

fastHan: A BERT-based Joint Many-Task Toolkit for Chinese NLP

Zhichao Geng, Hang Yan, Xipeng Qiu*, Xuanjing Huang

School of Computer Science, Fudan University

Key Laboratory of Intelligent Information Processing, Fudan University

{xpqiu, xjhuang}@fudan.edu.cn

Abstract

We present fastHan, an open-source toolkit for four basic tasks in Chinese natural language processing: Chinese word segmentation, Part-of-Speech tagging, named entity recognition, and dependency parsing. The kernel of fastHan is a joint many-task model based on a pruned BERT, which uses the first 8 layers in BERT. We also provide a 4-layer base version of model compressed from the 8-layer model. The joint-model is trained and evaluated in 13 corpora of four tasks, yielding near state-of-the-art (SOTA) performance in the dependency parsing task and SOTA performance in the other three tasks. In addition to its small size and excellent performance, fastHan is also very user-friendly. Implemented as a python package, fastHan allows users to easily download and use it. Users can get what they want with one line of code, even if they have little knowledge of deep learning. The project is released on Github¹.

1 Introduction

Recently, the need for Chinese natural language processing (NLP) has a dramatic increase for many downstream applications. There exist four basic tasks for Chinese NLP: Chinese word segmentation (CWS), Part-of-Speech (POS) tagging, named entity recognition (NER), and dependency parsing. These basic tasks are usually the cornerstones or provide useful features for other downstream tasks.

Chinese is different from English. There are no obvious boundaries between all characters in a Chinese sentence. CWS is to divide a Chinese character sequence into a sequence of words. POS tagging is to mark the POS of each word, like adjective, verb, noun, and so on. NER recognizes meaningful proprietary names (named entities) in the text, and

mainly studies the recognition of person names, place names, and organization names. Dependency parsing is based on the theory of dependency grammar to construct a dependency grammar tree of sentences. CWS is a character-level task while others are word-level tasks.

There are many pieces of research on how to perform multi-corpus training on these tasks and how to conduct multi-task joint training. There are differences between tasks and corpus. The task is a broader title, like CWS and POS tagging. Each task has several corpora, and each corpus has different criteria. A joint model can solve multiple tasks and multiple criteria at the same time. [Chen et al. \(2017\)](#) explore adversarial multi-criteria learning for CWS, proving more knowledge can be mined through training model in more corpora. [Ng and Low \(2004\)](#) use cross-label to label the POS so that POS tagging and CWS can be trained jointly. Results of the CWS task are contained in the output of the POS tagging task. [Wang et al. \(2013\)](#) firstly solve the CWS task by segmenting a sentence into a word lattice, then using the word lattice to solve the POS tagging task and dependency parsing task. [Yan et al. \(2020\)](#) uses the dependency arc of the “APP” label to jointly train the word-level dependency parsing task and character-level CWS task with the biaffine parser.

FastHan trains the BERT-based ([Devlin et al., 2018](#)) joint-model in 13 corpora to solve the above four tasks. Through multi-task learning, fastHan shares knowledge among the four tasks. There is a strong correlation between these tasks. For example, the model will perform better in the other three word-level tasks if its word segmentation ability is stronger. The joint training of POS tagging and dependency parsing can improve each other’s performance ([Zhang et al., 2020](#)) and so on. This shared information can improve fastHan’s performance on these tasks. Besides, training in more corpora can

^{*}Corresponding author

¹<https://github.com/fastnlp/fastHan>

obtain a larger vocabulary, which can reduce the situation of the model encountering character out of vocabulary. What’s more, the joint-model can greatly reduce the occupied space. Compared with training a model for each task, the joint-model can reduce the occupied space by four times.

FastHan has two versions of the kernel model, base and large. The large version of the model uses the first eight layers of BERT, and the base version of the model uses the Theseus strategy (Xu et al., 2020) to compress the large version of the model to four layers. To improve the performance of the model, fastHan has done a lot of optimization. For example, using the output of POS tagging to improve the performance on the dependency parsing task, using Theseus strategy to improve the performance of the base version model, and so on.

