

Fine-Tuning and Masked Language Models

*Larvatus prodeo [Masked, I go forward]
Descartes*

In the previous chapter we saw how to pretrain transformer language models, and how these pretrained models can be used as a tool for many kinds of NLP tasks, by casting the tasks as word prediction. The models we introduced in Chapter 10 to do this task are **causal or left-to-right transformer models**.

In this chapter we'll introduce a second paradigm for pretrained language models, called the **bidirectional transformer** encoder, trained via **masked language modeling**, a method that allows the model to see entire texts at a time, including both the right and left context. We'll introduce the most widely-used version of the masked language modeling architecture, the **BERT** model (Devlin et al., 2019).

We'll also introduce two important ideas that are often used with these masked language models. The first is the idea of **fine-tuning**. Fine-tuning is the process of taking the network learned by these pretrained models, and further training the model, often via an added neural net classifier that takes the top layer of the network as input, to perform some downstream task like named entity tagging or question answering or coreference. The intuition is that the pretraining phase learns a language model that instantiates rich representations of word meaning, that thus enables the model to more easily learn ('be fine-tuned to') the requirements of a downstream language understanding task. The pretrain-finetune paradigm is an instance of what is called **transfer learning** in machine learning: the method of acquiring knowledge from one task or domain, and then applying it (transferring it) to solve a new task.

The second idea that we introduce in this chapter is the idea of **contextual embeddings**: representations for words in context. The methods of Chapter 6 like word2vec or GloVe learned a single vector embedding for each unique word w in the vocabulary. By contrast, with contextual embeddings, such as those learned by masked language models like BERT, each word w will be represented by a different vector each time it appears in a different context. While the causal language models of Chapter 10 also use contextual embeddings, the embeddings created by masked language models **seem to function particularly well as representations**.

11.1 Bidirectional Transformer Encoders

Let's begin by introducing the bidirectional transformer encoder that underlies models like BERT and its descendants like **RoBERTa** (Liu et al., 2019) or **SpanBERT** (Joshi et al., 2020). In Chapter 10 we explored causal (left-to-right) transformers that can serve as the basis for powerful language models—models that can easily be applied to autoregressive generation problems such as contextual generation, summarization and machine translation. However, when applied to sequence classification and labeling problems causal models have obvious shortcomings since they

are based on an incremental, left-to-right processing of their inputs. If we want to assign the correct named-entity tag to each word in a sentence, or other sophisticated linguistic labels like the parse tags we'll introduce in later chapters, we'll want to be able to take into account information from the right context as we process each element. Fig. 11.1a, reproduced here from Chapter 10, illustrates the information flow in the purely left-to-right approach of Chapter 10. As can be seen, the hidden state computation at each point in time is based solely on the current and earlier elements of the input, ignoring potentially useful information located to the right of each tagging decision.

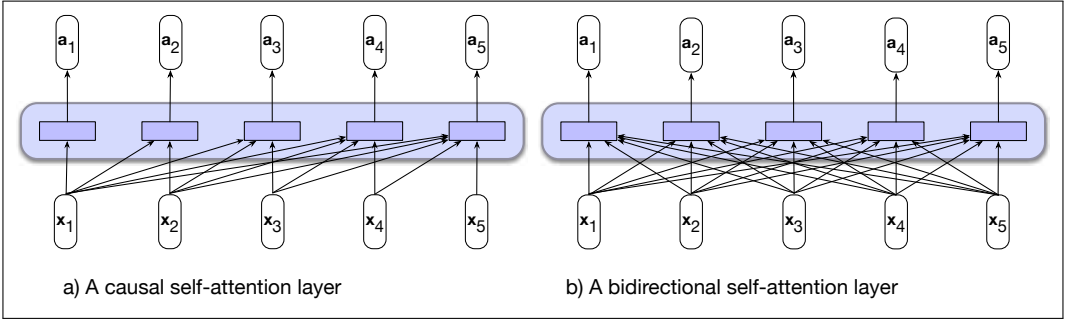


Figure 11.1 (a) The causal, backward looking, transformer model we saw in Chapter 10. Each output is computed independently of the others using only information seen earlier in the context. (b) Information flow in a bidirectional self-attention model. In processing each element of the sequence, the model attends to all inputs, both before and after the current one.

Bidirectional encoders overcome this limitation by allowing the self-attention mechanism to range over the entire input, as shown in Fig. 11.1b.

Why bidirectional encoders? The causal models of Chapter 10 are generative models, designed to easily generate the next token in a sequence. But the focus of bidirectional encoders is instead on computing contextualized representations of the input tokens. Bidirectional encoders use self-attention to map sequences of input embeddings (x_1, \dots, x_n) to sequences of output embeddings the same length (y_1, \dots, y_n) , where the output vectors have been contextualized using information from the entire input sequence. These output embeddings are contextualized representations of each input token that are generally useful across a range of downstream applications. The models of Chapter 10 are sometimes called decoder-only; the models of this chapter are sometimes called encoder-only, because they produce an encoding for each input token but generally aren't used to produce running text by decoding/sampling.

11.1.1 The architecture for bidirectional models

Bidirectional models use the same self-attention mechanism as causal models. The first step is to generate a set of key, query and value embeddings for each element of the input vector x through the use of learned weight matrices W^Q , W^K , and W^V . These weights project each input vector x_i into its specific role as a key, query, or value.

$$q_i = W^Q x_i; \quad k_i = W^K x_i; \quad v_i = W^V x_i \tag{11.1}$$

The output vector y_i corresponding to each input element x_i is a weighted sum of all

the input value vectors \mathbf{v} , as follows:

$$\mathbf{y}_i = \sum_{j=1}^n \alpha_{ij} \mathbf{v}_j \quad (11.2)$$

The α weights are computed via a softmax over the comparison scores between every element of an input sequence considered as a query and every other element as a key, where the comparison scores are computed using dot products.

$$\alpha_{ij} = \frac{\exp(\text{score}_{ij})}{\sum_{k=1}^n \exp(\text{score}_{ik})} \quad (11.3)$$

$$\text{score}_{ij} = \mathbf{q}_i \cdot \mathbf{k}_j \quad (11.4)$$

As with the models of Chapter 10, since each output vector, \mathbf{y}_i , is computed independently, the processing of an entire sequence can be parallelized via matrix operations. The first step is to pack the input embeddings \mathbf{x}_i into a matrix $\mathbf{X} \in \mathbb{R}^{N \times d_h}$. That is, each row of \mathbf{X} is the embedding of one token of the input. We then multiply \mathbf{X} by the key, query, and value weight matrices (all of dimensionality $d \times d$) to produce matrices $\mathbf{Q} \in \mathbb{R}^{N \times d}$, $\mathbf{K} \in \mathbb{R}^{N \times d}$, and $\mathbf{V} \in \mathbb{R}^{N \times d}$, containing all the key, query, and value vectors in a single step.

$$\mathbf{Q} = \mathbf{XW}^{\mathbf{Q}}; \mathbf{K} = \mathbf{XW}^{\mathbf{K}}; \mathbf{V} = \mathbf{XW}^{\mathbf{V}} \quad (11.5)$$

Given these matrices we can compute all the requisite query-key comparisons simultaneously by multiplying \mathbf{Q} and \mathbf{K}^T in a single operation. Fig. 11.2 illustrates the result of this operation for an input with length 5.

N	q1•k1	q1•k2	q1•k3	q1•k4	q1•k5
	q2•k1	q2•k2	q2•k3	q2•k4	q2•k5
	q3•k1	q3•k2	q3•k3	q3•k4	q3•k5
	q4•k1	q4•k2	q4•k3	q4•k4	q4•k5
	q5•k1	q5•k2	q5•k3	q5•k4	q5•k5
N					

Figure 11.2 The $N \times N$ \mathbf{QK}^T matrix showing the complete set of $q_i \cdot k_j$ comparisons.

Finally, we can scale these scores, take the softmax, and then multiply the result by \mathbf{V} resulting in a matrix of shape $N \times d$ where each row contains a contextualized output embedding corresponding to each token in the input.

$$\text{SelfAttention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{QK}^T}{\sqrt{d_k}}\right) \mathbf{V} \quad (11.6)$$

The key architecture difference is in bidirectional models we don't mask the future. As shown in Fig. 11.2, the full set of self-attention scores represented by \mathbf{QK}^T constitute an all-pairs comparison between the keys and queries for each element of the input. In the case of causal language models in Chapter 10, we masked the

upper triangular portion of this matrix (in Fig. ??) to eliminate information about future words since this would make the language modeling training task trivial. With bidirectional encoders we simply skip the mask, allowing the model to contextualize each token using *information from the entire input*.

