

RA-MTR: A Retrieval Augmented Multi-Task Reader based Approach for Inspirational Quote Extraction from Long Documents

Sayantana Adak and Animesh Mukherjee

Indian Institute of Technology Kharagpur, West Bengal – 721302

sayantanadak.skni@kgpian.iitkgp.ac.in, animeshm@cse.iitkgp.ac.in

Abstract

Inspirational quotes from famous individuals are often used to convey thoughts in news articles, essays, and everyday conversations. In this paper, we propose a novel context-based quote extraction system that aims to extract the most relevant quote from a long text. We formulate this quote extraction as an open domain question answering problem first by employing a vector-store based retriever and then applying a multi-task reader. We curate three context-based quote extraction datasets and introduce a novel multi-task framework RA-MTR that improves the state-of-the-art performance, achieving a maximum improvement of 5.08% in BoW F1-score.¹

1 Introduction

Inspirational quotes from famous individuals are powerful tools that convey wisdom and insight in a concise and often figurative manner. They provide a secondary voice that reinforces the author’s thoughts and beliefs (Liu et al., 2019). Context-aware quote extraction (also known as quote recommendation) is crucial in writing news articles, blogs, and summaries, as it helps to strengthen the expressed ideas. This process involves identifying phrases or sentences within a paragraph that are quotable and determining their relevance and quotability in a given context. Since “context” can be highly subjective, finding the most relevant quotes can be challenging due to the linguistic nuances involved. Figure 1 demonstrates a recommendation for a quotable phrase from a source paragraph, based on one context from the example of our dataset. It turns out that authors have to spend far too much time deciding what-to-quote from many source texts analyzing their context. Accordingly, it is in significant demand to

¹Code and Data are available at https://github.com/sayantana11995/Context_based_Quote_Extraction

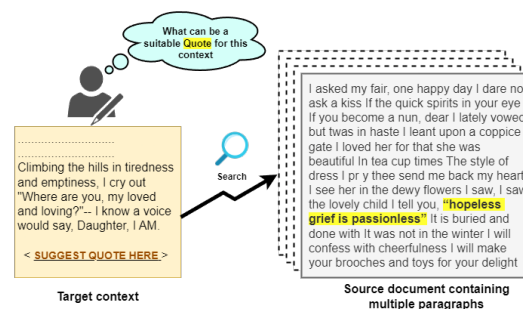


Figure 1: Example use-case of context-aware quote extraction from source document while composing an article. The highlighted portion from the source document can be a suitable quote for the target context in the left.

automate the process of extracting quotes from a text.

To tackle the task, Bendersky and Smith (2012) attempts to identify “quotable” phrase from books on the basis of linguistic and rhetorical properties. Unlike this, (Tan et al. (2015), Tan et al. (2016, 2018), Qi et al. (2022)), leverage “context” to select the most relevant quote from a list of quotes. (Lee et al. (2016); Wang et al. (2021)) use dialogue history as the *context*. The task of finding the most relevant quote itself remains challenging. Moreover, our task poses inherent difficulty, as we not only attempt to find the most relevant quote for a given context, but also extract the quote from a full source document (containing several hundreds of paragraphs). To the best of our knowledge, only MacLaughlin et al. (2021) attempts to extract context-aware quotes from text documents (US presidential speech transcripts). However, the length of the documents are considerably small (see Table 1 for details) and they only cover the political domain for quote extraction. In addition, none of the experimental dataset apart from Qi et al. (2022) is publicly available.

In this research, we focus on bridging the gap by rigorously curating three datasets for context aware quote extraction task, and presenting a novel

framework that can enhance the task of extracting quotes from a much longer text. Our contributions can be summarized as follows.

- To better extract quotes based on the context, we propose a Retrieval Augmented Multi-task Reader (RA-MTR) framework that utilizes a vector-store for initial retrieval followed by *Llama-3* based re-ranker, and a multi-task framework that leverages two training tasks tailored specifically for the quote extraction scenario.
- We curate two datasets for context-based quote extraction by adapting two existing quote recommendation resources—QuoteR, predominantly featuring literary quotes (Qi et al., 2022), and Quotus, which includes quotes from political speeches (Niculae et al., 2015). Additionally, we introduce a dataset centered on quotes from Mahatma Gandhi. Together, these datasets span diverse genres, and we make them publicly available to foster further research in this domain.
- We conduct rigorous experiments using RA-MTR and show that our framework outperforms the best-performing baseline by a maximum of **5.08%** in terms of BoW F1-score while considering the top-ranked paragraph as the location of the quote.
- We also perform analysis with the *multi-task reader* to demonstrate that our fine-tuned multi-task framework based on SpanBERT \oplus SpanBERT-CRF improves the quote predictions over a series of baselines. Our multi-task framework outperforms the standard BERT-based models by a large margin in a few shot settings. In particular, we see that even with eight data points from the target domain, our model beats BERT and SpanBERT by **14%** and **11%** in BoW F1-score respectively (see Table 5.)

