

# Harnessing Information Retrieval Techniques in Retrieval-Augmented Generation

Lixiang Li, Wendy Jiang, Chen Peng, Jie Zheng  
{li4256, jiang794, peng326, zheng795} @purdue.edu



Department of Computer Science

# Outline

- Introduction 4m
- Background: 4m
  - Vector Database
  - Marginalization and Embedding
  - Information Retrieval
  - Large Language Model (LLM)
- Methodology 2m
  - Re-ranking
- Experiment Setup 3m
  - Datasets, Metrics, benchmarks
- Experimental Results and Insights 3m

# Background

- Retrieval-Augmented Generation:
  - A technique that combines the strengths of retrieval-based and generation-based models.
- Improve Factual Accuracy
- Provide context information
- Provide up to date information

# Background

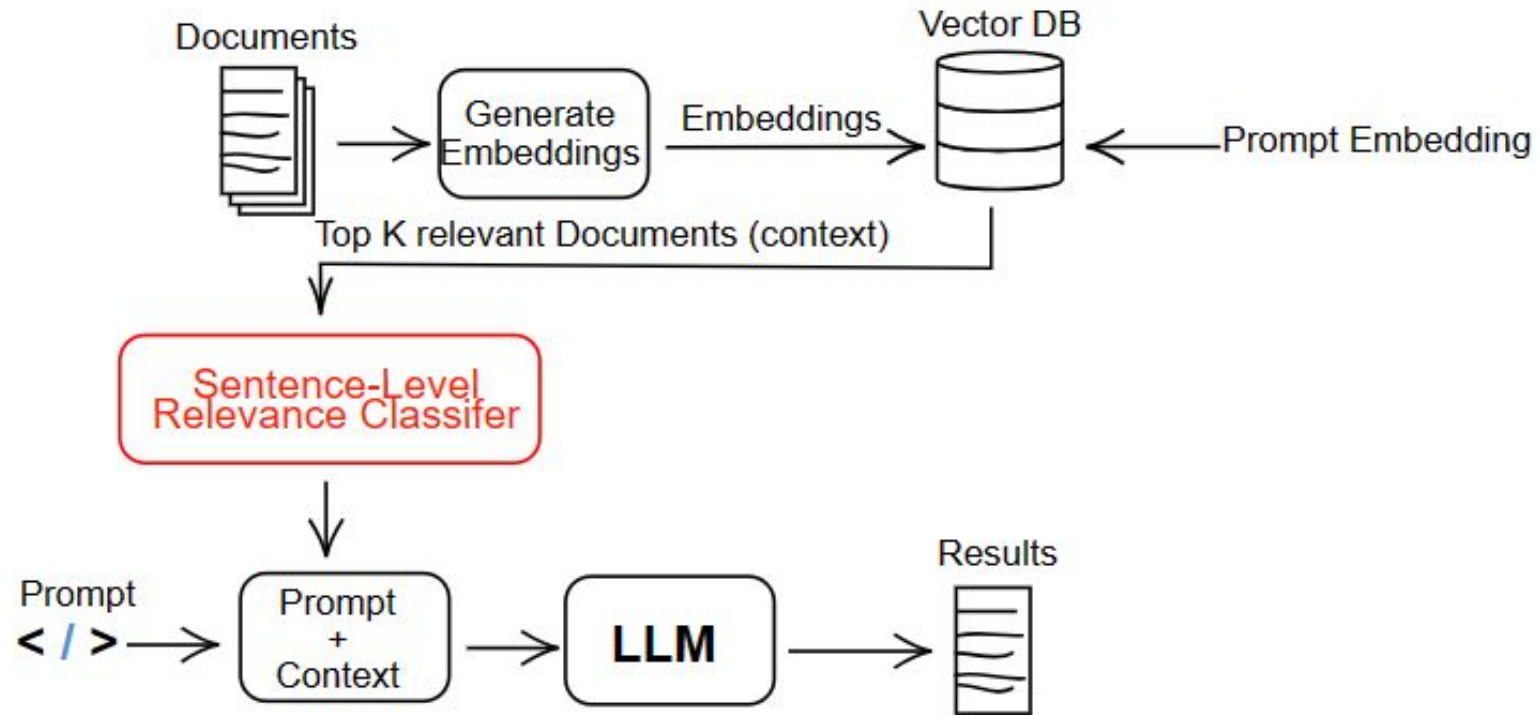
- Limitation:
  - LLM relies on highly relevant and specific information.
  - But traditional methods don't provide relevant and effective information
  - Traditional methods provides long top k documents which are too long to be effective
  - Truncation methods are not reliable enough

# Research Question

- How to design a model to help RAG to get the relevant and important information?
- We propose a model to identify most relevant information.

# Proposed Method

- 

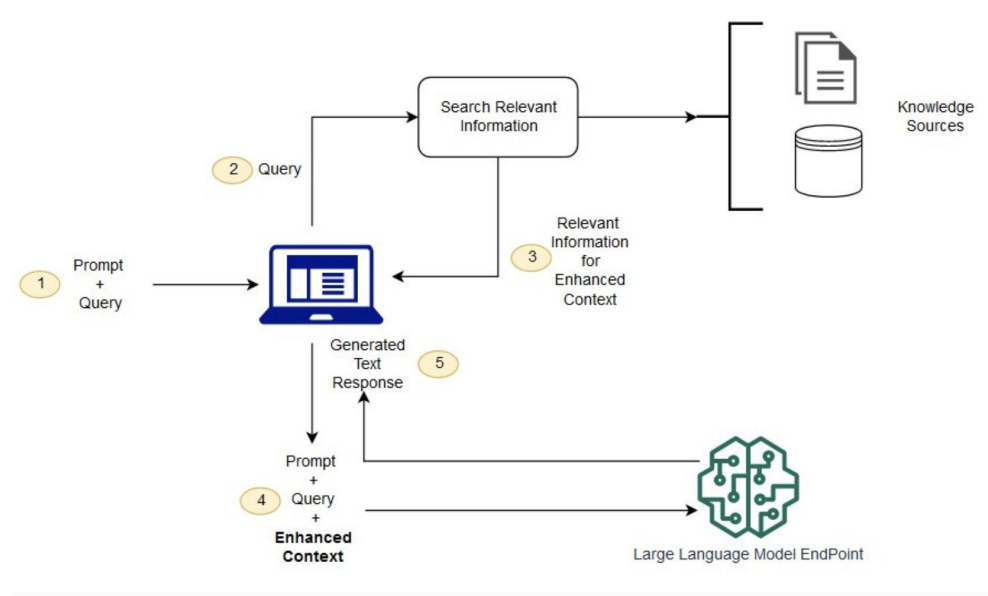


# Related Work

- Retrieval Methods in RAG
  - Can't perform well for long documents
- Truncation Methods in RAG
  - Not reliable since it's hard to determine the truncation cutoff.

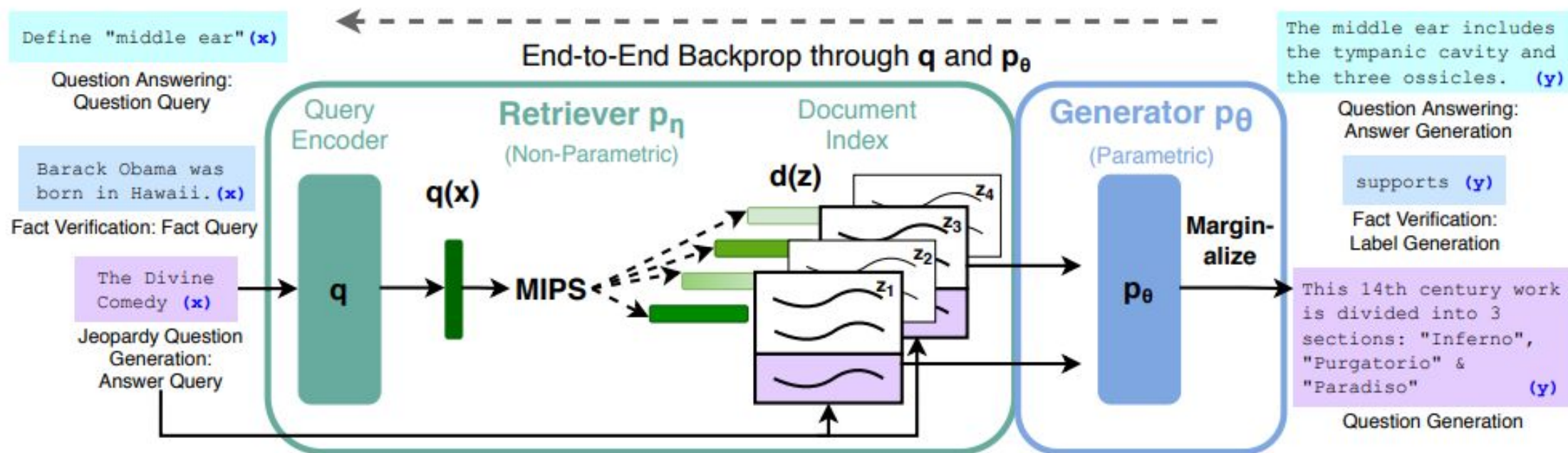
# Overview

- Benefits of RAG: improved model performance, information retrieval, and cost efficiency.
- Challenges: potential biases in datasets, computational complexity.





# RAG

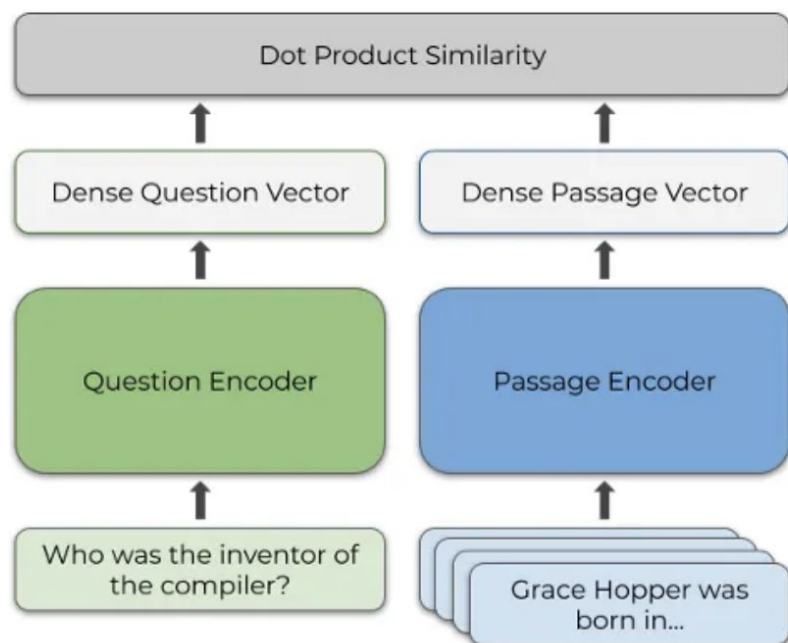


# RAG Sequence and RAG Token

$$p_{\text{RAG-Sequence}}(y|x) \approx \sum_{z \in \text{top-}k(p(\cdot|x))} p_{\eta}(z|x) p_{\theta}(y|x, z) = \sum_{z \in \text{top-}k(p(\cdot|x))} p_{\eta}(z|x) \prod_i^N p_{\theta}(y_i|x, z, y_{1:i-1})$$

$$p_{\text{RAG-Token}}(y|x) \approx \prod_i^N \sum_{z \in \text{top-}k(p(\cdot|x))} p_{\eta}(z|x) p_{\theta}(y_i|x, z, y_{1:i-1})$$

# RAG Retriever



Dense Passage Retrieval (DPR) architecture.

$$p_{\eta}(z|x) \propto \exp(\mathbf{d}(z)^{\top} \mathbf{q}(x))$$

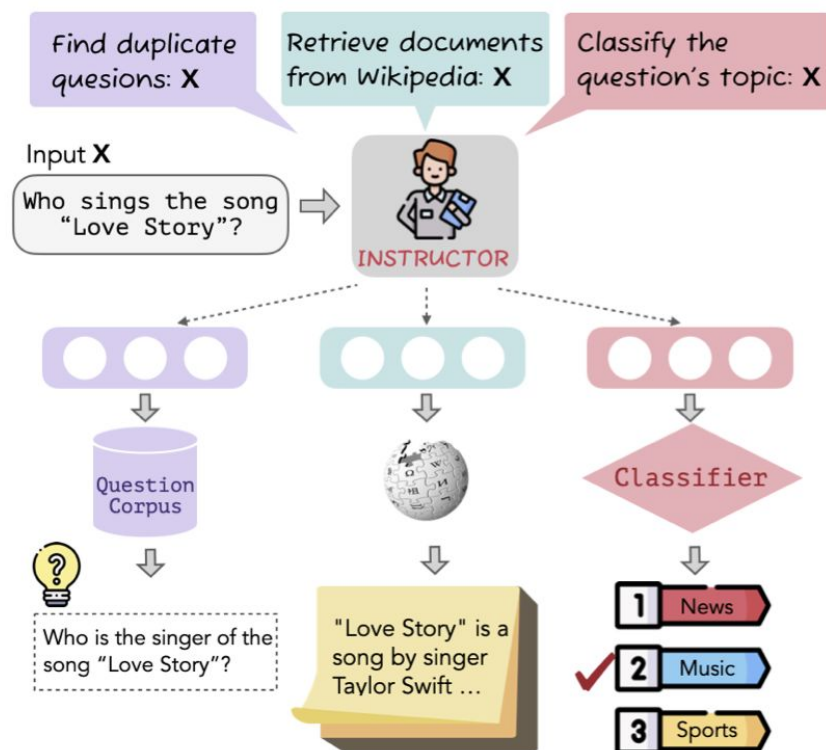
$$\mathbf{d}(z) = \text{BERT}_d(z), \quad \mathbf{q}(x) = \text{BERT}_q(x)$$

# Embedding Models

Overall								
Bitext Mining								
Classification								
Clustering								
Pair Classification								
Reranking								
Retrieval								
STS								
Summarization								
English								
Chinese								
French								
Polish								
Overall MTEB English leaderboard 🏆								
Metric: Various, refer to task tabs								
Languages: English								
Rank	Model	Model Size (Million Parameters)	Memory Usage (GB, fp32)	Embedding Dimensions	Max Tokens	Average (56 datasets)	Classification Average (12 datasets)	Clustering Average (11 datasets)
1	<a href="#">SFR-Embedding-Mistral</a>	7111	26.49	4096	32768	67.56	78.33	51.67
2	<a href="#">voyage-lite-02-instruct</a>	1220	4.54	1024	4000	67.13	79.25	52.42
3	<a href="#">GritLM-7B</a>	7242	26.98	4096	32768	66.76	79.46	50.61
4	<a href="#">e5-mistral-7b-instruct</a>	7111	26.49	4096	32768	66.63	78.47	50.26
5	<a href="#">google-gecko.text-embedding-ff</a>	1200	4.47	768	2048	66.31	81.17	47.48
6	<a href="#">GritLM-8v7B</a>	16702	173.08	1006	32768	65.66	78.53	50.11
40	<a href="#">instructor-large</a>	335	1.25	768	512	61.59	73.86	45.29

<https://huggingface.co/spaces/mteb/leaderboard>

# Instructor Model



```
import pickle
with open('NQDataset_with_ContentEmbeddings.pkl', 'wb') as f:
    pickle.dump(df_copy, f)
```

