

Measuring Context-Word Biases in Lexical Semantic Datasets

Anonymous ACL submission

Abstract

State-of-the-art pretrained contextualized models (PCM) eg. BERT use tasks such as WiC and WSD to evaluate their *word-in-context* representations. This inherently assumes that performance in these tasks reflect how well a model represents the coupled word and context semantics. We question this assumption by presenting the first quantitative analysis on the context-word interaction being tested in major contextual lexical semantic tasks. To achieve this, we run probing baselines on masked input, and propose measures to calculate and visualize the degree of context or word biases in existing datasets. The analysis was performed on both models and humans. Our findings demonstrate that models are usually not being tested for word-in-context semantics in the same way as humans are in these tasks, which helps us better understand the model-human gap. Specifically, to PCMs, most existing datasets fall into the extreme ends (the retrieval-based tasks exhibit strong target word bias while WiC-style tasks and WSD show strong context bias); In comparison, humans are less biased and achieve much better performance when both word and context are available than with masked input. We recommend our framework for understanding and controlling these biases for model interpretation and future task design.

1 Introduction

Meaning contextualization (i.e., identifying the correct meaning of a target word in linguistic context) is essential for understanding natural language, and has been the focus in many lexical semantic tasks. Pretrained contextualized models (PCMs) have brought large improvements in these tasks including WSD (Hadiwinoto et al., 2019; Loureiro and Jorge, 2019; Huang et al., 2019; Blevins and Zettlemoyer, 2020), WiC (Pilehvar and Camacho-Collados, 2019; Garí Soler et al., 2019) and entity

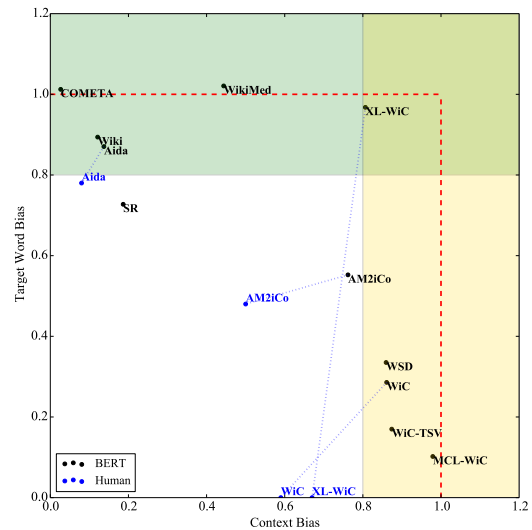


Figure 1: Plotting context and target word biases from BERT (black) and humans (blue) across popular context-aware lexical semantic datasets. The green shade and the yellow shade roughly indicate the areas for high target word bias and high context bias (>0.8). We would ideally want a dataset to lie 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.

linking (EL) (Wu et al., 2020; Broscheit, 2019).

These superior performances have been taken as proof that PCMs can successfully model *word-in-context* semantics. Many studies have investigated the process of lexical contextualization in these PCMs. Specifically, Vulić et al. (2020); Aina et al. (2019) found language models ‘contextualize’ words in higher layers while the type-level information is better kept in lower layers¹. Voita et al. (2019) point out different learning objectives affect the contextualization process, and Garí Soler

¹In this study, we do not perform layer-wise analysis as we fine-tune the PCMs to achieve the optimal performance in each task and we expect the relevant information is already surfaced to the last layer (Peters et al., 2019)

054 and Apidianaki (2021); Pimentel et al. (2020) show 101
055 PCMs can capture words’ ambiguity levels. 102

