# Dataset vs Reality: Understanding Model Performance from the Perspective of Information Need

Mengying Yu

School of Computer Science and Engineering, Nanyang Technological University
50 Nanyang Avenue, Singapore 639798
`yume0004@e.ntu.edu.sg`

Aixin Sun*

School of Computer Science and Engineering, Nanyang Technological University
50 Nanyang Avenue, Singapore 639798
`axsun@ntu.edu.sg`

**Abstract**

Deep learning technologies have brought us many models that outperform human beings on a few benchmarks. An interesting question is: *can these models well solve real-world problems with similar settings (*e.g., *identical input/output) to the benchmark datasets?* We argue that a model is trained to answer the *same information need* for which the training dataset is created. Although some datasets may share high structural similarities, *e.g.,* question-answer pairs for the question answering (QA) task and image-caption pairs for the image captioning (IC) task, they may represent different research tasks aiming for answering different information needs. To support our argument, we use the QA task and IC task as two case studies and compare their widely used benchmark datasets. From the perspective of *information need* in the context of information retrieval, we show the differences in the dataset creation processes, and the differences in morphosyntactic properties between datasets. The differences in these datasets can be attributed to the different information needs of the specific research tasks. We encourage all researchers to consider the information need the perspective of a research task before utilizing a dataset to train a model. Likewise, while creating a dataset, researchers may also incorporate the information need perspective as a factor to determine the degree to which the dataset accurately reflects the research task they intend to tackle.

## 1 Introduction

In the very first chapter of the Information Retrieval (IR) book, Manning *et al.* distinguish **information need** from **query**: "An information need is the topic about which the user desires to know more, and a query is what the user conveys to the computer in an attempt to communicate the information need" (Manning, Raghavan, & Schütze, 2008). When a query is entered into a search engine, the latter provides a list of potential documents that may be relevant to the user's search. The user then evaluates each document to
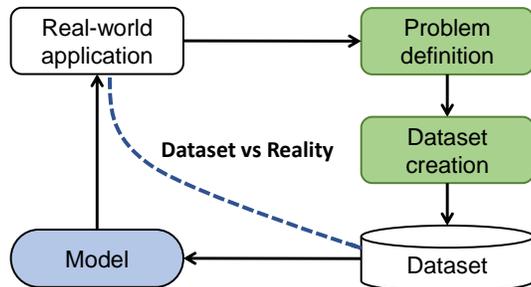
---

*Corresponding author.

Figure 1: Overview of dataset vs reality. Model learns from dataset only and is expected to address the practical task.

determine if it contains the information he/she is looking for.[1] Crafting queries that accurately capture the information need is crucial, and this principle also holds true when it comes to creating datasets for training models, as the ultimate goal is to address practical tasks in the real world.

As illustrated in Figure 1, before we can build a model to address a real-world problem, we need to formally formulate the problem by identifying its input/output, as well as any key constraints. To develop a model, the next essential step is to create a dataset. A dataset is used to simulate the practical task by providing inputs that are the same as, or resemble, those encountered in the real world, along with the expected outputs or ground truth labels. It is expected that the model trained on this dataset will be capable of addressing the practical problem at hand.

It is important to note that the model has no direct access to the practical problem and lacks a literal understanding of its problem definition. The model is trained to solve *the problem that is "defined" by the dataset*, utilizing the input and corresponding expected output provided in the dataset. In this sense, a dataset is analogous to a "query" in information retrieval, and the trained model plays the role of a search engine. Let us suppose that we have succeeded in training a model that is almost perfect using a dataset. Whether this trained model can effectively solve the practical problem *i.e.,* answering the "information need", depends largely on *to what extent the dataset truly reflects the real-world problem.*

In light of the above discussion, this paper aims to conduct a systematic comparison of datasets that share similar formats (*e.g.,* input, output) and are designed for similar or disparate objectives, yet have been frequently employed as benchmarks for identical tasks. From the perspective of information need, we hope to provide a different angle in understanding a research task and also the model performance obtained on these datasets, particularly for inexperienced researchers. For instance, on some datasets, well-trained models significantly outperform human beings; on some other datasets, models are far behind human's performance. Our findings may also be helpful in guiding new dataset creation.

As two case studies, we compare the datasets that are widely used for two research problems: question answering (QA) and image captioning (IC). Specifically, for QA, we study the differences between SQuAD2.0 and Natural Questions (NQ) datasets. For image captioning, we compare four datasets, namely, MS-COCO, Flickr30K, SBU Captions, and Conceptual Captions. Among the six datasets, the questions in SQuAD2.0 and the image captions in MS-COCO and Flickr30K are annotated by crowdworkers; the questions in NQ and the captions in SBU Captions and Conceptual Captions are from real users or content

---

[1]In our discussion, we refer to informational queries, and not navigational or transnational queries. For the latter types of queries, users often have clear expectations on the answers.

providers. In our analysis, we compare the creation processes of the datasets, and try to align the dataset creation with an information need. Again, we argue that if a model is developed on a dataset, then the model is trained to answer the same information need that the dataset was created for. We perform morphosyntactic analysis on datasets with similar input/output but are created for different information needs, to show their differences. Specifically, at the word level, we compare type token ratios for lexical richness and diversity. Based on WordNet, we also compare word specificity reflected by the depth of word path in WordNet. At the sentence level, we compare the similarity of the syntactic structure between sentences (*e.g.,* questions or image captions), computed by tree kernel similarity, and we report grammatical and/or spelling errors in sentences with an off-the-shelf language tool.

Our morphosyntactic analysis between similar datasets shows that: (i) The texts from real users or content providers have a higher type token ratio, indicating richer and more diverse lexicon usage. Accordingly, the texts by crowdworkers are relatively easier to be handled by models, thanks to the less diverse lexicon. Based on the depth of word path in WordNet, words from real users tend to be more specific compared to the crowdsourced words, unless a dataset is purposely collected to cover such objectives during creation. (ii) In terms of sentence structure, the crowdsourced texts share similar syntactic structures, and the tree kernel similarity between these texts is much larger than that from real users or content providers.

The variations observed in the morphosyntactic analysis indicate that deep learning models may have an easier time detecting patterns in crowdsourced texts as opposed to texts written by real users. This could be a reason why models are able to beat human performance on certain datasets, but not the other. However, as we show that the creation of these datasets may or may not share the same information need, the performance obtained on a dataset reflects the capability of deep learning models in addressing the corresponding information need that the dataset embodies. For example, the crowdsourced questions are conditioned on and are restricted to the context information that is provided during the annotation process. Therefore, these questions can only be used for testing the comprehension ability of a model or a human being. Questions of this nature do not align with the information needs typically sought by users of search engines. Conversely, captions sourced from crowdworkers provide a broad perspective of an image, making them very useful in facilitating image search through natural language.

