

# Revisiting subword tokenization: A case study on affixal negation in large language models

Anonymous ACL submission

## Abstract

In this work, we measure the impact of affixal negation on modern English large language models (LLMs). In affixal negation, the negated meaning is expressed through a negative morpheme, which is potentially challenging for LLMs as their tokenizers are often not morphologically plausible. We conduct extensive experiments using LLMs with different subword tokenization methods, which lead to several insights on the interaction between tokenization performance and negation sensitivity. Despite some interesting mismatches between tokenization accuracy and negation detection performance, we show that models can, on the whole, reliably recognize the meaning of affixal negation.

## 1 Introduction

Negation is central to language understanding but is not properly captured by modern NLP methods (Hossain et al., 2022; Truong et al., 2023, inter alia). While state-of-the-art large language models (LLMs) have improved negation understanding capabilities, challenges remain, such as the ability to correctly determine the enclosed scope of negation, or when negation interacts with other linguistic constructions like quantifiers (She et al., 2023; Truong et al., 2023). Negations in common NLP benchmarks are typically marked by separate negation cues such as *not*, *no*. However, in practice, negation can also be expressed through morphemes of words, i.e. by negative prefixes or suffixes such as in *uninteresting* or *effortless*.

While humans can identify affixal negation by leveraging morphological cues, NLP systems only rarely consider word-internal structure, beyond normalizing syntactic variation (Liu et al., 2012). Modern NLP methods such as language models employ subword tokenization, in which words are broken down into smaller units. This has an advantage of reducing vocabulary size, as well as learning

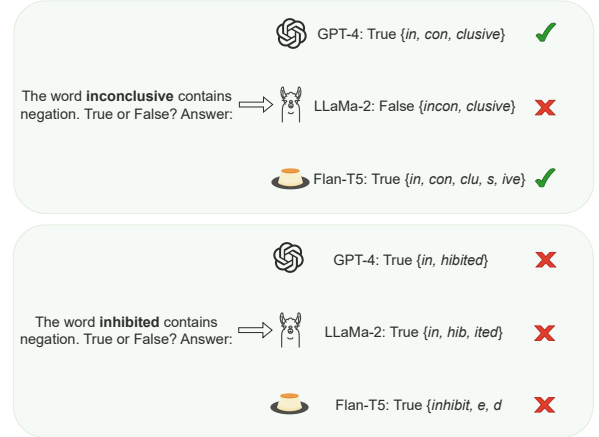


Figure 1: Example of our affixal negation prediction task, with the tokenization output for each model.

shared representation between words with similar subwords. The intent to improve such representation by making tokenization methods more linguistically sound has driven the invention of several morphology segmentation methods, such as Morfeessor (Grönroos et al., 2014). However, these have not been broadly adopted in modern LLMs as they do not scale well.

We hypothesize that current subword tokenization methods could lead to sub-optimal performance on language understanding tasks involving negation, because they do not correctly break words down morphologically. For instance, Table 1 demonstrates how different models employing different subword tokenization methods tokenize the word *anticlinal*. Another known challenge which could affect models is the high false positive rate in detecting affixal negations (Blanco and Moldovan, 2011), for example misinterpreting "de" in "deserve" as a negative affix while the word just coincidentally starts with "de" and should not be interpreted as negating "serve".

In this work, we analyze the impact of affixal negations on transformer-based language models, where two main tokenization methods are em-

Model	Type	Variant	Data source	Vocab. size	Output	NegMorph
BERT	BPE	WordPiece	books, wiki	30K	{anti, clin, al}	Correct
RoBERTa	BPE	Byte-level BPE	books, wiki	50K	{antic, l, inal}	Under-segmented
XLNet	ULM	SentencePiece	book, wiki, web text	32K	{anti, clin, al}	Correct
ALBERT	ULM	SentencePiece	book, wiki	32K	{anti, clin, al}	Correct
T5	BPE	SentencePiece	web text	32K	{anti, clin, al}	Correct
Llama-2	BPE	SentencePiece	web text, code, books, wiki, scientific publications	32K	{ant, ic, l, inal}	Over-segmented
GPT-2	BPE	Byte-level BPE	web text	50K	{antic, l, inal}	Under-segmented
GPT-4	BPE	Byte-level BPE	undisclosed	100K	{antic, l, inal}	Under-segmented

Table 1: Summary of different tokenizers used in our experiments. Output are tokenized version of the word “anticlinal” (model-specific special tokenization characters are removed for clarity purpose). All models are the base version unless specified otherwise.

played including Byte-pair encoding (BPE) (Gage, 1994; Sennrich et al., 2016), and Unigram language model (Unigram LM) (Kudo, 2018). We consider three research questions:

**RQ1: Are current subword tokenization methods able to produce morphologically-aligned tokens?** We analyzed the performance of various subword tokenization methods used in modern LMs. We find that most do not effectively produce morphologically-aligned tokens.

**RQ2: Are modern LMs aware of the presence of negation in affixal negations?** We design a negation prediction task to probe models’ awareness of affixal negation. We find that despite not performing well on the tokenization task, current LLMs can reliably infer the negated meaning of words with negative affix. For this task, there is only a weak positive correlation between tokenizer and classifier performance.

**RQ3: What are the impacts of affixal negation on downstream tasks?** As negation and sentiment are closely related, we measure the impact on downstream sentiment analysis task by looking at samples containing affixal negations from common datasets. Results show that models perform well on those samples, implying that the impact of affixal negation is minimal. However, there exists a bias in predicting negative sentiment for affixal negations.

## 2 Related work

There are two contrasting ways to construct a vocabulary for LMs using subword tokenization methods: BPE, which starts from a base character set, then merges those characters based on bigram frequency to form subword units (bottom-up) and un-

igram language models, which start from a large subword vocabulary, which is then reduced based on a regularization method (top-down). There are multiple variants of BPE, differing in how the base vocabulary is represented or how merging is done. WordPiece (Schuster and Nakajima, 2012) uses characters to represent the base vocabulary, then selects pairs that maximize the likelihood of training data, Byte-level BPE (Sennrich et al., 2016) uses bytes instead of Unicode to represent the base vocabulary; the merging is done based on the frequency count of bigrams. In contrast, the unigram language model (Kudo, 2018) starts from a large base vocabulary and iteratively trims down tokens based on unigram LM perplexity until a target vocabulary size is reached.