Overall, fastHan has advantages as follows:

Small size: The total parameter of the base model is 151MB, and the large model is 262MB.

High accuracy: The base version of the model achieved good results in all tasks, while the large version of the model approached SOTA in dependency parsing and achieved SOTA performance in other tasks.

Easy to use: Implemented as a python package, users just need one line of code to load the model or process sentences.

For developers of downstream applications, they do not need to do repetitive work for basic tasks and do not need to understand complex codes like BERT. Even if users have little knowledge of deep learning, by using fastHan they can get the results of SOTA performance quickly and easily. In addition, the smaller size can reduce the need for hardware, so that fastHan can be deployed on more platforms.

2 Proposed Model

The kernel of fastHan is a joint-model, which performs multi-task learning in 13 corpora of the four tasks. The architecture of the model is shown in Figure 1. For this model, sentences of different tasks are first added with corpus tags at the beginning of the sentence, and then the feature vectors are sent to the decoding layer through the BERT-based feature extraction layer. The decoding layer will use different decoders according to the current task: use conditional random field (CRF) to decode in the NER task; use MLP and CRF to decode

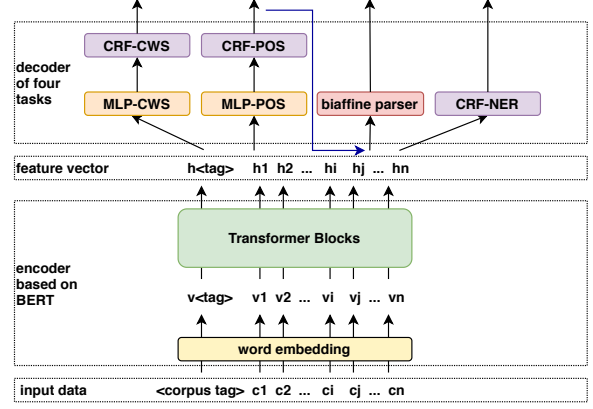


Figure 1: Architecture of the proposed model. Here input data are the input characters that have been vectorized.

in POS tagging and CWS task; use the output of POS tagging task combined with biaffine parser to decode in dependency parsing task.

Each task uses independent label set here, CWS uses label set $Y = \{B, M, E, S\}$; POS tagging uses cross-labels set based on $\{B, M, E, S\}$; NER uses cross-labels set based on $\{B, M, E, S, O\}$; dependency parsing uses arc head and arc label to represent dependency grammar tree.

2.1 BERT-based feature extraction layer

BERT (Devlin et al., 2018) is a language model trained in large-scale corpus. The pre-trained BERT can be used to encode the input sequence. We take the output of the last layer of transformer blocks as the feature vector of the sequence. The attention (Vaswani et al., 2017) mechanism of BERT can extract rich and semantic information related to the context. In addition, the calculation of attention is parallel in the entire sequence, which is faster than the feature extraction layer based on LSTM. In addition, we add layer pruning and corpus tags to BERT.

Layer Pruning: The original BERT has 12 layers of transformer blocks, which will occupy a lot of memory space. Even if the calculations in the transformer are parallel, the time cost of calculating 12 times is very high for these basic tasks. Inspired by Huang et al. (2019), we only use 4 or 8 layers. Our experiment found that using the first eight layers performs well on all tasks, and after compressing, four layers are enough for CWS, POS tagging, and NER.

Corpus Tags: Compared to linear projection layer, we use corpus tags to distinguish various tasks and corpora. Each corpus of each task corresponds to a specific corpus tag, and the embedding of these tags needs to be initialized and trained in fine-tuning. As shown in Figure 1, before inputting the sequence into BERT, we add the corpus tag to the head of the sequence. The attention mechanism will ensure that the vector of the corpus tag and the vector of each other position generate sufficiently complex calculations to bring the corpus and task information to each character.