We make two main contributions in this paper. First, we propose the perspective of information need to understand the relationship between a research task and its benchmark datasets. We argue that a model is trained to answer the same information need for which the training dataset is created. Second, in order to support our argument, we conduct experiments on datasets created with different information needs. Specifically, we conduct morphosyntactic analysis on datasets for two research tasks QA and IC as two case studies. Through this paper, we encourage researchers (particularly those who are new to a particular task) to look beyond the performance metric numbers obtained on some widely adopted datasets for a task, and to interpret the numbers through the lens of information need. When creating a dataset, researchers may also incorporate the information need perspective as a factor to determine the degree to which the dataset accurately reflects the research task they intend to tackle.

## 2 Related Work

The rapid development of deep learning technologies has led to significant performance improvement on various tasks from vision to language. On a few benchmarks, deep learning solutions outperform human

beings. Examples include SQuAD2.0 for the reading comprehension[2] task, and MS-COCO dataset for the image captioning[3] task. At the same time, researchers have also indicated that models have yet to outperform humans in the task itself (Sen & Saffari, 2020; Raji, Bender, Paullada, Denton, & Hanna, 2021).

Although recent solutions bring in huge performance gains on various tasks, there are questions about the real progress that has been made in some areas (Cremonesi & Jannach, 2021; Raji et al., 2021). The authors offer a detailed discussion on the limitations of using influential datasets to benchmark progress in machine learning (Raji et al., 2021). Progress made should be justified based on "how closely our evaluations hit the mark in appropriately characterizing the actual anticipated behaviour of the system in the real world or progress on stated motivations and goals for the field" (Raji et al., 2021). The authors also discuss the ineffectiveness of benchmarks that measure the general ability of machine learning models in visual understanding and language understanding. Furthermore, Miceli *et al.* conduct analyses from three areas: data quality, data work, and data documentation to illustrate that training models on incomplete or biased datasets may result in discriminatory outputs (Miceli, Posada, & Yang, 2022). Schlangen (2021) answered a similar question on *why models make better results on benchmark datasets constitute research progress*. In the discussion, Schlangen states that a dataset shall be verified to check "whether the provided input/output pairs can indeed be judged correct relative to the task (in its intensional description)". In our discussion, the perspective of information need can be considered as part of the "intensional description". In this paper, we offer a more concrete framework to interpret the relationships between real-world tasks, datasets, and models (see Figure 1), through the lens of information need. We target on the concrete differences between datasets that are structurally similar but are created for answering different information needs. We argue that a model is trained to answer the information need for which the dataset was originally created for. Hence, progress made measured by benchmark datasets may or may not translate to the progress made in addressing the practical problems depending on the degree of the dataset truly reflecting the practical setting.

In our paper, we use Question Answer (QA) and Image Captioning (IC) tasks and their datasets as two case studies. For QA datasets, (Rogers, Gardner, & Augenstein, 2023) offer a comprehensive survey and discussed the two different types of questions in datasets created for question answering and reading comprehension. The two types of questions are (i) information-seeking questions, and (ii) probing questions. The former is asked by users who do not know the answers, and the latter is for "testing the knowledge of another person or machine". In other words, the asker of a probing question knows the answer. The authors also note that the two classes of questions require different types of reasoning. The classification of question types well aligns with our perspective of information need. Nevertheless, we consider the perspective of information need is more general and can be applied to characterize datasets for many other tasks, by comparing the information need used for dataset creation and the information need for a practical task.

Regarding the image captioning case study, Torralba and Efros (2011) argue that using captured datasets to represent the visual world and focusing solely on beating numbers on benchmarks led researchers to lose sight of their original goal. They analyzed and compared multiple datasets in computer vision, such as ImageNet, MSRC, and PASCAL, and proposed that the datasets have selection bias, capture bias, and negative set bias. These factors cause the model performance to drop significantly when training the model on one dataset and testing on another dataset. In addition, the objects in the ImageNet dataset are primarily

---

[2]SQuAD2.0 leaderboard `https://rajpurkar.github.io/SQuAD-explorer/`
[3]MS-COCO leaderboard: `https://cocodataset.org/#captions-leaderboard`

located in the middle part of an image and are not occluded; hence the pictures do not well represent real-world images (Barbu et al., 2019). In this paper, our interest is not in image coverage or visual content. Rather, we are interested in the differences in the information need for dataset creation, reflected through the captions that come from crowdworkers and that come from content providers. The Flickr30K dataset, which is often used for image captioning, is annotated by crowdworkers, where the annotator's stereotypes and unwarranted inferences about the image can make the data biased when annotating (Van Miltenburg, 2016). In our analysis, we do not consider the bias within annotated captions, but compare them with the captions from image uploaders.

The differences between datasets created through crowdsourcing and datasets collected from real users or content providers, can also be considered as a kind of "bias" to the different information needs. Researchers in several fields have found that biases in datasets often cause errors or misunderstandings of models' capability. In visual question answering (VQA) tasks, VQA models often suffer from language prior due to suspicious correlations between answer occurrences and certain patterns of questions (Cadene, Dancette, Ben younes, Cord, & Parikh, 2019). For instance, color of a banana is always answered as "yellow" regardless of the input image containing a yellow or green banana. In sequential recommendation tasks, a recommender is evaluated using datasets with genuine sequential information. However, Woolridge *et al.* (Woolridge, Wilner, & Glick, 2021) analyze the timestamp information in commonly used datasets and found that they do not represent a meaningful sequential order, especially the MovieLens dataset. As a more generic perspective, we believe information need for dataset creation and that for the practical task is a good way to evaluate the degree of a dataset truly represents the real-world scenario.

## 3   Case Study 1: Question Answering Datasets

We compare two question answering datasets, SQuAD2.0 and Natural Questions (NQ) from Google. Due to their structure similarity, both datasets have been widely used for evaluating question answering solutions (Sen & Saffari, 2020).

### 3.1   Dataset Creation

**SQuAD2.0**   dataset is a large-scale question answering dataset created through crowdsourcing, on top of SQuAD1.1 (Rajpurkar, Zhang, Lopyrev, & Liang, 2016; Rajpurkar, Jia, & Liang, 2018). In the original paper (Rajpurkar et al., 2016), the authors clearly indicate that the formulation of the SQuAD dataset is *reading comprehension* (see Table 1 in (Rajpurkar et al., 2016)). Briefly, a collection of carefully selected 536 high quality Wikipedia articles are sampled, and their paragraphs are extracted. By reading the extracted paragraphs, crowdworkers are requested to create questions based on the content, and to highlight spans in the passage as answers. The creation of question-answer pairs was done through the Daemo platform (Gaikwad et al., 2015). The crowdworkers were encouraged to post questions *in their own words* and the copy-paste feature of the paragraph was disabled. SQuAD2.0 further includes 50,000 unanswerable questions, and the unanswerable question needs to be relevant to the paragraph content. When marking an unanswerable question, workers can view the corresponding paragraph.

