

# CED: Catalog Extraction from Documents

Tong Zhu<sup>1</sup>, Guoliang Zhang<sup>1</sup>, Zechang Li<sup>2</sup>, Zijian Yu<sup>1</sup>  
Junfei Ren<sup>1</sup>, Mengsong Wu<sup>1</sup>, Zhefeng Wang<sup>2</sup>  
Baoping Huai<sup>2</sup>, Pingfu Chao<sup>1</sup>, and Wenliang Chen<sup>1\*</sup>

<sup>1</sup> Institute of Artificial Intelligence, School of Computer Science and Technology,  
Soochow University, China

<sup>2</sup> Huawei Cloud, China

{tzhu7, glzhang, zjyu, jfrenjffren, mswumsw}@stu.suda.edu.cn

{lizechang1, wangzhefeng, huaibaoping}@huawei.com

{pfchao, wlchen}@suda.edu.cn

**Abstract.** Sentence-by-sentence information extraction from long documents is an exhausting and error-prone task. As the indicator of document skeleton, catalogs naturally chunk documents into segments and provide informative cascade semantics, which can help to reduce the search space. Despite their usefulness, catalogs are hard to be extracted without the assist from external knowledge. For documents that adhere to a specific template, regular expressions are practical to extract catalogs. However, handcrafted heuristics are not applicable when processing documents from different sources with diverse formats. To address this problem, we build a large manually annotated corpus, which is the first dataset for the Catalog Extraction from Documents (CED) task. Based on this corpus, we propose a transition-based framework for parsing documents into catalog trees. The experimental results demonstrate that our proposed method outperforms baseline systems and shows a good ability to transfer. We believe the CED task could fill the gap between raw text segments and information extraction tasks on extremely long documents. Data and code are available at <https://github.com/Spico197/CatalogExtraction>

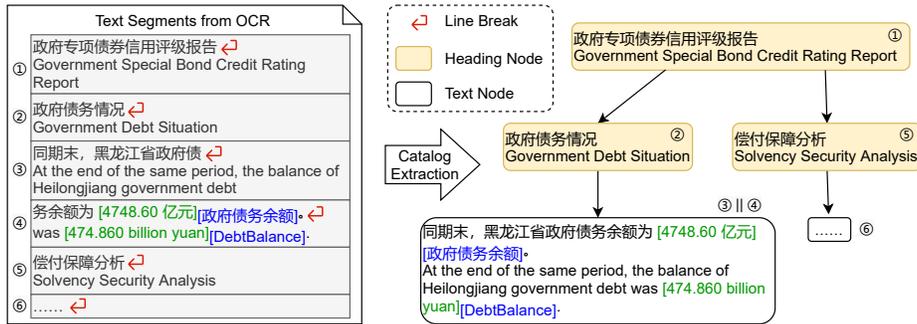
**Keywords:** Catalog Extraction · Information Extraction · Intelligent Document Processing.

## 1 Introduction

Information in long documents is usually sparsely distributed [13,21], so a preprocessing step that distills the structure is necessary to help reduce the search space for subsequent processes. Catalogs, as the skeleton of documents, can naturally locate coarse information by searching the leading section titles. As exemplified in Figure 1, the debt balance “474.860 billion yuan” appears in only one segment in the credit rating report that is 30 to 40 pages long. Taking the whole document

---

\* Corresponding author



**Fig. 1.** An example of catalog extraction. The text segments on the left are converted to a catalog tree on the right. The third and the fourth segments are concatenated after catalog extraction.

into Information Extraction (IE) systems is not practical in this condition. By searching the catalog tree, this entity can be located in the “Government Debt Situation” section with prior knowledge. Unfortunately, most documents are in plain text and do not contain catalogs in an easily accessible format. Thus, we propose the Catalog Extraction from Documents (CED) task as a preliminary step to any extremely long document-level IE tasks. In this manner, fine-grained entities, relations, and events can be further extracted within paragraphs instead of the entire document, which is pragmatic in document-level entity relationship extraction [12,15,14] and document-level event extraction [1].

Designing handcrafted heuristics may be a partial solution to the automatic catalog extraction problem. However, the performance is limited due to three major challenges: 1) Section titles vary across documents, and there are almost no common rules. For documents that are in the same format or inherited from the same template, the patterns of section titles are relatively fixed. Therefore, it is common to use regular expression matching to obtain the whole catalog. However, such handcrafted heuristics are not reusable when the formats of documents change, and researchers have to design new patterns from scratch, making catalog extraction laborious. 2) Catalogs have deep hierarchies with five- to six-level section headings. As the level of section headings deepens, titles become increasingly complex, and simple rule systems usually cannot handle fine-grained deep section headings well. 3) A complete sentence may be cut into multiple segments due to mistakes in data acquisition tools. For example, Optical Character Recognition (OCR) systems are commonly used for obtaining document texts. However, these systems often make mistakes, and sentences may be incorrectly cut into several segments by line breaks. These challenges increase the difficulties of using handcrafted rules.