Both methods assume that the input text uses spaces to separate words, which is not true for languages such as Chinese or Vietnamese. Therefore, a word segmentation step must be performed in advance. SentencePiece (Kudo and Richardson, 2018) was introduced to solve this problem by considering whitespace as part of words, essentially treating the whole input stream as the smallest unit to perform tokenization on. Then, either BPE or unigram LM can be applied to construct the vocabulary. Regardless of methods, they purely rely on statistical information and thus are not expected to produce morphologically-aligned subword tokens.

There have been efforts to build linguistically-sound word tokenization methods, most notably the Morfessor and its variants (Grönroos et al., 2014, 2020). Building morphology-aligned segmentation methods, especially in a multilingual setting, is an active line of research through recent SIGMORPHON shared tasks (Batsuren et al., 2022). These methods outperform general tokeniz-

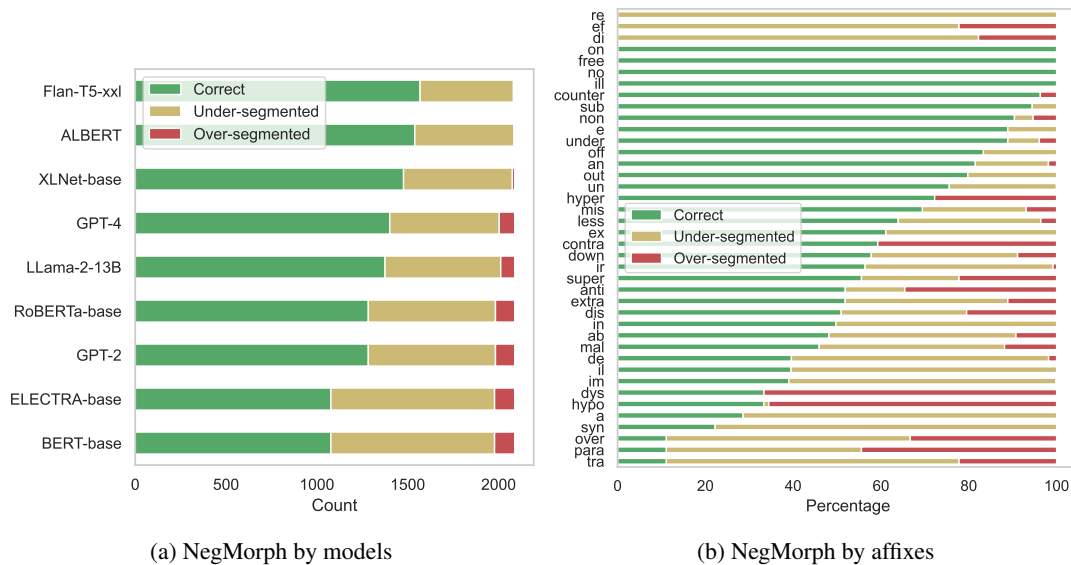


Figure 2: Word segmentation performance on the set of affixal negations (van Son et al., 2016).

ers in producing morphologically-aligned tokens, but their benefit on downstream tasks is often negligible (Domingo et al., 2019; Saleva and Lignos, 2021). In this work, we examine if morphologically correct tokenization is important for LLMs to deal with negation.

BERT and its variants have been shown to be insensitive to negation (Kassner and Schütze, 2020; Ettinger, 2020), affecting many downstream NLP tasks such as sentiment analysis, NLI, or QA (Hossain et al., 2020, 2022; Ravichander et al., 2022; Truong et al., 2022). Compared to previous models, current LLMs have improved negation handling ability, but still struggle with some unconventional types of negation and linguistic constructions (Truong et al., 2023). Here, we investigate the treatment of affixal negation in modern LMs, with the intuition that subword tokenization methods that don’t appropriately reflect this morphology will lead to misinterpretation of their semantics.

### 3 Experiment settings

We focus our analysis particularly on how affixal negations are represented in modern LLMs, designing probing tasks to test their awareness of negation, and the effect on downstream tasks.

#### 3.1 A lexicon of affixal negation

We use the lexicon created in van Son et al. (2016). The dataset contains a list of affixal negation and their non-negated counterparts (e.g. *unintended* - *intended*). For each affixal negation, the corresponding negative affix is also annotated. In total,

there are 2089 affixal negations, and 2055 non-negated words which are antonyms of the negations. These numbers are not equal because one word can have multiple corresponding negated counterparts, e.g. *intrusive* - *extrusive*, *unintrusive*.

#### 3.2 Tokenization methods

For each tokenizer type (along with their variants), we consider the most representative models that use them, based on their popularity. Although some models use the exact same tokenizer, it is worth investigating them as the difference in training corpora could lead to difference in tokenization results.

**BPE** We consider models using with different flavors of BPE. For WordPiece, we consider BERT (Devlin et al., 2019), ELECTRA (Clark et al., 2020); Byte-level BPE: RoBERTa (Liu et al., 2019), and GPT-family models including GPT-2 (Radford et al., 2019) and GPT-4 (OpenAI, 2023). For SentencePiece, we examine Flan-T5 (Chung et al., 2022) and Llama-2 (Touvron et al., 2023).

**Unigram LM** Models using unigram LM tokenization methods considered in this work are always used in combination with SentencePiece: XLNet (Yang et al., 2019), AIBERT (Lan et al., 2020).

#### 3.3 Morphologically-aligned segmentation for affixal negation

We consider a segmentation to be aligned with morphology (**Correct**) only if the negative affix matches with one of the produced tokens (e.g. *anti-climatic* → *anti*, *clima*, *tic*). Otherwise, it is either

**Under-segmented** if the negative affix is a substring of one of the produced tokens (e.g. *anticlima, tic*), or **Over-segmented** (e.g. *ant, i, clima, tic*). In a formal way, given an affixal negation word  $w$  having the negative affix  $a$ , if  $w$  is tokenized into  $T_k = \{t_i, t_{i+1}, \dots, t_n\}$  under tokenizer  $k$  then we define  $\text{NegMorph}_k(w)$  as follows:

$$\text{NegMorph}_k(w) = \begin{cases} \text{Correct} & \text{if } a \in T_k. \\ \text{Under-segmented} & \text{if } a \text{ is a substring of any } t_i \in T_k. \\ \text{Over-segmented} & \text{otherwise} \end{cases}$$

## 4 Findings

### 4.1 Current subword tokenization methods are morphologically suboptimal

As shown in Figure 2a, models employing the unigram LM method outperform those using BPE in producing morphologically correct tokens for affixal negations. This is in line with previous findings that the unigram LM produces subword units that align with morphology better than BPE (Bostrom and Durrett, 2020). Moreover, models that employ SentencePiece (T5, ALBERT, XLNet, LLaMa) outperform those that don't (BERT, RoBERTa, GPT-2). However, the best performing models only produce up to 75% correct NegMorph, showing moderate room for improvement. Most failed cases relate to under-segmentation.