Beyond this simple change, all of the other elements of the transformer architecture remain the same for bidirectional encoder models. Inputs to the model are segmented using subword tokenization and are combined with positional embeddings before being passed through a series of standard transformer blocks consisting of self-attention and feedforward layers augmented with residual connections and layer normalization, as shown in Fig. 11.3.

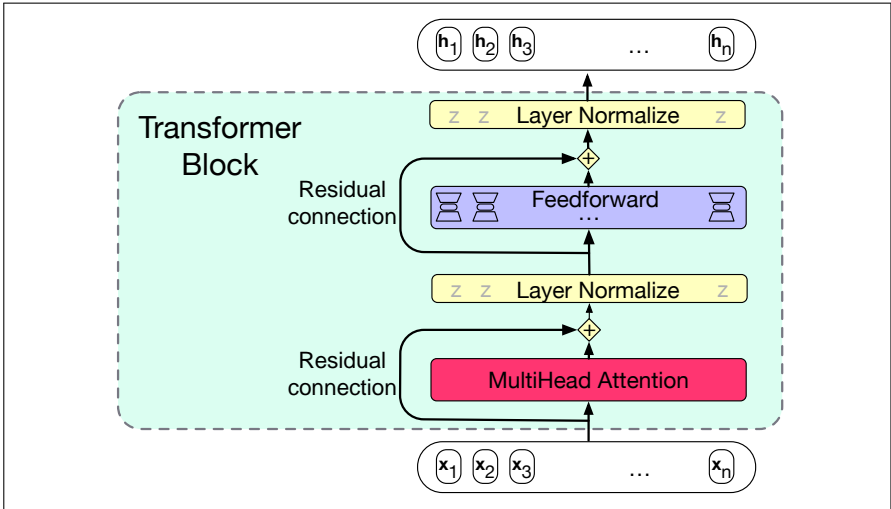


Figure 11.3 A transformer block showing all the layers.

To make this more concrete, the original English-only bidirectional transformer encoder model, BERT (Devlin et al., 2019), consisted of the following:

- An English-only subword vocabulary consisting of 30,000 tokens generated using the WordPiece algorithm (Schuster and Nakajima, 2012).
- Hidden layers of size of 768,
- 12 layers of transformer blocks, with 12 multihead attention layers each.
- The resulting model has about 100M parameters.

The larger multilingual XLM-RoBERTa model, trained on 100 languages, has

- A multilingual subword vocabulary with 250,000 tokens generated using the SentencePiece Unigram LM algorithm (Kudo and Richardson, 2018).
- 24 layers of transformer blocks, with 16 multihead attention layers each
- Hidden layers of size 1024
- The resulting model has about 550M parameters.

The use of WordPiece or SentencePiece Unigram LM tokenization (two of the large family of subword tokenization algorithms that includes the BPE algorithm we saw in Chapter 2) means that—like the large language models of Chapter 10—BERT and its descendants are based on subword tokens rather than words. Every input sentence first has to be tokenized, and then all further processing takes place on subword tokens rather than words. This will require, as we’ll see, that for some

NLP tasks that **require notions of words** (like named entity tagging, or parsing) we will occasionally need to map subwords back to words.

As with causal transformers, the size of the input layer dictates the complexity of the model. Both the time and memory requirements in a transformer grow quadratically with the length of the input. It's necessary, therefore, to set a fixed input length that is long enough to provide sufficient context for the model to function and yet still be computationally tractable. For BERT and XLR-RoBERTa, a fixed input size of 512 subword tokens was used.

11.2 Training Bidirectional Encoders

cloze task

We trained causal transformer language models in Chapter 10 by making them iteratively predict the next word in a text. But eliminating the causal mask makes the guess-the-next-word language modeling task trivial since the answer is now directly available from the context, so we're in need of a new training scheme. Fortunately, the traditional learning objective suggests an approach that can be used to train bidirectional encoders. Instead of trying to predict the next word, the model learns to perform a fill-in-the-blank task, technically called the **cloze task** (Taylor, 1953). To see this, let's return to the motivating example from Chapter 3. Instead of predicting which words are likely to come next in this example:

Please turn your homework ____ .

we're asked to predict a missing item given the rest of the sentence.

Please turn ____ homework in.

That is, given an input sequence with one or more elements missing, the learning task is to predict the missing elements. More precisely, during training the model is deprived of one or more elements of an input sequence and must generate a probability distribution over the vocabulary for each of the missing items. We then use the cross-entropy loss from each of the model's predictions to drive the learning process.

This approach can be generalized to any of a variety of methods that corrupt the training input and then asks the model to recover the original input. Examples of the kinds of manipulations that have been used include masks, substitutions, reorderings, deletions, and extraneous insertions into the training text.

11.2.1 Masking Words

Masked
Language
Modeling
MLM

The original approach to training bidirectional encoders is called **Masked Language Modeling** (MLM) (Devlin et al., 2019). As with the language model training methods we've already seen, **MLM** uses unannotated text from a large corpus. Here, the model is presented with a series of sentences from the training corpus where a random sample of tokens from each training sequence is selected for use in the learning task. Once chosen, a token is used in one of three ways:

- It is replaced with the unique vocabulary token [MASK].
- It is replaced with another token from the vocabulary, randomly sampled based on token unigram probabilities.
- It is left unchanged.

In BERT, 15% of the input tokens in a training sequence are sampled for learning. Of these, 80% are replaced with [MASK], 10% are replaced with randomly selected tokens, and the remaining 10% are left unchanged.

The MLM training objective is to predict the original inputs for each of the masked tokens using a bidirectional encoder of the kind described in the last section. The cross-entropy loss from these predictions drives the training process for all the parameters in the model. Note that all of the input tokens play a role in the self-attention process, but **only the sampled tokens are used for learning.**

More specifically, the original input sequence is first tokenized using a subword model. The sampled items which drive the learning process are chosen from among the set of tokenized inputs. Word embeddings for all of the tokens in the input are retrieved from the word embedding matrix and then combined with positional embeddings to form the input to the transformer.

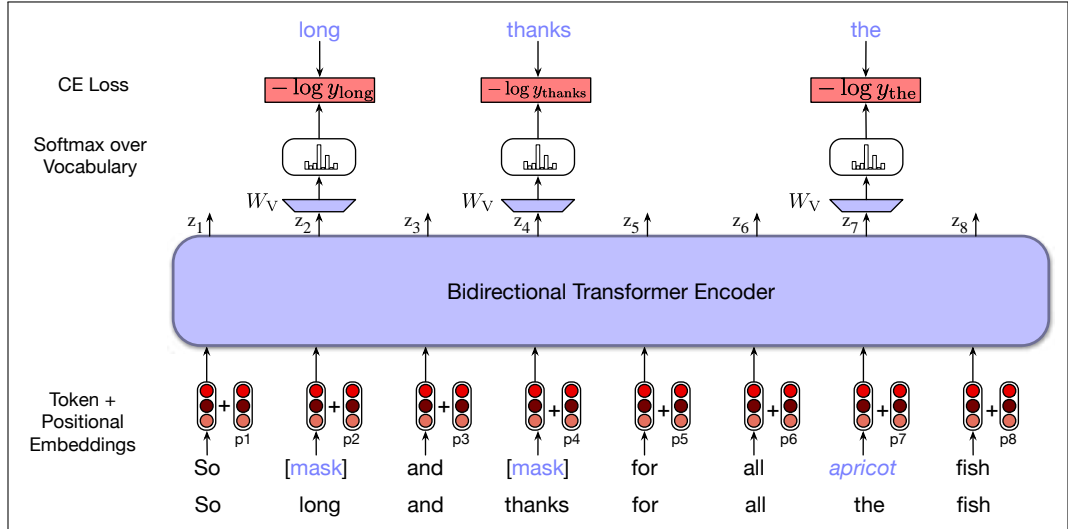


Figure 11.4 Masked language model training. In this example, three of the input tokens are selected, two of which are masked and the third is replaced with an unrelated word. The probabilities assigned by the model to these three items are used as the training loss. The other 5 words don't play a role in training loss. (In this and subsequent figures we display the input as words rather than subword tokens; the reader should keep in mind that BERT and similar models actually use subword tokens instead.)

