

A Treebank dataset tree types

Due to space constraints, we only show the top 20 tree types. This can be found in Table 7.

B Treebank dataset phrase lengths

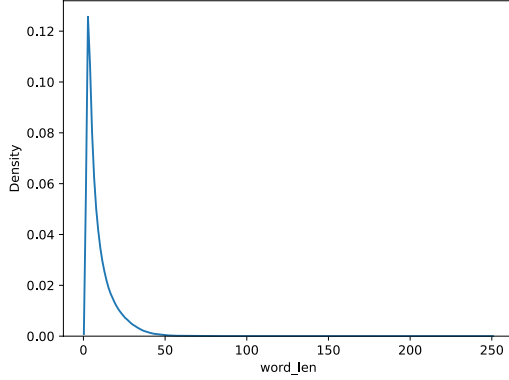


Figure 6: Length distribution of phrases mined from the treebank, in number of words. The modal length was 3 words, followed closely by 2 words. Few phrases contained more than 50 words.

C Probe learning curves

Learning curves of the approximative probes (across 10 folds) are shown in Figure 7.

D Length Correlation

The correlations of the phrase length (in words) and compositionality scores in Treebank are shown in Table 8.

E Error ratio of probes

Model/representation	Probe	Mean err. ratio (\downarrow)
BERT _{CLS}	ADD	0.4668
	W1	0.7806
	W2	0.3903
	LIN	0.3940
	AFF	0.3908
	MLP	0.3830
RoBERTa _{CLS}	ADD	0.4152
	W1	0.7946
	W2	0.2980
	LIN	0.3063
	AFF	0.3013
	MLP	0.3065
DeBERTa _{CLS}	ADD	0.7577
	W1	0.4661
	W2	0.7090
	LIN	0.6777
	AFF	0.9373
	MLP	0.5856
GPT-2 _{last}	ADD	0.4668
	W1	0.7806
	W2	0.3903
	LIN	0.3940
	AFF	0.3908
	MLP	0.3830
BERT _{AVG}	ADD	0.3873
	W1	0.8060
	W2	0.2167
	LIN	0.2327
	AFF	0.2098
	MLP	0.2283
RoBERTa _{AVG}	ADD	0.4504
	W1	0.8422
	W2	0.2431
	LIN	0.2471
	AFF	0.2095
	MLP	0.2181
DeBERTa _{AVG}	ADD	0.4472
	W1	0.8886
	W2	0.3202
	LIN	0.3143
	AFF	0.3044
	MLP	0.2952
GPT-2 _{AVG}	ADD	0.5013
	W1	0.9074
	W2	0.4226
	LIN	0.4041
	AFF	0.3475
	MLP	0.3554

Table 9: Error ratio ($\frac{\text{dist}_{\text{probe}}}{\text{dist}_{\text{control}}}$) for probes trained to predict representations from different model types. Mean across 10 folds.

Tree type	Count	Example
PP → IN NP	77716	((in) (american romance))
S → NP-SBJ VP	62948	((he) (said simultaneously, "i wish they were emeralds"))
NP → DT NN	40876	((the) (way))
NP → NP PP	35743	((the temporal organization) (of the dance))
S → NP-SBJ S <VP->	24467	((the partners) (said they already hold 15 % of all shares outstanding.))
VP → TO VP	21833	((to) (be the enemy))
PP-LOC → IN NP	18005	((in) (the marketplace))
NP → DT NP <JJ-NN>	14898	((a) (professional linguist))
VP → MD VP	13575	((could) (make up his mind))
VP → VB NP	11838	((evaluate) (the progress of therapy))
PP-TMP → IN NP	11032	((for) (almost a year))
PP-CLR → IN NP	10054	((from) (the most sympathetic angle))
NP → NNP NNP	9863	((honolulu) (harbor))
NP → JJ NNS	9477	((recent) (years))
VP → VBD VP	8356	((was) (salted))
SBAR → WHNP-1 S	8332	((what) (to look for))
SBAR → IN S	7848	((that) (it exceeds the company 's annual sales and its market capitalization))
NP-SBJ → DT NN	7600	((the) (rebound))
S → NP-SBJ-1 VP	7486	((draperies) (could be designed to serve structural purposes))
NP → NP SBAR	7317	((the " culture shock ") (they might encounter in remote overseas posts))

Table 7: Counts of the top 20 grammatical tree types found in the WSJ and Brown sections of the Penn Treebank, with some examples given.

