# GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding

**Alex Wang,**[1] **Amanpreet Singh,**[1] **Julian Michael,**[2] **Felix Hill,**[3]
**Omer Levy,**[2] **and Samuel R. Bowman**[1]

[1]New York University, New York, NY

[2]Paul G. Allen School of Computer Science & Engineering, University of Washington, Seattle, WA

[3]DeepMind, London, UK

{alexwang,amanpreet,bowman}@nyu.edu
{julianjm,omerlevy}@cs.washington.edu
felixhill@google.com

## Abstract

For natural language understanding (NLU) technology to be maximally useful, it must be able to process language in a way that is not exclusively tailored to a specific task, genre, or dataset. In pursuit of this objective, we introduce the General Language Understanding Evaluation (GLUE) benchmark, a collection of tools for evaluating and analyzing the performance of models across a diverse set of existing NLU tasks. By including tasks with limited training data, GLUE is designed to favor and encourage models that share general linguistic knowledge across tasks. GLUE also includes a hand-crafted diagnostic test suite that enables detailed linguistic analysis of models. We evaluate baselines based on current methods for transfer and representation learning and find that multi-task training on all our tasks yields better results than training a separate model for each task. However, the low absolute performance of our best model indicates the need for improved general NLU systems.

## 1 Introduction

The human ability to understand language is *general*, *flexible*, and *robust*. In contrast, most NLU models above the word level are designed for a specific task and struggle with out-of-domain data. If we aspire to develop models with understanding beyond the detection of superficial correspondences between inputs and outputs, then it is critical to develop a more unified model that can learn to execute a range of different linguistic tasks in different domains.

To facilitate research in this direction, we present the General Language Understanding Evaluation (GLUE, gluebenchmark.com) benchmark: a collection of NLU tasks including question answering, sentiment analysis, and textual entailment, and an associated online platform for model evaluation, comparison, and analysis. GLUE does not place any constraints on model architecture beyond the ability to process single-sentence and sentence-pair inputs and to make corresponding predictions. For some GLUE tasks, training data is plentiful, but for others it is limited or fails to match the genre of the test set. GLUE therefore favors models that can learn to represent linguistic knowledge in a way that facilitates sample-efficient learning and effective knowledge-transfer across tasks. While none of the datasets in GLUE were created from scratch for the benchmark, four of them feature privately-held test data, which will be used to ensure that the benchmark is used fairly.

To understand the types of knowledge learned by models and to encourage linguistic or semantically-meaningful solution strategies, GLUE also includes a set of hand-crafted analysis examples for probing trained models. This dataset is designed to highlight common phenomena, such as the use of world knowledge, logical operators, and lexical entailments, that models must grasp if they are to robustly solve the tasks.

To better understand the challenged posed by GLUE, we conduct experiments with simple baselines and state-of-the-art sentence representation models. We find that unified multi-task trained models slightly outperform comparable models trained on each task separately. Our best multi-task model makes use of ELMo (Peters et al., 2018), a recently proposed pre-training technique. However, this model still achieves a fairly low absolute score, indicating room for improved general NLU systems. Analysis with our diagnostic dataset reveals that our baseline models deal well with strong lexical signals but struggle with deeper logical structure.

In summary, we offer: (i) A suite of nine sentence or sentence-pair NLU tasks, built on established annotated datasets and selected to cover

a diverse range of text genres, dataset sizes, and degrees of difficulty. (ii) An online evaluation platform and leaderboard, based primarily on privately-held test data. The platform is model-agnostic, and can evaluate any method capable of producing results on all nine tasks. (iii) An expert-constructed diagnostic evaluation dataset. (iv) Baseline results for several major existing approaches to sentence representation learning.

## 2 Related Work

Collobert et al. (2011), one of the earliest works exploring deep learning for NLP, used a multi-task model with a shared sentence understanding component to jointly learn POS tagging, chunking, named entity recognition, and semantic role labeling. More recent work has explored using labels from core NLP tasks to supervise training of lower levels of deep neural networks (Søgaard and Goldberg, 2016; Hashimoto et al., 2016) and automatically learning cross-task sharing mechanisms for multi-task learning (Ruder et al., 2017).

Beyond multi-task learning, much work towards developing general NLU systems has focused on sentence-to-vector encoder functions (Le and Mikolov, 2014; Kiros et al., 2015, i.a.), leveraging unlabeled data (Hill et al., 2016; Peters et al., 2018), labeled data (Conneau and Kiela, 2018; McCann et al., 2017), and combinations of these (Collobert et al., 2011; Subramanian et al., 2018). In this line of work, a standard evaluation practice has emerged, recently codified as SentEval (Conneau et al., 2017; Conneau and Kiela, 2018). Like GLUE, SentEval relies on a set of existing classification tasks that involve either one or two sentences as inputs. Unlike GLUE, SentEval only evaluates sentence-to-vector encoders. Specifically, SentEval feeds the output of a pre-trained sentence encoder into lightweight task-specific models (typically linear classifiers) that are trained and tested on task-specific data.

SentEval is well-suited for evaluating sentence representations *in isolation*. However, cross-sentence contextualization and alignment, such as that yielded by methods like soft-attention, is instrumental in achieving state-of-the-art performance on tasks such as machine translation (Bahdanau et al., 2014; Vaswani et al., 2017), question answering (Seo et al., 2016; Xiong et al., 2016), and natural language inference (Rocktäschel et al., 2016) . GLUE is designed to facilitate the development of these methods: it is model-agnostic, allowing for any kind of representation or contextualization, including models that use no systematic vector or symbolic representations for sentences whatsoever. Indeed, among the baseline models we evaluate, the use of attention consistently leads to improved performance on GLUE.

GLUE also diverges from SentEval in the selection of evaluation tasks that are included in the suite. Many of the SentEval tasks are closely related to sentiment analysis, such as MR (Pang and Lee, 2005), SST (Socher et al., 2013), CR (Hu and Liu, 2004), and SUBJ (Pang and Lee, 2004). Other tasks are so close to being solved that evaluation on them is relatively uninformative, such as MPQA (Wiebe et al., 2005) and TREC question classification (Voorhees et al., 1999). In GLUE, we attempt to construct a benchmark that is both diverse and difficult.

In work which appeared after the initial launch of GLUE, McCann et al. (2018) introduce de-caNLP, which also scores NLP systems based on their performance on multiple datasets. Their benchmark recasts the ten evaluation tasks as question answering, converting tasks like summarization and text-to-SQL semantic parsing into question answering using automatic transformations. That benchmark lacks the leaderboard and error analysis toolkit of GLUE, but more importantly, we see it as pursuing a more ambitious but less immediately practical goal: While GLUE rewards methods that yield good performance on a circumscribed set of tasks using methods like those that are currently used for those tasks, their benchmark rewards systems that make progress toward their goal of unifying all of NLU under the rubric of question answering.