2 Related work

Quotability detection: Bendersky and Smith (2012) developed a quotable phrase extraction process that includes a supervised quotable phrase detection using lexical and syntactic features. Wang et al. (2021) introduced a transformation matrix that directly maps the query representations to quotation representations. MacLaughlin and Smith (2021) utilized BERT-based models for ranking the quotable paragraphs while evaluating on five dif-

ferent datasets. Voskarides et al. (2021) discussed challenges of retrieving news articles in the context of developing event-centric narratives.

Context based quote recommendation: Tan et al. (2015) proposed a learning-to-rank framework for quote recommendation. Tan et al. (2016) proposed a quote recommendation framework by learning the distributed meaning representations for the contexts and the quotes using LSTM. Lee et al. (2016) built a quote recommender system to predict quotes based on Twitter dialogues as context. Qi et al. (2022) built a large and the first publicly available dataset for quote recommendation. MacLaughlin et al. (2021) attempted to simultaneously rank the most quotable paragraphs and predict the most quotable spans from source transcripts modeling quote recommendation as an openQA problem.

The present work: We extend the work of MacLaughlin et al. (2021), by proposing a novel retriever augmented multi-task reader based quote extraction. The framework employs a *vector-store* based paragraph retriever followed by a decoder-only transformer based re-ranker and a novel multi-task based reader containing a sequence tagging module for identifying quotable phrases along with context aware span prediction. We curate three datasets of different genres and evaluate our approach. Our method outperforms all the previous baselines and generalizes better in a cross-domain few-shot setting.

3 Approach

We formalize the problem as an open-QA task, similar to the one described in MacLaughlin et al. (2021). Given a target context (T_C), and a source document (S_D) which consists of several paragraphs ($= \{P_1, P_2, \dots, P_n\}$), we require to first identify the most relevant list of paragraphs depending upon T_C , and then identify the most quotable phrase from the selected paragraphs. We propose the overall quote extraction approach consisting of a *retriever* (detailed in section 5.1) to select the most relevant paragraph followed by a *multi-task reader* (detailed in section 5.2) to extract a quote.

4 Dataset curation

In this section we present the details of the datasets first by discussing the quotes that we consider, followed by construction of *source paragraph* and *target context* for our experiments.

4.1 Training data

QuoteR: We primarily consider the English subset of the QuoteR dataset proposed by Qi et al. (2022), known to be the largest publicly available dataset for the quote recommendation task. The authors of the corresponding paper collected several quotes from the popular WikiQuote² project and search for the occurrences of these quotes in the Project Gutenberg corpus³, the BookCorpus (Zhu et al., 2015), and the OpenWebText corpus (Gokaslan and Cohen, 2019) respectively, and considered the preceding and the following 40 words of a particular quote to its left and right context respectively. After preprocessing, the authors provided a total of **6108** unique quotes and around **127k** contexts for those quotes.

4.2 Test data

Gandhi quotes: Websites such as mkgandhi⁴ has made the Collected Works of Mahatma Gandhi (CWMG) publicly available, which is a huge text corpus consisting of 100 volumes. We collect around a total of 800 Mahatma Gandhi quotes in English from Goodreads⁵ and the *mkgandhi* portal.

Quotus data: The authors in MacLaughlin et al. (2021) utilizes the *Quotus* (Niculae et al., 2015) dataset for their experiments. The dataset consists of two sets of texts – transcripts of US Presidential speeches and press conferences from 2009-2014 (*the source document*), and news articles that report on the speeches (*the target document*). The authors crawled the articles and transcripts from the provided links in the Quotus data, and preprocessed them to gather a significant amount of quote, contexts, and paragraphs. However, they did not make their dataset publicly available. We ourselves re-scraped the links from the source Quotus data.

Curating source paragraph and target context: From these three dataset (i.e., *QuoteR*, *Gandhi*, and *Quotus*) we get a list of quotes. However, to evaluate our *retriever* and *reader*, we require to curate the source paragraph and the target context for each of the quotes. We leverage the *Project Gutenberg* corpus to construct 4,889 quote-context-paragraph triples (containing 1,708 unique quotes) for QuoteR. We use *Gadhipedia*⁶ search engine to

²<https://en.wikiquote.org/>

³<https://www.gutenberg.org/>

⁴<https://mkgandhi.org/>

⁵https://www.goodreads.com/author/quotes/5810891.Mahatma_Gandhi?page=35

⁶<https://www.gandhipedia150.in>

curate 737 triples for the Gandhi quotes. For the Quotus, we utilize the *Quoting POTUS*⁷ website containing the news articles and align the quotes to source transcript for constructing 2,698 triples. The detailed steps and algorithms are provided in Appendix D.

