

Evaluating Commonsense in Pre-Trained Language Models

Xuhui Zhou,^{1*} Yue Zhang,² Leyang Cui,^{2,3} Dandan Huang²

¹University of Washington

²School of Engineering, Westlake University

³Zhejiang University

xuhuizh@uw.edu, {yue.zhang, cuileyang, huangdandan}@westlake.edu.cn

Abstract

Contextualized representations trained over large raw text data have given remarkable improvements for NLP tasks including question answering and reading comprehension. There have been works showing that syntactic, semantic and word sense knowledge are contained in such representations, which explains why they benefit such tasks. However, relatively little work has been done investigating commonsense knowledge contained in contextualized representations, which is crucial for human question answering and reading comprehension. We study the commonsense ability of GPT, BERT, XLNet, and RoBERTa by testing them on seven challenging benchmarks, finding that language modeling and its variants are effective objectives for promoting models' commonsense ability while bi-directional context and larger training set are bonuses. We additionally find that current models do poorly on tasks require more necessary inference steps. Finally, we test the robustness of models by making dual test cases, which are correlated so that the correct prediction of one sample should lead to correct prediction of the other. Interestingly, the models show confusion on these test cases, which suggests that they learn commonsense at the surface rather than the deep level. We release a test set, named CATs publicly, for future research.

Introduction

Contextualized representations trained over large-scale text data have given remarkable improvements to a wide range of NLP tasks, including natural language inference (Bowman et al. 2015), question answering (Rajpurkar, Jia, and Liang 2018) and reading comprehension (Lai et al. 2017). Giving new state-of-the-art results that approach or surpass human performance on several benchmark datasets, it is an interesting question what types of knowledge are learned in pre-trained contextualized representations in order to better understand how they benefit the NLP problems above. There has been work investigating the nature of syntactic (Liu et al. 2019a), semantic (Liu et al. 2019a) and word sense (Kim et al. 2019) knowledge contained in such contextualized representations, in particular BERT (Devlin et al. 2019), showing

that such knowledge can be effectively learned via language model (LM) pre-training over large scale data.

Commonsense knowledge spans “a huge portion of human experience, encompassing knowledge about the spatial, physical, social, temporal, and psychological aspects of typical everyday life.” (Liu and Singh 2004). Intuitively, such knowledge is at least as useful as semantic and syntactic knowledge in natural language inference, reading comprehension and coreference resolution. For example, the word “it” in the sentence “the dog cannot cross the street because it is too X” can refer to three different entities when the word “X” is “timid”, “wide” and “dark”, respectively, and resolving such ambiguity can require that a system has relevant commonsense knowledge beyond the sentence level. However, relatively little work has been conducted on systematically evaluating the nature of commonsense knowledge learned in contextualized representations.

We fill this gap by evaluating five state-of-the-art contextualized embedding models on seven commonsense benchmarks. The models include off-the-shelf embeddings¹ from GPT (Radford and Sutskever 2018), GPT2 (Radford et al. 2019), BERT (Devlin et al. 2019), XLNet (Yang et al. 2019) and RoBERTa (Liu et al. 2019b), and the benchmarks include Conjunction Acceptability, Sense Making (Wang et al. 2019), Winograd Schema Challenge (Levesque, Davis, and Morgenstern 2012), SWAG (Zellers et al. 2018), HellaSwag (Zellers et al. 2019), Sense Making with Reasoning (Wang et al. 2019), and Argument Reasoning Comprehension (Habernal et al. 2018). We evaluate commonsense knowledge contained in the above models by unifying the form of all the datasets and comparing LM perplexities on positive and negative samples (i.e., sentences that make sense and those that do not make sense, respectively). Commonsense contained in our data covers a wide range of subjects, from physical world knowledge to social conventions, from scientific domains to daily life scenes. We further categorize them by the difficulty level, namely the number of inference steps necessary in making sense.

We reframe the datasets in order to conduct both word- and sentence-level testing. For word-level testing, negative samples are drawn by replacing words from positive sam-

*Work done while at Westlake University
Copyright © 2020, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

