# What Makes My Model Perplexed?
# A Linguistic Investigation on Neural Language Models Perplexity

**Alessio Miaschi**[⋆◇]**, Dominique Brunato**[◇]**, Felice Dell'Orletta**[◇]**, Giulia Venturi**[◇]
[⋆]Department of Computer Science, University of Pisa
[◇]Istituto di Linguistica Computazionale "Antonio Zampolli", Pisa
ItaliaNLP Lab – *www.italianlp.it*
`alessio.miaschi@phd.unipi.it, {name.surname}@ilc.cnr.it`

## Abstract

This paper presents an investigation aimed at studying how the linguistic structure of a sentence affects the perplexity of two of the most popular Neural Language Models (NLMs), BERT and GPT-2. We first compare the sentence–level likelihood computed with BERT and the GPT-2's perplexity showing that the two metrics are correlated. In addition, we exploit linguistic features capturing a wide set of morpho-syntactic and syntactic phenomena showing how they contribute to predict the perplexity of the two NLMs.

## 1 Introduction and Motivation

Perplexity is one of the most standard metrics to assess the quality of a language model. It is also used in different scenarios, such as to classify formal and colloquial tweets (González, 2015), to detect the boundaries between varieties belonging to the same language family (Gamallo et al., 2017), to identify speech samples produced by subjects with cognitive and/or language diseases e.g. dementia, (Cohen and Pakhomov, 2020) or to assess whether it matches various human behavioural measures, such as gaze duration during reading (Demberg and Keller, 2008; Goodkind and Bicknell, 2018). With the recent success gained by Neural Language Models (NLMs) across a variety of NLP tasks, the notion of perplexity has started being investigated also to dig into issues related to the interpretability of contextual word representations, with the aim of understanding whether there is a relationship between this metric and the grammatical abilities implicitly encoded by a NLM (Gulordava et al., 2018; Marvin and Linzen, 2018; Kuncoro et al., 2019). In this context, Hu et al. (2020) and Warstadt et al. (2020) observed a dissociation between the perplexity of a NLM and its performance on targeted syntactic assessments probing the model's ability to encode a range of subtle syntactic phenomena.

These findings seem to be valid for models tested across languages (Mueller et al., 2020).

In this paper, we address this scenario but from a different perspective. Rather than studying the relation between the NLM's perplexity and its linguistic competences assessed on sentences undergoing controlled syntactic modifications, we focus on sentences representative of real usage. Our purpose indeed is to understand which linguistic phenomena of the input sentence may make perplexed a NLM and whether they can effectively predict the assigned perplexity score. To have a in-depth understanding of the relation between linguistic structure and perplexity, we rely on a wide spectrum of linguistic features modeling a variety of phenomena, specifically morpho-syntactic and syntactic ones. As we also intend to evaluate the possible influence of the NLM architecture on this relation, in all our experiments we consider two of the most popular NLMs, a traditional unidirectional one, i.e. GPT-2 (Radford et al., 2019), and a bidirectional model such as BERT (Devlin et al., 2019).

**Contributions** In this paper: (i) we showed that a sentence-level likelihood computed by masking each word sequentially for the BERT model has a robust correlation with GPT-2's perplexity scores; (ii) we verified whether it is possible to predict NLMs' perplexities using a wide set of linguistic features extracted by a sentence; (iii) we identified the linguistic properties of a sentence that mostly cause perplexity, reporting differences and similarities between the two models.

## 2 Our Approach

We defined two sets of experiments. The first consists in investigating the relationship between BERT and GPT-2 sentence-level perplexity (*PPL*) scores. To do so, we first computed BERT and GPT-2 *PPL* scores for sentences contained in the English Universal Dependencies (UD) treebank (Nivre et al., 2016) and we assessed their corre-

40

lation. In the second set of experiments, we studied whether a simple regression model that takes as input a wide range of linguistic features automatically extracted from each UD sentence is able to predict the two NLMs sentence-level perplexities.

To understand which linguistic phenomena contribute to the prediction of BERT and GPT-2 PPLs, and how these features differ between them, we performed an in-depth investigation training the regression model with one feature at a time.

## 2.1 Linguistic Features

The set of considered linguistic features is based on the ones described in Brunato et al. (2020) which are acquired from raw, morpho-syntactic and syntactic levels of annotation for a total of 78 features that can be categorised in 9 groups corresponding to different linguistic phenomena. A summary of the linguistic features is reported in Table 1, while the whole list is provided in Appendix A.

As shown in Table, these features model linguistic phenomena ranging from raw text one, to morpho–syntactic information and inflectional properties of verbs, to more complex aspects of sentence structure modeling global and local properties of the whole parsed tree and of specific subtrees, such as the order of subjects and objects with respect to the verb, the distribution of UD syntactic relations, also including features referring to the use of subordination and to the structure of verbal predicates.

