

# MCScript2.0: A Machine Comprehension Corpus Focused on Script Events and Participants

Simon Ostermann<sup>1</sup> Michael Roth<sup>1,2</sup> Manfred Pinkal<sup>1</sup>

<sup>1</sup>Saarland University <sup>2</sup>Stuttgart University  
{simono|pinkal}@coli.uni-saarland.de rothml@ims.uni-stuttgart.de

## Abstract

We introduce *MCScript2.0*, a machine comprehension corpus for the end-to-end evaluation of **script** knowledge. *MCScript2.0* contains approx. 20,000 questions on approx. 3,500 texts, crowdsourced based on a new collection process that results in challenging questions. Half of the questions cannot be answered from the reading texts, but require the use of commonsense and, in particular, script knowledge. We give a thorough analysis of our corpus and show that while the task is not challenging to humans, existing machine comprehension models fail to perform well on the data, even if they make use of a commonsense knowledge base. The dataset is available at [http://www.sfb1102.uni-saarland.de/?page\\_id=2582](http://www.sfb1102.uni-saarland.de/?page_id=2582)

## 1 Introduction

People have access to a wide range of commonsense knowledge that is naturally acquired during their lifetime. They make frequent use of such knowledge while speaking to each other, which makes communication highly efficient. The conversation in Example 1 illustrates this: For Max, it is natural to assume that Rachel paid during her restaurant visit, even if this fact was not mentioned by Rachel.

- (1) Rachel: “Yesterday, I went to this new Argentinian restaurant to have dinner. I enjoyed a tasty steak.”  
Max: “What did you pay?”

*Script knowledge* is one of the most important types of commonsense knowledge and subsumes such information (Schank and Abelson, 1977). It is defined as knowledge about everyday situations, also referred to as *scenarios*. It subsumes information about the actions that take place during such situations, and their typical temporal order, referred

to as *events*, as well as the persons and objects that typically play a role in the situation, referred to as *participants*. In the example, script knowledge subsumes the fact that the *paying* event is a part of the *eating in a restaurant* scenario.

Recent years have seen different approaches to learning script knowledge, centered around two strands: Work around narrative chains that are learned from large text collections (Chambers and Jurafsky, 2008, 2009), and the manual induction of script knowledge via crowdsourcing (Regneri et al., 2010; Wanzare et al., 2016). Script knowledge has been represented both symbolically (Jans et al., 2012; Pichotta and Mooney, 2014; Rudinger et al., 2015) and with neural models (Modi and Titov, 2014; Pichotta and Mooney, 2016). Scripts have been evaluated mostly intrinsically (Wanzare et al., 2017; Ostermann et al., 2017). An exception is *MCScript* (Ostermann et al., 2018a), a reading comprehension corpus with a focus on script knowledge, and a predecessor to the data set presented in this work. Previous work has shown, however, that script knowledge is not required for performing well on the data set (Ostermann et al., 2018b). Hence, to date, there exists no evaluation method that allows one to systematically assess the contribution of models of script knowledge to the task of automated text understanding.

Our work closes this gap: We present *MCScript2.0*, a reading comprehension corpus focused on script events and participants. It contains more than 3,400 texts about everyday scenarios, together with more than 19,000 multiple-choice questions on these texts. All data were collected via crowdsourcing. About half of the questions require the use of commonsense and script knowledge for finding the correct answer (like the question in Example 1), a notably higher number than in *MCScript*. We show that in comparison to *MCScript*, commonsense-based questions in *MCScript2.0* are

<b>T</b>	(...) We put our ingredients together to make sure they were at the right temperature, preheated the oven, and pulled out the proper utensils. We then prepared the batter using eggs and some other materials we purchased and then poured them into a pan. After baking the cake in the oven for the time the recipe told us to, we then double checked to make sure it was done by pushing a knife into the center. We saw some crumbs sticking to the knife when we pulled it out so we knew it was ready to eat !
<b>Q1</b>	<i>When did they put the pan in the oven and bake it according to the instructions?</i> After eating the cake. ✗ After mixing the batter. ✓
<b>Q2</b>	<i>What did they put in the oven?</i> The cake mix. ✓ Utensils. ✗

Figure 1: Example text fragment from MCScript2.0

also harder to answer, even for a model that makes use of a commonsense database. Thus, we argue that MCScript2.0 is the first resource which makes it possible to evaluate how far models are able to exploit script knowledge for automated text understanding.

Figure 1 shows a text snippet from a text in MCScript2.0, together with two questions with answer alternatives<sup>1</sup>. To find an answer for question 1, information about the temporal order of the steps for baking a cake is required: The cake is put in the oven after mixing the batter, and not after eating it—a piece of information not given in the text, since the event of putting the cake in the oven is not explicitly mentioned. Similarly, one needs script knowledge about which participants are typically involved in which events to know that the cake mix rather than the utensils is put into the oven. Both incorrect answer candidates are distractive: The utensils as well as the action of eating the cake are mentioned in the text, but wrong answers to the question. Our contributions are as follows:

- We present a new collecting method for challenging questions whose answers require commonsense knowledge and in particular script knowledge, as well as a new resource that was created with this method.

<sup>1</sup>More text samples are given in the Supplemental Material.

- We show that the task is simple for humans, but that existing benchmark models, including a top-scoring machine comprehension model that utilizes a resource for commonsense knowledge, struggle on the questions in MCScript2.0; especially on questions that require commonsense knowledge.
- We compare MCScript2.0 to MCScript, the first machine comprehension resource for evaluating models of script knowledge. We show that in comparison to MCScript, the number of questions that require script knowledge is increased by a large margin and that such questions are hard to answer. Consequently, we argue that our dataset provides a more robust basis for future research on text understanding models that use script knowledge.

## 2 Why another Machine Comprehension Dataset on Script Knowledge?

MCScript (Ostermann et al., 2018a) is the first machine comprehension dataset designed to evaluate script knowledge in an end-to-end machine comprehension application, and to our knowledge the only other existing extrinsic evaluation dataset for script knowledge. Recent research has shown, however, that commonsense knowledge is not required for good performance on the dataset (Ostermann et al., 2018b; Merkhofer et al., 2018).