2.2 CRF Decoder

We use conditional random field (CRF) (Lafferty et al., 2001) to do the final decoding work in POS tagging, CWS, and NER tasks. In CRF, the conditional probability of a label sequence can be formalized as:

$$P(Y|X) = \frac{1}{Z(x; \theta)} \exp\left(\sum_{t=1}^T \theta_1^T f_1(X, y_t) + \sum_{t=1}^{T-1} \theta_2^T f_2(X, y_t, y_{t+1})\right). \quad (1)$$

Compared to decoding using MLP only, CRF utilizes the information before the current position. Our experiments found that, when performing multi-task learning, compared to training on a single task, not using CRF will greatly reduce the performance of CWS and POS tagging. The performance improvement brought by CRF is worth the cost of its runtime.

2.3 Biaffine Parser with Output of POS tagging

This task refers to the work of Yan et al. (2020). Yan’s work uses the biaffine parser to solve both CWS and dependency parsing tasks. Compared to Yan’s work, our model will use the output of POS tagging for two reasons. First, dependency parsing has a large semantic and formal gap with other tasks, and sharing the parameter space with other tasks will reduce its performance. Our experimental results show that when the prediction of dependency parsing is independent of other tasks, the performance is worse than that of training dependency parsing only. And using the output of POS, dependency parsing can get more information, such as word segmentation and POS tagging labels. More importantly, users have the need to

obtain all information in one sentence. If running POS tagging and dependency parsing separately, the word segmentation results by the two tasks may conflict, and this contradiction cannot be resolved by engineering methods. Even if there is error propagation in this way, when the POS tagging accuracy is high enough, the impact is acceptable.

When predicting dependency parsing, we first add the POS tagging corpus tag at the head of the original sentence to get the POS tagging output. Then we add the corpus tag of dependency parsing at the head of the original sentence to get the feature vector. Then, using the word segmentation results from POS tagging to split the feature vector of dependency parsing by token. The feature vectors of characters in a token are averaged to represent the token. In addition, embedding is established for POS tagging labels, with the same dimension as the feature vector. The feature vector of each token is added to the embedding vector by position, and the result is input into the biaffine parser. During the training phase, the model uses golden POS tagging labels. Of course, the premise of using POS tagging output is that the corpus contains both dependency parsing and POS tagging information.

2.4 Training Strategy

We mainly applied three training strategies: warm-up, Theseus strategy, and task iteration strategy.

Warm-Up: The idea of warm-up comes from He et al. (2016). During fine-tuning, in order to prevent randomly initialized parameters from causing pre-trained BERT to oscillate, the learning rate increases linearly from 0 to the maximum. In order to a better fit, the learning rate decreases linearly to 0 after reaching the maximum value. Finally, the learning rate curve and the x-axis form a triangle.

Theseus Strategy: Theseus strategy (Xu et al., 2020) is a way to compress BERT, and we use it to train the base version of the model. As shown in Figure 2, after getting the large version of the model we use the module replacement strategy to train the four-layer base version of the model. The base version of the model is initialized with the first four layers of BERT, and its layer i is bound to the layer $2i - 1$ and $2i$ of the large version of the model. They are the corresponding modules. The training phase is divided into two parts. In the first part, we randomly choose whether to replace the corresponding module in the A version of the model, and we make a choice in each module. We

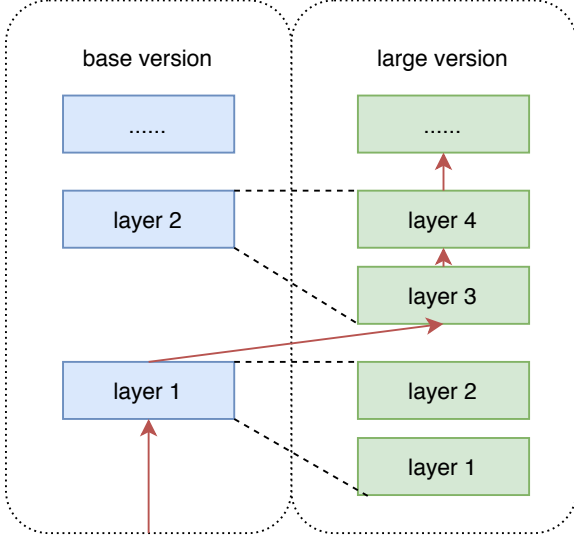


Figure 2: This diagram shows how to train the base version of the model using Theseus strategy. The red arrow represents a possible data path during training.

