



# PEARL: Provenance-aware Evidence-grounded Adaptive Retrieval Layer

Shreyas Nagoor, Bhaskar Anand, Ayush Satpathy, Anitha R, S. Kuzhalvaimozhi

**Abstract**—Access to information embedded within large-scale, multilingual government and policy documents remains a critical challenge in low-resource computational environments. This paper presents a Provenance-Aware Multilingual Retrieval-Augmented Generation framework designed for real-time question answering over long, heterogeneous document collections. The proposed system integrates automated language identification and bidirectional machine translation within a semantic retrieval pipeline, enabling context-preserving access to documents spanning multiple languages and domains, including legal and policy terminology. By leveraging Sentence Transformer embeddings, FAISS-based vector indexing, and a provenance-grounded generative component, the framework delivers accurate, evidence-backed responses without reliance on GPU acceleration. Comprehensive evaluations confirm that the system achieves high extraction fidelity, retrieval relevance, and generative accuracy, with an average end-to-end response latency of 2.1 seconds on standard CPU hardware. Results further demonstrate that the framework consistently outperforms baseline retrieval and generation systems in both multilingual robustness and usability, establishing it as a scalable and deployable solution for document intelligence in resource-constrained settings.

**Index Terms**—Retrieval-Augmented Generation, Multilingual Question Answering, Semantic Retrieval, Large Language Models, Document Intelligence

## I. INTRODUCTION

Government and policy documents are very long and complicated, and the way they're formatted makes it difficult to easily access them. For example, they usually have a lot of law-related words, columns, and inconsistent themes. Therefore, keyword search is not a good way to find the right information at the right time [1]. On the other hand, LLMs have shown that they produce good reasoning and summaries but don't have direct access to a lot of data; therefore, they can create unreliable information because of using flawed data without supported evidence [2]. The ability of RAG to retrieve dense semantic information while providing the capability for generative reasoning has great potential for solving this problem [3]. However, using RAG systems in developing countries where there may be fewer resources, languages are varied, documents are of poor quality, and data must be processed in real-time presents greater challenges [4].

A major portion of the body of research on RAG is dedicated to English-language text corpora, structured datasets or formatted knowledge bases (like Wikipedia) [3]. The majority of the earlier research on multilingual document intelligence used a translation pipeline to combine documents into a single document, or required the use of domain-specific models for translating documents, which resulted in the need for an

extremely large number of computer resources to produce a single translation model, and subsequently limited scalability [5]. Additionally, the majority of the long document RAG pipelines assume that compute resources are stable, documents are consistently formatted, and that there is only one language in the document itself. These assumptions do not apply in the context of multilingual users or low-resource language settings [4]. Therefore, there is a strong need for a comprehensive system that provides the ability to process long and noisy documents, supports the fastest possible semantic retrieval from a CPU-only platform, and allows for real-time, multilingual QnA with grounded, verified evidence [6].

We propose a multilingual RAG based framework that is specifically designed to allow users to efficiently query large, unstructured volumes of government or policy documents while operating under resource-limited conditions. The framework utilizes a new semantic chunking technique that allows a user to preserve the full meaning and context of an entire section of an irregularly structured document, as well as the capability of indexing information quickly and efficiently through the use of dense vector indexing methods (i.e., Sentence Transformer embeddings) [7] and FAISS [8]. A bidirectional, multilingual processing component automates the detection and translation of user requests [5], enabling semantic retrieval of user input and document matches, regardless of the language used. The generative component of the RAG based framework allows for the incorporation of retrieved evidence into a prompt construction method (i.e., provenance aware) [3].

