

TeleQnA: A Benchmark Dataset to Assess Large Language Models Telecommunications Knowledge

Ali Maatouk^{*†}, Fadhel Ayed^{*†}, Nicola Piovesan[†], Antonio De Domenico[†], Merouane Debbah[‡], Zhi-Quan Luo[§]

[†]Paris Research Center, Huawei Technologies, Boulogne-Billancourt, France

[‡]Khalifa University of Science and Technology, Abu Dhabi, UAE

[§]The Chinese University of Hong Kong, Shenzhen, China

Abstract—We introduce TeleQnA¹, the first benchmark dataset designed to evaluate the knowledge of Large Language Models (LLMs) in telecommunications. Comprising 10,000 questions and answers, this dataset draws from diverse sources, including standards and research articles. This paper outlines the automated question generation framework responsible for creating this dataset, along with how human input was integrated at various stages to ensure the quality of the questions. Afterwards, using the provided dataset, an evaluation is conducted to assess the capabilities of LLMs, including GPT-3.5 and GPT-4. The results highlight that these models struggle with complex standards-related questions but exhibit proficiency in addressing general telecom-related inquiries. Additionally, our results showcase how incorporating telecom knowledge context significantly enhances their performance, thus shedding light on the need for a specialized telecom foundation model. Finally, the dataset is shared with active telecom professionals, whose performance is subsequently benchmarked against that of the LLMs. The findings illustrate that LLMs can rival the performance of active professionals in telecom knowledge, thanks to their capacity to process vast amounts of information, underscoring the potential of LLMs within this domain. The dataset has been made publicly accessible on GitHub².

I. INTRODUCTION

LLMs have recently sparked a revolution in the realm of Natural Language Processing (NLP), elevating automatic text generation and interaction to high levels of performance. Currently, the advent of LLMs has begun to impact many other fields beyond NLP, such as medicine [1] and finance [2]. Much like these fields, the telecom industry stands on the brink of experiencing the potential impact LLMs could have on its landscape. Indeed, it is foreseen that LLMs may greatly facilitate a wide range of tasks, including troubleshooting anomalies, comprehending standards, and providing optimization recommendations for network enhancement [3]–[5].

The success of LLMs depends critically on benchmark datasets designed to assess their proficiency in specific domains. These datasets also play a pivotal role in determining the optimal architectural design for LLMs and guiding the pre-training procedure in these specialized fields. In the context of NLP applications, notable examples include the TriviaQA [1], HelloSwag [6], and SIQA [7] datasets, tailored to evaluate LLMs capabilities in reading comprehension, commonsense reasoning, and social intelligence, respectively. The same can

be observed in other domains, such as medicine and finance, where benchmark datasets like MultiMedQA [1] and FLUE [8] have been introduced to assess the proficiency of LLMs in these fields.

As LLMs find their way into the telecommunications industry, a clear and pressing issue arises—there is a notable absence of a benchmark dataset designed to evaluate these models’ proficiency in telecom. Consequently, there is an urgent need for such a dataset, as highlighted in various prior research (e.g., [9]). This paper aims to bridge this gap by introducing TeleQnA, the first benchmark dataset tailored specifically for evaluating LLMs’ telecom knowledge. The contributions of our paper are fourfold:

- First, we curate a comprehensive dataset that encompasses a vast corpus, spanning various disciplines within the telecommunications domain and taking on diverse forms, such as standards and research articles.
- Second, we introduce an automated Question and Answer (QnA) generation process based on the gathered material, involving two LLMs communicating with one another. We show how this architecture, complemented by human input at different stages, not only ensures the quality of the generated dataset but also facilitates its scalability to accommodate a wide array of multidisciplinary topics.
- Third, we proceed to assess the telecom knowledge of prominent LLMs: GPT-3.5 and GPT-4 [10]. We illustrate their proficiency in responding to general telecom inquiries and highlight their deficiencies in addressing intricate questions related to standards specifications. Additionally, we demonstrate how the incorporation of contextual information significantly enhances these models’ performance, particularly in areas where they typically face difficulties. All of these findings underscore the necessity for a specialized telecom foundation model to fully harness the potential of LLMs within the industry.
- Lastly, we assess the performance of professionals actively engaged in this field in comparison to that of LLMs. Our findings demonstrate that LLMs can compete with active professionals in telecom knowledge. This potential stems from their ability to process extensive volumes of information, highlighting the valuable role LLMs are poised to play in shaping the future of this domain.

II. DATASET SOURCES

When curating the dataset sources, our primary objective was to assemble a diverse and comprehensive collection of

^{*}Equal contribution.

¹Samples from this dataset can be found in Appendix A.