¹<https://github.com/huggingface/transformers>

Token-level		
CA	They broadcast an announcement, but a subway came into the station and I couldn't hear it.	✓
	They broadcast an announcement, before a subway came into the station and I couldn't hear it .	✗
WSC	The trophy doesn't fit into the brown suitcase because the trophy is too large.	✓
	The trophy doesn't fit into the brown suitcase because the suitcase is too large.	✗
SM	money can be used for buying cars	✓
	money can be used for buying stars	✗
Sentence-level		
SMR	¬“ he put an elephant into the fridge” (because) ← an elephant is much bigger than a fridge .	✓
	¬“ he put an elephant into the fridge ” (because) ← elephants are usually gray...	✗
	¬“ he put an elephant into the fridge ” (because) ← an elephant cannot eat a fridge .	✗
SWAG	Someone unlocks the door and they go in. → Someone leads the way in.	✓
	Someone unlocks the door and they go in. → Someone opens the door and walks out.	✗
	Someone unlocks the door and they go in. → Someone walks out of the driveway.	✗
	Someone unlocks the door and they go in. → Someone walks next to someone and sits on a pew.	✗
HellaSwag	A carved pumpkin with a light in it glows on a counter. Supplies for carving are then shown.	
	→ A woman cuts the top off the pumpkin, emptying the seeds.	✓
	→ she cuts down all the pieces and dumps them in a trash bin in the end.	✗
	→ she then carves the traced lines to cut out the design.	✗
ARCT	→ she tapes the top shut as the continue carving the pumpkin.	✗
	People can choose not to use Google ∧ Other search engines don't redirect to Google	
	→ Google is not a harmful monopoly	✓
	People can choose not to use Google ∧ All other search engines redirect to Google	
	→ Google is not a harmful monopoly	✗

Table 1: Example of reframed test instances corresponding to each of our test task. The key word is **bolded** in token-level tasks. \wedge , \neg , \leftarrow and \rightarrow are used for showing the logic flows and replaced by natural language in actual test data.

ples. We are concerned about nouns, verbs, adjectives, adverbs, pronouns and conjunctions, which reflect different aspects of commonsense. For example, while verbs such as “buy, throw, sell ...” are relatively more associated with event knowledge, conjunctions such as “because, but, so ...” are more associated with logical reasoning. For sentence-level testing, negative examples are drawn by replacing a full sub-sentences (such as a clause) with irrelevant or conflicting contents. Sentence-level tests concern more about commonsense inference.

From the results we have four salient observations. First, the pre-trained models give consistently better performances than random baselines, which demonstrates that language model pre-training is useful for learning commonsense knowledge. Second, models based on bi-directional contexts such as BERT, XLNet and RoBERTa are stronger in learning commonsense knowledge compared to those based on uni-directional contexts, such as GPT and GPT2. Third, more commonsense knowledge can be learned from larger training sets, which conforms well to the intuition. Fourth, the models have a certain degree of commonsense reasoning ability. However, as the number of necessary inference steps increase, the model performances drop, which shows that commonsense is still a big challenge that is not completely solved by pre-trained contextualized language models (LMs).

Finally, we further test the robustness of the five models by making dual test samples. Here a dual test sample is built by adding, deleting and replacing words in a test sam-

ple, or swapping two words in the sample, thereby resulting in a closely related test case. In theory, a model equipped with relevant commonsense should give consistent predictions on a pair of dual test cases. However, we find that none of the models are able to reach such consistency. Instead, the models are confused by the modification, tending to give the same predictions over a pair of dual samples despite they may have different gold labels. This further reveals that commonsense contained in the pre-trained models may remain in a surface level, without deep semantic comprehension. We publicly release our datasets, named commonsense ability tests (CATs), and the test script at GitHub.²

Tasks for Evaluating Commonsense

Commonsense ability can be broadly divided to two categories. First, a model with commonsense ability should have basic knowledge about the world, for example, *water always goes down*. Second, it should have the ability to reason over commonsense knowledge, such as *water always goes down because there is gravity on the earth and if you are injured, you should go to the hospital*. To comprehensively test different models' commonsense ability, we synthesize six challenging tasks by taking positive and negative samples from existing benchmarks, and further introduce a new task called Conjunction Acceptability (CA).

We reframe all the tasks into sentence-scoring tasks by substitution or concatenation. For example, we create posi-

²<https://github.com/XuhuiZhou/CATS>

Original:

Paul tried to call George on the phone, but he wasn't successful.

Who is he?

Candidate: A. Paul (correct) B. George

Reframed:

A. Paul tried to call George on the phone, but Paul wasn't successful. (Positive sample)

B. Paul tried to call George on the phone, but George wasn't successful. (Negative sample)

Table 2: Example of reframing a WSC question; Note that there can be additional negative samples.

tive and negative samples by replacing a pronoun in the sentence of a WSC question with the candidates to obtain a test instance as Table 2. A model is asked to score the sentences and we pick the sentence with the highest score as its prediction in a test instance. Below we introduce the data sources and reframed tasks in detail (the correct answer is **bolded**).

Sense Making (SM)

Introduced by Wang et al. (2019), this task tests whether a model can differentiate sense-making and non-sense-making statements. Given a pair of statements (i.e a test instance), it requires the model to choose the more sensible statement. One example is: ***I work 8 hours a day** / I work 25 hours a day.* This task conforms to our evaluation schema without a change. More examples are shown in the SM section of Table 1. The statements typically differ only in one key word which covers nouns, verbs, adjectives, and adverbs.