To address the CED task, we first construct a corpus with a total of 650 manually annotated documents. The corpus includes bid announcements, financial announcements, and credit rating reports. These three types of documents vary in length and catalog complexity. This corpus is able to serve as a benchmark

for the evaluation of CED systems. Among these three sources, bid announcements are the shortest in length with simple catalog structures, and financial announcements contain multifarious heading formats, while credit rating reports have deep and nested catalog structures. In addition, we collect documents from Wikipedia with catalog structures as a large-scale corpus for general model pre-training to enhance the transfer learning ability. These four types of data cover the first two challenges in catalog extraction. We also chunk sentences to simulate the incorrect segmentation problem observed in OCR systems, which covers the third challenge in CED.

Based on the constructed dataset, we design a transition-based framework for the CED task. The catalog tree is formulated as a stack and texts are encased in an input queue. These two buffers are used to help make action predictions, where each action stands for a control signal that manipulates the composition of a catalog tree. By constantly comparing the top element of the catalog stack with one text piece from the input queue, the catalog tree is constructed while action predictions are obtained. The final experimental results show that our method achieves promising results and outperforms other baseline systems. Besides, the model pre-trained on Wikipedia data is able to transfer the learned information to other domains when training data are limited.

Our contributions are summarized as follows:

- We propose a new task to extract catalogs from long documents.
- We build a manually annotated corpus for the CED task, together with a large-scale Wikipedia corpus with catalog structures for pre-training. The experimental results show the efficacy of low-resource transfer.
- We design a transition-based framework for the task. To the best of our knowledge, this is the first system that extracts catalogs from plain text segments without handcrafted patterns.

## 2 Related Work

Since CED is a new task that has not been widely studied, in this section, we mainly introduce approaches applied to similar tasks below.

**Parsing Problems:** Similar to other text-to-structure tasks, CED can be recognized as a parsing problem. A common practice to build syntactic parsers is biaffine-based frameworks with delicate decoding algorithms (e.g., CKY, Eisner, MST) to obtain global optima [4,19]. However, when the problem shifts from sentences to documents, former token-wise encoding and decoding methods become less applicable. As to documents, there are also many popular discourse parsing theories [8,11,6], which aim to extract the inner semantics among Elementary Discourse Units (EDU). However, the number of EDUs in current corpora is small. For instance, in the popular RST-DT corpus, the average number of EDU is only 55.6 per document [17]. When the number of EDUs grows larger, the transition-based method becomes a popular choice [7]. Our proposed CED task is based on naive catalog structures that are similar to syntactic structures, but some traditional parsing mechanisms are not suitable since one document

may contain thousands of segments. To this end, we utilize the transition-based method to deal with the CED task.

**Transition-based Applications:** The transition-based method parses texts to structured trees in a bottom-up style, which is fast and applicable for extremely long documents. Despite the successful applications in syntactic and discourse parsing [20,7,5], transition-based methods are widely used in information extraction tasks with particular actions, such as Chinese word segmentation [18], discontinuous named entity recognition [3] and event extraction [16]. Considering all the characteristics of the CED task, we propose a transition-based method to parse documents into catalog trees.

### 3 Dataset Construction

In this section, we introduce our constructed dataset, the ChCatExt. Specifically, we first elaborate on the pre-processing, annotation and post-processing methods, then we provide detailed data statistics.

#### 3.1 Processing & Annotation

We collect three types of documents to construct the proposed dataset, including bid announcements<sup>3</sup>, financial announcements<sup>4</sup> and credit rating reports<sup>5</sup>. We adopt Acrobat PDF reader to convert PDF files into docx format and use Office Word to make annotations. Annotators are required to: 1) remove running titles, footers (e.g., page numbers), tables and figures; 2) annotate all headings with different outline styles; and 3) merge mis-segmented paragraphs. To reduce the annotation bias, each document is assigned to two annotators, and an expert will check the annotations and make final decisions in case of any disagreement. Due to the length and structure variations, one document may take up to twenty minutes for an annotator to label.

After the annotation process, we use pandoc<sup>6</sup> and additional scripts to parse these files into program-friendly JSON objects. We insert a pseudo root node before each object to ensure that every document object has only one root. In real applications, documents are usually in PDF formats, which are immutable and often image-based. Using OCR tools to extract text contents from those files is a common practice. However, the OCR tools often split a natural sentence apart when a sentence is physically cut by line breaks or page turnings in PDF, as shown in Figure 1. To simulate real-world scenarios, we randomly sample some paragraphs with a probability of 50% and chunk them into segments. For heading strings, we chunk them into segments with lengths of 7 to 20 with jieba<sup>7</sup> assistance. This makes heading segmenting more natural, for example, “招标

<sup>3</sup> <http://ggzy.hebei.gov.cn/hbjyzz>

<sup>4</sup> <http://www.cninfo.com.cn>

<sup>5</sup> <https://www.chinaratings.com.cn> and <https://www.dfratings.com>

<sup>6</sup> <https://pandoc.org>

<sup>7</sup> <https://github.com/fxsjy/jieba>

**Table 1.** Data statistics. BidAnn refers to bid announcements, FinAnn is financial announcements and CreRat is credit rating reports. One node may contain multiple segments in its content, and we list the number of nodes here. Depth represents the depth of the document catalog tree (text nodes are also included). Length is obtained by counting the number of document characters.