056 While most these studies have focused on prob- 103
057 ing the innerworkings of the PCM feature space, 104
058 there is no systematic study to quantify the word- 105
059 context interaction (either learned by PCMs or in- 106
060 trinsic) across different lexical semantic tasks. On 107
061 one hand, these datasets often vary in their em- 108
062 phasis on context vs target words. For example, 109
063 we could expect tasks such as WSD and WiC to 110
064 rely more on context by design as the target words 111
065 are either given or the same in each input pair 112
066 ². On the other hand, models may find shortcuts 113
067 from datasets to avoid learning the complex word- 114
068 context interaction. **What is missing in the cur- 115
069 rent literature is an accurate quantification of 116
070 this word-context interplay being tested in each 117
071 task so that we can fully understand task goals 118
072 and model performance.** In particular, we need 119
073 to flag the situation where a model can solve a task 120
074 by relying solely on context or the target words. 121
075 Such heavy word or context reliance hinders a sci- 122
076 entific assessment of the models’ meaning contex- 123
077 tualization abilities as it essentially bypasses the 124
078 key word-context interaction challenge in human 125
079 understanding of lexical semantics. Therefore, we 126
080 refer to such heavy reliance on target words or 127
081 context in a contextual lexical semantic dataset as 128
082 target word biases or context biases ³.

083 This study presents an analysis framework to 129
084 quantify this context-word interaction by measur- 130
085 ing context and target word biases across lexical 131
086 semantic tasks. We first run controlled probing 132
087 baselines by masking the input to show the context 133
088 or the target word alone. Based on model’s perfor- 134
089 mance on these probing baselines, we calculate two 135
090 ratios that reflect how much of the model perfor- 136
091 mance in this dataset can be achieved from simply 137
092 relying on context alone or the target word alone, 138
093 i.e. the degree of context or target word biases (See 139
094 Figure 1 which will be discussed fully in Section 3). 140
095 The design of the probing baselines follows pre- 141
096 vious studies that applied input perturbation tech- 142
097 niques for model and task analysis in GLUE (Pham 143
098 et al., 2020), NLI (Poliak et al., 2018; Wang et al., 144
099 2018; Talman et al., 2021) and relation extraction 145
100 (Peng et al., 2020). While previous probing studies 146

²Notice that the exact amount of context/target word re-
liance in these tasks is to be tested as humans naturally use
both to make prediction.

³This is also in line with Gardner et al. (2021)’s claim that
all simple feature correlations from partial input are spurious.

usually assume no meaningful information from 101
corrupted input with no human verification, we pro- 102
vide fairer comparison with model performance by 103
collecting human judgment on the same masked in- 104
put in four tasks. Such comparison reveals whether 105
the biases are learned by models from the datasets 106
or are inherent in the tasks. 107

Our key findings are (1) the tasks can be clearly 108
divided into target-word-biased (the retrieval-based 109
tasks), and context-biases (WiC-style tasks and 110
WSD). Among the retrieval-based tasks, domain 111
affects ambiguity level and thus the target word 112
bias: models even achieve the best performance us- 113
ing target words alone in the medical domain. (2) 114
AM²ICO and Sense Retrieval show less extreme 115
model biases and challenge a model more to rep- 116
resent both the context and target words; and (3) 117
a similar trend of biases exists in humans but is 118
much less extreme, as humans find semantic judg- 119
ment more difficult on masked input and require 120
both word and context to do well in each task. This 121
analysis helps us better understand the nuanced 122
differences between models and humans in exist- 123
ing tasks, and we recommend the framework to 124
be applied when designing new datasets to check 125
whether word and context are required and whether 126
the models rely on the coupled word and context 127
semantics in a similar way to humans. 128

2 The Analysis Framework 129

2.1 Task Selection 130

We examine the following contextual lexical se- 131
mantic tasks, and for illustration, we list example 132
data for each task in Appendix A in the appendix. 133

Word Sense Disambiguation (WSD). WSD (Nav- 134
igli, 2009; Raganato et al., 2017) requires a model 135
to assign sense labels to target words in context 136
from a set of possible candidates for the target 137
words. Following the standard practice, we use 138
SemCor (sense-annotated texts created from Brown 139
Corpus) as the train set, Semeval2007 as dev, and 140
report accuracy on the concatenated ALL testset. 141