²<https://github.com/netop-team/TeleQnA>

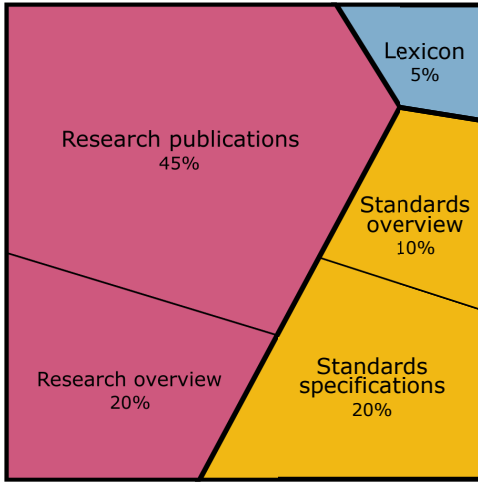


Figure 1: Distribution of the QnA dataset among the categories of the collected source materials.

open-access data related to the telecom industry. This involved a broad spectrum of content, encompassing standards and scholarly research materials. In total, our collection comprised around 25,000 pages, containing approximately 6 million words. This approach guarantees that the resultant QnA dataset offers an in-depth and well-rounded depiction of the telecom industry’s multifaceted aspects. Furthermore, this diversity serves as a valuable resource for evaluating the strengths of an LLM in various aspects, including standard specifications and general telecom expertise. The allocation of the QnA dataset among the three specified categories – research, standards, and lexicon – is illustrated in Fig. 1.

A. Standards Documents

Standards play a pivotal role in the field of telecommunications by ensuring that different technologies from various vendors can work together cohesively. These standards are established and maintained by recognized standardization bodies, including but not limited to 3GPP, IEEE, and ITU. Incorporating these documents when creating the QnA dataset is immensely valuable for capturing the intricacies of the telecom industry. Particularly, it allows delving deep into the technical facets and the present-day implementation of telecom technologies. With this in mind, we have integrated two distinct document sets into our raw data, all closely associated with standard activities.

- **Technical Specifications:** Technical specifications are detailed documents that define the technical standards for telecommunications systems. Given the sheer number of these documents, we have uniformly sampled a portion of them to ensure a non-biased representative selection of these documents across multiple standardization bodies.
- **Standards Overview:** To offer a more comprehensive perspective on standards, moving beyond the intricate technical specifications, we have gathered materials from various sources, encompassing standards summaries, review publications, and standards-related white papers.

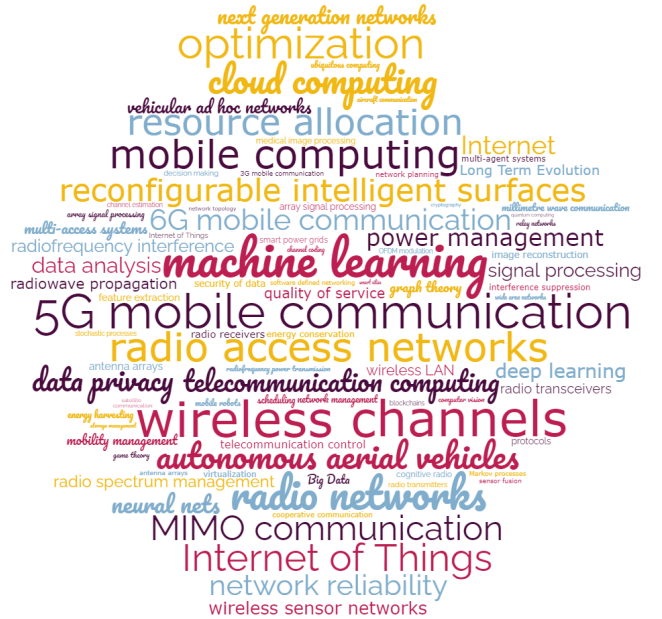


Figure 2: Word cloud illustrating the prominent keywords extracted from the collected research materials.

B. Research Material

The landscape of publishing venues within the telecommunications field is vast. It encompasses a wide array of publishers, and is presented in diverse formats such as journals and conferences, with each repository housing a plethora of valuable information. The sheer abundance and diversity of material available pose a considerable challenge when attempting to create a QnA dataset based on these scholarly works. Additionally, the varying levels of citation and peer review standards, coupled with our commitment to exclusively use open-access material, introduce another layer of complexity to the task of curating a reliable and representative dataset. To address these challenges, we have categorized our sources into two groups: *Research publications* and *Research overview*. The former category primarily comprises technical details sourced from research articles and in-depth technical open-access books, while the second includes content from review and survey venues. Within each of these categories, we have adhered to the subsequent set of guiding principles:

- **Diversity:** Our first guiding principle was the diversity of material. In particular, we considered a broad spectrum of outlets, encompassing well-established venues hosted by various publishers.
- **Relevance:** Within these outlets, we have chosen to identify the most influential material by considering their citation counts over the past 15 years. Our selection of this time frame is deliberate, as it allows us to incorporate the most recent and pertinent findings. On the other hand, the number of citations serves as an indicator of a manuscript popularity and impact over time. Furthermore, to infuse the dataset with the insights of industry experts, we have added material that telecom professionals have

identified as pivotal contributions to this collection.