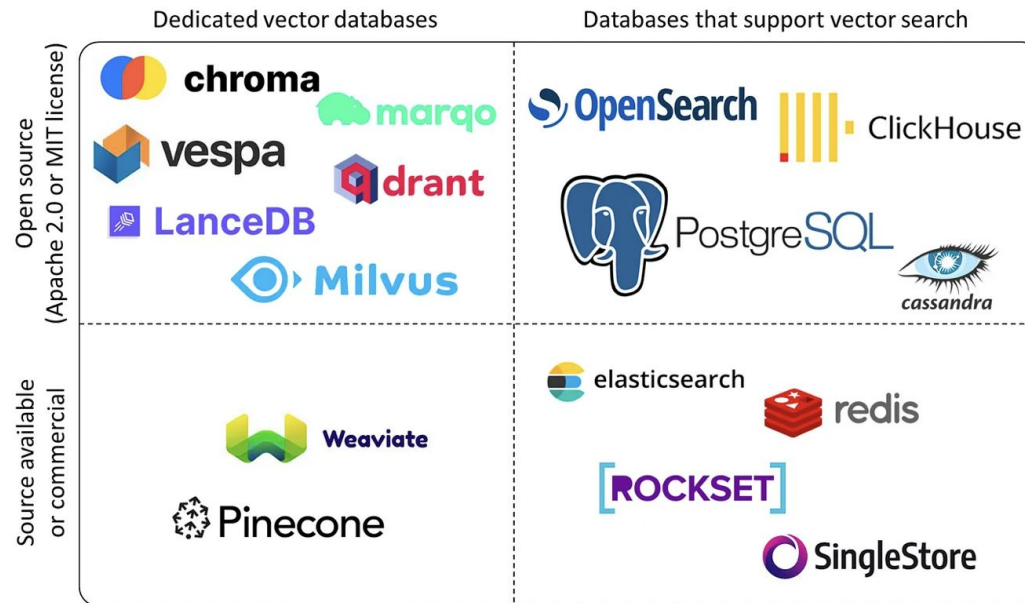
```
with open('NQDataset_with_ContentEmbeddings.pkl', 'rb') as f:
    loaded_df = pickle.load(f)
```

```
loaded_df.iloc[0]
```

```
Unnamed: 0      0
query          wolf of wall street number of f words
long_answer    Film  Year  Fuck count  Minutes  Uses / mi...
short_answer                    569
title      List of films that most frequently use the wor...
bert_title  list of films that most frequently use the wor...
abstract    The use of profanity in films has always been ...
content      This is a list of non-pornographic , English l...
url          https://en.wikipedia.org/wiki/List%20of%20film...
index                                             109430
content_embedding  [-0.024242813, 0.019909445, -0.042916078, 0.01...
Name: 0, dtype: object
```

# Storage

- Vector index (ex. Faiss) and vector database options
- Used Wikipedia dataset



The landscape of vector databases.



# Storage

**embeddings** ●

METRIC	DIMENSIONS	HOST
euclidean	768	https://embeddings-d7j98sx.svc.aped-4627-b74a.pinecone.io

CLOUD	REGION	VECTOR COUNT
AWS	us-east-1	<b>2,195</b>

**BROWSER** METRICS NAMESPACES (1)

Namespace  
( Default )

Query by Vector ▼

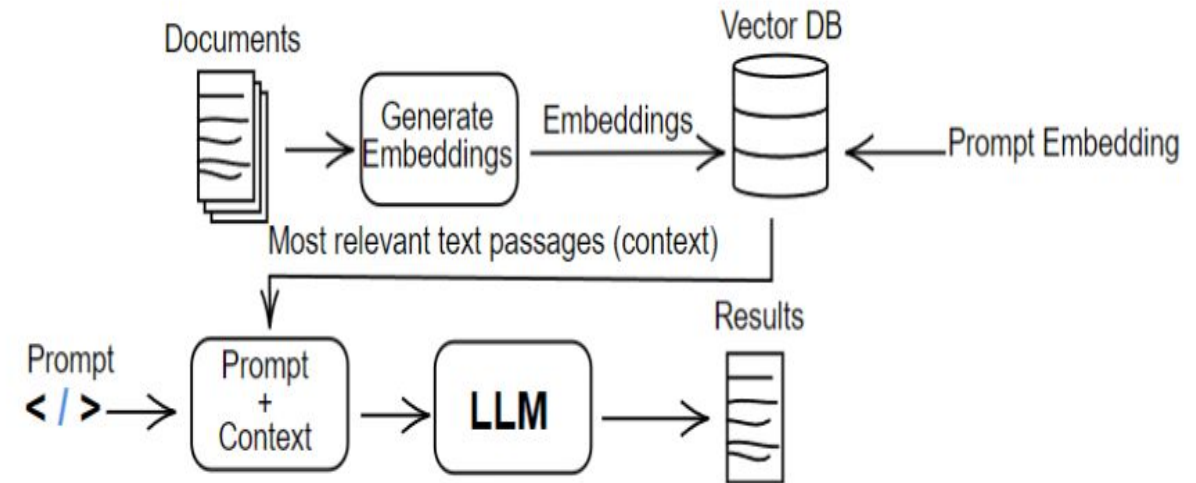
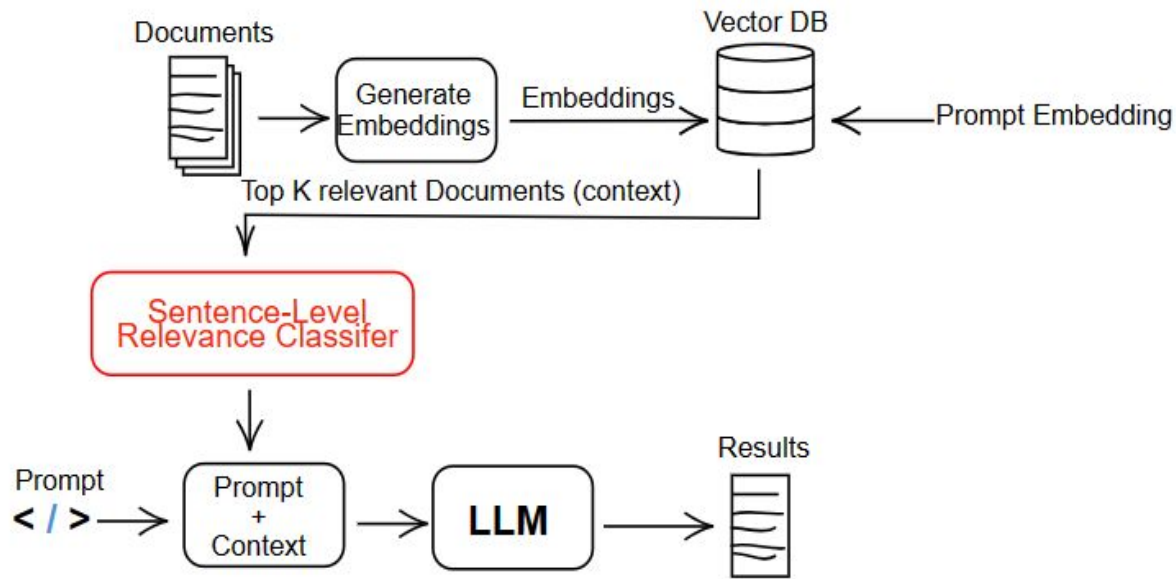
vector  
0.43,0.57,0.81,0.47,0.52,0.1,0.45,0.77,0.59,0.34,0.76,0.43

Top K\*  
10

Query

sending upsert requests: 100%  2355/2355 [01:23<00:00, 33.19it/s]  
{ 'upserted\_count' : 2355 }

# Proposed vs Existing Framework





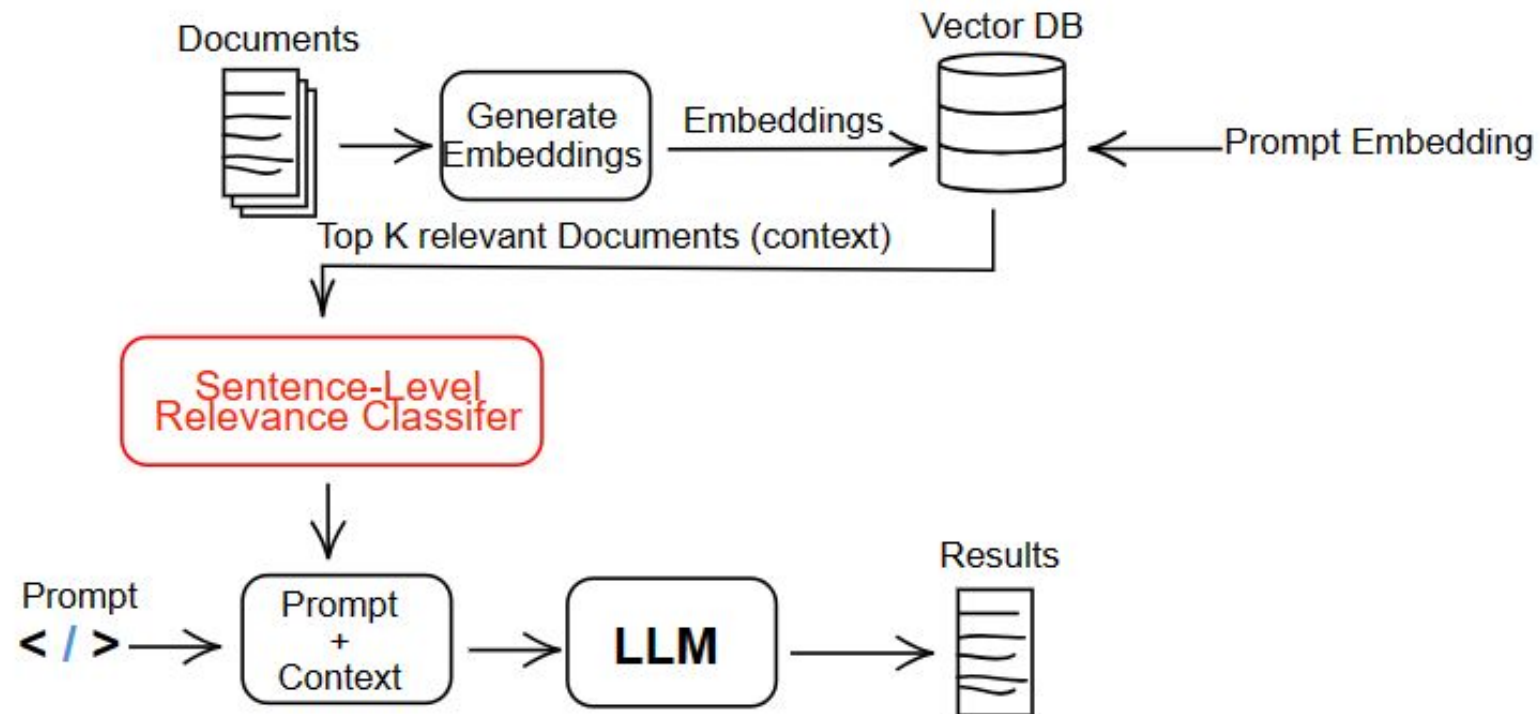
# Wikipedia Document Dataset



Unnamed: 0	query	long_answer	short_answer	title	bert_title	abstract	content	url	index	content_embedding
24	what are the toll roads called in mexico	This is a list of autopistas , or tolled ( quota ) highways , in Mexico . T	autopistas	List of Mexican aut	list of mexican autopistas	This is a list (	Many federal	https://en.wi	109631	4.96650993e-02 5.33997050e-02 4.75333607e-02 2.25149544e-02
25	what are the top five wine producing states	2016 production of still wine State Production ( gal ) Production	California Washington New York Penn	American wine	american wine	American wine	The North Ar	https://en.wi	98864	[[-7.31019536e-03 -1.26333470e-02 -1.06264362e-02 1.31341349e-02
26	who sings the theme song for living single	Living Single Season 1 DVD cover Created by Yvette Lee Bowser	performed by	Living Single	living single	Living Single	Throughout	https://en.wi	15276	[[-0.01875372 -0.00664845 -0.01411123 0.01214166 0.03976656 0.0153
27	what type of reproduction do whiptail lizards use	summer , and hatching approximately eight weeks later . The New M -		New Mexico whipt	new mexico whiptail	Cnemidoph	The New Mex	https://en.wi	103123	[[-2.50664949e-02 -7.14592077e-03 -2.15893276e-02 3.62700666e-03
28	where are the summer olympics held in 2012	The 2012 Summer Olympics , formally the Games of the XXX Olympic Olympiad		2012 Summer Olym	2012 summer olympics	The 2012 Su	Followed a bid headed by former Olympic champion Sebastian Coe and then Mayor of London Ken Livingstone, Lond			

# Methodology

- Framework:

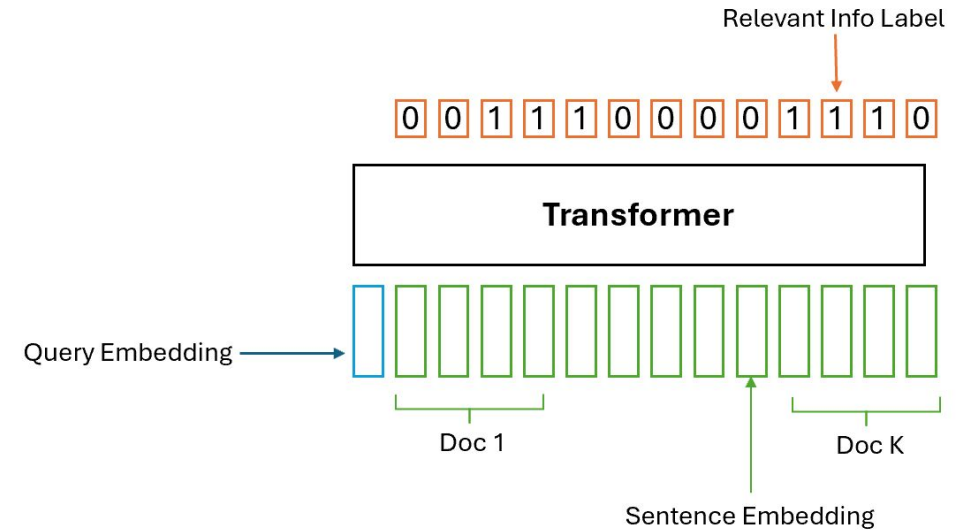


# Methodology

- Motivation:
  - RAG relies on highly relevant information.
  - Traditional methods like reranking and truncations have some limitations
- We propose a model to identify most relevant information.

# Proposed Method

- SLRC: Sentence-Level Relevance Classifier
- Input:
  - Query Embedding and sentences embedding of Top K documents
- Output:
  - Relevant info label: if the sentences is relevant with the query.



# Evaluation

## Natural Questions Benchmark

### Example 1

**Question:** what color was john wilkes booth's hair

**Wikipedia Page:** John\_Wilkes\_Booth

**Long answer:** Some critics called Booth “the handsomest man in America” and a “natural genius”, and noted his having an “astounding memory”; others were mixed in their estimation of his acting. He stood 5 feet 8 inches (1.73 m) tall, had jet-black hair, and was lean and athletic. Noted Civil War reporter George Alfred Townsend described him as a “muscular, perfect man” with “curling hair, like a Corinthian capital”.

