

---

# Teach Me to Explain: A Review of Datasets for Explainable Natural Language Processing

---

Sarah Wiegreffe\*

School of Interactive Computing  
Georgia Institute of Technology  
saw@gatech.edu

Ana Marasović\*

Allen Institute for AI  
University of Washington  
anam@allenai.org

## Abstract

1 Explainable Natural Language Processing (ExNLP) has increasingly focused on  
2 collecting human-annotated textual explanations. These explanations are used  
3 downstream in three ways: as data augmentation to improve performance on a  
4 predictive task, as supervision to train models to produce explanations for their  
5 predictions, and as a ground-truth to evaluate model-generated explanations. In  
6 this review, we identify 61 datasets with three predominant classes of textual expla-  
7 nations (highlights, free-text, and structured), organize the literature on annotating  
8 each type, identify strengths and shortcomings of existing collection methodologies,  
9 and give recommendations for collecting ExNLP datasets in the future.

## 10 1 Introduction

11 Interpreting supervised machine learning (ML) models is crucial for ensuring their reliability and  
12 trustworthiness in high-stakes scenarios. Models that produce justifications for their individual  
13 predictions (sometimes referred to as *local explanations*) can be inspected for the purposes of  
14 debugging, quantifying bias and fairness, understanding model behavior, and ascertaining robustness  
15 and privacy [80]. These benefits have led to the development of datasets that contain human  
16 justifications for the true label (overviewed in Tables 3–5). In particular, human justifications are used  
17 for three goals: (i) to aid models with additional training supervision [139], (ii) to train interpretable  
18 models that explain their own predictions [19], and (iii) to evaluate plausibility of model-generated  
19 explanations by measuring their agreement with human explanations [28].

20 Dataset collection is the most under-scrutinized component of the ML pipeline [90]—it is estimated  
21 that 92% of ML practitioners encounter data cascades, or downstream problems resulting from poor  
22 data quality [106]. It is important to constantly evaluate data collection practices critically and  
23 standardize them [12, 38, 92]. We expect that such examinations are particularly valuable when many  
24 related datasets are released contemporaneously and independently in a short period of time, as is the  
25 case with ExNLP datasets.

26 This survey aims to review and summarize the literature on collecting textual explanations, high-  
27 light what has been learned to date, and give recommendations for future dataset construction. It  
28 complements other explainable AI (XAI) surveys and critical retrospectives that focus on definitions,  
29 methods, and/or evaluation [32, 14, 75, 1, 100, 49, 41, 130, 25, 43, 79, 118, 11, 83, 52, 18], but  
30 not on datasets. We call such datasets ExNLP datasets, because modeling them for the three goals  
31 mentioned above requires NLP techniques. Datasets and methods for explaining fact checking [63]  
32 and reading comprehension [114] have been reviewed; we are the first to review all datasets with  
33 textual explanations regardless of task, comprehensively categorize them into three distinct classes,  
34 and provide critical retrospectives and best-practice recommendations.

---

\* Equal contributions.

Instance	Explanation
<i>Premise:</i> A white race dog wearing the number eight runs on the track. <i>Hypothesis:</i> A white race dog runs around his yard. <i>Label:</i> contradiction	<b>(highlight)</b> <i>Premise:</i> A white race dog wearing the number eight runs on the <b>track</b> . <i>Hypothesis:</i> A white race dog runs around his <b>yard</b> . <b>(free-text)</b> A race track is not usually in someone’s yard.
<i>Question:</i> Who sang the theme song from Russia With Love? <i>Paragraph:</i> ...The theme song was composed by Lionel Bart of Oliver! fame and sung by Matt Monro... <i>Answer:</i> Matt Monro	<b>(structured)</b> <i>Sentence selection:</i> (not shown) <i>Referential equality:</i> “the theme song from russia with love” (from question) = “The theme song” (from paragraph) <i>Entailment:</i> X was composed by Lionel Bart of Oliver! fame and sung by ANSWER. ⊢ ANSWER sung X

Table 1: Examples of explanation types discussed in §2. The first two rows show a highlight and free-text explanation for an E-SNLI instance [19]. The last row shows a (partial) structured explanation from QED for a NATURALQUESTIONS instance [68].

Instance with Highlight	Highlight Type Clarification
<i>Review:</i> this film is <b>extraordinarily horrendous</b> and I’m not going to waste any more words on it. <i>Label:</i> negative	<b>(-comprehensive)</b> <i>Review:</i> this film is [REDACTED] and I’m not going to waste any more words on it.
<i>Review:</i> this film is <b>extraordinarily horrendous</b> and I’m not going to <b>waste any more words on it</b> . <i>Label:</i> negative	<b>(comprehensive)</b> <i>Review:</i> this film is [REDACTED] and I’m not going to [REDACTED].
<i>Premise:</i> A shirtless man wearing white shorts. <i>Hypothesis:</i> A <b>man</b> in white shorts is <b>running on the sidewalk</b> . <i>Label:</i> neutral	<b>(-sufficient)</b> <i>Premise:</i> [REDACTED] <i>Hypothesis:</i> [REDACTED] <b>man</b> [REDACTED] <b>running on the sidewalk</b> .

Table 2: Examples of highlights differing in comprehensiveness and sufficiency (discussed in §2, §4).

35 We first define relevant EXNLP terminology (§2) and overview 61 existing datasets (§3), ac-  
36 companied with a live version of the tables as a website accepting community contributions:  
37 <https://exnlpdatasets.github.io>. We next analyze what can be learned from existing dataset  
38 collection methodologies. In particular, §4 discusses the traditional process of collecting explanations  
39 by instructing annotators to highlight parts of the input, and its discrepancies with evaluating model-  
40 generated highlight explanations. We also draw attention to how assumptions made in collection of  
41 free-text explanations (introduced in §2) influence their modeling, and call for better documentation  
42 of EXNLP data collection. In §5, we illustrate that not all template-like free-text explanations are  
43 unwanted, and call for embracing the structure of an explanation when appropriate. In §6, we present  
44 a proposal for controlling quality in explanation collection. Finally, §7 gathers recommendations  
45 from related subfields to further reduce data artifacts by increasing diversity of collected explanations.

## 46 2 Explainability Lexicon

47 An explanation can be described as a “three-place predicate: *someone* explains *something* to *someone*”  
48 [48]. The *something* being explained in machine learning systems are task labels: explanations are  
49 implicitly or explicitly designed to answer the question “why is [input] assigned [label]?”. However,  
50 collected explanations can vary in format. We identify three types in the EXNLP literature: *highlights*,  
51 *free-text*, and *structured* explanations. An example of each type is given in Table 1. Since a consensus  
52 on terminology has not yet been reached, we describe each type below.

53 **Highlights** are subsets of the input elements (words, phrases, or sentences) that explain a prediction.  
54 Lei et al. [71] coin them *extractive rationales*, or subsets of the input tokens of a textual task that  
55 satisfy two properties: (i) *compactness*, they are short and coherent, and (ii) *sufficiency*, they suffice  
56 for prediction as a substitute of the original text. Yu et al. [138] introduce a third criterion, (iii)  
57 *comprehensiveness*, that all the evidence that supports the prediction is selected, not just a sufficient  
58 set. Since the term “rationale” implies human-like intent, Jacovi and Goldberg [53] argue to call  
59 this type of explanation *highlights* to avoid inaccurately attributing human-like social behavior to AI  
60 systems. They are also called *evidence* in fact-checking and multi-document question answering (QA)  
61 [63]—a part of the source that refutes/supports the claim. To reiterate, highlights should be sufficient  
62 to explain a prediction and compact; if they are also comprehensive, we call them *comprehensive*  
63 *highlights*. Although the community has settled on criteria (i)–(iii) for highlights, the extent to

Dataset	Task	Granularity	Collection	# Instances
MOVIEREVIEWS [139]	sentiment classification	none	author	1,800
MOVIEREVIEWS <sub>c</sub> [28]	sentiment classification	none	crowd	200 <sup>‡</sup> ◇
SST [110]	sentiment classification	none	crowd	11,855◇
WIKIQA [133]	open-domain QA	sentence	crowd + authors	1,473
WIKIATTACK [21]	detecting personal attacks	none	students	1089◇
E-SNLI <sup>†</sup> [19]	natural language inference	none	crowd	~569K (1 or 3)
MULTIRC [58]	reading comprehension QA	sentences	crowd	5,825
FEVER [115]	verifying claims from text	sentences	crowd	~136K <sup>‡</sup>
HOTPOTQA [134]	reading comprehension QA	sentences	crowd	112,779
Hanselowski et al. [46]	verifying claims from text	sentences	crowd	6,422 (varies)
NATURALQUESTIONS [66]	reading comprehension QA	1 paragraph	crowd	n/a <sup>‡</sup> (1-5)
CoQA [101]	conversational QA	none	crowd	~127K (1 or 3)
COS-E v1.0 <sup>†</sup> [97]	commonsense QA	none	crowd	8,560
COS-E v1.11 <sup>†</sup> [97]	commonsense QA	none	crowd	10,962
BOOLQ <sub>c</sub> [28]	reading comprehension QA	none	crowd	199 <sup>‡</sup> ◇
EVIDENCEINFERENCE v1.0 [69]	evidence inference	none	experts	10,137
EVIDENCEINFERENCE v1.0 <sub>c</sub> [28]	evidence inference	none	experts	125 <sup>‡</sup>
EVIDENCEINFERENCE v2.0 [29]	evidence inference	none	experts	2,503
SciFACT [120]	verifying claims from text	1-3 sentences	experts	995 <sup>‡</sup> (1-3)
Kutlu et al. [65]	webpage relevance ranking	2-3 sentences	crowd	700 (15)
SCAT [136]	document-level machine translation	none	experts	~14K
ECTHR [23]	alleged legal violation prediction	paragraphs	auto + expert	~11K

Table 3: Overview of datasets with textual **highlights**. Values in parentheses indicate number of explanations collected per instance (if  $> 1$ ). DeYoung et al. [28] collected or recollected annotations for prior datasets (marked with the subscript  $c$ ). ◇ Collected  $> 1$  explanation per instance but only release 1. † Also contains free-text explanations. ‡ A subset of the original dataset that is annotated. It is not reported what subset of NATURALQUESTIONS has both a long and short answer.

64 which collected datasets (Table 3) reflect them varies greatly, as we will discuss in §4. Table 2 gives  
65 examples of sufficient vs. non-sufficient and comprehensive vs. non-comprehensive highlights.

66 **Free-text explanations** are free-form textual justifications that are not constrained to the words or  
67 modality of the input instance. They are thus more expressive and generally more readable than  
68 highlights. This makes them useful for explaining reasoning tasks where explanations must contain  
69 information outside the given input sentence or document [19, 125]. They are also called *textual* [60]  
70 or *natural language explanations* [19], terms that have been overloaded [95]. Synonyms, *free-form*  
71 [19] or *abstractive explanations* [84] do not emphasize that the explanation is textual.

72 Finally, **structured explanations** are explanations that are not entirely free-form although they are  
73 still written in natural language. For example, there may be constraints placed on the explanation-  
74 writing process, such as the required use of specific inference rules. We discuss the recent emergence  
75 of structured explanations in §5. Structured explanations do not have one common definition; we  
76 elaborate on dataset-specific designs in §3. An example is given in the bottom row of Table 1.

### 77 3 Overview of Existing Datasets

78 We overview currently available EXNLP datasets by explanation type: highlights (Table 3), free-text  
79 explanations (Table 4), and structured explanations (Table 5). To the best of our knowledge, all  
80 existing datasets are in English with the exception of SCAT [136]. The authors of 32.69% papers  
81 cited in Tables 3-5 do not report the dataset license in the paper or a repository, and 47.17% use  
82 *common* permissive licenses; for more information see Appendix B. See Appendix C for collection  
83 details.

84 For each dataset, we report the number of instances (input-label pairs) and the number of explanations  
85 per instance (if  $> 1$ ). The annotation procedure used to collect each dataset is reported as: crowd-  
86 annotated (“crowd”); automatically annotated through a web-scrape, database crawl, or merge of  
87 existing datasets (“auto”); or annotated by others (“experts”, “students”, or “authors”). Some authors  
88 perform semantic parsing on collected explanations (denoted with \*); we classify them by the dataset  
89 type before parsing and list the collection type as “crow + authors”. Tables 3-5 elucidate that the  
90 dominant collection paradigm ( $\geq 90\%$ ) is via human (crowd, student, author, or expert) annotation.