We argue that this is due to the way in which questions were collected. During the collection process, workers were not shown a text, but only a very short description of the text scenario. As a result, many questions ask about general aspects of the scenario, without referring to actual details. This leads to the problem that there are many questions with standardized answers, i.e. questions that can be answered irrespective of a concrete reading text. Examples 2 and 3 show two such cases, where the correct answer is almost exclusively *in the shower* and *on the stove*, independent of the text or even scenario.

- (2) Where did they wash their hair?
- (3) Where did they make the scrambled eggs?

Merkhofer et al. (2018) found that such information can essentially be learned from only the training data, using a simple logistic regression classifier and surface features regarding words in the text, question and answer candidates.

Also, many questions require vague inference over general commonsense knowledge rather than script knowledge. Example 4 illustrates this: The simple fact that planting a tree gets easier if you have help is not subsumed by script knowledge about planting a tree, but a more general type of commonsense knowledge.

- (4) *Text*: Once you know where to dig , select what type of tree you want. (...) Dig a hole large enough for the tree and roots . Place the tree in the hole and then fill the hole back up with dirt . (...)

**Q**: Would it have been easier to plant the tree if they had help?

yes ✓

no ✗

We inspected a random sample of 50 questions from the publicly available development set that were misclassified by the logistic model of [Merkhofer et al. \(2018\)](#). We found that for over 90% of the inspected questions, the use of script knowledge would be only marginally relevant.

We present a new data collection method and corpus that results in a larger number of challenging questions that require script knowledge. In particular, we define a revised question collection procedure, which ensures a large proportion of non-trivial commonsense questions.

### 3 Corpus Creation

Texts, questions, and answer candidates are required for a multiple choice machine comprehension dataset. Our data collection process for texts and answers is based on the MCScript data and the methods developed there, but with several crucial differences. Like [Ostermann et al. \(2018a\)](#), we create the data set via crowdsourcing. The question collection is revised to account for the shortcomings found with MCScript.

Similarly to [Ostermann et al. \(2018a\)](#), we are interested in questions that require inference over script knowledge for finding a correct answer. Creating such questions is challenging: When questions are collected by showing a reading text and asking crowdsourcing workers to write questions, their answer can usually be read off the text. The authors of MCScript thus decided to not show a reading text at all, but only a short summary of the text scenario. This resulted in too general questions, so we decided for a third option: We identi-

fied a number of *target sentences* in the reading text and guided workers to formulate questions about script-related details in these sentences. The target sentences were then hidden from the text, meaning that relevant information would have to be inferred from common sense during the answer collection and also in the task itself. In the following sections, we describe the three data collection steps in detail.

#### 3.1 Text Collection

As a starting point, we reused all texts from MCScript (2,119 texts on 110 scenarios) for our data set. To increase the topical coverage and diversity of the data set, we added texts for 90 new scenarios to our collection. As for MCScript, we selected topically different and plausible everyday scenarios of varying complexity, which were not too fine-grained (such as *opening a window*). The scenarios were taken from 3 sources: First, we extracted scenarios from several script collections ([Wanzare et al., 2016](#); [Regneri et al., 2010](#); [Singh et al., 2002](#)) that are not part of MCScript. Second, we inspected the *spinn3r* blog story corpus ([Burton et al., 2009](#)), a large corpus of narrative blog stories and identified additional scenarios in these stories. Third, we added new scenarios that are related to existing ones or that extend them.

We collected 20 texts per new scenario, using the same text collection method as [Ostermann et al. \(2018a\)](#): We asked workers to tell a story about a certain everyday scenario “as if talking to a child”. This instruction ensures that the resulting stories are simple in language and clearly structured. Texts collected this way have been found to explicitly mention many script events and participants ([Modi et al., 2016](#); [Ostermann et al., 2018a](#)). They are thus ideal to evaluate script-based inference.

#### 3.2 Question Collection

For the question collection, we followed [Ostermann et al. \(2018a\)](#) in telling workers that the data are collected for a reading comprehension task for children, in order to get linguistically simple and explicit questions. However, as mentioned above, we guide workers towards asking questions about target sentences rather than a complete text.

As target sentences, we selected every fourth sentence in a text. In order to avoid selecting target sentences with too much or too little content, we only considered sentences with less than 20 tokens,

First up was pitching the tent , then we set about making a fire to cook dinner.	Ask for: <b>the tent!</b> What did they pitch first? Ask for: <b>a fire!</b> What did they set in order to make dinner?
It took a while to collect all the wood we needed and get it roaring, but once we got the food on the fire we could finally relax. After soaking in nature and enjoying our campfire meal, night had fallen and it was time for bed. We curled up in our sleeping bags and told scary stories until we drifted off.	
The next morning we awoke early to stoke the fire.	Ask for: <b>the next morning!</b> When did they stoke the fire? Ask for: <b>the fire!</b> What did they
We then started breakfast. A day of hiking and fishing went by far too fast and after another night under the stars we packed up our gear and headed home.	

Figure 2: Screenshot of an item in the participant question collection experiment.

but that contained 2 or more noun phrases.<sup>2</sup>

In a series of pilot studies, we then showed the texts with highlighted target sentences to workers and asked them to write questions about these sentences. We however found, that in many cases, the written questions were too general or nonsensical.

We concluded that an even more structured task was required and decided to concentrate on questions of two types: (1) questions that ask about participants, and (2) questions about the temporal event structure of a scenario. Participants are usually instantiated by noun phrases (NPs), while events are described by verb phrases (VPs). We thus used *Stanford CoreNLP* (Manning et al., 2014) to extract both NPs and VPs in the target sentences and split up the experiment into two parts: In the first part, workers were required to write questions that *ask about the given noun phrase*. Figure 2 shows a screenshot of an item from the first part. The first column shows the reading text with the target sentence highlighted. The second columns shows all extracted phrases with a field for one question per phrase.<sup>3</sup> Full details of the experiment instructions are given in the Supplemental Material.

In the second part, we then asked workers to write a temporal question (when, how long, etc.) *using the given verb phrase*. We found that an exact repetition of the NP instructions for the second part (“ask about the given verb phrase”) resulted in unnatural questions, so we adapted the instructions. A screenshot of the VP experiment is given in the Supplemental Material.

We showed each text to two workers and asked

<sup>2</sup>All parameters were selected empirically, by testing different values and analyzing samples of the resulting data.

