

# TexSmart: A System for Enhanced Natural Language Understanding

Lemao Liu, Haisong Zhang, Haiyun Jiang, Yangming Li, Enbo Zhao, Kun Xu, Linfeng Song, Suncong Zheng, Botong Zhou, Jianchen Zhu, Xiao Feng, Tao Chen, Tao Yang, Dong Yu, Feng Zhang, Zhanhui Kang, Shuming Shi\*

Tencent AI

{texsmart, redmondliu, hansonzhang, shumingshi}@tencent.com

## Abstract

This paper introduces TexSmart, a text understanding system that supports fine-grained named entity recognition (NER) and enhanced semantic analysis functionalities. Compared to most previous publicly available text understanding systems and tools, TexSmart holds some unique features. First, the NER function of TexSmart supports over 1,000 entity types, while most other public tools typically support several to (at most) dozens of entity types. Second, TexSmart introduces new semantic analysis functions like semantic expansion and deep semantic representation, that are absent in most previous systems. Third, a spectrum of algorithms (from very fast algorithms to those that are relatively slow but more accurate) are implemented for one function in TexSmart, to fulfill the requirements of different academic and industrial applications. The adoption of unsupervised or weakly-supervised algorithms is especially emphasized, with the goal of easily updating our models to include fresh data with less human annotation efforts.<sup>1</sup>

## 1 Introduction

The long-term goal of natural language processing (NLP) is to help computers understand natural language as well as we do, which is one of the most fundamental and representative challenges for artificial intelligence. Natural language understanding includes a broad variety of tasks covering lexical analysis, syntactic analysis and semantic analysis. In this paper we introduce TexSmart, a new text understanding system that provides enhanced named entity recognition (NER) and semantic analysis functionalities besides standard NLP modules. Compared to most previous publicly-available text understanding systems (Loper and

Bird, 2002; OpenNLP; Manning et al., 2014; Gardner et al., 2018; Che et al., 2010; Qiu et al., 2013), TexSmart holds the following key characteristics:

- Fine-grained named entity recognition (NER)
- Enhanced semantic analysis
- A spectrum of algorithms implemented for one function, to fulfill the requirements of different academic and industrial applications

First, the fine-grained NER function of TexSmart supports over 1,000 entity types while most previous text understanding systems typically support several to (at most) dozens of coarse entity types (among which the most popular types are people, locations, and organizations). Large-scale fine-grained entity types are expected to provide richer semantic information for downstream NLP applications. Figure 1 shows a comparison between the NER results of a previous system and the fine-grained NER results of TexSmart. It is shown that TexSmart recognizes more entity types (e.g., work.movie) and finer-grained ones (e.g., loc.city vs. the general location type). Examples of entity types (and their important sub-types) which TexSmart is able to recognize include *people, locations, organizations, products, brands, creative work, time, numerical values, living creatures, food, drugs, diseases, academic disciplines, languages, celestial bodies, organs, events, activities, colors*, etc.

Second, TexSmart provides two advanced semantic analysis functionalities: semantic expansion, and deep semantic representation for a few entity types. These two functions are not available in most previous public text understanding systems. Semantic expansion suggests a list of related entities for an entity in the input sentence (as shown in Figure 1). It provides more information about the semantic meaning of an entity. Semantic expansion could also benefit upper-layer applications like web search (e.g., for query suggestion) and recommendation systems. For time and quantity entities, in addition to recognizing them from a sentence, TexSmart also tries to parse them into deep representations (as shown in Figure 1). This kind of deep representations is essential for some NLP applications. For example, when a chatbot is processing query

Project lead and chief architect

<sup>1</sup>TexSmart is available at <https://texsmart.qq.com/en>, and the long version of this paper can be found in the technical report (Zhang et al., 2020).

(a)			(b)			
No.	Entity	Type ID	No.	Entity	Type ID	Semantic Expansion
1	Marvel	person	1	Captain Marvel	work.movie	{ "related": ["Batman", "Superman", "Wonder Woman", "Green Lantern", "the flash", "Aquaman", "Spider-Man", "Green Arrow", "Supergirl", "Captain America"] }
2	Los Angeles	location	2	Los Angeles	loc.city	{ "related": ["Toronto", "Montreal", "Vancouver", "Ottawa", "Calgary", "London", "Paris", "Chicago", "Edmonton", "Boston"] }
3	24 months ago	time	3	24 months ago	time.generic	{ "value": [2019,3] }

Fine-grained NER
Deep Semantic Expression

Figure 1: Comparison between the NER results of a traditional text understanding system in (a) and the fine-grained NER and semantic analysis results provided by TexSmart in (b). The input sentence is “Captain Marvel was premiered in Los Angeles 24 months ago.”. The screenshot was taken in Mar. 2021.

“please book an air ticket to London at 4 pm the day after tomorrow”, it needs to know the exact time represented by “4 pm the day after tomorrow”.