Dataset	Task	Collection	# Instances
Jansen et al. [54]	science exam QA	authors	363
Ling et al. [74]	solving algebraic word problems	auto + crowd	~101K
Srivastava et al. [112]*	detecting phishing emails	crowd + authors	7 (30-35)
BABBLELABBLE [45]*	relation extraction	students + authors	200 <sup>‡‡</sup>
E-SNLI [19]	natural language inference	crowd	~569K (1 or 3)
LIAR-PLUS [3]	verifying claims from text	auto	12,836
COS-E v1.0 [97]	commonsense QA	crowd	8,560
COS-E v1.11 [97]	commonsense QA	crowd	10,962
SEN-MAKING [121]	commonsense validation	students + authors	2,021
CHANGEMYVIEW [9]	argument persuasiveness	crowd	37,718
WINOWHY [141]	pronoun coreference resolution	crowd	273 (5)
SBIC [108]	social bias inference	crowd	48,923 (1-3)
PUBHEALTH [62]	verifying claims from text	auto	11,832
Wang et al. [122]*	relation extraction	crowd + authors	373
Wang et al. [122]*	sentiment classification	crowd + authors	85
E- $\delta$ -NLI [17]	defeasible natural language inference	auto	92,298 (~8)
BDD-X <sup>††</sup> [60]	vehicle control for self-driving cars	crowd	~26K
VQA-E <sup>††</sup> [73]	visual QA	auto	~270K
VQA-X <sup>††</sup> [91]	visual QA	crowd	28,180 (1 or 3)
ACT-X <sup>††</sup> [91]	activity recognition	crowd	18,030 (3)
Ehsan et al. [33] <sup>††</sup>	playing arcade games	crowd	2000
VCR <sup>††</sup> [140]	visual commonsense reasoning	crowd	~290K
E-SNLI-VE <sup>††</sup> [31]	visual-textual entailment	crowd	11,335 (3) <sup>‡</sup>
ESPRIT <sup>††</sup> [98]	reasoning about qualitative physics	crowd	2441 (2)
VLEP <sup>††</sup> [70]	future event prediction	auto + crowd	28,726
EMU <sup>††</sup> [26]	reasoning about manipulated images	crowd	48K

Table 4: Overview of ExNLP datasets with **free-text explanations** for textual and visual-textual tasks (marked with <sup>††</sup> and placed in the lower part). Values in parentheses indicate number of explanations collected per instance (if > 1). <sup>‡</sup> A subset of the original dataset that is annotated. <sup>‡‡</sup> Subset publicly available. \* Authors semantically parse the collected explanations.

91 **Highlights** (Table 3) The granularity of highlights depends on the task they are collected for. The  
92 majority of authors do not place a restriction on granularity, allowing words, phrases, or sentences  
93 of the original input document to be selected. The coarsest granularity in Table 3 is one or more  
94 paragraphs in a longer document [66, 23]. We exclude datasets that include an associated document  
95 as evidence without specifying the location of the explanation within the document (namely document  
96 retrieval datasets). We exclude BEERADVOCATE [77] because it has been retracted.

97 Some highlights are re-purposed from annotations for a different task. For example, MULTIRC [58]  
98 contains sentence-level highlights that indicate justifications of answers to questions. However,  
99 they were originally collected for the authors to assess that each question in the dataset requires  
100 multi-sentence reasoning to answer. Another example is STANFORD SENTIMENT TREEBANK [SST; 110]  
101 which contains crowdsourced sentiment annotations for word phrases extracted from movie reviews  
102 [87]. Word phrases that have the same sentiment label as the review can be heuristically merged to  
103 get phrase-level highlights [22]. Other highlights in Table 3 are collected by instructing annotators.  
104 Instead of giving these instructions verbatim, their authors typically describe them concisely, e.g.,  
105 they say annotators are asked to highlight words justifying, constituting, indicating, supporting, or  
106 determining the label, or words that are essential, useful, or relevant for the label. The difference in  
107 wording of these instructions affects how people annotate explanations. In §4, we discuss how one  
108 difference in annotation instructions (requiring comprehensiveness or not) can be important.

109 **Free-Text Explanations** (Table 4) This is a popular explanation type for both textual and visual-  
110 textual tasks, shown in the first and second half of the table, respectively. Most free-text explanations  
111 are generally no more than a few sentences per instance. One exception is LIAR-PLUS [4], which  
112 contains the conclusion paragraphs of web-scraped human-written fact-checking summaries.

113 **Structured Explanations** (Table 5) Structured explanations take on dataset-specific forms. One  
114 common approach is to construct a chain of facts that detail the reasoning steps to reach an answer  
115 (“chains of facts”). Another is to place constraints on the textual explanations that annotators can  
116 write, such as requiring the use of certain variables in the input (“semi-structured text”).

Dataset	Task	Explanation Type	Collection	# Instances
WORLDTREE V1 [55]	science exam QA	explanation graphs	authors	1,680
OPENBOOKQA [78]	open-book science QA	1 fact from WORLDTREE	crowd	5,957
Yang et al. [132] <sup>††</sup>	action recognition	lists of relations + attributes	crowd	853
WORLDTREE V2 [129]	science exam QA	explanation graphs	experts	5,100
QED [68]	reading comp. QA	inference rules	authors	8,991
QASC [59]	science exam QA	2-fact chain	authors + crowd	9,980
EQASC [56]	science exam QA	2-fact chain	auto + crowd	9,980 (~10)
+ PERTURBED	science exam QA	2-fact chain	auto + crowd	n/a <sup>‡</sup>
EOBQA [56]	open-book science QA	2-fact chain	auto + crowd	n/a <sup>‡</sup>
Ye et al. [135]*	SQUAD QA	semi-structured text	crowd + authors	164
Ye et al. [135]*	NATURALQUESTIONS QA	semi-structured text	crowd + authors	109
R <sup>4</sup> C [51]	reading comp. QA	chains of facts	crowd	4,588 (3)
STRATEGYQA [40]	implicit reasoning QA	reasoning steps w/ highlights	crowd	2,780 (3)

Table 5: Overview of EXNLP datasets with **structured explanations** (§5). Values in parentheses indicate number of explanations collected per instance (if > 1). †† Visual-textual dataset. \* Authors semantically parse the collected explanations. ‡ Subset of instances annotated with explanations is not reported. Total # of explanations is 855 for EQASC PERTURBED and 998 for EOBQA.

117 The WORLDTREE datasets [55, 129] propose explaining elementary-school science questions with  
 118 a combination of chains of facts and semi-structured text, termed “explanation graphs”. The facts  
 119 are individual sentences written by the authors that are centered around a set of shared relations and  
 120 properties. Given the chain of facts for an instance (6.3 facts on average), the authors can construct  
 121 an explanation graph by linking shared words in the question, answer, and explanation.

122 OPENBOOKQA [OBQA; 78] uses single WORLDTREE facts to prime annotators to write QA pairs.  
 123 Similarly, each question in QASC [59] contains two associated science facts from a corpus selected  
 124 by human annotators who wrote the question. Jhamtani and Clark [56] extend OBQA and QASC with  
 125 two-fact chain explanation annotations, which are automatically extracted from a fact corpus and  
 126 validated with crowdsourcing. The resulting datasets, EQASC and EOBQA, contain multiple valid and  
 127 invalid explanations per instance, as well as perturbations for robustness testing (EQASC-PERTURBED).

128 A number of structured explanation datasets supplement datasets for reading comprehension. Ye  
 129 et al. [135] collect semi-structured explanations for NATURALQUESTIONS [66] and SQUAD [99]. They  
 130 require annotators to use phrases in both the input question and context, and limit them to a small set  
 131 of connecting expressions. Inoue et al. [51] collect R<sup>4</sup>C, fact chain explanations for HOTPOTQA [134].  
 132 Lamm et al. [68] collect explanations for NATURALQUESTIONS that follow a linguistically-motivated  
 133 form (see the example in Table 1). We discuss structured explanations further in §5.

#### 134 4 Link Between EXNLP Data, Modeling, and Evaluation Assumptions

135 All three parts of the machine learning pipeline (data collection, modeling, and evaluation) are  
 136 inextricably linked. In this section, we discuss what EXNLP modeling and evaluation research reveals  
 137 about the qualities of available EXNLP datasets, and how best to collect such datasets in the future.

138 Highlights are usually evaluated following two criteria: (i) *plausibility*, according to humans, how  
 139 well a highlight supports a predicted label [130, 28], and (ii) *faithfulness* or *fidelity*, how accurately a  
 140 highlight represents the model’s decision process [5, 124]. Human-annotated highlights (Table 2) are  
 141 used to measure the plausibility of model-produced highlights: the higher the agreement between the  
 142 two, the more plausible model highlights are considered. On the other hand, a highlight that is both  
 143 sufficient (implies the prediction, §2, first example in Table 2 in Appendix A) and comprehensive (its  
 144 complement in the input does *not* imply the prediction, §2, second example in Table 2) is regarded as  
 145 faithful to the prediction it explains [28, 22]. Since human-annotated highlights are used only for  
 146 evaluation of plausibility but not faithfulness, it might seem that the measurement and modeling of  
 147 faithfulness cannot influence how human-annotated explanations should be collected. This might  
 148 lead to collecting highlights that are not suitable for the goals (ii) and (iii) in §1.

149 Typical instructions for collecting highlights encourage sufficiency and compactness, but not compre-  
 150 hensiveness. For example, DeYoung et al. [28] deem MOVIEREVIEWS and EVIDENCEINFERENCE high-  
 151 lights non-comprehensive. Carton et al. [22] expect that FEVER highlights are non-comprehensive, in

152 contrast to DeYoung et al. [28]. Contrary to the characterization of both of these work, we observe that  
153 the E-SNLI authors collect non-comprehensive highlights, since they instruct annotators to highlight  
154 only words in the hypothesis (and not the premise) for neutral pairs, and consider contradiction/neutral  
155 explanations correct if at least one piece of evidence in the input is highlighted. Based on these  
156 discrepancies in characterization, we first conclude that post-hoc assessment of comprehensiveness  
157 from a general description of data collection is error-prone.

158 Alternatively, Carton et al. [22] empirically show that available human highlights are not necessarily  
159 sufficient nor comprehensive for predictions of highly accurate models. This suggests that the same  
160 might hold for gold labels, leading us to ask: are gold highlights in existing datasets flawed?

161 Let us first consider insufficiency. Highlighted input elements taken together have to reasonably  
162 indicate the label. Otherwise, a highlight is an invalid explanation. Consider two datasets whose  
163 sufficiency Carton et al. [22] found to be most concerning: neutral E-SNLI pairs and no-attack  
164 WIKIATTACK examples. Neutral E-SNLI cases are not justifiable by highlighting because they  
165 are obtained only as an intermediate step to collecting free-text explanations, and only free-text  
166 explanations truly justify a neutral label [19]. Table 2 shows one E-SNLI highlight that is not sufficient.  
167 No-attack WIKIATTACK examples are not explainable by highlighting because the absence of offensive  
168 content justifies the no-attack label, and this absence cannot be highlighted. We recommend (i)  
169 avoiding human-annotated highlights with low sufficiency when evaluating and collecting highlights,  
170 and (ii) assessing whether the true label can be explained by highlighting.

171 Consider a highlight that is non-comprehensive because it is redundant with its complement in the  
172 input (e.g., a word appears multiple times, but only one occurrence is highlighted). Highlighting  
173 only one occurrence of “great” is a valid justification, but quantifying faithfulness of this highlight  
174 is hard because the model might rightfully use the unhighlighted occurrence of “great” to make  
175 the prediction. Thus, comprehensiveness is modeled to make faithfulness evaluation feasible. Non-  
176 comprehensiveness of human highlights, however, hinders evaluating plausibility of comprehensive  
177 model highlights since model and human highlights do not match by design. To be able to eval-  
178 uate both plausibility and faithfulness, we should annotate comprehensive human highlights. We  
179 summarize these observations in Figure 2 in Appendix A.

180 Mutual influence of data and modeling assumptions also affects free-text explanations. For example,  
181 the E-SNLI guidelines have far more constraints than the COS-E guidelines, such as requiring self-  
182 contained explanations. Wiegrefe et al. [125] show that such data collection decisions can influence  
183 modeling assumptions. This is not an issue per se, but we should be cautious that EXNLP data  
184 collection decisions do not popularize explanation properties as *universally necessary* when they  
185 are not, e.g., that free-text explanations should be understandable without the original input or that  
186 highlights should be comprehensive. We believe this could be avoided with better documentation,  
187 e.g., with additions to a standard datasheet [38]. Explainability fact sheets have been proposed  
188 for models [111], but not for datasets. For example, an E-SNLI datasheet could note that self-  
189 contained explanations were required during data collection, but that this is not a necessary property  
190 of a valid free-text explanation. A dataset with comprehensive highlights should emphasize that  
191 comprehensiveness is required to simplify faithfulness evaluation.

## 192 Takeaways

- 193 1. It is important to precisely report how explanations were collected, e.g., by giving access to  
194 the annotation interface, screenshotting it, or giving the annotation instructions verbatim.
- 195 2. Sufficiency is necessary for highlights, and EXNLP researchers should avoid human-  
196 annotated highlights with low sufficiency for evaluating and developing highlights.
- 197 3. Comprehensiveness isn’t necessary for a valid highlight, it is a means to quantify faithfulness.
- 198 4. Non-comprehensive human-annotated highlights cannot be used to automatically evaluate  
199 plausibility of highlights that are constrained to be comprehensive. In this case, EXNLP  
200 researchers should collect and use comprehensive human-annotated highlights.
- 201 5. When deciding which datasets to use, EXNLP researchers should not make post-hoc esti-  
202 mates of comprehensiveness of human-annotated highlights from datasets’ general descrip-  
203 tions since that is error-prone.
- 204 6. EXNLP researchers should be careful to not popularize their data collection decisions as  
205 universally necessary. We advocate for documenting all constraints on collected explanations  
206 in a datasheet, highlighting whether each constraint is necessary for explanation to be valid

207 or not, and noting how each constraint might affect modeling and evaluation to the best of  
208 the author’s knowledge.

