

Pre-trained Models for Natural Language Processing: A Survey

Xipeng Qiu*, Tianxiang Sun, Yige Xu, Yunfan Shao, Ning Dai & Xuanjing Huang

*School of Computer Science, Fudan University, Shanghai 200433, China;
Shanghai Key Laboratory of Intelligent Information Processing, Shanghai 200433, China*

Recently, the emergence of pre-trained models (PTMs) has brought natural language processing (NLP) to a new era. In this survey, we provide a comprehensive review of PTMs for NLP. We first briefly introduce language representation learning and its research progress. Then we systematically categorize existing PTMs based on a taxonomy with four perspectives. Next, we describe how to adapt the knowledge of PTMs to the downstream tasks. Finally, we outline some potential directions of PTMs for future research. This survey is purposed to be a hands-on guide for understanding, using, and developing PTMs for various NLP tasks.

Deep Learning, Neural Network, Natural Language Processing, Pre-trained Model, Distributed Representation, Word Embedding, Self-Supervised Learning, Language Modelling

1 Introduction

With the development of deep learning, various neural networks have been widely used to solve Natural Language Processing (NLP) tasks, such as convolutional neural networks (CNNs) [75, 80, 45], recurrent neural networks (RNNs) [160, 100], graph-based neural networks (GNNs) [146, 161, 111] and attention mechanisms [6, 171]. One of the advantages of these neural models is their ability to alleviate the *feature engineering* problem. Non-neural NLP methods usually heavily rely on the discrete handcrafted features, while neural methods usually use low-dimensional and dense vectors (aka. *distributed representation*) to implicitly represent the syntactic or semantic features of the language. These representations are learned in specific NLP tasks. Therefore, neural methods make it easy for people to develop various NLP systems.

Despite the success of neural models for NLP tasks, the performance improvement may be less significant compared to the Computer Vision (CV) field. The main reason is that current datasets for most supervised NLP tasks are rather small (except machine translation). Deep neural networks usually

have a large number of parameters which make them overfit on these small training data and do not generalize well in practice. Therefore, the early neural models for many NLP tasks were relatively shallow and usually consisted of only 1~3 neural layers.

Recently, substantial work has shown that pre-trained models (PTMs) on the large corpus can learn universal language representations, which are beneficial for downstream NLP tasks and can avoid training a new model from scratch. With the development of computational power, the emergence of the deep models (i.e., Transformer [171]) and the constant enhancement of training skills, the architecture of PTMs has been advanced from shallow to deep. The *first-generation PTMs* aim to learn good word embeddings. Since these models themselves are no longer needed by downstream tasks, they are usually very shallow for computational efficiencies, such as Skip-Gram [116] and GloVe [120]. Although these pre-trained embeddings can capture semantic meanings of words, they are context-free and fail to capture higher-level concepts of text like syntactic structures, semantic roles, anaphora, etc.

* Corresponding author (email: xpqiu@fudan.edu.cn)

The *second-generation PTMs* focus on learning contextual word embeddings, such as CoVe [113], ELMo [122], OpenAI GPT [130] and BERT [32]. These learned encoders are still needed to represent words in context by downstream tasks. Besides, various pre-training tasks are also proposed to learn PTMs for different purposes.

The contributions of this survey can be summarized as follows:

1. *Comprehensive review.* We provide a comprehensive review of PTMs for NLP, including background knowledge, model architecture, pre-training tasks, various extensions, adaption approaches, and applications. We provide detailed descriptions of representative models, make the necessary comparison, and summarise the corresponding algorithms.
2. *New taxonomy.* We propose a taxonomy of PTMs for NLP, which categorizes existing PTMs from four different perspectives: 1) type of word representation; 2) architecture of PTMs; 3) type of pre-training tasks; 4) extensions for specific types of scenarios or inputs.
3. *Abundant resources.* We collect abundant resources on PTMs, including open-source systems, paper lists, etc.
4. *Future directions.* We discuss and analyze the limitations of existing PTMs. Also, we suggest possible future research directions.

The rest of the survey is organized as follows. Section 2 outlines the background concepts and commonly used notations of PTMs. Section 3 gives a brief overview of PTMs and clarifies the categorization of PTMs. Section 4 provides extensions of PTMs. Section 5 discusses how to transfer the knowledge of PTMs to downstream tasks. Section 6 gives the related resources on PTMs, including open-source systems, paper lists, etc. Section 7 presents a collection of applications across various NLP tasks. Section 8 discusses the current challenges and suggests future directions. Section 9 summarizes the paper.

2 Background

2.1 Language Representation Learning

As suggested by Bengio et al. [12], a good representation should express general-purpose priors that are not task-specific but would be likely to be useful for a learning machine to solve AI-tasks. When it comes to language, a good representation should capture the implicit linguistic rules and common sense knowledge hiding in text data, such as lexical meanings, syntactic structures, semantic roles, and even pragmatics.

The core idea of distributed representation is to describe the meaning of a piece of text by low-dimensional real-valued vectors. And each dimension of the vector has no corresponding sense while the whole represents a concrete concept. Figure 1 illustrates the generic neural architecture for NLP. There are two kinds of word embeddings: non-contextual and contextual embeddings. The difference between them is whether the embedding for a word dynamically changes according to the context it appears in.

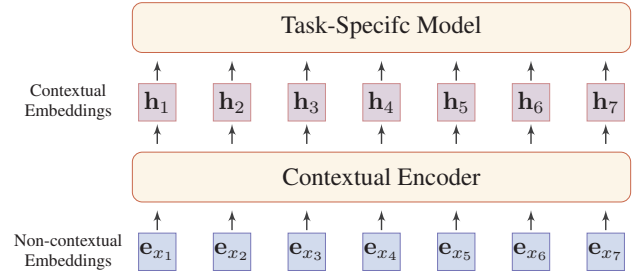


Figure 1: Generic Neural Architecture for NLP.

Non-contextual Embeddings To represent language, the first step is to map discrete language symbols into a distributed embedding space. Formally, for each word (or sub-word) x in a vocabulary \mathcal{V} , we map it to a vector $\mathbf{e}_x \in \mathbb{R}^{D_e}$ with a lookup table $\mathbf{E} \in \mathbb{R}^{D_e \times |\mathcal{V}|}$, where D_e is a hyper-parameter indicating the dimension of token embeddings. These embeddings are trained on task data along with other model parameters.

There are two main limitations to this kind of embeddings. The first issue is that the embeddings are static. The embedding for a word does is always the same regardless of its context. Therefore, these *non-contextual embeddings* fail to model polysemous words. The second issue is the out-of-vocabulary problem. To tackle this problem, character-level word representations or sub-word representations are widely used in many NLP tasks, such as CharCNN [81], FastText [13] and Byte-Pair Encoding (BPE) [141].

Contextual Embeddings To address the issue of polysemous and the context-dependent nature of words, we need distinguish the semantics of words in different contexts. Given a text x_1, x_2, \dots, x_T where each token $x_i \in \mathcal{V}$ is a word or sub-word, the contextual representation of x_i depends on the whole text.

$$[\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_T] = f_{\text{enc}}(x_1, x_2, \dots, x_T), \quad (1)$$

where $f_{\text{enc}}(\cdot)$ is neural encoder, which is described in Section 2.2, \mathbf{h}_i is called *contextual embedding* or *dynamical embedding* of token x_i because of the contextual information included in.

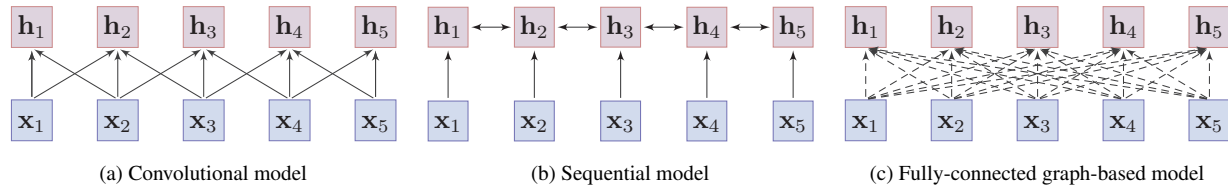


Figure 2: Neural Contextual Encoders

2.2 Neural Contextual Encoders

Most of the neural contextual encoders can be classified into three categories: convolutional models, sequential models, and graph-based models. Figure 2 illustrates the architecture of these models.

(1) *Convolutional models*. Convolutional models take the embeddings of words in the input sentence and capture the meaning of a word by aggregating the local information from its neighbors by convolution operations [80].

Convolutional models are usually easy to train and can capture the local contextual information.

(2) *Sequential models*. Sequential models usually adopt RNNs (such as LSTM [60] and GRU [21]) to capture the contextual representation of a word. In practice, bi-directional RNNs are used to collect information from both sides of a word, but its performance is often affected by the long-term dependency problem.

(3) *Graph-based models*. Different from the above models, graph-based models take the word as nodes and learn the contextual representation with a pre-defined linguistic structure between words, such as the syntactic structure [146, 161] or semantic relation [111].

Although the linguistic-aware graph structure can provide useful inductive bias, how to build a good graph structure is also a challenging problem. Besides, the structure depends heavily on expert knowledge or external NLP tools, such as the dependency parser.

In practice, a more straightforward way is to use a fully-connected graph to model the relation of every two words and let the model learn the structure by itself. Usually, the connection weights are dynamically computed by the self-attention mechanism, which implicitly indicates the connection between words.

A successful implementation of such an idea is the Transformer [171], which adopts the fully-connected self-attention architecture as well as other useful designs, such as positional embeddings, layer normalization, and residual connections.

Analysis Both convolutional and sequential models learn the contextual representation of the word with locality bias and are hard to capture the long-range interactions between

words. In contrast, Transformer can directly model the dependency between every two words in a sequence, which is more powerful and suitable to model the language.

However, due to its heavy structure and less model bias, Transformer usually requires a large training corpus and is easy to overfit on small or modestly-sized datasets [130, 49].

2.3 Why Pre-training?

With the development of deep learning, the number of model parameters has increased rapidly. The much larger dataset is needed to fully train model parameters and prevent overfitting. However, building large-scale labeled datasets is a great challenge for most NLP tasks due to the extremely expensive annotation costs, especially for syntax and semantically related tasks.

In contrast, large-scale unlabeled corpora are relatively easy to construct. To leverage the huge unlabeled text data, we can first learn a good representation from them and then use these representations for other tasks. Recent studies have demonstrated significant performance gains on many NLP tasks with the help of the representation extracted from the PTMs on the large unannotated corpora.

The advantages of pre-training can be summarized as follows:

1. Pre-training on the huge text corpus can learn universal language representations and help with the downstream tasks.
2. Pre-training provides a better model initialization, which usually leads to a better generalization performance and speeds up convergence on the target task.
3. Pre-training can be regarded as a kind of regularization to avoid overfitting on small data [39].

2.4 A Brief History of PTMs for NLP

Pre-training has always been an effective strategy to learn the parameters of deep neural networks, which are then fine-tuned on downstream tasks. As early as 2006, the breakthrough of deep learning came with greedy layer-wise unsupervised pre-training followed by supervised fine-tuning [58]. In CV, it

has been in practice to pre-train models on the huge ImageNet corpus, and then fine-tune further on smaller data for different tasks. This is much better than a random initialization because the model learns general image features which can then be used in various vision tasks.

In NLP, PTMs on large corpus have also been proved to be beneficial for the downstream NLP tasks, from the shallow word embedding to deep neural models.

2.4.1 Pre-trained word embeddings

Representing words as dense vectors has a long history [56]. The “modern” word embedding is introduced in pioneer work of neural network language model (NNLM) [11]. Collobert et al. [24] showed that the pre-trained word embedding on the unlabelled data can significantly improve many NLP tasks. To address the computational complexity, they learned word embeddings with *pairwise ranking* task instead of language modeling. Their work is the first attempt to obtain generic word embeddings useful for other tasks from unlabeled data. Mikolov et al. [116] showed that there is no need for deep neural networks to build good word embeddings. They propose two shallow architectures: Continuous Bag-of-Words (CBOW) and Skip-Gram (SG) models. Although the proposed models are simple and shallow, they can still learn the effective word embeddings capturing the latent syntactic and semantic similarities. Word2vec is one of the most popular implementations of these models and makes the pre-trained word embeddings accessible for different tasks in NLP. Besides, GloVe [120] is also a widely-used model for obtaining pre-trained word embeddings, which are computed by global word-word co-occurrence statistics from a corpus.

Although pre-trained word embeddings have been shown effective in NLP tasks, they are context-independent and mostly trained by shallow models. When used in a downstream task, the rest of the whole model still needs to be learned from scratch.

During the same time period, many researchers also try to learn embeddings of paragraph, sentence or document, such as paragraph vector [89], Skip-thought vectors [82], Context2Vec [114] and so on. Different from their modern successors, these sentence embedding models try to encode input sentences into a fixed-dimensional vector representation, rather than the contextual representation for each token.

2.4.2 Pre-trained contextual encoders

Since most NLP tasks are beyond word-level, it is natural to pre-train the neural encoders on sentence-level or higher. The output vectors of neural encoders are also called *contextual word embeddings* since they represent the word semantics

depending on its context.

McCann et al. [113] pre-trained a deep LSTM encoder from an attentional sequence-to-sequence model with machine translation (MT). The context vectors (CoVe) output by the pretrained encoder can improve the performance of a wide variety of common NLP tasks. Peters et al. [122] pre-trained 2-layer LSTM encoder with a bidirectional language model (BiLM), consisting of a forward LM and a backward LM. The contextual representations output by the pre-trained BiLM, ELMo (Embeddings from Language Models), are shown to bring large improvements on a broad range of NLP tasks. Akbik et al. [1] captured word meaning with contextual string embeddings pre-trained with character-level LM.

However, these PTMs are usually used as a feature extractor to produce the contextual word embeddings, which are fed into the main model for downstream tasks. Their parameters are fixed and the rest parameters of the main model are still trained from scratch.

Ramachandran et al. [134] found the seq2seq models can be significantly improved by unsupervised pre-training. The weights of both encoder and decoder are initialized with pre-trained weights of two language models and then fine-tuned with labeled data. ULMFiT (Universal Language Model Fine-tuning) [62] attempted to fine-tune pre-trained LM for text classification (TC) and achieved state-of-the-art results on six widely-used TC datasets. ULMFiT consists of 3 phases: 1) pre-training LM on general-domain data; 2) fine-tuning LM on target data; 3) fine-tuning on the target task. ULMFiT also investigates some effective fine-tuning strategies, including discriminative fine-tuning, slanted triangular learning rates, and gradual unfreezing. Since ULMFiT, fine-tuning has become the mainstream approach to adapt PTMs for the downstream tasks.

