Leveraging Fine-Tuned Retrieval-Augmented Generation with Long-Context Support: For 3GPP Standards

Omar Erak¹, Nouf Alabbasi¹, Omar Alhussein¹, Ismail Lotfi¹, Amr Hussein¹, Sami Muhaidat^{2,3}, Merouane Debbah²

¹KU 6G Research Centre, Department of Computer Science, Khalifa University, Abu Dhabi, UAE
²KU 6G Research Centre, Department of Computer and Information Engineering, Khalifa University, Abu Dhabi, UAE
³Department of Systems and Computer Engineering, Carleton University, Ottawa, Canada Emails: omarerak@ieee.org, 100064507@ku.ac.ae, omar.alhussein@ku.ac.ae, ismail.lotfi@ku.ac.ae, 100059484@ku.ac.ae, muhaidat@ieee.org, merouane.debbah@ku.ac.ae

Abstract—Recent studies show that large language models (LLMs) struggle with technical standards in telecommunications. We propose a fine-tuned retrieval-augmented generation (RAG) system based on the Phi-2 small language model (SLM) to serve as an oracle for communication networks. Our developed system leverages forward-looking semantic chunking to adaptively determine parsing breakpoints based on embedding similarity, enabling effective processing of diverse document formats. To handle the challenge of multiple similar contexts in technical standards, we employ a re-ranking algorithm to prioritize the most relevant retrieved chunks. Recognizing the limitations of Phi-2's small context window, we implement a recent technique, namely SelfExtend, to expand the context window during inference, which not only boosts the performance but also can accommodate a wider range of user queries and design requirements from customers to specialized technicians. For fine-tuning, we utilize the low-rank adaptation (LoRA) technique to enhance computational efficiency during training and enable effective fine-tuning on small datasets. Our comprehensive experiments demonstrate substantial improvements over existing question-answering approaches in the telecom domain, achieving performance that exceeds larger language models such as GPT-4 (which is about 880 times larger in size). This work presents a novel approach to leveraging SLMs for communication networks, offering a balance of efficiency and performance. This work can serve as a foundation towards agentic language models for networks.

Index Terms—6G networks, AGI, LLM, LoRA, RAG, retrieval

I. Introduction

LLMs have demonstrated impressive capabilities, from basic automation to complex decision-making [1]. They have proven their efficiency in a variety of tasks, including questions and answering (QnA), code generation, and other problems. Their ability to process natural language and generate humanlike responses makes them powerful tools in many fields, including telecommunications.

As telecom networks grow more complex and data-driven, large language models (LLMs) offer significant potential to enhance automation, optimize network management, and improve customer experiences [2]–[4]. However, to fully leverage

agentic LLM-based models in telecommunications, we need to develop models that deeply understand the nuances of telecom systems and possess comprehensive knowledge of telecom models. Building such specialized models is crucial for adapting LLMs to agentic roles where they can autonomously handle involved tasks, such as dynamic network optimization and predictive maintenance.

LLMs can be adapted to various tasks through fine-tuning or retrieval-augmented generation (RAG). Fine-tuning enhances the performance of LLMs by adjusting the model's internal knowledge through iterative training on specialized datasets. However, the telecom industry is rapidly evolving, rendering fine-tuning an expensive and inefficient approach that cannot easily keep up with such a fast-paced industry. Furthermore, once a model is fine-tuned, editing or forgetting a specific piece of information becomes challenging [5].

RAG, on the other hand, augments text generation with information retrieval, enabling models to produce more accurate and contextually aware responses. This approach allows flexible model adaptation and rapid integration of new information. RAG also grounds the model's response in the relevant retrieved context, reducing (yet not entirely eliminating) the risk of hallucination [6]. Telecom applications can benefit from real-time data to produce more accurate and up-to-date responses. Therefore, we believe LLM-based RAG systems can better enable emerging applications, such as dynamic network management, customer support, and predictive maintenance.