Third, a spectrum of algorithms is implemented for one task (e.g., part-of-speech tagging and NER) in TexSmart, to fulfill the requirements of different academic and industrial applications. On one side of the spectrum are the algorithms that are very fast but not necessarily the best in accuracy. On the opposite side are those that are relatively slow yet delivering state-of-the-art performance in terms of accuracy. Different application scenarios may have different requirements for efficiency and accuracy. Unfortunately, it is often very difficult or even impossible for a single algorithm to achieve the best in both speed and accuracy at the same time. With multiple algorithms implemented for one task, we have more chances to better fulfill the requirements of more applications.

One design principle of TexSmart is to put a lot of efforts into designing and implementing unsupervised or weakly-supervised algorithms for a task, based on large-scale structured, semi-structured, or unstructured data. The goal is to update our models easier to include fresh data with less human annotation efforts.

## 2 System Modules

Compared to most other public text understanding systems, TexSmart supports three unique modules, i.e., fine-grained NER, semantic expansion and deep semantic representation. Besides, traditional tasks supported by both TexSmart and many other systems include word segmentation, part-of-speech (POS) tagging, coarse-grained NER, constituency parsing, semantic role labeling, text classification and text matching. Below we first introduce the unique modules, and then describe the traditional tasks, followed by System Usage.

### 2.1 Key Modules

Since the implementation of fine-grained NER depends on semantic expansion, we first present semantic expansion, then fine-grained NER, and fi-

nally deep semantic representation.

#### 2.1.1 Semantic Expansion

Given an entity within a sentence, the semantic expansion module suggests a list of entities related to the given entity. For example in Figure 1, the suggestion results for “Captain Marvel” include “Spider-Man”, “Captain America”, and other related movies. Semantic expansion attaches additional information to an entity mention, which could be leveraged by upper-layer applications for better understanding the entity and the source sentence. Possible applications of the expansion results include web search (e.g., for query suggestion) and recommendation systems.

Semantic expansion task was firstly introduced in Han et al. (2020), and it was addressed by a neural method. However, this method is not as efficient as one expected for some industrial applications. Therefore, we propose a light-weight alternative approach in TexSmart for this task.

This approach includes two offline steps and two online ones, as illustrated in Figure 2. During the offline procedure, Hearst patterns are first applied to a large-scale text corpus to obtain a is-a map (or called a hyponym-to-hypernym map) (Hearst, 1992; Zhang et al., 2011). Then a clustering algorithm is employed to build a collection of term clusters from all the hyponyms, allowing a hyponym to belong to multiple clusters. Each term cluster is labeled by one or more hypernyms (or called type names). Term similarity scores used in the clustering algorithm are calculated by a combination of word embedding, distributional similarity, and pattern-based methods (Mikolov et al., 2013; Song et al., 2018; Shi et al., 2010).

During the online testing time, clusters containing the target entity mention are first retrieved by referring to the cluster collection. Generally, there may be multiple (ambiguous) clusters containing the target entity mention and thus it is necessary to pick the best cluster through disambiguation. Once the best cluster is chosen, its members (or instances) can be returned as the expansion results.

Now the core challenge is how to calculate the

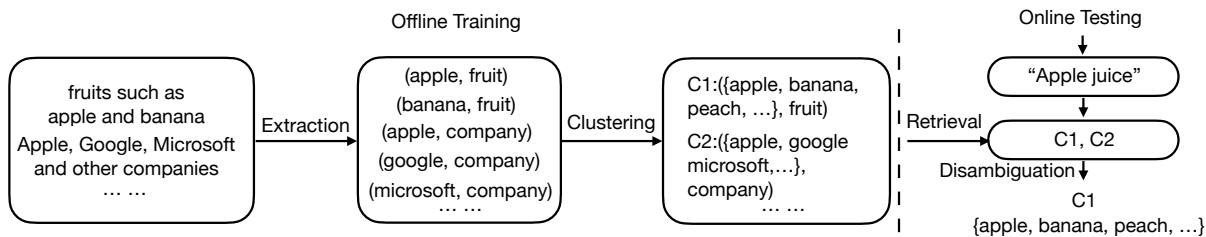


Figure 2: Key steps for semantic expansion: extraction, clustering, retrieval and disambiguation. The first two steps are conducted offline and the last two are performed online.

score of a cluster given an entity mention. We choose to compute the score as the average similarity score between a term in the cluster and a term in the context of the entity mention. Formally, suppose  $e$  is a mention in a sentence, context  $\mathbf{C} = \{c_1, c_2, \dots, c_m\}$  is a window of  $e$  within the sentence, and  $\mathbf{L} = \{e_1, e_2, \dots, e_n\}$  is a term cluster containing the entity mention (i.e.,  $e \in \mathbf{L}$ ). The cluster score is then calculated below:

$$\text{sim}(\mathbf{C}, \mathbf{L}; e) = \frac{1}{(m-1) \times (n-1)} \sum_{x \in \mathbf{C} \setminus \{e\}, y \in \mathbf{L} \setminus \{e\}} \cos(v_x, w_y) \quad (1)$$

where  $\mathbf{C} \setminus \{e\}$  means excluding a subset  $\{e\}$  from a set  $\mathbf{C}$ ,  $v_x$  denotes the input word embedding of  $x$ ,  $w_y$  denotes the output word embedding of  $y$  from a well-trained word embedding model, and  $\cos$  is the cosine similarity function.