More recently, the very deep PTMs have shown their powerful ability in learning universal language representations: e.g., OpenAI GPT (Generative Pre-training) [130] and BERT (Bidirectional Encoder Representation from Transformer) [32]. Besides LM, an increasing number of self-supervised tasks (see Section 3.1) are proposed to make the PTMs capturing more knowledge from large scale text corpora.

3 Overview of PTMs

The major differences between PTMs are the usages of contextual encoders, pre-training tasks, and purposes. We have briefly introduced the architectures of contextual encoders in Section 2.2. In this section, we focus on the description of pre-training tasks and give a taxonomy of PTMs.

3.1 Pre-training Tasks

The pre-training tasks are crucial for learning the universal representation of language. Usually, these pre-training tasks should be challenging and have substantial training data. In this section, we summarize the pre-training tasks into three categories¹⁾: supervised learning, unsupervised learning, and self-supervised learning.

1. *Supervised learning* is to learn a function that maps an input to an output based on training data consisting of input-output pairs.
2. *Unsupervised learning* is to find some intrinsic knowledge, such as clusters, densities, latent representations, from unlabeled data.
3. *Self-Supervised learning* (SSL) is a blend of supervised learning and unsupervised learning. The key idea of SSL is to predict any part of the input from other parts in some form. For example, the masked language model (MLM) is a self-supervised task that attempts to predict the masked words in a sentence given the rest words.

In CV, many PTMs are trained on large supervised training sets like ImageNet. However, in NLP field, the datasets of most supervised tasks are not large enough to train a good PTM. The only exception is machine translation (MT). A large-scale MT dataset, WMT 2017, consists of more than 7 million sentence pairs. Besides, MT is one of the most challenging tasks in NLP, and an encoder pretrained on MT can benefit a variety of downstream NLP tasks. As a successful PTM, CoVe [113] is an encoder pretrained on MT task and improves a wide variety of common NLP tasks: sentiment analysis (SST, IMDB), question classification (TREC), entailment (SNLI), and question answering (SQuAD).

The pre-training tasks widely-used in existing PTMs are listed as follows:

3.1.1 Language Modeling (LM)

The most common unsupervised task in NLP is probabilistic language modeling (LM), which is a classic probabilistic density estimation problem. Although LM is a general concept, in practice, LM often refers in particular to auto-regressive LM or unidirectional LM.

Given a text sequence $x_{1:T} = [x_1, x_2, \dots, x_T]$, its joint probability $p(x_{1:T})$ can be decomposed as

$$p(x_{1:T}) = \prod_{t=1}^T p(x_t | x_{0:t-1}), \quad (2)$$

¹⁾ Indeed, it is hard to clearly distinguish the unsupervised learning and self-supervised learning. For clarification, we refer “unsupervised learning” to the learning without human-annotated supervised labels”.

where x_0 is special token indicating the begin of sequence.

The conditional probability $p(x_t | x_{0:t-1})$ can be modeled by a probability distribution over the vocabulary given linguistic context $x_{0:t-1}$. The context $x_{0:t-1}$ is modeled by neural encoder $f_{\text{enc}}(\cdot)$, and the conditional probability is

$$p(x_t | x_{0:t-1}) = g_{\text{LM}}(f_{\text{enc}}(x_{0:t-1})), \quad (3)$$

where $g_{\text{LM}}(\cdot)$ is prediction layer.

Given a huge corpus, we can train the entire network with a maximum-likelihood estimation (MLE).

A drawback of unidirectional LM is that the representation of each token encodes only the leftward context tokens and itself. However, better contextual representations of text should encode contextual information from both directions. An improved solution is bidirectional LM (BiLM), which consists of two unidirectional LMs: a forward left-to-right LM and a backward right-to-left LM.

For BiLM, Baevski et al. [5] proposed a two-tower model that the forward tower operates the left-to-right LM and the backward tower operates the right-to-left LM.

3.1.2 Masked Language Modeling (MLM)

Masked language modeling (MLM) is first proposed by Taylor [165] in the literature, who referred this as a Cloze task. Devlin et al. [32] adapted this task as a novel pre-training task to overcome the drawback of the standard unidirectional LM. Loosely speaking, MLM first masks out some tokens from the input sentences and then trains the model to predict the masked tokens by the rest of tokens. However, this pre-training method will create a mismatch between the pre-training phase and the fine-tuning phase, because the mask token does not appear during the fine-tuning phase. Empirically, to deal with this issue, Devlin et al. [32] used a special [MASK] token 80% of the time, a random token 10% of the time and the original token 10% of the time to perform masking.

Sequence-to-Sequence MLM (Seq2Seq MLM) MLM is usually solved as classification problem. We feed the masked sequences to a neural encoder, whose output vectors are further fed into a softmax classifier to predict the masked token. Alternatively, we can use encoder-decoder (aka. sequence-to-sequence) architecture for MLM, in which the encoder is fed a masked sequence and the decoder sequentially produces the masked tokens in auto-regression fashion. We refer to this kind of MLM as sequence-to-sequence MLM (Seq2Seq MLM), which is used in MASS [147] and T5 [132]. Seq2Seq MLM can benefit the Seq2Seq-style downstream tasks, such as question answering, summarization and machine translation.

Enhanced Masked Language Modeling (E-MLM) Concurrently, there are multiple research proposing different enhanced versions of MLM to further improve on BERT. Instead of static masking, RoBERTa [105] improves BERT by dynamic masking.

UniLM [35, 7] extends the task of mask prediction on three types of language modeling tasks: unidirectional, bidirectional, and sequence-to-sequence prediction.

XLM [25] performs MLM on a concatenation of parallel bilingual sentence pairs, called translation language modeling (TLM).

To integrate structure information into pre-training, SpanBERT [72] replaces MLM with random contiguous words masking and Span Boundary Objective (SBO), which requires the system to predict masked spans based on span boundaries. Besides, StructBERT [180] introduces the span order recovery task to further incorporating language structures.

Another way to enrich MLM is to incorporate external knowledge. ERNIE(Baidu) [157, 158] proposes to mask phrases and entities instead of random spans. Knowledge is also introduced by incorporating entity embeddings as inputs, such as E-BERT [128], and ERNIE(THU) [199].

3.1.3 Permuted Language Modeling (PLM)

Despite the wide use of the MLM task in pre-training, Yang et al. [194] claimed that some special tokens used in pre-training of MLM, like [MASK], are absent when the model is applied on downstream tasks, leading to a gap between pre-training and fine-tuning. To overcome this issue, Permuted Language Modeling (PLM) [194] is a pre-training objective to replace MLM. In short, PLM is a language modeling task on a random permutation of input sequences. Given a sequence, a permutation is randomly sampled from all possible permutations. Then some of the tokens in the permuted sequence are chosen as the target and the model is trained to predict these targets, depending on the rest of tokens and the natural positions of targets. Note that this permutation does not affect the natural positions of sequences and only defines the order of token predictions. In practice, only the last few tokens in the permuted sequences are predicted, due to the slow convergence. And a special two-stream self-attention is introduced for target-aware representations.

3.1.4 Denoising Autoencoder (DAE)

Denoising autoencoder (DAE) takes a partially corrupted input and aims to recover the original undistorted input. Specific to language, a sequence-to-sequence model, such as the standard Transformer, is used to reconstruct the original text. There are several ways to corrupt text [93]:

(1) *Token Masking*: Randomly sampling tokens from the input and replacing them with [MASK] elements.

(2) *Token Deletion*: Randomly deleting tokens from the input. Different to token masking, the model need decide the positions of missing inputs.

(3) *Text Infilling*: Like SpanBERT, a number of text spans are sampled and replaced with a single [MASK] token. Each span length is drawn from a Poisson distribution ($\lambda = 3$). The model need predict how many tokens are missing from a span.

(4) *Sentence Permutation*: Dividing a document into sentences based on full stops, and shuffling these sentences in a random order.

(5) *Document Rotation*: Selecting a token uniformly at random, and rotating the document so that it begins with that token. The model need identify the real start position of the document.

3.1.5 Contrastive Learning (CTL)

Contrastive learning [140] assumes some observed pairs of text that are more semantically similar than randomly sampled text. A score function $s(x, y)$ for text pair (x, y) is learned to minimize the objective function:

$$\mathbb{E}_{x, y^+, y^-} \left[-\log \frac{\exp(s(x, y^+))}{\exp(s(x, y^+)) + \exp(s(x, y^-))} \right], \quad (4)$$

where (x, y^+) are a similar pair and y^- is presumably dissimilar to x . y^+ and y^- are typically called positive and negative sample. The score function $s(x, y)$ is often computed by a learnable neural encoder in two ways: $s(x, y) = f_{\text{enc}(x)}^T f_{\text{enc}(y)}$ or $s(x, y) = f_{\text{enc}}(x \oplus y)$.

The idea behind CTL is “learning by comparison”. Compared to LM, CTL usually has less computational complexity, and therefore is desirable alternative training criteria for PTMs.

Collobert et al. [24] proposed *pairwise ranking* task to distinguish real and fake phrases. The model need to predict a higher score for legal phrase than an incorrect phrase obtained by replacing its central word with a random words. Mnih and Kavukcuoglu [118] trained word embeddings efficiently with Noise-Contrastive Estimation (NCE) [51], which trains a binary classifier to distinguish real and fake samples. The idea of NCE is also used in the well-known word2vec embedding [116].

We briefly describes some recently proposed CTL tasks in the following paragraphs.

Deep InfoMax (DIM) Deep InfoMax (DIM) [59] is originally proposed for images, which improves the quality of the representation by maximizing the mutual information between an image representation and local regions of the image.

Kong et al. [83] applied DIM to language representation learning. The global representation of a sequence x is defined to be the hidden state of the first token (assumed to be a special start of sentence symbol) output by contextual encoder $f_{\text{enc}}(x)$. The objective of DIM is to assign a higher score for $f_{\text{enc}}(x_{i:j})^T f_{\text{enc}}(\hat{x}_{i:j})$ than $f_{\text{enc}}(\tilde{x}_{i:j})^T f_{\text{enc}}(\hat{x}_{i:j})$, where $x_{i:j}$ denotes an n -gram²⁾ span from i to j in x , $\hat{x}_{i:j}$ denotes a sentence masked at position i to j , and $\tilde{x}_{i:j}$ denotes a randomly-sampled negative n -gram from corpus.

Replaced Token Detection (RTD) Replaced Token Detection (RTD) is the same as NCE but predicts whether a token is replaced given its surrounding context.

CBOW with negative sampling (CBOW-NS) [116] can be viewed as a simple version of RTD, in which the negative samples are randomly sampled from vocabulary with simple proposal distribution.

ELECTRA [22] improves RTD by utilizing a generator to replacing some tokens of a sequence. A generator G and a discriminator D are trained following a two-stage procedure: (1) Train only the generator with MLM task for n_1 steps; (2) Initialize the weights of the discriminator with the weights of the generator. Then train the discriminator with a discriminative task for n_2 steps, keeping G frozen. Here the discriminative task indicates justifying whether the input token has been replaced by G or not. The generator is thrown after pre-training, and only the discriminator will be fine-tuned on downstream tasks.

RTD is also an alternative solution for the mismatch problem the network sees [MASK] during pre-training but not when being fine-tuned in downstream tasks.

Similarly, WKLM [188] replaces words on the entity-level instead of token-level. Concretely, WKLM replaces entity mentions with names of other entities of the same type and train the models to distinguish whether the entity has been replaced.

Next Sentence Prediction (NSP) Punctuations are the natural separators of text data. So, it is reasonable to construct pre-training methods by utilizing them. Next Sentence Prediction (NSP) [32] is just a great example for this. As its name suggests, NSP trains the model to distinguish whether two input sentences are continuous segments from the training corpus. Specifically, when choosing the sentences pair for each pre-training example, 50% of the time the second sentence is the actual next sentence of the first one, and 50% of the time it is a random sentence from the corpus. By doing so, it is capable to teach the model to understand the relationship between two input sentences and thus benefit downstream tasks that are sensitive to this information, such as Question

Answering and Natural Language Inference.

However, the necessity of the NSP task has been questioned by subsequent work [72, 194, 105, 86]. Yang et al. [194] found the impact of the NSP task unreliable, while Joshi et al. [72] found that single-sentence training without the NSP loss is superior to sentence-pair training with the NSP loss. Moreover, Liu et al. [105] conducted further analysis for the NSP task, which shows that when training with blocks of text from a single document, removing the NSP loss matches or slightly improves performance on downstream tasks.

Sentence Order Prediction (SOP) To better model inter-sentence coherence, ALBERT [86] replaces the NSP loss with a sentence order prediction (SOP) loss. As conjectured in Lan et al. [86], NSP conflates topic prediction and coherence prediction in a single task. However, topic prediction is easier to learn compared to coherence prediction, which allows the model to make predictions merely rely on topic learning. Different to NSP, SOP uses two consecutive segments from the same document as positive examples, and the same two consecutive segments but with their order swapped as negative examples. As a result, ALBERT consistently outperforms BERT on various downstream tasks.

StructBERT [180] and BERTje [29] also take SOP as their self-supervised learning task.

3.1.6 Others

Apart from the above tasks, there are many other tasks designated for specific tasks, such as sentiment label-aware MLM for sentiment analysis [78], gap sentence generation (GSG) for text summarization [197], disfluency detection [179] and so on. In addition, some auxiliary pre-training tasks are designed to incorporate factual knowledge, such as denoising entity auto-encoding (dEA) in ERNIE(THU) [199], entity linking (EL) in KnowBERT [123].

Furthermore, several tasks are introduced to obtain multi-modal pre-trained model. Typically, tasks like visual-based MLM, masked visual-feature modeling and visual-linguistic matching are widely used in multi-modal pre-training, such as VideoBERT [152], VisualBERT [95], ViLBERT [107] and so on.

3.2 Taxonomy of PTMs

To clarify the relations of existing PTMs for NLP, we build the taxonomy of PTMs, which categorizes existing PTMs from different perspective: (1) the type of word representation used by PTMs, (2) the backbone network used by PTMs, (3) the type of pre-training tasks used by PTMs, and (4) the PTMs designed for specific types of scenarios or inputs. Figure 3 shows

²⁾ n is drawn from a Gaussian distribution $\mathcal{N}(5, 1)$ clipped at 1 (minimum length) and 10 (maximum length).