Winograd Schema Challenge (WSC)

The Winograd Schema Challenge (WSC) dataset (Levesque, Davis, and Morgenstern 2012) consists 273 instances of the pronoun resolution problem. Each instance contains a sentence with a pronoun referring to one of nouns; the original question is to pick the correct noun. For our task, we transform the test as shown in Table 2. More examples are shown in the WSC section of Table 1. WSC is recognized as one of the most difficult commonsense datasets.

Conjunction Acceptability (CA)

As stated by LoBue and Yates (2011), logic-based commonsense knowledge is an important part of world knowledge in addition to content-based knowledge. We aim to probe a model's ability to understand the logic relations in the language by extracting 189 positive samples from the WSC dataset and replacing the conjunction manually with another conjunction to obtain a negative sample. We pair the positive and negative samples to obtain a test instance. For example, *The lawyer asked the witness a question, and the witness was reluctant to answer it / **The lawyer asked the witness a***

question, but the witness was reluctant to answer it. More examples are shown in the CA section of Table 1. This task using "because", "before", "when", "but", "and" to correspond to the Cause and Effect, Preconditions, Simultaneous Conditions, Contradiction, and Addition logic relations, respectively. It is complementary to the other token-level tasks which focus more on content-based knowledge.

SWAG

SWAG (Zellers et al. 2018) is a dataset with multiple choices questions about grounded situations. It questions models' understanding towards the relationship between two physical scenes. With the help of adversarial filtering (AF), Zellers et al. created a sufficiently large amount of questions automatically. For example, given *On stage, a woman takes a seat at the piano. She,* the question is to choose the following candidates: *A. sits on a bench as her sister plays with the doll B. smiles with someone as the music plays C. is in the crowd, watching the dancers D. **nervously sets her fingers on the keys.*** We obtain a positive or negative sample by concatenating the context and a candidate together (e.g *On stage, a woman takes a seat at the piano. She nervously sets her fingers on the keys*). There are one positive sample and three negative samples in a SWAG test instance. More examples are shown in the SWAG section of Table 1. By forcing the model to predict the next action, it requires inductive reasoning and temporal reasoning.

HellaSwag

HellaSwag (Zellers et al. 2019) is an argued version of SWAG with the same data format as SWAG, more inference steps and higher data quality. While HellaSwag also includes the dataset from WikiHow, we choose only the instances coming from ActivityNet to make the results comparable to the original SWAG dataset.

Sense Making with Reasoning (SMR)

Sense Making with Reasoning focuses on identifying the reason behind a statement (Wang et al. 2019) against commonsense. A model needs to understand that a specific statement (e.g *can is usually made of gold*) is against commonsense and to make a choice for the reason behind from three candidates (e.g *gold is too bright to make cans, **gold is too soft to make cans** and gold is too expensive to make cans*). We make a positive or negative sample by concatenating the statement and candidate reason together. For each test instance in SMR, there is a positive sample and two negative samples. More examples are shown in the SMR section of Table 1. This task is intuitively difficult since it requires a model to have deeper knowledge of with higher-level inference, which belongs to abductive reasoning.

Argument Reasoning Comprehension Task (ARCT)

Similar to SMR, Habernal et al. (2018) propose the ARCT dataset to test a model's abductive reasoning ability. Its domain lies in social topics such as search engine and LGBT rights, which is different from the daily-routine scenarios.

For example, given a reason R : *I find the idea that it is a sin to be born or live a life at all to be preposterous* and a claim C : *Christians have created a harmful atmosphere for gays*, this task is to pick the correct warrant W from two candidates: A . *being gay isn't considered a sin* B . ***being gay is considered a sin***, where $R \wedge W \rightarrow C$. We make a positive or negative sample by concatenating the reason, candidate warrant and claim together (e.g *I find the idea that it is a sin to be born or live a life at all to be preposterous and since being gay is considered a sin, Christians have created a harmful atmosphere for gays*). A test instance in ARCT contains a pair of positive and negative samples. More examples are shown in the ARCT section of Table 1. We further break this task into two variants, where ARCT1 represents the original dataset, ARCT2 represents an argumented dataset by adding negation to original instances to alleviate the statistical cues in the dataset (Niven and Kao 2019).

We integrated the above test sets into a commonsense ability test (CATs) benchmark, released for future research.

Pre-trained Models

We take six contextualized representation models that give the state-of-the-art performances on NLP benchmarks such as GLUE (Wang et al. 2018) and SQuAD (Rajpurkar, Jia, and Liang 2018). Off-the-shelf models are taken. Below we give the detailed settings.