Integrating RAG into telecommunication systems involves deploying LLM frameworks on user equipment and edge devices, a process which presents a significant challenge due to the involved computational intensity of LLMs, both in terms of training and inference costs. This challenge highlights the appeal of SLMs. These language models (LMs) offer computational and storage efficiency while maintaining adequate performance, suggesting their suitability for deployment on edge devices and possible enablement of on-device artificial intelligence (AI) [7]. Several small language models (SLMs) have been proposed in the literature, including Microsoft's Phi-

2 with 2.7B parameters [8] and Gemini Nano 2 with 3.2B parameters [9]. The Phi-2 model is currently considered as the state-of-the-art SLM as it is able to match or outperform models up to 25x larger such as Llama-2-70b which has 70B parameters [10].

Recent studies indicate that while state-of-the-art LLMs perform well on general telecommunications queries, they struggle with questions related to technical standards in the field [11], [12]. We believe this is mainly due to the fact that standard-type knowledge and system specifications do not exist in common research papers and other publications, which serve as the main learning sources for LLMs. For instance, with respect to telecom systems, the polysemy of abbreviations (e.g., SAP: service access sample vs. system application and protocol) can hinder the model from inferring a correct answer. Additionally, the LMs training on generalized knowledge can interfere with its performance in specialized domain tasks. For instance, certain protocols and methods in telecom-specific domains do not necessarily follow the generally followed-upon practices in broader contexts. A model trained on generalized knowledge may not adequately capture these domain-specific nuances and practices.

We propose a carefully developed Phi-2 based fine-tuned RAG system to serve as an oracle for communication networks. To the best of our knowledge, this is the first work to present a fine-tuned RAG system for communication networks. Previous works focus on presenting a frozen RAG framework or fine-tuning an LM. Our RAG system leverages a forward-looking semantic chunking (or parsing) strategy that adaptively determines breakpoints between sentences based on embedding similarity. This approach enables the system to effectively process documents with diverse formatting. In the Third Generation Partnership Project (3GPP) documents, a query can often relate to multiple similar contexts, as discussions and paragraphs on related topics may appear in various sections or be phrased similarly. Therefore, we utilize a re-ranking algorithm to further rank the retrieved chunks based on their relevance to the input query.

Since Phi-2 is an SLM, its performance is limited by its small context window, rendering it inefficient for certain tasks such as responding to open-ended or under-specified queries. This limitation is particularly relevant in telecom applications, where users can range from customers to specialized technicians. Therefore, we utilize a new technique, namely SelfExtend [13], to significantly extend the context window during inference.

Finally, we use Low-Rank Adaptation (LoRA) not only to enhance computational efficiency during training but also because it allows users to effectively fine-tune on small datasets. Our in-depth experiments demonstrate considerable improvements over existing QnA approaches in telecom, contributing to the ongoing advancement of the field.

The remainder of the paper is organized as follows. Section II provides an overview of related works. Section III discusses the proposed fine-tuned Phi-2 RAG system. Section IV provides our experimental results, and Section V concludes the

paper and discusses some insightful future research directions.

II. RELATED WORKS

Several benchmarking datasets have been proposed to enable LLM-based systems in the telecommunication domain, with mainly three different tasks: text classification, text summarization, and multiple choice questions (MCQs). In [14], the SPEC5G dataset was introduced with the objective of performing text classification and text summarization in telecom domain. Ericsson's team introduced TeleQuAD, a private 4,000 entry QnA dataset, and developed a proprietary TeleRoBERTa, a 124M bidirectional encoder representation from transformers (BERT)-based RAG system [15].

In [11], the introduction of the TeleQnA dataset marks a significant advancement in evaluating QnA tasks for telecommunications. The TeleQnA dataset contains 10,000 MCQs about telecommunication systems, curated from the various sources such as 3GPP and research papers. The dataset is verified by telecom human-in-the-loop experts. Another dataset, namely TSpec-LLM, is recently released [16]. The authors develop an automated framework to generate QnA pairs from 3GPP specifications, then test a naive RAG architecture to assess the quality of their dataset. In [17], Gajjar et al. introduce ORAN-Bench-13K, a dataset dedicated for the evaluation of open radio access networks (O-RAN) tasks. The dataset is based on 116 O-RAN specification documents and contains 13,000 pairs of MCQs, based on which ORANSight is developed, a RAG-based framework.