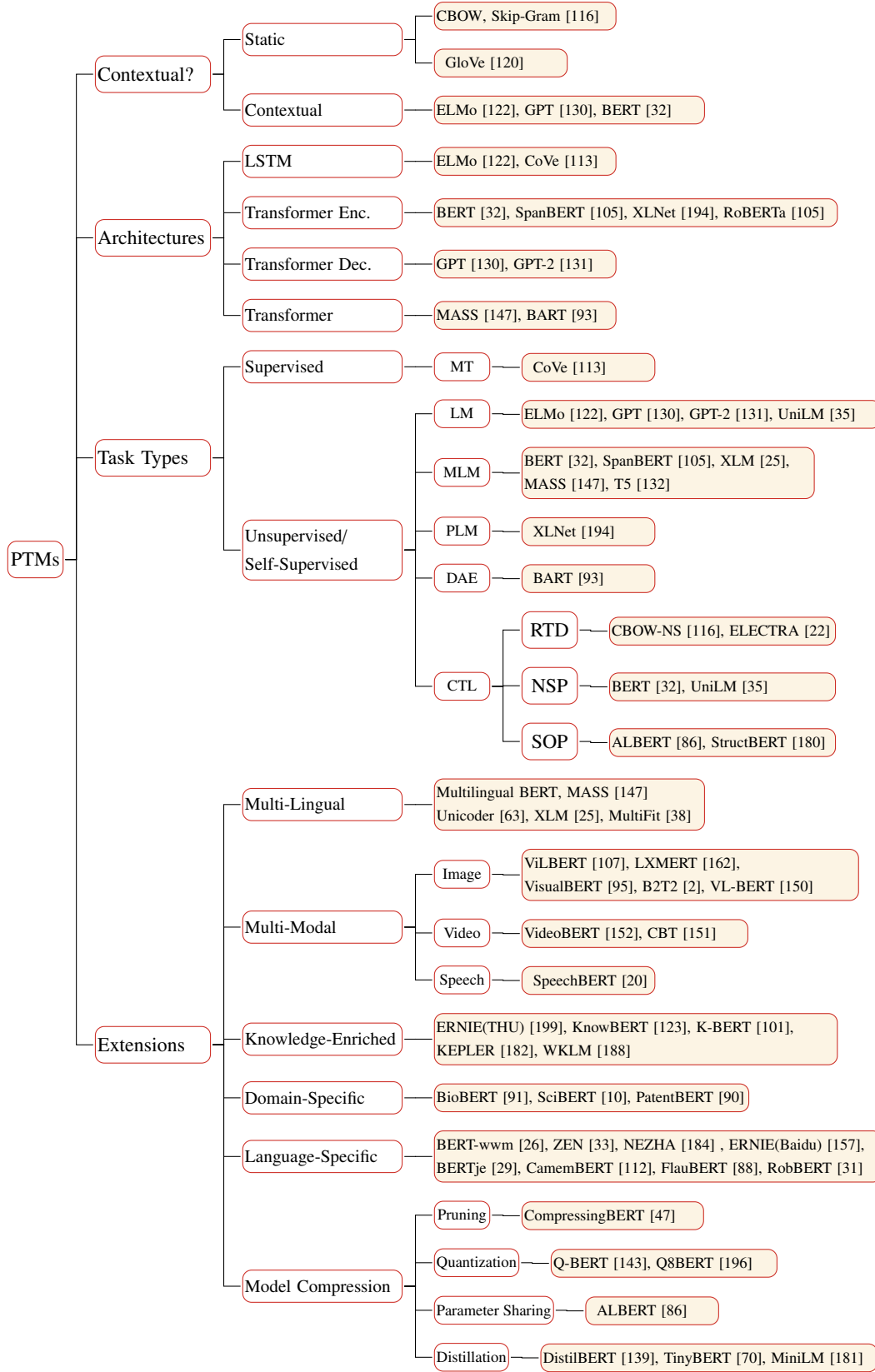


Figure 3: Taxonomy of PTMs with Representative Examples

Table 1: List of Representative PTMs

PTMs	Architecture [†]	Input	Pre-Training Task	Corpus	Params	GLUE [‡]	FT? [‡]
ELMo [122]	LSTM	Text	BiLM	WikiText-103			No
GPT [130]	Transformer Dec.	Text	LM	BookCorpus	117M	72.8	Yes
GPT-2 [131]	Transformer Dec.	Text	LM	WebText	117M ~ 1542M		No
BERT [32]	Transformer Enc.	Text	MLM & NSP	WikiEn+BookCorpus	110M ~ 340M	81.9*	Yes
InfoWord [83]	Transformer Enc.	Text	DIM+MLM	WikiEn+BookCorpus	=BERT	81.1*	Yes
RoBERTa [105]	Transformer Enc.	Text	MLM	BookCorpus+CC-News+OpenWebText+ STORIES	355M	88.5	Yes
XLNet [194]	Two-Stream Transformer Enc.	Text	PLM	WikiEn+ BookCorpus+Giga5+ClueWeb+Common Crawl	≈BERT	90.5 [§]	Yes
ELECTRA [22]	Transformer Enc.	Text	RTD+MLM	same to XL-Net	335M	88.6	Yes
UniLM [35]	Transformer Enc.	Text	MLM [‡] NSP	WikiEn+BookCorpus	340M	80.8	Yes
MASS [147]	Transformer	Text	Seq2Seq MLM	*Task-dependent			Yes
BART [93]	Transformer	Text	DAE	same to RoBERTa	110% of BERT	88.4*	Yes
T5 [132]	Transformer	Text	Seq2Seq MLM	Colossal Clean Crawled Corpus (C4)	220M ~ 11B	89.7*	Yes
ERNIE(THU) [199]	Transformer Enc.	Text+Entities	MLM+NSP+dEA	WikiEn+ Wikidata	114M	79.6	Yes
KnowBERT [123]	Transformer Enc.	Text	MLM+NSP+EL	WikiEn+WordNet/Wiki	253M ~ 523M		Yes
K-BERT [101]	Transformer Enc.	Text+Triples	MLM+NSP	WikiZh+WebtextZh+ CN-DBpedia+HowNet+MedicalKG	=BERT		Yes
KEPLER [182]	Transformer Enc.	Text	MLM+KE	WikiEn+Wikidata			Yes
WKLM [188]	Transformer Enc.	Text	MLM+ERD	WikiEn+Wikidata	=BERT		Yes

[†] “Transformer Enc.” and “Transformer Dec.” mean the encoder and decoder part of the standard Transformer architecture respectively. Their difference is that the decoder part uses masked self-attention with triangular matrix to prevent tokens from attending their future (right) positions. “Transformer” means the standard encoder-decoder architecture.

[‡] the averaged score on 9 tasks of GLUE benchmark (see Section 7.1).

* without WNLI task.

[§] indicates ensemble result.

[‡] means whether is model usually used in fine-tuning fashion.

◊ The MLM of UniLM is built on three versions of LMs: Unidirectional LM, Bidirectional LM, and Sequence-to-Sequence LM.

the taxonomy as well as some corresponding representative PTMs.

Besides, Table 1 distinguishes some representative PTMs in detail.

3.3 Model Analysis

Due to the great success of PTMs, it is important to understand what kinds of knowledge are captured by them, and how to induce knowledge from them. There is a wide range of literature analyzing linguistic knowledge and world knowledge stored in pre-trained non-contextual and contextual embeddings.

3.3.1 Non-Contextual Embeddings

Static word embeddings are first probed for kinds of knowledge. Mikolov et al. [117] found that word representations learned by neural network language models are able to capture linguistic regularities in language, and the relationship between words can be characterized by a relation-specific vector offset. Further analogy experiments [116] demonstrated that word vectors produced by skip-gram model can capture both syntactic and semantic word relationships, such as $\text{vec}(\text{“China”}) - \text{vec}(\text{“Beijing”}) \approx \text{vec}(\text{“Japan”}) - \text{vec}(\text{“Tokyo”})$. Besides, they find compositionality property of word vectors, for example, $\text{vec}(\text{“Germany”}) + \text{vec}(\text{“capital”})$ is close to $\text{vec}(\text{“Berlin”})$. Inspired by these work, Rubinstein

et al. [138] found that distributional word representations are good at predicting taxonomic properties (e.g., dog is an animal) but fail to learn attributive properties (e.g., swan is white). Similarly, Gupta et al. [50] showed that word2vec embeddings implicitly encode referential attributes of entities. The distributed word vectors along with a simple supervised model can learn to predict numeric and binary attributes of entities with a reasonable degree of accuracy.

3.3.2 Contextual Embeddings

A large number of studies have probed and induced different types of knowledge in contextual embeddings. In general, there are two types of knowledge: linguistic knowledge and world knowledge.

Linguistic Knowledge Tenney et al. [167], Liu et al. [99] designed a wide range of probing tasks for PTMs and found that BERT performs well on many syntactic tasks such as part-of-speech tagging and constituent labeling. However, BERT is not good enough at semantic and fine-grained syntactic tasks, compared with simple syntactic probing tasks. Furthermore, knowledge of subject-verb agreement [46] and semantic roles [40] are also confirmed to exist in BERT. Besides, Hewitt and Manning [55], Jawahar et al. [67], Kim et al. [79] proposed several methods to extract dependency trees and constituency trees from BERT, which proved the BERT’s ability

to encode sentence structure. Reif et al. [136] explored the geometry of internal representations in BERT and find some evidences: 1) linguistic features seem to be represented in separate semantic and syntactic subspaces; 2) attention matrices contain grammatical representations; 3) BERT distinguishes word senses at a very fine level.

World Knowledge Besides linguistic knowledge, PTMs may also store world knowledge presented in the training data. A straightforward method of probing world knowledge is to query BERT with “fill-in-the-blank” cloze statements, for example, “Dante was born in [MASK]”. Petroni et al. [125] constructed LAMA (Language Model Analysis) task by manually creating single-token cloze statements (queries) from several knowledge sources. Their experiments show that BERT contains world knowledge competitive with traditional information extraction methods. Since the simplicity of query generation procedure in LAMA, Jiang et al. [69] argued that LAMA is just measuring a lower bound for what language models know and propose more advanced methods to generate more efficient queries. Despite the surprising findings of LAMA, it has also been questioned by subsequent work [129, 77]. Similarly, several studies induce relational knowledge [14] and commonsense knowledge [28] from BERT for downstream tasks.

4 Extensions of PTMs

4.1 Knowledge-Enriched PTMs

PTMs usually learn universal language representation from general-purpose large-scale text corpora, but lack domain-specific knowledge. Incorporating domain knowledge from external knowledge bases into PTM has been shown to be effective. The external knowledge ranges from linguistic [87, 78, 123, 178], semantic [92], commonsense [48], factual [199, 123, 101, 188, 182], to domain-specific knowledge [54].

On the one hand, external knowledge can be injected during pre-training. Early studies [183, 202, 187, 190] focused on learning knowledge graph embeddings and word embedding jointly. Since BERT, some auxiliary pre-training tasks are designed to incorporate external knowledge into deep PTMs. LIBERT [87] (linguistically-informed BERT) incorporates linguistic knowledge via an additional linguistic constraints task. Ke et al. [78] integrated sentiment polarity of each word to extend the MLM to Label-Aware MLM (LA-MLM). As a result, their proposed model, SentiLR, achieves state-of-the-art performance on several sentence- and aspect-level sentiment classification tasks. Levine et al. [92] proposed SenseBERT, which is pre-trained to predict not only the masked tokens but

also their supersenses in WordNet. ERNIE(THU) [199] integrates entity embeddings pre-trained on a knowledge graph with corresponding entity mentions in the text to enhance the text representation. Similarly, KnowBERT [123] trains BERT jointly with an entity linking model to incorporate entity representation in an end-to-end fashion. Wang et al. [182] proposed KEPLER, which jointly optimizes knowledge embedding and language modeling objectives. These work inject structure information of knowledge graph via entity embedding. In contrast, K-BERT [101] explicitly injects related triples extracted from KG into the sentence to obtain an extended tree-form input for BERT. Moreover, Xiong et al. [188] adopted entity replacement identification to encourage the model to be more aware of factual knowledge. However, most of these methods update the parameters of PTMs when injecting knowledge, which may suffer from catastrophic forgetting when injecting multiple kinds of knowledge. To address this, K-Adapter [178] injects multiple kinds of knowledge by training different adapters independently for different pre-training tasks, which allows continual knowledge infusion.

On the other hand, one can incorporate external knowledge into pre-trained models without retraining them from scratch. As an example, K-BERT [101] allows injecting factual knowledge during fine-tuning on downstream tasks. Guan et al. [48] employed commonsense knowledge bases, ConceptNet and ATOMIC, to enhance GPT-2 for story generation. Yang et al. [192] proposed a knowledge-text fusion model to acquire related linguistic and factual knowledge for machine reading comprehension.

Besides, Logan IV et al. [106] and Hayashi et al. [53] extended language model to knowledge graph language model (KGLM) and latent relation language model (LRLM) respectively, both of which allow prediction conditioned on knowledge graph. These novel KG-conditioned language models show potential for pre-training.

4.2 Multi-Modal PTMs

Observing the success of PTMs across many NLP tasks, some research has focused on obtaining a cross-modal version of PTMs. A great majority of these models are designed for a general visual and linguistic feature encoding. And these models are pre-trained on some huge corpus of cross-modal data, such as videos with spoken words or images with captions, incorporating extended pre-training tasks to fully utilize the multi-modal feature. VideoBERT [152] and CBT [151] are joint video and text models. To obtain sequences of visual and linguistic tokens used for pre-training, the videos are pre-processed by CNN-based encoders and off-the-shelf speech recognition techniques, respectively. And a single Transformer encoder is trained on the processed data to learn

the vision-language representations for downstream tasks like video caption. Furthermore, UniViLM [109] proposes to bring in generation tasks to further pre-train the decoder using in downstream tasks.

Besides methods for video-language pre-training, several works introduce PTMs on image-text pairs, aiming to fit downstream tasks like visual question answering(VQA) and visual commonsense reasoning(VCR). Several proposed models adopt two separate encoders for image and text representation independently, such as ViLBERT [107] and LXMERT [162]. While other methods like VisualBERT [95], B2T2 [2], VL-BERT [150], Unicoder-VL [94] and UNITER [16] propose single-stream unified Transformer. Though these model architectures are different, similar pre-training tasks, such as MLM and image-text matching, are introduced in these approaches. And to better exploit visual elements, images are converted into sequences of regions by applying RoI or bounding box retrieval techniques before encoded by pre-trained Transformers.