## 3 Tasks

GLUE is centered on nine English sentence understanding tasks, which cover a broad range of domains, data quantities, and difficulties. As the goal of GLUE is to spur development of generalizable NLU systems, we design the benchmark such that good performance should require a model to share substantial knowledge (e.g., trained parameters) across all tasks, while still maintaining some task-specific components. Though it is possible to train a single model for each task and evaluate the resulting set of models on this benchmark, we expect that our inclusion of several data-scarce

| Corpus | \|Train\| | \|Dev\| | \|Test\| | Task | Metrics | Domain |
|--------|-------|------|-------|------|---------|--------|
| | | | | Single-Sentence Tasks | | |
| CoLA | 8.5k | 1k | **1k** | acceptability | Matthews corr. | misc. |
| SST-2 | 67k | 872 | 1.8k | sentiment | acc. | movie reviews |
| | | | | Similarity and Paraphrase Tasks | | |
| MRPC | 3.7k | 408 | 1.7k | paraphrase | acc./F1 | news |
| STS-B | 7k | 1.5k | 1.4k | sentence similarity | Pearson/Spearman corr. | misc. |
| QQP | 364k | 40k | **391k** | paraphrase | acc./F1 | social QA questions |
| | | | | Inference Tasks | | |
| MNLI | 393k | 20k | **20k** | NLI | matched acc./mismatched acc. | misc. |
| QNLI | 108k | 5.7k | 5.7k | QA/NLI | acc. | Wikipedia |
| RTE | 2.5k | 276 | 3k | NLI | acc. | misc. |
| WNLI | 634 | 71 | **146** | coreference/NLI | acc. | fiction books |

Table 1: Task descriptions and statistics. All tasks are single sentence or sentence pair classification, except STS-B, which is a regression task. MNLI has three classes; all other classification tasks have two. Test sets shown in bold use labels that have never been made public in any form.

tasks will ultimately render this approach uncompetitive. We describe the tasks below and in Table 1. Appendix A includes additional details. Unless otherwise mentioned, tasks are evaluated on accuracy and are balanced across classes.

## 3.1 Single-Sentence Tasks

**CoLA** The Corpus of Linguistic Acceptability[1] consists of English acceptability judgments drawn from books and journal articles on linguistic theory. Each example is a sequence of words annotated with whether it is a grammatical English sentence. Judgments of this particular kind are the primary form of evidence in syntactic theory (Schütze, 1996), so a machine learning system capable of predicting them reliably would offer potentially substantial evidence on questions of language learnability and innate bias. Following the authors, we use the Matthews correlation coefficient (Matthews, 1975) as the evaluation metric, which evaluates classifiers on unbalanced binary classification and ranges from -1 to 1, with 0 being the performance of uninformed guessing. We use the standard test set, for which we obtained private labels from the authors. We report a single performance number on the combination of the in- and out-of-domain sections of the test set.

**SST-2** The Stanford Sentiment Treebank (Socher et al., 2013) consists of sentences extracted from movie reviews and human annotations of their sentiment. Given a sentence, the task is to determine the sentiment of the sentence.

We use the two-way (positive/negative) class split, and use only sentence-level labels.

## 3.2 Similarity and Paraphrase Tasks

**MRPC** The Microsoft Research Paraphrase Corpus (Dolan and Brockett, 2005) is a corpus of sentence pairs automatically extracted from online news sources, with human annotations for whether the sentences in the pair are semantically equivalent. Because the classes are imbalanced (68% positive, 32% negative), we follow common practice and report both accuracy and F1 score.

**QQP** The Quora Question Pairs[2] dataset is a collection of question pairs from the community question-answering website Quora. Given two questions, the task is to determine whether they are semantically equivalent. As in MRPC, the class distribution in QQP is unbalanced (37% positive, 63% negative), so we report both accuracy and F1 score. We use the standard test set, for which we obtained private labels from the authors.

**STS-B** The Semantic Textual Similarity Benchmark (Cer et al., 2017) is a collection of sentence pairs drawn from news headlines, video and image captions, and natural language inference data. Each pair is human-annotated with a similarity score from 1 to 5; the task is to predict these scores. Follow common practice, we evaluate using Pearson and Spearman correlation coefficients.

---

[1] Available at: nyu-mll.github.io/CoLA

[2] data.quora.com/First-Quora-Dataset-Release-Question-Pairs

### 3.3 Inference Tasks

**MNLI** The Multi-Genre Natural Language Inference Corpus (Williams et al., 2018) is a crowd-sourced collection of sentence pairs with textual entailment annotations. Given a premise sentence and a hypothesis sentence, the task is to predict whether the premise entails the hypothesis (*entailment*), contradicts the hypothesis (*contradiction*), or neither (*neutral*). The premise sentences are gathered from ten different domains of text, including transcribed speech, fiction, and government reports. We use the standard test set, for which we obtained private labels from the authors, and evaluate on both the *matched* (in-domain) and *mismatched* (cross-domain) sections. We also use and recommend the SNLI corpus (Bowman et al., 2015) as 550k examples of auxiliary training data.

**QNLI** The Stanford Question Answering Dataset (Rajpurkar et al. 2016) is a question-answering dataset consisting of question-paragraph pairs, where one of the sentences in the paragraph (drawn from Wikipedia) contains the answer to the corresponding question (written by an annotator). We convert the task into sentence pair classification by forming a pair between each question and each sentence in the corresponding context, and filtering out pairs with low lexical overlap between the question and the context sentence. The task is to determine whether the context sentence contains the answer to the question. This modified version of the original task removes the requirement that the model select the exact answer, but also removes the simplifying assumptions that the answer is always present in the input and that lexical overlap is a reliable cue. This process of recasting existing datasets into NLI is similar to methods introduced in White et al. (2017). We call the converted dataset QNLI (Question-answering NLI).

**RTE** The Recognizing Textual Entailment (RTE) datasets come from a series of annual challenges on the task of textual entailment. We combine the data from RTE1 (Dagan et al., 2006), RTE2 (Bar Haim et al., 2006), RTE3 (Giampiccolo et al., 2007), and RTE5 (Bentivogli et al., 2009).[3] Examples are constructed based on news and Wikipedia text. We convert all datasets to a two-class split, where for three-class datasets we collapse *neutral* and *contradiction* into *not_entailment*, for consistency.

