

Developing a Benchmarking Framework for RAG Systems

Feedback Bericht

von Jan Albrecht

I. ZUSAMMENFASSUNG DER PRÄSENTATION

Ich habe in der Präsentation das Framework für RAG-Systeme vorgestellt und bin auf die Besonderheiten eingegangen, die bei der Umsetzung relevant sind. Darunter zählt sowohl der modulare Architektur eines RAG-Systems und die verschiedenen Aspekte der Evaluation. Ersteres fokussiert sich darauf, dass möglichst viele RAG-Varianten und IRA-Methoden (Iterativ, Rekursiv, Adaptiv) mit dem Framework abgebildet werden können. Die Evaluation kann in End-to-End Evaluation und Komponenten-Evaluation unterschieden werden. Während End-to-End Evaluation für den Endnutzer relevanter ist, kann Komponentenevaluation zur Fehleranalyse und Hyperparameter Tuning verwendet werden.

II. ZUSAMMENFASSUNG DER FEEDBACKS

Das Feedback war positiv und der Plan zur Umsetzung stimmig. Es wurde eingeworfen, den Fokus von gleichermaßen Framework und Benchmarks mehr auf das Framework zu rejustieren. Es kam die Befürchtung auf, dass mit zunehmender Arbeit für die Benchmarks, die eigentliche Umsetzung des Frameworks zu kurz kommt. Es soll ein hoher Fokus auf die Re-Konfigurierbarkeit gesetzt werden. Neben der potenziellen Evaluation durch typisch Task-Based Metriken wie Precision, Recall und F1-Score, sollten auch Hardware Metriken wie RAM, Inferenzzeit und CPU-Zeit verwendet werden, da diese in der Praxis eine wichtige Rolle spielen. Es wurde weiterhin angemerkt, dass das Paper der SWS-Abteilung "A Methodology for Evaluating RAG Systems: A Case Study On Configuration Dependency Validation" zu lesen und zu analysieren, wie der entwickelte Blueprint in der Umsetzung benutzt werden kann. Außerdem wurde darauf hingewiesen, dass SLURM für die Experimente eingesetzt werden kann.

III. ANGEPASSTE UMSETZUNG

Für die Anpassung der Umsetzung muss die Implementierung des Frameworks und das Evaluieren durch einen Software Engineering Task angepasst werden.

A. Angepasste Umsetzung des Frameworks

Ein Framework, welches nach dem Blueprint des durch die SWS-Abteilung veröffentlichten Paper aufgebaut ist, muss folgende Voraussetzungen haben:

- Easy integration of relevant context resources
- High coverage of RAG variants (including IRA methods and common special cases)
- High coverage of evaluation metrics (e.g. Precision, Recall, F1-Score)
- An LLM or Naive-RAG baseline
- Splitting the dataset into Validation and Holdout-Testset, with the Validation set for failure analysis and hyperparameter tuning and the Holdout-Testset for evaluating the generalization of the system.

Weiterhin wurden die Anforderungen gestellt, dass das Framework einfach zu rekonfigurieren sein soll, sodass die Evaluationen nach Fehleranalyse einfach zu replizieren sind. Replizierbarkeit und Interne Gültigkeit der Experimente sind Voraussetzungen für das Framework.

Das Framework wird neben Task-spezifischen Metriken auch Hardware-Metriken nutzen. Darunter zählen zum Beispiel Inferenzzeit, RAM und CPU-Zeit. Daher ist es wichtig, dass die Evaluierungen nicht parallel gestartet werden, sondern isoliert von anderen Prozessen und sequenziell auf einer dedizierten Hardware.

1) *Analyse bestehender Frameworks und Bibliotheken:*
Für die Neuausrichtung des Fokus, werde ich noch sicherstellen, dass ein solches Framework noch nicht existiert.

a) *AutoRAG[7]:* AutoRAG ist ein Tool für Hyperparameter-Tuning für RAG-Systeme. Der Nutzer erstellt eine Konfigurationsdatei, welche die sogenannten Nodes (Retriever, Generator, usw.) definiert und die Hyperparameter (e.g. Top-K, Models, Prompts, etc.) angibt. Danach läuft die Analyse und es wird die Konfiguration des besten Systems angegeben.