Moreover, several methods have explored the chance of PTMs on audio-text pairs, such as SpeechBERT [20]. This work tries to build an end-to-end Speech Question Answering(SQA) model by encoding audio and text with a single Transformer encoder, which is pre-trained with MLM on speech and text corpus and fine-tuned on Question Answering.

4.3 Model Compression

Since the pre-trained language models usually consist of at least hundreds of millions of parameters, they are difficult to be deployed on the on-line service in real-life applications and on resource-restricted devices. Model compression [15] is a potential approach to reduce the model size and increase computation efficiency.

There are four common ways to compress PTMs [42]: (1) pruning, which removes less important parameters, (2) weight quantization [36], which uses fewer bits to represent the parameters, (3) parameter sharing across similar model units, and (4) knowledge distillation [57], which trains a smaller student model that learns from intermediate outputs from the original model. Table 3 distinguishes some representative compressed PTMs in detail.

4.3.1 Model Pruning

Model pruning refers to removing part of neural network (e.g., weights, neurons, layers, channels, attention heads, etc.), thereby achieving the effects of reducing the model size and speeding up inference time.

Gordon et al. [47] explored the timing of pruning (e.g., pruning during pre-training, after downstream fine-tuning) and the

pruning regimes. Li and Eisner [96] compressed ELMo word token embeddings using variational information bottleneck. Michel et al. [115] and Voita et al. [174] tried to prune the entire self-attention heads in the transformer block.

4.3.2 Quantization

Quantization refers to the compression of higher precision parameters to lower precision. Works from Shen et al. [143] and Zafrir et al. [196] solely focus on this area. Note that quantization often requires compatible hardware.

4.3.3 Parameter Sharing

Another well-known approach to reduce the number of parameters is parameter sharing, which is widely used in CNNs, RNNs and Transformer [30]. ALBERT [86] uses *cross-layer parameter sharing* and *factorized embedding parameterization* to reduce the parameters of PTMs.

4.3.4 Knowledge Distillation

Knowledge distillation (KD) [57] is a compression technique in which a small model called *student model* is trained to reproduce the behaviors of a large model called *teacher model*. Here the teacher model can be an ensemble of many models and usually well pre-trained. Different to model compression, distillation techniques learn a small student model from a fixed teacher model through some optimization objectives, while compression techniques aiming at searching a sparser architecture.

Generally, distillation mechanisms can be divided into three types: distillation from soft target probabilities, distillation from other knowledge, and distillation to other structures:

(1) *Distillation from soft target probabilities*. Bucilua et al. [15] showed that making the student approximate the teacher model can transfer knowledge from teacher to student. A common method is approximating the logits of the teacher model. DistilBERT [139] trained the student model with a distillation loss over the soft target probabilities of the teacher as:

$$\mathcal{L}_{CE} = \sum_i t_i * \log(s_i), \quad (5)$$

where t_i and s_i are the probabilities estimated by the teacher model and the student respectively.

Distillation from soft target probabilities can also be used in task-specific models, such as information retrieval [108], and sequence labeling [168].

(2) *Distillation from other knowledge*. Distillation from soft target probabilities regards the teacher model as a black box and only focus on its outputs. Moreover, decomposing

Table 2: Comparison of Compressed PTMs

Method	Type	#Layer	Loss Function*	Speed Up	Params	Teacher	GLUE [‡]
BERT _{BASE} [32]	Baseline	12	CE _{MLM} + CE _{NSP}		110M		79.6
BERT _{LARGE} [32]		24	CE _{MLM} + CE _{NSP}		340M		81.9
Q-BERT [143]	Quantization	12	HAWQ + GWQ	-		BERT _{BASE}	≈ 99% BERT [°]
Q8BERT [196]		12	DQ + QAT	-		BERT _{BASE}	≈ 99% BERT
ALBERT [§] [86]	Param. Sharing	12	CE _{MLM} + CE _{SOP}	×5.6 ~ 0.3	12 ~ 235M		89.4 (ensemble)
DistilBERT [139]	Distillation	6	CE _{KD} + CE _{cosKD} + CE _{MLM}	×1.63	66M	BERT _{BASE}	77.0 (dev)
TinyBERT [†] [70]		4	MSE _{embed} + MSE _{attn} + MSE _{hidn} + CE _{KD}	×9.4	14.5M	BERT _{BASE}	76.5
BERT-PKD [156]		3 ~ 6	CE _{KD} + PT _{KD} + CE _{TASK}	×3.73 ~ 1.64	45.7 ~ 67 M	BERT _{BASE}	76.0 ~ 80.6 [#]
PD [170]		6	CE _{KD} + CE _{TASK} + CE _{MLM}	×2.0	67.5M	BERT _{BASE}	81.2 [#]
MobileBERT [‡] [159]		24	FMT + AT + PKT + CE _{KD} + CE _{MLM}	×4.0	25.3M	BERT _{LARGE}	79.7
MiniLM [181]		6	AT + AR	×1.99	66M	BERT _{BASE}	81.0 [°]
DualTrain [§] [201]		12	Dual Projection + CE _{MLM}	-	1.8 ~ 19.2M	BERT _{BASE}	75.8 ~ 81.9 [‡]

[‡] the averaged score on 8 tasks (without WNLI) of GLUE benchmark (see Section 7.1). Here MNLI-m and MNLI-mm are regarded as two different tasks. ‘dev’ indicates the result is on dev set. ‘ensemble’ indicates the result is from the ensemble model.

* ‘MLM’, ‘NSP’, and ‘SOP’ indicate pre-training objective (see Section 3.1). ‘HAWQ’, ‘GWQ’, ‘DQ’, and ‘QAT’ indicate Hessian AWARE Quantization, Group-wise Quantization, Quantization-Aware Training, and Dynamically Quantized, respectively. ‘KD’ means knowledge distillation. ‘TASK’ means task-specific loss. ‘FMT’, ‘AT’, and ‘PKT’ mean Feature Map Transfer, Attention Transfer, and Progressive Knowledge Transfer, respectively. ‘AR’ means Self-Attention value relation.

[§] The dimensionality of the hidden or embedding layers is reduced.

[†] used a smaller vocabulary.

[‡] Generally, the F1 score is usually used as the main metric of QQP task. But MiniLM reports the accuracy, which is incomparable to other works.

[°] Result on MNLI and SST-2 only.

[#] Result on the other tasks except STS-B and CoLA.

[‡] Result on MRPC, MNLI, and SST-2 only.

the teacher model and distilling more knowledge can bring improvement to the student model.

We summarize some representative models as Table 3.

Table 3: Distillation with other knowledge

Models	Teacher	Distilled Knowledge
TinyBERT [70]	BERT _{BASE}	Layer-to-Layer distillation with embedding outputs, hidden states, and self-attention distributions.
MobileBERT [159]	IB-BERT _{LARGE}	Layer-to-Layer distillation with soft target probabilities, hidden states, and self-attention distributions.
MiniLM [181]	BERT _{BASE}	Self-attention distributions and self-attention value relation.

Besides these representative models, other models distill knowledge through many approaches. Sun et al. [156] introduced a “*patient*” teacher-student mechanism, Liu et al. [102] explored the used of KD to improve a pre-trained multi-task deep neural network.

(3) *Distillation to other structures*. Generally, the structure of the student model is the same as the teacher model, except for a smaller layer size and a smaller hidden size. However, not only decreasing parameters but also simplifying model structures from Transformer to RNN [164] or CNN [18] can reduce the computational complexity.

Table 2 gives a comparison of some representative compressed PTMs.

3) <https://github.com/google-research/bert/blob/master/multilingual.md>

4.4 Domain-Specific PTMs

Most publicly available PTMs are trained on general domain corpus such as Wikipedia, which limits their applications to specific domains. Recently, some studies have proposed PTMs trained on specialty corpora, such as BioBERT [91] for biomedical text, SciBERT [10] for scientific text, ClinicalBERT [64, 3] for clinical text. In addition to pre-training a domain-specific PTM, some work attempts to adapt available pre-trained models to target applications, such as biomedical entity normalization [68], patent classification [90], progress notes classification and keyword extraction [163].

4.5 Multilingual and Language-Specific PTMs

Learning multilingual text representations shared across languages plays an important role for many cross-lingual NLP tasks. Most of early works focus on learning multilingual word embedding [41, 110, 145], which represents text from multiple languages in a single semantic space. However these methods usually need (weak) alignment between languages.

Multilingual BERT³⁾ (M-BERT) is pre-trained by MLM with the shared vocabulary and weights on Wikipedia text from the top 104 languages. Each training sample is a monolingual document, and there are no cross-lingual objectives specifically designed nor any cross-lingual data. Even so, M-BERT performs cross-lingual generalization surprisingly well [127]. K et al. [74] showed that the lexical overlap between languages plays a negligible role in the cross-lingual

success. MASS [147] pretrained on multiple languages also achieves significant improvement for unsupervised NMT.

XLM [25] improves M-BERT by incorporating a cross-lingual task, translation language modeling (TLM), which performs MLM on a concatenation of parallel bilingual sentence pairs. Unicoder [63] further propose three new cross-lingual pre-training tasks, including cross-lingual word recovery, cross-lingual paraphrase classification and cross-lingual masked language model.

Although multilingual PTMs perform well on many languages, recent work showed that PTMs trained on a single language significantly outperform the multilingual results [112, 88, 173].

For Chinese, which does not have explicit word boundaries, modeling larger granularity [26, 33, 184] and multi-granularity [157, 158] word representations have shown great success. Kuratov and Arkhipov [85] used transfer learning techniques to adapt a multilingual PTM to a monolingual PTM for Russian language. In addition, some monolingual PTMs pre-trained on language-specific corpus for French [112, 88], Finnish [173], Dutch [29, 31] have been released.

5 Adapting PTMs to Downstream Tasks

Although PTMs capture the general language knowledge from large corpus, how effectively adapting their knowledge to the downstream task is still a key problem.

5.1 Transfer Learning

Transfer learning [119] is to adapt the knowledge from a source task (or domain) to a target task (or domain). Figure 4 gives an illustration of transfer learning.

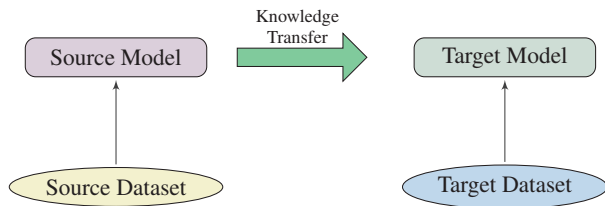


Figure 4: Transfer Learning

There are many types of transfer learning in NLP, such as domain adaptation, cross-lingual learning, multi-task learning, and so on. Adapting PTMs to downstream tasks is *sequential transfer learning* task, in which tasks are learned sequentially and the target task has labeled data.

5.2 How to Transfer?

To transfer the knowledge of a PTM to the downstream NLP tasks, we need to consider the following issues:

5.2.1 Choosing appropriate pre-training task, model architecture and corpus

Different PTMs usually have different effects on the same downstream task, since these PTMs are trained with various pre-training tasks, model architecture, and corpora.

(1) Currently, language model is the most popular pre-training task and can more efficiently solve a wide range of NLP problems [131]. However, different pre-training tasks have their own bias and give different effects for different tasks. For example, the NSP task [32] makes PTM understand the relationship between two sentences. Thus, the PTM can benefit downstream tasks such as Question Answering (QA) and Natural Language Inference (NLI).

(2) The architecture of PTM is also important for the downstream task. For example, although BERT helps with most natural language understanding tasks, it is hard to generate language.

(3) The data distribution of the downstream task should be approximate to PTMs. Currently, there are a large number of off-the-shelf PTMs, which can just as conveniently be used for various domain-specific or language-specific downstream tasks.

Therefore, given a target task, it is always a good solution to choose the PTMs trained with appropriate pre-training task, architecture, and corpus.

5.2.2 Choosing appropriate layers

Given a pre-trained deep model, different layers should capture different kinds of information, such as POS tagging, parsing, long-term dependencies, semantic roles, coreference, and so on. For RNN-based models, Belinkov et al. [9] and Melamud et al. [114] showed that representations learned from different layers in a multi-layer LSTM encoder benefit different tasks (e.g., predicting POS tags and understanding word sense). For transformer-based PTMs, Tenney et al. [166] found BERT represents the steps of the traditional NLP pipeline: basic syntactic information appears earlier in the network, while high-level semantic information appears at higher layers.

Let $\mathbf{H}^{(l)}$ ($1 \leq l \leq L$) denotes the l -th layer representation of the pre-trained model with L layers, and $g(\cdot)$ denote the task-specific model for the target task.

There are three ways to select the representation:

a) *Embedding Only*. One approach is to choose only the pre-trained static embeddings, while the rest of the model still needs to be trained from scratch for a new target task.

They fail to capture higher-level information that might be even more useful. Word embeddings are only useful in capturing semantic meanings of words but we also need to understand higher-level concepts like word sense.

b) *Top Layer*. The most simple and effective way is to feed the representation at the top layer into the task-specific model $g(\mathbf{H}^{(L)})$.

c) *All Layers*. A more flexible way is to automatic choose the best layer in a soft version, like ELMo [122]:

$$\mathbf{r}_t = \gamma \sum_{l=1}^L \alpha_l \mathbf{h}_t^{(l)}, \quad (6)$$

where α_l is the softmax-normalized weight for layer l and γ is a scalar to scale the vectors output by pre-trained model. The mixup representation is fed into the task-specific model $g(\mathbf{r}_t)$.

5.2.3 To tune or not to tune?

Currently, there are two common ways of model transfer: feature extraction (where the pre-trained parameters are frozen), and fine-tuning (where the pre-trained parameters are unfrozen and fine-tuned).

In feature extraction way, the pretrained models are regarded as off-the-shelf feature extractors. Moreover, it is important to expose the internal layers as they typically encode the most transferable representations [124].

Although both these two ways can significantly benefit most of NLP tasks, feature extraction way requires more complex task-specific architecture. Therefore, fine-tuning way is usually more general and convenient for many different downstream tasks than feature extraction way.

Table 4 gives some common combinations of adapting PTMs.

Table 4: Some common combinations of adapting PTMs.

Where	FT/FE? [†]	PTMs
Embedding Only	FT/FE	Word2vec [116], GloVe [120]
Top Layer	FT	BERT [32], RoBERTa [105]
Top Layer	FE	BERT [§] [203, 204]
All Layers	FE	ELMo [122]

[†] FT and FE mean Fine-tuning and Feature Extraction respectively.

