

## 601 A Datasheet for SITUATEDGEN

### 602 A.1 Motivation

603 **1. For what purpose was the dataset created? Was there a specific task in mind? Was there a**  
604 **specific gap that needed to be filled? Please provide a description.**

605 This dataset aims to probe the commonsense reasoning ability of generative language models through  
606 the lens of keyword generation tasks. The task requires machines to compose a pair of contrastive  
607 sentences with a given set of keywords containing geographical or temporal entities. Current models  
608 lack the ability to correctly reason for the relationship among these entities and thus generate sentences  
609 that contradict commonsense knowledge. We hope our dataset could stir more research to fill this gap  
610 of generative commonsense reasoning.

611 **2. Who created the dataset (e.g., which team, research group) and on behalf of which entity**  
612 **(e.g., company, institution, organization)?**

613 The dataset is created by Yunxiang Zhang and Xiaojun Wan on behalf of the Text Mining and  
614 Linguistic Computing Group, Wangxuan Institute of Computer Technology, Peking University. Most  
615 part of this paper is done when the first author is at Peking University before moving to University of  
616 Michigan.

617 **3. Who funded the creation of the dataset? If there is an associated grant, please provide the**  
618 **name of the grantor and the grant name and number.**

619 It is funded by the Text Mining and Linguistic Computing Group, Wangxuan Institute of Computer  
620 Technology, Peking University.

### 621 A.2 Composition

622 **1. What do the instances that comprise the dataset represent (e.g., documents, photos, people,**  
623 **countries)? Are there multiple types of instances (e.g., movies, users, and ratings; people and**  
624 **interactions between them; nodes and edges)? Please provide a description.**

625 The dataset is comprised of pure text data in English, presented in a Jsonline format. The file is  
626 composed of a list of instances containing input keywords and targeted outputs.

627 **2. How many instances are there in total (of each type, if appropriate)?**

628 Our dataset consists of 8,268 instances. Please refer to Table 2 for detailed information.

629 **3. Does the dataset contain all possible instances or is it a sample (not necessarily random)**  
630 **of instances from a larger set? If the dataset is a sample, then what is the larger set? Is the**  
631 **sample representative of the larger set (e.g., geographic coverage)? If so, please describe how**  
632 **this representativeness was validated/verified. If it is not representative of the larger set, please**  
633 **describe why not (e.g., to cover a more diverse range of instances, because instances were**  
634 **withheld or unavailable).**

635 This dataset does not cover all aspects of commonsense knowledge so it does not contain all possible  
636 instances. We focus on geographical and temporal commonsense in this work since they provide  
637 testbeds for evaluating machines' reasoning ability under different extra-linguistic contexts.

638 **4. What data does each instance consist of? "Raw" data (e.g., unprocessed text or images) or**  
639 **features? In either case, please provide a description.**

640 Each instance is a dictionary has the following fields:

- 641 • "keywords": a list of keywords as input
- 642 • "statement": a string concatenation of two sentences as the target generations

- 643 • “ids”: the origins of the commonsense statement (from which (train/dev/test) split of which  
644 source datasets/corpora) represented in the format of “{src\_dataset}::{split}::{id}”
- 645 • “statements”: a list of the two sentences in “statement” field.

646 **5. Is there a label or target associated with each instance? If so, please provide a description.**

647 Yes. It is represented as the “statement” filed in each instance.

648 **6. Is any information missing from individual instances? If so, please provide a description,  
649 explaining why this information is missing (e.g., because it was unavailable). This does not  
650 include intentionally removed information, but might include, e.g., redacted text.**

651 No. All instances are complete.

652 **7. Are relationships between individual instances made explicit (e.g., users’ movie ratings, social  
653 network links)? If so, please describe how these relationships are made explicit.**

654 Individual instances are independent of each other. The train/dev/test splits do not overlap in any  
655 single sentence.