4.3 Dataset statistics

Thus, overall we consider *three* datasets each from a different genre - (i) QuoteR - where most of the quotes are from novels, 2) Gandhi - solely based on the quotes of Mahatma Gandhi, and 3) Quotus - quotes from the political speech. The basic statistics of these three datasets are noted in Table 1. Figure 2 demonstrates the most prominent words present in the three datasets. The quotes in the QuoteR and Gandhi datasets contain positive words like “God”, “good”, “love”, “truth”, “petition” etc. The Quotus dataset on the other hand contains quotes having words “america”, “president” etc. We also compare our dataset with the dataset used in other similar works (see Table 2). We present the only dataset containing quote, context and source paragraph. These datasets will be made publicly available upon acceptance.

Dataset	# of unique quotes	# of quote-context-paragraph triples	Avg. # of tokens / quote	Median # of tokens / quote	Avg # of para / Src	Avg # of tokens / Src
QuoteR	1708	4889	13.51	11	551	98783.17
Gandhi	737	737	20.42	19	19.54	3812.47
Quotus	2698	2698	20.46	16	86.78	4631.55

Table 1: Statistics for the three datasets. For QuoteR we report the instances that we could find in the Gutenberg corpus.

Dataset	Context based	Context type	Source paragraph	Public
Bendersky and Smith (2012)	✗	writings	✗	✗
Tan et al. (2015)	✗	writings	✗	✗
Wang et al. (2021)	✓	dialogue	✗	✓
Qi et al. (2022)	✓	writings	✗	✓
MacLaughlin et al. (2021)	✓	writings	✓	✗
Our dataset	✓	writings	✓	✓

Table 2: Dataset comparison for other related tasks with ours.

5 Methodology

In this section, we describe the details of our methodology for quote extraction. We propose the overall quote extraction approach as an open-QA framework, which normally consists of a *retriever* and a *reader*. The retriever is essential for selecting the top paragraphs relevant to the context from the whole document. We employ a novel multi-task learning framework in the reader, which extracts

⁷<http://snap.stanford.edu/quotus/vis/>



Figure 2: Most prominent words present in the quotes across the three datasets.

the most quotable spans from the selected paragraphs and is discussed in detail below. The overall retriever-reader architecture RA-MTR is illustrated in the left part of Figure 3.

5.1 RA-MTR: Retriever module

Inspired by the RAG (Lewis et al., 2020) architecture, we employ a vector-store based retriever to initially retrieve top- k ($k = 20$)⁸ paragraph based on the given context. We utilize *langchain* API⁹, to split the source document into several chunks (we choose chunk-size of 1200¹⁰ characters and chunk-window of 100), followed by encoding of each chunk using *sentence-transformers*, and finally store the embeddings into a vector-store for efficient searching. We use ChromaDB¹¹ for storing the embeddings of the chunks. In parallel, the query context is also embedded using *sentence-transformers*. To extract the relevant inspirational quote from the source document we perform a similarity search by comparing the query context embedding and the embeddings in the vector-store to retrieve top- k chunks. The retrieved chunks are then passed to the more powerful sequence-to-sequence re-ranking module for further processing. *Fine-tuning paragraph re-ranking module:* After retrieving initial sets of candidate paragraphs, many past literature leveraged deep neural network based paragraph re-ranking modules to get a final ranked list of paragraphs. Works such as Dai and Callan (2019); Yilmaz et al. (2019); Nogueira et al. (2019) exploited BERT for paragraph/document retrieval tasks. Nogueira et al. (2020) used a T5-based encoder-decoder architecture for document ranking. We apply a more sophisticated decoder-only transformer based model *Llama-3*¹² to re-rank the

paragraph. Similar to Nogueira et al. (2020) we formulate the problem as a binary classification task, and the input prompt is:

Context: {c}
 Document: {d}
 Is the document relevant to the context? Answer yes/no:

where c and d are the context and paragraph texts, respectively. The model is fine-tuned to produce the words yes or no depending on whether the document is relevant or not to the query. That is, yes and no are the ‘target words’ (i.e., ground truth predictions in the sequence-to-sequence transformation). To generate training and test examples for the models, we iterate over each context and create (context, source paragraph, label) example triples for each paragraph in the corresponding source document. The label is yes if the author actually quoted from the paragraph (positive triple) and no (negative triple) otherwise. At inference time, to compute probabilities for each query–document pair (in a re-ranking setting), we retrieve the unprocessed next-token probabilities for the tokens yes and no. From these, we calculate the *yes – score* as follows.

