

PiC: A Phrase-in-Context Dataset for Phrase Understanding and Semantic Search

Thang M. Pham ^{†*}
thangpham@auburn.edu

Seunghyun Yoon [§]
syoon@adobe.com

Trung Bui [§]
bui@adobe.com

Anh Nguyen [†]
anh.ng8@gmail.com

[†]Auburn University [§]Adobe Research

Abstract

Since BERT [9], learning contextualized word embeddings has been a de-facto standard in NLP. However, the progress of learning contextualized *phrase* embeddings is hindered by the lack of a human-annotated benchmark that tests model understanding of phrase semantics given a context sentence or a paragraph (as opposed to phrases alone). To fill this gap, we propose PiC—a dataset of ~28K of noun phrases accompanied by their contextual Wikipedia pages and a suite of three tasks for training and evaluating phrase embeddings. We find that training on our dataset improves ranking models’ accuracy and remarkably pushes question-answering-style (QA) models (i.e., predicting the start and end index of the target phrase) near human-accuracy, which is 95% Exact Match (EM) on semantic search given a query phrase and a passage. Interestingly, we find evidence that such impressive performance is because the QA models learn to better capture the common meaning of a phrase *regardless* of its actual context. SotA models perform poorly—in differentiating between two different senses of the same phrase under two different context paragraphs (60% EM) and estimating the similarity between two different phrases given the same context (70% EM). Learning contextualized phrase embeddings remains an interesting, open challenge.

1 Introduction

Understanding phrases in context is a key to learning new vocabularies [22, 11], disambiguation [24], and many downstream tasks, including semantic search [10]. While contextualized *word* embeddings in BERT [9] have been a major success, the contextualized *phrase* embeddings [41] in existing systems mostly capture the common meaning of a phrase, i.e. without strong dependence on its context [41]. While there are *word*-sense disambiguation datasets, e.g. WiC [24], no such benchmarks exist for *phrases*. Existing phrase similarity benchmarks [23, 33, 2, 42, 39] compare phrases alone (without context) and some of them [23, 42] contain a large amount (~38% to 99%) of phrase pairs that share words, i.e. lexical overlap (see Table 1).

Others automatically generated the context for a phrase using GPT-2 [34] or by retrieving from Wikipedia [41]. Yet, there was no human verification of the realism of generated text [34] and no human annotation of how a phrase’s meaning changes w.r.t. the context [41]. All above drawbacks are limiting the evaluation of progress in phrase understanding.

*TP started his work during an internship at Adobe Research, and then continued at Auburn University.

To advance the development of contextualized phrase embeddings, we propose Phrase-in-Context (PiC), a suite of three tasks: (1) Phrase Similarity (PS), i.e. compare the semantic similarity of two phrases in the same context sentence (Fig. 1b); (2) Phrase Retrieval (PR), which is divided into PR-pass and PR-page (Fig. 1c), i.e. from a passage or a Wikipedia page, retrieve a phrase semantically-similar to a given query phrase; and (3) Phrase Sense Disambiguation (PSD), i.e. find the target phrase p semantically similar to the query phrase from a 2-paragraph document where p appears twice, each time in a different context paragraph that provides a *unique* meaning to p (Fig. 1e). Our $\sim 28\text{K}$ -example dataset is rigorously (a) *annotated* and *verified* by two groups of annotators: linguistics experts on [Upwork.com](https://www.upwork.com) and non-experts on Amazon Mechanical Turk (MTurk); and then (b) *tested* by models, linguists, and graduate students. Our contributions are:

1. We build PiC², the first, human-annotated benchmark suite of three tasks (PS, PR, and PSD) for evaluating and training contextualized phrase embeddings (Sec. 4). Compared to phrase similarity datasets, PS is the first to require models to rely on context. PSD is the first contextualized phrase disambiguation task (as opposed to word sense disambiguation). And PR is the first phrase-retrieval task where the query is a phrase (as opposed to questions).
2. After training on PR-pass, i.e. finding a phrase from a passage, Question Answering (QA) models perform at a near-human accuracy (92–94% vs. 95% EM). They also score high (84–89% EM) on PR-page, i.e. semantic phrase search in a single Wikipedia page (Sec. 5.4). This suggests our training set and resultant improved embeddings are useful for real-world semantic search and perhaps other downstream tasks (e.g. topic modeling as shown in [34]).
3. Interestingly, on PR-pass, harnessing these QA models’ phrase embeddings in a ranking approach (i.e. comparing the similarity between the query and *all* candidate phrases) yields much worse accuracy, i.e. 59% EM, at best (Sec. 5.5), which suggests that learning robust phrase embeddings remains an open question.
4. After training on PR-pass, state-of-the-art (SotA) models perform relatively well on PR-pass and even PR-page but, interestingly, not PSD (Sec. 5.6). Similarly, on PS, SotA models perform poorly (below 70% accuracy) in labeling phrase similarity given their context sentence (Sec. 5.1). That is, learning representations of phrases in context for poses a grand challenge for the community.

To our best knowledge, we present the first, human-annotated phrase-in-context benchmark and corresponding extensive evaluation of SotA models (of both ranking and QA approaches).

2 Related Work

We propose three different tasks (PS—phrase similarity; PR—phrase retrieval; and PSD—disambiguation), each related to a separate research area that we discuss below.

Phrase similarity First, most existing phrase similarity datasets (e.g. PPDB-annotated [35], PPDB-filtered [34], BiRD [2], and PAWS-short [34, 42]) contain a large percent of instances with lexical overlap between two paired phrases while our PS contains the least percent (6.09%; Table 1). Second, PS compares each pair of phrases in a context sentence while existing datasets only compare phrases alone (no context). Third, the phrases in PS are, on average, 2-token long (Table 1), comparable to that of other datasets. Fourth, unlike other datasets, PS contains exclusively *noun-phrases*, the most common phrase type according to the search-query statistics from Yahoo [36] (79.54%; Appendix D) and Adobe (internal Acrobat Pro data not shown), and each phrase contains ≥ 2 words.

Question answering (QA) Our phrase retrieval tasks—PR and PSD—follow the format of QA datasets except that our queries are *phrases* instead of questions and hence shorter (Table 2). Like SQuAD 1.1 [25] and HotpotQA [40], our documents and queries are extracted from Wikipedia articles. While our PR dataset is $\sim 3.5\times$ smaller than those two datasets, the paragraph document length in PR-pass and PSD is $\sim 2\times$ longer than those of SQuAD 1.1 and HotpotQA (Table 2). For our task, intuitively, the longer the document, the harder the task since there would be more candidates a model must compare with the query.

²Our dataset, evaluation and dataset-construction code are at <https://phrase-in-context.github.io>

Table 1: Our Phrase Similarity (PS) dataset is the largest in terms of the number of instances (i.e. phrase pairs) and the number of unique phrases. PS has the least percent of lexical-overlap instances. In addition, PS is the *only* human-annotated dataset that provides phrases in their context sentences.

	PS (ours)	WiC [24]	PPDB- annotated [35]	PPDB- filtered [34]	BiRD [2]	Turney [33]	PAWS- short [34]
All instances	10,208	7,466	3,000	15,532	3,345	2,180	1,214
Instances w/ lexical overlap	6.09%	100%	70.10%	97.93%	14.98%	0%	99.42%
Unique phrases	7,371	2,345	6,000	12,023	2,840	9,776	1,214
Mean length (in tokens)							
phrase ₁	2.06	1	3.67	2	2	2	9.52
phrase ₂	2.47	1	3.73	2	1.49	1	9.42
context sentence	22.54	8.40	0	0	0	0	0

Table 2: Our PR-pass, PR-page and PSD datasets are smaller in size compared to common QA datasets and contain shorter queries that are noun phrases instead of questions. However, our tasks require searching in much longer documents.

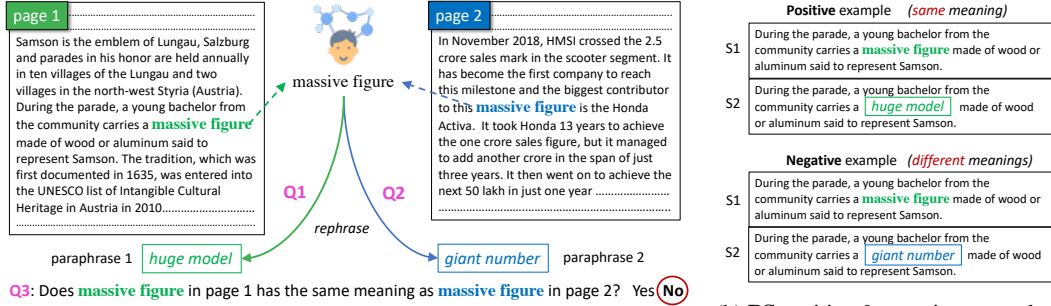
	PR-pass	PR-page	PSD	SQuAD 1.1 [25]	HotpotQA [40]
All instances	28,147	28,098	5,150	98,169	105,257
Unique queries/questions	27,055	27,016	5,100	97,888	105,249
Unique answers	13,458	13,423	2,416	72,469	57,259
Mean length of					
query/question (tokens)	2.42	2.42	2.46	11.42	20.03
answer (tokens)	2.17	2.17	2.06	3.46	2.35
sentence (tokens)	23.22	24.08	22.98	27.62	26.77
document (sentences)	10.26	119.32	20.41	5.10	4.14
document (tokens)	238.34	2,872.73	468.85	140.92	110.72

Sense disambiguation While word-sense disambiguation (WSD) is a long-standing problem in NLP, recently, SotA models have reached super-human accuracy (80% F_1) on the common English WSD [4]. Interestingly, these high-scoring models still struggle with rare senses that may be outside of the predefined sense inventories or have few training examples [6]. Without the need for predefined senses, WiC [24] poses disambiguation as a binary classification task where the goal is to predict whether the same target word in two different sentences carries the same or different meanings.

Compared to WiC PS is also a binary classification task, but with two major differences: (1) in WiC, the same target word appears in two different sentences while in PS, two different phrases appear in the same context sentence; (2) PS compares phrases composed of ≥ 2 words instead of a single word as in WiC and WSD. While word senses are defined in WordNet and BabelNet dictionaries [4], there are no English dictionaries of senses for multi-word noun phrases (*mNPs*). Thus, it is more challenging to acquire and learn the senses of *mNPs*, hence the importance of our PiC dataset. Like WiC, PSD tests disambiguating the meanings of the same *n*-gram in two different contexts. Yet, PSD is a phrase search task, which involves many more phrase comparisons per example than PS or WiC.