All these features have been shown to play a highly predictive role when leveraged by traditional learning models on a variety of classification problems, also including the development of probes as reported by Miaschi et al. (2020), who showed that these features can be effectively used to profile the knowledge encoded in the language representations of a pretrained NLM.

## 2.2 Models and Data

For our experiments, we rely on the pre-trained version of the two NLMs previously defined. BERT (Devlin et al., 2019) is a Transformer-based masked language model, pretrained on BookCorpus (Zhu et al., 2015) and English Wikipedia. GPT-2 (Radford et al., 2018) is a large transformer-based language model trained using the language modeling task (LM) on 8 million documents for a total of 40 GB of text.

We first computed GPT-2's sentence-level perplexities by dividing the sum of all sub-word con-

| Linguistic Feature |
| --- |
| **Raw Text Properties** |
| Sentence Length |
| Word Length |
| **Vocabulary Richness** |
| Type/Token Ratio for words and lemmas |
| **Morphosyntactic information** |
| Distibution of UD and language–specific POS |
| Lexical density |
| **Inflectional morphology** |
| Inflectional morphology of auxiliary verbs |
| **Verbal Predicate Structure** |
| Distribution of verbal heads and verbal roots |
| Verb arity and distribution of verbs by arity |
| **Global and Local Syntactic Tree Structures** |
| Depth of the whole syntactic tree |
| Average length of dependency links and of the longest link |
| Average length of prepositional chains and distribution by depth |
| Clause length |
| **Relative order of elements** |
| Order of subject and object |
| **Syntactic Relations** |
| Distribution of dependency relations |
| **Use of Subordination** |
| Distribution of subordinate and principal clauses |
| Average length of subordination chains and distribution by depth |
| Relative order of subordinate clauses |

Table 1: Linguistic Features used in the experiments.

ditional log-probabilities by the total number of words for each sentence in the UD dataset. On the other hand, since BERT masked language modeling task does not allow to compute well-formed probability distributions over sentences, we measure BERT sentence-level likelihood by masking each word sequentially and computing the probability as follows:

$$p(S) \approx \prod_{i=1}^{k} p(w_i|context)$$

where *context*, given the deep bidirectionality of the model, corresponds to $w_1, ..., w_{i-1}, w_{i+1}, ..., w_k$. The perplexity is then computed as follows:

$$PPL_S = e^{\left(\frac{p(S)}{N}\right)}$$

where *N* correspond to the length of sentence *S*. In order to uniform the terminology, in what follows we will refer to the BERT sentence-level likelihood as perplexity.

In order to evaluate our approach on gold annotated sentences, we relied on the three English Universal Dependencies (UD) treebanks: the English version of ParTUT (Sanguinetti and Bosco, 2015), the UD version of the GUM corpus (Zeldes, 2017) and of the English Web Treebank (EWT) (Silveira et al., 2014). Overall, the final dataset consists of 22,505 sentences.

| Lengths | $\rho$ score | # samples |
|---------|---------|-----------|
| All | 0.63 | 22,505 |
| n=10 | 0.66 | 847 |
| n=15 | 0.60 | 793 |
| n=20 | 0.64 | 643 |
| n=25 | 0.53 | 422 |
| n=30 | 0.54 | 277 |

Table 2: Spearman correlations between BERT and GPT-2 perplexities computed for all UD sentences (*All*) and sentences with fixed-length *n*.

# 3 A Linguistic Investigation on Perplexity

As a first step, we assessed whether there is a relationship between the perplexity of a traditional NLM and of a masked NLM. We thus calculated BERT and GPT-2 perplexity scores for each UD sentence and measured the correlation between them. Since *PPL* scores are highly affected by the length of the input sequence, we computed $\rho$ correlation coefficients also considering groups of sentences with fixed length. Specifically, we relied on Spearman correlation because we were interested in measuring how the variations in perplexity scores relate each other, rather than focusing on the actual *PPL* values. Results are reported in Table 2. As we can notice, even considering samples with fixed length, the two NLMs' perplexities exhibit moderate to substantial correlation (with $p < 0.001$), thus showing that BERT an GPT-2 do not diverge excessively in their ability of predicting the likelihood of the input sentences. Moreover, this allows us to confirm that, although the deep bidirectional structure of BERT does not permit to compute a well-formed probability distribution over a sentence (see Section 2.2), this metric could be considered as a valid approximation of the perplexity computed with a unidirectional NLM.

