

CarExpert: Leveraging Large Language Models for In-Car Conversational Question Answering

Md Rashad Al Hasan Rony¹, Christian Süß¹, Sinchana Ramakanth Bhat², Viju Sudhi², Julia Schneider⁴, Maximilian Vogel³, Roman Teucher², Ken E. Friedl¹, Soumya Sahoo²

¹BMW Group, ²Fraunhofer IAIS, ³BIG PICTURE GmbH, ⁴ONSEI GmbH
md-rashad-al-hasan.rony@bmw.de, christian.suess@bmw.de

Abstract

Large language models (LLMs) have demonstrated remarkable performance by following natural language instructions without fine-tuning them on domain-specific tasks and data. However, leveraging LLMs for domain-specific question answering suffers from severe limitations. The generated answer tends to hallucinate due to the training data collection time (when using off-the-shelf), complex user utterance and wrong retrieval (in retrieval-augmented generation). Furthermore, due to the lack of awareness about the domain and expected output, such LLMs may generate unexpected and unsafe answers that are not tailored to the target domain. In this paper, we propose CarExpert, an in-car retrieval-augmented conversational question-answering system leveraging LLMs for different tasks. Specifically, CarExpert employs LLMs to control the input, provide domain-specific documents to the extractive and generative answering components, and controls the output to ensure safe and domain-specific answers. A comprehensive empirical evaluation exhibits that CarExpert outperforms state-of-the-art LLMs in generating natural, safe and car-specific answers.

1 Introduction

Conversational question answering (CQA) has recently gained increased attention due to the advancements of Transformer-based (Vaswani et al., 2017) large language models (LLMs). These LLMs (Devlin et al., 2019; Brown et al., 2020; OpenAI, 2023; Touvron et al., 2023b) are nowadays widely adopted for performing question answering in both open-domain and domain-specific settings (Robinson and Wingate, 2023). As the source of additional knowledge conversational question answering systems are typically provided with text paragraphs (Kim et al., 2021; Rony et al., 2022c), and knowledge graphs (Rony et al., 2022b; Chaudhuri et al., 2021) for generating informative dialogues in a domain-specific setting, where such

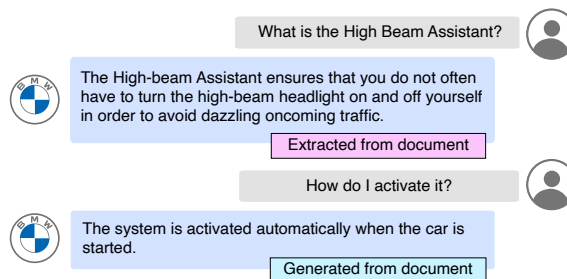


Figure 1: Illustration of a multi-turn in-car conversation between a user (in gray) and CarExpert (in blue).

systems typically engage in a multi-turn interaction with a user in form of speech or text. Figure 1 demonstrates a conversation between a user and a conversational question answering system (CarExpert) in a BMW car.

Leveraging LLMs end-to-end has several drawbacks (Liang et al., 2022; Srivastava et al., 2023; OpenAI, 2023). **Firstly**, the generated answer is often hallucinated as the knowledge from the pre-trained weights of LLMs is limited to their training data collection time (Ji et al., 2022). Furthermore, retrieval-augmented answer generation suffers from hallucination as well, due to wrong retrieval, complexity of the user utterance and retrieved document. **Secondly**, LLMs can be exploited using adversarial instructions that may lead the system to ingest malicious input and generate unsafe output (Perez and Ribeiro, 2022; Greshake et al., 2023). In the context of a car, the aforementioned downsides imply that the answer could lead to unsafe handling of the vehicle due to a lack of instructions, preservation, warning messages, or appropriate information; or by providing erroneous or confusing information.

Addressing the aforementioned issues, in this paper we propose CarExpert, an in-car conversational question-answering system, powered by LLMs. CarExpert is a modular, language model agnostic, easy to extend and controllable conversational question-answering system developed to work on

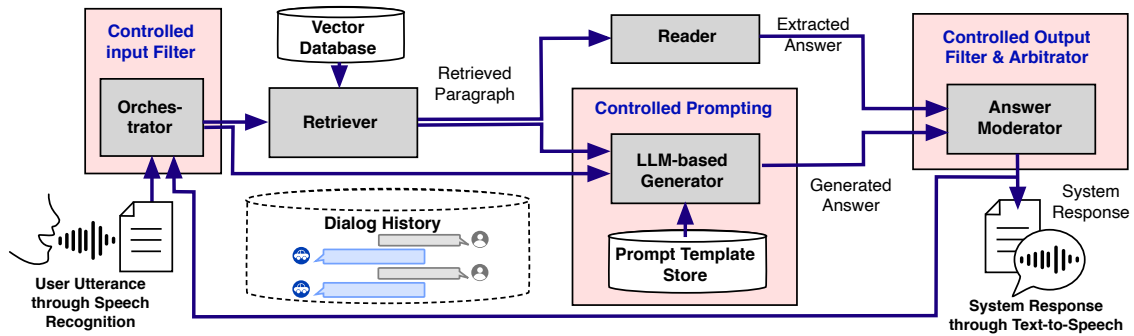


Figure 2: High level overview of the CarExpert system architecture.

the text level. On a high-level CarExpert performs question answering in two steps. First, given a user utterance it retrieves domain-specific relevant documents wherein the potential answer may exist. Second, for predicting the answer, CarExpert employs both extractive and generative answering mechanisms. Specifically, there are four sub-tasks involved in the overall process: 1) orchestration, 2) semantic search, 3) answer generation, and 4) answer moderation. Furthermore, CarExpert tackles unsafe scenarios by employing control mechanisms in three ways: i) in the *Orchestrator* using an input filter, ii) by defining prompts for controlling LLM-based answer generation, and iii) by an output filter in the *Answer Moderator*. Furthermore, CarExpert employs a heuristic during answer moderation to select answers from multiple models (extractive and generative) and provide the user with the potential best answer as the output. To facilitate voice-based user interaction in the car for real-life use, we encapsulate CarExpert with text-to-speech and speech-to-text services. Figure 2 depicts a high-level overview of the CarExpert architecture. Such modular design of CarExpert allows flexible integration to various types of interfaces such as web browser and mobile app (i.e., BMW App).

To assess the performance of CarExpert we conduct exhaustive evaluations (both qualitative and quantitative). An empirical evaluation exhibits that CarExpert outperforms off-the-shelf state-of-the-art LLMs in in car question answering. The contribution of this paper can be summarized as follows:

- We introduce CarExpert, a modular, language model agnostic, safe and controllable in-car conversational question answering system.
- A novel answer moderation heuristic for selecting a potential best answer from multiple possible outputs.

- A comprehensive empirical evaluation, demonstrating the effectiveness of CarExpert over the state-of-the-art LLMs for in-car conversational question answering.

2 Approach

CarExpert aims to generate domain-specific document-grounded answers. The task is divided into four sub-tasks: 1) Orchestration, 2) Semantic Search, 3) Answer Generation, and 4) Answer Moderation. We describe the sub-tasks below.

2.1 Orchestration

A prompt-based *Orchestrator* component is incorporated in CarExpert to tackle unsafe content and deal with multi-turn scenarios. Depending on the user utterance, CarExpert also can e.g. respond by saying that it does not have enough information or ask a clarification question, since the system is designed to only answer questions about the car. Thus the *Orchestrator* controls the input in CarExpert. The prompt used for this purpose is as follows:

Task: Given a question and paragraphs:

1. *For unsafe or harmful questions, politely decline to answer as they are out of context. Stop any further generation.*
2. *Flag any unsafe or harmful questions by politely stating that you cannot provide an answer. Stop any further generation.*
3. *If the question is safe and relevant, suggest a clarification question that demonstrates comprehension of the concept and incorporates information from the provided paragraphs. Start the question with "Do you mean".*
4. *If unsure about suggesting a specific clarification question, politely request more information to provide an accurate response. Stop any further generation.*

Question: {user utterance} *Paragraphs:* {paragraphs} *Answer:*

where, **user utterance** represent the current turn’s user utterance and **paragraphs** the top-3 retrieved documents obtained from the semantic search (discussed in Section §2.2).

2.2 Semantic Search

For efficient and fast semantic search of the relevant documents, CarExpert pre-processes data and parses clean contents from various curated sources (owners’ manuals, self-service FAQs, car configurator feature descriptions and press club publications) utilizing a data pipeline (more details in the Appendix A.1.1). The parsed data is utilized in two different ways. Firstly, we put humans in the loop to obtain high quality and domain expert annotated question-answer pairs for training an answer extraction model (discussed in Section 2.3.1). Secondly, the vector representation of the text is indexed only once as a pre-processing step to facilitate fast *Semantic Search* over a large set of text during the inference (see Figure 3). In the next step LLMs are fed with top-3 retrieved document for the answer generation. We use the terms ‘document’ and ‘paragraph’ interchangeably throughout this paper.