[§] BERT used as feature extractor.

5.3 Fine-Tuning Strategies

With the increase of the depth of PTMs, the representation captured by them makes the downstream task easier. Therefore, the task-specific layer of the whole model is simple. Since ULMFit and BERT, fine-tuning has become the main adaption method of PTMs. However, the process of fine-tuning is often

brittle: even with the same hyper-parameter values, distinct random seeds can lead to substantially different results [34].

Besides standard fine-tuning, there are also some useful fine-tuning strategies.

Two-stage fine-tuning An alternative solution is *two-stage transfer*, which introduces an intermediate stage between pre-training and fine-tuning. In the first stage, the PTM is transferred into a model fine-tuned by an intermediate task or corpus. In the second stage, the transferred model is fine-tuned to the target task. Sun et al. [154] showed that the “further pre-training” on the related-domain corpus can further improve the ability of BERT and achieved state-of-the-art performance on eight widely-studied text classification datasets. Phang et al. [126] and Garg et al. [44] introduced the intermediate supervised task related to the target task, which brings a large improvement for BERT, GPT, and ELMo.

Multi-task fine-tuning Liu et al. [103] fine-tuned BERT under the multi-task learning framework, which demonstrates that multi-task learning and pre-training are complementary technologies.

Fine-tuning with extra adaptation modules The main drawback of fine-tuning is its parameter inefficiency: every downstream task has its own fine-tuned parameters. Therefore, a better solution is to inject some fine-tunable adaptation modules into PTMs while the original parameters are fixed.

Stickland and Murray [149] equipped a single share BERT model with small additional task-specific adaptation modules, projected attention layers (PALs). The shared BERT with the PALs matches separately fine-tuned models on the GLUE benchmark with roughly 7 times fewer parameters. Similarly, Housby et al. [61] modified the architecture of pre-trained BERT by adding adapter modules. Adapter modules yield a compact and extensible model; they add only a few trainable parameters per task, and new tasks can be added without revisiting previous ones. The parameters of the original network remain fixed, yielding a high degree of parameter sharing.

Others Instead of fine-tuning all the layers simultaneously, *gradual unfreezing* [62] is also an effective method that gradually unfreeze layers of PTMs starting from the top layer. Chronopoulou et al. [19] proposed a simpler unfreezing method, *sequential unfreezing*, which first fine-tunes only the randomly-initialized task specific layers, and then unfreezes the hidden layers of PTM, and finally unfreezes the embedding layer.

Motivated by the success of widely-used ensemble models, Xu et al. [191] improved the fine-tuning of BERT with two effective mechanisms: *self-ensemble* and *self-distillation*.

Generally, the above works show that the utility of PTMs can be further stimulated by better fine-tuning strategies.

Table 5: Open-Source Implementations

System	Framework	PTMs	URL
word2vec	-	CBOW, Skip-Gram	https://github.com/tmikolov/word2vec
GloVe	-	Pre-trained word vectors	https://nlp.stanford.edu/projects/glove
FastText	-	Pre-trained word vectors	https://github.com/facebookresearch/fastText
Transformers	PyTorch & TF	BERT, GPT-2, RoBERTa, XLNet, etc.	https://github.com/huggingface/transformers
Fairseq	PyTorch	English LM, German LM, RoBERTa, etc.	https://github.com/pytorch/fairseq
Flair	PyTorch	BERT, ELMo, GPT, RoBERTa, XLNet, etc.	https://github.com/flairNLP/flair
AllenNLP [43]	PyTorch	ELMo, BERT, GPT-2, etc.	https://github.com/allenai/allennlp
FastNLP	PyTorch	BERT, RoBERTa, GPT, etc.	https://github.com/fastnlp/fastnlp
Chinese-BERT [26]	-	BERT, RoBERTa, etc. (for Chinese)	https://github.com/ymcui/Chinese-BERT-wwm
BERT [32]	TF	BERT, BERT-wwm	https://github.com/google-research/bert
RoBERTa [105]	PyTorch		https://github.com/pytorch/fairseq/tree/master/examples/roberta
XLNet [194]	TF		https://github.com/zihangdai/xlnet/
ALBERT [86]	TF		https://github.com/google-research/ALBERT
T5 [132]	TF		https://github.com/google-research/text-to-text-transfer-transformer
ERNIE(Baidu) [157, 158]	PaddlePaddle		https://github.com/PaddlePaddle/ERNIE

6 Resources of PTMs

6.1 Open-Source Implementations

There are many third-party implementations for PTMs systems available online with pretrained models. Table 5 summarizes popular ones.

6.2 Collections of Related Resources

Table 6 provides some repositories that list papers and other related resource of PTMs.

Table 6: Collections of Related Resources

Collection	URL
Papers List	https://github.com/thunlp/PLMpapers
Papers List	https://github.com/tomohideshibata/BERT-related-papers
Papers List	https://github.com/cedrickchee/awesome-bert-nlp
Bert Lang Street [†]	https://bertlang.unibocconi.it/
BertViz [172]	https://github.com/jessevig/bertviz

[†] a collection of BERT models with reported performances on different datasets, tasks and languages.

7 Applications

In this section, we summarize some applications of PTMs in several classic NLP tasks.

7.1 General Evaluation Benchmark

There is an essential issue for the NLP community that how can we evaluate PTMs in a comparable metric. Thus, large-scale-benchmark is necessary.

The General Language Understanding Evaluation (GLUE) benchmark [177] is a collection of nine natural language understanding tasks, including single-sentence classification tasks

(CoLA and SST-2), pairwise text classification tasks (MNLI, RTE, WNLI, QQP, and MRPC), text similarity task (STS-B), and relevant ranking task (QNLI). GLUE benchmark is well-designed for evaluating the robustness as well as generalization of models. GLUE does not provide the labels for the test set, but set up an evaluation server.

However, motivating by the fact that the progress in recent years has eroded headroom on the GLUE benchmark dramatically, a new benchmark called SuperGLUE [176] was presented. Compared to GLUE, SuperGLUE has more challenging tasks and more diverse task formats (e.g., coreference resolution and question answering).

State-of-the-art PTMs are listed in the corresponding leaderboard. ⁴⁾ ⁵⁾

7.2 Machine Translation

Machine Translation (MT) is an important task in the NLP community which has attracted many researchers. Almost all of Neural Machine Translation (NMT) models share the encoder-decoder framework, which first encodes input tokens to hidden representations by the encoder and then decodes output tokens in the target language from the decoder. Given the superb performance of BERT on other NLP tasks, it is natural to investigate how to incorporate pre-training techniques into NMT models.

Conneau and Lample [25] tried to initialize the entire encoder and decoder by a multilingual pre-trained BERT model and shows a significant improvement can be achieved on unsupervised MT and English-Romanian supervised MT. Edunov et al. [37] used ELMo to set the word embedding layer in the NMT model. This work shows performance improvements on English-Turkish and English-German NMT model by using a

⁴⁾ <https://gluebenchmark.com/>

⁵⁾ <https://super.gluebenchmark.com/>

pre-trained language model for source word embedding initialization. Similarly, Clinchant et al. [23] devised a series of different experiments for examining the best strategy to utilize BERT on the encoder part of NMT models. They achieve some improvement by using BERT as an initialization of the encoder. Also, they found that these models can get better performance on the out-of-domain dataset. Imamura and Sumita [65] proposed a two stages BERT fine-tuning method for NMT. At the first stage, the encoder is initialized by a pre-trained BERT model and they only train the decoder on the training set. At the second stage, the whole NMT model is jointly fine-tuned on the training set. By experiment, they show this approach can surpass the one stage fine-tuning method which directly fine-tunes the whole model. Apart from that, Zhu et al. [204] suggested using pre-trained BERT as an extra memory to facilitate NMT models. Concretely, they first encode the input tokens by a pre-trained BERT and use the output of the last layer as an extra memory. Then, the NMT model can access the memory via an extra attention module in each layer of both encoder and decoder. And they show a noticeable improvement on supervised, semi-supervised and unsupervised MT. Instead of only pre-training the encoder, Song et al. [147] proposed a masked sequence-to-sequence pre-training method (MASS) to pre-train the encoder and decoder jointly. In the experiment, this approach can surpass the BERT-style pre-training proposed by Conneau and Lample [25] both on unsupervised MT and English-Romanian supervised MT.

7.3 Question Answering

Question answering (QA), or a narrower concept machine reading comprehension (MRC), is an important application in the NLP community. From easy to hard, there are three types of QA tasks: single-round extractive QA (SQuAD) [133], multi-round generative QA (CoQA) [135], and multi-hop QA (HotpotQA) [193].

BERT creatively transforms the extractive QA task to the spans prediction task that predicts the starting span as well as the ending span of the answer [32]. After that, PTM as an encoder for predicting spans has become a competitive baseline. For extractive QA, Zhang et al. [200] proposed a retrospective reader architecture and initialize the encoder with PTM (e.g., ALBERT). For multi-round generative QA, Ju et al. [73] proposed a “PTM+Adversarial Training+Rationale Tagging+Knowledge Distillation” model. For multi-hop QA, Tu et al. [169] proposed an interpretable “Select, Answer, and Explain” (SAE) system that PTM acts as the encoder in the selection module.

Generally, encoder parameters in the proposed QA model are initialized through a PTM, and other parameters are randomly initialized. State-of-the-art models are listed in the corresponding leaderboard.^{6) 7) 8)}

7.4 Sentiment Analysis

BERT outperforms previous state-of-the-art models by simply fine-tuning on SST-2, which is a widely used dataset for sentiment analysis (SA) [32]. Bataa and Wu [8] utilized BERT with transfer learning techniques and achieve new state-of-the-art in Japanese SA.

Despite their success in simple sentiment classification, directly applying BERT to aspect-based sentiment analysis (ABSA), which is a fine-grained SA task, shows less significant improvement [153]. To better leverage the powerful representation of BERT, Sun et al. [153] constructed an auxiliary sentence by transforming ABSA from a single sentence classification task to a sentence pair classification task. Xu et al. [189] proposed post-training to adapt BERT from its source domain and tasks to the ABSA domain and tasks. Furthermore, Rietzler et al. [137] extended the work of [189] by analyzing the behavior of cross-domain post-training with ABSA performance. Karimi et al. [76] showed that the performance of post-trained BERT can be further improved via adversarial training. Song et al. [148] added an additional pooling module, which can be implemented as either LSTM or attention mechanism, to leverage BERT intermediate layers for ABSA. In addition, Li et al. [97] jointly learned aspect detection and sentiment classification towards end-to-end ABSA.

For sentiment transfer, Wu et al. [186] proposed “Mask and Infill” based on BERT. In the mask step, the model disentangle sentiment from content by masking sentiment tokens. In the infill step, it uses BERT along with a target sentiment embedding to infill the masked positions.

7.5 Summarization

Summarization, aiming at producing a shorter text which preserves the most meaning of a longer text, has attracted the attention of the NLP community in recent years. The task has been improved significantly since the widespread use of PTM. Zhong et al. [203] introduced transferable knowledge (e.g., BERT) for summarization and surpassed previous models. Zhang et al. [198] tries to pre-trained a document-level model that predicts sentences instead of words, and then apply it on downstream tasks such as summarization. More elaborately, Zhang et al. [197] designed a Gap Sentence Generation

6) <https://rajpurkar.github.io/SQuAD-explorer/>

7) <https://stanfordnlp.github.io/coqa/>

8) <https://hotpotqa.github.io/>

(GSG) task for pre-training, whose objective involves generating summary-like text from the input. Furthermore, Liu and Lapata [104] proposed BERTSUM. BERTSUM included a novel document-level encoder, and a general framework for both extractive summarization and abstractive summarization. In the encoder frame, BERTSUM extends BERT by inserting multiple [CLS] tokens to learn the sentence representations. For extractive summarization, BERTSUM stacks several inter-sentence Transformer layers. For abstractive summarization, BERTSUM proposes a two-staged fine-tuning approach using a new fine-tuning schedule.

7.6 Named Entity Recognition

Named Entity Recognition (NER) is a fundamental task in information extraction and plays an important role in many NLP downstream tasks. In deep learning, most of NER methods are in sequence-labeling framework. The entity information in a sentence will be transformed into the sequence of labels, and one label corresponds to one word. The model is used to predict the label of each word. Since ELMo and BERT have shown their power in NLP, there is much work about pre-trained models for NER.

TagLM [121] and ELMo [122] use a pre-trained language model’s last layer output and weighted-sum of each layer output as a part of word embedding. Devlin et al. [32] used the first BPE’s BERT representation to predict each word’s label without CRF. Pires et al. [127] realized zero-shot NER through multilingual BERT. Akbik et al. [1] used a pre-trained character-level language model to produce word-level embedding for NER. Liu et al. [98] used layer-wise pruning and dense connection to speed up ELMo’s inference on NER. Tsai et al. [168] leveraged knowledge distillation to run a small BERT for NER on a single CPU. Besides, BERT is also used on domain-specific NER, such as biomedicine [52, 91], etc.

8 Future Directions

Though PTMs have proven their power for various NLP tasks, challenges still exist due to the complexity of language. In this section, we suggest five future directions of PTMs.

(1) Upper Bound of PTMs Currently, PTMs have not yet reached its upper bound. Most of current PTMs can be further improved by more training steps and larger corpus.

The state of the art in NLP can be further advanced by increasing the depth of models, such as Megatron-LM [144] (8.3 billion parameters, 72 Transformer layers with a hidden size of 3072 and 32 attention heads) and Turing-NLG⁹⁾ (17

billion parameters, 78 Transformer layers with a hidden size of 4256 and 28 attention heads).

The general-purpose PTMs are always our pursuits for learning the intrinsic universal knowledge of language (even world knowledge), however, such PTMs usually need deeper architecture, larger corpus and challenging pre-training tasks, which further result in higher training cost. However, training huge models is also a challenging problem, which needs more sophisticated and efficient training techniques such as distributed training, mixed precision, gradient accumulation, etc. Therefore, a more practical direction is to design more efficient model architecture, self-supervised pre-training tasks, optimizers and training skills using existing hardware and software. ELECTRA [22] is a good solution towards this direction.

(2) Task-oriented Pre-training and Model Compression

In practice, different downstream tasks require the different abilities of PTMs. The discrepancy between PTMs and downstream tasks usually lies in two aspects: model architecture and data distribution. A larger discrepancy may result in that the benefit of PTMs may be insignificant. Although larger PTMs can usually lead to better performance, a practical problem is how to leverage these huge PTMs on special scenarios, such as low-capacity devices and low-latency applications. Therefore, we can carefully design the specific model architecture and pre-training tasks for downstream task or extract partial task-specific knowledge from existing PTMs.