656 **8. Are there recommended data splits (e.g., training, development/validation, testing)? If so,  
657 please provide a description of these splits, explaining the rationale behind them.**

658 Yes, see Tabel 2 for details. The splitting process makes sure that the train/dev/test splits do not  
659 overlap in any single sentence. See Appendix D.3 for details.

660 **9. Are there any errors, sources of noise, or redundancies in the dataset? If so, please provide a  
661 description.**

662 There is noise in the train and dev set. We manually filter out unqualified examples in the test set.  
663 See more analysis in Section 5.1 and Appendix E.2.

664 **10. Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g.,  
665 websites, tweets, other datasets)? If it links to or relies on external resources, a) are there  
666 guarantees that they will exist, and remain constant, over time; b) are there official archival  
667 versions of the complete dataset (i.e., including the external resources as they existed at the time  
668 the dataset was created); c) are there any restrictions (e.g., licenses, fees) associated with any of  
669 the external resources that might apply to a dataset consumer? Please provide descriptions of  
670 all external resources and any restrictions associated with them, as well as links or other access  
671 points, as appropriate.**

672 The SITUATEDGEN dataset is self-contained and we welcome practitioners to consider additional  
673 knowledge sources.

674 **11. Does the dataset contain data that might be considered confidential (e.g., data that is  
675 protected by legal privilege or by doctor–patient confidentiality, data that includes the content  
676 of individuals’ non-public communications)? If so, please provide a description.**

677 No.

678 **12. Does the dataset contain data that, if viewed directly, might be offensive, insulting, threaten-  
679 ing, or might otherwise cause anxiety? If so, please describe why.**

680 No.

### 681 A.3 Collection Process

682 **1. How was the data associated with each instance acquired? Was the data directly observable  
683 (e.g., raw text, movie ratings), reported by subjects (e.g., survey responses), or indirectly  
684 inferred/derived from other data (e.g., part-of-speech tags, model-based guesses for age or  
685 language)? If the data was reported by subjects or indirectly inferred/derived from other data,  
686 was the data validated/verified? If so, please describe how.**

687 The data is sourced from several commonsense related datasets and corpora. We design an automatic  
688 pipeline to convert and filter data into our desired format.

689 **2. What mechanisms or procedures were used to collect the data (e.g., hardware apparatuses  
690 or sensors, manual human curation, software programs, software APIs)? How were these  
691 mechanisms or procedures validated?**

692 We first convert instances from other datasets as commonsense statements. Then we match these  
693 statements into pairs and extract keywords from them. We further manually filter out invalid examples  
694 in the test set.

695 **3. Who was involved in the data collection process (e.g., students, crowdworkers, contractors)  
696 and how were they compensated (e.g., how much were crowdworkers paid)?**

697 We hired crowdworkers and compensated them with 0.1 yuan for each entry they checked, which is  
698 higher than the statutory minimum wage.

699 **4. Over what timeframe was the data collected? Does this timeframe match the creation  
700 timeframe of the data associated with the instances (e.g., recent crawl of old news articles)? If  
701 not, please describe the time frame in which the data associated with the instances was created.**

702 Our dataset was built in 2022 while the original source data is published between 2018-2021. Usually,  
703 commonsense statements are not changing over time.

704 **5. Were any ethical review processes conducted (e.g., by an institutional review board)? If so,  
705 please provide a description of these review processes, including the outcomes, as well as a link  
706 or other access point to any supporting documentation.**

707 No.

#### 708 **A.4 Preprocessing/cleaning/labeling**

709 **1. Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing,  
710 tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing  
711 of missing values)? If so, please provide a description. If not, you may skip the remaining  
712 questions in this section.**

713 Yes. We use templated-based and neural-based models to convert and filter the source data into our  
714 desired format. See details in Appendix D.

715 **2. Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to  
716 support unanticipated future uses)? If so, please provide a link or other access point to the  
717 “raw” data.** Yes. The raw data is available on the corresponding dataset websites (CREAK –  
718 <https://github.com/yasumasaonoe/creak>, OpenbookQA – [https://allenai.org/data/  
719 open-book-qa](https://allenai.org/data/open-book-qa), StrategyQA – <https://allenai.org/data/strategyqa>, CommonsenseQA –  
720 <https://www.tau-nlp.sites.tau.ac.il/commonsenseqa>, ARC – [https://allenai.org/  
721 data/arc](https://allenai.org/data/arc)).