- **Bias Minimization:** To mitigate potential biases within the dataset, we have also taken a deliberate approach to uniformly sample a number of manuscripts from these venues. This strategy enhances the dataset diversity, encompassing not only popular topics but also those that, while less prominent, remain highly pertinent to the field.
- **Multi-disciplinarity:** We have ensured that our sources extend beyond the realm of pure communications and encompass topics relevant to the field. This expanded scope includes material from various domains, such as machine learning and optimization. This approach stems from the recognition that an LLM tailored for the telecom sector should possess a degree of proficiency in interconnected domains. Although the majority of the sources consist of communication-oriented material, this diverse and multidisciplinary knowledge enriches the dataset, fortifying its comprehensive and interdisciplinary nature.

A word cloud representing the keywords extracted from the considered research material is illustrated in Fig. 2.

C. Telecom Lexicon

The telecommunications industry possesses a distinctive lexicon, characterized by an extensive array of technical terms, each holding significance in understanding the intricacies of communication networks. The incorporation of this lexicon into the generation of the QnA dataset holds significant value, as it provides a perfect avenue to evaluate the proficiency of an LLM in comprehending telecommunications terminology. To procure this source of data, we conducted a thorough examination of the standards documents and research papers we have collected and created a comprehensive dictionary of the technical terminology used therein.

III. DATASET CHARACTERISTICS

A. Question Type

With the source material prepared, our first course of action involved adopting a question type, and after careful consideration, we opted for a multiple-choice approach. This decision was made for several reasons:

- **Precision Assessment:** The multiple-choice format offers a straightforward means of gauging the accuracy of an LLM when presented with the dataset.
- **Complex Decision-Making:** The multiple-choice format allows us to probe deeper into the LLM’s decision-making capabilities. It enables us to assess whether the LLM can accurately discern situations where multiple choices may be correct, rather than just one.
- **Nuanced Discrimination:** By adopting a multiple-choice framework, we can examine the LLM’s ability to select the correct answer when presented with options that resemble the correct choice but are ultimately incorrect.

For the reasons mentioned above, multiple-choice approaches have gained popularity in the realm of machine learning benchmarks. For instance, they have been prominently utilized in datasets like MCTest and SWAG [11], [12].

B. Dataset Format

To maintain consistency across all questions, we have implemented a standardized format. Each question is represented in JSON format, comprising five distinct fields:

- **Question:** This field consists of a string that presents the question associated with a specific concept within the telecommunications domain.
- **Options:** This field comprises a set of strings representing the various answer options.
- **Answer:** This field contains a string that adheres to the format “option ID: Answer” and presents the correct response to the question. A single option is correct; however, options may include choices like “All of the Above” or “Both options 1 and 2”.
- **Explanation:** This field encompasses a string that clarifies the reasoning behind the correct answer.
- **Category:** This field includes a label identifying the source category depicted in Fig. 1.

IV. DATASET CREATION

To construct a comprehensive QnA dataset covering the multifaceted domain of telecommunications, a substantial number of questions is required. This task is further complicated by the specialized nature of telecom knowledge, demanding expertise to craft pertinent questions, answers, and explanations. Moreover, the telecom documents we collected often contain highly intricate information, making it infeasible for a team of human experts to generate questions and answers that comprehensively cover the diverse range of telecom subdomains. Adding to the complexity, we require these questions to be challenging to answer, and asking human experts to craft incorrect options that are also challenging across a wide array of subdomains is a task far from trivial. Given these challenges, adopting a crowdsourcing approach, similar to the one employed for the MCTest dataset [11], becomes impractical. To address these challenges, we leveraged OpenAI’s GPT-3.5’s API³ to facilitate the generation process. This involved developing an automated framework based on LLMs, which will be elaborated on in this section. An overview of the entire process can be found in Figure 3.

A. Preprocessing

Our initial step involved cleaning and formatting the collected data sources for seamless integration with the utilized LLMs. In parallel, we created acronyms extractors that identify the acronyms found in these documents and their corresponding definitions, allowing us to grasp the terminology employed within these materials.

B. LLM Agents

Our LLM-based framework revolves around two GPT-3.5 LLM agents: a generator and a validator, each with distinct roles as outlined below.

³The quality of the generation process was comparable between GPT-3.5 and GPT-4, which led us to favor the latter due to its lower cost.

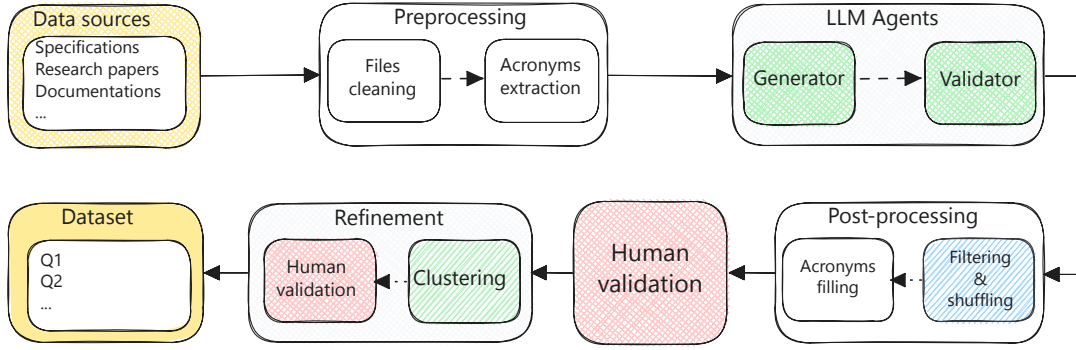


Figure 3: A high-level overview of the QnA generation process.