Besides, instead of training task-oriented PTMs from scratch, we can teach them with general-purpose existing PTMs by using techniques such as model compression (see Section 4.3). Although model compression is widely studied for CNNs in CV [17], compression for PTMs for NLP is just beginning. The fully-connected structure of Transformer also makes model compression more challenging.

(3) Architecture of PTMs Transformer has been proved to be an effective architecture for pre-training. However, the main limitation of Transformer is its computation complexity, which is quadratic to the input length. Limited by the memory of GPUs, most of current PTMs cannot deal with the sequence longer than 512 tokens. Breaking this limit needs to improve the architecture of Transformer, such as Transformer-XL [27]. Therefore, searching for more efficient model architecture for PTMs is important to capture longer-range contextual information.

The design of deep architecture is challenging, and we may seek help from some automatic methods, such as neural architecture search (NAS) [205].

9) <https://www.microsoft.com/en-us/research/blog/turing-nlg-a-17-billion-parameter-language-model-by-microsoft/>

(4) Knowledge Transfer Beyond Fine-tuning Currently, fine-tuning is the dominant method to transfer PTMs' knowledge to downstream tasks, but one deficiency is its parameter inefficiency: every downstream task has its own fine-tuned parameters. An improved solution is to fix the original parameters of PTMs and by adding small fine-tunable adaption modules for specific task [149, 61]. Thus, we can use a shared PTM to serve multiple downstream tasks. Indeed, mining knowledge from PTMs can be more flexible, such as feature extraction, knowledge distillation [195], data augmentation [185, 84], using PTMs as external knowledge [125], and so on. More efficient methods are expected.

(5) Interpretability and Reliability of PTMs Although PTMs reach impressive performance, their deep non-linear architecture makes the procedure of decision-making highly non-transparent.

Recently, explainable artificial intelligence (XAI) [4] has become a hotspot in general AI community. Unlike CNNs for images, interpreting PTMs is harder due to the complexities of both the Transformer-like architecture and language. Extensive efforts (see Section 3.3) have been made to analyze the linguistic and world knowledge included in PTMs, which help us understand these PTMs with some degree of transparency. However, much work on model analysis depends on the attention mechanism, and the effectiveness of attention for interpretability is still controversial [66, 142].

Besides, the reliability of PTMs is also becoming an issue of great concern with the extensive use of PTMs in production systems. The deep neural models are vulnerable to adversarial examples that can mislead a model to produce a specific wrong prediction with imperceptible perturbations from the original input. In CV, adversarial attacks and defenses have been widely studied. However, it is still challenging for text due to the discrete nature of languages. Generating of adversarial samples for text needs to possess such qualities: (1) imperceptible to human judges yet misleading to neural models; (2) fluent in grammar and semantically consistent with original inputs. Jin et al. [71] successfully attacked the fine-tuned BERT on text classification and textual entailment with adversarial examples. Wallace et al. [175] defined universal adversarial triggers that can induce a model to produce a specific-purpose prediction when concatenated to any input. Some triggers can even cause the GPT-2 model to generate racist text. The studies of adversarial attacks against PTMs help us understand their capabilities by fully exposing their vulnerabilities. Sun et al. [155] showed BERT is not robust on misspellings. Besides, adversarial defenses for PTMs are also promising, which improve the robustness of PTMs and make them be immune against adversarial attack.

Overall, as key components in many NLP applications, the

interpretability and reliability of PTMs still remain to be explored further in many respects, which helps us understand how PTMs work, and provides guide for better usage and further improvement.

9 Conclusion

In this survey, we conduct a comprehensive overview of PTMs for NLP, including the background knowledge, model architecture, pre-training tasks, various extensions, adaption approaches, applications, related resources of open-source systems and paper lists. Based on the current PTMs, we categorize them from four different perspectives: 1) type of word representation; 2) architecture of PTMs; 3) type of pre-training tasks; 4) extensions for specific types of scenarios or inputs. We also suggest several possible future research directions for PTMs.

References

- [1] Alan Akbik, Duncan Blythe, and Roland Vollgraf. Contextual string embeddings for sequence labeling. In *COLING*, pages 1638–1649, 2018.
- [2] Chris Alberti, Jeffrey Ling, Michael Collins, and David Reitter. Fusion of detected objects in text for visual question answering. In *EMNLP-IJCNLP*, pages 2131–2140, 2019.
- [3] Emily Alsentzer, John R. Murphy, Willie Boag, Wei-Hung Weng, Di Jin, Tristan Naumann, and Matthew B. A. McDermott. Publicly available clinical BERT embeddings. *arXiv preprint arXiv:1904.03323*, 2019.
- [4] Alejandro Barredo Arrieta, Natalia Díaz-Rodríguez, Javier Del Ser, Adrien Bannetot, Siham Tabik, Alberto Barbado, Salvador García, Sergio Gil-López, Daniel Molina, Richard Benjamins, et al. Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai. *Information Fusion*, 58:82–115, 2020.
- [5] Alexei Baevski, Sergey Edunov, Yinhan Liu, Luke Zettlemoyer, and Michael Auli. Cloze-driven pretraining of self-attention networks. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan, editors, *EMNLP-IJCNLP*, pages 5359–5368, 2019.
- [6] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- [7] Hangbo Bao, Li Dong, Furu Wei, Wenhui Wang, Nan Yang, Xiaodong Liu, Yu Wang, Songhao Piao, Jianfeng Gao, Ming Zhou, et al. Unilmv2: Pseudo-masked language models for unified language model pre-training. *arXiv preprint arXiv:2002.12804*, 2020.
- [8] Enkhbold Bataa and Joshua Wu. An investigation of transfer learning-based sentiment analysis in japanese. In *ACL*, 2019.

- [9] Yonatan Belinkov, Nadir Durrani, Fahim Dalvi, Hassan Sajjad, and James Glass. What do neural machine translation models learn about morphology? In *ACL*, pages 861–872, 2017.
- [10] Iz Beltagy, Kyle Lo, and Arman Cohan. SciBERT: A pre-trained language model for scientific text. In *EMNLP-IJCNLP*, pages 3613–3618, 2019.
- [11] Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Jauvin. A neural probabilistic language model. *Journal of machine learning research*, 3(Feb):1137–1155, 2003.
- [12] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828, 2013.
- [13] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. Enriching word vectors with subword information. *TACL*, 5:135–146, 2017.
- [14] Zied Bouraoui, José Camacho-Collados, and Steven Schockaert. Inducing relational knowledge from BERT. *arXiv preprint arXiv:1911.12753*, 2019.
- [15] Cristian Bucilua, Rich Caruana, and Alexandru Niculescu-Mizil. Model compression. In *KDD*, pages 535–541, 2006.
- [16] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. UNITER: learning universal image-text representations. *arXiv preprint arXiv:1909.11740*, 2019.
- [17] Yu Cheng, Duo Wang, Pan Zhou, and Tao Zhang. A survey of model compression and acceleration for deep neural networks. *arXiv preprint arXiv:1710.09282*, 2017.
- [18] Yew Ken Chia, Sam Witteveen, and Martin Andrews. Transformer to CNN: Label-scarce distillation for efficient text classification. *arXiv preprint arXiv:1909.03508*, 2019.
- [19] Alexandra Chronopoulou, Christos Baziotis, and Alexandros Potamianos. An embarrassingly simple approach for transfer learning from pretrained language models. In *NAACL-HLT*, pages 2089–2095, 2019.
- [20] Yung-Sung Chuang, Chi-Liang Liu, and Hung-yi Lee. SpeechBERT: Cross-modal pre-trained language model for end-to-end spoken question answering. *arXiv preprint arXiv:1910.11559*, 2019.
- [21] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014.
- [22] Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. ELECTRA: Pre-training text encoders as discriminators rather than generators. In *ICLR*, 2020.
- [23] Stephane Clinchant, Kweon Woo Jung, and Vassilina Nikoulina. On the use of BERT for neural machine translation. In *Proceedings of the 3rd Workshop on Neural Generation and Translation*, Hong Kong, November 2019.
- [24] Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel P. Kuksa. Natural language processing (almost) from scratch. *J. Mach. Learn. Res.*, 2011.
- [25] Alexis Conneau and Guillaume Lample. Cross-lingual language model pretraining. In *NeurIPS*, pages 7057–7067, 2019.
- [26] Yiming Cui, Wanxiang Che, Ting Liu, Bing Qin, Ziqing Yang, Shijin Wang, and Guoping Hu. Pre-training with whole word masking for chinese BERT. *arXiv preprint arXiv:1906.08101*, 2019.
- [27] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc Le, and Ruslan Salakhutdinov. Transformer-XL: Attentive language models beyond a fixed-length context. In *ACL*, pages 2978–2988, 2019.
- [28] Joe Davison, Joshua Feldman, and Alexander M. Rush. Commonsense knowledge mining from pretrained models. In *EMNLP-IJCNLP*, pages 1173–1178, 2019.
- [29] Wietse de Vries, Andreas van Cranenburgh, Arianna Bisazza, Tommaso Caselli, Gertjan van Noord, and Malvina Nissim. BERTje: A dutch BERT model. *arXiv preprint arXiv:1912.09582*, 2019.
- [30] Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Lukasz Kaiser. Universal transformers. In *ICLR*, 2019.
- [31] Pieter Delobelle, Thomas Winters, and Bettina Berendt. RoBERT: a dutch RoBERTa-based language model. *arXiv preprint arXiv:2001.06286*, 2020.
- [32] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, 2019.
- [33] Shizhe Diao, Jiaxin Bai, Yan Song, Tong Zhang, and Yonggang Wang. ZEN: pre-training chinese text encoder enhanced by n-gram representations. *arXiv preprint arXiv:1911.00720*, 2019.
- [34] Jesse Dodge, Gabriel Ilharco, Roy Schwartz, Ali Farhadi, Hannaneh Hajishirzi, and Noah Smith. Fine-tuning pretrained language models: Weight initializations, data orders, and early stopping. *arXiv preprint arXiv:2002.06305*, 2020.
- [35] Li Dong, Nan Yang, Wenhui Wang, Furu Wei, Xiaodong Liu, Yu Wang, Jianfeng Gao, Ming Zhou, and Hsiao-Wuen Hon. Unified language model pre-training for natural language understanding and generation. In *NeurIPS*, pages 13042–13054, 2019.
- [36] Zhen Dong, Zhewei Yao, Amir Gholami, Michael W Mahoney, and Kurt Keutzer. Hawq: Hessian aware quantization of neural networks with mixed-precision. In *ICCV*, pages 293–302, 2019.
- [37] Sergey Edunov, Alexei Baevski, and Michael Auli. Pre-trained language model representations for language generation. In Jill Burstein, Christy Doran, and Tamar Solorio, editors, *NAACL-HLT*, pages 4052–4059, 2019.
- [38] Julian Eisenschlos, Sebastian Ruder, Piotr Czapla, Marcin Kadras, Sylvain Gugger, and Jeremy Howard. MultiFiT: Efficient multi-lingual language model fine-tuning. In *EMNLP-IJCNLP*, pages 5701–5706, 2019.

- [39] Dumitru Erhan, Yoshua Bengio, Aaron C. Courville, Pierre-Antoine Manzagol, Pascal Vincent, and Samy Bengio. Why does unsupervised pre-training help deep learning? *J. Mach. Learn. Res.*, 11:625–660, 2010.
- [40] Allyson Ettinger. What BERT is not: Lessons from a new suite of psycholinguistic diagnostics for language models. *TACL*, 8:34–48, 2020.
- [41] Manaal Faruqui and Chris Dyer. Improving vector space word representations using multilingual correlation. In *EACL*, pages 462–471, 2014.
- [42] Prakhar Ganesh, Yao Chen, Xin Lou, Mohammad Ali Khan, Yin Yang, Deming Chen, Marianne Winslett, Hassan Sajjad, and Preslav Nakov. Compressing large-scale transformer-based models: A case study on bert. *arXiv preprint arXiv:2002.11985*, 2020.
- [43] Matt Gardner, Joel Grus, Mark Neumann, Oyvind Tafjord, Pradeep Dasigi, Nelson F. Liu, Matthew Peters, Michael Schmitz, and Luke S. Zettlemoyer. Allennlp: A deep semantic natural language processing platform. 2017.
- [44] Siddhant Garg, Thuy Vu, and Alessandro Moschitti. Tanda: Transfer and adapt pre-trained transformer models for answer sentence selection. *arXiv preprint arXiv:1911.04118*, 2019.
- [45] Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N Dauphin. Convolutional sequence to sequence learning. In *ICML*, pages 1243–1252, 2017.
- [46] Yoav Goldberg. Assessing BERT’s syntactic abilities. *arXiv preprint arXiv:1901.05287*, 2019.
- [47] Mitchell A Gordon, Kevin Duh, and Nicholas Andrews. Compressing BERT: Studying the effects of weight pruning on transfer learning. *arXiv preprint arXiv:2002.08307*, 2020.
- [48] Jian Guan, Fei Huang, Zhihao Zhao, Xiaoyan Zhu, and Minlie Huang. A knowledge-enhanced pretraining model for commonsense story generation. *arXiv preprint arXiv:2001.05139*, 2020.
- [49] Qipeng Guo, Xipeng Qiu, Pengfei Liu, Yunfan Shao, Xiangyang Xue, and Zheng Zhang. Star-transformer. In *NAACL-HLT*, pages 1315–1325, 2019.
- [50] Abhijeet Gupta, Gemma Boleda, Marco Baroni, and Sebastian Padó. Distributional vectors encode referential attributes. In *EMNLP*, pages 12–21, 2015.
- [51] Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *AISTATS*, pages 297–304, 2010.
- [52] Kai Hakala and Sampo Pyysalo. Biomedical named entity recognition with multilingual BERT. In *BioNLP Open Shared Tasks@EMNLP*, pages 56–61, 2019.
- [53] Hiroaki Hayashi, Zecong Hu, Chenyan Xiong, and Graham Neubig. Latent relation language models. *arXiv preprint arXiv:1908.07690*, 2019.
- [54] Bin He, Di Zhou, Jinghui Xiao, Xin Jiang, Qun Liu, Nicholas Jing Yuan, and Tong Xu. Integrating graph contextualized knowledge into pre-trained language models. *arXiv preprint arXiv:1912.00147*, 2019.
- [55] John Hewitt and Christopher D. Manning. A structural probe for finding syntax in word representations. In *NAACL-HLT*, pages 4129–4138, 2019.
- [56] GE Hinton, JL McClelland, and DE Rumelhart. Distributed representations. In *Parallel distributed processing: explorations in the microstructure of cognition, vol. 1: foundations*, pages 77–109. 1986.
- [57] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [58] Geoffrey E Hinton and Ruslan R Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 313 (5786):504–507, 2006.
- [59] R. Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Philip Bachman, Adam Trischler, and Yoshua Bengio. Learning deep representations by mutual information estimation and maximization. In *ICLR*, 2019. URL <https://openreview.net/forum?id=Bklr3j0cKX>.
- [60] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 1997.
- [61] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for NLP. In *ICML*, pages 2790–2799, 2019.
- [62] Jeremy Howard and Sebastian Ruder. Universal language model fine-tuning for text classification. In *ACL*, pages 328–339, 2018.
- [63] Haoyang Huang, Yaobo Liang, Nan Duan, Ming Gong, Linjun Shou, Daxin Jiang, and Ming Zhou. Unicoder: A universal language encoder by pre-training with multiple cross-lingual tasks. In *EMNLP-IJCNLP*, pages 2485–2494, 2019.
- [64] Kexin Huang, Jaan Allosa, and Rajesh Ranganath. ClinicalBERT: Modeling clinical notes and predicting hospital readmission. *arXiv preprint arXiv:1904.05342*, 2019.
- [65] Kenji Imamura and Eiichiro Sumita. Recycling a pre-trained BERT encoder for neural machine translation. In *Proceedings of the 3rd Workshop on Neural Generation and Translation*, Hong Kong, November 2019.
- [66] Sarthak Jain and Byron C Wallace. Attention is not explanation. In *NAACL-HLT*, pages 3543–3556, 2019.
- [67] Ganesh Jawahar, Benoît Sagot, and Djamé Seddah. What does BERT learn about the structure of language? In *ACL*, pages 3651–3657, 2019.
- [68] Zongcheng Ji, Qiang Wei, and Hua Xu. BERT-based ranking for biomedical entity normalization. *arXiv preprint arXiv:1908.03548*, 2019.
- [69] Zhengbao Jiang, Frank F. Xu, Jun Araki, and Graham Neubig. How can we know what language models know? *arXiv preprint arXiv:1911.12543*, 2019.
- [70] Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen,

