

# What does BERT learn about the structure of language?

Ganesh Jawahar    Benoît Sagot    Djamé Seddah  
Inria, France  
{firstname.lastname}@inria.fr

## Abstract

BERT is a recent language representation model that has surprisingly performed well in diverse language understanding benchmarks. This result indicates the possibility that BERT networks capture structural information about language. In this work, we provide novel support for this claim by performing a series of experiments to unpack the elements of English language structure learned by BERT. We first show that BERT's phrasal representation captures phrase-level information in the lower layers. We also show that BERT's intermediate layers encode a rich hierarchy of linguistic information, with surface features at the bottom, syntactic features in the middle and semantic features at the top. BERT turns out to require deeper layers when long-distance dependency information is required, e.g. to track subject-verb agreement. Finally, we show that BERT representations capture linguistic information in a compositional way that mimics classical, tree-like structures.

## 1 Introduction

BERT (Bidirectional Encoder Representations from Transformers) (Devlin et al., 2018) is a bidirectional variant of Transformer networks (Vaswani et al., 2017) trained to jointly predict a masked word from its context and to classify whether two sentences are consecutive or not. The trained model can be fine-tuned for downstream NLP tasks such as question answering and language inference without substantial modification. BERT outperforms previous state-of-the-art models in the eleven NLP tasks in the GLUE benchmark (Wang et al., 2018) by a significant margin. This remarkable result suggests that BERT could "learn" structural information about language.

Can we unveil the representations learned by BERT to proto-linguistics structures? Answering this question could not only help us understand

the reason behind the success of BERT but also its limitations, in turn guiding the design of improved architectures. This question falls under the topic of the interpretability of neural networks, a growing field in NLP (Belinkov and Glass, 2019). An important step forward in this direction is Goldberg (2019), which shows that BERT captures syntactic phenomena well when evaluated on its ability to track subject-verb agreement.

In this work, we perform a series of experiments to probe the nature of the representations learned by different layers of BERT.<sup>1</sup> We first show that the lower layers capture phrase-level information, which gets diluted in the upper layers. Second, we propose to use the probing tasks defined in Conneau et al. (2018) to show that BERT captures a rich hierarchy of linguistic information, with surface features in lower layers, syntactic features in middle layers and semantic features in higher layers. Third, we test the ability of BERT representations to track subject-verb agreement and find that BERT requires deeper layers for handling harder cases involving long-distance dependencies. Finally, we propose to use the recently introduced Tensor Product Decomposition Network (TPDN) (McCoy et al., 2019) to explore different hypotheses about the compositional nature of BERT's representation and find that BERT implicitly captures classical, tree-like structures.

## 2 BERT

BERT (Devlin et al., 2018) builds on Transformer networks (Vaswani et al., 2017) to pre-train bidirectional representations by conditioning on both left and right contexts jointly in all layers. The representations are jointly optimized by predicting randomly masked words in the input and classify-

<sup>1</sup>The code to reproduce our experiments is publicly accessible at [https://github.com/ganeshjawahar/interpret\\_bert](https://github.com/ganeshjawahar/interpret_bert)

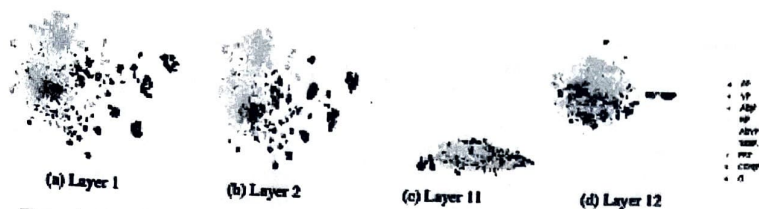


Figure 1: 2D t-SNE plot of span embeddings computed from the first and last two layers of BERT.

| layer | 1    | 2    | 3    | 4   | 5    | 6   | 7    | 8    | 9    | 10   | 11   | 12   |
|-------|------|------|------|-----|------|-----|------|------|------|------|------|------|
| NMI   | 0.38 | 0.37 | 0.35 | 0.3 | 0.24 | 0.2 | 0.19 | 0.16 | 0.17 | 0.18 | 0.16 | 0.19 |

Table 1: Clustering performance of span representations obtained from different layers of BERT.

ing whether the sentence follows a given sentence in the corpus or not. The authors of BERT claim that bidirectionality allows the model to swiftly adapt for a downstream task with little modification to the architecture. Indeed, BERT improved the state-of-the-art for a range of NLP benchmarks (Wang et al., 2018) by a significant margin.