2.3 Answer Generation

CarExpert employs both extractive and generative models to get answers for the same user utterance. The answer generation step is controlled by instructing the LLM using prompts and next by an *Answer Moderator* component. It selects the best answer based on an extraction ratio-based heuristic (discussed in Section 2.4). We describe the answer generation methods in the following sections.

2.3.1 LLM-based Answer generation

In this step, CarExpert takes off-the-shelf GPT-3.5-turbo and instructs it in a few-shot manner for answer generation based on the current user utterance, retrieved documents and the dialogue history. The probability distribution of generating a response can be formally defined as:

$$p(S_t|\mathcal{P}; \mathcal{H}; \mathcal{Q}) = \prod_{i=1}^n p(s_i|s_{<i}, \mathcal{P}; \mathcal{H}; \mathcal{Q}, \theta), \quad (1)$$

where S_t is the generated answer, \mathcal{P} is the prompt, \mathcal{H} is the dialogue history, \mathcal{Q} is the user utterance in the current turn, θ is model parameters, and n is the length of the response. Here, “;” indicates a concatenation operation between two

texts. Depending on the type of questions that the user may ask, the generation task is split into two major categories: 1) Abstractive Summarization and 2) Informal Talk. We design separate prompt templates for both the categories to handle various types of user utterances. We provide a brief description of both the categories below.

i. Abstractive Summarization: We design a prompt template to handle information seeking user utterances that can be answered from the semantic search results where the template aims to generate the answer in a natural sentence. The abstractive summarization template is as follows:

Task: Answer questions about the car given the following context and dialog. Answer always helpful. Answer in complete sentences. Don’t use more than two sentences. Extract the answer always from the context as literally as possible.

Dialogue 1:{example dialogue 1 }

Dialogue 6: Context: {top paragraphs , dialogue history} *User:*{user utterance} *System:*

where **example dialogue 1** is a variable that represents a complete multi-turn conversation. Each dialogue may contain 1 to 5 user-system utterance pairs. The variables **top paragraphs** and **dialogue history** represent top-3 paragraphs from the semantic search results and the complete dialogue history such as adjacent user-system pairs, respectively. Furthermore, **user utterance** indicates the current user utterance that the system needs to answer.

ii. Informal Talk: A conversational AI system not only deals with information-seeking utterances but also needs to tackle follow-up questions, clarifications, commands, etc. which makes the conversation engaging and natural. To tackle various forms of user utterances we design an *Informal Talk* template as follows:

Task: Answer the user feedback in a friendly and positive way. When asked about factual knowledge or about your opinion, just say that you can’t answer these questions. Please never answer a question with a factual statement. If a question is about something else than the car, you may append a ‘Please ask me something about the car’.

Dialogue 1:{example dialogue 1 }

Dialogue 20: User:{user utterance} *System:*

In the *Informal Talk* template we provide 20 example dialogues covering various forms of user utterance. This way both abstract summarization and informal talk templates leverages pre-trained

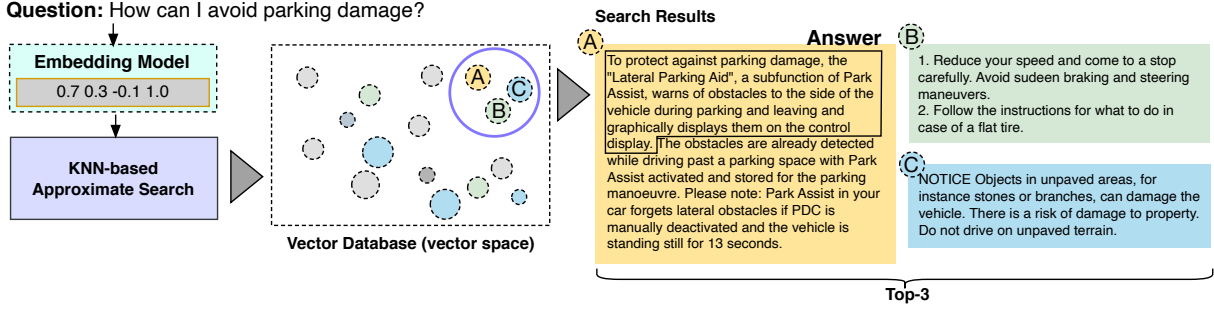


Figure 3: Semantic search during the inference (the vector space is depicted as a vector database for demonstration). The potential answer to the question is encapsulated in the box of retrieved document A.

large language model in a few-shot manner to generate natural and engaging dialogues. The prompt templates are stored in the *Prompt Template Store*.

2.3.2 Answer extraction

In CarExpert, we investigate two different answer extraction methods:

i. Machine Reading Comprehension Reader: Given a user utterance and a document the task of a MRC Reader model is to predict a continuous text span from the provided document that answers the user question. We fine-tune an Albert (Lan et al., 2020) model for the answer extraction task.

ii. LLM-based Reader: Engineering prompts is a popular way to instruct LLMs how to leverage their knowledge to solve downstream NLP tasks. In this approach, we leverage the pre-trained knowledge of LLMs, contained in their parameters to perform the same answer extraction task as the MRC Reader. However, in this case CarExpert does not need training data to perform the answer extraction. Specifically, in CarExpert we design a prompt that instructs the LLMs to perform answer extraction as literally as possible using both question and top-3 paragraphs from the semantic search results. The prompt template is as follows:

Task: Given the following question and paragraphs, extract exactly one continuous answer span from only one of the paragraphs.

Question: {user utterance} Paragraphs: {paragraphs} Answer:

During the inference, the variables **user utterance** and **paragraphs** are replaced with the actual user utterance and top three paragraphs retrieved from the semantic search.

2.4 Answer Moderation

An *Answer Moderator* component selects the best answer given the user utterance and potential answers (extractive and generative). We investigate the following two moderation techniques for answer moderation.

i. Cosine Similarity: This approach measures the semantic similarity between a user utterance and system response. The answer with a higher similarity score is selected as the system response. Formally, in this approach the answer selection can be defined as: $\max(\cosine(a_{ex}, \vec{Q}), \cosine(a_g, \vec{Q}))$, where a_{ex} , a_g , and \vec{Q} are the embedding representation of extracted answer, generated answer and user utterance.

ii. Extraction Score: This is a weighted Levenshtein distance-based heuristic that measures how syntactically close the system response is to the retrieved paragraphs. Formally, the Extraction Score (ES) can be defined as:

$$ES = \frac{1}{n} * \sum_{i=1}^n 1 - \frac{dist(x, y_i)}{\max(|x|, |y_i|)}, \quad (2)$$

where x is the generated answer, y_i is the i th paragraph and n is the number of paragraphs. The cost of edit operation is computed by $dist(\cdot)$. This moderation technique allows CarExpert to generate a controlled and document grounded answer by (i) grounding the system response to the retrieved documents, and (ii) filtering out incorrect and hallucinated responses. More details on the edit operations can be found in Appendix A.5.

3 Experimental Setup

Data: The reader and retriever models in CarExpert are fine-tuned and evaluated on car-specific

data from various sources (owners' manuals, self-service FAQs, car configurator feature descriptions and press club publications).

Baselines: We choose Dense Passage Retriever (DPR) (Karpukhin et al., 2020a), BM25 (Robertson et al., 2009), Sentence-transformer (Reimers and Gurevych, 2019) and SPLADE (Formal et al., 2022) as the baseline retriever. For answer generation we experiment with Albert (Lan et al., 2020) (extractive) and GPT-3.5-turbo¹ (generative) and Luminous-extended² (generative).

Metrics To measure the performance of the *Retriever* we use Mean Reciprocal Rank (MRR@3). For evaluating extractive *Reader*, we utilize token-level metrics, such as F1-Score and Exact Match (EM). Furthermore, we employ Cosine Similarity and METEOR (Banerjee and Lavie, 2005) to capture the similarity of generated answer against the reference response.

Further details of the datasets, hyper-parameter settings, and metrics can be found in the Appendix, in A.1, A.3 and A.4 respectively.

4 Experiments and Results

We conduct both qualitative and quantitative experiments to assess different parts contributing to the overall performance of CarExpert.

4.1 Quantitative Analysis

Table 2 and Table 3 demonstrate that the fine-tuned DPR and fine-tuned Reader perform better than the baseline models in the corresponding tasks. The performance improvement may attributed to their inherent capability of effectively learning and capturing the distribution and characteristics of the training data. In Table 2, we notice that a fine-tuned DPR outperforms a fine-tuned Sentence-transformer. The fine-tuned DPR model preforms in MRR@1 and hence we integrate DPR as the retriever used for semantic search in CarExpert.

