

# Don't Blame the Annotator: Bias Already Starts in the Annotation Instructions

Mihir Parmar<sup>1\*</sup> Swaroop Mishra<sup>1\*</sup> Mor Geva<sup>2</sup> Chitta Baral<sup>1</sup>

<sup>1</sup>Arizona State University <sup>2</sup>Allen Institute for AI

## Abstract

In recent years, progress in NLU has been driven by benchmarks. These benchmarks are typically collected by crowdsourcing, where annotators write examples based on annotation instructions crafted by dataset creators. In this work, we hypothesize that annotators pick up on patterns in the crowdsourcing instructions, which bias them to write similar examples that are then over-represented in the collected data. We study this form of bias, termed *instruction bias*, in 14 recent NLU benchmarks, showing that instruction examples often exhibit concrete patterns, which are propagated by crowdworkers to the collected data. This extends previous work (Geva et al., 2019) and raises a new concern of whether we are modeling the *dataset creator's instructions*, rather than the task. Through a series of experiments, we show that, indeed, instruction bias can lead to overestimation of model performance, and that models struggle to generalize beyond biases originating in the crowdsourcing instructions. We further analyze the influence of instruction bias in terms of pattern frequency and model size, and derive concrete recommendations for creating future NLU benchmarks.<sup>1</sup>

## 1 Introduction

Benchmarks have been proven pivotal for driving progress in Natural Language Understanding (NLU) in recent years (Rogers et al., 2021; Bach et al., 2022; Wang et al., 2022). Nowadays, NLU benchmarks are mostly created through crowdsourcing, where crowdworkers write examples following annotation instructions crafted by dataset creators (Callison-Burch and Dredze, 2010; Zheng et al., 2018; Suhr et al., 2021). The instructions typically include a short description of the task, along with several examples (Dasigi et al., 2019; Zhou et al., 2019; Sakaguchi et al., 2020).

Despite the vast success of this method, past studies have shown that data collected through crowdsourcing often exhibit various biases that lead to overestimation of model performance (Schwartz et al., 2017; Gururangan et al., 2018; Poliak et al., 2018; Tsuchiya, 2018; Le Bras et al., 2020; Hettiachchi et al., 2021). Such biases are often attributed to annotator-related biases, such as writing style and background knowledge (Gururangan et al., 2018; Geva et al., 2019).<sup>2</sup>

In this work, we propose that biases in crowdsourced NLU benchmarks often originate at an early stage in the data collection process of designing the annotation task. In particular, we hypothesize that task instructions provided by dataset creators, which serve as the guiding principles for annotators to complete the task, often influence crowdworkers to follow specific patterns, which are then propagated to the dataset and subsequently over-represented in the collected data. For instance,  $\sim 36\%$  of the instruction examples for the QUOREF dataset (Dasigi et al., 2019) start with “What is the name”, and this same pattern can be observed in  $\sim 59\%$  of the collected instances.

To test our hypothesis, we conduct a broad study of this form of bias, termed *instruction bias*, in 14 recent NLU benchmarks. We find that instruction bias is evident in most of these datasets, showing that  $\sim 72\%$  of instruction examples on average share a few clear patterns. Moreover, we find that these patterns are propagated by annotators to the collected data, covering  $\sim 61\%$  of the instances on average. This suggests that instruction examples play a critical role in the data collection process and the resulting example distribution.

High presence of instruction patterns in a dataset prevents it from representing the underlying task because a task and its associated reasoning have a larger scope than these patterns. For example co-reference resolution, temporal commonsense rea-

\*Equal Contribution

<sup>1</sup><https://github.com/swarooprm/instruction-bias>

<sup>2</sup>See more discussion on related work in App. A.

soning, and numerical reasoning are much broader tasks than the prevalent patterns in QUOREF (“*what is the name...*”), MC-TACO (“*how long...*”) and DROP (“*how many field goals...*”) datasets.

We investigate the effect of instruction bias on model performance, showing that performance is overestimated by instruction bias and that models often fail to generalize beyond instruction patterns. Moreover, we observe that a higher frequency of instruction patterns in the training set often increases the model performance gap on pattern and non-pattern examples and that large models are generally less sensitive to instruction bias.

In conclusion, our work shows that instruction bias widely exists in NLU benchmarks, which often leads to an overestimation of model performance. Based on our study, we derive concrete recommendations for monitoring and alleviating instruction bias in future NLU data collection efforts. From a broader perspective, our findings also have implications on the recent leaning-by-instructions paradigm (Efrat and Levy, 2020; Mishra et al., 2021), where crowdsourcing instructions are used in model training.

## 2 Instruction Bias in NLU Benchmarks

Instructions are the primary resource for educating crowdworkers on how to perform their task (Zheng et al., 2018). Bias in the instructions, dubbed *instruction bias*, could lead crowdworkers to propagate specific patterns to the collected data.