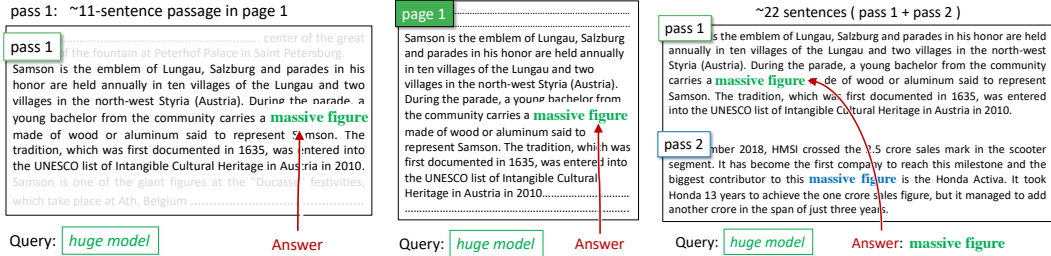
3 PiC Dataset Construction

We first collect a set of phrases with context and human annotations. Then, we derive the examples and labels for three main tasks: PS, PR, and PSD (Fig. 1). Our idea is to mine a set of triplets (p , page₁, page₂) from Wikipedia where the phrase p is a polysemous *mNP* that carries two different senses in two Wikipedia pages (e.g., “massive figure” means *a large number* in page₁ but *a huge physical shape* in page₂; Fig. 1a). Then, we ask linguistic experts to rephrase p into two paraphrases q_1 and q_2 , maintaining the two original senses of p in page₁ and page₂, respectively. The resultant set of 5-tuples (p , q_1 , q_2 , page₁, page₂) enables us to derive PS, PR-pass, PR-page, and PSD—the tasks that test (1) comparing the semantic similarity of two phrases given the same context sentence (Fig. 1b); (2) finding a semantically similar phrase in a document (PR-pass & PR-page; Fig. 1c); (3) disambiguating the senses of the same target *mNP* given two context paragraphs (PSD; Fig. 1e).



(a) Q1 & Q2 ask annotators to rephrase “massive figure” in page 1 and page 2. Q3 asks whether this phrase’s meaning is the same in both pages.

(b) PS positive & negative examples constructed using page 1 context (similarly, we repeat for page 2).



(c) A PR-pass example.

(d) A PR-page example.

(e) A PSD example.

Figure 1: Given a phrase, two associated Wikipedia pages, and expert annotations, i.e. answers to Q1, Q2, and Q3 (a), we are able to construct *two* pairs of positive and negative examples for PS (b), a PR-pass example (c), a PR-page example (d), and a PSD example *only if* the answer to Q3 is No (e).

3.1 Data Collection

As there are no English dictionaries that contain sense inventories for *mNPs*, the key **challenge** to our data collection is to mine *mNPs* that have (1) multiple senses; and (2) a Wikipedia context page for each sense. To do that, we take a Wikipedia dump and perform a **6-step** procedure that essentially extract all the *mNPs* that occur in more than one Wikipedia pages and that contain at least one polysemous word defined in the WiC dataset. From the triplets of (p , $page_1$, $page_2$), we programmatically narrow down to ~600K triplets where the context sentence of the *mNP* in $page_1$ is the most semantically dissimilar to the context sentence in $page_2$ (using SimCSE [12]). We continue filtering down to the top 19,500 triplets where $page_1$ and $page_2$ have *the most* semantically different lists of Wiki categories. That is, 19,500 triplets are estimated to yield ~15K annotated triplets (the target size based on our budget) after the human annotation process (as we allow annotators to ignore the cases they are not confident). See Appendix C for a detailed description of the data collection.

3.2 Data Annotation

Via Upwork, we hire 13 linguistics experts who are native English speakers at a rate of \$30/hour to annotate 15,021 out of 19,500 examples. For each phrase, we provide Upworkers with a triplet (p , $passage_1$, $passage_2$) where each $passage_i$ consists of 5 sentences centered at the phrase-containing sentence in the corresponding $page_i$. We ask them to answer three questions (Fig. 1a):

- Q1 Rephrase the target phrase p to a paraphrase q_1 such that its meaning is constant in $passage_1$.
- Q2 Similarly, rephrase p w.r.t. $passage_2$ to obtain a paraphrase q_2 .
- Q3 Answer Y/N if p has the same meaning in both contextual $passage_1$ and $passage_2$.

Upworkers are asked to provide paraphrases that (1) have at least two words and (2) minimize lexical overlap with each other and the target p . See the annotation guidelines given to Upworkers in [27] and a sample annotation assignment in [28]. After receiving annotations, we use LanguageTool [20]

to automatically find syntactical errors when the paraphrases are replaced by the original target phrase in the original passage and ask Upworkers to fix them. We also ask annotators to fix the remaining errors that we find by manual inspection.

3.3 Verifying Annotations

To verify the annotations obtained in Sec. 3.2 (i.e. $2 \times 15,021 = 30,042$ paraphrases; and 15,021 Y/N answers), first, we present the same Q1, Q2, and Q3 questions to 1,000 qualified MTurkers and ask whether they agree with the answers by Upwork annotators in Sec. 3.2. And then, for the cases that the MTurkers disagree with, we seek second opinions from 5 Upwork experts. After these two verification rounds, we discard all the examples where Upwork verifiers reject and arrive at the final 28,325 paraphrases and 13,413 Y/N labels (i.e. those annotations that *either* a MTurk or Upwork verifier endorses). See more details in Appendix H.

The total amount of annotation fee (including MTurk and Upwork) is around USD 25K.

4 Three Phrase Understanding Tasks

Using the human-annotated data, we construct three tasks of PS, PR, and PSD (as summarized in Fig. 1) for evaluating contextualized phrase-embeddings and semantic-search models.

4.1 Phrase Similarity (PS)

PS is a binary classification task with the goal of predicting whether two m NPs are semantically similar or not given *the same context* sentence. Given the annotated data, a positive example is a triplet of (an original phrase p , a paraphrase q_1 , an original page₁’s sentence that contains p). To create a negative example, from the same triplets, we select only those where the paraphrase q_2 holds a *different* meaning than q_1 given the page₁ context of q_1 (i.e., when the answer to Q3 is No; see Fig. 1b). In total, we obtain 5,104 negative examples. Then, we randomly select 5,104 positive examples to form a class-balanced PS dataset.

4.2 Phrase Retrieval (PR)

That is, PR is a task of finding in a given document d a phrase p that is semantically similar to the given query phrase, which is the paraphrase q_1 (the answer by annotators to Q1) or q_2 (the answer to Q2). We release two versions of PR: **PR-pass** and **PR-page**, i.e. datasets of triplets (query q_1 , target phrase p , document d) where d is a random 11-sentence passage that contains p (Fig. 1c) or an entire Wikipedia page (Fig. 1d). While PR-pass contains 28,147 examples, PR-page contains slightly fewer examples (28,098) as we remove those examples whose Wikipedia pages coincidentally also contain exactly the query phrase (in addition to the target phrase). Both datasets are split into $\sim 20K/3K/5K$ for train/dev/test, respectively.

4.3 Phrase Sense Disambiguation (PSD)

From the verified annotations in Sec. 3.3, there are in total 5,150 phrases p such that both annotators and verifiers agree to hold *different* meanings across the two context Wikipedia pages (i.e., No answer to Q3 in Fig. 1a). To create a PSD example, given a phrase p from the above 5,150, we extract from its associated page₁ and page₂ two corresponding ~ 11 -sentence paragraphs (as in PR) and concatenate them (separated by an empty line) into a single document (see Fig. 1e). The task is to find the occurrence of the target phrase p that is semantically similar to a paraphrase q (e.g., the first “massive figure” in pass₁ but not the second one in Fig. 1e example) in this new document.

5 Experiments and Results

Here, we test SotA models on PS, PR-pass, PR-page, and PSD to (1) assess whether existing models are able to leverage context to improve accuracy; and (2) compute the best model accuracy, quantifying the headroom for future research.

Non-contextualized and Contextualized phrase embeddings Besides testing SotA trained BERT-based classifiers, we also test the ranking approach based on SotA contextualized embeddings, i.e. comparing the cosine similarity between the query and candidate phrases. Following [41], we compute a *contextualized* phrase embedding by feeding the entire phrase-containing sentence (e.g. S_1 in Fig. 1b) into a model, e.g. BERT, and taking the mean pooling of the last-layer embeddings over the words of the given phrase only. For non-contextualized phrase embeddings, we repeat the same process but input to the model only the phrase (without the surrounding words in the sentence).

Models We choose SotA models in phrase similarity (PhraseBERT [34]), sentence similarity (USE-v5 [8], SentenceBERT [26], and SimCSE [12]), question-answering (Longformer [3], DensePhrase [21]), and contextualized embeddings (SpanBERT [18], BERT [9]). For DensePhrase, we use their Phrase-Encoder (as opposed to the Query-Encoder) to compute phrase embeddings. USE-v5 is only available via public APIs [32], which do not support extraction of contextualized embeddings.

5.1 Phrase Similarity: Incorporating context into phrase embeddings improves accuracy

RQ: *Does incorporating context improve the phrase similarity accuracy on PS?*

Experiment We split the PS dataset 70/10/20 for train/dev/test and test two approaches: (1) using the cosine similarity score between two pre-trained phrase-embeddings (with and without context) to predict phrase similarity; (2) training BERT-based binary classifiers directly using PS training set. We use 6 backbone BERT models that are all “base” version unless specified otherwise (Table 3).

Approach 1: Cosine similarity First, we test how *pre-trained* phrase embeddings alone (without finetuning or extra weights) can be leveraged to solve PS. For each PS example of two phrases, we compute their non-contextualized phrase embeddings and compute their cosine similarity score. To evaluate the pre-trained embeddings on PS, we follow [38] and tune the binary-classification threshold T to maximize the training-set accuracy, and then use the same optimal T to report the test-set accuracy. We repeat the experiment for *contextualized* phrase embeddings.

Approach 2: BERT-based classifiers To complement Approach 1, we test Approach 2, i.e. building a binary classifier by adding two extra MLP layers on top of the pre-trained embeddings used in Approach 1. For a phrase pair, we concatenate the two 768-D phrase embeddings from BERT_{base} into a 1,536-D vector, and then place one ReLU layer (256 units) and a 2-output linear classification layer with softmax on top. Following [34], we finetune these models for a maximum of 100 epochs (with early stopping and patience of 10 epochs) on the train set. See Appendix A for more training details.

Results First, without the context, all models perform lower than 50% (i.e. the random chance; Table 3a & c), suggesting that PS requires models to rely on context. Second, incorporating context information into phrase embeddings substantially improves mean model accuracy on PS for both Approach 1 (from 48.92% to 62.62%; Table 3b vs. a) and Approach 2 (Table 3d vs. c). Interestingly, BERT and PhraseBERT consistently benefit the most from the added context information (e.g., +41.06 and +37.29 gains). Third, while starting from the same backbone models, Approach 2 yields consistently higher accuracy than Approach 1 (Table 3; 66.57 vs. 62.62 mean accuracy), which is expected as Approach 2 models have more capacity and the backbones are allowed to be finetuned on PS. See Figs. A1–A4 for qualitative PS predictions from the PhraseBERT-based classifier.

5.2 Human Baselines and Upperbound (95% Exact Match) on Phrase Retrieval

To interpret the progress of machine phrase understanding on PR, here, we establish multiple human baselines for both non-experts and linguistics experts (with and without training them).

Experiment We recruit participants and have them perform one or two tests per person. A test consists of 20 PR-pass examples. That is, PR-pass documents are 11-sentence long and are feasible for a person to read in minutes (compared to reading an entire Wikipedia page). We test three groups of people: (1) 21 graduate students at our institution (1 test per person); (2) five Upwork experts (1 test per person); and (3) another five Upwork experts (2 tests per person, i.e., for a total of $2 \times 5 = 10$ tests). The students in Group 1 volunteer to help our study unpaid while the Upworkers (Group 2 and 3) are hired using the same procedure described in Sec. 3.2.