Fig. 11.4 illustrates this approach with a simple example. Here, *long*, *thanks* and *the* have been sampled from the training sequence, with the first two masked and *the* replaced with the randomly sampled token *apricot*. The resulting embeddings are passed through a stack of bidirectional transformer blocks. To produce a probability distribution over the vocabulary for each of the masked tokens, the output vector z_i from the final transformer layer for each masked token i is multiplied by a learned set of classification weights $W_V \in \mathbb{R}^{|V| \times d_h}$ and then through a softmax to yield the required predictions over the vocabulary.

$$y_i = \text{softmax}(W_V z_i)$$

With a predicted probability distribution for each masked item, we can use cross-entropy to compute the loss for each masked item—the negative log probability assigned to the actual masked word, as shown in Fig. 11.4. More formally, for a given vector of input tokens in a sentence or batch be \mathbf{x} , let the set of tokens that are

masked be M , the version of that sentence with some tokens replaced by masks be \mathbf{x}^{mask} , and the sequence of output vectors be \mathbf{z} . For a given input token x_i , such as the word *long* in Fig. 11.4, the loss is the probability of the correct word *long*, given \mathbf{x}^{mask} (as summarized in the single output vector \mathbf{z}_i):

$$L_{MLM}(x_i) = -\log P(x_i|\mathbf{z}_i)$$

The gradients that form the basis for the weight updates are based on the average loss over the sampled learning items from a single training sequence (or batch of sequences).

$$L_{MLM} = -\frac{1}{|M|} \sum_{i \in M} \log P(x_i|\mathbf{z}_i)$$

Note that only the tokens in M play a role in learning; the other words play no role in the loss function, so in that sense BERT and its descendents are inefficient; only 15% of the input samples in the training data are actually used for training weights.¹

11.2.2 Next Sentence Prediction

The focus of mask-based learning is on predicting words from surrounding contexts with the goal of producing effective word-level representations. However, an important class of applications involves determining the relationship between pairs of sentences. These include tasks like paraphrase detection (detecting if two sentences have similar meanings), entailment (detecting if the meanings of two sentences entail or contradict each other) or discourse coherence (deciding if two neighboring sentences form a coherent discourse).

To capture the kind of knowledge required for applications such as these, some models in the BERT family include a second learning objective called **Next Sentence Prediction** (NSP). In this task, the model is presented with pairs of sentences and is asked to predict whether each pair consists of an actual pair of adjacent sentences from the training corpus or a pair of unrelated sentences. In BERT, 50% of the training pairs consisted of positive pairs, and in the other 50% the second sentence of a pair was randomly selected from elsewhere in the corpus. The NSP loss is based on how well the model can distinguish true pairs from random pairs.

To facilitate NSP training, BERT introduces two new tokens to the input representation (tokens that will prove useful for fine-tuning as well). After tokenizing the input with the subword model, the token [CLS] is prepended to the input sentence pair, and the token [SEP] is placed between the sentences and after the final token of the second sentence. Finally, embeddings representing the first and second segments of the input are added to the word and positional embeddings to allow the model to more easily distinguish the input sentences.

During training, the output vector from the final layer associated with the [CLS] token represents the next sentence prediction. As with the MLM objective, a learned set of classification weights $\mathbf{W}_{\text{NSP}} \in \mathbb{R}^{2 \times d_h}$ is used to produce a two-class prediction from the raw [CLS] vector.

$$y_i = \text{softmax}(\mathbf{W}_{\text{NSP}} h_i)$$

¹ There are members of the BERT family like ELECTRA that do use all examples for training (Clark et al., 2020).

Cross entropy is used to compute the NSP loss for each sentence pair presented to the model. Fig. 11.5 illustrates the overall NSP training setup. In BERT, the NSP loss was used in conjunction with the MLM training objective to form final loss.

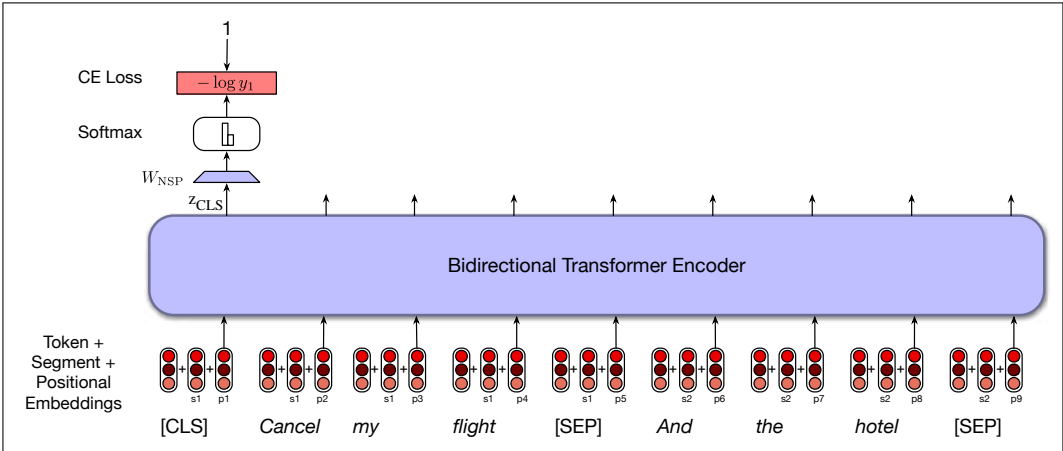


Figure 11.5 An example of the NSP loss calculation.

11.2.3 Training Regimes

BERT and other early transformer-based language models were trained on about 3.3 billion words (a combination of English Wikipedia and a corpus of book texts called BooksCorpus (Zhu et al., 2015) that is no longer used for intellectual property reasons). Modern masked language models are now trained on much larger datasets of web text, filtered a bit, and augmented by higher-quality data like Wikipedia, the same as those we discussed for the causal large language models of Chapter 10. Multilingual models similarly use webtext and multilingual Wikipedia. For example, the XLM-R model was trained on about 300 billion tokens in 100 languages, taken from the web via Common Crawl (<https://commoncrawl.org/>).

To train the original BERT models, pairs of text segments were selected from the training corpus according to the next sentence prediction 50/50 scheme. Pairs were sampled so that their combined length was less than the 512 token input. Tokens within these sentence pairs were then masked using the MLM approach with the combined loss from the MLM and NSP objectives used for a final loss. Approximately 40 passes (epochs) over the training data was required for the model to converge.

Some models, like the RoBERTa model, drop the next sentence prediction objective, and therefore change the training regime a bit. Instead of sampling pairs of sentence, the input is simply a series of contiguous sentences. If the document runs out before 512 tokens are reached, an extra separator token is added, and sentences from the next document are packed in, until we reach a total of 512 tokens. Usually large batch sizes are used, between 8K and 32K tokens.

Multilingual models have an additional decision to make: what data to use to build the vocabulary? Recall that all language models use subword tokenization (BPE or SentencePiece Unigram LM are the two most common algorithms). What text should be used to learn this multilingual tokenization, given that it's easier to get much more text in some languages than others? One option would be to create this vocabulary-learning dataset by sampling sentences from our training data (perhaps

web text from Common Crawl), randomly. In that case we will choose a lot of sentences from languages like languages with lots of web representation like English, and the tokens will be biased toward rare English tokens instead of creating frequent tokens from languages with less data. Instead, it is common to divide the training data into subcorpora of N different languages, compute the number of sentences n_i of each language i , and readjust these probabilities so as to upweight the probability of less-represented languages (Lample and Conneau, 2019). The new probability of selecting a sentence from each of the N languages (whose prior frequency is n_i) is $\{q_i\}_{i=1\dots N}$, where:

$$q_i = \frac{p_i^\alpha}{\sum_{j=1}^N p_j^\alpha} \quad \text{with} \quad p_i = \frac{n_i}{\sum_{k=1}^N n_k} \quad (11.7)$$

Recall from (??) in Chapter 6 that an α value between 0 and 1 will give higher weight to lower probability samples. Conneau et al. (2020) show that $\alpha = 0.3$ works well to give rare languages more inclusion in the tokenization, resulting in better multilingual performance overall.

The result of this pretraining process consists of both learned word embeddings, as well as all the parameters of the bidirectional encoder that are used to produce contextual embeddings for novel inputs.

For many purposes, a pretrained multilingual model is more practical than a monolingual model, since it avoids the need to build many (100!) separate monolingual models. And multilingual models can improve performance on low-resourced languages by leveraging linguistic information from a similar language in the training data that happens to have more resources. Nonetheless, when the number of languages grows very large, multilingual models exhibit what has been called the **curse of multilinguality** (Conneau et al., 2020): the performance on each language degrades compared to a model training on fewer languages. Another problem with multilingual models is that they ‘have an accent’: grammatical structures in higher-resource languages (often English) bleed into lower-resource languages; the vast amount of English language in training makes the model’s representations for low-resource languages slightly more English-like (Papadimitriou et al., 2023).