From Table 4 we observe that GPT-3.5-turbo performs better than the Luminous-extended model since the former is a larger model and hence offers better representations and generalization.

Table 5 exhibits that *Extraction Score* does a better job in moderating and selecting the best answer which aligns better to the retrieved documents. CarExpert incorporate the *Extraction Score*-based

heuristic for answer moderation. The *Extraction Score* technique is described in Appendix A.5.

4.2 Qualitative Analysis

Table 1 demonstrates a qualitative comparison between CarExpert (with document) and GPT-3.5-turbo (with and without document) of answer generation. When provided with the document we instruct both the models to answer from the provided documents. In the first case, without any documents provided GPT-3.5-turbo could not answer the question, where with the document it generated a very long answer. Furthermore, when answering it is referring to a specific paragraph such as ". . .The first paragraph mentions. . .", which is irrelevant to the user. CarExpert in this case correctly generated the expected answer. In the second case, we asked the system about how to mount a child seat. Off-the-shelf GPT-3.5-turbo generated generic answer from its pre-trained knowledge, which includes unnecessary detail such as ". . .Read the instruction. . .", and is not tailored to the target car brand. On the contrary, although GPT-3.5-turbo generated a better answer, it includes additional irrelevant and lengthy details which are not suitable for in-car CQA. Still adding irrelevant information (right column, 3rd row: item 1., 6. & 7.). Overall, in both the cases, CarExpert exhibits precise answer prediction then off-the-shelf GPT-3.5-turbo with and without documents. Although, CarExpert leverages GPT-3.5-turbo for the answer generation, carefully designed prompts in CarExpert helped the system to generate precise answers. Precise answers are suitable for real-time use in the car, where the user may find an unnecessary detailed answer (which GPT-3.5-turbo generated) very exhausting. More lemon- and cherry-picked examples can be found in Appendix D.

5 Discussions and Potential Impact

CarExpert is built in a modular fashion, which allows for expansion and adaptability to diverse industrial use cases. Furthermore, the proposed architecture enables the system to maintain, modify and scale the data more effectively. Moreover, a pipeline approach such as CarExpert improves the overall interpretability and debugging of a system. Finally, the introduced system is controllable and domain-specific as it allows for explicit control over the design and behavior of each of the

¹<https://openai.com/>

²<https://www.aleph-alpha.com/>

User	CarExpert (with paragraphs)	GPT-3.5-turbo (w/o paragraphs)	GPT-3.5-turbo (with paragraphs)
Does the car have a parking assist function?	Yes, Park Assist supports the following functions: When parallel parking: Parking in reverse, parallel to the road. When perpendicular parking: Parking forward or in reverse, perpendicular to the road. When driving out of parallel parking spaces. When driving out of perpendicular parking spaces.	I'm sorry, but I would need more specific information about the car you're referring to. There are many different car models with various features, and the availability of parking assist functions can vary between manufacturers and models. Can you provide me with the make and model of the car you're interested in?	Yes, the car has a parking assist function. The first paragraph mentions that Park Assist supports functions like parallel parking, perpendicular parking, and driving out of both types of parking spaces. The second paragraph further elaborates on the features of the Automatic Parking Assistant, which includes [...]
How to mount child seats?	To install the child restraint system in the vehicle, proceed as follows: 1. Mount the child restraint system in accordance with the instructions from the child seat manufacturer. 2. Make sure that the child restraint system anchors engage correctly on both sides in the lower seat mounts. 3. Pull out the seat belt strap completely. 4. Secure the child restraint system with the seat belt. 5. Allow the seat belt strap to be pulled in, then pull it tight against the child restraint system.	Mounting child seats properly is crucial for the safety of your child while traveling in a vehicle. The specific instructions for mounting a child seat may vary depending on the make and model of the seat and the vehicle. However, here is a general guide on how to mount child seats: 1. Read the instructions: Carefully read the instruction manual provided with your child seat. It will provide specific details about the installation process and any additional requirements. 2. Choose the appropriate seat: [...]	To mount child seats, follow these steps: 1. Read and follow the instructions provided by the child seat manufacturer. 2. Ensure that the child restraint system anchors engage correctly on both sides in the lower seat mounts. 3. Pull out the seat belt strap compl. [...] 6. Refer strictly to the installation instructions supplied with the child restraint system for specific guidelines. 7. Each rear seat position is equipped with a head restraint. [...]

Table 1: Dialog examples with answers of CarExpert compared to GPT-3.5-turbo stand alone and GPT-3.5-turbo with relevant paragraphs provided through prompting. Longer texts are truncated and replaced by [...] for the demonstration purpose.

Retriever	MRR@1	MRR@3
BM25	0.26	0.31
Fine-tuned DPR	0.47	0.52
Fine-tuned Sentence-transformer	0.42	0.49
SPLADE	0.44	0.53

Table 2: Performance comparison of retriever models.

Reader	F1	EM
Pre-trained Albert-large	0.31	0.01
Fine-tuned Albert-large	0.60	0.21
GPT-3.5-turbo	0.51	0.14
Luminous-extended	0.36	0.01

Table 3: Evaluation results on the module: *Reader*.

modules such as *Orchestrator* and answer generation. We anticipate that CarExpert will aid other industrial use cases leverage LLMs in developing fine-grained and regulated conversational question answering systems.

6 Related Works

Large Language Models: Large language model (LLM) such as GPT-3 (Brown et al., 2020), PaLM (Chowdhery et al., 2022), LaMDA (Thoppi-

Generator	Cos. Sim.	METEOR
GPT-3.5-turbo	0.68	0.38
Luminous-extended	0.52	0.14

Table 4: Performance of *LLM-based Generator* models.

Answer Moderator	Accuracy
Cosine Similarity	0.82
Extraction Score	0.87

Table 5: Performance of *Answer Moderator* approaches.

lan et al., 2022a), LLaMA (Touvron et al., 2023a) and GPT-4 (OpenAI, 2023) are capable of performing complex downstream tasks without being trained for that tasks. A different line of recent research focuses on controlling the behaviour of LLMs such as NeMo-Guardrails³. Inspired by humans capabilities of following instructions in natural language, recent research works fine-tuned LLMs so that it can understand instructions in a zero-shot or few-shot settings and perform a given task following the language instruction (Wei et al., 2022; Taori et al., 2023; Brown et al., 2020; Rony et al., 2022a; Schick and Schütze, 2021; Prasad

³ <https://developer.nvidia.com/nemo>

et al., 2023). In CarExpert, prompt-guided LLMs are employed to control various tasks of the answer generation process.

Conversational Question Answering: Recent advancements of LLMs significantly improved multi-turn question answering systems (Chowdhery et al., 2022; Thoppilan et al., 2022b; Zaib et al., 2021). However, in multi-task objectives these models lack robustness (Liang et al., 2022; Srivastava et al., 2023). A different line of work (Daull et al., 2023) emphasised on the needs for hybrid approaches to take advantage of multiple learning models to better handle the limitations. Architectural compositions such as LLM + semantic information retrieval (de Jong et al., 2023; Borgeaud et al., 2022), LLM + instruction tuning module (Khattab et al., 2022), LLM + Router (Xu et al., 2023), cascaded LLMs (Dohan et al., 2022), LLM + RLHF/RLAIF (Ouyang et al., 2022; Bai et al., 2022). Despite significant progress over time, CQA systems still struggle with long-standing issues like hallucination, the ability to scale models and data, and formal reasoning.

7 Conclusion

We have introduced CarExpert, a new and controlled in-car conversational question-answering system powered by LLMs. Specifically, CarExpert employed semantic search to restrict the system generated answer within the car domain and incorporated LLMs to predict natural, controlled and safe answers. Furthermore, to tackle hallucinated answers, CarExpert proposed an Extraction Score-based Answer Moderator. We anticipate that the proposed approach can not only be applicable for the in-car question answering but also be easily extendable and adapted for other domain-specific settings. In future, we plan to integrate multi-task models to handle multiple task using a single LLM and reduce error propagation in the system.

Limitations

While our modular framework offers considerable flexibility in employing diverse models and aligning them with specific tasks and objectives, it comes with few challenges as well. One major drawback is the difficulty in jointly optimizing and fine-tuning the individual components toward a common objective. When optimized independently, each module may overfit to certain tasks and subsequently propagate errors due to intricate inter-

actions, ultimately impacting the overall system performance. Furthermore, given our reliance on LLMs, occasional hallucinations may occur despite our efforts to maintain control. Moreover, our system may struggle with handling highly complex and ambiguous queries, potentially requiring external resolution modules. In future, we intend to tackle the existing issues to develop a more robust conversational question answering system.