Once established the correlation between the perplexities of the two NLMs, we performed a second experiment to investigate (i) if the considered set of linguistic features plays a role in predicting their perplexity and (ii) which are the features that contribute more to the prediction task. To do so, we trained a LinearSVR model that predicts perplexity's scores using our set of linguistic properties as input features. Since most of them refer to syntactic properties of sentence that are strongly correlated with its length, we considered as a baseline a SVR model that takes sentence length as input and outputs BERT/GPT-2 sentence's perplexity. Regression results deriving by considering both the
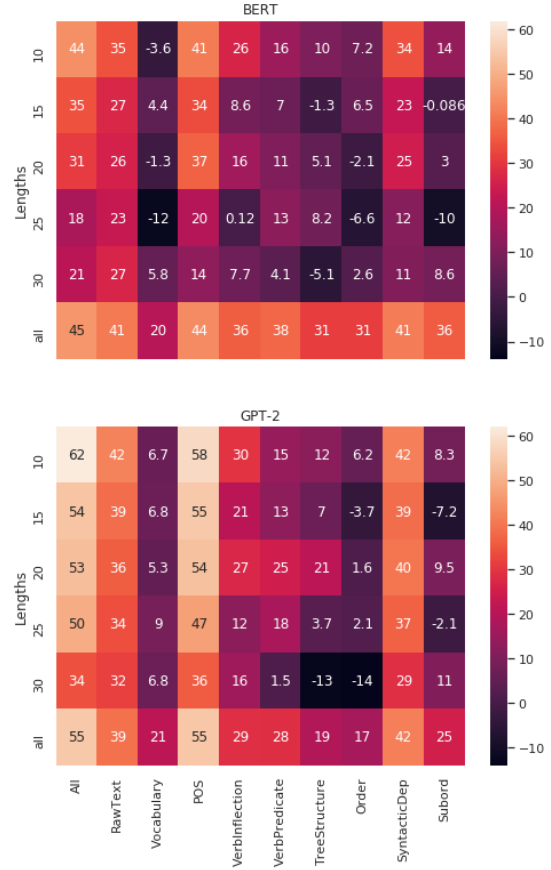


Figure 1: BERT and GPT-2 $\rho$ scores (multiplied by 100) obtained with the LinearSVR model using linguistic features, for the whole UD dataset and groups of sentences with fixed length.

whole set (*All*) and each of the 9 groups of linguistic features separately are reported in Figure 1. As a general remark, for the whole UD dataset, we can observe that the results considering both all and the 9 groups of linguistic features outperform the results obtained by the baseline, i.e. $\rho$=0.38 for BERT and 0.22 for GPT-2 respectively. This demonstrates that the considered features are able to model aspects involved in NLM's perplexity that go beyond the simple length of sentence. This is particularly the case of GPT-2, suggesting that the probability assigned to a sentence by a traditional NLM is more explainable in terms of linguistic phenomena mainly affecting morpho-syntactic and syntactic structure. Consequently, the baseline score is higher for BERT. If we consider the scores obtained for each group of sentences with fixed length, we can see that higher scores are obtained for groups containing shorter sentences, for both NLMs. This is quite expected since in these sentences the possible output space is smaller for almost all features,

| | Sentence length = All | | | Sentence length = 16 | |
|---|---|---|---|---|---|
| lexical_density | 0.4 | 0.38 (1) | lexical_density | 0.29 | 0.38 (1) |
| %_upos_PRON | 0.38 | 0.35 (2) | %_upos_PROPN | 0.25 | 0.29 (4) |
| verbal_heads | 0.37 | 0.26 (5) | %_xpos_NNP | 0.25 | 0.25 (6) |
| %_dep_root | 0.37 | 0.22 (10) | %_dep_compound | 0.2 | 0.21 (7) |
| sent_length | 0.37 | 0.22 (8) | char_per_tok | 0.19 | 0.33 (2) |
| avg_verb_edges | 0.35 | 0.22 (9) | %_upos_PRON | 0.19 | 0.31 (3) |
| parse_depth | 0.34 | 0.17 (24) | %_xpos_PRP | 0.16 | 0.26 (5) |
| max_links_len | 0.32 | 0.18 (19) | %_upos_AUX | 0.15 | 0.12 (13) |
| %_dep_nsubj | 0.31 | 0.22 (11) | %_dep_mark | 0.13 | 0.16 (8) |
| char_per_tok | 0.31 | 0.34 (3) | verbal_heads | 0.13 | 0.14 (10) |
| %_subj_pre | 0.3 | 0.17 (25) | %_xpos_VB | 0.11 | 0.14 (12) |
| clause_length | 0.3 | 0.13 (37) | %_aux_mood_Ind | 0.09 | 0.078 (25) |
| %_upos_AUX | 0.3 | 0.21 (12) | %_dep_punct | 0.086 | -0.078 (74) |
| %_verbal_root | 0.29 | 0.18 (18) | %_upos_PUNCT | 0.086 | -0.13 (77) |
| %_xpos_PRP | 0.29 | 0.29 (4) | %_dep_det | 0.082 | 0.057 (37) |
| avg_links_len | 0.28 | 0.13 (35) | %_dep_nsubj | 0.08 | 0.16 (9) |
| %_aux_form_Fin | 0.28 | 0.18 (21) | %_dep_advmod | 0.077 | 0.072 (30) |
| avg_subord_chain | 0.27 | 0.2 (14) | %_upos_DET | 0.077 | 0.068 (32) |
| %_subord_prop | 0.26 | 0.18 (17) | %_dep_aux | 0.074 | 0.099 (18) |
| %_upos_VERB | 0.26 | 0.18 (22) | %_upos_VERB | 0.074 | 0.087 (21) |
| | BERT | GPT-2 | | BERT | GPT-2 |

