# Literature Review: Active Retrieval Augmented Generation (Jiang et al., 2023)

Elias Ulf Hörnberg | 384928 | elias.hornberg@epfl.ch
AB-EH-ME

## 1 Summary

The document elaborates on a method known as Forward-Looking Active Retrieval augmented generation (FLARE), aimed at enhancing the capabilities of large language models (LMs). Traditional LMs, despite their proficiency in language understanding and generation, often struggle with producing factually accurate outputs due to their inclination to generate content based on learned patterns rather than factual correctness. Prior methods typically involve a single retrieval step before the generation process, which restricts their ability to adapt as the context develops during longer text generation tasks.

FLARE addresses these limitations by introducing a dynamic, continuous retrieval mechanism throughout the generation process. Unlike conventional models that only retrieve once at the beginning, FLARE actively decides when and what information to retrieve next based on both the current context and the predicted needs of the upcoming content. This method integrates the retrieval process more deeply into the generation workflow, allowing the model to adjust its knowledge base continually and ensure the relevance and accuracy of the information being generated.

The document claims that FLARE achieves superior performance across a variety of complex tasks that require extensive and accurate text generation, such as multihop question answering, commonsense reasoning, long-form question answering, and open-domain summarization. These results are presented as part of a comprehensive evaluation that benchmarks FLARE against existing models that either do not use retrieval enhancements or use less dynamic, single-time retrieval methods.

## 2 Strengths

Innovative Integration of Retrieval and Generation: FLARE's integration of a continuous retrieval mechanism within the generative process allows it to dynamically adjust to the evolving information needs. The paper emphasizes this by noting that *"FLARE iteratively generates a temporary next sentence, uses it as the query to retrieve relevant documents if it contains low-probability tokens, and regenerates the next sentence"*. This continual adaptation is a significant enhancement over static retrieval methods.

Improved Factual Accuracy and Relevance: By continuously updating its knowledge base during the generation process, FLARE ensures that the content it generates is not only relevant but also factually correct. This is especially important in applications such as academic writing, technical documentation, and other professional settings where accuracy is critical.

Broad Applicability and Scalability: The method's successful application across diverse NLP tasks demonstrates its scalability and adaptability. This versatility is crucial for the development of generalized AI systems capable of performing well across various domains without needing task-specific adjustments.

Enhanced Long-form Text Generation:The paper specifically points out FLARE's proficiency in long-form text generation, which has been a challenging area for many LMs, noting its ability to maintain coherence and integrate relevant information throughout extended passages.

## 3 Weaknesses

Increased Computational Requirements: The ongoing retrieval mechanism, while beneficial, also introduces significant computational overheads. This could limit FLARE's use in environments where computing resources are limited or rapid response times are necessary.

Dependency on External Knowledge Sources: FLARE's performance heavily relies on the quality of external databases. The paper indicates potential

limitations by stating that *"the effectiveness of our method could be significantly reduced in scenarios where the external knowledge bases are limited"*.

Risk of Overfitting and Reduced Generalizability: The complexity of FLARE's architecture might lead to overfitting on specific tasks or datasets. The paper discusses the challenges of generalizability, noting that *"FLARE did not provide significant gains over not using retrieval in some tested scenarios"*, reflecting potential overfitting issues.

Potential for Information Overload: The active retrieval process could lead to excessive information retrieval, cluttering the generated content with unnecessary details. The paper subtly hints at this risk, suggesting that managing the balance between detail and conciseness remains a challenge.

# References

Zhengbao Jiang, Frank F. Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. Active retrieval augmented generation.