11.3 Contextual Embeddings

contextual embeddings

Given a pretrained language model and a novel input sentence, we can think of the sequence of model outputs as constituting **contextual embeddings** for each token in the input. These contextual embeddings are vectors representing some aspect of the meaning of a token in context, and can be used for any task requiring the meaning of tokens or words. More formally, given a sequence of input tokens x_1, \dots, x_n , we can use the output vector \mathbf{z}_i from the final layer of the model as a representation of the meaning of token x_i in the context of sentence x_1, \dots, x_n . Or instead of just using the vector \mathbf{z}_i from the final layer of the model, it’s common to compute a representation for x_i by averaging the output tokens \mathbf{z}_i from each of the last four layers of the model.

Just as we used static embeddings like word2vec in Chapter 6 to represent the meaning of words, we can use contextual embeddings as representations of word meanings in context for any task that might require a model of word meaning. Where static embeddings represent the meaning of word *types* (vocabulary entries), contextual embeddings represent the meaning of word *instances*: instances of a particular

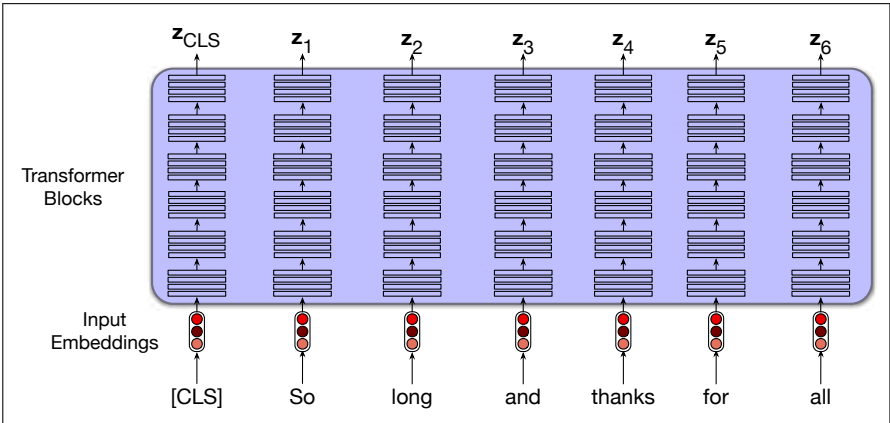


Figure 11.6 The output of a BERT-style model is a contextual embedding vector z_i for each input token x_i .

word type in a particular context. Thus where word2vec had a single vector for each word type, contextual embeddings provide a single vector for each instance of that word type in its sentential context. Contextual embeddings can thus be used for tasks like measuring the semantic similarity of two words in context, and are useful in linguistic tasks that require models of word meaning.

11.3.1 Contextual Embeddings and Word Sense

ambiguous Words are **ambiguous**: the same word can be used to mean different things. In Chapter 6 we saw that the word “mouse” can mean (1) a small rodent, or (2) a hand-operated device to control a cursor. The word “bank” can mean: (1) a financial institution or (2) a sloping mound. We say that the words ‘mouse’ or ‘bank’ are **polysemous** (from Greek ‘many senses’, *poly-* ‘many’ + *sema*, ‘sign, mark’).²

word sense A **sense** (or **word sense**) is a discrete representation of one aspect of the meaning of a word. We can represent each sense with a superscript: **bank**¹ and **bank**², **mouse**¹ and **mouse**². These senses can be found listed in online thesauruses (or thesauri) like **WordNet** (Fellbaum, 1998), which has datasets in many languages listing the senses of many words. In context, it’s easy to see the different meanings:

- mouse**¹ : a *mouse* controlling a computer system in 1968.
mouse² : a quiet animal like a *mouse*
bank¹ : ...a *bank* can hold the investments in a custodial account ...
bank² : ...as agriculture burgeons on the east *bank*, the river ...

This fact that context disambiguates the senses of *mouse* and *bank* above can also be visualized geometrically. Fig. 11.7 shows a two-dimensional project of many instances of the BERT embeddings of the word *die* in English and German. Each point in the graph represents the use of *die* in one input sentence. We can clearly see at least two different English senses of *die* (the singular of *dice* and the verb *to die*, as well as the German article, in the BERT embedding space.

Thus while thesauruses like WordNet give discrete lists of senses, embeddings (whether static or contextual) offer a continuous high-dimensional model of meaning

² The word **polysemy** itself is ambiguous; you may see it used in a different way, to refer only to cases where a word’s senses are related in some structured way, reserving the word **homonymy** to mean sense ambiguities with no relation between the senses (Haber and Poesio, 2020). Here we will use ‘polysemy’ to mean any kind of sense ambiguity, and ‘structured polysemy’ for polysemy with sense relations.

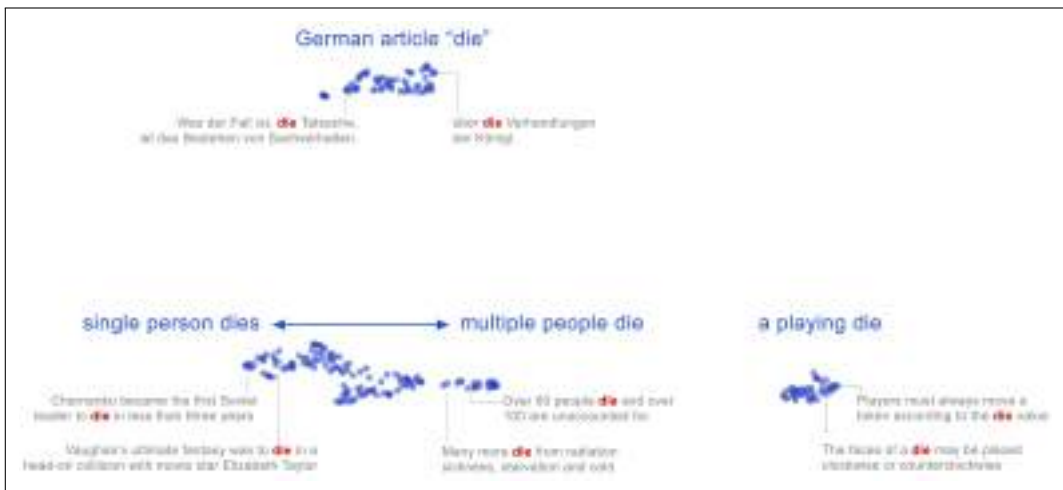


Figure 11.7 Each blue dot shows a BERT contextual embedding for the word *die* from different sentences in English and German, projected into two dimensions with the UMAP algorithm. The German and English meanings and the different English senses fall into different clusters. Some sample points are shown with the contextual sentence they came from. Figure from [Coenen et al. \(2019\)](#).

that, although it can be clustered, doesn't divide up into fully discrete senses.

Word Sense Disambiguation

word sense
disambiguation
WSD

The task of selecting the correct sense for a word is called **word sense disambiguation**, or **WSD**. WSD algorithms take as input a word in context and a fixed inventory of potential word senses (like the ones in WordNet) and outputs the correct word sense in context. Fig. 11.8 sketches out the task.

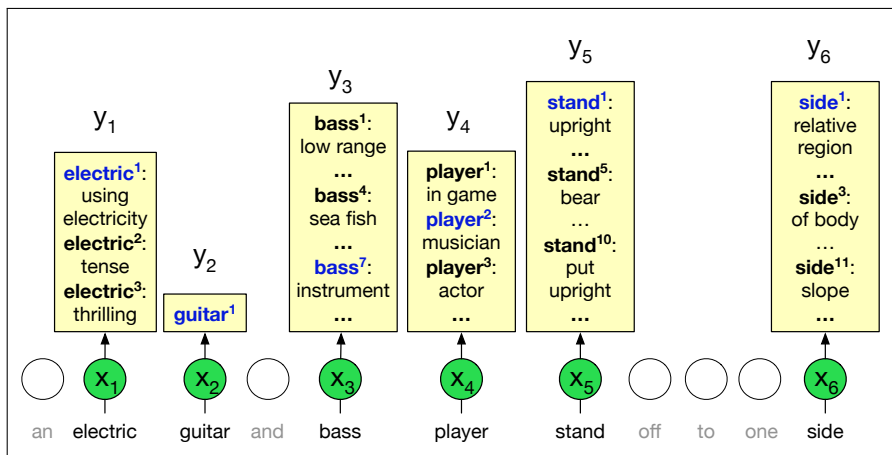


Figure 11.8 The all-words WSD task, mapping from input words (x) to WordNet senses (y). Figure inspired by [Chaplot and Salakhutdinov \(2018\)](#).