722 **3. Is the software that was used to preprocess/clean/label the data available? If so, please  
723 provide a link or other access point.**

724 Yes. Please see [https://github.com/yunx-z/situated\\_gen](https://github.com/yunx-z/situated_gen).

#### 725 **A.5 Uses**

726 **1. Has the dataset been used for any tasks already? If so, please provide a description.**

727 Not yet.

728 **2. Is there a repository that links to any or all papers or systems that use the dataset? If so,  
729 please provide a link or other access point.**

730 There has not been such a paper or system yet.

731 **3. What (other) tasks could the dataset be used for?**

732 It can be used to develop better language models for commonsense reasoning. It can be used to  
733 evaluate language models, especially their understanding of commonsense knowledge. It could  
734 potentially benefit many downstream applications such as document summarization [44], story  
735 writing [51] and dialogue response generation [31].

736 **4. Is there anything about the composition of the dataset or the way it was collected and**  
737 **preprocessed/cleaned/labeled that might impact future uses? For example, is there anything**  
738 **that a dataset consumer might need to know to avoid uses that could result in unfair treatment**  
739 **of individuals or groups (e.g., stereotyping, quality of service issues) or other risks or harms**  
740 **(e.g., legal risks, financial harms)? If so, please provide a description. Is there anything a**  
741 **dataset consumer could do to mitigate these risks or harms?**

742 The dataset has very low risks of containing sentences with toxicity and offensiveness. Since our data  
743 is sourced from existing datasets, we may inherit geographical biases [16] that result in an uneven  
744 distribution of commonsense knowledge about western and non-western regions. The commonsense  
745 statements may not sound familiar to people who live in locations that are poorly represented in the  
746 source datasets. Therefore, models developed on our dataset may preserve biases learned from the  
747 annotators of the source datasets. We note that pretrained language models may also inherit the bias  
748 in the massive pretraining data. It is important that interested parties carefully address those biases  
749 before deploying the model to real-world settings.

750 **5. Are there tasks for which the dataset should not be used? If so, please provide a description.**

751 The dataset can only be used for research purposes.

752 **A.6 Distribution**

753 **1. Will the dataset be distributed to third parties outside of the entity (e.g., company, institution,**  
754 **organization) on behalf of which the dataset was created? If so, please provide a description.**

755 The dataset is already publicly available.

756 **2. How will the dataset will be distributed (e.g., tarball on website, API, GitHub)? Does the**  
757 **dataset have a digital object identifier (DOI)?**

758 The dataset is available at [https://github.com/yunx-z/situated\\_gen](https://github.com/yunx-z/situated_gen).

759 **3. When will the dataset be distributed?**

760 It has already been distributed.

761 **4. Will the dataset be distributed under a copyright or other intellectual property (IP) license,**  
762 **and/or under applicable terms of use(ToU)? If so, please describe this license and/or ToU, and**  
763 **provide a link or other access point to, or otherwise reproduce, any relevant licensing terms or**  
764 **ToU, as well as any fees associated with these restrictions.**

765 This dataset is licensed under the Creative Commons Attribution-NonCommercial-ShareAlike 4.0  
766 International License (CC BY-NC-SA 4.0). The full text of the license can be accessed at the  
767 following link: <https://creativecommons.org/licenses/by-nc-sa/4.0/>.

