

Fine-Grained Element Identification in Complaint Text of Internet Fraud

Tong Liu^{1*}, Siyuan Wang^{1*}, Jingchao Fu¹, Lei Chen¹, Zhongyu Wei^{1†},
Yaqi Liu², Heng Ye², Liaosa Xu², Weiqiang Wang², Xuanjing Huang¹

¹School of Data Science, Fudan University, China; ²Ant Group, China

{18210980056, wangsy18, chenl18, zywei, xjhuang}@fudan.edu.cn

fuuuuugcn@gmail.com; {yaqiliu.lyq, daokun.yh, liaosa.xls, weiqiang.wwq}@antgroup.com

ABSTRACT

Existing system dealing with online complaint provides a final decision without explanations. We propose to analyse the complaint text of internet fraud in a fine-grained manner. Considering the complaint text includes multiple clauses with various functions, we propose to identify the role of each clause and classify them into different types of fraud element. We construct a large labeled dataset originated from a real finance service platform. We build an element identification model on top of BERT and propose additional two modules to utilize the context of complaint text for better element label classification, namely, global context encoder and label refiner. Experimental results show the effectiveness of our model.

CCS CONCEPTS

• Applied computing → Secure online transactions.

KEYWORDS

complaint text mining, fraud element identification, dataset

ACM Reference Format:

Tong Liu¹ [1], Siyuan Wang¹ [1], Jingchao Fu¹, Lei Chen¹, Zhongyu Wei¹ [2], Yaqi Liu², Heng Ye², Liaosa Xu², Weiqiang Wang², Xuanjing Huang¹. 2021. Fine-Grained Element Identification in Complaint Text of Internet Fraud. In *Proceedings of the 30th ACM International Conference on Information and Knowledge Management (CIKM '21)*, November 1–5, 2021, Virtual Event, QLD, Australia. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/3459637.3482108>

1 INTRODUCTION

In the field of e-commercial business, various payment platforms provide an easy way of capital transferring but also exposes huge threat of Internet fraud. Every year, financial companies receive thousands of fraud complaints, from imprudence remittance to mendacious business contract. Although complaint text is usually full of sore and loss, it is valuable for organizations to understand incident mechanism and avoid potential risks [1, 5]. Handling fraud

* Tong Liu and Siyuan Wang have equal contribution.

† Corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
CIKM '21, November 1–5, 2021, Virtual Event, QLD, Australia

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-8446-9/21/11...\$15.00

<https://doi.org/10.1145/3459637.3482108>

complaints is of necessity to resolve victims' burning issues and recognizing fraud elements enables financial organizations and law enforcement agencies to detect security flaws and put an end to unfair and misleading swindle practices.

However, processing numerous complaint texts throws huge challenge for Internet finance companies. Different from credit and trading record, fraud complaints are composed of informal and unstructured text that takes painstaking efforts to understand and analyze. Besides, the inspect outcome is usually short (accept further investigation or not), which might make victims indignant for they cannot figure out why their loss never comes back.



Figure 1: Ideal pipelines of fraud complaints processing.

Figure 1 introduces a more well-grounded scheme for complaint text processing and fraud judgement. When a user realizes being swindled and submits complaints to the transfer platform, the third-party payment service provider gets involved in and segments the complaint text into pieces so as to inspect fraud elements inside. The text on the right briefly shows the three fundamental elements to ascertain a fraud: (1) the fraudster deliberately *fabricate* an illusory circumstance to ask for money, (2) the *remittance* is material and (3) the injured party *realizes* the falsehood afterwards and asks for compensation. With fraud elements automatically detected, it is more convenient for inspectors to check whether the statement is valid and more interpretable for users to comprehend the inspect results. Once the fraud case is established, the third-party and legal institution will follow the hints of fraud elements to investigate the transaction history and then present a fairer arbitrament.

To this end, we propose a novel task to identify fraud elements in a complaint paragraph. We design an annotation criterion, split paragraphs into clauses and ask undergraduates to annotate each clause. Moreover, we analyze clause distribution in complaint paragraphs and explore the connection between successional clauses while finding that position and global coherence are influential for identifying clause role. Therefore, we propose a hierarchical architecture which integrate context information to obtain more accurate prediction of clauses. Our contributions are of three-folds:

- We propose a novel task and construct a dataset for fine-grained fraud element identification on clause level to further analyse internet fraud issues.

