

CPT: A Pre-Trained Unbalanced Transformer for Both Chinese Language Understanding and Generation

Yunfan Shao,¹ Zhichao Geng,¹ Yitao Liu,¹ Junqi Dai,¹

Fei Yang,² Li Zhe,² Hujun Bao,² Xipeng Qiu^{*1}

¹ School of Computer Science, Fudan University

² Shanghai Key Laboratory of Intelligent Information Processing, Fudan University

² Zhejiang Lab, Hangzhou, China

{yfshao19,zcgeng20,xpqiu}@fudan.edu.cn, {yangf,lizhe}@zhejianglab.com, bao@cad.zju.edu.cn

Abstract

In this paper, we take the advantage of previous pre-trained models (PTMs) and propose a novel Chinese Pre-trained Unbalanced Transformer (CPT). Different from previous Chinese PTMs, CPT is designed for both natural language understanding (NLU) and natural language generation (NLG) tasks. CPT consists of three parts: a shared encoder, an understanding decoder, and a generation decoder. Two specific decoders with a shared encoder are pre-trained with masked language modeling (MLM) and denoising auto-encoding (DAE) tasks, respectively. With the partially shared architecture and multi-task pre-training, CPT can (1) learn specific knowledge of both NLU or NLG tasks with two decoders and (2) be fine-tuned flexibly that fully exploits the potential of the model. Moreover, the unbalanced Transformer saves the computational and storage cost, which makes CPT competitive and greatly accelerates the inference of text generation. Experimental results on a wide range of Chinese NLU and NLG tasks show the effectiveness of CPT¹.

Introduction

Recently, large-scale pre-trained models (PTMs) have become backbone models for many natural language processing (NLP) tasks (Qiu et al. 2020b). However, existing PTMs are usually trained with different architectures and pre-training tasks. When applying PTMs to a downstream task, we should choose a suitable one as the backbone model according to its pre-training nature. For example, we usually select BERT or RoBERTa (Devlin et al. 2019; Liu et al. 2019) as the backbone model for natural language understanding (NLU) tasks, and BART or GPT (Lewis et al. 2020; Radford 2018) for natural language generation (NLG) tasks. With the success of PTMs in English, many works have been done to train the counterparts for Chinese (Cui et al. 2019a; Sun et al. 2019; Wei et al. 2019; Zhang et al. 2020, 2021; Zeng et al. 2021). However, these Chinese PTMs usually follow the settings of English PTMs, which makes these models focus on either language understanding or language generation, limiting their application to a much wider range of Chinese NLP tasks. Therefore, it is attractive to pre-train a joint model for both NLU and NLG tasks.

^{*}Corresponding Author.

Copyright © 2022, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

¹Code is available at <https://github.com/fastnlp/CPT>

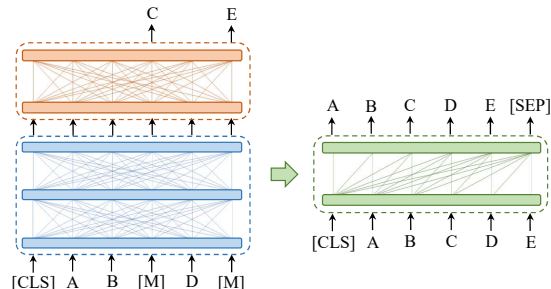


Figure 1: CPT: An unbalanced Transformer encoder-decoder pre-trained with MLM and DAE jointly, which is suitable for both NLU and NLG tasks.

Few works attempt to fuse NLU and NLG into a unified model. UniLMs (Dong et al. 2019; Bao et al. 2020) and GLM (Du et al. 2021) adapt a unified Transformer encoder for both understanding and generation; however, their architectures restrict them to employ more flexible pre-training tasks, such as denoising auto-encoding (DAE) used in BART, a widely successful pre-training task for NLG. PALM (Bi et al. 2020) adopts the standard Transformer and adds an auxiliary masked language modeling (MLM) task to enhance the understanding ability; however, it still focuses on language generation tasks.

In this paper, we propose **CPT**, a novel Chinese Pre-trained Unbalanced Transformer for both NLU and NLG tasks. The architecture of CPT is very concise (as shown in Figure 1), which divides a full Transformer encoder-decoder into three parts: 1) a shared encoder to capture the common representation; 2) a decoder for understanding, which uses full self-attention and is pre-trained with masked language modeling (MLM); 3) a decoder for generation, which adopts masked self-attention and is pre-trained with the DAE task. By multi-task pre-training, CPT is able to improve the performance on both language understanding and generation, respectively.

The main properties of CPT are as follows:

1. CPT can be regarded as two separated PTMs with a shared encoder. Two specific decoders are pre-trained with MLM and DAE tasks, respectively. Each decoder can learn the specific knowledge on either NLU or NLG tasks, while

	BERT RoBERTa	ZEN NEZHA ERNIE-1.0/2.0	PanGu- α	CPM	CPM-2	BART	CPT
# Params	Base - 110M Large - 340M	\approx BERT	32Layers - 2.6B 40Layers - 13.1B 64Layers - 207.0B	Small - 110M Medium - 340M Large - 2.6B	Base - 11B MOE - 198B	Base - 139M Large - 406M	Base - 121M Large - 393M
Arch.	Transformer Encoder	Transformer Encoder Variant	Transformer Decoder	Transformer Decoder	Full Transformer	Full Transformer	Unbalanced Full Transformer
PreTrain. Task	MLM	MLM	LM	LM	Seq2Seq MLM	DAE	MLM+DAE
Tok. Masking Prediction	Char Word Char	Char - Char	Word/Char - Word/Char	Word/Char - Word/Char	Word/Char - Word/Char	Char Word Char	Char Word Char
NLU	✓	✓	✗	✗	✗	✗	✓
NLG	✗	✗	✓	✓	✓	✓	✓

Table 1: Summary of some representative Chinese PTMs. “# Params” refers to the number of parameters. “Arch.” refers to the model architecture. “LM” refers to language modeling in auto-regression fashion, while “Seq2Seq MLM” refers to masked language modeling in Seq2Seq fashion. “Tok.”, “Masking” and “Prediction” refer to the tokenization, masking and prediction granularity of the model, respectively. “✓” means “could be directly used to”. And “✗” means “need to be adapted to”.

the shared encoder learns the common knowledge for universal language representation.