Source	#Docs	Avg.Length	Avg.#Nodes			Avg.Depth
			Heading	Text	Total	
BidAnn	100	1,756.76	8.04	30.61	38.65	3.00
FinAnn	300	3,504.22	12.09	52.31	64.40	3.79
CreRat	250	15,003.81	27.70	81.07	108.77	4.59
Total ChCatExt	650	7,658.30	17.47	60.03	77.50	3.98
Wiki	214,989	1,960.41	11.07	19.34	30.41	3.86

公告” will be split into “招标 (zhao biao)” and “公告 (gong gao)” instead of “招 (zhao)” and “标公告 (biao gong gao)”. For other normal texts, we split them into random target lengths between 70 and 100. Since the workflow is rather complicated, we will open-source all the processing scripts to help further development.

In addition to the above manually annotated data, we collect 665,355 documents from Wikipedia<sup>8</sup> for model pre-training. Most of these documents are shallow in catalog structures and short in text lengths. We keep documents with a catalog depth from 2 to 4 to reach higher data complexity, so that these documents are more similar to the manually annotated ones. After that, 214,989 documents are obtained. We chunk these documents in the same way as the manually annotated ones to simulate OCR segmentation.

### 3.2 Data Statistics

Table 1 lists the statistics of the whole dataset. Among the three types, BidAnn has the shortest length and the shallowest structure, and the headings are similar to each other. FinAnn is more complex in structure than BidAnn and contains more nodes. Moreover, there are many forms of headings in FinAnn without obvious patterns, which increases the difficulty of catalog extraction. CreRat is the most sophisticated one among all types of data. Its average length is 8.5 times longer than BidAnn while the average depth is 4.59. However, it contains fewer variations in headings, which may be easier for models to locate. Compared to manually annotated domain-specific data, Wiki is easy to obtain. The structure depth is similar to that of FinAnn while its length is 1.5k shorter. Because of the large size, Wiki is well suited for model pre-training and parameter initializing.

It is worth noting that leaf nodes can be heading or normal texts in catalog trees. Since normal texts cannot lead a section, all texts are leaf nodes in catalogs. However, headings could also be leaf nodes if the leading section has no children. Such a phenomenon appears in approximately 24% of documents. Therefore one

<sup>8</sup> <https://dumps.wikimedia.org/zhwiki/20211220/>

**Table 2.** An example of transition-based catalog tree construction. Elements in **red bold** represent the current stack top  $s$ , and elements in **blue underline** represent the input text  $q$ . \$ means the terminal of  $\mathcal{Q}$  and the finale of action prediction.

Step	Catalog Tree Stack $\mathcal{S}$	Input Queue $\mathcal{Q}$	Predicted Action
1	<b>Root</b>	<u>Credit Rating Report</u> , ...	Sub-Heading
2	Root [ <b>Credit Rating Report</b> ]	<u>Debt Situation</u> , ...	Sub-Heading
3	Root [ Credit Rating Report [ <b>Debt Situation</b> ] ]	<u>The balance</u> , ...	Sub-Text
4	Root [ Credit Rating Report [ Debt Situation [ <b>The balance</b> ] ] ]	<u>was 474 billion yuan.</u> , ...	Concat
5	Root [ Credit Rating Report [ Debt Situation [ <b>The balance was 474 billion yuan.</b> ] ] ]	<u>Security Analysis</u> , ...	Reduce
6	Root [ Credit Rating Report [ <b>Debt Situation</b> [ The balance was 474 billion yuan. ] ] ]	<u>Security Analysis</u> , ...	Reduce
7	Root [ <b>Credit Rating Report</b> [ Debt Situation [ The balance was 474 billion yuan. ] ] ]	<u>Security Analysis</u> , ...	Sub-Heading
8	Root [ Credit Rating Report [ Debt Situation [ The balance was 474 billion yuan. ] <b>Security Analysis</b> ] ]	<u>Texts</u> , \$	Sub-Text
9	Root [ Credit Rating Report [ Debt Situation [ The balance was 474 billion yuan. ] Security Analysis [ <b>Texts</b> ] ] ]	\$	\$

node cannot be recognized as a text node simply by the number of children, which makes the CED task more complicated.

## 4 Transition-based Catalog Extraction

In this section, we introduce details of our proposed TRAnSition-based Catalog trEe constRuction method *TRACER*. We first describe the transition-based process, and then introduce the model architecture.

### 4.1 Actions & Transition Process

The transition-based method is designed for parsing trees from extremely long texts. Since the average length of our CreRat documents is approximately 15k Chinese characters, popular global optimized tree algorithms are apparently too costly to be utilized here.

Action design plays an important role in our transition-based method. There are two buffers here: 1) the input queue  $\mathcal{Q}$  providing one text segment  $q$  at each time; and 2) the tree buffer  $\mathcal{S}$  that records the final catalog tree, where the current stack top points to  $s$ . Actions are obtained by comparing  $s$  and  $q$

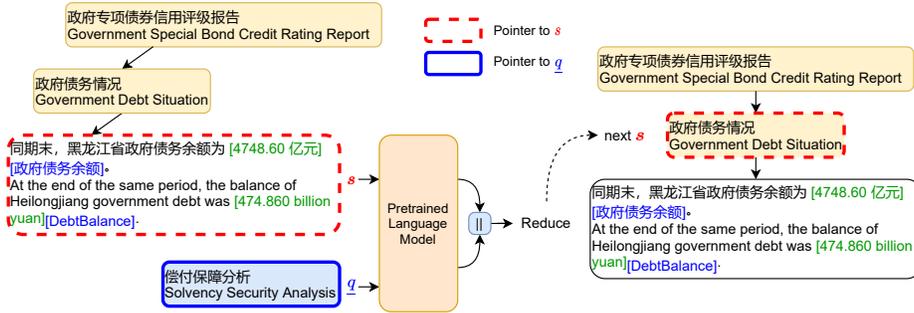


Fig. 2. Framework of the transition-based catalog extraction.

continuously, which results in the buffer changing. As the comparison process continues, actions compose a control sequence to build the target catalog tree simultaneously.