## 209 5 Rise of Structured Explanations

210 The merit of free-text explanations is their expressivity, which can come at the costs of underspecifica-  
211 tion and inconsistency due to the difficulty of quality control (stressed by the creators of two popular  
212 free-text explanation datasets: E-SNLI and COS-E). In this section, we highlight and challenge one  
213 prior approach to overcoming these difficulties: discarding template-like free-text explanations.

214 We gather crowdsourcing guidelines for the above-mentioned datasets in Tables 6-7 in Appendix and  
215 compare them. We observe two notable similarities between the guidelines for the above-mentioned  
216 datasets. First, both asked annotators to first highlight input words and then formulate a free-text  
217 explanation from them, to control quality. Second, template-like explanations are discarded because  
218 they are deemed uninformative. The E-SNLI authors assembled a list of 56 templates (e.g., “There  
219 is ⟨hypothesis⟩”) to identify explanations whose edit distance to one of the templates is  $<10$ . They  
220 re-annotate the detected template-like explanations (11% in the entire dataset). The COS-E authors  
221 discard sentences “⟨answer⟩ is the only option that is correct/obvious” (the only given example  
222 of a template). Template explanations concern researchers because they can result in artifact-like  
223 behaviors in certain modeling architectures. For example, a model which predicts a task output from  
224 a generated explanation can produce explanations that are plausible to a human user and give the  
225 impression of making label predictions on the basis of this explanation. However, it is possible that  
226 the model learns to ignore the semantics of the explanation and instead makes predictions based on  
227 the explanation’s template type [64, 53]. In this case, the semantic interpretation of the explanation  
228 (that of a human reader) is not faithful (an accurate representation of the model’s decision process).

229 Despite re-annotating, Camburu et al. [20] report that E-SNLI explanations still largely follow  
230 28 label-specific templates (e.g., an entailment template “X is another form of Y”) even after re-  
231 annotation. Brahman et al. [17] report that models trained on gold E-SNLI explanations generate  
232 template-like explanations for the defeasible NLI task. These findings lead us to ask: what are the  
233 differences between templates considered uninformative and filtered out, and those identified by  
234 Camburu et al. [20], Brahman et al. [17] that remain after filtering? Are *all* template-like explanations  
235 uninformative?

236 Although prior work indicates that template-like explanations are undesirable, most recently, struc-  
237 tured explanations have been intentionally collected (see Table 5; §3). What these studies share is  
238 that they acknowledge structure as *inherent* to explaining the tasks they investigate. Related work  
239 [GLUCOSE; 82] takes the matter further, arguing that explanations should not be entirely free-form.  
240 Following GLUCOSE, we recommend running pilot studies to explore how people define and generate  
241 explanations for a task *before* collecting free-text explanations for it. If they reveal that informative  
242 human explanations are naturally structured, incorporating the structure in the annotation scheme is  
243 useful since the structure is natural to explaining the task. This turned out to be the case with NLI;  
244 Camburu et al. [20] report: “Explanations in E-SNLI largely follow a set of label-specific templates.  
245 This is a *natural consequence of the task* and dataset”. We recommend embracing the structure when  
246 possible, but also encourage creators of datasets with template-like explanations to highlight in a  
247 dataset datasheet (§4) that template structure can influence downstream modeling decisions. There  
248 is no all-encompassing definition of explanation, and researchers could consult domain experts or  
249 follow literature from other fields to define an appropriate explanation in a task-specific manner, such  
250 as in GLUCOSE [82]. For conceptualization of explanations in different fields see Tiddi et al. [116].

251 Finally, what if pilot studies do not reveal any obvious structure to human explanations of a task?  
252 Then we need to do our best to control the quality of free-text explanations because low dataset quality  
253 is a bottleneck to building high-quality models. COS-E is collected with notably less annotation  
254 constraints and quality controls than E-SNLI, and has annotation issues that some have deemed make  
255 the dataset unusable [84]; see examples in Table 7 of Appendix A. As exemplars of quality control,  
256 we point the reader to the annotation guidelines of VCR [140] in Table 8 and GLUCOSE [81]. In §6  
257 and §7, we give further task-agnostic recommendations for collecting high-quality EXNLP datasets,  
258 applicable to all three explanation types.

259 **Takeaways**

- 260 1. ExNLP researchers should study how people define and generate explanations for the task  
261 before collecting free-text explanations.
- 262 2. If pilot studies show that explanations are naturally structured, embrace the structure.
- 263 3. There is no all-encompassing definition of explanation. Thus, ExNLP researchers could con-  
264 sult domain experts or follow literature from other fields to define an appropriate explanation  
265 form, and these matters should be open for debate on a given task.

## 266 6 Increasing Explanation Quality

267 When asked to write free-text sentences from scratch for a table-to-text annotation task outside  
268 ExNLP, Parikh et al. [89] note that crowdworkers produce “vanilla targets that lack [linguistic]  
269 variety”. Lack of variety can result in annotation artifacts, which are prevalent in the popular SNLI  
270 [15] and MNLI [126] datasets [94, 44, 117], among others [39]. These authors demonstrate the  
271 harms of such artifacts: models can overfit to them, leading to both performance over-estimation and  
272 problematic generalization behaviors.

273 Artifacts can occur from poor-quality annotations and inattentive annotators, both of which have been  
274 on the rise on crowdsourcing platforms [24, 6, 84]. To mitigate artifacts, both increased **diversity of**  
275 **annotators** and **quality control** are needed. We focus on quality control here and diversity in §7

### 276 6.1 A Two-Stage Collect-And-Edit Approach

277 While ad-hoc methods can improve quality [19, 140, 81], an effective and generalizable method is to  
278 collect annotations in two stages. A two-stage methodology has been applied by a small minority  
279 of ExNLP dataset papers [56, 141, 140], who first compile explanation candidates automatically or  
280 from crowdworkers, and secondly perform quality-control by having other crowdworkers assess the  
281 quality of the collected explanations (we term this COLLECT-AND-JUDGE). Judging improves the  
282 overall quality of the final dataset by removing low-quality instances, and additionally allows authors  
283 to release quality ratings for each instance.

284 Outside ExNLP, Parikh et al. [89] use an extended version of this approach (that we term COLLECT-  
285 AND-EDIT): they generate a noisy automatically-extracted dataset for the table-to-text generation task,  
286 and then ask annotators to edit the datapoints. Bowman et al. [16] use this approach to re-collect  
287 NLI hypotheses, and find, crucially, that having annotators edit rather than create hypotheses reduces  
288 artifacts in a subset of MNLI. In XAI, Kutlu et al. [65] collect highlight explanations for Web page  
289 ranking with annotator editing. We advocate expanding the COLLECT-AND-JUDGE approach for  
290 explanation collection to COLLECT-AND-EDIT. This has potential to increase linguistic diversity via  
291 multiple annotators per-instance, reduce individual annotator biases, and perform quality control.  
292 Through a case study of two multimodal free-text explanation datasets, we will demonstrate that  
293 collecting explanations automatically without human editing (or at least judging) can lead to artifacts.

294 E-SNLI-VE [31] and VQA-E [73] are two visual-textual datasets for entailment and question-  
295 answering, respectively. E-SNLI-VE combines annotations of two datasets: (i) SNLI-VE [128],  
296 collected by replacing the textual premises of SNLI [15] with FLICKR30K images [137], and (ii) E-  
297 SNLI [19], a dataset of crowdsourced explanations for SNLI. This procedure is possible because every  
298 SNLI premise was originally the caption of a FLICKR30K photo. However, since SNLI’s hypotheses  
299 were collected from crowdworkers who did not see the original images, the photo replacement process  
300 results in a significant number of errors [119]. Do et al. [31] re-annotate labels and explanations for  
301 the neutral pairs in the validation and test sets of SNLI-VE. However, it has been argued that the  
302 dataset remains low-quality for training models due to artifacts in the entailment and the neutral class’  
303 training sets [76]. With a full EDIT approach, we expect that these artifacts would be significantly  
304 reduced, and the resulting dataset could have quality on-par with E-SNLI. Similarly, the VQA-E  
305 dataset [73] converts image captions from the VQA v2.0 dataset [42] into explanations, but a notably  
306 lower plausibility compared to a carefully-crowdsourced VCR explanations is reported in [76].

307 Both E-SNLI-VE and VQA-E present novel and cost-effective ways to produce large ExNLP datasets  
308 for new tasks, but also show the quality tradeoffs of automatic collection. Strategies such as  
309 crowdsourced judging and editing, even on a small subset, can reveal and mitigate such issues.



## 310 6.2 Teach and Test the Underlying Task

311 In order to both create and judge explanations, annotators must understand the underlying task and  
312 label-set well. In most cases, this necessitates teaching and testing the task. Prior work outside  
313 of ExNLP has noted the difficulty of scaling annotation to crowdworkers for complex linguistic  
314 tasks [103, 34, 96, 82]. To increase annotation quality, these works provide intensive training  
315 to crowdworkers, including personal feedback. Since label understanding is a prerequisite for  
316 explanation collection, task designers should consider relatively inexpensive strategies such as  
317 qualification tasks and checker questions. This need is correlated with the difficulty and domain-  
318 specificity of the task, as elaborated above.

319 Similarly, people cannot explain all tasks equally well and even after intensive training they might  
320 struggle to explain tasks such as deception detection and recidivism prediction [86]. Human explana-  
321 tions for such tasks might be limited in serving the three goals outlined in §1.

## 322 6.3 Addressing Ambiguity

323 Data collectors often collect explanations post-hoc, i.e., annotators are asked to explain labels assigned  
324 by a system or other annotators. The underlying assumption is that the explainer believes the assigned  
325 label to be correct or at least likely (there is no task ambiguity). However, this assumption has been  
326 shown to be inaccurate (among others) for relation extraction [7], natural language inference [93, 85],  
327 and complement coercion [34], and the extent to which it is true likely varies by task, instance, and  
328 annotator. If an annotator is uncertain about a label, their explanation may be at best a hypothesis and  
329 at worst a guess. HCI research encourages leaving room for ambiguity rather than forcing raters into  
330 binary decisions, which can result in poor or inaccurate labels [105].

331 To ensure explanations reflect human decisions as closely as possible, it is ideal to collect both labels  
332 and explanations from the same annotators. Given that this is not always possible, including a checker  
333 question to assess whether an explanation annotator agrees with a label is a good alternative.

### 334 Takeaways

- 335 1. Using a COLLECT-AND-EDIT method can reduce individual annotator biases, perform quality  
336 control, and potentially reduce dataset artifacts.
- 337 2. Teaching and testing the underlying task and addressing ambiguity can improve data quality.

## 338 7 Increasing Explanation Diversity

339 Beyond quality control, increasing annotation diversity is another task-agnostic means to mitigate  
340 artifacts and collect more representative data. We elaborate on suggestions from related work (inside  
341 and outside ExNLP) here.

### 342 7.1 Use a Large Set of Annotators

343 Collecting representative data entails ensuring that a handful of annotators do not dominate data  
344 collection. Outside ExNLP, Geva et al. [39] report that recruiting only a small pool of annotators  
345 (1 annotator per 100–1000 examples) allows models to overfit on annotator characteristics. Such  
346 small annotator pools exist in ExNLP—for instance, E-SNLI reports an average of 860 explanations  
347 written per worker. The occurrence of the incorrect explanation “rivers flow trough valleys” for  
348 529 different instances in COS-E v1.11 is likely attributed to a single annotator. Al Kuwatly et al.  
349 [2] find that demographic attributes can predict annotation differences. Similarly, Davidson et al.  
350 [27], Sap et al. [107] show that annotators often consider African-American English writing to be  
351 disproportionately offensive.<sup>2</sup> A lack of annotator representation concerns ExNLP for three reasons:  
352 explanations depend on socio-cultural background [61], annotator traits should not be predictable  
353 [39], and the subjectivity of explaining leaves room for social bias to emerge.

354 On most platforms, annotators are not restricted to a specific number of instances. Verifying that no  
355 worker has annotated an excessively large portion of the dataset in addition to strategies from Geva

---

<sup>2</sup>In another related study, 82% of annotators reported their race as white [108]. This is a likely explanation for the disproportionate annotation.

356 et al. [39] can help mitigate annotator bias. More elaborate methods for increasing annotator diversity  
357 include collecting demographic attributes or modeling annotators as a graph [2, 123].