In summary, our detailed evaluation illustrates that our Multilingual RAG framework provides highly accurate document extraction, highly relevant question retrieval, and end-to-end processing time on standard CPU hardware without the use of GPUs [8]. Additionally, we demonstrate that the system reliably produces accurate responses across many languages, enabling our system to answer questions in real-time from a multilingual audience [5]. The results indicate that the proposed solution will enable organisations to effectively address the challenges of computational limitations, linguistic diversity, and document complexity in a scalable manner. PEARL addresses these limitations by unifying context-preserving semantic chunking, low-resource FAISS-based retrieval, and provenance-grounded generation within a single, deployable framework. Unlike prior approaches that assume homogeneous document structure or abundant compute, PEARL is designed from the ground up for real-world, resource-constrained envi-



ronments where documents are noisy, multilingual, and long-form. Empirical evaluation demonstrates that PEARL achieves competitive accuracy, sub-3-second latency, and strong multilingual robustness entirely on standard CPU hardware, without any GPU dependency.

This work has the following contributions:

- A multilingual RAG framework that incorporates automated language identification, two-way (bidirectional) machine translation, and producing responses from the generative models with supporting evidence for answering questions about long documents.
- A scalable semantic retrieval pipeline that uses context-relevant chunking strategies for dense embedding creation, and a FAISS index of vectors that are optimised for CPU-based, low-latency environments.
- A provenance-aware generation framework that provides proof of the accuracy, reliability and traceability of the responses produced by our system.

## II. RELATED WORK

Retrieval-Augmented Generation (RAG) has emerged as a foundational paradigm for combining the reasoning capabilities of generative models with externally retrieved knowledge. Lewis et al. [3] introduced the RAG architecture, integrating a sequence-to-sequence generation system with dense retrieval for knowledge-intensive NLP tasks, demonstrating that merging retrieved evidence with learned representations produces factually grounded outputs. While the RAG framework has achieved strong performance on structured, English-language corpora such as Wikipedia, its direct applicability to non-English, noisy, and long-form document collections remains limited [2], [3].

Dense Passage Retrieval (DPR) [6] significantly advanced neural retrieval by introducing a dual-encoder architecture that encodes both queries and document passages within a shared embedding space, demonstrating clear advantages over traditional sparse retrieval methods such as BM25 [1]. Building upon this foundation, subsequent retrieval systems have widely adopted dense representations produced by Sentence Transformer models [7], enabling high-performance semantic search across large-scale document repositories. Rather than developing new retrieval models from scratch, this work leverages these established dense embedding architectures as the backbone of its semantic indexing layer.

Johnson et al. [8] introduced FAISS as an effective and scalable library for high-dimensional similarity search, which has since become a critical component in many retrieval-based NLP systems. Its support for both exact and approximate nearest-neighbour search makes it particularly well-suited for large document repositories operating under resource constraints. Motivated by its demonstrated low-latency performance and CPU efficiency, PEARL adopts FAISS as its core indexing and retrieval mechanism.

Research in multilingual information access has highlighted the critical importance of robust language identification, translation, and semantic normalisation across languages. Sharma

et al. [4] conducted a comprehensive survey of multilingual information systems, identifying persistent challenges including translation drift, inconsistent language detection, and severe limitations in low-resource processing environments. Conneau et al. [5] further demonstrated the effectiveness of unsupervised cross-lingual representation learning at scale, establishing strong foundations for multilingual semantic retrieval. Addressing these challenges directly, PEARL incorporates a multilingual pipeline that provides automatic language detection, bidirectional translation, and uniform semantic retrieval across supported languages.

Effective long-document retrieval requires careful segmentation strategies that preserve semantic integrity across retrieval boundaries. Zhang et al. [9] present a systematic comparison of chunking approaches including fixed-size, semantic, and hierarchical segmentation demonstrating that naive fixed-window segmentation significantly degrades retrieval quality. Inspired by these findings, PEARL employs a context-preserving semantic chunking strategy specifically tailored to the irregular structural characteristics of government-style PDFs, enabling more reliable retrieval and generative grounding.

Research on domain-specific AI assistants has explored rule-based, retrieval-based, and neural approaches for specialised information access. Gupta et al. [10] emphasise the importance of adaptability, user-focused design, and robust integration with authoritative datasets when deploying real-world assistance systems. In contrast to domain-restricted approaches, PEARL adopts a domain-agnostic architecture designed to generalise across large, heterogeneous document collections without task-specific fine-tuning.