In this work, we investigate the linguistic structure implicitly learned by BERT’s representations. We use the PyTorch implementation of BERT, which hosts the models trained by (Devlin et al., 2018). All our experiments are based on the bert-base-uncased variant,<sup>2</sup> which consists of 12 layers, each having a hidden size of 768 and 12 attention heads (110M parameters). In all our experiments, we seek the activation of the first input token (‘[CLS]’) (which summarizes the information from the actual tokens using a self-attention mechanism) at every layer to compute BERT representation, unless otherwise stated.

### 3 Phrasal Syntax

Peters et al. (2018) have shown that the representations underlying LSTM-based language models (Hochreiter and Schmidhuber, 1997) can capture phrase-level (or span-level) information.<sup>3</sup> It remains unclear if this holds true for models not trained with a traditional language modeling objective, such as BERT. Even if it does, would the information be present in multiple layers of the model? To investigate this question we extract span representations from each layer of BERT.

<sup>2</sup>We obtained similar results in preliminary experiments with the bert-large-uncased variant.

<sup>3</sup>Peters et al. (2018) experimented with ELMo-style CNN and Transformer but did not report this finding for these models.

Following Peters et al. (2018), for a token sequence  $s_1, \dots, s_j$ , we compute the span representation  $s_{(s_1, s_j), l}$  at layer  $l$  by concatenating the first  $(h_{s_1, l})$  and last hidden vector  $(h_{s_j, l})$ , along with their element-wise product and difference. We randomly pick 3000 labeled chunks and 500 spans not labeled as chunks from the CoNLL 2000 chunking dataset (Sang and Buchholz, 2000).

As shown in Figure 1, we visualize the span representations obtained from multiple layers using t-SNE (Maaten and Hinton, 2008), a non-linear dimensionality reduction algorithm for visualizing high-dimensional data. We observe that BERT mostly captures phrase-level information in the lower layers and that this information gets gradually diluted in higher layers. The span representations from the lower layers map chunks (e.g. ‘to demonstrate’) that project their underlying category (e.g. VP) together. We further quantify this claim by performing a  $k$ -means clustering on span representations with  $k = 10$ , i.e. the number of distinct chunk types. Evaluating the resulting clusters using the Normalized Mutual Information (NMI) metric shows again that the lower layers encode phrasal information better than higher layers (cf. Table 1).

### 4 Probing Tasks

Probing (or diagnostic) tasks (Adi et al., 2017; Hupkes et al., 2018; Conneau et al., 2018) help in unearthing the linguistic features possibly encoded in neural models. This is achieved by setting up an auxiliary classification task where the final output of a model is used as features to predict a linguistic phenomenon of interest. If the auxiliary classifier can predict a linguistic prop-