Here, we study instruction bias in NLU benchmarks, focusing on two research questions: (a) Do crowdsourcing instructions exhibit patterns that annotators can pick up on? and (b) Are such patterns propagated by crowdworkers to the collected data? In our study, we use the instructions of 14 recent NLU benchmarks:<sup>3</sup> (1) CLARIQ (Aliannejadi et al., 2020), (2) COSMOSQA (Huang et al., 2019), (3) DROP (Dua et al., 2019), (4) DUORC (Saha et al., 2018), and (5) HOTPOTQA (Yang et al., 2018) (6) HYBRIDQA (Chen et al., 2020), (7) MC-TACO (Zhou et al., 2019), (8) MULTIRC (Khashabi et al., 2018), (9) PIQA (Bisk et al., 2020), (10) QASC (Khot et al., 2020), (11) QUOREF (Dasigi et al., 2019), (12) ROPES (Lin et al., 2019), (13) SCIQA (Welbl et al., 2017), (14) WINOGRANDE (Sakaguchi et al., 2020). These benchmarks were created through different crowdsourcing protocols

<sup>3</sup>The instructions were obtained from Mishra et al. (2021), who have collected those from the dataset authors.

| Dataset        | Pattern  | % Ins. | % $S_{\text{train}}$ | % $S_{\text{test}}$ |
|----------------|--|--------|----------------------|---------------------|
| CLARIQ         | [Are Would<br> Do] you   | 72.2   | 85.1                 | 89                  |
| COSMOSQA       | What AUX   | 87.5   | 45.1                 | 38.4                |
| DROP           | How many<br>[field goals<br>  yards  <br>points  <br>touchdowns] | 70     | 62.5                 | 62.5                |
| DUORC          | [How old  <br>How   What  <br>Who] AUX                           | 70     | 85.1                 | 84                  |
| HOTPOTQA       | [In   Of  <br>From   _ ]<br>[Which What]<br>AUX                  | 87.5   | 53.8                 | 54.2                |
| HYBRIDQA       | Which AUX  | 29.4   | 25.7                 | 15.1                |
| MC-TACO        | How long AUX   | 100    | -                    | 87.6                |
|                | What AUX   | 100    | -                    | 90.1                |
|                | How often AUX  | 100    | -                    | 85.3                |
|                | AUX...<br>[still always<br> by the time]                         | 100    | -                    | 67.3                |
|                | When did /<br>What time  | 100    | -                    | 83.4                |
| MULTIRC        | What AUX   | 14.3   | 38.4                 | 41.5                |
| PIQA           | How [do  <br>can]  | 66.7   | 43.7                 | 42.9                |
| QASC           | What AUX   | 57.1   | 49.3                 | 47                  |
| QUOREF         | What is the<br>[ _   full  <br>real  first<br> last] name        | 36.4   | 57                   | 60                  |
| ROPES          | Which AUX  | 42.9   | 74.1                 | 20.7                |
| SCIQA          | What AUX   | 100    | 83.6                 | 84.5                |
| WINOGRANDE     | [because  <br>so   while  <br>since   but]<br>... the            | 73.7   | 63.4                 | 63.1                |
| <b>Average</b> |  | 71.8   | 59                   | 62                  |

Table 1: Patterns in instruction examples (Ins.) and their propagation to the train ( $S_{\text{train}}$ ) and test ( $S_{\text{test}}$ ) sets of NLU datasets.  $\text{AUX} \in \{\text{am, is, are, was, were, has, have, had, do, does, did, will, would, can, could, may, might, shall, should, must}\}$ , and  $\_$  is an empty string. For MC-TACO, each row corresponds to a different data subset (see App. B).

to evaluate diverse tasks (Mishra et al., 2021) (see task and dataset statistics in App. B, C).

### 2.1 Patterns in Crowdsourcing Instructions

Our goal is to quantify biases in example instructions that propagate to collected data instances. In this study, we focus on an intuitive form of bias

of recurring word patterns, which crowdworkers can easily pick up on. To find such patterns, we manually analyze the instruction examples of each dataset to find a *dominant pattern*, using the following procedure: (a) identifying repeating patterns of  $n \geq 2$  words, (b) merging patterns that are semantically similar or have a significant word overlap, and (c) selecting the most frequent pattern as the dominant pattern (an example is provided in App. D).

Tab. 1 shows the dominant pattern in the instruction examples of each dataset. On average, 71.8% of the instruction examples used to create a dataset exhibit the same dominant pattern, and for 10 out of 14 datasets, the dominant pattern covers more than half of the instruction examples. This suggests that crowdsourcing instructions indeed demonstrate a small set of repeating “shallow” patterns. Moreover, the short length of the patterns (2-4 words) and the typically low number of instruction examples (App. B) make the patterns easily visible to crowdworkers, who can end up following them.

Notably, our results are an underestimation of the actual instruction bias, since (a) we only consider the dominant pattern for each dataset (b) our manual analysis over instruction examples has a preference to short patterns (c) we do not consider paraphrased patterns (beyond the shallow paraphrases which are visible in annotation instructions).

## 2.2 Instruction Bias Propagation to Datasets

We now turn to investigate whether patterns in instruction examples are further propagated by crowdworkers to the collected data. To this end, we analyze the train and test sets of each benchmark<sup>4</sup> to find the same patterns, using simple string matching. To account for syntactic modifications in identified patterns, we also consider synonym words where appropriate and match the paraphrased version of each pattern.

Tab. 1 shows the results. Across all datasets, instruction patterns are ubiquitous in the collected data, occurring in 60.5% of the instances on average, with similar presence in training (59%) and test (62%) examples. While the dominant pattern’s frequency in the data is typically not higher than in the instructions, for CLARIQ, DUORC, MULTIRC, QUOREF and ROPES, the pattern frequency was amplified by the crowdworkers. Interestingly, these datasets used a relatively large number of instruc-

tion examples (App. B), which implies that more examples do not necessarily alleviate the propagation of instruction bias. Example data instances with instruction patterns are provided in App. E.

