On the Diversity and Limits of Human Explanations

Chenhao Tan

Department of Computer Science & Harris School of Public Policy
The University of Chicago
chenhao@uchicago.edu

Abstract

A growing effort in NLP aims to build datasets of human explanations. However, it remains unclear whether these datasets serve their intended goals. This problem is exacerbated by the fact that the term *explanation* is *overloaded* and refers to a broad range of notions with different properties and ramifications. Our goal is to provide an overview of the diversity of explanations, discuss human limitations in providing explanations, and ultimately provide implications for collecting and using human explanations in NLP.

Inspired by prior work in psychology and cognitive sciences, we group existing human explanations in NLP into three categories: proximal mechanism, evidence, and procedure. These three types differ in nature and have implications for the resultant explanations. For instance, procedure is not considered explanation in psychology and connects with a rich body of work on learning from instructions. The diversity of explanations is further evidenced by proxy questions that are needed for annotators to interpret and answer "why is [input] assigned [label]". Finally, giving explanations may require different, often deeper, understandings than predictions, which casts doubt on whether humans can provide valid explanations in some tasks.

1 Introduction

With the growing interest in explainable NLP systems, the NLP community have become increasingly interested in building datasets of human explanations. These human explanations can *ideally* capture human reasoning of why a (correct) label is chosen. If this is indeed the case, they are hypothesized to aid models with additional supervision, train models that explain their own predictions, and evaluate machine-generated explanations (Wiegreffe and Marasović, 2021). In fact, DeYoung et al. (2020) already developed a leaderboard, where the

implicit assumption is that humans can provide *valid* explanations and these explanations can in turn be *uniformly* considered as groundtruths.

However, are these assumptions satisfied and can human explanations serve these goals? In this work, we aim to introduce prior relevant literature in psychology to the NLP community and argue against abusing the term explanations and prematurely assuming that human explanations provide valid reasoning for inferring a label.

First, we point out the rich diversity in what the NLP community refer to as explanations and how researchers collect them. The term "explanation" is overloaded in the NLP and AI community: it often refers to many distinct concepts and outcomes.¹ For example, procedural instructions are different from explanations that attempt to convey proximal causal mechanisms. The diversity of explanations is further evidenced by the variety of proxy questions that researchers ask to collect explanations, e.g., "highlight the important words that would tell someone to see the movie" vs. "highlight ALL words that reflect a sentiment". These proxy questions are necessary because the question of "why is [input] assigned [label]" is too open-ended. It follows that these "human explanations" are supposed to answer different questions in the first place and may not all be used for the same goals, e.g., serving as groundtruth labels.

In addition to the diversity, we highlight two insights from psychology on whether humans can provide valid explanations: 1) prediction does not entail explanation (Wilson and Keil, 1998), i.e., although humans may be able to provide valid labels, they may not be able to provide explanations that capture the reasoning process needed or used to infer a label; 2) everyday explanations are necessarily incomplete (Keil, 2006; Lombrozo, 2006), because they seldom capture the complete deduc-

¹Broniatowski et al. (2021) further argues that interpretation and explainability are distinct concepts.

tive processes from a set of axioms to a statement.

In summary, *not* all explanations are equal and humans may *not* always be able to provide valid explanations. We encourage the NLP community to embrace the complex and intriguing phenomena behind human explanations instead of simply viewing explanations as another set of uniform labels. A better understanding and characterization of human explanations will inform how to collect and use human explanations in NLP.

2 Types of Human Explanations in NLP

To understand whether datasets of human explanations can serve their intended goals, we first connect current human explanations in NLP with existing psychology literature to examine the use of the term "explanation" in NLP. We adapt the categorization in Lombrozo (2006) and group human explanations in NLP into the following three categories based on the conveyed information:

- **Proximal mechanisms**. This type of explanation attempts to provide the mechanism behind the predicted label, i.e., how to infer the label from the text, and match efficient cause in Lombrozo (2006). We created E1 in Table 1 to illustrate this type of explanation. Note that E1 does not provide the complete mechanism. For instance, it does not define "year" or "temporal modifier", or make clear that "1997" is a "year". Neither does it cover the axioms of logic. This is a common property of human explanations: they are known to cover partial/proximal mechanisms rather than the complete deduction from natural laws and empirical conditions (Hempel and Oppenheim, 1948).
- Evidence. This type of explanation includes the relevant tokens in the input (e.g., E2 in Table 1) and directly maps to highlights in Wiegreffe and Marasović (2021). However, it does not map to any existing definitions of explanations in the psychology literature since the evidence does not provide any information on *how* evidence leads to the label. In other words, evidence alone does not *explain*.
- **Procedure**. Unlike proximal mechanisms, this type of explanation provides step-by-step rules or procedures that one can directly follow, e.g., E3 in Table 1. They are more ex-

plicit and unambiguous than proximal mechanisms. In fact, one can write a rule based on E3 to find marriage relation, but one cannot easily do that with E1. Furthermore, the procedures are grounded to the input, so it is related to *formal cause*, "the form or properties that make something what it is" (Lombrozo, 2006), which is definitional and does not convey the underlying mechanisms. Procedural instructions are only possible for some tasks, while proximal mechanisms are the most common form of everyday explanations.

These three categories empirically capture all the explanations discussed in NLP literature. Lombrozo (2006) also discuss two other categories, *final causes* (the goal) and *material causes* (the constituting substance). For instance, a final cause to "why [input] is assigned [label]" can be that "this label is provided to train a classifier". These two categories have been less relevant for NLP.