To solve the mentioned challenges, actions are designed to distinguish between headings and texts. Our actions can also capture the difference between headings from adjacent depth levels. In this way, we construct the catalog tree without regard to its depth and complexity. Additionally, we propose an additional action for text segment concatenation. Based on these facts, we design 4 actions as follows:

- Sub-Heading: current input text  $q$  is a child heading node of  $s$ ;
- Sub-Text: current input text  $q$  is a child text node of  $s$ ;
- Concat: current input text  $q$  is the latter part of  $s$  and their contents should be concatenated;
- Reduce: the level of  $q$  is above or at the same level as  $s$ , and  $s$  should be updated to its parent node.

An example is provided in Table 2. To start the prediction, a *Root* node is given in advance. The first heading *Credit Rating Report* is regarded as a child of *Root*. Then, *Debt Situation* becomes another heading node. After that, the **Sub-Text** action suggests that *The balance* is the child node of *Debt Situation* as the body text. Action **Concat** concatenates two body text. Next, action **Reduce** leads to the second layer from the third one. We can eventually build a catalog tree with such a sequence of actions. Furthermore, we present two constraints to avoid illegal results. The first one is that the action between *Root* node and the first input  $q$  can only be Sub-Heading or Sub-Text; Another constraint restricts text nodes to be leaf nodes in the tree, and only Reduce and Concat actions are allowed when  $s$  is not a heading. If the predicted action is illegal, we take the second-best prediction as the final result.

## 4.2 Model Architecture

As Figure 2 shows, the given inputs  $s$  and  $q$  are encoded via a pre-trained language model (PLM). Here, we use a light version of Chinese whole word masking

RoBERTa (RBT3) [2] to obtain encoded representations  $\mathbf{s}$  and  $\mathbf{q}$ . After concatenation,  $\mathbf{g} = \mathbf{s}||\mathbf{q}$  is fed into Feed-Forward Networks (FFN). The FFN is composed of two linear transform layers with ReLU activation function and dropout. Finally we adopt the softmax function to obtain the predicted probabilities as shown in Equation 1.

$$\begin{aligned} \mathbf{o} &= \text{FFN}(\mathbf{g}), \\ p(\mathcal{A}|s, q) &= \text{softmax}(\mathbf{o}), \end{aligned} \tag{1}$$

where  $\mathcal{A}$  denotes all the action candidates. In this way, we can capture the implicit semantic relationship between two nodes.

During prediction, we take the action with maximal probability  $p$  as the predicted result:

$$a_i = \underset{a \in \mathcal{A}}{\text{argmax}} p(\mathcal{A}|s, q),$$

where  $a_i \in \mathcal{A}$  is the predicted action. As discussed in § 4.1, we use two extra constraints to help force decoding legal action results. If  $a_i$  is an illegal action, we sort the predicted probabilities in reverse order, and then find the legal result with the highest probability.

As for training, we take cross entropy as the loss function to help update the model parameters:

$$\mathcal{L} = - \sum_i \mathbb{I}_{y_a=a_i} \log p(a_i|s, q),$$

where  $\mathbb{I}$  is the indicator function,  $y_a$  is the gold action, and  $a_i \in \mathcal{A}$  is the predicted action.

## 5 Experiments

### 5.1 Datasets

We further split the datasets into train, development, and test sets with a proportion of 8:1:1 for training. To fully utilize the scale advantage of the Wiki corpus, we use it to train the model for 40k steps and subsequently save the PLM parameters for transferring experiments.

### 5.2 Evaluation Metrics

We use the overall micro F1 score on predicted tree nodes to evaluate performances. Each node in a tree can be formulated as a tuple: (level, type, content), where level refers to the depth of the current node, type refers to the node type (either *Heading* or *Text*), and content refers to the string that the node carries. The F1 score can be obtained by comparing gold tuples and predicted tuples.

$$P = \frac{N_r}{N_p}, \quad R = \frac{N_r}{N_g}, \quad F1 = \frac{2PR}{P + R},$$

where  $N_r$  denotes the number of correctly matched tuples,  $N_g$  represents the number of gold tuples and  $N_p$  denotes the number of predicted tuples.

**Table 3.** Main results on ChCatExt. The scores are calculated via the method described in § 5.2. Please beware the overall scores are NOT the average of Heading and Text scores. Heading and Text scores are obtained from a subset of predicted tuple results, where all the node types are “Heading” or “Text”. The overall scores are calculated from the universal set, so they are often lower than the Heading and Text scores. WikiBert represents the PLM that is trained on the wiki corpus in advance.