| Layer | SentLen<br>(Surface) | WC<br>(Surface) | TreeDepth<br>(Syntactic) | TopConst<br>(Syntactic) | BShift<br>(Syntactic) | Tense<br>(Semantic) | SubjNum<br>(Semantic) | ObjNum<br>(Semantic) | SOMO<br>(Semantic) | CoordInv<br>(Semantic) |
|-------|----------------------|-----------------|--------------------------|-------------------------|-----------------------|---------------------|-----------------------|----------------------|--------------------|------------------------|
| 1     | 93.9 (2.8)           | 24.9 (24.8)     | 35.9 (6.1)               | 63.6 (9.0)              | 50.3 (0.3)            | 82.2 (18.4)         | 77.6 (10.2)           | 76.7 (26.3)          | 49.9 (10.1)        | 33.9 (3.9)             |
| 2     | 95.9 (5.4)           | 43.0 (64.8)     | 40.6 (11.5)              | 71.3 (16.1)             | 53.8 (5.8)            | 85.9 (23.5)         | 82.5 (15.3)           | 80.6 (17.1)          | 53.8 (4.4)         | 38.5 (8.5)             |
| 3     | 96.3 (5.9)           | 66.5 (69.0)     | 39.7 (10.4)              | 71.3 (18.3)             | 64.9 (14.9)           | 86.6 (23.6)         | 82.9 (14.4)           | 89.3 (16.6)          | 55.8 (5.9)         | 39.3 (9.3)             |
| 4     | 94.2 (2.3)           | 68.8 (69.6)     | 39.4 (10.8)              | 71.3 (18.3)             | 74.4 (24.5)           | 87.6 (25.2)         | 81.9 (13.4)           | 81.4 (10.1)          | 59.0 (8.5)         | 36.1 (8.1)             |
| 5     | 92.0 (6.5)           | 69.2 (69.0)     | 40.6 (11.8)              | 81.3 (30.8)             | 81.4 (31.4)           | 89.5 (26.7)         | 85.1 (19.4)           | 81.2 (18.6)          | 60.2 (10.3)        | 64.1 (14.1)            |
| 6     | 88.4 (3.0)           | 63.5 (63.4)     | 41.3 (13.8)              | 83.3 (36.6)             | 82.9 (32.9)           | 89.8 (27.6)         | 86.1 (21.9)           | 82.0 (20.1)          | 60.7 (10.2)        | 71.1 (21.2)            |
| 7     | 83.7 (7.7)           | 56.9 (56.7)     | 40.1 (12.0)              | 84.8 (39.5)             | 83.0 (32.9)           | 89.9 (27.5)         | 87.4 (22.2)           | 82.3 (21.4)          | 61.6 (11.7)        | 74.3 (24.9)            |
| 8     | 82.9 (8.1)           | 51.1 (51.0)     | 39.2 (10.3)              | 84.0 (39.5)             | 83.9 (33.9)           | 89.9 (27.6)         | 87.5 (22.2)           | 81.2 (19.7)          | 62.1 (12.2)        | 76.4 (26.4)            |
| 9     | 80.1 (11.1)          | 47.9 (47.8)     | 38.5 (10.8)              | 83.1 (39.1)             | 87.0 (37.1)           | 90.0 (28.0)         | 87.6 (22.4)           | 81.8 (20.5)          | 63.4 (13.4)        | 76.7 (26.9)            |
| 10    | 77.0 (14.8)          | 43.4 (43.2)     | 38.1 (9.9)               | 81.7 (39.1)             | 86.7 (36.7)           | 89.7 (27.8)         | 87.1 (22.4)           | 80.5 (19.9)          | 63.3 (12.7)        | 78.4 (28.1)            |
| 11    | 73.9 (17.8)          | 42.8 (42.7)     | 36.3 (7.9)               | 80.3 (39.1)             | 86.8 (36.8)           | 89.9 (27.8)         | 85.7 (21.9)           | 78.9 (18.6)          | 64.4 (14.5)        | 73.6 (27.9)            |
| 12    | 68.5 (21.4)          | 49.1 (49.2)     | 34.7 (6.9)               | 76.5 (31.2)             | 86.4 (36.4)           | 89.5 (27.7)         | 84.0 (26.2)           | 78.7 (18.4)          | 69.3 (15.5)        | 74.9 (25.4)            |

Table 2: Probing task performance for each BERT layer. The value within the parentheses corresponds to the difference in performance of trained vs. untrained BERT.

