

# Breaking It Down: Domain-Aware Semantic Segmentation for Retrieval Augmented Generation

Aparajitha Allamraju, Maitreya Prafulla Chitale, Hiranmai Sri Adibhatla,  
Rahul Mishra, and Manish Shrivastava

International Institute of Information Technology, Hyderabad, India

**Abstract.** Document chunking is a crucial component of Retrieval-Augmented Generation (RAG), as it directly affects the retrieval of relevant and precise context. Conventional fixed-length and recursive splitters often produce arbitrary, incoherent segments that fail to preserve semantic structure. Although semantic chunking has gained traction, its influence on generation quality remains underexplored. This paper introduces two efficient semantic chunking methods, Projected Similarity Chunking (PSC) and Metric Fusion Chunking (MFC), trained on PubMed data using three different embedding models. We further present an evaluation framework that measures the effect of chunking on both retrieval and generation by augmenting PubMedQA with full-text PubMed Central articles. Our results show substantial retrieval improvements ( $\uparrow 24x$  with PSC) in MRR and higher Hits@ $k$  on PubMedQA. We provide a comprehensive analysis, including statistical significance and response-time comparisons with common chunking libraries. Despite being trained on a single domain, PSC and MFC also generalize well, achieving strong out-of-domain generation performance across multiple datasets. Overall, our findings confirm that our semantic chunkers, especially PSC, consistently deliver superior performance.

**Keywords:** Semantic Chunking · RAG Evaluation · Question Answering

## 1 Introduction

Retrieval-Augmented Generation (RAG) [8] has become a key approach for enhancing Large Language Models (LLMs) in high-stakes domains such as finance, law, and healthcare. By integrating external, domain-relevant context, RAG improves the factual grounding and reliability of model outputs. However, its effectiveness depends heavily on the quality of the retrieval stage.

A central factor in retrieval performance is chunking [2], which determines how documents are segmented for indexing. Effective chunking preserves semantic coherence, enabling the retrieval of contextually complete and relevant segments and improving downstream tasks like question answering. Yet mainstream

RAG frameworks largely rely on fixed-length, sentence-based, or paragraph-based methods, which frequently yield incoherent or overly broad chunks that weaken retrieval precision.

We propose two semantically aware, domain-trained chunking methods, Projected Similarity Chunking (PSC) and Metric Fusion Chunking (MFC), which use distinct boundary-detection mechanisms and are trained on augmented biomedical corpora. Our results show substantial improvements in both retrieval and generation compared to fixed-length, sentence-level, semantic, and recursive baselines. Furthermore, the chunkers generalize well to out-of-domain datasets, demonstrating robustness beyond the biomedical domain. All code, trained models, and augmentation scripts will be released upon acceptance.

## 2 Background

Chunking is a critical preprocessing step in Retrieval-Augmented Generation (RAG), playing a fundamental role in improving both the efficiency and relevance of the retrieved context for downstream tasks. Existing approaches typically based on fixed-size segments or rigid predefined rules, often fail to capture semantic boundaries effectively. While benchmarks exist for evaluating retrieval [15] and generation [5, 14, 17, 18], along with established metrics for both tasks [1], most studies evaluate the retriever and generator modules independently. The influence of chunking on retrieval and generation performance, however, remains comparatively underexplored.

Semantic chunking addresses boundary segmentation by producing coherent, contextually aligned text segments instead of arbitrary fixed-length chunks. Some existing library implementations [7] attempt to solve the limitations of uniform chunking by employing embedding-based strategies. In such methods, candidate chunks are identified by measuring the cosine similarity between adjacent sentences (or merged segments) and splitting when the similarity falls below a threshold, indicating a topic shift. However, cosine similarity has known shortcomings for downstream tasks such as information retrieval [23] and often lacks robust, domain-specific performance, as our experiments demonstrate. The cost of use of the library implementation of semantic chunking was elaborated in [11] and concluded that the existing semantic chunker may not always be worth the computation cost. Additionally, recent research [22] has explored predicting optimal chunk granularity, but to our knowledge, no recent work has proposed new semantic chunking algorithms or systematically examined their joint impact on retrieval and generation.

We propose an approach that better aligns with human judgment than existing baselines, particularly because the gold data used for evaluation, both for our models and the baselines, is human-authored. Our methods not only deliver higher retrieval and generation quality but also achieve measurable computational efficiency gains compared to current chunking strategies.

### 3 Methodology

In this section, we detail the process of training and employing our semantic chunkers using domain-specific documents. The hypothesis behind our methods is that two sentences similar in ideas would occur together while two sentences with not so close ideas would occur in different paragraphs and thereby be in separate chunks. Instead of just calculating distance between sentence embeddings, we train our models to be able to amplify the distances and give us scores which we then use to chunk. We train the chunkers using positive and negative samples created from the PQA-U dataset as described in Section 4.1.

#### 3.1 Semantic Chunking

The two methods that we introduce are distinctive in how they define the semantic boundaries of a chunk, ultimately determining the separation points between two chunks. In both the methods, we employ a 1:1 negative sampling of data as mentioned in Section 4.1, to ensure that the learned representations and boundaries are robust. We choose simple light-weight approaches and models to ensure that performance is optimized for downstream tasks and the chunking itself doesn't take much time or compute. The shared methodology between both of our chunking strategies is explained below.

Given a sentence pair  $S_i$  and  $S_j$ , we encode them using state-of-the-art sentence transformers [12]: **all-MiniLM-L6-v2**(MiniLM) [20], **all-mpnet-base-v2**(mpnet) [13], and **e5-large-v2**(e5) [19] to have a total of three chunking models per strategy.