### 2.1.2 Fine-Grained NER

Generally, it is challenging to build a fine-grained NER system. Xu et al. (2020) create a fine-grained NER dataset for Chinese, but the number of its types is less than 20. A knowledge base (such as Freebase (Bollacker et al., 2008)) is utilized in Ling and Weld (2012) as distant supervision to obtain a training dataset for fine-grained NER. However, this dataset only includes about one hundred types whereas TexSmart supports up to one thousand types. Moreover, the fine-grained NER module in TexSmart does not rely on any knowledge bases and thus can be readily extended to other languages for which there is no knowledge base available.

**Ontology** To establish fine-grained NER in TexSmart, we need to define an ontology of entity types. The TexSmart ontology was built in a semi-automatic way, based on the term clusters in Figure 2. Please note that each term cluster is labeled by one or more hypernyms as type names of the cluster. We first conduct a simple statistics over the term clusters to get a list of popular type names (i.e., those having a lot of corresponding term clusters). Then we manually create one or more *formal types* from one popular type name and add the type name to the name list of the formal types. For example, formal type “work.movie” is manually built

from type name “movie”, and the word “movie” is added to the name list of “work.movie”. As another example, formal types “language.human\_lang” and “language.programming” are manually built from type name “language”, and the word “language” is added to the name lists of both the two formal types. Each formal type is also assigned with a *sample instance list* in addition to a name list. Instances can be chosen manually from the clusters corresponding to the names of the formal type. To reduce manual efforts, the sample instance list for every type is often quite short. The supertype/subtype relation between the formal types are also specified manually. As a result, we obtain a type hierarchy containing about 1,000 formal types, each assigned with a standard id (e.g., work.movie), a list of names (e.g., “movie” and “film”), and a short list of example instances (e.g., “Star Wars”). The TexSmart ontology is available on the download page<sup>2</sup>. Figure 3 shows a sub-tree (with type id “loc.generic” as the root) sampled from the entire ontology.

**Unsupervised method** The unsupervised fine-grained NER method works in two steps. First, run the semantic expansion algorithm (referring to the previous subsection) to get the best cluster for the entity mention. Second, derive an entity type from the cluster.

For the best cluster obtained in the first step, it contains a list of terms as instances and is also labeled with a list of hypernyms (or type names). The final entity type id for the cluster is determined by a type scoring algorithm. The candidate types are those in the TexSmart ontology whose name lists contain at least one hypernym of the cluster. Please note that each entity type in the TexSmart ontology has been assigned with a name list and a sample instance list. Therefore the score of a candidate entity type can be calculated according to the information of the entity type and cluster.

This unsupervised method has a major drawback: It cannot recognize unknown entity mentions (i.e., entity mentions that are not in any of our term clusters).

<sup>2</sup><https://ai.tencent.com/ailab/nlp/textsmart/en/download.html>

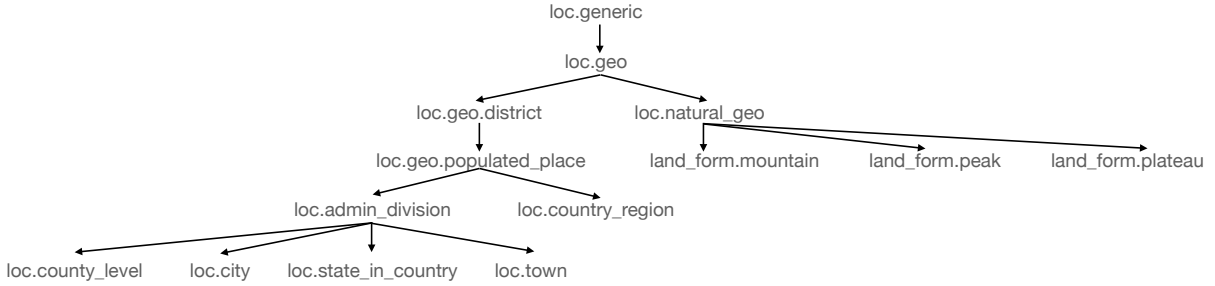


Figure 3: A sub-tree of the TextSmart ontology, with “loc.generic” as the root

**Hybrid method** In order to address the above issue, we propose a hybrid method for fine-grained NER. Its key idea is to combine the results of the unsupervised method and those of a coarse-grained NER model. We train a coarse-grained NER model in a supervised manner using an off-the-shelf training dataset (for example, Ontonotes dataset (Weischedel et al., 2013)). Given the supervised and unsupervised results, the combination policy is as follows: If the fine-grained type is compatible with the coarse type, i.e., the fine-grained one is a subtype of the coarse one, the fine-grained type is returned; otherwise the coarse type is chosen.

For example, assume that the entity mention “apple” in the sentence “...apple juice...” is determined as “food.fruit” by the unsupervised method and “food.generic” by the supervised model. The hybrid approach returns “food.fruit” according to the above policy. However, if the unsupervised method returns “org.company”, the hybrid approach will return “food.generic” because the two types returned by the supervised method and the unsupervised method are not compatible.