- Two separated decoders enable CPT to adapt to various downstream tasks flexibly. For example, CPT could be fine-tuned with at least five modes for classification tasks (as shown in Figure 2), which exploits the full potential of CPT. Thus, we could choose a suitable fine-tuning mode based on the attributes and characteristics of downstream tasks.
- The overall architecture of CPT is an unbalance Transformer. To make the computational cost and the size of CPT comparable with popular PTMs, such as BERT and BART, we use a novel architecture consisting of a deeper shared encoder and two shallower decoders. Especially, the shallow generation decoder greatly accelerates the inference of text generation.

We conduct experiments on various language understanding and text generation tasks, including datasets for text classification, sequence labeling, machine reading comprehension, summarization, data-to-text generation, etc. Results show that CPT could achieve competitive results with state-of-the-art on these datasets.

Related Work

PTMs towards both NLU and NLG

Recently, there are some efforts to combine language understanding and generation into a single pre-trained model. UniLM (Dong et al. 2019) pre-trained with an ensemble of attention masks, which allows the model to be used for both generative and classification tasks. A difference is that all parameters of UniLM are shared between generation and discrimination, whereas CPT uses two separated decoders. Thus, CPT can utilize the DAE pre-training task which is proven to be effective for NLG tasks (Lewis et al. 2020).

PALM (Bi et al. 2020) is a pre-trained model focusing on conditional generation. To force the encoder to comprehend the meaning of the given context, MLM is added to pre-train the encoder. In contrast, CPT has an individual decoder for MLM which can avoid the negative effects brought by DAE. Therefore CPT also has good performance on NLU tasks.

More recently, ERNIE 3.0 (Sun et al. 2021b) also uses a universal encoder and several task-specific decoders, but it adopts Transformer-XL as the backbone and its generative pre-training task is left-to-right LM with a special masked attention matrix. Different from ERNIE 3.0, CPT adopts the encoder-decoder architecture and is more suitable for sequence-to-sequence (Seq2Seq) tasks.

Chinese PTMs

Many attempts have been conducted to pre-train the Chinese counterparts of PTMs.

The first line of works follows BERT and uses MLM with whole word masking strategy to pre-train Transformer encoder, such as Chinese versions of BERT and RoBERTa (Cui et al. 2019a), NEZHA (Wei et al. 2019), ZEN (Diao et al. 2020). Some of them add special features of Chinese characters or words to further boost the performance of NLU tasks, such as ERNIE 1.0/2.0 (Sun et al. 2019, 2020), ChineseBERT (Sun et al. 2021c). However, these PTMs could not be adopted to text generation directly.

The second line of works follows GPT and uses the left-to-right LM task to pre-train a Transformer decoder, such as CPM (Zhang et al. 2021) and PanGu (Zeng et al. 2021). Although large-scale PTMs with tens of billions parameters have been released recently, the huge computation and storage cost hinders their applications.

The third line of works aims to pre-train the full Transformer encoder-decoder. CPM-2 (Zhang et al. 2021) follows T5 (Raffel et al. 2020) and adopts a Seq2Seq MLM pre-training task, which predicts the masked tokens in a Seq2Seq

fashion. Although BART (Lewis et al. 2020) has achieved widely success on conditional text generation tasks, such as text summarization (Dou et al. 2021; Liu and Liu 2021) and dialogue system (Lin et al. 2020), it still lacks corresponding Chinese versions².

Different from the above Chinese PTMs, CPT is a pre-trained unbalanced Transformer with MLM and DAE tasks, which is capable of achieving competitive results on both NLU and NLG tasks. Besides, CPT is parameter efficient compared to these large-scale models. Table 1 compares different Chinese PTMs.

Model Architecture

As shown in Figure 1, The architecture of CPT is a variant of the full Transformer and consists of three parts:

1. **Shared Encoder** (S-Enc): a Transformer encoder with fully-connected self-attention, which is designed to capture the common semantic representation for both language understanding and generation.
2. **Understanding Decoder** (U-Dec): a shallow Transformer encoder with fully-connected self-attention, which is designed for NLU tasks. The input of U-Dec is the output of S-Enc.
3. **Generation Decoder** (G-Dec): a Transformer decoder with masked self-attention, which is designed for generation tasks with auto-regressive fashion. G-Dec utilizes the output of S-Enc with cross-attention.

With the two specific decoders, CPT can be used flexibly. For example, CPT can be easily fine-tuned for NLU tasks using just S-Enc and U-Dec, and can be regarded as the standard Transformer encoder; while for NLG tasks, CPT adopts S-Enc and G-Dec, and forms a Transformer encoder-decoder. With different combinations, CPT is able to be effectively applied on various downstream tasks, which fully exploits the pre-trained parameters and obtains competitive performance. More combinations and use cases will be discussed in **Fine-Tuning** Section.

Different from most PTMs with encoder-decoders, we exploit a deep-shallow framework for shared encoder and decoders. More specifically, we use a deeper encoder and two shallow decoders for CPT. We assume that a shallow decoder retains the performance on text generation and reduces decoding time, which has proven to be effective for neural machine translation (Kasai et al. 2021) and spell checking (Sun et al. 2021a).

The deep-shallow setup makes CPT more general for both understanding and generative tasks with minor parameter overheads. It also accelerates the inference of CPT for text generation as the G-Dec is a light decoder.

Pre-Training

To make CPT good at both NLU and NLG tasks, we introduce two pre-training tasks.

1. **Masked Language Modeling** (MLM): We pre-train the parameters of S-Enc and U-Dec with MLM (Devlin et al.

2019; Cui et al. 2019a). Given a sentence, we randomly replace some tokens with the [MASK] token and train S-Enc and U-Dec to predict the masked tokens. Following Cui et al. (2019a), we adopt Whole Word Masking (WWM) to replace the tokens. Compared to randomly token masking, WWM is more suitable for inducing semantic information carried by words and spans.

2. **Denoising Auto-Encoding** (DAE): We pre-train the parameters of S-Enc and G-Dec by reconstructing the original document based on the corrupted input. According to the studies of BART (Lewis et al. 2020), we corrupted the input by two effective ways. 1) **Token Infilling**: a Whole Word Masking (WWM) strategy with single mask replacement. First, a number of words are sampled based on the segmentation. Then, each selected word is replaced with a single [MASK] token, regardless of how many tokens it consists; and 2) **Sentence Permutation**: sentences are extracted from a document based on punctuation, and shuffled in a random order.