GPT (Radford and Sutskever 2018) is a uni-directional transformer LM trained on 800M tokens of BookCorpus (Zhu et al. 2015). Given a text sequence $\mathbf{x} = [x_1, \dots, x_T]$, GPT works in a way similar to conventional auto-regressive (AR) LM:

$$\max_{\theta} \log p_{\theta}(\mathbf{x}) = \sum_{t=1}^T \log p_{\theta}(x_t | \mathbf{x}_{<t}),$$

where $\mathbf{x}_{<t} = [x_1, \dots, x_{t-1}]$. The model has dimension of hidden states $H = 768$, attention head numbers $A = 12$, number of layers $L = 12$ and total parameter size $P = 110M$.

GPT2 (Radford et al. 2019) works similarly as GPT with a few modifications on the hyperparameters. In particular, GPT2 optimizes the layer normalization, expands the vocabulary size to 50,257, increases the context size from 512 to 1024 tokens, and optimizes with a larger batchsize of 512. In addition, GPT2 is pre-trained on WebText, which was created from scraping web pages. The dataset roughly contains 8 million documents (40 GB). We study GPT2-base and GPT2-medium, with model size $H = 768, A = 12, L = 12, P = 117M$ and $H = 1024, A = 16, L = 24, P = 345M$, respectively, where the definitions of H, L and A are the same as for GPT.

BERT (Devlin et al. 2019) jointly trains on a masked language modeling task and a next sentence prediction task (NSP). The model is trained on the BookCorpus and English Wikipedia, a total of approximately 3300M tokens. BERT is designed with the following objective:

$$\max_{\theta} \log p_{\theta}(\bar{\mathbf{x}} | \tilde{\mathbf{x}}) \approx \sum_{t=1}^T m_t \log p_{\theta}(x_t | \tilde{\mathbf{x}}),$$

where $\tilde{\mathbf{x}}$ is a corrupted version of text sequence \mathbf{x} , and $\bar{\mathbf{x}}$ is masked tokens. $m_t = 1$ if token x_t belongs to $\bar{\mathbf{x}}$.

Here we consider BERT-base and BERT-large, with $H = 768, A = 12, L = 12, P = 117M$ and $H = 1024, A = 16, L = 24, P = 340M$, respectively, where the definitions of H, L and A are the same as for GPT.

XLNet (Yang et al. 2019) is trained with a permutation-based language modeling objective to capture bidirectional contexts while retain the benefits of AR models. Specifically, they let \mathcal{Z}_T be the set of all possible permutations of the length- T sequence $\mathbf{x} = [x_1, \dots, x_T]$:

$$\max_{\theta} \mathbf{E}_{\mathbf{z} \sim \mathcal{Z}_T} \left[\sum_{t=1}^T m_t \log p_{\theta}(x_{z_t} | \tilde{\mathbf{x}}_{\mathbf{z}_{<t}}) \right],$$

where z_t and $\mathbf{z}_{<t}$ are the t -th element and the first $t - 1$ elements of a permutation $\mathbf{z} \in \mathcal{Z}_T$, respectively. In this way, XLNet ensures that any specific token x_t in \mathbf{x} has seen all the tokens before or after it.

We consider XLNet-base and XLNet-large, whose model sizes are $H = 768, A = 12, L = 12, P = 117M$ and $H = 1024, A = 16, L = 24, P = 340M$, respectively, where the definitions of H, L and A are the same as for GPT. Note that XLNet-base is trained with the same data as BERT, while XLNet-large is trained with a larger dataset that consists of 32.98B subword pieces coming from Wiki, BookCorpus, Giga5, ClueWeb, and Common Crawl.

RoBERTa (Liu et al. 2019b) has the same architecture as BERT but is trained with dynamic masking, FULL-SENTENCES without NSP loss, a larger batch-size and a larger vocabulary size. Given the optimized design choice, one key difference of RoBERTa with other models is its large training dataset, which consists of BookCorpus, CC-NEWS, OpenWebText, and STORIES. With a total 160GB text, RoBERTa has access to more potential knowledge than the other models.

Experimental Design

The CAT datasets are applicable to any model that has a method to score a sentence. They fit with the pre-trained models above, which are by nature language models. We derive the score of a sentence below with uni-directional-context LMs and bi-directional-context LMs, respectively.