The resulting embeddings  $E_i$  and  $E_j$  undergo a linear transformation. The training enables the model to specifically tune the embedding space for boundary detection.

In **Projected Similarity Chunking** (PSC), the dot-product similarity between embeddings is obtained and the resultant final scalar similarity score is used to determine the boundary

In **Metric Fusion Chunking** (MFC), we employ a combined metric of evaluating similarity between two given embeddings. We employ an equally weighted combination of Dot Product Similarity, Euclidean Distance, and Manhattan Distance which is passed through a single neural layer to obtain the final similarity score.

The scalar similarity score  $S_{i,j}$  is normalized into a probability range  $(0, 1)$  using sigmoid. This value represents the model's final prediction. If the output probability is greater than or equal to 0.5 (determined empirically), the individual model predicts that the two sentences occur in the same section. Otherwise, if the probability is below 0.5, it predicts that the sentences should not occur one consecutively and we use that to draw our segment boundary or create the chunks.

## 4 Experiments

This section presents experiments conducted using the proposed model on our augmented PubMedQA dataset and the RAGBench [3] datasets. It is organized into five subsections: dataset, evaluation metrics, baselines, experimental setup, and evaluation setup.

### 4.1 Dataset

Currently, no publicly available datasets exist for jointly evaluating retrieval and generation performance. To address this gap, we repurposed the PubMedQA (PQA) dataset [6] by augmenting it with corresponding full-text articles from PubMed Central (PMC). The full texts were downloaded in XML format following the official instructions <sup>1</sup>.

PubMedQA is a carefully curated biomedical question-answering dataset derived from PubMed articles whose titles are phrased as questions. The dataset is divided into labeled (PQA-L), unlabeled (PQA-U), and artificial (PQA-A) subsets. PQA-L and PQA-U are human-annotated, whereas PQA-A is automatically labeled. In the original dataset, the title serves as the question, the context consists of selected sentences from the abstract that support the conclusion, and the conclusion itself serves as the long-form answer. Upon downloading the corresponding full-text XML files, we observed that not all PQA entries have available articles on PMC. Specifically, from the original 1k samples in PQA-L, only 115 have full-text matches; PQA-U has 6.45k matches instead of 61.2k; and PQA-A has 51.2k matches from an original 211.3k.

The XML files were cleaned to remove figures, tables, captions, references, appendices, and other supplementary materials, retaining only the abstract and main article text. For our experiments, we kept the original question from the dataset and used the cleaned full text. The original “context” sentences were used to evaluate retriever performance, while the generated answers were compared to the original long answers from the conclusions.

We used PQA-A’s full-text data to train our semantic chunking models and reserved PQA-L for evaluation. Training data was prepared by splitting PQA-A articles into pairs. To help the model distinguish between sequentially related and unrelated sentences, we applied negative sampling. For each positive pair (sentences occurring consecutively), one negative pair was added, yielding a balanced set of 50.3% positive and 49.7% negative samples. Positive pairs were defined as sentences grouped by the human author within the same section of the text. Negative pairs were created by randomly selecting sentences from different sections within the same document that never co-occur in any section. This design ensures that the model learns to identify contextually related sentences based on human-defined structure. The final dataset contained 93M sentence pairs for training our models.

---

<sup>1</sup> <https://www.ncbi.nlm.nih.gov/research/bionlp/RESTful/pmcqa.cgi/>

The original dataset columns included: *pubid*, *question*, *context*, *long\_answer*, and *final\_decision*. In our augmented version, the schema is updated to include *full\_text*. We do not use *final\_decision*, focusing exclusively on long-form answering.

By introducing full-text articles alongside the original question answer pairs, our augmented PQA dataset enables a unified framework for evaluating both retrieval quality and answer generation in a controlled setting. The curated contexts allow us to measure retriever precision, while the long-form conclusions serve as a robust reference for generation metrics, making it uniquely suited for end-to-end RAG evaluation.

## 4.2 Evaluation metrics

We evaluate the system at two levels: retrieval quality and answer generation. For the retriever, we employ Hits@k and Mean Reciprocal Rate (MRR) to measure the effectiveness and ranking efficiency of the retrieved chunks. For the generator, we evaluate the outputs using standard generation metrics such as BLEU [10], ROUGE [9], and BERTScore [21]. Additionally, we measure query response time and time-to-first-token to assess the responsiveness of our chunker models.

## 4.3 Baselines

In our experiments, we compare several popular chunking strategies against our proposed methods to establish a baseline for performance. All chunking strategies are implemented using LangChain<sup>2</sup> library. The sentences and chunks derived are then encoded using one of three sentence transformer [12] models: all-MiniLM-L6-v2 [20], all-mpnet-base-v2 [13], and e5-large-v2 [19]. The performance of each chunker is evaluated based on its impact on retrieval and generation tasks.

**Character Chunking** (Fixed-Length) (Char) splits text into fixed-size blocks, often ignoring natural linguistic boundaries and reducing coherence [2, 4, 8]. We adopt a 1000-token chunk size with 200-token overlap, chosen because two consecutive chunks typically capture a complete semantic unit while avoiding unnecessary redundancy. **Sentence Chunking** (Sent) groups a fixed number of full sentences, offering better linguistic consistency than character-level splitting. We use three sentences with a one-sentence overlap, roughly matching the 1000 and 200-token configuration, though this produces variable chunk lengths that complicate uniform processing. **Recursive Chunking** (Rec) segments documents hierarchically (e.g., paragraphs to sentences). While more structured, it requires additional computation to identify segmentation levels. We use the same 1000-token and 200-token setups for comparability. **Semantic Chunking** (Sem) uses sentence embeddings to detect semantic shifts and form meaningfully aligned chunks. Although conceptually strong, LangChain’s default SemanticChunker under-performs in our evaluations relative to the proposed methods.