<sup>3</sup>If the noun phrase was part of a prepositional phrase or a construction of the form “NP of NP”, we took the whole phrase instead, because it is more natural to ask for the complete phrase. In order to avoid redundancy, we only looked at NPs that had no other NPs as parents. We also excluded noun phrases that referred to the narrator (*I, me* etc.).

them to write one question per VP or NP. Workers were only allowed to work on either the VP or the NP part, since the instructions could easily be confused. In order to exclude yes/no questions, we did not accept inputs starting with an auxiliary or modal verb. Also, all questions needed to contain at least 4 words. We asked workers to use *they* to refer to the protagonist of the story and other types of mentions (e.g. pronouns like *I, you, we* or the word *narrator*) were not accepted.

### 3.3 Answer Collection

For collecting answer candidates we hid the target sentences from the texts and showed them with up to 12 questions, to keep the workload at an acceptable level. If there were more questions for a text, we selected 12 questions at random.

Since the target sentences are hidden in the texts, it can be expected that some questions cannot be answered from the text anymore. However, the necessary information for finding an answer might be inferred from script knowledge, so workers were explicitly told that they might need commonsense to find an answer. Some answers can still be read off the text, if other parts of the texts contain the same information as the hidden target sentences. For other questions, neither the text nor script knowledge provides sufficient information for finding an answer.

As for the creation of MCScript, workers first had to conduct a 4-way classification for each question to account for these cases: *text-based* (answer is in the text), *script-based* (answer can be inferred from script knowledge), *unfitting* (question doesn’t make sense), *unknown* (answer is not known). Having such class annotations is not only useful for evaluation, but it also sensitizes workers for the fact that they are explicitly allowed to use background knowledge.

In the experiment, workers were also instructed



<i>text-based</i>	<i>script-based</i>	<i>text-or-script</i>	<i>unfitting</i>	<i>unknown</i>
9,357	12,433	2,403	3,240	6,457
<b>total answerable: 24,193</b>			<b>total not answerable: 9,397</b>	

Table 1: Distribution of question labels, before validation.

to write both a correct and a plausible incorrect answer for questions labeled as *text-based* or *script-based*. We follow (Ostermann et al., 2018a) and require workers to write an alternative question if the labels *unfitting* or *unknown* are used, in order to level out the workload.

We presented each question to 5 workers, resulting in 5 judgements and up to 5 incorrect and correct answer candidates per question. For the final data set, we considered questions with a majority vote (3 out of 5) on *text-based* or *script-based*. We also included questions without a majority vote, but for which at least 3 workers assigned one of *text-based* or *script-based*. In that case, we assigned the new label *text-or-script* and also accepted the question for the final data set. This seemed reasonable, since at least 3 workers wrote answers for the question, meaning it could still be used in the final data collection. The remaining questions were discarded.

### 3.4 Answer Candidate Selection

In a last step, we selected one correct and one incorrect answer from all possible candidates per question for the data set. To choose the most plausible correct answer candidate, we adapt the procedure from Ostermann et al. (2018a): We normalize all correct answers (lowercasing, normalizing numbers<sup>4</sup>, deleting stopwords<sup>5</sup>) and then merge candidates that are contained in another candidate, and candidate pairs with a Levenshtein (1966) distance of less than 3. The most frequent candidate is then selected as correct answer. If there was no clear majority, we selected a candidate at random.

To select an incorrect answer candidate, we adapt the *adversarial filtering* algorithm from Zellers et al. (2018). Our implementation uses a simple classifier that utilizes shallow surface features. The algorithm selects the incorrect answer candidate from the set of possible candidates that is most difficult for the classifier, i.e. an incorrect answer that is hard to tell apart from the correct

answer (e.g. the incorrect answers in Figure 1: *eating* and *utensils* are also mentioned in the text). By picking incorrect answers with the adversarial filtering method, the dataset becomes robust against surface-oriented methods.

Practically, the algorithm starts with a random assignment, i.e. a random incorrect answer candidate per question. This assignment is refined iteratively, such that the most difficult candidate is selected. In each iteration, the algorithm splits the data into a random training part and a test part. The classifier is trained on the training part and then used to classify *all* possible candidates in the test part. The assignment of answer candidates in the test data is then changed such that the most difficult incorrect answer candidate per question is picked as incorrect answer. After several iterations through the whole dataset, the number of changed answer candidates usually stagnates and the algorithm converges.

For MCScript2.0, we use the logistic classifier mentioned in Section 2, which only uses surface features and is thus well suited for the filtering algorithm. Implementation details and pseudocode are given in the Supplemental Material.

## 4 Corpus Analysis

### 4.1 General Statistics

In total, MCScript2.0 comprises 19,821 questions on 3,487 texts, i.e. 5.7 questions on average per text. The average length of texts, questions and answers is 164.4 tokens, 8.2 tokens and 3.4 tokens, respectively.

In the data collection process, we crowdsourced 1,800 new texts, resulting in a total of 3,919 texts for 200 scenarios. On average, there are 1.98 target sentences per text. In the question collection, we gathered 42,132 questions that were then used for the answer collection. For 8,242 questions, there was no clear majority on the question label. Table 1 shows the label distribution on the remaining 33,890 questions. 24,193 of these could be answered, i.e. 71%.

To increase data quality, we conducted a manual validation of the data. Four student assistants re-

<sup>4</sup>We used *text2num*, <https://github.com/ghewgill/text2num>.

<sup>5</sup>*and, or, to, the, a*

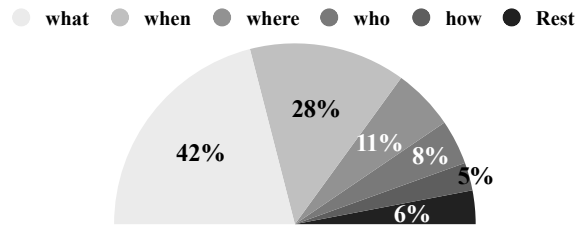


Figure 3: Distribution of question types

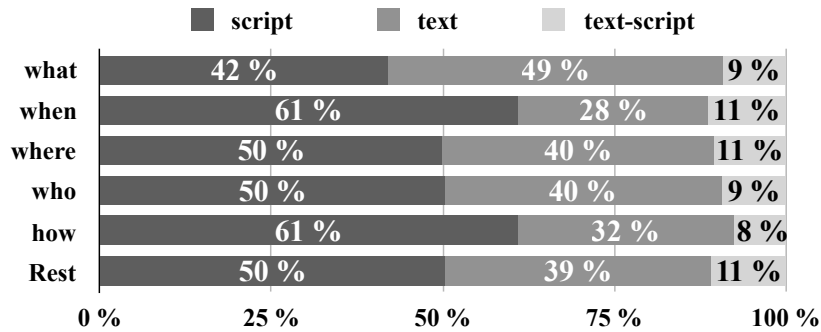


Figure 4: Proportion of labels per question type.

placed erroneous answers and deleted nonsensical questions, question duplicates and incoherent texts.