An analysis of what types of negative affixes are hard to tokenize is provided in Figure 2b, and their most frequent incorrect tokenizations are shown in Figure 3. Some common affixes that are incorrectly tokenized are  $il \rightarrow ill$  (*illicit, illogical*),  $ir \rightarrow irre$  (*irresolute, irreposibly, irregular*),  $a \rightarrow as$  (*asymmetric*),  $at$  (*atypically*). Overall, we see that some affixes can be wrongly tokenized in a wide range of ways (represented by the large number of substacks), showing that current tokenization methods are inefficient. Overcoming this problem would greatly reduce the vocabulary size of LLMs, as well as improve word representations.

### 4.2 Negative affix is a signal for negation, but word knowledge is essential

We design a binary classification task on the lexicon described in Section 3.1 to probe the ability of models to understand affixal negation, denoted *Affix*. The prompts are captured below.

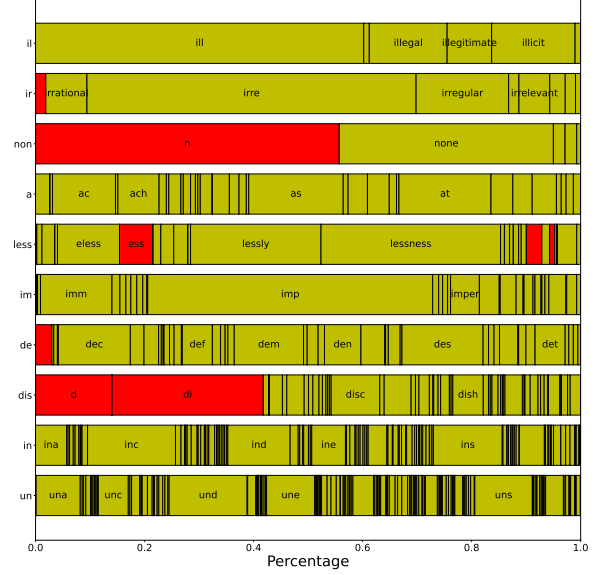


Figure 3: Top 10 most frequent affixes in the dataset and the distribution of tokens that they are wrongly tokenized into. Yellow bar denotes Under-segmented, while Red bar denotes Over-segmented.

#### Affix (zero-shot)

The word {word} contains negation. True or False?  
Answer:

#### Affix (few-shot)

A word contains negation if it has a negated meaning, usually expressed through a negative prefix (such as un, in) or suffix (such as less).

The word unwell contains negation. True or False?  
Answer: True  
Explanation: decentralize is created by prepending the root word centralize with the negative prefix de.

The word deserve contains negation. True or False?  
Answer: False  
Explanation: deserve just coincidentally starts with de.

The word {word} contains negation. True or False?  
Answer:

For a few-shot prompt, we provide an explicit instruction to explain what negation means in this context, as well as two demonstrating samples, to prevent any ambiguity (such as confusion with negative sentiment).

We evaluate three state-of-the-art LLMs in a zero-, and few-shot manner and breakdown the results into two categories: Neg (only affixal negations), Non-neg (only non-negated words). Results are summarized in Figure 5 (full results in Table 3).

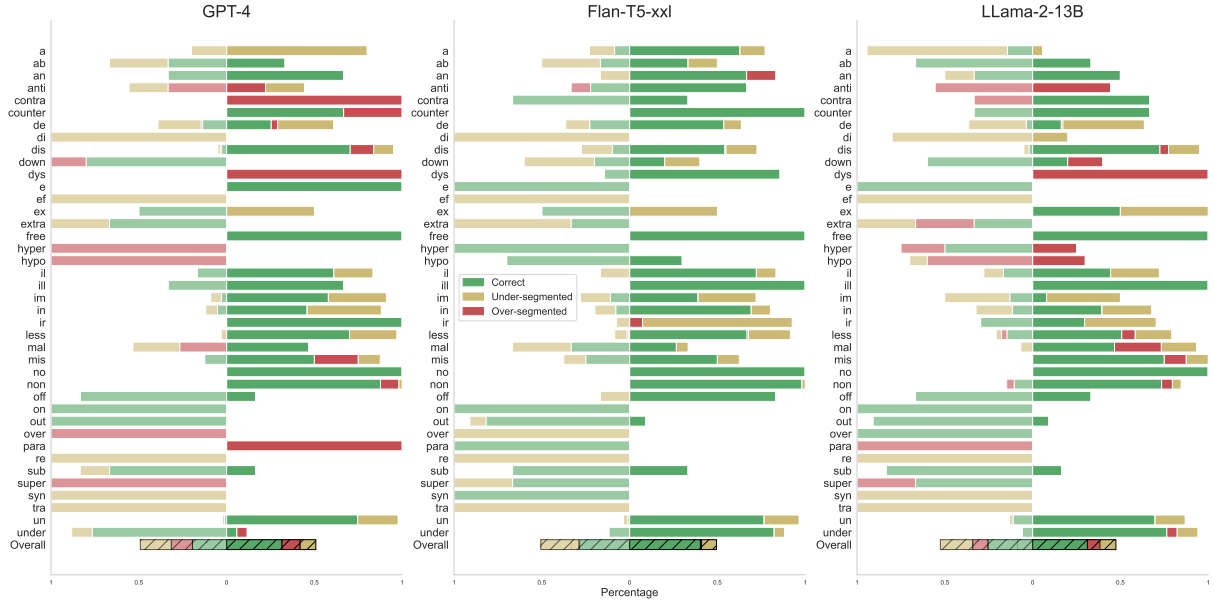


Figure 4: Fraction of correct/incorrect prediction on the Affix (fewshot) task, breakdown by affixes. The left greyed-out side of each subplot corresponds to wrong predictions.