The SQuAD2.0 is widely used to train and evaluate question-answering models. For example, the Retro-Reader model (Z. Zhang, Yang, & Zhao, 2021) composed of the sketchy reading module and the intensive reading module has an $F_1$ score of 92.978. Google research uses the SQuAD2.0 dataset to verify that
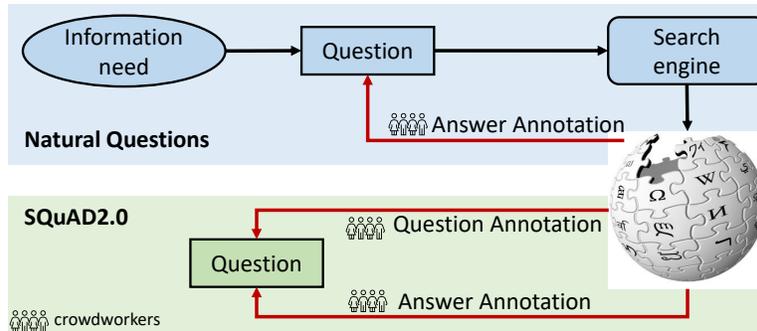
Figure 2: Information needs and data annotation. In NQ, user asks a question to learn more about the topic. In SQuAD dataset, crowdworkers ask questions to test others' (or machine's) understanding of the given context, *e.g.,* a Wikipedia article.

ALBERT achieves higher performance in natural language understanding with fewer parameters than BERT-large (Lan et al., 2019). At the time of writing, a good number of submissions report better performance than humans by either exact match or $F_1$ measure on its leaderboard.

**Natural Questions** dataset is large scale training data for QA problems. Similar to SQuAD2.0, the answers to the questions are also from Wikipedia articles. One key difference is that the questions are real queries (anonymized and aggregated) issued to the Google search engine (Kwiatkowski et al., 2019). The queries are filtered heuristically to select natural questions. Wikipedia pages that appear in the top 5 search results are given to annotators to select a long answer, a short answer, or to mark null if answers cannot be found. The long answer refers to a bounding box containing the answer on the Wikipedia page, typically a paragraph or table. The short answer is a span or a set of spans, typically entities. Because the questions are real queries, the annotator also needs to judge whether a given question is good or bad; a bad question is incomprehensible and does not express the message clearly.

Long and short answer tasks are evaluated separately. Among the state-of-the-art models, Poolingformer achieves excellent results in both tasks, with $F_1$ scores of 0.798 and 0.616 for long and short tasks, respectively (H. Zhang et al., 2021). The Reflection Net (X. Wang, Shou, Gong, Duan, & Jiang, 2020) achieves the top ranking for the short answer task at the time of writing.[4]

## 3.2 Information Need

Despite the structural similarity between the two datasets, *e.g.,* both contain question-answer pairs, and all their answers are from Wikipedia pages, we see very different performances reported on leaderboards. We now compare their sources of questions, from the perspective of information need.

Because the questions in Natural Questions are real queries, we can in general assume that each query represents an information need *for which a searcher is willing to know more*, see Figure 2. Before issuing this query, the searcher does not know the answer. In this sense, the composition of this query is not restricted by or conditioned on the answer. For reading comprehension task, the annotator is given a piece of text, and based on which to ask questions and to annotate answers. Both questions and answers are annotated.

---

[4]QA Leaderboard `https://ai.google.com/research/NaturalQuestions` As our focus in this paper is not the technical details of specific models, we refer readers to their original papers for technical details.

Table 1: Example questions from SQuAD2.0 and Natrual Questions (NQ) with the same (short) answers. Example questions from SQuAD are restricted to the context for evaluating comprehension. For not knowning the answer, NQ questions use more generic terms like "when" instead of "what season".

| Dataset | Example questions with the answer "football" |
|---------|-----------------------------------------------|
| SQuAD | 1. What sports activity is featured in The Times on Mondays? |
| | 2. What did Nigeria win a Summer Olympics gold medal for? |
| | 3. What is London's most popular athletic sport? |
| NQ | 1. where does the term muffed punt come from |
| | 2. what is the most famous sport in russia |
| | 3. what sports did jackie robinson play besides baseball |

| Dataset | Example questions with the answer "winter" |
|---------|---------------------------------------------|
| SQuAD | 1. Which season is the most dry in Oklahoma? |
| | 2. What season is characterized as short in Charleston? |
| | 3. In what season does New Delhi's air pollution worsen? |
| | 4. During what season is DST usually not observed because of the detriments of dark mornings? |
| | 5. The Piedmont is colder than the coast in what season? |
| NQ | 1. when do purple martins migrate to south america |
| | 2. if you live in the southern hemisphere what season do you have in august |
| | 3. when is the newest episode of steven universe coming out |
| | 4. when do deer lose their antlers in california |
| | 5. when is there going to be a new steven universe episode |

The information need in this case is more aligned to "*which question-answer pairs better evaluate another reader's understanding of this piece of text*". The questions are constrained by the given text passage and/or answers in the passage.

Table 1 lists sample questions from both datasets that lead to the same short answer *i.e.,* "football" and "winter", respectively. The first question for football in SQuAD2.0 well reflects its constrained context *i.e.,* "The Times on Monday". The question is not very meaningful if the context is not given. In the contrast, in Natural Questions, football could be an answer to a "where" question. For the "winter" answer, questions in SQuAD2.0 mostly indicate "what season" while the questions in Natural Questions mostly use "when" as the question issuer may not know whether the answer is a season or a month or even a specific date. We argue that questions reflect information needs, and differences in information needs in the two datasets lead to differences in the questions. As the result, the models developed on these two datasets are unlikely to fit into the same real-world application scenarios.

Next, we quantify the differences between questions in the two datasets through word-level and sentence-level analysis. In our analysis, we consider all the questions in both trainning and development, for both datasets.

## 3.3 Word-Level Analysis

Considering questions in both train and development sets, the two datasets SQuAD2.0 and Natural Questions contain 142,192 and 315,203 questions respectively, reported in Table 2.