Propagation of instruction bias to the test set raises concerns regarding its reliability for evaluation of the task and the reasoning abilities it estimates, which we address next.

## 3 Effect on Model Learning

We saw that patterns in crowdsourcing instructions contaminate NLU datasets. In this section, we investigate the effect of this on model performance.

Let  $\mathcal{S}_{\text{train}}$  be the set of training examples, we denote by  $\mathcal{S}_{\text{train}}^p$  and  $\mathcal{S}_{\text{train}}^{-p}$  its disjoint subsets of examples with and without instruction patterns, respectively. We use the similar notation for the set of test examples  $\mathcal{S}_{\text{test}}$ . To analyze the influence of instruction bias on model performance, we fine-tune models on training examples with increasing levels of the instruction pattern. Namely, for  $k \in \{0, 20, 50, 75, 100\}$ , we randomly sample a set  $\mathcal{S}_{\text{train}\%k}^p$  of  $k\%$  of examples from  $\mathcal{S}_{\text{train}}^p$  and train the model on the union  $\mathcal{S}_{\text{train}\%k}^p \cup \mathcal{S}_{\text{train}}^{-p}$ .

Another important question to consider is to what extent models generalize from instruction patterns to the downstream task. To this end, we train models on  $\mathcal{S}_{\text{train}}^p$  and measure their performance on both  $\mathcal{S}_{\text{test}}^{-p}$  and  $\mathcal{S}_{\text{test}}^p$ .

### 3.1 Experimental Setting

**Datasets** We select seven datasets: (1) CLARIQ, (2) DROP, (3) MULTIRC, (4) PIQA, (5) QUOREF, (6) ROPES, and (7) SCIQA. These datasets cover a variety of tasks, different types and levels of instruction bias (Tab. 1), and are different in size (App. C), which allows us to analyze various aspects of instruction bias on model learning.

**Models** We evaluate multiple strong models on these datasets. For all datasets, we use the T5-base and T5-large models<sup>5</sup> (Raffel et al., 2020), except for DROP, where we use Numnet+ (Ran et al., 2019), a RoBERTa model (Liu et al., 2019) with specialized output heads for numerical reasoning. Numnet+ has 355M parameters, which is closer to T5-base (220M) than to T5-large (770M) in size.

**Evaluation** We evaluate model performance using the standard  $F_1$  evaluation score for all datasets.

<sup>4</sup>For some benchmarks, we analyze the validation set in absence of explicit test set.

<sup>5</sup>We use HuggingFace models with default parameters.

|                | Base                          |                                  |         | Large                         |                                  |         |
|----------------|-------------------------------|----------------------------------|---------|-------------------------------|----------------------------------|---------|
|                | $\mathcal{S}_{\text{test}}^p$ | $\mathcal{S}_{\text{test}}^{-p}$ |         | $\mathcal{S}_{\text{test}}^p$ | $\mathcal{S}_{\text{test}}^{-p}$ |         |
| CLARIQ         | 30.7                          | 25.9                             | 15.6% ↓ | 30                            | 27.7                             | 7.7% ↓  |
| DROP           | 77.3                          | 78.7                             | 1.8% ↑  |                               |                                  |         |
| MULTIRC        | 44.6                          | 38.8                             | 13% ↓   | 41.9                          | 43.4                             | 3.6% ↑  |
| PIQA           | 20.8                          | 19.6                             | 5.8% ↓  | 21.9                          | 20.1                             | 8.2% ↓  |
| QUOREF         | 85.8                          | 71.7                             | 16.4% ↓ | 91.9                          | 81.4                             | 11.4% ↓ |
| ROPES          | 61.5                          | 44.5                             | 27.6% ↓ | 55.2                          | 59.2                             | 7.2% ↑  |
| SCIQA          | 80.4                          | 80.3                             | 0.1% ↓  | 82.3                          | 82.5                             | 0.2% ↑  |
| <b>Average</b> | 57.3                          | 51.4                             | 10.3% ↓ | 53.8                          | 52.4                             | 2.6% ↓  |

Table 2: Performance on  $\mathcal{S}_{\text{test}}^p$  vs.  $\mathcal{S}_{\text{test}}^{-p}$  of models trained on  $\mathcal{S}_{\text{train}}$ .

|                | Base                          |                                  |         | Large                         |                                  |         |
|----------------|-------------------------------|----------------------------------|---------|-------------------------------|----------------------------------|---------|
|                | $\mathcal{S}_{\text{test}}^p$ | $\mathcal{S}_{\text{test}}^{-p}$ |         | $\mathcal{S}_{\text{test}}^p$ | $\mathcal{S}_{\text{test}}^{-p}$ |         |
| CLARIQ         | 29.7                          | 25.9                             | 12.8% ↓ | 28.8                          | 23.9                             | 17% ↓   |
| DROP           | 58.2                          | 6.3                              | 89.2% ↓ |                               |                                  |         |
| MULTIRC        | 42.7                          | 31.9                             | 25.3% ↓ | 43.8                          | 37.1                             | 15.3% ↓ |
| PIQA           | 20.8                          | 15                               | 27.9% ↓ | 21.9                          | 15.3                             | 30.1% ↓ |
| QUOREF         | 85.8                          | 64.8                             | 24.5% ↓ | 91.1                          | 76                               | 16.6% ↓ |
| ROPES          | 60.3                          | 45.6                             | 24.4% ↓ | 57.1                          | 57.7                             | 1.1% ↑  |
| SCIQA          | 80.6                          | 80.4                             | 0.3% ↓  | 82.5                          | 82.4                             | 0.1% ↓  |
| <b>Average</b> | 49.7                          | 37.7                             | 24.1% ↓ | 54.2                          | 48.7                             | 10.2% ↓ |

Table 3: Performance on  $\mathcal{S}_{\text{test}}^p$  vs.  $\mathcal{S}_{\text{test}}^{-p}$  of models trained on data instances containing instruction patterns ( $\mathcal{S}_{\text{train}}^p$ ).