## 358 7.2 Multiple Annotations Per Instance

359 HCI research has long considered the ideal of crowdsourcing a single ground-truth as a “myth” that  
360 fails to account for the diversity of human thought and experience [8]. Similarly, ExNLP researchers  
361 should not assume there is always one correct explanation. Many of the assessments crowdworkers  
362 are asked to make when writing explanations are subjective in nature, and there are many different  
363 models of explanation based on a user’s cognitive biases, social expectations, and socio-cultural  
364 background [79]. Prasad et al. [95] present a theoretical argument to illustrate that there are multiple  
365 ways to highlight input words to explain an annotated sentiment label. Camburu et al. [19] find a low  
366 inter-annotator BLEU score [88] between free-text explanations collected for E-SNLI test instances.

367 If a dataset contains only one explanation when multiple are plausible, a plausible model explanation  
368 can be penalized unfairly for not agreeing with it. We expect that modeling multiple explanations can  
369 also be a useful learning signal. Some existing datasets contain multiple explanations per instance  
370 (last column of Tables 3-5). Future ExNLP data collections should do the same if there is subjectivity  
371 in the task or diversity of correct explanations (which can be measured via inter-annotator agreement).  
372 If annotators exhibit low agreement between explanations deemed as plausible, this can reveal a  
373 diversity of correct explanations for the task, which should be considered in modeling and evaluation.

## 374 7.3 Get Ahead: Add Contrastive and Negative Explanations

375 The machine learning community has championed modeling *contrastive explanations* that justify  
376 why a prediction was made *instead of* another, to align more closely with human explanation  
377 [30, 47, 79]. Most recently, methods have been proposed in NLP to produce contrastive edits of the  
378 input as explanations [104, 131, 127, 53]. Outside of ExNLP, datasets with contrastive edits have  
379 been collected to assess and improve robustness of NLP models [57, 37, 72] and might be used for  
380 explainability too.

381 Just as highlights are not sufficiently intelligible for complex tasks, the same might hold for contrastive  
382 input edits. To the best of our knowledge, there is no dataset that contains contrastive free-text or  
383 structured explanations. These could take the form of (i) collecting explanations that answer the  
384 question “why...instead of...”, or (ii) collecting explanations for other labels besides the gold label,  
385 to be used as an additional training signal. A related annotation paradigm is to collect *negative*  
386 *explanations*, i.e., explanations that are invalid for an (input, gold label) pair. Such examples can  
387 improve ExNLP models by providing supervision of what is *not* a correct explanation [109]. A  
388 human JUDGE or EDIT phase automatically gives negative explanations: the low-scoring instances  
389 (former) or instances pre-editing (latter) [56, 141].

## 390 Takeaways

- 391 1. To increase annotation diversity, a large set of annotators, multiple annotations per instance,  
392 and collecting explanations that are most useful to the needs of end-users are important.
- 393 2. Reporting inter-annotator agreement with plausibility of annotated explanations is useful  
394 to know whether there is a natural diversity of explanations for the task and should the  
395 diversity be considered for modeling and evaluation.

## 396 8 Conclusions

397 We have presented a review of existing datasets for ExNLP research, highlighted discrepancies in  
398 data collection that can have downstream modeling effects, and synthesized the literature both inside  
399 and outside ExNLP into a set of recommendations for future data collection.

400 We note that a majority of the work reviewed in this paper has originated in the last 1-2 years,  
401 indicating an explosion of interest in collecting datasets for ExNLP. We provide reflections for current  
402 and future data collectors in an effort to promote standardization and consistency. This paper also  
403 serves as a starting resource for newcomers to ExNLP, and, we hope, a starting point for further  
404 discussions.

## References

- 405
- 406 [1] Amina Adadi and Mohammed Berrada. Peeking inside the black-box: A survey on explainable  
407 artificial intelligence (XAI). *IEEE Access*, 6:52138–52160, 2018. doi: 10.1109/ACCESS.  
408 2018.2870052.
- 409 [2] Hala Al Kuwatly, Maximilian Wich, and Georg Groh. Identifying and measuring annotator  
410 bias based on annotators’ demographic characteristics. In *Proceedings of the Fourth Workshop  
411 on Online Abuse and Harms*, pages 184–190, Online, November 2020. Association for  
412 Computational Linguistics. doi: 10.18653/v1/2020.alw-1.21. URL <https://www.aclweb.org/anthology/2020.alw-1.21>
- 413
- 414 [3] Tariq Alhindi, Savvas Petridis, and Smaranda Muresan. Where is your evidence: Im-  
415 proving fact-checking by justification modeling. In *Proceedings of the First Workshop  
416 on Fact Extraction and VERification (FEVER)*, pages 85–90, Brussels, Belgium, Novem-  
417 ber 2018. Association for Computational Linguistics. doi: 10.18653/v1/W18-5513. URL  
418 <https://www.aclweb.org/anthology/W18-5513>.
- 419 [4] Tariq Alhindi, Smaranda Muresan, and Daniel Preotiuc-Pietro. Fact vs. opinion: the role  
420 of argumentation features in news classification. In *Proceedings of the 28th International  
421 Conference on Computational Linguistics*, pages 6139–6149, Barcelona, Spain (Online),  
422 December 2020. International Committee on Computational Linguistics. URL <https://www.aclweb.org/anthology/2020.coling-main.540>.
- 423
- 424 [5] David Alvarez-Melis and T. Jaakkola. Towards robust interpretability with self-explaining  
425 neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.  
426 URL <https://arxiv.org/abs/1806.07538>.
- 427 [6] Antonio Alonso Arechar and David Rand. Turking in the time of covid. PsyArXiv, 2020.  
428 URL <https://psyarxiv.com/vktqu>.
- 429 [7] Lora Aroyo and Chris Welty. Crowd truth: Harnessing disagreement in crowdsourcing a  
430 relation extraction gold standard. *WebSci2013. ACM*, 2013(2013), 2013.
- 431 [8] Lora Aroyo and Chris Welty. Truth is a lie: Crowd truth and the seven myths of human  
432 annotation. *AI Magazine*, 36(1):15–24, 2015. URL [https://ojs.aaai.org//index.php/  
433 aimagazine/article/view/2564](https://ojs.aaai.org//index.php/aimagazine/article/view/2564).
- 434 [9] David Atkinson, Kumar Bhargav Srinivasan, and Chenhao Tan. What gets echoed? un-  
435 derstanding the “pointers” in explanations of persuasive arguments. In *Proceedings of the  
436 2019 Conference on Empirical Methods in Natural Language Processing and the 9th In-  
437 ternational Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages  
438 2911–2921, Hong Kong, China, November 2019. Association for Computational Linguistics.  
439 doi: 10.18653/v1/D19-1289. URL <https://www.aclweb.org/anthology/D19-1289>.
- 440 [10] Yujia Bao, Shiyu Chang, Mo Yu, and Regina Barzilay. Deriving machine attention from human  
441 rationales. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language  
442 Processing*, pages 1903–1913, Brussels, Belgium, October-November 2018. Association for  
443 Computational Linguistics. doi: 10.18653/v1/D18-1216. URL [https://www.aclweb.org/  
444 anthology/D18-1216](https://www.aclweb.org/anthology/D18-1216).
- 445 [11] Alejandro Barredo Arrieta, Natalia Díaz-Rodríguez, Javier Del Ser, Adrien Bénézet, Si-  
446 ham Tabik, Alberto Barbado, Salvador Garcia, Sergio Gil-Lopez, Daniel Molina, Richard  
447 Benjamins, Raja Chatila, and Francisco Herrera. Explainable artificial intelligence (XAI):  
448 Concepts, taxonomies, opportunities and challenges toward responsible ai. *Information Fusion*,  
449 58:82–115, 2020. ISSN 1566-2535. doi: <https://doi.org/10.1016/j.inffus.2019.12.012>. URL  
450 <https://www.sciencedirect.com/science/article/pii/S1566253519308103>.
- 451 [12] Emily M. Bender and Batya Friedman. Data statements for natural language processing:  
452 Toward mitigating system bias and enabling better science. *Transactions of the Association  
453 for Computational Linguistics*, 6:587–604, 2018. doi: 10.1162/tacl\_a\_00041. URL [https:  
454 //www.aclweb.org/anthology/Q18-1041](https://www.aclweb.org/anthology/Q18-1041).
- 455 [13] Chandra Bhagavatula, Ronan Le Bras, Chaitanya Malaviya, Keisuke Sakaguchi, Ari Holtzman,  
456 Hannah Rashkin, Doug Downey, Wen-tau Yih, and Yejin Choi. Abductive commonsense  
457 reasoning. In *International Conference on Learning Representations (ICLR)*, 2020. URL  
458 <https://arxiv.org/abs/1908.05739>.

- 459 [14] Or Biran and Courtenay Cotton. Explanation and justification in machine learning: A survey.  
460 In *IJCAI Workshop on Explainable AI (XAI)*, pages 8–13, 2017. URL [http://www.cs.columbia.edu/~orb/papers/xai\\_survey\\_paper\\_2017.pdf](http://www.cs.columbia.edu/~orb/papers/xai_survey_paper_2017.pdf).
- 462 [15] Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. A  
463 large annotated corpus for learning natural language inference. In *Proceedings of the 2015  
464 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon,  
465 Portugal, September 2015. Association for Computational Linguistics. doi: 10.18653/v1/  
466 D15-1075. URL <https://www.aclweb.org/anthology/D15-1075>.
- 467 [16] Samuel R. Bowman, Jennimaria Palomaki, Livio Baldini Soares, and Emily Pitler. New  
468 protocols and negative results for textual entailment data collection. In *Proceedings of  
469 the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*,  
470 pages 8203–8214, Online, November 2020. Association for Computational Linguistics. doi:  
471 10.18653/v1/2020.emnlp-main.658. URL [https://www.aclweb.org/anthology/2020.  
472 emnlp-main.658](https://www.aclweb.org/anthology/2020.emnlp-main.658).
- 473 [17] Faeze Brahman, Vered Shwartz, Rachel Rudinger, and Yejin Choi. Learning to rationalize  
474 for nonmonotonic reasoning with distant supervision. In *the AAAI Conference on Artificial  
475 Intelligence*, 2021. URL <https://arxiv.org/abs/2012.08012>.
- 476 [18] Nadia Burkart and Marco F. Huber. A survey on the explainability of supervised machine  
477 learning. *The Journal of Artificial Intelligence Research (JAIR)*, 70, 2021. doi: [https://www.  
478 jair.org/index.php/jair/article/view/12228](https://www.jair.org/index.php/jair/article/view/12228).
- 479 [19] Oana-Maria Camburu, Tim Rocktäschel, Thomas Lukasiewicz, and Phil Blunsom. e-SNLI:  
480 Natural language inference with natural language explanations. In *Advances in Neural In-  
481 formation Processing Systems (NeurIPS)*, 2018. URL [https://papers.nips.cc/paper/  
482 2018/hash/4c7a167bb329bd92580a99ce422d6fa6-Abstract.html](https://papers.nips.cc/paper/2018/hash/4c7a167bb329bd92580a99ce422d6fa6-Abstract.html).
- 483 [20] Oana-Maria Camburu, Brendan Shillingford, Pasquale Minervini, Thomas Lukasiewicz, and  
484 Phil Blunsom. Make up your mind! adversarial generation of inconsistent natural language  
485 explanations. In *Proceedings of the 58th Annual Meeting of the Association for Computational  
486 Linguistics*, pages 4157–4165, Online, July 2020. Association for Computational Linguistics.  
487 doi: 10.18653/v1/2020.acl-main.382. URL [https://www.aclweb.org/anthology/2020.  
488 acl-main.382](https://www.aclweb.org/anthology/2020.acl-main.382).
- 489 [21] Samuel Carton, Qiaozhu Mei, and Paul Resnick. Extractive adversarial networks: High-recall  
490 explanations for identifying personal attacks in social media posts. In *Proceedings of the  
491 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3497–3507,  
492 Brussels, Belgium, October–November 2018. Association for Computational Linguistics. doi:  
493 10.18653/v1/D18-1386. URL <https://www.aclweb.org/anthology/D18-1386>.
- 494 [22] Samuel Carton, Anirudh Rathore, and Chenhao Tan. Evaluating and characterizing hu-  
495 man rationales. In *Proceedings of the 2020 Conference on Empirical Methods in Nat-  
496 ural Language Processing (EMNLP)*, pages 9294–9307, Online, November 2020. Ass-  
497 ociation for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.747. URL  
498 <https://www.aclweb.org/anthology/2020.emnlp-main.747>.
- 499 [23] Ilias Chalkidis, Manos Fergadiotis, Dimitrios Tsarapatsanis, Nikolaos Aletras, Ion Androu-  
500 topoulos, and Prodromos Malakasiotis. Paragraph-level rationale extraction through regular-  
501 ization: A case study on European court of human rights cases. In *Proceedings of the 2021  
502 Conference of the North American Chapter of the Association for Computational Linguistics:  
503 Human Language Technologies*, pages 226–241, Online, June 2021. Association for Computa-  
504 tional Linguistics. URL <https://www.aclweb.org/anthology/2021.naacl-main.22>.
- 505 [24] Michael Chmielewski and Sarah C Kucker. An MTurk crisis? Shifts in data quality and the  
506 impact on study results. *Social Psychological and Personality Science*, 11(4):464–473, 2020.  
507 URL <https://journals.sagepub.com/doi/abs/10.1177/1948550619875149>.
- 508 [25] Miruna-Adriana Clinciu and Helen Hastie. A survey of explainable AI terminology. In *Proceed-  
509 ings of the 1st Workshop on Interactive Natural Language Technology for Explainable Artificial  
510 Intelligence (NLXAI 2019)*, pages 8–13. Association for Computational Linguistics, 2019.  
511 doi: 10.18653/v1/W19-8403. URL <https://www.aclweb.org/anthology/W19-8403>.
- 512 [26] Jeff Da, M. Forbes, Rowan Zellers, Anthony Zheng, Jena D. Hwang, Antoine Bosselut, and  
513 Yejin Choi. Edited media understanding: Reasoning about implications of manipulated images.  
514 arXiv:2012.04726, 2020. URL <https://arxiv.org/abs/2012.04726>.