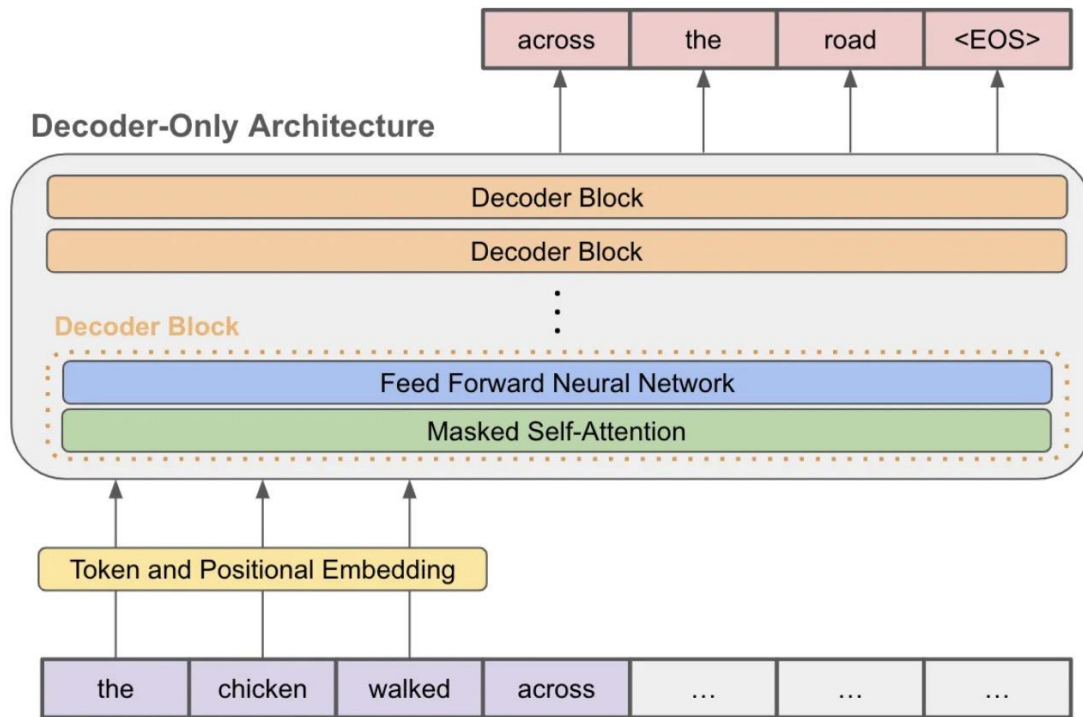
**Short answer:** jet-black

### Example 1

**Question:** who played will on as the world turns **Long answer:** William “Will” Harold Ryan Munson is a fictional character on the CBS soap opera As the World Turns. He was portrayed by Jesse Soffer on recurring basis from September 2004 to March 2005, after which he got a contract as a regular. Soffer left the show on April 4, 2008 and made a brief return in July 2010. **Judgment:** Correct. **Justification:** It is clear beyond a reasonable doubt that the answer is correct.

<https://storage.googleapis.com/gweb-research2023-media/pubtools/pdf/1f7b46b5378d757553d3e92ead36bda2e4254244.pdf>

# LLM (gpt-2)

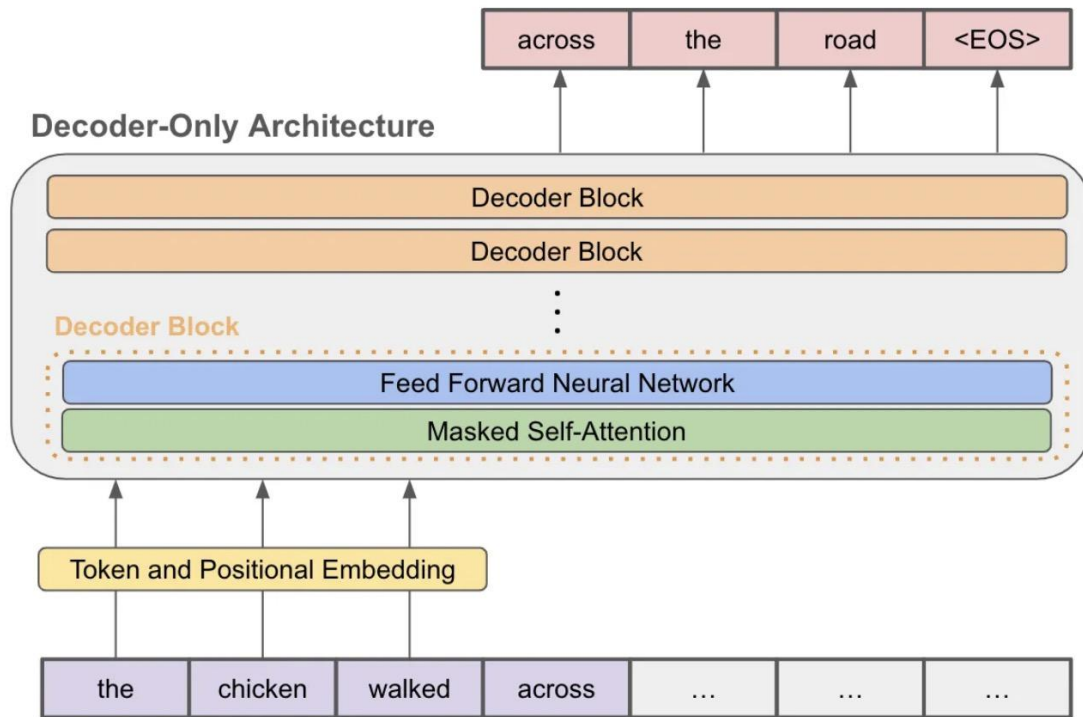


Decoder-only transformer architecture (created by author)

GPT-2 is a decoder-only transformer architecture. It removes the following components from the transformer:

- The entire encoder module
- All encoder-decoder self-attention modules in the decoder

# LLM (gpt-2)



Decoder-only transformer architecture (created by author)

After these components have been removed, each layer of the decoder simply consists of

- a masked self-attention layer followed by a feed forward neural network. Stacking several of such layers on top of each other forms a deep, decoder-only transformer architecture, such as those used for GPT or GPT-2.

# Evaluation Metrics

We choose two metrics to measure the similarity of the machine-translated text to a set of high quality reference translations. Both metrics range from 0 to 1.

- BLEU\_SCORE: A precision based measure.
- ROUGE\_SCORE: The harmonic mean of recall and precision.
- Perplexity score ?



# Evaluation Benchmarks

We choose the RAG and GPT-2 as the benchmarks where GPT-2 is pre-trained on large scale datasets.

We compare the performance of:

1. Pretrained GPT-2 model
2. Pretrained GPT-2 model + RAG prompts (top-3 documents)
3. Fine-tuned GPT-2 model using RAG prompts (top-3 documents)
4. LLaMa 7B + RAG prompts
5. LLaMa 7B

# Evaluation settings

We customize the dataset by selecting 2335 query&answer pairs from NQDataset. The dataset is splitted into two parts:

- Train Split: 1335 query&answer pairs
- Test Split: 1000 query&answer pairs

What's more, the indices of top-3 document is added for each query&answer pairs.

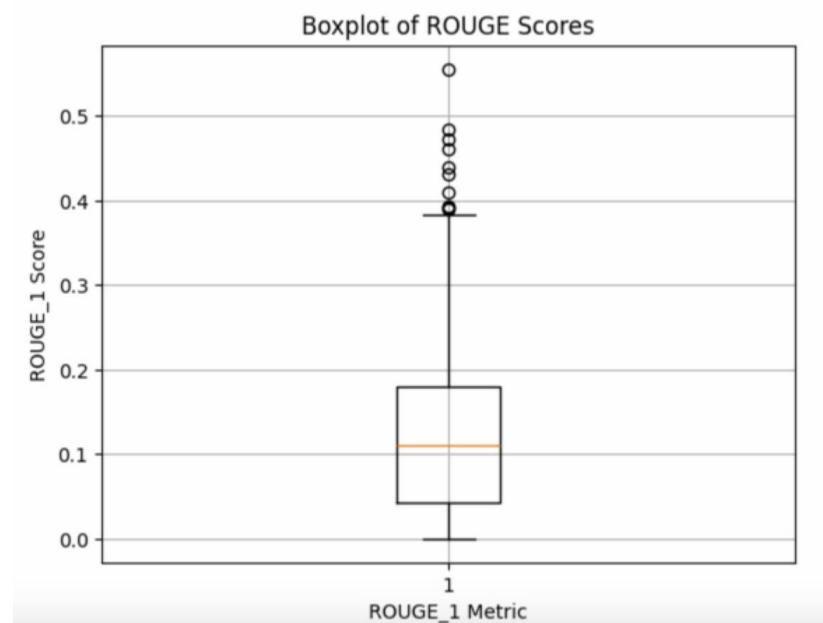
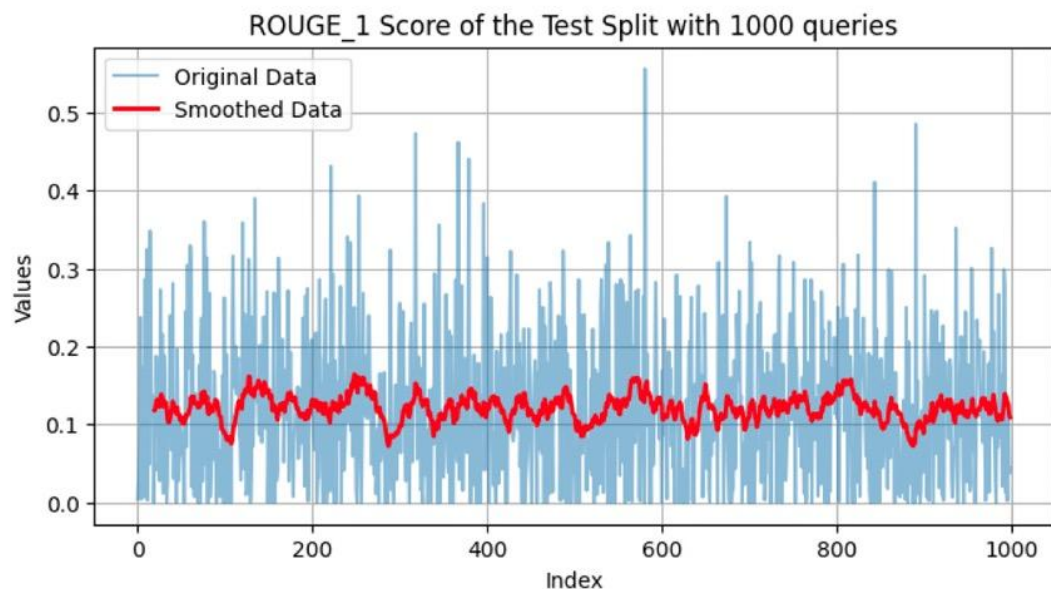
query	long_answer	Topk_indices
wolf of wall street number of f words	Film Year Fuck count Minutes Uses / minute Source	[ 0 178 684]

Example of our customized dataset

# Evaluation Results — ROUGE of EXP1

The following results show the ROUGE 1 score and BLEU score of the pre-trained GPT-2 model on the test split of our customized dataset.

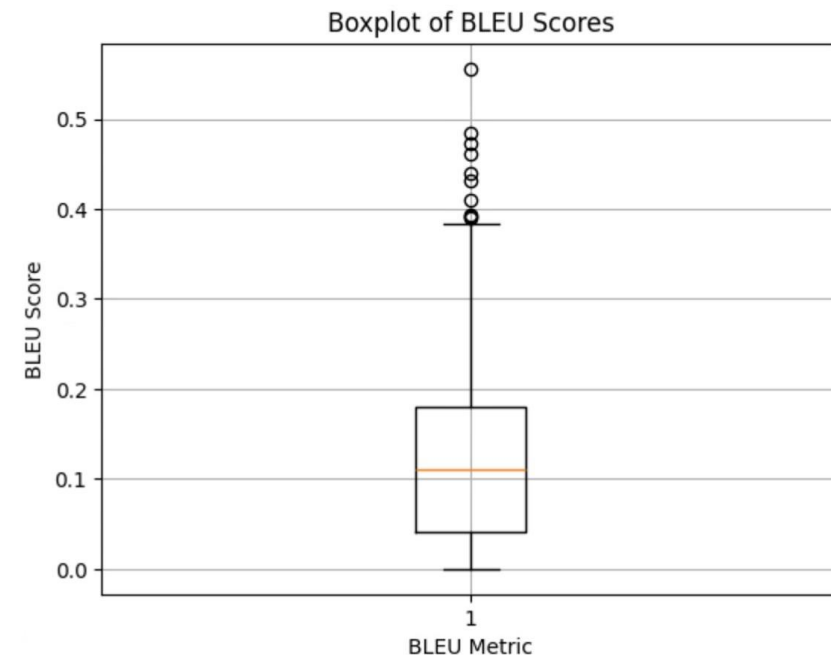
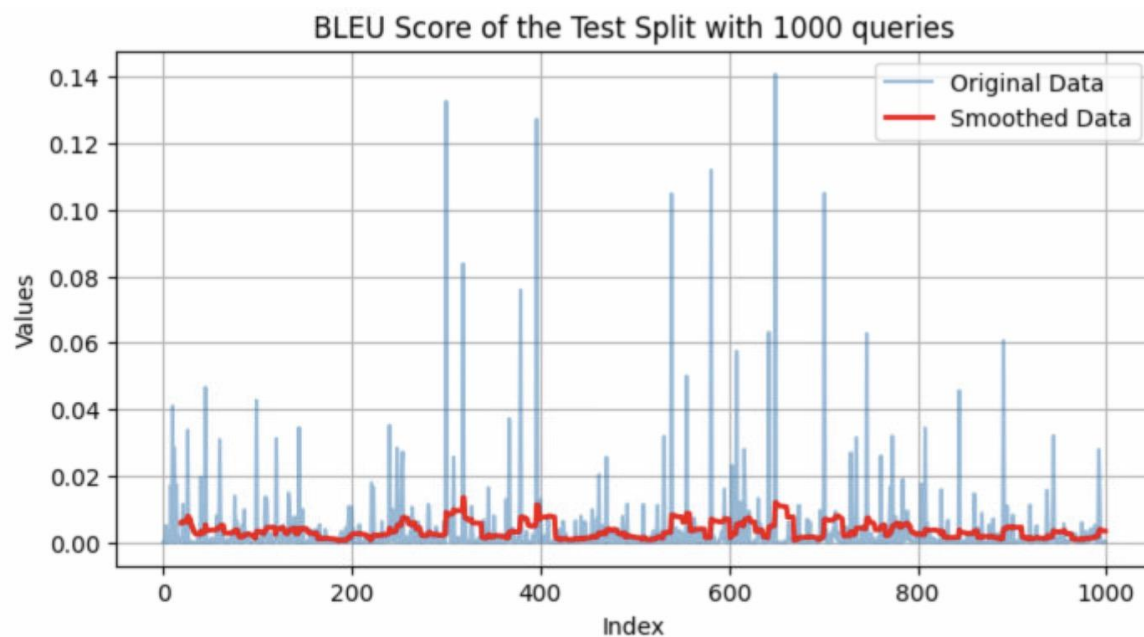
It shows that the overall performance of pre-trained GPT-2 is bad.



