# Is Word Segmentation Necessary for Deep Learning of ChineseRepresentations

## 1 Introduction

**Firstly, word data sparsity inevitably leads to overfitting and the ubiquity of OOV words limits the model's learning capacity**



Table 1: Word statistics of Chinese Treebank.



Figure 4: Effects of removing training instances containing OOV words.

**Secondly, the state-of-the-art word segmenta- tion performance is far from perfect, the errors of which would bias downstream NLP tasks.**

| Corpora | Yiǒu | Míng | reaches | the final |
|---|---|---|---|---|
| CTB | 拥有 | 明 | 进入 | 总决赛 |
| PKU | 拥 | 明 | 进入 | 总决赛 |

Table 2: CTB and PKU have different segmentation criteria (Chen et al., 2017c).

**Thirdly, if we ask the fundamental problem of how much benefit word segmentation may provide, it is all about how much additional semantic infor- mation is present in a labeled CWS dataset.**

the answer to this question remains unclear

**Before neural network models became popular, there were discussions on whether CWS is nec- essary and how much improvement it can bring about.**

## 2 Related Work

## 3 Experimental Results

### 3.1 Language Modeling

| model | dimension | ppl |
|---|---|---|
| word | 512 | 199.9 |
| char | 512 | 193.0 |
| word | 2048 | 182.1 |
| char | 2048 | 170.9 |
| hybrid (word+char) | 1024+1024 | 175.7 |
| hybrid (word+char) | 2048+1024 | 177.1 |
| hybrid (word+char) | 2048+2048 | 176.2 |
| hybrid (char only) | 2048 | 171.6 |

### 3.2 Machine Translation



Table 5: Results of different models on the Ch-En machine translation task. Results of Mixed RNN are taken from Li et al (2017). BLEU point on CWS task reaches 96.7.

### 3.3 Sentence Matching/Paraphrase



Table 6: Results on the LCQMC and BQ corpora.

### 3.4 Text Classification



Table 7: Results on the validation and the test set for text classification.

### Domain Adaptation Ability



| | train | domain test | |
|---|---|---|---|
| model | acc | | proportion of sen containing OOV |
| word-based | 81.28% | | 11.79% |
| char-based | 83.33% | | 0.56% |
| | train | test domain | |
| model | acc | | proportion of sen containing OOV |
| word-based | 67.32% | | 7.93% |
| char-based | 67.93% | | 46.85% |

Table 8: Domain adaptation of the word-based model and the char-based model.

## 4 Analysis

### 4.1 Data Sparsity



Figure 2: Effects of data sparsity on the char-based model and the word-based model.

### 4.2 Out-of-Vocabulary Words



Figure 4: Effects of removing training instances containing OOV words.

best Vietnam database (CAIS, 2000). Using Jieba, the most widely-used open-sourced Chinese word segmentation system, to seg- ment the CTB, we end up with a dataset consist- ing of 615,194 words with 50,266 distinct words. Among the 50,266 distinct words, 24,458 words appear only once, amounting to 48.7% of the total vocabulary, yet they only take up 4.0% of the entire corpus. If we increase the frequency bar to 4, we get 38,889 words appearing less or equal to 4 times, which constitute 77.4% of the vocabulary but only 10.1% of the corpus.

### 4.3 Overfitting



Figure1: Effects of dropout rates on the char-based model and the word-based model.

### 4.4 Visualization



Figure 3: Sentence matching between two Chinese sentences with char-based model and word-based model.

## 5 Conclusion

**Conclusion**

Through direct comparisons between these two types of models, we find that char- based models consistently outperform word- based models.

**Reasons**

it is because word-based models are more vulnerable to data sparsity and the presence of out- of-vocabulary (OOV) words, and thus more prone to overfitting.

---

1. 目前绝多在用 subword model 而 char-based 和 word-based 可以作一下对比 (这不同代名上的)

