher forcing learning approach. For the encoder-decoder multi-head attention in each of the six decoder layers, the inputs are the fact description embedding from the BERT encoder, the charge embedding, the aggregated law article embedding, and the output from the last decoder layer (or the court opinion input for the first layer). The charge embedding size is 768 and is also randomly initialized. The other hyper-parameters on the decoder side use the default values of the original Transformer architecture.

As described before, in our model, the BERT encoder and the decoder use two different optimizers. Their learning rates are different. We use two Adam optimizers with  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$  for the encoder and the decoder, respectively, but they have different learning rates and warmup-steps. In Equation 2, we set  $learnRate = 2e-3$  and  $warmUp = 30,000$  for the BERT encoder, and we set  $learnRate = 0.05$

and  $warmUp = 15,000$  for the decoder. For the whole model, the batch size is 32 and dropout is 0.1.

## Evaluation Results and Analysis

### Experiment Result

Table 2 and Table 3 present the comparison result on the two data sets. B-4 means BLEU-4 value. R-1, R-2, R-L are the F1 values of ROUGE-1, ROUGE-2 and ROUGE-L, respectively. The result shows that our proposed approach performs clearly better than all the other methods. To verify if the performance improvement is statistically significant, we conducted *t-test* (Rice, 2006) between our model and the best baseline (i.e. DCA). The *t-test* results show that the performance improvement is statistically significant at the level of  $p=0.01$  on both data sets, and for all the four evaluation metrics. The scores on data set 1 are lower than that on data set 2, and the reason is that the average length of court view in data set 1 are longer than that of data set 2, and data set 1 covers more types of criminal cases.

From the two tables we can also notice that even the Random+charge&law method has a relatively high score, which means that the expressions of court view with the same charge label and law articles are similar, with many overlapped n-grams. We can also notice that the BM25+charge&law method, which is based on retrieving the case having similar fact description and also having the same charge label and law articles, performs pretty well. The seq2seq model with attention mechanism performs better than the BM25+charge&law model. Adding charge and law information to the Seq2seq-LSTM-attention model will improve the performance more.

### Ablation Test

The underlying ideas behind our model are the following two main hypotheses: 1. Adding charge and law article information in the decoding process will improve court view generation performance, since it will provide more charge-discrimination information and article related information. 2. Our generation model with BERT’s pre-trained model as the encoder and utilizing two separate optimizers for encoder and decoder will improve the court view generation performance. The results in Table 2 and 3 have proven that these two features together make our model outperform all other approaches. In this ablation test, we want to see the effect of each feature. The ablation test was conducted on the data set 1. The results are presented in Table 4. The results show that removing either charge or law article information from the model will decrease its performance, which validates our hypothesis 1. If both of them are not included in our model, then it still outperforms the seq2seq-LSTM-attention model (result line 5 in Table 4 vs. result line 5 in Table 2), which demonstrates that our hypothesis 2 is valid. This test also shows that adding charge and law article

Method	R-1	R-2	R-L	B-4
Random	24.1	5.1	23.6	5.2
Random+charge&law	52.3	27.3	47.8	23.4
BM25	61.3	41.4	58.1	37.8
BM25+charge&law	65.3	46.0	62.2	41.1
Seq2seq-LSTM-attention	66.9	48.1	63.6	41.8
Seq2seq-LSTM-attention+charge&law	68.4	50.1	65.2	42.9
DCA	69.7	51.4	66.8	44.6
Our approach	<b>71.7</b>	<b>53.2</b>	<b>68.5</b>	<b>46.7</b>

Table 2. Algorithm comparison result on data set 1

Method	R-1	R-2	R-L	B-4
Random	26.5	6.2	25.1	6.4
Random+charge&law	54.7	29.8	50.2	25.7
BM25	63.5	43.7	60.3	40.1
BM25+charge&law	67.8	48.2	64.5	43.4
Seq2seq-LSTM-attention	69.1	50.3	65.9	44.0
Seq2seq-LSTM-attention+charge&law	70.9	52.5	67.7	45.8
DCA	71.9	52.8	68.3	46.9
Our approach	<b>73.8</b>	<b>54.6</b>	<b>70.6</b>	<b>48.7</b>

Table 3. Algorithm comparison result on data set 2

Method	R-1	R-2	R-L	B-4
Our approach	71.7	53.2	68.5	46.7
- without charge	70.3	52.2	67.1	45.2
- without law	69.9	51.7	66.6	44.5
- without charge&law in the final output layer	70.5	52.4	67.4	45.4
- without charge & law	69.2	50.6	65.9	43.7

Table 4. Ablation test result using data set 1

embeddings to the final linear and softmax layer will improve the generation performance (result line 5 in Table 4), since it provides more context for decoder to choose the right token.