# Evaluation Results — BLEU of EXP1

The following results show the ROUGE 1 score and BLEU score of the pre-trained GPT-2 model on the test split of our customized dataset.

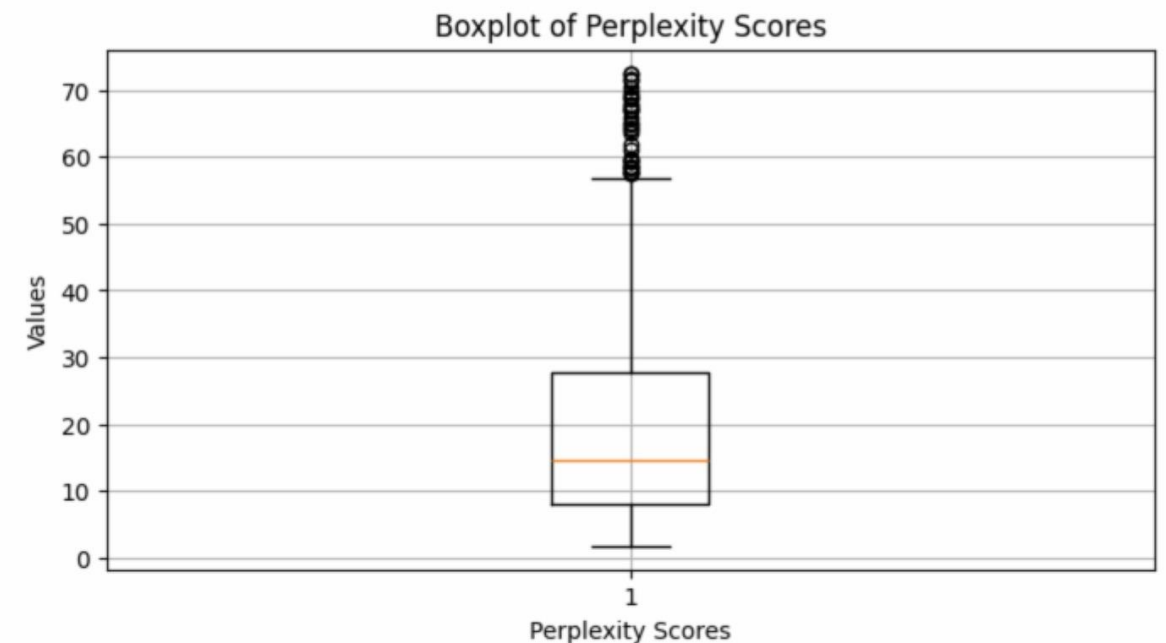
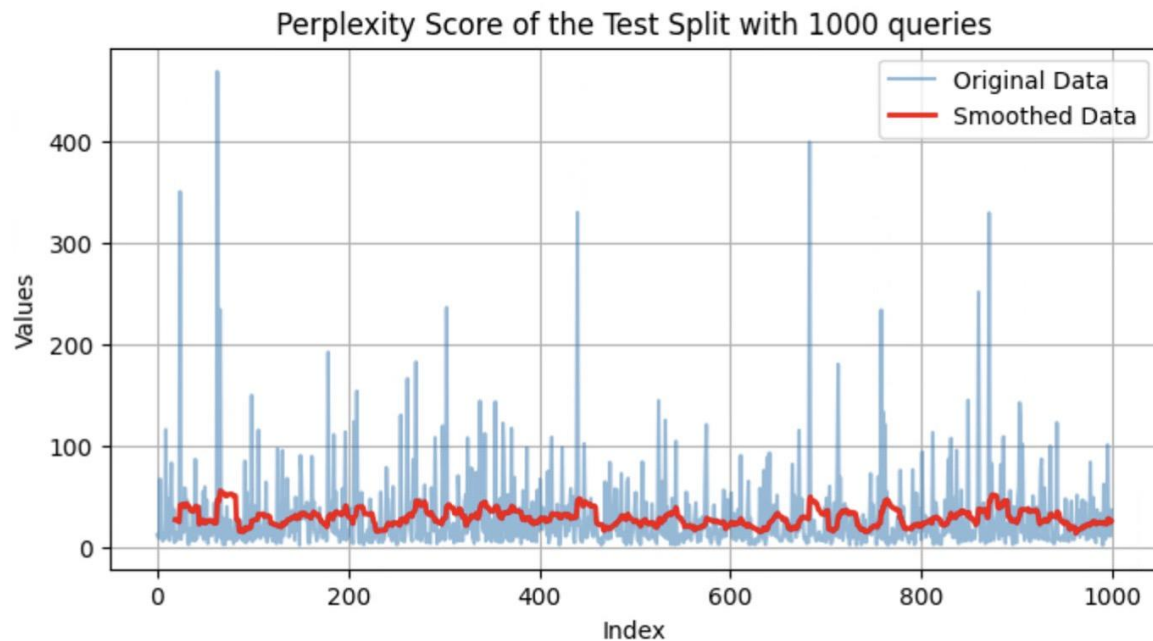
It shows that the overall performance of pre-trained GPT-2 is bad.



# Evaluation Results — Perplexity of EXP1

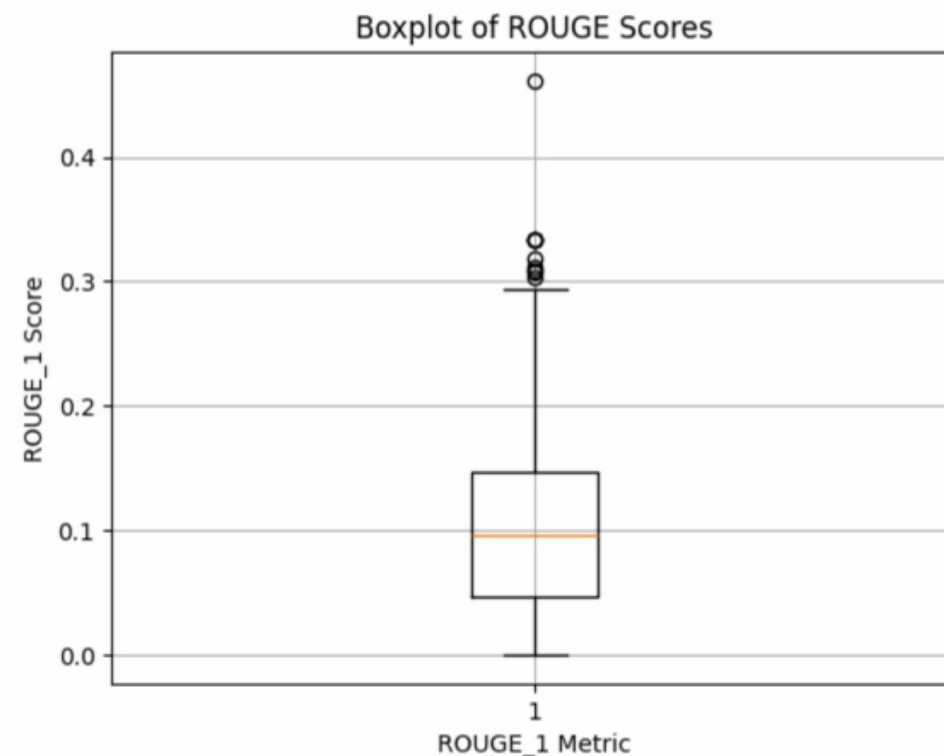
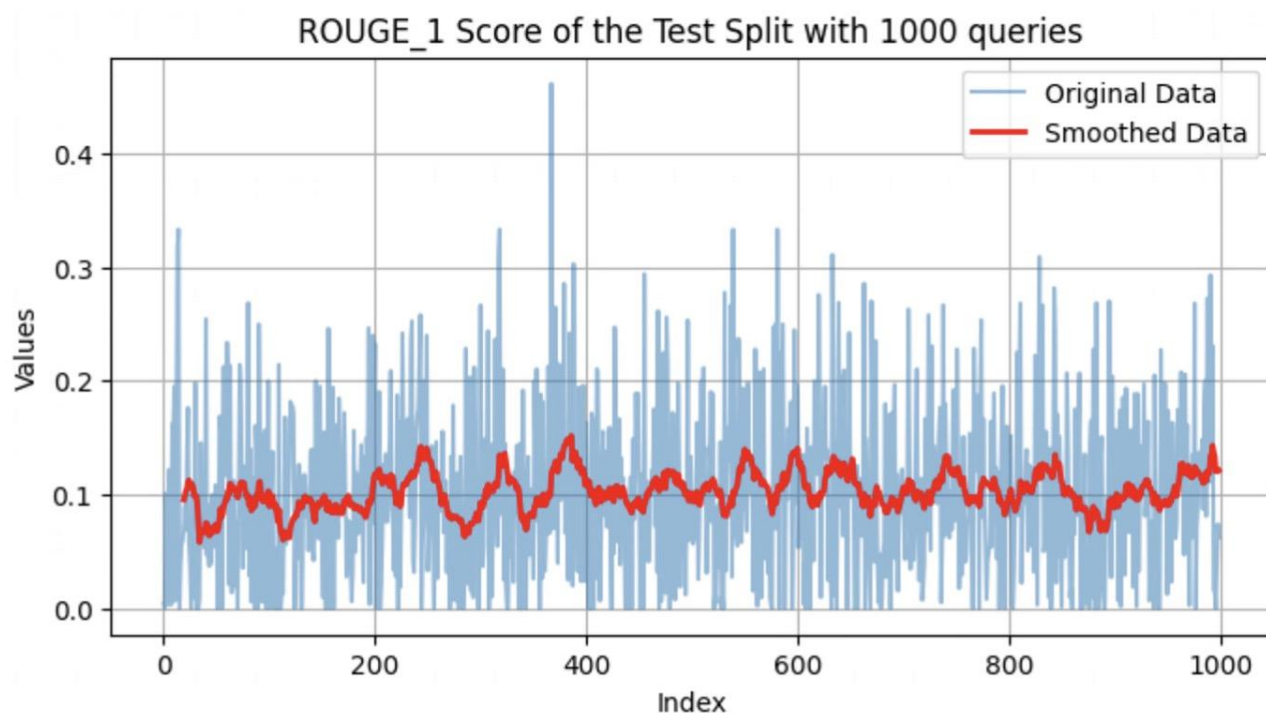
The following results show the ROUGE 1 score and BLEU score of the pre-trained GPT-2 model on the test split of our customized dataset.

It shows that the overall performance of pre-trained GPT-2 is bad.



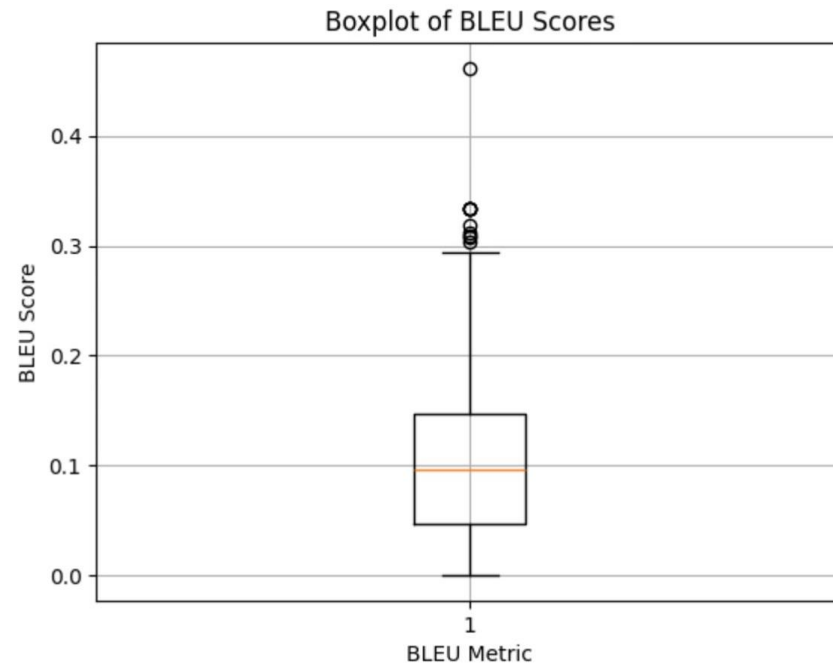
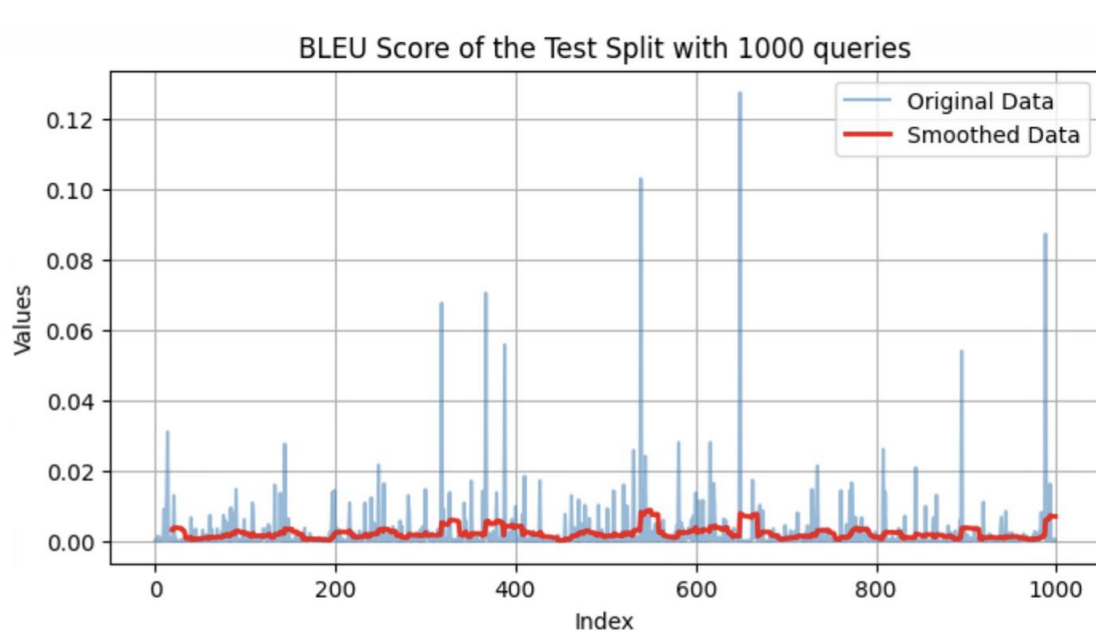
# Evaluation Results — ROUGE of EXP2

In terms of the bad performance of pretrained GPT-2 mode, we use RAG prompts (top-3 documents) to help improve the scores.