Table 2: Comparison of two Question-Answering datasets at word and sentence levels. NQ has much higher type token ratio (or higher lexical diversity), and much lower similarity between sentence structures.

| Dimension | SQuAD2.0 | Natural Questions |
|---|---|---|
| **Dataset** | | |
| Number of questions | 142,192 | 315,203 |
| Avg question length | 11.257 | 9.370 |
| **Word** | | |
| #Tokens | 1,600,676 | 2,953,389 |
| #Token types | 45,901 | 59,356 |
| Std. Type Token Ratio | 8.860 | 11.749 |
| Word path depth (stem) | 8.186 (8.371) | 8.371 (8.535) |
| **Sentence** | | |
| Tree Kernel Similarity | 0.088±0.002 | 0.057±0.003 |
| Tree Kernel w/o tokens | 0.267±0.008 | 0.241±0.009 |
| Error types | 533 | 742 |
| Error per question | 0.043 | 0.079 |

The average length of the questions (in the number of tokens by SpaCy) is 11.257 and 9.370, respectively. Figure 3 plots the distribution of question lengths of the two datasets. Nearly half of the questions in Natural Questions contain 8 words, and the distribution of question length in SQuAD2.0 is relatively more smooth. Here, question sampling and selection process in dataset construction could be a possible reason for the differences in the two distributions. Hence this result is reported for reference only. Next, we compare the questions in terms of type token ratio and word specificity.

**Type Token Ratio.** TTR is a commonly used measure for lexical density and richness (Richards, 1987). We use the following formula to compute TTR where the number of token types is the number of unique tokens.

$$TTR = \frac{\#token\ types}{\#tokens} \times 100$$

According to Heaps' law, there is a nonlinear relationship between token type and tokens (Manning et al., 2008). Therefore, it is not appropriate to directly compare the TTR of the two datasets, due to their different sizes. Instead, we calculate the Standard Type Token Ratio - STTR (Richards, 1987). The basic idea is to divide a corpus into multiple chunks of equal size, calculate the TTR values of each chunk respectively, and then average these TTR values to obtain its STTR value. If the number of tokens contained in a chunk is fewer than the predefined number, the chunk is discarded. We use SpaCy to get tokens, and calculate the TTR for every 100,000 tokens on each dataset.

The STTR of the Natural Questions is 11.749, while the value of the SQuAD2.0 dataset is 8.860 (see Table 2). The larger STTR suggests that the lexical density and richness of real questions raised by users is much higher than that created by annotators. In general, high lexical diversity also makes a dataset more challenging.
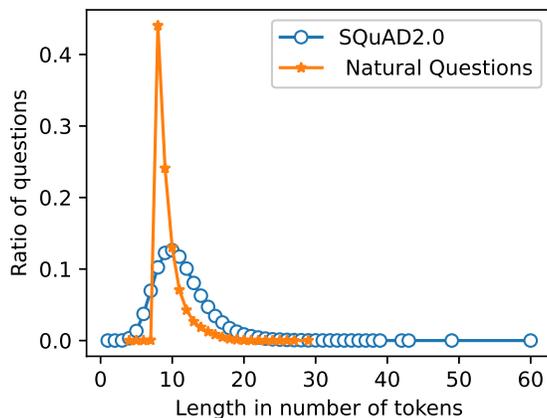
Figure 3: Question length distribution in the two datasets. Question sampling and selection process in dataset construction could be a possible reason for the distribution differences.

**Word Specificity.** So far, our analysis of the questions shows that Natural Questions has a higher type token ratio, suggesting Natural Questions to be more diverse and challenging. Next, we use WordNet (Miller, 1995) to measure word specificity. WordNet is a large lexical database of English where words are organized by their meanings (Miller, 1995). A few relations are defined in WordNet between words or more specifically synsets. In our analysis, we use the super-subordinate relation *i.e.,* hypernymy, to measure the word specificity. Hypernymy is a relation linking more general synsets to specific synsets. Specifically, we use the hypernym_paths() function in the WordNet library to calculate the average path length of nouns that appear in the questions from each dataset.

The mean depth values of nouns with and without stemming are listed in Table 2. Note that the average depth of nouns in the Natural Questions dataset is deeper than that in the SQuAD2.0 dataset. This result suggests that the nouns used in the questions asked by real users are more specific than those annotated by annotators. We believe this result is also consistent with the high diversity reflected by the type token ratio.

## 3.4 Sentence-Level Analysis

We compare the questions between the two datasets from two perspectives: tree kernel similarity, and grammar errors.

**Tree Kernel Similarity.** We use tree kernel similarity to quantify the structural similarity between questions in each dataset (Collins & Duffy, 2002). The central idea of tree kernel is to count the number of common subtrees between two constituency parse trees. Figure 4 gives an illustration with two simple parse trees. Each parse tree has 6 subtrees and the pair has 3 common subtrees. The tree kernel similarity is the ratio between the number of common subtrees and the normalized number of subtrees in the pair.

We use StanfordCoreNLP to obtain constituency parse trees of questions. Then, we randomly select 10,000 parse trees and use the tree kernel to calculate their similarity. We compute two sets of values with and without the leaf-level tokens. The version with leaf-level token means that we consider the complete parse tree of a sentence. The version without leaf-level tokens is to focus solely on the sentence syntactic structure without considering the actual words, *i.e.,* the POS tags are the leaf-level nodes in a parse tree
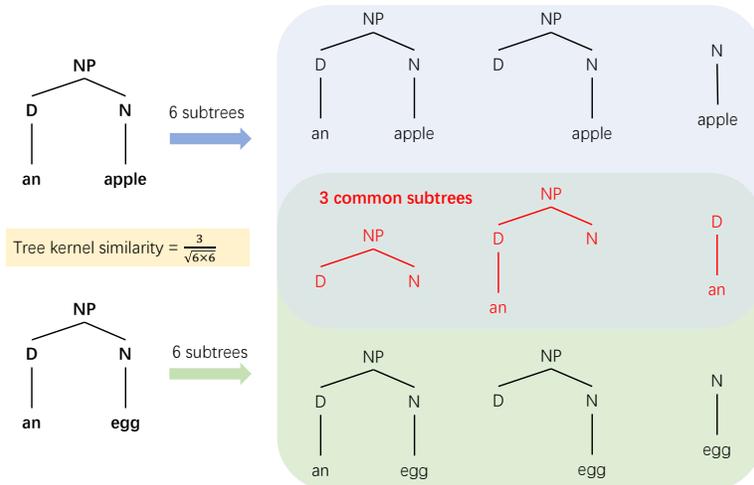
Figure 4: Illustration of tree kernel similarity between two constituency parse trees. Each parse tree has 6 subtrees and the pair has 3 common subtrees. The tree kernel similarity is the ratio between the number of common subtrees and the normalized number of subtrees in the pair, shown in the yellow box.

instead of the words. The average tree kernel similarity of the questions (with leaf-level token) in SQuAD2.0 is 0.088, while that of the Natural Questions dataset is 0.057. This reflects that when crowdworkers annotate questions, they are more inclined to use similar sentence patterns, leading to similar question structures in the dataset. The similarities without considering leaf-level tokens show a similar trend with 0.267 vs 0.241. In short, the questions raised by search engine users are more diverse in terms of syntax structure.