Table 3: Accuracy (%) of state-of-the-art BERT-based models on the PS test set. Contextualized phrase embeddings (“Phrase + Context”) yield substantially higher performance on PS than non-contextualized embeddings (“Phrase”), e.g. a remarkable gain of **+41.06** from 28.57% for BERT.

Model	Approach 1: Cosine similarity		Approach 2: BERT-based classifiers	
	(a) Phrase	(b) Phrase + Context	(c) Phrase	(d) Phrase + Context
PhraseBERT [34]	48.48	63.23 (+14.75)	28.40	65.69 (+37.29)
BERT [9]	48.76	64.32 (+15.56)	28.57	69.63 (+41.06)
SpanBERT [18]	49.58	62.61 (+13.04)	45.79	68.04 (+22.25)
SpanBERT _{Large} [18]	49.76	63.70 (+13.94)	37.57	67.10 (+29.53)
SentenceBERT [26]	48.81	61.32 (+12.51)	27.36	62.02 (+34.66)
SimCSE [12]	48.14	60.51 (+12.37)	36.07	66.95 (+30.88)
mean \pm std	48.92 \pm 0.63	62.62 \pm 1.45	33.96 \pm 7.22	66.57 \pm 2.59

Table 4: Best QA models reach near human-upperbound Exact Match on PR-pass. Yet, ranking models based on phrase embeddings significantly *underperform* QA models.

Accuracy of human groups and models	EM (%)
Group 1: 20 Non-experts (w/o training)	73.60 \pm 7.90
Group 2: 05 Experts (w/o training)	82.00 \pm 12.00
Group 3: 05 Experts (w/ training)	90.50 \pm 3.70
Best human accuracy (4 people)—Upperbound	95.00 \pm 0.00
Best untrained, ranking model (BERT)	47.44
Best PR-trained, ranking model (PhraseBERT)	59.02
Best PR-trained, QA model (Longformer _{Large})	94.28

Results First, we find a large, expected ~ 9 -point gap between non-experts and experts (Table 4; 73.60% vs. 82.00%). Second, we train experts in Group 3 by asking them to do a preliminary tests and giving them feedback before testing them. We find the training to substantially boost expert accuracy further (from 82.00% to 90.50%). Importantly, we find the Human Exact Match (EM) Upperbound to be 95%, i.e. the highest scores that 4 people (among all groups) make. Upon manual inspection of the submissions of these best performers, we find their incorrect answers sometimes partially overlap with the groundtruth and are sometimes reasonable. In other cases, the best performers find acceptable answers but that do not overlap at all with the groundtruth labels in PR. That is, we estimate a 5% of noise in the annotations of PR.

5.3 Phrase Retrieval: In ranking, context only helps BERT embeddings but not others

One way to evaluate the quality of SotA phrase embeddings is by testing:

RQ: *How well do phrase embeddings perform in the **ranking** approach on PR?*

Approaching PR by ranking is a challenging and meaningful embedding test because the embedding of the query is compared against that of all phrase candidates (extracted by tokenizing the document), which can include syntactically-incorrect phrases, meaningless phrases or rare phrases. Such out-of-distribution challenge appears less often in PS or WiC [24], i.e. a binary classification setting.

Experiment As described in Sec. 4.2, the PR train/dev/test splits are 20,147/3K/5K examples and we only use the 5K-example test set to test the models in this ranking experiment (no training). To construct a list of candidate phrases, we split each PR document into multiple sentences (using NLTK sentence splitter) and tokenize each sentence into tokens (using NLTK tokenizer) and build an exhaustive list of n -grams (here, $n \in \{2, 3\}$ only for computational tractability). For every example, we add the groundtruth phrase (which can be longer than 3 words) to the list of candidates (since we are only interested in testing phrase embeddings, not the phrase extractor).

Results We report top-1/3/5 accuracy and MRR@5 on the PR-pass test set in Table 5a. First, for most SotA embeddings, incorporating context sentence into the phrase embeddings *hurts* the accuracy

Table 5: **Ranking** accuracy (%) on **PR-pass** using the state-of-the-art pretrained phrase embeddings (a) and those finetuned on PR-pass via QA-style training (b). See Appendix E for the results on PR-page. Δ (e.g. **-3.62**) denotes the differences between the Top-1 accuracy in the contextualized (“Phrase + Context”) vs. the non-contextualized (“Phrase”) setting.

Model	Phrase				Phrase + Context			
	Top-1	Top-3	Top-5	MRR@5	Top-1 (Δ)	Top-3	Top-5	MRR@5
(a) Pre-trained embeddings								
PhraseBERT [34]	36.62	66.96	75.90	52.20	33.00 (-3.62)	49.60	56.70	41.90
BERT [9]	29.80	47.90	55.40	39.50	47.44 (+17.64)	65.78	73.30	57.30
BERT _{Large} [9]	23.76	38.52	45.40	31.70	42.80 (+19.04)	58.90	64.90	51.30
SpanBERT [18]	20.88	31.04	35.20	26.40	14.40 (-6.48)	30.46	39.80	23.40
SentenceBERT [26]	22.30	50.64	60.60	36.80	25.14 (+2.84)	39.52	46.20	32.90
SimCSE [12]	28.10	53.70	64.60	41.60	32.40 (+4.30)	53.44	62.80	43.70
USE-v5 [8]	43.36	70.12	78.90	57.30	n/a	n/a	n/a	n/a
DensePhrase [21]	32.24	51.30	60.50	42.60	31.50 (-0.74)	46.30	53.80	39.70
(b) PR-pass-trained QA models’ phrase embeddings								
PhraseBERT [34]	59.02	81.58	87.90	70.60	24.98 (-34.04)	37.78	43.90	32.00
BERT [9]	50.10	66.16	71.40	58.60	20.34 (-29.76)	31.40	37.10	26.50
BERT _{Large} [9]	32.70	42.40	45.90	37.80	11.40 (-21.30)	17.00	20.50	14.60
SpanBERT [18]	15.22	22.88	26.60	19.40	8.92 (-6.30)	13.56	16.60	11.60
SentenceBERT [26]	53.14	74.86	80.70	64.20	20.12 (-33.02)	30.04	34.90	25.60
SimCSE [12]	50.96	76.70	83.40	64.00	37.70 (-13.26)	52.38	58.90	45.60
(c) Differences between after vs. before finetuning, i.e. the 6 models in (b) vs. those in (a)								
mean differences	+16.61				-11.95			

rather than helps *except for* BERT embeddings. That is, interestingly, for all BERT embeddings (base and large), the accuracy increases substantially (**+17.64** and **+19.04**; Table 5a) when the one-sentence context is the input. For some reason, most models that started from BERT but were later finetuned lost the capability to leverage the context information, e.g., PhraseBERT, DensePhrase, and SpanBERT in Table 5a.

Second, the best top-1 accuracy scores on PR-pass for non-contextualized (USE-v5; 43.36%) and contextualized (BERT; 47.44%) phrase embeddings are substantially lower than the non-expert baselines (73.60%; Table 4) and Human Upperbound (95%). Future work is required to learn more robust, phrase embeddings for ranking. See Figs. A8–A9 for qualitative examples.

5.4 Phrase Retrieval: Question-Answering models reach near-human accuracy

Consistent with [41], our ranking results in Sec. 5.3 reveal that there exists a large headroom for improving both non-contextualized and contextualized phrase embeddings. Yet, because ranking is a naive approach and QA models [17, 9] are the SotA approach on many QA tasks [25], there, here we train QA models directly on the train set of PR-pass and PR-page in order to answer:

RQ: *How well do SotA semantic-search models perform on PR-pass and PR-page?*

Experiment We take the SotA embeddings tested in Sec. 5.3 and add a linear classification layer on top and finetune each entire classifier on the train set of PR-pass or PR-page for 2 epochs using the default HuggingFace hyperparameters (see Appendix B for finetuning details). Following the standard setup of BERT architectures for QA tasks [17, 9], each QA model predicts the start and end index of the target phrase. Additionally, since PR-page documents are much longer than a typical QA paragraph (Table 2), we also test training Longformer [3], which has a sequence length of 4,096, sufficient for an entire Wikipedia page. We take the models of smallest dev loss and report their test-set performance in Table 6.

Table 6: **Test-set performance (%) of QA models** on PR-pass (a), PR-page (b), and PSD (c). When trained directly on PR-pass (a) and PR-page (b), SotA QA models perform well. However, testing the PR-pass-trained models on PSD shows a significant drop in accuracy (c). That is, QA models tend to understand a *single sense* of a phrase in context well (high PR-pass, PR-page, and PSD EM scores). Yet, they are not able to differentiate two senses of the same phrase (e.g., here, PhraseBERT accuracy drops **-40.90** points between EM+loc vs. EM scores).

Model	(a) PR-pass		(b) PR-page		(c) PSD			
	EM	F ₁	EM	F ₁	EM	F ₁	EM+loc	F ₁ +loc
PhraseBERT [34]	93.42	94.97	85.24	87.19	91.81	93.17	50.91 (-40.90)	51.10
BERT [9]	93.26	94.65	85.64	87.77	92.39	93.64	52.01 (-40.38)	52.21
BERT _{Large} [9]	93.64	95.16	87.36	89.52	93.55	94.73	54.93 (-38.62)	55.16
SpanBERT [18]	93.50	95.02	87.28	87.66	92.17	93.49	54.80 (-37.37)	55.01
SentenceBERT [26]	93.24	94.54	84.66	86.89	92.52	93.83	52.56 (-39.96)	52.80
SimCSE [12]	92.90	94.51	85.68	87.66	92.17	93.58	53.50 (-38.67)	53.72
Longformer [3]	94.26	95.58	89.54	91.15	95.36	96.14	60.52 (-34.84)	60.69
Longformer _{Large} [3]	94.28	95.53	87.58	89.32	95.63	96.27	60.10 (-35.53)	60.19
mean	93.56	95.00	86.92	88.85	93.20	94.36	54.92 (-38.28)	55.11
± std	0.49	0.42	1.93	1.73	1.51	1.23	3.59	3.56

Results On PR-pass, in contrast to the poor performance of ranking models (Sec. 5.3), our PR-pass-trained QA models perform impressively at a near-upperbound level (~ 93 – 94% EM; Table 6a) surpassing the accuracy of trained experts (90.50% EM). Surprisingly, on PR-page where the documents are substantially longer (around $12\times$) than the documents of PR-pass, QA models’ accuracy only drop slightly (from $\sim 94\%$ to ~ 85 – 89% EM; Table 6b). Note that in a full Wikipedia page of PR-page, there might be phrases that could be considered correct but are *not* labeled groundtruth according to our annotations. This remarkable result suggests that training on PR-pass can enable high-performing models on real-world semantic search and potentially other NLP downstream tasks.