The work in [18] evaluates the performance of various zero shot LLMs in a few tasks in the telecommunications domain, including a QnA task. Among their findings, they note that though LLM models such as Zephyr, and Mistral perform outstandingly in the tasks, their performance still comes strikingly short when compared to GPT-3.5 or GPT-4. The authors in [19] demonstrate the effectiveness of fine-tuning different SLMs for the telecom domain. The work in [12] proposes Telco-RAG, a framework specialized for MCQ answering for telecom applications, tailored to the specific requirements of telecom standards, particularly 3GPP documents. Their contribution focuses on modifying the RAG framework by employing a router, using generated candidate answers to enhance retrieval quality, and appending the definitions of acronyms and technical terms to the user's query.

To the best of the authors' knowledge, this is the first work to present a RAG architecture with fine-tuned SLM generator. Additionally, critical components such as SelfExtend for handling long contexts, re-ranking for enhancing retrieval accuracy, and semantic chunking to preserve contextual coherence have not been proposed in the relevant literature.

III. PROPOSED ARCHITECTURE

A. General Overview

RAG integrates four key components: a chunking mechanism to segment information, an embedding model to encode the information in a latent space, a retriever to fetch relevant context, and a generator to produce responses. In this work,

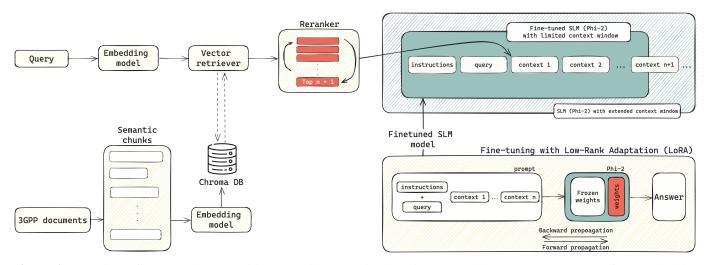


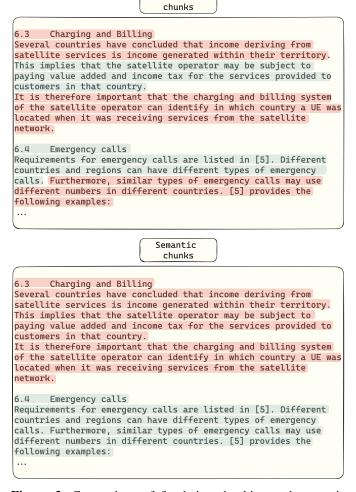
Figure 1. Overview of proposed RAG architecture with semantic chunking, extended context support, and fine-tuned Phi-2 SLM integration for 3GPP document processing.

we present an advanced RAG framework, illustrated in Fig. 1, that optimizes each component to enhance performance and adapt the model to the telecom domain. In our architecture, 3GPP documents are first chunked using a semantic chunker then embedded and stored in a vector database (DB). During inference, the user's query is embedded and then passed to the vector retriever. The retriever performs a vector similarity search in embedding space to return the nearest relevant neighbors from the indexed corpus. The retrieved chunks are then passed to a re-ranking algorithm that returns an ordered sublist containing the most relevant chunks.

The Phi-2 generator model is then given the chunks, the query, and a set of instructions to provide an answer. Here, we efficiently fine-tune the generator using LoRA. The model initially is fine-tuned with n contexts based on retrieved information and the prompt. During inference, we employ the self-extend technique to dynamically expand the context window, accommodating retrieved contexts beyond the initial n limit. In what follows, we highlight the key components of the proposed framework.

B. Semantic Chunking Strategy

The method of splitting or chunking is critical as improper chunks can lead to inaccurate representations in the embedding space. Chunking can be implemented through either fixed-size segments or adaptable segments based on specific criteria. While fixed-size chunking can yield reasonable results and is computationally efficient, it often creates blocks of text that do not consider the content or context. To mitigate this issue, we leverage semantic chunking to split the 3GPP documents [20]. Semantic chunking adaptively determines breakpoints between sentences based on embedding similarity. This way the meaning of the text is preserved through the logical breaks where sentences are semantically connected, rather than arbitrarily cutting the text at fixed intervals. Fig. 2 shows an illustrative example from a 3GPP document.