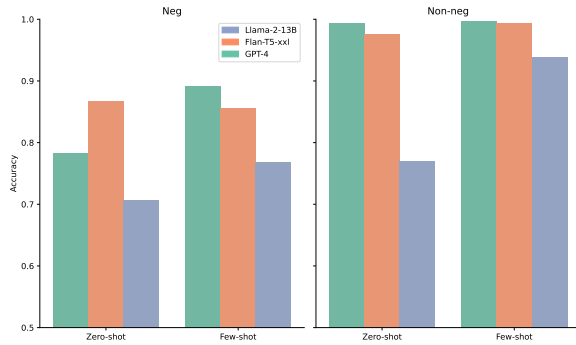


Figure 5: Zero- and Few-shot results on the affixal negation prediction task.

Overall, we find that the performance on Neg is much lower compared to Non-neg, where the best models achieve near-perfect performance.

For zero-shot setting, surprisingly, Flan-T5 outperforms both Llama-2 and GPT-4 on this prediction task, despite being smallest in size. After adding more explicit instruction and examples (Affix (fewshot)), we observe large increases in performance for GPT-4 and Llama-2, and little to no difference for Flan-T5. Whereas for the non-negated subset, all models have near-perfect performance, with GPT-4 slightly outperforming Flan-T5. Llama-2 performance for this task is much lower compared to the other two.

We further breakdown the results based on affixes. Figure 4 illustrates the percentage of correct/incorrect prediction for each affix, divided by

NegMorph categories. Compared to the relatively high results for Neg in Table 3, we have a clearer view on the actual performance of models. On average, we see that models made errors equally as likely for all affixes (as shown by the last Overall bar, where the percentages of incorrect and correct predictions are roughly 50%). From the figure, we can also observe that the correct/incorrect prediction distribution is similar across models (especially between GPT-4 and Flan-T5), showing that they tend to make the same errors. Moreover, the percentage of correct segmentation is larger in cases where they made correct prediction for the Affix task. However, calculating the Pearson’s coefficient between NegMorph and Accuracy on Neg set did not yield any statistically significant correlation.

**Hyphenated words** To make sure that the negative affixes are not further broken down by tokenizers, we convert words into their “hyphenated” form (e.g. *unintended* → *un-intended*). From Figure 6, we see that this greatly increases the performance of different tokenizers on the NegMorph metric (by as much as 32%). Compared to the normal setting, the accuracy of all models also increases on the Affix task, suggesting a positive correlation between NegMorph and Accuracy. Llama-2 benefited the most from this setting, having the largest increases in both Accuracy and NegMorph.



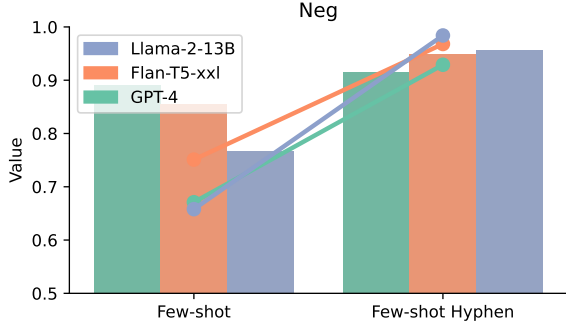


Figure 6: Results of Few-shot and Few-shot Hyphen on the affixal negation prediction task. Bars denote the accuracy on the prediction task, while Dots denote the Correct NegMorph scores for the segmentation task.

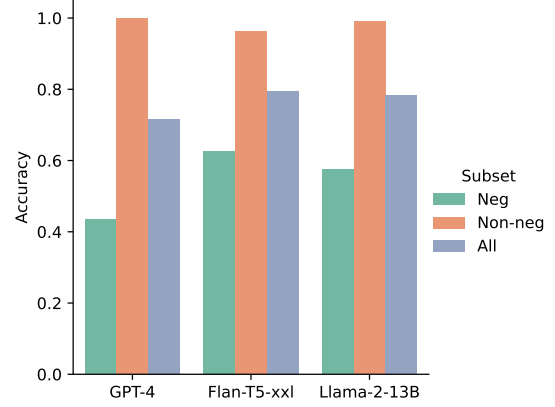


Figure 7: Accuracy on the affixal nonce words prediction task.

**Nonce words** Nonce words are words that look and sound like real words, but are created for a single-purpose use and not recognized as words within a language (e.g. *roagly*). To measure the effect of negative affixes on word semantics, we construct a list of “affixal nonce words” by prepending or appending negative affixes to a list of nonce words. We collect a list of adjective nonce words from [Cremers \(2022\)](#). For affixes, we used the list of 40 negative affixes provided in [van Son et al. \(2016\)](#) and collected 40 non-negative affixes (e.g. *auto-*, *bi-*, *-ism*, *-ful*).<sup>1</sup> For each nonce word, we prepend (or append) the affixes to form “affixal nonce word”. In total, the set consists of 11 nonce words  $\times$  80 affixes = 880 samples, evenly distributed between negated (e.g. *dis-roagly*) and non-negated (e.g. *auto-roagly*). We adopt the Affix (few-shot) prompt and add an instruction to prevent models from refusing to answer the questions because of invalid words (full prompt in Appendix B). Similarly, we also report the results of two subsets of negative affixes (Neg) and non-negative affixes (Non-neg) in Figure 7. For the Neg set, we find that for the performance of all models is relatively low, despite them being able to correctly tokenize the negative affixes. Whereas for the Non-neg set, performances are near-perfect for all models, similar to the previous Affix-Hyphen task. Looking at the results, however, we found that most errors made by the models are when the negative affixes are ambiguous, i.e. their meaning depends on which words they are attached to (e.g. *a-*, *di-*, *ef-*, *para-*, *re-*). This reveals an important insight that whether something is considered to be a negation should be

judged with context (which is parametric knowledge about words in this case).

**Non-negated words with tokens homonymous with negative affixes** To explore the false positive problem raised in [Blanco and Moldovan \(2011\)](#), we collect words from the Non-neg subset and WordNet ([Miller, 1995](#)) which do not have negated meaning, but have negative prefix/suffix as the first/last subword token. We tokenize WordNet using the Flan-T5 tokenizer and select all words that start/end with the negative prefixes/suffixes, then subtract all words in the list of affixal negations. We manually go through the extracted list again to remove errors, resulting in a set of 330 words<sup>2</sup>. Following the same affixal negation prediction task, we find that Flan-T5 has very good performance (0.958 accuracy), showing that it can synthesize information from all subword tokens instead of only relying on the negative affixes. Most errors come from the “uni-” cases, where model tokenized them into “un-” (e.g. *unidirectional*, *univalent*).