5.5 QA-style training improves *non*-contextualized but not contextualized phrase embeddings

As the QA models trained on PR-pass and PR-page perform impressively (Sec. 5.4), almost $1.5\times$ better than the ranking models based on pre-trained embeddings, an interesting question is:

RQ: *Does finetuning following the QA style improve the non- or contextualized phrase embeddings?*

This is important to understand because (1) the impressive QA-models’ performance gain may come from the extra linear-classification layer (not necessarily from the finetuned embeddings); (2) it is under-explored how finetuning improves the contextualization of phrase embeddings.

Experiment We extract the phrase embeddings (both non-contextualized and contextualized) from the PR-pass-trained QA models from Sec. 5.4 (i.e. discarding the linear classification layer) and test them on the PR-pass ranking experiments (as in Sec. 5.3).

Results After training on PR-pass, the *non*-contextualized, phrase embeddings improve substantially for most models at an average gain of **+16.61** in top-1 accuracy (e.g., PhraseBERT top-1 accuracy increases from 36.62% to 59.02%; Table 5b). This result shows that training on PR-pass improves non-contextualized phrase embeddings. In stark contrast, after finetuning on PR-pass via QA-style training, the ranking scores of *contextualized* phrase embeddings surprisingly drop significantly, **-11.95** points on average (Table 5c).

In sum, we are observing a consistent trend that the contextualized phrase embeddings of the original pre-trained BERT (both base and large) are remarkably beneficial for retrieval (i.e. PR). However, after finetuning, e.g. on PR-pass or using other techniques (e.g. in PhraseBERT or SentenceBERT), such benefits of leveraging context disappears. Aligned with Yu and Ettinger [41], we find that incorporating context effectively into phrase embeddings is an open research direction.

5.6 Phrase Sense Disambiguation: Best models also perform poorly

We find that SotA PR-pass-trained QA models reach superhuman accuracy on PR-pass, i.e. finding a phrase of the same meaning (Sec. 5.4). Yet, PR-pass only tests models’ understanding of a *single sense* of the target phrase at a time. It is interesting to study:

RQ: *Do PR-pass-trained QA models understand phrases deeply enough to separate two different senses of the same target phrase?*

Experiment To do that, here we take the best QA models trained on PR-pass and test them on PSD. Note that, PSD has the same task format as PR-pass (see Fig. 1c–e) except that the document is twice as long and contains **two occurrences of the same target phrase**. We only test the QA models trained on PR-pass but not the ranking models (Sec. 5.3) since their accuracy is expected to be worse.

Results Although the PR-pass-trained QA models are never trained on PSD, they interestingly frequently find one occurrence of the target phrase (mean of 93.20% EM; Table 6c). However, they mostly locate the **target phrase in the wrong context passage** with high confidence scores. That is, if we consider also the correctness of the location of the predicted phrase, their EM+loc³ accuracy drops significantly to an average of 54.92%. We also find that directly finetuning on a 2K-example subset of PSD only slightly improved the EM+loc to an average of 64.30% on a 3K-example PSD test set (Appendix F). Note that we estimate the Human Upperbound on PSD to be also 95%, i.e. the same as that of PR-pass. See qualitative examples and predictions of Longformer (i.e. the best model tested) in Figs. A5–A7.

In sum, there is a large headroom for future research on PSD. Consistent with the findings in Sec. 5.5, QA models are *not* yet capable of leveraging surrounding words to understand phrases in context.

6 Discussion and Conclusion

Limitations Our dataset is currently limited to multi-word, English noun-phrases. Furthermore, it is expected to contain around 5% of errors on PR-pass (i.e. the best human performance is 95% EM). On PR-page, there may be more than one correct target phrase; however, we only label one phrase as the correct answer per document. We use only phrases that contain at least one WiC word.

While WiC and English WSD rely exclusively on dictionaries [24] to obtain word senses and example sentences, our data collection approach depends on Wikipedia, WiC, NLP models, and humans for annotation. It is possible to extend to other types of phrases in the future. In sum, we present PiC, the first 3-task suite for evaluating phrases in context. QA models can obtain high accuracy on semantic search after training on our PR-pass and PR-page datasets. Yet, their capability is limited to finding a semantically-similar phrase given a single context that contains the target phrase (in PR-pass). The results on PS and PSD show that SotA phrase embeddings are still limited in encoding a phrase and its context well enough to understand a phrase sense given a specific context. It is interesting future work to improve these models for disambiguating the senses of a phrase in context (PS and PSD).

Acknowledgments and Disclosure of Funding

We thank Qi Li, Peijie Chen, Hai Phan, Giang Nguyen, and Naman Bansal from Auburn University for helpful feedback on the early results. We also thank graduate students from AN’s Deep Learning class for helping us with the human study. We are grateful for the valuable support and feedback from Phat Nguyen. AN is supported by NaphCare Foundations, Adobe gifts, and NSF grants (1850117, 2145767).

References

- [1] Tarik Arici, Hayreddin Ceker, and Ismail Baha Tutar. 2020. MULTI-SPAN QUESTION ANSWERING USING SPAN-IMAGE NETWORK. *preprint*. 14

³For a PSD example, if the predicted span does not intersect at all with the groundtruth span, the EM+loc and F₁+loc scores would be 0. If they intersect, the two scores would be equal to EM and F₁, respectively.

- [2] Shima Asaadi, Saif Mohammad, and Svetlana Kiritchenko. 2019. **Big BiRD: A Large, Fine-Grained, Bigram Relatedness Dataset for Examining Semantic Composition**. 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 505–516, Minneapolis, Minnesota. Association for Computational Linguistics. 1, 2, 3
- [3] Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*. 6, 8, 9, 18
- [4] Michele Bevilacqua, Tommaso Pasini, Alessandro Raganato, Roberto Navigli, et al. 2021. Recent trends in word sense disambiguation: A survey. In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*. International Joint Conference on Artificial Intelligence, Inc. 3
- [5] Steven Bird, Ewan Klein, and Edward Loper. 2009. *Natural Language Processing with Python: Analyzing Text with the Natural Language Toolkit*. O'Reilly, Beijing. 15
- [6] Terra Blevins, Mandar Joshi, and Luke Zettlemoyer. 2021. **FEWS: Large-Scale, Low-Shot Word Sense Disambiguation with the Dictionary**. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 455–465, Online. Association for Computational Linguistics. 3
- [7] Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. **A large annotated corpus for learning natural language inference**. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon, Portugal. Association for Computational Linguistics. 24
- [8] Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, et al. 2018. Universal sentence encoder. *arXiv preprint arXiv:1803.11175*. 6, 8, 17
- [9] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*. 1, 6, 7, 8, 9, 14, 17, 18
- [10] Lev Finkelstein, Evgeniy Gabrilovich, Yossi Matias, Ehud Rivlin, Zach Solan, Gadi Wolfman, and Eytan Ruppín. 2001. Placing search in context: The concept revisited. In *Proceedings of the 10th international conference on World Wide Web*, pages 406–414. 1
- [11] Ute Fischer. 1994. Learning words from context and dictionaries: An experimental comparison. *Applied Psycholinguistics*, 15(4):551–574. 1
- [12] Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. **SimCSE: Simple Contrastive Learning of Sentence Embeddings**. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6894–6910, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics. 4, 6, 7, 8, 9, 15, 17, 18
- [13] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé III, and Kate Crawford. 2021. Datasheets for datasets. *Communications of the ACM*, 64(12):86–92. 26
- [14] Princeton NLP Group. 2022. `princeton-nlp/sup-simcse-roberta-large` · Hugging Face. <https://huggingface.co/princeton-nlp/sup-simcse-roberta-large>. (Accessed on 06/08/2022). 15
- [15] Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. 2020. **spaCy: Industrial-strength Natural Language Processing in Python**. <https://spacy.io/>. 15, 17
- [16] Huggingface. 2022. `super_glue` · Datasets at Hugging Face. https://huggingface.co/datasets/super_glue/viewer/wic. (Accessed on 06/08/2022). 15
- [17] Huggingface. 2022. `transformers/examples/pytorch/question-answering` at main · huggingface/transformers. <https://github.com/huggingface/transformers/tree/main/examples/pytorch/question-answering>. (Accessed on 06/09/2022). 8

- [18] Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S Weld, Luke Zettlemoyer, and Omer Levy. 2020. Spanbert: Improving pre-training by representing and predicting spans. *Transactions of the Association for Computational Linguistics*, 8:64–77. 6, 7, 8, 9, 18
- [19] Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. **Natural Questions: A Benchmark for Question Answering Research**. *Transactions of the Association for Computational Linguistics*, 7:452–466. 24
- [20] LanguageTool. 2022. LanguageTool - Online Grammar, Style & Spell Checker. <https://languagetool.org/>. (Accessed on 06/08/2022). 4
- [21] Jinhyuk Lee, Mujeen Sung, Jaewoo Kang, and Danqi Chen. 2021. Learning Dense Representations of Phrases at Scale. In *Association for Computational Linguistics (ACL)*. 6, 8
- [22] William E Nagy, Patricia A Herman, and Richard C Anderson. 1985. Learning words from context. *Reading research quarterly*, pages 233–253. 1
- [23] Ellie Pavlick, Pushpendre Rastogi, Juri Ganitkevitch, Benjamin Van Durme, and Chris Callison-Burch. 2015. **PPDB 2.0: Better paraphrase ranking, fine-grained entailment relations, word embeddings, and style classification**. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 425–430, Beijing, China. Association for Computational Linguistics. 1
- [24] Mohammad Taher Pilehvar and Jose Camacho-Collados. 2019. **WiC: the Word-in-Context Dataset for Evaluating Context-Sensitive Meaning Representations**. 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 1267–1273, Minneapolis, Minnesota. Association for Computational Linguistics. 1, 3, 7, 10, 24
- [25] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. **SQuAD: 100,000+ Questions for Machine Comprehension of Text**. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas. Association for Computational Linguistics. 2, 3, 8
- [26] Nils Reimers and Iryna Gurevych. 2019. **Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks**. 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 3982–3992, Hong Kong, China. Association for Computational Linguistics. 6, 7, 8, 9, 18
- [27] PiC Team. 2021. upwork_annotation_guidelines.pdf. https://auburn.edu/~tmp0038/upwork_annotation_guidelines.pdf. (Accessed on 06/08/2022). 4
- [28] PiC Team. 2021. upwork_samples.pdf. https://auburn.edu/~tmp0038/upwork_samples.pdf. (Accessed on 06/08/2022). 4
- [29] PiC Team. 2022. upwork_samples.pdf. https://auburn.edu/~tmp0038/upwork_verification_samples.pdf. (Accessed on 06/08/2022). 24
- [30] Wikimedia Team. 2021. API:Categories - MediaWiki. <https://www.mediawiki.org/wiki/API:Categories>. (Accessed on 06/08/2022). 15
- [31] Wikimedia Team. 2021. Wikimedia Downloads. <https://dumps.wikimedia.org/>. (Downloaded on November 1st, 2021). 15
- [32] TensorFlow. 2022. Universal Sentence Encoder | TensorFlow Hub. https://www.tensorflow.org/hub/tutorials/semantic_similarity_with_tf_hub_universal_encoder. (Accessed on 08/11/2022). 6
- [33] Peter D Turney. 2012. Domain and function: A dual-space model of semantic relations and compositions. *Journal of artificial intelligence research*, 44:533–585. 1, 3