Although both unsupervised and hybrid methods are described on top of the ontology manually defined above, they can actually be used for other ontologies such as those in FIGER and Ontonotes datasets, because most type names in these ontologies can be covered by our clusters obtained in semantic expansion as long as the training data is sufficient. In this sense, both methods are general in practice.

### 2.1.3 Deep Semantic Representation

For a time or quantity entity within a sentence, TextSmart can analyze its potential structured representation, so as to further derive its precise semantic meaning. For example in Figure 1, the deep semantic representation given by TextSmart for “24 months ago” is a structured string with a precise date in JSON format: {"value": [2019, 3]} if the screenshot time was Mar. 2021. Deep semantic representation is important for applications like task-oriented chatbots, where the precise meanings of some entities are required. So far, most public text understanding tools do not provide such a fea-

ture. As a result, applications using these tools have to implement deep semantic representation by themselves.

Some NLP toolkits make use of regular expressions or supervised sequence tagging methods to recognize time and quantity entities. However, it is difficult for those methods to derive structured or deep semantic information of entities. To overcome this problem, time and quantity entities are parsed in TextSmart by Context Free Grammar (CFG), which is more expressive than regular expressions. Its key idea is similar to that in Shi et al. (2015) and can be described as follows: First, CFG grammar rules are manually written according to possible natural language expressions of a specific entity type. Second, the Earley algorithm (Earley, 1970) is employed to parse a piece of text to obtain semantic trees of entities. Finally, deep semantic representations of entities are derived from the semantic trees.

## 2.2 Other Modules

**Word Segmentation** In order to support different application scenarios, TextSmart provides word segmentation results of two granularity levels: word level (or basic level), and phrase level. For phrase-level segmentation, some phrases (especially noun phrases) may contained as a unit. An unsupervised algorithm is implemented in TextSmart for both English and Chinese word segmentation. We choose an unsupervised method over supervised ones due to two reasons. First, it is at least 10 times faster. Second, its accuracy is good enough for most applications.

**Part-of-Speech Tagging** Part-of-Speech (POS) denotes the syntactic role of each word in a sentence, also known as word classes or syntactic categories and it is helpful for many downstream text understanding tasks such as parsing (Huang, 2008; Chen and Manning, 2014; Liu et al., 2018a). We implement three models among many popular ones for part-of-speech tagging (Ratnaparkhi, 1996; Huang et al., 2015; Li et al., 2021b): Log-linear based model (Ratnaparkhi, 1996), conditional random field (CRF) based model (Lafferty et al., 2001) and deep neural network (DNN) based model (Akbik

et al., 2018; Liu et al., 2019). We denote them as: `log_linear`, `crf` and `dnn`, respectively.

**Coarse-grained NER** The difference between fine-grained and coarse-grained NERs is that the former involves more entity types with a finer granularity. We implement coarse-grained NER using supervised learning methods, including conditional random field (CRF) (Lafferty et al., 2001) based and deep neural network (DNN) based models (Akbiik et al., 2018; Liu et al., 2019; Li et al., 2020).

**Constituency Parsing** We implement the constituency parsing model based on the work (Kitaev and Klein, 2018). Kitaev and Klein (2018) build the parser by combining a sentence encoder with a chart decoder based on the self-attention mechanism. Different from work (Kitaev and Klein, 2018), we use pre-trained BERT model as the text encoder to extract features to support the subsequent decoder-based parsing. Our model achieves excellent performance and has low search complexity.

**Semantic Role Labeling** Semantic role labeling (also called shallow semantic parsing) tries to assign role labels to words or phrases in a sentence. TexSmart takes a sequence labeling model with BERT as the text encoder for semantic role labeling similar to Shi and Lin (2019). TexSmart supports semantic role labeling on both Chinese and English texts.

**Text Classification** Text Classification aims to assign a semantic label for an input text among a predefined label set. Text Classification is a classical task in NLP and it has been widely used in many applications, such as spam filtering, sentiment analysis and question classification. The predefined label set in TexSmart is available on the web page.<sup>3</sup>

**Text Matching** We implement two text matching algorithms in TexSmart: Linkage and ESIM (Chen et al., 2017). Linkage is an unsupervised algorithm designed by ourselves that incorporates synonymy information and word embedding knowledge to compute semantic similarity. Different from the previous models with complicated network architectures, ESIM carefully designs the sequential model with both local and global inference based on chain LSTMs and outperforms the counterparts.

### 3 System Usage

Two ways are available to use TexSmart: Calling the HTTP API directly, or downloading one version of the offline SDK. Note that for the same input text, the results from the HTTP API and the SDK may be slightly different, because the HTTP API employs a larger knowledge base and supports more

<sup>3</sup>[https://ai.tencent.com/ailab/nlp/textsmart/table\\_html/tc\\_label\\_set.html](https://ai.tencent.com/ailab/nlp/textsmart/table_html/tc_label_set.html).

text understanding tasks and algorithms. The detailed comparison between the SDK and the HTTP API is available in <https://ai.tencent.com/ailab/nlp/textsmart/en/instructions.html>.

