

# On the Origin of Hallucinations in Conversational Models: Is it the Datasets or the Models?

Anonymous ACL submission

## Abstract

Although knowledge-grounded conversational models are able to generate fluent responses, they are known to suffer from producing factually invalid statements, a phenomenon commonly called hallucination. In this work, we investigate the underlying causes of this phenomenon: is hallucination due to the training data, or to the models? We conduct a comprehensive human study on both existing knowledge-grounded conversational benchmarks and several state-of-the-art models. Our study reveals that the standard benchmarks consist of more than 60% hallucinated responses, leading to models that not only hallucinate but even amplify hallucinations. We hope these insights will show the way forward towards building hallucination-free conversational models.

## 1 Introduction

In recent years, knowledge-grounded conversational models, powered by large pre-trained language models (Radford et al., 2019; Brown et al., 2020; Raffel et al., 2020), achieved a remarkable performance in generating fluent and natural-looking responses. However, these systems exhibit an undesirable phenomenon commonly called *hallucination*. Hallucinated outputs are characterized as being unfaithful to some externally provided knowledge (Dziri et al., 2021b; Rashkin et al., 2021). The common belief in the literature is that researchers need to fix the models in order to fix hallucinations (Shuster et al., 2021; Mielke et al., 2020; Dziri et al., 2021a; Rashkin et al., 2021) and no attempt has been made so far to audit the conversational benchmarks to the best of our knowledge. Nonetheless, it is not yet well-understood why conversational models have a propensity to hallucinate; is it because conversational benchmarks are noisy and contain factually incorrect sentences or does it stem from potential shortcomings in models’ architectures and/or training procedures? In this work,

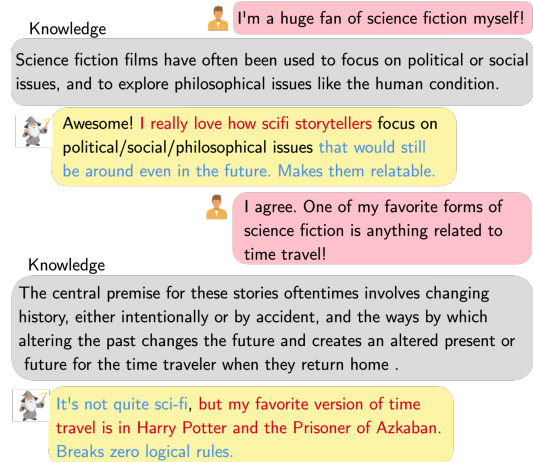


Figure 1: An example of a hallucinated conversation from the Wizard of Wikipedia dataset (Dinan et al., 2018). The wizard (yellow) is hallucinating information that cannot be inferred from the knowledge-snippet: hallucinated subjective content (red) and hallucinated objective content (blue).

we investigate both existing benchmarks and generated responses of prominent conversational models to shed light on the origins of hallucinations.

On one hand, knowledge-grounded conversational benchmarks may contain hallucinations *due to error-prone collection protocols and lack of adequate quality control*. Existing dialogue systems are typically trained on corpora crowd-sourced through online platforms (Dinan et al., 2018; Zhou et al., 2018; Gopalakrishnan et al., 2019; Moon et al., 2019). With loose incentive to come up with faithfully-grounded utterances on the provided knowledge, crowdworkers may ignore knowledge-snippets altogether, use their personal knowledge or sometimes assume a fictional persona, resulting in conversations that are rife with subjective content and unverified factual knowledge. Figure 1 shows a hallucinated conversation from the Wizard of Wikipedia dataset (Dinan et al., 2018),

On the other hand, neural conversational models are not *necessarily designed to generate faithful outputs, but to mimic the distributional properties of the data*. This kind of optimization will likely

push the models to replicate and even amplify the hallucination behaviour at test time (Bender et al., 2021). The presence of even few hallucinated responses may skew the data distribution in a way that curbs the model’s ability to generate faithful responses (Kang and Hashimoto, 2020).

In-depth understanding of the various sources of hallucination and how they manifest themselves can help researchers enforce faithfulness in conversational models and thus reduce hallucinations. In this work, we take a step closer to gain such an understanding via a systematic study where we identify and categorize hallucinations in the widely-used benchmarks, measure their frequency, and overall negative impact on generated responses as judged by human evaluators. Specifically, drawing insights from the linguistic coding system for discourse phenomena called Verbal Response Modes (VRM, Stiles, 1992), we manually annotate conversations from the three widely-used knowledge-grounded conversational benchmarks: Wizard of Wikipedia (Dinan et al., 2018), CMU-DoG (Zhou et al., 2018) and TopicalChat (Gopalakrishnan et al., 2019). Our analysis reveals that more than 60% of the responses are hallucinated in the three datasets, with major hallucination modes that manifest principally through the expression of subjective information (e.g., thoughts, beliefs, feelings, intentions, personal experiences) and the expression of unsupported objective factual information.

Further, to understand if neural conversational models make this hallucination more severe, we annotate responses generated by several state-of-the-art models, including ones that are designed to alleviate hallucinations. We find that the generated responses consist of an even larger portion of hallucinations, in comparison with the training data. Our findings, thus, question the quality of current knowledge-grounded conversational datasets and challenge the robustness of existing conversational models.