In practice, We first use a Chinese Word Segmentation (CWS) tool to split the sentences into words. Then, we select 15% of the words and mask the corresponding characters. For the masked characters, we follow the setup of BERT to (1) replace 80% of them with a special [MASK] token, (2) replace 10% of them by random tokens, (3) keep the rest 10% of them unchanged.

Finally, we train CPT with two pre-training tasks under a multi-task learning framework. Thus, CPT can learn for both understanding and generation, and can easily deal with downstream NLU or NLG tasks.

Fine-Tuning

PTMs are usually fine-tuned in only few ways for a given downstream task. For example, for sentence-level classification, we fine-tune BERT by taking the top-layer output of [CLS] token as the representation of the whole sentence, while fine-tune GPT by using the representation of the last token of the sequence.

Thanks to the separated understanding and generation decoders, CPT can be fine-tuned in multiple patterns. For a given downstream task, one could choose the most suitable way to fully stimulate the potential of CPT to achieve competitive results.

Fine-Tuning for Sentence-Level Classification

When incorporating external classifiers, CPT have three fine-tuning modes for sentence-level classification (As shown in Figure 2 (a), (b) and (c)).

1. CPT_u : a BERT-style mode. The sentence representation is from U-Dec module only, which is usually the first state of [CLS] token.
2. CPT_g : a BART-style mode. The same input is fed into the S-Enc and G-Dec, and the representation from the final output token [SEP] from G-Dec is used.
3. CPT_{ug} : The same input is fed into the S-Enc and G-Dec, and the final representation is the concatenation of the first output of U-Dec and the final output of G-Dec.

²Besides CPT, we also provide a Chinese BART as a byproduct.

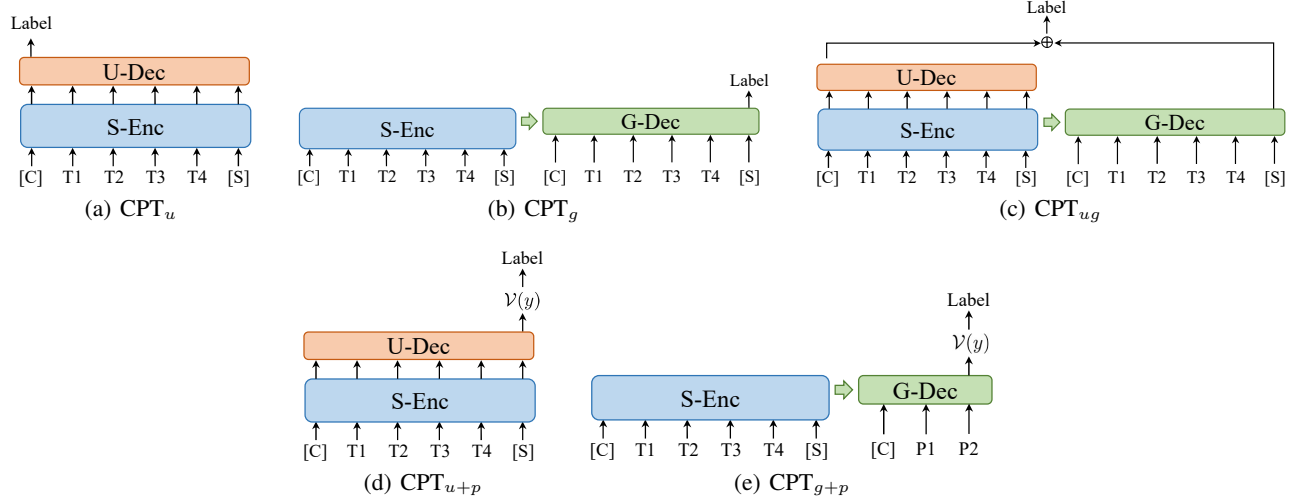


Figure 2: Five ways to fine-tune CPT for text classification. “T1-4” and “P1-2” refer to text input x and prompt tokens, respectively. $\mathcal{V}(y)$ is the mapping function that maps the language model predictions to the label. [C] and [S] are abbreviations for [CLS] and [SEP], respectively.

Recently, a powerful and attractive framework, prompt-based learning (Schick and Schütze 2021; Gao, Fisch, and Chen 2021; Liu et al. 2021), is also able to boost the performance of PTMs. By defining prompting templates and reformulating the classification tasks into a generative fashion, the framework utilizes PTMs to generate words corresponding to task labels. The generative patterns are so close to the pre-training tasks of PTMs that they have the ability of few-shot or even zero-shot learning.

The prompt-based methods could also be applied on CPT with more flexibly fashions since CPT has two decoders. As shown in Figure 2(d) and (e), we construct prompts and convert the task into an generation task with CPT by the following two modes:

1. CPT_{u+p} : A MLM task. We manually construct an input template and assign a word to each task label. CPT is fine-tuned to predict the word at the masked positions, which will be mapped to the task labels. Since a word may be tokenized into multiple character tokens, the predicted distributions at masked positions are averaged to get the predicted distribution of labels.
2. CPT_{g+p} : Conditional text generation. We encode the input text with S-Enc and train CPT to generate prompt text initialized with corresponding labels by teacher forcing. For inference, we first construct the prompt text for each label. Then, the perplexity of each prompt text is calculated. Finally, the prediction is assign to the label with the highest corresponding perplexity.

Fine-Tuning for Sequence Labeling

For sequence labeling, each token needs a representation for token-level classification. Similar to sequence-level classification, we leverage PTMs to obtain high quality token representations and then put the representations to a trainable classifier to assign labels for these tokens. Thus, similar to

sentence-level classification, we can fine-tune CPT for sequence labeling as CPT_u , CPT_g and CPT_{ug} , using (1) U-Dec only, (2) G-Dec only, or (3) both U-Dec and G-Dec. Figure 3 shows two examples for sequence labeling.

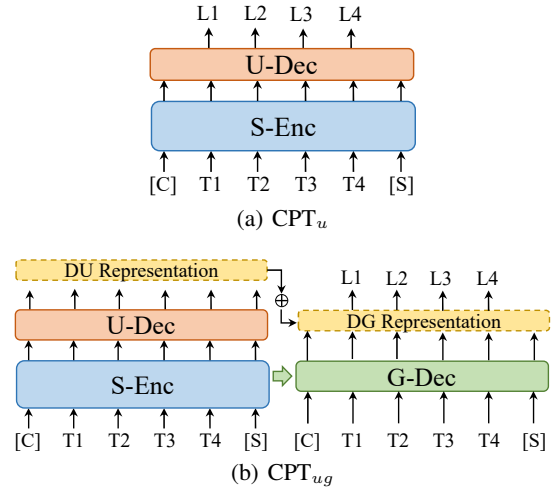


Figure 3: Two examples of fine-tuning CPT for sequence labeling. “T1-4” and “L1-4” refer to text input x and token labels, respectively.