Das Tool nutzt nur End-to-End Evaluation und ist auch schwierig anpassbar. Laut den Entwicklern im Discord sind auch fortgeschrittene Methoden (IRA) derzeit nicht möglich. Es bietet auch keine Möglichkeit ein Baseline-LLM zu evaluieren.

b) *Llama-Index[8] & Langchain[2]:* Llama-Index und Langchain sind umfangreiche Bibliotheken für die Entwicklung von RAG-Pipelines. Sie bieten jedoch keine zeiteffiziente Möglichkeit zum Benchmarken vieler RAG-Systeme.

c) *RaLLe*[3]: RaLLe ist eine Bibliothek zur zeiteffizienten Entwicklung von RAG-Systemen. Der letzte Commit ist über ein Jahr alt und die Bibliothek scheint sehr restriktiv. Unter anderem ist es nur möglich .tsv-Dateien mit einer vorgegebenen Struktur zu nutzen. Fortgeschrittene Methoden sind nicht möglich.

d) *FlashRAG*[6]: FlashRAG ist ein Python Toolkit, das eine ähnliche Motivation hat, wie ich sie hier versuche umzusetzen. Es bietet eine Konfigurationsdatei als Start, es hat eine Auswahl von Task-spezifischen Metriken und eine Vielzahl an RAG-Komponenten. Es lassen sich viele RAG-Varianten und einige Spezialfälle damit umsetzen (IRA, Self-RAG[1], FLARE[5], IRCOT[13]).

e) *RAGLAB*[14]: RAGLAB ist das neueste Framework und bietet ähnliche Funktionalitäten wie dessen Vorgänger. Es führt auch "Fair Comparison" ein, einen faireren Vergleich für LLM's und Retriever zum Beispiel das Entfernen von Special Tokens im Training oder dem Seedsetzen in Retrieval und Generator Pipelines.

f) *Haystack*[9] & *FastRAG*[4]: Haystack ist ähnlich wie Llama-Index und LangChain eine umfangreiche Bibliothek für die Entwicklung von RAG-Pipelines. Sie bietet Pipeline Definition mittels YAML oder JSON Dateien und ist stark modular aufgebaut. Es lassen sich IRA Methoden damit entwickeln. FastRAG ist kompatibel mit Haystack und hat eine Vielzahl an RAG-Komponenten entwickelt, mit denen es auch möglich ist Spezialfälle aus Papers zu entwickeln. Darunter zählen zum Beispiel REPLUG[12], PLAID[10], ColBERT[11] oder quantized LLM's.

2) *Einordnung bestehender Frameworks und Bibliotheken*: Bei jeder Analyse wurden zwei Punkte betrachtet:

- Bietet das untersuchte Framework bereits die Lösung für das in der Masterarbeit untersuchte Thema?
- Kann ich das Framework als Fundament für mein Framework nutzen?

Keines der Frameworks bietet eine vollständige Pipeline für das Entwickeln von RAG-Varianten durch iteratives Re-Konfigurieren und systematischer replizierbarer Evaluierung mit Holdout-Testset. Eine vordefinierte Untersuchung durch Einführung einer LLM-Baseline, ob ein RAG-System eine Verbesserung darstellt ist bei keines der Frameworks vorhanden. Weiterhin misst keines der Frameworks den Einfluss der Komplexität auf Metriken wie Inferenzzeit, RAM-Usage oder CPU-Time. Damit ist der erste Punkt beantwortet. Für den zweiten Punkt können einige Frameworks herausgefiltert werden. RaLLe und AutoRAG sind sehr restriktiv und es lassen sich nur eine begrenzte Anzahl von RAG-Varianten umsetzen. Llama-Index und LangChain sind ebenfalls restriktiv und nicht zeiteffizient genug, um viele ver-

schiedene RAG-Varianten zu testen. RAGLAB bietet keine Pipeline-Definition durch Konfigurationsdateien wie YAML oder JSON an. Dies könnte die Entwicklungsgeschwindigkeit stark reduzieren.