### 4.3 Impact on downstream tasks

One main drawback of our probing task is that the words lack context. Negation is a context-dependant concept, that is, what is considered a negation could differ depending on the context of use. Investigating the impact of affixal negation on downstream tasks is essential.

<sup>1</sup>We collected the affixes from <https://litinfocus.com/120-root-words-prefixes-and-suffixes-pdf-list/>

<sup>2</sup>We didn’t consider other models as the list of words would be different between models.

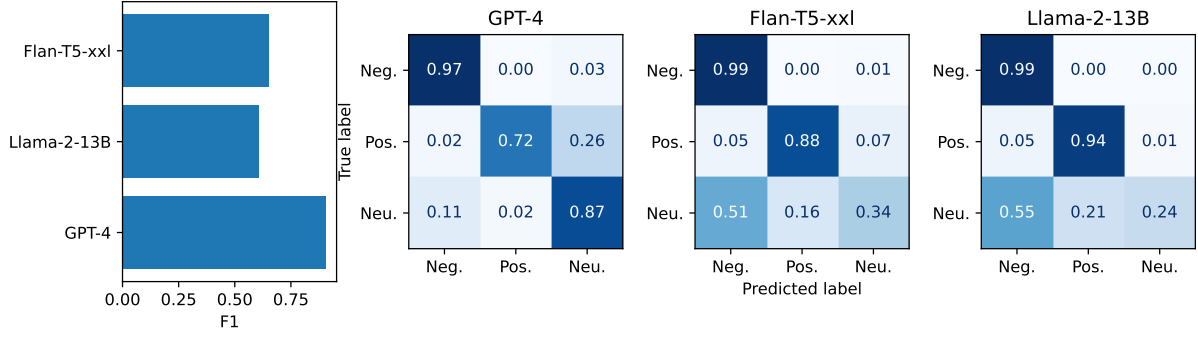


Figure 8: Performance on word-level sentiment task

### 4.3.1 Sentiment analysis

Previous works have shown that negation is a strong indicator for negative sentiment (Wiegand et al., 2010). Furthermore, the fact that sentiment analysis is part of many NLP benchmarks could create a bias in models which leads to negation being conflated with negative sentiment. For instance, the word *incredible* is constructed by prepending the root word *credible* with the negative affix *in-*, meaning “not credible” but is used to express a positive meaning. This inspired us to extend our analysis to the downstream sentiment analysis task. We evaluate LLMs few-shot performance in two settings of word- and sentence-level sentiment analysis (full prompts in Appendix C).

**Word-level sentiment** Using SentiWordNet 3.0 (Baccianella et al., 2010), we map the sentiment to the affixal negation lexicon from Section 3.1. After that, two graduate researchers went over the list to determine the final labels (positive, negative, neutral). In general, we find that GPT-4 outperforms Flan-T5 and Llama-2 on this word-level task. All models have almost perfect performance on predicting negative words, but struggle with the other two classes. In particular, we find Flan-T5 and Llama-2 overpredict Negative for Neutral words, while GPT-4 often mistakes Positive for Neutral.

**Sentence-level sentiment** For this task, we look at common sentence-level sentiment analysis datasets including SST-2 (Socher et al., 2013), and Rotten Tomatoes (RT) (Pang and Lee, 2005). One drawback of this evaluation is that samples tend to contain many sentiment signals, making it hard to gauge the effect of affixal negations.

We consider 3 settings, (1) Affix: only samples containing affixal negation; (2) Non affix: only samples without affixal negation; (3) Replace affix:

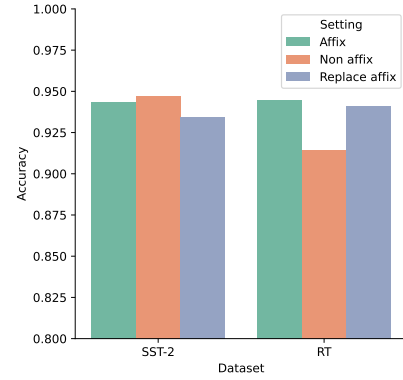


Figure 9: Accuracy on sentence-level sentiment analysis task. Results are average across 3 models.

similar to Affix, but we replace all instances of affixal negations with equivalent syntactic negations, i.e. *not* + word (*uninteresting* → *not interesting*), and summarize the results in Figure 9. Note that the numbers of samples in Non affix are much larger than Affix for both datasets. Overall, we can conclude that affixal negation is a strong signal to guide models’ prediction. We observe good performance for Affix in both datasets, where the accuracy are comparable to Non Affix in SST-2 and higher in RT. Attempting to replace affixal negations would slightly decrease the performance of models in both datasets. This suggests that affixal negation is actually a stronger sentiment cue compared to syntactic negation. We further report class-wise performance of the Affix set in Figure 10. Accuracy on samples having Negative sentiment is higher than Positive, once again showing that affixal negation is a strong cue for predicting negative sentiment.

## 5 A look into token attribution

We perform an interpretation analysis to give insight into what drives models’ predictions. For this analysis, we chose the Flan-T5-xxl model as GPT-

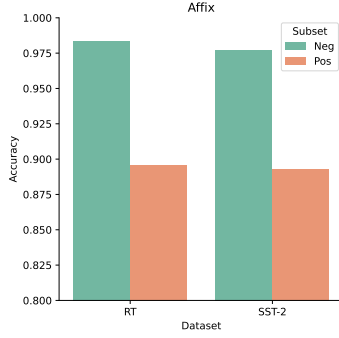


Figure 10: Accuracy of Neg/Pos class of the Affix set. Results are average across 3 models.

4 predictions are not interpretable. We calculate the attribution for each token corresponding to the predictions using the Integrated Gradient method (Sundararajan et al., 2017), with probability as the scoring function, implemented in *Inseq* (Sarti et al., 2023). Overall, we observe high attribution scores from relevant tokens, such as the subword tokens of the target words, showing that models know where to pay attention to when making inference.