WSD can be a useful analytic tool for text analysis in the humanities and social sciences, and word senses can play a role in model interpretability for word representations. Word senses also have interesting distributional properties. For example a word often is used in roughly the same sense through a discourse, an observation called the **one sense per discourse** rule ([Gale et al., 1992](#)).

one sense per
discourse

The best performing WSD algorithm is a simple 1-nearest-neighbor algorithm using contextual word embeddings, due to [Melamud et al. \(2016\)](#) and [Peters et al. \(2018\)](#). At training time we pass each sentence in some sense-labeled dataset (like the SemCore or SenseEval datasets in various languages) through any contextual embedding (e.g., BERT) resulting in a contextual embedding for each labeled token. (There are various ways to compute this contextual embedding v_i for a token i ; for BERT it is common to pool multiple layers by summing the vector representations of i from the last four BERT layers). Then for each sense s of any word in the corpus, for each of the n tokens of that sense, we average their n contextual representations v_i to produce a contextual **sense embedding** \mathbf{v}_s for s :

$$\mathbf{v}_s = \frac{1}{n} \sum_i \mathbf{v}_i \quad \forall \mathbf{v}_i \in \text{tokens}(s) \tag{11.8}$$

At test time, given a token of a target word t in context, we compute its contextual embedding \mathbf{t} and choose its nearest neighbor sense from the training set, i.e., the sense whose sense embedding has the highest cosine with \mathbf{t} :

$$\text{sense}(t) = \underset{s \in \text{senses}(t)}{\operatorname{argmax}} \operatorname{cosine}(\mathbf{t}, \mathbf{v}_s) \tag{11.9}$$

Fig. 11.9 illustrates the model.

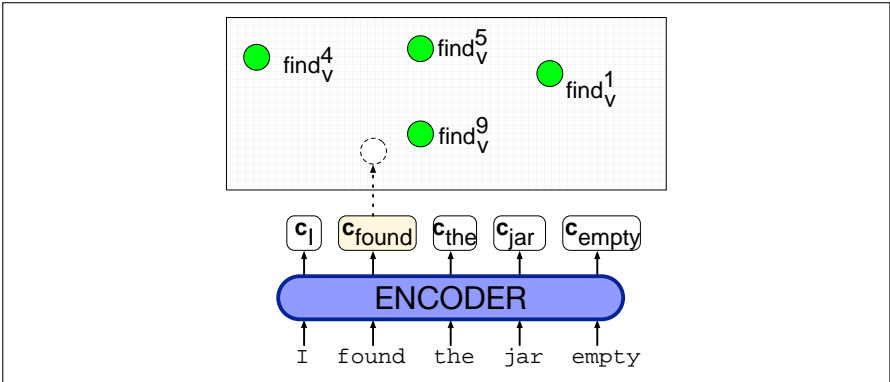


Figure 11.9 The nearest-neighbor algorithm for WSD. In green are the contextual embeddings precomputed for each sense of each word; here we just show a few of the senses for *find*. A contextual embedding is computed for the target word *found*, and then the nearest neighbor sense (in this case **find_v⁹**) is chosen. Figure inspired by [Loureiro and Jorge \(2019\)](#).

11.3.2 Contextual Embeddings and Word Similarity

In Chapter 6 we introduced the idea that we could measure the similarity of two words by considering how close they are geometrically, by using the cosine as a similarity function. The idea of meaning similarity is also clear geometrically in the meaning clusters in Fig. 11.7; the representation of a word which has a particular sense in a context is closer to other instances of the same sense of the word. Thus we often measure the similarity between two instances of two words in context (or two instances of the same word in two different contexts) by using the cosine between their contextual embeddings.

Usually some transformations to the embeddings are required before computing cosine. This is because contextual embeddings (whether from masked language

anisotropy

models or from autoregressive ones) have the property that the vectors for all words are extremely similar. If we look at the embeddings from the final layer of BERT or other models, embeddings for instances of any two randomly chosen words will have extremely high cosines that can be quite close to 1, meaning all word vectors tend to point in the same direction. The property of vectors in a system all tending to point in the same direction is known as **anisotropy**. [Ethayarajh \(2019\)](#) defines the **anisotropy** of a model as the expected cosine similarity of any pair of words in a corpus. The word ‘isotropy’ means uniformity in all directions, so in an isotropic model, the collection of vectors should point in all directions and the expected cosine between a pair of random embeddings would be zero. [Timkey and van Schijndel \(2021\)](#) show that one cause of anisotropy is that cosine measures are dominated by a small number of dimensions of the contextual embedding whose values are very different than the others: these **rogue dimensions** have very large magnitudes and very high variance.

[Timkey and van Schijndel \(2021\)](#) shows that we can make the embeddings more isotropic by standardizing (z-scoring) the vectors, i.e., subtracting the mean and dividing by the variance. Given a set C of all the embeddings in some corpus, each with dimensionality d (i.e., $\mathbf{x} \in \mathbb{R}^d$), the mean vector $\mu \in \mathbb{R}^d$ is:

$$\mu = \frac{1}{|C|} \sum_{\mathbf{x} \in C} \mathbf{x} \quad (11.10)$$

The standard deviation in each dimension $\sigma \in \mathbb{R}^d$ is:

$$\sigma = \sqrt{\frac{1}{|C|} \sum_{\mathbf{x} \in C} (\mathbf{x} - \mu)^2} \quad (11.11)$$

Then each word vector \mathbf{x} is replaced by a standardized version \mathbf{z} :

$$\mathbf{z} = \frac{\mathbf{x} - \mu}{\sigma} \quad (11.12)$$

One problem with cosine that is not solved by standardization is that cosine tends to underestimate human judgments on similarity of word meaning for very frequent words ([Zhou et al., 2022](#)).

In the next section we’ll see the most common use of contextual representations: as representations of words or even entire sentences that can be the inputs to classifiers in the fine-tuning process for downstream NLP applications.

11.4 Fine-Tuning Language Models

fine-tuning

The power of pretrained language models lies in their ability to extract generalizations from large amounts of text—generalizations that are useful for myriad downstream applications. There are two ways to make practical use of the generalizations. One way is to use natural language to **prompt** the model, putting it in a state where it contextually generates what we want. We’ll introduce prompting in Chapter 12. An alternative is to create interfaces from pretrained language models to downstream applications through a process called **fine-tuning**. In fine-tuning, we create applications on top of pretrained models by adding a small set of application-specific parameters. The fine-tuning process consists of using labeled data about

the application to train these additional application-specific parameters. Typically, this training will either freeze or make only minimal adjustments to the pretrained language model parameters.

The following sections introduce fine-tuning methods for the most common applications including sequence classification, sequence labeling, sentence-pair inference, and span-based operations.

11.4.1 Sequence Classification

Sequence classification applications often represent an input sequence with a single consolidated representation. With RNNs, we used the hidden layer associated with the final input element to stand for the entire sequence. A similar approach is used with transformers. An additional vector is added to the model to stand for the entire sequence. This vector is sometimes called the **sentence embedding** since it refers to the entire sequence, although the term ‘sentence embedding’ is also used in other ways. In BERT, the [CLS] token plays the role of this embedding. This unique token is added to the vocabulary and is prepended to the start of all input sequences, both during pretraining and encoding. The output vector in the final layer of the model for the [CLS] input represents the entire input sequence and serves as the input to a **classifier head**, a logistic regression or neural network classifier that makes the relevant decision.

As an example, let’s return to the problem of sentiment classification. A simple approach to fine-tuning a classifier for this application involves learning a set of weights, \mathbf{W}_C , to map the output vector for the [CLS] token— \mathbf{z}_{CLS} —to a set of scores over the possible sentiment classes. Assuming a three-way sentiment classification task (positive, negative, neutral) and dimensionality d_h for the size of the language model hidden layers gives $\mathbf{W}_C \in \mathbb{R}^{3 \times d_h}$. Classification of unseen documents proceeds by passing the input text through the pretrained language model to generate \mathbf{z}_{CLS} , multiplying it by \mathbf{W}_C , and finally passing the resulting vector through a softmax.

$$\mathbf{y} = \text{softmax}(\mathbf{W}_C \mathbf{z}_{CLS}) \tag{11.13}$$

Finetuning the values in \mathbf{W}_C requires supervised training data consisting of input sequences labeled with the appropriate class. Training proceeds in the usual way; cross-entropy loss between the softmax output and the correct answer is used to drive the learning that produces \mathbf{W}_C .