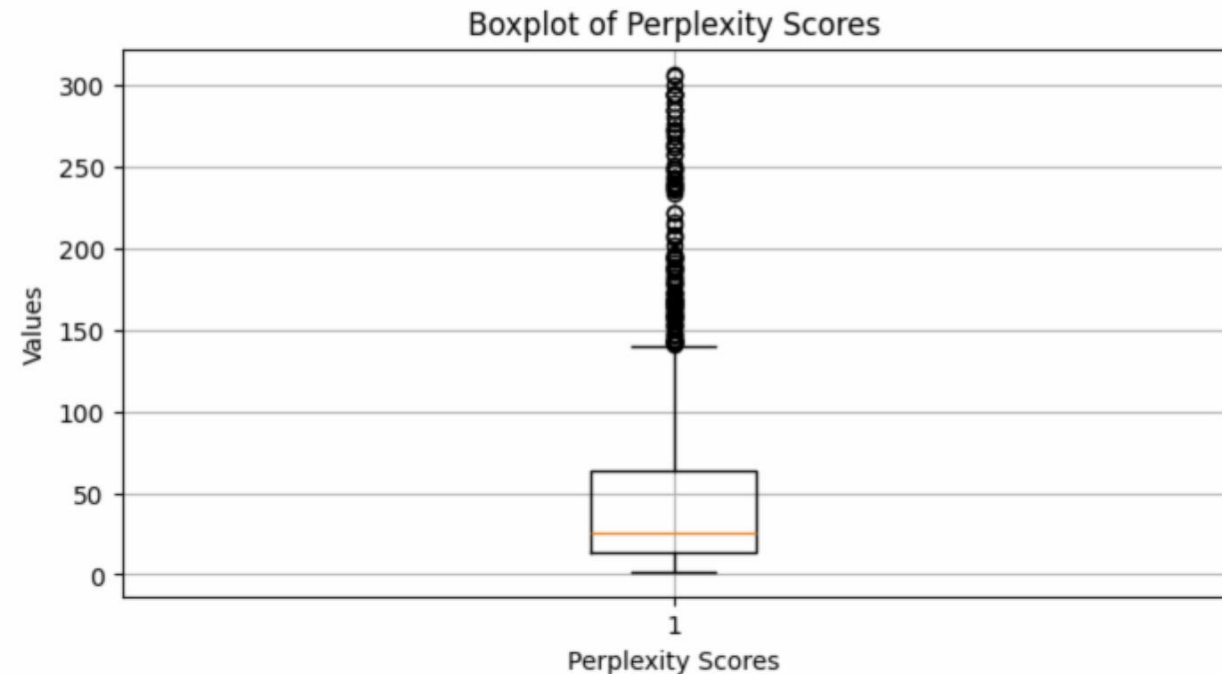
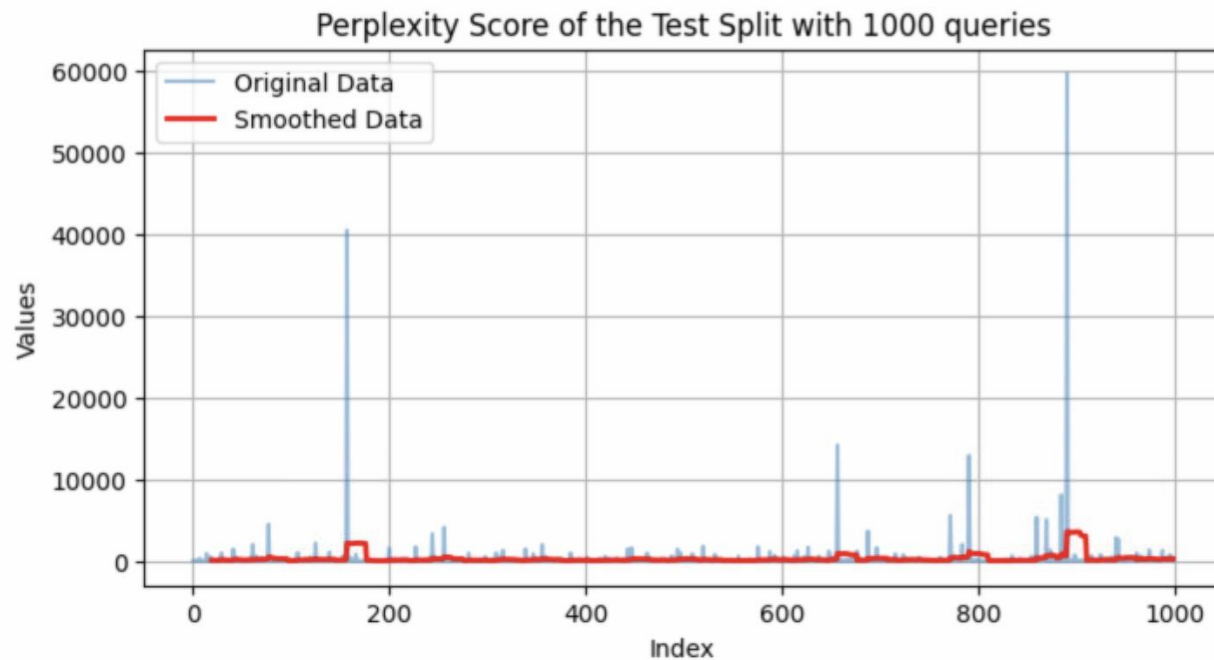
# Evaluation Results — BLEU of EXP2

In terms of the bad performance of pretrained GPT-2 mode, we use RAG prompts (top-3 documents) to help improve the scores.



# Evaluation Results — Perplexity of EXP2

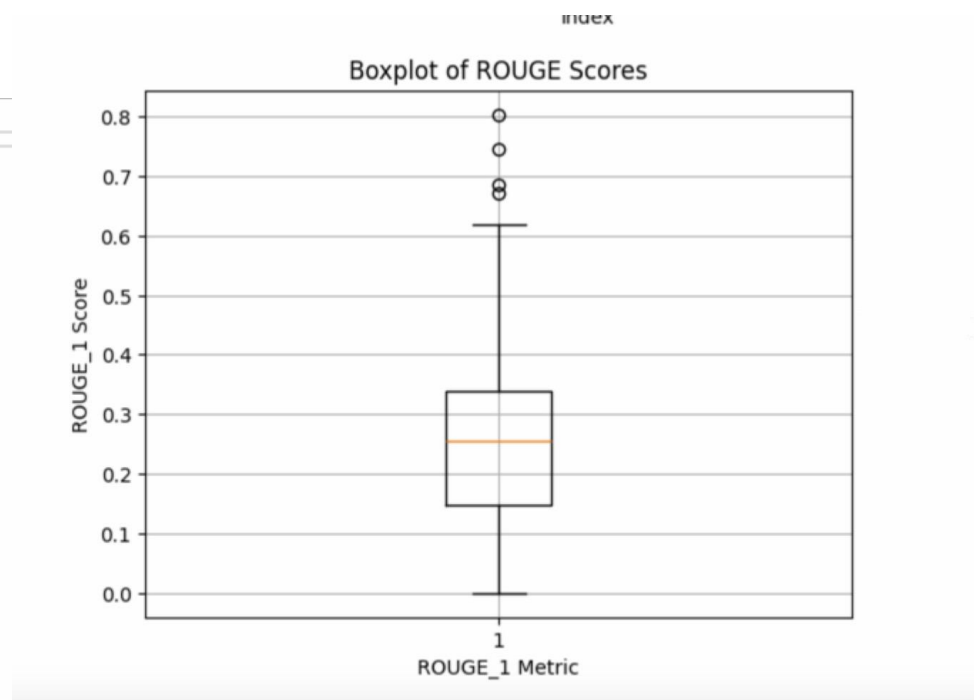
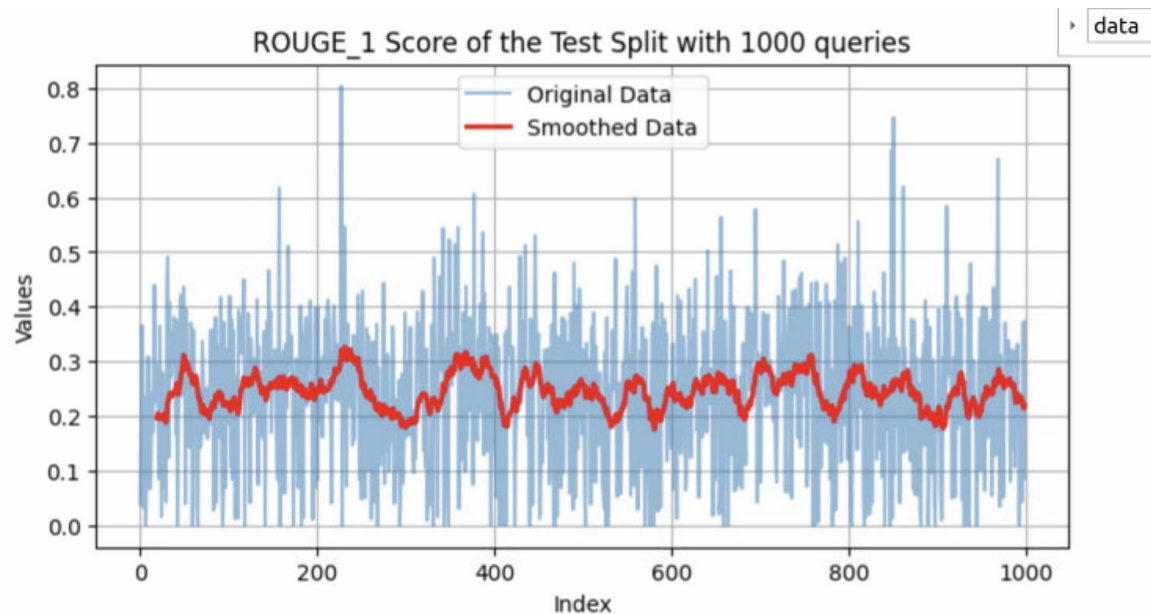
In terms of the bad performance of pretrained GPT-2 mode, we use RAG prompts (top-3 documents) to help improve the scores.





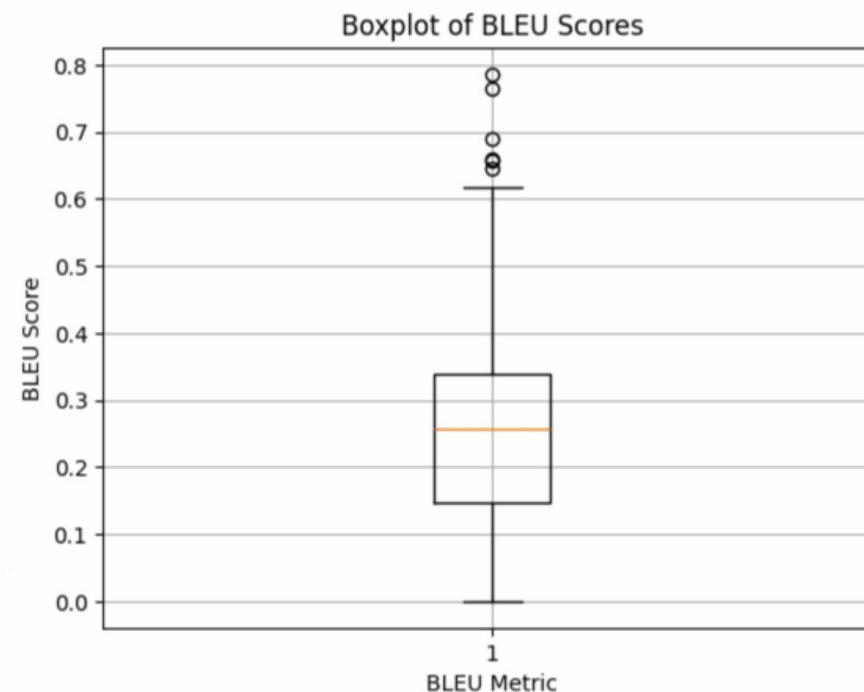
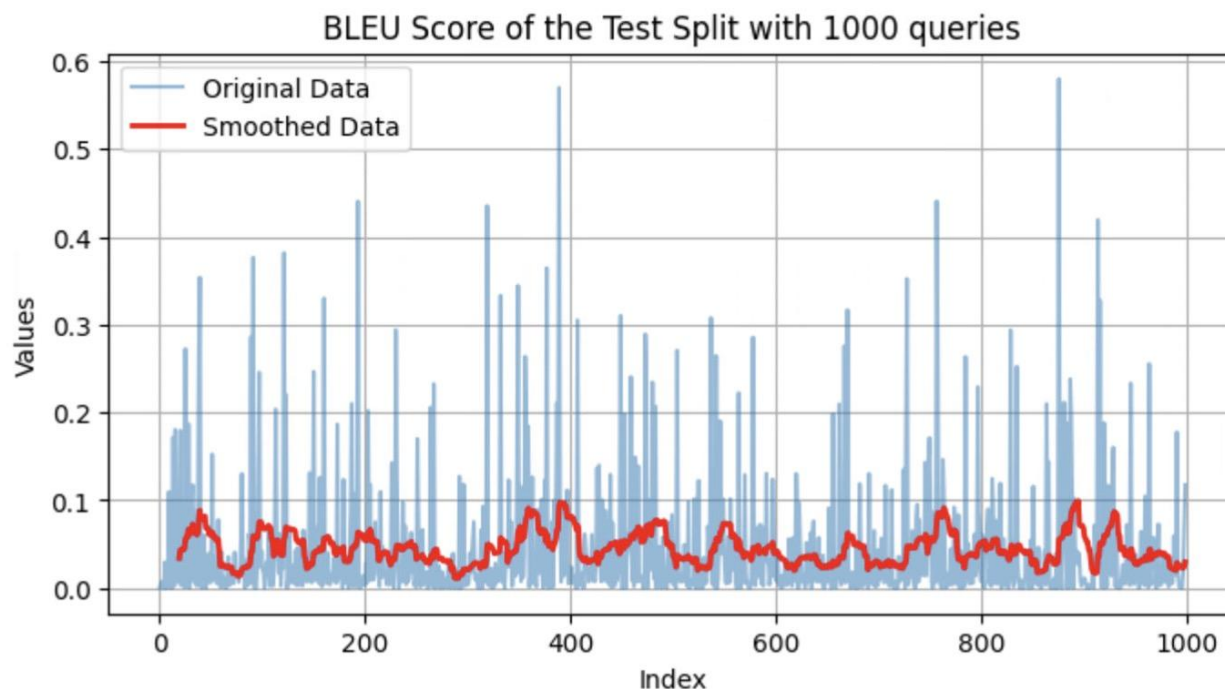
# Evaluation Results — ROUGE of EXP3

In terms of the bad performance of pretrained GPT-2 mode, we use a classifier to select the most relevant sentences from top-3 documents as the prompt to help improve the scores.