Formally, suppose the sentence S of n words $S = \{w_1, \dots, w_{k-1}, w_k, w_{k+1}, \dots, w_n\}$. We define the score of a sentence as:

$$Score(S) = \frac{\sum_{k=1}^n \log(P_{\theta}(w_k | context_k))}{n},$$

where the denominator n is for alleviating the influence of the sentence length to models' prediction, especially in sentence-level tasks. For a uni-directional model, $context_k = S_{<k} \equiv \{w_1, \dots, w_{k-1}\}$. The numerator becomes $\sum_{k=1}^n \log(P_{\theta}(w_k | S_{<k}))$, which is factorized from $\log(P_{\theta}(w_1, \dots, w_{k-1}, w_k, w_{k+1}, \dots, w_n))$. This is essentially a LM. For a bi-directional model, the $context_k = S_{-k}$, which represents the S with the k -th word being removed. In particular, the k -th word can be removed with being replaced by a special token '[MASK]' in BERT. The

	CA	WSC	SM	SMR	SWAG	HellaSwag	ARCT1	ARCT2	Average
RANDOM	0.500	0.500	0.500	0.333	0.250	0.250	0.500	0.500	0.416
GPT	0.830	0.558	0.735	0.354	0.592	0.263	0.472	0.528	0.542
GPT2-base	0.787	0.512	0.705	0.355	0.503	0.300	0.466	0.509	0.517
GPT2-medium	0.885	0.568	0.746	0.385	0.591	0.338	0.462	0.527	0.563
BERT-base	0.891	0.523	0.697	0.419	0.625	0.373	0.477	0.503	0.563
BERT-large	0.934	0.625	0.694	0.444	0.696	0.393	0.468	0.517	0.596
XLNet-base	0.809	0.544	0.662	0.374	0.494	0.381	0.516	0.526	0.543
XLNet-large	0.891	0.636	0.583	0.394	0.662	0.435	0.563	0.570	0.591
RoBERTa-base	0.901	0.623	0.750	0.423	0.712	0.414	0.501	0.537	0.565
RoBERTa-large	0.962	0.694	0.792	0.512	0.769	0.5	0.606	0.599	0.679
HUMAN	0.993	0.920	0.991	0.975	0.880	0.945	0.909	0.909	0.945

Table 3: Accuracy for each pre-trained contextualizer on each test set. The rightmost column shows the average of accuracy score of each model.

numerator $\sum_{k=1}^n \log(P_{\theta}(w_k|S_{-k}))$ can also be factorized from $\log(P_{\theta}(w_1, \dots, w_{k-1}, w_k, w_{k+1}, \dots, w_n))$ under the assumption that w_k is independent of the successive words (i.e. $w_{k+1}, w_{k+2}, \dots, w_n$), which is the bi-directional-context LM.

Intuitively, $P_{\theta}(w_k|context_k)$ can be interpreted as how probable a word w_k is given the $context_k: S_{<k}$ or S_{-k} . For example, let $S_{-k} = He\ put\ an\ [MASK]\ into\ the\ fridge,$ $w_{k1} = elephant$ and $w_{k2} = turkey$. $P_{\theta}(w_{k2}|S_{-k})$ should have a relatively larger value since filling in the “elephant” in the first case results in an improper sentence, which is against commonsense.

As introduced earlier (Table 1), all CATs tasks consist of instances with positive and negative sentences. After we score each sample in a test instance, the models predict the positive sample simply by taking the highest score in the instance.

Commonsense Tests Results

Table 3 shows the model performances with random choices as the baseline. Take WSC for example, the random baseline is 0.5, the human is 0.920 and all the models range between 0.512 and 0.694 with RoBERTa-large giving the best result of 0.694. Except for the ARCT task, all tested models demonstrate stronger performances than RANDOM, which indicates that the models all have varying degrees of commonsense. However, for most of the tasks, all of models are well below human performance.

Uni-directional Vs Bi-directional LM

We compare uni-directional (GPT, GPT2-base, GPT2-medium) and bi-directional models (Bert-base, Bert-large, XLNet-base, XLNet-large, RoBERTa-base, and RoBERTa-large). Picking the strongest model from each group, RoBERTa-large outperforms GPT2-medium by a large margin for every task. As mentioned before, RoBERTa-large has the same parameter size as GPT2-medium. However, RoBERTa-large is trained with much more data than GPT2-medium.

From Figure 1, we can see that except for the SM task, both BERT-large and XLNet-large outperform GPT2-

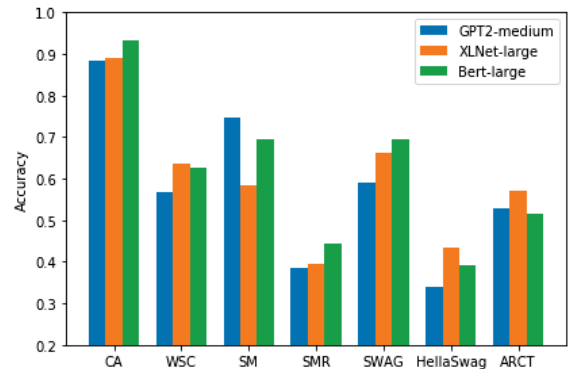


Figure 1: Comparison between bidirectional and unidirectional models among different tasks.

medium while BERT-large is trained with a smaller dataset than GPT2-medium. This indicates that bi-directional context can be more useful for learning commonsense. Intuitively, the models with bi-directional context can make more sentence-level inference. While only the preceeding words receive sufficient context in a uni-directional model, every word has the full context for bi-directional models. Table 4 shows examples where RoBERTa-large makes the correct prediction but GPT2-medium does not, we can see that the key tokens, which are considered to be the most influential part in making the correct prediction, lie in the middle of the sentence. This can be the main reason why bi-directional context is important for models’ commonsense ability.