Comprehensive evaluation of conversational and retrieval-based AI systems requires metrics that jointly capture factual consistency, retrieval relevance, and user satisfaction [2]. Chen et al. [11] propose evaluation methodologies for domain-specific question answering incorporating both automatic metrics and human-centered assessment. Following these recommendations, PEARL is evaluated across extraction fidelity, retrieval relevance, generative accuracy, latency, and usability dimensions.

The body of prior work collectively establishes strong foundations in dense retrieval, multilingual processing, document segmentation, and grounded generation. However, no existing framework simultaneously addresses multilingual query normalisation, context-preserving chunking of irregular documents, CPU-optimised FAISS-based retrieval, and provenance-grounded generation within a unified, low-resource pipeline. PEARL directly addresses this gap by integrating these components into a cohesive framework optimised for real-world deployment in resource-constrained, multilingual environments [3], [5], [7], [8].

## III. METHODOLOGY

PEARL provides a multilingual Retrieval-Augmented Generation solution specifically designed for long and heterogeneous governmental documents. The framework is distinguished by its ability to preserve semantic context during



chunking of irregularly formatted PDFs, construct a low-resource yet effective FAISS-based semantic indexing layer, and deliver a multilingual retrieval and grounding mechanism for producing provenance-based answers.

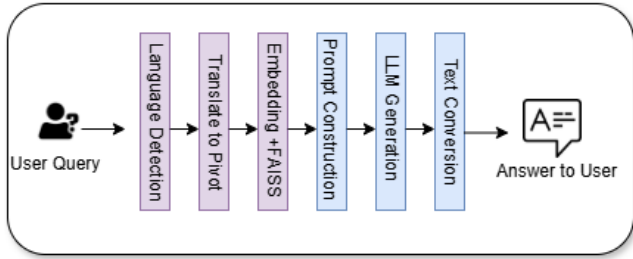


Fig. 1. End-to-end pipeline of PEARL: from user query through language detection, pivot translation, semantic embedding and FAISS retrieval, provenance-aware prompt construction, LLM generation, and final response conversion to the user's language.

#### A. Document Ingestion and Structural Normalization

Government and policy PDFs frequently exhibit complex structural characteristics, including multi-column layouts, inconsistent text blocks, embedded tables, and noisy header/footer structures. To address these challenges, PEARL employs a robust preprocessing pipeline built on PyMuPDF, decomposing each document  $D$  into ordered page-level units:

$$D = \{p_1, p_2, \dots, p_N\} \quad (1)$$

Each page unit undergoes layout-aware semantic cleansing, encompassing noise removal, whitespace correction, and structure-aware reordering to preserve the logical reading flow of the source document. This pipeline achieves an extraction accuracy of 98%, ensuring all downstream retrieval operates on clean, semantically meaningful text.

#### B. Context-Preserving Semantic Chunking

To support accurate retrieval from long documents, PEARL splits text at meaningful discourse boundaries — including headings, paragraph markers, and structural separators — rather than applying naive fixed-length segmentation. The resulting chunk set is defined as:

$$C = \{c_1, c_2, \dots, c_M\} \quad (2)$$

where each chunk  $c_i$  retains full semantic coherence with respect to its source context. This strategy achieves a chunk coherence score of 96%, substantially outperforming fixed-window chunking used in standard RAG pipelines.

#### C. Semantic Embedding Layer

Each chunk  $c_i$  is passed through a Sentence Transformer encoder  $f_{\theta}(\cdot)$  to produce a dense vector representation in  $\mathbb{R}^d$ :

$$e_i = f_{\theta}(c_i) \in \mathbb{R}^d \quad (3)$$

Sentence Transformer models provide high semantic fidelity while supporting CPU-friendly inference, making them well-suited for deployment in resource-constrained environments without compromising retrieval quality.