## 2 Hallucinations in Benchmarks

We conduct study on three English crowdsourced knowledge-grounded conversational benchmarks: Wizard of Wikipedia (WoW), CMU-DoG and TopicalChat. These datasets consist of dialogues between two speakers, where the goal is to communicate information about particular topics while speakers are presented with a knowledge snippet relevant to the current turn. More details about

these datasets are provided in Appendix A.

### 2.1 Response Classification Taxonomy

We sample 200 random knowledge-grounded responses from training sets of each benchmark, and annotate each response based on whether it can be inferred exclusively from the knowledge-snippet. Inspired by the BEGIN taxonomy (Dziri et al., 2021b) of response classification, we annotate each turn with labels as follows:<sup>1</sup> **Entailment**: a response is fully supported by the knowledge. **Hallucination**: a response’s factual correctness cannot be fully verified from the knowledge-snippet (even if the response is true in the real world). **Partial Hallucination**: part of the response is hallucinated while the rest is entailed. **Generic**: a response that is vague and does not convey any factual information. **Uncooperative**: an entailed response that does not follow the principles of conversational cooperation according to Gricean maxims (Grice, 1989). We provide more details of these classes and their examples in Appendix C.

**(Q1) How much hallucination exists in the benchmarks?** For each dialogue turn, we solicit judgements from two linguists who are experts in the task in order to obtain high-quality annotations. The inter-annotator agreement measured through average Krippendorff’s alpha coefficient is 0.91 which indicates high agreement. Figure 2 shows the breakdown of different conversational phenomenon in the three benchmarks and Table 5 (Appendix D) depicts exemplars of hallucinated responses. Surprisingly, the three benchmarks are fraught with hallucinations. TopicalChat contains 62% responses that are purely hallucinated against only 12% responses that are fully entailing the source knowledge. Our analysis shows similar trends in the CMU\_DoG dataset. On the other hand, in WoW, hallucinated responses are largely mixed with faithful content (37.7% v.s. 22.5% fully hallucinated responses), which amounts to 60.2% hallucinations in total. These findings raise the alarm on the quality of the widely used dialogue datasets in the community.

### 2.2 The Linguistic Nature of Hallucinations

To understand the linguistic nature of hallucinations, we further annotate responses based on a

<sup>1</sup>We omit the label contradiction and off-topic from BEGIN as we consider it a subcategory of hallucination.

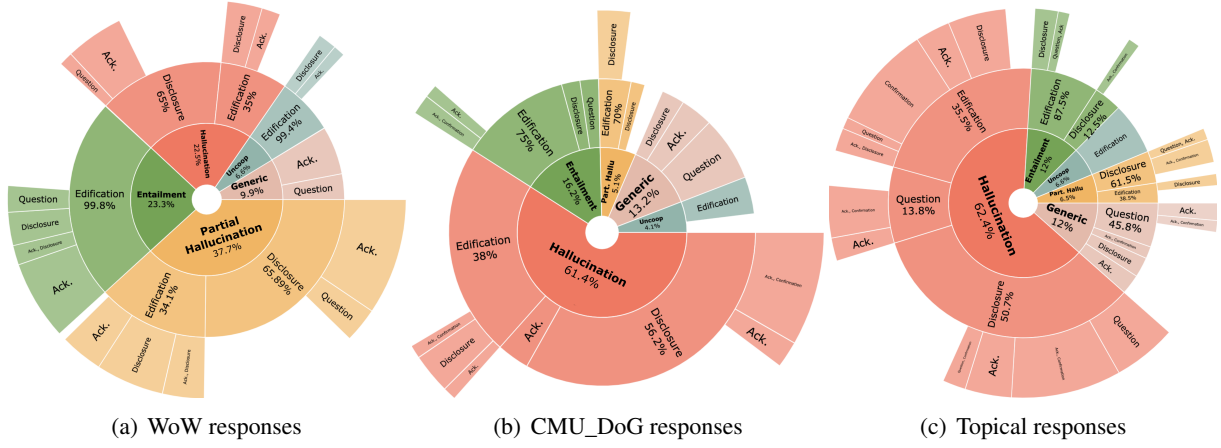


Figure 2: BEGIN and VRM breakdown of responses from WoW, CMU\_DoG and TopicalChat. The inner circle shows the breakdown of BEGIN classes and the outer shows the VRM types in each BEGIN type: Hallucination (red), Entailment (green), Partial Hallucination (yellow), Generic (pink), and Uncooperative (blue).

VRM Type	Description	Example
Disclosure	Reveal the speaker’s subjective opinions, personal experience, thoughts, feelings, wishes, and intentions.	“I think science fiction is an amazing genre. Future science, technology they’re all interesting.”
Edification	Concerns information that is, in principle, objective.	“Recycling includes items like metal and plastic.”
Advisement	Corresponds to guiding the behaviour of the addressee through: commands, requests, suggestions, advice, permission, prohibition.	“You should be patient and persistent to succeed.”
Confirmation	Compares the speaker’s experience with the other’s by expressing shared ideas/memories/beliefs, or by agreement/disagreement	“I agree that love encompasses a variety of different emotional and mental states.”
Question	Concerns requesting information or guidance.	“What is your favorite song?”
Acknowledge	Expresses no content, it conveys only receipt of communication from the other’s speaker.	“Mmm. OK,...”, “Yeah, ...”, “Hello, ...”

