

Development of a Platform for Scalable Launch of Chatbots Based on RAG Approach

The project aims to develop a scalable platform for deploying chatbots that utilize the Retrieval-Augmented Generation (RAG) approach. The RAG approach combines retrieval of relevant documents with generative models to produce accurate and contextually appropriate responses. The platform will enable efficient scaling to handle multiple chatbots and high user demand ensuring reliability and performance across diverse applications.

1 Introduction

Table 1 shows a comparative analysis of various recent solutions for developing scalable platforms for launching chatbots based on the RAG approach. The table summarizes the strengths and weaknesses of each solution according to the project's quality criteria.

Table 1: Comparative analysis of solutions for scalable RAG-based chatbot platforms

Solution	Strengths	Weaknesses
LangChain framework [1]	Modular components for building RAG pipelines; simplifies integration with various LLMs and vector stores	May have limitations in customization for specific enterprise use cases; relatively new with evolving features
Haystack framework [2]	Open-source and highly customizable; supports various backends for retrieval and generation	Steeper learning curve; requires more setup and configuration
Meta's RAG implementation [3]	State-of-the-art performance; well-researched architecture; open-source code available	Complex implementation; requires significant computational resources; less focus on deployment scalability

Continued on next page

Table 1 – *Comparative analysis of solutions for scalable RAG-based chatbot platforms*

Solution	Strengths	Weaknesses
Microsoft’s DeepSpeed Chat [4]	Optimized for training large models efficiently; supports distributed training and inference	Primarily focused on training efficiency; may require adaptation for RAG and deployment; steep learning curve
OpenAI’s ChatGPT with Retrieval Plugin [5]	Easy to integrate and deploy; leverages powerful LLM capabilities; supports retrieval augmentation	Dependent on external API; limited customization
Using Kubernetes with ONNX Runtime [6]	Provides robust scaling and orchestration; optimized inference with ONNX; supports containerized deployments	Adds complexity in deployment and management; requires expertise in DevOps and ML model optimization

References

- [1] Harrison Chase. Langchain: A framework for developing applications powered by language models. 2022.
- [2] deepset. Haystack: End-to-end neural question answering at scale. 2020.
- [3] Patrick Lewis, Ethan Perez, Aleksandra Piktus, and et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Advances in Neural Information Processing Systems*, volume 33, pages 9459–9474, 2020.
- [4] Shaojie Ren, Shaden Kumar, Minjia Zhang, et al. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. *Proceedings of the 26th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, pages 1–14, 2021.
- [5] OpenAI. Retrieval plugin for chatgpt, 2023.
- [6] Microsoft. Scaling machine learning inference with onnx runtime and kubernetes, 2022.
- [7] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. Dense passage retrieval for open-domain question answering. *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6769–6781, 2020.
- [8] Gautier Izacard and Edouard Grave. Leveraging passage retrieval with generative models for open domain question answering. *arXiv preprint arXiv:2007.01282*, 2021.

- [9] Ben Wang and Aran Komatsuzaki. Gpt-j: An open-source autoregressive language model. 2022.
- [10] Siddharth Gudibande, Yasaman Khazaeni, and et al. False negatives in retrieval-based chatbots. *arXiv preprint arXiv:2104.07669*, 2021.
- [11] Shuguang Sun, Zhiwei Cao, Hong Zhu, and Jinhui Zhao. A survey of optimization methods for deep learning models training. *Journal of Systems Architecture*, 117:102112, 2021.
- [12] Pranav Rajpurkar, Robin Jia, and Percy Liang. Know what you don’t know: Unanswerable questions for squad. *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, pages 784–789, 2018.
- [13] Tianyu Gao, Xiao Liu, Shunzhong Zheng, Hao Zhu, and Pengcheng Li. Rethinkdenoising: Efficient and accurate training of large transformer models with denoising objectives. *arXiv preprint arXiv:2106.16241*, 2021.
- [14] Yida Feng, Ziqi Yang, and et al. A survey of deep learning approaches for document retrieval. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, pages 2107–2120, 2020.
- [15] Yang Liu, Ming Zhou, et al. Is chatgpt a good translator? a preliminary study. *arXiv preprint arXiv:2301.08745*, 2023.
- [16] Urvashi Khandelwal, Angela Fan, Dan Jurafsky, and Luke Zettlemoyer. knn-lm: Improving language models with nearest neighbor retrieval. *Transactions of the Association for Computational Linguistics*, 9:727–742, 2021.
- [17] Hugo Touvron et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [18] Stephen Roller, Emily Dinan, Naman Goyal, Da Ju, Mary Williamson, Yinhan Liu, Jing Xu, Myle Ott, Kurt Shuster, Eric M. Smith, et al. Recipes for building an open-domain chatbot. *arXiv preprint arXiv:2004.13637*, 2021.
- [19] Jinhua Su, Dongdong Ren, et al. One teacher is enough? pre-trained language model distillation for neural machine translation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*, pages 136–146, 2022.
- [20] Li Yang, Yan Wang, et al. Model compression with two-stage multi-teacher knowledge distillation for web question answering system. *IEEE Access*, 8:173334–173347, 2020.