Figure 2: BERT and GPT-2 $\rho$ scores obtained with the LinearSVR model, for the whole UD dataset and 16 token-long sentences. Scores are reported for the 20 top-ranked features for BERT. Numbers in brackets correspond to the relative in the GPT-2 ranking.

thus making them more predictive. Also in this case, the impact of the linguistic features is always higher for the prediction of GPT-2's perplexity.

A more in-depth analysis of these results shows that the distribution of the morpho-syntactic characteristics of a sentence (*POS*) and of the syntactic dependency relations (*SyntacticDep*) are the two most predictive sources of linguistic information. As Figure 1 reports, this holds for the two NLM models and it remains constant throughout all the groups of sentences with fixed lengths. Interestingly, if we consider the whole set of sentences, the effect of the morpho-syntactic information on the prediction of GPT-2's perplexity is exactly the same of that of the whole set of linguistic features. For some sentence lengths (15, 20, 30) the scores obtained using only this type of information outperform even those obtained considering the whole set of features. Note that this last remark is true also in the prediction of BERT's perplexity. As expected the other most predictive group is the one (*RawText*) that includes the length of sentence.

### 3.1 Focus on the contribution of individual features

To investigate more in depth which linguistic phenomena are more involved in the perplexity of the two models, we trained the LinearSVR model using each individual feature at a time. This was done for both the whole dataset and the subset of sentences (i.e. 758 sentences) having a length of 16 tokens,

which corresponds to the mean sentence length of the UD dataset. A subset of results is reported in Figure 2, while the whole results are provided in Appendix B. As we can see in the left-side of the heatmap, the two models share many features in the first ten positions, thus showing that the two NLM architectures are made perplexed by similar linguistic characteristics of a sentence. In particular, for both of them, the two most predictive features correspond to the lexical density and the presence of pronouns confirming the highly predictive power of morpho-syntactic information. They are followed by features related to the presence of verbs and to their internal structure (i.e. *verbal_heads* and *avg_verb_edges*), and, as it was expected, by the length of the sentence. Despite these similarities, we can see that the scores obtained by the regression model to predict BERT's perplexity are on average higher than GPT-2's scores. Considering that we obtained higher scores using all (or groups of) features in the prediction of GPT-2' perplexity (see Figure 1), this latter result may suggest that the interaction among features is less relevant in the prediction of BERT's perplexity. Differences among the two models concern features that are highly sensitive to sentence length, which result to be more predictive of BERT's perplexity. This is the case of syntactic features capturing global and local aspects of sentence structure, i.e. the depth of the whole syntactic tree (*parse_depth*), the maximum length of dependency links (*max_links_len*) and the length of verbal clauses (*clause_length*). Also, the canonical order of nuclear sentence elements such as pre-verbal subjects contribute more to predict BERT's than GPT-2's perplexity. Instead, the distribution of proper nouns (*%_upos_PROPN*), in particular in their singular form (*%_xpos_NNP*), the length of token (*char_per_tok*) and vocabulary richness are more predictive of GPT-2's perplexity. Although we cannot say from ranking results whether features highly ranked are positively or negatively correlated with perplexity, we can hypothesize that knowing the distribution of tokens belonging to open lexical categories (e.g. proper nouns vs determiners) make the perplexity easier to identify.

The right-side heatmap shows the top-ranked features used to predict the two models perplexity for sentences 16-token long. As expected, when sentence length is controlled, the role of other features less related to length becomes predominant.

In particular, morpho-syntactic information is still highly predictive for the two models, with lexical parts-of-speech showing to be relevant not only for GPT-2's but also of BERT's perplexity.

## 4 Conclusion

In this paper we proposed an investigation of the linguistic phenomena characterizing the perplexity of a undirectional and a bidirectional Neural Language Model, GPT-2 and BERT. We first reported robust correlations between GPT-2's perplexity and the sentence-level likelihood computed with BERT. This is a quite prominent result, especially considering that these two metrics are differently computed as a consequence of the two NLMs architectures.

Interestingly, we found the effectiveness of linguistic features modelling a wide set of morpho-syntactic and syntactic phenomena in predicting the perplexity of the two NLMs, especially for shorter sentences. Despite similar trends, we observed some differences between the two NLMs both at the level of regression accuracy and in the rankings of the features exploited in the prediction of perplexity. GPT-2's perplexity is better captured by the considered features and it resulted to be more affected by lexical parts-of-speech and features capturing the vocabulary richness of a sentence. On the contrary, BERT's perplexity seems to be best predicted by syntactic features highly sensitive to sentence length.