768 **5. Have any third parties imposed IP-based or other restrictions on the data associated with the**  
769 **instances? If so, please describe these restrictions, and provide a link or other access point to,**  
770 **or otherwise reproduce, any relevant licensing terms, as well as any fees associated with these**  
771 **restrictions.**

772 No.

773 **6. Do any export controls or other regulatory restrictions apply to the dataset or to individual**  
774 **instances? If so, please describe these restrictions, and provide a link or other access point to,**  
775 **or otherwise reproduce, any supporting documentation.**

776 No.

## 777 **A.7 Maintenance**

### 778 **1. Who will be supporting/hosting/maintaining the dataset?**

779 The first author, Yunxiang Zhang, is hosting and maintaining the dataset.

### 780 **2. How can the owner/curator/manager of the dataset be contacted (e.g., email address)?**

781 Email: yunxiang@umich.edu

### 782 **3. Is there an erratum? If so, please provide a link or other access point.**

783 No.

### 784 **4. Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete instances)? If so, please describe how often, by whom, and how updates will be communicated to dataset consumers (e.g., mailing list, GitHub)?**

787 We are interested to collect more data using our automatic pipelines and conduct manual filtering as  
788 future work. We also welcome interested parties to point out errors in the dataset via contact email or  
789 github issues so we could correct them. If there is a plan for systematic updates, we will announce it  
790 at the earliest opportunity.

### 791 **5. If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so? If so, please provide a description. Will these contributions be validated/verified? If so, please describe how. If not, why not? Is there a process for communicating/distributing these contributions to dataset consumers? If so, please provide a description.**

795 People can use this repository following the licenses and cite our paper.

## 796 **B Limitations**

797 Since our dataset is derived from existing commonsense benchmarks, we may inherit their annotation  
798 artifacts [18] and contain certain types of spurious lexical patterns (e.g., “A lived in B”). We could  
799 also conduct an extra manual evaluation on the machine generations, so as to gauge its correlation  
800 with automatic metrics, though this has been verified by [25] on the original generative commonsense  
801 reasoning task. Recently, a lot of work has developed new retrieval-augmented commonsense  
802 text generation models [54, 19], which could also be included as baseline models for a more  
803 comprehensive benchmark.

## 804 **C Ethics Statement**

805 Our data is built upon publicly available datasets and we will follow their licenses when releasing our  
806 data. There is no explicit detail that leaks an annotator’s personal information. The dataset has very  
807 low risks of containing sentences with toxicity and offensiveness. Since our data is sourced from  
808 existing datasets, we may inherit geographical biases [16] that result in an uneven distribution of  
809 commonsense knowledge about western and non-western regions. The commonsense statements may  
810 not sound familiar to people who live in locations that are poorly represented in the source datasets.  
811 Therefore, models developed on our dataset may preserve biases learned from the annotators of the  
812 source datasets. We note that pretrained language models may also inherit the bias in the massive  
813 pretraining data. It is important that interested parties carefully address those biases before deploying  
814 the model to real-world settings.

Table 5: Source dataset examples. **Correct answers** are in bold and underlined.

Dataset	Size	Format	Raw Data → Statement Conversion Example
CREAK [34]	13,418	True/False statement	In the calendar year, May comes after April and before June. ( <b><u>True</u></b> /False) → In the calendar year, May comes after April and before June.
StrategyQA [17]	5,111	Yes/No Question	Are more watermelons grown in Texas than in Antarctica? ( <b><u>Yes</u></b> /No) → More watermelons are grown in Texas than in Antarctica.
CommonsenseQA [46]	12,247	Multiple-choice Question	Where in Southern Europe would you find many canals? (A) Michigan (B) New York (C) Amsterdam ( <b><u>D</u></b> ) Venice (E) Sydney → You would find many canals in Venice, Southern Europe.
ARC [11]	7,787	Multiple-choice Question	How long does it take for Earth to rotate on its axis seven times? (A) one day ( <b><u>B</u></b> ) one week (C) one month (D) one year → It takes one week for Earth to rotate on its axis seven times.
OpenbookQA [30]	6,493	Commonsense Statement	You wear shorts in the summer. → You wear shorts in the summer.