- **Generator:** The generator is an LLM instance responsible for producing multiple-choice questions. To accomplish this, it is supplied with a segment of a telecommunications document, selected from those discussed in Section II. Subsequently, it is instructed to generate a set of multiple-choice questions based on the given context, adhering to the formatting guidelines. Additionally, it is directed to ensure that the questions are self-contained. Furthermore, it is prompted to provide a rationale for the correct answer, typically by referencing sentences within the document to substantiate the choice of the answer.
- **Validator:** The questions and options generated by the generator, together with the contextual information (i.e., the relevant segment of the document), are given as input to another GPT-3.5 LLM, referred to as the validator. To prevent any biasing, the answer and the rationale behind the correct answer are intentionally omitted when providing input to the validator. The validator’s primary function is to select an option among those provided solely based on its knowledge and the contextual information provided. If the validator picks the option designated as correct by the generator, the question is retained. However, if the validator selects an incorrect option, the question is discarded. This automated layer serves to enhance the overall accuracy and mitigate any potential instances of misinformation or fabrication that the generator may introduce.

C. Post-Processing

Moving forward in the process, we transition into the post-processing stage. Following the initial validation step, this phase focuses on the following key aspects:

- **Filtering:** We developed a tool designed to identify and filter out questions that are not self-contained, meaning they refer to the text from which are taken from by including particular terms such as “figure” and “fig”.
- **Shuffling:** Our experiments have revealed a tendency of GPT-3.5 to place the correct answer in option 1. To mitigate this bias, we shuffled the options to ensure equal likelihood for all choices.
- **Acronyms:** We designed acronym detectors that identify the use of acronyms in the questions and map them to their respective definitions.

D. Human Validation - First Stage

Having validated the questions using the validator and completed the necessary post-processing steps, we now introduce the first stage that incorporates human intervention. Specifically, a telecom expert is assigned the task of carefully reviewing each question, the provided answer options, the explanation given by the question generator for the answer, and the contextual information used. The objective is to make an informed decision regarding whether or not to discard a question. The primary objectives of this human-in-the-loop process are twofold:

- **Ensure question correctness:** The expert’s role is to verify the accuracy of the questions, their associated options, and the designated correct answer, effectively eliminating improperly generated questions.
- **Ensure that the questions are self-contained:** The questions should not revolve around content that is of localized interest, such as experimental results specific to a proposed algorithm.

E. Refinement

The collected telecom sources inherently exhibit overlaps across the various venues. This overlap gives rise to redundancy within the QnA dataset, a challenge that necessitates a refinement of the dataset. To address this issue systematically, we employed Open AI’s Ada v2 text embeddings model to compute embeddings for all the questions in the QnA dataset along with their corresponding explanations. These embeddings serve as compact numerical representations, encapsulating the essence of phrases, thus enabling us to gauge the semantic similarity between any two questions by measuring the proximity of their respective embeddings in this high dimensional space. Having computed these embeddings, we then applied the K-Means clustering algorithm to group questions into clusters, in order to streamline the subsequent human intervention step. In fact, the concluding phase encompasses a secondary human validation that serves a dual purpose:

- **Eliminating redundant questions:** The primary goal of the intervention is to eliminate duplicate or semantically-redundant questions within each cluster.
- **Final iteration of quality checking:** The second objective entails conducting a final round of verification, similar to the first human validation stage.

	Run 1	Run 2	Run 3	Run 4	Mean	Std
50 questions	63.35	65.63	65.10	63.79	64.46	0.92
25 questions	65.19	67.56	66.90	66.30	66.48	0.87
10 questions	67.64	64.32	66.94	67.20	66.52	1.29
5 questions	65.45	67.55	68.52	66.85	67.09	1.11
1 question	66.43	67.48	68.10	66.46	67.11	0.70

Table 1: Illustration of GPT-3.5’s accuracy (%) across multiple runs and question batch sizes.

V. PERFORMANCE EVALUATION

In this section, we conduct a comprehensive evaluation of GPT-3.5 and GPT-4 using the TeleQnA dataset. Our objective is to assess their proficiency in telecom knowledge while also offering insights into their strengths and weaknesses. Additionally, we compare the performance of telecom professionals to these language models, establishing a benchmark for reference. As a performance measure, we define the accuracy as the percentage of questions for which the entity at hand selected the option marked as correct in the dataset.