### 3.2 Results

**Model performance is overestimated by instruction bias.** We start by comparing the performance on  $\mathcal{S}_{\text{test}}^p$  and  $\mathcal{S}_{\text{test}}^{-p}$  of models trained on the full training set (Tab. 2). The average performance across all datasets is higher on examples that exhibit instruction patterns by  $\sim 10\%$  and  $\sim 3\%$  for the base and large models, respectively. Specifically, the base models’ performance is lower on  $\mathcal{S}_{\text{test}}^{-p}$  for all datasets except DROP, in some cases by a dramatic gap of  $> 15\%$  (e.g. 27.6% in ROPES and 16.4% in QUOREF). In contrast, results for the large models vary across datasets, while the performance gap is generally smaller in magnitude. This shows that model performance is often overestimated by instructions bias, but large models are generally less sensitive to instruction patterns, which might be attributed to their larger capacity to capture knowledge and skills during pre-training.

Next, we analyze how the performance gap between  $\mathcal{S}_{\text{test}}^p$  and  $\mathcal{S}_{\text{test}}^{-p}$  changes for increasing levels of instruction pattern in the training set (Fig. 1). Considering the base models, for most datasets (DROP, QUOREF, PIQA, CLARIQ), a higher presence of instruction patterns in the training set widens the

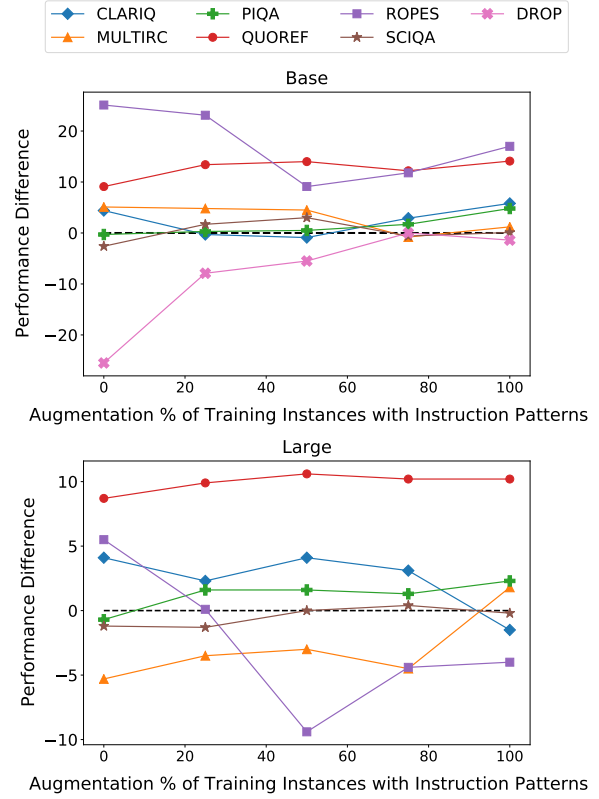


Figure 1: Performance difference on  $\mathcal{S}_{\text{test}}^p$  and  $\mathcal{S}_{\text{test}}^{-p}$ , for increasing levels of examples with instruction patterns in the training set.

performance gap (e.g. by  $> 20$  points in DROP and by  $\sim 5$  in PIQA). MULTIRC and ROPES, do not show clear trends, which might be due to their relatively small size (App. C). Moving to the large models, changes in performance difference are smaller in magnitude (as in Tab. 2). Nonetheless, there is a marginal increase in the performance gap for QUOREF, PIQA, and SCIQA when having more instruction patterns in the training data.

Interestingly, performance gaps on QUOREF are consistently large across all experiments (Tab. 2, Fig. 1). A possible explanation is the relatively long and frequent instruction pattern (Tab. 1), which draws a clear boundary between instruction patterns and non-patterns.

**Models often fail to generalize beyond instruction patterns.** Tab. 3 shows the performance on  $\mathcal{S}_{\text{test}}^p$  and  $\mathcal{S}_{\text{test}}^{-p}$  when training only on examples with instruction patterns. Across all experiments, we observe large performance gaps, reaching to  $\sim 89\%$  in DROP and  $> 15\%$  in both base and large models for PIQA, MULTIRC, and QUOREF. As in previous results, the performance gap is lower for the large models compared to the base ones, reiterating



that they are less sensitive to instruction patterns. Overall, our results indicate that models trained only on examples with instruction patterns fail to generalize to other task examples. This further stresses that instruction bias should be monitored and avoided during data collection.

## 4 Conclusions and Discussion

We identify a prominent source of bias in crowd-sourced NLU datasets, called *instruction bias*, which originates in annotation instructions written by dataset creators. We study this bias in 14 NLU benchmarks, showing that instruction examples used to create NLU benchmarks often exhibit clear patterns that are propagated by annotators to the collected data. In addition, we investigate the effect of instruction bias on model performance, showing that instruction patterns can lead to overestimation of model performance as well as limit the ability of models to generalize to other task examples. These findings also have implications on the recently popular leaning-by-instructions paradigm (Efrat and Levy, 2020; Mishra et al., 2021), where crowdsourcing instructions are utilized as a signal for model training.