**WNLI** The Winograd Schema Challenge (Levesque et al., 2011) is a reading comprehension task in which a system must read a sentence with a pronoun and select the referent of that pronoun from a list of choices. The examples are manually constructed to foil simple statistical methods: Each one is contingent on contextual information provided by a single word or phrase in the sentence. To convert the problem into sentence pair classification, we construct sentence pairs by replacing the ambiguous pronoun with each possible referent. The task is to predict if the sentence with the pronoun substituted is entailed by the original sentence. We use a small evaluation set consisting of new examples derived from fiction books[4] that was shared privately by the authors of the original corpus. While the included training set is balanced between two classes, the test set is imbalanced between them (35% entailment, 65% not entailment). As with QNLI, each example is evaluated separately, so there is not a systematic correspondence between a model's score on this task and its score on the unconverted original task. We call converted dataset WNLI (Winograd NLI).

### 3.4 Evaluation

The GLUE benchmark follows the same evaluation model as SemEval and Kaggle. To evaluate a system on the benchmark, one must run the system on the provided test data for the tasks, then upload the results to the website for scoring. The benchmark site then shows per-task scores, as well as a macro-average of those scores to determine a system's position on the leaderboard. For tasks with multiple metrics (e.g., accuracy and F1), we use an unweighted average of the metrics as the score for the task when computing the overall macro-average. The website also provides fine- and coarse-grained results on the diagnostic dataset. See Appendix C for details.

### 3.5 Data and Bias

We do not endorse the use of the task training sets for any specific *non-research* use. They do not cover every dialect of English one may wish to handle, nor languages other than English. As all

---

[3]RTE4 is not publicly available, while RTE6 and RTE7 do not fit the standard NLI task.

[4]See similar examples at `cs.nyu.edu/faculty/davise/papers/WinogradSchemas/WS.html`

| Tags | Sentence 1 | Sentence 2 | Fwd | Bwd |
|---|---|---|---|---|
| *Lexical Entailment (Lexical Semantics), Downward Monotone (Logic)* | The timing of the meeting has not been set, according to a Starbucks spokesperson. | The timing of the meeting has not been considered, according to a Starbucks spokesperson. | N | E |
| *Universal Quantifiers (Logic)* | Our deepest sympathies are with all those affected by this accident. | Our deepest sympathies are with a victim who was affected by this accident. | E | N |
| *Quantifiers (Lexical Semantics), Double Negation (Logic)* | I have never seen a hummingbird not flying. | I have never seen a hummingbird. | N | E |

Table 2: Examples from the diagnostic set. *Fwd* denotes the label when sentence 1 is the premise; *Bwd* is the label when sentence 2 is the premise. Labels are *entailment* (E), *neutral* (N), or *contradiction* (C). Examples are tagged with the phenomena they demonstrate, and each phenomenon belongs to one of four broad categories (in parentheses). See Table 5 in Appendix A for a complete tag taxonomy.

of them contain text or annotations that were collected in uncontrolled settings, they contain evidence of stereotypes and biases that one may not wish one's system to learn (Rudinger et al., 2017).

## 4   Diagnostic Dataset

Drawing inspiration from the FraCaS suite (Cooper et al., 1996) and the recent Build-It-Break-It competition (Ettinger et al., 2017), we include a small, manually-curated test set (with private labels) for the analysis of system performance. While the main benchmark mostly reflects an application-driven distribution of examples, our diagnostic dataset highlights a pre-defined set of modeling-relevant phenomena.

Each example in the diagnostic dataset is an NLI sentence pair with fine-grained tags for the phenomena it demonstrates. The NLI task is well-suited to this kind of analysis, as it can straightforwardly evaluate the full set of skills involved in (ungrounded) sentence understanding, from the resolution of syntactic ambiguity to pragmatic reasoning with world knowledge. We ensure that the data is reasonably diverse by producing examples for a wide variety of linguistic phenomena, and basing our examples on naturally-occurring sentences from several domains. This approaches differs from that of FraCaS, which was designed to test linguistic theories with a minimal and uniform set of examples. A sample from our dataset is shown in Table 2, and a full list of linguistic categories is in Table 5 in the appendix.

**Domains**   We construct sentence pairs based on text from four domains: News (articles linked from the front page), Reddit (threads linked from the Front Page), Wikipedia (Featured Articles), and academic papers from recent ACL conferences. We include 100 sentence pairs constructed from each source and 150 artificially-constructed sentence pairs for 550 total.

**Annotation Process**   We begin with a target set of phenomena, based roughly on those used in the FraCaS suite (Cooper et al., 1996). We construct each example by locating a sentence that can be easily made to demonstrate a target phenomenon, and editing it in two ways to produce an appropriate sentence pair. We make minimal modifications so as to maintain high lexical and structural overlap within each sentence pair and limit superficial cues. We then label the inference relationships between the sentences, considering each sentence alternatively as the premise, producing two labeled examples for each pair (1100 total). Where possible, we produce several pairs with different labels for a single source sentence, to have minimal sets of sentence pairs that are lexically and structurally very similar but correspond to different entailment relationships. The resulting labels are 42% *entailment*, 35% *neutral*, and 23% *contradiction*.

**Evaluation**   Since the class distribution in the diagnostic set is not balanced, we use $R_3$ (Gorodkin, 2004), a three-class generalization of the Matthews correlation coefficient, for evaluation.

In light of recent work showing that crowd-sourced data often contains artifacts which can be exploited to perform well without solving the intended task (Schwartz et al., 2017; Gururangan et al., 2018; Poliak et al., 2018; Tsuchiya, 2018), we audit the data for such artifacts. We reproduce the methodology of Gururangan et al. (2018), training two fastText classifiers (Joulin et al., 2016) to predict entailment labels on SNLI

and MNLI using only the hypothesis as input. Testing the trained classifiers on the diagnostic data, we obtain accuracies close to chance, 32.7% and 36.4% respectively, showing that the data does not suffer from artifacts of this kind.

To establish human baseline performance on the diagnostic set, we have six NLP researchers annotate 50 sentence pairs (100 entailment examples) randomly sampled from the diagnostic set. Inter-annotator agreement is high, with a Fleiss's $\kappa$ of 0.73. The average $R_3$ score among the annotators is 0.80, much higher than any of the baseline systems described in Section 5.

**Intended Use**   Because these analysis examples are hand-picked to address certain phenomena, we expect that they will not be representative of the distribution of language as a whole, even in the targeted domains. However, NLI is a task with no natural input distribution. We deliberately select sentences that we hope will be able to provide insight into what models are doing, what phenomena they catch on to, and where are they limited. This means that the raw performance numbers on the analysis set should be taken with a grain of salt. The set is provided not as a benchmark, but as an analysis tool to paint in broad strokes the kinds of phenomena a model may or may not capture, and to provide a set of examples that can serve for error analysis, qualitative model comparison, and development of adversarial examples that expose a model's weaknesses.