## 815 D Additional Details of Dataset Collection

### 816 D.1 Commonsense Statement Collection

817 We briefly introduce the nature of each source dataset in Section 4.1.

- 818 • **CREAK** [34] is a commonsense fact verification dataset featuring entity commonsense,  
819 which includes 13,418 true or false statements about entity knowledge written by crowd-  
820 workers.
- 821 • **StrategyQA** [17] is a commonsense question answering dataset that requires multi-hop  
822 implicit reasoning. It consists of 5,111 questions whose answers are either Yes or No.  
823 Machines need to decompose a question into multiple atomic questions to arrive at an  
824 answer.
- 825 • **CommonsenseQA** [46] is a commonsense question answering dataset of 12,247 five-way  
826 multiple-choice questions with a focus on knowledge in everyday life.
- 827 • **ARC** [11] is a commonsense question answering dataset. It has 7,787 four-way multiple-  
828 choice natural science questions collected from grade-school standardized tests.
- 829 • **OpenbookQA** [30] is a commonsense question answering dataset that simulates openbook  
830 test. The data set is made up of 5,957 multiple-choice questions, accompanied by 6,493  
831 commonsense statements about science facts. Since there is a significant overlap between  
832 the knowledge in questions and statements, we only use the statements data for simplicity.

833 We now detail the specific preprocessing method for each source dataset to convert them (i.e.,  
834 question-answer pairs) into statements.

- 835 • If the raw data comes in the statement format (CREAK and OpenbookQA), we obtain the  
836 true statements (part of CREAK and all of OpenbookQA) without extra processing.
- 837 • If the raw data comes in Yes/No question format (StrategyQA), we leverage a POS-rule-  
838 based open-sourced system `question_to_statement`<sup>8</sup> to transform a pair of question and  
839 Yes/No answer into a statement.

<sup>8</sup>[https://github.com/SunnyWay/question\\_to\\_statement](https://github.com/SunnyWay/question_to_statement)

840 • If the raw data comes in multiple-choice format (CommonsenseQA and ARC), we utilize  
841 a neural model to convert a pair of question and correct choice  $(q, a)$  into a statement in a  
842 sequence-to-sequence fashion. Concretely, we use the QA-to-statement model checkpoint  
843 released by [36], which is a BART [22] model finetuned on QA2D [13], a dataset of  
844 human-annotated statements for QA pairs.

845 Converting QA pair to statement is not a difficult task for pretrained seq2seq models. We observe that  
846 the generated statements are mostly fluent and faithful to the input. Additionally, we have manually  
847 filtered out unnatural examples in the test set. We summarize the basic information of these datasets  
848 and provide an example of statement conversion for each dataset in Table 5.

## 849 D.2 Antithesis Mining

850 **Keyword Masking.** We use entities and other nouns as the keywords of sentences because as  
851 a pilot study, we only consider the relationships between spatio-temporal contexts and nouns and  
852 ignore the influence of other part-of-speech categories such as verbs, adjectives, and prepositions.  
853 We use the same NER tagger in Section 4.2 to extract entities. We leverage spaCy<sup>9</sup> to extract all the  
854 nouns (including proper nouns) from a sentence. We merge the entities and nouns as keywords after  
855 removing duplicates. In particular, if a noun and an entity partly overlap (e.g., “month” and “a lunar  
856 month”), we retain the entity when deduplicating.

857 **Masked Sentence Similarity Matching.** We use the pretrained language model  
858 all-MiniLM-L6-v2<sup>10</sup> released by SentenceTransformers [41] to obtain high-quality embed-  
859 dings of keyword-masked sentences. We calculate the cosine similarity to pair highly similar masked  
860 sentences. Computing the similarity of all possible sentence pairs requires  $\mathcal{O}(n^2)$  time complexity.  
861 To accelerate this process, we use the paraphrase\_mining API of SentenceTransformers [41].