- We formulate the task as a form of sequence labeling as plentiful analysis indicates that position and relation are significant features for complaint clauses.
- We build our model on top of BERT with a global context encoder to capture the textual context and a label-refining mechanism to utilize the label context. Experiment results on the dataset show the effectiveness of our proposed model.

2 DATASET

2.1 Annotation Framework

With the support of an Internet finance service company, we collect a considerable Chinese complaint corpus with 7 categories of fraud elements: *content fabrication* (CF), *identity fabrication* (IF), *remittance excuse* (RE), *contact platform* (CP), *fraud realization* (FR), *user demand* (UD) and *non-fraudulent statement* (NONE).

The annotation process consists of three steps. First, we split each complaint paragraph into clauses according to Chinese punctuation marks (comma, semicolon and space). Second, each clause is assigned to 2 undergraduates to acquire a unique label of fraud element. If any two annotators disagree with each other, the instance will be checked by a third annotator following the majority rule. If any two cannot reach an agreement, the instance will be discarded. Finally, we calculate the Cohen’s kappa coefficient between two annotators to assess annotation quality and get an averaged kappa coefficient of 78%. Overall, we construct a fraud complaint dataset containing 41,103 paragraphs and 197,878 labeled clauses.

2.2 Dataset Analysis

Different from existing complaint-related work [6, 7, 21] that regards every single complaint text as independent, our corpus contains a hierarchical relationship between clauses and paragraphs. Therefore, except for categorical statistics, we further explore the distribution and relation of clauses in specific paragraph.

Categorical Statistics. Table 1 shows the statistics of each category. The proportion of each category is quite uneven which makes the task challenging. *Non-fraudulent statements* account for a large proportion since clauses exclusive of a specific fraud element will be regard as non-fraudulent. The elements of *content fabrication*, *remittance excuse* and *fraud realization* have large amounts which are the critical elements for a fraud case. The number of *identity fabrication* elements is relatively small, but as an important supplement of factuality modification, they are essential for companies and legal institutions to locate credulous people and scam artists. Although the elements of *contact platform* and *user demand* only take up a minor space, they are indicative for the third-party companies to determine whether they have the authority of intervention and stimulate them to provide better customer service.

Table 1: Statistics of the fraud complaint dataset.

Statistics	CF	IF	RE	CP	FR	UD	NONE
# of clauses	39,207	2,739	19,546	7,882	35,289	2,608	90,607
Proportion	19.81%	1.38%	9.88%	3.98%	17.83%	1.32%	45.79%
Avg. length of clauses	12.0	12.0	11.4	10.5	9.9	10.0	8.5
Vocabulary size	12,646	2,337	6,041	3,597	7,382	1,595	17,067
Clause novelty	0.323	0.853	0.309	0.456	0.209	0.612	0.188

Besides, we measure the semantic richness of each clause by calculating their average length. As shown in Table 1, the elements of *fabricating content* and *identity* are longer which confirms that they contain more abundant information to depict the fraud action. *Non-fraudulent statements* are the shortest because they are more casual and untargeted. We also compute vocabulary size for each category and divide the size by the number of clauses to obtain clause novelty. Since *identity fabrication* and *user demand* are barely mentioned, their repetitive use of words is quite small. For frequently emerged elements such as *content fabrication* and *fraud realization*, their novelty difference reflects that fraudsters practice various deception, but most of the cheated are suffering alike.

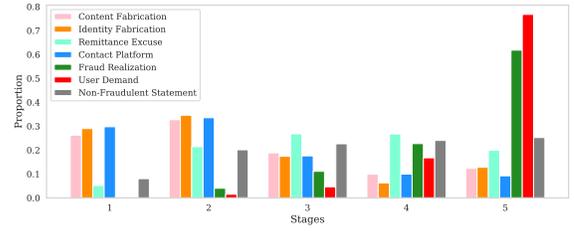


Figure 2: Distribution of fraud elements in different stages of complaint paragraph.