| Layer | 0 (1.5) | 1 (5.2) | 2 (7.7) | 3 (10.5) | 4 (13.3) |
|-------|---------|---------|---------|----------|----------|
| 1     | 99.89   | 40.43   | 23.22   | 21.46    | 20       |
| 2     | 92.01   | 42.6    | 23.84   | 24.78    | 26.02    |
| 3     | 92.77   | 47.05   | 28.77   | 27.22    | 29.56    |
| 4     | 94.39   | 52.97   | 31.02   | 29.13    | 30.99    |
| 5     | 94.98   | 63.12   | 43.64   | 36.61    | 36.11    |
| 6     | 95.45   | 67.28   | 46.93   | 38.22    | 36.46    |
| 7     | 95.52   | 72.44   | 51.03   | 43.5     | 41.86    |
| 8     | 95.68   | 78.66   | 58.74   | 48.88    | 45.49    |
| 9     | 95.34   | 73.84   | 57.96   | 50.34    | 48.85    |
| 10    | 94.08   | 69.21   | 51.5    | 45.26    | 48.59    |
| 11    | 94.33   | 66.62   | 51.89   | 46.09    | 42.45    |
| 12    | 94.06   | 62.78   | 51.07   | 46.04    | 46.17    |

Table 3: Subject-verb agreement scores for each BERT layer. The last five columns correspond to the number of nouns intervening between the subject and the verb (attractors) in test instances. The average distance between the subject and the verb is enclosed in parentheses next to each attractor category.

erty well, then the original model likely encodes that property. In this work, we use probing tasks to assess individual model layers in their ability to encode different types of linguistic features. We evaluate each layer of BERT using ten probing sentence-level datasets/tasks created by Conneau et al. (2018), which are grouped into three categories. Surface tasks probe for sentence length (SentLen) and for the presence of words in the sentence (WC). Syntactic tasks test for sensitivity to word order (BShift), the depth of the syntactic tree (TreeDepth) and the sequence of top-level constituents in the syntax tree (TopConst). Semantic tasks check for the tense (Tense), the subject (resp. direct object) number in the main clause (SubjNum, resp. ObjNum), the sensitivity to random replacement of a noun/verb (SOMO) and the random swapping of coordinated clausal conjuncts (CoordInv). We use the SentEval toolkit (Conneau and Kiela, 2018) along with the recommended hyperparameter space to search for the best probing classifier. As random encoders can

surprisingly encode a lot of lexical and structural information (Zhang and Bowman, 2018), we also evaluate the untrained version of BERT, obtained by setting all model weights to a random number.

Table 2 shows that BERT embeds a rich hierarchy of linguistic signals: surface information at the bottom, syntactic information in the middle, semantic information at the top. BERT has also surpassed the previously published results for two tasks: BShift and CoordInv. We find that the untrained version of BERT corresponding to the higher layers outperforms the trained version in the task of predicting sentence length (SentLen). This could indicate that untrained models contain sufficient information to predict a basic surface feature such as sentence length, whereas training the model results in the model storing more complex information, at the expense of its ability to predict such basic surface features.

## 5 Subject-Verb Agreement

Subject-verb agreement is a proxy task to probe whether a neural model encodes syntactic structure (Linzen et al., 2016). The task of predicting the verb number becomes harder when there are more nouns with opposite number (attractors) intervening between the subject and the verb. Goldberg (2019) has shown that BERT learns syntactic phenomenon surprisingly well using various stimuli for subject-verb agreement. We extend his work by performing the test on each layer of BERT and controlling for the number of attractors. In our study, we use the stimuli created by Linzen et al. (2016) and the SentEval toolkit (Conneau and Kiela, 2018) to build the binary classifier with the recommended hyperparameter space, using as features the activations from the (masked) verb at hand.