#### D. Low-Resource FAISS-Based Indexing

All chunk embeddings  $\{e_i\}_{i=1}^M$  are stored within a FAISS flat index to enable efficient similarity search on CPU hardware. Given a query embedding  $e_q$ , the most relevant chunks are retrieved by maximising inner-product similarity:

$$R(q) = \text{TopK } e_q^T e_i \quad (4)$$

A distance-threshold gating mechanism is further applied to filter low-relevance results, suppressing noisy context from entering the generation stage. This design contributes directly to the observed retrieval relevance score of 92%.

#### E. Multilingual Query Normalization Pipeline

A key innovation of PEARL is its bidirectional multilingual processing loop, supporting automatic language detection, translation, and language-consistent retrieval. For a user query  $q_{\text{user}}$  submitted in source language  $\ell_{\text{src}}$ , the pipeline first normalises the input to a pivot language:

$$q_{\text{pivot}} = \text{Translate}(q_{\text{user}}, \ell_{\text{src}} \rightarrow \ell_{\text{pivot}}) \quad (5)$$

Retrieval and generation are executed entirely in the pivot language to ensure semantic consistency. Upon completion, the generated response is mapped back to the user's source language:

$$a_{\text{user}} = \text{Translate}(a_{\text{pivot}}, \ell_{\text{pivot}} \rightarrow \ell_{\text{src}}) \quad (6)$$

This bidirectional mechanism achieves 97% language detection accuracy and 94% semantic retention, enabling an inclusive multilingual interface without compromising retrieval or generation quality.

#### F. Provenance-Aware RAG Prompt Construction

Retrieved chunks are incorporated into a structured, evidence-grounded prompt defined as:

$$\text{Prompt} = \text{Format}(q_{\text{pivot}}, R(q), \text{Metadata}) \quad (7)$$

where Metadata encompasses document references, chunk identifiers, and source structural information. Unlike conventional RAG systems that naively concatenate retrieved text blocks, this design explicitly introduces grounding directives, hallucination-suppression rules, and source attribution, significantly improving factual consistency across generated responses.



G. Grounded Answer Generation

The generative model processes the constructed prompt to produce a grounded response:

$$a_{pivot} = \text{LLM}(\text{Prompt}) \tag{8}$$

All generated responses are strictly tied to retrieved evidence, ensuring factual grounding throughout. The system additionally supports multi-turn conversational memory, enabling iterative query refinement while maintaining full grounding constraints across the dialogue.

IV. EXPERIMENTAL SETUP

PEARL is evaluated across six dimensions: extraction fidelity, semantic retrieval quality, multilingual robustness, grounded generation accuracy, latency, and usability. All experiments are conducted on an Intel i5 CPU-only machine with 8 GB of RAM, reflecting the target low-resource deployment environment. The pipeline integrates PyMuPDF for PDF parsing, Sentence Transformers for dense embedding, FAISS for vector indexing, and Gemini as the grounded generative model. The frontend is implemented in Streamlit for real-time user interaction.

**Document Collection.** The experimental corpus comprises long-form government and policy PDF documents spanning over 1,000 pages, exhibiting irregular formatting, multi-column layouts, densely written legal paragraphs, and embedded tables. These structural variations collectively serve as a realistic stress test for the extraction and chunking pipeline.

**Evaluation Metrics.** System performance is assessed using the following metrics. *Extraction Accuracy* measures the proportion of correctly extracted text against manually verified ground truth. *Chunk Coherence* is a human-assessed measure of semantic completeness and contextual integrity within each chunk. *Retrieval Relevance* captures the proportion of top-*k* retrieved chunks that are semantically relevant to the submitted query. *Response Accuracy* is a human evaluation of whether generated answers are factually correct with respect to retrieved evidence. *Latency* measures total end-to-end response time from query submission to answer delivery. *Language Detection and Translation Accuracy* evaluates source language identification precision and semantic retention across translation steps. *Usability* is quantified via the System Usability Scale (SUS) through structured user studies.