We conclude with the following recommendations: (1) Crowdsourcing instructions should be diverse; this could be achieved, for example, by having a large number of instructive examples, periodically sampling examples from previously collected data, or rephrasing examples using neural models. (2) Word patterns in collected instances should be analyzed during data collection, as well as possible correspondence to instruction examples. (3) Correlation between model performance and input patterns should be checked, when evaluating models. We hope our work will bring more attention to developing better representation of reasoning tasks beyond the benchmarks containing instruction bias.

## Acknowledgments

We thank Daniel Khashabi from Allen Institute for AI for his valuable feedback on an early stage of this work. This research was supported in part by the Computer Science Scholarship granted by the Séphora Berrebi Foundation.

## References

- Mohammad Aliannejadi, Julia Kiseleva, Aleksandr Chuklin, Jeff Dalton, and Mikhail Burtsev. 2020. Convai3: Generating clarifying questions for open-domain dialogue systems (clariq). *arXiv preprint arXiv:2009.11352*.
- Anjana Arunkumar, Swaroop Mishra, Bhavdeep Sachdeva, Chitta Baral, and Chris Bryan. 2020. Real-time visual feedback for educative benchmark creation: A human-and-metric-in-the-loop workflow. *NeurIPS 2020 Workshop HAMLETS*.
- Stephen H Bach, Victor Sanh, Zheng-Xin Yong, Albert Webson, Colin Raffel, Nihal V Nayak, Abheesht Sharma, Taewoon Kim, M Saiful Bari, Thibault Fèvre, et al. 2022. Promptsource: An integrated development environment and repository for natural language prompts. *arXiv preprint arXiv:2202.01279*.
- Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. 2020. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 7432–7439.
- Chris Callison-Burch and Mark Dredze. 2010. [Creating speech and language data with Amazon’s Mechanical Turk](#). In *Proceedings of the NAACL HLT 2010 Workshop on Creating Speech and Language Data with Amazon’s Mechanical Turk*, pages 1–12, Los Angeles. Association for Computational Linguistics.
- Joseph Chee Chang, Saleema Amershi, and Ece Kamar. 2017. Revolt: Collaborative crowdsourcing for labeling machine learning datasets. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*, pages 2334–2346.
- Wenhu Chen, Hanwen Zha, Zhiyu Chen, Wenhan Xiong, Hong Wang, and William Yang Wang. 2020. [HybridQA: A dataset of multi-hop question answering over tabular and textual data](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1026–1036, Online. Association for Computational Linguistics.
- Jonathan H Clark, Eunsol Choi, Michael Collins, Dan Garrette, Tom Kwiatkowski, Vitaly Nikolaev, and Jennimaria Palomaki. 2020. Tydi qa: A benchmark for information-seeking question answering in typologically diverse languages. *Transactions of the Association for Computational Linguistics*, 8:454–470.
- Pradeep Dasigi, Nelson F. Liu, Ana Marasović, Noah A. Smith, and Matt Gardner. 2019. [Quoref: A reading comprehension dataset with questions requiring coreferential reasoning](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5925–5932, Hong Kong, China. Association for Computational Linguistics.

- Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019. [DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2368–2378, Minneapolis, Minnesota. Association for Computational Linguistics.
- Avia Efrat and Omer Levy. 2020. The turking test: Can language models understand instructions? *arXiv preprint arXiv:2010.11982*.
- Carsten Eickhoff. 2018. Cognitive biases in crowdsourcing. In *Proceedings of the eleventh ACM international conference on web search and data mining*, pages 162–170.
- Mor Geva, Yoav Goldberg, and Jonathan Berant. 2019. [Are we modeling the task or the annotator? an investigation of annotator bias in natural language understanding datasets](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1161–1166, Hong Kong, China. Association for Computational Linguistics.
- Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies. *Transactions of the Association for Computational Linguistics*, 9:346–361.
- Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel Bowman, and Noah A. Smith. 2018. [Annotation artifacts in natural language inference data](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 107–112, New Orleans, Louisiana. Association for Computational Linguistics.
- Danula Hettiachchi, Mark Sanderson, Jorge Goncalves, Simo Hosio, Gabriella Kazai, Matthew Lease, Mike Schaekermann, and Emine Yilmaz. 2021. Investigating and mitigating biases in crowdsourced data. In *Companion Publication of the 2021 Conference on Computer Supported Cooperative Work and Social Computing*, pages 331–334.
- Xiao Hu, Haobo Wang, Anirudh Vegesana, Somesh Dube, Kaiwen Yu, Gore Kao, Shuo-Han Chen, Yung-Hsiang Lu, George K Thiruvathukal, and Ming Yin. 2020. Crowdsourcing detection of sampling biases in image datasets. In *Proceedings of The Web Conference 2020*, pages 2955–2961.
- Lifu Huang, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2019. [Cosmos QA: Machine reading comprehension with contextual commonsense reasoning](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2391–2401, Hong Kong, China. Association for Computational Linguistics.
- Rabeeh Karimi Mahabadi, Yonatan Belinkov, and James Henderson. 2020. [End-to-end bias mitigation by modelling biases in corpora](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8706–8716, Online. Association for Computational Linguistics.
- Daniel Khashabi, Snigdha Chaturvedi, Michael Roth, Shyam Upadhyay, and Dan Roth. 2018. Looking beyond the surface: A challenge set for reading comprehension over multiple sentences. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 252–262.
- Tushar Khot, Peter Clark, Michal Guerquin, Peter Jansen, and Ashish Sabharwal. 2020. QASC: A dataset for question answering via sentence composition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 8082–8090.
- Walter S Lasecki, Mitchell Gordon, Danai Koutra, Malte F Jung, Steven P Dow, and Jeffrey P Bigham. 2014. Glance: Rapidly coding behavioral video with the crowd. In *Proceedings of the 27th annual ACM symposium on User interface software and technology*, pages 551–562.
- Ronan Le Bras, Swabha Swayamdipta, Chandra Bhagavatula, Rowan Zellers, Matthew E Peters, Ashish Sabharwal, and Yejin Choi. 2020. Adversarial filters of dataset biases. In *ICML*.
- Kevin Lin, Oyvind Tafjord, Peter Clark, and Matt Gardner. 2019. [Reasoning over paragraph effects in situations](#). In *Proceedings of the 2nd Workshop on Machine Reading for Question Answering*, pages 58–62, Hong Kong, China. Association for Computational Linguistics.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Rabeeh Karimi Mahabadi, Yonatan Belinkov, and James Henderson. 2021. Variational information bottleneck for effective low-resource fine-tuning. *ICLR*.
- Swaroop Mishra, Anjana Arunkumar, Bhavdeep Sachdeva, Chris Bryan, and Chitta Baral. 2020. Dqi: Measuring data quality in nlp. *arXiv preprint arXiv:2005.00816*.
- Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2021. Cross-task generalization via natural language crowdsourcing instructions. *ACL*.