- 515 [27] Thomas Davidson, Debasmita Bhattacharya, and Ingmar Weber. Racial bias in hate speech  
516 and abusive language detection datasets. In *Proceedings of the Third Workshop on Abusive*  
517 *Language Online*, pages 25–35, Florence, Italy, August 2019. Association for Computational  
518 Linguistics. doi: 10.18653/v1/W19-3504. URL [https://www.aclweb.org/anthology/  
519 W19-3504](https://www.aclweb.org/anthology/W19-3504).
- 520 [28] Jay DeYoung, Sarthak Jain, Nazneen Fatema Rajani, Eric Lehman, Caiming Xiong, Richard  
521 Socher, and Byron C. Wallace. ERASER: A benchmark to evaluate rationalized NLP mod-  
522 els. In *Proceedings of the 58th Annual Meeting of the Association for Computational Lin-*  
523 *guistics*, pages 4443–4458, Online, July 2020. Association for Computational Linguistics.  
524 doi: 10.18653/v1/2020.acl-main.408. URL [https://www.aclweb.org/anthology/2020.  
525 acl-main.408](https://www.aclweb.org/anthology/2020.acl-main.408).
- 526 [29] Jay DeYoung, Eric Lehman, Benjamin Nye, Iain Marshall, and Byron C. Wallace. Evidence  
527 inference 2.0: More data, better models. In *Proceedings of the 19th SIGBioMed Workshop*  
528 *on Biomedical Language Processing*, pages 123–132, Online, July 2020. Association for  
529 Computational Linguistics. doi: 10.18653/v1/2020.bionlp-1.13. URL [https://www.aclweb.  
530 org/anthology/2020.bionlp-1.13](https://www.aclweb.org/anthology/2020.bionlp-1.13).
- 531 [30] Amit Dhurandhar, Pin-Yu Chen, Ronny Luss, Chun-Chen Tu, Paishun Ting, Karthikeyan  
532 Shanmugam, and Payel Das. Explanations based on the missing: Towards contrastive expla-  
533 nations with pertinent negatives. In *Proceedings of the 32nd International Conference on*  
534 *Neural Information Processing Systems*, pages 590–601, 2018. URL [https://proceedings.  
535 neurips.cc/paper/2018/file/c5ff2543b53f4cc0ad3819a36752467b-Paper.pdf](https://proceedings.neurips.cc/paper/2018/file/c5ff2543b53f4cc0ad3819a36752467b-Paper.pdf).
- 536 [31] Virginie Do, Oana-Maria Camburu, Zeynep Akata, and Thomas Lukasiewicz. e-SNLI-VE-  
537 2.0: Corrected Visual-Textual Entailment with Natural Language Explanations. In *IEEE*  
538 *CVPR Workshop on Fair, Data Efficient and Trusted Computer Vision*, 2020. URL [https:  
539 //arxiv.org/abs/2004.03744](https://arxiv.org/abs/2004.03744).
- 540 [32] Finale Doshi-Velez and Been Kim. Towards a rigorous science of interpretable machine  
541 learning. arXiv:1702.08608, 2017. URL <https://arxiv.org/abs/1702.08608>.
- 542 [33] Upol Ehsan, Pradyumna Tambwekar, Larry Chan, Brent Harrison, and Mark O. Riedl. Auto-  
543 mated rationale generation: A technique for explainable AI and its effects on human percep-  
544 tions. In *Proceedings of the Conference of Intelligent User Interfaces (ACM IUI)*, 2019. URL  
545 <https://arxiv.org/abs/1901.03729>.
- 546 [34] Yanai Elazar, Victoria Basmov, Shauli Ravfogel, Yoav Goldberg, and Reut Tsarfaty. The  
547 extraordinary failure of complement coercion crowdsourcing. In *Proceedings of the First*  
548 *Workshop on Insights from Negative Results in NLP*, pages 106–116, Online, November  
549 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.insights-1.17. URL  
550 <https://www.aclweb.org/anthology/2020.insights-1.17>.
- 551 [35] Angela Fan, Yacine Jernite, Ethan Perez, David Grangier, Jason Weston, and Michael Auli.  
552 ELI5: Long form question answering. In *Proceedings of the 57th Annual Meeting of the*  
553 *Association for Computational Linguistics*, pages 3558–3567, Florence, Italy, July 2019.  
554 Association for Computational Linguistics. doi: 10.18653/v1/P19-1346. URL [https://www.  
555 aclweb.org/anthology/P19-1346](https://www.aclweb.org/anthology/P19-1346).
- 556 [36] Maxwell Forbes, Jena D. Hwang, Vered Shwartz, Maarten Sap, and Yejin Choi. Social  
557 chemistry 101: Learning to reason about social and moral norms. In *Proceedings of the 2020*  
558 *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 653–670,  
559 Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.  
560 emnlp-main.48. URL <https://www.aclweb.org/anthology/2020.emnlp-main.48>.
- 561 [37] Matt Gardner, Yoav Artzi, Victoria Basmov, Jonathan Berant, Ben Bogin, Sihao Chen,  
562 Pradeep Dasigi, Dheeru Dua, Yanai Elazar, Ananth Gottumukkala, Nitish Gupta, Han-  
563 naneh Hajishirzi, Gabriel Ilharco, Daniel Khoshdel, Kevin Lin, Jiangming Liu, Nelson F.  
564 Liu, Phoebe Mulcaire, Qiang Ning, Sameer Singh, Noah A. Smith, Sanjay Subrama-  
565 nian, Reut Tsarfaty, Eric Wallace, Ally Zhang, and Ben Zhou. Evaluating models’ local  
566 decision boundaries via contrast sets. In *Findings of the Association for Computa-*  
567 *tional Linguistics: EMNLP 2020*, pages 1307–1323, Online, November 2020. Associa-  
568 tion for Computational Linguistics. doi: 10.18653/v1/2020.findings-emnlp.117. URL  
569 <https://www.aclweb.org/anthology/2020.findings-emnlp.117>.

- 570 [38] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna M.  
571 Wallach, Hal Daumé, and Kate Crawford. Datasheets for datasets. In *Proceedings of the 5th*  
572 *Workshop on Fairness, Accountability, and Transparency in Machine Learning*, 2018. URL  
573 [https://www.fatml.org/media/documents/datasheets\\_for\\_datasets.pdf](https://www.fatml.org/media/documents/datasheets_for_datasets.pdf)
- 574 [39] Mor Geva, Yoav Goldberg, and Jonathan Berant. Are we modeling the task or the annotator?  
575 an investigation of annotator bias in natural language understanding datasets. In *Proceedings*  
576 *of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th*  
577 *International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages  
578 1161–1166, Hong Kong, China, November 2019. Association for Computational Linguistics.  
579 doi: 10.18653/v1/D19-1107. URL <https://www.aclweb.org/anthology/D19-1107>.
- 580 [40] Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. Did  
581 aristotle use a laptop? a question answering benchmark with implicit reasoning strategies.  
582 *Transactions of the Association for Computational Linguistics*, 2021. URL [https://arxiv](https://arxiv.org/pdf/2101.02235.pdf)  
583 [.org/pdf/2101.02235.pdf](https://arxiv.org/pdf/2101.02235.pdf).
- 584 [41] Leilani H. Gilpin, David Bau, B. Yuan, A. Bajwa, M. Specter, and Lalana Kagal. Explaining  
585 explanations: An overview of interpretability of machine learning. *2018 IEEE 5th International*  
586 *Conference on Data Science and Advanced Analytics (DSAA)*, pages 80–89, 2018. URL  
587 <https://arxiv.org/pdf/1806.00069.pdf>.
- 588 [42] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and D. Parikh. Making the V  
589 in VQA Matter: Elevating the Role of Image Understanding in Visual Question Answering.  
590 In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6325–  
591 6334, 2017. URL <https://arxiv.org/abs/1612.00837>.
- 592 [43] Riccardo Guidotti, A. Monreale, F. Turini, D. Pedreschi, and F. Giannotti. A survey of methods  
593 for explaining black box models. *ACM Computing Surveys (CSUR)*, 51:1 – 42, 2019. URL  
594 <https://dl.acm.org/doi/pdf/10.1145/3236009>.
- 595 [44] Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel Bowman, and  
596 Noah A. Smith. Annotation artifacts in natural language inference data. In *Proceedings of*  
597 *the 2018 Conference of the North American Chapter of the Association for Computational*  
598 *Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 107–112, New  
599 Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/  
600 N18-2017. URL <https://www.aclweb.org/anthology/N18-2017>.
- 601 [45] Braden Hancock, Paroma Varma, Stephanie Wang, Martin Bringmann, Percy Liang, and  
602 Christopher Ré. Training classifiers with natural language explanations. In *Proceedings of*  
603 *the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long*  
604 *Papers)*, pages 1884–1895, Melbourne, Australia, July 2018. Association for Computational  
605 Linguistics. doi: 10.18653/v1/P18-1175. URL [https://www.aclweb.org/anthology/](https://www.aclweb.org/anthology/P18-1175)  
606 [P18-1175](https://www.aclweb.org/anthology/P18-1175).
- 607 [46] Andreas Hanselowski, Christian Stab, Claudia Schulz, Zile Li, and Iryna Gurevych. A  
608 richly annotated corpus for different tasks in automated fact-checking. In *Proceedings of*  
609 *the 23rd Conference on Computational Natural Language Learning (CoNLL)*, pages 493–  
610 503, Hong Kong, China, November 2019. Association for Computational Linguistics. doi:  
611 10.18653/v1/K19-1046. URL <https://www.aclweb.org/anthology/K19-1046>.
- 612 [47] Lisa Anne Hendricks, Ronghang Hu, Trevor Darrell, and Zeynep Akata. Generating counter-  
613 factual explanations with natural language. In *International Conference on Machine Learning*  
614 *(ICML)*, 2018.
- 615 [48] Denis J Hilton. Conversational processes and causal explanation. *Psy-*  
616 *chological Bulletin*, 107(1):65, 1990. URL [https://www.researchgate](https://www.researchgate.net/profile/Denis_Hilton/publication/232543382_Conversational_processes_and_causal_explanation/links/00b7d519bd8fa613f1000000/Conversational-processes-and-causal-explanation.pdf)  
617 [.net/profile/Denis\\_Hilton/publication/232543382\\_Conversational\\_](https://www.researchgate.net/profile/Denis_Hilton/publication/232543382_Conversational_processes_and_causal_explanation/links/00b7d519bd8fa613f1000000/Conversational-processes-and-causal-explanation.pdf)  
618 [processes\\_and\\_causal\\_explanation/links/00b7d519bd8fa613f1000000/](https://www.researchgate.net/profile/Denis_Hilton/publication/232543382_Conversational_processes_and_causal_explanation/links/00b7d519bd8fa613f1000000/Conversational-processes-and-causal-explanation.pdf)  
619 [Conversational-processes-and-causal-explanation.pdf](https://www.researchgate.net/profile/Denis_Hilton/publication/232543382_Conversational_processes_and_causal_explanation/links/00b7d519bd8fa613f1000000/Conversational-processes-and-causal-explanation.pdf)
- 620 [49] Robert R Hoffman, Shane T Mueller, Gary Klein, and Jordan Litman. Metrics for explainable  
621 AI: Challenges and prospects. arXiv:1812.04608, 2018. URL [https://arxiv.org/abs/](https://arxiv.org/abs/1812.04608)  
622 [1812.04608](https://arxiv.org/abs/1812.04608).
- 623 [50] Lifu Huang, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Cosmos QA: Machine  
624 reading comprehension with contextual commonsense reasoning. In *Proceedings of the 2019*