During validation, 152 texts were found to be incoherent and discarded (along with all questions). Additionally, 3,388 questions were deleted because they were nonsensical or duplicates. 1,620 correct and 2,977 incorrect answers were exchanged, resp., because they were inappropriate. If a question deletion resulted in texts without any questions, or if a text did not have any answerable questions, the text was discarded, too.

After question validation, the final dataset comprises 9,935 questions that are labeled as script-based, 7,908 as text-based, and 1,978 as text-or-script.

## 4.2 Questions

Figure 3 gives the distribution over question types, which we extracted by looking at the first word in the question. The largest number of questions are *what* questions, most of which ask about participants of a script. *When* questions make up the second largest group, asking for temporal event structure. During the VP question experiment, some workers ignored that we asked for temporal questions only, which resulted in a number of *how* questions.

MCScript2.0 contains 50% questions labeled as script-based, which is a notably larger amount as compared to the approximately 27% of questions in MCScript labeled as script-based. The number of

script-based questions varies between the question types, as can be seen in Figure 4. While *when* and *how* questions require script knowledge for finding an answer in more than 60% of cases, less than half of *what* questions do so. A simple explanation for this could be that *when* or *how* questions typically ask for events, while *what* questions ask for participants. Events are usually referred only once in a text, i.e. with the hiding of the respective event mention, the needed information has to be inferred. Participants in contrast tend to appear more often throughout a story.

Example 5 below illustrates this. Question 1 was originally asked about a sentence in which the plates are set for the dinner guests. The guests still appear in another sentence, so the answer can be inferred from the text. For question 2, in contrast, script knowledge is required for finding an answer: The event of *bringing the items to the table* is not mentioned anymore, so the information that this happens typically after counting plates and silverware needs to be inferred.

- (5) **T:** (...) I was told that there would be 13 or 14 guests. First I counted out 14 spoons, then the same number of salad forks, dinner forks, and knives. (...) I set each place with one napkin, one dinner fork, one salad fork, one spoon, and one knife. (...)

**Q1:** Who are the plate and cup for?  
dinner guests ✓ the neighbor ✗

**Q2:** When did they bring the items over to

- the table?
- after counting them ✓
- after placing them on the table ✗

## 5 Experiments

We test three benchmark models on MCScript2.0 that were also evaluated on MCScript, so a direct comparison is possible. For the experiments, we split the data into a training set (14,191 questions on 2,500 texts), a development set (2,020 questions on 355 texts) and a test set (3,610 questions on 632 texts). All texts of 5 randomly chosen scenarios were assigned completely to the test set, so a part of the test scenarios are unseen during training.

### 5.1 Models

#### Logistic Regression Classifier

As first model, we reimplemented the logistic regression classifier proposed by Merkhofer et al. (2018), which was also used in the adversarial filtering algorithm. The classifier employs 3 types of features: (1) Length features, encoding the length of the text, answer and questions on the word and character level, (2) overlap features, encoding the amount of literal overlap between text, question, and answers, and (3) binary lexical patterns encoding the presence or absence of words or combinations of words in answer, text and question.

#### Attentive Reader

As second model, we implement an attentive reader (Hermann et al., 2015). We adopt the formulation by Ostermann et al. (2018a) (originally by Chen et al. (2016)). All tokens in text, question and answers are represented with word embeddings. Bidirectional gated recurrent units (GRUs, Cho et al. (2014)) process the text, question and answers and transform them into sequences of contextualized hidden states. The text is represented as a weighted average of the hidden states with a bilinear attention formulation, and another bilinear weight matrix is used to compute a scalar as score for each answer. For a formalization, we refer to Ostermann et al. (2018a) and Chen et al. (2016).

#### Three-way Attentive Network (TriAN)

As third model, we use a three-way attentive network (Wang et al., 2018), the best-scoring model of the shared task on MCScript<sup>6</sup>. Various types of in-

<sup>6</sup>Code available at <https://github.com/intfloat/commonsense-rc>

formation are employed to represent tokens: Word embeddings, part of speech tags, named entity embeddings, and word count/overlap features, similar to the logistic classifier. Three bidirectional LSTM (Hochreiter and Schmidhuber, 1997) modules are used to encode text, question and answers. The resulting hidden representations are reweighted with three attention matrices and then summed into vectors using three self-attention layers.

Additionally, token representations are enhanced with *ConceptNet* (Speer et al., 2017) relations as a form of induced commonsense knowledge. ConceptNet is a large database of commonsense facts, represented as triples of two entities with a predicate. Relevant ConceptNet relations between words in the answer and the text are queried from the database and represented with relation embeddings, which are learned end-to-end during training and appended to the text token representations.

In contrast to Wang et al. (2018), we use the non-ensemble version of TriAN without pretraining on *RACE* (Lai et al., 2017), for better comparability to the other models.

### 5.2 Human Upper Bound

To assess human performance, 5 student assistants performed the reading comprehension task on 60 texts each. To assess agreement, 20 texts were annotated by all students. The annotators reached averaged pairwise agreement of 96.3% and an average accuracy of 97.4%, which shows that this is a simple task for humans.

### 5.3 Results

**Overall Performance.** Table 2 gives details about the performance of the 3 benchmark models on the test set, and on script-based ( $acc_{scr}$ ) and text-based ( $acc_{txt}$ ) questions in the test set. As can be seen, the logistic model scores worst, presumably because it has been used for the adversarial filtering algorithm and the data are thus most challenging for this model. TriAN performs best, clearly outperforming the attentive reader. TriAN is apparently superior in its way of text processing, since it employs a richer text representation and exploits attention mechanisms on more levels, which is reflected by a higher accuracy on text-based questions. In contrast, script-based questions seem to be challenging for TriAN. This is interesting, because it shows that ConceptNet alone cannot provide sufficient information for answering the kind of questions that can be found in MCScript2.0.