$$yes - score_{(c,d_i)} = \frac{p(yes|P)}{p(yes|P) + p(no|P)} \quad (1)$$

where, c is the context, d_i is the i^{th} document and P is the prompt. Similarly, as baseline, we also fine-tune encoder-decoder based models (T5, FLAN-T5) for the re-ranking task using similar approach. *Sampling hard negatives:* Out of all the negative triples obtained we select n hard samples for training to make the model more robust. We explore two different hard negative sampling methods - a) select the paragraphs that are closest next to the positive paragraph, and b) select the top ranked paragraphs (other than the positive one) returned by BM25 retriever model. However, we observe that both choices produce similar results (see results section for details).

⁸Increasing k did not change the performance too much.

⁹https://python.langchain.com/docs/modules/data_connection/

¹⁰We find maximum length of context + paragraph is 1005.

¹¹<https://pypi.org/project/chromadb/>

¹²We apply the chat model from huggingface *meta-llama/Meta-Llama-3-8B-Instruct*

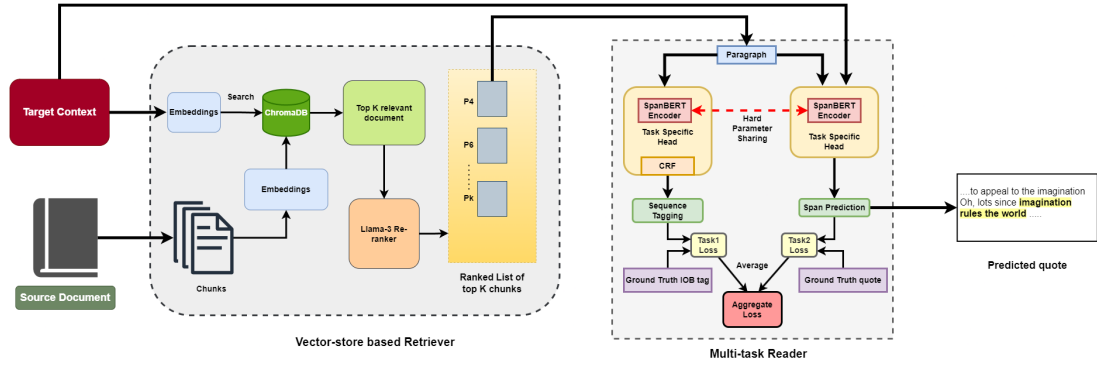


Figure 3: The RA-MTR architecture.

5.2 RA-MTR: Multi-task reader module

Motivation for multi-task training: Unlike normal spans of text, quotes have certain inherent special properties or some figurative language that make them unique (Bendersky and Smith, 2012). We believe that identifying such special occurrences of phrases is essential for quote prediction from paragraphs. We cast this as a sequence tagging, i.e., marking only the portions of a text that can be recounted as quotable. We attempt to optimize two tasks in parallel - quotable sequence tagging (using SpanBERT-CRF) and context aware span prediction (using SpanBERT). In a paragraph, there can be multiple spans of text which will be relevant to the context. However, not every relevant span is quotable. The span prediction module is essentially a variant of a question-answering module, which might not be good enough to identify quotability of the answer. Also, many of the quotes are only subparts of a sentence (e.g., “He travels fastest who travels alone,...”) while few of the quotes consist of more than one sentence (e.g., “In this world there are only two tragedies. One is not getting what one wants, and the other is getting it.”). To mitigate this gap, we use a specific sequence identification module (SpanBERT-CRF) to find quotable sequences.

Span prediction from paragraph: We train the span prediction model using context-quote-paragraph triple as the training example. Similar to MacLaughlin et al. (2021), we utilize the span-level BERT architecture, which receives the packed sequence of the context and paragraph as input. By utilizing the final hidden vector $T_i \in \mathbb{R}^h$ as the representation for each wordpiece in a given paragraph, we follow the standard approach of casting span prediction as two classification tasks, i.e., separately predicting the start and end of the span

(Devlin et al., 2019). To this purpose, we introduce a start vector, $S \in \mathbb{R}^h$, and an end vector, $E \in \mathbb{R}^h$. The probability of a word w being the start of the quoted span is the dot product $S \cdot T_w$ followed by a softmax over all wordpieces in the example. We follow the same approach for calculating the probability of word w being the end of the span using E . The loss is calculated as the average NLL (Negative log-likelihoods) of the correct start and end positions, i.e., the tokens in the paragraph the author actually quoted. Following Devlin et al. (2019), at prediction time, the score of the span from position i to j is $S \cdot T_i + E \cdot T_j$. We consider a span valid if $j > i$ and i and j occur in the paragraph portion of the input. We retain a mapping from wordpieces to original tokens for prediction.