Table 1: The definitions of the VRM types with examples.

linguistic coding system for discourse phenomena, dubbed Verbal Response Modes (VRM; Stiles 1992). Concretely, we label a turn with the following speech acts: **Disclosure**, **Edification**, **Advisement**, **Confirmation**, **Question** and **Acknowledgement (Ack.)**. Table 1 displays the definition for each VRM type along with examples.

**(Q2) What are the hallucination strategies used in human-human data?** Figure 2 shows the VRM breakdown for each BEGIN category in the three benchmarks. We make the following observations: The majority of hallucinations belong to *disclosure* (i.e., subjective information) in all benchmarks (65.9%, 56.2% and 50.7% in WoW, CMU\_DoG and TopicalChat respectively). Although the strategy of sharing subjective information such as thoughts, opinions and feelings is natural in conversations, it often comes at a cost of ignoring the knowledge snippet in these datasets. Moreover, *edification* is also a common phenomenon in hallucinated responses, suggesting that humans not only discuss subjective informa-

tion but also bring extra unsupported facts, either true or false. Other linguistic modes are also associated with hallucinations such as acknowledging unsupported claims or asking irrelevant questions. Conversely, entailment responses have high percentage of edification ( $> 70\%$ ) with information inferred from the knowledge snippet.

### 3 Hallucination Amplification in Models

Next, we investigate how much models amplify the hallucination phenomenon at inference time. We consider a range of representative models:

- **GPT2** (Radford et al., 2019; Wolf et al., 2019) is an autoregressive model which takes as input a concatenation of the knowledge and the history.
- **DoHA** (Prabhumoye et al., 2021) builds a BART-based conversational model (Lewis et al., 2020) for knowledge-grounding, with a two-view attention mechanism to handle separately the encoded document and the history during generation.
- **CTRL** (Rashkin et al., 2021) augments the GPT2 model with control tokens (Keskar et al., 2019) that

	Model	ROUGE $\uparrow$	Hallucination Rate $\downarrow$			Entailment Rate $\uparrow$		
			Full	Partial	Overall	Entail.	Uncoop.	Overall
WoW	Gold	–	22.5	37.7	60.2	23.2	6.6	29.8
	GPT2	19.1	66.0	15.2	81.2	11.7	3.6	15.3
	DoHA	21.5	39.6	28.9	68.5	12.7	7.1	19.8
	CTRL	24.4	29.0	34.0	<b>63.0</b>	12.2	17.3	29.5
CMU DoG	Gold	–	61.4	5.1	66.5	16.2	4.1	20.3
	GPT2	13.7	75.0	6.0	81.0	5.0	5.5	10.5
	DoHA	15.4	62.5	10.0	72.5	8.1	5.9	14.0
	CTRL	19.3	62.3	6.0	<b>68.3</b>	9.1	12.4	19.8
Topical	Gold	–	62.4	6.5	68.9	12.0	6.6	18.6
	GPT2	13.2	71.5	8.5	80.5	6.0	5.5	11.5
	DoHA	16.3	52.5	26.5	79.0	9.0	5.0	14.0
	CTRL	18.3	47.5	29.5	<b>77.0</b>	8.5	10.0	18.5

Table 2: Amplification of models on the test data from WoW and CMU\_DoG and TopicalChat. ‘Entail.’ and ‘Uncoop.’ mean entailment and uncooperative.

guide the generation towards less subjective and more entailed content.

We fine-tune each model on the benchmarks and use nucleus sampling (Holtzman et al., 2019) with  $p = 0.6$  for decoding (more implementation details are in Appendix B). As seen in Table 2, CTRL is the best model followed by DoHA based on the ROUGE score. Table 4 in Appendix F shows a sample of generated responses. Similar to the analysis in §2, we task the two linguists to analyze model-generated responses for 200 randomly-selected test samples from each benchmark. We seek to answer the following questions:

**(Q3) Do state-of-the-art conversational models amplify hallucination?** Table 2 shows the degree of amplification across different models trained on the three benchmarks. Contrasting this with human gold responses, the models not only hallucinate but also amplify the percentage of hallucinations. For example, GPT2 amplifies full hallucination by 21% in WoW, 14.5% in CMU\_DoG and 11.6% in TopicalChat. Conversely, it reduces entailment by 14.5%, 9.8% and 7.1% respectively. This suggests that hallucination patterns are easier to learn than entailment. Among the three, CTRL hallucinates less followed by DoHA. Overall, these results demonstrate that hallucination is not only a reflection of training data issues, but also a consequence of the weaknesses of models. We hypothesize that there are multiple factors that can contribute to the models’ deficiencies including teacher forcing (Ranzato et al., 2016), maximum likelihood estimation (Kang and Hashimoto, 2020), bias in pre-trained LMs (Nadeem et al., 2021) and decoding strategies (Shuster et al., 2021). We leave investigating the role of each factors to hallucination amplification for future work.

