

Measuring Context-Word Biases in Lexical Semantic Datasets

Anonymous ACL submission

Abstract

001 State-of-the-art contextualized models eg.
002 BERT use tasks such as WiC and WSD to
003 evaluate their *word-in-context* representations.
004 This inherently assumes that performance in
005 these tasks reflect how well a model represents
006 the coupled word and context semantics. We
007 question this assumption by presenting the first
008 quantitative analysis on the context-word inter-
009 action required and being tested in major con-
010 textual lexical semantic tasks, taking into ac-
011 count that tasks can be inherently biased and
012 models can learn spurious correlations from
013 datasets. To achieve this, we run probing base-
014 lines on masked input, based on which we then
015 propose measures to calculate the degree of
016 context or word biases in a dataset, and plot
017 existing datasets on a continuum. The anal-
018 ysis were performed on both models and hu-
019 mans to decouple biases inherent to the tasks
020 and biases learned from the datasets. We found
021 that, (1) to models, most existing datasets fall
022 into the extreme ends of the continuum: the
023 retrieval-based tasks and especially the ones
024 in the medical domain (eg. COMETA) ex-
025 hibit strong target word bias while WiC-style
026 tasks and WSD show strong context bias; (2)
027 AM²ICO and Sense Retrieval show less ex-
028 treme model biases and challenge a model
029 more to represent both the context and target
030 words. (3) A similar trend of biases exists in
031 humans but humans are much less biased com-
032 pared with models as humans found seman-
033 tic judgments more difficult with the masked
034 input, indicating models are learning spuri-
035 ous correlations. This study demonstrates that
036 with heavy context or target word biases, mod-
037 els are usually not being tested for word-in-
038 context representations as such in these tasks
039 and results are therefore open to misinterpreta-
040 tion. We recommend our framework as a san-
041 ity check for context and target word biases in
042 future task design and model interpretation in
043 lexical semantics.

1 Introduction

044
045 Meaning contextualization (i.e., identifying the cor-
046 rect meaning of a target word in linguistic context)
047 is essential for understanding natural language,
048 and has been the focus in many lexical semantic
049 tasks. Pretrained contextualized models (PCMs)
050 have brought large improvements in these tasks in-
051 cluding WSD (Hadiwinoto et al., 2019; Loureiro
052 and Jorge, 2019; Huang et al., 2019; Blevins and
053 Zettlemoyer, 2020), WiC (Pilehvar and Camacho-
054 Collados, 2019; Garí Soler et al., 2019) and entity
055 linking (EL) (Wu et al., 2020; Broscheit, 2019).

056 These superior performances have been taken
057 as proof that PCMs can successfully model *word-*
058 *in-context* semantics. However, on one hand, the
059 evaluation benchmarks often vary in their emphasis
060 on context vs target words. For example, we could
061 expect tasks such as WSD and WiC to rely more
062 on context by design as the target words are either
063 given or the same in each input pair. Notice that
064 the exact amount of context/target word reliance in
065 these tasks is to be tested as humans naturally use
066 both to make prediction. On the other hand, models
067 may find shortcuts from datasets to avoid learning
068 the complex word-context interaction. **What is**
069 **missing in the current literature is an accurate**
070 **quantification of this word-context interplay re-**
071 **quired and being tested in each task so that we**
072 **can fully understand task goals and model per-**
073 **formance.** In particular, we need to flag heavy
074 word and context reliance where a model can solve
075 a task by relying solely on context or the target
076 words. Such heavy word or context reliance hinders
077 a scientific understanding of the models' meaning
078 contextualization abilities as it essentially bypasses
079 the key word-context interaction challenge in mean-
080 ing contextualization, which requires the modeling
081 of both target words and their contexts (Words are
082 frequently ambiguous, but so are contexts. In "*I*
083 *like XX.*", *XX* could have a number of meanings).

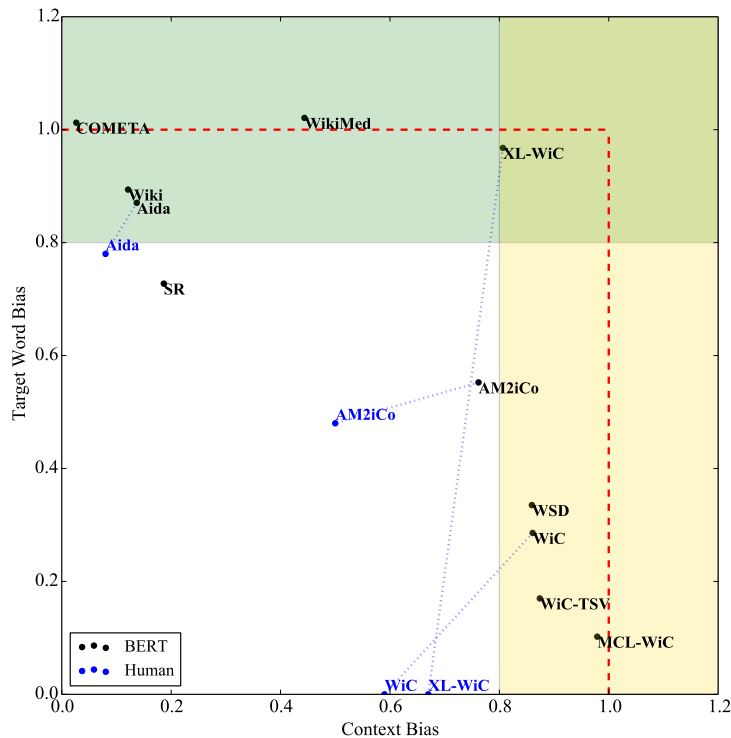


Figure 1: Plotting context and target word biases from BERT (blue) and humans (black) across popular context-aware lexical semantic datasets. The green shade and the yellow shade roughly indicate the areas for high target word bias and high context bias (>0.8). We would ideally want a dataset to lie towards the bottom left corner which is bias-free. The dashed red lines indicate 1.0 context (right) and 1.0 target word bias (top), implying a dataset is in effect dealt with by relying on target words alone or context alone.