## References

Dominique Brunato, Andrea Cimino, Felice Dell'Orletta, Giulia Venturi, and Simonetta Montemagni. 2020. Profiling-ud: a tool for linguistic profiling of texts. In *Proceedings of The 12th Language Resources and Evaluation Conference*, pages 7147–7153, Marseille, France. European Language Resources Association.

Trevor Cohen and Serguei Pakhomov. 2020. A tale of two perplexities: Sensitivity of neural language models to lexical retrieval deficits in dementia of the Alzheimer's type. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1946–1957, Online. Association for Computational Linguistics.

V. Demberg and Frank Keller. 2008. Data from eye-tracking corpora as evidence for theories of syntactic processing complexity. *Cognition*, 109:193–210.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186.

Pablo Gamallo, Jose Ramom Pichel, and Iñaki Alegria. 2017. A perplexity-based method for similar languages discrimination. In *VarDial2017 workshop at EACL 2017. Proceedings of the Fourth Workshop on NLP for Similar Languages, Varieties and Dialects, pages 109–114,Valencia, Spain, April 3, 2017. c©2017 Association for Computational Linguistics (http://web.science.mq.edu.au/ smalmasi/vardial4/index.html)*.

M. González. 2015. An analysis of twitter corpora and the differences between formal and colloquial tweets. In *TweetMT@SEPLN*.

Adam Goodkind and Klinton Bicknell. 2018. Predictive power of word surprisal for reading times is a linear function of language model quality. In *Proceedings of the 8th Workshop on Cognitive Modeling and Computational Linguistics, CMCL 2018, Salt Lake City, Utah, USA, January 7, 2018*, pages 10–18. Association for Computational Linguistics.

Kristina Gulordava, Piotr Bojanowski, Edouard Grave, Tal Linzen, and Marco Baroni. 2018. Colorless green recurrent networks dream hierarchically. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1195–1205, New Orleans, Louisiana. Association for Computational Linguistics.

Jennifer Hu, Jon Gauthier, Peng Qian, Ethan Wilcox, and Roger Levy. 2020. A systematic assessment of syntactic generalization in neural language models. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1725–1744, Online. Association for Computational Linguistics.

Adhiguna Kuncoro, Chris Dyer, Laura Rimell, Stephen Clark, and Phil Blunsom. 2019. Scalable syntax-aware language models using knowledge distillation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3472–3484, Florence, Italy. Association for Computational Linguistics.

Rebecca Marvin and Tal Linzen. 2018. Targeted syntactic evaluation of language models. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1192–1202.

Alessio Miaschi, Dominique Brunato, Felice Dell'Orletta, and Giulia Venturi. 2020. Linguistic profiling of a neural language model. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 745–756,

Barcelona, Spain (Online). International Committee on Computational Linguistics.

Aaron Mueller, Garrett Nicolai, Panayiota Petrou-Zeniou, Natalia Talmina, and Tal Linzen. 2020. Cross-linguistic syntactic evaluation of word prediction models. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5523–5539, Online. Association for Computational Linguistics.

Joakim Nivre, Marie-Catherine De Marneffe, Filip Ginter, Yoav Goldberg, Jan Hajic, Christopher D Manning, Ryan McDonald, Slav Petrov, Sampo Pyysalo, Natalia Silveira, et al. 2016. Universal dependencies v1: A multilingual treebank collection. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 1659–1666.

Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *Technical report*.

Manuela Sanguinetti and Cristina Bosco. 2015. Parttut: The turin university parallel treebank. In *Harmonization and Development of Resources and Tools for Italian Natural Language Processing within the PARLI Project*, pages 51–69. Springer.

Natalia Silveira, Timothy Dozat, Marie-Catherine De Marneffe, Samuel R Bowman, Miriam Connor, John Bauer, and Christopher D Manning. 2014. A gold standard dependency corpus for english. In *LREC*, pages 2897–2904.

Alex Warstadt, Alicia Parrish, Haokun Liu, Anhad Mohananey, Wei Peng, Sheng-Fu Wang, and Samuel R. Bowman. 2020. BLiMP: The benchmark of linguistic minimal pairs for english. *Transactions of the Association for Computational Linguistics*, 8:377–392.

Amir Zeldes. 2017. The GUM corpus: Creating multilayer resources in the classroom. *Language Resources and Evaluation*, 51(3):581–612.

Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *The IEEE International Conference on Computer Vision (ICCV)*.