Quotability detection as sequence tagging: Scheible et al. (2016); Pareti et al. (2013) framed the quotation detection task as sequence tagging. Portelli et al. (2021) used sequence labeling for the adverse drug events (ADE) detection from a given text. Along similar lines, we employ SpanBERT neural model combined with Conditional Random Field (CRF) to identify quotable phrases or words. Each example from the dataset is accompanied by a paragraph, the start and end character positions of the quote in that paragraph. Using this information, we first convert this into the commonly used IOB (Inside, Outside, Beginning) schema using *Spacy*. Consider the example in Figure 1, every word except the bold portion (i.e., the quote) should be marked as ‘O’. The word ‘hopeless’ in the quote should be labeled as ‘B’ and the rest of the quote should be labeled as ‘I’. The BIO tagging is illustrated in Figure 4. Since BERT-based models generally employ wordpiece tokenizers to tackle the out-of-vocabulary words, which actually break such words into multiple subwords, we require to

```

... . you, "hopeless grief is passionless" It is buried . . .
.. 'O' 'O' 'O' 'B' 'I' 'I' 'I' 'O' 'O' 'O' 'O' ..

```

Figure 4: Example of BIO tagging.

Method	Dataset					
	QuoteR (test)		Gandhi		Quotus	
	Top-1	Top-3	Top-1	Top-3	Top-1	Top-3
DrQA (Chen et al., 2017)	7.19	8.22	6.20	8.38	4.26	5.43
ParagraphRanker (Lee et al., 2018)	16.58	21.45	12.17	14.35	11.31	15.11
BM25 + (MacLaughlin et al., 2021) (Positive only settings)	31.78	37.37	23.70	26.60	32.58	37.68
Contriever + FiD	32.81	39.22	23.78	28.25	33.36	39.43
BM25 + BERT-base + MTR*	37.20	46.28	25.17	27.30	36.15	41.12
BM25 + T5-base + MTR*	34.17	45.21	21.31	23.45	37.13	39.25
BM25 + T5-large + MTR*	39.07	48.29	28.13	30.67	39.29	41.11
BM25 + FLAN-T5-large + MTR*	43.12	51.43	34.46	42.30	42.26	48.21
Vector-store based retriever + LLM reader (Llama-3-8b-instruct)	14.81	18.76	17.53	23.28	31.38	39.77
RA-MTR (ours)	45.74	57.25	38.71	49.38	43.40	53.45

Table 3: The result (F1 score) for the quote extraction using the different baseline models and our RA-MTR approach. For a fair comparison, we took the results from the positive-only settings of MacLaughlin et al. (2021). Note that all the fine-tuned models are only trained on the QuoteR training data. *MTR: Our fine-tuned multi-task reader. Results using different LLMs are reported in Table 10 in Appendix B.

decide on a consistent IOB schema for the subwords. We set a rule for the sub-labels which are consistent with the IOB schema: words labeled as B generate a series of subwords labeled as [B, I, . . . , I], while words labeled as I (or O) generate a series of identical I (or O) sub-labels. For example, if the word ‘resisted’ has the label B, then its corresponding wordpieces - [‘Resist’, ‘##ed’] would get labeled as [B, I].

Multi-task training: To take advantage of both the span prediction model and the quotable phrase identification model, we adopt a multi-task based framework where we have two independent models and they share the same transformer encoder. The span prediction model tries to match the start and end token of the quote in the paragraph, whereas the quotable phrase identification model tries to predict the ‘B’, ‘I’, and ‘O’ labels for each token. During the backpropagation, we average the losses from the two models. The right part in the Figure 3 demonstrates the architecture of the multi-task framework.

6 Experiments

In this section, we discuss the experiments that we conduct and the details of the experimental setups.

Fine-tuning paragraph re-ranking: We pass the packed input of context and paragraphs to the retriever model. Out of **4889** data points in the QuoteR dataset, and we select 80% for training, 10% each for dev and test. We choose to fine-tune

the *Llama-3-8b-instruct* model for the paragraph ranking task. For model implementation details and hyper-parameters see Appendix F.

Fine-tuning reader: We fine-tune the reader by randomly selecting 80% QuoteR data for training, 10% each for dev and test (see Appendix F for implementation details). To test the generalizability of the model in a few-shot setting, we consider random training samples $\in \{4, 8, 16, 32, 64\}$ from the other two datasets (i.e., Gandhi and Quotus) for further fine-tuning with a slightly lower learning rate ($1e^{-5}$), and test on the remaining data samples for the respective datasets.