Therefore, we refer to such heavy reliance on target words or context in a contextual lexical semantic dataset as target word biases or context biases. This is also in line with Gardner et al. (2021)’s claim that all simple feature correlations based on partial input are spurious.

This study presents an analysis framework to quantify this context-word interaction by measuring context and target word biases. We first run controlled probing baselines by masking the input to show the context or the target word alone. Based on model’s performance on these probing baselines, we calculate two ratios that reflect how much of the model performance in this dataset can be achieved from simply relying on context alone or the target word alone, i.e. the degree of context or target word biases (See Figure 1 which will be discussed fully in Section 3). The design of the probing baselines follows previous studies that applied input permutation techniques for model and task analysis in GLUE (Pham et al., 2020), NLI (Poliak et al., 2018; Wang et al., 2018; Talman et al.,

2021) and relation extraction (Peng et al., 2020). While previous probing studies usually assume no meaningful information from corrupted input with no human verification, we provide fairer comparison with model performance by collecting human judgment on the same masked input in four tasks. Such comparison reveals whether the biases are learned spuriously by models from the datasets or are inherent in the tasks.

2 The Analysis Framework

2.1 Task Selection

We examine a number of popular context-aware lexical semantic tasks. For illustration, we list example data for each task in Table 4 in the appendix.

Word Sense Disambiguation (WSD). WSD (Navigli, 2009; Raganato et al., 2017) requires a model to assign a sense label to a target word in context from a set of possible candidates for the target word. Following the standard practice, we use SemCor as the train set, Semeval2007 as dev, and report

126 accuracy results on the concatenated ALL testset.

127 **The WiC-style Tasks (WiC, WiC-TSV, MCL-**
128 **WiC and XL-WiC).** To alleviate WSD’s require-
129 ment for a sense inventory, WiC (Pilehvar and
130 Camacho-Collados, 2019) presents a pairwise clas-
131 sification task where each pair consists of two word-
132 in-context instances. The model needs to judge
133 whether the target words in a pair have the same
134 contextual meanings. WiC-TSV (Breit et al., 2021)
135 extends the WiC framework to multiple domains
136 and settings. This study adopts the combined set-
137 ting where each input consists of a word in context
138 instance paired with a definition and a hypernym,
139 and the task is to judge whether the sense intended
140 by the target word in context matches the one de-
141 scribed by the definition and is the hyponym of the
142 hypernym. The WiC-style tasks have also been ex-
143 tended to the multilingual and crosslingual settings
144 in MCL-WiC (Martelli et al., 2021), XL-WiC (Ra-
145 ganato et al., 2020) and more recently in AM²ICO
146 (Liu et al., 2021). MCL-WiC provides test sets
147 for five languages with full gold annotation scores.
148 However, MCL-WiC only covers training data in
149 English. To ensure the analysis will be testing
150 the same data distribution during both training and
151 testing, we will only use the English dataset of
152 MCL-WiC. XL-WiC extends WiC to 12 languages.
153 While most languages in this task do not have train-
154 ing data, we perform analysis on its German dataset
155 which does contain both train (50k) and test data
156 (20k). AM²ICO covers 14 datasets, each of which
157 pairs English word-in-context instances with word-
158 in-context instances in a target language. In this
159 study, we perform analysis on the English-Chinese
160 dataset which contain 13k train and 1k test data ¹.

161 **Sense Retrieval (SR).** Based on WSD with the
162 same train and test data, SR (Loureiro and Jorge,
163 2019) requires a model to retrieve a correct entry
164 from the full sense inventory of all words from
165 WordNet (Miller, 1998).

166 **AIDA and Wikification.** An important applica-
167 tion scenario for testing meaning contextualization
168 is Entity Linking (EL). EL maps a mention (an en-
169 tity in its context) to a knowledge base (KB) which
170 is usually Wikipedia in the general domain. The
171 target word and its context help solve name vari-

172 ations and lexical ambiguity, which are the main
173 challenges in EL (Shen et al., 2014). In addition,
174 the context itself can help learn better representa-
175 tions for rare or new entities (Schick and Schütze,
176 2019; Ji et al., 2017). We test on two popular
177 Wikipedia-based EL benchmarks: AIDA (Hoffart
178 et al., 2011) and Wikification (Wiki) (Ratinov et al.,
179 2011; Bunescu and Paşca, 2006). AIDA provides
180 manual annotations of entities with Wikipedia and
181 YAGO2 labels for 946, 216 and 231 articles as train,
182 dev and test sets respectively. The Wiki Dataset
183 is based on the hyperlinks from Wikipedia. We
184 randomly sampled 50k sentences from Wikipedia
185 as the test and another 50k as the dev set. The rest
186 is used for training. For both AIDA and Wiki, the
187 search space is the full Wikipedia entity list.

188 **WikiMed and COMETA.** To test domain ef-
189 fects, we evaluate on two medical EL tasks. We
190 use the WikiMed corpus (Vashishth et al., 2020),
191 an automatically extracted medical subset from
192 Wikipedia, for medical wikification. Each men-
193 tion is mapped to a Wikipedia page linked to a
194 concept in UMLS (Bodenreider, 2004), a massive
195 medical concept KB. We define the search space as
196 the Wikipedia entities covered in UMLS. With the
197 same Wikipedia ontology but a different domain
198 subset, WikiMed can be directly compared with
199 Wiki for assessing domain influence. We also test
200 on COMETA (Basaldella et al., 2020), a medical
201 EL task in social media. COMETA consists of 20k
202 English biomedical entity mentions from online
203 posts in Reddit. The expert-annotated labels are
204 linked to SNOMED CT (Donnelly et al., 2006),
205 another widely-used medical KB.

206 We report accuracy for WSD and all the WiC
207 style tasks, and accuracy@1 for retrieval-based
208 tasks including Wiki, AIDA, etc.