- Linlin Li, Fang Wang, and Qun Liu. TinyBERT: Distilling BERT for natural language understanding. *arXiv preprint arXiv:1909.10351*, 2019.
- [71] Di Jin, Zhijing Jin, Joey Tianyi Zhou, and Peter Szolovits. Is bert really robust? natural language attack on text classification and entailment. *arXiv preprint arXiv:1907.11932*, 2019.
- [72] Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S. Weld, Luke Zettlemoyer, and Omer Levy. SpanBERT: Improving pre-training by representing and predicting spans. *arXiv preprint arXiv:1907.10529*, 2019.
- [73] Ying Ju, Fubang Zhao, Shijie Chen, Bowen Zheng, Xuefeng Yang, and Yunfeng Liu. Technical report on conversational question answering. *arXiv preprint arXiv:1909.10772*, 2019.
- [74] Karthikeyan K, Zihan Wang, Stephen Mayhew, and Dan Roth. Cross-lingual ability of multilingual BERT: An empirical study. In *ICLR*, 2020. URL <https://openreview.net/forum?id=HJeT3yrtDr>.
- [75] Nal Kalchbrenner, Edward Grefenstette, and Phil Blunsom. A convolutional neural network for modelling sentences. *arXiv preprint arXiv:1404.2188*, 2014.
- [76] Akbar Karimi, Leonardo Rossi, Andrea Prati, and Katharina Full. Adversarial training for aspect-based sentiment analysis with BERT. *arXiv preprint arXiv:2001.11316*, 2020.
- [77] Nora Kassner and Hinrich Schütze. Negated LAMA: birds cannot fly. *arXiv preprint arXiv:1911.03343*, 2019.
- [78] Pei Ke, Haozhe Ji, Siyang Liu, Xiaoyan Zhu, and Minlie Huang. Sentir: Linguistic knowledge enhanced language representation for sentiment analysis. *arXiv preprint arXiv:1911.02493*, 2019.
- [79] Taeuk Kim, Jihun Choi, Daniel Edmiston, and Sang-goo Lee. Are pre-trained language models aware of phrases? simple but strong baselines for grammar induction. *arXiv preprint arXiv:2002.00737*, 2020.
- [80] Yoon Kim. Convolutional neural networks for sentence classification. In *EMNLP*, pages 1746–1751, 2014.
- [81] Yoon Kim, Yacine Jernite, David Sontag, and Alexander M Rush. Character-aware neural language models. In *AAAI*, 2016.
- [82] Ryan Kiros, Yukun Zhu, Ruslan R Salakhutdinov, Richard Zemel, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. Skip-thought vectors. In *NeurIPS*, pages 3294–3302, 2015.
- [83] Lingpeng Kong, Cyprien de Masson d’Autume, Lei Yu, Wang Ling, Zihang Dai, and Dani Yogatama. A mutual information maximization perspective of language representation learning. In *ICLR*, 2019.
- [84] Varun Kumar, Ashutosh Choudhary, and Eunah Cho. Data augmentation using pre-trained transformer models. *arXiv preprint arXiv:2003.02245*, 2020.
- [85] Yuri Kuratov and Mikhail Arkhipov. Adaptation of deep bidirectional multilingual transformers for russian language. *arXiv preprint arXiv:1905.07213*, 2019.
- [86] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. ALBERT: A lite BERT for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942*, 2019.
- [87] Anne Lauscher, Ivan Vulic, Edoardo Maria Ponti, Anna Korhonen, and Goran Glavas. Informing unsupervised pre-training with external linguistic knowledge. *arXiv preprint arXiv:1909.02339*, 2019.
- [88] Hang Le, Loïc Vial, Jibril Frej, Vincent Segonne, Maximin Coavoux, Benjamin Lecouteux, Alexandre Allauzen, Benoît Crabbé, Laurent Besacier, and Didier Schwab. FlauBERT: Unsupervised language model pre-training for french. *arXiv preprint arXiv:1912.05372*, 2019.
- [89] Quoc Le and Tomas Mikolov. Distributed representations of sentences and documents. In *ICML*, pages 1188–1196, 2014.
- [90] Jieh-Sheng Lee and Jieh Hsiang. PatentBERT: Patent classification with fine-tuning a pre-trained BERT model. *arXiv preprint arXiv:1906.02124*, 2019.
- [91] Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *arXiv preprint arXiv:1901.08746*, 2019.
- [92] Yoav Levine, Barak Lenz, Or Dagan, Dan Padnos, Or Sharir, Shai Shalev-Shwartz, Amnon Shashua, and Yoav Shoham. SenseBERT: Driving some sense into BERT. *arXiv preprint arXiv:1908.05646*, 2019.
- [93] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. BART: denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*, 2019.
- [94] Gen Li, Nan Duan, Yuejian Fang, Daxin Jiang, and Ming Zhou. Unicoder-vl: A universal encoder for vision and language by cross-modal pre-training. *arXiv preprint arXiv:1908.06066*, 2019.
- [95] Liunan Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. VisualBERT: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019.
- [96] Xiang Lisa Li and Jason Eisner. Specializing word embeddings (for parsing) by information bottleneck. In *EMNLP-IJCNLP*, pages 2744–2754, 2019.
- [97] Xin Li, Lidong Bing, Wenxuan Zhang, and Wai Lam. Exploiting BERT for end-to-end aspect-based sentiment analysis. In *W-NUT@EMNLP*, 2019.
- [98] Liyuan Liu, Xiang Ren, Jingbo Shang, Xiaotao Gu, Jian Peng, and Jiawei Han. Efficient contextualized representation: Language model pruning for sequence labeling. In *EMNLP*, pages 1215–1225, 2018.
- [99] Nelson F. Liu, Matt Gardner, Yonatan Belinkov, Matthew E.

- Peters, and Noah A. Smith. Linguistic knowledge and transferability of contextual representations. In *NAACL-HLT*, pages 1073–1094, 2019.
- [100] Pengfei Liu, Xipeng Qiu, and Xuanjing Huang. Recurrent neural network for text classification with multi-task learning. *arXiv preprint arXiv:1605.05101*, 2016.
- [101] Weijie Liu, Peng Zhou, Zhe Zhao, Zhiruo Wang, Qi Ju, Haotang Deng, and Ping Wang. K-BERT: Enabling language representation with knowledge graph. *arXiv preprint arXiv:1909.07606*, 2019.
- [102] Xiaodong Liu, Pengcheng He, Weizhu Chen, and Jianfeng Gao. Improving multi-task deep neural networks via knowledge distillation for natural language understanding. *arXiv preprint arXiv:1904.09482*, 2019.
- [103] Xiaodong Liu, Pengcheng He, Weizhu Chen, and Jianfeng Gao. Multi-task deep neural networks for natural language understanding. *arXiv preprint arXiv:1901.11504*, 2019.
- [104] Yang Liu and Mirella Lapata. Text summarization with pre-trained encoders. *arXiv preprint arXiv:1908.08345*, 2019.
- [105] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. RoBERTa: A robustly optimized BERT pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- [106] Robert L. Logan IV, Nelson F. Liu, Matthew E. Peters, Matt Gardner, and Sameer Singh. Barack’s wife hillary: Using knowledge graphs for fact-aware language modeling. In *ACL*, 2019.
- [107] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. ViLBERT: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *NeurIPS*, pages 13–23, 2019.
- [108] Wenhao Lu, Jian Jiao, and Ruofei Zhang. TwinBERT: Distilling knowledge to twin-structured BERT models for efficient retrieval. *arXiv preprint arXiv:2002.06275*, 2020.
- [109] Huaishao Luo, Lei Ji, Botian Shi, Haoyang Huang, Nan Duan, Tianrui Li, Xilin Chen, and Ming Zhou. Univlm: A unified video and language pre-training model for multimodal understanding and generation. *arXiv preprint arXiv:2002.06353*, 2020.
- [110] Minh-Thang Luong, Hieu Pham, and Christopher D Manning. Bilingual word representations with monolingual quality in mind. In *Proceedings of the 1st Workshop on Vector Space Modeling for Natural Language Processing*, pages 151–159, 2015.
- [111] Diego Marcheggiani, Joost Bastings, and Ivan Titov. Exploiting semantics in neural machine translation with graph convolutional networks. In *NAACL-HLT*, pages 486–492, 2018.
- [112] Louis Martin, Benjamin Müller, Pedro Javier Ortiz Suárez, Yoann Dupont, Laurent Romary, Éric Villemonte de la Clergerie, Djamé Seddah, and Benoît Sagot. CamemBERT: a tasty french language model. *arXiv preprint arXiv:1911.03894*, 2019.
- [113] Bryan McCann, James Bradbury, Caiming Xiong, and Richard Socher. Learned in translation: Contextualized word vectors. In *NeurIPS*, 2017.
- [114] Oren Melamud, Jacob Goldberger, and Ido Dagan. Context2Vec: Learning generic context embedding with bidirectional LSTM. In *CoNLL*, pages 51–61, 2016.
- [115] Paul Michel, Omer Levy, and Graham Neubig. Are sixteen heads really better than one? In *NeurIPS*, pages 14014–14024, 2019.
- [116] Tomas Mikolov, Ilya Sutskever, Kai Chen, Gregory S. Corrado, and Jeffrey Dean. Distributed representations of words and phrases and their compositionality. In *NeurIPS*, 2013.
- [117] Tomas Mikolov, Wen-tau Yih, and Geoffrey Zweig. Linguistic regularities in continuous space word representations. In *ACL*, pages 746–751, 2013.
- [118] Andriy Mnih and Koray Kavukcuoglu. Learning word embeddings efficiently with noise-contrastive estimation. In *NeurIPS*, pages 2265–2273, 2013.
- [119] Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10):1345–1359, 2009.
- [120] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. GloVe: Global vectors for word representation. In *EMNLP*, 2014.
- [121] Matthew E. Peters, Waleed Ammar, Chandra Bhagavatula, and Russell Power. Semi-supervised sequence tagging with bidirectional language models. In *ACL*, pages 1756–1765, 2017.
- [122] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. In *NAACL-HLT*, 2018.
- [123] Matthew E. Peters, Mark Neumann, Robert L. Logan IV, Roy Schwartz, Vidur Joshi, Sameer Singh, and Noah A. Smith. Knowledge enhanced contextual word representations. In *EMNLP-IJCNLP*, 2019.
- [124] Matthew E. Peters, Sebastian Ruder, and Noah A. Smith. To tune or not to tune? adapting pretrained representations to diverse tasks. In *Proceedings of the 4th Workshop on Representation Learning for NLP, RepL4NLP@ACL 2019, Florence, Italy, August 2, 2019*, pages 7–14, 2019.
- [125] Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick S. H. Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander H. Miller. Language models as knowledge bases? In *EMNLP-IJCNLP*, pages 2463–2473, 2019.
- [126] Jason Phang, Thibault FÉvry, and Samuel R Bowman. Sentence encoders on STILTs: Supplementary training on intermediate labeled-data tasks. *arXiv preprint arXiv:1811.01088*, 2018.
- [127] Telmo Pires, Eva Schlinger, and Dan Garrette. How multilin-