- Maryam M Najafabadi, Flavio Villanustre, Taghi M Khoshgoftaar, Naeem Seliya, Randall Wald, and Edin Muharemagic. 2015. Deep learning applications and challenges in big data analytics. *Journal of big data*, 2(1):1–21.
- Nikita Nangia, Saku Sugawara, Harsh Trivedi, Alex Warstadt, Clara Vania, and Samuel R. Bowman. 2021. [What ingredients make for an effective crowdsourcing protocol for difficult NLU data collection tasks?](#) In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1221–1235, Online. Association for Computational Linguistics.
- Adam Poliak, Jason Naradowsky, Aparajita Haldar, Rachel Rudinger, and Benjamin Van Durme. 2018. Hypothesis only baselines in natural language inference. In *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics*, pages 180–191.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. [Exploring the limits of transfer learning with a unified text-to-text transformer](#). *Journal of Machine Learning Research*, 21(140):1–67.
- Hossein A Rahmani and Jie Yang. 2021. Demographic biases of crowd workers in key opinion leaders finding. *arXiv preprint arXiv:2110.09248*.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. [Know what you don’t know: Unanswerable questions for SQuAD](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 784–789, Melbourne, Australia. Association for Computational Linguistics.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. [SQuAD: 100,000+ questions for machine comprehension of text](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas. Association for Computational Linguistics.
- Qiu Ran, Yankai Lin, Peng Li, Jie Zhou, and Zhiyuan Liu. 2019. NumNet: Machine reading comprehension with numerical reasoning. In *Proceedings of EMNLP*.
- Anna Rogers, Matt Gardner, and Isabelle Augenstein. 2021. Qa dataset explosion: A taxonomy of nlp resources for question answering and reading comprehension. *arXiv preprint arXiv:2107.12708*.
- Amrita Saha, Rahul Aralikkatte, Mitesh M. Khapra, and Karthik Sankaranarayanan. 2018. [DuoRC: Towards complex language understanding with paraphrased reading comprehension](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1683–1693, Melbourne, Australia. Association for Computational Linguistics.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2020. Winogrande: An adversarial winograd schema challenge at scale. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 8732–8740.
- Roy Schwartz, Maarten Sap, Ioannis Konstas, Leila Zilles, Yejin Choi, and Noah A Smith. 2017. The effect of different writing tasks on linguistic style: A case study of the roc story cloze task. In *Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017)*, pages 15–25.
- Alane Suhr, Clara Vania, Nikita Nangia, Maarten Sap, Mark Yatskar, Samuel R. Bowman, and Yoav Artzi. 2021. [Crowdsourcing beyond annotation: Case studies in benchmark data collection](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing: Tutorial Abstracts*, pages 1–6, Punta Cana, Dominican Republic & Online. Association for Computational Linguistics.
- Swabha Swayamdipta, Roy Schwartz, Nicholas Lourie, Yizhong Wang, Hannaneh Hajishirzi, Noah A Smith, and Yejin Choi. 2020. Dataset cartography: Mapping and diagnosing datasets with training dynamics. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9275–9293.
- Masatoshi Tsuchiya. 2018. [Performance impact caused by hidden bias of training data for recognizing textual entailment](#). In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).
- Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Anjana Arunkumar, Arjun Ashok, Arut Selvan Dhanasekaran, Atharva Naik, David Stap, et al. 2022. Benchmarking generalization via in-context instructions on 1,600+ language tasks. *arXiv preprint arXiv:2204.07705*.
- Johannes Welbl, Nelson F. Liu, and Matt Gardner. 2017. [Crowdsourcing multiple choice science questions](#). In *Proceedings of the 3rd Workshop on Noisy User-generated Text*, pages 94–106, Copenhagen, Denmark. Association for Computational Linguistics.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D Manning. 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2369–2380.
- Feifei Zheng, Ruoling Tao, Holger R Maier, Linda See, Dragan Savic, Tuqiao Zhang, Qiuwen Chen,

Thaine H Assumpção, Pan Yang, Bardia Heidari, et al. 2018. Crowdsourcing methods for data collection in geophysics: State of the art, issues, and future directions. *Reviews of Geophysics*, 56(4):698–740.