Acknowledgement

We would like to thank Dr. Hans-Joerg Voegel, Dr. Robert Bruckmeier, and Dr. Peter Lehnert from the BMW Group in Munich, Germany for their support in this work. We would like to extend our thank to Dr. Nicolas Flores-Herr, Dr. Joachim Koehler, Alexander Arno Weber and the Fraunhofer IAIS team for the helpful discussions and contributions to this work, and the members who contributed to this project from BIG PICTURE GmbH and ONSEI GmbH.

References

- Ben Athiwaratkun, Andrew Wilson, and Anima Anandkumar. 2018. [Probabilistic FastText for multi-sense word embeddings](#). In *Proc. of ACL*, pages 1–11, Melbourne, Australia. Association for Computational Linguistics.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, John Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, E Perez, Jamie Kerr, Jared Mueller, Jeff Ladish, J Landau, Kamal Ndousse, et al. 2022. [Constitutional ai: Harmlessness from ai feedback](#). *ArXiv preprint*, abs/2212.08073.
- Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72.
- Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George van den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, Diego de Las Casas, Aurelia Guy, Jacob Menick, Roman Ring, Tom Hennigan, Saffron Huang, Loren Maggiore, Chris Jones, Albin Cassirer, Andy Brock, Michela Paganini, Geoffrey Irving, Oriol Vinyals, Simon Osindero, Karen Simonyan, Jack W. Rae, Erich Elsen, and Laurent Sifre. 2022. [Improving language models by retrieving from](#)

- trillions of tokens. In *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pages 2206–2240. PMLR.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared 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 M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher 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 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Debanjan Chaudhuri, Md Rashad Al Hasan Rony, and Jens Lehmann. 2021. Grounding dialogue systems via knowledge graph aware decoding with pre-trained transformers. In *The Semantic Web: 18th International Conference, ESWC 2021, Virtual Event, June 6–10, 2021, Proceedings 18*, pages 323–339. Springer.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, and Hyung Won Chung. 2022. [Palm: Scaling language modeling with pathways](#). *ArXiv preprint*, abs/2204.02311.
- Xavier Daull, Patrice Bellot, Emmanuel Bruno, Vincent Martin, and Elisabeth Murisasco. 2023. [Complex qa and language models hybrid architectures, survey](#). *ArXiv preprint*, abs/2302.09051.
- Michiel de Jong, Yury Zemlyanskiy, Joshua Ainslie, Nicholas FitzGerald, Sumit Sanghai, Fei Sha, and William Cohen. 2023. [FiDO: Fusion-in-decoder optimized for stronger performance and faster inference](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 11534–11547, Toronto, Canada. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proc. of NAACL-HLT*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- David Dohan, Winnie Xu, Aitor Lewkowycz, Jacob Austin, David Bieber, Raphael Gontijo Lopes, Yuhuai Wu, Henryk Michalewski, Rif A. Saurous, Jascha Narain Sohl-Dickstein, Kevin Murphy, and Charles Sutton. 2022. [Language model cascades](#). *ArXiv preprint*, abs/2207.10342.
- Thibault Formal, Carlos Lassance, Benjamin Piwowarski, and Stéphane Clinchant. 2022. [From distillation to hard negative sampling: Making sparse neural ir models more effective](#).
- Thibault Formal, Benjamin Piwowarski, and Stéphane Clinchant. 2021. [Splade: Sparse lexical and expansion model for first stage ranking](#).
- Kai Greshake, Sahar Abdelnabi, Shailesh Mishra, Christoph Endres, Thorsten Holz, and Mario Fritz. 2023. [More than you’ve asked for: A comprehensive analysis of novel prompt injection threats to application-integrated large language models](#). *CoRR*, abs/2302.12173.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Wenliang Dai, Andrea Madotto, and Pascale Fung. 2022. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55:1 – 38.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020a. [Dense passage retrieval for open-domain question answering](#). In *Proc. of EMNLP*, pages 6769–6781, Online. Association for Computational Linguistics.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020b. [Dense passage retrieval for open-domain question answering](#). In *Proc. of EMNLP*, pages 6769–6781, Online. Association for Computational Linguistics.
- O. Khattab, Keshav Santhanam, Xiang Lisa Li, David Leo Wright Hall, Percy Liang, Christopher Potts, and Matei A. Zaharia. 2022. [Demonstrate-search-predict: Composing retrieval and language models for knowledge-intensive nlp](#). *ArXiv preprint*, abs/2212.14024.
- Boeun Kim, Dohaeng Lee, Sihyung Kim, Yejin Lee, Jin-Xia Huang, Oh-Woog Kwon, and Harksoo Kim. 2021. [Document-grounded goal-oriented dialogue systems on pre-trained language model with diverse input representation](#). In *Proceedings of the 1st Workshop on Document-grounded Dialogue and Conversational Question Answering (DialDoc 2021)*, pages 98–102, Online. Association for Computational Linguistics.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2020. [ALBERT: A lite BERT for self-supervised learning of language representations](#). In *Proc. of ICLR*. OpenReview.net.
- Jey Han Lau and Timothy Baldwin. 2016. [An empirical evaluation of doc2vec with practical insights into document embedding generation](#). In *Proceedings of the 1st Workshop on Representation Learning for NLP*, pages 78–86, Berlin, Germany. Association for Computational Linguistics.
- Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby Yan,

- Ce Zhang, Christian Cosgrove, Christopher D. Manning, and Yuta Koreeda. 2022. Holistic evaluation of language models. *Annals of the New York Academy of Sciences*.
- OpenAI. 2023. [Gpt-4 technical report](#).
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Gray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. [Training language models to follow instructions with human feedback](#). In *Advances in Neural Information Processing Systems*.
- Fábio Perez and Ian Ribeiro. 2022. Ignore previous prompt: Attack techniques for language models. *ArXiv*, abs/2211.09527.
- Archiki Prasad, Peter Hase, Xiang Zhou, and Mohit Bansal. 2023. [GrIPS: Gradient-free, edit-based instruction search for prompting large language models](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 3845–3864, Dubrovnik, Croatia. Association for Computational Linguistics.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. [SQuAD: 100,000+ questions for machine comprehension of text](#). In *Proc. of EMNLP*, pages 2383–2392, Austin, Texas. Association for Computational Linguistics.
- Nils Reimers and Iryna Gurevych. 2019. [Sentence-BERT: Sentence embeddings using Siamese BERT-networks](#). In *Proc. of EMNLP*, pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.
- Stephen Robertson, Hugo Zaragoza, et al. 2009. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends® in Information Retrieval*, 3(4):333–389.
- Joshua Robinson and David Wingate. 2023. [Leveraging large language models for multiple choice question answering](#). In *The Eleventh International Conference on Learning Representations*.
- Md Rashad Al Hasan Rony, Debanjan Chaudhuri, Ricardo Usbeck, and Jens Lehmann. 2022a. [Tree-kgqa: An unsupervised approach for question answering over knowledge graphs](#). *IEEE Access*, 10:50467–50478.
- Md Rashad Al Hasan Rony, Ricardo Usbeck, and Jens Lehmann. 2022b. [DialoKG: Knowledge-structure aware task-oriented dialogue generation](#). In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 2557–2571, Seattle, United States. Association for Computational Linguistics.
- Md Rashad Al Hasan Rony, Ying Zuo, Liubov Kovrigina, Roman Teucher, and Jens Lehmann. 2022c. [Climate bot: A machine reading comprehension system for climate change question answering](#). In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, pages 5249–5252. International Joint Conferences on Artificial Intelligence Organization. AI for Good - Demos.
- Timo Schick and Hinrich Schütze. 2021. [Few-shot text generation with natural language instructions](#). In *Proc. of EMNLP*, pages 390–402, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adria Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda Askell, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S. Iyer, Anders Johan Andreassen, Andrea Madotto, Andrea Santilli, Andreas Stuhlmüller, Andrew M. Dai, Andrew La, Andrew Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabassum, Arul Menezes, Arun Kirubaranjan, Asher Mullokandov, Ashish Sabharwal, Austin Herrick, Avia Efrat, Aykut Erdem, Ayla Karakaş, B. Ryan Roberts, Bao Sheng Loe, Barret Zoph, Bartłomiej Bojanowski, Batuhan Özyurt, Behnam Hedayatnia, Behnam Neyshabur, Benjamin Inden, Benno Stein, Berk Ekmekci, Bill Yuchen Lin, Blake Howald, Bryan Orinon, Cameron Diao, Cameron Dour, Catherine Stinson, Cedrick Argueta, Cesar Ferri, Chandan Singh, Charles Rathkopf, Chenlin Meng, Chitta Baral, Chiyu Wu, Chris Callison-Burch, Christopher Waites, Christian Voigt, Christopher D Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera, Clemencia Siro, Colin Raffel, Courtney Ashcraft, Cristina Garbacea, Damien Sileo, Dan Garrette, Dan Hendrycks, Dan Kilman, Dan Roth, C. Daniel Freeman, Daniel Khachabi, Daniel Levy, Daniel Moseguí González, Danielle Perszyk, Danny Hernandez, Danqi Chen, Daphne Ippolito, Dar Gilboa, David Dohan, David Drakard, David Jurgen, Debajyoti Datta, Deep Ganguli, Denis Emelin, Denis Kleyko, Deniz Yuret, Derek Chen, Derek Tam, Dieuwke Hupkes, Diganta Misra, Dilyar Buzan, Dimitri Coelho Mollo, Diyi Yang, Dong-Ho Lee, Dylan Schrader, Ekaterina Shutova, Ekin Dogus Cubuk, Elad Segal, Eleanor Hagerman, Elizabeth Barnes, Elizabeth Donoway, Ellie Pavlick, Emanuele Rodolà, Emma Lam, Eric Chu, Eric Tang, Erkut Erdem, Ernie Chang, Ethan A Chi, Ethan Dyer, Ethan Jerzak, Ethan Kim, Eunice Engefu Manyasi, Evgenii Zheltonozhskii, Fanyue Xia, Fatemeh Siar, Fernando Martínez-Plumed, Francesca Happé, Francois Chollet, Frieda Rong, Gaurav Mishra, Genta In-dra Winata, Gerard de Melo, Germán Kruszewski,