---

<sup>2</sup> <https://www.langchain.com>

#### 4.4 Experimental setup

We trained our two chunker architectures, PSC and MFC, on three selected sentence embedding models, MiniLM, mpnet, and E5, using the augmented PQA-A dataset. The augmented PQA-A dataset contains approximately 93M training samples.

All model training was performed on a single NVIDIA RTX 6000 GPU with 48 GB of VRAM. Each embedding–chunker combination was trained for 5 epochs, and empirical evaluation indicated that models trained up to epoch 4 achieved the best overall performance. We used the BCEWithLogitsLoss objective function, which combines a sigmoid activation with binary cross-entropy loss.

Training time for each model combination was under 12 hours, resulting in a total cumulative GPU usage of roughly 72 hours across all experiments.

For generation and inference, we used a pair of NVIDIA RTX 4090 GPUs (24 GB VRAM each). The PQA-L dataset, comprising 115 documents, was used for in-domain evaluation, requiring about 1 hour per model. Out-of-domain evaluation on RagBench averaged 14 hours per model. This setup enabled consistent, large-scale experimentation while maintaining practical computational requirements.

#### 4.5 Evaluation Setup

We evaluate the impact of chunking strategies on both retrieval quality and answer generation. Using the augmented PQA-L dataset, gold-standard chunks and answers enable direct comparison against retrieved contexts and generated outputs. Alongside our PSC and MFC chunkers, we benchmark four baselines: Character, Sentence, Semantic, and Recursive Chunkers.

For out-of-domain testing, we use 12 RagBench [3] datasets to measure generation quality across chunkers. Since RagBench lacks gold context labels, we evaluate only generated answers, which though LLM-produced, are reliable for benchmarking. This setup allows us to test how well domain-trained chunkers generalize to varied document types.

Retriever evaluation is conducted solely on the PubMedQA extension using MRR and Hits@k. PSC and MFC produce naturally variable chunk lengths (averaging 471 tokens), which we find improves retrieval and downstream answer generation compared to fixed-size segmentation.

For generation, we retrieve the top-5 chunks and pass them, along with the query, to a Llama-3.1-8B [16] instruction-tuned model with fixed decoding parameters across all datasets. Outputs are evaluated against gold answers using ROUGE [9], BLEU [10], and BERTScore [21]. Keeping the retriever and generator constant ensures all performance differences arise solely from the chunking strategy.

Each dataset, chunker, and encoder configuration follows a unified workflow: the document is segmented, its embeddings are indexed, the most relevant chunks are retrieved for a given query, and a generator model produces the final answer.

Model	Chunker	Hits@3	Hits@5	MRR	Query Time	TTFT
MinILM	Char	0.0	0.0	0.0	0.01	5.91
	Sent	0.0	0.0	0.0	0.01	5.91
	Sem	0.00*	0.00*	0.01	0.01	0.56
	Rec	0.0	0.0	0.0	0.01	0.19
PSC	PSC	<b>0.10</b>	<b>0.13</b>	<b>0.30</b>	<b>0.00</b>	<b>0.12</b>
	MFC	0.09	0.12	0.26	0.01	0.13
mpnet	Char	0.0	0.0	0.0	0.02	5.86
	Sent	0.0	0.0	0.0	0.01	0.34
	Sem	0.00*	0.00*	0.01	0.01	0.62
	Rec	0.0	0.0	0.0	0.01	0.19
e5	PSC	<b>0.11</b>	<b>0.15</b>	<b>0.31</b>	0.01	<b>0.13</b>
	MFC	0.10	0.14	0.29	0.01	0.15
e5	Char	0.0	0.0	0.0	0.02	5.58
	Sent	0.0	0.0	0.0	0.01	0.34
	Sem	0.00*	0.00*	0.01	0.01	0.63
	Rec	0.0	0.0	0.0	0.01	0.19
PSC	PSC	<b>0.13</b>	<b>0.16</b>	<b>0.36</b>	0.01	<b>0.14</b>
	MFC	0.12	0.15	0.34	0.01	0.20

Table 1: Retrieval performance measured with Hits@3, Hits@5 and MRR across the three embedding models, show significant improvement, on using our chunkers. \* indicates infinitesimal non-zero values. The p-test values between PSC and Sem across the three embeddings and metrics were in the order of  $e^{-11}$  and t-test values greater than 7. Query Time and Time to First Token (TTFT) indicate that PSC is the fastest.

This integrated design removes redundant pipelines, enforces consistent evaluation across setups, and offers a scalable framework for systematically comparing RAG components.

## 5 Analysis

Our trained chunkers significantly outperform commonly used methods, MRR improves from 0.0086 (recursive) to 0.2914 (PSC in e5) and Hits@5 from 0.0028 to 0.1234 (PSC in e5), as shown in Table 1. The library implementation of semantic chunker under performs in retrieval (MRR: 0.0130 vs. 0.2914 for PSC) due to its reliance on cosine similarity for chunk boundary detection, which often results in suboptimal segmentation. These results confirm that chunk quality strongly influences retrieval effectiveness.