Implications. This categorization allows us to think about what kind of explanations are desired for NLP systems and help clarify how to use them appropriately. First, proximal mechanisms are best aligned with human intuitions of explanations, especially in terms of hinting at causal mechanisms. However, they can be difficult to collect for NLP tasks. For example, Table 2 shows example explanations in E-SNLI that fail to convey any proximal mechanisms: they either repeat the hypothesis or express invalid mechanisms ("the bike competition" does not entail "bikes on a stone road"). See further discussions on the challenges in collecting explanations in §4. Furthermore, they may be difficult to use for supervising or evaluating a model.

Second, evidence by definition provides little information about the mechanisms behind a label, but it can be potentially useful as additional supervision or groundtruths. We will further elaborate on the nature of evidence in different tasks in §3. However, it may be useful to the community to use clear terminology (e.g., evidence or rationale (Lei et al., 2020; Carton et al., 2020)) to avoid lumping everything into "explanation".

Finally, procedures are essentially instructions, and Keil (2006) explicitly distinguishes explanations from simple procedural knowledge: "Knowing how to operate an automated teller machine or make an international phone call might not entail having any understanding of how either system

Instance	Label	Explanation
----------	-------	-------------

Task: whether the query is supported or refuted by the preceding texts (Thorne et al., 2018)

S1: No Way Out is the debut studio album by American hip hop recording artist, songwriter and record producer Puff Daddy. S2: It was released on July 1, 1997, by Supports his Bad Boy record label. The label's official crediting as "The Family"...... Query: 1997 was the year No Way Out was released.

Proximal mechanism (E1): "It" in S2 refers to "No Way Out" and "on July 1, 1997" is the temporal modifier of "release", we can thus infer 1997 was the year that No Way Out was released.

Evidence (E2): S1, S2

Task: whether person 1 is married to person 2 (Hancock et al., 2018)

True

Tom Brady and his wife Gisele Bündchen were spotted in New York City on Monday amid rumors of Brady's alleged role in Deflategate **Procedure (E3):** The words "and his wife" are between person 1 and person 2.

Table 1: Types of human explanations and corresponding examples. "S1:" and "S2:" were added to facilitate writing the explanation, which also shows the non-triviality of writing explanations.

P: Men in green hats appear to be attending a gay pride festival. H: Men are attending a festival. E: The men are attending the festival.

P: Several bikers on a stone road, with spectators watching. H: The bikers are on a stone road. E: That there are spectators watching the bikers on a stone road implies there is a bike competition.

Table 2: Examples from E-SNLI (Camburu et al., 2018). P: premise, H:: hypothesis, E: explanation.

works". Another reason to clarify the procedure category is that it would be useful to engage with a rich body of work on learning from instructions when human explanations are procedural (Goldwasser and Roth, 2014; Matuszek et al., 2013).

We would like to emphasize that procedures or instructions are powerful and can potentially benefit many NLP problems (e.g., relation extraction). At the same time, it is useful to point out that procedures are different from proximal mechanisms.

3 Proxy Questions Used to Collect Human Explanations

Although explanations are supposed to answer "why is [input] assigned [label]" (Wiegreffe and Marasović, 2021), this literal form is too openended and may not induce "useful" human explanations. As a result, proxy questions are often necessary for collecting human explanations. These proxy questions further demonstrate the diversity of human explanations beyond the types of explanation. Here we discuss these proxy questions for collecting evidence. See the appendix for discussions on proximal mechanisms and procedures.

To collect evidence (highlights), researchers adopt diverse questions for relatively simple single-text classification tasks (see Table 3). Consider

the seemingly straightforward case of sentiment analysis, "why is the sentiment of a review positive/negative". A review can present both positive and negative sentiments (Aithal and Tan, 2021), so the label often comes from one sentiment outweighing the other. However, in practice, researchers often ask annotators to identify **only** words supporting the label. Critical wording differences remain in their questions: Zaidan et al. (2007) ask for *the most important* words and phrases that would *tell someone to see the movie*, while Sen et al. (2020) requires *all* words *reflecting the sentiment*. Two key differences arise: 1) "the most important" vs. "all"; 2) "telling someone to see the movie" vs. "reflecting the sentiment".

In contrast, personal attack detection poses a task where the negative class ("no personal attack") by definition points to the lack of evidence in the text. It follows that the questions that researchers can ask almost exclusively apply to the positive class (i.e., "highlight sections of comments that they considered to constitute personal attacks").

In comparison, researchers approach evidence more uniformly for document-query classification tasks. They generally use similar proxy questions (e.g., Thorne et al. (2018) and Hanselowski et al. (2019) ask almost the same questions) and ask people to select sentences instead of words. That said, intriguing differences still exist: 1) Lehman et al. (2019) simply ask annotators to provide accompanying rationales; 2) Thorne et al. (2018) aim for "strong" reasons, which likely induces different interpretations among annotators; 3) Khashabi et al. (2018) collect questions, answer, and sentence indices at the same time, among which sentence indices can be used to find the corresponding sentences as evidence. It remains unclear how these

Reference	Task	Questions		
Zaidan et al. (2007)	Sentiment analysis	Evidence: single-text classification To justify why a review is positive, highlight the most important words and phrases that would tell someone to see the movie. To justify why a review is negative, highlight words and phrases that would tell someone not to see the movie.		
Sen et al. (2020) Carton et al. (2018)	•	Label the sentiment and highlight ALL words that reflect this sentiment. Highlight sections of comments that they considered to constitute personal attacks.		
	Evidence: document-query classification			
Lehman et al. (2019)		Generators were also asked to provide answers and accompanying rationales to the prompts that they provided.		
Thorne et al. (2018)		If I was given only the selected sentences, do I have strong reason to believe the claim is true (supported) or stronger reason to believe the claim is false (refuted).		
Khashabi et al. (2018)	Question answering	Ask them (participants) for a correct answer and for the sentence indices required to answer the question.		