Fixed size

Figure 2. Comparison of fixed-size chunking and semantic chunking applied to an excerpt from a 3GPP document.

```
prompt:
<context>
Provide a correct answer to a multiple choice question on
wireless communications and standards. Use only one option
from A, B, C, D or E.
<question>
<options>
Instructions:
Step 1. Read the context carefully to identify the most
relevant information.
Step 2. Pay special attention to key terms in the question,
such as "main," "primary," "most important," etc. These terms
indicate that the answer should focus on the most significant
aspect mentioned in the context.
Step 3. Focus on keywords and phrases in the context that
directly relate to the question being asked.
Step 4. Match the context to the options provided and choose
the option that is most explicitly supported by the context
and aligns with the key terms in the question.
Step 5. If the context does not provide clear support for any
option, choose the option that logically fits the context and
the question.
Step 6. Pay attention to abbreviations related to wireless
communications and 3GPP standards
Answer:
```

Figure 3. Prompt structure that includes retrieved context and instructions.

Given that telecom documents don't follow a strict format, leveraging semantic chunking preserves essential context while minimizing information fragmentation and irrelevant grouping.

C. Embedding Model

We utilize *bge-small-en-v1.5*, an open-source embedding model [21]. It is optimized for balancing efficiency and accuracy in text embedding tasks. The model is trained using contrastive learning on a large-scale dataset [22] and it creates vector representations that capture semantic relationships between elements. This method effectively reduces the distance between similar pairs while increasing it between dissimilar ones which helps refine the model's ability to discriminate between relevant features. To ensure faster performance at runtime, these embeddings are calculated and stored in Chroma dB, an AI-native vector database designed to efficiently handle high-dimensional embeddings which allows for faster similarity search compared to traditional databases [23].

D. Retrieval with Re-ranking

The performance of RAG is highly dependent on the relevance and quality of the retrieved context. This work employs a cross-encoder re-ranker, a widely adopted semantic re-ranking method to ensure that the most relevant contexts are ranked first [24]. Unlike bi-encoders, which embed each chunk independently, cross-encoders process pairs of text to calculate the similarity between them. This approach, allows it to fully capture the interactions and relationships between the query and each chunk of context. We specifically use the *ms-marco-MiniLM-L-6-v2* model due to its balance between efficiency and performance, both essential for telecom tasks [24]

E. Extending the Context Window with SelfExtend

However, SLMs typically struggle to generalize effectively to input sequences longer than those encountered during training. This presents as a challenge during inference with long contexts. The context window of SLMs are often short. For example, Phi-2 has a context window of 2048 tokens. Semantic chunking does not guarantee a fixed chunk size, and therefore, might exceed the context window of the SLMs.

Furthermore, supplementing the context with tables from telecom documents is crucial given the significance of the information they hold. Therefore, in order to support enhanced performance, and to allow for future expansions, we utilize SelfExtend [13].

This method leverages the inherent capabilities of LLMs to handle extended contexts without the need for fine-tuning. SelfExtend achieves this by implementing a bi-level attention mechanism: grouped attention for capturing dependencies between distant tokens, and neighbor attention for focusing on adjacent tokens. These attentions are computed using the model's existing self-attention during inference. By making use of SelfExtend, we are able to extend Phi-2's context window to 8192 tokens.

F. The Generator: Fine-tuned Phi-2 with Multiple Contexts

A substantial portion of telecom key terms and special language is confined to specification documents and white papers, neither of which LMs are heavily trained on. We fine-tune Phi-2 to overcome this and to adapt the model to recognize telecom terminology.

It should also be emphasized that fine-tuning the model is not intended to expand its knowledge base, but rather to enhance its ability to discern important details within the context and respond in the correct format. To accommodate limited resources, gradient accumulation and LoRA are utilized for fine-tuning [25]. Gradient accumulation is a technique that allows the model to effectively handle large batch sizes by minimizing the memory needed for storing gradients. It does this by processing several small batches and accumulating the gradients from each batch before updating the weights rather than calculating and updating the weights after each batch.

Fine-tuning all parameters in a model is impractical in our domain due to resource constraints. Moreover, for smaller or specialized datasets, this approach risks overfitting and poor generalization, potentially yielding diminishing returns.