**The WiC-style Tasks (WiC, WiC-TSV, MCL- 142
WiC and XL-WiC).** To alleviate WSD’s require- 143
ment for a sense inventory, WiC (Pilehvar and 144
Camacho-Collados, 2019) presents a pairwise clas- 145
sification task where each pair consists of two word- 146
in-context instances. The model needs to judge 147
whether the target words in a pair have the same 148
contextual meanings. WiC-TSV (Breit et al., 2021) 149
extends the WiC framework to multiple domains 150

and settings. This study adopts the combined setting where each input consists of a word in context instance paired with a definition and a hypernym, and the task is to judge whether the sense intended by the target word in context matches the one described by the definition and is the hyponym of the hypernym. The WiC-style tasks have also been extended to the multilingual and crosslingual settings in MCL-WiC (Martelli et al., 2021), XL-WiC (Raganato et al., 2020) and more recently in AM²ICO (Liu et al., 2021). MCL-WiC provides test sets for five languages with full gold annotation scores. However, MCL-WiC only covers training data in English. To ensure the analysis will be testing the same data distribution during both training and testing, we will only use the English dataset of MCL-WiC. XL-WiC extends WiC to 12 languages. While most languages in this task do not have training data, we perform analysis on its German dataset which does contain both train (50k) and test data (20k). AM²ICO covers 14 datasets, each of which pairs English word-in-context instances with word-in-context instances in a target language. In this study, we perform analysis on the English-Chinese dataset which contain 13k train and 1k test data ⁴.

Sense Retrieval (SR). With the same train and test data as WSD, SR (Loureiro and Jorge, 2019) requires a model to retrieve a correct entry from the full sense inventory of WordNet (Miller, 1998).

AIDA and Wikification. An important application scenario for testing meaning contextualization is Entity Linking (EL). EL maps a mention (an entity in its context) to a knowledge base (KB) which is usually Wikipedia in the general domain. The target word and its context help solve name variations and lexical ambiguity, which are the main challenges in EL (Shen et al., 2014). In addition, the context itself can help learn better representations for rare or new entities (Schick and Schütze, 2019; Ji et al., 2017). We test on two popular Wikipedia-based EL benchmarks: AIDA (Hoffart et al., 2011) and Wikification (Wiki) (Ratinov et al., 2011; Bunescu and Paşca, 2006). AIDA provides manual annotations of entities with Wikipedia and YAGO2 labels for 946, 216 and 231 articles as train, dev and test sets respectively. The Wiki Dataset is based on the hyperlinks from Wikipedia. We randomly sampled 50k sentences from Wikipedia as the test and another 50k as the dev set. The rest

⁴We performed the analysis on other datasets of AM²ICO and found the trend is similar

is used for training. For both AIDA and Wiki, the search space is the full Wikipedia entity list.

WikiMed and COMETA. To test domain effects, we evaluate on two medical EL tasks. We use the WikiMed corpus (Vashishth et al., 2020), an automatically extracted medical subset from Wikipedia, for medical wikification. Each mention is mapped to a Wikipedia page linked to a concept in UMLS (Bodenreider, 2004), a massive medical concept KB. We define the search space as the Wikipedia entities covered in UMLS. With the same Wikipedia ontology but a different domain subset, WikiMed can be directly compared with Wiki for assessing domain influence. We also test on COMETA (Basaldella et al., 2020), a medical EL task in social media. COMETA consists of 20k English biomedical entity mentions from online posts in Reddit. The expert-annotated labels are linked to SNOMED CT (Donnelly et al., 2006), another widely-used medical KB.

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

2.2 Probing Baselines

Context vs. Word: For the main experiment, we design the WORD baseline where we input only the target word ⁵ to the model, and the CONTEXT baseline where the target word is replaced with a [MASK] token in the input. The model is then trained and tested on the perturbed input. A high performance in CONTEXT or WORD will indicate strong context or target word bias. Example baseline input is shown in Table 1. **Lower Bound:** 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

⁵In the surveyed tasks, a target word can show different surface variations of number, case and etc. Eg., *breed*, *breeds*.

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 dataset⁶. In each setup, an annotator is assigned a randomly sampled 100 examples from the test set of each task⁷ 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 be⁸. 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 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 $Bias^{M_C}$ and $Bias^{M_W}$ (as calculated in Equation (1) and Equation (2)) to measure the model’s context and

target word biases in a dataset. $Bias^{M_C}$ 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. $Bias^{M_W}$ is calculated in the same way as $Bias^{M_C}$ 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. $Bias^{M_C}$ and $Bias^{M_W}$ 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 $Bias^{M_C}$ 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.