- Giambattista Parascandolo, Giorgio Mariani, Gloria Xinyue Wang, Gonzalo Jaimovitch-Lopez, Gregor Betz, Guy Gur-Ari, Hana Galijasevic, Hannah Kim, Hannah Rashkin, Hannaneh Hajishirzi, Harsh Mehta, Hayden Bogar, Henry Francis Anthony Shevlin, Hinrich Schuetze, Hiromu Yakura, Hongming Zhang, Hugh Mee Wong, Ian Ng, Isaac Noble, Jaap Jumelet, Jack Geissinger, Jackson Kernion, Jacob Hilton, Jaehoon Lee, Jaime Fernández Fisac, James B Simon, James Koppel, James Zheng, James Zou, Jan Kocon, Jana Thompson, Janelle Wingfield, Jared Kaplan, Jarema Radom, Jascha Sohl-Dickstein, Jason Phang, Jason Wei, Jason Yosinski, Jekaterina Novikova, Jelle Bosscher, Jennifer Marsh, Jeremy Kim, Jeroen Taal, Jesse Engel, Jesujoba Alabi, Jiacheng Xu, Jiaming Song, Jillian Tang, Joan Waweru, John Burden, John Miller, John U. Balis, Jonathan Batchelder, Jonathan Berant, Jörg Frohberg, Jos Rozen, Jose Hernandez-Orallo, Joseph Boudeman, Joseph Guerr, Joseph Jones, Joshua B. Tenenbaum, Joshua S. Rule, Joyce Chua, Kamil Kanclerz, Karen Livescu, Karl Krauth, Karthik Gopalakrishnan, Katerina Ignatyeva, Katja Markert, Kaustubh Dhole, Kevin Gimpel, Kevin Omondi, Kory Wallace Mathewson, Kristen Chiafullo, Ksenia Shkaruta, Kumar Shridhar, Kyle McDonell, Kyle Richardson, Laria Reynolds, Leo Gao, Li Zhang, Liam Dugan, Lianhui Qin, Lidia Contreras-Ochando, Louis-Philippe Morency, Luca Moschella, Lucas Lam, Lucy Noble, Ludwig Schmidt, Luheng He, Luis Oliveros-Colón, Luke Metz, Lütfi Kerem Senel, Maarten Bosma, Maarten Sap, Maartje Ter Hoeve, Maheen Farooqi, Manaal Faruqui, Mantas Mazeika, Marco Baturan, Marco Marelli, Marco Maru, Maria Jose Ramirez-Quintana, Marie Tolkiehn, Mario Giulianelli, Martha Lewis, Martin Potthast, Matthew L Leavitt, Matthias Hagen, Mátyás Schubert, Medina Orduna Baitemirova, Melody Arnaud, Melvin McElrath, Michael Andrew Yee, Michael Cohen, Michael Gu, Michael Ivanitskiy, Michael Starritt, Michael Strube, Michał Śwędrowski, Michele Bevilacqua, Michihiro Yasunaga, Mihir Kale, Mike Cain, Mimeo Xu, Mirac Suzgun, Mitch Walker, Mo Tiwari, Mohit Bansal, Moin Amnaseri, Mor Geva, Mozhdheh Gheini, Mukund Varma T, Nanyun Peng, Nathan Andrew Chi, Nayeon Lee, Neta Gur-Ari Krakover, Nicholas Cameron, Nicholas Roberts, Nick Doiron, Nicole Martinez, Nikita Nangia, Niklas Deckers, Niklas Muennighoff, Nitish Shirish Keskar, Niveditha S. Iyer, Noah Constant, Noah Fiedel, Nuan Wen, Oliver Zhang, Omar Agha, Omar Elbaghdadi, Omer Levy, Owain Evans, Pablo Antonio Moreno Casares, Parth Doshi, Pascale Fung, Paul Pu Liang, Paul Vicol, Pegah Alipoormolabashi, Peiyuan Liao, Percy Liang, Peter W Chang, Peter Eckersley, Phu Mon Htut, Pinyu Hwang, Piotr Miłkowski, Piyush Patil, Pouya Pezeshkpour, Priti Oli, Qiaozhu Mei, Qing Lyu, Qinlang Chen, Rabin Banjade, Rachel Etta Rudolph, Raefer Gabriel, Rahel Habacker, Ramon Risco, Raphaël Millière, Rhythm Garg, Richard Barnes, Rif A. Saurous, Riku Arakawa, Robbe Raymaekers, Robert Frank, Rohan Sikand, Roman Novak, Roman Sitelew, Roman Le Bras, Rosanne Liu, Rowan Jacobs, Rui Zhang, Russ Salakhutdinov, Ryan Andrew Chi, Seungjae Ryan Lee, Ryan Stovall, Ryan Teehan, Rylan Yang, Sahib Singh, Saif M. Mohammad, Sajant Anand, Sam Dillavou, Sam Shleifer, Sam Wiseman, Samuel Gruetter, Samuel R. Bowman, Samuel Stern Schoenholz, Sanghyun Han, Sanjeev Kwatra, Sarah A. Rous, Sarik Ghazarian, Sayan Ghosh, Sean Casey, Sebastian Bischoff, Sebastian Gehrmann, Sebastian Schuster, Sepideh Sadeghi, Shadi Hamdan, Sharon Zhou, Shashank Srivastava, Sherry Shi, Shikhar Singh, Shima Asaadi, Shixiang Shane Gu, Shubh Pachchigar, Shubham Toshniwal, Shyam Upadhyay, Shyamolima Shammie Debnath, Siamak Shakeri, Simon Thormeyer, Simone Melzi, Siva Reddy, Sneha Priscilla Makini, Soo-Hwan Lee, Spencer Torene, Sriharsha Hatwar, Stanislas Dehaene, Stefan Divic, Stefano Ermon, Stella Biderman, Stephanie Lin, Stephen Prasad, Steven Piantadosi, Stuart Shieber, Summer Mishnerghi, Svetlana Kiritchenko, Swaroop Mishra, Tal Linzen, Tal Schuster, Tao Li, Tao Yu, Tariq Ali, Tatsunori Hashimoto, Te-Lin Wu, Théo Desbordes, Theodore Rothschild, Thomas Phan, Tianle Wang, Tiberius Nkinyili, Timo Schick, Timofei Kornev, Titus Tunduny, Tobias Gerstenberg, Trenton Chang, Trishala Neeraj, Tushar Khot, Tyler Shultz, Uri Shaham, Vedant Misra, Vera Demberg, Victoria Nyamai, Vikas Raunak, Vinay Venkatesh Ramasesh, vinay uday prabhu, Vishakh Padmakumar, Vivek Srikumar, William Fedus, William Saunders, William Zhang, Wout Vossen, Xiang Ren, Xiaoyu Tong, Xinran Zhao, Xinyi Wu, Xudong Shen, Yadollah Yaghoobzadeh, Yair Lakretz, Yangqiu Song, Yasaman Bahri, Yejin Choi, Yichi Yang, Yiding Hao, Yifu Chen, Yonatan Belinkov, Yu Hou, Yufang Hou, Yuntao Bai, Zachary Seid, Zhuoye Zhao, Zijian Wang, Zijie J. Wang, Zirui Wang, and Ziyi Wu. 2023. [Beyond the imitation game: Quantifying and extrapolating the capabilities of language models.](#) *Transactions on Machine Learning Research*.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. [Stanford alpaca: An instruction-following llama model.](#) https://github.com/tatsu-lab/stanford_alpaca.
- Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al. 2022a. [Lamda: Language models for dialog applications.](#) *ArXiv preprint*, abs/2201.08239.
- Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam M. Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, and Quoc Le. 2022b. [Lamda: Language models for dialog applications.](#) *ArXiv preprint*, abs/2201.08239.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. [Llama: Open and efficient foundation language models.](#)

Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Daniel M. Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony S. Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel M. Kloumann, A. V. Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashii Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Arulien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. Llama 2: Open foundation and fine-tuned chat models.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008.

Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. 2022. [Finetuned language models are zero-shot learners](#). In *Proc. of ICLR*. OpenReview.net.

Haike Xu, Zongyu Lin, Jing Zhou, Yanan Zheng, and Zhilin Yang. 2023. [A universal discriminator for zero-shot generalization](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10559–10575, Toronto, Canada. Association for Computational Linguistics.

Munazza Zaib, Wei Emma Zhang, Quan Z. Sheng, Adnan Mahmood, and Yang Zhang. 2021. Conversational question answering: a survey. *Knowledge and Information Systems*, 64:3151 – 3195.

Jianjin Zhang, Zheng Liu, Weihao Han, Shitao Xiao, Ruicheng Zheng, Yingxia Shao, Hao Sun, Hanqing Zhu, Premkumar Srinivasan, Weiwei Deng, et al. 2022. Uni-retriever: Towards learning the unified embedding based retriever in bing sponsored search. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 4493–4501.

A Appendix

A.1 Data

Sources: Information sources were comprised of the following documents:

- *Owners’ manual:* Detailed descriptions of functional features and step-by-step instructions on their usage for the target car. Including information about safe usage of the car as well as warnings to prevent unsafe situations and handling.
- *Self-service:* A collection of frequently asked questions and answers about cars and their features (language: English UK and US).
- *Car Configurator:* Description of configuring the car’s appearance and its technical features.
- *Press Club:* A collection of target car specific articles published as press publications.

Table 6 illustrates number of paragraphs and the median word count of each information source. While the owners’ manual has the highest number of relatively short paragraphs, Press Club contains longer paragraphs in smaller quantities. Figure 4 depicts the distribution of word count in one paragraph for the data sources.

Data Sources	# Para.	Median
Owners’ Manual	3,537	38
Self Service	312	70
Car Configurator	150	52
Press Club	125	228

Table 6: Overview of number of paragraph and median word count per paragraph for each source document.

Training and Evaluation Data: We constructed a set of in-house annotated data by randomly sampling paragraphs from different data sources. The train/dev/test splits are shown in table 7. The eval-

	# Queries	# Para.
In-house Train ⁺	757	278
In-house Dev ⁺	176	80
In-house Test ⁺	66	40
Evaluation*	60	40

Table 7: Data statistics for in-house data and human-annotated evaluation data. ⁺ indicates data used for training and evaluating the *Retriever* and *Reader* models. * indicates data used for evaluating the individual modules and the system as a whole.

uation set contains, 60 multi-turn dialogues (33%

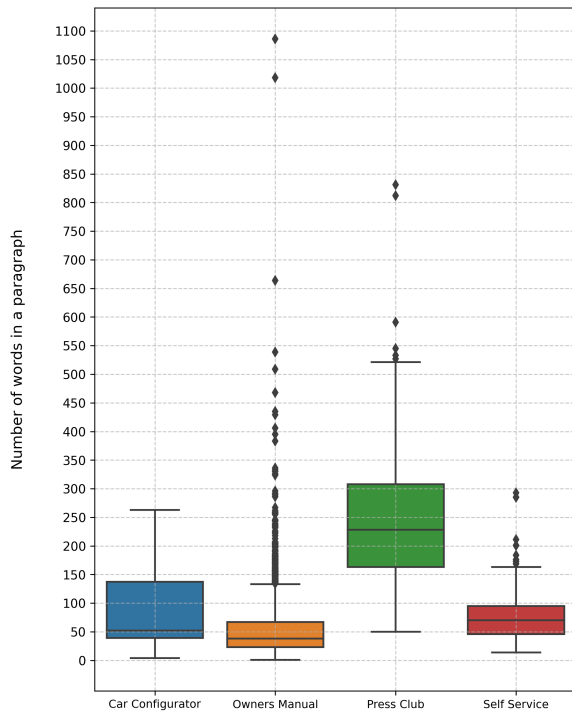


Figure 4: Word count distributions per paragraph.

with 2 turns, 33% with 4 turns and 33% with 6 turns), curated from 40 different paragraphs for randomly sampled document collection. We ensured that at least one dialog is crafted for every paragraph in this evaluation set. The human-annotation process for collecting these data are described in Section § C.

A.1.1 Data Processing Pipeline

The data processing pipeline in CarExpert takes data in various format (such as unstructured text, PDF, Excel, CSV, XML) and transforms them into SQuAD (Rajpurkar et al., 2016) format. SQuAD is a widely used question answering dataset format. The paragraphs in the SQuAD format are then converted into vectors, obtained from the Sentence-transformer and stored them in a vector database to facilitate quick semantic search (retrieval) given a user query.

A.2 Baselines

The baseline models used for comparing each components are as follows:

Retriever: (i) *Sparse embeddings:* BM25 (Robertson et al., 2009) (ii) *Static embedding models:* FastText (Athiwaratkun et al., 2018) and doc2vec (Lau and Baldwin, 2016) (iii) *Contextual embedding models:* Dense Passage

Retrieval (DPR) (Karpukhin et al., 2020b) and Sentence-transformers (Zhang et al., 2022) (iv) *Hybrid embedding models:* SPLADE (Formal et al., 2021).

Reader: (i) *Encoder-based (LM) models:* pre-trained reader models including Albert-large (Lan et al., 2020) (ii) *Decoder-based (LLM) models:* GPT-3.5-turbo and Luminous-extended.

Generator: (i) GPT-3.5-turbo (ii) Luminous-extended.

A.3 Hyper-parameter Settings

We describe the hyper-parameters used in different components of the CarExpert below.

Retriever: We fine-tune the DPR model by employing a query encoder: facebook/dpr-question_encoder-multiset-base and facebook/dpr-ctx_encoder-multiset-base as the paragraph encoder. We continued training for 10 epochs with a batch size of 8, warm-up steps of 6, and one hard negative sample per data point. We further fine-tuned the Sentence-transformer model all-MiniLM-L6-v2 with a batch size of 16 for 1 epoch, combining the objective of reducing Masked Language Modelling (MLM) and Next Sentence Prediction (NSP) loss.

Reader: As the reader model, we fine-tuned Albert-large (Lan et al., 2020) as the base model. For the LLM-based reader, we used GPT-3.5-turbo and Luminous-extended models. In both cases, we set a temperature of 0 to facilitate deterministic text generation, as well as a presence penalty of 0, top- p sampling rate of 0 and repetition penalty of 1.

Generator: For the LLM-based answer generation, we use GPT-3.5-turbo and Luminous-extended with a temperature of 0.8, top- p sampling rate 0.4, repetition penalty 1 and presence penalty of 0.6. These settings allow for a more flexible answer generation, in contrast to the LLM-based reader.

A.4 Metrics

For quantitative evaluation of the system components and the system as a whole, we relied upon the following metrics.

Retriever: (i) *Mean Reciprocal Rank (MRR)* for the top-3 paragraphs calculates the average reciprocal rank of the first relevant document across

multiple queries. The focus is on the rank of the first relevant document.

Reader: (i) *F1-Score* considers both precision (how many predicted words are correct) and recall (how many correct words are predicted). (ii) *Exact Match* (EM) measures the percentage of predicted answers that exactly match the ground truth answers. It is a strict metric that demands the model response to be identical to the ground truth.

Type of token	INS	DEL	SUB
Default	1.0	1.0	1.0
Stop words	0.5	0.5	0.5
Input tokens	0.5	-	0.1
Reference tokens	-	2.0	-

Table 8: Insertion costs (INS), Deletion costs (DEL) and Substitution costs (SUB) for different types of tokens.

Generator: (i) *Cosine Similarity* between the system response and the human annotated response. (ii) *METEOR* (Banerjee and Lavie, 2005) provides a single score reflecting the overall quality and fluency of the generated response against the human annotated response.

Answer Moderator (i) *Accuracy* of correctly yielding the extracted or the generated response as annotated by the human annotators.

System as a whole: (i) *Cosine Similarity* between the final system response and the expected system response. (ii) *Component Contributions* revealing if the system yields more extractive responses or generative results.