**(Q4) What are the hallucination strategies used by models?** Surprisingly, different models use different strategies for hallucination. While DoHA and GPT2 predominantly rely on and amplify *disclosure*, CTRL relies on *edification*. This is because CTRL is trained explicitly to avoid pronouns (a crucial ingredient for disclosure) and to generate entailed responses. As a side-effect, it ends up amplifying uncooperative responses (to 260%, 300% and 151% as seen in Table 2). Full results of all models and datasets are in Appendix E.

## 4 Related Work

The problem of hallucination in neural language generation has been receiving an increased attention from the NLP community in recent years, including machine translation (Raunak et al., 2021; Wang and Sennrich, 2020) and summarization (Durmus et al., 2020; Lewis et al., 2020; Kang and Hashimoto, 2020). Hallucinations in knowledge-grounded neural dialogue generation is a relatively new research problem (Roller et al., 2021; Mielke et al., 2020; Shuster et al., 2021; Dziri et al., 2021a; Rashkin et al., 2021). While existing methods focus on addressing hallucinations in models, we focus on investigating their root causes and analyzing their nature. Closest to our work are Dziri et al. (2021b) and Santhanam et al. (2021) who introduce testbeds for quantifying groundedness in dialogue systems including hallucinations, whereas we resort to a much finer-grained manual analysis on multiple benchmarks and models.

## 5 Conclusion

Our investigations demonstrate empirically that hallucination is a prevalent issue in both dialog benchmarks and models. Our analysis on three widely used benchmarks reveals that they are rife with hallucinations (more than half of the data), and the most common strategies people use are *disclosure* and *edification*. Moreover, we show that conversational models trained on these benchmarks not only hallucinate but also amplify hallucinations, even the models that were designed to alleviate this issue. These results indicate that the data is not the only responsible factor but also the models. The community shall urgently reconsider the quality of the datasets as well as the the models. Following checklists for responsible data collection (Bender and Friedman, 2018; Rogers et al., 2021) could circumvent some of these problems in future datasets.



## Impact Statement

Our analytical study reveals that a large portion of standard knowledge-grounded dialogue benchmarks is hallucinated, leading us to reflect on the potential harm of low-quality data releases for conversational models. In recent years, the conversational AI market has seen a proliferation of a variety of applications—which are powered by large pre-trained LMs—that span across a broad range of domains, such as customer support, education, e-commerce, health, entertainment, etc (Vakulenko et al., 2021). Ensuring that these systems are trustworthy is key to deploy systems at a large scale in real-world application, especially in high-stake domains (Sambasivan et al., 2021). However, even if we come up with a model that is robust enough against hallucination, it will be ultimately bounded by the data quality.

## References

- Emily M. Bender and Batya Friedman. 2018. [Data statements for natural language processing: Toward mitigating system bias and enabling better science](#). *Transactions of the Association for Computational Linguistics*, 6:587–604.
- Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 610–623.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Emily Dinan, Stephen Roller, Kurt Shuster, Angela Fan, Michael Auli, and Jason Weston. 2018. Wizard of wikipedia: Knowledge-powered conversational agents. In *International Conference on Learning Representations*.
- Esin Durmus, He He, and Mona Diab. 2020. [FEQA: A question answering evaluation framework for faithfulness assessment in abstractive summarization](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5055–5070, Online. Association for Computational Linguistics.
- Nouha Dziri, Andrea Madotto, Osmar Zaiane, and Avishek Joey Bose. 2021a. [Neural path hunter: Reducing hallucination in dialogue systems via path grounding](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 2197–2214, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Nouha Dziri, Hannah Rashkin, Tal Linzen, and David Reitter. 2021b. Evaluating groundedness in dialogue systems: The begin benchmark. *arXiv preprint arXiv:2105.00071*.
- Karthik Gopalakrishnan, Behnam Hedayatnia, Qinglang Chen, Anna Gottardi, Sanjeev Kwatra, Anu Venkatesh, Raefer Gabriel, and Dilek Hakkani-Tür. 2019. [Topical-Chat: Towards Knowledge-Grounded Open-Domain Conversations](#). In *Proc. Interspeech 2019*, pages 1891–1895.
- Paul Grice. 1989. *Studies in the Way of Words*. Harvard University Press.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2019. The curious case of neural text degeneration. In *International Conference on Learning Representations*.
- Daniel Kang and Tatsunori B. Hashimoto. 2020. [Improved natural language generation via loss truncation](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 718–731, Online. Association for Computational Linguistics.
- Nitish Shirish Keskar, Bryan McCann, Lav R Varshney, Caiming Xiong, and Richard Socher. 2019. Ctrl: A conditional transformer language model for controllable generation. *arXiv preprint arXiv:1909.05858*.
- Diederik P. Kingma and Jimmy Ba. 2015. [Adam: A method for stochastic optimization](#). In *ICLR (Poster)*.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. [BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.
- Sabrina J Mielke, Arthur Szlam, Y-Lan Boureau, and Emily Dinan. 2020. Linguistic calibration through metacognition: aligning dialogue agent responses with expected correctness. *arXiv preprint arXiv:2012.14983*.