2.2 Probing Baselines 209

210 **Context vs. Word:** For the main experiment, we
211 design the WORD baseline where we input only
212 the target word ² to the model, and the CONTEXT
213 baseline where the target word is replaced with a
214 [MASK] token in the input. The model is then
215 trained and tested on the perturbed input. A high
216 performance in CONTEXT or WORD will indicate
217 strong context or target word bias. Example base-
218 line input is shown in Table 1. **Lower Bound:**

¹We performed the analysis on other datasets of AM²ICO and found the trend is similar

²In the surveyed tasks, a target word can show different surface variations of number, case and etc. Eg., *breed*, *breeds*.

Apart from a RANDOM baseline, we also set up a LABEL baseline where all the input is masked and the learning is only from the label distribution in the task. Notice that training the LABEL baseline is preferable to simply counting label occurrences in the data as the former can work with both continuous and categorical label space. All the probing baselines are compared with model performance on the full input (FULL). We refer to model M’s performance in WORD, CONTEXT, LABEL and FULL as M_W , M_C , M_L and M_{Full} respectively.

Human Evaluation: To measure the inherent task biases, we collect human judgment (HUM) for a subset (WiC, XL-WiC, AM²ICO and AIDA) as being representative of the tasks described in Section 2.1 and feasible given resources for annotation. WiC, XL-WiC and AM²ICO cover WiC-style datasets in different languages; AIDA is chosen as a representative retrieval-based task. We follow the quality control procedures in Pilehvar and Camacho-Collados (2019); Liu et al. (2021) to recruit two different annotators for each baseline input from CONTEXT, WORD and for FULL input in each task. The annotators are recruited from Prolific. They have graduate degrees and are fluent or native in the language of the dataset. In each setup, an annotator is assigned a randomly sampled 100 examples from the test set of each task³ and there is a 50 example overlap between the two annotators for agreement calculation. The annotators are asked to perform meaning judgment in WiC, XL-WiC and AM²ICO, and to find the corresponding Wikipedia pages for entities for AIDA. For CONTEXT input where the target words are masked, annotators are encouraged to first guess what the target words could be. As to the WORD input, annotators are asked to think of the most representative meaning of the out-of-context words when performing the tasks. As the pairs of input are always the same word by design in WiC and XL-WiC, we assume humans will give true judgment for all the examples and therefore will score 0.5 on WORD input in WiC and XL-WiC. As to human’s LABEL baseline performance, while humans are not given any prior indication of how the task labels will be distributed, it is reasonable to expect that an annotator will give a random choice between the available labels or stick with one label

³We cannot use the test set for WiC and XL-WiC as the test labels are undisclosed. As the dev set comes from the same distribution of the test, we use dev to estimate human performance in these two tasks.

when there is no input. Therefore, we approximate the LABEL human baseline as being 0.5 for WiC, XL-WiC and AM²ICO, and 0 for AIDA.

2.3 Calculating the Bias Measures

Based on a model M ’s performance on the full input and on the baseline input, we propose $Bias^{MC}$ and $Bias^{MW}$ (as calculated in Equation (1) and Equation (2)) to measure the model’s context and target word biases in a dataset. $Bias^{MC}$ is the ratio of M_C to M_{Full} with the LABEL performance M_L deducted from both M_C and M_{Full} . M_L has to be deducted as it is unrelated to the input. Otherwise, the ratio will give an inflated bias measurement. $Bias^{MW}$ is calculated in the same way as $Bias^{MC}$ except that we replace M_C with M_W in the equation. The two measures can also be seen as M_C and M_W under min-max normalization where the min value is M_L and the max value is M_{Full} , and therefore the normalized values can be fairly compared across datasets. $Bias^{MC}$ and $Bias^{MW}$ reflect how much of what a model has learned from the input in a dataset can be achieved from context alone or target word alone, which will give us indicators of the degree of context and target word biases in the dataset. These bias indicators will in turn tell us how important the masked part of the input is. For example, we can interpret a $Bias^{MC}$ of 0.9 as 90% of what the model has learned from the full input can be achieved from the context alone. The 10% gap can be gained from adding the masked target word and since this gap is small with a high context bias, we can conclude that the model can do pretty well just from the context alone and it is not learning much from the target word.

$$Bias^{MC} = \frac{(M_C - M_L)}{(M_{Full} - M_L)} \quad (1)$$

$$Bias^{MW} = \frac{(M_W - M_L)}{(M_{Full} - M_L)} \quad (2)$$

Like models, humans can also be biased as they can also use their prior knowledge or biases (eg. humans can guess the typical meaning of a word without knowing the context) to make predictions based on partial input (Gardner et al., 2021). To measure how much humans can perform on the baseline input will help us understand the biases inherent in a task. We therefore calculate the context and target word bias scores for humans in the same way.

Input	Sentence1	Sentence2	BERT	HUMAN
FULL	Google represents a new [breed] of entrepreneurs .	The [breed] of tulip .	F	F
CONTEXT	Google represents a new [MASK] of entrepreneurs .	The [MASK] of tulip .	F	T
WORD	breed	breed	T	-
GUESSEDDWORD	Google represents a new [type] of entrepreneurs .	The [type] of tulip .	F	T

Table 1: Example input of FULL, CONTEXT and WORD in WiC. Target words are in brackets and the original WiC label for the FULL example is F. GUESSEDDWORD shows human-elicited target words based on CONTEXT. Comparing CONTEXT and GUESSEDDWORD also shows BERT’s contextual bias in WiC as BERT is not sensitive to the target word change.

2.4 Experiment setup

The underlying model for our main experiments is BERT (Devlin et al., 2019), one of the most successful PCMs that offer dynamic contextual word representations as bidirectional hidden layers from a transformer architecture. To ensure the general trend of our findings are consistent across different models, we also performed the analysis using ROBERTA (Liu et al., 2019), which improves upon BERT by optimized design decisions during training.