Methods	Heading			Text			Overall		
	P	R	F1	P	R	F1	P	R	F1
Pipeline	88.637	86.595	87.601	81.627	82.475	82.047	76.837	77.338	77.085
Tagging	87.456	88.241	87.846	81.079	81.611	81.344	77.746	78.800	78.269
TRACER	<b>90.634</b>	<b>90.341</b>	<b>90.486</b>	83.031	<b>85.673</b>	<b>84.328</b>	<b>81.017</b>	<b>83.818</b>	<b>82.390</b>
w/o Constraints	89.911	89.713	89.811	82.491	84.948	83.698	80.216	83.035	81.596
TRACER w/ WikiBert	88.671	89.785	89.221	<b>83.308</b>	85.025	84.156	80.820	83.357	82.063

### 5.3 Baselines

Few studies focus on the catalog extraction task, thus we propose two baselines for objective comparisons.

**1) Classification Pipeline:** The catalog extraction task can be formulated in two steps: segment concatenation and tree prediction. For the first step, we take the text pairs as input and adopt the [CLS] token representation to predict the concatenation results. Suppose the depth of a tree is limited, the depth level can be regarded as a classification task with  $\text{MaxHeadingDepth} + 1$  labels, where “1” stands for the text node label. We use PLM with TextCNN [9] to make level predictions.

**2) Tagging:** Inheriting the idea of two-step classification from above, the whole task can be formulated as a tagging task. The segment concatenation sub-task reflects the BIO tagging scheme, and the level depth and node type are tagging labels. We use PLM with LSTM and CRF to address this tagging task.

### 5.4 Experiment Settings

Experiments are conducted with an NVIDIA TitanXp GPU. We use RBT3<sup>9</sup>, a Chinese RoBERTa variation, as the PLM. We use AdamW [10] to optimize the model with a learning rate of  $2e-5$ . Models are trained for 10 epochs. The training batch size is 20, and the dropout rate is 0.5. We take 5 trials with different random seeds for each experiment and report average results on the test set with the best model evaluated on the development set. For the classification pipeline and the tagging baselines, we set the maximal heading depth to 8.

### 5.5 Main Results

From Table 3, we find that our proposed TRACER outperforms the classification pipeline and tagging baselines by 5.305% and 4.121% overall F1 scores. The

<sup>9</sup> <https://huggingface.co/hfl/rbt3>

**Table 4.** Transferring F1 results: train on the source set, evaluate on the target set. Training on Wiki is the process of obtaining the WikiBert, so results in TRACER w/ WikiBert are absent since the experiments are duplicates.

tgt ↓ src →	TRACER				TRACER w/ WikiBert			
	BidAnn	FinAnn	CreRat	Wiki	BidAnn	FinAnn	CreRat	Wiki
BidAnn	88.076	<b>25.557</b>	8.400	2.703	<b>88.200</b>	25.260	<b>11.741</b>	-
FinAnn	7.391	<b>69.249</b>	15.543	11.388	<b>8.100</b>	68.588	<b>20.174</b>	-
CreRat	2.361	14.420	<b>92.790</b>	14.029	<b>7.000</b>	<b>30.821</b>	92.290	-

**Table 5.** Transfer F1 results: train on k documents from the source set, evaluate on the whole target set.

tgt ↓ src →	TRACER								
	BidAnn			FinAnn			CreRat		
	k=3	5	10	k=3	5	10	k=3	5	10
BidAnn	63.033	10.969	7.242	0.713	20.798	9.164	0.000	11.490	<b>39.264</b>
FinAnn	63.460	<b>17.758</b>	10.613	0.815	28.177	11.755	0.047	11.337	48.543
CreRat	77.259	14.363	14.845	1.725	25.110	<b>12.636</b>	3.255	10.768	70.277
TRACER w/ WikiBert									
BidAnn	<b>66.644</b>	<b>14.578</b>	<b>18.719</b>	<b>2.355</b>	<b>21.482</b>	<b>18.385</b>	<b>1.024</b>	<b>12.781</b>	30.626
FinAnn	<b>67.040</b>	15.509	<b>15.467</b>	<b>3.515</b>	<b>32.568</b>	<b>14.125</b>	<b>1.765</b>	<b>25.285</b>	<b>56.192</b>
CreRat	<b>79.029</b>	<b>16.424</b>	<b>14.936</b>	<b>4.517</b>	<b>27.528</b>	12.398	<b>18.659</b>	<b>19.238</b>	67.775

pipeline method requires two separate steps to reveal catalog trees, which may accumulate errors across modules and lead to an overall performance decline. Although the tagging method is a stronger baseline than the pipeline one, it still cannot match TRACER. The reason may be the granularities that these methods focus on. The pipeline and the tagging methods directly predict the depth level for each node, while TRACER pays attention to the structural relationships between each node pair. Besides, since the two baselines need a set of predefined node depth labels, TRACER is more flexible and can predict deeper and more complex structures.

As discussed in § 4.1, we use two additional constraints to prevent TRACER from generating illegal trees. The significance of these constraints is presented in the last line of Table 3. If we remove them, the overall F1 score drops 0.794%. The decline is expected, but the variation is small, which shows the robustness of the TRACER model design.