| Role scheme \ Layer | 1      | 2      | 3      | 4      | 5      | 6      | 7      | 8      | 9      | 10     | 11     | 12     |
|---------------------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|
| Left-to-right       | 0.0003 | 0.0007 | 0.0008 | 0.0034 | 0.0058 | 0.0087 | 0.0201 | 0.0179 | 0.0284 | 0.0428 | 0.0562 | 0.0305 |
| Right-to-left       | 0.0004 | 0.0007 | 0.0007 | 0.0032 | 0.0060 | 0.0099 | 0.0233 | 0.0303 | 0.0337 | 0.0486 | 0.0481 | 0.0339 |
| Bag-of-words        | 0.0006 | 0.0009 | 0.0012 | 0.0039 | 0.0066 | 0.0108 | 0.0251 | 0.0221 | 0.0355 | 0.0307 | 0.0422 | 0.0348 |
| Bidirectional       | 0.0025 | 0.0030 | 0.0034 | 0.0053 | 0.0079 | 0.0106 | 0.0226 | 0.0201 | 0.0311 | 0.0453 | 0.0391 | 0.0334 |
| Tree                | 0.0005 | 0.0009 | 0.0011 | 0.0037 | 0.0055 | 0.0081 | 0.0179 | 0.0155 | 0.0249 | 0.0263 | 0.0389 | 0.0278 |
| Tree (random)       | 0.0005 | 0.0009 | 0.0011 | 0.0038 | 0.0063 | 0.0099 | 0.0237 | 0.0214 | 0.0338 | 0.0406 | 0.0415 | 0.0340 |

Table 4: Mean squared error between TPDN and BERT representation for a given layer and role scheme on SNLI test instances. Each number corresponds to the average across five random initializations.

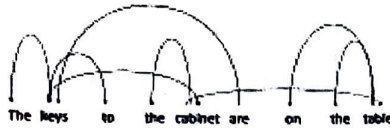


Figure 2: Dependency parse tree induced from attention head #11 in layer #2 using gold root ‘are’ as starting node for maximum spanning tree algorithm.

Results in Table 3 show that the middle layers perform well in most cases, which supports the result in Section 4 where the syntactic features were shown to be captured well in the middle layers. Interestingly, as the number of attractors increases, one of the higher BERT layers (#8) is able to handle the long-distance dependency problems caused by the longer sequence of words intervening between the subject and the verb, better than the lower layer (#7). This highlights the need for BERT to have deeper layers to perform competitively on NLP tasks.

## 6 Compositional Structure

Can we understand the compositional nature of representation learned by BERT, if any? To investigate this question, we use Tensor Product Decomposition Networks (TPDN) (McCoy et al., 2019), which explicitly compose the input token (“filler”) representations based on the role scheme selected beforehand using tensor product sum. For instance, a role scheme for a word can be based on the path from the root node to itself in the syntax tree (e.g. ‘LR’ denotes the right child of left child of root). The authors assume that, for a given role scheme, if a TPDN can be trained well to approximate the representation learned by a neural model, then that role scheme likely specifies the compositionality implicitly learned by the model. For each BERT layer, we work with five different role schemes. Each word’s role is computed based on its left-to-right index, its right-to-left index, an ordered pair containing its left-to-right and

right-to-left indices, its position in a syntactic tree (formatted version of the Stanford PCFG Parser (Klein and Manning, 2003) with no unary nodes and no labels) and an index common to all the words in the sentence (bag-of-words), which ignores its position. Additionally, we also define a role scheme based on random binary trees.

Following McCoy et al. (2019), we train our TPDN model on the premise sentences in the SNLI corpus (Bowman et al., 2015). We initialize the filler embeddings of the TPDN with the pre-trained word embeddings from BERT’s input layer, freeze it, learn a linear projection on top of it and use a Mean Squared Error (MSE) loss function. Other trainable parameters include the role embeddings and a linear projection on top of tensor product sum to match the embedding size of BERT. Table 4 displays the MSE between representation from pretrained BERT and representation from TPDN trained to approximate BERT. We discover that BERT implicitly implements a tree-based scheme, as a TPDN model following that scheme best approximates BERT’s representation at most layers. This result is remarkable, as BERT encodes classical, tree-like structures despite relying purely on attention mechanisms.

Motivated by this study, we perform a case study on dependency trees induced from self attention weight following the work done by Raganato and Tiedemann (2018). Figure 2 displays the dependencies inferred from an example sentence by obtaining self attention weights for every word pairs from attention head #11 in layer #2, fixing the gold root as the starting node and invoking the Chu-Liu-Edmonds algorithm (Chu and Liu, 1967). We observe that determiner-noun dependencies (“the keys”, “the cabinet” and “the table”) and subject-verb dependency (“keys” and “are”) are captured accurately. Surprisingly, the predicate-argument structure seems to be partly modeled as shown by the chain of dependencies between “key”, “cabinet” and “table”.