Positional Distribution. Since clauses in the same paragraph are mutually complementary, it is worthwhile to investigate how fraud elements distribute among complaint paragraphs. Given a paragraph, we record the serial number of each clause and divide by the sequence length to obtain the relative position. Further we segment the whole paragraph into 5 equidistant stages. The proportion of different stages for each fraud element is shown in Figure 2. The elements of *content fabrication*, *identity fabrication* and *contact platform* are more likely to appear in the early stage as they come straight to describe the origin, development and transition of the fraud case. The towering bars of *fraud realization* and *user demand* in the rear stages illustrate victims tend to state outcomes and make requests at the end. *Non-fraudulent statements* are quite evenly dispersed in different stages for they have the effect of lubrication to make the whole paragraph clear and coherent.

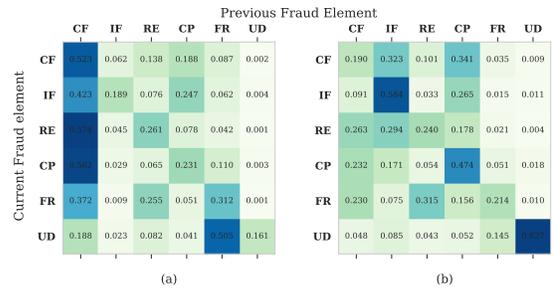


Figure 3: Probabilities of successional fraud element pairs: (a) the original result; (b) the balanced result.

Ordinal Relation. We investigate all the possible combinations of successional element pair [16]. Since non-fraudulent statements lack of explicit semantics, we filter them out. For each of other

categories, we calculate the proportion of the adjacent parent. Similarly, we remove the effect of class imbalance by dividing the prior possibility of the previous fraud element before computing the proportion. The original and balanced results are shown in Figure 3. Before relieving the imbalance effect, it is inevitable that *content fabrication* is more likely to appear before other elements due to its vast existence. Even so, the *user demand* is less likely to follow *content fabrication* which means an absence of consistency between them. However, with a nearly same proportion, the *fraud realization* is less possible to become the prerequisite except for user demand. After normalizing with prior distribution, the elements of *content fabrication*, *identity fabrication* and *contact platform* are the most possible previous statements. Interestingly, both (a) and (b) in Figure 3 have a deep-colored diagonal indicating that clauses frequently succeed the previous of the same type and there exists a strong semantic transitivity between consecutive clauses.

With all the findings, we claim that the position and relation are significant attributes for complaint clauses, and it is reasonable to treat fraud element identification as a sequence labeling task.

3 TASK AND MODEL

This section presents the task of fraud element identification in complaint text and our classification model. The overall framework is shown in Figure 4, which is based on a pre-trained BERT model [4]. It consists of a Local Clause Encoder, a Global Context Encoder and a Label Refiner. The Global Context Encoder is designed to capture the textual context across the complaint text and Label Refiner is proposed to utilize label context.

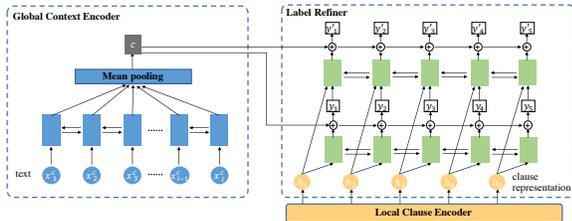


Figure 4: The overall architecture of our proposed model.

3.1 Task Definition

We first introduce some notations in our task:

- $x = (x_1, \dots, x_n)$: an complaint text with n clauses, where x_i is the i -th clause;
- $x_i = (x_{i,1}, \dots, x_{i,m})$: a clause with m tokens, where $x_{i,j}$ is the j -th token in x_i ;
- y_i : the fraud element of the i -th clause x_i in the complaint text x .

We formulate the task as a sequence-level classification task that assigns a element type y_i for each clause x_i in the complaint text x .

3.2 Our Model

In our model, the Local Clause Encoder first utilizes BERT to encode each clause in a complaint text x and take the representations of [CLS] symbol as the clause representations (h_1, \dots, h_n). Then we introduce the Global Context Encoder and the Label Refiner.