FlashRAG und Haystack+FastRAG wären mögliche Kandidaten als Basis für das geplante Framework. Ich werde Haystack und FastRAG nutzen, da Haystack stark modular aufgebaut ist und in Zukunft eine UI bieten wird um RAG-Systeme per Drag&Drop zu erstellen. Die daraus resultierenden RAG-Varianten können in YAML's exportiert werden und für mein Framework benutzt werden.

B. Angepasste Umsetzung der Benchmarks

Der Fokus auf die Umsetzung des Frameworks bedeutet, dass der Umfang der Benchmark-Experimente reduziert werden muss. Anstelle eines eigenen Experiments werde ich mich auf bereits bestehenden Experimente fokussieren. Mögliche Experimente sind Configuration Validation mit dem CTest Datensatz und Vulnerability Detection mit den Datensätzen SySeVR, diversful oder bifi.

REFERENCES

- [1] Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-rag: Learning to retrieve, generate, and critique through self-reflection. <http://arxiv.org/pdf/2310.11511v1>.
- [2] Harrison Chase. LangChain. <https://github.com/langchain-ai/langchain>, October 2022.
- [3] Yasuto Hoshi, Daisuke Miyashita, Youyang Ng, Kento Tatsuno, Yasuhiro Morioka, Osamu Torii, and Jun Deguchi. Ralle: A framework for developing and evaluating retrieval-augmented large language models. <https://arxiv.org/abs/2308.10633>, 2023.
- [4] Peter Izsak, Moshe Berchansky, Daniel Fleischer, and Ronen Lapedon. fastRAG: Efficient Retrieval Augmentation and Generation Framework. <https://github.com/IntelLabs/fastrag>, February 2023.
- [5] Zhengbao Jiang, Frank F. Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. Active retrieval augmented generation. <http://arxiv.org/pdf/2305.06983v2>.
- [6] Jiajie Jin, Yutao Zhu, Xinyu Yang, Chenghao Zhang, and Zhicheng Dou. Flashrag: A modular toolkit for efficient retrieval-augmented generation research. *CoRR*, abs/2405.13576, 2024.
- [7] Jeffrey Kim and Bobb Kim. AutoRAG. <https://github.com/Marker-Inc-Korea/AutoRAG>, February 2024.
- [8] Jerry Liu. LlamaIndex. https://github.com/jerryjliu/llama_index, 11 2022.
- [9] Malte Pietsch, Timo Möller, Bogdan Kostic, Julian Risch, Massimiliano Pippi, Mayank Jobanputra, Sara Zanzottera, Silvano Cerza, Vladimir Blagojevic, Thomas Stadelmann, Tanay Soni, and Sebastian Lee. Haystack: the end-to-end NLP framework for pragmatic builders. <https://github.com/deepset-ai/haystack>, November 2019.
- [10] Keshav Santhanam, Omar Khattab, Christopher Potts, and Matei Zaharia. Plaid: An efficient engine for late interaction retrieval, 2022.
- [11] Keshav Santhanam, Omar Khattab, Jon Saad-Falcon, Christopher Potts, and Matei Zaharia. Colbertv2: Effective and efficient retrieval via lightweight late interaction, 2022.
- [12] Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettlemoyer, and Wen tau Yih. Replug: Retrieval-augmented black-box language models, 2023.
- [13] Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions. <http://arxiv.org/pdf/2212.10509v2>.

- [14] Xuanwang Zhang, Yunze Song, Yidong Wang, Shuyun Tang, et al. RAGLAB: A modular and research-oriented unified framework for retrieval-augmented generation. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Association for Computational Linguistics, December 2024.