A. Batch Size and Variance Study

Our numerical analysis begins by examining the variance of the accuracy achieved by the LLMs when responding to the TeleQnA dataset. Specifically, we are interested in tracking how the accuracy evolves with each iteration of LLM questioning. Additionally, we investigate the impact of sending questions in random batches of varying sizes ($B = 1, 5, 10, 25,$ and 50) to the LLM. The former investigation aims to show the number of iterations required to obtain a statistically reliable accuracy score for the LLM’s telecom knowledge. Meanwhile, the latter investigation seeks to identify trends when querying the LLM’s API with larger question batches, potentially streamlining the process by sending questions in batches rather than one by one. To obtain these findings, we uniformly sampled 1000 questions and conducted accuracy assessments across multiple runs and with different batch sizes. The results of these experiments are presented in Table I for GPT-3.5⁴.

The results demonstrate the consistency of accuracy across multiple runs, with a standard deviation of approximately 1%. Given that the standard deviation decreases with the number of questions considered, it can be inferred that this deviation is even smaller when applied to the entire dataset. Therefore, we can conclude that conducting a single run is sufficiently reliable for assessing the proficiency of the LLM. Conversely, we observe a decline in mean accuracy when questions are submitted in larger batches. This decrease may be attributed to the inherent diversity of questions stemming from various sub-domains within telecommunications. When questions are batched together, the LLM encounters a broader spectrum of topics, which may result in less specialized responses and, therefore, a potential drop in accuracy. However, we note that this decline in performance remains modest when transitioning from a batch size of $B = 1$ to $B = 5$. Thus, we have chosen to

⁴GPT-4 exhibited similar trends, and therefore, we have opted to exclude the specific details.

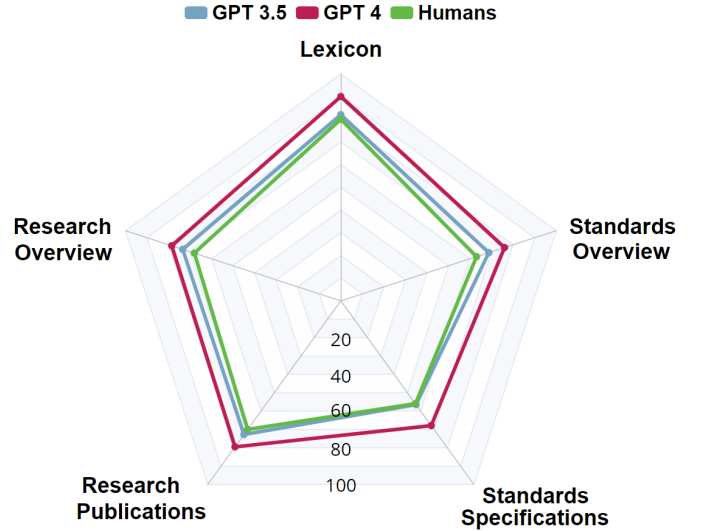


Figure 4: Accuracy (%) comparison among GPT-3.5, GPT-4, and active professionals in all five categories.

utilize this batch size in the rest of our experiments to reduce the number of queries needed for assessment.

B. Performance Benchmarks

We evaluate the performance of GPT-3.5 and GPT-4 across the various categories of the TeleQnA dataset. The results are reported in Fig. 4, and the details can be found in Appendix B. As anticipated, GPT-4 consistently outperforms GPT-3.5, demonstrating around 7% improvement across all categories. Notably, LLMs exhibit exceptional performance in the lexicon category, which encompasses general telecom knowledge and terminology, achieving approximately 87% accuracy for GPT-4. Conversely, these models face challenges when confronted with more intricate questions related to standards, with the highest performing model, GPT-4, achieving a modest 64% accuracy in this domain. In summary, GPT-3.5 averaged an accuracy of 67%, while GPT-4 achieved an accuracy of 74%. These results demonstrate that these models possess a solid foundation in general telecom expertise. However, to attain higher accuracy in responding to complex inquiries, further adaptations to the telecom domain are necessary.

In our final step, we conducted a performance benchmark comparing active professionals to LLMs. To accomplish this, we designed surveys consisting of ten questions, with two questions representing each of the five categories. These surveys were then distributed to thirty active professionals in the telecom field, working in various domains such as standardization, signal processing, optical networks, etc. These professionals were explicitly instructed not to utilize search engines or external references, relying solely on their personal knowledge. The results reveal that LLMs and active professionals exhibit comparable performance in general telecom knowledge. However, in the case of intricate questions related to research and standards, LLMs demonstrate the capability to rival these professionals. This is attributed to LLMs’ ability to digest and memorize complex and intricate documents.

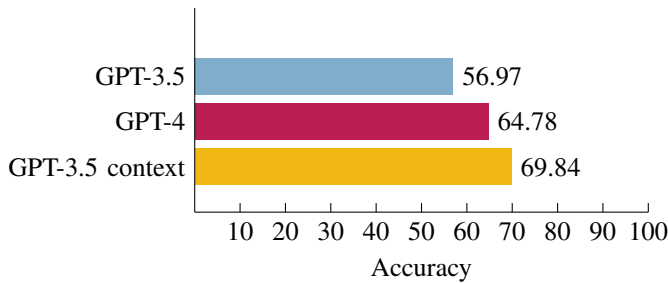


Figure 5: Accuracy (%) comparison among GPT-3.5, GPT-4, and GPT-3.5 with context in the standards specifications category.