Interestingly, the PLM trained on the Wiki corpus does not bring performance improvements as expected. This may be due to the different data distributions between Wikipedia and our manually annotated ChCatExt. The following transferring analysis section § 5.6 contains more results with WikiBert.

## 5.6 Analysis of Transfer Ability

One of our motivations for building a model to solve the CED task is that we want to provide a general model that fits all kinds of documents. Therefore, we

**Table 6.** Transfer F1 results: train on the source set, further train on k target documents, evaluate on the target set.

src ↓ tgt →	TRACER								
	BidAnn			FinAnn			CreRat		
	k=3	5	10	k=3	5	10	k=3	5	10
BidAnn	-	87.995	74.630	<b>26.640</b>	-	29.607	<b>37.991</b>	<b>56.658</b>	-
FinAnn	-	87.991	75.921	24.502	-	<b>35.672</b>	<b>38.560</b>	<b>68.287</b>	-
CreRat	-	<b>88.923</b>	79.061	27.988	-	43.066	52.139	72.954	-
	TRACER w/ WikiBert								
BidAnn	-	<b>91.400</b>	<b>76.626</b>	25.709	-	<b>29.729</b>	29.762	53.406	-
FinAnn	-	<b>93.608</b>	<b>76.106</b>	28.035	-	33.698	36.217	65.825	-
CreRat	-	88.777	<b>81.020</b>		-	<b>45.488</b>	<b>57.580</b>	<b>73.519</b>	-
				<b>32.345</b>					

**Table 7.** Transfer F1 results: concatenate the source set with k target documents, train on the merged set, evaluate on the target set.

src ↓ tgt →	TRACER								
	BidAnn			FinAnn			CreRat		
	k=3	5	10	k=3	5	10	k=3	5	10
BidAnn	-	<b>80.924</b>	73.703	27.237	-	29.528	<b>45.813</b>	<b>56.273</b>	-
FinAnn	-	<b>88.902</b>	76.137	24.800	-	31.989	32.173	22.583	-
CreRat	-	<b>88.310</b>	<b>82.551</b>	<b>31.768</b>	-	<b>45.847</b>	<b>61.107</b>	<b>73.933</b>	-
	TRACER w/ WikiBert								
BidAnn	-	78.647	<b>79.606</b>	<b>28.243</b>	-	<b>33.735</b>	45.226	55.887	-
FinAnn	-	83.227	<b>76.556</b>	<b>30.823</b>	-	<b>34.878</b>	<b>36.559</b>	<b>62.070</b>	-
CreRat	-	83.823	76.442	29.587	-	35.086	59.217	71.713	-

conduct transfer experiments under different settings to find the interplay among different sources with diverse structure complexities. The results are listed in Table 4 to 7.

We first train models on three separate source datasets and make direct predictions on target datasets. From the left part of Table 4, we can obtain a rough intuition of the data distribution. The model trained on BidAnn makes poor predictions on FinAnn & CreRat, and gets only 7.391% and 2.361% F1 scores, which also conforms with former discussions in § 3.2. BidAnn is the easiest one among the three sources of datasets, so the generalization ability is less robust. FinAnn is shallower in structure, but it contains more variations. The model trained on FinAnn only obtains a 69.249% F1 score evaluated on FinAnn itself. However, it gets better results on BidAnn (25.557%) and CreRat (14.420%) than the others. The model trained on CreRat gets 92.790% on itself. However, it does not generalize well on the other two sources. We also provide the zero-shot cross-domain results from Wiki to the other three subsets. Although the results are poor under the zero-shot setting, the pre-trained WikiBert shows great transfer ability. Comparing results horizontally in Table 4, we find that the pre-trained WikiBert could provide good generalization and outperforms

the vanilla TRACER among 6 out of 9 transferring data pairs. The other 3 pairs’ results are very close and competitive.

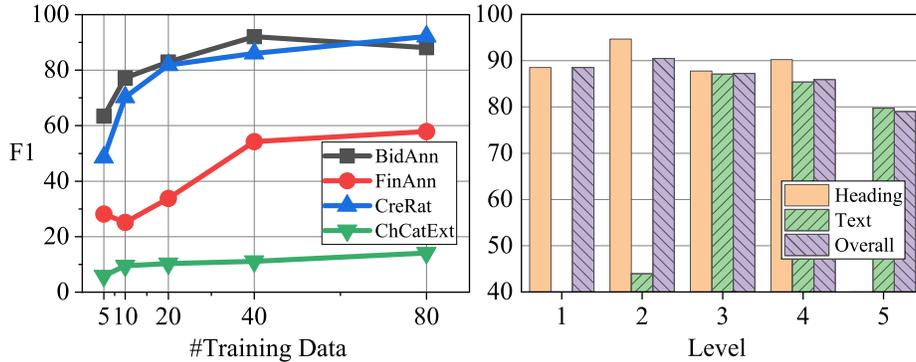
To further investigate the generalization ability of pre-training on the Wiki corpus, we take an extreme condition into consideration, where only a few documents are available to train a model. In this case, as shown in Table 5, we train models with only  $k$  source documents and calculate the final evaluation results on the whole target test set. Each model is evaluated on the original source development set to select the best model and then the best model makes final predictions on the target test set. TRACER w/ WikiBert outperforms vanilla TRACER among 23 out of 27 transferring pairs. There is no obvious upward trend when increasing  $k$  from 3 to 10, which is unexpected and suggests that the model may suffer from overfitting problems on such extremely small training sets.