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

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

Like models, humans can also be biased as they can use their prior knowledge (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.

2.4 Experiment setup

The underlying model for our main experiments is BERT (Devlin et al., 2019), the most popular PCM that offers 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.

⁶Notice the annotator profiles in this study may be different from the original annotation scheme which was not always clearly specified. Therefore results on FULL in this study may be different from figures originally reported.

⁷We 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.

⁸We provide an example annotation guideline in appendix G. Human has lower but still reasonable agreement in the probing baselines (where there is naturally less information) than with FULL input (appendix F).

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	-
GUESSEWORD	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. GUESSEWORD shows human-elicited target words based on CONTEXT. Comparing CONTEXT and GUESSEWORD also shows BERT’s contextual bias in WiC as BERT is not sensitive to the target word change.

We adopt standard model finetuning setups in each task. We use the base uncased variant of BERT⁹ for general domain experiments and PUB-MEDBERT (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 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 7 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

⁹All PCM configurations are listed in Appendix D. We also conducted experiments with ROBERTA (Liu et al., 2019) and reported the results in Appendix E

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.

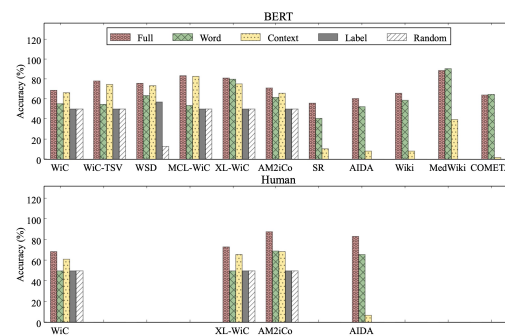


Figure 2: BERT and human performance on probing baselines across popular context-aware lexical semantic tasks. For the retrieval-based tasks, we report @1 accuracy, and the LABEL and RANDOM baselines are not visible as they are close to 0.

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 score highly on 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 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.

	SemCor	Wiki	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

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 alongside 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 the bottom right corner with relatively high context bias; AIDA lies near the top left corner with high

	Input	Retrieved concept entry	Result
FULL	Formerly many more species were attributed to "Miltonia", ... including [Miltoniopsis] and Oncidium ...	miltonia: miltonia is an orchid genus comprising twelve epiphyte species and eight natural hybrids.	Wrong
WORD	Miltoniopsis	miltoniopsis: miltoniopsis is a genus of orchids native to costa rica and etc.	Correct

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

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.

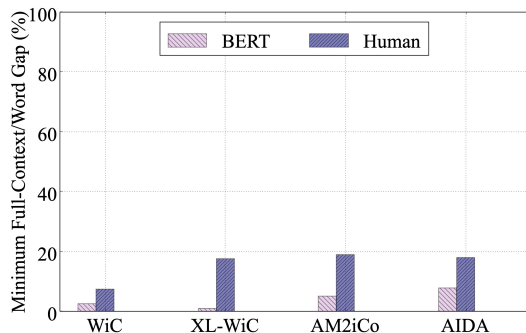


Figure 3: The minimum gap between FULL and CONTEXT or WORD, i.e. $\min(\text{FULL-CONTEXT}, \text{FULL-WORD})$ with BERT and human performance. A small gap will indicate strong bias.

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 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 for all WiC-style tasks¹⁰ in Table 4, and we found XL-WiC has the lowest label entropy as 0.2084, and on average a word pair has the same label for 87% of the examples it appears in the dataset. Therefore, the model learns correlation between the word itself and the label without needing context for disambiguation.

	WiC	AM ² iCo	XL-WiC	MCL-WiC	WiC-TSV
LE	0.4022	0.5383	0.2084	0.2976	0.4212

Table 4: Label Entropy (LE) across WiC style tasks

Finally, the fact that models cannot achieve a similarly large jump of performance from masked input to FULL like humans could indicate the word-context interaction is particularly challenging for models and this might eventually explain the model-human gap. We take AM²iCo as an example that explicitly requires word-context interaction (Figure 5). While BERT achieves comparable results with humans in CONTEXT and WORD, a significantly larger human-model gap is found in FULL, indicating the word-context interaction is what the model lacks the most to achieve human-like performance.