- [34] Shufan Wang, Laure Thompson, and Mohit Iyyer. 2021. **Phrase-BERT: Improved Phrase Embeddings from BERT with an Application to Corpus Exploration**. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 10837–10851, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics. 1, 2, 3, 6, 7, 8, 9, 17, 18
- [35] John Wieting, Mohit Bansal, Kevin Gimpel, and Karen Livescu. 2015. **From Paraphrase Database to Compositional Paraphrase Model and Back**. *Transactions of the Association for Computational Linguistics*, 3:345–358. 2, 3
- [36] Yahoo. 2022. Webscope | Yahoo Labs. <https://webscope.sandbox.yahoo.com/catalog.php?datatype=1&did=66>. (Accessed on 08/10/2022). 2
- [37] Yahoo. 2022. Webscope | Yahoo Labs. <https://webscope.sandbox.yahoo.com/catalog.php?datatype=1&did=66>. (Accessed on 08/10/2022). 17
- [38] Ronghui Yang. 2022. `arcface-pytorch/test.py` at master · ronghuaiyang/arcface-pytorch. <https://github.com/ronghuaiyang/arcface-pytorch/blob/master/test.py>. (Accessed on 06/09/2022). 6
- [39] Yinfei Yang, Yuan Zhang, Chris Tar, and Jason Baldridge. 2019. PAWS-X: A Cross-lingual Adversarial Dataset for Paraphrase Identification. In *Proc. of EMNLP*. 1
- [40] Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. **HotpotQA: A Dataset for Diverse, Explainable Multi-hop Question Answering**. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2369–2380, Brussels, Belgium. Association for Computational Linguistics. 2, 3
- [41] Lang Yu and Allyson Ettinger. 2020. **Assessing Phrasal Representation and Composition in Transformers**. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4896–4907, Online. Association for Computational Linguistics. 1, 6, 8, 9
- [42] Yuan Zhang, Jason Baldridge, and Luheng He. 2019. PAWS: Paraphrase Adversaries from Word Scrambling. In *Proc. of NAACL*. 1, 2

Appendix for: PiC: A Phrase-in-Context Dataset for Phrase Understanding and Semantic Search

This dataset is a joint work between Adobe Research and Auburn University.

Authors

- Thang M. Pham – Auburn University thangpham@auburn.edu
- David S. Yoon – Adobe Research syoon@adobe.com
- Trung Bui – Adobe Research bui@adobe.com
- Anh Nguyen – Auburn University anh.ng8@gmail.com

A Training models on Phrase Similarity

Hyperparameters We train each *BERT-based classifier* for a maximum of 100 epochs with early stopping monitored on validation accuracy (patience of 10 epochs). We use a batch size of 200 and Adam optimizer with learning rate $\alpha = 0.0001$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, and $\epsilon = 10^{-8}$.

Training time On average, with early stopping, training a single model using one V100 GPU takes ~ 5 and ~ 8 mins for non-context and context settings, respectively.

B Training QA models on Phrase Retrieval

We finetune each QA model that consists of a linear layer on top of a pretrained model selected in Sec. 5 to predict the start and end indices of answers (as the common setup in BERT QA models [9, 1]). The format of a tokenized input is “[CLS] query [SEP] document [SEP]” with maximum sequence length of 4,096 for Longformer_{Base} and Longformer_{Large} and 512 for the remaining models. If the document exceeds the maximum sequence length, it is split into smaller features for prediction and thus start and end indices with the highest confidence scores are selected.

Hyperparameters We follow HuggingFace scheme to finetune the QA models for 2 epochs using Adam optimizer with learning rate $\alpha = 0.00003$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$. The batch size varies from 1 to 8 for each model: On one V100 GPU, the “base” models can handle 8 examples while the “large” BERT models can only fit 2–4 examples into 16GB of memory. For Longformer_{Large}, we use an A100 GPU to feed one PR-page example into the model. We take the smallest dev-loss models from the training and report their test-set results.

Training time On average, training a single QA model for 2 epochs using one A100 GPU takes ~ 20 mins for base models and ~ 9.5 hours for Longformer_{Large}.

C Data collection

From a Wikipedia dump, we perform a 6-step procedure (summarized in Table A1) for mining a list of *mNPs* sorted descendingly by their likelihood of containing multiple senses. The most polysemous 19,500 *mNPs* are then passed to experts for annotation (Sec. 3.2) and others for verification (Sec. 3.3).

Step 1: Download Wiki articles We download a Wikipedia dump file [31] that contains $\sim 15.78\text{M}$ Wikipedia articles and filter out all empty pages to arrive at $\sim 6.27\text{M}$ non-empty articles.

Step 2: Extract phrases We use NLTK sentence splitter [5] to split each Wikipedia article into multiple sentences. And then we use SpaCy [15] to extract noun phrases and proper nouns. For each phrase, we remove stopwords (those among the 179 stopwords in NLTK v3.6.5) and preceding non-alphanumeric characters. We remove stopwords because they tend to create more pairs of phrases with lexical overlap, rendering the phrase similarity task easier. We then remove *unigram* phrases to arrive at $\sim 286.78\text{M}$ *mNPs*. For each *mNP*, we construct a 3-tuple (phrase, sentence, metadata), i.e. the phrase, its container sentence, and metadata for identifying the Wikipedia webpage (hereafter, page).

Step 3: Remove phrases of a single context We further remove all phrases that (1) contain non-ASCII characters (e.g. “phaenná nâsos”, which are non-English); and (2) appear only once, i.e. keeping those that occur in multiple sentences since we look for *polysemous mNPs*, which have multiple senses and contexts. After this step, $\sim 17.96\text{M}$ phrases remain.

While some phrases with non-ASCII characters are also commonly used in English (e.g., “déjà vu”), we find only 2.48% of phrases at this stage contain non-ASCII characters, and 29% of them are common in English. In short, we are removing only 0.72% of the English phrases that contain non-ASCII characters at Step 3.

Step 4: Find phrases of polysemous words To increase the chance of collecting polysemous *mNPs*, we only keep *mNPs* that have at least one word in the list of 2,345 unique multiple-sense words of WiC [16], arriving at $\sim 6.5\text{M}$ *mNPs*, each appearing in ≥ 2 sentences and in ≥ 1 Wikipedia pages. We empirically find that Step 4 is important and substantially increases our chance of finding polysemous *mNPs* (compared to skipping Step 4).

Step 5: Find phrases in distinct contexts

We observe that a *mNP* is likely to be polysemous when (a) its context sentences are semantically different; and (b) its context Wikipedia pages are of dissimilar categories (e.g. “massive figure” in finance vs. history; Fig. A6).

To implement this filter, we form all possible triplets (phrase, sentence₁, sentence₂) from the list of context sentences of each *mNP*⁴. We compute the cosine similarity of two sentences at the CLS embedding space of a SimCSE [12] provided on HuggingFace [14]. To find triplets where the two sentences are semantically dissimilar, we keep only the triplets where (sentence₁, sentence₂) has a low cosine similarity, i.e. $\in [-0.3, 0.2]$ and the length difference of the two sentences is < 4 words (as two sentences of substantially different lengths often have a low cosine similarity regardless of their semantic differences). As the result, there are $\sim 600\text{K}$ triplets remaining after this step.

We further re-rank these $\sim 600\text{K}$ descendingly by the dissimilarity of the lists of Wikipedia categories⁵ of the context pages that contain sentence₁ and sentence₂. That is, we treat each Wikipedia page’s comma-separated list of categories as an input text to SimCSE and sort the $\sim 600\text{K}$ descendingly by the cosine similarity of the resultant embeddings.

Step 6: Select data for annotation Before asking annotators to label our sorted phrases we perform final filtering by removing proper nouns and phrases whose Wikipedia documents contain missing words.

We perform final filtering to ensure the data given to annotators is in proper format. That is, from $\sim 600\text{K}$ phrases, we filter down to $\sim 475\text{K}$ phrases by applying two filters: (1) Remove all phrases that are proper nouns (i.e. POS tagging returns PROPN) since proper nouns often refer to a single identity and thus unambiguous; (2) Remove all phrases that have a newline character and all phrases whose context Wikipedia page contains missing words (i.e. errors in the Wikipedia dump).

As the result, we obtain a list of $\sim 475\text{K}$ phrases sorted by their estimate chance of carrying two different senses. After manual inspection, we take the top 19,500 triplets of the format (phrase, page₁,

⁴For computational tractability, we only keep at most 32 context sentences per *mNP* where each sentence’s length in words is $\in [5, 25]$.

⁵We use the provided Wikipedia API [30] to obtain the categories for each article as the dump file has no category-related information.

Table A1: Summary of our 3-stage data construction. p, s, m, d, q, l denote target phrase, sentence, metadata, document, query, and label, respectively.

	Remaining #	Data type	Description
Sec. 3.1 Data Collection			
Step 1: Download Wiki articles	~6.27M	articles	Remove ~9.51M empty articles.
Step 2: Extract phrases	~286.78M	(p, s, m)	Extract noun phrases and proper nouns along with their context sentences from Wikipedia articles.
Step 3: Remove phrases of a single context	~17.96M	$(p, [s_1, \dots, s_n], m)$	For each phrase, gather all sentences where that phrase is used.
Step 4: Find phrases of polysemous words	~6.5M	$(p, [s_1, \dots, s_n], m)$	Filter those phrases that do not contain WiC words.
Step 5: Find phrases in distinct contexts			Sort by X_i and apply filters to find pairs of sentences where their phrase potentially has different meanings.
- Sort and filter by semantic dissimilarity	~600K	(p, s_1, s_2, m)	X_1 : cosine similarity scores of sentences embeddings.
- Sort by domain dissimilarity	~600K	(p, s_1, s_2, m)	X_2 : cosine similarity scores of domain embeddings i.e., use categories of each article to get embeddings.
Step 6: Select data for expert annotation	19,500	(p, d_1, d_2)	Remove proper nouns and phrases with missing information and select top 19,500 examples for annotation.
Sec. 3.2 Data Annotations			
	30,042	(p, d, q)	Create a query i.e., paraphrase from the given phrase in each context document.
	15,021	(p, d_1, d_2, l)	Create a Yes/No label for each pair of documents.
Sec. 3.3 Verifying Annotations			
Round 1: MTurk verifier	22,496	(p, d, q)	Verify queries and Yes/No label by MTurkers.
	10,043	(p, d_1, d_2, l)	
Round 2: Upwork verifiers	28,325	(p, d, q)	Verify instances rejected in Round 1.
	13,413	(p, d_1, d_2, l)	

page₂)—i.e. a phrase p and its two context Wikipedia pages where p is the most likely to have two different senses (e.g., see “massive figure” in Fig. 1a)—and hire linguistic experts to annotate them.

Our manual inspection involves taking 1,000 random triplets and manually read them. We find that at least ~30% of the 1,000-triplet subset contain a polysemous target phrase p and two Wikipedia pages that give p two unique meanings. We perform this manual inspection repeatedly throughout the process of inventing and refining the data collection process in order to arrive at the final list of steps as presented in this paper.