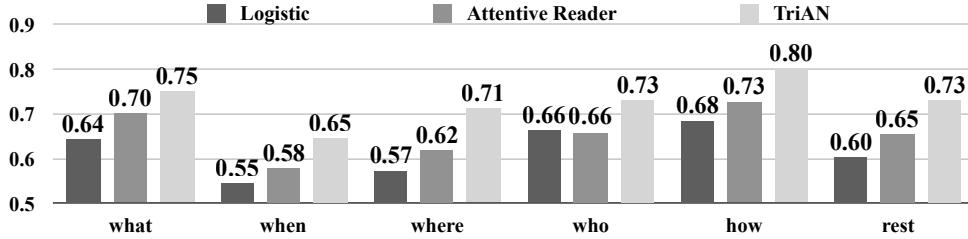


Figure 5: Performance of the models on question types.

	acc	acc <sub>scr</sub>	acc <sub>txt</sub>
<i>Logistic Model</i>	0.61	0.56	0.67
<i>Attentive Reader</i>	0.65	0.63	0.68
<i>TriAN</i>	<b>0.72</b>	<b>0.67</b>	<b>0.78</b>
<i>Humans</i>	0.97		

Table 2: Accuracy on test set, and on script/text-based questions (acc<sub>scr</sub>, acc<sub>txt</sub>) on **MCScript2.0**. The maximum per column is printed in **bold**.

**Comparison to MCScript.** Since the same models were used for MCScript, a comparison of their performance on both datasets is possible. Results on MCScript are given in Table 3.<sup>7</sup> As can be seen, the performance of all three models is worse on MCScript2.0, showing that the dataset is generally more challenging. In contrast to MCScript, script-based questions in MCScript2.0 are clearly harder to answer than text-based questions: All models perform worse on script-based questions as compared to text-based questions. In comparison to MCScript, the performance of TriAN is **12%** lower. This indicates that the new mode of question collection and the answer selection via adversarial filtering resolve some of the difficulties with MCScript.

To assess whether the performance difference to MCScript is due to the 90 new scenarios being more challenging, we additionally evaluated the models on these scenarios only. We found no performance difference on the new vs. old scenarios.

**Influence of Adversarial Filtering.** To find out how large the influence of the new question collection method and the answer selection via adversarial filtering is, we conducted an additional experiment: We applied the answer selection method of Ostermann et al. (2018a) to our data set to create

<sup>7</sup>For the attentive reader, numbers were taken from (Ostermann et al., 2018b). The other models were retrained (and in the case of the logistic model re-implemented), since no exact numbers on script/text-based questions were published.

	acc	acc <sub>scr</sub>	acc <sub>txt</sub>
<i>Logistic Model</i>	0.79	0.76	<b>0.81</b>
<i>Attentive Reader</i>	0.72	0.75	0.71
<i>TriAN</i>	<b>0.80</b>	<b>0.79</b>	<b>0.81</b>
<i>Humans</i>	0.98		

Table 3: Accuracy on the test set and on script/text-based questions (acc<sub>scr</sub>, acc<sub>txt</sub>) on **MCScript**. The maximum per column is printed in **bold**.

an alternative version of the data that is not based on adversarial filtering. Correct answers were selected to have the lowest possible overlap with the reading text. Incorrect answers were selected using the majority voting technique described in Section 3.4.

We found that the adversarial filtering accounts for around two thirds of the total accuracy difference of TriAN as compared to MCScript, i.e. one third of the difference can be attributed to the new question collection. This means that both modifications together add to the larger difficulty of MCScript2.0.

**Question Types.** Figure 5 shows the performance of the models on single question types, as identified in Section 4. It is clear that *when* questions are most challenging for all models. The logistic classifier performs almost at chance level. As far as TriAN is concerned, we found that many cases of errors ask for the typical temporal order of events, as Example 6 illustrates:

- (6) **Q:** When did they put the nozzle in their tank?  
before filling up with gas. ✓  
after filling up with gas. ✗

The event of *put the nozzle in the tank* is not mentioned in the shown version of the text, so it is not possible to read off the text when the event actually took place.



*How* questions are the least difficult questions. This can be explained with the fact that many *how* questions ask for numbers that are mentioned in the text (e.g. *How long did they stay in the sauna?* or *How many slices did they place onto the paper plate?*). The answer to such questions can often be found with a simple text lookup. Another part of *how* questions asks for the typical duration of an activity. These questions often have similar answers irrespective of the scenario, since most of the narrations in MCScript2.0 span a rather short time period. Such answers can easily be memorized by the models.

Especially for TriAN, *what* and *who* questions seem to be easy. This could be explained with the fact that ConceptNet contains lots of information about entities and their relations to each other, apparently also covering some information about script participants, which seems to be useful for these question types.

## 6 Related Work

Recent years have seen a number of datasets that evaluate commonsense inference. Like our corpus, most of these data sets choose a machine comprehension setting. The data sets can be roughly classified along their text domain:

**News Texts.** Two recently published machine comprehension data sets that require commonsense inference are based on news texts. First, *NewsQA* (Trischler et al., 2017) is a dataset of newswire texts from CNN with questions and answers written by crowdsourcing workers. During data collection, full texts were not shown to workers as a basis for question formulation, but only the text’s title and a short summary, to avoid literal repetitions and support the generation of non-trivial questions requiring background knowledge. Second, *ReCoRD* (Zhang et al., 2018) contains cloze-style questions on newswire texts that were not crowdsourced, but automatically extracted by pruning a named entity in a larger passage from the text.

**Web Texts.** Other corpora use web documents. An example is *TriviaQA* (Joshi et al., 2017), a corpus that contains automatically collected question-answer pairs from 14 trivia and quiz-league websites, together with web-crawled evidence documents from *Wikipedia* and *Bing*. While a majority of questions require world knowledge for finding the correct answer, it is mostly factual knowledge. *CommonsenseQA* (Talmor et al., 2018) con-

tains a total of over 9000 multiple-choice questions that were crowdsourced based on knowledge base triples extracted from ConceptNet. Texts were only added post-hoc, by querying a web search engine based on the formulated question, and by adding the retrieved evidence texts to the questions and answers.