Metrics: Since the setup for our span prediction task is identical to QA, we evaluate the span-level models using the two popular QA metrics – (i) exact match (EM), and (ii) macro-averaged bag-of-words (BoW) F1. EM measures if the predicted span exactly matches the positive quote, and BoW-F1 measures their average word overlap.

Baselines: Both the retriever and the reader can have many variants which serve as ideal baselines. In the retriever part we use vanilla BM25 as a first baseline. Apart from the simple BM25 retriever, we employ encoder-based (BERT), encoder-decoder based (T5, FLAN-T5) document re-ranking to improve paragraph selection.

For the *reader* part, as primitive baselines, we consider using the first and last sentences of each paragraph as potential quotes. To further explore, we also fine-tune the BERT and the SpanBERT pretrained models on the BERT question answering architecture. We keep the same hyperparameter settings as the multi-task framework. Additionally, we also observe the ability of different medium sized open-source LLMs such as FLAN-T5-large¹³, FLAN-T5-XL¹⁴, Bloomz-3b¹⁵, Falcon-7b¹⁶, Llama-3-8b¹⁷ models for zero-shot context-aware quote extraction task. For the detailed methodology, refer to Appendix E.

7 Results

Performance of RA-MTR: To examine the efficacy of our entire pipeline, we conduct an end-to-end prediction from our approach. In the *retriever-reader* based (baseline) approach, we first provide

¹³<https://huggingface.co/google/flan-t5-large>

¹⁴<https://huggingface.co/google/flan-t5-xl>

¹⁵<https://huggingface.co/bigscience/bloomz-3b>

¹⁶<https://huggingface.co/tiiuae/falcon-7b>

¹⁷<https://huggingface.co/meta-llama/>

Meta-Llama-3-8B-Instruct

Method	EM	BoW-F1
First sentence	0.55	6.31
Last sentence	1.08	5.88
BERT-base	69.1±0.5	76.2±0.9
BERT-large	71±0.3	77.9±0.3
SpanBERT-base	71.7±0.6	77.7±0.4
SpanBERT-large	72.3±0.6	79.2±0.4
Multi-task using SpanBERT-base (Ours)	73±0.8	78.2 ±0.3
Multi-task using SpanBERT-large (Ours)	77±0.4	86.1 ±0.2

Table 4: Reader performance on the QuoteR dataset. We provide the positive paragraph to predict the quote span.

Test data	# training samples	Methods		
		BERT	SpanBERT	Multi-task (Ours)
Gandhi	8	27.71	30.30	41.32
	16	32.60	32.65	50.29
	32	38.12	36.85	62.91
	64	43.54	44.65	72.08
Quotus	8	33.97	36.82	40.58
	16	37.90	41.84	49.33
	32	40.56	44.80	55.08
	64	47.86	51.27	59.12

Table 5: Few-shot inference performance on the 1) Gandhi and 2) Quotus datasets. We have used the BoW F1-score as the metric for comparison here.

the context and the list of paragraphs segmented from a particular book to the paragraph retrieval module. We initially get a list of 20 top-ranked paragraphs relevant to the context from the RAG model and then re-rank these using the Llama-3 model. We take the top three paragraphs further and sequentially pass them with respect to the context to our multi-task quote extraction module. This span prediction module within the multi-task framework predicts the top three quotable spans, each from one corresponding paragraph. We measure the BoW F1-score for these three predictions with respect to the ground truth quote and report the scores for 1) top-1 prediction - score when we compare the ground truth with the predicted span from the top 1 out of the three ranked paragraphs, and 2) top-3 predictions - best score when we compare the ground truth with all the three predictions. Table 3 demonstrates the result for the end-to-end quote prediction RA-MTR. We compare the performance of our pipeline with two commonly used baselines for open-domain question answering tasks – DrQA (Chen et al., 2017) and ParagraphRanker (Lee et al., 2018). In addition, we employ Contriever (Izacard et al., 2021) for paragraph retrieval and fine-tune a Fusion in Decoder (FiD) (Izacard and Grave, 2021) model as the reader. We also compare RA-MTR against MacLaughlin et al. (2021), which focuses on context-based quote extraction. Lastly, we compare our model with present day LLM variants. RA-MTR by far outperforms all the baselines.