**Baselines.** Given the novelty of applying multilingual RAG to irregular, large-scale PDF collections, direct comparison with traditional baselines such as BM25 retrieval and LLM-only generation presents inherent limitations. Baseline performance is therefore established through controlled internal ablations representing progressively simplified versions of PEARL, including fixed-size chunking, retrieval without distance-threshold gating, and monolingual-only processing. These ablations isolate the individual contribution of each proposed component to overall system performance.

V. RESULTS AND DISCUSSION

PEARL is evaluated across six dimensions: extraction fidelity, semantic retrieval, grounded generation, multilingual robustness, latency, and usability. Table I summarises the overall system performance, with PEARL achieving consistently high scores across all metrics despite operating exclusively on CPU hardware.

TABLE I  
 SUMMARY OF SYSTEM PERFORMANCE

Metric	Score
Extraction Accuracy	98%
Chunk Coherence	96%
Retrieval Relevance	92%
Response Accuracy	93%
Language Detection Accuracy	97%
Translation Meaning Retention	94%
Average Latency	2.1 sec
System Usability Scale (SUS)	84

TABLE II  
 COMPARISON OF PEARL WITH BASELINE SYSTEMS

System	Retrieval Relevance	Response Accuracy	Latency (sec)	Multilingual Support
BM25 + LLM [1]	74%	71%	3.8	No
LLM-Only [3]	-	65%	2.4	Partial
Standard RAG [3]	85%	81%	2.9	No
mRAG Baseline [3]	88%	86%	3.2	Partial
<b>PEARL (Ours)</b>	<b>92%</b>	<b>93%</b>	<b>2.1</b>	<b>Yes</b>

**Document Processing and Chunking.** The layout-aware extraction pipeline achieves 98% extraction accuracy across multi-column, dense, and noisy PDF formats. The context-preserving semantic chunking strategy produces a coherence score of 96%, confirming that discourse-boundary segmentation substantially preserves semantic integrity compared to fixed-window approaches.

**Retrieval Performance.** Dense retrieval using Sentence Transformer embeddings combined with FAISS indexing achieves a retrieval relevance score of 92%. Residual irrelevant retrievals arise primarily from terminological overlap across document sections. PEARL indexes over 10,000 embeddings within two minutes, with per-query retrieval consistently completing in under one second, confirming real-time viability on low-resource hardware.

**Grounded Generation.** The Gemini LLM achieves 93% response accuracy against human-verified ground truth. Observed errors are attributable to ambiguous or underspecified content within source documents rather than model hallucination, validating the effectiveness of the provenance-aware prompt construction strategy. Multi-turn dialogue support maintains contextual continuity without degrading retrieval reliability.

**Multilingual Robustness.** The bidirectional multilingual pipeline achieves 97% language detection accuracy and 94% semantic retention across supported languages. Response quality remains consistent across all tested languages, confirming



that the normalisation process introduces no measurable degradation in retrieval or generation performance.

**Latency.** PEARL achieves an average end-to-end response latency of 2.1 seconds, with 95% of queries resolved within three seconds. This encompasses language detection, translation, semantic retrieval, grounded generation, and response rendering demonstrating real-time usability on standard CPU hardware without GPU acceleration.

**Usability.** A user study involving both student and non-expert participants yielded a System Usability Scale (SUS) score of 84, corresponding to a “Good” usability rating. Participants highlighted the multilingual interface, natural conversational flow, and transparent source attribution as the three primary strengths of the system.

**Analysis.** The consistent performance across all evaluation dimensions validates the four core design decisions of PEARL: context-preserving semantic chunking, distance-threshold gated retrieval, bidirectional multilingual normalisation, and provenance-aware prompt construction. Collectively, these components demonstrate that high-quality grounded generation is achievable in low-resource, multilingual settings without reliance on GPU infrastructure.