freeze the parameters of the large version of the model when using gradients to update parameters. The replacement probability p is initialized to 0.5 and decreases linearly to 0. In the second part, We only use the base version of the model, no replacement.

Task Iteration Strategy: In training phase, the samples of each step comes from the same task, and the task selection of every step requires a strategy. We tried three strategies.

The first is to iterate one by one in sequence, skipping tasks that have already been iterated. The disadvantage of this method is that in the end the model is only trained on the task with the largest amount of data, resulting in performance bias on different tasks.

The second is to iterate by probability according to the remaining data of each task. The probability of each task being selected can be formalized as:

$$p_i = \frac{n_i}{\sum_{j=1}^T n_j}. \quad (2)$$

Where p_i represents the probability of the i -th task being selected, n_i represents the current remaining amount of data for the i -th task. When the data volume of each task is very different, the task with fewer data will be iterated at the end of each epoch. However, this will conflict with the warm-up strategy, because the learning rate of tasks with a small amount of data is either too large or too small in each epoch.

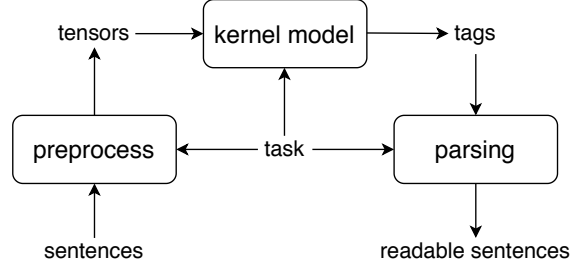


Figure 3: The workflow of fastHan. At each stage, fastHan needs to act according to the task currently being performed.

The third strategy is to divide all tasks into k shares as much as possible, where k is the data amount of the task with the least data. In this way, each epoch can be divided into k big iterations. Within each big iteration, $\frac{1}{k}$ of the data of each task is taken, and then each task is iterated one by one in sequence. This can solve the shortcomings of the first two strategies.

3 fastHan

FastHan is a Chinese NLP toolkit based on the above model, developed based on fastNLP and PyTorch. FastNLP is a modularized and extensible NLP framework that can quickly implement NLP tasks and build complex models.

To install fastHan, users just need one line of command:

```
pip install fastHan
```

3.1 Workflow

When FastHan initializes, it first loads the pre-trained model parameters from the file system. Then, fastHan uses the parameters to initialize BERT and randomly initialize a kernel model. At last, fastHan load the pre-trained parameters to the kernel model.

After initialization, FastHan’s workflow is shown in Figure 3. In the preprocessing stage, fastHan first adds a corpus tag to the head of each sentence according to the current task and then uses the vocabulary to convert the sentence into a batch of vectors as well as padding. FastHan is robust and does not preprocess the original sentence redundantly, such as removing stop words, processing numbers, and English words.

In the parsing phase, fastHan first converts the label sequence into character form and then parses it. The final result will divide the original sentence


```

from fastHan import FastHan
model=FastHan()
sentences=['我爱踢球。','林丹是冠军']
answer=model(sentences, 'Parsing')
for i,sentence in enumerate(answer):
    print(i)
    for token in sentence:
        print(token, token.pos, token.head, token.head_label)
answer_ner=model(sentences, 'NER')
print(answer_ner)

```

```

0
我 PN 2 nsubj
爱 VV 0 root
踢球 VV 2 ccomp
。 PU 2 punct
1
林丹 NR 2 top
是 VC 0 root
冠军 NN 2 attr
[[[], [['林丹', 'NR']]]

```

Figure 4: An example of using fastHan. On the left is the code entered by the user, and on the right is the corresponding output.

into several tokens, each token has some attributes. FastHan defines the Sentence and Token classes to store the output results, which is highly readable.