A.5 Answer Moderator

Edit Operations in Extraction score: Table 8 demonstrates the edit operation cost used in Extraction Score. Note that when the system *deletes* any reference token, it receives a maximum penalty. Eventually, the distance is normalized to a consistent scale using the maximum absolute value.

B Ablation Studies

B.1 Retriever

We performed an extensive ablation study on different types of retriever (sparse, static, contextual, and hybrid) on both in-house and human-annotated evaluation datasets.

The retriever scores from the traditional BM25 and the static models are significantly lower, as expected, than the rest of the candidates. We observe that our datasets are reasonably hard for the retrievers which rely upon just the frequencies or associations between query-document pairs, essentially failing to yield meaningful contextual representations. The fine-tuned DPR performs the best on the human-annotated evaluation set, while the fine-tuned Sentence-transformer model performs the best on the in-house test set. It is also worth noting that the off-the-shelf SPLADE model performs almost as good as the fine-tuned contextual models. This could be attributed to how hybrid models are trained to combine the best of both worlds from the sparse and dense representations.

B.2 System as a whole

Table 10 demonstrates the experimental results of CarExpert with various system configurations. The component-wise evaluation presented earlier in Table 2 through 5) motivated us to conduct this elaborate study, within a scope with (i) fine-tuned DPR and fine-tuned Sentence Transformer models as *Retriever*, (ii) fine-tuned *Reader* and GPT-3.5-turbo based *Reader*, (iii) GPT-3.5-turbo as the *Generator*, and (iv) both answer moderation techniques.

It is evident from the results that the *Extraction Score* based *Answer Moderator* always prefers extractive responses than the generative responses when compared to the Cosine Similarity-based counterpart. For instance, the configurations **C01** and **C03** differ only by the Answer Moderator, however there is a significant increase in the contribution of extractive responses from 23% to 52%. This moderation technique helps our model to stay controllable regardless of the nature of the user utterances. The best share of extractive responses is obtained from **C03**.

We also observe how different retriever models affect the overall system response. For instance, the configurations **C04** and **C08** differ only by the retrievers used, however, with a significant difference in the similarity between the system response and reference response. In future, we intend to explore other sophisticated metrics that measure more nuanced aspects of language generation. In addition, we hypothesize that the cosine-similarity-based system evaluation might be biased towards the cosine similarity-based arbitration method as they may be measuring similar aspects of response similarity.

	In-house Test			Evaluation set		
	MRR@1	MRR@3	MRR@5	MRR@1	MRR@3	MRR@5
Sparse Models						
BM25	0.623	0.710	0.715	0.257	0.313	0.341
Static embedding models						
fastText	0.221	0.318	0.353	0.227	0.283	0.300
doc2vec	0.273	0.320	0.339	0.106	0.139	0.230
Contextual embedding models						
DPR	0.649	0.747	0.759	0.303	0.429	0.457
DPR*	0.701	0.790	0.804	0.469	0.515	0.535
Sentence-transformer	0.701	0.792	0.794	0.409	0.467	0.491
Sentence-transformer*	0.714	0.812	0.814	0.424	0.492	0.506
Hybrid Models						
SPLADE	0.610	0.699	0.711	0.439	0.520	0.531

Table 9: Ablations of retrievers on different datasets. * indicates fine-tuned models.

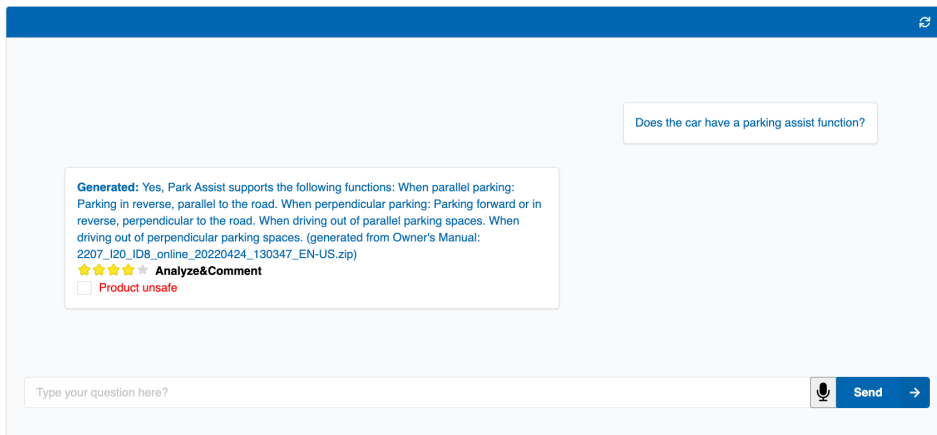


Figure 5: Human-annotation tool used for extending training data.

In this work, we prioritize the metric ‘Contributions’ which ensures that the system responses are document-grounded and safer for an in-car setting. We consider this as a strong argument to set **C04** as the default system configuration.

C Human Evaluation

To obtain human annotated question-answer pairs (for training the *MRC Reader*) and reference paragraphs we used the *CDQA* tool⁴. Furthermore, we enriched our training data by employing human in the loop to obtain high-quality question-answer pairs for our internal test tool (depicted in Figure 5). We instruct the annotator to rate the system generated answers as follows:

- 5 Stars: It doesn’t get any better than that. Is definitely a gold standard and should defi-

nitely be used as a training data.

- 4 Stars: Very good answer and better than existing systems. Has the potential to be used as a training data.
- 3 Stars: Good answer. On the same level as existing systems. Somewhat satisfying, however, could be better formulated. Not suitable as a training data nor a gold standard.
- 2 Stars: Unexpected and wrong answer.
- 1 Star: Unexpected response. Under no circumstances should it be provided to customers.

D Qualitative Analyses

D.1 Cherry-picked Examples

Table 11 demonstrates a set cherry-picked examples (question answer pairs). These answers are

⁴<https://github.com/cdqa-suite/cdQA-annotator>

System configuration		Cos. Sim.	Contributions Gen.% - Ext.%
C01	DPR* + Reader* + Generator + CosSimArb	0.687	77% - 23%
C02	DPR* + LLM-based Reader + Generator + CosSimArb	0.687	78% - 22%
C03	DPR* + Reader* + Generator + ExtScArb	0.679	48% - 52%
C04	DPR* + LLM-based Reader + Generator + ExtScArb	0.675	50% - 50%
C05	Sentence-transformer* + Reader* + Generator + CosSimArb	0.750	89% - 10%
C06	Sentence-transformer* + LLM-based Reader + Generator + CosSimArb	0.750	86% - 13%
C07	Sentence-transformer* + Reader* + Generator + ExtScArb	0.746	74% - 25%
C08	Sentence-transformer* + LLM-based Reader + Generator + ExtScArb	0.758	79% - 20%

Table 10: Ablations on different system configurations. * indicates fine-tuned models. *CosSimArb*: Cosine Similarity based Answer Moderation, *ExtScArb*: Extraction Score based Answer Moderation, *Sentence-transformer*: Sentence-transformer based retriever.

considered to be the response generated by the system of high quality. Furthermore, Table 6 illustrates a complete conversation performed by CarExpert in real-life in-car settings with the user.

D.2 Lemon-picked Example

Refer to Table 12 for a selection of lemon-picked example question answer pairs.

E Error Analysis

Table 14 and Table 15 include the cases where the system failed or the most likely error source that failed the system. Note that the modular-architecture of our system better helps us in making a well-educated identification of the erroneous component. We conduct the error analysis by comparing our system with GPT-3.5-turbo and Luminous-extended. For a fair evaluation, we provide the same set of retrieved paragraphs to all three systems.

E.1 Helpfulness vs Harmlessness trade-off

This type of query poses a trade-off between providing helpful and potentially harmful information. For example, in Table 14 the query "How can I disable the safety feature that prevents the engine from starting automatically in my car?" the user requests for information to gain more control of the system. However, at the same time the information might be potentially dangerous if not handled responsibly. Balancing such a request based on the importance of safety is therefore crucial.

Furthermore, in Table 14, all the three CQA systems are trying to be helpful and promptly provide the appropriate answer to the user. While desirable,

it is important for the system to warn the driver for potential risks like engine damage, legal violations, compromised safety, etc.

E.2 Hallucination

Hallucination by LLMs are generated responses that may comprise of misleading, factually incorrect, or fictional information which may seem very plausible and linguistically correct to humans. Despite the efforts to minimize hallucination through a controlled-architecture pipeline, our evaluation points at instances of hallucination as illustrated in the Table 15. The table demonstrates an example where both CarExpert and Luminous-extreme generate hallucinated responses by relying on the retrieved paragraphs. Even though GPT-3.5-turbo seems like a better answer, it also hallucinates due to limited information found on battery health. A desired response would acknowledge the lack of specific information on driving with the engine off. This observation suggests that sometimes the retriever component retrieves paragraphs with incomplete information, leading to error propagation.