862 **Rule-based Filtering.** We devise the following rules to filter invalid sentence pairs based on  
863 iterative observation of the data:

- 864 • The masked sentence similarity exceeds a certain threshold<sup>11</sup>, which indicates parallel  
865 sentence structure of antithesis.
- 866 • The number of masked keywords ([UNK]) of every single sentence should not be more  
867 than 5 and less than 2, which controls for a reasonable difficulty of the keyword-to-text  
868 generation task.
- 869 • Any entity in one sentence does not appear in the other sentence within a pair (including  
870 the deformation of entity words, such as singular/plural form, upper/lower case, etc.). This  
871 is to avoid both sentences expressing the information of the same entity, while contrastive  
872 sentences should describe two opposite things.
- 873 • Both of the two sentences contain either GEO entities or TEMP entities (GEO+GEO or  
874 TEMP+TEMP), which avoids sentences comparing GEO context to a non-parallel TEMP  
875 context (GEO+TEMP).

## 876 D.3 Dataset Splitting

877 We treat dataset splitting as a community structure [7] discovery problem. Community structure  
878 refers to a group of tightly connected nodes that have a high density of internal connections and  
879 a low density of external connections. We regard a single sentence as a node in the graph. If two  
880 single sentences can be matched into a pair of contrastive sentences, an undirected edge will connect  
881 the corresponding nodes of these two single sentences. In this way, we obtain an undirected graph

---

<sup>9</sup>[https://spacy.io/models/en#en\\_core\\_web\\_sm](https://spacy.io/models/en#en_core_web_sm)

<sup>10</sup><https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2>

<sup>11</sup>We set the threshold as 0.8 via manual inspection.

882 describing the dataset structure. A subset of a dataset (such as a training set) is equivalent to a  
883 subgraph containing all sentence pairs (edges) and single sentences (nodes) of that subset.

884 In order to prevent the same sentence from appearing across different sets, we require that the  
885 subgraph node sets of the training set, validation set, and test set are disjoint. We use a community  
886 structure detection algorithm to meet this requirement. We use the community as the basic unit  
887 of dataset splitting, putting all the edges (sentence pairs) in one community into a certain dataset  
888 split. Connecting edges between communities (two vertices belonging to different communities) are  
889 removed. We note that sentences with similar syntactic structures tend to be connected to each other  
890 in the graph and thus fall into the same community, which ensures the syntactic variability between  
891 train/dev/test splits.

892 We use the Louvain [7] community structure detection algorithm<sup>12</sup> and divide our graph into 79  
893 communities. The largest community contains 3,273 edges, accounting for about 26% of the total  
894 data. We remove edges connecting different communities and then randomly divide the communities  
895 of contrastive sentence pairs into training set, validation set or test set.

## 896 E Dataset Quality Analysis

### 897 E.1 Manual Filtering of the Test Set

898 To ensure the high quality of the dataset, we manually filter out invalid examples in the test set that  
899 are not fluent antitheses or context-dependent. This process is important for the very high human  
900 performance shown in Table 3. Table 6 shows the instructions for annotators. We first ask two  
901 graduate students with proficiency in English to annotate 100 examples as valid or invalid. They agree  
902 with each other (i.e., give the same label) on 88% of examples. The inter-annotator agreement in  
903 terms of Cohen’s Kappa [12] is 0.76, which indicates substantial agreement [21]. Since the agreement  
904 ratio is satisfactory, we ask one of the annotators to complete the rest of the filtering process.

### 905 E.2 Error Cases Analysis

906 In Section 5.1, we annotate 100 random examples for whether it is actually 1) (fluent) antithesis  
907 and 2) context-dependent. Below, we analyze the bad cases in detail, including non-contrastive and  
908 non-context-dependent sentence pairs.

909 The main explanation that accounts for the production of non-contrastive sentence pair is that the  
910 remaining verbs after keyword masking may have lexical ambiguity, e.g. “play” in “*Slaves **play** a role  
911 in the history of the united states.*” and “*A team sport **played** mostly in Canada is Lacrosse.*” Although  
912 the pretrained language models could infer the meaning of a word according to its context [14], the  
913 contexts are lost after keyword masking. As a result, two sentences with different syntactic structures  
914 are matched together, thus violating the antithesis rule. This poses a limitation of our antithesis  
915 mining algorithm.