Generation results on PubMedQA (Table 2) show PSC and MFC performing competitively and often surpassing baselines. MFC with the E5 embedding model achieves the strongest BLEU, ROUGE, and BERTScore, while PSC and MFC with mpnet also achieve high ROUGE scores. This indicates that our chunkers yield contextually aligned outputs.

	Embedding	Chunker	BLEU	ROUGE1	ROUGE2	ROUGEL	BERTScore
MinILM	Char	<b>28.10</b>	<b>30.63</b>	<b>13.16</b>	<b>22.26</b>	<b>87.43</b>	
	Sent	21.84	22.64	8.37	15.97	84.72	
	Sem	24.73	27.19	10.67	19.75	86.57	
	Rec	22.68	24.33	9.08	17.53	85.73	
mpnet	PSC	24.20	27.08	10.00	19.45	86.92	
	MFC	22.55	25.56	9.08	18.07	86.48	
	Char	23.31	25.64	9.81	18.52	86.80	
	Sent	23.19	25.08	9.50	17.39	86.45	
e5	Sem	22.01	25.30	<b>10.50</b>	18.44	86.83	
	Rec	20.22	23.26	8.73	16.95	84.17	
	PSC	23.14	<b>26.94</b>	10.21	19.68	85.91	
	MFC	<b>23.60</b>	26.61	9.83	<b>19.71</b>	<b>87.38</b>	
	Char	24.06	26.32	10.11	18.77	86.53	
	Sent	19.42	21.66	8.32	15.14	85.92	
	Sem	22.19	23.88	8.70	17.19	85.93	
	Rec	22.10	24.78	10.09	18.22	86.56	
	PSC	23.41	26.52	9.59	18.77	87.42	
	MFC	<b>25.34</b>	<b>27.68</b>	<b>10.31</b>	<b>19.94</b>	<b>87.46</b>	

Table 2: Evaluation metrics (BLEU, ROUGE-1, ROUGE-2, ROUGE-L, and BERTScore) for different embedding models and chunking strategies. Results show Character chunking with miniLM, MFC with mpnet and E5 perform the best. The p-test values between PSC and Char across the three embeddings and five metrics were less than 0.3 and t-test values around than 1.

Evaluation reveals a mismatch between conventional metrics and actual generation quality. Standard metrics like BLEU, ROUGE, and BERTScore prioritize lexical overlap over semantic nuance and factual accuracy. Consequently, they fail to capture the superiority of PSC and MFC chunkers, which produce concise, context-focused responses free from the verbosity and drift found in other methods.

A two-tailed independent samples t-test was conducted to evaluate the statistical significance of the performance metrics between all the chunkers. The difference between the reported means of MRR, Hits@3, Hits@5, BLEU and ROUGE was found to be statistically significant, as indicated by a p-value and t-value in the captions of Tables 1 and 2. This suggests that the improvement of the PSC Chunker over the next best chunker is not due to random chance with a confidence of 99.9% for retrieval, and 90% for generation.

For out-of-domain tests, the generation results with BLEU and ROUGE-L scores on every dataset from RAGBench are summarized in Tables 3.1, 3.2 and 3.3. Despite our chunkers being trained only on PubMed dataset, we can see

a clear improvement in the performance across all datasets with our chunkers, especially PSC.