# A  Appendix A

| **fRaw Text Properties** |
| :--- |
| [sent_length]: average length of sentences in a document, calculated in terms of the number of words per sentence |
| [char_per_tok]: average number of characters per word (excluded punctuation) |

| **Vocabulary Richness** |
| :--- |
| [ttr_lemma]: Type/Token Ratio (TTR) calculated with respect to the lemmata in a sentence. It ranges between 1 (high lexical variety) and 0 (low vocabulary richness) |
| [ttr_form]: Type/Token Ratio (TTR) calculated with respect to the word forms in a sentence. It ranges between 1 (high lexical variety) and 0 (low vocabulary richness) |

| **Morphosyntactic information** |
| :--- |
| [%_upos_*]: distribution of the part-of-speech categories defined in the Universal POS tags, as detailed at the following link: `https://universaldependencies.org/u/pos/index.html` |
| [%_xpos_*]: distribution of the part-of-speech categories defined in the Penn Treebank POS tags, as detailed at the following link: `https://www.ling.upenn.edu/courses/Fall_2003/ling001/penn_treebank_pos.html` |
| [lexical_density]: the value corresponds to the ratio between content words (nouns, proper nouns, verbs, adjectives, adverbs) over the total number of words in a sentence |

| **Inflectional morphology** |
| :--- |
| [%_aux_tense_*]: distribution of auxiliary verbs according to their tense: `https://universaldependencies.org/u/feat/Tense.html` |
| [%_aux_mood_dist_*]: distribution of auxiliary verbs according to their moods: `https://universaldependencies.org/u/feat/Mood.html` |
| [%_aux_form_*]: distribution of auxiliary verbs according to their forms: `https://universaldependencies.org/u/feat/VerbForm.html` |
| [verbs_gender_dist_*]: distribution of verbs according to the gender of participle forms, for the languages that have this features: `https://universaldependencies.org/u/feat/Gender.html` |
| [%_aux_num_pers_*]: distribution of auxiliary verbs according to their number and person: `https://universaldependencies.org/u/feat/Person.html` |

| **Verbal Predicate Structure** |
| :--- |
| [verbal_head]: average distribution of verbal heads in the document, out ot the total of heads. |
| [%_verbal_roo]: average distribution of roots headed by a lemma tagged as verb, out of the total of sentence roots; |
| [avg_verb_edges]: verbal arity, calculated as the average number of instantiated dependency links (covering both arguments and modifiers) sharing the same verbal head, excluding punctuation and auxiliaries bearing the syntactic role of copula according to the UD scheme |
| [verbal_arity*]: distribution of verbs for arity class (e.g. verbs with arity 1, 2, ...) |

| **Global and Local Syntactic Tree Structures** |
| :--- |
| [parse_depth]: mean of the maximum tree depths of the sentence. The maximum depth is calculated as the longest path (in terms of occurring dependency links) from the root of the dependency tree to some leaf |
| [clause_length]: average clause length, calculated in terms of the average number of tokens per clause, where a clause is defined as the ratio between the number of tokens in a sentence and the number of either verbal or copular head |
| [avg_links_len]: average number of words occurring linearly between each syntactic head and its dependent (excluding punctuation dependencies) |
| [max_links_len]: the value of the longest dependency link in the document, calculated in number of tokens |
| [prep_1]: distribution of prepositional chains 1-complement long. A prepositional chain is calculated as the number of embedded prepositional complements dependent on a noun |

| **Relative order of elements** |
| :--- |
| [%_obj_post]: distribution of objects following the verb |
| [%_subj_pre]: distribution of subjects preceding the verb |

| **Syntactic Relations** |
| :--- |
| [%_dep_*]: average distribution of the 37 universal syntactic relations used in UD (`https://universaldependencies.org/u/dep/index.html`) |

| **Use of Subordination** |
| :--- |
| [principal_prop_dist]: distribution of principal clauses |
| [%_subord_prop]: distribution of subordinate clauses, as defined in the UD scheme: `https://universaldependencies.org/u/overview/complex-syntax.html#subordination` |
| [subord_post]: distribution of subordinate clauses following the main clause |
| [avg_subord_chain]: average length of subordinate chains, where a subordinate 'chain' is calculated as the number of subordinate clauses embedded on a first subordinate clause |
| [subord_1]: distribution of subordinate chains 1-clause long |

Table 3: Linguistic features used in the experiments.