397	Seungwhan Moon, Pararth Shah, Anuj Kumar, and Rajen Subba. 2019. <a href="#">OpenDialKG: Explainable conversational reasoning with attention-based walks over knowledge graphs</a> . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 845–854, Florence, Italy. Association for Computational Linguistics.	455
398		456
399		
400		457
401		458
402		459
403		460
404	Moin Nadeem, Anna Bethke, and Siva Reddy. 2021. <a href="#">StereoSet: Measuring stereotypical bias in pre-trained language models</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 5356–5371, Online. Association for Computational Linguistics.	461
405		462
406		463
407		464
408		465
409		466
410		467
411		468
412	Shrimai Prabhumoye, Kazuma Hashimoto, Yingbo Zhou, Alan W Black, and Ruslan Salakhutdinov. 2021. <a href="#">Focused attention improves document-grounded generation</a> . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 4274–4287, Online. Association for Computational Linguistics.	469
413		470
414		
415		471
416		472
417		473
418		474
419		475
420	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9.	476
421		
422		477
423		478
424	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. <a href="#">Exploring the limits of transfer learning with a unified text-to-text transformer</a> . <i>Journal of Machine Learning Research</i> , 21(140):1–67.	479
425		480
426		481
427		482
428		483
429		
430	Marc’Aurelio Ranzato, Sumit Chopra, Michael Auli, and Wojciech Zaremba. 2016. <a href="#">Sequence level training with recurrent neural networks</a> . In <i>4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings</i> .	484
431		485
432		486
433		487
434		488
435		
436	Hannah Rashkin, David Reitter, Gaurav Singh Tomar, and Dipanjan Das. 2021. <a href="#">Increasing faithfulness in knowledge-grounded dialogue with controllable features</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 704–718, Online. Association for Computational Linguistics.	491
437		492
438		493
439		494
440		495
441		496
442		497
443		498
444		499
445	Vikas Raunak, Arul Menezes, and Marcin Junczys-Dowmunt. 2021. <a href="#">The curious case of hallucinations in neural machine translation</a> . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1172–1183, Online. Association for Computational Linguistics.	500
446		
447		501
448		502
449		503
450		504
451		505
452	Anna Rogers, Timothy Baldwin, and Kobi Leins. 2021. ‘just what do you think you’re doing, dave?’ a checklist for responsible data use in nlp. In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 4821–4833.	506
453		507
454		508
		509
		510
	Stephen Roller, Emily Dinan, Naman Goyal, Da Ju, Mary Williamson, Yinhan Liu, Jing Xu, Myle Ott, Eric Michael Smith, Y-Lan Boureau, and Jason Weston. 2021. <a href="#">Recipes for building an open-domain chatbot</a> . In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 300–325, Online. Association for Computational Linguistics.	
	Nithya Sambasivan, Shivani Kapania, Hannah Highfill, Diana Akrong, Praveen Paritosh, and Lora M Aroyo. 2021. “everyone wants to do the model work, not the data work”: Data cascades in high-stakes ai. In <i>proceedings of the 2021 CHI Conference on Human Factors in Computing Systems</i> , pages 1–15.	
	Sashank Santhanam, Behnam Hedayatnia, Spandana Gella, Aishwarya Padmakumar, Seokhwan Kim, Yang Liu, and Dilek Hakkani-Tur. 2021. Rome was built in 1776: A case study on factual correctness in knowledge-grounded response generation. <i>arXiv preprint arXiv:2110.05456</i> .	
	Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela, and Jason Weston. 2021. <a href="#">Retrieval augmentation reduces hallucination in conversation</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 3784–3803, Punta Cana, Dominican Republic. Association for Computational Linguistics.	
	Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: a simple way to prevent neural networks from overfitting. <i>The journal of machine learning research</i> , 15(1):1929–1958.	
	William B Stiles. 1992. <i>Describing talk: A taxonomy of verbal response modes</i> . Sage Publications.	
	Svitlana Vakulenko, Evangelos Kanoulas, and Maarten de Rijke. 2021. A large-scale analysis of mixed initiative in information-seeking dialogues for conversational search. <i>arXiv preprint arXiv:2104.07096</i> .	
	Chaojun Wang and Rico Sennrich. 2020. <a href="#">On exposure bias, hallucination and domain shift in neural machine translation</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 3544–3552, Online. Association for Computational Linguistics.	
	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. <a href="#">Trans-formers: State-of-the-art natural language processing</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing</i> .	

511 *System Demonstrations*, pages 38–45, Online. Asso-  
512 ciation for Computational Linguistics.

513 Thomas Wolf, Victor Sanh, Julien Chaumond, and  
514 Clement Delangue. 2019. Transfertransfo: A  
515 transfer learning approach for neural network  
516 based conversational agents. *arXiv preprint*  
517 *arXiv:1901.08149*.

518 Kangyan Zhou, Shrimai Prabhumoye, and Alan W  
519 Black. 2018. [A dataset for document grounded con-](#)  
520 [versations](#). In *Proceedings of the 2018 Conference*  
521 *on Empirical Methods in Natural Language Process-*  
522 *ing*, pages 708–713, Brussels, Belgium. Association  
523 for Computational Linguistics.