Fine-Tuning for Machine Reading Comprehension

Machine Reading Comprehension requires the model to predict an answer span shown in the passage for a given question. A typical fine-tuning pattern is to train PTMs to predict the start and end positions of the span in the passage. The prediction is based on the tokens of the passage. Thus, CPT_u , CPT_g and CPT_{ug} can be fine-tuned, similar to sequence-labeling.

Figure 4 shows the example of CPT_u .

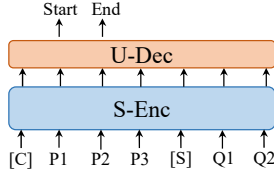


Figure 4: Example of fine-tuning CPT_u for Machine Reading Comprehension. “P1-3”, “Q1-2”, “A1-5” refer to passages, questions and answers, respectively.

Fine-Tuning for Conditional Generation

Apart from NLU tasks, CPT can do text generation efficiently. As shown in Figure 5, we simply fine-tune CPT_g with S-Enc and G-Dec modules on text generation tasks, similar to the usage of other auto-regressive PTMs (Lewis et al. 2020).

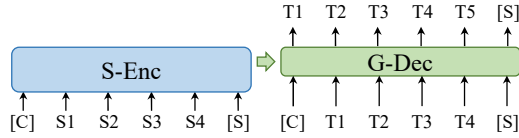


Figure 5: Example of fine-tuning CPT_g for Conditional Generation. “S1-4” and “T1-4” refer to input and target sequences, respectively.

Experiments

Pre-Training Setups

We implement two versions of CPT, namely, *base* and *large*, respectively consisting of 14/28 Transformer layers with 10/20 layers for shared encoder and 2/4 layers for each task specific decoder. And the hidden units and attention heads per layer for base and large versions are 768/1,024 and 12/16, respectively. The total number of layers activated for a given task is always equal to 12/24, which makes our model comparable with base/large-size of BERT and its variants (RoBERTa, ERNIE 1.0/2.0, etc).

We train our models on the open source large-scale raw text, Chinese Wikipedia and a part of WuDaoCorpus. The training data contains 200GB cleaned text ranges from different domains. We use Jieba to segment Chinese words for Whole Word Masking and use WordPiece tokenizer inherited from BERT to split input text into tokens. We use Adam to train the models for 500k steps, with the batch size of 2048, the learning rate of $1e-4$, $\beta_1 = 0.9$, $\beta_2 = 0.98$, weight decay of 0.01. We warmup the learning rate for first 10,000 steps then do linear decay. In addition, a Chinese BART is pre-trained with the same corpora, tokenization and hyperparameters as a baseline.

Evaluation Tasks

To evaluate the effectiveness of our model, we conduct experiments on various NLP datasets across different understanding and generation tasks, with details illustrated below.

Classification We evaluate the model on the Chinese Language Understanding Evaluation Benchmark (CLUE) (Xu et al. 2020), which contains text classification **TNEWS**, **IFLYTEK**, natural language inference (NLI), **OCNLI**, sentence pair matching (SPM) **AFQMC**, and coreference resolution (CoRE) CLUEWSC 2020 (**WSC**.) key word recognition (KwRE) **CSL**. We conduct data augmentation **CSL** as Zhang and Li (2020) performed, and evaluate **TNEWS** on version 1.1 test set. Accuracy is used for these datasets.

Sequence Labeling We evaluate our model on Chinese word segmentation (**CWS**) and named entity recognition (**NER**), which are two representative sequence labeling tasks. We use two datasets from **SIGHAN2005** (Emerson 2005) for CWS, which are **MSR**, **PKU**. And for NER, **MSRA** (Levov 2006), **OntoNotes**³ are used. We use the same dataset pre-processing and split methods as in previous work (Li et al. 2021, 2020; Qiu et al. 2020a). And F1 scores are reported.

MRC Span based machine reading comprehension (MRC) dataset CMRC 2018 (**CMRC**) (Cui et al. 2019b) and Traditional Chinese MRC dataset **DRCD** (Shao et al. 2018) are used. We follow the data processing in Cui et al. (2019a, 2020) and transform the text from **DRCD** is transformed to Simplified Chinese. The Exact Match (EM) scores are reported.

Text Generation We use two abstractive summarization datasets, **LCSTS** (Hu, Chen, and Zhu 2015) and **CSL**⁴, and a data-to-text generation dataset, **ADGEN** (Shao et al. 2019) to evaluate the text generation ability of our model. Among them, **LCSTS** is a large corpus of Chinese short text summarization dataset constructed from Sina Weibo, consisting of 2 million real Chinese short texts with short summaries. And **CSL** is an academic domain text summarization dataset, constructed from abstract and titles from publications in computer science domain. And **ADGEN** is a data-to-text dataset that requires models to generate long text for advertisement based on some keywords. And we evaluate PTMs on test sets of **LCSTS** and **ADGEN** and the development set of **CSL**. The character-level Rouge-L is used to evaluate the summarization results. For **ADGEN**, we follow Zhang et al. (2021) to use **BLEU-4**.

Compared PTMs

We compare CPT with a series of state-of-the-art PTMs for either natural language understanding or text generation. The details are as follows.

PTMs for NLU PTMs with the Transformer Encoder structure and pre-trained with MLM usually perform well in NLU tasks, such as the Chinese versions of BERT and RoBERTa (Cui et al. 2019a), NEZHA (Wei et al. 2019), ERNIE 2.0 (Sun et al. 2020), MacBERT (Cui et al. 2020). Unless otherwise specified, we use BERT and RoBERTa to refer to BERT-wwm-ext and RoBERTa-wwm-ext, respectively.