**Offline Toolkit (SDK)** So far the SDK supports Linux, Windows, and Windows Subsystem for Linux (WSL). Mac OS support will be added in v0.3.0. Programming languages supported include C, C++, Python (both version 2 and version 3) and Java (version  $\geq 1.6.0$ ). Example codes for using the SDK with different programming languages are in the `./examples` sub-folder. For example, the Python codes in `./examples/python/en_nlu_example1.py` show how to use the TexSmart SDK to process an English sentence. The C++ codes in `./examples/c_cpp/src/nlu_cpp_example1.cc` show how to use the SDK to analyze both an English sentence and a Chinese sentence.

**HTTP API** The HTTP API of TexSmart contains two parts: the text understanding API and the text matching API. The text understanding API can be accessed via HTTP-POST and the URL is available on the web page.<sup>4</sup> The text matching API is used to calculate the similarity between a pair of sentences. Similar to the text understanding API, the text matching API also supports access via HTTP-POST and the URL is available on the web page.<sup>5</sup>

## 4 System Evaluation

### 4.1 Settings

**Semantic Expansion** The performance of semantic expansion are evaluated based on human annotation. We first select at random 5,000 `<sentence, entity mention>` pairs (called SE pairs) from our test set of NER (to make sure that the entities selected are correct). Then our semantic expansion algorithm is applied to the SE pairs to generate a related-entity list for each pair. Top nine expansion results of each SE pair are then judged by human annotators in terms of quality and relatedness, with each result annotated by two annotators. For each result, a label of 2, 1, or 0 is assigned by each annotator. The three labels mean “highly related”, “slightly related”, and “not related” respectively. In calculating evaluation scores, the three labels are normalized to scores 100, 50, and 0 respectively. As there is no context for each expanded entity, it is challenging for human to annotate its ground-truth label. In fact, the overall disagreement rate between two annotators is 23.5%. To measure the quality of our model, we report the average score according to both annotators.

**Fine-grained NER** Ling and Weld (2012) provide a test set for fine-grained NER evaluation.

<sup>4</sup><https://textsmart.qq.com/api>

<sup>5</sup>[https://textsmart.qq.com/api/match\\_text](https://textsmart.qq.com/api/match_text).

	SE		FGNER	
	ZH	EN	Base	Hybrid
Quality	79.5	80.5	45.9	53.8

Table 1: Semantic expansion (SE) and fine-grained NER (FGNER) evaluation results. SE is evaluated by human annotators and FGNER is evaluated by a variant of F1 score. Base denotes the supervised coarse NER model.

However, this dataset only contains about 400 sentences. In addition, it misses some important entities during human annotation, which is a common issue in building a dataset for evaluating fine-grained NER (Li et al., 2021a). Therefore, we create a larger fine-grained NER dataset, based on the Ontonotes 5.0 dataset. We ask three human annotators to label fine-grained types for each coarse-labeled entity. Since human annotators do not need to identify mentions from scratch, it would mitigate the missing entities issue to some extent. Furthermore, because it is too costly for three human annotators to annotate types from the entire ontology, we instead take a sub-ontology for human annotation which combines all types from Ling and Weld (2012) and Gillick et al. (2014), including 140 types in total. Due to ambiguous entities, there are indeed some disagreement annotations among three annotators but their overall agreement rate is respectful, i.e., the averaged pair-wise agreement rate is about 87.1% in terms of Mi-F1 scores.

	Parsing		SRL	
	EN	ZH	EN	ZH
F1	95.42	92.25	86.7	82.1
Sents/sec	16.60	16.00	10.2	11.5

Table 2: Evaluation results for constituency parsing and SRL. The decoding speed in is measured upon a GPU P40 machine.

To set the hybrid method for fine-grained NER, we select LUA (Li et al., 2020) as the coarse-grained NER model, which is trained on Ontonotes 5.0 training dataset (Weischedel et al., 2013). To compare fine-grained NER against coarse-grained NER, we report a variant of F1 measure for evaluation which only differs from standard F1 in matching count accumulation: if an output type is a fine-grained type and it exactly matches a gold fine-grained type, the matching count accumulates 1; if an output is a coarse grained type and it is compatible with a gold fine-grained type, the matching count accumulates 0.5.

**POS Tagging** We evaluate three POS tagging algorithms: log-linear, CRF, and DNN. They are all trained on the standard training datasets from

PTB for English and CTB 9.0 for Chinese. We use their corresponding test sets to evaluate all the models.

**Coarse-grained NER** To ensure better generalization to industrial applications, we combine several public training sets together for English NER. They are CoNLL2003 (Sang and De Meulder, 2003), BTC (Derczynski et al., 2016), GMB (Bos et al., 2017), SEC\_FILING (Alvarado et al., 2015), WikiGold (Balasuriya et al., 2009; Nothman et al., 2013), and WNUT17 (Derczynski et al., 2017). Since the label set for all these datasets are slightly different, we only maintain three common labels (Person, Location and Organization) for training and testing. For Chinese, we create a NER dataset including about 80 thousand sentences labeled with 12 entity types, by following a similar guideline to that of the Ontonotes dataset. We randomly split it into a training set and a test set with ratio of 3:1. We evaluate two algorithms for coarse-grained NER: CRF and DNN. For DNN, we implement the RoBERTa-CRF and Flair models. As we found RoBERTa-CRF performs better on the Chinese dataset while Flair is better on the English dataset, we report results of RoBERTa-CRF for Chinese and Flair for English in our experiments.