## A Datasets

We conduct our analysis on the following datasets:

**Wizard of Wikipedia:** consists of dialogues between a “wizard” and an “apprentice”, where the goal of the wizard is to communicate information about a particular topic. The Apprentice, in turn, is expected to seek information about the topic. At each turn, the wizard is presented with a knowledge snippet (typically a sentence) from Wikipedia while the apprentice is not; and the Wizard is allowed to form an utterance that does not use the evidence. We omit data points in which the wizard did not explicitly select a passage as evidence for the response. The dataset consists of 82722 grounded-responses in train, 8800 valid and 8690 test.

**CMU\_DoG:** consists of conversations focusing only on the movie domain. Each response is grounded on a document (typically a section from Wikipedia). Workers are asked to either persuade the other speakers to watch the movie based on the knowledge or to discuss the content of the document with them. In total, there are 78136 grounded responses in train, 13800 in valid and 13796 in test.

**TopicalChat:** consists of dialogues conversing around a variety of topics. Workers are provided with relevant facts from Reddit, Wikipedia and news articles. The collection protocol consists of two scenarios: symmetric and asymmetric. In the symmetric scenario, workers will have access to the same source knowledge and in the asymmetric scenario, they will have access to different sources. In total, the dataset has 292215 grounded responses in train, 23601 in valid and 23623 in test.

## B Implementation Details

**GPT2:** We implement this model using the Pytorch Huggingface Transformers library (Wolf et al., 2020) and the Pytorch-lightning library<sup>2</sup>. During training, we use the Adam optimizer (Kingma and Ba, 2015) with Dropout (Srivastava et al., 2014) on a batch size of 32 with a learning rate of  $6.25 \times 10^{-5}$  that is linearly decayed. The maximum dialogue history length is set to 3 utterances. The model early-stops at epoch {6, 10, 10} respectively for WoW, CMU\_DoG and TopicalChat. The average runtime is {1.5, 3, 3}

<sup>2</sup><https://github.com/PyTorchLightning/pytorch-lightning>

hours for WoW, CMU\_DoG and TopicalChat respectively.

**DoHA:** We use the code and the pre-trained model on CMU\_DoG that are publicly available by the authors at their Github’s account<sup>3</sup>. For WoW and TopicalChat, we follow closely the authors’ training procedure described in (Prabhumoye et al., 2021) and we trained two models on both datasets. The average runtime of these models is {5, 10} hours for WoW and TopicalChat respectively.

**CTRL:** The code is not publicly available for this model. We were able to reproduce the results ourselves by following training details in the paper and having multiple discussions with the authors. Similar to GPT2, we implement this model using the Pytorch Huggingface Transformers library and the Pytorch-lightning library.

For each dataset, we save the best model based on the validation set. Training for all models is done on an Nvidia V100 GPU 32GB and for inference, we use nucleus sampling with  $p=0.6$ .

## C Definitions of the BEGIN Taxonomy

Here are more detailed definitions of each of the BEGIN classes. Examples can be found in Table 3:

- **Entailment:** The response can be verified as true based solely on the knowledge-snippet. Any factual information which it contains can be found in or derived from the knowledge provided.
- **Hallucination:** The response cannot be verified as true based solely on the knowledge-snippet. It is comprised of information which cannot be found in or derived from the knowledge provided.
- **Partial Hallucination:** Part of the response can be verified as true using the knowledge-snippet (i.e. is entailed), while another part of the response cannot (i.e. is hallucinated).
- **Generic:** The response is vague and does not contain factual claims. Such responses are often used simply to convey receptiveness to what the interlocutor has to say.
- **Uncooperative:** The response is entailed but violates the Gricean maxims of conversation

<sup>3</sup><https://bit.ly/3bBup2M>



614 (Grice, 1989). These responses may be per-  
615 ceived as rude, purposefully misleading, or  
616 showing a general unwillingness to cooperate  
617 with the interlocutor for effective communica-  
618 tion.

## 619 **D Hallucinated Human-Human** 620 **Responses**

621 Table 5 contains hallucinated gold responses from  
622 WoW, CMU\_DoG and TopicalChat.

## 623 **E Breakdown of BEGIN and VRM in** 624 **Machine-generated Responses**

625 Figure 3, 4 and 5 display the distribution of BEGIN  
626 and VRM in GPT2, DoHA and CTRL trained on  
627 the three benchmark.

## 628 **F Machine-generated Responses**

629 Table 4 contains a sample of generated responses  
630 from GPT2, DoHA and CTRL on the WoW and  
631 CMU\_DoG.



Figure 3: Breakdown of BEGIN classes and VRM speech acts on WoW machine-generated responses.