Table 3: Questions that prior work uses to collect human explanations. We include the short version of the guidelines here for space reasons. Refer to the appendix for the full text of the relevant annotation guidelines.

differences in annotation processes and question phrasings affect the collected human explanations.

Implications. Our observation on proxy questions aligns with dataset-specific designs discussed in Wiegreffe and Marasović (2021). We emphasize that these different forms of questions entail different properties of the collected human explanations, as evidenced by Carton et al. (2020). For example, the lack of evidence in the negative class in personal attack classification likely requires special strategies in using human explanations to train a model and evaluate machine rationales. Sentencelevel and token-level annotations also lead to substantially different outcomes, at least in the forms of explanations. We believe that it is important for the NLP community to investigate the effect of proxy questions and use the collected explanations with care, rather than lumping all datasets under the umbrella of explanations.

We also recommend all researchers to provide detailed annotation guidelines used to collect human explanations. As the area of collecting human explanation is nascent, the goal is not to promote consistent and uniform annotation guidelines but to encourage the community to pay attention to the different underlying questions and characterize the resultant diverse properties of human explanations.

4 Can Humans Provide Explanations?

In order for human explanations to serve as additional supervision in training models and evaluate machine-generated explanations, human explanations need to provide valid mechanisms for

a correct label. Finally, we discuss challenges for humans to provide explanations of such qualities.

Conceptual framework. We situate our discussion in the psychological framework provided by Wilson and Keil (1998) to highlight what may be required to explain. Wilson and Keil (1998) examines where explanation falls in three central notions: prediction, understanding, and theories. They argue that these three notions "form a progression of increasing sophistication and depth with explanations falling between understanding and theories". For instance, we may be able to predict that a car will start when we turn the ignition switch, but few of us are able to explain in detail why this is so. In contrast, if a person is able to explain in detail why a car starts when you turn on the ignition switch, they can likely predict what will happen if various parts of the engine are damaged or removed.

These three central notions are also essential in machine learning. Traditional label annotation is concerned with prediction, however, being able to predict does not entail being able to explain.

Emulation vs. discovery. Next, we gradually unfold the practical challenges in collecting valid explanations from humans. The first challenge lies in whether humans can predict, i.e., assign the correct label. We highlight two types of tasks for AI: emulation vs. discovery (Lai et al., 2020). In *emulation* tasks, models are trained to emulate human intelligence and labels are often crowdsourced. Labels, however, can also derive from external (social/biological) processes, e.g., the popularity of a tweet and the effect of a medical treatment. Mod-

els can thus discover patterns that humans may not recognize in these *discovery* tasks. While most NLP tasks such as NLI and QA are emulation tasks, many NLP problems, especially when concerned with social interaction, are discovery tasks, ranging from identifying memorable movie quotes to predicting the popularity of messages (Danescu-Niculescu-Mizil et al., 2012; Tan et al., 2014).

Aligning with our discussion on explanation and prediction, most datasets of human explanations in NLP assume that humans are able to predict and are on emulation tasks. However, we note exceptions such as explanations of actions in gaming (Ehsan et al., 2019), where humans may often choose suboptimal actions (labels).

Cognitive challenges in providing valid explana-

tions. Even conditioned on that humans can predict the label, humans may not be able to provide valid explanations for at least two reasons. First, as Wilson and Keil (1998) suggests, explanation requires more depth than prediction. For instance, we may possess some notions of common sense (e.g., one should not slap a stranger), but it is unclear whether we can explain common sense in detail (e.g., why one should not slap a stranger through theory of morality), similar to the car ignition example. One may argue that theory of morality may not be what NLP researchers seek, but it is critical to consider the desiderata of human explanations, with the limits in mind.

Second, explanation often requires people to report their *subjective* mental processes, i.e., how our minds arrive at a particular judgement, rather than following *objective* consensual guidelines such as annotating logical entailment. However, classic work by Nisbett and Wilson (1977) suggests that our verbal reports on our mental processes can be highly inaccurate. For instance, in admission decisions, legitimate information can be used to justify preferences based on illegitimate factors such as race (Norton et al., 2006). Many studies on implicit bias also reinforces that we are not aware of our biases and thus cannot include them (i.e., the actual reasoning in our mind) in our explanations (Greenwald et al., 1998).

Explanations are necessarily incomplete. Finally, there are indeed cases where we believe that humans can provide valid mechanisms. For instance, some question answering tasks boil down to logical inference from evidence to query. In these

cases, NLP researchers need to recognize that human explanations are necessarily incomplete: people do not start from a set of axioms and present all the deductive steps (Keil, 2006; Lombrozo, 2006). Therefore, even for simple tasks such as natural language inference, we may simply give explanations such as repeating the hypothesis without presenting any axiom or deduction required to infer the label.

Implications. We cannot assume that humans are capable of providing explanations that contain valuable proximal mechanisms. The very fact that humans can still provide explanations for incorrect labels and tasks where they do not perform well suggests that one should be skeptical about whether human explanations can be used to train models as additional supervision or evaluate machine-generated explanations as groundtruths.