**Fictional Texts.** *NarrativeQA* (Kočíský et al., 2018) is a reading comprehension dataset that largely differs from other corpora by means of text length. Instead of providing short reading texts, models have to process complete books or movie scripts and answer very complex questions.

Because of their domains, the aforementioned data sets evaluate a very broad notion of commonsense, including e.g. physical knowledge (for trivia texts) and knowledge about political facts (for newswire texts). However, none of them explicitly tackle script knowledge.

## 7 Conclusion

We presented MCScript2.0, a new machine comprehension dataset with a focus on challenging inference questions that require script knowledge or commonsense knowledge for finding the correct answer. Our new question collection procedure results in about half of the questions in MCScript2.0 requiring such inference, which is a much larger amount compared to a previous dataset.

We evaluate several benchmark models on MCScript2.0 and show that even a model that makes use of ConceptNet as a source for commonsense knowledge struggles to answer many question in our corpus. MCScript2.0 forms the basis of a shared task organized at the COIN workshop.<sup>8</sup>

## Acknowledgments

We thank our student assistants Leonie Harter, David Meier, Christine Schäfer and Georg Seiler for the help with data validation, and Kathryn Chapman, Srestha Ghosh, Trang Hoang, Ben Posner and Miriam Schulz for help with assessing the human upper bound. We also thank the numerous workers on MTurk for their good work and Carina Silberer and the reviewers for their helpful comments on the paper. This research was funded by the German Research Foundation (DFG) as part of SFB 1102 Information Density and Linguistic Encoding.

<sup>8</sup><https://coinnlp.github.io/>

## References

- Kevin Burton, Akshay Java, Ian Soboroff, et al. 2009. The ICWSM 2009 Spinn3r Dataset. In *Third Annual Conference on Weblogs and Social Media (ICWSM 2009)*.
- Nathanael Chambers and Dan Jurafsky. 2008. Unsupervised Learning of Narrative Event Chains. In *Proceedings of the 46th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 789–797. Association for Computational Linguistics.
- Nathanael Chambers and Dan Jurafsky. 2009. Unsupervised Learning of Narrative Schemas and their Participants. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP*, pages 602–610. Association for Computational Linguistics.
- Danqi Chen, Jason Bolton, and Christopher D. Manning. 2016. A Thorough Examination of the CNN/Daily Mail Reading Comprehension Task. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 1, pages 2358–2367.
- Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*.
- Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend. In *Advances in Neural Information Processing Systems*, pages 1693–1701.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.
- Bram Jans, Steven Bethard, Ivan Vulić, and Marie Francine Moens. 2012. Skip N-grams and Ranking Functions for Predicting Script Events. In *Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics*, EACL '12, pages 336–344. Association for Computational Linguistics.
- Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. 2017. TriviaQA: A Large Scale Distantly Supervised Challenge Dataset for Reading Comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*, pages 1601–1611. Association for Computational Linguistics.
- Tomáš Kočiský, Jonathan Schwarz, Phil Blunsom, Chris Dyer, Karl Moritz Hermann, Gábor Melis, and Edward Grefenstette. 2018. The NarrativeQA Reading Comprehension Challenge. *Transactions of the Association of Computational Linguistics*, 6:317–328.
- Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. 2017. RACE: Large-scale Reading Comprehension Dataset From Examinations. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 785–794.
- Vladimir I. Levenshtein. 1966. Binary codes capable of correcting deletions, insertions, and reversals. In *Soviet Physics Doklady*, volume 10, pages 707–710.
- Christopher D. Manning, Mihai Surdeanu, John Bauer, Jenny Rose Finkel, Steven Bethard, and David McClosky. 2014. The Stanford CoreNLP Natural Language Processing Toolkit. In *ACL (System Demonstrations)*, pages 55–60.
- Elizabeth M. Merkhofer, John Henderson, David Bloom, Laura Strickhart, and Guido Zarrella. 2018. MITRE at SemEval-2018 Task 11: Commonsense Reasoning without Commonsense Knowledge. In *Proceedings of the 12th International Workshop on Semantic Evaluations (SemEval-2018)*, pages 1078–1082.
- Ashutosh Modi, Tatjana Anikina, Simon Ostermann, and Manfred Pinkal. 2016. InScript: Narrative texts annotated with script information. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016)*, pages 3485–3493. European Language Resources Association (ELRA).
- Ashutosh Modi and Ivan Titov. 2014. Inducing Neural Models of Script Knowledge. In *Proceedings of the Conference on Computational Natural Language Learning (CoNLL)*, pages 49–57.
- Simon Ostermann, Ashutosh Modi, Michael Roth, Stefan Thater, and Manfred Pinkal. 2018a. MCScript: A Novel Dataset for Assessing Machine Comprehension Using Script Knowledge. In *Proceedings of the 11th International Conference on Language Resources and Evaluation (LREC 2018)*, pages 3567–3574.
- Simon Ostermann, Michael Roth, Ashutosh Modi, Stefan Thater, and Manfred Pinkal. 2018b. SemEval-2018 Task 11: Machine Comprehension using Commonsense Knowledge. In *Proceedings of The 12th International Workshop on Semantic Evaluation*, pages 747–757.
- Simon Ostermann, Michael Roth, Stefan Thater, and Manfred Pinkal. 2017. Aligning Script Events with Narrative Texts. *Proceedings of the 6th Joint Conference on Lexical and Computational Semantics (\*SEM 2017)*, pages 128–134.
- Karl Pichotta and Raymond J. Mooney. 2014. Statistical script learning with multi-argument events. In *Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics (EACL 2014)*, pages 220–229.