## 5   Baselines

We evaluate a simple multi-task learning model trained on the benchmark tasks, as well as several more sophisticated variants based on recent pre-training methods, as baselines. We briefly describe them here. See Appendix B for details. We implement our models in the AllenNLP library (Gardner et al., 2017).

**Architecture**   Our simplest baseline architecture is based on sentence-to-vector encoders, and sets aside GLUE's ability to evaluate models with more complex structures. Taking inspiration from Conneau et al. (2017), the model uses a two-layer, 1500D (per direction) BiLSTM with max pooling and 300D GloVe word embeddings (840B Common Crawl version; Pennington et al., 2014). For single-sentence tasks, we encode the sentence and pass the resulting vector to a classifier. For sentence-pair tasks, we encode sentences independently to produce vectors $u, v$, and pass $[u; v; |u - v|; u * v]$ to a classifier. The classifier is an MLP with a 512D hidden layer.

We also consider a variant of our model which for sentence pair tasks uses an attention mechanism inspired by Seo et al. (2016) between all pairs of words, followed by a second BiLSTM with max pooling. By explicitly modeling the interaction between sentences, these models fall outside the sentence-to-vector paradigm.

**Pre-Training**   We augment our base model with two recent methods for pre-training: ELMo and CoVe. We use existing trained models for both.

ELMo uses a pair of two-layer neural language models (one forward, one backward) trained on the Billion Word Benchmark (Chelba et al., 2013). Each word is represented by a contextual embedding, produced by taking a linear combination of the corresponding hidden states of each layer of the two models. We follow the authors' recommendations[5] and use ELMo embeddings in place of any other embeddings.

CoVe (McCann et al., 2017) uses a sequence-to-sequence model with a two-layer BiLSTM encoder trained for English-to-German translation. The CoVe vector of a word is the corresponding hidden state of the top-layer LSTM. As in the original work, we concatenate the CoVe vectors to the GloVe word embeddings.

**Training**   We train our models with the BiLSTM sentence encoder and post-attention BiLSTMs shared across tasks, and classifiers trained separately for each task. For each training update, we sample a task to train with a probability proportional to the number of training examples for each task. We train our models with Adam (Kingma and Ba, 2014) with initial learning rate $10^{-3}$ and batch size 128. We use the macro-average score as the validation metric and stop training when the learning rate drops below $10^{-5}$ or performance does not improve after 5 validation checks.

We also train a set of single-task models, which are configured and trained identically, but share no parameters. While this is generally an effective model for the tasks under study, to allow for fair comparisons with the multi-task analogs we do not tune parameter or training settings for each

---

[5] `github.com/allenai/allennlp/blob/master/tutorials/how_to/elmo.md`

| Model | Avg | Single Sentence | | Similarity and Paraphrase | | | Natural Language Inference | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | CoLA | SST-2 | MRPC | QQP | STS-B | MNLI | QNLI | RTE | WNLI |
| Single-Task Training | | | | | | | | | | |
| BiLSTM | 62.0 | 15.7 | 85.9 | 69.3/79.4 | 81.7/61.4 | 66.0/62.8 | 70.3/70.8 | 60.8 | 52.8 | 62.3 |
| +ELMo | 66.2 | **35.0** | 90.2 | 69.0/80.8 | 85.7/65.6 | 64.0/60.2 | 72.9/73.4 | 69.4 | 50.1 | **65.1** |
| +CoVe | 62.4 | 14.5 | 88.5 | 73.4/81.4 | 83.3/59.4 | 67.2/64.1 | 64.5/64.8 | 64.8 | 53.5 | 61.6 |
| +Attn | 60.0 | 15.7 | 85.9 | 68.5/80.3 | 83.5/62.9 | 59.3/55.8 | 74.2/73.8 | 51.9 | 51.9 | 55.5 |
| +Attn, ELMo | 64.8 | **35.0** | 90.2 | 68.8/80.2 | 86.5/66.1 | 55.5/52.5 | 76.9/76.7 | 61.1 | 50.4 | **65.1** |
| +Attn, CoVe | 60.8 | 14.5 | 88.5 | 68.6/79.7 | 84.1/60.1 | 57.2/53.6 | 71.6/71.5 | 53.8 | 52.7 | 64.4 |
| Multi-Task Training | | | | | | | | | | |
| BiLSTM | 63.5 | 24.0 | 85.8 | 71.9/82.1 | 80.2/59.1 | 68.8/67.0 | 65.8/66.0 | 71.1 | 46.8 | 63.7 |
| +ELMo | 64.8 | 27.5 | 89.6 | 76.2/83.5 | 78.5/57.8 | 67.0/65.9 | 67.1/68.0 | 66.7 | 55.7 | 62.3 |
| +CoVe | 62.2 | 16.2 | 84.3 | 71.8/80.0 | 82.0/59.1 | 68.0/67.1 | 65.3/65.9 | 70.4 | 44.2 | **65.1** |
| +Attn | 65.7 | 0.0 | 85.0 | 75.1/**83.7** | 84.3/63.6 | 73.9/71.8 | 72.2/72.1 | 82.1 | **61.7** | 63.7 |
| +Attn, ELMo | **69.0** | 18.9 | **91.6** | **77.3**/83.5 | 85.3/63.3 | 72.8/71.1 | 75.6/75.9 | 81.7 | 61.2 | **65.1** |
| +Attn, CoVe | 64.3 | 19.4 | 83.6 | 75.2/83.0 | 84.9/61.1 | 72.3/71.1 | 69.9/68.7 | 78.9 | 38.3 | **65.1** |
| Pre-Trained Sentence Representation Models | | | | | | | | | | |
| CBoW | 58.9 | 0.0 | 80.0 | 73.4/81.5 | 79.1/51.4 | 61.2/58.7 | 56.0/56.4 | 75.1 | 54.1 | 62.3 |
| Skip-Thought | 61.5 | 0.0 | 81.8 | 71.7/80.8 | 82.2/56.4 | 71.8/69.7 | 62.9/62.8 | 74.7 | 53.1 | **65.1** |
| InferSent | 64.7 | 4.5 | 85.1 | 74.1/81.2 | 81.7/59.1 | 75.9/75.3 | 66.1/65.7 | 79.8 | 58.0 | **65.1** |
| DisSent | 62.1 | 4.9 | 83.7 | 74.1/81.7 | 82.6/59.5 | 66.1/64.8 | 58.7/59.1 | 75.2 | 56.4 | **65.1** |
| GenSen | 66.6 | 7.7 | 83.1 | 76.6/83.0 | 82.9/59.8 | **79.3/79.2** | 71.4/71.3 | **82.3** | 59.2 | **65.1** |

Table 3: Baseline performance on the GLUE tasks. For MNLI, we report accuracy on the matched and mismatched test sets. For MRPC and Quora, we report accuracy and F1. For STS-B, we report Pearson and Spearman correlation. For CoLA, we report Matthews correlation. For all other tasks we report accuracy. All values are scaled by 100. A similar table is presented on the online platform.

task, so these single-task models do not generally represent the state of the art for each task.