**Constituency Parsing** We conduct parsing experiments on both English and Chinese datasets. For English task, we use WSJ sections in Penn Treebank (PTB) (Marcus et al., 1993), and we follow the standard splits: the training data ranges from section 2 to section 21; the development data is section 24; and the test data is section 23. For Chinese task, we use the Penn Chinese Treebank (CTB) of the version 5.1 (Xue et al., 2005). The training data includes the articles 001-270 and articles 440-1151; the development data is the articles 301- 325; and the test data is the articles 271-300.

**SRL** Semantic role labeling experiments are conducted on both English and Chinese datasets. We use the CoNLL 2012 datasets (Pradhan et al., 2013) and follow the standard splits for the training, development and test sets. The network parameters of our model are initialized using RoBERTa. The batch size is set to 32 and the learning rate is  $5 \times 10^{-5}$ .

**Text Matching** Two text matching algorithms are evaluated: ESIM and Linkage. The datasets used in evaluating English text matching are MRPC<sup>6</sup> and QUORA<sup>7</sup>. For Chinese text matching, four datasets are involved: LCQMC (Liu et al., 2018b), AFQMC (Xu et al., 2020), BQ\_CORPUS (Chen et al., 2018), and PAWS-zh (Zhang et al., 2019). We evaluate the quality

<sup>6</sup><https://www.microsoft.com/en-us/download/details.aspx?id=52398>.

<sup>7</sup><https://www.quora.com/q/quoradata/First-Quora-Dataset-Release-Question-Pairs>.

	POS Tagging						Coarse-grained NER			
	Log-linear		CRF		DNN		CRF		DNN	
	EN	ZH	EN	ZH	EN	ZH	EN	ZH	EN	ZH
F1	96.76	93.94	96.50	93.73	97.04	98.08	73.24	67.26	83.12	75.23
Sents/sec	3.9K		1.3K		149		1.1K		107	

Table 3: Evaluation results for some POS Tagging and coarse-grained NER algorithms in TexSmart on both English (EN) and Chinese (ZH) datasets. The English and Chinese NER datasets are labeled with 3 and 12 entity types respectively.

Algorithms	Sents/Sec	English		Chinese			
		MRPC	QUORA	LCQMC	AFQMC	BQ_CORPUS	PAWS-zh
ESIM	861	-	-	82.63	51.30	71.05	61.55
Linkage	1973	82.18	74.94	79.26	48.66	71.23	62.30

Table 4: Text matching evaluation results. ESIM is a supervised algorithm and it is trained on an in-house labeled dataset only for Chinese. Linkage is an unsupervised algorithm and it is trained for both English and Chinese.

and speed for both ESIM and Linkage algorithms in terms of F1 score and sentences per second, respectively. Since we have not trained the English version of ESIM yet, the corresponding evaluation results are not reported.

## 4.2 Evaluation Results

Table 1 shows the evaluation results of semantic expansion and fine-grained NER. For semantic expansion, it is shown that TexSmart achieves an accuracy of about 80.0 on both English and Chinese datasets. It is a pretty good performance. For fine-grained NER, it is observed that the hybrid approach performs much better than the supervised model (LUA).

Evaluation results for constituency parsing and semantic role labeling are summarized in Table 2. For constituency parsing, the F1 scores on the English and Chinese test sets are 95.42 and 92.25, respectively. The decoding speed depends on the input sentence length. It can process 16.6 and 16.0 sentences per second on our test sets. For SRL, the F1 scores on the English and Chinese test sets are 86.7 and 82.1 respectively and it processes about 10 sentences per second. The speed may be not efficient enough for some applications. As future work, we plan to design more efficient syntactic parsing and SRL algorithms.

The evaluation results for POS Tagging and coarse-grained NER are listed in Table 3. The speed values in this table are measured in sentences per second and they are measured upon a machine with `Platinum 8255C CPU @ 2.50GHz`. Please note that the speed results for Log-linear and CRF are obtained using one single thread, while the speed results for DNN are on 6 threads.

It is clear from the POS tagging results that the three algorithms form a spectrum. On one side of

the spectrum is the log-linear algorithm, which is very fast but less accurate than the DNN algorithm. On the opposite side is the DNN algorithm, which achieves the best accuracy but are much slower than the other two algorithms. The CRF algorithm is in the middle of the spectrum.

Also from Table 3, we can see that the two coarse-grained NER algorithms form another spectrum. The CRF algorithm is on the high-speed side, while the DNN algorithm is on the high-accuracy side. Note that for DNN methods in this table, we employ a data augmentation method to improve their generalization abilities and a knowledge distillation method to speed up its inference (Hinton et al., 2015).