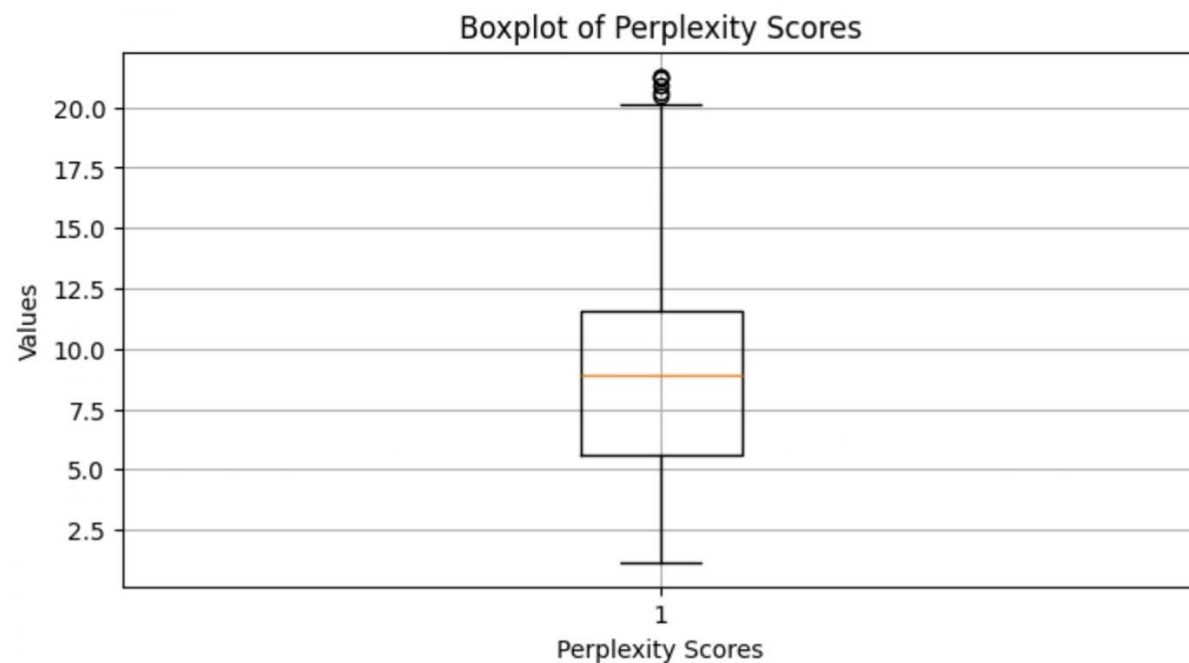
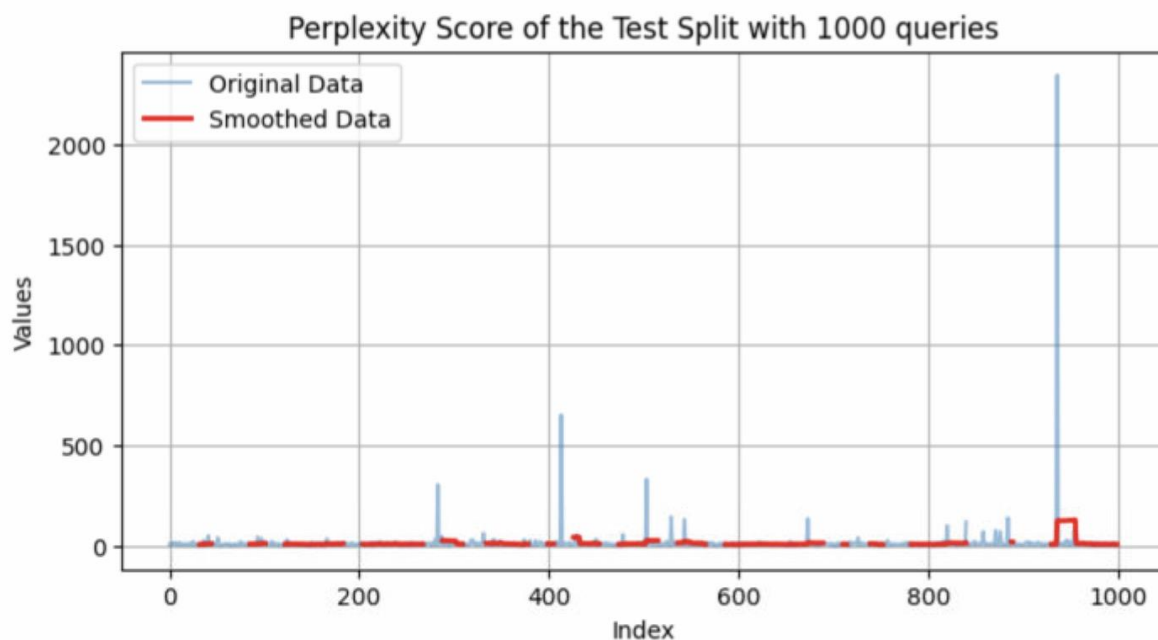
# Evaluation Results — BLEU of EXP3

In terms of the bad performance of pretrained GPT-2 mode, we use a classifier to select the most relevant sentences from top-3 documents as the prompt to help improve the scores.

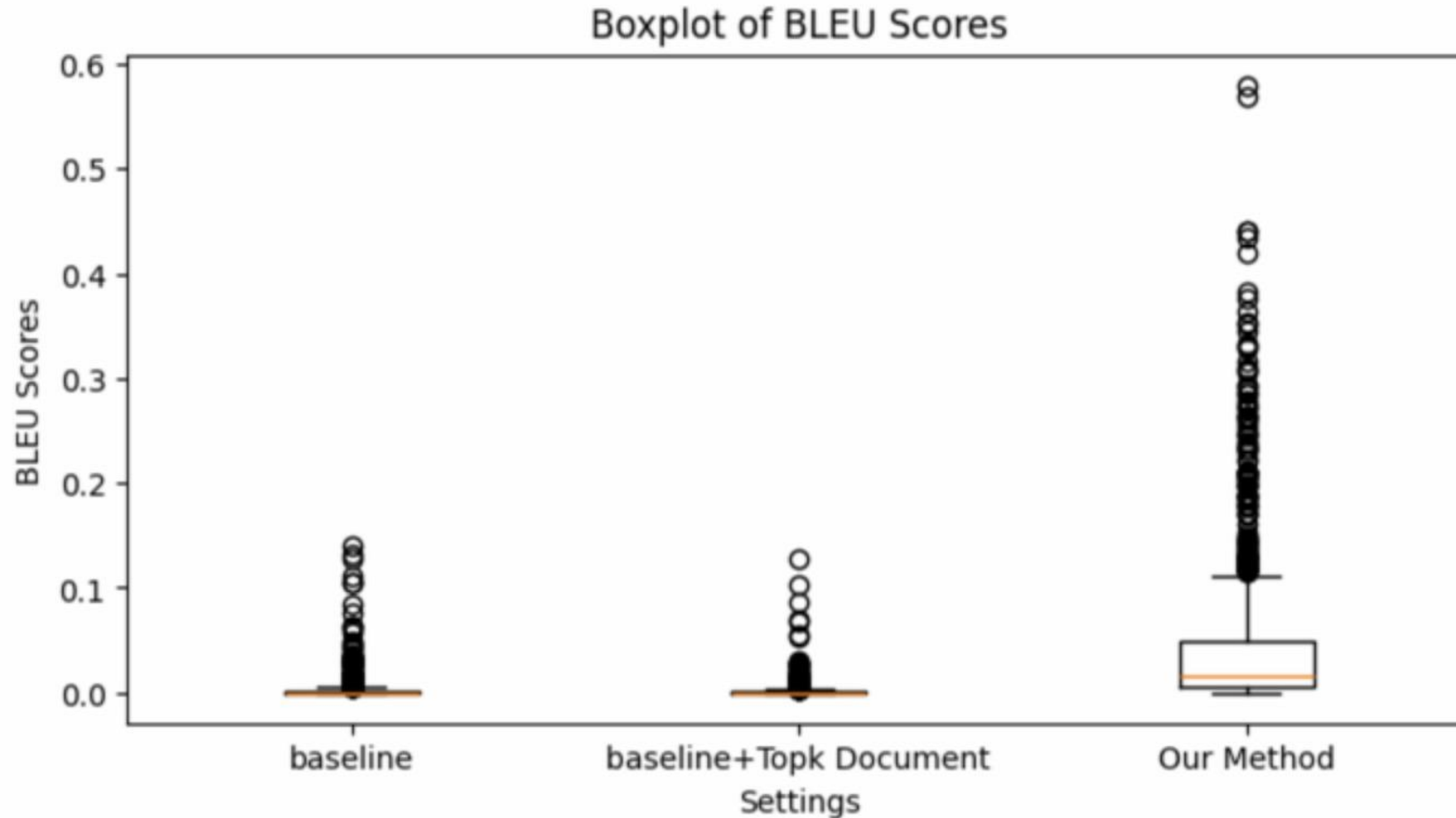


# Evaluation Results — Perplexity of EXP3

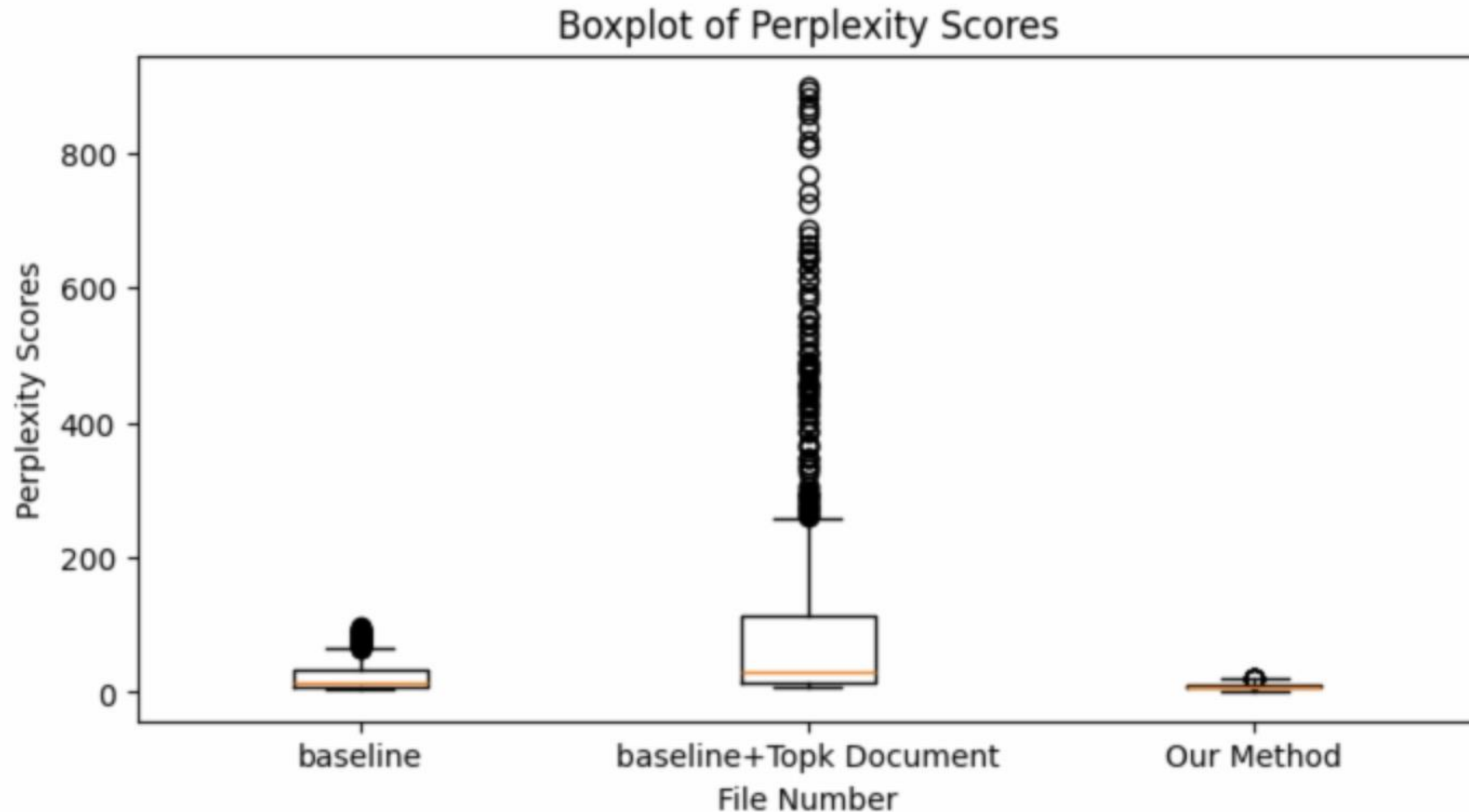
In terms of the bad performance of pretrained GPT-2 mode, we use a classifier to select the most relevant sentences from top-3 documents as the prompt to help improve the scores.



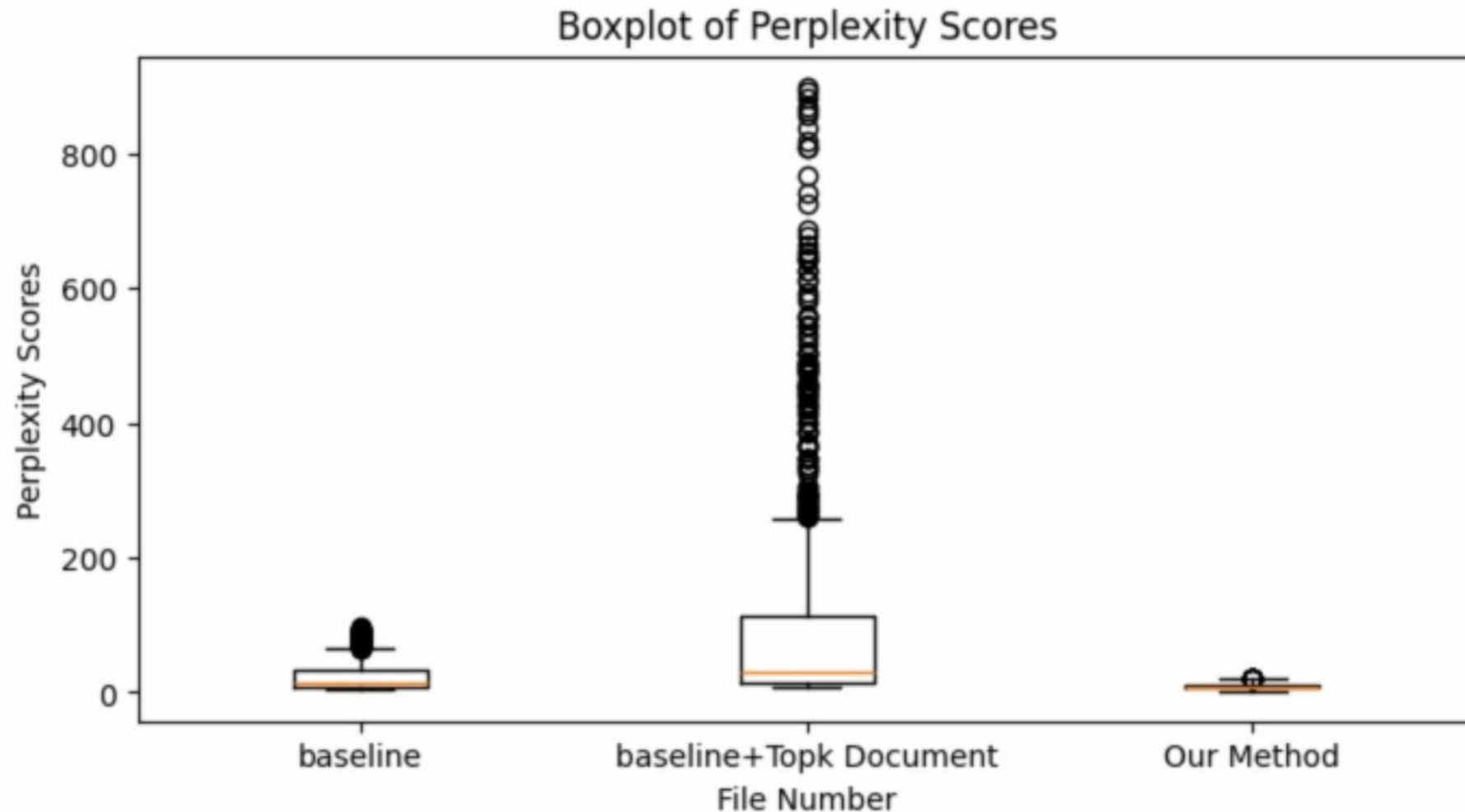
# Evaluation Results — BLEU Comparison of GPT2