### Length of Court Opinion

We wanted to see the relation between the length of the court view and the model performance. Figure 4 shows how the BLEU-4 score changes with the court view length. At length of 15, we have the best performance, and the score decreases as the length grows, with score of 48 at length of 30 and going down to 32 at length of 50. This is not surprising, because when the text to be generated is long, it is hard to capture the correct source information when generating a token, not matter the model is based on Transformer or other seq2seq models.

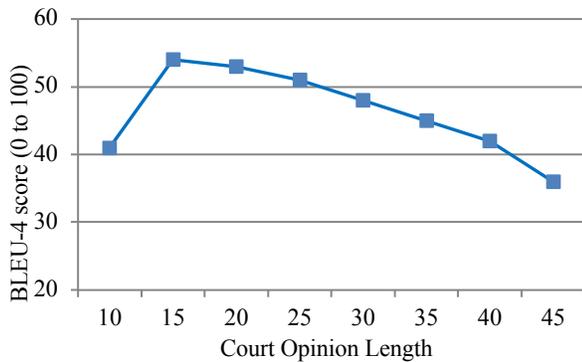


Figure 4. Model performance vs. court opinion length. The length is based on the number of tokens

### Effect of Normalizing Numbers and Entities

For numbers and entities appearing in facts description, unlike previous studies converting them to just two symbols, in this study, we differentiate the entities into four symbols and the converting numbers into 20 symbols, by analyzing the training data. We expect this change will make the fact description more distinctive for certain cases. By analyzing some examples, we found that this do make the generated court view more accurate. For example, there is a case about a government employee conducting the crime of corruption with a huge amount of money. We found that if we convert the names and numbers as other studies do, or do not do any conversion, the generated court view will have a BLEU-4 value of 0.33. In contrast, after the conversion using our approach, the score changes to 0.45. The reason is that without conversion, the view is generated as a crime of regular professional embezzlement, but with the conversion the generated view correctly tells that this is about corruption with public economy involving “数额特别巨大” (huge amount of money), and the applied laws and charge are also different.

### Limitation and Conclusion

Below we discuss some limitations, and then give a conclusion of this study.

**Legal fact points.** Court opinion generation is more challenging than the general text generation problem, since the court view contain legal logic and law related information. Our model has employed law article embedding model and charge labels to make the court view more law and charge discriminative, and the view generation process directed by related law information. And by using BERT, instead of the general encoder part of a Transformer model, our approach can take advantage of the large amount of case judgment documents available to train a legal-aware model. But this is not enough, and one direction of improving our model is to incorporate the legal fact points of law articles into the model. Each law article explicitly or implicitly defines a set

of fact points and related conditions. We can exploit a legal knowledge base, which contains the fact points and reasoning logic derived semi-automatically from law articles and legal documents, to achieve this. These fact points can be used to extract facts from the fact description, and then they can be incorporated into the decoding side to direct the court view generation.

**Possible harm.** One application of this model is to be deployed in some online legal services or devices located in some rural areas, to help individuals seeking legal advices, especially the low-income and poor people. But it does have limitations and possible harms. For example: 1. Some rural areas in China do not have internet yet, or people cannot access internet due to their financial limitation, and therefore they cannot access this system. This may cause more inequality in terms of access to legal services. 2. The legal opinions generated by this model may not be always correct. It is possible that the wrong (or even correct) opinion may intimidate a user, especially when the opinion implies that the user would lose the case if he/she takes it to the court. This may discourage the user from seeking further legal advices or actions.