Ben Zhou, Daniel Khashabi, Qiang Ning, and Dan Roth. 2019. “going on a vacation” takes longer than “going for a walk”: A study of temporal commonsense understanding. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3363–3369, Hong Kong, China. Association for Computational Linguistics.



## A Biases in NLU Benchmarks

Crowdsourcing has been a widely adapted approach to create large scale datasets such as SQUAD 1.1 (Rajpurkar et al., 2016, 2018), DROP (Dua et al., 2019), QUOREF (Dasigi et al., 2019) and many more (Najafabadi et al., 2015; Callison-Burch and Dredze, 2010; Lasecki et al., 2014; Zheng et al., 2018; Chang et al., 2017). Many past works investigate different types of bias in crowdsourcing datasets such as cognitive bias (Eickhoff, 2018), annotator bias (Gururangan et al., 2018; Geva et al., 2019), sampling bias (Hu et al., 2020), demographic bias (Rahmani and Yang, 2021) and others (Hettiachchi et al., 2021). Many works on bias in NLU benchmarks focus on biases resulting from the crowdsourcing annotations, and how annotator-specific patterns create biases in data (Geva et al., 2019).

To mitigate the bias, prior works have focused on priming crowdsourcing annotators with minimal information to increase their imagination (Geva et al., 2021; Clark et al., 2020) to avoid recurring patterns. Arunkumar et al. (2020) develops a real time feedback and metric-in-the loop (Mishra et al., 2020) workflow to educate crowdworkers in controlling dataset biases. Nangia et al. (2021) provides an iterative protocol with expert assessments for crowdsourcing data collection to increase difficulty of instances. (Swayamdipta et al., 2020) introduces dataset map as a model-based tool to characterize and diagnose datasets. Also, Karimi Mahabadi et al. (2020); Mahabadi et al. (2021) propose learning strategies to train neural models, which are more robust to such biases and transfer better to out-of-domain datasets.

In this work, we show that biases exhibited by annotators start from the crowdsourcing instructions designed by dataset creators.

## B Tasks and Instruction Statistics

Tab. 4 describes the tasks of all NLU datasets used in our work, and provides the number of examples present in crowdsourcing instructions of each dataset. From Tab. 4, we can observe that our analysis involves a wide range of different tasks. Also, we believe that the lower number of examples in crowdsourcing instructions might be limiting the imagination of annotators while creating samples, resulting in instruction bias.

| Dataset    | Task                    | # of Examples |
|------------|-------------------------|---------------|
| CLARIQ     | Ambiguous QA            | 18            |
| COSMOSQA   | Commonsense Reasoning   | 8             |
| DROP       | Numerical Reasoning     | 10            |
| DUORC      | Paraphrased RC          | 10            |
| HOTPOTQA   | Multi-hop QA            | 8             |
| HYBRIDQA   | QA                      | 17            |
| MC-TACO    | Event Duration          | 3             |
|            | Event Ordering          | 2             |
|            | Frequency               | 2             |
|            | Stationary              | 2             |
|            | Absolute Point          | 2             |
| MULTIRC    | Complex QA              | 7             |
| PIQA       | Physical Interaction QA | 6             |
| QASC       | Complex QA              | 7             |
| QUOREF     | Coreference QA          | 11            |
| ROPES      | RC                      | 14            |
| SCIQA      | Science-based QA        | 6             |
| WINOGRANDE | Commonsense Reasoning   | 19            |
| Average    |                         | 8.4           |

Table 4: Tasks of each dataset and number of examples in crowdsourcing instruction of each dataset. RC: Reading Comprehension, QA: Question Answering.

## C Dataset Statistics

Tab. 5 describes the statistics of train and evaluation sets of datasets used in our experiments. Here, we can observe that each selected dataset differs in terms of number of training samples, % of instruction patterns, and tasks.

## D Pattern Extraction Method

Here, we describe an example to show how we extract the dominant pattern from the crowdsourcing instructions and subsequently identify the same pattern in the dataset. We try to find recurring word patterns such as “Are you...”, “how many points...”, “Was... still...”, “since... the...”.

For example, MC-TACO (event duration) has 3 examples in crowdsourcing instructions: (1) how long did Jack play basketball?, (2) how long did he do his homework?, and (3) how long did it take for him to get the Visa? In step (a), we analyze examples manually and find *dominant pattern*. Here, we can see that all examples contain tri-gram pattern, i.e., “how long did”. In step (b), we try to generate more possible patterns that are semantically similar to the *dominant pattern* or have a significant