- 625 *Conference on Empirical Methods in Natural Language Processing and the 9th International*  
626 *Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2391–2401,  
627 Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.  
628 18653/v1/D19-1243. URL <https://www.aclweb.org/anthology/D19-1243>.
- [51] Naoya Inoue, Pontus Stenetorp, and Kentaro Inui. R4C: A benchmark for evaluating RC  
629 systems to get the right answer for the right reason. In *Proceedings of the 58th Annual*  
630 *Meeting of the Association for Computational Linguistics*, pages 6740–6750, Online, July  
631 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.602. URL  
632 <https://www.aclweb.org/anthology/2020.acl-main.602>.
- [52] Alon Jacovi and Yoav Goldberg. Towards faithfully interpretable NLP systems: How should  
634 we define and evaluate faithfulness? In *Proceedings of the 58th Annual Meeting of the*  
635 *Association for Computational Linguistics*, pages 4198–4205, Online, July 2020. Association  
636 for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.386. URL <https://www.aclweb.org/anthology/2020.acl-main.386>.
- [53] Alon Jacovi and Yoav Goldberg. Aligning faithful interpretations with their social attribution.  
639 *Transactions of the Association for Computational Linguistics*, 2021. URL <https://arxiv.org/abs/2006.01067>.
- [54] Peter Jansen, Niranjan Balasubramanian, Mihai Surdeanu, and Peter Clark. What’s in an  
642 explanation? characterizing knowledge and inference requirements for elementary science  
643 exams. In *Proceedings of COLING 2016, the 26th International Conference on Computational*  
644 *Linguistics: Technical Papers*, pages 2956–2965, Osaka, Japan, December 2016. The COLING  
645 2016 Organizing Committee. URL <https://www.aclweb.org/anthology/C16-1278>.
- [55] Peter Jansen, Elizabeth Wainwright, Steven Marmorstein, and Clayton Morrison. WorldTree:  
647 A corpus of explanation graphs for elementary science questions supporting multi-hop in-  
648 ference. In *Proceedings of the Eleventh International Conference on Language Resources*  
649 *and Evaluation (LREC 2018)*, Miyazaki, Japan, May 2018. European Language Resources  
650 Association (ELRA). URL <https://www.aclweb.org/anthology/L18-1433>.
- [56] Harsh Jhamtani and Peter Clark. Learning to explain: Datasets and models for identifying  
652 valid reasoning chains in multihop question-answering. In *Proceedings of the 2020 Con-*  
653 *ference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 137–150,  
654 Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.  
655 emnlp-main.10. URL <https://www.aclweb.org/anthology/2020.emnlp-main.10>.
- [57] Divyansh Kaushik, Eduard Hovy, and Zachary Lipton. Learning the difference that makes a  
657 difference with counterfactually-augmented data. In *International Conference on Learning*  
658 *Representations (ICLR)*, 2020. URL <https://openreview.net/pdf?id=SkIgsONFvr>.
- [58] Daniel Khashabi, Snigdha Chaturvedi, Michael Roth, Shyam Upadhyay, and Dan Roth. Look-  
660 ing beyond the surface: A challenge set for reading comprehension over multiple sentences.  
661 In *Proceedings of the 2018 Conference of the North American Chapter of the Association for*  
662 *Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages  
663 252–262, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi:  
664 10.18653/v1/N18-1023. URL <https://www.aclweb.org/anthology/N18-1023>.
- [59] Tushar Khot, Peter Clark, Michal Guerquin, Peter Jansen, and Ashish Sabharwal. QASC:  
666 A dataset for question answering via sentence composition. In *Proceedings of the AAAI*  
667 *Conference on Artificial Intelligence*, volume 34, pages 8082–8090, 2020. URL <https://arxiv.org/pdf/1910.11473.pdf>.
- [60] Jinkyu Kim, Anna Rohrbach, Trevor Darrell, John F. Canny, and Zeynep Akata. Textual  
670 Explanations for Self-Driving Vehicles. In *Proceedings of the European Conference on*  
671 *Computer Vision (ECCV)*, 2018. URL <https://arxiv.org/abs/1807.11546>.
- [61] Hana Kopecká and Jose M Such. Explainable AI for Cultural Minds. In *Workshop on*  
673 *Dialogue, Explanation and Argumentation for Human-Agent Interaction (DEXA HAI) at the*  
674 *24th European Conference on Artificial Intelligence (ECAI)*, 2020. URL [https://kclpure.kcl.ac.uk/portal/files/134728815/DEXA\\_aug\\_crc.pdf](https://kclpure.kcl.ac.uk/portal/files/134728815/DEXA_aug_crc.pdf).
- [62] Neema Kotonya and Francesca Toni. Explainable automated fact-checking for public  
677 health claims. In *Proceedings of the 2020 Conference on Empirical Methods in Nat-*  
678 *ural Language Processing (EMNLP)*, pages 7740–7754, Online, November 2020. Asso-  
679

- 680 ciation for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.623. URL  
681 <https://www.aclweb.org/anthology/2020.emnlp-main.623>.
- 682 [63] Neema Kotonya and Francesca Toni. Explainable automated fact-checking: A survey. In  
683 *Proceedings of the 28th International Conference on Computational Linguistics*, pages 5430–  
684 5443, Barcelona, Spain (Online), December 2020. International Committee on Computational  
685 Linguistics. doi: 10.18653/v1/2020.coling-main.474. URL <https://www.aclweb.org/anthology/2020.coling-main.474>.
- 686
- 687 [64] Sawan Kumar and Partha Talukdar. NILE : Natural language inference with faithful  
688 natural language explanations. In *Proceedings of the 58th Annual Meeting of the As-*  
689 *sociation for Computational Linguistics*, pages 8730–8742, Online, July 2020. Associa-  
690 tion for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.771. URL <https://www.aclweb.org/anthology/2020.acl-main.771>.
- 691
- 692 [65] Mucahid Kutlu, Tyler McDonnell, Matthew Lease, and Tamer Elsayed. Annotator rationales  
693 for labeling tasks in crowdsourcing. *Journal of Artificial Intelligence Research*, 69:143–189,  
694 2020. URL [https://www.ischool.utexas.edu/~ml/papers/kutlu\\_jair20.pdf](https://www.ischool.utexas.edu/~ml/papers/kutlu_jair20.pdf).
- 695 [66] Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh,  
696 Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova,  
697 Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc  
698 Le, and Slav Petrov. Natural questions: A benchmark for question answering research.  
699 *Transactions of the Association for Computational Linguistics*, 7:452–466, March 2019. doi:  
700 10.1162/tacl\_a\_00276. URL <https://www.aclweb.org/anthology/Q19-1026>.
- 701 [67] Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. RACE: Large-scale  
702 ReAding comprehension dataset from examinations. In *Proceedings of the 2017 Conference on*  
703 *Empirical Methods in Natural Language Processing*, pages 785–794, Copenhagen, Denmark,  
704 September 2017. Association for Computational Linguistics. doi: 10.18653/v1/D17-1082.  
705 URL <https://www.aclweb.org/anthology/D17-1082>.
- 706 [68] Matthew Lamm, Jennimaria Palomaki, Chris Alberti, Daniel Andor, Eunsol Choi, Livio Baldini  
707 Soares, and Michael Collins. QED: A framework and dataset for explanations in question  
708 answering. *arXiv preprint arXiv:2009.06354*, 2020. URL [https://arxiv.org/pdf/2009.](https://arxiv.org/pdf/2009.06354.pdf)  
709 [06354.pdf](https://arxiv.org/pdf/2009.06354.pdf).
- 710 [69] Eric Lehman, Jay DeYoung, Regina Barzilay, and Byron C. Wallace. Inferring which medical  
711 treatments work from reports of clinical trials. In *Proceedings of the 2019 Conference of the*  
712 *North American Chapter of the Association for Computational Linguistics: Human Language*  
713 *Technologies, Volume 1 (Long and Short Papers)*, pages 3705–3717, Minneapolis, Minnesota,  
714 June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1371. URL  
715 <https://www.aclweb.org/anthology/N19-1371>.
- 716 [70] Jie Lei, Licheng Yu, Tamara Berg, and Mohit Bansal. What is more likely to happen next?  
717 video-and-language future event prediction. In *Proceedings of the 2020 Conference on Empiri-*  
718 *cal Methods in Natural Language Processing (EMNLP)*, pages 8769–8784, Online, November  
719 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.706.  
720 URL <https://www.aclweb.org/anthology/2020.emnlp-main.706>.
- 721 [71] Tao Lei, Regina Barzilay, and Tommi Jaakkola. Rationalizing neural predictions. In *Proceed-*  
722 *ings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages  
723 107–117, Austin, Texas, November 2016. Association for Computational Linguistics. doi:  
724 10.18653/v1/D16-1011. URL <https://www.aclweb.org/anthology/D16-1011>.
- 725 [72] Chuanrong Li, Lin Shengshuo, Zeyu Liu, Xinyi Wu, Xuhui Zhou, and Shane Steinert-Threlkeld.  
726 Linguistically-informed transformations (LIT): A method for automatically generating contrast  
727 sets. In *Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural*  
728 *Networks for NLP*, pages 126–135, Online, November 2020. Association for Computational  
729 Linguistics. doi: 10.18653/v1/2020.blackboxnlp-1.12. URL <https://www.aclweb.org/anthology/2020.blackboxnlp-1.12>.
- 730
- 731 [73] Qing Li, Qingyi Tao, Shafiq R. Joty, Jianfei Cai, and Jiebo Luo. VQA-E: Explaining, Elaborat-  
732 ing, and Enhancing Your Answers for Visual Questions. In *Proceedings of the European Con-*  
733 *ference on Computer Vision (ECCV)*, 2018. URL <https://arxiv.org/abs/1803.07464>.



- 734 [74] Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. Program induction by rationale  
735 generation: Learning to solve and explain algebraic word problems. In *Proceedings of the 55th*  
736 *Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*,  
737 pages 158–167, Vancouver, Canada, July 2017. Association for Computational Linguistics.  
738 doi: 10.18653/v1/P17-1015. URL <https://www.aclweb.org/anthology/P17-1015>.
- 739 [75] Zachary C Lipton. The mythos of model interpretability. *Queue*, 16(3):31–57, 2018. URL  
740 <https://dl.acm.org/doi/pdf/10.1145/3236386.3241340>.
- 741 [76] Ana Marasović, Chandra Bhagavatula, Jae sung Park, Ronan Le Bras, Noah A. Smith, and Yejin  
742 Choi. Natural language rationales with full-stack visual reasoning: From pixels to semantic  
743 frames to commonsense graphs. In *Findings of the Association for Computational Linguistics:*  
744 *EMNLP 2020*, pages 2810–2829, Online, November 2020. Association for Computational  
745 Linguistics. doi: 10.18653/v1/2020.findings-emnlp.253. URL [https://www.aclweb.org/](https://www.aclweb.org/anthology/2020.findings-emnlp.253)  
746 [anthology/2020.findings-emnlp.253](https://www.aclweb.org/anthology/2020.findings-emnlp.253).
- 747 [77] Julian McAuley, Jure Leskovec, and Dan Jurafsky. Learning attitudes and attributes from  
748 multi-aspect reviews. In *2012 IEEE 12th International Conference on Data Mining, 2012*.  
749 URL <https://ieeexplore.ieee.org/document/6413815>.
- 750 [78] Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor  
751 conduct electricity? a new dataset for open book question answering. In *Proceedings of the*  
752 *2018 Conference on Empirical Methods in Natural Language Processing*, pages 2381–2391,  
753 Brussels, Belgium, October–November 2018. Association for Computational Linguistics. doi:  
754 10.18653/v1/D18-1260. URL <https://www.aclweb.org/anthology/D18-1260>.
- 755 [79] Tim Miller. Explanation in artificial intelligence: Insights from the social sciences. *Artificial*  
756 *intelligence*, 267:1–38, 2019. URL <https://arxiv.org/pdf/1706.07269.pdf>.
- 757 [80] Christoph Molnar. *Interpretable Machine Learning*. 2019. [https://christophm.github.](https://christophm.github.io/interpretable-ml-book/)  
758 [io/interpretable-ml-book/](https://christophm.github.io/interpretable-ml-book/).
- 759 [81] Lori Moon, Lauren Berkowitz, Jennifer Chu-Carroll, and Nasrin Mostafazadeh. Details of  
760 data collection and crowd management for glucose (generalized and contextualized story ex-  
761 planations). *Github*, 2020. URL [https://github.com/ElementalCognition/glucose/](https://github.com/ElementalCognition/glucose/blob/master/data_collection_quality.pdf)  
762 [blob/master/data\\_collection\\_quality.pdf](https://github.com/ElementalCognition/glucose/blob/master/data_collection_quality.pdf).
- 763 [82] Nasrin Mostafazadeh, Aditya Kalyanpur, Lori Moon, David Buchanan, Lauren Berkowitz,  
764 Or Biran, and Jennifer Chu-Carroll. GLUCOSE: Generalized and Contextualized story  
765 explanations. In *Proceedings of the 2020 Conference on Empirical Methods in Natural*  
766 *Language Processing (EMNLP)*, pages 4569–4586, Online, November 2020. Association for  
767 Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.370. URL [https://www.](https://www.aclweb.org/anthology/2020.emnlp-main.370)  
768 [aclweb.org/anthology/2020.emnlp-main.370](https://www.aclweb.org/anthology/2020.emnlp-main.370).
- 769 [83] W. James Murdoch, Chandan Singh, Karl Kumbier, Reza Abbasi-Asl, and Bin Yu. Definitions,  
770 methods, and applications in interpretable machine learning. *Proceedings of the National*  
771 *Academy of Sciences*, 116(44):22071–22080, 2019. ISSN 0027-8424. doi: 10.1073/pnas.  
772 1900654116. URL <https://www.pnas.org/content/116/44/22071>.
- 773 [84] Sharan Narang, Colin Raffel, Katherine Lee, Adam Roberts, Noah Fiedel, and Karishma  
774 Malkan. WT5?! Training Text-to-Text Models to Explain their Predictions. arXiv:2004.14546,  
775 2020. URL <https://arxiv.org/abs/2004.14546>.
- 776 [85] Yixin Nie, Xiang Zhou, and Mohit Bansal. What can we learn from collective human opinions  
777 on natural language inference data? In *Proceedings of the 2020 Conference on Empirical*  
778 *Methods in Natural Language Processing (EMNLP)*, pages 9131–9143, Online, November  
779 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.734.  
780 URL <https://www.aclweb.org/anthology/2020.emnlp-main.734>.
- 781 [86] R. Nisbett and T. Wilson. Telling more than we can know: Verbal reports on mental processes.  
782 *Psychological Review*, 84:231–259, 1977.
- 783 [87] Bo Pang and Lillian Lee. Seeing stars: Exploiting class relationships for sentiment cate-  
784 gorization with respect to rating scales. In *Proceedings of the 43rd Annual Meeting of the*  
785 *Association for Computational Linguistics (ACL’05)*, pages 115–124, Ann Arbor, Michigan,  
786 June 2005. Association for Computational Linguistics. doi: 10.3115/1219840.1219855. URL  
787 <https://www.aclweb.org/anthology/P05-1015>.