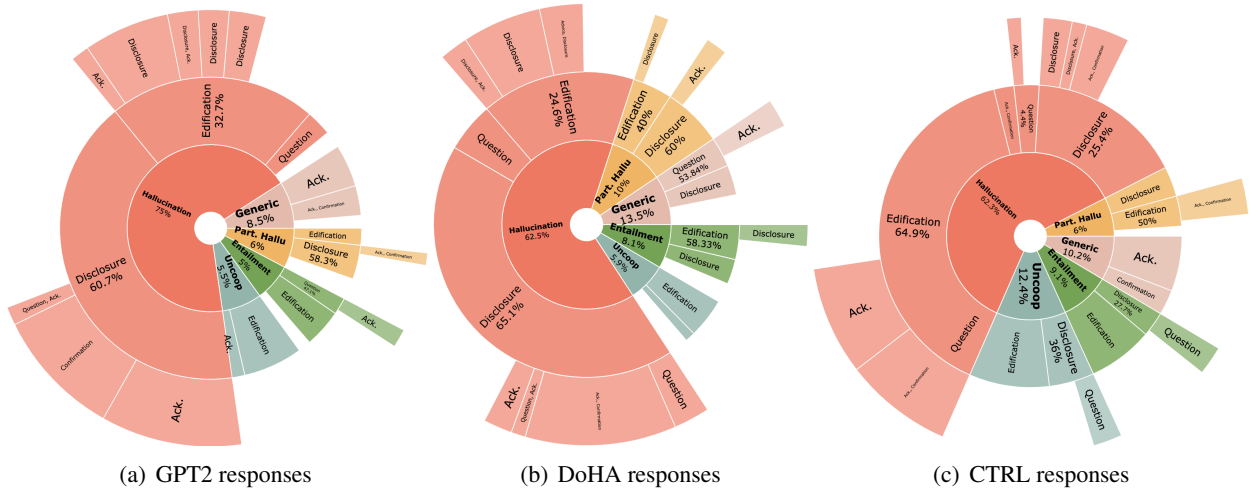


Figure 4: Breakdown of BEGIN classes and VRM speech acts on CMU\_DoG machine-generated responses.

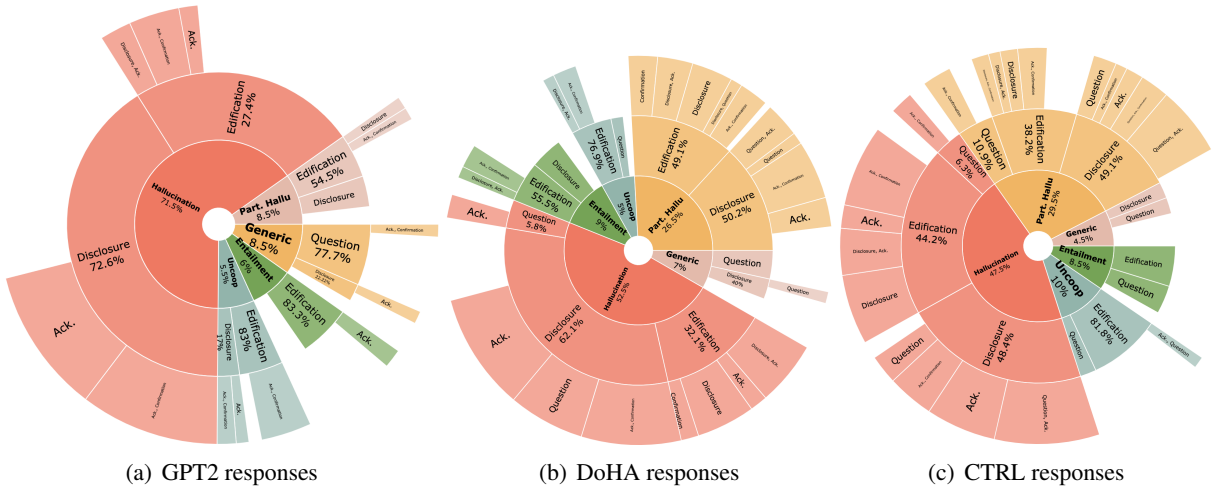


Figure 5: Breakdown of BEGIN classes and VRM speech acts on Topical machine-generated responses.

BEGIN	VRM	Example
Entailment	Disclosure	<p>Knowledge: A dragon is a legendary creature, typically scaled or fire-spewing and with serpentine, reptilian or avian traits, that features in the myths of many cultures around world.</p> <p>History: Dragons are so fascinating, I wonder where they originated from.</p> <p>Response: I'm not sure, but I know that it is a legendary creature featured in myths of many cultures around the world!</p>
Hallucination	Edification Disclosure	<p>Knowledge: The central premise for these stories oftentimes involves changing history, either intentionally or by accident, and the ways by which altering the past changes the future and creates an altered present or future for the time traveler when they return home.</p> <p>History: One of my favorite forms of science fiction is anything related to time travel! I find it fascinating.</p> <p>Response: It's not quite sci-fi, but my favorite version of time travel is in Harry Potter and the Prisoner of Azkaban. Breaks zero logical rules.</p>
Partial hallucination	Confirmation	<p>Knowledge: Love encompasses a variety of different emotional and mental states, typically strongly and positively experienced, ranging from the deepest interpersonal affection to the simplest pleasure.</p> <p>History: Romance can be small acts, like making breakfast for your significant other. But it's portrayed as grand gestures which are unattainable.</p> <p>Response: I agree it has been portrayed as impossible actions . Love consists of a variety of emotional and mental states.</p>
Generic	Question	<p>Knowledge: Hard rock is a loosely defined subgenre of rock music that began in the mid-1960s, with the garage, psychedelic and blues rock movements.</p> <p>History: yea i have heard of the American hard rock band Van Halen.</p> <p>Response: Are you a fan of hard rock?.</p>
Uncooperative	Edification	<p>Knowledge: Marketing is used to create, keep and satisfy the customer.</p> <p>History: Well, I wish I could. I think my training was a little off. I can't find a job and thinking about going back to school to be a counselor. What do you do?</p> <p>Response: Marketing is used to create a customer and then retain the customer and keep them happy.</p>

Table 3: Examples from Wizard of Wikipedia (Dinan et al., 2018) showing the BEGIN breakdown and different VRM linguistic phenomena for each response.