- Karl Pichotta and Raymond J. Mooney. 2016. Learning Statistical Scripts with LSTM Recurrent Neural Networks. *Proceedings of the 30th AAAI Conference on Artificial Intelligence (AAAI-16)*, pages 2800–2806.
- Michaela Regneri, Alexander Koller, and Manfred Pinkal. 2010. Learning Script Knowledge with Web Experiments. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, pages 979–988. Association for Computational Linguistics.
- Rachel Rudinger, Pushpendre Rastogi, Francis Ferraro, and Benjamin Van Durme. 2015. Script induction as language modeling. *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1681–1686.
- Roger C. Schank and Robert P. Abelson. 1977. *Scripts, Plans, Goals, and Understanding*. Lawrence Erlbaum Associates.
- Push Singh, Thomas Lin, Erik T. Mueller, Grace Lim, Travell Perkins, and Wan Li Zhu. 2002. Open Mind Common Sense: Knowledge Acquisition from the General Public. In *On the move to Meaningful Internet Systems 2002: CoopIS, DOA, and ODBASE*, pages 1223–1237. Springer.
- Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence (AAAI-17)*, pages 4444–4451.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2018. CommonsenseQA: A Question Answering Challenge Targeting Commonsense Knowledge. *arXiv preprint arXiv:1811.00937*.
- Adam Trischler, Tong Wang, Xingdi Yuan, Justin Harris, Alessandro Sordani, Philip Bachman, and Kaheer Suleman. 2017. NewsQA: A Machine Comprehension Dataset. In *Proceedings of the 2nd Workshop on Representation Learning for NLP*, pages 191–200.
- Liang Wang, Meng Sun, Wei Zhao, Kewei Shen, and Jingming Liu. 2018. Yuanfudao at SemEval-2018 Task 11: Three-way Attention and Relational Knowledge for Commonsense Machine Comprehension. In *Proceedings of the 12th International Workshop on Semantic Evaluations (SemEval-2018)*, pages 758–762.
- Lilian Wanzare, Alessandra Zarcone, Stefan Thater, and Manfred Pinkal. 2016. DeScript: A Crowdsourced Database for the Acquisition of High-quality Script Knowledge. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016)*, pages 3494–3501. European Language Resources Association (ELRA).
- Lilian Wanzare, Alessandra Zarcone, Stefan Thater, and Manfred Pinkal. 2017. Inducing Script Structure from Crowdsourced Event Descriptions via Semi-Supervised Clustering. In *Proceedings of the 2nd Workshop on Linking Models of Lexical, Sentential and Discourse-level Semantics*, pages 1–11. Association for Computational Linguistics.
- Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. 2018. SWAG: A Large-Scale Adversarial Dataset for Grounded Commonsense Inference. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 93–104.
- Sheng Zhang, Xiaodong Liu, Jingjing Liu, Jianfeng Gao, Kevin Duh, and Benjamin Van Durme. 2018. ReCoRD: Bridging the Gap between Human and Machine Commonsense Reading Comprehension. *arXiv preprint arXiv:1810.12885*.

## A Supplemental Material

### A.1 Additional Data Sample

- (7) **T:** I am at work . I have a guest sit at the bar . The ordered themselves a beer . I check that he is of age , and that his license is valid . I then go to the beer cooler , and grab a nice cold mug , and fill it up with beer . I place a napkin down and set the beer on top in front of the bar guest . I present him the check and tell him no rush , whenever he is ready . He then places his cash with the receipt . I go to cash him out , offer to be right back with his change , and he responds with , ” Keep the change ” . I like nights like this .

**Q:** Why did they receive a nice tip?  
the customer was happy with the service ✓  
the customer was in a rush ✗

**Q:** When was the check printed?  
after the order ✓  
before the order ✗

**Q:** What did they create at the computer and print?  
the check ✓  
change ✗

- (8) **T:** I wanted to throw a Bachelorette Party for my best friend . She lives in Dallas , but she wants to have her party in New Orleans for a girls weekend . The first thing we did was talk about the theme of the party . We decided on the theme of “ Something Blue ” . We would have all colors of blue and activities that have titles with the word blue for the whole weekend . She gave me a list of 20 girls . I created an invitation that had blue and included a picture of her . I also included an itinerary of our weekend activities with all of our fun “ blue ” titles , to set the fun mood . I sealed them before hand writing the addresses and adding a stamp . Next , they were off to the post office , so everyone could be invited to our fun weekend .

**Q:** What was printed out?  
itinerary ✓

invitations to a weddingy ✗

**Q:** When was each invitation placed into their blue envelope?

Before handwriting addresses ✓

After adding stamps ✗

**Q:** Where did she place the invitations?

Post office ✓

Dallas ✗

- (9) **T:** The restaurant was terrible again and I probably should not have given it another chance . The management at the store level is obviously not paying attention to me so it is time to right to headquarters . I opened the word processing program on my computer and opened a new document . I went all the way to the right side and entered my street address on one line and the city , state and zip code below that . Next I entered the date and then moved all the way to the left and entered the street address of the restaurant headquarters and the city , state and zip code of the headquarters . I started the letter with Dear Sir and on the next lines , proceeded to explain the problems I had been having with this particular location , it ’s service and food . I explained that I had tried to resolve it at store level but had been unsuccessful . On the final line , I went all the way to right and entered ‘ Sincerely ’ and hit return a couple times , then added my name below that . I folded it and put it in an addressed and stamped envelope , and mailed it to the company headquarters .

**Q:** What did they print out?

The letter ✓

The receipt from the restaurant ✗

**Q:** When did they sign above their printed name?

After the letter was printed ✓

After putting the letter in the envelope. ✗

**Q:** What they did they sign?

The letter ✓

The receipt at the restaurant ✗

### A.2 Crowdsourcing Details

All data were collected using Amazon Mechanical Turk<sup>9</sup> . We paid \$0.50 for each item in the text and answer collection experiment. For the question collection experiment, we paid \$0.50 per item, if the text contained 4 or more target sentences,

<sup>9</sup><https://www.mturk.com>



and \$0.30 per item if fewer target sentences were highlighted.

### A.3 Implementation Details

For implementation details and preprocessing of the logistic model and TriAN, we follow (Merkhofer et al., 2018) and (Wang et al., 2018), respectively. NLTK<sup>10</sup> was used as preprocessing tool for the Attentive Reader.

The learning rate was tuned to 0.002 and 0.1 for TriAN and the attentive reader, resp. and the hidden size for both models to 64. As in the original formulation, dropout was set to 0.5 for the attentive reader and to 0.4 for TriAN. Batch size was set to 32 and both models were trained for 50 epochs. The model with the best accuracy on the dev data was used for testing.