To further compare the syntactic structure similarity between the two datasets, we conduct the unpaired t-test between the pairwise similarity values sampled from the two datasets. Specifically, we randomly select 100 questions and compute their pairwise similarity within each dataset. Then we compute the $p$-value between the two sets of similarity values and obtain $p < 0.001$. We also computed the $p$-value for sample sizes of 1,000 and 10,000 questions, and in both settings $p < 0.001$. This set of results indicates significant differences in the syntactic structure similarity distributions between the two datasets.

**Grammar Accuracy.** Lastly, we use LanguageTool[5] to detect grammatical errors in questions from both datasets. LanguageTool is an open-source rule-based grammar tool for comprehensive grammar and spelling error detection (Naber, 2003). Its API also serves OpenOffice's spell checker.

Because the two datasets are formatted differently, we filter out errors related to white space and upper/lowercase. As reported in Table 2, the number of error types in Natural Questions is significantly higher than that in SQuAD2.0, as expected. Due to the different sizes of the two databases, we calculate the average number of errors per question for a fair comparison. The average number of errors per question in SQuAD2.0 is 0.043, and in Natural Questions is 0.079. In short, questions searched by users on search engines contain more grammar and/or spelling errors than questions annotated by annotators.

---

[5] https://github.com/languagetool-org/languagetool

Table 3: Example image captions that contain keyword "football" from the four datasets. Crowdsourced captions provide an objective description of an image, and content providers provide additional context information that complements the image.

| Dataset | Captions are annotated through crowdsourcing |
|---|---|
| **MS-COCO** | 1. A man holding a football on a grassy field |
| | 2. A man is running while carrying a football |
| **Flickr30K** | 1. Two football players chase after a ball. |
| | 2. A man tries to catch a football on grass surrounded by American flags. |

| Dataset | Captions are from content providers |
|---|---|
| **SBU Captions** | 1. Alexander is blowing the candles on his football themed cake (made by Mom, of course!) |
| | 2. Flag football by the girls and Cheerleading by the boys. Seniors vs Juniors. The boys were deffenetly more interesting than the girls :D |
| **C-Captions** | 1. football player celebrates scoring his side 's second goal of the game with teammates |
| | 2. rugby player of american football team during the match at american football team |

# 4 Case Study 2: Image Captioning Datasets

For image captioning datasets analysis, we consider MS-COCO, Flickr30K, SBU Captions, and Conceptual Captions. These datasets are not created for the same purpose but have all been used for image captioning tasks. Specifically, image captions in MS-COCO and Flickr30K are by crowdworkers, and captions in the remaining two datasets are collected together with the image, or simply image creators. The basic statistics of the dataset regarding the number of images and captions are listed in Table 5.

## 4.1 Dataset Creation

**MS-COCO** (Microsoft Common Objects in COntext) is a large scale image dataset that can be used for image classification, object segmentation, recognition in context, and image captioning (Lin et al., 2014).[6] Here, we are only interested in the image captions (Chen et al., 2015) which are collected through crowdsourcing. The authors of the COCO image collection carefully selected entry-level categories, *i.e.,* category labels that people often use to describe objects, and then searched images from Flickr using a combination of object categories. To label images with captions, crowdworkers were informed to describe all the important parts of the scene, and not to describe unimportant details or things that might have happened in the future or past. Each image contains about five captions. MS-COCO establishes a leaderboard for the captioning task. The Oscar (Object-Semantics Aligned Pre-training) method (Li et al., 2020) is currently ranked first with a BLUE-4 value of 41.7.

**Flickr30K** dataset was created for the purpose of constructing a large scale visual denotation graph with images of everyday activities (Young, Lai, Hodosh, & Hockenmaier, 2014). Each image in this dataset has five captions from five independent crowdworkers "who are not familiar with the specific entities and

---
[6]We used the 2017 version from `https://cocodataset.org/\#download`

11

circumstances depicted" (Young et al., 2014). Flickr30K has been used to evaluate and compare image-captioning models (You, Jin, Wang, Fang, & Luo, 2016a). Captions in other languages have been extended from this dataset (Elliott, Frank, Sima'an, & Specia, 2016; Peng & Li, 2016).

**SBU Captions**   dataset contains 1 million images with relevant captions from Flickr (Ordonez, Kulkarni, & Berg, 2011). Captions are from image uploaders. During data collection, images are filtered to ensure that their text is visually descriptive. Furthermore, the text needs to contain at least one preposition for the spatial relationship of the objects in the image. The SBU Captions dataset was initially used in retrieval tasks (Hodosh, Young, & Hockenmaier, 2013) and later used for image caption generation (Kiros, Salakhutdinov, & Zemel, 2014).

**Conceptual Captions**   dataset contains about 3.3M image and description pairs (Sharma, Ding, Goodman, & Soricut, 2018). Its images and original descriptions were obtained from Web pages. Image and caption pairs were extracted, filtered, and transformed through an automatic pipeline (Chambers et al., 2010). The image descriptions are from the Alt-text elements in HTML, with careful filtering. Google AI provides a leaderboard for the image captioning task on Conceptual Captions.[7]

**Object Distributions.**   Images in all datasets are from the public domain; specifically, images in Conceptual Captions are from web pages and images in the other three datasets are from Flickr. Prior to performing the morphosyntactic analysis, it is imperative to take into account the distribution of objects present in the images across the four datasets. To identify the objects depicted in an image, we utilize the Faster R-CNN image object detection model (Ren, He, Girshick, & Sun, 2017). To be more specific, the Faster R-CNN model used in our experiments is a pre-trained model with the MS COCO dataset.[8] Given the impracticality of manually verifying object detection accuracy across all datasets, we validate the detected objects through the use of ground truth captions present in each of the four datasets. That is, we consider an object is captured in an image if the object is detected by the Faster R-CNN model, and the name of the object is mentioned in the image caption.

Subsequent to extracting the objects from the images, we sample a subset of commonly occurring objects present in all four datasets to investigate potential differences in their relative frequencies across the datasets. Specifically, we use the top 80% most frequent objects from the Flickr30K dataset as the initial pool of sampled objects. If an object in this pool is present in all four datasets, it is included in the final set of sampled objects. Consequently, we obtain a uniform set of objects that are present in all four datasets, albeit with varying frequencies of occurrence in each individual dataset. Based on this common list, we conduct paired t-test and Table 4 reports the $p$-values between datasets. Observe that the $p$-values between MS-COCO, Flickr30K, and Conceptual Captions datasets are much higher than 0.05, indicating these three datasets share similar object distributions for the sampled objects. The SBU Captions dataset on the other hand has a $p < 0.001$ compared to any other dataset. The reason is that the SBU Captions dataset contains more scenery photos which resulted in a very different object distribution. It is important to note that this experiment does not provide comprehensive coverage of the range of objects present in each dataset. Instead, the results obtained from our analysis of the common pool of sampled objects indicate that three of the four datasets exhibit remarkably similar object distributions.