Biases in the data collection While there are many filtering steps in our data collection above, most of them are data cleaning filters that are typically needed in a regular NLP dataset construction.

We recognize that there are *three key filters* in our system that impose strong biases:

1. In Step 4, we use only phrases that contain one word in the WiC. That is, we find Step 4 to substantially increase our chance of finding triplets with a polysemous target phrase. We have added this note in the Data Collection description. It is possible to remove Step 4, but that would require a larger human annotation effort to reach the same 15K labeled triplets.
2. In Step 5, we rely on SimCSE to find target phrases that are placed in two sentences of dissimilar meaning.
3. In Step 5, we rely on SimCSE to find target phrases that are placed in two Wikipedia pages of distinct topics.

D Statistics for search queries in Yahoo Search Query dataset

We analyze 4,496 user queries released in the Yahoo Search Query Log To Entities dataset [37] and use SpaCy tokenizer [15] to classify them into 4 main categories: Noun phrases, verb phrases, URLs and others. As a result, noun phrases are the most common query type from users with 3,576 queries ($\sim 79.54\%$) followed by URLs with 675 queries ($\sim 15.01\%$) while verb phrases and other types are less preferred by users. Moreover, the average length of the real user queries is ~ 1.60 which is quite close to our PS task with ~ 2.27 .

Table A2: Statistics of Yahoo queries across different query types.

Query type	# queries	Percentage (%)
Noun phrases	3,576	79.54
Verb phrases	148	3.29
URLs	675	15.01
Others	97	2.16
Total	4,496	100.00

E Quantitative results on PR-page

Table A3: **Ranking** accuracy on **PR-page** using the state-of-the-art pretrained phrase embeddings (a) and those finetuned on PR-pass via QA-style training (b).

Model	Phrase				Phrase + Context			
	Top-1	Top-3	Top-5	MRR@5	Top-1	Top-3	Top-5	MRR@5
(a) Pre-trained embeddings								
BERT [9]	20.70	34.30	41.00	28.20	35.40 (+14.70)	52.10	59.10	44.50
USE-v5 [8]	32.20	52.70	60.80	43.20	n/a	n/a	n/a	n/a
(b) PR-pass-trained QA models' phrase embeddings								
PhraseBERT [34]	49.40	69.40	76.70	60.10	14.70	21.60	26.10	18.70
SimCSE [12]	44.20	66.60	73.50	55.70	24.60	37.80	43.20	31.70

F Finetuning on PSD does not substantially improve accuracy

As PSD has only 5,150 examples, we use all examples for testing in Sec. 5.6 and find the best PR-trained QA models to perform poorly. To further understand the challenge of PSD, here, we ask:

RQ: *Does training on PR-pass and finetuning on PSD improve accuracy on PSD?*

Experiment We take the PR-pass-trained QA models and further finetune them on a subset of PSD to measure how training directly on PSD improves QA models. We split PSD into 1,650/500/3,000 examples for train/dev/test sets, respectively, and finetune the PR-pass-trained QA models on this PSD train set. For comparison with the results in Sec. 5.4, we use the same set of hyperparameters as when finetuning on PR-pass in Sec. 5.4. Below, we report the test-set results of the lowest dev-loss models.

Results On the PSD-3K test set, all models perform poorly at a mean EM score of 55.49% (Table A4a; mean). Interestingly, finetuning the original models using the 2,150 examples (hereafter, PSD-2K) instead of PR-pass decreases accuracy, on average by **-7.21** points. An explanation is that 1,650 PSD training examples are too few for the finetuning to be effective. Indeed, finetuning the PR-pass-trained QA models further on PSD-2K increases the scores for all models by **+8.81** on average (Table A4c; mean). The best model is Longformer_{Base} [3] (Table A4; 70.00 EM), which is still substantially lower than the human upperbound of 95%.

Table A4: Performance of **QA models** on 3,000 PSD **test** examples. (a) and (b) models are **finetuned** only on PR-pass and 2,150 PSD examples (PSD-2K), respectively. (c) models are finetuned on PR-pass first and then finetuned on PSD-2K. All models are “base” unless otherwise specified. The definitions of EM+loc and F₁+loc are in Table 6’s caption.

Models finetuned on	(a) PR-pass		(b) PSD-2K		(c) PR-pass + PSD-2K	
	EM+loc	F ₁ +loc	EM+loc	F ₁ +loc	EM+loc	F ₁ +loc
PhraseBERT [34]	51.10	51.27	34.33 (-16.77)	34.73	55.90 (+4.80)	56.27
BERT [9]	51.90	52.11	45.67 (-6.23)	46.22	63.20 (+11.30)	63.43
BERT _{Large} [9]	56.20	56.43	53.13 (-3.07)	53.85	66.87 (+10.67)	67.11
SpanBERT [18]	55.17	55.35	43.97 (-11.20)	45.33	68.63 (+13.46)	68.93
SentenceBERT [26]	53.17	53.40	39.73 (-13.44)	40.50	59.10 (+5.93)	59.43
SimCSE [12]	54.37	54.56	40.83 (-13.54)	41.28	63.63 (+9.26)	63.82
Longformer [3]	61.07	61.23	63.03 (+1.96)	63.79	70.00 (+8.93)	70.08
Longformer _{Large} [3]	60.90	60.94	65.53 (+4.63)	66.45	67.03 (+6.13)	67.23
mean	55.49	55.66	48.28 (-7.21)	49.02	64.30 (+8.81)	64.54
± std	3.78	3.74	11.26	11.36	4.85	4.78

G Qualitative examples for PS, PR-pass, PR-page and PSD

PS example. Groundtruth: “positive”	
P ₁	moderate speed
P ₂	steady pace
S ₁	Deforestation due to logging and land conversion has likely caused the population to decline at a moderate speed.
S ₂	Deforestation due to logging and land conversion has likely caused the population to decline at a steady pace.

Figure A1: PhraseBERT-based classifier **correctly** predicts “positive” given two phrases P₁ and P₂ with and without the presence of context S₁ and S₂. Here, to humans, the phrases are non-polysemous and have the same meaning.

PS example. Groundtruth: “negative”	
P ₁	greatest emphasis
P ₂	highest stress
S ₁	However, the rock art had the greatest emphasis on domesticated cattle.
S ₂	However, the rock art had the highest stress on domesticated cattle.

Figure A2: PhraseBERT-based classifier **correctly** predicts “negative” given two phrases P₁ and P₂ with and without the presence of context S₁ and S₂. Here, to humans, the two phrases are non-ambiguously carrying different meanings.

PS example. Groundtruth: “positive”	
P ₁	unique image
P ₂	uncommon style
S ₁	Bayliss has been praised for her unique image and tendency to change up songs.
S ₂	Bayliss has been praised for her uncommon style and tendency to change up songs.

Figure A3: PS case that requires context to determine similarity. Without context, a PhraseBERT-based classifier incorrectly thinks P₁ and P₂ are different. Yet, it changes the prediction to “positive”, i.e. thinking two phrases have the same meaning, when the context is taken into account.

PS example. Groundtruth: “negative”	
P ₁	permanent post
P ₂	stable location
S ₁	His assistant, John Carver took over as caretaker manager, managing one win, but was not considered for the permanent post, and left in September 2004.
S ₂	His assistant, John Carver took over as caretaker manager, managing one win, but was not considered for the stable location, and left in September 2004.

Figure A4: PS case that requires context to determine similarity. Without context, PhraseBERT-based classifier incorrectly thinks P₁ and P₂ carry the same meaning. Yet, it correctly changes the prediction to “negative” when the context is taken into account.

PSD example.	
d	<p>Bubble memory is a type of non-volatile computer memory that uses a thin film of a magnetic material to hold small magnetized areas, known as "bubbles" or "domains", each storing one bit of data. The material is arranged to form a series of parallel tracks that the bubbles can move along under the action of an external magnetic field. The bubbles are read by moving them to the edge of the material where they can be read by a conventional magnetic pickup, and then rewritten on the far edge to keep the memory cycling through the material. In operation, bubble memories are similar to delay line memory systems. Bubble memory started out as a promising technology in the 1970s, offering memory density of an order similar to hard drives but performance more comparable to core memory while lacking any moving parts. This led many to consider it a contender for a "universal memory" that could be used for all storage needs. The introduction of dramatically faster semiconductor memory chips pushed bubble into the slow end of the scale, and equally dramatic improvements in hard drive capacity made it uncompetitive in price terms. Bubble memory was used for some time in the 1970s and 80s where its non-moving nature was desirable for maintenance or shock-proofing reasons. The introduction of Flash RAM and similar technologies rendered even this niche uncompetitive, and bubble disappeared entirely by the late 1980s. History. Precursors.</p> <p>The Inkerman stone, of which the building is made, was mined near Sevastopol and transported by barges. No convenient mooring facilities existed at that time, so the barges had to anchor in the harbor and the load was moved to the shore by boats and then transported to the construction site across the steppe. During the first year of construction, the builders concentrated on the basic structure at the expense of various facilities and decorations. At the end of 1816, the lighthouse looked like a conic 36-metre-high stone tower with a wooden 3.3-metre-high decagonal lantern. The lighthouse became operational in 1817 after its lighting system had been repaired. Three houses were built next to the tower to accommodate the lighthouse personnel and for storage needs. However, cold and humid winters of the Tarkhanut Peninsula, however, made these houses nearly unsuitable for living. In 1862, the lighting system was upgraded, and the spread of light reached 12.4 miles. In 1873, the construction resumed along with cleaning efforts of the surrounding areas. The building was finished and painted white. In 1876, an additional telegraph spot was built near the tower.</p>
q ₁	storehouse purposes Groundtruth: storage needs & Prediction: storage needs (confidence: 0.99)
q ₂	data caching Groundtruth: storage needs & Prediction: storage needs (confidence: 0.99)

Figure A5: Given document *d*, our Longformer_{Large} QA model trained on PR-pass correctly retrieves storage needs in the second paragraph for the query *q*₁ “storehouse purposes” but fails to retrieve the answer when the query *q*₂ is “data caching”. The predicted answer for *q*₂ should be storage needs (i.e. in the first passage) since this phrase relates to caching data digitally in computers while storage needs refers to physically storing objects.