In most cases of real-world applications, a few target documents are available. Supposing we want to transfer models from source sets to target sets with  $k$  target documents available, there are two possible methods to utilize such data. The first one is to train on the source set, and then further train with  $k$  target documents; the other one is to concatenate the source set and  $k$  targets into a new train set. We conduct experiments under these two settings. The results are presented in Table 6 and 7. Comparing the vanilla TRACER model results, we find that concatenating has 10 out of 18 pairs that outperform the further training method. From  $k=3$  to 10, there are 2, 3, and 5 pairs that show better results, indicating that the concatenation method is better as  $k$  increases. WikiBert has different effects under these two settings. In the further training method, WikiBert is more powerful (11 out of 18 pairs), while it is less useful in the concatenation method (8 out of 18 pairs).

Overall, we find that: 1) WikiBert achieves good performances, especially when the training set is small; 2) If there are  $k$  target documents available besides the source set, WikiBert is not a must, and concatenating the source set with  $k$  targets to make a new train set may produce better results.

## 5.7 Analysis on the Number of Training Data

The left part of Figure 3 shows the average results on each separate dataset with different training data scales. Although BidAnn is the smallest data, the model still gets a 63.460% F1 score and surpasses the other datasets. Interestingly, a decline is observed in BidAnn when the number of training documents increases from 40 to 80. We take it as a normal fluctuation representing a performance saturation since the performance standard deviation is 4.950% when the training data scale is 40. Besides, we find that TRACER has good performance on CreRat. This indicates that TRACER performs well in datasets with deeply nested documents if the catalog heading forms are less varied. In contrast, TRACER is lower in performance on FinAnn than BidAnn and CreRat, and it is more data-hungry than other data sources. For ChCatExt, the merged dataset, performance grows slowly with the increase of training data scale, and more data are needed to be fully trained. Comparing the overall F1 performance of 82.390%



**Fig. 3.** F1 scores with different numbers training data scales (left) and levels (right). The scale means the number of documents participating in training. The results with different levels are evaluated on ChCatExt.

on the whole ChCatExt, the small scale of the training set may lead to a bad generalization.

### 5.8 Analysis on Different Depth

From the right bar plot of Figure 3, it is interesting to see the F1 scores are 0% in level 1 text and level 5 heading. This is mainly due to the golden data distribution characteristics that there are no text nodes in level 1, and there are few headings in deeper levels, leading to zero performances. The F1 score on level 2 text is only 43.938%, which is very low compared to the level 3 text result. Considering that there are only 6.092% of text nodes among all the level 2 nodes, this indicates that TRACER may be not robust enough. Combining the above factors, we find that the overall performance increases from level 1 to 2 and then decreases as the level grows deeper. To reduce the performance decline with deeper levels, additional historical information needs to be considered in future work.

## 6 Conclusion and Future Discussion

In this paper, we build a large dataset for automatic catalog extraction, including three domain-specific subsets with human annotations and large-scale Wikipedia documents with automatically annotated structures. Based on this dataset, we design a transition-based method to help address the task and get promising results. We pre-train our model on the Wikipedia documents and conduct experiments to evaluate the transfer learning ability. We expect that this task and new data could boost the development of Intelligent Document Processing.

We also find some imperfections from the experimental results. Due to the distribution gaps, pre-training on Wikipedia documents does not bring performance improvements on the domain-specific subsets, although it is proven to be

useful under the low-resource transferring settings. Besides, the current model only compares two single nodes each time and misses the global structural histories. Better encoding strategies may need to be discovered to help the model deal with deeper structure predictions. We leave these improvements to future work.

## Acknowledgments

This work was supported by the National Natural Science Foundation of China (Grant No. 61936010) and Provincial Key Laboratory for Computer Information Processing Technology, Soochow University. This work was also supported by the Priority Academic Program Development of Jiangsu Higher Education Institutions, and the joint research project of Huawei Cloud and Soochow University. We would also like to thank the anonymous reviewers for their valuable comments.

## References

1. Chen, Y., Liu, S., Zhang, X., Liu, K., Zhao, J.: Automatically labeled data generation for large scale event extraction. In: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 409–419. Association for Computational Linguistics, Vancouver, Canada (Jul 2017). <https://doi.org/10.18653/v1/P17-1038>, <https://aclanthology.org/P17-1038>
2. Cui, Y., Che, W., Liu, T., Qin, B., Yang, Z.: Pre-training with whole word masking for chinese bert. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* **29**, 3504–3514 (2021). <https://doi.org/10.1109/TASLP.2021.3124365>
3. Dai, X., Karimi, S., Hachey, B., Paris, C.: An effective transition-based model for discontinuous NER. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics. pp. 5860–5870. Association for Computational Linguistics, Online (Jul 2020). <https://doi.org/10.18653/v1/2020.acl-main.520>, <https://aclanthology.org/2020.acl-main.520>
4. Dozat, T., Manning, C.D.: Deep biaffine attention for neural dependency parsing. In: 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24–26, 2017, Conference Track Proceedings. OpenReview.net (2017), <https://openreview.net/forum?id=Hk95PK9le>
5. Dyer, C., Ballesteros, M., Ling, W., Matthews, A., Smith, N.A.: Transition-based dependency parsing with stack long short-term memory. In: Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers). pp. 334–343. Association for Computational Linguistics, Beijing, China (Jul 2015). <https://doi.org/10.3115/v1/P15-1033>, <https://aclanthology.org/P15-1033>
6. Halliday, M.A.K., Matthiessen, C.M., Halliday, M., Matthiessen, C.: An introduction to functional grammar. Routledge (2014)
7. Ji, Y., Eisenstein, J.: Representation learning for text-level discourse parsing. In: Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 13–24. Association for Computational Linguistics, Baltimore, Maryland (Jun 2014). <https://doi.org/10.3115/v1/P14-1002>, <https://aclanthology.org/P14-1002>