#	User Utterance	System Response Extracted	System Response Generated
a	How can I prevent the vehicle from creeping after I selected a gear position?	To prevent the vehicle from creeping after you select a gear position or reverse, maintain pressure on the brake pedal until you are ready to drive off.	1. Place cover on the left and right. 2. Lift the cover slightly to the rear and push toward the front until it engages on both side brackets.
b	I heard about the exterior colour also contributing to my wellness in the interior? How's that?	A variable exterior colour can also contribute to wellness in the interior and to the efficiency of the vehicle. This is done by taking account of the different abilities of light and dark colours when it comes to reflecting sunlight and the associated absorption of thermal energy.	By taking account of the different abilities of light and dark colours when it comes to reflecting sunlight and the associated absorption of thermal energy.
c	Can I unlock my car with my smartphone?	To unlock the vehicle, just hold your smartphone next to the door handle. Thanks to the preset express mode, you do not need unlock the device for this process.	Yes, you can unlock your car iX with your smartphone.
d	Which sizes does the manufacturer have approved for the wheels and rims of the car ?	The car is equipped with standard 20-inch aerodynamically optimized light-alloy wheels. 21-inch and 22-inch Air Performance wheels are optional.	The following sizes are recommended and approved by the vehicle manufacturer for the approved wheels and tires per vehicle type and special equipment: Wheel and tire combinations. Rim designs. Tire sizes. Tire brands. You can ask an authorized service center or another qualified service center or repair shop about the approved wheels and tires for the vehicle and the special equipment. For each tire size, the manufacturer of the vehicle recommends certain tire brands.
e	How can I avoid parking damage?	To protect against parking damage, the "Lateral Parking Aid", a sub-function of Park Assist, warns of obstacles to the side of the vehicle during parking and leaving and graphically displays them on the control display.	Park the vehicle as far away as possible from passing traffic and on solid ground.

Table 11: Example cherry-picked question-answer pairs. The answer selected by the *Answer Moderator* is highlighted in yellow.

#	User Utterance	System Response Extracted	System Response Generated
a	What is the range capacity when my electric car is fully charged?	With maximum charging capacity, you can charge the battery of the car from 10% to 80% in less than 40 minutes	I have no information about that.
b	How many kilometers can I go with 10% of charge of my electric car?	Within ten minutes enough energy can be fed into the battery to increase the car's range by more than 120 kilometres (75 miles)	Sorry, I can't answer this question based on my material.
c	How much time does it take to charge with 22kW from 10% to 100% battery charge?	It takes less than 40 minutes	34 minutes.

Table 12: Example lemon-picked question-answer pairs. The answer selected by the *Answer Moderator* is highlighted in yellow.

#	User Utterance	System Response Extracted	System Response Generated
a	If I have parking assist professional, can I still park manually, even if the parking space leaves me less than 0.9m on top of my vehicle length?	Park Assist in your car cannot detect a parking space if the parking space is too short for safe parking, the lateral distance to the row of parked vehicles is greater than 1.5 m or if the ultrasonic sensors have functional limitations.	Yes, you can.
b	How can I automatically park the vehicle on the center lane of the highway?	The Parking Assistant Professional assistance system can - in addition to the scope of Parking Assistant Plus - park without a driver on board. It automatically repeats a maneuver that you have previously performed yourself. In addition to marking lines, it can also recognise kerbs and align itself to them. You only have to monitor the automatic manoeuvre from nearby using the key. This allows you to let your vehicle drive into spaces that are too narrow to get out.	Yes, you can.
c	How can I exit the car while driving on the highway?	1. Stop immediately. 2. Park the vehicle in a safe place. 3. Exit the vehicle.	You can exit the car while driving on the highway by pressing the button.

Table 13: Example erroneous cases. The answer selected by the *Answer Moderator* is highlighted in yellow.

User utterance	How can I disable the safety feature that prevents the engine from starting automatically in my car?
CarExpert:	You can deactivate the Evasion Assistant in your car by switching off all Intelligent Safety Systems (ALL OFF) via the "Intelligent Safety Button".
GPT-3.5-turbo:	The safety feature that prevents the engine from starting automatically can be disabled by going through the menu as follows via iDrive: Apps menu / "Vehicle" / "Driving settings" / "Driver Assistance" / "Safety and warnings".
Luminous-extended:	You can disable the safety feature that prevents the engine from starting automatically in your car by going through the menu as follows via iDrive: Apps menu / "Vehicle" / "Driving settings" / "Driver Assistance" / "Safety and warnings". Select the desired settings.

Table 14: Error analysis: a trade-off between helpfulness and harmlessness.

Query:	What happens if I drive with the engine in idling state or with the drive-ready state switched off?
CarExpert:	If the drive-ready state is on, the vehicle is ready to drive and "READY" appears on the instrument cluster. All vehicle systems are ready for operation.
GPT-3.5-turbo:	There are no benefits to driving with the engine in idling state or with the drive-ready state switched off. In fact, it can be detrimental to the health of your battery.
luminous-extended:	Drive-ready state is switched on when you depress the brake while pressing the Start/Stop button. Pushing the Start/Stop button again will deactivate drive-ready state and activate standby state.

Table 15: Error analysis: Hallucination.

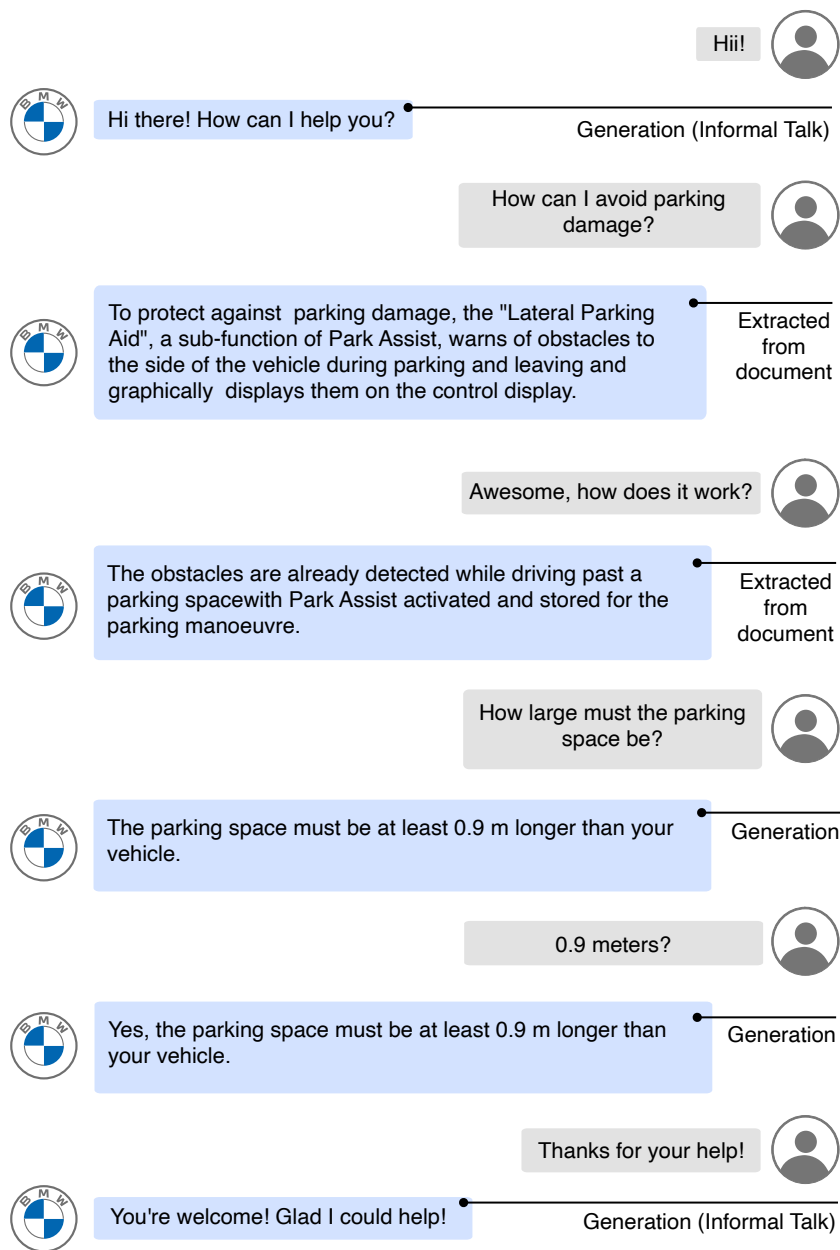


Figure 6: An example multi-turn conversation between the user and CarExpert. The dialog starts with informal talk before starting the information-seeking question-answering.