- gual is multilingual BERT? *arXiv preprint arXiv:1906.01502*, 2019. URL <http://arxiv.org/abs/1906.01502>.
- [128] Nina Pörner, Ulli Waltinger, and Hinrich Schütze. BERT is not a knowledge base (yet): Factual knowledge vs. name-based reasoning in unsupervised QA. *arXiv preprint arXiv:1911.03681*, 2019.
 - [129] Nina Pörner, Ulli Waltinger, and Hinrich Schütze. BERT is not a knowledge base (yet): Factual knowledge vs. name-based reasoning in unsupervised QA. *CoRR*, abs/1911.03681, 2019.
 - [130] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. 2018. URL <https://s3-us-west-2.amazonaws.com/openai-assets/researchcovers/languageunsupervised/languageunderstandingpaper.pdf>.
 - [131] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI Blog*, 2019.
 - [132] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv preprint arXiv:1910.10683*, 2019.
 - [133] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100, 000+ questions for machine comprehension of text. In Jian Su, Xavier Carreras, and Kevin Duh, editors, *EMNLP*, pages 2383–2392, 2016.
 - [134] Prajit Ramachandran, Peter J Liu, and Quoc Le. Unsupervised pretraining for sequence to sequence learning. In *EMNLP*, pages 383–391, 2017.
 - [135] Siva Reddy, Danqi Chen, and Christopher D. Manning. Coqa: A conversational question answering challenge. *TACL*, 7: 249–266, 2019.
 - [136] Emily Reif, Ann Yuan, Martin Wattenberg, Fernanda B Viégas, Andy Coenen, Adam Pearce, and Been Kim. Visualizing and measuring the geometry of BERT. In *NeurIPS*, pages 8592–8600, 2019.
 - [137] Alexander Rietzler, Sebastian Stabinger, Paul Opitz, and Stefan Engl. Adapt or get left behind: Domain adaptation through BERT language model finetuning for aspect-target sentiment classification. *arXiv preprint arXiv:1908.11860*, 2019.
 - [138] Dana Rubinstein, Effi Levi, Roy Schwartz, and Ari Rappoport. How well do distributional models capture different types of semantic knowledge? In *ACL*, pages 726–730, 2015.
 - [139] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
 - [140] Nikunj Saunshi, Orestis Plevrakis, Sanjeev Arora, Mikhail Khodak, and Hrishikesh Khandeparkar. A theoretical analysis of contrastive unsupervised representation learning. In *ICML*, pages 5628–5637, 2019.
 - [141] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. In *ACL*, 2016.
 - [142] Sofia Serrano and Noah A Smith. Is attention interpretable? In *ACL*, pages 2931–2951, 2019.
 - [143] Sheng Shen, Zhen Dong, Jiayu Ye, Linjian Ma, Zhewei Yao, Amir Gholami, Michael W Mahoney, and Kurt Keutzer. Q-BERT: Hessian based ultra low precision quantization of BERT. In *AAAI*, 2020.
 - [144] Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. Megatron-LM: Training multi-billion parameter language models using gpu model parallelism. *arXiv preprint arXiv:1909.08053*, 2019.
 - [145] Karan Singla, Doğan Can, and Shrikanth Narayanan. A multi-task approach to learning multilingual representations. In *ACL*, pages 214–220, 2018.
 - [146] Richard Socher, Alex Perelygin, Jean Y Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *EMNLP*, pages 1631–1642. *ACL*, 2013.
 - [147] Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. MASS: masked sequence to sequence pre-training for language generation. In *ICML*, volume 97 of *Proceedings of Machine Learning Research*, pages 5926–5936, 2019.
 - [148] Youwei Song, Jiahai Wang, Zhiwei Liang, Zhiyue Liu, and Tao Jiang. Utilizing BERT intermediate layers for aspect based sentiment analysis and natural language inference. *arXiv preprint arXiv:2002.04815*, 2020.
 - [149] Asa Cooper Stickland and Iain Murray. BERT and PALs: Projected attention layers for efficient adaptation in multi-task learning. In *ICML*, pages 5986–5995, 2019.
 - [150] Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. VL-BERT: pre-training of generic visual-linguistic representations. *arXiv preprint arXiv:1908.08530*, 2019.
 - [151] Chen Sun, Fabien Baradel, Kevin Murphy, and Cordelia Schmid. Contrastive bidirectional transformer for temporal representation learning. *arXiv preprint arXiv:1906.05743*, 2019.
 - [152] Chen Sun, Austin Myers, Carl Vondrick, Kevin Murphy, and Cordelia Schmid. VideoBERT: A joint model for video and language representation learning. *arXiv preprint arXiv:1904.01766*, 2019.
 - [153] Chi Sun, Luyao Huang, and Xipeng Qiu. Utilizing BERT for aspect-based sentiment analysis via constructing auxiliary sentence. In *NAACL-HLT*, 2019.
 - [154] Chi Sun, Xipeng Qiu, Yige Xu, and Xuanjing Huang. How to fine-tune BERT for text classification? In *China National Conference on Chinese Computational Linguistics*, pages 194–206, 2019.
 - [155] Lichao Sun, Kazuma Hashimoto, Wenpeng Yin, Akari Asai, Jia Li, Philip Yu, and Caiming Xiong. Adv-bert: Bert is not

- robust on misspellings! generating nature adversarial samples on bert. *arXiv preprint arXiv:2003.04985*, 2020.
- [156] Siqi Sun, Yu Cheng, Zhe Gan, and Jingjing Liu. Patient knowledge distillation for BERT model compression. In *EMNLP-IJCNLP*, pages 4323–4332, 2019.
- [157] Yu Sun, Shuohuan Wang, Yukun Li, Shikun Feng, Xuyi Chen, Han Zhang, Xin Tian, Danxiang Zhu, Hao Tian, and Hua Wu. ERNIE: enhanced representation through knowledge integration. *arXiv preprint arXiv:1904.09223*, 2019.
- [158] Yu Sun, Shuohuan Wang, Yukun Li, Shikun Feng, Hao Tian, Hua Wu, and Haifeng Wang. ERNIE 2.0: A continual pre-training framework for language understanding. *arXiv preprint arXiv:1907.12412*, 2019.
- [159] Zhiqing Sun, Hongkun Yu, Xiaodan Song, Renjie Liu, Yiming Yang, and Denny Zhou. MobileBERT: Task-agnostic compression of BERT by progressive knowledge transfer. 2019. URL <https://openreview.net/pdf?id=SJxjVaNKwB>.
- [160] Ilya Sutskever, Oriol Vinyals, and Quoc VV Le. Sequence to sequence learning with neural networks. In *NeurIPS*, pages 3104–3112, 2014.
- [161] Kai Sheng Tai, Richard Socher, and Christopher D. Manning. Improved semantic representations from tree-structured long short-term memory networks. In *ACL*, pages 1556–1566, 2015.
- [162] Hao Tan and Mohit Bansal. LXMERT: learning cross-modality encoder representations from transformers. In *EMNLP-IJCNLP*, pages 5099–5110, 2019.
- [163] Matthew Tang, Priyanka Gandhi, Md Ahsanul Kabir, Christopher Zou, Jordyn Blakey, and Xiao Luo. Progress notes classification and keyword extraction using attention-based deep learning models with BERT. *arXiv preprint arXiv:1910.05786*, 2019.
- [164] Raphael Tang, Yao Lu, Linqing Liu, Lili Mou, Olga Vechtomova, and Jimmy Lin. Distilling task-specific knowledge from BERT into simple neural networks. *arXiv preprint arXiv:1903.12136*, 2019.
- [165] Wilson L. Taylor. “cloze procedure”: A new tool for measuring readability. *Journalism Quarterly*, 30(4):415–433, 1953.
- [166] Ian Tenney, Dipanjan Das, and Ellie Pavlick. BERT rediscovers the classical NLP pipeline. In Anna Korhonen, David R. Traum, and Lluís Màrquez, editors, *ACL*, pages 4593–4601, 2019.
- [167] Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang, Adam Poliak, R. Thomas McCoy, Najoung Kim, Benjamin Van Durme, Samuel R. Bowman, Dipanjan Das, and Ellie Pavlick. What do you learn from context? probing for sentence structure in contextualized word representations. In *ICLR*, 2019.
- [168] Henry Tsai, Jason Riesa, Melvin Johnson, Naveen Arivazhagan, Xin Li, and Amelia Archer. Small and practical BERT models for sequence labeling. In *EMNLP-IJCNLP*, pages 3632–3636, 2019.
- [169] Ming Tu, Kevin Huang, Guangtao Wang, Jing Huang, Xiaodong He, and Bowen Zhou. Select, answer and explain: Interpretable multi-hop reading comprehension over multiple documents. In *AAAI*, 2020.
- [170] Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Well-read students learn better: The impact of student initialization on knowledge distillation. *arXiv preprint arXiv:1908.08962*, 2019.
- [171] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017.
- [172] Jesse Vig. A multiscale visualization of attention in the transformer model. *arXiv preprint arXiv:1906.05714*, 2019. URL <https://arxiv.org/abs/1906.05714>.
- [173] Antti Virtanen, Jenna Kanerva, Rami Ilo, Jouni Luoma, Juhani Luotolahti, Tapio Salakoski, Filip Ginter, and Sampo Pyysalo. Multilingual is not enough: BERT for finnish. *arXiv preprint arXiv:1912.07076*, 2019.
- [174] Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned. In *ACL*, pages 5797–5808, 2019.
- [175] Eric Wallace, Shi Feng, Nikhil Kandpal, Matt Gardner, and Sameer Singh. Universal adversarial triggers for attacking and analyzing NLP. In *EMNLP-IJCNLP*, pages 2153–2162, 2019.
- [176] Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. SuperGLUE: A stickier benchmark for general-purpose language understanding systems. In *NeurIPS*, pages 3261–3275, 2019.
- [177] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *ICLR*, 2019.
- [178] Ruize Wang, Duyu Tang, Nan Duan, Zhongyu Wei, Xuanjing Huang, Jianshu Ji, Guihong Cao, Daxin Jiang, and Ming Zhou. K-adapter: Infusing knowledge into pre-trained models with adapters. *arXiv preprint arXiv:2002.01808*, 2020.
- [179] Shaolei Wang, Wanxiang Che, Qi Liu, Pengda Qin, Ting Liu, and William Yang Wang. Multi-task self-supervised learning for disfluency detection. *arXiv preprint arXiv:1908.05378*, 2019.
- [180] Wei Wang, Bin Bi, Ming Yan, Chen Wu, Zuyi Bao, Liwei Peng, and Luo Si. StructBERT: Incorporating language structures into pre-training for deep language understanding. In *ICLR*, 2020.
- [181] Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. MiniLM: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. *arXiv preprint arXiv:2002.10957*, 2020.
- [182] Xiaozhi Wang, Tianyu Gao, Zhaocheng Zhu, Zhiyuan Liu, Juanzi Li, and Jian Tang. KEPLER: A unified model for knowledge embedding and pre-trained language representa-

- tion. *arXiv preprint arXiv:1911.06136*, 2019.
- [183] Zhen Wang, Jianwen Zhang, Jianlin Feng, and Zheng Chen. Knowledge graph and text jointly embedding. In *EMNLP*, pages 1591–1601, 2014.
- [184] Junqiu Wei, Xiaozhe Ren, Xiaoguang Li, Wenyong Huang, Yi Liao, Yasheng Wang, Jiashu Lin, Xin Jiang, Xiao Chen, and Qun Liu. NEZHA: Neural contextualized representation for chinese language understanding. *arXiv preprint arXiv:1909.00204*, 2019.
- [185] Xing Wu, Shangwen Lv, Liangjun Zang, Jizhong Han, and Songlin Hu. Conditional BERT contextual augmentation. In *International Conference on Computational Science*, pages 84–95, 2019.
- [186] Xing Wu, Tao Zhang, Liangjun Zang, Jizhong Han, and Songlin Hu. "mask and infill" : Applying masked language model to sentiment transfer. *arXiv preprint arXiv:1908.08039*, 2019.
- [187] Ruobing Xie, Zhiyuan Liu, Jia Jia, Huanbo Luan, and Maosong Sun. Representation learning of knowledge graphs with entity descriptions. In *IJCAI*, 2016.
- [188] Wenhan Xiong, Jingfei Du, William Yang Wang, and Veselin Stoyanov. Pretrained encyclopedia: Weakly supervised knowledge-pretrained language model. In *ICLR*, 2020.
- [189] Hu Xu, Bing Liu, Lei Shu, and Philip S. Yu. BERT post-training for review reading comprehension and aspect-based sentiment analysis. In *NAACL-HLT*, 2019.
- [190] Jiacheng Xu, Xipeng Qiu, Kan Chen, and Xuanjing Huang. Knowledge graph representation with jointly structural and textual encoding. In *IJCAI*, pages 1318–1324, 2017.
- [191] Yige Xu, Xipeng Qiu, Ligao Zhou, and Xuanjing Huang. Improving BERT fine-tuning via self-ensemble and self-distillation. *arXiv preprint arXiv:2002.10345*, 2020.
- [192] An Yang, Quan Wang, Jing Liu, Kai Liu, Yajuan Lyu, Hua Wu, Qiaoqiao She, and Sujian Li. Enhancing pre-trained language representations with rich knowledge for machine reading comprehension. In *ACL*, pages 2346–2357, 2019.
- [193] Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In *EMNLP*, pages 2369–2380, 2018.
- [194] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. XLNet: Generalized autoregressive pretraining for language understanding. In *NeurIPS*, pages 5754–5764, 2019.
- [195] Ziqing Yang, Yiming Cui, Zhipeng Chen, Wanxiang Che, Ting Liu, Shijin Wang, and Guoping Hu. Textbrewer: An open-source knowledge distillation toolkit for natural language processing. *arXiv preprint arXiv:2002.12620*, 2020.
- [196] Ofir Zafrir, Guy Boudoukh, Peter Izsak, and Moshe Wasserblat. Q8BERT: Quantized 8bit BERT. *arXiv preprint arXiv:1910.06188*, 2019.
- [197] Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter J Liu. PEGASUS: Pre-training with extracted gap-sentences for abstractive summarization. *arXiv preprint arXiv:1912.08777*, 2019.
- [198] Xingxing Zhang, Furu Wei, and Ming Zhou. HIBERT: Document level pre-training of hierarchical bidirectional transformers for document summarization. In *ACL*, pages 5059–5069, 2019.
- [199] Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang, Maosong Sun, and Qun Liu. ERNIE: enhanced language representation with informative entities. In *ACL*, 2019.
- [200] Zhuosheng Zhang, Junjie Yang, and Hai Zhao. Retrospective reader for machine reading comprehension. *arXiv preprint arXiv:2001.09694*, 2020.
- [201] Sanqiang Zhao, Raghu Gupta, Yang Song, and Denny Zhou. Extreme language model compression with optimal subwords and shared projections. *arXiv preprint arXiv:1909.11687*, 2019.
- [202] Huaping Zhong, Jianwen Zhang, Zhen Wang, Hai Wan, and Zheng Chen. Aligning knowledge and text embeddings by entity descriptions. In *EMNLP*, pages 267–272, 2015.
- [203] Ming Zhong, Pengfei Liu, Danqing Wang, Xipeng Qiu, and Xuanjing Huang. Searching for effective neural extractive summarization: What works and what’s next. In *ACL*, pages 1049–1058, 2019.
- [204] Jinhua Zhu, Yingce Xia, Lijun Wu, Di He, Tao Qin, Wengang Zhou, Houqiang Li, and Tieyan Liu. Incorporating BERT into neural machine translation. In *ICLR*, 2020.
- [205] Barret Zoph and Quoc V Le. Neural architecture search with reinforcement learning. *arXiv preprint arXiv:1611.01578*, 2016.