	CONTEXT	WORD	FULL
Human	69	68.5	87.9
BERT	66	61	71
Human - BERT	3	7.5	16.9

Table 5: Model-human gap in CONTEXT, WORD and FULL in AM²iCo

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,

¹⁰We disregard words that only occur once

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 GUESSEDDWORD 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 GUESSEDDWORD 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 GUESSEDDWORD cases where the human-elicited target words change the labels (We show more examples in Table 8), 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 Implications for future dataset design

We recommend this analysis framework in future dataset design and result interpretation for contextual word representation evaluation. In particular, we recommend (1) creating probing baselines by masking the context and word (if relevant), and (2) providing a sample to humans (details in section 2.2, (3) and comparing human and model performance of full input vs the masked baseline/s, and then calculate bias indicators. In terms of task design, we would ideally want both models and humans to show low baseline performance and

thus low bias measures. When interpreting the results, apart from evaluating model performance on the FULL input, we should also ensure the model shows a human-like gap in performance (between FULL and the baseline(s)) on the same data.

5 Conclusion and Limitations

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 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 relying on different cues 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.

One limitation of this study is that we do not have large-scaled quantitative evidence to pinpoint the cues the models rely on from partial input¹¹. Possible future directions will be to design such ablation studies to identify any spurious correlations the models have learned and introduce adversarial examples that penalize sole reliance on context or target words in both task design and model training.

¹¹It could be that models like BERT (trained with masked language modelling) are genuinely better than humans in exploiting partial input

References

Laura Aina, Kristina Gulordava, and Gemma Boleda. 2019. [Putting words in context: LSTM language models and lexical ambiguity](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3342–3348, Florence, Italy. Association for Computational Linguistics.

Marco Basaldella, Fangyu Liu, Ehsan Shareghi, and 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.

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.

Olivier Bodenreider. 2004. The unified medical language system (umls): integrating biomedical terminology. *Nucleic acids research*, 32(suppl_1):D267–D270.

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.

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.

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.

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.

Kevin Donnelly et al. 2006. Snomed-ct: The advanced terminology and coding system for ehealth. *Studies in health technology and informatics*, 121:279.

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.

Aina Garí Soler and Marianna Apidianaki. 2021. [Let’s play mono-poly: BERT can reveal words’ polysemy level and partitionability into senses](#). *Transactions of the Association for Computational Linguistics*, 9:825–844.

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.

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](#).

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.

Johannes Hoffart, Mohamed Amir Yosef, Ilaria Bordino, Hagen Fürstenu, 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.

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.

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.