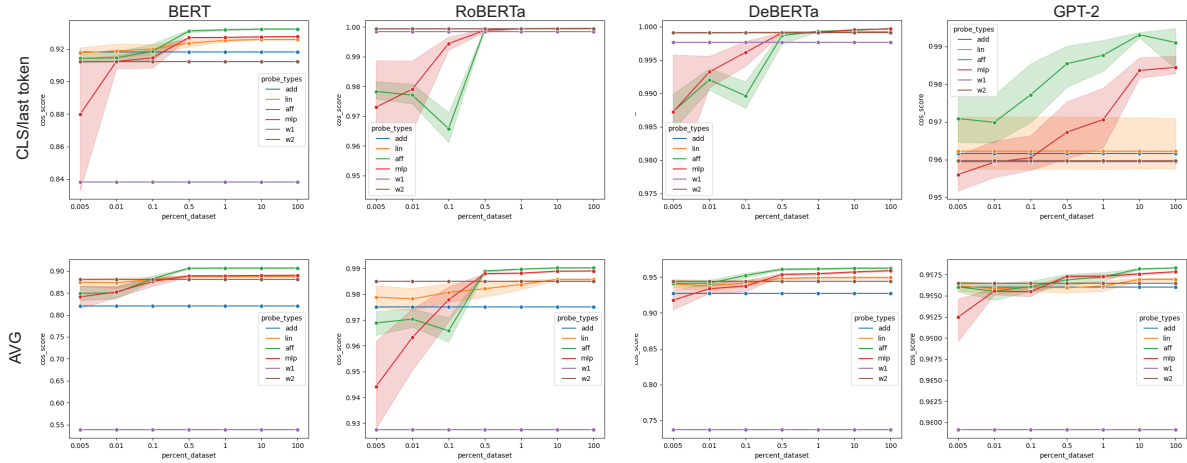


Figure 7: Learning curves of approximative probes trained on differing percentages of train data.

Model and representation	Spearman ρ	p-val
BERT _{CLS}	-0.0700	0.0
RoBERTa _{CLS}	0.1659	0.0
DeBERTa _{CLS}	0.1166	0.0
BERT _{AVG}	0.7143	0.0
RoBERTa _{AVG}	0.7086	0.0
DeBERTa _{AVG}	0.7866	0.0

Table 8: Spearman ρ correlation between phrase length (in words) and compositionality score in the treebank.

F Annotation setup and instructions

Annotators were recruited from a population of graduate students. Initially, 6 annotators completed the pilot experiment, which consisted of 101 examples. The subset of three annotators with highest agreement was asked if they would like to complete the full study. One annotator in the highest-agreement group could not continue to the full study, so this annotator was excluded, and the next group with highest agreement was chosen. The agreement values in [subsubsection 6.1.1](#) are for the final group of annotators chosen.

The experiment was implemented on the Qualtrics platform, and participants were first presented with a consent form, linking to more background information on the study, and informing them that their participation was entirely voluntary. After agreeing to the terms, participants were shown some examples and went through 3 practice questions. The example given are shown in [Figure 8](#), and the annotation interface is shown in [Figure 9](#) and [Figure 10](#). After completing the practice section, annotators began annotating the real examples, which followed the same interface as the practice examples.

Annotators were all located in the United States, paid approximately \$15 per hour for their work.

Examples

The following examples illustrate some examples of compositionality. Compositionality means that you can understand the meaning of a phrase from its parts.

Ivory tower

- This means that someone or something is out of touch with ordinary people. It doesn't mean "a white tower", and you wouldn't know what this means unless you came across it before, so it should be marked as **non-compositional**.

Balance sheet

- This means a spreadsheet that someone calculates ("balances") their finances on. Its meaning can be inferred once you know what it is, but it might not be obvious right away, so it should be marked as **somewhat compositional**.

Brown dog

- This is a dog which is brown. You can fully figure out the meaning just from the two words, so it is **fully compositional**.

←

→

Figure 8: Examples of compositionality judgments shown to annotators

Consider the following phrase:
Raining cats and dogs

If the phrase has both a literal and idiomatic meaning, please consider the idiomatic meaning. E.g. "raining cats and dogs" could mean literal cats and dogs falling out of the sky, but you should consider the usual meaning. Some of the phrases are quite rare, so you should search up the meaning if you don't know it.

Please consider the two parts of this phrase individually:
Raining
Cats and dogs

Consider the most typical meaning of the two parts of the phrase.

How well can you understand the phrase by combining the most typical meaning of the two parts of the phrase?

1. Not at all - you cannot understand the phrase from the typical meanings of its two parts.

2. Somewhat - you can understand the phrase somewhat. You may have to guess its meaning, or one of its parts is used in an atypical or figurative way.

3. Fully - you can completely understand the phrase by understanding the typical meanings of its two parts.

1 - Not compositional

2 - Somewhat compositional

3 - Fully compositional

Figure 9: First page of annotation interface for a practice phrase

Knowing the final meaning of the phrase **Raining cats and dogs**, how much do you think **Raining** contributes to the final meaning?

1 - Not at all

2 - Somewhat

3 - A great deal or fully

Knowing the final meaning of the phrase **Raining cats and dogs**, how much do you think **Cats and dogs** contributes to the final meaning?

1 - Not at all

2 - Somewhat

3 - A great deal or fully

Figure 10: Second page of annotation interface for a practice phrase