Table 4 shows the performance of two algorithms for text matching. We can see from this table that, in terms of speed, both algorithms are fairly efficient. Please note that the speed is measured in sentences per second using one single CPU from a machine with `Platinum 8255C CPU @ 2.50GHz`. In terms of accuracy, their performance comparison depends on the dataset being used. ESIM performs apparently better on the first two datasets, while slightly worse on the last one. Applications may need to test on their datasets before making decision between the two algorithms.

## 5 Conclusion

In this paper we have presented TexSmart, a text understanding system that supports fine-grained NER, enhanced semantic analysis, as well as some common text understanding functionalities. We have introduced the main functions of TexSmart and key algorithms for implementing the functions. We have also reported some evaluation results on major modules of TexSmart.

## References

- Alan Akbik, Duncan Blythe, and Roland Vollgraf. 2018. Contextual string embeddings for sequence labeling. In *COLING 2018, 27th International Conference on Computational Linguistics*, pages 1638–1649.
- Julio Cesar Salinas Alvarado, Karin Verspoor, and Timothy Baldwin. 2015. Domain adaption of named entity recognition to support credit risk assessment. In *Proceedings of the Australasian Language Technology Association Workshop 2015*, pages 84–90.
- Dominic Balasuriya, Nicky Ringland, Joel Nothman, Tara Murphy, and James R Curran. 2009. Named entity recognition in wikipedia. In *Proceedings of the 2009 Workshop on The People’s Web Meets NLP: Collaboratively Constructed Semantic Resources (People’s Web)*, pages 10–18.
- Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. 2008. Freebase: a collaboratively created graph database for structuring human knowledge. In *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*, pages 1247–1250.
- Johan Bos, Valerio Basile, Kilian Evang, Noortje Venhuizen, and Johannes Bjerva. 2017. The groningen meaning bank. In Nancy Ide and James Pustejovsky, editors, *Handbook of Linguistic Annotation*, volume 2, pages 463–496. Springer.
- Wanxiang Che, Zhenghua Li, and Ting Liu. 2010. Ltp: A chinese language technology platform. In *Coling 2010: Demonstrations*, pages 13–16.
- Danqi Chen and Christopher Manning. 2014. A fast and accurate dependency parser using neural networks. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 740–750.
- Jing Chen, Qingcai Chen, Xin Liu, Haijun Yang, Daohe Lu, and Buzhou Tang. 2018. The BQ corpus: A large-scale domain-specific Chinese corpus for sentence semantic equivalence identification. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*.
- Qian Chen, Xiao-Dan Zhu, Z. Ling, Si Wei, Hui Jiang, and Diana Inkpen. 2017. Enhanced lstm for natural language inference. In *ACL*.
- Leon Derczynski, Kalina Bontcheva, and Ian Roberts. 2016. Broad twitter corpus: A diverse named entity recognition resource. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 1169–1179.
- Leon Derczynski, Eric Nichols, Marieke van Erp, and Nut Limsopatham. 2017. Results of the wnut2017 shared task on novel and emerging entity recognition. In *Proceedings of the 3rd Workshop on Noisy User-generated Text*, pages 140–147.
- Jay Earley. 1970. An efficient context-free parsing algorithm. *Communications of the ACM*, 13(2):94–102.
- Matt Gardner, Joel Grus, Mark Neumann, Oyvind Tafjord, Pradeep Dasigi, Nelson F. Liu, Matthew Peters, Michael Schmitz, and Luke Zettlemoyer. 2018. AllenNLP: A deep semantic natural language processing platform. In *Proceedings of Workshop for NLP Open Source Software (NLP-OSS)*, pages 1–6, Melbourne, Australia. Association for Computational Linguistics.
- Dan Gillick, Nevena Lazic, Kuzman Ganchev, Jesse Kirchner, and David Huynh. 2014. Context-dependent fine-grained entity type tagging. *arXiv preprint arXiv:1412.1820*.
- Jialong Han, Aixin Sun, Haisong Zhang, Chenliang Li, and Shuming Shi. 2020. Case: Context-aware semantic expansion. In *AAAI*, pages 7871–7878.
- Marti A Hearst. 1992. Automatic acquisition of hyponyms from large text corpora. In *Coling 1992 volume 2: The 15th international conference on computational linguistics*.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*.
- Liang Huang. 2008. Forest reranking: Discriminative parsing with non-local features. In *Proceedings of ACL-08: HLT*, pages 586–594.
- Zhiheng Huang, Wei Xu, and Kai Yu. 2015. Bidirectional lstm-crf models for sequence tagging. In *Proceedings of ACL*.
- Nikita Kitaev and D. Klein. 2018. Constituency parsing with a self-attentive encoder. In *ACL*.
- John Lafferty, Andrew McCallum, and Fernando CN Pereira. 2001. Conditional random fields: Probabilistic models for segmenting and labeling sequence data.
- Yangming Li, Lemao Liu, and Shuming Shi. 2020. Segmenting natural language sentences via lexical unit analysis. *arXiv preprint arXiv:2012.05418*.
- Yangming Li, Lemao Liu, and Shuming Shi. 2021a. Empirical analysis of unlabeled entity problem in named entity recognition. In *Proceedings of ICLR*.
- Yangming Li, Lemao Liu, and Kaisheng Yao. 2021b. Neural sequence segmentation as determining the leftmost segments. In *Proceedings of NAACL*.
- Xiao Ling and Daniel S Weld. 2012. Fine-grained entity recognition. In *AAAI*, volume 12, pages 94–100.