			__negative			</s>											
			__The			0.026			0.046								
			__sentiment			0.062			0.104								
__positive			</s>			__of			0.033			0.042					
__The			0.029			0.049			__the			0.04			0.029		
__sentiment			0.092			0.115			__word			0.071			0.036		
__of			0.028			0.052			__un			0.202			0.036		
__the			0.034			0.037			interest			0.165			0.05		
__word			0.076			0.049			ing			0.054			0.022		
__interesting			0.198			0.042			__is			0.029			0.019		
__is			0.058			0.023											
									__negative			</s>					
			__The			0.024			0.046								
			__sentiment			0.061			0.104								
			__of			0.021			0.045								
			__the			0.04			0.03								
			__word			0.052			0.034								
			__un			0.12			0.034								
			i			0.072			0.02								
			ntended			0.266			0.064								
			__is			0.039			0.021								

Figure 11: Token attribution of selected samples on word-level sentiment prediction task. Only parts of the prompts are shown for clarity purpose.

### Negative affixes have flipping sentiment effect

In Section 4.3.1, we showed that models tend to overpredict negative sentiment on the list of affixal negations. Through the saliency heatmap in Figure 11, we can see high token attributions of the negative affixes that changed the sentiment of the root words (either positive or neutral) into negative.

This is in-line with previous finding that negation flips sentiment direction (Tigges et al., 2023). This effect could be the main cause for the low performance on Neutral class observed in our word-level sentiment analysis task. When applying to the negation prediction task, however, we did not observe a similar effect and did not see any clear pattern for token attributions.

### Correct tokenization is not essential for negation awareness

Through many experiments, we have shown that overall, correct tokenization leads to better awareness of models to the presence of negation. This effect, however is not significant. By comparing token attributions between 3 cases of NegMorph (Figure 12), we saw that models are able to combine information from relevant subword tokens corresponding to a word to make the correct inference.

__True </s>			__Fal </s>			__True </s>		
__The	0.035	0.03	__The	0.029	0.021	__The	0.036	0.035
__word	0.043	0.049	__word	0.041	0.03	__word	0.045	0.052
__anti	0.053	0.091	__male	0.09	0.042	__	0.021	0.033
clin	0.078	0.073	vol	0.077	0.059	d	0.023	0.037
al	0.029	0.039	ence	0.036	0.026	issent	0.1	0.09
__contains	0.061	0.079	__contains	0.058	0.03	__contains	0.068	0.046
__neg	0.13	0.042	__neg	0.188	0.066	__neg	0.124	0.047
ation	0.054	0.023	ation	0.053	0.031	ation	0.051	0.028
.	0.05	0.028	.	0.041	0.034	.	0.05	0.032

Figure 12: Token attribution of selected sample samples on negation prediction task. Three subplots correspond to Correct, Under-segmented, and Over-segmented case respectively. Only parts of the prompts are shown for clarity purpose.

## 6 Conclusion

In this work, we conduct an in-depth analysis on how well modern LLMs handle affixal negation, a type of negation where morphology is essential to understanding word semantics. We have shown that there is significant room to improve current tokenization methods in producing morphologically-aligned tokens. Despite that, the effect of morphologically incorrect tokenization on the ability of models to understand word meaning in downstream tasks including sentiment analysis is minimal. Regardless, designing better subword tokenization methods may have many immediate benefits such as reducing vocabulary size, learning better word representations, and improving models' interpretability.



## 7 Limitations

**Prompting** As this work involves experiments using LLMs, there is always a possibility that the prompts we used are not optimal (and also, the problem of reproducibility). We attempted to reuse prompts templates from existing works where possible and strove to design prompts that are intuitive and specific otherwise.

**Multilinguality** Morphology is a language-dependent problem. We recognize that the lack of investigation in other languages other than English is a drawback of this work.

**Broader impact** Given that our focus is on presenting and analysing the problem of poor treatment of affixal negation in LLMs, we did not propose any immediate solutions to improve the status quo. The finding on the impact on downstream tasks could be limited by the lack of samples (both in size and meaningful patterns) in the test data.

## References

Stefano Baccianella, Andrea Esuli, and Fabrizio Sebastiani. 2010. [SentiWordNet 3.0: An enhanced lexical resource for sentiment analysis and opinion mining](#). In *Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC'10)*, Valletta, Malta. European Language Resources Association (ELRA).

Khuyagbaatar Batsuren, Gábor Bella, Aryaman Arora, Viktor Martinovic, Kyle Gorman, Zdeněk Žabokrtský, Amarsanaa Ganbold, Šárka Dohnalová, Magda Ševčíková, Kateřina Pelegrinová, Fausto Giunchiglia, Ryan Cotterell, and Ekaterina Vylomova. 2022. [The SIGMORPHON 2022 shared task on morpheme segmentation](#). In *Proceedings of the 19th SIGMORPHON Workshop on Computational Research in Phonetics, Phonology, and Morphology*, pages 103–116, Seattle, Washington. Association for Computational Linguistics.

Eduardo Blanco and Dan Moldovan. 2011. Some issues on detecting negation from text. In *Twenty-Fourth International FLAIRS Conference*.

Kaj Bostrom and Greg Durrett. 2020. [Byte pair encoding is suboptimal for language model pretraining](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4617–4624, Online. Association for Computational Linguistics.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. [Scaling instruction-finetuned language models](#). *ArXiv preprint*, abs/2210.11416.

Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. 2020. [ELECTRA: pre-training text encoders as discriminators rather than generators](#). In *8th International Conference on Learning Representations (ICLR 2020)*.

Alexandre Cremers. 2022. Interpreting gradable adjectives: rational reasoning or simple heuristics? *Empirical Issues in Syntax and Semantics*, 14:31–61.

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, Minneapolis, Minnesota. Association for Computational Linguistics.

Miguel Domingo, Mercedes García-Martínez, Alexandre Helle, Francisco Casacuberta, and Manuel Heranz. 2019. How much does tokenization affect neural machine translation? In *International Conference on Computational Linguistics and Intelligent Text Processing*, pages 545–554. Springer.

Allyson Ettinger. 2020. [What BERT is not: Lessons from a new suite of psycholinguistic diagnostics for language models](#). *Transactions of the Association for Computational Linguistics*, 8:34–48.

Philip Gage. 1994. A new algorithm for data compression. *C Users Journal*, 12(2):23–38.

Stig-Arne Grönroos, Sami Virpioja, and Mikko Kurimo. 2020. [Morfessor EM+Prune: Improved sub-word segmentation with expectation maximization and pruning](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 3944–3953, Marseille, France. European Language Resources Association.