We adopt standard model finetuning setups in each task. We use the base uncased variant of BERT⁴ for general domain experiments and PUBMEDBERT (Gu et al., 2020) for the medical tasks. For WSD, we use GLOSSBERT (Huang et al., 2019) that learns a sentence-gloss pair classification model based on BERT. For the WiC-style tasks, we follow the SuperGlue (Wang et al., 2019) practices to concatenate BERT’s last layer of [CLS] and the target words’ token representations for each input pair, followed by a linear classifier. For the retrieval-based tasks including SR and EL, we adopt a bi-encoder architecture to model query and target candidates with BERT (Wu et al., 2020). For the query, we insert [and] to mark the start and end positions of the target word in context. Each target candidate is reformatted as “[CLS]Name || Description[SEP]”. Name is an entity title (EL) or synset lemmas from WordNet (SR). Description is the first sentence in an entity’s Wikipedia page (Wiki & WikiMed), a gloss (SR), or n/a (COMETA). The model learns to draw closer the true query-target pairs’ representations using triplet loss with triplet miners during finetuning (Liu et al., 2020). For each experiment, we perform grid search for the learning rate in $[1e-5, 2e-5, 3e-5]$ and select models with early

⁴All PCM configurations are listed in Appendix D. We also conducted experiments with ROBERTA (Liu et al., 2019) and reported the results in Appendix E

stopping on the dev set. We also run all the models with three random seeds and select the models with the best performance on the dev set. The performance across random seeds are stable as shown by small standard deviations which can be referred to in Table 5 in the appendix.

3 Main Results and Discussion

We report BERT’s baseline performance in Figure 2, based on which we calculate $Bias^{BERT_C}$ and $Bias^{BERT_W}$ for each dataset and plot the results (black dots) in Figure 1 (We also report ROBERTA biases in Appendix E and found a similar trend). For comparison, we plot human baseline performance and biases alongside the model performance in each figure.

3.1 Model biases

Models can learn extreme context or target word biases from the datasets. One obvious observation from Figure 1 is that, probed with BERT, most of the datasets lie close to the dashed red lines: tasks such as WiC and MCL-WiC lie towards the right and are close to the vertical red line which indicates 1.0 context bias; the retrieval-based tasks such as WikiMed and Wiki lie towards the top and are close to or even surpass the horizontal red line which indicates 1.0 target word bias. This pattern indicates that BERT can learn a tremendous amount from these datasets by relying only on the target words or only on the context. In other words, context or target words can be much ignored when the model learns to solve the tasks. It is therefore questionable how much word-context interaction, which requires the modeling of both word and context representations, is actually learned by BERT when applied to these tasks.

Moreover, the datasets tend to concentrate in two corners. That is, models usually learn strong bias from either context or the target word: the retrieval-based datasets (eg. Wiki) lie in the top

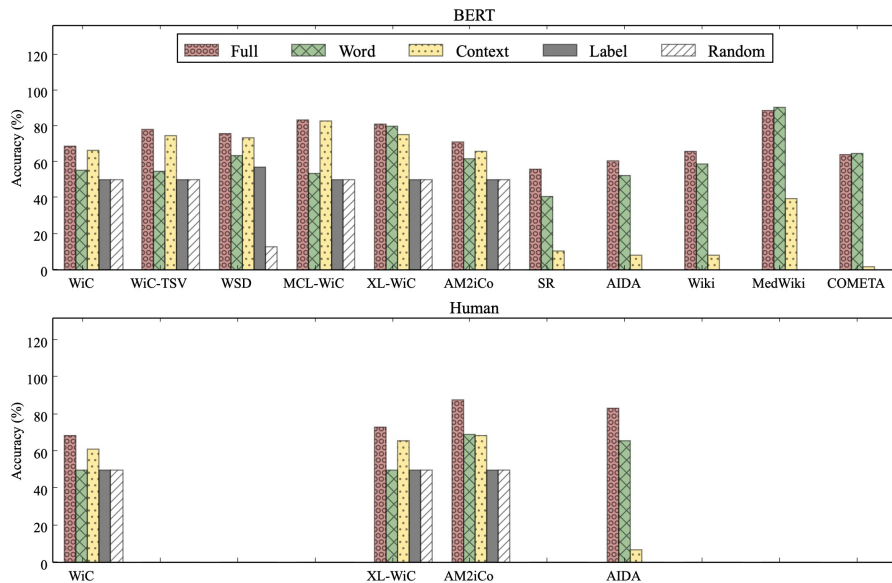


Figure 2: BERT and human performance on probing baselines across popular context-aware lexical semantic tasks. For the retrieval-based tasks, we report @1 accuracy, and the LABEL and RANDOM baselines are not visible as they are close to 0.

left corner, showing large target word bias and low context bias; the WiC style datasets and WSD lie in the bottom right corner with large context bias and low target word bias. XL-WiC is an exception as it contains both strong context and target word biases. We will come back to this later in Section 3.2 where we compare model and human performance.

AM²iCo and SR are closest to testing word-context interaction from models. There are few existing datasets that in effect require the modeling of the context-word interaction, which should result in both low context and target word biases. SR and AM²iCo can be seen as two such datasets which, in Figure 1, can be found further inside of the red lines towards the bias-free left bottom corner. This is because these two tasks are designed to require balanced attention over context and target words. In SR, a system needs to model the target words in order to retrieve all the possible senses associated with the word, and because there is plenty of ambiguity in the dataset, context is also crucial to identify the correct sense. AM²iCo was specifically designed to include adversarial examples to penalize models that rely only on the context, and therefore elicits the lowest context bias from models among the WiC-style tasks. As such, SR and AM²iCo are the closest tasks that we have to test word-context interaction.

Domains affect lexical ambiguity and the target word bias.