### A.4 Adversarial Filtering

**Data:** data set  $D$ , a randomly initialized assignment  $S$ , and a classifier  $C$

**Result:**  $\hat{S}$

**repeat**

split the data into test batches of size  $b$ , such that each batch contains all questions for  $b$  texts;

**for**  $D_{test}$  in batches **do**

$D_{train} \leftarrow D \setminus D_{test}$ ;

$\mathcal{D}_{train} \leftarrow compile(D_{train})$ ;

train  $C$  on  $\mathcal{D}_{train}$ ;

**for all instances**

$\langle T_i, Q_i, a_i^+, \langle a_{i,0}^-, \dots, a_{i,j}^- \rangle \rangle$  in  $D_{test}$

**do**

use  $C$  to classify all incorrect

answer candidates  $a_{i,0}^-, \dots, a_{i,j}^-$ ;

set  $s_{i,j}$  to the index of the answer

candidate with the highest

probability of being correct;

**end**

**end**

**until** number of changed assignments

stagnates or increases;

**Algorithm 1:** Adversarial Filtering for MC-Script2.0

Formally, let a dataset be defined as a list of tuples  $\langle t_i, q_i, a_i^+, \langle a_{i,0}^-, \dots, a_{i,j}^- \rangle \rangle$ , where  $t_i$  is a reading text,  $q_i$  is a question on the text,  $a_i^+$  is the correct answer (as selected via majority vote, s. last Section)

<sup>10</sup><https://www.nltk.org>

and  $\langle a_{i,0}^-, \dots, a_{i,j}^- \rangle$  is a list of 3 to 5 incorrect answer candidates<sup>11</sup>. The aim of the algorithm is to find an assignment  $\hat{S} = \{s_{0,0}, \dots, s_{i,j}\}$ , where each  $s_{i,j}$  is the index of the most difficult answer candidate in  $\langle a_{i,0}^-, \dots, a_{i,j}^- \rangle$ .

A dataset that is *compiled* with the assignment  $S$  is a list of instances  $\langle t_i, q_i, a_i^+, a_i^- \rangle$ , such that there is only one incorrect answer candidate per question, according to the indices given by  $S$ .

Once the algorithm converges,  $\hat{S}$  is used to compile the final version of the dataset,  $\hat{D}$ , which contains incorrect answer candidates that are most likely to be correct.

For the batch size we tried values in  $\{50, 100, 250, 500\}$ , but we found that for all values, the performance of the classifier would drop close to chance level after one iteration only. We set  $b = 250$ , since the performance was closest to chance after convergence with that setting. Also, we defined that the algorithm converges if the number of changed assignments since the last iteration is  $\leq 50$ .

<sup>11</sup>Note that since there are several questions per text, the value of  $t_i$  may appear in several instances.

## A.5 Screenshot of the VP-based Question Collection Experiment.

<p>packed up the tent and the tools from the last trip, we put everything into the car and set out on our way. After a four hour drive we arrived at the campgrounds. We found our spot and unpacked the car.</p>	
<p>First up was pitching the tent, then we set about making a fire to cook dinner.</p>	<p>Answer: <b>pitching the tent</b>          Question: _____  <b>When did they pitch the tent?</b>          Answer: <b>making a fire</b>          Question: _____  <b>When did they make a fire?</b>          Answer: <b>to cook dinner</b>          Question: _____  <b>How long did they cook dinner?</b></p>
<p>It took a while to collect all the wood we needed and get it roaring, but once we got the food on the fire we could finally relax. After soaking in nature and enjoying our campfire meal, night had fallen and it was time for bed. We curled up in our sleeping bags and told scary stories until we drifted off.</p>	
<p>The next morning we awoke early to stoke the fire.</p>	<p>Answer: <b>awoke early</b>          Question: _____  <b>When did they awake early?</b>          Answer: <b>stoke the fire</b>          Question: _____  <b>How long did they</b></p>
<p>We then started breakfast. A day of hiking and fishing went by far too fast and after another night under the stars we packed up our gear and headed home.</p>	

## A.6 Screenshot of the Instructions for the NP-based Question Collection Experiment.

**Instructions**

In this experiment, we collect reading comprehension questions for children, based on stories about everyday activities. We want to test if the children understand certain parts of the text, so we are looking for questions about specific phrases in the text.

**Your Task:**

Please, first read the full text that appears on the left-hand side. Then click the **Continue** button. A number of text fields will appear on the right-hand side, next to a highlighted sentence. Please fill each text field with a question that can be answered with the given phrase, in the context of the highlighted sentence. Then click **Continue**, and the text fields for the second highlighted sentence will appear.

Note that relevant information from the text should be repeated in the question, so that the children have a chance to find an answer. In the example, *What did they make?* would be a bad question, because it is too general, while *What did they light to make dinner?* is specific enough.

**Important remarks:**

1. Please stick to the following rules, otherwise we reserve the right to reject your work:
  - Do not write yes/no questions. We will check for these questions automatically.
  - Do not write questions that are shorter than 4 words. We will check for these questions automatically.
  - If persons are mentioned by name ("John", "Michael") or function ("the waiter", "the guest"), please refer to them in the same way in the question. If you want to refer to the protagonist (usually "I", "we"), always use "they". Do not use I/you/we. We will check for these pronouns automatically.
  - Please write the questions in grammatical and coherent English. We will manually check a major number of HITS.
2. When you start the experiment, there will be no *Submit* button, but don't worry! You just need to go through the 2 sentences, and the *Submit* button will appear right afterwards. Once you finished one sentence, you can no longer edit your questions.
3. You can access both these instructions and the examples below during the experiment by clicking on one of the two buttons at the top of the experiment page.

**Example**

We spent a week getting our gear together for our big camping trip. After making sure everything was packed and the tent had no holes from the last trip , we put everything into the car and set out on our way. After a four hour drive we arrived at the campgrounds. We found our spot and unpacked the car.		
<b>First up was pitching the tent , then we set about making a fire to cook dinner.</b>	Ask for <b>the tent!</b>	What did they pitch after arriving?
	Ask for <b>a fire!</b>	What did they light to make dinner?
It took a while to collect all the wood we needed and get it roaring, but once we got the food on the fire we could finally relax. After soaking in nature and enjoying our campfire meal, night had fallen and it was time for bed. We curled up in our sleeping bags and told scary stories until we drifted off.		
<b>The next morning we awoke early to stoke the fire.</b>	Ask for <b>the next morning!</b>	When did they wake up?
	Ask for <b>the fire!</b>	What did they stake in the morning?
We then started breakfast. A day of hiking and fishing went by far too fast and after another night under the stars we packed up our gear and headed home.		

Got it. Start the experiment!