## 7 Related Work

Peters et al. (2018) studies how the choice of neural architecture such as CNNs, Transformers and RNNs used for language model pretraining affects the downstream task accuracy and the qualitative properties of the contextualized word representations that are learned. They conclude that all architectures learn high quality representations that outperform standard word embeddings such as GloVe (Pennington et al., 2014) for challenging NLP tasks. They also show that these architectures hierarchically structure linguistic information, such that morphological, (local) syntactic and (longer range) semantic information tend to be represented in, respectively, the word embedding layer, lower contextual layers and upper layers. In our work, we observe that such hierarchy exists as well for BERT models that are not trained using the standard language modelling objective. Goldberg (2019) shows that the BERT model captures syntactic information well for subject-verb agreement. We build on this work by performing the test on each layer of BERT controlling for the number of attractors and then show that BERT requires deeper layers for handling harder cases involving long-distance dependency information.

Tenney et al. (2019) is a contemporaneous work that introduces a novel edge probing task to investigate how contextual word representations encode sentence structure across a range of syntactic, semantic, local and long-range phenomena. They conclude that contextual word representations trained on language modeling and machine translation encode syntactic phenomena strongly, but offer comparably small improvements on semantic tasks over a non-contextual baseline. Their result using BERT model on capturing linguistic hierarchy confirms our probing task results although using a set of relatively simple probing tasks. Liu et al. (2019) is another contemporaneous work that studies the features of language captured/missed by contextualized vectors, transferability across different layers of the model and the impact of pretraining on the linguistic knowledge and transferability. They find that (i) contextualized word embeddings do not capture fine-grained linguistic knowledge, (ii) higher layers of RNN to be task-specific (with no such pattern for a transformer) and (iii) pretraining on a closely related task yields better performance than language model pretraining. Hewitt and Manning (2019) is

a very recent work which showed that we can recover parse trees from the linear transformation of contextual word representation consistently, better than with non-contextual baselines. They focused mainly on syntactic structure while our work additionally experimented with linear structures (left-to-right, right-to-left) to show that the compositionality modelling underlying BERT mimics traditional syntactic analysis.

The recent burst of papers around these questions illustrates the importance of interpreting contextualized word embedding models and our work complements the growing literature with additional evidences about the ability of BERT in learning syntactic structures.

## 8 Conclusion

With our experiments, which contribute to a currently bubbling line of work on neural network interpretability, we have shown that BERT does capture structural properties of the English language. Our results therefore confirm those of Goldberg (2019); Hewitt and Manning (2019); Liu et al. (2019); Tenney et al. (2019) on BERT who demonstrated that span representations constructed from those models can encode rich syntactic phenomena. We have shown that phrasal representations learned by BERT reflect phrase-level information and that BERT composes a hierarchy of linguistic signals ranging from surface to semantic features. We have also shown that BERT requires deeper layers to model long-range dependency information. Finally, we have shown that BERT's internal representations reflect a compositional modelling that shares parallels with traditional syntactic analysis. It would be interesting to see if our results transfer to other domains with higher variability in syntactic structures (such as noisy user generated content) and with higher word order flexibility as experienced in some morphologically-rich languages.

## Acknowledgments

We thank Grzegorz Chrupala and our anonymous reviewers for providing insightful comments and suggestions. This work was funded by the ANR projects ParSITi (ANR-16-CE33-0021), SoSweet (ANR15-CE38-0011-01) and the French-Israeli PHC Maimonide cooperation program.

**Q1. Based on the above paper, specifically section 4 and Table 2, comment of the separation of linguistic information present in different layers. Why do you think the model collects SOMO feature information in the last layer?**

**[10 Marks]**

**Q2. Provide a critical negative analysis of this paper. (Which means tell me the gaps and mistakes in the paper that make the conclusion untrustworthy).**

**[10 Marks]**

**(Ganesh Jawahar is an Alumnus of IIITH. He did his Master's with Prof. Vasudeva Varma)**