Multi-task reader performance: We show the

results for span prediction using various methods in Table 4 for the QuoteR dataset. The results in the first two rows are from two very primitive baselines. Scores in the next two rows are from only the span selection models, which (MacLaughlin et al., 2021) has considered. We can clearly see that our multi-task based approach outperforms the other methods. The improvements are significant with $p < 0.05$ as per the Mann-Whitney U test (Mann and Whitney, 1947) and experiment with our three datasets. In Table 5, we demonstrate the few-shot performances on the Gandhi and the Quotus data for the quote prediction task. We can observe that, in the few-shot settings, the multi-task framework performs much better than the simple span prediction models that are normally used for the QA tasks. In fact, with only 8 data samples from the target domain, our model beats BERT and SpanBERT by **14%** and **11%** for the Gandhi data, and by **7%** and **4%** for the Quotus data respectively. We can infer from these results that the addition of the quotable phrase identification task actually helps the model learn the linguistic properties of the quotes much better than the simple span prediction model. Further, the multi-task framework generalizes particularly well in the cross-domain setting even with the training and test paragraphs coming from different genres.

Method	EM	BoW-F1
BERT span prediction	19.20±0.30	30.90±0.80
SpanBERT span prediction	18.30±0.50	29.70±0.40
Multi-task (Ours)	22.00±0.80	38.20±0.70

Table 6: Results for the quote extraction in absence of the context (for QuoteR dataset).

Performance of the sequence tagger: We analyze the output generated by the sequence tagger head from our multi-task framework. Note that this was an auxiliary task to improve the main task of span prediction. The sequence tagger head typically predicts ‘B’, ‘I’, or ‘O’ tags for each token, and the prediction is independent of the context. We apply the model to instances in the QuoteR test dataset and extract the predicted labels from the sequence tagger head. We find that the model correctly predicts the BIO labels for 48.1% of the instances. In 20.7% of the cases, the model predicts multiple BIO labels within a single paragraph, indicating that one paragraph can contain multiple instances of quotable phrases.

Context (in)dependence: We conduct an ablation experiment to observe the impact of context for the quote prediction in the multi-task setting. We re-

move the whole context from the inputs in the test set for the quote prediction models while experimenting with the QuoteR dataset¹⁸. Table 6 clearly shows that the baseline models’ performances are drastically reduced, whereas our multi-task framework outperforms the two baselines. As the sequence tagging task is independent of the context we observe that the multi-task framework performs better in the absence of context while the two other models that are highly context-sensitive. We can infer that the linguistic boundary identification for the quotes in terms of the BIO markers enhances the performance and makes it robust to the absence of context. This is one of the prime strengths of the multi-task framework.

8 Analysis of retriever

8.1 Quantitative analysis of retriever

We attempt to measure how much re-ranking is effective in ranking the paragraphs based on the context. We use the curated set for paragraph ranking (the QuoteR test set, and other two datasets including the source) to perform the evaluation. As baselines we use BM25 based, sentence-bert similarity based and contriever based retriever. Further we use Flan-T5 and Llama-3 based re-ranking to evaluate the importance of re-ranking. From the retriever we first take top-20 (which we use during the main experiment) retrieved paragraphs from the source and then use re-ranking to measure Precision@k ($k \in \{1, 3\}$). The result is reported in Table 7. We can observe that the re-ranking drastically improves the retriever performance.

Method	Dataset					
	QuoteR		Gandhi		Quotus	
	P@1	P@3	P@1	P@3	P@1	P@3
BM25	0.15	0.18	0.09	0.16	0.13	0.19
sentence-bert similarity	0.24	0.29	0.14	0.17	0.23	0.27
contriever-based	0.27	0.33	0.15	0.19	0.25	0.3
sentence-bert + Flan-T5 re-ranking	0.43	0.52	0.38	0.42	0.44	0.53
sentence-bert + Llama-3 re-ranking	0.48	0.57	0.41	0.44	0.47	0.54

Table 7: Analysis of retriever performance along with different reranker.

8.2 Qualitative analysis of different re-ranking methods

In this section, we attempt to analyse how well our vector-store based retriever performed. As we cannot directly compare the retrieved chunks with

¹⁸The results from the other datasets show similar trends and hence are not shown.

the positive paragraph in our dataset (due to variable word length), we measure the average *Jaccard* similarity between the top predicted chunk with the positive paragraph in our dataset for a specific context. We present the results in Figure 5. We observe that, using Llama-3 based re-ranking, the similarity significantly improved for all the three datasets.

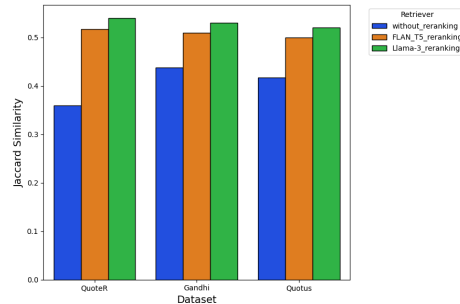


Figure 5: Average *Jaccard* similarity between top predicted chunk and positive paragraph for a specific context.