| Dataset      | Train                        |                                |                                   | Test                        |                               |                                  |
|--------------|------------------------------|--------------------------------|-----------------------------------|-----------------------------|-------------------------------|----------------------------------|
|              | $\mathcal{S}_{\text{train}}$ | $\mathcal{S}_{\text{train}}^p$ | $\mathcal{S}_{\text{train}}^{-p}$ | $\mathcal{S}_{\text{test}}$ | $\mathcal{S}_{\text{test}}^p$ | $\mathcal{S}_{\text{test}}^{-p}$ |
| CLARIQ       | 8566                         | 7286 85.1%                     | 1280 14.9%                        | 4499                        | 4006 89%                      | 493 11%                          |
| DROP         | 77409                        | 48422 62.5%                    | 28987 37.5%                       | 9536                        | 5960 62.5%                    | 3576 37.3%                       |
| MULTIRC      | 5131                         | 1972 38.4%                     | 3159 61.6%                        | 953                         | 395 41.5%                     | 558 58.6%                        |
| PIQA         | 17171                        | 7508 43.7%                     | 9663 56.3%                        | 3268                        | 1401 42.9%                    | 1867 57.1%                       |
| QUOREF       | 19399                        | 11052 57%                      | 8347 43%                          | 2418                        | 1451 60%                      | 967 40%                          |
| ROPES        | 1412                         | 1046 74.1%                     | 366 25.9%                         | 203                         | 42 20.7%                      | 161 79.3%                        |
| SCIQA        | 11679                        | 9765 83.61%                    | 1914 16.4%                        | 1000                        | 845 84.5%                     | 155 15.5%                        |
| <b>Total</b> | 140767                       | 87051 61.8%                    | 53716 38.2%                       | 21877                       | 14100 64.5%                   | 7777 35.6%                       |

Table 5: Statistics of number of train and test examples with and without instruction patterns.  $\mathcal{S}_{\text{train}}$ : set of examples in train set,  $\mathcal{S}_{\text{train}}^p$ : set of examples in train set with instruction pattern,  $\mathcal{S}_{\text{train}}^{-p}$ : set of examples in train set without instruction pattern,  $\mathcal{S}_{\text{test}}$ : set of examples in test set,  $\mathcal{S}_{\text{test}}^p$ : set of examples in test set with instruction pattern,  $\mathcal{S}_{\text{test}}^{-p}$ : set of examples in test set without instruction pattern.

word overlap. Here, “*how long did*” can be “*how long was*”, “*how long does*”, etc. (i.e, How long AUX). In step (c), we look for all these possible patterns in datasets using simple word-matching techniques.

## E Examples

Tab. 6 provides dataset, instruction patterns and corresponding examples of data instances that exhibit the instruction patterns.

| Dataset     | Pattern  | Examples   |
|-------------|--|--|
| CLARIQ      | [Are Would Do] you                                   | Are you looking for a specific web site?   |
|             |  | What kind of train are you looking for?  |
|             |  | Do you want to watch news videos or read the news?   |
|             |  | Would you like the location of the ritz carlton lake las vegas?  |
| COSMOSQA    | What AUX   | What may happen after the young man makes his call?  |
|             |  | What might happen if you have him for the whole day?   |
|             |  | What's a possible reason the writer doesn't look disabled on the outside?  |
| DROP        | How many [field goals   yards   points   touchdowns] | How many touchdowns did Jones have?  |
|             |  | How many field goals did Kris Brown kick   |
|             |  | How many yards was the longest touchdown of the game?  |
|             |  | After Akers 32-yard field goal, how many points behind was Washington?   |
| HOTPOTQA    | [in of from _] [Which What] AUX                      | Which franchise was founded in 1978, Chuck E. Cheese's or Jet's Pizza?   |
|             |  | Busan, in the area surrounding the mountain of Geumjeongsan, is the second most populated city in which country? |
|             |  | What is the name of the third album from singer Selena Quintanilla-Pérez?  |
| MC-TACO     | How long AUX   | How long was his mother ill?   |
|             | What AUX   | What did the government decide after the 9/11 attack?  |
|             | How often AUX  | How often would one family be able to do something like this?  |
|             | AUX... [still always]                                | Will electronic espionage always be happening in the U.S.?   |
|             | When did / What time                                 | Is she still gone?<br>What time did the planes crash into the World Trade Center?<br>When did Durer die?         |
| MULTIRC     | Which AUX  | What was Poe's first published work?   |
|             |  | What is the full name of the person described?   |
|             |  | What kind of career does Christie Brinkley have?   |
| PIQA        | How [do   can]                                       | How do I make orange icing if I have store-bought white frosting?  |
|             |  | How can I make popsicles for dogs?   |
|             |  | Are you nervous about giving a speech or doing something? How can you calm yourself?                             |
| QUOREF      | What is the [full   real   first   last] name        | What is the first name of the person who purchases a revolver?   |
|             |  | What is the full name of the person who is calmly asked to leave?  |
|             |  | What was the name of the house where Appleton Water Tower was built?   |
|             |  | What is the last name of the person who convinces the girls to help him look for the treasure?                   |
| ROPES       | Which AUX  | Which area would be less likely to experience a drought and have better chance at a new growth?                  |
|             |  | Which hair spray brand should Greg buy to be environmentally friendly?   |
|             |  | Which markalong was produced asexually?  |
| SCIQA       | What AUX   | What are by far the most common type of invertebrate?  |
|             |  | What do waves deposit to form sandbars and barrier islands?  |
|             |  | What is the term for the total kinetic energy of moving particles of matter?                                     |
| WINO-GRANDE | [because   so   while   since   but] ... the         | The dog didn't like its collar but was okay with its leash because the _ was loose on it.                        |
|             |  | Hunter took Benjamin's clothes to the laundromat, since _ had the day off that day.                              |
|             |  | James sang his song at the top of his voice so as to be heard over the noise but the _ is louder.                |

Table 6: Examples of data instances from original dataset that contain instruction patterns. AUX  $\in$  {am, is, are, was, were, has, have, had, do, does, did, will, would, can, could, may, might, shall, should, must}. \_ : <blank>.