The retrieval-based tasks in this study offer comparison between two domains, general vs medical, by comparing Wiki/AIDA and WikiMed. The target word bias is increased in the medical domain where relying on the target words alone gives the best performance (i.e. COMETA and WikiMed both have > 1.0 target word bias). Such divergence across domains is arguably caused by the different degrees of lexical ambiguity in these tasks. In particular, domain could reduce ambiguity (Magnini et al., 2002; Koeling et al., 2005), and therefore affect the importance of the context and therefore the target word bias. As a quantitative measure for lexical ambiguity, we calculate average sense entropy across all words in each task’s training data, see Table 2. Confirming our hypothesis, sense entropy (lexical ambiguity) in a task does roughly correlate with the model’s target word bias: the medical domain tasks (WikiMed and COMETA) contain the lowest lexical ambiguity as reflected by the lowest sense entropy, and therefore missing context in these two tasks will not bring so much negative impact on the model performance, resulting in the highest target word biases; whereas higher sense entropy and thus higher lexical ambiguity (eg. Wiki and then SR) will necessarily require context along-

	SemCor	Wikification	AIDA	WikiMed	COMETA
Sense Entropy	0.2102	0.060	0.0438	0.026	0.0004
$Bias^{BERT_w}$	0.7274	0.8939	0.8705	1.0208	1.0124
$Bias^{RoBERTa_w}$	0.7315	0.8994	0.8319	0.9957	1.1798

Table 2: Target Word Bias and Sense Entropy across retrieval-based tasks

side the target word, which leads to lower target word biases.

Context can harm model performance in Medical EL. We notice that the model’s target word bias in COMETA and WikiMed can go beyond 1.0, indicating that the model learning is dominated entirely by the target words with the context being useless or even harmful. This comes as a surprise as medical EL has been treated as a contextual lexical semantic task where the context is usually provided in the hope for higher modeling accuracy. We examined the errors from FULL as compared with WORD, and we found that the model tends to get distracted by related context words. Table 3 shows an example where the retrieval model selects the entry that is closer to a context word (“Miltonia”) than to the target word (“Miltoniopsis”), but in fact knowing the target word alone in this case is sufficient to retrieve the correct label. This indicates that the model has not learned a good strategy to incorporate word and context representations from the datasets (i.e. not knowing when to focus on the context and when to focus on the target words).

3.2 Human vs Model

There are inherent task biases. Our first finding is that humans show a similar trend of biases in the tasks in comparison to model biases (except for XL-WiC). This is evident from Figure 1 where, with the human bias indicators, WiC still lies near bottom right corner with relatively high context bias; AIDA lies near top left corner with high target word bias and AM²ICO remains in the middle. This confirms that there are some degrees of biases inherent in the task design so that humans can also rely on either target words or context alone to perform the task to some extent.

Humans are less biased than models.

That being said, the second finding and the more important one is that humans exhibit overall much weaker biases in comparison with models in all the four tasks. If we compare human performance with model performance in Figure 2, we can see the

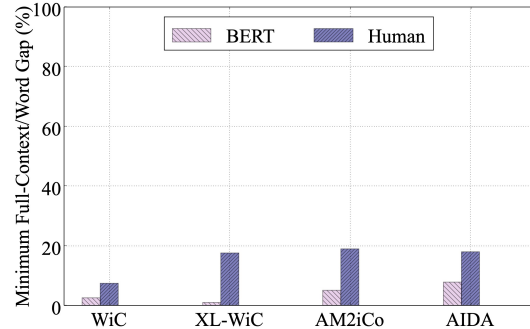


Figure 3: The minimum gap between FULL and CONTEXT or WORD, i.e. $\min(\text{FULL-CONTEXT}, \text{FULL-WORD})$ with BERT and human performance. A small gap will indicate strong bias.

CONTEXT and WORD baseline scores are lower in comparison to FULL from human performance. For clearer comparison, we calculate and plot the minimum gap between FULL to either of the two baselines in Figure 3, and we can see substantial difference between humans and models where humans exhibit much larger gaps across the four tasks. The much larger gaps from humans also result in all the four tasks moving further towards the left-bottom “bias-free” corner as shown Figure 1. In other words, humans are more likely than models to rely on both word and context as the absence of either part will lead to much more negative impact for humans when performing these tasks.

The most dramatic difference is in XL-WiC where the model’s strong target word bias disappears in humans. The task of XL-WiC by nature should not leak any information from the target word alone (hence 0 target word bias for humans) as the input pair will always contain the same target word. The high target word bias from models comes from the fact the dataset does not contain sufficient ambiguous cases where the same word pair can have both true and false labels dependent on the contexts. We confirm this by calculating the per-word average label entropy of the training data as 0.09, and on average a word pair has the same label for 94% of the context examples it appears

Baseline Input	Retrieved concept entry	Result
FULL	Formerly many more species were attributed to “Miltonia”, ... including [Miltoniopsis] and Oncidium ...	Wrong
WORD	Miltoniopsis	Correct

Table 3: Error analysis on FULL and WORD BERT predictions on WikiMed.

in the dataset. Therefore, the model learns correlation between the word itself and the label without needing context for disambiguation.

Target words are important in WiC for humans.

The much lower context bias from humans in tasks such as WiC suggests that the absence of the target words drastically decreases performance. In fact, human CONTEXT baseline (0.61) is even worse than BERT (0.65) as shown in Figure 2. This may also come as a surprise, considering that target words are always the same and only the context is different in each pair of input. We examined human response in CONTEXT and found that humans can guess another valid target word based on the context, which gives a different prediction. Table 1 shows such an example. While the original WiC label of the input is F, our annotator gave T for the CONTEXT input, guessing the target word is *type*. This is a reasonable prediction as *type* fits the contexts and does hold its meaning across the two sentences. We refer to this new example with human-elicited target words as GUESSEWORD input. The same annotator was able to give the WiC label F when we reveal the original target word (*breed*) which has the specific meaning of *species* in sentence1 and *personality* in sentence2 (see the FULL input in Table 1). BERT however still predicts F regardless of the target word change in this GUESSEWORD example.