- 788 [88] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic  
789 evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association*  
790 *for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA, July 2002.  
791 Association for Computational Linguistics. doi: 10.3115/1073083.1073135. URL <https://www.aclweb.org/anthology/P02-1040>.  
792
- 793 [89] Ankur Parikh, Xuezhi Wang, Sebastian Gehrmann, Manaal Faruqui, Bhuwan Dhingra, Diyi  
794 Yang, and Dipanjan Das. ToTTo: A controlled table-to-text generation dataset. In *Pro-*  
795 *ceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*  
796 *(EMNLP)*, pages 1173–1186, Online, November 2020. Association for Computational Linguistics.  
797 doi: 10.18653/v1/2020.emnlp-main.89. URL [https://www.aclweb.org/anthology/](https://www.aclweb.org/anthology/2020.emnlp-main.89)  
798 [2020.emnlp-main.89](https://www.aclweb.org/anthology/2020.emnlp-main.89).
- 799 [90] Praveen Paritosh. Achieving data excellence. In *NeurIPS 2020 Crowd Science Workshop*,  
800 2020. URL [https://neurips.cc/virtual/2020/public/workshop\\_16111.html](https://neurips.cc/virtual/2020/public/workshop_16111.html). In-  
801 vited talk.
- 802 [91] Dong Huk Park, Lisa Anne Hendricks, Zeynep Akata, Anna Rohrbach, Bernt  
803 Schiele, Trevor Darrell, and Marcus Rohrbach. Multimodal explanations: Justi-  
804 fying decisions and pointing to the evidence. In *Proceedings of the IEEE Con-*  
805 *ference on Computer Vision and Pattern Recognition*, pages 8779–8788, 2018.  
806 URL [https://openaccess.thecvf.com/content\\_cvpr\\_2018/papers/Park\\_](https://openaccess.thecvf.com/content_cvpr_2018/papers/Park_Multimodal_Explanations_Justifying_CVPR_2018_paper.pdf)  
807 [Multimodal\\_Explanations\\_Justifying\\_CVPR\\_2018\\_paper.pdf](https://openaccess.thecvf.com/content_cvpr_2018/papers/Park_Multimodal_Explanations_Justifying_CVPR_2018_paper.pdf).
- 808 [92] Amandalynne Paullada, Inioluwa Deborah Raji, Emily M. Bender, Emily L. Denton, and  
809 A. Hanna. Data and its (dis)contents: A survey of dataset development and use in machine  
810 learning research. In *The ML-Retrospectives, Surveys & Meta-Analyses NeurIPS 2020 Work-*  
811 *shop*, 2020. URL <https://arxiv.org/abs/2012.05345>.
- 812 [93] Ellie Pavlick and Tom Kwiatkowski. Inherent disagreements in human textual inferences.  
813 *Transactions of the Association for Computational Linguistics*, 7:677–694, March 2019. doi:  
814 10.1162/tacl\_a\_00293. URL <https://www.aclweb.org/anthology/Q19-1043>.
- 815 [94] Adam Poliak, Jason Naradowsky, Aparajita Haldar, Rachel Rudinger, and Benjamin  
816 Van Durme. Hypothesis only baselines in natural language inference. In *Proceedings of the Sev-*  
817 *enth Joint Conference on Lexical and Computational Semantics*, pages 180–191, New Orleans,  
818 Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/S18-2023.  
819 URL <https://www.aclweb.org/anthology/S18-2023>.
- 820 [95] Grusha Prasad, Yixin Nie, Mohit Bansal, Robin Jia, Douwe Kiela, and Adina Williams. To  
821 what extent do human explanations of model behavior align with actual model behavior?  
822 arXiv:2012.13354, 2020. URL <https://arxiv.org/abs/2012.13354>.
- 823 [96] Valentina Pyatkin, Ayal Klein, Reut Tsarfaty, and Ido Dagan. QADiscourse - Discourse  
824 Relations as QA Pairs: Representation, Crowdsourcing and Baselines. In *Proceedings of*  
825 *the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*,  
826 pages 2804–2819, Online, November 2020. Association for Computational Linguistics. doi:  
827 10.18653/v1/2020.emnlp-main.224. URL [https://www.aclweb.org/anthology/2020.](https://www.aclweb.org/anthology/2020.emnlp-main.224)  
828 [emnlp-main.224](https://www.aclweb.org/anthology/2020.emnlp-main.224).
- 829 [97] Nazneen Fatema Rajani, Bryan McCann, Caiming Xiong, and Richard Socher. Explain  
830 yourself! leveraging language models for commonsense reasoning. In *Proceedings of the 57th*  
831 *Annual Meeting of the Association for Computational Linguistics*, pages 4932–4942, Florence,  
832 Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1487. URL  
833 <https://www.aclweb.org/anthology/P19-1487>.
- 834 [98] Nazneen Fatema Rajani, Rui Zhang, Yi Chern Tan, Stephan Zheng, Jeremy Weiss, Aadit Vyas,  
835 Abhijit Gupta, Caiming Xiong, Richard Socher, and Dragomir Radev. ESPRIT: Explaining  
836 solutions to physical reasoning tasks. In *Proceedings of the 58th Annual Meeting of the*  
837 *Association for Computational Linguistics*, pages 7906–7917, Online, July 2020. Association  
838 for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.706. URL [https://www.](https://www.aclweb.org/anthology/2020.acl-main.706)  
839 [aclweb.org/anthology/2020.acl-main.706](https://www.aclweb.org/anthology/2020.acl-main.706).
- 840 [99] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+  
841 questions for machine comprehension of text. In *Proceedings of the 2016 Conference on*

- 842 *Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas, November  
843 2016. Association for Computational Linguistics. doi: 10.18653/v1/D16-1264. URL  
844 <https://www.aclweb.org/anthology/D16-1264>.
- 845 [100] Gabriëlle Ras, Marcel van Gerven, and Pim Haselager. *Explanation Methods in Deep Learning:  
846 Users, Values, Concerns and Challenges*, pages 19–36. Springer International Publishing,  
847 Cham, 2018. ISBN 978-3-319-98131-4. doi: 10.1007/978-3-319-98131-4\_2. URL [https://doi.org/10.1007/978-3-319-98131-4\\_2](https://doi.org/10.1007/978-3-319-98131-4_2).
- 849 [101] Siva Reddy, Danqi Chen, and Christopher D. Manning. CoQA: A conversational question  
850 answering challenge. *Transactions of the Association for Computational Linguistics*, 7:  
851 249–266, March 2019. doi: 10.1162/tacl\_a\_00266. URL <https://www.aclweb.org/anthology/Q19-1016>.
- 853 [102] Matthew Richardson, Christopher J.C. Burges, and Erin Renshaw. MCTest: A challenge  
854 dataset for the open-domain machine comprehension of text. In *Proceedings of the 2013  
855 Conference on Empirical Methods in Natural Language Processing*, pages 193–203, Seattle,  
856 Washington, USA, October 2013. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/D13-1020>.
- 858 [103] Paul Roit, Ayal Klein, Daniela Stepanov, Jonathan Mamou, Julian Michael, Gabriel Stanovsky,  
859 Luke Zettlemoyer, and Ido Dagan. Controlled crowdsourcing for high-quality QA-SRL  
860 annotation. In *Proceedings of the 58th Annual Meeting of the Association for Computational  
861 Linguistics*, pages 7008–7013, Online, July 2020. Association for Computational Linguistics.  
862 doi: 10.18653/v1/2020.acl-main.626. URL <https://www.aclweb.org/anthology/2020.acl-main.626>.
- 864 [104] Alexis Ross, Ana Marasović, and Matthew E Peters. Explaining nlp models via minimal  
865 contrastive editing (mice). *arXiv preprint arXiv:2012.13985*, 2020. URL <https://arxiv.org/pdf/2012.13985.pdf>.
- 867 [105] Nithya Sambasivan. Human-data interaction in ai. In *PAIR Symposium*, 2020. URL <https://www.youtube.com/watch?v=cjRF5a4eo2Y&t=83s>. Invited talk.
- 869 [106] Nithya Sambasivan, Shivani Kapania, Hannah Highfill, Diana Akrong, Praveen Kumar  
870 Paritosh, and Lora Mois Aroyo. "everyone wants to do the model  
871 work, not the data work": Data cascades in high-stakes ai. In *Proceedings  
872 of the 2021 CHI Conference on Human Factors in Computing Systems*, 2021.  
873 URL <https://storage.googleapis.com/pub-tools-public-publication-data/pdf/0d556e45afc54afeb2eb6b51a9bc1827b9961ff4.pdf>.
- 875 [107] Maarten Sap, Dallas Card, Saadia Gabriel, Yejin Choi, and Noah A. Smith. The risk of racial  
876 bias in hate speech detection. In *Proceedings of the 57th Annual Meeting of the Association  
877 for Computational Linguistics*, pages 1668–1678, Florence, Italy, July 2019. Association for  
878 Computational Linguistics. doi: 10.18653/v1/P19-1163. URL <https://www.aclweb.org/anthology/P19-1163>.
- 880 [108] Maarten Sap, Saadia Gabriel, Lianhui Qin, Dan Jurafsky, Noah A. Smith, and Yejin Choi.  
881 Social bias frames: Reasoning about social and power implications of language. In *Proceedings  
882 of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5477–  
883 5490, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.486. URL <https://www.aclweb.org/anthology/2020.acl-main.486>.
- 885 [109] Hendrik Schuff, Heike Adel, and Ngoc Thang Vu. F1 is Not Enough! Models and Evaluation  
886 Towards User-Centered Explainable Question Answering. In *Proceedings of the 2020 Confer-  
887 ence on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7076–7095,  
888 Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.575. URL <https://www.aclweb.org/anthology/2020.emnlp-main.575>.
- 890 [110] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng,  
891 and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment  
892 treebank. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language  
893 Processing*, pages 1631–1642, Seattle, Washington, USA, October 2013. Association for  
894 Computational Linguistics. URL <https://www.aclweb.org/anthology/D13-1170>.
- 895 [111] Kacper Sokol and Peter Flach. Explainability fact sheets: a framework for systematic as-  
896 sessment of explainable approaches. In *Proceedings of the 2020 Conference on Fairness*,