Note that incomplete explanations can still be very useful for NLP. We believe that recognizing and characterizing this incompleteness (e.g., which proximal mechanism is more salient to humans) is critical for understanding and leveraging human explanations for the intended goals in NLP. To summarize, we argue that human explanations are necessarily incomplete and it is important to understand and characterize this incompleteness, which can inform how we can leverage it for the intended goals in NLP.

5 Conclusion

Explanations represent a fascinating phenomenon and are actively studied in psychology, cognitive science, and other social sciences. While the growing interest in explanations from the NLP community is exciting, we encourage the community to view this as an opportunity to understand how humans approach explanations and contribute to understanding and exploring the explanation processes. This will in turn inform how to collect and use human explanations in NLP. A modest proposal is that it is useful to examine and characterize human explanations before assuming that all explanations are equal and chasing a leaderboard.

Acknowledgments

We thank anonymous reviewers for their feedback, and members of the Chicago Human+AI Lab for their insightful suggestions. This work is supported in part by research awards from Amazon, IBM, Salesforce, and NSF IIS-2040989, 2125116, 2126602.

References

- Madhusudhan Aithal and Chenhao Tan. 2021. On positivity bias in negative reviews. In *Proceedings of ACL*.
- Tariq Alhindi, Savvas Petridis, and Smaranda Muresan. 2018. Where is your evidence: Improving fact-checking by justification modeling. In *Proceedings of the First Workshop on Fact Extraction and VERification (FEVER)*, pages 85–90, Brussels, Belgium. Association for Computational Linguistics.
- Faeze Brahman, Vered Shwartz, Rachel Rudinger, and Yejin Choi. 2020. Learning to rationalize for nonmonotonic reasoning with distant supervision. *arXiv preprint arXiv:2012.08012*.
- David A Broniatowski et al. 2021. Psychological foundations of explainability and interpretability in artificial intelligence.
- Oana-Maria Camburu, Tim Rocktäschel, Thomas Lukasiewicz, and Phil Blunsom. 2018. e-snli: Natural language inference with natural language explanations. In *Proceedings of NeurIPS*.
- Samuel Carton, Qiaozhu Mei, and Paul Resnick. 2018. Extractive adversarial networks: High-recall explanations for identifying personal attacks in social media posts. In *Proceedings of EMNLP*.
- Samuel Carton, Anirudh Rathore, and Chenhao Tan. 2020. Evaluating and characterizing human rationales. In *Proceedings of EMNLP*.
- Jeff Da, Maxwell Forbes, Rowan Zellers, Anthony Zheng, Jena D Hwang, Antoine Bosselut, and Yejin Choi. 2020. Edited media understanding: Reasoning about implications of manipulated images. arXiv preprint arXiv:2012.04726.
- Cristian Danescu-Niculescu-Mizil, Justin Cheng, Jon Kleinberg, and Lillian Lee. 2012. You had me at hello: How phrasing affects memorability. In *Proceedings of the ACL*.
- Jay DeYoung, Sarthak Jain, Nazneen Fatema Rajani, Eric Lehman, Caiming Xiong, Richard Socher, and Byron C Wallace. 2020. Eraser: A benchmark to evaluate rationalized nlp models. In *Proceedings of ACL*.
- Virginie Do, Oana-Maria Camburu, Zeynep Akata, and Thomas Lukasiewicz. 2020. e-snli-ve-2.0: Corrected visual-textual entailment with natural language explanations. In *IEEE CVPR Workshop on Fair, Data Efficient and Trusted Computer Vision*.
- Upol Ehsan, Pradyumna Tambwekar, Larry Chan, Brent Harrison, and Mark O Riedl. 2019. Automated rationale generation: a technique for explainable ai and its effects on human perceptions. In *Proceedings of the 24th International Conference on Intelligent User Interfaces*, pages 263–274.

- Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies. *arXiv preprint arXiv:2101.02235*.
- Dan Goldwasser and Dan Roth. 2014. Learning from natural instructions. *Machine learning*, 94(2):205–232.
- Anthony G Greenwald, Debbie E McGhee, and Jordan LK Schwartz. 1998. Measuring individual differences in implicit cognition: the implicit association test. *Journal of personality and social psychology*, 74(6):1464.
- Braden Hancock, Paroma Varma, Stephanie Wang, Martin Bringmann, Percy Liang, and Christopher Ré. 2018. Training classifiers with natural language explanations. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1884–1895, Melbourne, Australia. Association for Computational Linguistics.
- Andreas Hanselowski, Christian Stab, Claudia Schulz, Zile Li, and Iryna Gurevych. 2019. A richly annotated corpus for different tasks in automated fact-checking. *arXiv* preprint arXiv:1911.01214.
- Carl G Hempel and Paul Oppenheim. 1948. Studies in the logic of explanation. *Philosophy of science*, 15(2):135–175.
- Naoya Inoue, Pontus Stenetorp, and Kentaro Inui. 2020. R4C: A benchmark for evaluating RC systems to get the right answer for the right reason. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6740–6750, Online. Association for Computational Linguistics.
- Peter Jansen, Niranjan Balasubramanian, Mihai Surdeanu, and Peter Clark. 2016. What's in an explanation? characterizing knowledge and inference requirements for elementary science exams. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 2956–2965.
- Peter Jansen, Elizabeth Wainwright, Steven Marmorstein, and Clayton Morrison. 2018. WorldTree: A corpus of explanation graphs for elementary science questions supporting multi-hop inference. In Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018), Miyazaki, Japan. European Language Resources Association (ELRA).
- Harsh Jhamtani and Peter Clark. 2020. Learning to explain: Datasets and models for identifying valid reasoning chains in multihop question-answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 137–150, Online. Association for Computational Linguistics.