Scale of Training Data

A larger training dataset intuitively allows a model to have access to more commonsense knowledge, thus performs better in our tests. Trained with by far the most data, RoBERTa is the winner for every task. Most of the models are in fact trained on a subset of the dataset used to train RoBERTa. However, larger dataset do not always work when the model capacity is limited with regard to commonsense. For example, GPT2-base underperforms GPT for many tasks in our dataset, which suggests GPT2-base un-

Token-level:

A. **Sam pulled up a chair to the piano, but the chair was broken, so he had to stand instead.**

B. Sam pulled up a chair to the piano, but the piano was broken, so he had to stand instead.

Sentence-level:

A. **Comments sections permit a reader to analyze many different perspectives in one place, and since I want to see all these ideas, even stupid ones, Comment sections have not failed.**

B. Comments sections permit a reader to analyze many different perspectives in one place, but since I don't want to see all these stupid ideas, Comment sections have not failed.

Table 4: Examples that a bi-directional model (RoBERTa-large) predicts correctly while a uni-directional model (GPT-2-medium) makes an incorrect prediction; the correct answer is **bolded**; the key tokens are underlined.

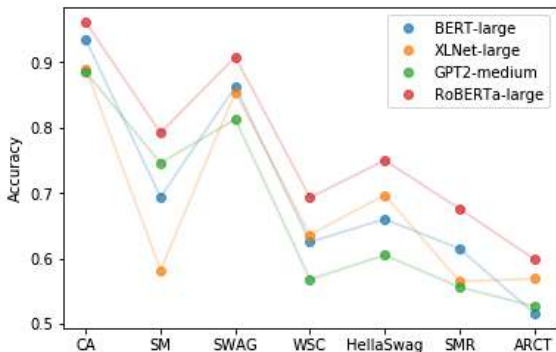


Figure 2: Models performances when the number of inference step (IS) increases. Tasks are ranked according to their IS in an increasing order from left to right.

derfits the WebText dataset with regard to commonsense. The fact that RoBERTa-base has the same parameter size as GPT2-base, yet benefits from the larger dataset suggests that bi-directional models have larger representative power in commonsense ability.

Number of Inference Steps

Similar to humans, the model performance can intuitively drop when commonsense inference becomes more complicated. To verify this intuition, we pick 100 sentences randomly from each test dataset and annotate the number of required inference steps (IS) of each instance manually. The inference step of each test dataset is defined as the average of the number of the turns of reasoning necessary for the instances from the test dataset. We choose to answer the question by counting the logical operations that exist in an instance. For example, for the sentence

Original:

A. People usually like wealth. B. People hardly like wealth.

Dual:

A. People usually hate wealth. B. People hardly hate wealth.

Table 5: Example of a robust test case; it contains a test instance from the original test set with a dual test instance created manually. When the key word changes from ‘like’ to ‘hate’, the correct answer switches from A to B. This is an unique feature of our robustness test sets.

They add a lot to the piece and I look forward to reading comments, but since comments sections always distract me from my work, Comment sections have failed., the logic chain is $(They\ add\ a\ lot\ to\ the\ piece \wedge I\ look\ forward\ to\ reading\ comments) \wedge comments\ sections\ always\ distract\ me\ from\ my\ work \rightarrow Comment\ sections\ have\ failed.$ Thus, this instance needs three inference steps.

In this way, we obtain the Inference Step (IS) for seven test datasets. Each instance is labeled by two expert annotators, and the inter-annotator agreement is 93%. The final IS is the average from both annotators. Figure 2 shows the results³ on the test cases with different IS. There is a decrease of performances as IS increases. SWAG and HellaSwag fall out the trend, which may suggest that the models have stronger commonsense ability in temporal reasoning.

Generally speaking, all of our tested models outperform the random baselines except for the ARCT task, which suggests that despite of using different modeling schemas, language modeling stands as an effective objective for extracting commonsense knowledge from large, raw texts. For each task, the overall performance increases with a larger model parameter size, a more sophisticated model design, and larger training data.

Robustness Test

The robustness of models in commonsense reasoning is an important perspective in evaluating deep commonsense ability. Intuitively, a person can reason whether a statement makes sense or not because he has consistent knowledge. If the statement changes slightly, for example, changing a key word, that person should still make the correct judgement.