Data	Knowledge & History	Responses
WoW	<p><b>Knowledge:</b> Cinematography (also called "Direction of Photography") is the science or art of motion-picture photography by recording light or other electromagnetic radiation, either electronically by means of an image sensor, or chemically by means of a light-sensitive material such as film stock.</p> <p><b>History:</b> Hi buddy, What you think about cinematography?</p>	<p><b>GPT2:</b> Cinematography is another title I'm currently in school for. I've been doing the art of photography for over 15 years now. (Hallucination)</p> <p><b>DoHA:</b> I do enjoy the science and art of motion-picture photography. (Partial Hallucination)</p> <p><b>CTRL:</b> You're right. Cinematography is an art of motion -picture photography by recording light or other electromagnetic radiation, either electronically by means of an image sensor, or chemically by means of a light-sensitive material such as film stock. What about you ? (Entailment)</p>
CMU_DoG	<p><b>Knowledge:</b> Mean Girls is a 2004 American teen comedy film directed by Mark Waters and written by Tina Fey. The film is partially based on Rosalind Wiseman's 2002 non-fiction self-help book Queen Bees and Wannabes, which describes female high school social cliques and the damaging effects they can have on girls. Genre: Comedy. MovieName: Mean Girls rating: Rotten Tomatoes: 84 % and average: 6.9/10.</p> <p><b>History:</b> I think Rachel McAdams had an even better role as Regina George however!</p>	<p><b>GPT2:</b> Mean Girls was a classic. (Hallucination)</p> <p><b>DoHA:</b> Yes I agree! she was great and she was better than Mark Waters! (Hallucination)</p> <p><b>CTRL:</b> I see. Mean Girls is a comedy movie that was directed by Mark Waters in 2004. (Entailment)</p>

Table 4: Generated responses from different models based on Wizard of Wikipedia (Dinan et al., 2018) and CMU\_DoG (Zhou et al., 2018) test samples.

Data	Knowledge & History	Hallucinated Gold Response
WoW	<p><b>Knowledge:</b> Science fiction often explores the potential consequences of scientific and other innovations, and has been called a ‘literature of ideas’.</p> <p><b>History:</b> And that’s difficult to do when dealing with time travel . I actually haven’t seen the latest Harry Potter movies. Guess it’s time to check them out!</p>	If you really want a look at the potential negative consequences of scientific innovation, what you should check out is the TV show Fringe. Incredibly well written.
CMU_DoG	<p><b>Knowledge:</b> Movie: The Social Network. In October 2003, 19-year-old Harvard University student Mark Zuckerberg is dumped by his girlfriend Erica Albright. Returning to his dorm, Zuckerberg writes an insulting entry about Albright on his LiveJournal blog and then creates a campus website called Facemash by hacking into college databases to steal photos of female students, then allowing site visitors to rate their attractiveness. After traffic to the site crashes parts of Harvard’s computer network, Zuckerberg is given six months of academic probation. However, Facemash’s popularity attracts the attention of Harvard upperclassmen and twins Cameron and Tyler Winklevoss and their business partner Divya Narendra. The trio invites Zuckerberg to work on Harvard Connection, a social network featuring the exclusive nature of Harvard students and aimed at dating.</p> <p><b>History:</b> The movie is The Social Network. I personally do not like Facebook as a company.</p>	The movie portrays the founding of social networking website Facebook and the resulting lawsuits. It even has Justin Timberlake in it, I don’t think I’ve ever seen him act.
TopicalChat	<p><b>Knowledge:</b> Wikipedia: first paragraph in <a href="https://en.wikipedia.org/wiki/Google">https://en.wikipedia.org/wiki/Google</a></p> <p><b>Reddit facts:</b> A single Google search requires more computing power than it took to send Neil Armstrong and eleven other astronauts to the moon. Google Maps calculates traffic by tracking how fast Android devices are moving on the road instead of hiring people to mow the lawns around their headquarters. Google uses hundreds of live goats. On 16th August 2013, Google went down for about five minutes, and took 40% of web traffic with it. When there is a disputed border, Google maps tailors its maps to the claims of each country where the Internet browser is located.</p> <p><b>History:</b> Google provides online related services and products, which includes online ads, search engine and cloud computing.</p>	Yeah, their services are good. I ’m just not a fan of intrusive they can be on our personal lives.

Table 5: Hallucinated responses from different benchmarks: Wikipedia (Dinan et al., 2018), CMU\_DoG (Zhou et al., 2018) and TopicalChat (Gopalakrishnan et al., 2019). Text highlighted in red indicates hallucinated content.