Global Context Encoder. Considering a complaint text contains several fine-grained clauses, we introduce a global context encoder [17] to capture the global textual context information c to help identify the element of clauses in the complaint text x . We concatenate all clauses in the complaint text $x^c = \text{concat}([x_1, \dots, x_n])$ and take the embeddings of them as the input of the global encoder, and use a bi-directional GRU [9] to encode the whole text sequence. Then a mean pooling is conducted over encoding states of all time steps to get the textual context representation as Eq.(1), where E_w is the word embedding matrix.

$$\begin{aligned} h_i^c &= \text{BiGRU}(E_w x_i^c, h_{i-1}^c) \\ c &= \text{mean}(h_1^c, \dots, h_n^c, \dots) \end{aligned} \quad (1)$$

Label Refiner. After encoding the clauses and getting the textual context representation of the complaint text, we identify the fraud element of each clause. Instead of directly feeding [CLS] representation from BERT into a output layer, we take the clause representations (h_1, \dots, h_n) in the complaint text x as the input and employ a BiGRU to learn the dependencies between the clauses and get their hidden states as ($s_1^1, s_2^1, \dots, s_n^1$). To be aware of the global textual context information, we concatenate the state s_i^1 of clause x_i with the textual context representation c [20] to get the classification probability through a feed-forward network as Eq.(2).

$$P(y_i|x_i, x) = \text{Softmax}(W_1[s_i^1; c]) \quad (2)$$

Then we propose a *Label-Refining* mechanism to utilize the label context information to refine the identification results. As shown in Figure 4, after the first BiGRU layer for classification, we take the concatenation of each clause representation h_i and the previous classification probability $P(y_i|x_i, x)$ as the input and employ another BiGRU to again encode the clauses sequence in the complaint text x . Then we get the refined probability $P(y'_i|x_i, x)$ as Eq.(3).

$$\begin{aligned} s_i^2 &= \text{BiGRU}([P(y_i|x_i, x); h_i], s_{i-1}^2) \\ P(y'_i|x_i, x) &= \text{Softmax}(W_2[s_i^2; c]) \end{aligned} \quad (3)$$

4 EXPERIMENTS

4.1 Experimental Setup

We randomly split the dataset into training, validation and test set by the proportion of 8:1:1. The numbers of complaint paragraphs and clauses in each data split are (32882, 4110, 4111) and (158521, 19731, 19626), respectively.

We adopt BERT-base, Chinese model as backbone. The global encoder is a two-layer GRU while the GRUs in label refiner are both of one-layer, and the size of hidden units in all GRUs is set as 256. Mini-batch of size 32 is taken and the dropout rate is 0.3. We first fine-tune BERT for the classification task as the base model for 4 epochs using AdamW optimizer with learning rate of $2e-5$. Then we load the fine-tuned BERT as the local clause encoder and freeze its parameters, and then train our model for another 10 epochs using Adam optimizer with learning rate $2e-4$.

4.2 Overall Performance

We compare *Our Model* with some baseline and state-of-the-art models, including *SVM* [19], *BiGRU* [13], *BERT* [4], *BERT-wwm-ext* [3], *BERT+BiGRU* and *RoBERTa* [14]. *Our Model-GC* and *Our*

Model-LR are ablation tests of our model, without the global context encoder and the *label-refining* mechanism, respectively.

Table 2: Evaluation results for element identification in complaint text of different models.

Model	Accuracy(%)	Precision(%)	Recall(%)	F1-score(%)
<i>SVM</i>	77.22	79.38	71.41	75.19
<i>BiGRU</i>	80.51	79.46	77.31	78.37
<i>BERT</i>	83.53	83.66	80.64	82.13
<i>BERT-www-ext</i>	83.27	82.80	81.63	82.21
<i>BERT+BiGRU</i>	82.70	82.38	80.44	81.40
<i>RoBERTa</i>	83.40	82.39	81.97	82.18
<i>Our Model-LR</i>	84.11	83.28	82.31	82.79
<i>Our Model-GC</i>	84.30	83.73	82.35	83.04
<i>Our Model</i>	84.47	84.11	82.30	83.19

We take Accuracy, Precision, Recall, and F1-score as evaluation metrics, and results are shown in the Table 2. We have several findings as follows: (1) *Our Model* outperforms others in terms of all metrics except recall. This indicates the effectiveness of our model considering both textual context and label context information. (2) *BiGRU* performs worse than all other BERT-based models, which demonstrate the effectiveness of pre-trained models. (3) *BERT*, *BERT-www-ext* and *RoBERTa* achieve similar performance, which means that *BERT* is enough to capture clause representation for classification. Thus we adopt *BERT* as our base model. (4) The improvements of *Our Model* compared to *Our Model-LR* and *Our Model-GC* respectively reveal the effectiveness of the global context encoder and *label-refining* mechanism for our element identification task.