2. 分字可以用 transformer

3. 评价的标准还是不够了解。
   例如: auc 是 accuracy 吗?
   ppl 是 perplexity

4. CTB6.0 Chinese TreeBank.
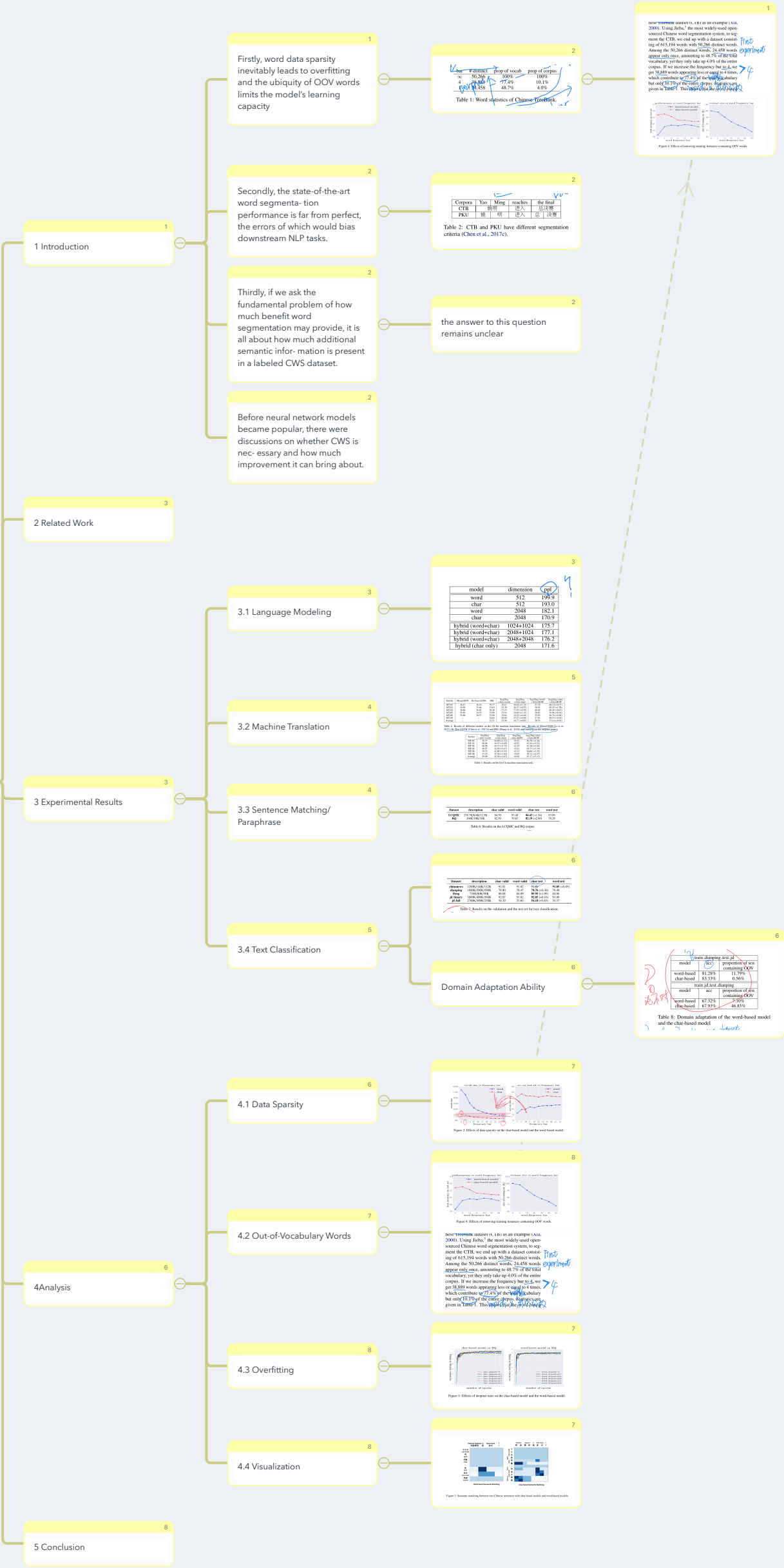   segmented, POS-targeted, bracketed Chinese corpus.
   来自约是 新闻. 文章杂志等

最后提到了 对于 oov 不同的 frequency bar，实验结果 先升后降.
因为 frequency bar 小，那么 infrequent 但 接近 frequency bar
的词 会被归于词库中. 而不利 `为 OV。同时 这种词的进入.
会加重 data sparity 的问题，使很多 有可能有 词库代表性的词，
即（频率低 有特征 ）的词，~~进入词库~~. 影响考虑扑火合效果。
无法进入词库

nese Treebank dataset (CTB) as an example (Xia,
2000). Using Jieba,[3] the most widely-used open-
sourced Chinese word segmentation system, to seg-
ment the CTB, we end up with a dataset consist-
ing of 615,194 words with 50,266 distinct words.
Among the 50,266 distinct words, 24,458 words
appear only once, amounting to 48.7% of the total
vocabulary, yet they only take up 4.0% of the entire
corpus. If we increase the frequency bar to 4, we
get 38,889 words appearing less or equal to 4 times,
which contribute to 77.4% of the total vocabulary
but only 10.1% of the entire corpus. Statistics are
given in Table 1. This shows that the word-based

First
experiment

>4