Embed Chunker	CovidQA	CUAD		DelucionQA		EManual		ExpertQA		FinQA		HAGRID		HotpotQA		MS Marco		PubMedQA		TAT-QA		TechQA			
	BLEU	R-L BLEU	BLEU	R-L BLEU	BLEU	R-L BLEU	BLEU	BLEU	R-L BLEU	BLEU	R-L BLEU	BLEU	R-L BLEU	BLEU	R-L BLEU	BLEU	R-L BLEU	BLEU	R-L BLEU	BLEU	R-L BLEU	BLEU	R-L BLEU		
Minilm	Char	1.95	3.16	9.86	9.82	2.91	4.0	3.98	6.22	16.34	12.49	2.13	3.44	2.63	4.17	1.03	1.94	1.63	2.53	8.49	10.21	3.07	5.01	13.48	12.05
	Sent	2.51	4.22	5.83	6.61	5.73	9.28	5.52	8.65	9.29	8.73	4.11	6.52	2.97	4.95	1.26	2.48	3.16	5.15	9.61	11.95	3.51	6.04	9.51	10.59
	Sem	3.48	5.44	5.64	5.82	7.35	10.93	8.32	11.13	10.32	8.59	4.64	6.72	4.59	6.95	2.07	3.88	4.33	6.34	14.36	15.90	4.83	7.12	5.67	5.91
	Rec	4.62	7.04	10.73	10.61	9.43	13.36	9.78	12.75	17.17	13.01	6.03	8.79	4.57	6.81	2.03	3.85	5.01	7.33	14.77	16.43	4.69	7.20	14.38	12.64
mpnet	PSC	<b>6.31</b>	<b>8.83</b>	<b>11.34</b>	<b>10.75</b>	<b>15.61</b>	<b>17.25</b>	<b>13.85</b>	<b>14.44</b>	20.83	14.32	<b>9.76</b>	<b>12.30</b>	<b>7.22</b>	<b>9.71</b>	4.00	7.27	<b>8.32</b>	<b>10.71</b>	<b>24.33</b>	<b>21.38</b>	<b>6.40</b>	<b>8.82</b>	<b>18.01</b>	<b>13.91</b>
	MFC	6.16	8.63	9.90	9.78	15.67	17.12	13.40	13.69	<b>20.93</b>	<b>14.33</b>	9.15	11.64	7.06	9.64	<b>4.06</b>	<b>7.40</b>	8.17	<b>10.71</b>	23.72	21.18	5.83	8.21	17.19	13.34
	Char	1.92	3.16	9.71	9.73	2.84	4.82	3.97	6.26	16.11	12.43	2.10	3.39	2.60	4.15	1.03	1.96	1.63	2.53	8.40	10.13	2.90	4.75	13.18	12.02
	Sent	2.48	4.20	5.90	6.69	5.64	9.24	5.45	8.71	9.20	8.63	4.19	6.49	2.97	4.96	1.23	2.43	3.14	5.13	9.55	11.89	3.53	6.07	9.50	10.72
E5	Sem	3.47	5.38	5.7	6.11	6.91	10.57	8.18	11.46	10.49	8.66	4.19	6.08	4.55	7.00	2.00	3.74	4.24	6.29	14.09	15.66	4.39	6.62	5.59	5.95
	Rec	4.58	6.97	10.72	10.35	9.4	13.49	9.21	12.42	17.07	13.05	5.89	<b>8.55</b>	4.62	6.93	2.02	3.82	5.01	7.37	14.68	16.38	<b>4.48</b>	<b>6.88</b>	14.07	12.56
	PSC	<b>6.76</b>	<b>9.39</b>	<b>13.14</b>	<b>11.65</b>	<b>12.61</b>	<b>15.78</b>	18.10	<b>17.88</b>	<b>19.28</b>	<b>13.69</b>	<b>6.16</b>	8.50	<b>7.02</b>	<b>9.68</b>	<b>3.69</b>	6.64	<b>8.33</b>	<b>10.91</b>	<b>24.47</b>	<b>21.62</b>	4.19	6.36	<b>15.78</b>	<b>12.82</b>
	MFC	6.66	9.23	12.9	11.5	12.67	15.53	<b>18.61</b>	17.84	18.47	13.30	5.48	7.67	6.65	9.34	3.65	<b>6.72</b>	7.99	10.52	23.91	21.49	3.89	5.96	15.39	12.61
Minilm	Char	1.80	3.01	9.44	9.32	2.79	4.71	3.86	6.04	15.96	12.06	2.07	3.46	2.14	3.40	1.04	1.97	1.58	2.44	8.43	10.18	3.25	5.36	12.52	12.04
	Sent	2.46	4.23	6.18	7.02	5.73	9.3	5.60	8.56	9.27	8.63	4.22	6.80	2.81	4.67	1.23	2.44	3.15	5.16	9.69	12.05	3.52	6.12	9.65	10.99
	Sem	3.30	5.33	6.29	6.05	7.67	11.15	8.50	11.69	10.23	8.56	4.55	6.73	4.24	6.45	1.74	3.26	5.27	7.59	14.46	16.05	4.80	7.24	4.92	5.54
	Rec	4.65	7.34	11.0	10.61	9.59	13.58	9.40	12.58	17.22	12.81	6.00	8.97	4.51	6.74	1.92	3.63	5.01	7.33	14.95	16.67	4.56	7.09	<b>13.87</b>	<b>12.82</b>
mpnet	PSC	<b>6.31</b>	<b>9.24</b>	<b>13.96</b>	<b>12.23</b>	<b>6.52</b>	<b>9.27</b>	<b>6.14</b>	<b>8.75</b>	18.16	13.17	<b>9.25</b>	<b>11.99</b>	<b>6.65</b>	<b>9.27</b>	3.21	6.05	<b>7.42</b>	<b>10.12</b>	<b>25.00</b>	<b>22.20</b>	<b>6.81</b>	<b>9.40</b>	13.39	11.77
	MFC	6.20	9.09	9.99	9.32	5.42	8.22	6.06	8.68	<b>19.17</b>	<b>13.45</b>	6.97	9.56	6.14	8.71	<b>3.73</b>	<b>6.90</b>	6.92	9.49	24.49	22.10	5.25	7.77	10.65	9.83

Table 3.1: Evaluation metrics BLEU and ROUGE-L (R-L) for different embedding models and chunking strategies on 12 datasets from RAGBench. The results in bold show that our chunkers perform better across the board despite being trained only one domain data.