- Frank C Keil. 2006. Explanation and understanding. *Annu. Rev. Psychol.*, 57:227–254.
- Daniel Khashabi, Snigdha Chaturvedi, Michael Roth, Shyam Upadhyay, and Dan Roth. 2018. Looking beyond the surface: A challenge set for reading comprehension over multiple sentences. In *Proceedings of NAACL*.
- Tushar Khot, Peter Clark, Michal Guerquin, Peter Jansen, and Ashish Sabharwal. 2020. Qasc: A dataset for question answering via sentence composition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 8082–8090.
- Jinkyu Kim, Anna Rohrbach, Trevor Darrell, John Canny, and Zeynep Akata. 2018. Textual explanations for self-driving vehicles. In *Proceedings of the European conference on computer vision (ECCV)*, pages 563–578.
- Neema Kotonya and Francesca Toni. 2020. Explainable automated fact-checking for public health claims. In *EMNLP*.
- Mucahid Kutlu, Tyler McDonnell, Matthew Lease, and Tamer Elsayed. 2020. Annotator rationales for labeling tasks in crowdsourcing. *Journal of Artificial Intelligence Research*, 69:143–189.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453–466.
- Vivian Lai, Han Liu, and Chenhao Tan. 2020. "why is 'chicago' deceptive?" towards building model-driven tutorials for humans. In *Proceedings of CHI*.
- Matthew Lamm, Jennimaria Palomaki, Chris Alberti, Daniel Andor, Eunsol Choi, Livio Baldini Soares, and Michael Collins. 2020. Qed: A framework and dataset for explanations in question answering. arXiv preprint arXiv:2009.06354.
- Eric Lehman, Jay DeYoung, Regina Barzilay, and Byron C Wallace. 2019. Inferring which medical treatments work from reports of clinical trials. In *NAACL*.
- Jie Lei, Licheng Yu, Tamara Berg, and Mohit Bansal. 2020. What is more likely to happen next? video-and-language future event prediction. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8769–8784, Online. Association for Computational Linguistics.
- Qing Li, Qingyi Tao, Shafiq Joty, Jianfei Cai, and Jiebo Luo. 2018. VQA-E: Explaining, elaborating, and enhancing your answers for visual questions. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 552–567.

- Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. 2017. Program induction by rationale generation: Learning to solve and explain algebraic word problems. *arXiv* preprint arXiv:1705.04146.
- Tania Lombrozo. 2006. The structure and function of explanations. *Trends in cognitive sciences*, 10(10):464–470.
- Cynthia Matuszek, Evan Herbst, Luke Zettlemoyer, and Dieter Fox. 2013. Learning to parse natural language commands to a robot control system. In *Experimental robotics*, pages 403–415. Springer.
- Richard E Nisbett and Timothy D Wilson. 1977. Telling more than we can know: verbal reports on mental processes. *Psychological review*, 84(3):231.
- Michael I Norton, Samuel R Sommers, Joseph A Vandello, and John M Darley. 2006. Mixed motives and racial bias: The impact of legitimate and illegitimate criteria on decision making. *Psychology, Public Policy, and Law*, 12(1):36.
- Dong Huk Park, Lisa Anne Hendricks, Zeynep Akata, Anna Rohrbach, Bernt Schiele, Trevor Darrell, and Marcus Rohrbach. 2018. Multimodal explanations: Justifying decisions and pointing to the evidence. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8779–8788.
- Nazneen Fatema Rajani, Bryan McCann, Caiming Xiong, and Richard Socher. 2019. Explain yourself! leveraging language models for commonsense reasoning.
- Nazneen Fatema Rajani, Rui Zhang, Yi Chern Tan, Stephan Zheng, Jeremy Weiss, Aadit Vyas, Abhijit Gupta, Caiming Xiong, Richard Socher, and Dragomir Radev. 2020. ESPRIT: Explaining solutions to physical reasoning tasks. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7906–7917, Online. Association for Computational Linguistics.
- Maarten Sap, Saadia Gabriel, Lianhui Qin, Dan Jurafsky, Noah A Smith, and Yejin Choi. 2020. Social bias frames: Reasoning about social and power implications of language. In *ACL*.
- Cansu Sen, Thomas Hartvigsen, Biao Yin, Xiangnan Kong, and Elke Rundensteiner. 2020. Human attention maps for text classification: Do humans and neural networks focus on the same words? In *Proceedings of ACL*, pages 4596–4608.
- Shashank Srivastava, Igor Labutov, and Tom Mitchell. 2017. Joint concept learning and semantic parsing from natural language explanations. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1527–1536, Copenhagen, Denmark. Association for Computational Linguistics.

Chenhao Tan, Lillian Lee, and Bo Pang. 2014. The effect of wording on message propagation: Topicand author-controlled natural experiments on twitter. In *Proceedings of ACL*.

James Thorne, Andreas Vlachos, Christos Arpit Mittal. Christodoulopoulos, and 2018 FEVER: a large-scale dataset for fact extraction and VERification. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 809-819, New Orleans, Louisiana. Association for Computational Linguistics.

David Wadden, Shanchuan Lin, Kyle Lo, Lucy Lu Wang, Madeleine van Zuylen, Arman Cohan, and Hannaneh Hajishirzi. 2020. Fact or fiction: verifying scientific claims. *arXiv preprint arXiv:2004.14974*.

Ziqi Wang, Yujia Qin, Wenxuan Zhou, Jun Yan, Qinyuan Ye, Leonardo Neves, Zhiyuan Liu, and Xiang Ren. 2020. Learning from explanations with neural execution tree.