The concept underlying LoRA is that pre-trained LMs possess a low 'intrinsic dimension' [26]. That is, the model's essential information is concentrated in a smaller subspace, even if the overall parameter space is high-dimensional. LoRA harnesses this observation and focus the model updates on a smaller only a subset of the learning parameters. During fine-tuning, weight updates W_{new} are represented as

$$W_{new} = W_0 + \Delta W,\tag{1}$$

where $W_0 \in \mathcal{R}^{d \times k}$ are the initial pre-trained weights, and $\Delta W \in \mathcal{R}^{d \times k}$ represents the change in weights. Computing

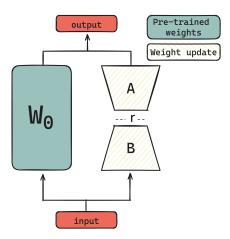


Figure 4. Schematic illustration of the Low-Rank Adaptation (LoRA) technique for efficient fine-tuning of neural networks with low-rank matrices (e.g., LMs).

 ΔW is computationally expensive. In LoRA, trainable parameters ΔW are expressed as a product of two low-rank matrices, $B \times A$, where $B \in \mathcal{R}^{d \times r}$ and $A \in \mathcal{R}^{r \times k}$, with rank $r \ll \min(d,k)$. Thus, with LoRA, weight updates W_{new} are computed as

$$W_{new} = W_0 + (\frac{\alpha}{r})B \times A, \tag{2}$$

where α is a scaling factor, reflecting how important the weight updates are to the initial pre-trained weights. As a result of this matrix decomposition, only $d \times r + r \times k$ parameters need to be updated.

G. Prompt Engineering

Prompt engineering is crucial for optimizing the performance of LLMs like Phi-2 in telecom applications. In this work, we supplement the inputted question with chunks of relevant context, and a set of instructions for the SLMs to follow, forming the prompt. In exemplary prompt is outlined in 3. This approach helps unify the format of the SLM's output, focus its attention on the appearance of critical terms, and encourage it to reply on the context rather than its prior knowledge. This prompt results in focused outputs, ensuring accurate and relevant responses for complex telecom inquiries.

IV. EXPERIMENTAL RESULTS

A. Settings

The RAG framework relies on around 550 3GPP documents up to Release 18. For fine-tuning, we utilize the TeleQnA dataset which contains 10,000 multiple-choice questions [11]. The testing set features another 2,000 multiple-choice questions that focus on 3GPP standards. For fine-tuning, we set the following parameters: weight decay of 0.01, batch size of 32, dropout rate of 0.05, and learning rate of 10^{-4} . For LoRA, we set the rank to 32 and the alpha to 64. For semantic chunking, there are two hyper-parameters, namely the breakpoint percentile threshold and the buffer size. The former represents the percentile of cosine dissimilarity that

must be exceeded between a group of sentences and the next to form a node. The latter determines the number of sentences to group together when evaluating semantic similarity. We set the breakpoint percentile threshold to 90, and the buffer size to 3.

For the vector retriever, for each query, we first retrieve 150 chunks, from which the re-ranker would return the top 15 most relevant chunks. The aforementioned hyper-parameters are not necessarily optimized, but rather set based on qualitative judgment and limited exploration. Fine-tuning the Phi-2 generator (with multiple contexts) relies on three retrieved-context chunks along with the set of instructions and the respective query. For reproducibility and reuse, our source code is made publicly available [27].

B. Results and Analysis

TABLE I: Accuracy comparison of our fine-tuned Phi-2 model against baseline models, both with and without retrieved context.

	Accu	Accuracy	
Model	w/o Context	w/ Context	
Phi-2	49.95%	71.35%	
gpt-4o	61.30%	69.30%	
finetuned Phi-2 (LoRA)	49.10%	80.30%	

In Table I, we benchmark our developed framework against three other solutions: base Phi-2 (2.7B), GPT-40 mini (8B), and GPT40 (1.76T). As expected, the base Phi-2 model shows poor performance of about 49.95% accuracy when tested on the dataset when no context is provided. Although the performance is notably improved to 71.35% when the base model is supplemented by retrieved context, it is still outperformed by our proposed model. Our fine-tuned RAG model is better aligned with the required task, enabling it to leverage the retrieved context better and produce more accurate recommendations.