# B    Appendix B

| Sentence length = All | | |
| --- | --- | --- |
| | BERT | GPT-2 |
| lexical_density | 0.4 | 0.38 (1) |
| %_upos_PRON | 0.38 | 0.35 (2) |
| verbal_heads | 0.37 | 0.26 (5) |
| %_dep_root | 0.37 | 0.22 (10) |
| sent_length | 0.37 | 0.22 (8) |
| avg_verb_edges | 0.35 | 0.22 (9) |
| parse_depth | 0.34 | 0.17 (24) |
| max_links_len | 0.32 | 0.18 (19) |
| %_dep_nsubj | 0.31 | 0.22 (11) |
| char_per_tok | 0.31 | 0.34 (3) |
| %_subj_pre | 0.3 | 0.17 (25) |
| clause_length | 0.3 | 0.13 (37) |
| %_upos_AUX | 0.3 | 0.21 (12) |
| %_verbal_root | 0.29 | 0.18 (18) |
| %_xpos_PRP | 0.29 | 0.29 (4) |
| avg_links_len | 0.28 | 0.13 (35) |
| %_aux_form_Fin | 0.28 | 0.18 (21) |
| avg_subord_chain | 0.27 | 0.2 (14) |
| %_subord_prop | 0.26 | 0.18 (17) |
| %_upos_VERB | 0.26 | 0.18 (22) |
| %_upos_DET | 0.25 | 0.16 (29) |
| %_dep_det | 0.24 | 0.15 (30) |
| subord_post | 0.24 | 0.17 (23) |
| %_dep_mark | 0.23 | 0.2 (13) |
| %_aux_mood_Ind | 0.22 | 0.15 (31) |
| %_upos_PROPN | 0.22 | 0.25 (6) |
| %_xpos_NNP | 0.22 | 0.22 (7) |
| %_dep_advmod | 0.22 | 0.14 (34) |
| %_xpos_DT | 0.21 | 0.16 (28) |
| %_dep_aux | 0.21 | 0.16 (27) |
| %_upos_ADP | 0.21 | 0.096 (56) |
| %_dep_obl | 0.2 | 0.12 (43) |
| %_xpos_IN | 0.2 | 0.11 (48) |
| %_upos_ADV | 0.2 | 0.13 (36) |
| subord_1 | 0.19 | 0.11 (47) |
| principal_prop_dist | 0.19 | 0.075 (62) |
| %_dep_case | 0.19 | 0.053 (65) |
| %_xpos_VB | 0.19 | 0.18 (20) |
| %_upos_PART | 0.19 | 0.12 (39) |
| %_aux_num_pers_+ | 0.18 | 0.14 (33) |
| verbal_arity_3 | 0.18 | 0.12 (40) |
| %_obj_post | 0.18 | 0.1 (53) |
| %_xpos_RB | 0.18 | 0.12 (42) |
| %_dep_obj | 0.18 | 0.1 (52) |
| %_aux_tense_Pres | 0.17 | 0.095 (57) |
| %_upos_SCONJ | 0.15 | 0.12 (44) |
| %_dep_cc | 0.14 | 0.093 (58) |
| %_dep_cop | 0.14 | 0.12 (46) |
| %_upos_CCONJ | 0.14 | 0.1 (51) |
| verbal_arity_4 | 0.14 | 0.098 (54) |
| verbal_arity_2 | 0.14 | 0.083 (60) |
| %_xpos_CC | 0.14 | 0.096 (55) |
| %_xpos_VBP | 0.14 | 0.11 (50) |
| %_dep_advcl | 0.13 | 0.11 (49) |
| %_dep_conj | 0.13 | 0.082 (61) |
| %_xpos_TO | 0.13 | 0.12 (41) |
| ttr_lemma | 0.12 | 0.2 (15) |
| %_aux_num_pers_Sing+3 | 0.12 | 0.071 (63) |
| %_dep_nmod | 0.11 | -0.024 (77) |
| prep_1 | 0.11 | 0.03 (70) |
| %_xpos_VBD | 0.096 | 0.071 (64) |
| %_dep_nmod:poss | 0.092 | 0.047 (67) |
| ttr_form | 0.088 | 0.17 (26) |
| %_xpos_VBZ | 0.079 | 0.0069 (72) |
| %_xpos_VBN | 0.066 | -0.02 (76) |
| %_dep_compound | 0.05 | 0.19 (16) |
| %_upos_NOUN | 0.048 | 0.15 (32) |
| %_xpos_NN | 0.029 | 0.084 (59) |
| %_dep_punct | -0.0064 | -0.013 (73) |
| %_upos_PUNCT | -0.0071 | 0.034 (69) |
| %_upos_NUM | -0.016 | 0.04 (68) |
| %_xpos_NNS | -0.058 | 0.013 (71) |
| %_xpos_, | -0.067 | 0.13 (38) |
| %_dep_amod | -0.081 | 0.049 (66) |
| %_xpos_JJ | -0.085 | -0.019 (75) |
| %_upos_ADJ | -0.11 | -0.017 (74) |
| %_xpos_. | -0.13 | 0.12 (45) |