A key difference from what we’ve seen earlier with neural classifiers is that this loss can be used to not only learn the weights of the classifier, but also to update the weights for the pretrained language model itself. In practice, reasonable classification performance is typically achieved with only minimal changes to the language model parameters, often limited to updates over the final few layers of the transformer. Fig. 11.10 illustrates this overall approach to sequence classification.

11.4.2 Pair-Wise Sequence Classification

As mentioned in Section 11.2.2, an important type of problem involves the classification of pairs of input sequences. Practical applications that fall into this class include paraphrase detection (are the two sentences paraphrases of each other?), logical entailment (does sentence A logically entail sentence B?), and discourse coherence (how coherent is sentence B as a follow-on to sentence A?).

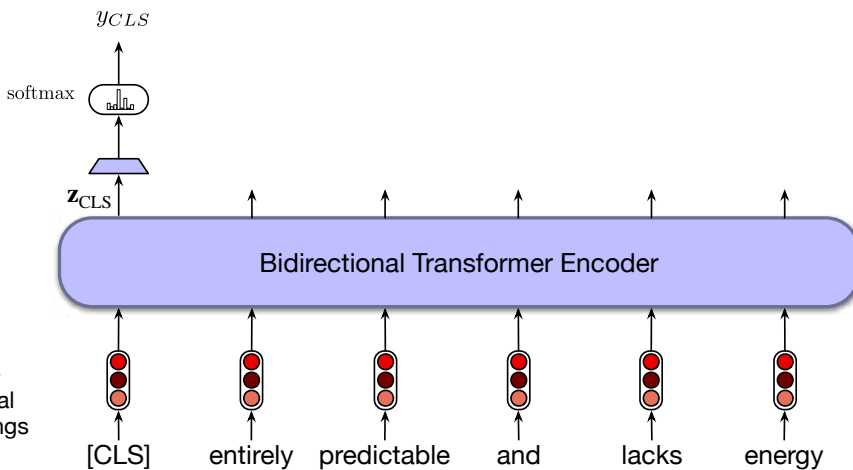


Figure 11.10 Sequence classification with a bidirectional transformer encoder. The output vector for the [CLS] token serves as input to a simple classifier.

Fine-tuning an application for one of these tasks proceeds just as with pretraining using the NSP objective. During fine-tuning, pairs of labeled sentences from the supervised training data are presented to the model, and run through all the layers of the model to produce the \mathbf{z} outputs for each input token. As with sequence classification, the output vector associated with the prepended [CLS] token represents the model’s view of the input pair. And as with NSP training, the two inputs are separated by the [SEP] token. To perform classification, the [CLS] vector is multiplied by a set of learning classification weights and passed through a softmax to generate label predictions, which are then used to update the weights.

As an example, let’s consider an entailment classification task with the Multi-Genre Natural Language Inference (MultiNLI) dataset (Williams et al., 2018). In the task of **natural language inference** or **NLI**, also called **recognizing textual entailment**, a model is presented with a pair of sentences and must classify the relationship between their meanings. For example in the MultiNLI corpus, pairs of sentences are given one of 3 labels: *entails*, *contradicts* and *neutral*. These labels describe a relationship between the meaning of the first sentence (the premise) and the meaning of the second sentence (the hypothesis). Here are representative examples of each class from the corpus:

- **Neutral**
 - a: Jon walked back to the town to the smithy.
 - b: Jon traveled back to his hometown.
- **Contradicts**
 - a: Tourist Information offices can be very helpful.
 - b: Tourist Information offices are never of any help.
- **Entails**
 - a: I’m confused.
 - b: Not all of it is very clear to me.

A relationship of *contradicts* means that the premise contradicts the hypothesis; *entails* means that the premise entails the hypothesis; *neutral* means that neither is necessarily true. The meaning of these labels is looser than strict logical entailment

or contradiction indicating that a typical human reading the sentences would most likely interpret the meanings in this way.

To fine-tune a classifier for the MultiNLI task, we pass the premise/hypothesis pairs through a bidirectional encoder as described above and use the output vector for the [CLS] token as the input to the classification head. As with ordinary sequence classification, this head provides the input to a three-way classifier that can be trained on the MultiNLI training corpus.

11.4.3 Sequence Labelling

Sequence labelling tasks, such as part-of-speech tagging or BIO-based named entity recognition, follow the same basic classification approach. Here, the final output vector corresponding to *each input token* is passed to a classifier that produces a softmax distribution over the possible set of tags. Again, assuming a simple classifier consisting of a single feedforward layer followed by a softmax, the set of weights to be learned for this additional layer is $\mathbf{W_K} \in \mathbb{R}^{k \times d_h}$, where k is the number of possible tags for the task. As with RNNs, a greedy approach, where the argmax tag for each token is taken as a likely answer, can be used to generate the final output tag sequence. Fig. 11.11 illustrates an example of this approach.

y_i = softmax(W_K z_i) (11.14)

t_i = argmax_k(y_i) (11.15)

Alternatively, the distribution over labels provided by the softmax for each input token can be passed to a conditional random field (CRF) layer which can take global tag-level transitions into account.

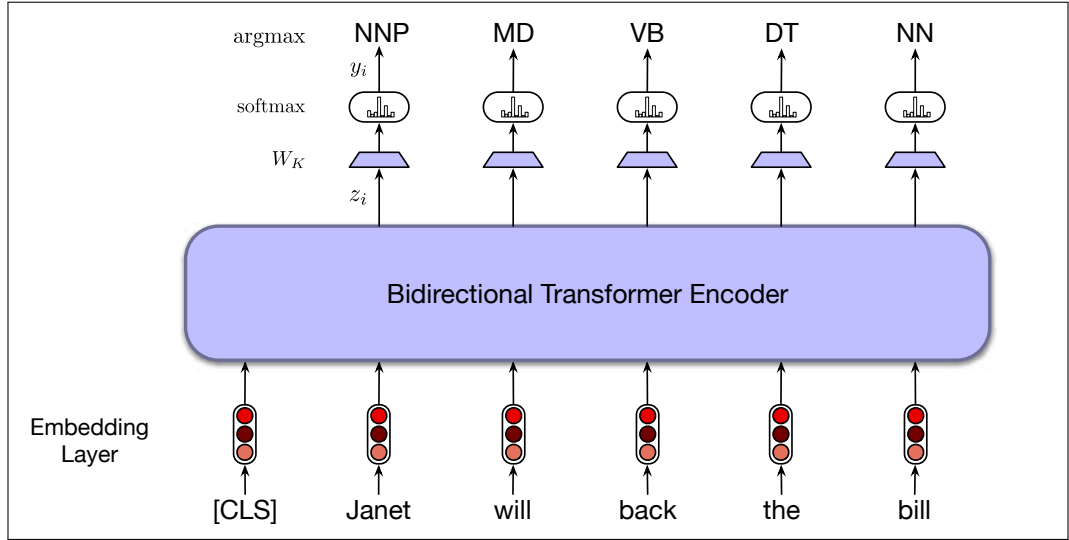


Figure 11.11 Sequence labeling for part-of-speech tagging with a bidirectional transformer encoder. The output vector for each input token is passed to a simple k-way classifier.

A complication with this approach arises from the use of subword tokenization such as WordPiece, SentencePiece Unigram LM or Byte Pair Encoding. Supervised training data for tasks like named entity recognition (NER) is typically in the form of BIO tags associated with text segmented at the word level. For example the following sentence containing two named entities:

[[LOC Mt. Sanitas](#)] is in [[LOC Sunshine Canyon](#)] .
 would have the following set of per-word BIO tags.

(11.16) *Mt. Sanitas is in Sunshine Canyon .*
 B-LOC I-LOC O O B-LOC I-LOC O

Unfortunately, the WordPiece tokenization for this sentence yields the following sequence of tokens which doesn't align directly with BIO tags in the ground truth annotation:

'Mt', '.', 'San', '##itas', 'is', 'in', 'Sunshine', 'Canyon' '.'

To deal with this misalignment, we need a way to assign BIO tags to subword tokens during training and a corresponding way to recover word-level tags from subwords during decoding. For training, we can just assign the gold-standard tag associated with each word to all of the subword tokens derived from it.

For decoding, the simplest approach is to use the argmax BIO tag associated with the first subword token of a word. Thus, in our example, the BIO tag assigned to "Mt" would be assigned to "Mt." and the tag assigned to "San" would be assigned to "Sanitas", effectively ignoring the information in the tags assigned to "." and "##itas". More complex approaches combine the distribution of tag probabilities across the subwords in an attempt to find an optimal word-level tag.