761	Rob Koeling, Diana McCarthy, and John Carroll. 2005.	Thang M Pham, Trung Bui, Long Mai, and Anh	817
762	Domain-specific sense distributions and predomi-	Nguyen. 2020. Out of order: How important is	818
763	nant sense acquisition. In <i>Proceedings of Human</i>	the sequential order of words in a sentence in nat-	819
764	<i>Language Technology Conference and Conference</i>	ural language understanding tasks? <i>arXiv preprint</i>	820
765	<i>on Empirical Methods in Natural Language Process-</i>	<i>arXiv:2012.15180.</i>	821
766	<i>ing</i> , pages 419–426, Vancouver, British Columbia,		
767	Canada. Association for Computational Linguistics.		
768	Fangyu Liu, Ehsan Shareghi, Zaiqiao Meng, Marco	Mohammad Taher Pilehvar and Jose Camacho-	822
769	Basaldella, and Nigel Collier. 2020. Self-alignment	Collados. 2019. WiC: the word-in-context dataset	823
770	pre-training for biomedical entity representations.	for evaluating context-sensitive meaning representa-	824
771	<i>arXiv preprint arXiv:2010.11784.</i>	tions. In <i>Proceedings of NAACL-HLT 2019</i> , pages	825
772		1267–1273.	826
773	Qianchu Liu, Edoardo M. Ponti, Diana McCarthy,	Tiago Pimentel, Rowan Hall Maudslay, Damian Blasi,	827
774	Ivan Vulić, and Anna Korhonen. 2021. Am2ico:	and Ryan Cotterell. 2020. Speakers fill lexical se-	828
775	Evaluating word meaning in context across low-	mantic gaps with context. In <i>Proceedings of the</i>	829
776	resourcelanguages with adversarial examples.	<i>2020 Conference on Empirical Methods in Natural</i>	830
777		<i>Language Processing (EMNLP)</i> , pages 4004–4015,	831
778	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	Online. Association for Computational Linguistics.	832
779	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,		
780	Luke Zettlemoyer, and Veselin Stoyanov. 2019.	Adam Poliak, Jason Naradowsky, Aparajita Haldar,	833
781	Roberta: A robustly optimized bert pretraining ap-	Rachel Rudinger, and Benjamin Van Durme. 2018.	834
782	proach. <i>arXiv preprint arXiv:1907.11692.</i>	Hypothesis only baselines in natural language in-	835
783		ference. In <i>Proceedings of the Seventh Joint Con-</i>	836
784	Daniel Loureiro and Alípio Jorge. 2019. Language	<i>ference on Lexical and Computational Semantics</i> ,	837
785	modelling makes sense: Propagating representa-	pages 180–191, New Orleans, Louisiana. Associa-	838
786	tions through WordNet for full-coverage word sense	tion for Computational Linguistics.	839
787	disambiguation. In <i>Proceedings of the 57th Annual</i>		
788	<i>Meeting of the Association for Computational Lin-</i>	Alessandro Raganato, José Camacho-Collados, and	840
789	<i>guistics</i> , pages 5682–5691, Florence, Italy. Associa-	Roberto Navigli. 2017. Word sense disambiguation:	841
790	tion for Computational Linguistics.	A unified evaluation framework and empirical com-	842
791		parison. In <i>Proceedings of EACL 2017</i> , pages 99–	843
792	Bernardo Magnini, Carlo Strapparava, Giovanni Pez-	110.	844
793	zulo, and Alfio Gliozzo. 2002. The role of domain	Alessandro Raganato, Tommaso Pasini, Jose Camacho-	845
794	information in word sense disambiguation. <i>Natural</i>	Collados, and Mohammad Taher Pilehvar. 2020.	846
795	<i>Language Engineering</i> , 8(4):359–373.	XL-WiC: A multilingual benchmark for evaluating	847
796		semantic contextualization. In <i>Proceedings of the</i>	848
797	Federico Martelli, Najla Kalach, Gabriele Tola, and	<i>2020 Conference on Empirical Methods in Natural</i>	849
798	Roberto Navigli. 2021. SemEval-2021 Task 2: Mul-	<i>Language Processing (EMNLP)</i> , pages 7193–7206,	850
799	tilingual and Cross-lingual Word-in-Context Disam-	Online. Association for Computational Linguistics.	851
800	biguation (MCL-WiC). In <i>Proceedings of the Fif-</i>		
801	<i>teenth Workshop on Semantic Evaluation (SemEval-</i>	Lev Ratinov, Dan Roth, Doug Downey, and Mike An-	852
802	<i>2021)</i> .	derison. 2011. Local and global algorithms for dis-	853
803	George A Miller. 1998. <i>WordNet: An electronic lexical</i>	ambiguation to Wikipedia. In <i>Proceedings of the</i>	854
804	<i>database.</i> MIT press.	<i>49th Annual Meeting of the Association for Com-</i>	855
805		<i>putational Linguistics: Human Language Technol-</i>	856
806	Roberto Navigli. 2009. Word sense disambiguation: A	<i>gies</i> , pages 1375–1384, Portland, Oregon, USA. As-	857
807	survey. <i>ACM Computing Surveys</i> , 41(2):1–69.	sociation for Computational Linguistics.	858
808			
809	Hao Peng, Tianyu Gao, Xu Han, Yankai Lin, Peng	Timo Schick and Hinrich Schütze. 2019. Attentive	859
810	Li, Zhiyuan Liu, Maosong Sun, and Jie Zhou. 2020.	mimicking: Better word embeddings by attending	860
811	Learning from Context or Names? An Empirical	to informative contexts. In <i>Proceedings of the 2019</i>	861
812	Study on Neural Relation Extraction. In <i>Proceed-</i>	<i>Conference of the North American Chapter of the</i>	862
813	<i>ings of the 2020 Conference on Empirical Methods</i>	<i>Association for Computational Linguistics: Human</i>	863
814	<i>in Natural Language Processing (EMNLP)</i> , pages	<i>Language Technologies, Volume 1 (Long and Short</i>	864
815	3661–3672, Online. Association for Computational	<i>Papers)</i> , pages 489–494, Minneapolis, Minnesota.	865
816	Linguistics.	Association for Computational Linguistics.	866
817			
818	Matthew E. Peters, Sebastian Ruder, and Noah A.	Wei Shen, Jianyong Wang, and Jiawei Han. 2014. En-	867
819	Smith. 2019. To tune or not to tune? adapting pre-	entity linking with a knowledge base: Issues, tech-	868
820	trained representations to diverse tasks. In <i>Proceed-</i>	niques, and solutions. <i>IEEE Transactions on Knowl-</i>	869
821	<i>ings of the 4th Workshop on Representation Learn-</i>	<i>edge and Data Engineering</i> , 27(2):443–460.	870
822	<i>ing for NLP (RepL4NLP-2019)</i> , pages 7–14, Flo-		
823	rence, Italy. Association for Computational Linguis-	Aarne Talman, Marianna Apidianaki, Stergios	871
824	tics.	Chatzikyriakidis, and Jörg Tiedemann. 2021. Nli	872