Stig-Arne Grönroos, Sami Virpioja, Peter Smit, and Mikko Kurimo. 2014. [Morfessor FlatCat: An HMM-based method for unsupervised and semi-supervised learning of morphology](#). In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, pages 1177–1185, Dublin, Ireland. Dublin City University and Association for Computational Linguistics.

Md Mosharaf Hossain, Dhivya Chinnappa, and Eduardo Blanco. 2022. [An analysis of negation in natural language understanding corpora](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 716–723, Dublin, Ireland. Association for Computational Linguistics.

Md Mosharaf Hossain, Venelin Kovatchev, Pranoy Dutta, Tiffany Kao, Elizabeth Wei, and Eduardo Blanco. 2020. [An analysis of natural language inference benchmarks through the lens of negation](#). In *Proceedings of the 2020 Conference on Empirical*

579	<i>Methods in Natural Language Processing (EMNLP)</i> ,	8729–8755, Abu Dhabi, United Arab Emirates. As-	634
580	pages 9106–9118, Online. Association for Computa-	sociation for Computational Linguistics.	635
581			
582	Nora Kassner and Hinrich Schütze. 2020. <a href="#">Negated and</a>	Jonne Saleva and Constantine Lignos. 2021. <a href="#">The ef-</a>	636
583	<a href="#">misprimed probes for pretrained language models:</a>	<a href="#">fectiveness of morphology-aware segmentation in</a>	637
584	<a href="#">Birds can talk, but cannot fly.</a> In <i>Proceedings of the</i>	<a href="#">low-resource neural machine translation.</a> In <i>Proceed-</i>	638
585	<i>58th Annual Meeting of the Association for Computa-</i>	<i>ings of the 16th Conference of the European Chap-</i>	639
586	<i>tational Linguistics</i> , pages 7811–7818, Online. Asso-	<i>ter of the Association for Computational Linguistics:</i>	640
587	ciation for Computational Linguistics.	<i>Student Research Workshop</i> , pages 164–174, Online.	641
		Association for Computational Linguistics.	642
588	Taku Kudo. 2018. <a href="#">Subword regularization: Improv-</a>	Gabriele Sarti, Nils Feldhus, Ludwig Sickert, and Os-	643
589	<a href="#">ing neural network translation models with multiple</a>	kar van der Wal. 2023. <a href="#">Inseq: An interpretability</a>	644
590	<a href="#">subword candidates.</a> In <i>Proceedings of the 56th An-</i>	<a href="#">toolkit for sequence generation models.</a> In <i>Proceed-</i>	645
591	<i>annual Meeting of the Association for Computational</i>	<i>ings of the 61st Annual Meeting of the Association</i>	646
592	<i>Linguistics (Volume 1: Long Papers)</i> , pages 66–75,	<i>for Computational Linguistics (Volume 3: System</i>	647
593	Melbourne, Australia. Association for Computational	<i>Demonstrations)</i> , pages 421–435, Toronto, Canada.	648
594	Linguistics.	Association for Computational Linguistics.	649
595	Taku Kudo and John Richardson. 2018. <a href="#">SentencePiece:</a>	Mike Schuster and Kaisuke Nakajima. 2012. Japanese	650
596	<a href="#">A simple and language independent subword tok-</a>	and korean voice search. In <i>International Conference</i>	651
597	<a href="#">enizer and detokenizer for neural text processing.</a> In	<i>on Acoustics, Speech and Signal Processing</i> , pages	652
598	<i>Proceedings of the 2018 Conference on Empirical</i>	5149–5152.	653
599	<i>Methods in Natural Language Processing: System</i>		
600	<i>Demonstrations</i> , pages 66–71, Brussels, Belgium.	Rico Sennrich, Barry Haddow, and Alexandra Birch.	654
601	Association for Computational Linguistics.	2016. <a href="#">Neural machine translation of rare words with</a>	655
		<a href="#">subword units.</a> In <i>Proceedings of the 54th Annual</i>	656
602	Zhenzhong Lan, Mingda Chen, Sebastian Goodman,	<i>Meeting of the Association for Computational Lin-</i>	657
603	Kevin Gimpel, Piyush Sharma, and Radu Soricut.	<i>guistics (Volume 1: Long Papers)</i> , pages 1715–1725,	658
604	2020. <a href="#">ALBERT: A lite BERT for self-supervised</a>	Berlin, Germany. Association for Computational Lin-	659
605	<a href="#">learning of language representations.</a> In <i>8th Inter-</i>	guistics.	660
606	<i>national Conference on Learning Representations</i>		
607	<i>(ICLR 2020)</i> .	Jingyuan S. She, Christopher Potts, Samuel R. Bowman,	661
		and Atticus Geiger. 2023. <a href="#">ScoNe: Benchmarking</a>	662
608	Haibin Liu, Tom Christiansen, William A Baumgart-	<a href="#">negation reasoning in language models with fine-</a>	663
609	ner, and Karin Verspoor. 2012. <a href="#">Biolemmatizer: a</a>	<a href="#">tuning and in-context learning.</a> In <i>Proceedings of the</i>	664
610	<a href="#">lemmatization tool for morphological processing of</a>	<i>61st Annual Meeting of the Association for Computa-</i>	665
611	<a href="#">biomedical text.</a> <i>Journal of Biomedical Semantics</i> ,	<i>tational Linguistics (Volume 2: Short Papers)</i> , pages	666
612	3:3.	1803–1821, Toronto, Canada. Association for Com-	667
		putational Linguistics.	668
613	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	Richard Socher, Alex Perelygin, Jean Wu, Jason	669
614	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,	Chuang, Christopher D. Manning, Andrew Ng, and	670
615	Luke Zettlemoyer, and Veselin Stoyanov. 2019.	Christopher Potts. 2013. <a href="#">Recursive deep models for</a>	671
616	<a href="#">RoBERTa: A robustly optimized BERT pretraining</a>	<a href="#">semantic compositionality over a sentiment treebank.</a>	672
617	<a href="#">approach.</a> <i>ArXiv preprint</i> , abs/1907.11692.	In <i>Proceedings of the 2013 Conference on Empiri-</i>	673
618	George A Miller. 1995. WordNet: a lexical database	<i>cal Methods in Natural Language Processing</i> , pages	674
619	for English. <i>Communications of the ACM</i> , 38(11):39–	1631–1642, Seattle, Washington, USA. Association	675
620	41.	for Computational Linguistics.	676
621	OpenAI. 2023. <a href="#">GPT-4 technical report.</a>	Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017.	677
622	Bo Pang and Lillian Lee. 2005. Seeing stars: Exploiting	Axiomatic attribution for deep networks. In <i>Interna-</i>	678
623	class relationships for sentiment categorization with	<i>tional conference on machine learning</i> , pages 3319–	679
624	respect to rating scales. In <i>Proceedings of the ACL.</i>	3328. PMLR.	680
625	Alec Radford, Jeff Wu, Rewon Child, David Luan,	Curt Tigges, Oskar John Hollinsworth, Atticus Geiger,	681
626	Dario Amodei, and Ilya Sutskever. 2019. Language	and Neel Nanda. 2023. <a href="#">Linear representations of</a>	682
627	models are unsupervised multitask learners. OpenAI	<a href="#">sentiment in large language models.</a>	683
628	blog.		
629	Abhilasha Ravichander, Matt Gardner, and Ana Maraso-	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	684
630	vic. 2022. <a href="#">CONDAQA: A contrastive reading com-</a>	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	685
631	<a href="#">prehension dataset for reasoning about negation.</a> In	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	686
632	<i>Proceedings of the 2022 Conference on Empirical</i>	Bhosale, et al. 2023. LLaMA 2: Open founda-	687
633	<i>Methods in Natural Language Processing</i> , pages	tion and fine-tuned chat models. <i>arXiv preprint</i>	688
		<i>arXiv:2307.09288.</i>	689