**Sentence Representation Models** Finally, we evaluate the following trained sentence-to-vector encoder models using our benchmark: average bag-of-words using GloVe embeddings (CBoW), Skip-Thought (Kiros et al., 2015), InferSent (Conneau et al., 2017), DisSent (Nie et al., 2017), and GenSen (Subramanian et al., 2018). See Appendix B for additional details. For these models, we only train task-specific classifiers on the representations they produce.

## 6 Benchmark Results

We train three runs of each model and evaluate the run with the best macro-average development set performance. For single-task and sentence representation models, we evaluate the best run for each individual task. We present performance on the main benchmark tasks in Table 3.

In most cases, using multi-task training over single-task training yields better overall scores, particularly among the parameter-rich attention models. Attention generally hurts performance in single task training, but helps in multi-task train-

ing. We see a consistent improvement in using ELMo embeddings in place of GloVe or CoVe embeddings, particularly for single-sentence tasks. Using CoVe slightly improves on GloVe for single task training but not for multi-task training.

Among the pre-trained sentence representation models, we observe fairly consistent gains by moving from CBoW to Skip-Thought to Infersent and GenSen. Relative to the models trained directly on the GLUE tasks, InferSent is competitive and GenSen outperforms all but the two best.

Looking at results per task, we find that the sentence representation models substantially underperform on CoLA compared to the models directly trained on the task. Similarly, with the exception of InferSent, the sentence representation models are outperformed on SST by our BiLSTM and its non-CoVe variants. These discrepancies indicate a need for better transfer methods for generalizing outside of the tasks a model was trained on and for task diversity in evaluation methods, as we have sought to do with GLUE. On the other hand, for STS-B, there is a significant gap between the models trained directly on the task and the best sentence representation model, which we interpret as indicating the necessity of using trans-

fer learning methods trained on data outside of the GLUE benchmark in order to solve it. Finally, there are tasks for which no model does particularly well. On WNLI, no model exceeds most-frequent-class guessing (65.1%). On RTE and in aggregate, even our best baselines leave room for improvement. These early results indicate that solving GLUE is beyond the capabilities of current models and methods, and that training on auxiliary tasks seems a necessary and promising direction.

## 7 Analysis

We analyze the baselines by evaluating each model's MNLI classifier on the diagnostic set to get a better sense of their linguistic capabilities. Results are presented in Table 4.

**Coarse Categories** Overall performance is low for all models: The highest total score of 28 still denotes poor absolute performance. Performance tends to be higher on Predicate-Argument Structure and lower on Knowledge, though numbers are not closely comparable across categories. Unlike on the main benchmark, the multi-task models are almost always outperformed by their single-task counterparts. This is perhaps unsurprising, since with our simple multi-task training regime, there is likely some destructive interference between MNLI and the other tasks. The models trained on the GLUE tasks largely outperform the pretrained sentence representation models, with the exception of GenSen. Using attention has a greater influence on diagnostic scores than using ELMo or CoVe, which we take to indicate that attention is especially important for generalization in NLI.

**Fine-Grained Subcategories** Most models handle universal quantification relatively well. Looking at relevant examples, it seems that catching on to lexical cues such as "all" often suffices for good performance. Similarly, lexical cues often provide good signal in examples of morphological negation.

We also observe weaknesses that vary between models. Double negation is especially difficult for the GLUE-trained models that only use GloVe embeddings. This is ameliorated by ELMo, and to some degree CoVe, perhaps because the translation and language modeling objectives teach models that phrases like "not bad" and "okay" have similar distributions. Also, while attention improves overall results, attention models tend to struggle with downward monotonicity. Examining their predictions, we found that the models were sensitive to hypernym/hyponym substitutions as signals of entailment, but predicted it in the wrong direction (as if the substituted word was in an upward monotone context). Restrictivity examples, which often depend on nuances of quantifier scope, are especially difficult for all models.

Overall, there is evidence that going beyond sentence-to-vector representations, e.g. with an attention mechanism, might aid performance on out-of-domain data, and that transfer methods like ELMo and CoVe encode linguistic information specific to their supervision signal. However, increased representational capacity may lead to overfitting, such as the failure of attention models in downward monotone contexts. We expect that our platform and diagnostic dataset will be useful for similar analyses in the future, so that model designers can better understand their models' generalization behavior and implicit knowledge.

## 8 Conclusion

We introduce GLUE, a platform and collection of resources for evaluating and analyzing natural language understanding systems. We find that, in aggregate, models trained jointly on our tasks see better performance than the combined performance of models trained for each task separately. We confirm the utility of attention mechanisms and transfer learning methods such as ELMo in NLU systems, which combine to outperform the best sentence representation models on the GLUE benchmark, but still leave room for improvement. When evaluating these models on our diagnostic dataset, we find that they fail (often spectacularly) on many linguistic phenomena, suggesting possible avenues for future work. In sum, the question of how to design general-purpose NLU models remains unanswered, and we believe that GLUE can provide fertile soil for addressing this challenge.