Notably, our developed system also outperforms GPT-40 with and without context. Since GPT-40 is a significantly more generalized model, it is not well-aligned to the specialized

TABLE II: Performance comparison of various configurations of the fine-tuned Phi-2 model with RAG and additional components; the table uses the following acronyms: SE for SelfExtend, RR for Rerank, SC for Semantic Chunking, and MC for Multiple Context.

Model	Accuracy
FT Phi-2 + RAG	72.10%
FT Phi-2 + RAG + RR	77.35%
FT Phi-2 + RAG + RR + SC	41.20%
FT Phi-2 (MC) + RAG + RR + SC + SE	80.30%

task at hand. Qualitative analysis reveals that GPT-4o's *a priori* knowledge from other domains sometimes interferes with our domain-specific task, albeit when the relevant context is present in the prompt.

Table II represents a brief ablation study, where we analyze the impact of each added component on the predictive accuracy. The re-ranking algorithm has significant positive effects, adding 5% in accuracy. This underscores the importance of prioritizing retrieved-context chunks. When implemented on its own, semantic chunking degrades the performance. However, when paired with SelfExtend, the accuracy is significantly increased to 80.30% which is an increase of around 8% compared to the base model. Although semantic chunking produces more semantically coherent chunks (as shown in Fig. 2), it often results in longer chunks, and SelfExtend helps incorporate them into complete semantic units.

V. CONCLUSIONS AND FUTURE WORKS

The significance of SLMs has not gone unrecognized in the industry, as a shift from LLMs to SLMs can be evidenced by the recent directions of several big firms. They are especially pertinent in the telecom industry given the constraints often encountered in AI-driven telecom systems, such as the need for efficient deployment on edge devices with limited computational power. Coupled with RAG, SLMs have the potential to be a dominant tool in the industry. In this paper, we develop an fine-tuned Phi-2-based RAG system to serve as an oracle for telecommunication networks. The proposed system integrates semantic chunking and re-ranking to improve the relevance and accuracy of retrieved contexts. Moreover, we fine-tune the generator with a carefully designed prompt and retrieved contexts to adapt it to the problem domain. Additionally, we utilize SelfExtend to significantly extend the model's context window, enabling it to process longer sequences without fine-tuning. Experimental results show that our approach competes with larger state-of-theart LLMs and also offers significant efficiency advantages, making it suitable for deployment on edge devices. Our approach prioritizes transparency and explainability, offering a framework that allows for scrutiny, understanding, and further development in this field. Maintaining transparency through open-source models is particularly important in telecom, where it fosters trust and facilitates community contributions.

We believe the developed system can serve as a foundation for other downstream telecommunication tasks. Future work can also explore optimizing the embedding model, better integration of structured data such as tables and graphs, and further enhancements to the RAG framework to continue advancing the capabilities of LLMs in telecom applications.

REFERENCES

- [1] S. Bubeck *et al.*, "Sparks of artificial general intelligence: Early experiments with GPT-4," *CoRR*, vol. abs/2303.12712, 2023.
- [2] A. Maatouk, N. Piovesan, F. Ayed, A. D. Domenico, and M. Debbah, "Large language models for telecom: Forthcoming impact on the industry," *IEEE Commun. Mag*, pp. 1–7, 2024.