³<https://catalog.ldc.upenn.edu/LDC2011T03>

⁴<https://github.com/CLUEbenchmark/CLGE>

PTMs for NLG For text generation, we compare CPT with generative Transformers ranging from normal size to large scale, including BART (Lewis et al. 2020), mT5 (Xue et al. 2021), CPM-2 (Zhang et al. 2021), and models with pre-trained encoders. BART is a sequence-to-sequence model pre-trained with DAE task. Due to the missing of Chinese version, we train a Chinese BART as the baseline on the same corpus and tokenization of CPT. mT5 is a multilingual variant of T5 pre-trained on over 101 languages, including Chinese. CPM-2 is a large-scale encoder-decoder model with 11 billion parameters, pre-trained in multiple stages with large-scale Chinese and bilingual data. We also report generative models adopted from Transformer encoders such as RoBERTa and ERNIE 2.0, to further evaluate the effectiveness generative pre-training.

Main Results

To fully release the potential of our model, we fine-tune CPT for NLU tasks in different ways as mentioned in **Fine-Tuning** Section, denoted as CPT_u , CPT_g and CPT_{ug} , CPT_{u+p} and CPT_{g+p} , respectively. We use (B) and (L) to distinguish base and large version of PTMs, respectively.

Classification Table 2 shows the development set results of CLUE Benchmark of different fine-tuning modes. As a result, CPT_u (B) achieves a 74.6 on average, surpassing other baselines and fine-tuning patterns on base version of CPT. Besides, CPT_{ug} (L) obtains a averaged accuracy 76.2, which is better than RoBERTa (L) by a large margin. Therefore, we choose CPT_u (B) and CPT_{ug} (L) as the most suitable fine-tuning patterns to do the classification. And we evaluate them on the test sets, reported in Table 3.

For prompt-based fine-tuning (Table 2), we find that directly fine-tuning without prompt works well on some datasets, with the small gaps between CPT_u , CPT_g and CPT_{ug} . Moreover, CPT_{u+p} achieves good results on some datasets that even outperform methods without prompt tuning. However, the accuracy of prompt-base methods on other datasets drops a lot. As there are many factors that affect prompt tuning performance including prompt design, choices of words for labels, etc. Manually designed prompts may be suboptimal. Besides, we find that CPT_{g+p} degenerates obviously on TNEWS and IFLYTEK. Both datasets have more than 3 classes, which contains 15 and 112 labels, respectively. Moreover, these labels are hard to represented by a single character. In practice we assign words with up to 7 characters to a label. We presume that the large number of labels and the multi-token issue hinders CPT_{g+p} to generate correctly.

Table 3 reports the performance of CPT on classification tasks and the comparison with previous representative Chinese PTMs. We report accuracy on the test sets of these datasets. Among the fine-tuned CPTs, we choose base version CPT_u and large version CPT_{ug} as they obtain the best results on development sets. Base size CPT consistently outperforms BERT, RoBERTa and ERNIE. Moreover, large size CPT achieves a 74.5 averaged score, outperforming RoBERTa (L) with a large margin.

Sequence Labeling The CPT is fine-tuned as CPT_u , CPT_g and CPT_{ug} and evaluated on development sets. We find that

Models	TNEWS	IFLYTEK	OCNLI	AFQMC	CSL	WSC	AVG
BERT (B)	56.8	58.9	75.4	72.0	82.3	83.2	71.4
RoBERTa (B)	57.5	59.4	76.5	74.4	86.1	88.8	73.8
CPT_u (B)	58.4	60.5	76.4	75.1	86.1	91.1	74.6
CPT_g (B)	57.3	60.4	76.3	71.4	86.4	87.2	73.2
CPT_{ug} (B)	57.4	61.9	76.8	70.6	86.3	89.8	73.8
CPT_{g+p} (B)	54.9	25.4	76.6	73.7	86.9	79.9	66.2
CPT_{u+p} (B)	58.4	61.6	76.6	75.1	86.9	79.9	73.1
RoBERTa (L)	58.3	61.7	78.5	75.4	86.3	89.5	75.0
CPT_u (L)	58.8	61.8	79.5	75.9	86.5	92.1	75.8
CPT_g (L)	59.1	61.7	79.9	75.8	86.9	91.8	75.9
CPT_{ug} (L)	59.2	62.4	79.8	75.8	86.6	93.4	76.2
CPT_{g+p} (L)	54.5	29.2	79.8	75.4	87.1	89.5	69.2
CPT_{u+p} (L)	59.0	61.2	79.6	75.4	87.3	87.8	75.1

Table 2: Accuracy results on dev set of CLUE Benchmark. We fine-tune CPT with five different ways as shown in Figure 2. (B) and (L) refer to base-size and large-size of PTMs, respectively.

Models	TNEWS	IFLYTEK	OCNLI	AFQMC	CSL	WSC	AVG
BERT (B)	58.6	59.4	73.2	74.1	84.2	74.5	70.7
RoBERTa (B)	59.5	60.3	73.9	74.0	84.7	76.9	71.5
CPT_u (B)	59.2	60.5	73.4	74.4	85.5	81.4	72.4
RoBERTa (L)	58.9	63.0	76.4	76.6	82.1	74.6	71.9
CPT_{ug} (L)	59.2	62.4	78.4	75.0	85.5	86.2	74.5

Table 3: Results on CLUE benchmarks. For all tasks we report accuracy on test sets.

CPT_u constantly obtains the best development results. We conjecture that CWS and NER have more dependency on local syntax than complex semantics used for text generation. Thus, CPT_u is more suitable for CWS and NER with its bidirectional fully connected self-attention. As a result, we report the test set results of CPT_u to compare with other PTMs.

	CWS		NER	
	MSR	PKU	MSRA	OntoNotes
BERT (B)	98.24	96.50	95.13	81.73
ERNIE 2.0* (B)	-	-	93.80	-
RoBERTa (B)	98.14	96.15	95.23	81.52
CPT_u (B)	98.29	96.58	95.78	82.08
ERNIE 2.0* (L)	-	-	95.00	-
RoBERTa (L)	98.42	96.37	95.20	81.78
CPT_u (L)	98.51	96.70	96.20	83.08

Table 4: Results on sequence labeling datasets. The F1 scores on test sets are reported. Models with * indicate the results are from Sun et al. (2020).

We compare our model with other state-of-the-art methods on sequence labeling datasets. As shown in Table 4, CPT_u (L) achieves the highest performance and exceed the BERT (L), RoBERTa (L) and ERNIE (L) on all sequence labeling tasks,

both CWS and NER. And CPT_u (B) obtains a comparable results, surpassing base versions of BERT and RoBERTa.