We aim to test the robustness of the five models by making dual test samples. A dual instance to the original instance should test the same commonsense knowledge point or largely relevant to the original one. In this way, we expect that the model can demonstrate consistency in the decision. One example is shown in Table 5, which choosing A in the original instance should lead to choosing B in the dual case (See Figure 3 for more examples).

³The performances on tasks with more than one negative sample are transformed to binary-choice scales.

	Add	Del	Swap	Sub
RANDOM	0.50	0.50	0.50	0.50
GPT	0.13	0.17	0.51	0.19
GPT2-base	0.20	0.23	0.45	0.22
GPT2-medium	0.24	0.24	0.52	0.22
BERT-base	0.26	0.15	0.50	0.29
BERT-large	0.26	0.25	0.56	0.24
XLNet-base	0.36	0.16	0.41	0.26
XLNet-large	0.36	0.39	0.37	0.26
RoBERTa-base	0.20	0.27	0.47	0.35
RoBERTa-large	0.29	0.33	0.56	0.42

Table 6: Portion of consistent cases of each method for each contextualizer. Add stands for adding key words in the test sample; Del stands for deleting key words in the test sample; Swap stands for swapping the position of words in the test sample; Sub stands for replacing key words in the test sample. The best contextualizer for each method is **bolded**.

We consider multiple ways to construct a dual test instance. Particularly, a dual test instance is built by methods: adding, deleting and replacing words in a test sample, or swapping two words in the sample, thereby resulting in a closely related test instance. All of our dual test instances are constructed from the original commonsense test data.

We construct 75 dual instances for each method above over WSC, SM, and ARCT, keeping the instances from each dataset approximately equivalent in order to evaluate the influence of different duality methods to the models. We then pair each dual instance with the original instance to form a new test case. If the model gives the correct or wrong prediction for both of the instances in this case, we recognize it as a *consistent* case.

The results are shown in Table 6. In theory, a model equipped with relevant commonsense should give consistent predictions on a pair of dual test case. However, we find that none of the models reach consistency. In fact, their consistency is well below the random baselines except for the Swap method.

To better investigate the reason behind the poor consistency, we look at inconsistent cases from the pre-trained model (i.e RoBERTa-large). Similar to Trinh and Le (2018), we investigate how the model makes decision between two candidate sentences $S_{correct}$ and $S_{incorrect}$ where they have the same number of words. In particular, we look at:

$$q_k = \log\left(\frac{P_\theta(w_k|S_{correct} - \{w_k\})}{P_\theta(w_k|S_{incorrect} - \{w_k\})}\right),$$

where $1 \leq k \leq n$. It follows that the choice between $S_{correct}$ and $S_{incorrect}$ is made by the value $Q = \sum_k q_k$ being bigger than 0 or not. Visualizing the value of each q_k provides more insights into the decisions of the model.

From Figure 3, we can tell that the model is confused by the modification, tending to give the same predictions over a pair of dual samples despite that they have different gold labels, especially for Sub, Add and Del. This further reveals that the commonsense knowledge contained in the

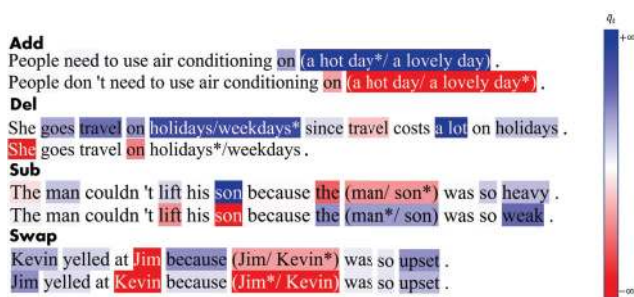


Figure 3: Samples of questions from Add, Del, Sub and Swap predicted correctly for the original instance but incorrectly for the dual instance. Note that a sentence here represents a test instance with a pair of positive and negative samples, represented by (.../...). Here we mark the correct prediction by an asterisk and display the normalized q_t by coloring its corresponding word.

pre-trained models may remain in a surface level, without deep semantic comprehension.

Related Work

Liu et al. (2019) evaluate BERT (Devlin et al. 2019), GPT (Radford and Sutskever 2018), and ELMo (Peters et al. 2018) on a variety of linguistics tasks. Their results suggest that the features generated by pre-trained contextualizer are sufficient for high performance on a board set of tasks but models fail on tasks requiring fine-grained linguistics knowledge. Tenney et al. (2019) evaluate similar models on a variety of sub-sentence linguistic analysis tasks. Their results suggest that contextualized word representation encode both syntax and semantics. Our work is in line in the sense that contextualized representation encode rich knowledge to be 'probed'. However, we focus on evaluating the commonsense in those representations. To our best knowledge, this is the first work to systematically evaluate commonsense in pre-trained models.

Our evaluation method is similar to Trinh and Le (2018), who make use of LM to score a sentence. However, they focus on Winograd schema questions with only self-trained recurrent LMs while we test five models' commonsense with seven diverse tasks.