873 data sanity check: Assessing the effect of data
 874 corruption on model performance. *arXiv preprint*
 875 *arXiv:2104.04751*.

876 Shikhar Vashishth, Denis Newman-Griffis, Rishabh
 877 Joshi, Ritam Dutt, and Carolyn Rose. 2020. Im-
 878 proving broad-coverage medical entity linking with
 879 semantic type prediction and large-scale datasets.
 880 *arXiv preprint arXiv:2005.00460*.

881 Elena Voita, Rico Sennrich, and Ivan Titov. 2019. The
 882 bottom-up evolution of representations in the trans-
 883 former: A study with machine translation and lan-
 884 guage modeling objectives. In *Proceedings of the*
 885 *2019 Conference on Empirical Methods in Natural*
 886 *Language Processing and the 9th International*
 887 *Joint Conference on Natural Language Processing*
 888 *(EMNLP-IJCNLP)*, pages 4396–4406, Hong Kong,
 889 China. Association for Computational Linguistics.

890 Ivan Vulić, Edoardo Maria Ponti, Robert Litschko,
 891 Goran Glavaš, and Anna Korhonen. 2020. **Probing**
 892 **pretrained language models for lexical semantics**. In
 893 *Proceedings of the 2020 Conference on Empirical*
 894 *Methods in Natural Language Processing (EMNLP)*,
 895 pages 7222–7240, Online. Association for Computa-
 896 tional Linguistics.

897 Alex Wang, Yada Pruksachatkun, Nikita Nangia,
 898 Amanpreet Singh, Julian Michael, Felix Hill, Omer
 899 Levy, and Samuel R Bowman. 2019. Superglue: A
 900 stickier benchmark for general-purpose language un-
 901 derstanding systems. *Advances in Neural Informa-*
 902 *tion Processing Systems*, 32.

903 Alex Wang, Amanpreet Singh, Julian Michael, Fe-
 904elix Hill, Omer Levy, and Samuel Bowman. 2018.
 905 **GLUE: A multi-task benchmark and analysis plat-**
 906 **form for natural language understanding**. In *Pro-*
 907 *ceedings of the 2018 EMNLP Workshop Black-*
 908 *boxNLP: Analyzing and Interpreting Neural Net-*
 909 *works for NLP*, pages 353–355, Brussels, Belgium.
 910 Association for Computational Linguistics.