As qualitative analysis on the human-model discrepancy on CONTEXT, we examined 20 cases where annotators did not predict WiC labels (from the corresponding FULL input) while BERT did. In 11 cases, humans guessed other valid target words to justify their predictions. We then perform preliminary analysis to test BERT on all the 11 GUESSEWORD cases where the human-elicited target words change the labels (We show more examples in Table 6), and found that for 7 out of 11, BERT is insensitive to the changed target words and maintains its original prediction. This suggests

BERT does not appreciate the same word-context interaction as humans, and is making prediction mainly based on contexts rather than modeling contextual lexical semantics in WiC.

4 Conclusion

This study presented an analysis framework to disentangle and quantify context-word interplay in application of popular contextual lexical semantic benchmarks. With our proposed bias measures, we plot datasets on a continuum, and we found that, to models, most existing datasets lie on the two ends with excessive biases (WiC-style tasks and WSD are heavily context-biased while retrieval-based tasks are heavily target-word-biased) that essentially bypass the key challenges in word-context interaction. SR and AM²ICO have been identified as two tasks that have less extreme biases and therefore can better test the representation of both word and context, and we call for more tasks that challenge models to do so. In addition, we identify that the degree of lexical ambiguity as a byproduct of domain affects target word bias (medical>general) in retrieval-based tasks. Most importantly, we differentiate biases spuriously learned by models and task-inherent biases by collecting human responses on the same baseline input. We found that models’ heavy context and target word biases are not attested to the same extent in humans who usually need both context and target words to perform well in the tasks. This suggests that models are learning spurious correlations instead of modeling contextual lexical semantics as intended by the tasks. Our paper highlights the importance of understanding these biases in existing datasets and encourages future dataset and model design to control for these biases and to focus more on testing the challenging word-context interaction in context-sensitive lexical semantics. Possible future directions will be to include adversarial examples that penalize sole reliance on context or target words in both task design and model training.

References

- 598
599 Marco Basaldella, Fangyu Liu, Ehsan Shareghi, and
600 Nigel Collier. 2020. [COMETA: A corpus for medical entity linking in the social media](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3122–3137, Online. Association for Computational Linguistics. 605
- 606 Terra Blevins and Luke Zettlemoyer. 2020. [Moving down the long tail of word sense disambiguation with gloss informed bi-encoders](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1006–1017, Online. Association for Computational Linguistics. 611
- 612 Olivier Bodenreider. 2004. The unified medical language system (umls): integrating biomedical terminology. *Nucleic acids research*, 32(suppl_1):D267–D270. 615
- 616 Anna Breit, Artem Revenko, Kiamehr Rezaee, Mohammad Taher Pilehvar, and Jose Camacho-Collados. 2021. [WiC-TSV: An evaluation benchmark for target sense verification of words in context](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1635–1645, Online. Association for Computational Linguistics. 623
- 624 Samuel Broscheit. 2019. [Investigating entity knowledge in BERT with simple neural end-to-end entity linking](#). In *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)*, pages 677–685, Hong Kong, China. Association for Computational Linguistics. 629
- 630 Razvan Bunescu and Marius Paşca. 2006. [Using encyclopedic knowledge for named entity disambiguation](#). In *11th Conference of the European Chapter of the Association for Computational Linguistics*, Trento, Italy. Association for Computational Linguistics. 635
- 636 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. 644
- 645 Kevin Donnelly et al. 2006. Snomed-ct: The advanced terminology and coding system for ehealth. *Studies in health technology and informatics*, 121:279. 647
- 648 Matt Gardner, William Merrill, Jesse Dodge, Matthew Peters, Alexis Ross, Sameer Singh, and Noah A. Smith. 2021. [Competency problems: On finding and removing artifacts in language data](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1801–1813, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics. 654
655
- Aina Garí Soler, Marianna Apidianaki, and Alexandre Allauzen. 2019. [Word usage similarity estimation with sentence representations and automatic substitutes](#). In *Proceedings of the Eighth Joint Conference on Lexical and Computational Semantics (*SEM 2019)*, pages 9–21, Minneapolis, Minnesota. Association for Computational Linguistics. 656
657
658
659
660
661
662
- Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2020. [Domain-specific language model pretraining for biomedical natural language processing](#). 663
664
665
666
667
- Christian Hadiwinoto, Hwee Tou Ng, and Wee Chung Gan. 2019. [Improved word sense disambiguation using pre-trained contextualized word representations](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5297–5306, Hong Kong, China. Association for Computational Linguistics. 668
669
670
671
672
673
674
675
676
- Johannes Hoffart, Mohamed Amir Yosef, Ilaria Bordino, Hagen Fürstenau, Manfred Pinkal, Marc Spaniol, Bilyana Taneva, Stefan Thater, and Gerhard Weikum. 2011. [Robust disambiguation of named entities in text](#). In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*, pages 782–792, Edinburgh, Scotland, UK. Association for Computational Linguistics. 677
678
679
680
681
682
683
684
- Luyao Huang, Chi Sun, Xipeng Qiu, and Xuanjing Huang. 2019. [GlossBERT: BERT for word sense disambiguation with gloss knowledge](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3509–3514, Hong Kong, China. Association for Computational Linguistics. 685
686
687
688
689
690
691
692
693
- Yangfeng Ji, Chenhao Tan, Sebastian Martschat, Yejin Choi, and Noah A. Smith. 2017. [Dynamic entity representations in neural language models](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1830–1839, Copenhagen, Denmark. Association for Computational Linguistics. 694
695
696
697
698
699
700
- Rob Koeling, Diana McCarthy, and John Carroll. 2005. [Domain-specific sense distributions and predominant sense acquisition](#). In *Proceedings of Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing*, pages 419–426, Vancouver, British Columbia, Canada. Association for Computational Linguistics. 701
702
703
704
705
706
707
- Fangyu Liu, Ehsan Shareghi, Zaiqiao Meng, Marco Basaldella, and Nigel Collier. 2020. [Self-alignment pre-training for biomedical entity representations](#). *arXiv preprint arXiv:2010.11784*. 708
709
710
711