Furthermore, it is crucial to recognize the challenge faced by professionals when responding to these questions, as they encompass a broad range of telecom subdomains that these individuals may not be necessarily actively engaged with in their work. Considering all factors, our results underscore the significant promise that LLMs hold within this domain, as demonstrated by their competitiveness within this extensive and comprehensive dataset.

C. Influence of Context

Until this stage of our benchmarking process, we have been querying the LLMs without accompanying contextual information for the questions. Nevertheless, in this subsection, our goal is to investigate how supplementing questions with additional context affects the accuracy of these models.

To accomplish this, we have focused on the standards specifications sources, encompassing thousands of technical standards pages. Our selection of this category is driven by the fact that this is where LLMs have exhibited the lowest performance. With this in mind, we segmented these pages into approximately 500-word segments before generating embeddings for each segment using OpenAI’s Ada v2 text embeddings. Moreover, we employed the same OpenAI model to create embeddings for the questions and corresponding options belonging to this category. Following that, we constructed a distance matrix between the embeddings of each question-options pair and those of each segment. The next step involved querying the LLM by supplying batches of five questions-options pairs, and additionally, as context, the top-3 closest segments to the five question-option pairs based on the distance matrix. The outcomes of these experiments are illustrated in Fig. 5.

The 22.5% relative accuracy gain highlights the significant enhancement in performance achieved by incorporating contextual information, demonstrating how even less advanced models like GPT-3.5 can match the performance of state-of-the-art GPT-4 model. This underscores the necessity for a specialized telecom language models, fine-tuned or trained specifically on telecom-related data. Developing such a foundation model has the potential to push the boundaries of LLMs performance in the telecom domain, paving the way for a wide range of use cases that demand telecom knowledge and expertise.

VI. CONCLUSIONS, LIMITATIONS, AND FUTURE DIRECTIONS

A. Conclusions

In this paper, we introduced TeleQnA, the first open-source benchmark dataset tailored specifically to the telecom industry. With a collection of 10,000 questions and answers drawn from various sources in the field, TeleQnA serves as an evaluation tool for assessing the knowledge of LLMs in the telecommunications domain. Throughout this paper, we have detailed the development process of TeleQnA and conducted an extensive evaluation of GPT-3.5 and GPT-4 using TeleQnA. The results have revealed interesting insights: although these models excel in addressing general telecom-related questions, they face challenges when tackling complex inquiries. Importantly, our findings emphasize the critical need for specialized telecom foundation models.

B. Limitations

While our framework’s design is aimed at minimizing errors in the questions, it remains dependent on human judgment to determine whether a generated question should be retained or discarded. This introduces an element of subjectivity, particularly in cases where questions offer multiple valid options, albeit with varying degrees of correctness. In such instances, human subjectivity comes into play. For instance, let us consider the following question:

Question 1: What advantages does a data center-enabled HAP (High Altitude Platform) offer from an energy perspective?

- *Option 1:* Saves cooling energy and harvests solar energy
- *Option 2:* Uses renewable energy and offsets carbon emissions
- *Option 3:* Deploys large-scale solar panels and batteries
- *Option 4:* Requires less energy for communication links

Answer: *Option 1:* Saves cooling energy and harvests solar energy

Explanation: The data center-enabled HAP saves cooling energy due to its location at the naturally low temperature stratosphere and harvests solar energy using large solar panels.

Category: Research publications

An expert might discern that Option 1 is more accurate than Option 2. In fact, although Option 2 is indeed an advantage offered by data center-enabled HAP, Option 1 provides more insights into the fact that high-altitude platforms are deployed in the low temperature stratosphere, while also considering the fact that solar energy is used (i.e., a renewable energy). We eliminated some of these questions to ensure the utmost objectivity and correctness, but deliberately retained others to challenge the LLM and assess its capability to discern between options exhibiting differing degrees of correctness.

Another important consideration is that a substantial portion of the generated questions draw from research papers to encompass a wide range of subjects. While this enriches the dataset, it introduces an element of subjectivity, as authors in

their papers may present their own perspectives and interpretations of phenomena, which can be inherently subjective. While our aim has been to filter out such questions, it is important to acknowledge that, in the end, the process involves an element of human subjectivity. All of these factors collectively indicate that, despite our efforts to optimize the process extensively, we must acknowledge that the dataset generation is not flawless, as it continues to rely on human input.