Conclusion

We studied the commonsense knowledge and reasoning ability of pre-trained contextualizers with a suite of seven diverse probing tasks, showing that large-scale pre-trained contextualized representation has a certain degree of commonsense knowledge, but there is still a quite large gap between the current state-of-the-art representation models and robust human-level commonsense reasoning, which may require more breakthrough in modeling. We release our test sets, named CATs, publicly.

Acknowledgments

We would like to thank the anonymous reviewers for their insightful comments, and Mr. Cunxiang Wang for his help on the collection of the data. Yue Zhang is the corresponding author.

References

- Bowman, S. R.; Angeli, G.; Potts, C.; and Manning, C. D. 2015. A large annotated corpus for learning natural language inference. In *EMNLP*.
- Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 4171–4186.
- Habernal, I.; Wachsmuth, H.; Gurevych, I.; and Stein, B. 2018. The argument reasoning comprehension task: Identification and reconstruction of implicit warrants. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, 1930–1940.
- Kim, N.; Patel, R.; Poliak, A.; Xia, P.; Wang, A.; McCoy, T.; Tenney, I.; Ross, A.; Linzen, T.; Van Durme, B.; Bowman, S. R.; and Pavlick, E. 2019. Probing what different NLP tasks teach machines about function word comprehension. In *Proceedings of the Eighth Joint Conference on Lexical and Computational Semantics (*SEM 2019)*, 235–249.
- Lai, G.; Xie, Q.; Liu, H.; Yang, Y.; and Hovy, E. H. 2017. Race: Large-scale reading comprehension dataset from examinations. In *EMNLP*.
- Levesque, H. J.; Davis, E.; and Morgenstern, L. 2012. The winograd schema challenge. In *Proceedings of the Thirteenth International Conference on Principles of Knowledge Representation and Reasoning, KR'12*, 552–561. AAAI Press.
- Liu, H., and Singh, P. 2004. Conceptnet – a practical commonsense reasoning tool-kit. *BT Technology Journal* 22:211–226.
- Liu, N. F.; Gardner, M.; Belinkov, Y.; Peters, M. E.; and Smith, N. A. 2019a. Linguistic knowledge and transferability of contextual 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)*, 1073–1094.
- Liu, Y.; Ott, M.; Goyal, N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; and Stoyanov, V. 2019b. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- LoBue, P., and Yates, A. 2011. Types of common-sense knowledge needed for recognizing textual entailment. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, 329–334.
- Niven, T., and Kao, H.-Y. 2019. Probing neural network comprehension of natural language arguments. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 4658–4664.
- Peters, M.; Neumann, M.; Iyyer, M.; Gardner, M.; Clark, C.; Lee, K.; and Zettlemoyer, L. 2018. Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, 2227–2237.
- Radford, A., and Sutskever, I. 2018. Improving language understanding by generative pre-training.
- Radford, A.; Wu, J.; Child, R.; Luan, D.; Amodei, D.; and Sutskever, I. 2019. Language models are unsupervised multitask learners.
- Rajpurkar, P.; Jia, R.; and Liang, P. 2018. Know what you don't know: Unanswerable questions for SQuAD. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, 784–789.
- Tenney, I.; Xia, P.; Chen, B.; Wang, A.; Poliak, A.; McCoy, R. T.; Kim, N.; Durme, B. V.; Bowman, S. R.; Das, D.; and Pavlick, E. 2019. What do you learn from context? probing for sentence structure in contextualized word representations. *ArXiv abs/1905.06316*.
- Trinh, T. H., and Le, Q. V. 2018. A simple method for commonsense reasoning. *CoRR abs/1806.02847*.
- Wang, A.; Singh, A.; Michael, J.; Hill, F.; Levy, O.; and Bowman, S. 2018. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, 353–355.
- Wang, C.; Liang, S.; Zhang, Y.; Li, X.; and Gao, T. 2019. Does it make sense? and why? a pilot study for sense making and explanation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 4020–4026.
- Yang, Z.; Dai, Z.; Yang, Y.; Carbonell, J. G.; Salakhutdinov, R.; and Le, Q. V. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. *CoRR abs/1906.08237*.
- Zellers, R.; Bisk, Y.; Schwartz, R.; and Choi, Y. 2018. SWAG: A large-scale adversarial dataset for grounded commonsense inference. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 93–104.
- Zellers, R.; Holtzman, A.; Bisk, Y.; Farhadi, A.; and Choi, Y. 2019. HellaSwag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 4791–4800.
- Zhu, Y.; Kiros, R.; Zemel, R. S.; Salakhutdinov, R.; Urtasun, R.; Torralba, A.; and Fidler, S. 2015. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. *2015 IEEE International Conference on Computer Vision (ICCV)* 19–27.