911 Ledell Wu, Fabio Petroni, Martin Josifoski, Sebastian
 912 Riedel, and Luke Zettlemoyer. 2020. **Scalable zero-**
 913 **shot entity linking with dense entity retrieval**. In
 914 *Proceedings of the 2020 Conference on Empirical*
 915 *Methods in Natural Language Processing (EMNLP)*,
 916 pages 6397–6407, Online. Association for Computa-
 917 tional Linguistics.

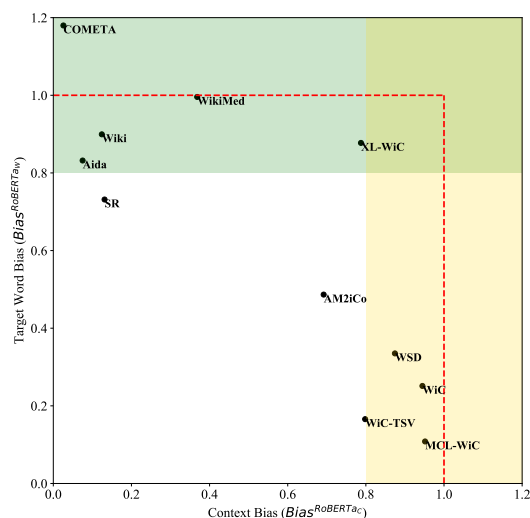


Figure 4: Plotting context and target word biases when applying ROBERTA across popular context-aware lexical semantic datasets. The green shade and the yellow shade roughly indicate the areas for high target word bias and high context bias (0.8). The dashed red lines indicate 1.0 context (right) and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

A Task examples 918

Table 6 lists example input and labels for tasks surveyed in this study. 919 920

B Dev performance 921

Table 7 shows BERT biases calculated over three runs on the dev set with standard deviation reported. 922 923 924

C Examples of the context bias in WiC 925

See Table 8 for two examples where the model relies solely on the context to make the prediction. 926 927

D Model configurations 928

ALL PCMs are from <https://huggingface.co/>. 929
 Model configurations are listed in Table 9. 930

931

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 6: 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 7: 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
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
FULL	His [treatment] of the race question is badly biased .	His [treatment] of space borrows from Italian architecture .	F	F
CONTEXT	His [MASK] of the race question is badly biased .	His [MASK] of space borrows from Italian architecture .	F	T
GUESSEWORD	His [understanding] of the race question is badly biased .	His [understanding] of space borrows from Italian architecture .	T	F
FULL	I could see it in the [distance] .	The [distance] from New York to Chicago .	F	F
CONTEXT	I could see it in the [MASK] .	The [MASK] from New York to Chicago .	F	T
GUESSEWORD	I could see it in the [train] .	The [train] from New York to Chicago .	T	F

Table 8: 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.

E ROBERTA Performance (Figure 4)

F Agreement in WiC-style tasks (Table 10)

G Annotation Guideline

Figure 5 shows an example annotation guideline for the CONTEXT experiment in WiC.

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 9: Model details in our experiments

The task asks you to decide whether a target word has the same contextual meaning or not in two different contexts.				
In this particular task, the target words are masked as [X] and you should rely on the contexts to make a prediction as best as you can.				
Open the 'Data' spreadsheet. Each row is a pair of contexts for a word [X].				
In the 'label' column, write 't' for pairs where you think the meanings of [X] are likely to be the same, and 'f' otherwise. If you really cannot make a decision based on the context, just write 't'.				
Below are two examples in English:				
Example no	context1	context2	label	word
1	[X] clothes .	She [X] her blouses .	t	iron
2	Do you [X] sugar in your coffee ?	A reading was [X] of the earth 's tremors .	f	take
The basic form for [X] was 'starch' in example 1 and 'take' in example 2. Although you cannot see the target words, you can still make a guess as to whether the two contextual meanings are the same.				
The actual target word that you think of might be different to the one in the original sentence, that doesn't matter. If you can, use 'word' column to indicate the word you think of in its basic form.				
Thank you!				

Figure 5: An annotation guideline for conducting the CONTEXT baseline of humans in WiC.

	AM ² iCo	XL-WiC	WiC
CONTEXT	94.0	88	76
FULL	87.9	66	64

Table 10: Human agreement in CONTEXT and FULL in WiC-style tasks