Sarah Wiegreffe and Ana Marasović. 2021. Teach me to explain: A review of datasets for explainable nlp.

Robert A Wilson and Frank Keil. 1998. The shadows and shallows of explanation. *Minds and machines*, 8(1):137–159.

Zhengnan Xie, Sebastian Thiem, Jaycie Martin, Elizabeth Wainwright, Steven Marmorstein, and Peter Jansen. 2020. WorldTree v2: A corpus of science-domain structured explanations and inference patterns supporting multi-hop inference. In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 5456–5473, Marseille, France. European Language Resources Association.

Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov, and Christopher D Manning. 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600*.

Qinyuan Ye, Xiao Huang, Elizabeth Boschee, and Xiang Ren. 2020. Teaching machine comprehension with compositional explanations. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1599–1615, Online. Association for Computational Linguistics.

Omar Zaidan, Jason Eisner, and Christine Piatko. 2007. Using "annotator rationales" to improve machine learning for text categorization. In *Human language technologies 2007: The conference of the North American chapter of the association for computational linguistics; proceedings of the main conference*, pages 260–267.

Hongming Zhang, Xinran Zhao, and Yangqiu Song. 2020. Winowhy: A deep diagnosis of essential commonsense knowledge for answering winograd schema challenge. In *ACL*.

A Proxy Questions for Proximal Mechanisms and Procedure

Proximal mechanisms. In collecting proximal mechanisms, studies are more likely to ask explicitly versions of "why is [input] assigned [label]", compared to the case of evidence. However, they often need to provide structured guidelines. For example, Camburu et al. (2018) and Rajani et al. (2019) discussed the need to enforce word overlap as a way to improve the quality of human rationales. The specific requirements are quite different (see Table 6 and Table 7). There are also specific formulations of explanations, e.g., "What aspect/stereotype/characteristic of this group (often un-fairly assumed) is referenced or implied by this post?" in Sap et al. (2020). Finally, it is common that we cannot infer the exact questions asked (8/18 papers that collect explanations in free text).

Procedures. We cannot identify the exact questions in three of five papers for explicitly step-by-step procedures, which reflects the importance of reporting detailed annotation guidelines. As researchers collect step-by-step guidelines, Ye et al. (2020) and Geva et al. (2021) adopt very different decomposition for their problems (see Table 12).

B Detailed Proxy Questions

Table 4-12 show the instructions we find in prior work that detail the proxy questions. Camburu et al. (2018) and Rajani et al. (2019) collect both evidence and proximal mechanism. We include them in the tables for proximal mechanisms. Also, for question answering tasks, the difference between procedure and proximal mechanism can be subtle. We consider the collected explanations *procedure* if they aim to explicitly provide step-by-step guides directly grounded in the input.

Reference	Task	Questions and guidelines
Zaidan et al. (2007)	sentiment analysis	Each review was intended to give either a positive or a negative overall recommendation. You will be asked to justify why a review is positive or negative. To justify why a review is positive, highlight the most important words and phrases that would tell someone to see the movie. To justify why a review is negative, highlight words and phrases that would tell someone not to see the movie. These words and phrases are called rationales. You can highlight the rationales as you notice them, which should result in several rationales per review. Do your best to mark enough rationales to provide convincing support for the class of interest. You do not need to go out of your way to mark everything. You are probably doing too much work if you find yourself go- ing back to a paragraph to look for even more rationales in it. Furthermore, it is perfectly acceptable to skim through sections that you feel would not contain many
Sen et al. (2020)	sentiment anal-	rationales, such as a re- viewer's plot summary, even if that might cause you to miss a rationale here and there. 1. Read the review and decide the sentiment of this review
	ysis	(positive or negative). Mark your selection.2. Highlight ALL words that reflect this sentiment. Click on a word to highlight it. Click again to undo.3. If multiple words refect this sentiment, please highlight them all.
Carton et al. (2018)	Personal attack detection	40 undergraduate students used Brat (Stenetorp et al., 2012) to highlight sections of comments that they considered to constitute personal attacks.
Lehman et al. (2019)	Question answering	Prompt Generation: Question answering & Prompt creators were instructed to identify a snippet, in a given full-text article, that reports a relationship between an intervention, comparator, and outcome. Generators were also asked to provide answers and accompanying rationales to the prompts that they provided; such supporting evidence is important for this task and domain. The annotator was also asked to mark a snippet of text supporting their response. Annotators also had the option to mark prompts as invalid, e.g., if the prompt did not seem answerable on the basis of the article.
Thorne et al. (2018)	fact verification (QA)	If I was given only the selected sentences, do I have strong reason to believe the claim is true (supported) or stronger reason to believe the claim is false (refuted). If I'm not certain, what additional information (dictionary) do I have to add to reach this conclusion. In the annotation interface, all sentences from the introductory section of the page for the main entity of the claim and of every linked entity in those sentences were provided as a default source of evidence (left-hand side in Fig. 2). We did not set a hard time limit for the task, but the annotators were advised not to spend more than 2-3 minutes per claim.