712	Qianchu Liu, Edoardo M. Ponti, Diana McCarthy, Ivan Vulić, and Anna Korhonen. 2021. Am2ico: Evaluating word meaning in context across low-resourcelanguages with adversarial examples.	Alessandro Raganato, José Camacho-Collados, and Roberto Navigli. 2017. Word sense disambiguation: A unified evaluation framework and empirical comparison. In <i>Proceedings of EACL 2017</i> , pages 99–110.	767 768 769 770 771
716	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. <i>arXiv preprint arXiv:1907.11692</i> .	Alessandro Raganato, Tommaso Pasini, Jose Camacho-Collados, and Mohammad Taher Pilehvar. 2020. XL-WiC: A multilingual benchmark for evaluating semantic contextualization. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 7193–7206, Online. Association for Computational Linguistics.	772 773 774 775 776 777 778
721	Daniel Loureiro and Alípio Jorge. 2019. Language modelling makes sense: Propagating representations through WordNet for full-coverage word sense disambiguation. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 5682–5691, Florence, Italy. Association for Computational Linguistics.	Lev Ratinov, Dan Roth, Doug Downey, and Mike Anderson. 2011. Local and global algorithms for disambiguation to Wikipedia. In <i>Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1375–1384, Portland, Oregon, USA. Association for Computational Linguistics.	779 780 781 782 783 784 785
728	Bernardo Magnini, Carlo Strapparava, Giovanni Pezzulo, and Alfio Gliozzo. 2002. The role of domain information in word sense disambiguation. <i>Natural Language Engineering</i> , 8(4):359–373.	Timo Schick and Hinrich Schütze. 2019. Attentive mimicking: Better word embeddings by attending to informative contexts. In <i>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)</i> , pages 489–494, Minneapolis, Minnesota. Association for Computational Linguistics.	786 787 788 789 790 791 792 793
732	Federico Martelli, Najla Kalach, Gabriele Tola, and Roberto Navigli. 2021. SemEval-2021 Task 2: Multilingual and Cross-lingual Word-in-Context Disambiguation (MCL-WiC). In <i>Proceedings of the Fifteenth Workshop on Semantic Evaluation (SemEval-2021)</i> .	Wei Shen, Jianyong Wang, and Jiawei Han. 2014. Entity linking with a knowledge base: Issues, techniques, and solutions. <i>IEEE Transactions on Knowledge and Data Engineering</i> , 27(2):443–460.	794 795 796 797
738	George A Miller. 1998. <i>WordNet: An electronic lexical database</i> . MIT press.	Aarne Talman, Marianna Apidianaki, Stergios Chatzikyriakidis, and Jörg Tiedemann. 2021. Nli data sanity check: Assessing the effect of data corruption on model performance. <i>arXiv preprint arXiv:2104.04751</i> .	798 799 800 801 802
740	Roberto Navigli. 2009. Word sense disambiguation: A survey. <i>ACM Computing Surveys</i> , 41(2):1–69.	Shikhar Vashishth, Denis Newman-Griffis, Rishabh Joshi, Ritam Dutt, and Carolyn Rose. 2020. Improving broad-coverage medical entity linking with semantic type prediction and large-scale datasets. <i>arXiv preprint arXiv:2005.00460</i> .	803 804 805 806 807
742	Hao Peng, Tianyu Gao, Xu Han, Yankai Lin, Peng Li, Zhiyuan Liu, Maosong Sun, and Jie Zhou. 2020. Learning from Context or Names? An Empirical Study on Neural Relation Extraction. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 3661–3672, Online. Association for Computational Linguistics.	Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2019. Superglue: A stickier benchmark for general-purpose language understanding systems. <i>Advances in Neural Information Processing Systems</i> , 32.	808 809 810 811 812 813
748	Thang M Pham, Trung Bui, Long Mai, and Anh Nguyen. 2020. Out of order: How important is the sequential order of words in a sentence in natural language understanding tasks? <i>arXiv preprint arXiv:2012.15180</i> .	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In <i>Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP</i> , pages 353–355, Brussels, Belgium. Association for Computational Linguistics.	814 815 816 817 818 819 820 821
755	Mohammad Taher Pilehvar and Jose Camacho-Collados. 2019. WiC: the word-in-context dataset for evaluating context-sensitive meaning representations. In <i>Proceedings of NAACL-HLT 2019</i> , pages 1267–1273.		
760	Adam Poliak, Jason Naradowsky, Aparajita Haldar, Rachel Rudinger, and Benjamin Van Durme. 2018. Hypothesis only baselines in natural language inference. In <i>Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics</i> , pages 180–191, New Orleans, Louisiana. Association for Computational Linguistics.		

822
823
824
825
826
827
828

Ledell Wu, Fabio Petroni, Martin Josifoski, Sebastian Riedel, and Luke Zettlemoyer. 2020. [Scalable zero-shot entity linking with dense entity retrieval](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6397–6407, Online. Association for Computational Linguistics.

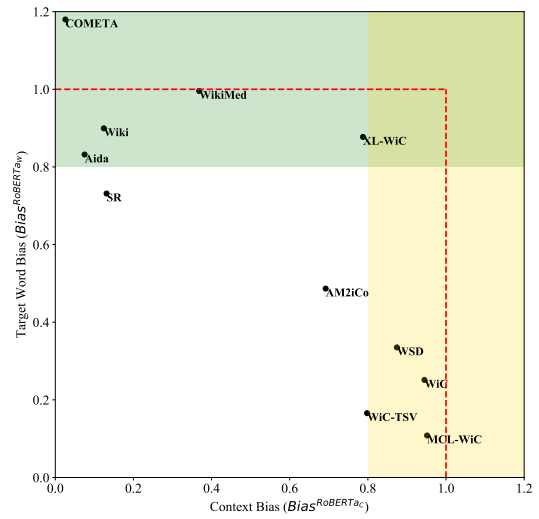


Figure 4: Plotting context and target word biases when applying ROBERTA across popular context-aware lexical semantic datasets. The green shade and the yellow shade roughly indicate the areas for high target word bias and high context bias (0.8). The dashed red lines indicate 1.0 context (right) and 1.0 target word bias (top), implying the model only requires the target words alone or context alone in this dataset.

