

# Hybrid Retrieval-Augmented Generation Approach for LLMs Query Response Enhancement

Pouria Omrani\*<sup>‡</sup>, Alireza Hosseini<sup>†‡</sup>, Kiana Hooshanfar<sup>†‡</sup>, Zahra Ebrahimian<sup>†‡</sup>, Ramin Toosi<sup>†‡</sup>, Mohammad Ali Akhaee<sup>†</sup>

<sup>‡</sup>Faculty of Electrical Engineering, K. N. Toosi University of Technology, Tehran, Iran

<sup>†</sup>School of Electrical and Computer Engineering, College of Engineering, University of Tehran, Tehran, Iran

<sup>‡</sup>Adak Vira Iranian Rahjoo Company, Tehran, Iran

{pouria.omrani@ieec.org, arhosseini77@ut.ac.ir, k.hooshanfar@ut.ac.ir, z.ebrahimian@ut.ac.ir, r.toosi@ut.ac.ir, akhaee@ut.ac.ir}

**Abstract**— In the domain of Natural Language Processing (NLP), the integration of Large Language Models (LLMs) with Retrieval-Augmented Generation (RAG) represents a significant advancement towards enhancing the depth and relevance of model-generated responses. This paper introduces a novel hybrid RAG framework that synergizes the Sentence-Window and Parent-Child methodologies with an innovative re-ranking mechanism, aimed at optimizing the query response capabilities of LLMs. By leveraging external knowledge sources more effectively, the proposed method enriches LLM outputs with greater accuracy, relevance, and information fidelity. We subject our hybrid model to rigorous evaluation against benchmark datasets and metrics, demonstrating its superior performance over existing state-of-the-art RAG techniques. The results highlight our method's enhanced ability to generate responses that are not only contextually appropriate but also demonstrate a high degree of faithfulness to the source material, thereby setting a new standard for query response enhancement in LLMs. Our study underscores the potential of hybrid RAG models in refining the interaction between LLMs and external knowledge, paving the way for future research in the field of NLP.

**Keywords**— LLM, Generative AI, NLP, Retrieval Augmented Generation (RAG).

# رویکرد ترکیبی تولید تقویت شده با بازیابی برای بهبود پاسخ دهی پرسش ها در مدل های زبانی بزرگ

پوریا عمرانی<sup>۱</sup>، علیرضا حسینی<sup>۱</sup>، کیانا هوشانفر<sup>۱</sup>، زهرا ابراهیمیان<sup>۱</sup>، رامین طوسی<sup>۱</sup>، محمدعلی اخایی<sup>۲</sup>

<sup>۱</sup> دانشکده مهندسی برق، دانشگاه صنعتی خواجه نصیرالدین طوسی، تهران، ایران - شرکت آداک ویرا ایرانیان رهجو، تهران، ایران  
pouria.omrani@ieee.org\_ arhosseini77@ut.ac.ir\_ k.hooshanfar@ut.ac.ir\_ z.ebrahimian@ut.ac.ir\_ r.toosi@ut.ac.ir

<sup>۲</sup> گروه مهندسی برق و کامپیوتر، دانشکده مهندسی، دانشگاه تهران، تهران، ایران  
akhaee@ut.ac.ir

## چکیده

در حوزه پردازش زبان طبیعی (NLP)، ادغام مدل های زبانی بزرگ (LLMs) با تولید تقویت شده توسط بازیابی (RAG) نشان دهنده پیشرفت قابل توجهی در افزایش عمق و ارتباط پاسخ های تولید شده توسط مدل است. این مقاله چهارچوب RAG ترکیبی نوآورانه ای را معرفی می کند که روش های Parent-Child و Sentence-Window را با مکانیزم بازترتیب بندی تلفیق می کند، با هدف بهینه سازی قابلیت های پاسخ گویی به پرسش ها توسط LLM ها. با استفاده مؤثرتر از منابع دانش خارجی، روش پیشنهادی خروجی های LLM را با دقت، ارتباط و وفاداری اطلاعاتی بیشتری غنی سازی می کند. ما مدل ترکیبی خود را با استفاده از مجموعه داده ها و معیارهای مرجع مورد ارزیابی قرار می دهیم و عملکرد برتر آن را نسبت به تکنیک های RAG پیشرفته موجود نشان می دهیم. نتایج، توانایی بهبود یافته روش ما را در تولید پاسخ هایی که نه تنها مناسب با زمینه هستند، بلکه وفاداری بالایی به مواد منبع نیز نشان می دهند، برجسته می کند و بدین ترتیب یک استاندارد جدید برای بهبود پاسخ گویی در LLM ها ایجاد می کند. مطالعه ما پتانسیل مدل های RAG ترکیبی را در پالایش تعامل بین LLM ها و دانش خارجی مورد تأکید قرار می دهد و راه را برای تحقیقات آینده در زمینه NLP هموار می سازد.

**کلمات کلیدی:** مدل زبانی بزرگ، هوش مصنوعی مولد، پردازش زبان طبیعی، تولید تقویت شده توسط بازیابی