| Sentence length = 16 | | |
| --- | --- | --- |
| | BERT | GPT-2 |
| lexical_density | 0.29 | 0.38 (1) |
| %_upos_PROPN | 0.25 | 0.29 (4) |
| %_xpos_NNP | 0.25 | 0.25 (6) |
| %_dep_compound | 0.2 | 0.21 (7) |
| char_per_tok | 0.19 | 0.33 (2) |
| %_upos_PRON | 0.19 | 0.31 (3) |
| %_xpos_PRP | 0.16 | 0.26 (5) |
| %_upos_AUX | 0.15 | 0.12 (13) |
| %_dep_mark | 0.13 | 0.16 (8) |
| verbal_heads | 0.13 | 0.14 (10) |
| %_xpos_VB | 0.11 | 0.14 (12) |
| %_aux_mood_Ind | 0.09 | 0.078 (25) |
| %_dep_punct | 0.086 | -0.078 (74) |
| %_upos_PUNCT | 0.086 | -0.13 (77) |
| %_dep_det | 0.082 | 0.057 (37) |
| %_dep_nsubj | 0.08 | 0.16 (9) |
| %_dep_advmod | 0.077 | 0.072 (30) |
| %_upos_DET | 0.077 | 0.068 (32) |
| %_dep_aux | 0.074 | 0.099 (18) |
| %_upos_VERB | 0.074 | 0.087 (21) |
| subord_post | 0.073 | 0.097 (19) |
| avg_subord_chain | 0.072 | 0.12 (14) |
| %_aux_num_pers_+ | 0.07 | 0.079 (23) |
| %_upos_PART | 0.066 | -0.039 (68) |
| %_aux_tense_Pres | 0.06 | 0.025 (45) |
| avg_verb_edges | 0.058 | 0.03 (44) |
| %_subord_prop | 0.057 | 0.067 (33) |
| verbal_arity_4 | 0.05 | -0.052 (71) |
| %_xpos_DT | 0.049 | 0.11 (15) |
| %_dep_obl | 0.046 | -0.0082 (53) |
| %_upos_SCONJ | 0.043 | 0.031 (43) |
| clause_length | 0.043 | 0.079 (24) |
| %_dep_cop | 0.042 | 0.043 (41) |
| %_upos_ADV | 0.042 | 0.077 (27) |
| %_xpos_TO | 0.04 | 0.088 (20) |
| %_upos_NOUN | 0.04 | 0.14 (11) |
| avg_links_len | 0.037 | -0.033 (64) |
| %_aux_form_Fin | 0.035 | 0.059 (36) |
| %_upos_NUM | 0.034 | 0.084 (22) |
| max_links_len | 0.03 | 0.049 (38) |
| %_xpos_VBP | 0.026 | 0.068 (31) |
| %_xpos_NN | 0.018 | 0.022 (47) |
| %_xpos_RB | 0.014 | 0.077 (26) |
| %_xpos_. | 0.013 | 0.077 (28) |
| %_xpos_, | 0.012 | 0.02 (50) |
| %_dep_nmod | 0.0084 | 0.1 (17) |
| %_verbal_root | 0.0051 | 0.061 (35) |
| %_dep_obj | 0.0036 | -0.054 (72) |
| ttr_form | 0.00076 | 0.076 (29) |
| %_dep_conj | -0.00065 | -0.0059 (52) |
| %_subj_pre | -0.012 | -0.051 (70) |
| %_xpos_VBZ | -0.017 | -0.011 (56) |
| %_dep_case | -0.019 | 0.022 (46) |
| %_dep_advcl | -0.02 | -0.0088 (54) |
| verbal_arity_3 | -0.02 | -0.01 (55) |
| %_dep_root | -0.025 | -0.017 (57) |
| prep_1 | -0.026 | 0.042 (42) |
| %_aux_num_pers_Sing+3 | -0.03 | -0.032 (63) |
| %_xpos_NNS | -0.031 | 0.047 (39) |
| subord_1 | -0.035 | -0.065 (73) |
| %_upos_ADJ | -0.036 | -0.04 (69) |
| parse_depth | -0.037 | -0.036 (67) |
| sent_length | -0.039 | -0.031 (62) |
| verbal_arity_2 | -0.041 | 0.02 (49) |
| %_dep_amod | -0.047 | 0.11 (16) |
| %_obj_post | -0.051 | -0.024 (60) |
| ttr_lemma | -0.06 | 0.061 (34) |
| %_dep_nmod:poss | -0.061 | -0.12 (75) |
| %_upos_CCONJ | -0.065 | 0.021 (48) |
| %_xpos_CC | -0.076 | -0.035 (66) |
| principal_prop_dist | -0.077 | -0.021 (58) |
| %_xpos_IN | -0.083 | -0.03 (61) |
| %_xpos_VBD | -0.091 | -0.022 (59) |
| %_dep_cc | -0.091 | 0.012 (51) |
| %_upos_ADP | -0.098 | -0.12 (76) |
| %_xpos_VBN | -0.11 | 0.044 (40) |
| %_xpos_JJ | -0.11 | -0.035 (65) |

Figure 3: BERT and GPT-2 $\rho$ scores obtained with the LinearSVR model using one feature at a time, for the whole UD dataset and sentences with lengths = 16.