MRC Table 5 shows the experimental results on MRC tasks, which also indicates the effectiveness of CPT. We report the Exact Match (EM) score on CMRC dev set, DRCD dev and test sets. We try and evaluate CPT_u , CPT_u and CPT_u on the development sets of these datasets and choose the pattern that acquires the best results to report. As a conclusion, CPT_u obtains comparable or higher results compared to previous systems that are widely used, such as RoBERTa, MacBERT, ERNIE and NEZHA. Moreover, CPT_u consistently outperforms other strong baselines by a large margin, with 72.3 EM score on the CMRC development set and 91.1 EM on the DRCD test set.

	CMRC 2018	DRCD	
	Dev	Dev	Test
RoBERTa (B)	67.9	85.9	85.2
MacBERT* (B)	68.2	89.2	88.7
ERNIE 2.0* (B)	69.1	88.5	88.0
NEZHA* (B)	67.8	-	-
CPT_u (B)	68.8	89.0	89.0
RoBERTa (L)	70.6	89.1	88.9
MacBERT* (L)	70.1	90.8	90.9
ERNIE 2.0* (L)	71.5	89.7	89.0
NEZHA* (L)	68.1	-	-
CPT_u (L)	72.3	91.0	91.1

Table 5: Results on MRC datasets. Exact Match (EM) scores are reported. Models with * indicate the results from the corresponding work.

Text Generation Table 6 compares the performance of our model on generation datasets with other strong methods. The character-level Rouge-L is used to evaluate the summarization results. For ADGEN, we follow Zhang et al. (2021) to use BLEU-4.

As a conclusion, CPT_g achieves competitive performance on text generation compared with other methods, such as mT5, CPM-2, BART. In addition, compared with other pre-trained encoders (RoBERTa and ERNIE 2.0), CPT_g improves the generation score with the NLG enhanced pre-training. When compared with pre-trained mT5 and CPM-2, CPT_g acquires better results on both base and large versions. We assume the difference of pre-training tasks that lead to the performance gaps. Both mT5 and CPM-2 exploit a T5 style masked span generation as their pre-training task, while CPT is pre-trained with DAE, which shows the effectiveness of DAE for text generation pre-training. In addition, while BART (L) have a slightly better results on CSL and ADGEN, CPT_g and BART have similar results on text generation. The shallow decoder of CPT_g may affect the performance on long text generation. However, the performance gaps are still small. And we believe the pre-training of the shallow decoder closes the gaps.

Moreover, because of the shallow decoder, CPT could generate texts more efficiently (Figure 6), which could be faster

Models	LCSTS (Rouge-L)	CSL (Rouge-L)	ADGEN (BLEU-4)
mT5 (S)	33.5	56.7	10.2
BART (B)	37.8	62.1	9.9
CPT_g (B)	38.2	63.0	9.8
CPM-2 [†]	35.9	-	10.6
mT5 (B)	36.5	61.8	-
ERNIE 2.0* (L)	41.4	-	-
RoBERTa* (L)	41.0	-	-
BART (L)	40.6	64.2	10.0
CPT_g (L)	42.0	63.7	10.7

Table 6: Results on text generation datasets. The small(base) version of mT5 has almost the same parameters as the base(large) version of other PTMs. CPM-2 has a much larger number of parameters than other large size PTMs. Models with * and [†] indicate the results are from Sun et al. (2021b) and Zhang et al. (2021), respectively.

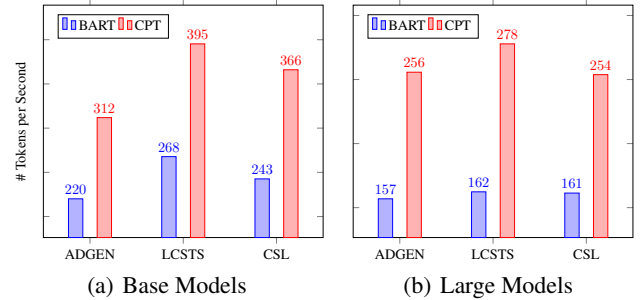


Figure 6: Inference throughput for BART and CPT. It is measured on the same parts of datasets that the models are evaluated. The beam size is 4 and the batch size is 8.

than other depth symmetric encoder-decoder Transformers with the same number of layers of the encoder and the decoder. As BART and CPT have similar number of parameters in both base and large versions. On all generation dataset, the decoding speed of CPT surpass BART with a large margin. Our model achieves $1.4 \times \sim 1.5 \times$ speedup compared with BART and still maintain comparable generation results in base size. And CPT (L) has up to $1.7 \times$ relative speedup compared to BART (L). As a conclusion, the shallow G-Dec is able to speed up the generation with minor performance loss.

Conclusion

In this paper, we propose CPT, a novel Chinese PTM for both language understanding and generation. With the flexible design, CPT can be assembled and disassembled in various fashions, which could fully exploit the potential of CPT. Experimental results on a wide range of Chinese NLU and NLG tasks show the effectiveness of CPT.

In future work, we will introduce more specific designs according to Chinese properties, such as better tokenization, pre-training tasks and model architectures.