**Fairness.** Fairness concern has been raised in the AI community in the past few years. If the training data reflect human biases, the outcome of the system is going to be biased. By incorporating the law article information in our model, our model can mitigate this type of bias in the output to certain degree. The description and information in the law articles are objective and usually not biased. In the decoding part of the model, compared to a model without the law article embedding component, the law information will make the generated view more objective, the language more law related, and the view less subjective to the entities (e.g. minority or majority) mentioned in the source description. Our model uses BERT as the encoder, and this gives the model the ability to learn semantic information from a huge amount of legal data. If the data set is large enough and representative, it can reduce the bias, but it still cannot completely solve this problem.

**Conclusion.** In this study, we proposed a new method for automatically generating court opinion from the fact description of a legal case. Our model exploits the charge and law article information in the generation process, in addition to the fact description text. We use the BERT model as the encoder and a Transformer architecture as decoder. Our model employs two separate optimizers for training the encoder and decoder. The experiments on two data sets of Chinese legal cases show that our approach outperforms other approaches. One of the future studies in this direction will be to incorporate the predictions of charge labels and law articles into the court view generation process, which will be a jointly learning task.

## References

- Celikyilmaz, A; Bosselut, A; He, X and Choi, Y. 2018. Deep communicating agents for abstractive summarization. *The 2018 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL 2018)*.
- Chalkidis, I; Androutsopoulos, I; Aletras, N; Neural Legal Judgment Prediction in English, *The 57th Annual Meeting of the Association for Computational Linguistics (ACL 2019)*
- Chen, Y.-L.; Liu, Y. and Ho, W. 2013. a text mining approach to assist the general public retrieval of legal documents. *The Journal of the Association for Information Science and Technology (JASIST 2013)*
- Devlin, J.; Chang, M.-W; Lee, K; and Toutanova, K. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. *The 2019 Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL 2019)*.
- Du, W.; Li, B.; Yang, M.; Qu, Q.; Shen, Y.; 2019. A Multi-Task Learning Approach for Answer Selection: A Study and a Chinese Law Dataset, The Thirty-third AAAI Conference on Artificial Intelligence (AAAI 2019)
- Edunov, S.; Baevski, A. and Auli M. 2019. Pre-trained language model representations for language generation. *The 2019 Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL 2019)*.
- Lan, Z.; Chen, M.; Goodman, S., et al. 2019. ALBERT: A Lite BERT for Self-supervised Learning of Language Representations, <https://arxiv.org/abs/1909.11942>, 2019
- Lin, C.-Y. and Hovy, E. 2003. Automatic evaluation of summaries using n-gram co-occurrence statistics. *The 2003 Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL 2003)*.
- Liu, C.-L. and Hsieh, C.-D. 2006. Exploring phrase-based classification of judicial documents for criminal charges in chinese. *In Foundations of Intelligent Systems, 16th International Symposium (ISMIS 2006)*
- Liu, C.-L. and Liao, T.-M. 2005. Classifying criminal charges in chinese for web-based legal services. *In Web Technologies Research and Development - APWeb 2005, 7th Asia-Pacific web conference*.
- Liu, Y. and Lapata, M. 2019. Text Summarization with Pre-trained Encoders, Conference on Empirical Methods in Natural Language Processing (EMNLP 2019)
- Liu, C.-L, Chang, C.-T. and Ho, J.-H. 2004. Case instance generation and refinement for case-based criminal summary judgments in Chinese, *Journal of Information Science and Engineering*, 20(4):783-800
- Liu, Y.; Ott, M. and et al., 2019, RoBERTa: A Robustly Optimized BERT Pretraining Approach, <https://arxiv.org/abs/1907.11692>
- Liu, Y.-H.; Chen, Y.-L. and Ho, W.-L. 2015. Predicting associated statutes for legal problems. *Information Processing and Management (IPM 2015)*
- Locke, D and Zuccon, G. 2018. A Test Collection for Evaluating Legal Case Law Search, *International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2018)*
- Luo, B.; Feng, Y.; Xu, J.; Zhang, X.; Zhao, D. 2017. Learning to predict charges for criminal cases with legal basis. Conference on Empirical Methods in Natural Language Processing (EMNLP 2017)
- Nallapati, R.; Zhou, B.; Nogueira, C.; Gulcehre, C.; Xiang, B. 2016. Abstractive text summarization using sequence-to-sequence RNNs and beyond. *Conference on Computational Natural Language Learning (CoNLL 2016)*
- Nallapati, R. and Manning, C.D. 2018. Legal Docket-Entry Classification: Where Machine Learning stumbles, Conference on Empirical Methods in Natural Language Processing (EMNLP 2018)
- Narayan, S.; Cohen, S.B and Lapata, M. 2018. Ranking sentences for extractive summarization with reinforcement learning. *The 2018 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL 2018)*
- Peters, M.; Neumann, M.; Iyyer, M. and et al. 2018. Deep contextualized word representations, *The 2018 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL 2018)*
- Papineni, K.; Roukos, S.; Ward, T.; Zhu, W.-J. 2002. Bleu: a method for automatic evaluation of machine translation. The 40th Annual Meeting of the Association for Computational Linguistics (ACL 2002)
- Paulus, R.; Xiong, C.; Socher, R. 2017. A Deep Reinforced Model for Abstractive Summarization, 2017, <https://arxiv.org/abs/1705.04304>
- Radford, A.; Wu, J.; et al. 2019. Language models are unsupervised multitask learners. <https://www.semanticscholar.org/paper/Language-Models-are-Unsupervised-Multitask-Learners-Radford-Wu/9405cc0d6169988371b2755e573cc28650d14dfe>
- Raghav, K.; Reddy, P.K; and Reddy, V.B. 2016. Analyzing the extraction of relevant legal judgments using paragraph-level and citation information. *In AI4J Artificial Intelligence for Justice*.
- Rice, J. A. 2006. *Mathematical Statistics and Data Analysis* (3rd ed.). Duxbury Advanced
- Robertson, S. and Zaragoza, H. 2009. The Probabilistic Relevance Framework: BM25 and Beyond, *Foundations and Trends in Information Retrieval 2009*
- Rothe, S.; Narayan, S. and Severyn, A. 2019. Leveraging pre-trained checkpoints for sequence generation tasks. *arXiv preprint arXiv:1907.12461*, 2019
- Rush, A.M; Chopra, S. and Weston, J. 2015. A neural attention model for abstractive sentence summarization. Conference on Empirical Methods in Natural Language Processing (EMNLP 2015)
- See, A.; Liu, P.J. and Manning, C.D. 2017. Get to the point: Summarization with pointer generator networks. The 55th