Embed Chunker	CovidQA				CUAD				DelucionQA				EManual				ExpertQA				FinQA			
	R1	R2	BScore	R1	R2	BScore	R1	R2	BScore	R1	R2	BScore	R1	R2	BScore	R1	R2	BScore	R1	R2	BScore			
Minilm	Char	3.78	2.68	<b>83.30</b>	15.14	5.42	79.23	5.67	4.26	<b>84.81</b>	7.72	5.24	<b>84.13</b>	23.49	9.80	<b>82.79</b>	4.55	2.40	79.78					
	Sent	4.95	3.55	82.67	9.88	3.69	<b>80.44</b>	10.97	7.93	83.98	10.89	7.12	80.33	15.30	6.61	82.16	8.62	4.70	<b>80.13</b>					
	Sem	6.60	4.55	83.19	8.89	3.25	79.71	13.20	9.36	84.75	14.95	9.07	83.51	16.00	6.95	82.66	9.13	4.85	79.71					
	Rec	8.61	5.95	83.10	16.47	<b>5.93</b>	79.22	16.65	11.56	83.86	17.32	10.59	83.29	24.42	10.11	82.75	11.77	6.33	79.47					
mpnet	PSC	<b>11.03</b>	<b>6.87</b>	<b>82.46</b>	<b>16.87</b>	5.89	79.84	<b>23.16</b>	<b>13.68</b>	82.97	<b>21.11</b>	<b>10.79</b>	82.61	<b>27.42</b>	10.56	82.60	<b>17.10</b>	<b>9.10</b>	79.88					
	MFC	10.79	6.72	82.55	15.12	5.31	79.81	23.05	13.47	82.93	20.16	9.94	82.72	27.39	<b>10.67</b>	82.51	16.15	8.53	79.85					
	Char	3.75	2.68	<b>83.19</b>	15.10	5.38	78.87	5.55	4.14	<b>85.08</b>	7.70	5.27	<b>84.33</b>	23.20	9.59	82.78	4.46	2.37	<b>79.74</b>					
	Sent	4.94	3.54	82.49	10.01	3.58	<b>80.17</b>	10.87	7.92	84.13	10.82	7.19	83.18	15.20	6.54	82.03	8.70	4.59	79.72					
E5	Sem	6.56	4.47	83.00	9.15	3.44	79.75	12.62	9.17	84.80	14.96	9.57	83.71	16.21	7.02	82.71	8.29	4.35	79.57					
	Rec	8.50	5.74	82.97	16.27	5.78	78.99	16.66	11.57	83.88	16.52	10.18	83.31	24.27	9.98	82.78	11.57	6.16	79.29					
	PSC	<b>11.75</b>	<b>7.40</b>	82.55	<b>18.65</b>	<b>6.43</b>	79.79	<b>20.49</b>	<b>13.01</b>	83.72	26.46	12.72	81.77	<b>25.98</b>	<b>10.52</b>	<b>82.81</b>	<b>11.71</b>	<b>6.12</b>	79.54					
	MFC	11.57	7.20	82.44	18.49	6.27	79.70	20.07	12.61	83.65	<b>26.75</b>	<b>12.91</b>	81.91	25.23	10.32	<b>82.81</b>	10.53	5.48	79.54					
Minilm	Char	3.56	2.61	<b>83.64</b>	14.56	5.16	79.19	5.43	4.07	<b>84.97</b>	7.46	5.14	84.17	22.74	9.40	82.64	4.47	2.47	79.93					
	Sent	4.94	3.62	82.88	10.62	3.84	<b>80.05</b>	10.98	8.00	84.03	11.08	7.06	82.79	15.15	6.57	81.90	8.91	4.95	<b>80.7</b>					
	Sem	6.37	4.53	83.58	9.39	3.39	79.55	13.60	9.51	84.67	15.41	9.58	83.46	15.80	6.93	82.57	9.02	4.91	79.81					
	Rec	8.78	6.23	83.48	16.70	5.84	79.25	<b>16.88</b>	<b>11.69</b>	83.87	<b>16.91</b>	<b>10.38</b>	83.44	24.00	9.76	<b>82.67</b>	11.92	6.60	79.57					
mpnet	PSC	<b>11.29</b>	<b>7.43</b>	82.99	<b>19.61</b>	<b>6.57</b>	79.47	11.27	7.89	84.62	11.27	7.30	<b>84.18</b>	24.87	10.07	82.61	<b>16.42</b>	<b>8.94</b>	79.99					
	MFC	11.10	7.34	82.99	14.66	5.06	79.49	9.82	7.11	84.76	11.22	7.10	83.79	<b>25.58</b>	<b>10.29</b>	82.66	12.94	7.03	79.86					

Table 3.2: Evaluation metrics ROUGE1 (R1), ROUGE2 (R2) and BERTScore (BScore) on 6 datasets from RAGBench.

We also evaluate the response time of the various chunkers by calculating the Query Response Time for retrieval and Time to First Token(TTFT) for generation. Since the models compute the similarity between sentences for chunking, there is a slight increase of 0.7s against Character chunker and 0.4s against Semantic chunker in indexing time on an average. As indexing is an offline one-time operation, we do not believe that the increase is of importance in this study. The