---

[7]Conceptual Captions https://ai.google.com/research/ConceptualCaptions/
[8]https://github.com/open-mmlab/mmdetection/tree/master/configs/faster_rcnn

Table 4: The *p*-values of comparing object distributions in image captioning datasets. SBU Captions has a different object distribution from the other three datasets from having many scenery photos.

| p-value | MS-COCO | Flickr30K | SBU Captions | Conceptual Captions |
|---|---|---|---|---|
| MS-COCO | – | 0.140 | < 0.001 | 0.208 |
| Flickr30K | 0.140 | – | < 0.001 | 0.494 |
| SBU Captions | < 0.001 | < 0.001 | – | < 0.001 |
| Conceptual Captions | 0.208 | 0.494 | < 0.001 | – |

## 4.2 Information Need

All four datasets contain image-caption pairs and all have been used to develop and evaluate image caption generation. The captions are from two different sources: through crowdsourcing in MS-COCO and Flickr30K, by content providers in SBU Captions and Conceptual Captions. Although we do not find a paper comparing the four datasets in one set of experiments, the results reported in different papers with the same measures clearly indicate that SBU Captions and Conceptual Captions are more challenging datasets compared to the other two datasets (Li et al., 2020; You, Jin, Wang, Fang, & Luo, 2016b; E. K. Wang, Zhang, Wang, Wu, & Chen, 2019; Vinyals, Toshev, Bengio, & Erhan, 2015; Changpinyo, Sharma, Ding, & Soricut, 2021).

Again, in this case study, we start with example captions from the datasets, and from which we explain the possible information needs in dataset creation. Table 3 lists example captions from the four datasets that contain a keyword "football". As expected, crowdworkers aim to objectively describe an image, similar to "translating" an image into a piece of text. Based on the description, if we have seen a similar scene before, we are able to visualize the picture in our minds. The *information need* during the annotation process could be "*how someone else would describe this image*". The captions generated in this way do not provide additional information on top of the visual content in a given picture, if we ignore the modality difference. In this sense, this information need is a good fit for a typical "image search by natural language" task.

Captions by content providers, *e.g.,* SBU and Conceptual captions, provide *additional context* to an image. Visual content and its accompanying caption complement each other and together tell the full story about the scene or event. For the same reason, captions cannot be fully discovered solely from the visual content of the image. As a result, image captioning models trained on such a dataset are less likely to give very meaningful results, if evaluated against the captions from content providers. Given that both the images and captions are sourced from content providers, there is no distinct "information need" linking an image to its corresponding caption in this setting, beyond the provision of supplementary information to enhance the overall presentation.

## 4.3 Word-Level Analysis

Similarly to the QA datasets, in this section, we compare the four image captioning datasets at the word level. In terms of caption length distribution, the captions in the four datasets range from 10.3 words on average in Conceptual Captions to 15.4 words in SBU Captions. Again, because the caption lengths are heavily affected by data filtering and instructions to crowdworkers, we show the length distributions in Figure 5 for comprehensiveness and reference purpose only.

Table 5: Analysis on MS-COCO, Flickr30K, SBU Captions and Conceptual Captions (Conceptual) datasets. Datasets by content providers have higher type token ratio and lower similarity between sentence syntactic structures.

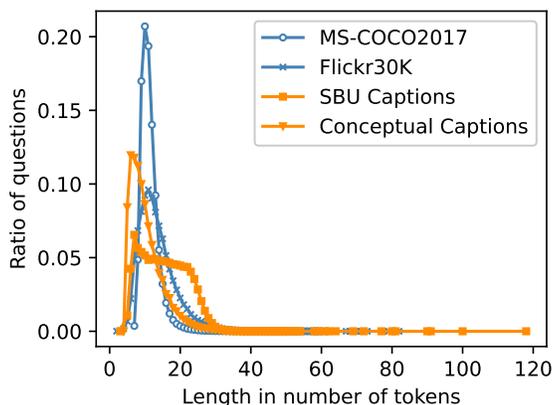| Source of Caption | Crowdsourced | | Content provider | |
|---|---|---|---|---|
| **Dimension** | **MS-COCO** | **Flickr30K** | **SBU Captions** | **Conceptual** |
| **Dataset** | | | | |
| Number of images | 132K | 31K | 1M | 3.3M |
| Number of captions | 616,767 | 158,915 | 1,000,000 | 3,334,173 |
| Average caption length | 11.342 | 13.495 | 15.388 | 10.319 |
| **Word** | | | | |
| #Tokens | 6,995,575 | 2,144,562 | 15,387,646 | 34,405,207 |
| #Token types | 27,837 | 18,377 | 255,126 | 48,939 |
| Std. Type Token Ratio | 3.818 | 5.000 | 11.421 | 8.575 |
| Word path depth (stem) | 8.536 (8.614) | 8.451 (8.526) | 8.667 (8.791) | 8.408 (8.530) |
| **Sentence** | | | | |
| Tree Kernel Similarity | 0.102±0.004 | 0.113±0.003 | 0.060±0.002 | 0.051±0.002 |
| Tree Kernel Similarity w/o tokens | 0.341±0.008 | 0.348±0.008 | 0.257±0.008 | 0.302±0.009 |
| Error types | 589 | 383 | 986 | 1167 |
| Error per caption | 0.077 | 0.051 | 0.069 | 0.052 |



Figure 5: Caption length distribution in four datasets. The distribution is provided for reference purpose because the captions in a dataset is heavily affected by data filtering during dataset creation.

The standard type token ratios of MS-COCO, Flickr30K, SBU Captions, and Conceptual Captions datasets are 3.818, 5, 11.421, and 8.575 respectively, listed in Table 5. Clearly, the captions collected from content providers contain richer and more diverse lexicons. In particular, SBU Captions has comparable STTR with Natural Questions, and both their texts are from real users. Note that, the number of captions in SBU Captions datasets is not the largest but its number of token types is significantly larger than others (see Table 5).

With the highest STTR, it is expected that the words in SBU Captions are more specific, with the deepest word path by WordNet. MS-COCO has a relatively deep word path as well because the images are purposely collected to cover concrete objects commonly seen in our daily life.

## 4.4 Sentence-Level Analysis

MS-COCO and Flickr30K share similar tree kernel similarity of 0.102 and 0.113 respectively. Both similarity values are much higher than that in SBU and Conceptual Captions for 0.06 and 0.051 respectively. We also conduct the unpaired t-test on the pairwise tree kernel similarities computed from the randomly sampled captions in each dataset, and all $p$-values obtained are $< 0.001$ between crowedsourced dataset and the dataset using user provided content. Again, tree kernel similarity measures similarity between sentence syntactic structures in the form of a parse tree. If not considering tokens in tree kernel similarity computation, the similarities from the two crowdsourced datasets remain much larger than the two datasets where data are collected from content providers. In short, captions written by content providers are richer in terms of sentence structure.