11.5 Advanced: Span-based Masking

For many NLP applications, the natural unit of interest may be larger than a single word (or token). Question answering, syntactic parsing, coreference and semantic role labeling applications all involve the identification and classification of longer phrases. This suggests that a span-oriented masked learning objective might provide improved performance on such tasks.

11.5.1 Masking Spans

A span is a contiguous sequence of one or more words selected from a training text, prior to subword tokenization. In span-based masking, a set of randomly selected spans from a training sequence are chosen. In the SpanBERT work that originated this technique ([Joshi et al., 2020](#)), a span length is first chosen by sampling from a geometric distribution that is biased towards shorter spans and with an upper bound of 10. Given this span length, a starting location consistent with the desired span length and the length of the input is sampled uniformly.

Once a span is chosen for masking, all the tokens within the span are substituted according to the same regime used in BERT: 80% of the time the span elements are substituted with the [MASK] token, 10% of the time they are replaced by randomly sampled tokens from the vocabulary, and 10% of the time they are left as is. Note that this substitution process is done at the span level—all the tokens in a given span are substituted using the same method. As with BERT, the total token substitution is limited to 15% of the training sequence input. Having selected and masked the training span, the input is passed through the standard transformer architecture to generate contextualized representations of the input tokens.

Downstream span-based applications rely on span representations derived from the tokens within the span, as well as the start and end points, or the boundaries, of a span. Representations for these boundaries are typically derived from the first and last tokens of a span, the tokens immediately preceding and following the span, or some combination of them. The SpanBERT learning objective augments the MLM objective with a boundary oriented component called the Span Boundary Objective (SBO). The SBO relies on a model’s ability to predict the tokens within a masked span from the tokens immediately preceding and following the span.

Let the sequence of output from the transformer encoder for the n input tokens s_1, \dots, x_n be z_1, \dots, z_n . A token x_i in a masked span of tokens (x_s, \dots, x_e) , i.e., starting with token x_s and ending with token x_e , is represented by concatenating 3 embeddings. The first two are the embeddings of two external boundary tokens x_{s-1} and x_{e+1} , i.e., the token preceding x_s , the token following x_e . The third embedding that is concatenated is the relative position embedding of the target token \mathbf{p}_{i-s+1} . The position embeddings p_1, p_2, \dots represent relative positions of the tokens with respect to the left boundary token x_{s-1} .

$$L(x) = L_{MLM}(x) + L_{SBO}(x) \tag{11.17}$$

$$L_{SBO}(x_i) = -\log P(x_i | x_{s-1}, x_{e+1}, p_{i-s+1}) \tag{11.18}$$

This probability for token x_i is formed by passing the concatenation of these embeddings through a 2-layer feedforward network to get the probability distribution over the whole vocabulary at i :

$$\mathbf{s}_i = \text{FFN}([\mathbf{z}_{s-1}; \mathbf{z}_{e+1}; \mathbf{p}_{i-s+1}]) \tag{11.19}$$

$$\mathbf{y}_i = \text{softmax}(\mathbf{W}_V \mathbf{s}_i) \tag{11.20}$$

We then use \mathbf{s}_i , the output of the vector representation of token i in the span, to predict the token x_i by reshaping it and passing it through a softmax to get a probability distribution \mathbf{y}_i over the vocabulary, and select from it the probability for input token x_i .

The final loss is the sum of the BERT MLM loss and the SBO loss.

Fig. 11.12 illustrates this with one of our earlier examples. Here the span selected is *and thanks for* which spans from position 3 to 5. The total loss associated with the masked token *thanks* is the sum of the cross-entropy loss generated from the prediction of *thanks* from the output \mathbf{z}_4 , plus the cross-entropy loss from the prediction of *thanks* from the output vectors from the left external boundary \mathbf{z}_2 , the right external boundary \mathbf{z}_6 , and the embedding for relative position 2 in the span.

11.5.2 Fine-tuning for Span-Based Applications

Span-oriented applications operate in a middle ground between sequence level and token level tasks. That is, in span-oriented applications the focus is on generating and operating with representations of contiguous sequences of tokens. Typical operations include identifying spans of interest, classifying spans according to some labeling scheme, and determining relations among discovered spans. Applications include named entity recognition, question answering, syntactic parsing, semantic role labeling and coreference resolution.

Formally, given an input sequence x consisting of T tokens, (x_1, x_2, \dots, x_T) , a span is a contiguous sequence of tokens with start i and end j such that $1 \leq i \leq j \leq T$. This formulation results in a total set of spans equal to $\frac{T(T+1)}{2}$. For practical

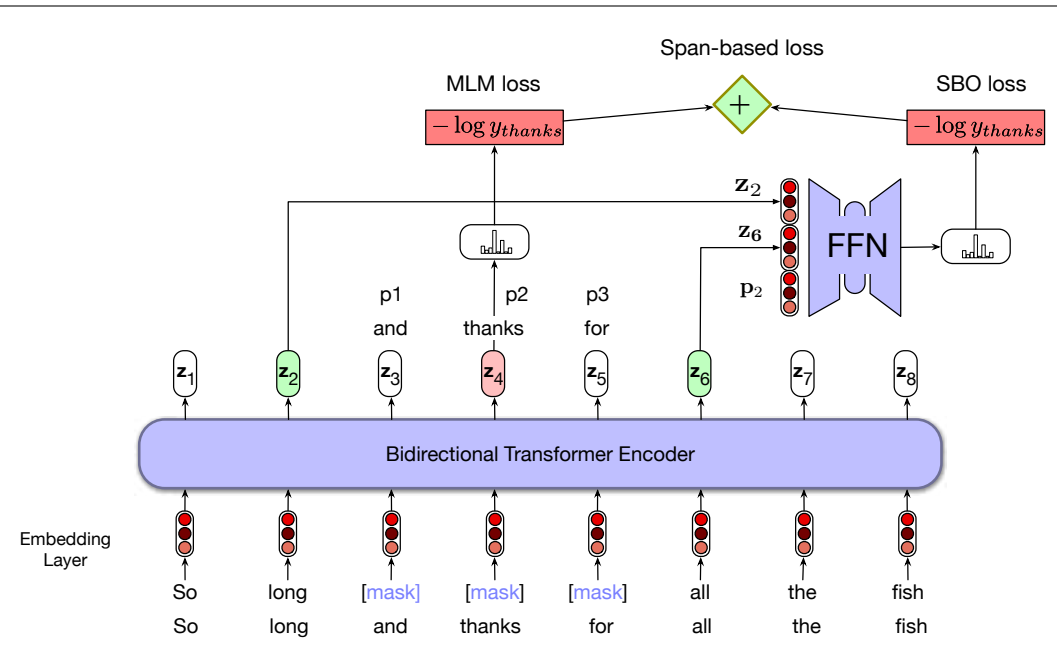


Figure 11.12 Span-based language model training. In this example, a span of length 3 is selected for training and all of the words in the span are masked. The figure illustrates the loss computed for word *thanks*; the loss for the entire span is the sum of the loss for the three words in the span.

purposes, span-based models often impose an application-specific length limit L , so the legal spans are limited to those where $j - i < L$. In the following, we'll refer to the enumerated set of legal spans in x as $S(x)$.

The first step in fine-tuning a pretrained language model for a span-based application is using the contextualized input embeddings from the model to generate representations for all the spans in the input. Most schemes for representing spans make use of two primary components: representations of the span boundaries and summary representations of the contents of each span. To compute a unified span representation, we concatenate the boundary representations with the summary representation.

In the simplest possible approach, we can use the contextual embeddings of the start and end tokens of a span as the boundaries, and the average of the output embeddings within the span as the summary representation.

$$\mathbf{g}_{ij} = \frac{1}{(j-i)+1} \sum_{k=i}^j \mathbf{z}_k \quad (11.21)$$

$$\text{spanRep}_{ij} = [\mathbf{z}_i; \mathbf{z}_j; \mathbf{g}_{i,j}] \quad (11.22)$$

A weakness of this approach is that it doesn't distinguish the use of a word's embedding as the beginning of a span from its use as the end of one. Therefore, more elaborate schemes for representing the span boundaries involve learned representations for start and end points through the use of two distinct feedforward networks:

$$\mathbf{s}_i = \text{FFN}_{\text{start}}(\mathbf{z}_i) \quad (11.23)$$

$$\mathbf{e}_j = \text{FFN}_{\text{end}}(\mathbf{z}_j) \quad (11.24)$$

$$\text{spanRep}_{ij} = [\mathbf{s}_i; \mathbf{e}_j; \mathbf{g}_{i,j}] \quad (11.25)$$

Similarly, a simple average of the vectors in a span is unlikely to be an optimal representation of a span since it treats all of a span’s embeddings as equally important. For many applications, a more useful representation would be centered around the head of the phrase corresponding to the span. One method for getting at such information in the absence of a syntactic parse is to use a standard self-attention layer to generate a span representation.

$$\mathbf{g}_{ij} = \text{SelfAttention}(\mathbf{z}_{i:j}) \tag{11.26}$$

Now, given span representations \mathbf{g} for each span in $S(x)$, classifiers can be fine-tuned to generate application-specific scores for various span-oriented tasks: binary span identification (is this a legitimate span of interest or not?), span classification (what kind of span is this?), and span relation classification (how are these two spans related?).