| Model | Coarse-Grained | | | | | UQuant | MNeg | Fine-Grained 2Neg | Coref | Restr | Down |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | All | LS | PAS | L | K | | | | | | |
| *Single-Task Training* | | | | | | | | | | | |
| BiLSTM | 21 | 25 | 24 | 16 | 16 | 70 | <u>53</u> | 4 | 21 | -15 | **12** |
| +ELMo | 20 | 20 | 21 | 14 | 17 | 70 | 20 | **42** | 33 | -26 | -3 |
| +CoVe | 21 | 19 | 23 | 20 | <u>18</u> | 71 | 47 | -1 | 33 | -15 | 8 |
| +Attn | 25 | 24 | 30 | 20 | 14 | 50 | 47 | 21 | **38** | -8 | -3 |
| +Attn, ELMo | **28** | **30** | **35** | **23** | 14 | **85** | 20 | **42** | 33 | -26 | -3 |
| +Attn, CoVe | 24 | 29 | 29 | 18 | 12 | 77 | 50 | 1 | 18 | <u>-1</u> | **12** |
| *Multi-Task Training* | | | | | | | | | | | |
| BiLSTM | 19 | 16 | 22 | 16 | 17 | 71 | 35 | -8 | 26 | **<u>0</u>** | 8 |
| +ELMo | 19 | 15 | 21 | 17 | **21** | 70 | **60** | 15 | 26 | **<u>0</u>** | **12** |
| +CoVe | 17 | 15 | 21 | 14 | 16 | 50 | 31 | -8 | 25 | -15 | **12** |
| +Attn | <u>25</u> | 23 | <u>32</u> | <u>19</u> | 16 | 58 | 26 | -5 | 28 | -1 | -20 |
| +Attn, ELMo | 23 | <u>24</u> | 30 | 17 | 13 | <u>78</u> | 27 | <u>37</u> | 30 | -15 | -20 |
| +Attn, CoVe | 20 | 16 | 25 | 15 | 17 | <u>78</u> | 37 | 14 | <u>31</u> | -15 | 8 |
| *Pre-Trained Sentence Representation Models* | | | | | | | | | | | |
| CBoW | 9 | 6 | 13 | 5 | 10 | 3 | 0 | <u>13</u> | 28 | <u>-15</u> | -11 |
| Skip-Thought | 12 | 2 | 23 | 11 | 9 | 61 | 6 | -2 | <u>30</u> | <u>-15</u> | 0 |
| InferSent | 18 | 20 | 20 | <u>15</u> | 14 | 77 | 50 | -20 | 15 | <u>-15</u> | -9 |
| DisSent | 16 | 16 | 19 | 13 | <u>15</u> | 70 | 43 | -11 | 20 | -36 | -09 |
| GenSen | <u>20</u> | <u>28</u> | <u>26</u> | 14 | 12 | <u>78</u> | <u>57</u> | 2 | 21 | <u>-15</u> | **12** |

Table 4: Results on the diagnostic set. We report $R_3$ coefficients between gold and predicted labels, scaled by 100. The coarse-grained categories (left) are *Lexical Semantics* (**LS**), *Predicate-Argument Structure* (**PAS**), *Logic* (**L**), and *Knowledge and Common Sense* (**K**). Our example fine-grained categories (right) are *Universal Quantification* (**UQuant**), *Morphological Negation* (**MNeg**), *Double Negation* (**2Neg**), *Anaphora/Coreference* (**Coref**), *Restrictivity* (**Restr**), and *Downward Monotone* (**Down**).

from AdeptMind.

# References

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. In *ICLR 2015*.

Roy Bar Haim, Ido Dagan, Bill Dolan, Lisa Ferro, Danilo Giampiccolo, Bernardo Magnini, and Idan Szpektor. 2006. The second PASCAL recognising textual entailment challenge.

Luisa Bentivogli, Ido Dagan, Hoa Trang Dang, Danilo Giampiccolo, and Bernardo Magnini. 2009. The fifth PASCAL recognizing textual entailment challenge.

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. Association for Computational Linguistics.

Daniel Cer, Mona Diab, Eneko Agirre, Inigo Lopez-Gazpio, and Lucia Specia. 2017. Semeval-2017 task 1: Semantic textual similarity-multilingual and cross-lingual focused evaluation. In *11th International Workshop on Semantic Evaluations*.

Ciprian Chelba, Tomas Mikolov, Mike Schuster, Qi Ge, Thorsten Brants, Phillipp Koehn, and Tony Robinson. 2013. One billion word benchmark for measuring progress in statistical language modeling. *arXiv preprint 1312.3005*.

Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel Kuksa. 2011. Natural language processing (almost) from scratch. *Journal of Machine Learning Research*, 12(Aug):2493–2537.

Alexis Conneau and Douwe Kiela. 2018. SentEval: An evaluation toolkit for universal sentence representations. In *LREC 2018*.

Alexis Conneau, Douwe Kiela, Holger Schwenk, Loïc Barrault, and Antoine Bordes. 2017. Supervised learning of universal sentence representations from natural language inference data. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017, Copenhagen, Denmark, September 9-11, 2017*, pages 681–691.

Robin Cooper, Dick Crouch, Jan Van Eijck, Chris Fox, Josef Van Genabith, Jan Jaspars, Hans Kamp, David Milward, Manfred Pinkal, Massimo Poesio, Steve Pulman, Ted Briscoe, Holger Maier, and Karsten Konrad. 1996. Using the framework. Technical report, The FraCaS Consortium.

Ido Dagan, Oren Glickman, and Bernardo Magnini. 2006. The PASCAL recognising textual entailment challenge. In *Machine learning challenges. evaluating predictive uncertainty, visual object classification, and recognising tectual entailment*, pages 177–190. Springer.

William B Dolan and Chris Brockett. 2005. Automatically constructing a corpus of sentential paraphrases. In *Proceedings of IWP*.

Allyson Ettinger, Sudha Rao, Hal Daumé III, and Emily M Bender. 2017. Towards linguistically generalizable NLP systems: A workshop and shared task. In *First Workshop on Building Linguistically Generalizable NLP Systems*.

Matt Gardner, Joel Grus, Mark Neumann, Oyvind Tafjord, Pradeep Dasigi, Nelson F. Liu, Matthew Peters, Michael Schmitz, and Luke S. Zettlemoyer. 2017. Allennlp: A deep semantic natural language processing platform.

Danilo Giampiccolo, Bernardo Magnini, Ido Dagan, and Bill Dolan. 2007. The third PASCAL recognizing textual entailment challenge. In *Proceedings of the ACL-PASCAL workshop on textual entailment and paraphrasing*, pages 1–9. Association for Computational Linguistics.

Jan Gorodkin. 2004. Comparing two k-category assignments by a k-category correlation coefficient. *Comput. Biol. Chem.*, 28(5-6):367–374.

Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel R. Bowman, and Noah A. Smith. 2018. Annotation artifacts in natural language inference data. In *Proceedings of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.

Kazuma Hashimoto, Caiming Xiong, Yoshimasa Tsuruoka, and Richard Socher. 2016. A joint many-task model: Growing a neural network for multiple nlp tasks. In *Proceedings of EMNLP 2017*.

Felix Hill, Kyunghyun Cho, and Anna Korhonen. 2016. Learning distributed representations of sentences from unlabelled data. *Proceedings of NAACL 2016*.

Minqing Hu and Bing Liu. 2004. Mining and summarizing customer reviews. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 168–177. ACM.

Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2016. Bag of tricks for efficient text classification. *arXiv preprint 1607.01759*.

Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. In *ICLR 2015*.

Ryan Kiros, Yukun Zhu, Ruslan R Salakhutdinov, Richard Zemel, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Skip-Thought vectors. In *Advances in neural information processing systems*, pages 3294–3302.

Quoc Le and Tomas Mikolov. 2014. Distributed representations of sentences and documents. In *Proceedings of the 31st International Conference on Machine Learning*, volume 32 of *Proceedings of Machine Learning Research*, pages 1188–1196, Bejing, China. PMLR.