Embed Chunker	HAGRID				HotpotQA				MS Marco				PubMedQA				TAT-QA				TechQA				
	R1	R2	BScore	R1	R2	BScore	R1	R2	BScore	R1	R2	BScore	R1	R2	BScore	R1	R2	BScore	R1	R2	BScore	R1	R2	BScore	
Minilm	Char	5.09	3.58	<b>83.38</b>	2.12	1.67	82.53	3.18	2.31	84.00	14.81	8.91	<b>85.87</b>	6.21	3.30	<b>80.12</b>	20.95	9.31	81.18						
	Sent	5.99	4.23	82.00	2.69	2.12	81.30	6.24	4.44	82.97	17.01	10.27	84.69	7.55	4.07	79.49	17.70	8.29	81.05						
	Sem	8.63	6.00	82.96	4.28	3.30	<b>82.55</b>	8.21	5.82	<b>84.04</b>	23.17	13.83	85.50	9.13	4.78	80.00	9.97	4.57	<b>81.23</b>						
	Rec	8.50	5.83	82.68	4.19	3.29	82.44	9.44	6.61	83.70	23.79	14.21	85.02	9.01	4.87	80.03	22.09	9.71	81.13						
	PSC	<b>12.53</b>	<b>8.01</b>	82.27	8.15	6.03	81.72	<b>14.68</b>	8.98	83.34	<b>33.38</b>	<b>18.02</b>	84.63	<b>11.39</b>	<b>5.93</b>	79.90	<b>25.16</b>	<b>10.34</b>	81.19						
	MFC	12.37	7.92	82.25	<b>8.26</b>	<b>6.09</b>	81.70	14.46	<b>9.00</b>	83.37	32.92	17.90	84.71	10.51	5.46	79.90	24.02	9.98	81.14						
mpnet	Char	5.06	3.57	<b>83.42</b>	2.13	1.68	82.61	3.17	2.30	84.01	14.68	8.84	<b>85.86</b>	5.87	3.13	<b>80.03</b>	20.82	9.32	81.03						
	Sent	6.01	4.23	82.06	2.62	2.07	81.30	6.23	4.44	82.98	16.90	10.20	84.71	7.68	4.04	79.05	17.69	8.34	80.92						
	Sem	8.62	6.01	83.04	4.13	3.22	<b>82.62</b>	8.06	5.74	<b>84.03</b>	22.83	13.66	85.51	8.35	4.38	79.92	9.89	4.66	<b>81.19</b>						
	Rec	8.61	5.89	82.75	4.18	3.26	82.41	9.44	6.61	83.74	23.70	14.17	85.04	<b>8.67</b>	<b>4.60</b>	79.73	21.80	9.73	80.94						
	PSC	<b>12.37</b>	<b>8.06</b>	82.47	7.42	5.61	82.16	<b>14.75</b>	<b>9.24</b>	83.46	<b>33.64</b>	<b>18.39</b>	84.69	8.01	4.20	79.90	<b>22.69</b>	<b>9.74</b>	80.96						
	MFC	11.84	7.80	82.56	<b>7.52</b>	<b>5.64</b>	82.11	14.25	8.98	83.53	33.20	18.27	84.77	7.48	3.92	79.95	22.34	9.66	81.07						
E5	Char	4.17	2.93	<b>83.45</b>	2.14	1.68	82.67	3.08	2.24	<b>84.03</b>	14.73	8.89	<b>85.88</b>	6.59	3.59	<b>80.24</b>	20.42	9.37	80.95						
	Sent	5.66	3.99	82.08	2.64	2.11	81.35	6.23	4.45	82.97	17.15	10.37	84.73	7.54	4.17	79.58	18.25	8.57	80.60						
	Sem	7.99	5.58	83.09	3.60	2.81	<b>82.70</b>	9.89	6.90	83.94	23.34	13.99	85.53	9.07	4.87	80.08	9.03	4.35	<b>81.17</b>						
	Rec	8.38	5.77	82.79	3.95	3.11	82.49	9.42	6.59	83.73	24.06	14.43	85.03	8.82	4.83	80.11	<b>22.05</b>	<b>10.00</b>	80.98						
	PSC	<b>11.79</b>	<b>7.71</b>	82.61	6.65	5.13	82.32	<b>13.39</b>	<b>8.78</b>	83.64	<b>34.32</b>	<b>18.93</b>	84.73	<b>12.06</b>	<b>6.36</b>	80.03	20.38	9.09	80.87						
	MFC	11.08	7.34	82.72	<b>7.71</b>	<b>5.81</b>	82.11	12.55	8.31	83.78	33.89	18.91	84.84	9.78	5.23	80.14	16.97	7.49	80.98						

Table 3.3: Evaluation metrics ROUGE1 (R1), ROUGE2 (R2) and BERTScore (BScore) on remaining 6 datasets from RAGBench.

results in Table 1 indicate a clear improvement in both the retrieval and generation time with PSC chunker. This shows that the chunking strategy is not only improving retrieval and generation but also lowering the time taken to retrieve relevant chunks and generate an answer. This reiterates our claim for using light and simple models for fast and improved chunking.

## 6 Conclusion

Our approach demonstrates a substantial improvement in retrieval effectiveness, with the MFC chunker trained on the E5 embedding model achieving an MRR increase of up to 24 times over baseline methods. This performance gain highlights the power of domain-specific semantic chunking in optimizing retrieval-augmented generation (RAG) pipelines. By producing chunks that naturally align with human-authored section boundaries, our method preserves contextual coherence and facilitates more accurate retrieval, which in turn benefits downstream generation quality.

Beyond PubMedQA, our results on out-of-domain datasets indicate strong potential for generalization, suggesting that the principles of semantic chunking we present are transferable across domains. This work not only provides a concrete, reproducible methodology for chunking optimization but also introduces a fresh perspective on document segmentation, a critical yet underexplored component in RAG systems. We believe that our findings can serve as a foundation for future research in retrieval-aware chunking strategies, adaptive chunk size determination, and domain-adaptive retrieval optimization. The lightweight models introduced in this paper can be integrated into complex models resulting in more sophisticated approaches.