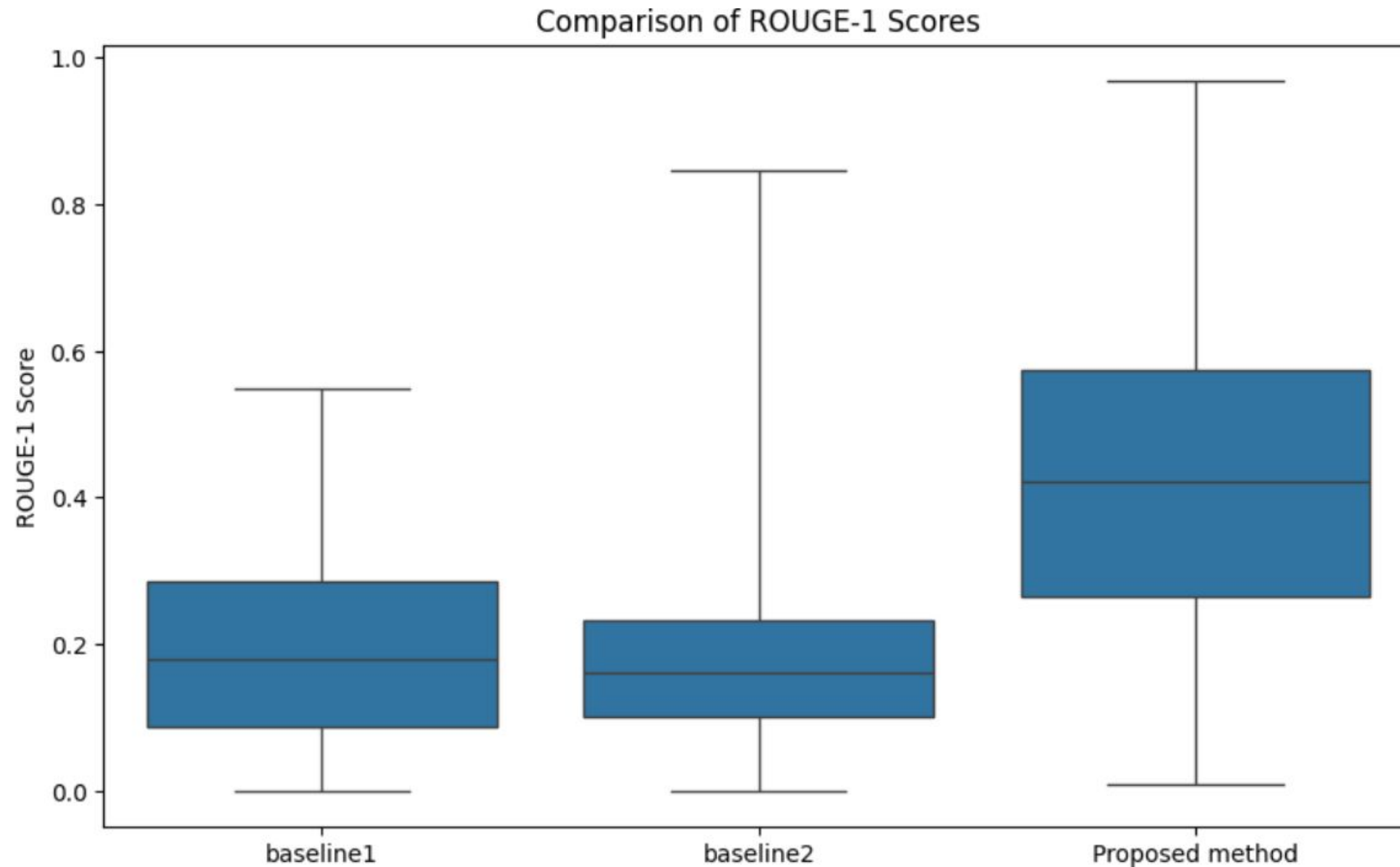
# Evaluation Results — ROUGE Comparison of GPT2



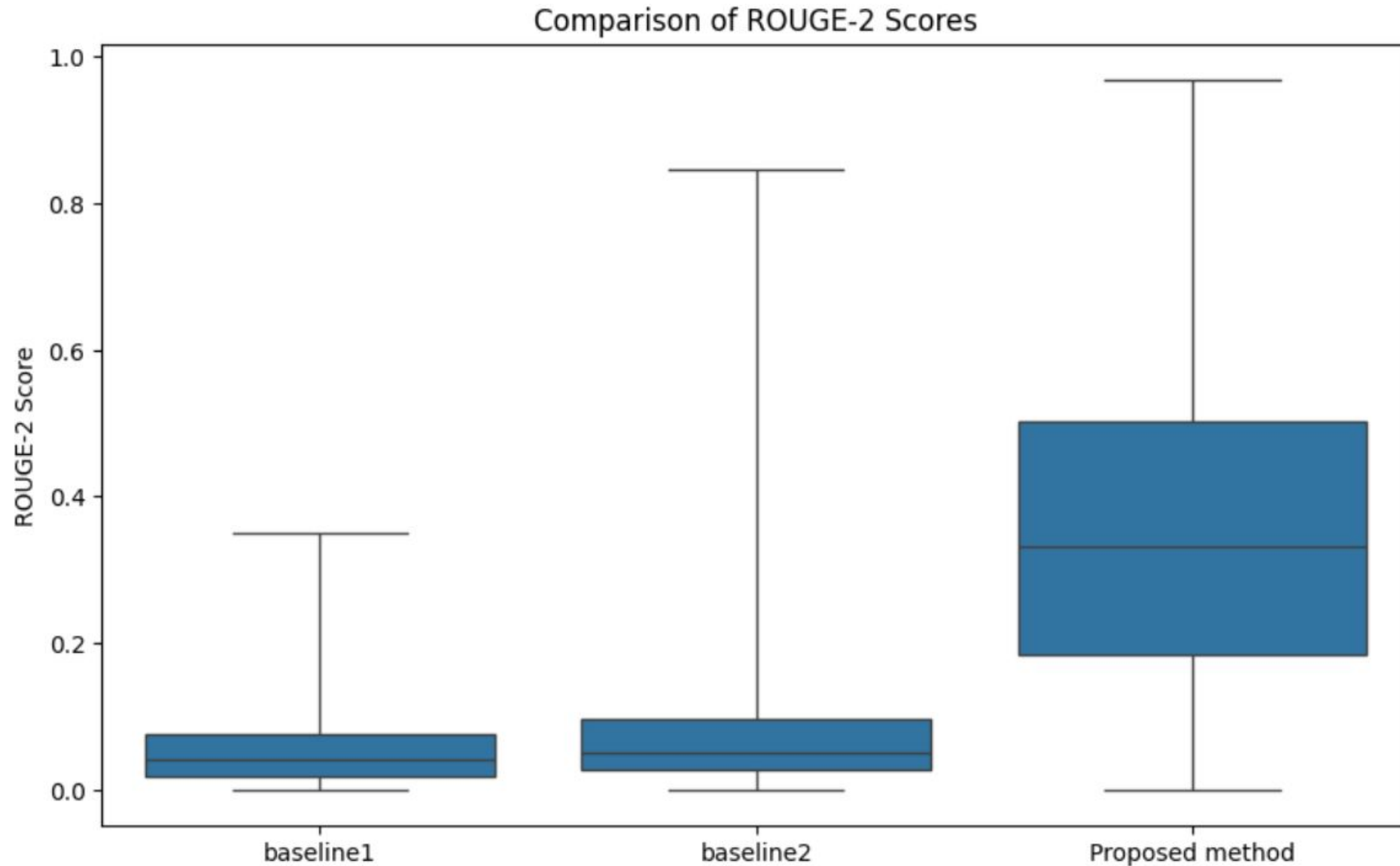
# Evaluation Results — Perplexity Comparison of GPT2