References

- Bao, H.; Dong, L.; Wei, F.; Wang, W.; Yang, N.; Liu, X.; Wang, Y.; Gao, J.; Piao, S.; Zhou, M.; and Hon, H. 2020. UniLMv2: Pseudo-Masked Language Models for Unified Language Model Pre-Training. In *ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, 642–652. PMLR.
- Bi, B.; Li, C.; Wu, C.; Yan, M.; Wang, W.; Huang, S.; Huang, F.; and Si, L. 2020. PALM: Pre-training an Autoencoding&Autoregressive Language Model for Context-conditioned Generation. In Webber, B.; Cohn, T.; He, Y.; and Liu, Y., eds., *EMNLP 2020, Online, November 16-20, 2020*, 8681–8691.
- Cui, Y.; Che, W.; Liu, T.; Qin, B.; Wang, S.; and Hu, G. 2020. Revisiting Pre-Trained Models for Chinese Natural Language Processing. In Cohn, T.; He, Y.; and Liu, Y., eds., *EMNLP 2020, Online Event, 16-20 November 2020*, 657–668.
- Cui, Y.; Che, W.; Liu, T.; Qin, B.; Yang, Z.; Wang, S.; and Hu, G. 2019a. Pre-Training with Whole Word Masking for Chinese BERT. *CoRR*, abs/1906.08101.
- Cui, Y.; Liu, T.; Che, W.; Xiao, L.; Chen, Z.; Ma, W.; Wang, S.; and Hu, G. 2019b. A Span-Extraction Dataset for Chinese Machine Reading Comprehension. In Inui, K.; Jiang, J.; Ng, V.; and Wan, X., eds., *EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, 5882–5888.
- Devlin, J.; Chang, M.; Lee, K.; and Toutanova, K. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In Burstein, J.; Doran, C.; and Solorio, T., eds., *NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, 4171–4186.
- Diao, S.; Bai, J.; Song, Y.; Zhang, T.; and Wang, Y. 2020. ZEN: Pre-training Chinese Text Encoder Enhanced by N-gram Representations. In Cohn, T.; He, Y.; and Liu, Y., eds., *EMNLP 2020, Online Event, 16-20 November 2020*, 4729–4740.
- Dong, L.; Yang, N.; Wang, W.; Wei, F.; Liu, X.; Wang, Y.; Gao, J.; Zhou, M.; and Hon, H. 2019. Unified Language Model Pre-training for Natural Language Understanding and Generation. In Wallach, H. M.; Larochelle, H.; Beygelzimer, A.; d’Alché-Buc, F.; Fox, E. B.; and Garnett, R., eds., *NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, 13042–13054.
- Dou, Z.; Liu, P.; Hayashi, H.; Jiang, Z.; and Neubig, G. 2021. GSum: A General Framework for Guided Neural Abstractive Summarization. In Toutanova, K.; Rumshisky, A.; Zettlemoyer, L.; Hakkani-Tür, D.; Beltagy, I.; Bethard, S.; Cotterell, R.; Chakraborty, T.; and Zhou, Y., eds., *NAACL-HLT 2021, Online, June 6-11, 2021*, 4830–4842.
- Du, Z.; Qian, Y.; Liu, X.; Ding, M.; Qiu, J.; Yang, Z.; and Tang, J. 2021. All NLP Tasks Are Generation Tasks: A General Pretraining Framework. *CoRR*, abs/2103.10360.
- Emerson, T. 2005. The Second International Chinese Word Segmentation Bakeoff. In *SIGHAN@IJCNLP 2005, Jeju Island, Korea, 14-15, 2005*. ACL.
- Gao, T.; Fisch, A.; and Chen, D. 2021. Making Pre-trained Language Models Better Few-shot Learners. In Zong, C.; Xia, F.; Li, W.; and Navigli, R., eds., *ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, 3816–3830.
- Hu, B.; Chen, Q.; and Zhu, F. 2015. LCSTS: A Large Scale Chinese Short Text Summarization Dataset. In Márquez, L.; Callison-Burch, C.; Su, J.; Pighin, D.; and Marton, Y., eds., *EMNLP 2015, Lisbon, Portugal, September 17-21, 2015*, 1967–1972. The Association for Computational Linguistics.
- Kasai, J.; Pappas, N.; Peng, H.; Cross, J.; and Smith, N. A. 2021. Deep Encoder, Shallow Decoder: Reevaluating Non-autoregressive Machine Translation. In *ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.
- Levow, G. 2006. The Third International Chinese Language Processing Bakeoff: Word Segmentation and Named Entity Recognition. In Ng, H. T.; and Kwong, O. O. Y., eds., *SIGHAN@COLING/ACL 2006, Sydney, Australia, July 22-23, 2006*, 108–117.
- Lewis, M.; Liu, Y.; Goyal, N.; Ghazvininejad, M.; Mohamed, A.; Levy, O.; Stoyanov, V.; and Zettlemoyer, L. 2020. BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension. In Jurafsky, D.; Chai, J.; Schluter, N.; and Tetraault, J. R., eds., *ACL 2020, Online, July 5-10, 2020*, 7871–7880.
- Li, X.; Shao, Y.; Sun, T.; Yan, H.; Qiu, X.; and Huang, X. 2021. Accelerating BERT Inference for Sequence Labeling via Early-Exit. In Zong, C.; Xia, F.; Li, W.; and Navigli, R., eds., *ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, 189–199.
- Li, X.; Yan, H.; Qiu, X.; and Huang, X. 2020. FLAT: Chinese NER Using Flat-Lattice Transformer. In Jurafsky, D.; Chai, J.; Schluter, N.; and Tetraault, J. R., eds., *ACL 2020, Online, July 5-10, 2020*, 6836–6842.
- Lin, Z.; Madotto, A.; Winata, G. I.; and Fung, P. 2020. MinTL: Minimalist Transfer Learning for Task-Oriented Dialogue Systems. In Webber, B.; Cohn, T.; He, Y.; and Liu, Y., eds., *EMNLP 2020, Online, November 16-20, 2020*, 3391–3405.
- Liu, P.; Yuan, W.; Fu, J.; Jiang, Z.; Hayashi, H.; and Neubig, G. 2021. Pre-train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing. *CoRR*, abs/2107.13586.
- Liu, Y.; and Liu, P. 2021. SimCLS: A Simple Framework for Contrastive Learning of Abstractive Summarization. In Zong, C.; Xia, F.; Li, W.; and Navigli, R., eds., *ACL/IJCNLP 2021, (Volume 2: Short Papers), Virtual Event, August 1-6, 2021*, 1065–1072.
- Liu, Y.; Ott, M.; Goyal, N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; and Stoyanov, V. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *CoRR*, abs/1907.11692.
- Qiu, X.; Pei, H.; Yan, H.; and Huang, X. 2020a. A Concise Model for Multi-Criteria Chinese Word Segmentation with Transformer Encoder. In Cohn, T.; He, Y.; and Liu, Y., eds., *Findings of EMNLP 2020*, 2887–2897.