PSD example.		
<i>d</i>	<p>In the libretto, Delilah is portrayed as a seductive "femme fatale", but the music played during her parts invokes sympathy for her. The 1949 biblical drama "Samson and Delilah", directed by Cecil B. DeMille and starring Victor Mature and Hedy Lamarr in the titular roles, was widely praised by critics for its cinematography, lead performances, costumes, sets, and innovative special effects. It became the highest-grossing film of 1950, and was nominated for five Academy Awards, winning two. According to "Variety", the film portrays Samson as a stereotypical "handsome but dumb hulk of muscle". Samson has been especially honored in Russian artwork because the Russians defeated the Swedes in the Battle of Poltava on the feast day of St. Sampson, whose name is homophonous with Samson's. The lion slain by Samson was interpreted to represent Sweden, as a result of the lion's placement on the Swedish coat of arms. In 1735, C. B. Rastrelli's bronze statue of Samson slaying the lion was placed in the center of the great cascade of the fountain at Peterhof Palace in Saint Petersburg. Samson is the emblem of Lungau, Salzburg and parades in his honor are held annually in ten villages of the Lungau and two villages in the north-west Styria (Austria). During the parade, a young bachelor from the community carries a massive figure made of wood or aluminum said to represent Samson. The tradition, which was first documented in 1635, was entered into the UNESCO list of Intangible Cultural Heritage in Austria in 2010. Samson is one of the giant figures at the "Ducasse" festivities, which take place at Ath, Belgium.</p> <p>On September 22, 2015, Honda announced that they had sold over 1 million Activas in five months in the Indian market, from April to August. Honda launched their 5th generation of Honda Activa in 2018, and the sixth-generation Honda Activa 6G have been launched in India with prices starting at 63,912 (ex-showroom, Delhi). Milestones. In April, 2014, "The Economic Times" reported the Honda Activa to be the best selling two wheeler in India, outselling the Hero Splendor. During the month of September 2013, 141,996 Honda Activa scooters were sold, nearly equal to Honda's entire annual sales in North America. The 110cc Activa is the company's biggest seller, by far. It is responsible for over 2,00,000 sales units each month. In November 2018, HMSI crossed the 2.5 crore sales mark in the scooter segment. It has become the first company to reach this milestone and the biggest contributor to this massive figure is the Honda Activa. It took Honda 13 years to achieve the one crore sales figure, but it managed to add another crore in the span of just three years. It then went on to achieve the next 50 lakh in just one year.</p>	
<i>q</i> ₁	huge model	Groundtruth: massive figure & Prediction: massive figure (confidence: 0.99)
<i>q</i> ₂	giant number	Groundtruth: massive figure & Prediction: massive figure (confidence: 0.99)

Figure A6: Given document *d*, Longformer_{Large} model trained with QA approach on PR-pass correctly retrieves **massive figure** in the second paragraph for the query *q*₂ "giant number" but *fails* to retrieve the answer when the query *q*₁ is "huge model". The predicted answer for *q*₁ should be **massive figure** in the first passage since this phrase relates to a physical shape instead of a number.

PSD example.	
d	<p>Eva held ambitions to replace Hortensio Quijano for the 1951 election, although her poor health kept her from this. Nonetheless many were concerned that her agenda would be pushed through. In march of 1951 the government arrested several retired army officers due to their dissent and disapproval of Perón's administration. This raised tensions among the rest of the army, although action did not occur. By September tensions had risen among the military due to the unrivalled power of the Peronist regime. On September 28, 1951, during the election, Menéndez led the military uprising in an attempt to overthrow the government. He led a core of officers, commanding a division, and left Campo de Mayo bound for the Casa Rosada. Resolve for the uprising, especially among the non-commissioned officers and enlisted men, was not strong enough. They were not prepared to fight their own countrymen. The uprising was over as soon as opposition was encountered, almost completely bloodless. Perón admired the loyalty of the troops and pardoned all those involved.</p> <p>The design uses a similar standard to the JVX in terms of distortion reduction with crossbraces and 27 cells but that's where the similarity ends. Petra was built from the ground up with entirely new panel shaping and trim. Petra has a highly elliptical planform and very high sweep. NZ Aerosports say she has a high roll rate, a long recovery arc and high maximal glide ratio. She is said to deliver unrivalled power in the turn, plane out and flare. Petra has a long list of World Records, National and International titles to back that up. She had an impressive debut at the PD Big Boy Pants event in July 2011, with Nick Batsch setting a new distance world record of 222.45m (729ft). One month later Nick took out the Pink Open in Klatovy and the FAI World Cup also; first in distance, speed and overall. He also won the 2011 US CP nationals on Petra. Patrick Boulongne came 2nd in the European Championships and 6th overall at the World Cup with Petra in his first competition with her. He went on to win the 2011 French Canopy Piloting Nationals.</p>
q₁	incomparable energy Groundtruth: unrivalled power & Prediction: unrivalled power (confidence: 0.99)
q₂	indomitable strength Groundtruth: unrivalled power & Prediction: unrivalled power (confidence: 0.99)

Figure A7: Given document **d**, Longformer_{Large} model trained via the QA approach on PR-pass correctly retrieves **unrivalled power** in the first paragraph for the query **q₂** “indomitable strength” but *fails* to retrieve the answer when the query **q₁** is “incomparable energy”. The predicted answer for **q₁** should be **unrivalled power** in the second passage since the second passage changes “unrivalled power” meaning to a competition strength instead of military power.

PR-pass example.		Groundtruth: common thought
d	<p>As the medical corps grew in size there was also specialization evolving. Physicians surfaced that specialized in disease, surgery, wound dressing and even veterinary medicine. Veterinary physicians were there to tend to livestock for agricultural purposes as well as combat purposes. The Cavalry was known for their use of horses in combat and scouting purposes. Because of the type of injuries that would have been commonly seen, surgery was a somewhat common occurrence. Tools such as scissors, knives and arrow extractors have been found in remains. In fact, Roman surgery was quite intuitive, in contrast to common thought of ancient surgery. The Roman military surgeons used a cocktail of plants, which created a sedative similar to modern anesthesia. Written documentation also showed surgeons would use oxidation from a metal such as copper and scrape it into wounds, which provided an antibacterial effect; however, this method was most likely more toxic than providing an actual benefit. Doctors had the knowledge to clean their surgical instruments with hot water after each use. Wounds were dressed, and dead tissue was removed when bandages were changed.</p>	
q	prevalent theory	
R	0.882 common thought 0.855 common thought of 0.702 fact 0.698 to common thought 0.675 common occurrence	

Figure A8: A ranking model based on the phrase embeddings of the PR-pass-trained PhraseBERT QA model correctly ranks and retrieves the most semantically relevant answer “common thought” as the top-1 prediction in the retrieval list **R** for the query “prevalent theory” in a PR-pass example (which contains a document **d** and a query **q**).

PR-page example.		Groundtruth: continued risk
d	<p>... Following a United Nations agreement between Indonesia and Portugal, a UN-supervised referendum held on 30 August 1999 offered a choice between autonomy within Indonesia and full independence. The people of East Timor voted overwhelmingly for independence. An Australian-led and Indonesian-sanctioned peacekeeping force, INTERFET, was sent into the territory to restore order following a violent 'scorched-earth' policy carried out by pro-integration militia and supported by elements of the Indonesian military. In response to Australia's involvement, Indonesia abrogated the 1995 security pact, asserting that Australia's actions in East Timor were inconsistent with 'both the letter and spirit of the agreement'. Official meetings were cancelled or delayed, including the Indonesia-Australia Ministerial Dialogue, which would not reconvene until March 2003. INTERFET was later replaced by a UN force of international police, UNTAET, which formed a detachment to investigate alleged atrocities. "Tampa" affair and the War on Terror. The relationship came under strain in August 2001 during the "Tampa" affair, when Australia refused permission for the Norwegian freighter ship MV "Tampa" to enter Australian waters while carrying Afghan asylum seekers that it had rescued from a distressed fishing vessel in international waters. The Indonesian Search and Rescue Agency did not immediately respond to requests from Australia to receive the vessel. When the ship entered Australian territorial waters after being refused permission, Australia attempted without success to persuade Indonesia to accept the asylum seekers. Norway also refused to accept the asylum seekers and reported Australia to international maritime authorities. The incident prompted closer coordination between Indonesian and Australian authorities, including regional conferences on people smuggling, trafficking in persons and other transnational crime. In 2002, a terrorist attack in Kuta, Bali killed 202 people, including 88 Australians, and injured a further 240. Jemaah Islamiyah, a violent Islamist group, claimed responsibility for the attack, allegedly in retaliation for Australia's support for East Timorese independence and the War on Terror. A subsequent attack in 2005 resulted in the deaths of a further 20 people, including 15 Indonesians and 4 Australians. The 2003 Marriott Hotel bombing was also perceived as targeted at Western interests in Indonesia; Al Qaeda claimed the attack was carried out by a Jemaah Islamiyah suicide bomber in response to actions of the United States and its allies, including Australia. A 2004 attack on the Australian embassy in Jakarta by Jemaah Islamiyah resulted in the deaths of nine Indonesians. The following year, Indonesian diplomatic and consular premises in Australia received a number of hoax and threat messages. Since then, both the United States and Australian governments have issued warnings against travel to Indonesia, advising their citizens of a continued risk of attacks. These incidents prompted greater cooperation between law enforcement agencies in the two countries, building on a 1999 agreement on drug trafficking and money laundering. The Australian Federal Police's Jakarta Regional Cooperation Team provided assistance to the Indonesian National Police, and has contributed to the Jakarta Centre for Law Enforcement Cooperation. This relationship has attracted criticism, particularly following the arrest and sentencing of the Bali Nine, a group of nine Australians arrested in Denpasar while attempting to smuggle heroin from Indonesia to Australia. The 2005 conviction of Schapelle Corby for attempting to smuggle drugs to Bali also attracted significant attention in the Australian media. The 2004 Indian Ocean earthquake prompted a significant humanitarian response from Australia, including a \$1 billion aid package from the federal government, a further \$17.45 million contribution from state and territory governments, and the commitment of 900 Australian Defence Force personnel to relief efforts in northern Sumatra and Aceh. A telethon broadcast on Australia's three major commercial television networks called "" generated pledges of more than \$10 million, contributing to total private aid of \$140 million. The Eighth "Australia-Indonesia Ministerial Forum" (AIMF) was held in Bali on 29 June 2006 and was attended by five Australian and eleven Indonesian ministers. A key outcome was support for the conclusion of a security agreement, later realised as the Lombok Agreement, providing a framework for the development of the security relationship by the end of 2006 on defence, law enforcement, counter-terrorism, intelligence, maritime security, aviation safety, WMD non-proliferation, and bilateral nuclear cooperation for peaceful purposes. Australia-Indonesia-East Timor Trilateral Ministerial Meetings occurred three times to September 2006. Recent relations. 2010 President Susilo Bambang Yudhoyono visited Australia in April 2010, and became the second Indonesian leader to address federal parliament: Finally, I look forward to a day in the near future. The day when policy makers, academicians, journalists and other opinion leaders all over the world take a good look at the things we are doing so well together. And they will say: these two used to be worlds apart. But they now have a fair dinkum of a partnership. ...</p>	
q	sustained threat	
R	0.830 threat . 0.802 potential threat 0.800 threat reached 0.787 threat as 0.787 threat to	

Figure A9: A ranking model based on the non-contextualized embeddings of USE-v5 fails to retrieve the correct answer “continued risk” for the query “sustained threat” in the PR-page example (which contains a document *d* and a query *q*). The top-5 phrases retrieved (R) contains the word “threat” but have no identifier conveying the “continued” or “sustained” sense. Here, the Wikipedia page is truncated to fit into a single manuscript page.

H Verifying annotations

There are two common methods for evaluation of dataset quality: (1) Verify only a small, random subset [24] to estimate the quality of the full dataset or (2) verifying the entire dataset with multiple annotators and use the inter-annotator agreement (IAA) to control quality [7, 19]. The first approach for approximation is budget-friendly but it remains unknown whether the rest of examples are at high quality, while IAA is more desired but annotating thousands of instances can be prohibitively slow and costly.