8. Kamp, H., Reyle, U.: From discourse to logic: Introduction to modeltheoretic semantics of natural language, formal logic and discourse representation theory, vol. 42. Springer Science & Business Media (2013)
9. Kim, Y.: Convolutional neural networks for sentence classification. In: Moschitti, A., Pang, B., Daelemans, W. (eds.) Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL. pp. 1746–1751. ACL (2014). <https://doi.org/10.3115/v1/d14-1181>, <https://doi.org/10.3115/v1/d14-1181>
10. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. In: 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net (2019), <https://openreview.net/forum?id=Bkg6RiCqY7>
11. MANN, W.C., THOMPSON, S.A.: Rhetorical structure theory: Toward a functional theory of text organization. *Text - Interdisciplinary Journal for the Study of Discourse* **8**(3), 243–281 (1988). <https://doi.org/doi:10.1515/text.1.1988.8.3.243>, <https://doi.org/10.1515/text.1.1988.8.3.243>
12. Peng, N., Poon, H., Quirk, C., Toutanova, K., Yih, W.t.: Cross-sentence n-ary relation extraction with graph LSTMs. *Transactions of the Association for Computational Linguistics* **5**, 101–115 (2017). [https://doi.org/10.1162/tacl\\_a.00049](https://doi.org/10.1162/tacl_a.00049), <https://aclanthology.org/Q17-1008>
13. Yang, H., Chen, Y., Liu, K., Xiao, Y., Zhao, J.: DCFEE: A document-level Chinese financial event extraction system based on automatically labeled training data. In: Proceedings of ACL 2018, System Demonstrations. pp. 50–55. Association for Computational Linguistics, Melbourne, Australia (Jul 2018). <https://doi.org/10.18653/v1/P18-4009>, <https://aclanthology.org/P18-4009>
14. Yao, Y., Ye, D., Li, P., Han, X., Lin, Y., Liu, Z., Liu, Z., Huang, L., Zhou, J., Sun, M.: DocRED: A large-scale document-level relation extraction dataset. In: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics. pp. 764–777. Association for Computational Linguistics, Florence, Italy (Jul 2019). <https://doi.org/10.18653/v1/P19-1074>, <https://aclanthology.org/P19-1074>
15. Yu, D., Sun, K., Cardie, C., Yu, D.: Dialogue-based relation extraction. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics. pp. 4927–4940. Association for Computational Linguistics, Online (Jul 2020). <https://doi.org/10.18653/v1/2020.acl-main.444>, <https://aclanthology.org/2020.acl-main.444>
16. Zhang, J., Qin, Y., Zhang, Y., Liu, M., Ji, D.: Extracting entities and events as a single task using a transition-based neural model. In: Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19. pp. 5422–5428. International Joint Conferences on Artificial Intelligence Organization (7 2019). <https://doi.org/10.24963/ijcai.2019/753>, <https://doi.org/10.24963/ijcai.2019/753>
17. Zhang, L., Xing, Y., Kong, F., Li, P., Zhou, G.: A Top-down Neural Architecture towards Text-level Parsing of Discourse Rhetorical Structure. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics. pp. 6386–6395. Association for Computational Linguistics (2020). <https://doi.org/10.18653/v1/2020.acl-main.569>, <https://aclanthology.org/2020.acl-main.569>

18. Zhang, M., Zhang, Y., Fu, G.: Transition-based neural word segmentation. In: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 421–431 (2016)
19. Zhang, Y., Li, Z., Zhang, M.: Efficient second-order TreeCRF for neural dependency parsing. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics. pp. 3295–3305. Association for Computational Linguistics, Online (Jul 2020). <https://doi.org/10.18653/v1/2020.acl-main.302>, <https://aclanthology.org/2020.acl-main.302>
20. Zhu, M., Zhang, Y., Chen, W., Zhang, M., Zhu, J.: Fast and accurate shift-reduce constituent parsing. In: Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 434–443. Association for Computational Linguistics, Sofia, Bulgaria (Aug 2013), <https://aclanthology.org/P13-1043>
21. Zhu, T., Qu, X., Chen, W., Wang, Z., Huai, B., Yuan, N., Zhang, M.: Efficient document-level event extraction via pseudo-trigger-aware pruned complete graph. In: Raedt, L.D. (ed.) Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22. pp. 4552–4558. International Joint Conferences on Artificial Intelligence Organization (7 2022). <https://doi.org/10.24963/ijcai.2022/632>, <https://doi.org/10.24963/ijcai.2022/632>, main Track