- Qiu, X.; Sun, T.; Xu, Y.; Shao, Y.; Dai, N.; and Huang, X. 2020b. Pre-trained Models for Natural Language Processing: A Survey. *SCIENCE CHINA Technological Sciences*, 63(10): 1872–1897.
- Radford, A. 2018. Improving Language Understanding by Generative Pre-Training.
- Raffel, C.; Shazeer, N.; Roberts, A.; Lee, K.; Narang, S.; Matena, M.; Zhou, Y.; Li, W.; and Liu, P. J. 2020. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *J. Mach. Learn. Res.*, 21: 140:1–140:67.
- Schick, T.; and Schütze, H. 2021. It’s Not Just Size That Matters: Small Language Models Are Also Few-Shot Learners. In Toutanova, K.; Rumshisky, A.; Zettlemoyer, L.; Hakkani-Tür, D.; Beltagy, I.; Bethard, S.; Cotterell, R.; Chakraborty, T.; and Zhou, Y., eds., *NAACL-HLT 2021, Online, June 6-11, 2021*, 2339–2352.
- Shao, C.; Liu, T.; Lai, Y.; Tseng, Y.; and Tsai, S. 2018. DRCD: a Chinese Machine Reading Comprehension Dataset. *CoRR*, abs/1806.00920.
- Shao, Z.; Huang, M.; Wen, J.; Xu, W.; and Zhu, X. 2019. Long and Diverse Text Generation with Planning-based Hierarchical Variational Model. In Inui, K.; Jiang, J.; Ng, V.; and Wan, X., eds., *EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, 3255–3266.
- Sun, X.; Ge, T.; Wei, F.; and Wang, H. 2021a. Instantaneous Grammatical Error Correction with Shallow Aggressive Decoding. In Zong, C.; Xia, F.; Li, W.; and Navigli, R., eds., *ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, 5937–5947.
- Sun, Y.; Wang, S.; Feng, S.; Ding, S.; Pang, C.; Shang, J.; Liu, J.; Chen, X.; Zhao, Y.; Lu, Y.; Liu, W.; Wu, Z.; Gong, W.; Liang, J.; Shang, Z.; Sun, P.; Liu, W.; Ouyang, X.; Yu, D.; Tian, H.; Wu, H.; and Wang, H. 2021b. ERNIE 3.0: Large-scale Knowledge Enhanced Pre-training for Language Understanding and Generation. *CoRR*, abs/2107.02137.
- Sun, Y.; Wang, S.; Li, Y.; Feng, S.; Chen, X.; Zhang, H.; Tian, X.; Zhu, D.; Tian, H.; and Wu, H. 2019. ERNIE: Enhanced Representation through Knowledge Integration. *CoRR*, abs/1904.09223.
- Sun, Y.; Wang, S.; Li, Y.; Feng, S.; Tian, H.; Wu, H.; and Wang, H. 2020. ERNIE 2.0: A Continual Pre-Training Framework for Language Understanding. In *AAAI 2020, New York, NY, USA, February 7-12, 2020*, 8968–8975. AAAI Press.
- Sun, Z.; Li, X.; Sun, X.; Meng, Y.; Ao, X.; He, Q.; Wu, F.; and Li, J. 2021c. ChineseBERT: Chinese Pretraining Enhanced by Glyph and Pinyin Information. *CoRR*, abs/2106.16038.
- Wei, J.; Ren, X.; Li, X.; Huang, W.; Liao, Y.; Wang, Y.; Lin, J.; Jiang, X.; Chen, X.; and Liu, Q. 2019. NEZHA: Neural Contextualized Representation for Chinese Language Understanding. *CoRR*, abs/1909.00204.
- Xu, L.; Hu, H.; Zhang, X.; Li, L.; Cao, C.; Li, Y.; Xu, Y.; Sun, K.; Yu, D.; Yu, C.; Tian, Y.; Dong, Q.; Liu, W.; Shi, B.; Cui, Y.; Li, J.; Zeng, J.; Wang, R.; Xie, W.; Li, Y.; Patterson, Y.; Tian, Z.; Zhang, Y.; Zhou, H.; Liu, S.; Zhao, Z.; Zhao, Q.; Yue, C.; Zhang, X.; Yang, Z.; Richardson, K.; and Lan, Z. 2020. CLUE: A Chinese Language Understanding Evaluation Benchmark. In Scott, D.; Bel, N.; and Zong, C., eds., *COLING 2020, Barcelona, Spain (Online), December 8-13, 2020*, 4762–4772. International Committee on Computational Linguistics.
- Xue, L.; Constant, N.; Roberts, A.; Kale, M.; Al-Rfou, R.; Siddhant, A.; Barua, A.; and Raffel, C. 2021. mT5: A Massively Multilingual Pre-trained Text-to-Text Transformer. In Toutanova, K.; Rumshisky, A.; Zettlemoyer, L.; Hakkani-Tür, D.; Beltagy, I.; Bethard, S.; Cotterell, R.; Chakraborty, T.; and Zhou, Y., eds., *NAACL-HLT 2021, Online, June 6-11, 2021*, 483–498.
- Zeng, W.; Ren, X.; Su, T.; Wang, H.; Liao, Y.; Wang, Z.; Jiang, X.; Yang, Z.; Wang, K.; Zhang, X.; Li, C.; Gong, Z.; Yao, Y.; Huang, X.; Wang, J.; Yu, J.; Guo, Q.; Yu, Y.; Zhang, Y.; Wang, J.; Tao, H.; Yan, D.; Yi, Z.; Peng, F.; Jiang, F.; Zhang, H.; Deng, L.; Zhang, Y.; Lin, Z.; Zhang, C.; Zhang, S.; Guo, M.; Gu, S.; Fan, G.; Wang, Y.; Jin, X.; Liu, Q.; and Tian, Y. 2021. PanGu- α : Large-scale Autoregressive Pretrained Chinese Language Models with Auto-parallel Computation. *CoRR*, abs/2104.12369.
- Zhang, X.; and Li, H. 2020. AMBERT: A Pre-trained Language Model with Multi-Grained Tokenization. *CoRR*, abs/2008.11869.
- Zhang, Z.; Gu, Y.; Han, X.; Chen, S.; Xiao, C.; Sun, Z.; Yao, Y.; Qi, F.; Guan, J.; Ke, P.; Cai, Y.; Zeng, G.; Tan, Z.; Liu, Z.; Huang, M.; Han, W.; Liu, Y.; Zhu, X.; and Sun, M. 2021. CPM-2: Large-scale Cost-effective Pre-trained Language Models. *CoRR*, abs/2106.10715.
- Zhang, Z.; Han, X.; Zhou, H.; Ke, P.; Gu, Y.; Ye, D.; Qin, Y.; Su, Y.; Ji, H.; Guan, J.; Qi, F.; Wang, X.; Zheng, Y.; Zeng, G.; Cao, H.; Chen, S.; Li, D.; Sun, Z.; Liu, Z.; Huang, M.; Han, W.; Tang, J.; Li, J.; Zhu, X.; and Sun, M. 2020. CPM: A Large-scale Generative Chinese Pre-trained Language Model. *CoRR*, abs/2012.00413.