We also applied the same LanguageTool on the four image captioning datasets to detect grammar and spelling errors. However, due to differences in dataset filtering, the errors detected do not show a clear pattern in these experiments. The errors per caption are reported in Table 5 last row, where errors related to upper/lower cases and white spaces are not counted. If considering such errors, then the average number of errors per caption in SBU Captions and Conceptual Captions datasets is significantly higher.

## 5 Discussion

Nowadays, leaderboards from large scale datasets have attracted significant attention from many researchers. Although improvements in leaderboard performance are indicative of progress made on a dataset, it is unclear whether these gains arise from the model's ability to address similar real-world tasks or from the model's improved fit to the peculiarities of the dataset itself. In this paper, we compare datasets through the lens of information need to which each dataset was designed to answer. We show that, even if two datasets demonstrate very similar structural similarity, and a model can be easily adapted from one dataset to another, the two datasets may represent two completely different kinds of real-world problems, because they are designed to answer two different information needs. Therefore, the trained models are meant to answer different information needs. As authors indicated in their dataset paper, SQuAD was designed for reading comprehension (Rajpurkar et al., 2016) and Natural Questions dataset is to provide a large scale QA dataset (Kwiatkowski et al., 2019) where the questions are asked by users who want to know more about the topics. The datasets were indeed created to answer two completely different information needs. We argue that the models trained mean to answer their corresponding information needs as well.

Our detailed analysis of the datasets covers both word-level diversity and sentence-level diversity among them. Specifically, datasets that have been annotated by crowdworkers exhibit significantly lower standard type token ratio values compared to datasets that have been collected from real-world settings. Consequently, a trained model has a greater likelihood of fitting a crowdsourced dataset more effectively. At the sentence level, the similarity of sentence structure among annotated datasets is higher than sentences collected from real users. For image captions from content providers, they cover additional context information about the image which is not observable from the image's visual content. Again, the similar sentence structure makes crowdsourced datasets relatively easier to model. In short, the differences between datasets revealed in our analysis partially explain the high and low performances obtained on different datasets *e.g.,* from their leaderboards. Our analysis results are unsurprising given that crowdworkers are required to adhere to specific dataset creation instructions in order to maintain data quality. However, our key focus remains the information need for the dataset creation. Note that, while the dimensions employed in our analysis may have potential implications for predicting task difficulty (Mishra, Bhattacharyya, & Carl, 2013), it is not within the scope of our present investigation to delve into the issue of task difficulty prediction. We leverage the outcomes of our analysis to substantiate our contentions regarding the distinct information needs and the degree to which a dataset can effectively capture the pertinent information needs. In the event that two datasets are not devised with the purpose of addressing identical information needs, any direct comparison of model performance would be rendered incongruous, since the models in question would have been trained to tackle divergent real-world problems.

We remark that our analysis and discussion are not to discourage dataset creation through crowdsourcing. For instance, the objective descriptions in MS-COCO and Flickr30K effectively facilitate natural language search on images. Apart from the person who creates an image, individuals who view the image typically possess a comparable objective perception of its visual elements. In this context, the captions sourced from the crowd serve as an accurate reflection of the real-world problem scenario, as they correspond to a congruent information need. The extent to which the images contained in a dataset accurately depict those encountered in the real world is a distinct topic from the matter under discussion.

Lastly, we acknowledge that understanding the relationship between model, dataset, and a practical task is complex and complicated. In this paper, we put forth a proposition to employ the information need perspective from an Information Retrieval framework as a means to comprehend their interconnections. Nevertheless, we believe information need could be one of many possible perspectives to approach this complex topic. Moreover, question answering and image captioning are just two of the many NLP and NLU tasks. The extent to which the information need perspective can be readily utilized to elucidate other tasks remains inadequately investigated.

# 6   Conclusion

Recent years have witnessed a trend in developing datasets to enable many more interesting studies on various problems. This is evidenced by the resource papers in SIGIR, and the datasets and benchmark track in NeurIPS conferences. At the same time, there is also a trend to reconsider model performance, as models outperform human benchmarks on datasets for various tasks. In this paper, we discuss dataset vs reality from the perspective of information need. We support our discussion with word-level and sentence-level analysis of datasets that are similar in format and have been used for similar tasks. We believe that all

datasets are created as the result of great effort from researchers, and we are not coming to a conclusion on which dataset should or should not be used. We are also not meant to discourage dataset creation through crowdsourcing. Our aspiration is that our analysis will provide insight into how to improve the comprehension of model performance on datasets, particularly with regards to the degree of congruence between the information need of the dataset and that of a real-world application. Stated differently, we urge researchers to contemplate the specific real-world problem that a dataset genuinely embodies, in terms of the corresponding information need, before incorporating it into experiments. Similarly, we advocate for researchers to carefully assess whether the methodology employed in the creation of a dataset is capable of guaranteeing that the information need reflected in the resultant dataset genuinely aligns with that of the practical task.

# References

Barbu, A., Mayo, D., Alverio, J., Luo, W., Wang, C., Gutfreund, D., . . . Katz, B. (2019). Objectnet: A large-scale bias-controlled dataset for pushing the limits of object recognition models. In *Advances in neural information processing systems* (Vol. 32). Curran Associates, Inc.

Cadene, R., Dancette, C., Ben younes, H., Cord, M., & Parikh, D. (2019). Rubi: Reducing unimodal biases for visual question answering. In *Advances in neural information processing systems* (pp. 839–850). Vancouver, Canada: Curran Associates, Inc.

Chambers, C., Raniwala, A., Perry, F., Adams, S., Henry, R. R., Bradshaw, R., & Weizenbaum, N. (2010). Flumejava: easy, efficient data-parallel pipelines. *ACM Sigplan Notices*, *45*(6), 363–375.

Changpinyo, S., Sharma, P., Ding, N., & Soricut, R. (2021, June). Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *Proceedings of the ieee/cvf conference on computer vision and pattern recognition (cvpr)* (p. 3558-3568).

Chen, X., Fang, H., Lin, T.-Y., Vedantam, R., Gupta, S., Dollár, P., & Zitnick, C. L. (2015). Microsoft coco captions: Data collection and evaluation server. *ArXiv*, *abs/1504.00325*.

Collins, M., & Duffy, N. (2002). New ranking algorithms for parsing and tagging: Kernels over discrete structures, and the voted perceptron. In *Acl* (pp. 263–270).