4.3 Further Analysis

Results of different fraud element. We investigate the performance of different fraud element types of our model in terms of Precision, Recall and F1-score in Table 3. We can see that for clauses of Fraud Realization, User Demand and Non-Fraudulent Statement, our model achieves a good performance. From the Figure 2 which shows clauses of Fraud Realization and User Demand concentrated in the last stage of the complaint text, our model is easier to capture this pattern so that it performs better for these two elements. As for Non-Fraudulent Statement, we assume the good performance is owing to its large proportion in the dataset.

Table 3: Evaluation results for complaint text classification of different fraud element types.

Fraud Element	Precision(%)	Recall(%)	F1-score(%)
CF	82.68	82.33	82.50
IF	83.90	82.49	83.19
RE	75.85	63.74	69.27
CP	85.90	83.23	84.54
FR	86.07	89.16	87.59
UD	88.43	86.99	87.70
NONE	85.92	88.15	87.02

Error Analysis. We compute the confusion matrix to analyse which pairs of fraud elements are more confused to be classified. We also compare the confusion matrices of *Our Model-LR* and *Our Model* to see which elements are improved by the label refining

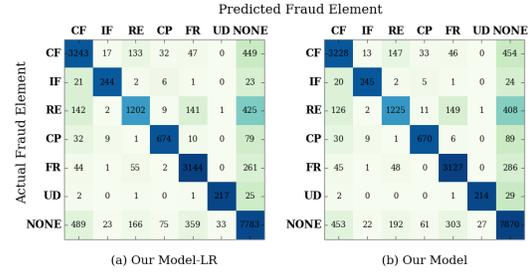


Figure 5: The confusion matrices for classification results of *Our Model-LR* (a) and *Our Model* (b).

mechanism. As shown in Figure 5, we can see that clauses of Non-Fraudulent Statement are easier to be confused with all other elements, because it makes up the majority of the imbalanced dataset. And the pairs of (Identity Fabrication, Content Fabrication), (Remittance Excuse, Content Fabrication) and (Remittance Excuse, Fraud Realization) are likely to be classified incorrectly. From comparison between the diagonal elements of matrix (a) and (b), we find that the classification of clauses in Remittance Excuse, Non-Fraudulent Statement and Identity Fabrication are indeed improved, and the corresponding confused numbers are decreased, for example, (Remittance Excuse, Content Fabrication) is reduced from 142 to 126.

5 RELATED WORK

Recently many complaint classification tasks have been widely studied [12]. Researchers categorize complaints in order to figure out the complaint reasons [21], identify downstream companies [7], and explore user demands [10]. Filgueiras et al. [6] analyze the complaints about food safety and economic surveillance to help government structuralize complaint letters and decide the fine-grained department ultimately responsible for the complaint. In our work, rather than assigning labels to the whole complaint text, we map each component of the paragraph and list all the fraud elements to provide more interpretable judicial outcomes.

Traditional methods for text classification utilize manually crafted features, and employ machine learning algorithms such as Naive Bayes [15] and Support Vector Machines [2] to obtain category bounds. With the popularity of deep learning, RNN, CNN [11], and pre-trained language models [4, 8, 14, 18] refresh the performance of text classification. In this paper, we build our model on top of BERT and propose a fine-grained classification schema for element identification using textual and label context information.

6 CONCLUSION

In this paper, we aim at analyzing complaint text of internet fraud in a fine-grained level. We first propose an annotation scheme to distinguish various types of fraud elements and construct a dataset as benchmark. We build a classification model on top of BERT and use context information of complaint text to better identify element types. Experiment results confirm the effectiveness of our model. In the future, we will explore how to utilize domain knowledge for complaint text modeling. We will also utilize the fine-grained element type information for better complaint text classification.