C. Future Directions

In light of these limitations, one thing becomes abundantly clear: the development and advancement of extensive and comprehensive telecom knowledge datasets for evaluating LLMs necessitate a collaborative effort spanning the entire industry. We are receptive to all forms of feedback, especially in cases where certain questions may elicit nuanced perspectives from active community members regarding their correctness. In the end, we view TeleQnA as an initial step toward such a collaborative endeavor. Researchers also have the potential to create their own datasets, and by combining these diverse datasets with TeleQnA, one can see a future where an expansive and all-encompassing dataset, akin to those successfully crafted by the machine learning community for their applications of interest, can be created. This, in turn, paves the way for the widespread adoption of high-performing LLMs tailored specifically for telecom applications.

VII. ACKNOWLEDGMENTS

We express our gratitude to the active telecom professionals for their generous contribution of time and participation in our surveys designed to evaluate telecom professionals' knowledge using TeleQnA.

REFERENCES

- [1] K. Singhal *et al.*, "Large language models encode clinical knowledge," *Nature*, vol. 620, no. 7972, pp. 172–180, Aug 2023. [Online]. Available: <https://doi.org/10.1038/s41586-023-06291-2>
- [2] S. Wu *et al.*, "BloombergGPT: A Large Language Model for Finance," *arXiv preprint arXiv:2303.17564*, May 2023.
- [3] A. Maatouk *et al.*, "Large language models for telecom: Forthcoming impact on the industry," *arXiv preprint arXiv:2308.06013*, August 2023.
- [4] Altman Solon, "The opportunity in generative AI for Telecom," White Paper, September 2023. [Online]. Available: <https://pages.awscloud.com/GLOBAL-other-DL-generative-ai-for-telecom-whitepaper-2023-learn.html>
- [5] L. Bariah *et al.*, "Understanding Telecom Language Through Large Language Models," *arXiv preprint arXiv:2306.07933*, 2023.
- [6] R. Zellers *et al.*, "HellaSwag: Can a machine really finish your sentence?" in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, Jul. 2019, pp. 4791–4800.
- [7] M. Sap *et al.*, "Social IQa: Commonsense reasoning about social interactions," in *EMNLP-IJCNLP 2019*, Nov. 2019, pp. 4463–4473.
- [8] R. Shah *et al.*, "When FLUE meets FLANG: Benchmarks and large pretrained language model for financial domain," in *EMNLP 2022*, Dec. 2022, pp. 2322–2335.
- [9] M. Kotaru, "Adapting foundation models for information synthesis of wireless communication specifications," *arXiv preprint arXiv:2308.04033*, August 2023.
- [10] OpenAI, "GPT-4 Technical Report," *arXiv preprint arXiv:2303.08774*, March 2023.
- [11] M. Richardson, "Mctest: A challenge dataset for the open-domain machine comprehension of text," in *EMNLP 2013*, October 2013.
- [12] R. Zellers *et al.*, "Swag: A large-scale adversarial dataset for grounded commonsense inference," in *EMNLP 2018*, October 2018.

Ali Maatouk (Member, IEEE) is a Researcher with Huawei Technologies, France. His research interests include the mathematical modeling and optimization of communication systems, machine learning, and the notion of information freshness.

Fadhel Ayed is a Senior Researcher with Huawei Technologies, France. His research interests include resource-efficient machine learning, stochastic processes, and mathematical foundations of deep learning.

Nicola Piovesan (Member, IEEE) is a Senior Researcher with Huawei Technologies, France. His research interests include energy sustainability, energy efficiency optimization, and machine learning in wireless communication systems.

Antonio De Domenico (Member, IEEE) is a Senior Researcher with Huawei Technologies, France. His research interests include heterogeneous wireless networks, machine learning, and green communications.

Mérouane Debbah (Fellow, IEEE) is a Professor at Khalifa University of Science and Technology in Abu Dhabi. His research interests include Large Language Models, distributed AI systems for networks and semantic communications.

Zhi-Quan Luo (Fellow, IEEE) is is Vice President (Academic) of the Chinese University of Hong Kong (Shenzhen). His research interests lie in optimization algorithms for signal processing, wireless communication and data analytics.

APPENDIX A
SAMPLES OF TELEQNA

Question 1: What does EIRP stand for?

- *Option 1:* Effective Isotropic Radio Power
- *Option 2:* Equivalent Isotropic Radiated Power
- *Option 3:* Efficient Isotropic Radiation Propagation
- *Option 4:* Effective Infrared Radiation Power
- *Option 5:* Equivalent Infrared Radiated Power

Answer: *Option 2:* Equivalent Isotropic Radiated Power

Explanation: EIRP stands for Equivalent Isotropic Radiated Power.

Category: [Lexicon](#)

Question 2: What is a Heterogeneous Network?

- *Option 1:* A network consisting of multiple cells with different characteristics
- *Option 2:* A network that provides services to the user in a managed way
- *Option 3:* A network that changes the radio access mode or system used for bearer services
- *Option 4:* A network that broadcasts a specific indicator and identity
- *Option 5:* A network where the subscriber's personal service environment is controlled

Answer: *Option 1:* A network consisting of multiple cells with different characteristics

Explanation: A heterogeneous network is a 3GPP access network with multiple cells having different characteristics.