To ground this discussion, let’s return to named entity recognition (NER). Given a scheme for representing spans and a set of named entity types, a span-based approach to NER is a straightforward classification problem where each span in an input is assigned a class label. More formally, given an input sequence x_1, \dots, x_n , we want to assign a label y , from the set of valid NER labels, to each of the spans in $S(x)$. Since most of the spans in a given input will not be named entities we’ll add the label NULL to the set of types in Y .

$$y_{ij} = \text{softmax}(\text{FFN}(\text{spanRep}_{ij})) \tag{11.27}$$

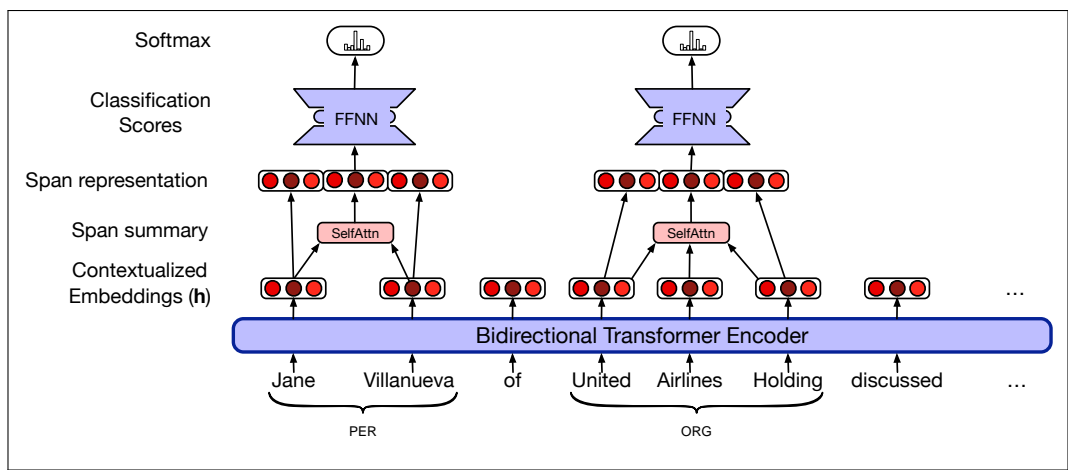


Figure 11.13 A span-oriented approach to named entity classification. The figure only illustrates the computation for 2 spans corresponding to ground truth named entities. In reality, the network scores all of the $\frac{T(T+1)}{2}$ spans in the text. That is, all the unigrams, bigrams, trigrams, etc. up to the length limit.

With this approach, fine-tuning entails using supervised training data to learn the parameters of the final classifier, as well as the weights used to generate the boundary representations, and the weights in the self-attention layer that generates the span content representation. During training, the model’s predictions for all spans are compared to their gold-standard labels and cross-entropy loss is used to drive the training.

During decoding, each span is scored using a softmax over the final classifier output to generate a distribution over the possible labels, with the argmax score for each span taken as the correct answer. Fig. 11.13 illustrates this approach with an

example. A variation on this scheme designed to improve precision adds a calibrated threshold to the labeling of a span as anything other than NULL.

There are two significant advantages to a span-based approach to NER over a BIO-based per-word labeling approach. The first advantage is that BIO-based approaches are prone to a labeling mis-match problem. That is, every label in a longer named entity must be correct for an output to be judged correct. Returning to the example in Fig. 11.13, the following labeling would be judged entirely wrong due to the incorrect label on the first item. Span-based approaches only have to make one classification for each span.

(11.28) *Jane Villanueva of United Airlines Holding discussed ...*
 B-PER I-PER O I-ORG I-ORG I-ORG O

The second advantage to span-based approaches is that they naturally accommodate embedded named entities. For example, in this example both *United Airlines* and *United Airlines Holding* are legitimate named entities. The BIO approach has no way of encoding this embedded structure. But the span-based approach can naturally label both since the spans are labeled separately.

11.6 Summary

This chapter has introduced the topic of transfer learning from pretrained language models. Here's a summary of the main points that we covered:

- Bidirectional encoders can be used to generate contextualized representations of input embeddings using the entire input context.
- Pretrained language models based on bidirectional encoders can be learned using a masked language model objective where a model is trained to guess the missing information from an input.
- Pretrained language models can be fine-tuned for specific applications by adding lightweight classifier layers on top of the outputs of the pretrained model.

Bibliographical and Historical Notes

- Chaplot, D. S. and R. Salakhutdinov. 2018. Knowledge-based word sense disambiguation using topic models. *AAAI*.
- Clark, K., M.-T. Luong, Q. V. Le, and C. D. Manning. 2020. [Electra: Pre-training text encoders as discriminators rather than generators](#). *ICLR*.
- Coenen, A., E. Reif, A. Yuan, B. Kim, A. Pearce, F. Viégas, and M. Wattenberg. 2019. Visualizing and measuring the geometry of bert. *NeurIPS*.
- Conneau, A., K. Khandelwal, N. Goyal, V. Chaudhary, G. Wenzek, F. Guzmán, E. Grave, M. Ott, L. Zettlemoyer, and V. Stoyanov. 2020. [Unsupervised cross-lingual representation learning at scale](#). *ACL*.
- Devlin, J., M.-W. Chang, K. Lee, and K. Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). *NAACL HLT*.
- Ethayarajh, K. 2019. [How contextual are contextualized word representations? Comparing the geometry of BERT, ELMo, and GPT-2 embeddings](#). *EMNLP*.
- Fellbaum, C., editor. 1998. *WordNet: An Electronic Lexical Database*. MIT Press.
- Gale, W. A., K. W. Church, and D. Yarowsky. 1992. [One sense per discourse](#). *HLT*.
- Haber, J. and M. Poesio. 2020. [Assessing polyseme sense similarity through co-predication acceptability and contextualised embedding distance](#). **SEM*.
- Joshi, M., D. Chen, Y. Liu, D. S. Weld, L. Zettlemoyer, and O. Levy. 2020. [SpanBERT: Improving pre-training by representing and predicting spans](#). *TACL*, 8:64–77.
- Kudo, T. and J. Richardson. 2018. [SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing](#). *EMNLP*.
- Lample, G. and A. Conneau. 2019. Cross-lingual language model pretraining. *NeurIPS*, volume 32.
- Liu, Y., M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov. 2019. [RoBERTa: A robustly optimized BERT pretraining approach](#). ArXiv preprint arXiv:1907.11692.
- Loureiro, D. and A. Jorge. 2019. [Language modelling makes sense: Propagating representations through WordNet for full-coverage word sense disambiguation](#). *ACL*.
- Melamud, O., J. Goldberger, and I. Dagan. 2016. [context2vec: Learning generic context embedding with bidirectional LSTM](#). *CoNLL*.
- Papadimitriou, I., K. Lopez, and D. Jurafsky. 2023. [Multilingual BERT has an accent: Evaluating English influences on fluency in multilingual models](#). *EACL Findings*.
- Peters, M., M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, and L. Zettlemoyer. 2018. [Deep contextualized word representations](#). *NAACL HLT*.
- Schuster, M. and K. Nakajima. 2012. [Japanese and Korean voice search](#). *ICASSP*.
- Taylor, W. L. 1953. Cloze procedure: A new tool for measuring readability. *Journalism Quarterly*, 30:415–433.
- Timkey, W. and M. van Schijndel. 2021. [All bark and no bite: Rogue dimensions in transformer language models obscure representational quality](#). *EMNLP*.
- Williams, A., N. Nangia, and S. Bowman. 2018. [A broad-coverage challenge corpus for sentence understanding through inference](#). *NAACL HLT*.
- Zhou, K., K. Ethayarajh, D. Card, and D. Jurafsky. 2022. [Problems with cosine as a measure of embedding similarity for high frequency words](#). *ACL*.
- Zhu, Y., R. Kiros, R. Zemel, R. Salakhutdinov, R. Urtasun, A. Torralba, and S. Fidler. 2015. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. *IEEE International Conference on Computer Vision*.