REFERENCES

- [1] Carol Brennan, Tania Sourdin, Jane Williams, Naomi Burstyner, and Chris Gill. 2017. Consumer vulnerability and complaint handling: Challenges, opportunities and dispute system design. *International journal of consumer studies* 41, 6 (2017), 638–646.
- [2] Corinna Cortes and Vladimir Vapnik. 1995. Support-vector networks. *Machine learning* 20, 3 (1995), 273–297.
- [3] Yiming Cui, Wanxiang Che, Ting Liu, Bing Qin, Ziqing Yang, Shijin Wang, and Guoping Hu. 2019. Pre-training with whole word masking for chinese bert. *arXiv preprint arXiv:1906.08101* (2019).
- [4] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [5] Eric Ferri. 2018. The Evolving Practice Of Complaint Management. *Bloomberg Law* (2018), 1–8.
- [6] Joao Filgueiras, Luis Barbosa, Gil Rocha, Henrique Lopes Cardoso, Luis Paulo Reis, Joao Pedro Machado, and Ana Maria Oliveira. 2019. Complaint Analysis and Classification for Economic and Food Safety. In *Proceedings of the Second Workshop on Economics and Natural Language Processing*. 51–60.
- [7] Yaakov HaCohen-Kerner, Rakefet Dilmon, Maor Hone, and Matanya Aharon Ben-Basan. 2019. Automatic classification of complaint letters according to service provider categories. *Information Processing & Management* 56, 6 (2019), 102102.
- [8] Jeremy Howard and Sebastian Ruder. 2018. Universal language model fine-tuning for text classification. *arXiv preprint arXiv:1801.06146* (2018).
- [9] Zhiheng Huang, Wei Xu, and Kai Yu. 2015. Bidirectional LSTM-CRF models for sequence tagging. *arXiv preprint arXiv:1508.01991* (2015).
- [10] Junegak Joung, Kiwook Jung, Sanghyun Ko, and Kwangsoo Kim. 2019. Customer complaints analysis using text mining and outcome-driven innovation method for market-oriented product development. *Sustainability* 11, 1 (2019), 40.
- [11] Yoon Kim. 2014. Convolutional neural networks for sentence classification. *arXiv preprint arXiv:1408.5882* (2014).
- [12] William Kos, MP Schraagen, MJS Brinkhuis, and FJ Bex. 2017. Classification in a Skewed Online Trade Fraud Complaint Corpus. In *Preproceedings of the 29th Benelux Conference on Artificial Intelligence November 8–9, 2017 in Groningen, The Netherlands*. 172–183.
- [13] Gang Liu and Jiabao Guo. 2019. Bidirectional LSTM with attention mechanism and convolutional layer for text classification. *Neurocomputing* 337 (2019), 325–338.
- [14] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).
- [15] Andrew McCallum, Kamal Nigam, et al. 1998. A comparison of event models for naive bayes text classification. In *AAAI-98 workshop on learning for text categorization*, Vol. 752. Citeseer, 41–48.
- [16] Dunja Mladenic and Marko Grobelnik. 1998. Word sequences as features in text-learning. In *In Proceedings of the 17th Electrotechnical and Computer Science Conference (ERK98)*. Citeseer.
- [17] Guocheng Niu, Hengru Xu, Bolei He, Xinyan Xiao, Hua Wu, and GAO Sheng. 2019. Enhancing Local Feature Extraction with Global Representation for Neural Text Classification. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. 496–506.
- [18] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training.
- [19] Johan AK Suykens and Joos Vandewalle. 1999. Least squares support vector machine classifiers. *Neural processing letters* 9, 3 (1999), 293–300.
- [20] Koji Tanaka, Junya Takayama, and Yuki Arase. 2019. Dialogue-Act Prediction of Future Responses Based on Conversation History. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: Student Research Workshop*. 197–202.
- [21] Xuesong Tong, Bin Wu, Shuyang Wang, and Jinna Lv. 2018. A complaint text classification model based on character-level convolutional network. In *2018 IEEE 9th International Conference on Software Engineering and Service Science (ICSESS)*. IEEE, 507–511.