- 897     *Accountability, and Transparency*, pages 56–67, 2020. URL [https://dl.acm.org/doi/](https://dl.acm.org/doi/pdf/10.1145/3351095.3372870)  
898     [pdf/10.1145/3351095.3372870](https://dl.acm.org/doi/pdf/10.1145/3351095.3372870).
- 899 [112] Shashank Srivastava, Igor Labutov, and Tom Mitchell. Joint concept learning and semantic  
900 parsing from natural language explanations. In *Proceedings of the 2017 Conference on*  
901 *Empirical Methods in Natural Language Processing*, pages 1527–1536, Copenhagen, Denmark,  
902 September 2017. Association for Computational Linguistics. doi: 10.18653/v1/D17-1161.  
903 URL <https://www.aclweb.org/anthology/D17-1161>.
- 904 [113] Julia Strout, Ye Zhang, and Raymond Mooney. Do human rationales improve machine  
905 explanations? In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and*  
906 *Interpreting Neural Networks for NLP*, pages 56–62, Florence, Italy, August 2019. Association  
907 for Computational Linguistics. doi: 10.18653/v1/W19-4807. URL [https://www.aclweb](https://www.aclweb.org/anthology/W19-4807)  
908 [org/anthology/W19-4807](https://www.aclweb.org/anthology/W19-4807).
- 909 [114] Moganarangan Thayaparan, Marco Valentino, and André Freitas. A survey on explainability  
910 in machine reading comprehension. arXiv:2010.00389, 2020. URL [https://arxiv.org/](https://arxiv.org/abs/2010.00389)  
911 [abs/2010.00389](https://arxiv.org/abs/2010.00389).
- 912 [115] James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. FEVER: a  
913 large-scale dataset for fact extraction and VERification. In *Proceedings of the 2018 Conference*  
914 *of the North American Chapter of the Association for Computational Linguistics: Human*  
915 *Language Technologies, Volume 1 (Long Papers)*, pages 809–819, New Orleans, Louisiana,  
916 June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-1074. URL  
917 <https://www.aclweb.org/anthology/N18-1074>.
- 918 [116] Ilaria Tiddi, M. d’Aquin, and E. Motta. An ontology design pattern to define explanations.  
919 In *Proceedings of the 8th International Conference on Knowledge Capture*, 2015. URL  
920 <http://oro.open.ac.uk/44321/>.
- 921 [117] Masatoshi Tsuchiya. Performance impact caused by hidden bias of training data for recognizing  
922 textual entailment. In *Proceedings of the Eleventh International Conference on Language*  
923 *Resources and Evaluation (LREC 2018)*, Miyazaki, Japan, May 2018. European Language  
924 Resources Association (ELRA). URL <https://www.aclweb.org/anthology/L18-1239>.
- 925 [118] Sahil Verma, John P. Dickerson, and Keegan Hines. Counterfactual explanations for machine  
926 learning: A review. arXiv:2010.10596, 2020. URL <https://arxiv.org/abs/2010.10596>.
- 927 [119] Hoa Trong Vu, Claudio Greco, Aliia Erofeeva, Somayeh Jafaritazehjan, Guido Linders, Marc  
928 Tanti, Alberto Testoni, Raffaella Bernardi, and Albert Gatt. Grounded textual entailment. In  
929 *Proceedings of the 27th International Conference on Computational Linguistics*, pages 2354–  
930 2368, Santa Fe, New Mexico, USA, August 2018. Association for Computational Linguistics.  
931 URL <https://www.aclweb.org/anthology/C18-1199>.
- 932 [120] David Wadden, Shanchuan Lin, Kyle Lo, Lucy Lu Wang, Madeleine van Zuylen, Arman  
933 Cohan, and Hannaneh Hajishirzi. Fact or fiction: Verifying scientific claims. In *Proceedings*  
934 *of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*,  
935 pages 7534–7550, Online, November 2020. Association for Computational Linguistics. doi:  
936 10.18653/v1/2020.emnlp-main.609. URL [https://www.aclweb.org/anthology/2020.](https://www.aclweb.org/anthology/2020.emnlp-main.609)  
937 [emnlp-main.609](https://www.aclweb.org/anthology/2020.emnlp-main.609).
- 938 [121] Cunxiang Wang, Shuailong Liang, Yue Zhang, Xiaonan Li, and Tian Gao. Does it make sense?  
939 and why? a pilot study for sense making and explanation. In *Proceedings of the 57th Annual*  
940 *Meeting of the Association for Computational Linguistics*, pages 4020–4026, Florence, Italy,  
941 July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1393. URL  
942 <https://www.aclweb.org/anthology/P19-1393>.
- 943 [122] Ziqi Wang, Yujia Qin, Wenxuan Zhou, Jun Yan, Qinyuan Ye, Leonardo Neves, Zhiyuan Liu,  
944 and Xiang Ren. Learning from explanations with neural execution tree. In *Proceedings*  
945 *of the International Conference on Learning Representations (ICLR)*, 2020. URL <https://arxiv.org/pdf/1911.01352.pdf>.  
946
- 947 [123] Maximilian Wich, Hala Al Kuwatly, and Georg Groh. Investigating annotator bias with a  
948 graph-based approach. In *Proceedings of the Fourth Workshop on Online Abuse and Harms*,  
949 pages 191–199, Online, November 2020. Association for Computational Linguistics. doi: 10.  
950 18653/v1/2020.alw-1.22. URL <https://www.aclweb.org/anthology/2020.alw-1.22>.

- 951 [124] Sarah Wiegrefe and Yuval Pinter. Attention is not not explanation. In *Proceedings of*  
952 *the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th*  
953 *International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages  
954 11–20, Hong Kong, China, November 2019. Association for Computational Linguistics. doi:  
955 10.18653/v1/D19-1002. URL <https://www.aclweb.org/anthology/D19-1002>.
- 956 [125] Sarah Wiegrefe, Ana Marasović, and Noah A. Smith. Measuring association between labels  
957 and free-text rationales. arXiv:2010.12762, 2020. URL <https://arxiv.org/abs/2010>  
958 [12762](https://arxiv.org/abs/2010.12762).
- 959 [126] Adina Williams, Nikita Nangia, and Samuel Bowman. A broad-coverage challenge corpus  
960 for sentence understanding through inference. In *Proceedings of the 2018 Conference of*  
961 *the North American Chapter of the Association for Computational Linguistics: Human*  
962 *Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122, New Orleans, Louisiana,  
963 June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-1101. URL  
964 <https://www.aclweb.org/anthology/N18-1101>.
- 965 [127] Tongshuang Wu, Marco Túlio Ribeiro, J. Heer, and Daniel S. Weld. Polyjuice: Automated,  
966 general-purpose counterfactual generation. arXiv:2101.00288, 2021. URL <https://arxiv>  
967 [org/abs/2101.00288](https://arxiv.org/abs/2101.00288).
- 968 [128] Ning Xie, Farley Lai, Derek Doran, and Asim Kadav. Visual entailment: A novel task for  
969 fine-grained image understanding. *arXiv preprint arXiv:1901.06706*, 2019. URL [https:](https://arxiv)  
970 [/arxiv.org/pdf/1901.06706.pdf](https://arxiv.org/pdf/1901.06706.pdf).
- 971 [129] Zhengnan Xie, Sebastian Thiem, Jaycie Martin, Elizabeth Wainwright, Steven Marmorstein,  
972 and Peter Jansen. WorldTree v2: A corpus of science-domain structured explanations and  
973 inference patterns supporting multi-hop inference. In *Proceedings of the 12th Language*  
974 *Resources and Evaluation Conference*, pages 5456–5473, Marseille, France, May 2020.  
975 European Language Resources Association. ISBN 979-10-95546-34-4. URL <https://www>  
976 [aclweb.org/anthology/2020.lrec-1.671](https://www.aclweb.org/anthology/2020.lrec-1.671).
- 977 [130] Fan Yang, Mengnan Du, and X. Hu. Evaluating explanation without ground truth in inter-  
978 pretable machine learning. arXiv:1907.06831, 2019. URL <https://arxiv.org/abs/1907>  
979 [06831](https://arxiv.org/abs/1907.06831).
- 980 [131] Linyi Yang, Eoin Kenny, Tin Lok James Ng, Yi Yang, Barry Smyth, and Ruihai Dong. Gener-  
981 ating plausible counterfactual explanations for deep transformers in financial text classification.  
982 In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6150–  
983 6160, Barcelona, Spain (Online), December 2020. International Committee on Computational  
984 Linguistics. URL <https://www.aclweb.org/anthology/2020.coling-main.541>.
- 985 [132] Shaohua Yang, Qiaozhi Gao, Sari Sadiya, and Joyce Chai. Commonsense justification  
986 for action explanation. In *Proceedings of the 2018 Conference on Empirical Methods in*  
987 *Natural Language Processing*, pages 2627–2637, Brussels, Belgium, October-November  
988 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1283. URL  
989 <https://www.aclweb.org/anthology/D18-1283>.
- 990 [133] Yi Yang, Wen-tau Yih, and Christopher Meek. WikiQA: A challenge dataset for open-domain  
991 question answering. In *Proceedings of the 2015 Conference on Empirical Methods in Natural*  
992 *Language Processing*, pages 2013–2018, Lisbon, Portugal, September 2015. Association for  
993 Computational Linguistics. doi: 10.18653/v1/D15-1237. URL <https://www.aclweb.org/>  
994 [anthology/D15-1237](https://www.aclweb.org/anthology/D15-1237).
- 995 [134] Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov,  
996 and Christopher D. Manning. HotpotQA: A dataset for diverse, explainable multi-hop question  
997 answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language*  
998 *Processing*, pages 2369–2380, Brussels, Belgium, October-November 2018. Association for  
999 Computational Linguistics. doi: 10.18653/v1/D18-1259. URL <https://www.aclweb.org/>  
1000 [anthology/D18-1259](https://www.aclweb.org/anthology/D18-1259).
- 1001 [135] Qinyuan Ye, Xiao Huang, Elizabeth Boschee, and Xiang Ren. Teaching machine com-  
1002 prehension with compositional explanations. In *Findings of the Association for Compu-*  
1003 *tational Linguistics: EMNLP 2020*, pages 1599–1615, Online, November 2020. Association  
1004 for Computational Linguistics. doi: 10.18653/v1/2020.findings-emnlp.145. URL  
1005 <https://www.aclweb.org/anthology/2020.findings-emnlp.145>.

- 1006 [136] Kayo Yin, Patrick Fernandes, Danish Pruthi, Aditi Chaudhary, André F. T. Martins, and  
1007 Graham Neubig. Do context-aware translation models pay the right attention? In *Proceedings*  
1008 *of the 59th Annual Meeting of the Association for Computational Linguistics*, 2021. URL  
1009 <https://arxiv.org/abs/2105.06977>.
- 1010 [137] Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions  
1011 to visual denotations: New similarity metrics for semantic inference over event descriptions.  
1012 *Transactions of the Association for Computational Linguistics*, 2:67–78, 2014. doi: 10.1162/  
1013 tacl\_a\_00166. URL <https://www.aclweb.org/anthology/Q14-1006>.
- 1014 [138] Mo Yu, Shiyu Chang, Yang Zhang, and Tommi Jaakkola. Rethinking cooperative ratio-  
1015 nalization: Introspective extraction and complement control. In *Proceedings of the 2019*  
1016 *Conference on Empirical Methods in Natural Language Processing and the 9th Interna-*  
1017 *tional Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4094–  
1018 4103, Hong Kong, China, November 2019. Association for Computational Linguistics. doi:  
1019 10.18653/v1/D19-1420. URL <https://www.aclweb.org/anthology/D19-1420>.
- 1020 [139] Omar Zaidan, Jason Eisner, and Christine Piatko. Using “annotator rationales” to improve  
1021 machine learning for text categorization. In *Human Language Technologies 2007: The*  
1022 *Conference of the North American Chapter of the Association for Computational Linguistics;*  
1023 *Proceedings of the Main Conference*, pages 260–267, Rochester, New York, April 2007.  
1024 Association for Computational Linguistics. URL [https://www.aclweb.org/anthology/](https://www.aclweb.org/anthology/N07-1033)  
1025 [N07-1033](https://www.aclweb.org/anthology/N07-1033).
- 1026 [140] Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. From recognition to cognition:  
1027 Visual commonsense reasoning. *2019 IEEE/CVF Conference on Computer Vision and Pattern*  
1028 *Recognition (CVPR)*, 2019. URL <https://ieeexplore.ieee.org/document/8953217>.
- 1029 [141] Hongming Zhang, Xinran Zhao, and Yangqiu Song. WinoWhy: A deep diagnosis of essential  
1030 commonsense knowledge for answering Winograd schema challenge. In *Proceedings of*  
1031 *the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5736–  
1032 5745, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.  
1033 acl-main.508. URL <https://www.aclweb.org/anthology/2020.acl-main.508>.
- 1034 [142] Ye Zhang, Iain Marshall, and Byron C. Wallace. Rationale-augmented convolutional neural  
1035 networks for text classification. In *Proceedings of the 2016 Conference on Empirical Methods*  
1036 *in Natural Language Processing*, pages 795–804, Austin, Texas, November 2016. Association  
1037 for Computational Linguistics. doi: 10.18653/v1/D16-1076. URL [https://www.aclweb](https://www.aclweb.org/anthology/D16-1076)  
1038 [.org/anthology/D16-1076](https://www.aclweb.org/anthology/D16-1076).
- 1039 [143] Yuhao Zhang, Victor Zhong, Danqi Chen, Gabor Angeli, and Christopher D. Manning.  
1040 Position-aware attention and supervised data improve slot filling. In *Proceedings of the*  
1041 *2017 Conference on Empirical Methods in Natural Language Processing*, pages 35–45,  
1042 Copenhagen, Denmark, September 2017. Association for Computational Linguistics. doi:  
1043 10.18653/v1/D17-1004. URL <https://www.aclweb.org/anthology/D17-1004>.