A Task examples

829

Table 4 lists example input and labels for tasks surveyed in this study.

830

831

B Dev performance

832

Table 5 shows BERT biases calculated over three runs on the dev set with standard deviation reported.

833

834

835

C Examples of the context bias in WiC

836

See Table 6 for two examples where the model relies solely on the context to make the prediction.

837

838

D Model configurations

839

ALL PCMs are from <https://huggingface.co/>. Model configurations are listed in Table 7.

840

841

842

E ROBERTA Performance (Figure 4)

843

Task	Input	Label	Label Space	Metrics
WiC	Room and [board]. He nailed [boards] across the windows.	F	T or F	Acc
WiC-TSV	I spent my [spring] holidays in Morocco. the season of growth; season, time of the year	T	T or F	Acc
MCL-WiC	Bolivia holds a key [play] in any process... A musical [play] on the same subject...	F	T or F	Acc
XL-WiC	Herr [Starke] wollte uns kein Interview geben. Das kann ich dir aber sagen: Wenn die Frau [Starke] kommt...	T	T or F	Acc
AM ² ICo	...航天员训练及[阿波罗]中飞船... ...the six [Apollo] Moon landings...	T	T or F	Acc
WSD	The [art] of change-ringing is peculiar to the English...	art : a superior skill that you can learn by study and practice and observation	art : the creation of beautiful or significant things art : the products of human creativity; works of art collectively ...(all possible meanings of <i>art</i>)	F1
SR	The [art] of change-ringing is peculiar to the English...	art : a superior skill that you can learn by study and practice and observation	art : a superior skill that you can learn by study and practice and observation door : a swinging or sliding barrier that will close the entrance... ... PLUS all other entries in WordNet	Acc
Wiki	an additional [Hash] literal syntax using colons for symbol keys...	hash table : in computing , a hash table (hash map) is a data structure...	hash table : in computing , a hash table (hash map) is a data structure ... united kingdom : the United Kingdom of Great Britain and Northern Ireland... ... (all entries in Wikipedia)	Acc@1
WikiMed	The flowers produce pollen, but no nectar. Various bees and flies visit the flowers looking in vain for nectar, for instance [sweat bees] in the genera “Lasiglossum” and “Halictus”...	halictidae : the Halictidae is the second largest family of Apoidea bees.	halictidae : the Halictidae is the second largest family of Apoidea bees. eomecon : eomecon is a monotypic genus of flowering plants in the poppy family... ... (all entries in the medical section of Wikipedia)	Acc@1
COMETA	I am [spacey] because I am thinking and daydreaming about my obsession.	dizziness (finding)	dizziness (finding) large intestine ...PLUS all other entries in SNOMED CT	Acc@1

Table 4: Examples for a selection of context-sensitive lexical semantic tasks surveyed in this thesis. Acc: accuracy; ρ : Spearman’s correlation; r : Pearson’s correlation; P&R: precision and recall.

	WiC	WiC-TSV	WSD	MCL-WiC	XL-WiC	AM ² iCo	SR	AIDA	Wiki	MedWiki	COMETA
$Bias^{BERT_W}$	0.473 (0.016)	0.266 (0.043)	0.346 (0.015)	0.122 (0.007)	0.903 (0.002)	0.665 (0.008)	0.648 (0.012)	0.910 (0.007)	0.946 (0.002)	1.024 (0.022)	1.017 (0.034)
$Bias^{BERT_C}$	1.055 (0.017)	0.890 (0.028)	0.874 (0.020)	0.864 (0.043)	0.844 (0.002)	0.768 (0.016)	0.237 (0.011)	0.241 (0.015)	0.308 (0.003)	0.447 (0.010)	0.028 (0.010)

Table 5: Average context and target word biases over three runs with three different random seeds on the dev set in each dataset. Standard deviation is reported in the parenthesis.

Input	Sentence1	Sentence2	BERT	HUM
FULL	[Misdirect] the letter .	The pedestrian [misdirected] the out - of - town driver .	F	F
CONTEXT	[MASK] the letter .	The pedestrian [MASK] the out - of - town driver .	F	T
GUESSEWORD	[Ignore] the letter .	The pedestrian [ignored] the out - of - town driver .	F	T
FULL	[Kill] the engine .	He [kills] the ball .	F	F
CONTEXT	[MASK] the engine	He [MASK] the ball .	F	T
GUESSEWORD	[Hit] the engine .	He [hits] the ball .	F	T

Table 6: Example input of WORD, CONTEXT and FULL in WiC. The original WiC label for these examples is F. GUESSEWORD contains human-elicited target words that flip the label. Comparing CONTEXT and GUESSEWORD also shows BERT’s contextual bias in WiC as BERT is not sensitive to the target word change.

Model	Variant name in Huggingface	Parameters	Pretraining corpus
BERT	bert-base-uncased	12-layer, 768-hidden, 12-heads, 110M parameters	Lowercased Wikipedia + BookCorpus
PUBMEDBERT	microsoft/BiomedNLP-PubMedBERT-base-uncased-abstract-fulltext	12-layer, 768-hidden, 12-heads, 110M parameters	Lowercased abstracts from PubMed and full-text articles from PubMedCentral
DEBERTA	microsoft/deberta-large	24-layer, 1024-hidden, 16-heads, 400M parameters	Wikipedia + BookCorpus + OPENWEBTEXT (public Reddit content) + STORIES

Table 7: Model details in our experiments