Hector J Levesque, Ernest Davis, and Leora Morgenstern. 2011. The Winograd schema challenge. In *Aaai spring symposium: Logical formalizations of commonsense reasoning*, volume 46, page 47.

Brian W Matthews. 1975. Comparison of the predicted and observed secondary structure of t4 phage lysozyme. *Biochimica et Biophysica Acta (BBA)-Protein Structure*, 405(2):442–451.

Bryan McCann, James Bradbury, Caiming Xiong, and Richard Socher. 2017. Learned in translation: Contextualized word vectors. In *Advances in Neural Information Processing Systems*, pages 6297–6308.

Bryan McCann, Nitish Shirish Keskar, Caiming Xiong, and Richard Socher. 2018. The natural language decathlon: Multitask learning as question answering. *arXiv preprint 1806.08730*.

Allen Nie, Erin D Bennett, and Noah D Goodman. 2017. Dissent: Sentence representation learning from explicit discourse relations. *arXiv preprint 1710.04334*.

Bo Pang and Lillian Lee. 2004. A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts. In *Proceedings of the 42nd annual meeting on Association for Computational Linguistics*, page 271. Association for Computational Linguistics.

Bo Pang and Lillian Lee. 2005. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. In *Proceedings of the 43rd annual meeting on association for computational linguistics*, pages 115–124. Association for Computational Linguistics.

Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. GloVe: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543.

Matthew E Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proceedings of NAACL 2018*.

Adam Poliak, Jason Naradowsky, Aparajita Haldar, Rachel Rudinger, and Benjamin Van Durme. 2018. Hypothesis only baselines in natural language inference. *arXiv preprint:1805.01042*.

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. Association for Computational Linguistics.

Tim Rocktäschel, Edward Grefenstette, Moritz Hermann, Karl, Tomáš Kočiský, and Phil Blunsom. 2016. Reasoning about entailment with neural attention. In *Proceedings of the International Conference on Learning Representations (ICLR)*.

Sebastian Ruder, Joachim Bingel, Isabelle Augenstein, and Anders Søgaard. 2017. Sluice networks: Learning what to share between loosely related tasks. *arXiv preprint 1705.08142*.

Rachel Rudinger, Chandler May, and Benjamin Van Durme. 2017. Social bias in elicited natural language inferences. In *Proceedings of the First ACL Workshop on Ethics in Natural Language Processing*, pages 74–79.

Carson T Schütze. 1996. *The empirical base of linguistics: Grammaticality judgments and linguistic methodology*. University of Chicago Press.

Roy Schwartz, Maarten Sap, Ioannis Konstas, Li Zilles, Yejin Choi, and Noah A. Smith. 2017. The effect of different writing tasks on linguistic style: A case study of the ROC story cloze task. In *Proceedings of CoNLL*.

Minjoon Seo, Aniruddha Kembhavi, Ali Farhadi, and Hannaneh Hajishirzi. 2016. Bidirectional attention flow for machine comprehension. In *ICLR 2017*.

Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages 1631–1642.

Anders Søgaard and Yoav Goldberg. 2016. Deep multi-task learning with low level tasks supervised at lower layers. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, volume 2, pages 231–235.

Sandeep Subramanian, Adam Trischler, Yoshua Bengio, and Christopher J. Pal. 2018. Learning general purpose distributed sentence representations via large scale multi-task learning. In *Proceedings of ICLR*.

Masatoshi Tsuchiya. 2018. Performance Impact Caused by Hidden Bias of Training Data for Recognizing Textual Entailment. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 6000–6010.

Ellen M Voorhees et al. 1999. The trec-8 question answering track report. In *Trec*, volume 99, pages 77–82.

Aaron Steven White, Pushpendre Rastogi, Kevin Duh, and Benjamin Van Durme. 2017. Inference is everything: Recasting semantic resources into a unified evaluation framework. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, volume 1, pages 996–1005.

Janyce Wiebe, Theresa Wilson, and Claire Cardie. 2005. Annotating expressions of opinions and emotions in language. *Language resources and evaluation*, 39(2-3):165–210.

Adina Williams, Nikita Nangia, and Samuel R. Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In *Proceedings of NAACL 2018*.

Caiming Xiong, Victor Zhong, and Richard Socher. 2016. Dynamic coattention networks for question answering. In *ICLR 2017*.

Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *Proceedings of the IEEE international conference on computer vision*, pages 19–27.

# A  Additional Data Details

## A.1  Dataset Construction

**QNLI**  To construct a balanced dataset, we select all pairs in which the most similar sentence to the question was *not* the answer sentence, as well as an equal amount of cases in which the correct sentence was the most similar to the question, but another distracting sentence was a close second. Our similarity metric is based on CBoW representations with pre-trained GloVe embeddings. This approach to converting pre-existing datasets into NLI format is closely related to recent work by White et al. (2017), as well as to the original motivation for textual entailment presented by Dagan et al. (2006). Both argue that many NLP tasks can be productively reduced to textual entailment.

## A.2  Diagnostic Data

We show the full label set used to tag the diagnostic set in Table 5.

| Coarse-Grained Categories | Fine-Grained Categories |
|---|---|
| Lexical Semantics | Lexical Entailment, Morphological Negation, Factivity, Symmetry/Collectivity, Redundancy, Named Entities, Quantifiers |
| Predicate-Argument Structure | Core Arguments, Prepositional Phrases, Ellipsis/Implicits, Anaphora/Coreference Active/Passive, Nominalization, Genitives/Partitives, Datives, Relative Clauses, Coordination Scope, Intersectivity, Restrictivity |
| Logic | Negation, Double Negation, Intervals/Numbers, Conjunction, Disjunction, Conditionals, Universal, Existential, Temporal, Upward Monotone, Downward Monotone, Non-Monotone |
| Knowledge | Common Sense, World Knowledge |

Table 5: The types of linguistic phenomena annotated in the diagnostic dataset, organized under four major categories.

## B  Additional Baseline Details

### B.1  Attention Mechanism

We implement our attention mechanism as follows: given two sequences of hidden states $u_1, u_2, \ldots, u_M$ and $v_1, v_2, \ldots, v_N$, we first compute matrix $H$ where $H_{ij} = u_i \cdot v_j$. For each $u_i$, we get attention weights $\alpha_i$ by taking a softmax over the $i^{th}$ row of $H$, and get the corresponding context vector $\tilde{v}_i = \sum_j \alpha_{ij} v_j$ by taking the attention-weighted sum of the $v_j$. We pass a second BiLSTM with max pooling over the sequence $[u_1; \tilde{v}_1], \ldots [u_M; \tilde{v}_M]$ to produce $u'$. We process the $v_j$ vectors analogously to obtain $v'$. Finally, we feed $[u'; v'; |u' - v'|; u' * v']$ into a classifier.