Reference	Task	Questions and guidelines
Khashabi et al. (2018)	Question answering	We show each paragraph to 5 turkers and ask them to write 3-5 questions such that: (1) the question is answerable from the pas- sage, and (2) only those questions are allowed whose answer cannot be determined from a single sentence. We clarify this point by providing example paragraphs and questions. In order to encourage turkers to write meaningful questions that fit our criteria, we additionally ask them for a correct answer and for the sentence indices required to answer the question.
Yang et al. (2018)	Question answering	Workers provide the supporting facts (cannot infer the exact question)
Hanselowski et al. (2019)	tion	Mechanical Turk to annotate whether an ETS (evidence text snippets) agrees with the claim, refutes it, or has no stance towards the claim. An ETS was only considered to express a stance if it explicitly referred to the claim and either expressed support for it or refuted it. In all other cases, the ETS was considered as having no stance. FGE annotation. We filtered out ETSs with no stance, as they do not contain supporting or refuting FGE. If an ETS was annotated as supporting the claim, the crowd workers selected only sup-porting sentences; if the ETS was annotated as refuting the claim, only refuting sentences were selected.
Kwiatkowski et al. (201	9)Question answering	Long Answer Identification: For good questions only, annotators select the earliest HTML bounding box containing enough information for a reader to completely infer the answer to the question. Bounding boxes can be paragraphs, tables, list items, or whole lists. Alternatively, annotators mark "no answer" if the page does not answer the question, or if the information is present but not contained in a single one of the allowed elements.
Wadden et al. (2020)	Fact verification	ca-An evidence set is a collection of sentences from the abstract that provide support or contradiction for the given claim. To decide whether a collection of sentences is an evidence set, ask yourself, "If I were shown only these sentences, could I reasonably conclude that the claim is true (or false)"? 1) Evidence sets should be minimal. If you can remove a sentence from the evidence set and the remaining sentences are sufficient for support / contradiction, you should remove it. 2) There may be multiple evidence sets in a given abstract. See more at https://scifact.s3-us-west-2.amazonaws.com/doc/evidence-annotation-instructions.pdf
Kutlu et al. (2020)	relevance assessment	Please copy and paste text 2-3 sentences from the webpage which you believe support your decision. For instance, if you selected Highly Relevant, paste some text that you feel clearly satisfies the given query. If you selected Definitely not relevant, copy and paste some text that shows that the page has nothing to do with the query. If there is no text on the page or images led you to your decision, please type "The text did not help me with my decision".

Reference	Task	Questions and guidelines
Jansen et al. (2016)	science QA	For each question, we create gold explanations that describe the inference needed to arrive at the correct answer. Our goal is to derive an explanation corpus that is grounded in grade-appropriate resources. Accordingly, we use two elementary study guides, a science dictionary for elementary students, and the Simple English Wiktionary as relevant corpora. For each question, we retrieve relevant sentences from these corpora and use them directly, or use small variations when necessary. If relevant sentences were not located, then these were constructed using simple, straightforward, and grade-level appropriate language. Approximately 18% of questions required specialized domain knowledge (e.g. spatial, mathematical, or other abstract forms) that did not easily lend itself to simple verbal description, which we removed from consideration. This resulted in a total of 363 gold explanations.
Rajani et al. (2019)	Question answering	Turkers are prompted with the following question: "Why is the predicted output the most appropriate answer?" Annotators were in- structed to highlight relevant words in the question that justifies the ground-truth answer choice and to provide a brief open-ended explanation based on the highlighted justification could serve as the commonsense reasoning behind the question. Annotators cannot move forward if they do not highlight any relevant words in the question or if the length of explanations is less than 4 words. We also check that the explanation is not a sub- string of the question or the answer choices without any other extra words. We collect these ex- planations from only one annotator per example, so we also perform some post-collection checks to catch examples that are not caught by our previ- ous filters. We filter out explanations that could be classified as a template. For example, explanations of the form " <answer> is the only option that is [correct—obvious]" are deleted and then reannotated.</answer>
Sap et al. (2020)	social bias	What aspect/stereotype/characteristic of this group (often unfairly assumed) is referenced or implied by this post? — Use simple phrases and do not copy paste from the post.

Table 6: Questions that prior work uses to solicit human explanations for **proximal mechanisms** (in free text).

Re	eference	Task	Questions and guidelines
	e et al. (2020)	Natual language inference	We encouraged the annotators to focus on the non-obvious elements that induce the given relation, and not on the parts of the premise that are repeated identically in the hypothesis. For entailment, we required justifications of all the parts of the hypothesis that do not appear in the premise. For neutral and contradictory pairs, while we encouraged stating all the elements that contribute to the relation, we consider an explanation correct, if at least one element is stated. Finally, we asked the annotators to provide self-contained explanations, as opposed to sentences that would make sense only after reading the premise and hypothesis. We did in-browser checks to ensure that each explanation contained at least three tokens and that it was not a copy of the premise or hypothesis. We further guided the annotators to provide adequate answers by asking them to proceed in two steps. First, we require them to highlight words from the premise and/or hypothesis that they consider essential for the given relation. Secondly, annotators had to formulate the explanation using the words that they highlighted. However, using exact spelling might push annotators to formulate grammatically incorrect sentences, therefore we only required half of the highlighted words to be used with the same spelling. For entailment pairs, we required at least one word in the premise to be highlighted. For contradiction pairs, we required highlighting at least one word in both the premise and the hypothesis. For neutral pairs, we only allowed highlighting words in the hypothesis, in order to strongly emphasize the asymmetry in this relation and to prevent workers from confusing the premise with the hypothesis. We believe these label-specific constraints helped in putting the annotator into the correct mindset, and additionally gave us a means to filter incorrect explanations. Finally, we also checked that the annotators used other words that were not highlighted, as we believe a correct explanation would need to articulate a link between t
	o et al. (2020) im et al. (2018)	self-driving cars	We provide a driving video and ask a human annotator in Amazon Mechanical Turk to imagine herself being a driving instructor. Note that we specifically select human annotators who are familiar with US driving rules. The annotator has to describe what the driver is doing (especially when the behavior changes) and why, from a point of view of a driving instructor. Each described action has to be accompanied with a start and end time-stamp. The annotator may stop the video, forward and backward through it while searching for the activities that are interesting and

justifiable.