# Evaluation Results - Rouge1 Comparison on LLaMa2

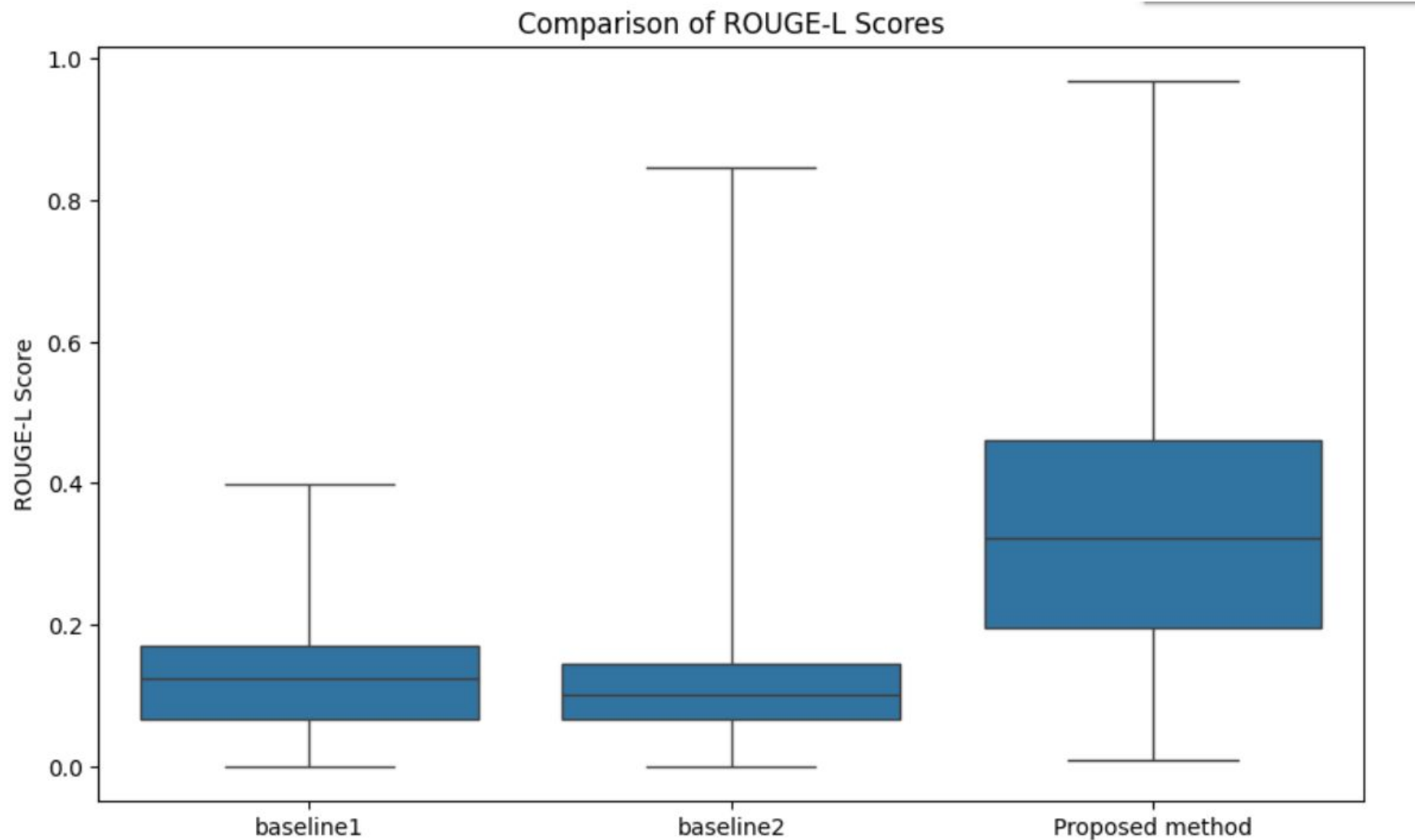


# Evaluation Results - Rouge2 Comparison on LLaMa2

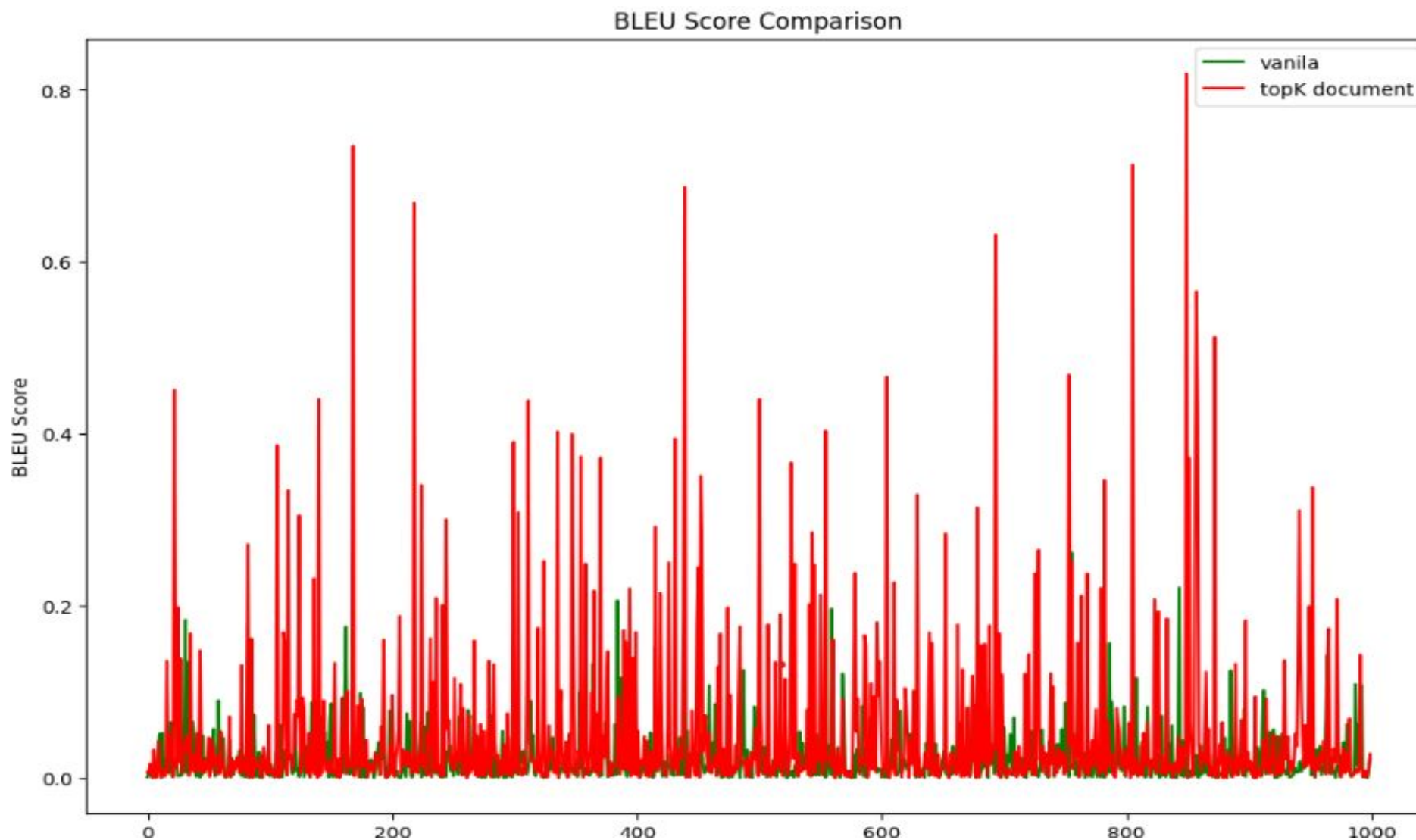




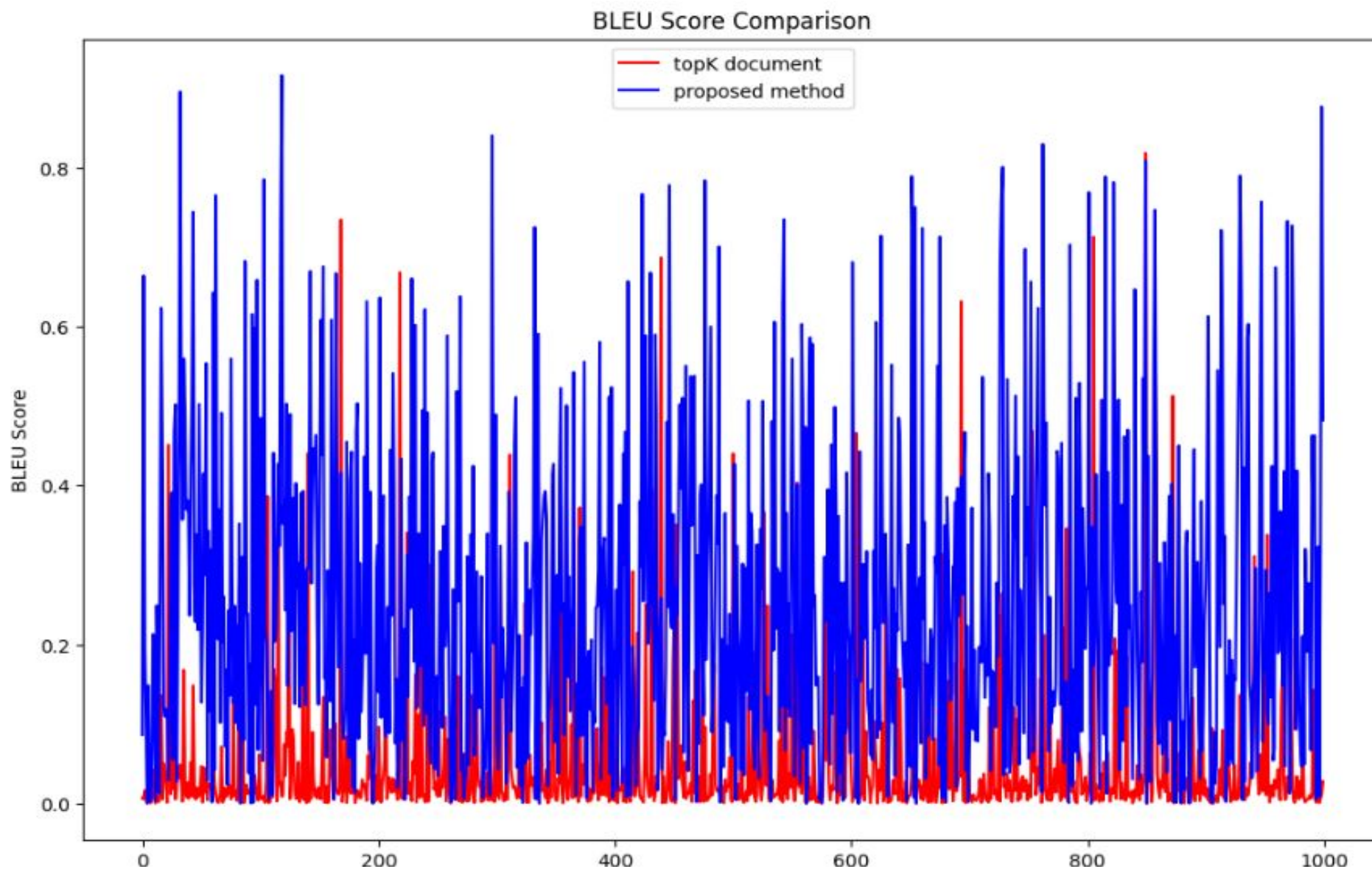
# Evaluation Results - RougeL Comparison on LLaMa2



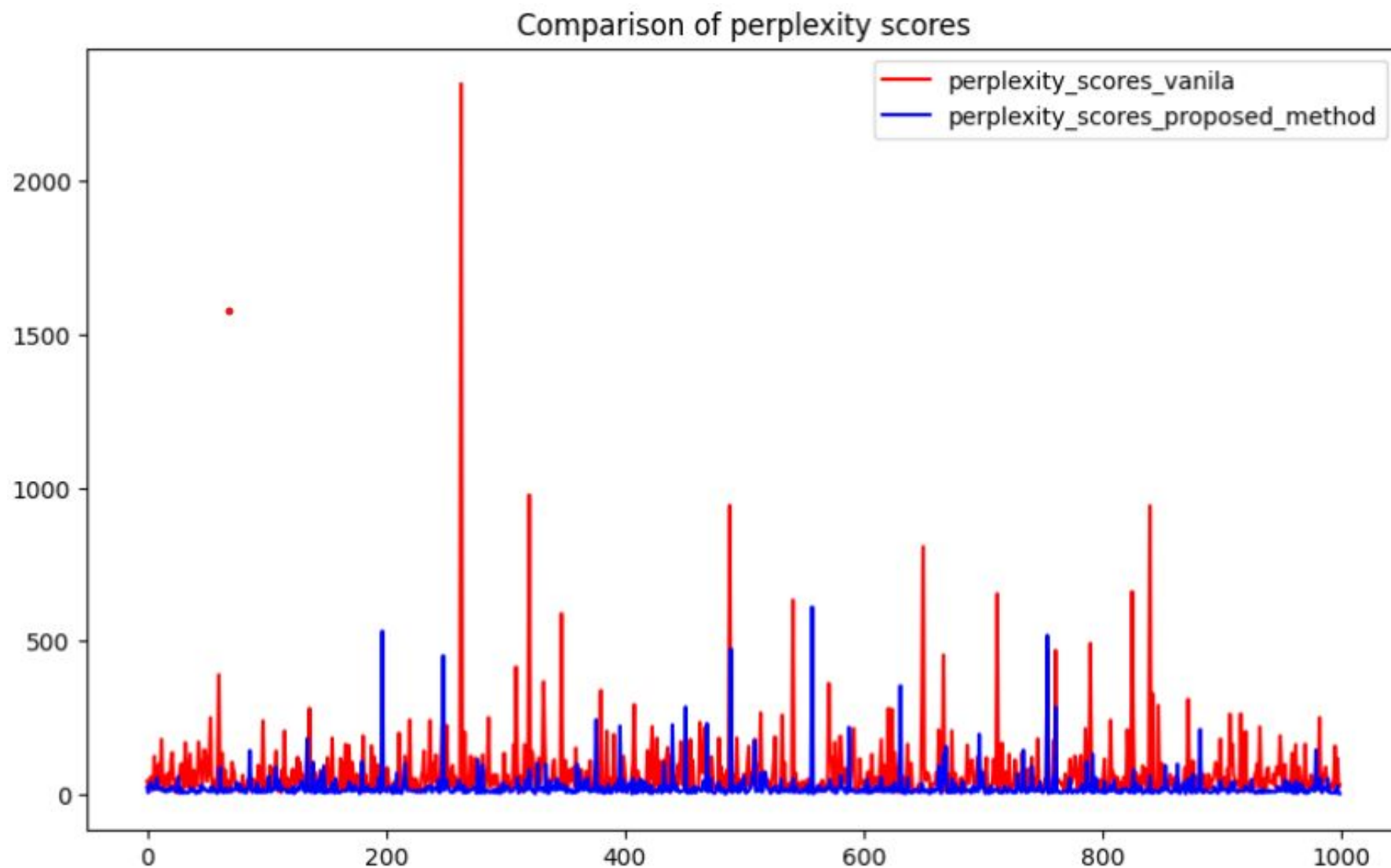
# Evaluation Results - BLEU Score Comparison on LLaMa2



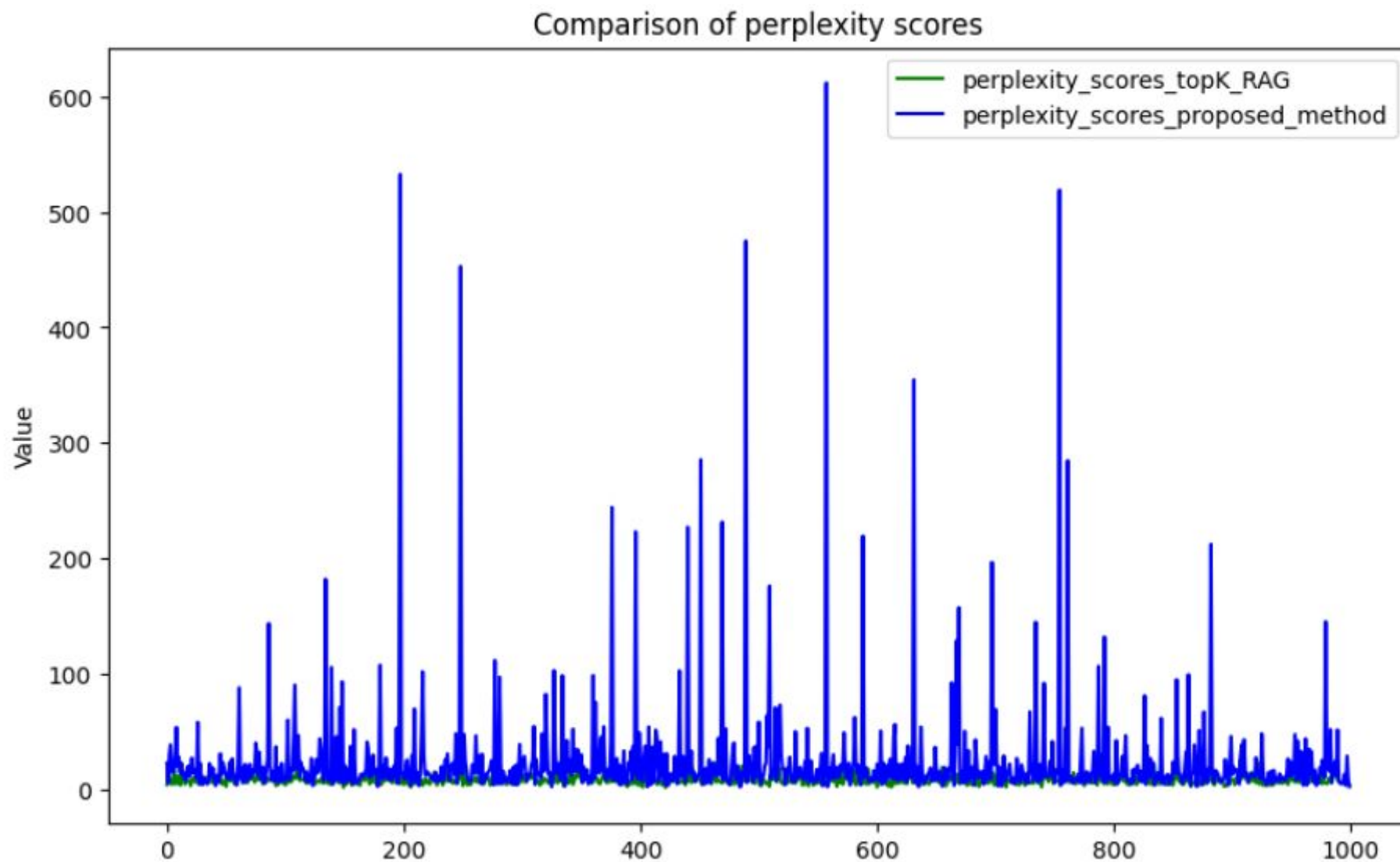
# Evaluation Results - BLEU Score Comparison on LLaMa2



# Evaluation Results - Perplexity Score Comparison on LLaMa2



# Evaluation Results - Perplexity Score Comparison on LLaMa2



# Difficulties and Limitations

- (1) We don't have enough GPU to run the inference using full document as context
- (2) It takes long time to generate the Wikipedia embeddings due to large size of wikipedia documents
- (3) The context size of GPT2 is limited and we can not fit the entire top document into the prompt of GPT2; LLaMa model has long context windows, but it takes GPUs to run experiments and we do not have enough resource to run large scale experiments compared to industry
- (4) On the evaluation of truthfulness, usefulness and trustworthiness, we lack the man power to evaluate the all the outputs (ROUGE, BLUE and Perplexity do not guarantee truthfulness, usefulness and trustworthiness)
- (5) The NQdataset evaluation might still not be comprehensive to encompass a wide range of human knowledge

# Summary of Contribution

(1) Proposed method on better extracting relevant information for RAG system

(2) Extensive evaluation on proposed method on open source model including GPT2 and LLaMa2; And the proposed method achieve better performance in ROUGE, BLUE and Perplexity compared to two baselines

# Future Work

- Evaluation reveals that while LLMs exhibit a certain degree of noise robustness
- Still struggle significantly in terms of negative rejection, information integration, and dealing with false information

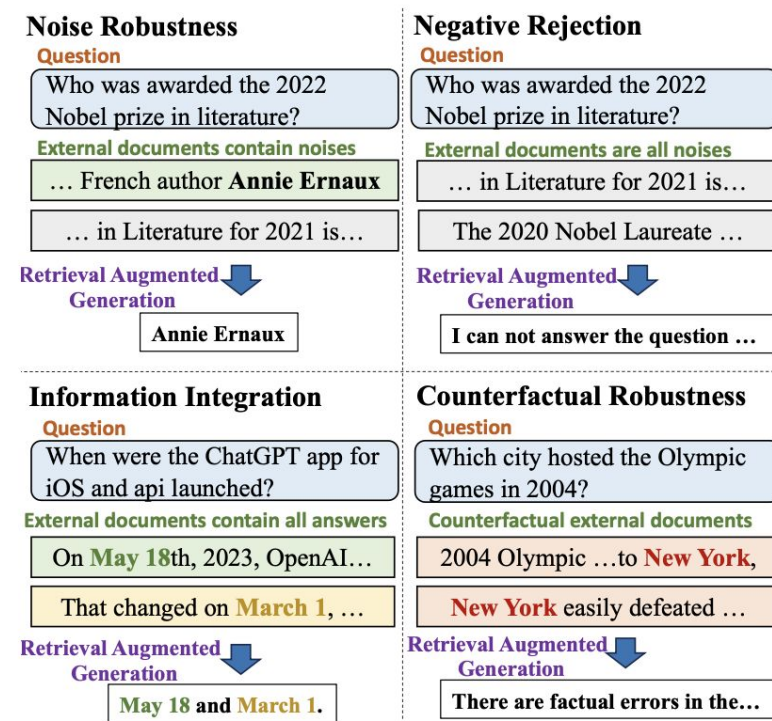


Figure 1: Illustration of 4 kinds of abilities required for retrieval-augmented generation of LLMs.



# Future Work

- Proposes iterative retrieval-generation collaborative framework
- Able to leverage both parametric and non-parametric knowledge, but also helps to find the correct reasoning path through retrieval-generation interactions