- Lemao Liu, Muhua Zhu, and Shuming Shi. 2018a. Improving sequence-to-sequence constituency parsing. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32.
- Xin Liu, Qingcai Chen, Chong Deng, Huajun Zeng, Jing Chen, Dongfang Li, and Buzhou Tang. 2018b. LCQMC: a large-scale Chinese question matching corpus. In *Proceedings of the 27th International Conference on Computational Linguistics*.
- 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*.
- Edward Loper and Steven Bird. 2002. Nltk: the natural language toolkit. In *Proceedings of the ACL-02 Workshop on Effective tools and methodologies for teaching natural language processing and computational linguistics*.
- Christopher D Manning, Mihai Surdeanu, John Bauer, Jenny Rose Finkel, Steven Bethard, and David McClosky. 2014. The stanford corenlp natural language processing toolkit. In *Proceedings of 52nd annual meeting of the association for computational linguistics: system demonstrations*, pages 55–60.
- Mitchell Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. 1993. Building a large annotated corpus of english: The penn treebank.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems*, 26:3111–3119.
- Joel Nothman, Nicky Ringland, Will Radford, Tara Murphy, and James R Curran. 2013. Learning multilingual named entity recognition from wikipedia. *Artificial Intelligence*, 194:151–175.
- OpenNLP. <https://opennlp.apache.org>.
- Sameer Pradhan, Alessandro Moschitti, Nianwen Xue, Hwee Tou Ng, Anders Björkelund, Olga Uryupina, Yuchen Zhang, and Zhi Zhong. 2013. Towards robust linguistic analysis using ontonotes. In *Proceedings of the Seventeenth Conference on Computational Natural Language Learning*, pages 143–152.
- Xipeng Qiu, Qi Zhang, and Xuan-Jing Huang. 2013. Fudannlp: A toolkit for chinese natural language processing. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 49–54.
- Adwait Ratnaparkhi. 1996. A maximum entropy model for part-of-speech tagging. In *Conference on empirical methods in natural language processing*.
- Erik F Sang and Fien De Meulder. 2003. Introduction to the conll-2003 shared task: Language-independent named entity recognition. *arXiv preprint cs/0306050*.
- Peng Shi and Jimmy Lin. 2019. Simple bert models for relation extraction and semantic role labeling. *arXiv preprint arXiv:1904.05255*.
- Shuming Shi, Yuehui Wang, Chin-Yew Lin, Xiaojiang Liu, and Yong Rui. 2015. Automatically solving number word problems by semantic parsing and reasoning. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1132–1142.
- Shuming Shi, Huibin Zhang, Xiaojie Yuan, and Ji-Rong Wen. 2010. Corpus-based semantic class mining: distributional vs. pattern-based approaches. In *Proceedings of the 23rd International Conference on Computational Linguistics (Coling 2010)*, pages 993–1001.
- Yan Song, Shuming Shi, Jing Li, and Haisong Zhang. 2018. Directional skip-gram: Explicitly distinguishing left and right context for word embeddings. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 175–180.
- Ralph Weischedel, Martha Palmer, Mitchell Marcus, Eduard Hovy, Sameer Pradhan, Lance Ramshaw, Nianwen Xue, Ann Taylor, Jeff Kaufman, Michelle Franchini, et al. 2013. Ontonotes release 5.0 ldc2013t19. *Linguistic Data Consortium, Philadelphia, PA*, 23.
- Liang Xu, Xuanwei Zhang, Lu Li, Hai Hu, Chenjie Cao, Weitang Liu, Junyi Li, Yudong Li, Kai Sun, Yechen Xu, et al. 2020. Clue: A chinese language understanding evaluation benchmark. *arXiv preprint arXiv:2004.05986*.
- Naiwen Xue, Fei Xia, Fu-Dong Chiou, and Marta Palmer. 2005. The penn chinese treebank: Phrase structure annotation of a large corpus. *Natural language engineering*, 11(2):207.
- Fan Zhang, Shuming Shi, Jing Liu, Shuqi Sun, and Chin-Yew Lin. 2011. Nonlinear evidence fusion and propagation for hyponymy relation mining. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 1159–1168.
- Haisong Zhang, Lemao Liu, Haiyun Jiang, Yangming Li, Enbo Zhao, Kun Xu, Linfeng Song, Suncong Zheng, Botong Zhou, Jianchen Zhu, et al. 2020. Textsmart: A text understanding system for fine-grained ner and enhanced semantic analysis. *arXiv preprint arXiv:2012.15639*.

Yuan Zhang, Jason Baldridge, and Luheng He. 2019.  
PAWS: paraphrase adversaries from word scrambling. *CoRR*, abs/1904.01130.