### B.2  Training

We train our models with the BiLSTM sentence encoder and post-attention BiLSTMs shared across tasks, and classifiers trained separately for each task. For each training update, we sample a task to train with a probability proportional to the number of training examples for each task. We scale each task's loss inversely proportional to the number of examples for that task, which we found to improve overall performance. We train our models with Adam (Kingma and Ba, 2014) with initial learning rate $10^{-3}$, batch size 128, and gradient clipping. We use macro-average score over all tasks as our validation metric, and perform a validation check every 10k updates. We divide the learning rate by 5 whenever validation performance does not improve. We stop training when the learning rate drops below $10^{-5}$ or performance does not improve after 5 validation checks.

### B.3  Sentence Representation Models

We evaluate the following sentence representation models:

1. CBoW, the average of the GloVe embeddings of the tokens in the sentence.

2. Skip-Thought (Kiros et al., 2015), a sequence-to-sequence(s) model trained to generate the previous and next sentences given the middle sentence. We use the original pre-trained model[6] trained on sequences of sentences from the Toronto Book Corpus (Zhu et al. 2015, TBC).

3. InferSent (Conneau et al., 2017), a BiLSTM with max-pooling trained on MNLI and SNLI.

4. DisSent (Nie et al., 2017), a BiLSTM with max-pooling trained to predict the discourse marker (*because*, *so*, etc.) relating two sentences on data derived from TBC. We use the variant trained for eight-way classification.

5. GenSen (Subramanian et al., 2018), a sequence-to-sequence model trained on a variety of supervised and unsupervised objectives. We use the variant of the model trained on both MNLI and SNLI, the Skip-Thought objective on TBC, and a constituency parsing objective on the Billion Word Benchmark.

## C  Benchmark Website Details

GLUE's online platform is built using React, Redux and TypeScript. We use Google Firebase for data storage and Google Cloud Functions to host and run our grading script when a submission is made. Figure 1 shows the visual presentation of our baselines on the leaderboard.
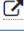
---

[6] github.com/ryankiros/skip-thoughts

| | PRIMARY | | | AUXILIARY | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| RankName | Model | URL | Score | CoLA | SST-2 | MRPC | STS-B | QQP | MNLI-m | MNLI-mm | QNLI | RTE | WNLI |
| 1   GLUE Baselines | BiLSTM+ELMo+Attn | ⬀ | 68.9 | 18.9 | 91.6 | 77.3/83.5 | 72.8/71.1 | 83.5/63.3 | 75.6 | 75.9 | 81.7 | 61.2 | 65.1 |
| | GenSen | ⬀ | 66.6 | 7.7 | 83.1 | 76.6/83.0 | 79.3/79.2 | 82.9/59.8 | 71.4 | 71.3 | 82.3 | 59.2 | 65.1 |
| | Single Task BiLSTM+ELMo | ⬀ | 66.2 | 35.0 | 90.2 | 69.0/80.8 | 64.0/60.2 | 85.7/65.6 | 72.9 | 73.4 | 69.4 | 50.1 | 65.1 |
| | BiLSTM+Attn | | 65.7 | 0.0 | 85.0 | 75.1/83.7 | 73.9/71.8 | 84.3/63.6 | 72.2 | 72.1 | 82.1 | 61.7 | 63.7 |
| | BiLSTM+ELMo | ⬀ | 64.9 | 27.5 | 89.6 | 76.2/83.5 | 67.0/65.9 | 78.5/57.8 | 67.1 | 68.0 | 66.7 | 55.7 | 62.3 |
| | Single Task BiLSTM+ELMo+Attn | ⬀ | 64.8 | 35.0 | 90.2 | 68.8/80.2 | 55.5/52.5 | 86.5/66.1 | 76.9 | 76.7 | 61.1 | 50.3 | 65.1 |
| | InferSent | ⬀ | 64.7 | 4.5 | 85.1 | 74.1/81.2 | 75.9/75.3 | 81.7/59.1 | 66.1 | 65.7 | 79.8 | 58.0 | 65.1 |
| | BiLSTM+CoVe+Attn | ⬀ | 64.3 | 19.4 | 83.6 | 75.2/83.0 | 72.3/71.1 | 84.9/61.1 | 69.9 | 68.7 | 78.9 | 38.3 | 65.1 |
| | BiLSTM | | 63.5 | 24.0 | 85.8 | 71.9/82.1 | 68.8/67.0 | 80.2/59.1 | 65.8 | 66.0 | 71.1 | 46.8 | 63.7 |
| | Single Task BiLSTM+CoVe | ⬀ | 62.4 | 14.5 | 88.5 | 73.4/81.4 | 67.2/64.1 | 83.3/59.4 | 64.5 | 64.8 | 64.8 | 53.5 | 61.6 |
| | BiLSTM+CoVe | ⬀ | 62.2 | 16.2 | 84.3 | 71.8/80.0 | 68.0/67.1 | 82.0/59.1 | 65.3 | 65.9 | 70.4 | 44.2 | 65.1 |
| | DisSent | ⬀ | 62.1 | 4.9 | 83.7 | 74.1/81.7 | 66.1/64.8 | 82.6/59.5 | 58.7 | 59.1 | 75.2 | 56.4 | 65.1 |
| | Single Task BiLSTM | | 62.0 | 15.7 | 85.9 | 69.3/79.4 | 66.0/62.8 | 81.7/61.4 | 70.3 | 70.8 | 60.8 | 52.8 | 62.3 |
| | Skip-Thought | ⬀ | 61.5 | 0.0 | 81.8 | 71.7/80.8 | 71.8/69.7 | 82.2/56.4 | 62.9 | 62.8 | 74.7 | 53.1 | 65.1 |
| | Single Task BiLSTM+CoVe+Attn | ⬀ | 60.8 | 14.5 | 88.5 | 68.6/79.7 | 57.2/53.6 | 84.1/60.1 | 71.6 | 71.5 | 53.8 | 52.7 | 64.4 |
| | Single Task BiLSTM+Attn | | 60.0 | 15.7 | 85.9 | 68.5/80.3 | 59.3/55.8 | 83.5/62.9 | 74.2 | 73.8 | 51.9 | 51.9 | 55.5 |
| | CBOW | | 58.9 | 0.0 | 80.0 | 73.4/81.5 | 61.2/58.7 | 79.1/51.4 | 56.0 | 56.4 | 75.1 | 54.1 | 62.3 |

Figure 1: The benchmark website leaderboard. An expanded view shows additional details about each submission, including a brief prose description and parameter count.