916 In addition, 7% of the sentence pairs are antitheses yet not context-dependent. Take the following  
917 sentence pair as an example: “*You could find millions of brownstone in New York City.*”<sup>13</sup> “*One can  
918 find a Holiday Inn inside the United States.*”. After swapping the GEO entity “New York City” and  
919 “United States” in these two sentences, they still conform to commonsense. The reason for this  
920 phenomenon is that New York City is part of the United States, and thus the “brownstone” related  
921 to New York will also be related to the United States. However, we would like to point out that  
922 contextual dependence is not an absolutely strict condition. Although this example still holds after  
923 swapping the GEO entities, it is not the optimal answer, because “brownstone” is more a typical thing  
924 in New York City and thus more suitable for a match with “New York City”.

---

<sup>12</sup><https://github.com/shobbrook/communities>

<sup>13</sup>As background knowledge, there are many historical buildings in New York City whose facades are made of brown sandstone, see <https://bungalow.com/articles/what-exactly-is-a-brownstone>.

Table 6: Annotator instructions for manual filtering of our dataset.

**Goal:** The objective of our project is to generate high-quality contrastive sentence pairs (antithesis) that incorporate geographical and temporal contexts. These sentence pairs will serve as a means to evaluate machines’ commonsense reasoning abilities under different extra-linguistic contexts. We aim to create sentences that require a deep understanding of real-world geographical and temporal entities but can be reasonably confirmed without resorting to external sources like Google or Wikipedia.

**Instructions:** We show a set of keywords and a pair of sentences containing these keywords. Your task is to determine whether this sentence pair satisfies *all* of the following criteria:

1. The sentence pair includes all of the given keywords.
2. Each sentence has at least one entity related to geography or time.
3. Each sentence is fluent and adheres to commonsense knowledge.
4. The two sentences have similar syntactic structures and create a contradiction in semantics.
  - Intuitively, the qualifying two sentences can be connected into a coherent sentence via a conjunction word such as “while”, “yet”, and “whereas” (e.g., “*July is summer in the United States, while July is winter in Australia.*”).
5. Swapping any of the geographical or temporal entities between the two sentences could lead to a contradiction with commonsense yet grammatical correctness.
  - For example, for the sentence pair “*July is summer in China. July is winter in Australia.*”, if the two geographical entities “China” and “Australia” are swapped, the resulting sentences do not adhere to commonsense anymore: “*July is summer in Australia. July is winter in China.*”

**Examples:**

Keywords: *morning, night, sunrise, sunset*

Sentence 1: "The sky is bright with the sunrise in the early morning."

Sentence 2: "The sky is dark with the sunset in the late night."

Criterion 1: Both sentences include the keywords "morning" and "night."

Criterion 2: Each sentence contains a geographical or temporal entity ("sunrise" and "sunset") related to the context.

Criterion 3: Both sentences are fluent and adhere to commonsense knowledge.

Criterion 4: The sentences have a similar syntactic structure and create a semantic contradiction: "The sky is bright with the sunrise in the early morning, while the sky is dark with the sunset in the late night."

Criterion 5: Swapping the temporal entities "early morning" and "late night" would result in a contradiction: "The sky is bright with the sunrise in the late night, while the sky is dark with the sunset in the early morning."

This example demonstrates how the sentence pairs satisfy the specified criteria of the task.

---

## 925 F Experimental Setup

### 926 F.1 Baseline Models

927 We use HuggingFace [50] implementations of the BART and T5 models. For the decoding method,  
928 we adopt the standard beam search with a beam size of 4 for all baseline models. As for checkpoint  
929 selection, we save a checkpoint for each epoch and select the checkpoint with the highest ROUGE-2  
930 on the validation set. Other default hyperparameters are shown in Table 7.

931 Table 8 shows an example of GPT prompt format, consisting of a fixed instruction (“*Generate a*  
932 *pair of contrastive sentences with the given set of keywords.*”) and a few in-context demonstrations  
933 (“*Keywords:  $c_1, \dots, c_k \setminus n$  Sentences:  $s_1 s_2$* ”).