Category: [Lexicon](#)

Question 3: Which communication technique uses Light Diodes as transmitters?

- *Option 1:* Large-Scale Antenna Arrays
- *Option 2:* Free Space Optical Communications
- *Option 3:* Heterogeneous Networks
- *Option 4:* Cooperative Relay Communications
- *Option 5:* Cognitive Radio Communications

Answer: *Option 2:* Free Space Optical Communications

Explanation: Free Space Optical Communications use Light Diodes as transmitters.

Category: [Research overview](#)

Question 4: What is the main advantage of using SDR (software-defined radio) techniques?

- *Option 1:* They provide higher computational power for signal processing.
- *Option 2:* They are more cost-effective compared to traditional communication systems.
- *Option 3:* They offer flexibility to update receiver and transmitter functionalities by modifying software code.
- *Option 4:* They result in lower latency in signal processing.
- *Option 5:* None of the above.

Answer: *Option 3:* They offer flexibility to update receiver and transmitter functionalities by modifying software code.

Explanation: SDR platforms are designed to be highly flexible, where all the receiver and transmitter functionalities

can be updated by a simple modification of the software code.

Category: [Research overview](#)

Question 5: What is the maximum number of eigenmodes that the MIMO channel can support? (nt is the number of transmit antennas, nr is the number of receive antennas)

- *Option 1:* nt
- *Option 2:* nr
- *Option 3:* min(nt, nr)
- *Option 4:* max(nt, nr)

Answer: *Option 3:* min(nt, nr)

Explanation: The maximum number of eigenmodes that the MIMO channel can support is min(nt, nr).

Category: [Research publications](#)

Question 6: What are the parameters governing the wireless channel propagation?

- *Option 1:* Positions of the transmitter and receiver antenna elements
- *Option 2:* Carrier frequency
- *Option 3:* Nature and positions of scattering objects in the environment
- *Option 4:* All of the above

Answer: *Option 4:* All of the above

Explanation: The parameters governing the wireless channel propagation include the positions of the transmitter (Tx) and receiver (Rx) antenna elements, the carrier frequency, as well as the nature and positions of scattering objects in the environment.

Category: [Research publications](#)

Question 7: What frequency band does Bluetooth use? [Bluetooth - Overview]

- *Option 1:* 2.4 GHz
- *Option 2:* 5 GHz
- *Option 3:* 40 MHz
- *Option 4:* 1 MHz

Answer: *Option 1:* 2.4 GHz

Explanation: Bluetooth uses the 2.4 GHz licence-free ISM band for transmission.

Category: [Standards overview](#)

Question 8: Which pairs of wires are used in 10/100Base-T? [IEEE 802.3]

- *Option 1:* Pair 1 and pair 2
- *Option 2:* Pair 2 and pair 3
- *Option 3:* Pair 3 and pair 4
- *Option 4:* Pair 4 and pair 1

Answer: *Option 2:* Pair 2 and pair 3

Explanation: In 10/100Base-T, pair 2 (orange/white and orange) and pair 3 (green/white and green) are used.

Category: [Standards overview](#)

Question 9: What is the RRC buffer size for a UE? [3GPP Release 17]

- *Option 1:* 45 MB
- *Option 2:* 45 KB

- *Option 3*: 45 GB
- *Option 4*: 45 TB
- *Option 4*: 4500 KB

Answer: *Option 2*: 45 KB

Explanation: The RRC buffer size for a UE is 45 Kbytes.

Category: [Standards specifications](#)

Question 10: What are the required inputs for the Ndrf_MLModelManagement_Delete service operation? [3GPP Release 18]

- *Option 1*: Storage Transaction Identifier
- *Option 2*: Unique ML Model identifier(s)
- *Option 3*: Model file address(es)
- *Option 4*: Storage Transaction Identifier or Unique ML Model identifier(s)
- *Option 4*: Storage Transaction Identifier and Unique ML Model identifier(s)

Answer: *option 4*: Storage Transaction Identifier or Unique ML Model identifier(s)

Explanation: The required inputs for the Ndrf_MLModelManagement_Delete service operation are either the Storage Transaction Identifier or the Unique ML Model identifier(s).

Category: [Standards specifications](#)

APPENDIX B DETAILED BENCHMARKING RESULTS

In this appendix, we report the results obtained from our experiments involving GPT-3.5 and GPT-4 queries using TeleQnA. These results were obtained by randomly sampling batches of five questions drawn from TeleQnA, which were then provided to the large language model, and subsequently compared to the actual, correct answers.

	GPT-3.5	GPT-4	Humans
Lexicon (500 questions)	82.20	86.80	80.33
Research overview (2000 questions)	68.50	76.25	63.66
Research publications (4500 questions)	70.42	77.62	68.33
Standards overview (1000 questions)	64.00	74.40	61.66
Standards specifications (2000 questions)	56.97	64.78	56.33
Overall accuracy (10000 questions)	67.29	74.91	64.86

Table II: Illustration of GPT-3.5, GPT-4, and active professionals accuracy (%) across the various TeleQnA categories.