Cremonesi, P., & Jannach, D. (2021, Nov.). Progress in recommender systems research: Crisis? what crisis? *AI Magazine*, *42*(3), 43-54. doi: 10.1609/aimag.v42i3.18145

Elliott, D., Frank, S., Sima'an, K., & Specia, L. (2016). Multi30K: Multilingual English-German image descriptions. In *Proceedings of the 5th workshop on vision and language* (pp. 70–74). Berlin, Germany: Association for Computational Linguistics. doi: 10.18653/v1/W16-3210

Gaikwad, S., Morina, D., Nistala, R., Agarwal, M., Cossette, A., Bhanu, R., . . . others (2015). Daemo: A self-governed crowdsourcing marketplace. In *Acm symposium on user interface software & technology* (pp. 101–102).

Hodosh, M., Young, P., & Hockenmaier, J. (2013). Framing image description as a ranking task: Data, models and evaluation metrics. *Journal of Artificial Intelligence Research*, *47*, 853–899.

Kiros, R., Salakhutdinov, R., & Zemel, R. S. (2014). Unifying visual-semantic embeddings with multimodal neural language models. *arXiv:1411.2539*.

Kwiatkowski, T., Palomaki, J., Redfield, O., Collins, M., Parikh, A., Alberti, C., . . . others (2019). Natural questions: a benchmark for question answering research. *Trans. ACL*, *7*, 453–466.

Lan, Z., Chen, M., Goodman, S., Gimpel, K., Sharma, P., & Soricut, R. (2019). Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942*.

Li, X., Yin, X., Li, C., Zhang, P., Hu, X., Zhang, L., . . . others (2020). Oscar: Object-semantics aligned pre-training for vision-language tasks. In *European conference on computer vision* (pp. 121–137).

Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., . . . Zitnick, C. L. (2014). Microsoft coco: Common objects in context. In *European conference on computer vision* (pp. 740–755).

Manning, C. D., Raghavan, P., & Schütze, H. (2008). *Introduction to information retrieval*. Cambridge, UK: Cambridge University Press. Retrieved from `http://nlp.stanford.edu/IR-book/information-retrieval-book.html`

Miceli, M., Posada, J., & Yang, T. (2022). Studying up machine learning data: Why talk about bias when we mean power? In *Acm group*.

Miller, G. A. (1995). Wordnet: a lexical database for english. *Communications of the ACM*, *38*(11), 39–41.

Mishra, A., Bhattacharyya, P., & Carl, M. (2013). Automatically predicting sentence translation difficulty. In *Proceedings of the 51st annual meeting of the association for computational linguistics, ACL, volume 2: Short papers* (pp. 346–351). The Association for Computer Linguistics.

Naber, D. (2003). *A rule-based style and grammar checker*. GRIN Verlag.

Ordonez, V., Kulkarni, G., & Berg, T. (2011). Im2text: Describing images using 1 million captioned photographs. *Advances in neural information processing systems*, *24*, 1143–1151.

Peng, H., & Li, N. (2016). Generating chinese captions for flickr30k images. (http://vision.soic.indiana.edu/b657/sp2016/projects/penghao/paper.pdf)

Raji, I. D., Bender, E. M., Paullada, A., Denton, E., & Hanna, A. (2021). Ai and the everything in the whole wide world benchmark. In *Neurips 2021 benchmarks and datasets track*.

Rajpurkar, P., Jia, R., & Liang, P. (2018). Know what you don't know: Unanswerable questions for SQuAD. In *Acl (volume 2: Short papers)* (pp. 784–789). Melbourne, Australia: Association for Computational Linguistics.

Rajpurkar, P., Zhang, J., Lopyrev, K., & Liang, P. (2016). SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 conference on empirical methods in natural language processing* (pp. 2383–2392). Austin, Texas: Association for Computational Linguistics.

Ren, S., He, K., Girshick, R., & Sun, J. (2017, Jun). Faster r-cnn: Towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.

Richards, B. (1987). Type/token ratios: What do they really tell us? *Journal of child language*, *14*(2), 201–209.

Rogers, A., Gardner, M., & Augenstein, I. (2023, feb). Qa dataset explosion: A taxonomy of nlp resources for question answering and reading comprehension. *ACM Comput. Surv.*, *55*(10). doi: 10.1145/3560260

Schlangen, D. (2021). Targeting the benchmark: On methodology in current natural language processing research. In C. Zong, F. Xia, W. Li, & R. Navigli (Eds.), *Proceedings of the 59th annual meeting of the association for computational linguistics and the 11th international joint conference on natural language processing, ACL/IJCNLP, volume 2: Short papers* (pp. 670–674). Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-short.85

Sen, P., & Saffari, A. (2020). What do models learn from question answering datasets? In *Emnlp* (pp. 2429–2438). Association for Computational Linguistics.

Sharma, P., Ding, N., Goodman, S., & Soricut, R. (2018). Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Acl (volume 1: Long papers)* (pp. 2556–2565).

Torralba, A., & Efros, A. A. (2011). Unbiased look at dataset bias. In *Cvpr* (pp. 1521–1528).

Van Miltenburg, E. (2016). Stereotyping and bias in the flickr30k dataset. *arXiv preprint arXiv:1605.06083*.

Vinyals, O., Toshev, A., Bengio, S., & Erhan, D. (2015, June). Show and tell: A neural image caption generator. In *Proceedings of the ieee conference on computer vision and pattern recognition (cvpr)*.

Wang, E. K., Zhang, X., Wang, F., Wu, T.-Y., & Chen, C.-M. (2019). Multilayer dense attention model for image caption. *IEEE Access*, *7*, 66358-66368. doi: 10.1109/ACCESS.2019.2917771

Wang, X., Shou, L., Gong, M., Duan, N., & Jiang, D. (2020). No answer is better than wrong answer: A reflection model for document level machine reading comprehension. In *Findings of emnlp*.

Woolridge, D., Wilner, S., & Glick, M. (2021). Sequence or pseudo-sequence? an analysis of sequential

recommendation datasets. In *Perspectives on the evaluation of recommender systems workshop, co-located with acm recsys.*

You, Q., Jin, H., Wang, Z., Fang, C., & Luo, J. (2016a). Image captioning with semantic attention. In *Cvpr* (pp. 4651–4659).

You, Q., Jin, H., Wang, Z., Fang, C., & Luo, J. (2016b, June). Image captioning with semantic attention. In *Proceedings of the ieee conference on computer vision and pattern recognition (cvpr).*

Young, P., Lai, A., Hodosh, M., & Hockenmaier, J. (2014). From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Trans. ACL*, *2*, 67–78. doi: 10.1162/tacl_a_00166

Zhang, H., Gong, Y., Shen, Y., Li, W., Lv, J., Duan, N., & Chen, W. (2021). Poolingformer: Long document modeling with pooling attention. In *Icml.*

Zhang, Z., Yang, J., & Zhao, H. (2021). Retrospective reader for machine reading comprehension. In *Aaai.*