### 934 F.2 Evaluation Metrics

935 We use the standard implementation of BLEU, ROUGE, METEOR, CIDEr, and SPICE in  
936 `pycocoevalcap`<sup>14</sup>. As recommended, we adopt the Recall score of BERTScore<sup>15</sup> and the hash code  
937 for evaluation setting is “`roberta-large_L17_no-idf_version=0.3.12(hug_trans=4.21.3)-rescaled_fast-`

<sup>14</sup><https://github.com/salaniz/pycocoevalcap>

<sup>15</sup>[https://github.com/Tiiiger/bert\\_score](https://github.com/Tiiiger/bert_score)

Table 7: Hyper-parameter settings for all baseline models.

Parameter	Value
epoch	10
batch size	32
beam size	4
max input length	64
max output length	128
learning rate	3e-5
warm-up steps	500

Table 8: An example of InstructGPT prompt format. We only show two in-context demonstrations here for brevity.

---

Generate a pair of contrastive sentences with the given set of keywords.

Keywords: Kansas, steakhouses, New York City, city, pizzerias

Sentences: Kansas city is known for its steakhouses. New York City is known for its pizzerias.

...

Keywords: seven days, one day, 1,440 minutes, a week

Sentences: There are 1,440 minutes in one day. There are seven days a week.

Keywords: axis, one day, one month, Earth, Moon

Sentences:

---

tokenizer”. In addition, we design and implement MATCH to evaluate how well the machines solve the challenge of situated semantic matching (Section 3.2). We now define the keyword matching accuracy MATCH based on mathematical notations introduced in Section 3.1.

$t = (t_1, \dots, t_k)$ ,  $t_i \in \{0, 1\}$  indicates that each keyword  $c_i$  appears in which sentence in the answer pair  $y^{true} = \{s_1^{true}, s_2^{true}\}$ . In other words, if  $c_i$  *should* appear in  $s_1$ , then  $t_i = 0$ ; if  $c_i$  *should* appear in  $s_2$ , then  $t_i = 1$ .  $p = (p_1, \dots, p_k)$ ,  $p_i \in \{-1, 0, 1\}$  indicates that each keyword  $c_i$  appears in which sentence in the output pair  $y^{pred} = \{s_1^{pred}, s_2^{pred}\}$ . In other words, if  $c_i$  *actually* appear in  $s_1$ , then  $p_i = 0$ ; if  $c_i$  *actually* appear in  $s_2$ , then  $p_i = 1$ ; if  $c_i$  does not *actually* appear in both  $s_1$  and  $s_2$ , then  $p_i = -1$ <sup>16</sup>. We define the matching accuracy of a sentence pair  $\text{match}(y^{true}, y^{pred})$  as the proportion of correctly matched keywords, which is calculated as  $\frac{1}{k} \max(\sum_{i=1}^k \mathbb{1}_{t_i=p_i}, \sum_{i=1}^k \mathbb{1}_{1-t_i=p_i}) \in [0, 1]$ . Here  $\mathbb{1}$  is the indicator function. The formula includes both  $1 - t$  and  $t$  in a symmetric way because the sentence pair is unordered. For the whole test set, we take the average matching accuracy of all examples as MATCH.

We illustrate the computing process of matching accuracy with a simple example. Given [July, China, winter, Australia, summer, July], the answer could be “*July is summer in China. July is winter in Australia.*” So  $t = (0, 0, 1, 1, 0, 1)$ . If the generated output is “*July is summer in Australia. July is winter in China.*”, then  $p = (0, 1, 1, 0, 0, 1)$ . As a result, the matching accuracy is  $4/6 = 0.67$ .

As for the implementation, we utilize NLTK<sup>17</sup> to split the output into two sentences. In particular, if there is only one sentence in the output, we append an empty string as the second one; if there are more than two sentences, we only take the former two sentences into consideration. We lemmatize the sentence before determining keyword appearance.

---

<sup>16</sup>By defining  $p_i = -1$ , MATCH can also reflect the coverage of keywords in the output.

<sup>17</sup><https://www.nltk.org/>