Annual Meeting of the Association for Computational Linguistics (*ACL 2017*)

Song, K.; Tan, X.; Qin, T.; Lu, J.; Liu, T.Y. 2019. MASS: Masked Sequence to Sequence Pre-training for Language Generation, <https://arxiv.org/pdf/1905.02450.pdf>, 2019

Sutskever, L.; Vinyals, O. and Le, Q.V. 2014. Sequence to sequence learning with neural networks. *Conference on Neural Information Processing Systems (NIPS 2014)*

Vaswani, A.; Shazeer, N.; Parmar, N., and et al. 2017. Attention is all you need. *Conference on Neural Information Processing Systems (NIPS 2017)*

Wang, P.; Yang, Z.; Niu, S.; Zhang, Y.; Zhang, L.; Niu, S. 2019. Modeling Dynamic Pairwise Attention for Crime Classification over Legal Articles, *International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2019)*

Yang, Z.; Dai, Z.; Yang, Y.; Carbonell, J.; Salakhutdinov, R.; Le, Q.V. 2019. XLNet: Generalized Autoregressive Pre-training for Language Understanding, *Conference on Neural Information Processing Systems (NIPS 2019)*

Yang, W.; Jia, W.; Zhou, X.; and Luo, Y. 2019. Legal Judgment Prediction via Multi-Perspective Bi-Feedback Network, *International Joint Conferences on Artificial Intelligence (IJCAI 2019)*.

Ye, H.; Jiang, X.; Luo, Z.; Chao, W. 2018. Interpretable Charge Predictions for Criminal Cases: Learning to Generate Court Views from Fact Descriptions, *The 2018 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL 2018)*

Zhong, H.; Guo, Z.; Tu, C., et al. 2018. Legal Judgment Prediction via Topological Learning, *Conference on Empirical Methods in Natural Language Processing (EMNLP 2018)*