Reference	Task	Questions and guidelines
Zhang et al. (2020)	coreference resolution	Given a context and a pronoun reference relationship, write how you would decide the selected candidate is more likely to be referred than the other candidate using natural language. Don't try to be overly formal, simply write what you think. In the first phase, we ask annotators to provide reasons for all WSC questions. Detailed instructions are provided such that annotators can fully understand the task1. As each question may have multiple plausible reasons, for each question, we invite five annotators to provide reasons based on their own judgments. A screenshot of the survey is shown in Figure 3. As a result, we collect 1,365 reasons. As the quality of some given reasons might not be satisfying, we introduce the second round annotation to evaluate the quality of collected reasons. In the second phase, for each reason, we invite five annotators to verify whether they think the reason is reasonable or not2. If at least four annotators think the reason is plausible, we will accept that reason. As a result, we identify 992 valid reasons.
Lei et al. (2020) Da et al. (2020)	future event pre- diction harm of manip- ulated images	we also require them to provide a rationale as to why it is more or less likely For each question we require annotators to provide both an answer to the question and a rationale (e.g. the physical change in the image edit that alludes to their answer). This is critical, as the rationales prevent models from guessing a response such as "would be harmful" without providing the proper reasoning for their response. We ask annotators to explicitly separate the rationale from the response by using the word "because" or "since" (however, we find that the vast majority of annotators naturally do this, without being explicitly prompted).
Ehsan et al. (2019)	gaming	"Please explain your action". During this time, the player's microphone automatically turns on and the player is asked to explain their most recent action while a speech-to-text library automatically transcribes the explanation real-time.

Table 8: Questions that prior work uses to solicit human explanations for proximal mechanisms (in free text).

Reference	Task	Questions and guidelines
Ling et al. (2017)	algebraic prob- lems	cannot infer the exact question
Alhindi et al. (2018)	fact verification	we cannot infer the exact question automatically extracting for each claim the justification that humans have provided in the fact-checking article associated with the claim. Most of the articles end with a summary that has a headline "our ruling" or "summing up"
Kotonya and Toni (2020)	fact verification	automatically scraped from the website, we cannot infer the exact question
Wang et al. (2020)	sentiment analysis & relation extraction	Turkers are prompted with a list of selected predicates (see Appendix) and several examples of NL explanations. We cannot infer the exact question
Brahman et al. (2020)	natural lan- guage inference	automatically generated. We cannot infer the exact question
Li et al. (2018)	visual QA	automatically generated, We cannot infer the exact question
Park et al. (2018)	visual QA	During data annotation, we ask the annotators to complete the sentence "I can tell the person is doing (action) because" where the action is the ground truth activity label. However, We cannot infer the exact question in VQA-X.
Rajani et al. (2020)	physics reason- ing	We cannot infer the exact question

Table 9: Questions that prior work uses to solicit human explanations for proximal mechanisms (in free text).

Reference	Task	Questions and guidelines
Jansen et al. (2018)	science QA	Specific interfaces were designed. For a given question, annotators identified the central concept the question was testing, as well as the inference required to correctly answer the question, then began progressively constructing the explanation graph. Sentences in the graph were added by querying the tablestore based on key- words, which retrieved both single sentences/table rows, as well as entire explanations that had been previously annotated. If any knowledge required to build an explanation did not exist in the table store, this was added to an appropriate table, then added to the explanation.
Xie et al. (2020)	science QA	similar to Jansen et al. (2018)
Khot et al. (2020)	question answering	
Jhamtani and Clark (2020)	question answering	
Inoue et al. (2020)	question answering	

Table 10: Questions that prior work uses to solicit human explanations in **proximal mechanisms** (in structured explanations).

Reference	Task	Questions and guidelines
Srivastava et al. (2017) Hancock et al. (2018)	ing	The screenshot includes both "explanations" and "instructions", however, we cannot infer the exact question we cannot infer the exact question
Hallcock et al. (2018)	tion	we cannot lifter the exact question

Table 11: Questions that prior work uses to solicit human explanations for procedure (in free text).

Reference	Task		Questions and guidelines
Lamm et al. (2020)	question swering	an-	referential equality, we cannot infer the exact question
Ye et al. (2020)	question swering	an-	Please read carefully to get accepted! (1) You're not required to answer the question. The answer is already provided and marked in red. Read examples below carefully to learn about what we want! (2) Identify important short phrases that appear both in the question and in the context. Important: The two appearances of the phrase should be exactly the same (trivial differences like plural form or past tense are still acceptable). Important: Write sentences like Y is "Switzerland". Make sure there is no typo in what you quote. (3) Explain how you locate the answer with the phrases you marked; Only use the suggested expressions in the table in the bottom.
Geva et al. (2021)	question swering	an-	1) Creative question writing: Given a term (e.g., silk), a description of the term, and an expected answer (yes or no), the task is to write a strategy question about the term with the expected answer, and the facts required to answer the question. 2) Strategy question decomposition: Given a strategy question, a yes/no answer, and a set of facts, the task is to write the steps needed to answer the question. 3) Evidence matching: Given a question and its de- composition (a list of single-step questions), the task is to find evidence paragraphs on Wikipedia for each retrieval step. Operation steps that do not require retrieval are marked as operation.

Table 12: Questions that prior work uses to solicit human explanations for procedure (in structured explanations).