We propose a **hybrid approach** to evaluate (leveraging both linguistic experts and non-experts) and ensure high quality for 30,042 queries and 15,021 Yes/No answers at lower cost compared to IAA via two rounds:

1. First, we ask around 1,000 highly qualified freelancers on Amazon Mechanical Turk (MTurk verifiers) to verify whether the *query* annotated by our Upwork annotators is interchangeable i.e. has the same meaning with the given *phrase* in *paragraph*. To verify Yes/No answers, MTurk verifiers need to read two short paragraphs containing the same phrase like Upwork annotators to make decisions. We do not show answers to the MTurk verifiers to avoid biases.
2. Second, we continue hiring 5 Upwork verifiers who are writing experts to double-check those instances rejected by MTurk verifiers from the previous round and only discard an example if the Upwork verifiers agree with MTurk verifiers.

H.1 Round 1: Verification by MTurk non-experts

We use AMT platform to recruit more than 1,000 MTurk verifiers. Also, we use Gorilla (gorilla.sc) to develop user interface to collect answers from participants because (1) Gorilla provides easy-to-use tools to build graphical interface, (2) it is straightforward to monitor and discard results from unqualified participants and (3) we can easily share the experiment with MTurk verifiers via a link. Per 30 verified answers in around ~ 20 minutes, the verification process costs us \$5.6 (AMT fees included) and 1 token to Gorilla to a single MTurk verifier.

Participants are given detailed instructions along with 5 practice samples to get familiar with the task (Fig. A10). They need to pass an evaluation checkpoint including 6 questions randomly sampled from our verified question bank in order to start working with sets of 30 questions. With this approach, all examples in the dataset are verified once and as a result, 22,496/30,042 queries ($\sim 74.88\%$) and 10,043/15,021 Yes/No answers ($\sim 66.86\%$) accepted by MTurkers are considered high quality since they are annotated by a writing expert and confirmed by a qualified English native speaker. The remaining 7,546 queries and 4,978 Yes/No answers rejected that are passed to another group of 5 writing experts for confirmation.

H.2 Round 2: Verification by Upwork experts

We hired another set of 5 writing experts from Upwork (Upwork verifiers) with an hourly rate of \$25-40/hour to verify 12,524 examples rejected by MTurk verifiers, i.e., at an average cost of approximately \$0.26 per example. See a sample assignment given to an Upwork expert in [29].

We rely on IAA to decide whether to accept or reject an example. Specifically, we use the *same question types* as shown to MTurk verifiers in the previous step and see whether these Upwork verifiers agree with the Upwork annotators to keep this example or with MTurk verifiers to reject it. We find that the agreement between the first- and third-round annotators are 5,829 (out of 7,546) paraphrases and 3,370 (out of 4,978) Yes/No answers in total and thus the total high-quality queries and Yes/No answers we achieve are 28,325 and 13,413, respectively.

Instructions

In this study, you will work with 2 types of Yes/No questions:

Type 1: Same paragraphs, different phrases

P1: 2020, The cancellation of the conference in 2020 due to the COVID-19 pandemic led to an online series of lectures entitled Skeptical Inquirer Presents. These sessions included presentations by well known figures in the skeptical community and opportunities for viewers to ask questions. Conference details.

P2: 2020, The cancellation of the conference in 2020 due to the COVID-19 pandemic led to an online series of lectures entitled Skeptical Inquirer Presents. These sessions included presentations by Omicron variant in the skeptical community and opportunities for viewers to ask questions. Conference details.

Question: In both passages, are well known figures and Omicron variant interchangeable (i.e., having the same meaning)?

Type 2: Different paragraphs, same phrases

P1: In 1949, a memorial for the 442 Regimental Central Postal Directory was incorporated and remembered for the Japanese-American soldiers who had fallen during World War II. Every year during the Obon Festival, families gather to upkeep their relatives' tombstones and to visit the spirits. Biddy Mason, nurse and philanthropist, was one of the well known figures to be buried at the cemetery in 1991. There is a section called the "Shower's Rest" in which 400 carnival workers and a circus performers are buried by a memorial that is decorated with a lion.

P2: 2020, The cancellation of the conference in 2020 due to the COVID-19 pandemic led to an online series of lectures entitled Skeptical Inquirer Presents. These sessions included presentations by well known figures in the skeptical community and opportunities for viewers to ask questions. Conference details.

Question: In both passages, does well known figures have the same meaning?

Continue

Training

P1: This library was supposedly founded in 1945, but has started work in current object in 1947. During 1953-1956 it has played the role of the national library since the Kosovo National Library was closed. Academy of Sciences and Art is a necessary institution for the education system that is placed in Pristina. This institution was founded in 1975 as the Association of Science and Arts of Kosovo.

P2: This library was supposedly founded in 1945, but has started work in current object in 1947. During 1953-1956 it has played the role of the national library since the Kosovo National Library was closed. Academy of Sciences and Art is a big supermarket for the education system that is placed in Pristina. This institution was founded in 1975 as the Association of Science and Arts of Kosovo.

Question 3: In both passages, are necessary institution and big supermarket interchangeable (i.e., having the same meaning)?

No Yes

Your progress: 0/36 trials completed

P1: HotDocs transforms documents and graphical (PDF) forms into document-generation templates and deploys of these templates to various server environments. Document modeling in HotDocs can range from variable insertions to the formation and insertions of complex, computed variables. Business logic consisting of IF/THEN statements and REPEAT loops can be built into the template to control the inclusion or exclusion of language blocks. HotDocs includes a variety of other scripting instructions and sets of pre-packaged functions using boolean logic.

P2: HotDocs transforms documents and graphical (PDF) forms into document-generation templates and deploys of these templates to various server environments. Document modeling in HotDocs can range from variable insertions to the formation and insertions of complex, computed variables. Business logic consisting of IF/THEN statements and REPEAT loops can be built into the template to control the inclusion or exclusion of linguistic sections. HotDocs includes a variety of other scripting instructions and sets of pre-packaged functions using boolean logic.

Q: In both passages, are language blocks and linguistic sections interchangeable (i.e., having the same meaning)?

No Yes

Note: Yes/No buttons will be displayed in 20 seconds. Please read the given contents carefully before answering question.

Training

P1: Mercury is poured over the dirt with bare hands. The method leaves much gold undetected, and therefore some miners are using metal detectors. The mercury pollution in the area is both an environment problem, and a health hazard. Most of the gold gets transported to the North of Paramaribo where the gold buyers are located.

P2: Tools such as KLEE, Cloud9, and Otter take this approach by implementing models for file system operations, sockets, IPC, etc. Forking the entire system state. Symbolic execution tools based on virtual machines solve the environment problem by forking the entire VM state. For example, in S2E each state is an independent VM snapshot that can be executed separately.

Question 1: In both passages, does environment problem have the same meaning?

No Yes

Training

P1: This library was supposedly founded in 1945, but has started work in current object in 1947. During 1953-1956 it has played the role of the national library since the Kosovo National Library was closed. Academy of Sciences and Art is a necessary institution for the education system that is placed in Pristina. This institution was founded in 1975 as the Association of Science and Arts of Kosovo.

P2: This library was supposedly founded in 1945, but has started work in current object in 1947. During 1953-1956 it has played the role of the national library since the Kosovo National Library was closed. Academy of Sciences and Art is a big supermarket for the education system that is placed in Pristina. This institution was founded in 1975 as the Association of Science and Arts of Kosovo.

Question 3: In both passages, are necessary institution and big supermarket interchangeable (i.e., having the same meaning)?

NO. The phrase necessary institution means a required establishment of the livery stable which does not totally refer to a big supermarket.

Yes Continue

Training

P1: Mercury is poured over the dirt with bare hands. The method leaves much gold undetected, and therefore some miners are using metal detectors. The mercury pollution in the area is both an environment problem, and a health hazard. Most of the gold gets transported to the North of Paramaribo where the gold buyers are located.

P2: Tools such as KLEE, Cloud9, and Otter take this approach by implementing models for file system operations, sockets, IPC, etc. Forking the entire system state. Symbolic execution tools based on virtual machines solve the environment problem by forking the entire VM state. For example, in S2E each state is an independent VM snapshot that can be executed separately.

Question 1: In both passages, does environment problem have the same meaning?

NO. The phrase environment problem in the first passage mentions an issue of the physical environment in which we are living while in the second passage, it means the issue of the digital environment of an operating system. Thus, it should have different meanings.

No Continue

Figure A10: Gorilla layouts shown to MTurkers to verify annotations in the first round.

I Data Sheet

We follow the documentation template provided by Gebru et al. 2021 [13].

I.1 Motivation

For what purpose was the dataset created? Understanding phrases in context plays a vital role in solving many Natural Language Understanding (NLU) tasks such as question answering or reading comprehension. While there are *word-sense* disambiguation datasets like WiC, no such benchmarks exist for *phrases*. Existing phrase benchmarks compare only phrases without context and some of them contain numerous phrase pairs that have lexical overlap. The major drawback is no human annotation of how a phrase’s meaning changes w.r.t the context. This motivates us to construct a Phrase-in-Context benchmark to drive the development of contextualized phrase embeddings in NLU.

Who created the dataset (e.g., which team, research group) and on behalf of which entity (e.g., company, institution, organization)? Auburn University and Adobe Research.

I.2 Composition/collection process/preprocessing/cleaning/labeling and uses

We describe the data construction process, annotation and verification methods in our paper (See Sec. 3 and Sec. 4).

I.3 Distribution

Will the dataset be distributed to third parties outside the entity (e.g., company, institution, organization) on behalf of which the dataset was created? We release three datasets PS, PR (including PR-pass and PR-page) and PSD to the public.

How will the dataset be distributed (e.g., tarball on website, API, GitHub)? The datasets are released and can be viewed and downloaded on HuggingFace <https://huggingface.co/PiC> or on our website <https://phrase-in-context.github.io>.

When will the dataset be distributed? It has been released in July 2022.

What is the dataset format and how it can be read? We use JSON - a widely used data format for PiC dataset and follow a scheme of HuggingFace datasets to host it. Three datasets PS, PR and PSD in the PiC dataset are loaded as follows:

```
1 # The following pip command is to install the HuggingFace library
2 "datasets": pip3 install datasets
3
4 from datasets import load_dataset
5
6 ps      = load_dataset("PiC/phrase_similarity")
7 pr_pass = load_dataset("PiC/phrase_retrieval", "PR-pass")
8 pr_page = load_dataset("PiC/phrase_retrieval", "PR-page")
9 psd     = load_dataset("PiC/phrase_sense_disambiguation")
```

Will the dataset be distributed under a copyright or other intellectual property (IP) license, and/or under applicable terms of use (ToU)? Our dataset is distributed under the CC-BY-NC 4.0 license.

I.4 Maintenance

How can the owner/curator/manager of the dataset be contacted (e.g., email address)? Thang Pham (thangpham@auburn.edu) and Anh Nguyen (anh.ng8@gmail.com) will be responsible for maintenance.

Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete instances)? Yes. If we include more tasks or find any errors, we will correct the dataset. It will be updated on our website and also HuggingFace.

If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so? They can contact us via email for the contribution.