3.2 Usage

As shown in Figure 4, fastHan is very easy to use. It only needs one line of code to initialize, where users can choose to use the base or large version of the model. When using it for the first time, fastHan will automatically download the parameters.

When calling fastHan, users need to select the task to be performed. The POS tagging task contains the information of CWS, while the dependency parsing task contains the POS tagging information. The information on the NER task is independent of other tasks. POS tagging and dependency parsing tasks use CTB label sets, while the NER task uses the MSRA label set.

The input of FastHan can be a string or a list of strings. In the output of fastHan, users can access sentences and tokens in sentences by index. Each token has four attributes, used to represent information of POS, NER, and dependency parsing.

3.3 Other Functions

FastHan can use the `set_cws_style` function to change the word segmentation style. Each CWS corpus has different granularity and coverage. By changing the corpus tag, fastHan will segment words in the style of the corresponding corpus.

Besides, fastHan can use the `set_device` function to change the device where the model is located. By using GPU, fastHan can greatly improve the execution speed.

4 Evaluation

We evaluate fastHan in terms of accuracy and execution speed. Some hyper-parameters of fastHan

Hyper-Parameter	Value
max learning-rate	3e-5
training epoch	10
drop out	0.1
batch size	32
weight decay	0.01
betas	(0.9,0.999)

Table 1: Some important hyper-parameters of fastHan when training the kernel model. Betas is a hyper-parameter for the optimizer.

are shown in Table 1. When training fastHan, we use AdamW optimizer. We use ReLU for all activation functions.

4.1 Accuracy Test

The accuracy test is performed on the test set of training data. We refer to the CWS corpora used by (Huang et al., 2019), including PKU, MSR, AS, CITYU (Emerson, 2005), CTB6 (Xue et al., 2005), SXU (Jin and Chen, 2008), UD, CNC, WTB (Wang et al., 2014) and ZX (Zhang et al., 2014). More details can be found in Huang et al. (2019). For POS tagging and dependency parsing, we use the Penn Chinese Treebank 9.0 (CTB-9) (Xue et al., 2005). For NER, we use MSRA and OntoNotes.

We conduct an additional set of experiments to make the base version of fastHan trained on each task separately. The final result is shown in Table 2. Both base and large models perform satisfactorily. It can be found that multi-task learning greatly improves fastHan’s performance on all tasks. The large version of fastHan outperforms the current best model in tasks other than dependency parsing.

Model	CWS F	Dependency F_{dep}	Parsing F_{ldep}	POS F	NER MSRA F	NER OntoNotes F
SOTA models	97.1	85.66	81.71	93.15	95.25	79.92
fastHan base trained separately	97.15	80.2	75.12	94.27	92.2	80.3
fastHan base trained jointly	97.27	81.22	76.71	94.88	94.33	82.86
fastHan large trained jointly	97.41	85.52	81.38	95.66	95.50	83.82

Table 2: The results of fastHan’s accuracy result. The score of CWS is the average of 10 corpora. When training dependency parsing separately, the biaffine parser use the same architecture as Yan et al. (2020). SOTA models are best-performing work we know for each task. They came from Huang et al. (2019), Yan et al. (2020), Meng et al. (2019), Diao et al. (2019) and Jie and Lu (2019) in order.

Models	CPU	GPU
fastHan base	25-55	22-111
fastHan large	14-28	21-97

Table 3: Speed test for fastHan. The numbers in the table represent the average number of sentences processed per second. The smaller value is the processing speed of dependency parsing, and the larger value is the processing speed of other tasks.