- Thinh Hung Truong, Timothy Baldwin, Karin Verspoor, and Trevor Cohn. 2023. [Language models are not naysayers: an analysis of language models on negation benchmarks](#). In *Proceedings of the 12th Joint Conference on Lexical and Computational Semantics (\*SEM 2023)*, pages 101–114, Toronto, Canada. Association for Computational Linguistics.
- Thinh Hung Truong, Yulia Otmakhova, Timothy Baldwin, Trevor Cohn, Jey Han Lau, and Karin Verspoor. 2022. [Not another negation benchmark: The NaN-NLI test suite for sub-clausal negation](#). In *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 883–894, Online only. Association for Computational Linguistics.
- Chantal van Son, Emiel van Miltenburg, and Roser Morante. 2016. [Building a dictionary of affixal negations](#). In *Proceedings of the Workshop on Extra-Propositional Aspects of Meaning in Computational Linguistics (ExProM)*, pages 49–56, Osaka, Japan. The COLING 2016 Organizing Committee.
- Michael Wiegand, Alexandra Balahur, Benjamin Roth, Dietrich Klakow, and Andrés Montoyo. 2010. [A survey on the role of negation in sentiment analysis](#). In *Proceedings of the Workshop on Negation and Speculation in Natural Language Processing*, pages 60–68, Uppsala, Sweden. University of Antwerp.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime G. Carbonell, Ruslan Salakhutdinov, and Quoc V. Le. 2019. [XLNet: Generalized autoregressive pretraining for language understanding](#). In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 5754–5764.

## A Models endpoints

- GPT-4: We access GPT-4 through the official API with the name `gpt-4`. Note that this is different from the GPT-4 turbo model with the name `gpt-4-1106-preview`.
- Llama-2-13B: We used the official instruction fine-tuned Llama-2-13B available on the HuggingFace hub with the name: `meta-llama/Llama-2-13b-chat-hf`.
- Flan-T5-xxl: We used the official xxl version (11.3B) of the Flan-T5 model available on the HuggingFace hub with the name: `google/flan-t5-xxl`.

## B Details of Affixal Nonce word prediction task

**List of nonce words** *roagly, vibble, drok, scrop, plard, hif, tepable, plawic, bluth, sprat, flurf*

**List of non-negative affixes** Prefix: *ambi-, aqu-, ast-, aud-, auto-, bi-, bio-, cent-, circum-, co-, cred-, cycl-, dec-, dia-, equ-, geo-, grad-, hydro-, inter-, medi-, mega-, min-, micro-, pan-, semi-, tele-, uni-, tri-*. Suffix: *-able, -al, -ance, -ful, -ian, -ic, -tic, -ile, -ism, -ist, -junct, -ly*

### Nonce

A nonce word is a word occurring, invented, or used just for a particular occasion, or a word with a special meaning used for a special occasion. Infer whether the given nonce word contains negation or not.

A word contains negation if it has a negated meaning, usually expressed through a negative prefix (such as *un-, in-*) or suffix (such as *less*).

The word *unwell* contains negation. True or False?

Answer: True

Explanation: *decentralize* is created by prepending the root word *centralize* with the negative prefix *de-*.

The word *deserve* contains negation. True or False?

Answer: False

Explanation: *deserve* just coincidentally starts with *de-*.

The word {word} contains negation. True or False?

Answer:

## C Prompts for sentiment analysis

### Word-level sentiment

{Few-shot samples}

The sentiment of the word {word} is positive, negative, or neutral.

Answer:

### Sentence-level sentiment

{Few-shot samples}

{sentence}  
Question: Is this sentence positive or negative?

Answer:

## D Full results

Model	Neg. Nonce	Non-neg. Nonce	All
GPT-4	0.434	1	0.717
Llama-2-13B	0.575	0.991	0.783
Flan-T5-xxl	0.627	0.964	0.795

Table 2: Affixal nonce words prediction task

Model	Neg	Accuracy Non-neg	All	NegMorph Correct
<i>Affix (zero-shot)</i>				
GPT-4	0.783	0.994	0.888	0.671
Llama-2-13B	0.707	0.770	0.738	0.658
Flan-T5-xxl	0.867	0.976	0.921	0.751
<i>Affix (fewshot)</i>				
GPT-4	0.890 (+0.107)	0.997 (+0.003)	0.943 (+0.055)	0.670
Llama-2-13B	0.767 (+0.060)	0.938 (+0.168)	0.852 (+0.114)	0.658
Flan-T5-xxl	0.855 (-0.012)	0.993 (+0.017)	0.924 (+0.003)	0.750
<i>Affix (fewshot)-Hyphen</i>				
GPT-4	0.916 (+0.133)	-	-	0.929 (+0.258)
Llama-2-13B	0.956 (+0.249)	-	-	0.984 (+0.326)
Flan-T5-xxl	0.948 (+0.081)	-	-	0.968 (+0.217)

Table 3: Results of our affixal negation prediction task. (+/- denote the change compared to the *Affix (zero-shot)* setting