## VI. CONCLUSION AND FUTURE WORK

The goal of this article was to introduce an innovative multilingual Retrieval-Augmented Generation (RAG) framework that will allow for efficient and accurate retrieval of information from long documents such as government and policy documents. Through the incorporation of various modules, such as context preserving semantic chunking, low-resource FAISS-based indexing, and multilingual query normalization, we provide new opportunities for researchers to access information that would otherwise be lost due to time or lack of access.

The proposed methodology has been tested in several ways and the results of these experiments indicate that our approach performs extremely well in each area of performance assessed. Our findings suggest that our proposed methodology has a relatively low average latency of only 2.1 seconds from the time that a question is asked until a relevant answer is obtained.

These results indicate that our approach to multilingual RAG provides a practical and valuable solution for any organization that needs multilingual retrieval capabilities on their literature.

The system will be improved by future development initiatives in multiple ways. To begin with, adding new hallucination detection and uncertainty estimation tooling will improve reliability and ensure higher levels of trust-worthiness. Next, the minority language aspects of the multilingual component must also be increased in order to enhance the overall capability of the system. A third area of development will include integrating offline or edge-aware usage components into the system architecture. Utilising quantised embedding representations or light-weight, low-latency language model alternatives will allow the system to be utilised in environments with limited connectivity, increasing the ability of individuals and organisations to access it, thereby widening

its accessibility. Lastly, the use cases of multimodal RAG processing will expand significantly, for example, through the inclusion of image-based document rendering or table extraction capabilities from a scanned document (using optical character recognition technology) or from a PDF file produced by a scanning device.

As demonstrated, a multilingual RAG system can be cost-effective and efficient in terms of scale. Consequently, the multilingual RAG Framework can serve as a foundation for future generation document intelligence applications and subsequent innovation in document management within organisations.

## REFERENCES

- [1] S. Robertson and H. Zaragoza, “The probabilistic relevance framework: BM25 and beyond,” *Foundations and Trends in Information Retrieval*, vol. 3, no. 4, pp. 333–389, 2009.
- [2] Z. Ji, N. Lee, R. Frieske, T. Yu, D. Su, Y. Xu, E. Ishii, Y. J. Bang, A. Madotto, and P. Fung, “Survey of hallucination in natural language generation,” *ACM Computing Surveys*, 2023.
- [3] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W.-t. Yih, T. Rocktäschel, S. Riedel, and D. Kiela, “Retrieval-augmented generation for knowledge-intensive NLP tasks,” in *Advances in Neural Information Processing Systems*, 2020.
- [4] R. Sharma, K. Patel, and A. Kumar, “Multilingual information systems: Challenges and solutions,” *Journal of Information Systems Research*, 2023.
- [5] A. Conneau, K. Khandelwal, N. Goyal, V. Chaudhary, G. Wenzek, F. Guzmán, E. Grave, M. Ott, L. Zettlemoyer, and V. Stoyanov, “Unsupervised cross-lingual representation learning at scale,” in *Proceedings of ACL*, 2020.
- [6] V. Karpukhin, B. Oguz, S. Min, P. Lewis, L. Wu, S. Edunov, D. Chen, and W.-t. Yih, “Dense passage retrieval for open-domain question answering,” in *Proceedings of EMNLP*, 2020.
- [7] N. Reimers and I. Gurevych, “Sentence-BERT: Sentence embeddings using siamese BERT-networks,” in *Proceedings of EMNLP*, 2019.
- [8] J. Johnson, M. Douze, and H. Jegou, “Billion-scale similarity search with GPUs,” *IEEE Transactions on Big Data*, 2019.
- [9] L. Zhang, Q. Wang, and Y. Liu, “Document chunking strategies for large-scale information retrieval,” in *Proceedings of EMNLP*, 2022.
- [10] M. Gupta, P. Singh, and S. Verma, “Ai-powered domain-specific assistants: A systematic review of deployment strategies,” *AI Applications Review*, 2023.
- [11] X. Chen, M. Rodriguez, and K. Johnson, “Evaluation frameworks for conversational AI in domain-specific applications,” *Transactions on Interactive Intelligent Systems*, 2023.