4.2 Speed Test

The speed test was performed on a personal computer configured with Intel Core i5-9400f + NVIDIA GeForce GTX 1660ti. The model was speed-tested on the first 800 sentences of the CTB CWS corpus, with an average of 45.2 characters per sentence and a batch size of 8.

The result is shown in Table 3. Dependency parsing runs slower, and the other tasks run at about the same speed. The base model with GPU performs poorly in dependency parsing because the use of POS tagging output in dependency parsing requires a lot of CPU calculations, and the acceleration effect of GPU is less than the burden of information transfer.

5 Conclusion

We presented fastHan, a BERT-based toolkit for CWS, NER, POS, and dependency parsing in Chinese NLP. After our optimization, fastHan has the characteristics of high accuracy, small size, and ease of use.

In the future, fastHan will continue to improve, with better performance and more functions.

References

Xinchi Chen, Zhan Shi, Xipeng Qiu, and XuanJing Huang. 2017. Adversarial multi-criteria learning for

chinese word segmentation. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1193–1203.

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*.

Shizhe Diao, Jiaxin Bai, Yan Song, Tong Zhang, and Yonggang Wang. 2019. Zen: pre-training chinese text encoder enhanced by n-gram representations. *arXiv preprint arXiv:1911.00720*.

Thomas Emerson. 2005. The second international chinese word segmentation bakeoff. In *Proceedings of the fourth SIGHAN workshop on Chinese language Processing*.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.

Weipeng Huang, Xingyi Cheng, Kunlong Chen, Taifeng Wang, and Wei Chu. 2019. Toward fast and accurate neural chinese word segmentation with multi-criteria learning. *arXiv preprint arXiv:1903.04190*.

Zhanming Jie and Wei Lu. 2019. Dependency-guided lstm-crf for named entity recognition. *arXiv preprint arXiv:1909.10148*.

Guangjin Jin and Xiao Chen. 2008. The fourth international chinese language processing bakeoff: Chinese word segmentation, named entity recognition and chinese pos tagging. In *Proceedings of the sixth SIGHAN workshop on Chinese language processing*.

John Lafferty, Andrew McCallum, and Fernando CN Pereira. 2001. Conditional random fields: Probabilistic models for segmenting and labeling sequence data.

Yuxian Meng, Wei Wu, Fei Wang, Xiaoya Li, Ping Nie, Fan Yin, Muyu Li, Qinghong Han, Xiaofei Sun, and Jiwei Li. 2019. Glyce: Glyph-vectors for chinese

- character representations. In *Advances in Neural Information Processing Systems*, pages 2746–2757.
- Hwee Tou Ng and Jin Kiat Low. 2004. Chinese part-of-speech tagging: One-at-a-time or all-at-once? word-based or character-based? In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 277–284.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.
- William Yang Wang, Lingpeng Kong, Kathryn Mazaitis, and William Cohen. 2014. Dependency parsing for weibo: An efficient probabilistic logic programming approach. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1152–1158.
- Zhiguo Wang, Chengqing Zong, and Nianwen Xue. 2013. A lattice-based framework for joint chinese word segmentation, pos tagging and parsing. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 623–627.
- Canwen Xu, Wangchunshu Zhou, Tao Ge, Furu Wei, and Ming Zhou. 2020. Bert-of-theseus: Compressing bert by progressive module replacing. *arXiv preprint arXiv:2002.02925*.
- 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.
- Hang Yan, Xipeng Qiu, and Xuanjing Huang. 2020. A graph-based model for joint chinese word segmentation and dependency parsing. *Transactions of the Association for Computational Linguistics*, 8:78–92.
- Meishan Zhang, Yue Zhang, Wanxiang Che, and Ting Liu. 2014. Type-supervised domain adaptation for joint segmentation and pos-tagging. In *Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics*, pages 588–597.
- Yu Zhang, Zhenghua Li, Houquan Zhou, and Min Zhang. 2020. Is pos tagging necessary or even helpful for neural dependency parsing? *arXiv preprint arXiv:2003.03204*.