- [3] R. Zhang, H. Du, Y. Liu, D. Niyato, J. Kang, S. Sun, X. Shen, and H. V. Poor, "Interactive AI with retrieval-augmented generation for next generation networking," *IEEE Netw.*, pp. 1–1, 2024.
- [4] O. Erak, O. Alhussein, S. Naser, N. Alabbasi, D. Mi, and S. Muhaidat, "Large language model-driven curriculum design for mobile networks," *CoRR*, vol. abs/2405.18039, 2024.
- [5] Y. Yao, X. Xu, and Y. Liu, "Large language model unlearning," CoRR, vol. abs/2310.10683, 2024.
- [6] L. Huang, W. Yu, W. Ma, W. Zhong, Z. Feng, H. Wang, Q. Chen, W. Peng, X. Feng, B. Qin, and T. Liu, "A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions," *CoRR*, vol. abs/2311.05232, 2023.
- [7] N. Piovesan, A. D. Domenico, and F. Ayed, "Telecom language models: Must they be large?" CoRR, vol. abs/2403.04666, 2024.
- [8] M. Javaheripi and S. Bubeck, "Phi-2: The surprising power of small language models," https://www.microsoft.com/en-us/research/blog/ phi-2-the-surprising-power-of-small-language-models/, Dec. 2023, accessed: 2024-8-19.
- [9] G. Team et al., "Gemini: A family of highly capable multimodal models," CoRR, vol. abs/2312.11805, 2024.
- [10] H. Touvron et al., "Llama 2: Open foundation and fine-tuned chat models," CoRR, vol. abs/2307.09288, 2023.
- [11] A. Maatouk, F. Ayed, N. Piovesan, A. D. Domenico, M. Debbah, and Z.-Q. Luo, "TeleQnA: A benchmark dataset to assess large language models telecommunications knowledge," *CoRR*, vol. abs/2310.15051, 2023.
- [12] A.-L. Bornea, F. Ayed, A. D. Domenico, N. Piovesan, and A. Maatouk, "Telco-RAG: Navigating the challenges of retrieval-augmented language models for telecommunications," *CoRR*, vol. abs/2404.15939, 2024.
- [13] H. Jin, X. Han, J. Yang, Z. Jiang, Z. Liu, C.-Y. Chang, H. Chen, and X. Hu, "LLM Maybe LongLM: Self-Extend LLM Context Window Without Tuning," CoRR, vol. abs/2401.01325, 2024.
- [14] I. Karim, K. S. Mubasshir, M. M. Rahman, and E. Bertino, "SPEC5G: A dataset for 5G cellular network protocol analysis," in *Proc. IJCNLP-AACL*, 2023, pp. 20–38.
- [15] A. Karapantelakis, M. Thakur, A. Nikou, F. Moradi, C. Orlog, F. Gaim, H. Holm, D. D. Nimara, and V. Huang, "Using large language models to understand telecom standards," *CoRR*, vol. abs/2404.02929, 2024.
- [16] R. Nikbakht, M. Benzaghta, and G. Geraci, "TSpec-LLM: An opensource dataset for LLM understanding of 3GPP specifications," CoRR, vol. abs/2406.01768, 2024.
- [17] P. Gajjar and V. K. Shah, "ORAN-Bench-13K: An open source benchmark for assessing LLMs in open radio access networks," CoRR, vol. abs/2407.06245, 2024.
- [18] T. Ahmed, N. Piovesan, A. D. Domenico, and S. Choudhury, "Linguistic intelligence in large language models for telecommunications," *CoRR*, vol. abs/2402.15818, 2024.
- [19] L. Bariah, H. Zou, Q. Zhao, B. Mouhouche, F. Bader, and M. Debbah, "Understanding telecom language through large language models," in Proc. IEEE Globecom. 2023, pp. 6542–6547.
- Proc. IEEE Globecom, 2023, pp. 6542–6547.
 [20] J. Liu, "LlamaIndex," 2022. [Online]. Available: https://github.com/jerryjliu/llama_index
- [21] S. Xiao, Z. Liu, P. Zhang, N. Muennighoff, D. Lian, and J.-Y. Nie, "C-Pack: Packaged resources to advance general chinese embedding," *CoRR*, vol. abs/2309.07597, 2024.
- [22] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *Proc. ICML*, 2020, pp. 1597–1607.
- [23] Chroma, "ChromaDB," https://www.trychroma.com/, accessed: 08/05/2024.
- [24] Inferless, "Ms marco: ms-marco-minilm-l-6-v2," https://huggingface.co/ cross-encoder/ms-marco-MiniLM-L-6-v2, 2023, accessed: 2024-08-18.
- [25] E. J. Hu, yelong shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen, "LoRA: Low-rank adaptation of large language models," in *Proc. ICLR*, 2022.
- [26] A. Aghajanyan, S. Gupta, and L. Zettlemoyer, "Intrinsic dimensionality explains the effectiveness of language model fine-tuning," in *Proc. ACL-IJCNLP*, 2021, pp. 7319–7328.
- [27] N. Alabbasi and O. Erak, "Specializing large language models for telecom networks," https://github.com/Nouf-Alabbasi/oKUmura_AI_ Telecom_challenge, 2024.