## References

1. Es, S., James, J., Espinosa Anke, L., Schockaert, S.: RAGAs: Automated evaluation of retrieval augmented generation. In: Aletras, N., De Clercq, O. (eds.) Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations. pp. 150–158. Association for Computational Linguistics, St. Julians, Malta (Mar 2024), <https://aclanthology.org/2024.eacl-demo.16/>
2. Finardi, P., Avila, L., Castaldoni, R., Gengo, P., Larcher, C., Piau, M., Costa, P., Caridá, V.: The chronicles of rag: The retriever, the chunk and the generator. arXiv preprint arXiv:2401.07883 (2024)
3. Friel, R., Belyi, M., Sanyal, A.: Ragbench: Explainable benchmark for retrieval-augmented generation systems. arXiv preprint arXiv:2407.11005 (2024)
4. Gao, Y., Xiong, Y., Gao, X., Jia, K., Pan, J., Bi, Y., Dai, Y., Sun, J., Wang, H., Wang, H.: Retrieval-augmented generation for large language models: A survey. arXiv preprint arXiv:2312.10997 **2**(1) (2023)
5. Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., Steinhardt, J.: Measuring massive multitask language understanding. In: International Conference on Learning Representations (2021), <https://openreview.net/forum?id=d7KBjmI3GmQ>
6. Jin, Q., Dhingra, B., Liu, Z., Cohen, W., Lu, X.: Pubmedqa: A dataset for biomedical research question answering. In: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP). pp. 2567–2577 (2019)
7. Kamradt, G.: Semantic chunking. <https://github.com/charlespwd/project-title> (2024)
8. Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W.t., Rocktäschel, T., Riedel, S., Kiela, D.: Retrieval-augmented generation for knowledge-intensive nlp tasks. In: Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., Lin, H. (eds.) Advances in Neural Information Processing Systems. vol. 33, pp. 9459–9474. Curran Associates, Inc. (2020), [https://proceedings.neurips.cc/paper\\_files/paper/2020/file/6b493230205f780e1bc26945df7481e5-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/6b493230205f780e1bc26945df7481e5-Paper.pdf)
9. Lin, C.Y.: Rouge: A package for automatic evaluation of summaries. In: Text summarization branches out. pp. 74–81 (2004)
10. Papineni, K., Roukos, S., Ward, T., Zhu, W.J.: Bleu: a method for automatic evaluation of machine translation. In: Proceedings of the 40th annual meeting of the Association for Computational Linguistics. pp. 311–318 (2002)
11. Qu, R., Tu, R., Bao, F.S.: Is semantic chunking worth the computational cost? In: Chiruzzo, L., Ritter, A., Wang, L. (eds.) Findings of the Association for Computational Linguistics: NAACL 2025. pp. 2155–2177. Association for Computational Linguistics, Albuquerque, New Mexico (Apr 2025). <https://doi.org/10.18653/v1/2025.findings-naacl.114>, <https://aclanthology.org/2025.findings-naacl.114>
12. Reimers, N., Gurevych, I.: Sentence-BERT: Sentence embeddings using Siamese BERT-networks. In: Inui, K., Jiang, J., Ng, V., Wan, X. (eds.) Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP). pp. 3982–3992. Association for Computational Linguistics

- tics, Hong Kong, China (Nov 2019). <https://doi.org/10.18653/v1/D19-1410>, <https://aclanthology.org/D19-1410/>
13. Song, K., Tan, X., Qin, T., Lu, J., Liu, T.Y.: Mpnet: Masked and permuted pre-training for language understanding. *Advances in neural information processing systems* **33**, 16857–16867 (2020)
  14. Su, L., Duan, N., Cui, E., Ji, L., Wu, C., Luo, H., Liu, Y., Zhong, M., Bharti, T., Sacheti, A.: GEM: A general evaluation benchmark for multimodal tasks. In: Zong, C., Xia, F., Li, W., Navigli, R. (eds.) *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. pp. 2594–2603. Association for Computational Linguistics, Online (Aug 2021). <https://doi.org/10.18653/v1/2021.findings-acl.229>, <https://aclanthology.org/2021.findings-acl.229>
  15. Thakur, N., Reimers, N., Rücklé, A., Srivastava, A., Gurevych, I.: BEIR: A heterogeneous benchmark for zero-shot evaluation of information retrieval models. In: *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)* (2021), <https://openreview.net/forum?id=wCu6T5xFjeJ>
  16. Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al.: Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971 (2023), <https://doi.org/10.48550/arXiv.2302.13971>
  17. Wang, A., Pruksachatkun, Y., Nangia, N., Singh, A., Michael, J., Hill, F., Levy, O., Bowman, S.: SuperGLUE: A stickier benchmark for general-purpose language understanding systems. *Advances in neural information processing systems* **32** (2019)
  18. Wang, A., Singh, A., Michael, J., Hill, F., Levy, O., Bowman, S.R.: GLUE: A multi-task benchmark and analysis platform for natural language understanding. In: *International Conference on Learning Representations* (2019), <https://openreview.net/forum?id=rJ4km2R5t7>
  19. Wang, L., Yang, N., Huang, X., Jiao, B., Yang, L., Jiang, D., Majumder, R., Wei, F.: Text embeddings by weakly-supervised contrastive pre-training. arXiv preprint arXiv:2212.03533 (2022)
  20. Wang, W., Wei, F., Dong, L., Bao, H., Yang, N., Zhou, M.: Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. *Advances in neural information processing systems* **33**, 5776–5788 (2020)
  21. Zhang\*, T., Kishore\*, V., Wu\*, F., Weinberger, K.Q., Artzi, Y.: Bertscore: Evaluating text generation with bert. In: *International Conference on Learning Representations* (2020), <https://openreview.net/forum?id=SkeHuCVFDr>
  22. Zhong, Z., Liu, H., Cui, X., Zhang, X., Qin, Z.: Mix-of-granularity: Optimize the chunking granularity for retrieval-augmented generation. In: Rambow, O., Wanner, L., Apidianaki, M., Al-Khalifa, H., Eugenio, B.D., Schockaert, S. (eds.) *Proceedings of the 31st International Conference on Computational Linguistics*. pp. 5756–5774. Association for Computational Linguistics, Abu Dhabi, UAE (Jan 2025), <https://aclanthology.org/2025.coling-main.384/>
  23. Zhou, K., Ethayarajh, K., Card, D., Jurafsky, D.: Problems with cosine as a measure of embedding similarity for high frequency words. In: Muresan, S., Nakov, P., Villavicencio, A. (eds.) *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. pp. 401–423. Association for Computational Linguistics, Dublin, Ireland (May 2022). <https://doi.org/10.18653/v1/2022.acl-short.45>, <https://aclanthology.org/2022.acl-short.45>