9 Analysis of the multi-task reader output

Analysis of top predicted quotes: Since there may be multiple quotes in a given paragraph for a given context, we also look at the top five predicted spans from our multi-task framework for each of the paragraphs in the test set. We manually annotate the relevance of the predicted spans for the top five predictions. We had two annotators, and each of them was provided with a set of context and the top five predicted spans. They were required to mark 1 if the predicted span is semantically coherent with the context, and 0 otherwise. In the case of ambiguity, a third annotator was involved to adjudicate. We obtain an inter-annotator agreement of Cohen’s $\kappa = 0.64$. We take the final relevancy (i.e., 0 or 1) based on majority vote. We achieve a high MAP@5 score of 0.78, indicating that our multi-task framework retrieved ~ 3.9 (on average) meaningful recommendations among the top five recommendations.

Error in the sequence tagger: We review some instances where the model failed to predict the correct BIO labels. A specific example is depicted in Figure 6, where the true quote is highlighted in green, while the predicted quotes are highlighted in yellow. Although the true and the predicted quotes come from different portions of the paragraph, they both are highly quotable as per human experts. We observe many such cases of (pseudo) errors that manifest due to the absence of valid ad-

ditional ground truth quotes.

Error in the multi-task reader: We examine the

inward representation, and a creative energy constantly fed by susceptibility to the veriest minutiae of experience, which it reproduces and constructs in fresh and fresh wholes not the habitual confusion of probable fact with the fictions of fancy and transient inclination, but a breadth of ideal association which informs every material object, every incidental fact, with far reaching memories and stored residues of passion, bringing into new light the less obvious relations of human existence. imagination, feeling and the whole inward life are being constantly shaped by our actions. **experience gives new character to the inward life, and at the same time determines its motives** and its inclinations. The muscles develop as they are used what has been once done it is easier to do again. in the same way, **our deeds influence our lives, and compel us to repeat our actions**. at least this is george elliot s opinion, and one she is fond of re affirming. after arthur had wronged hetty, his life was changed, and of this change wrought in his character by his conduct, george elliot says, **our deeds determine us, as much as we determine our deeds** and until we know what

Figure 6: Example of an instance where the sequence tagger wrongly predicts the BIO labels. True and predicted quotes are highlighted in green and yellow respectively.

predicted quotes, which do not entirely match with the ground truth quotes. We observe that in 72% the model predicts a sub-phrase of the original quote. For instance, while the actual quote is ‘Our Father, which art in heaven, hallowed be thy name’, the corresponding predicted quote is ‘which art in heaven, Hallowed be thy Name’. In a few cases, the model over-predicts, i.e., predicts a span containing the true quote and some phrases surrounding it. For example, the actual quote ‘Money begets money’, is predicted as ‘Money begets money and its offspring’.

10 Additional results

Ablation experiments: We conduct ablation experiments for the effectiveness of different components in the RA-MTR framework. We report the results for the QuoteR test set in Table 8 (see Appendix 8 for more results).

Vector space analysis: Apparently the two tasks

Method	Top-1 F1	Top-3 F1
RA-MTR	45.74	57.25
w/o paragraph ranking	36.33	44.84
w/o CRF in the sequence tagger	41.26	49.57

Table 8: Results of RA-MTR without auxiliary components for the QuoteR test set.

presented inside the multi-task framework may seem similar. To understand how the two tasks are different we conduct an experiment by passing 400 test examples (i.e., context-paragraph-quote triples) from QuoteR data as input to the trained multi-task framework and extracting the last hidden representations of the two task specific heads. We measure the average cosine similarity between these two representations and observe a very small similarity of 0.23. We also plot the cosine simi-

larity heatmap and T-SNE plot between these two representations (see Figure 7). From the T-SNE plot we observe a negligible overlap, which indicates the tasks are indeed not the same, rather the complementary nature of the two tasks assist each other resulting in improvement of both ‘quotability detection and recommendation’.

Additional user evaluation: We already presented

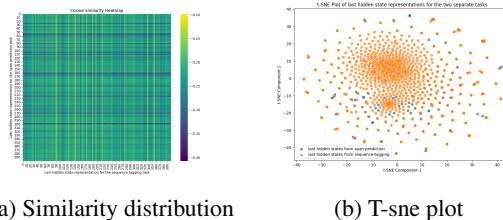


Figure 7: Vector space analysis of the two tasks presented in multi-task reader.

a set of human evaluation in section 9 to analyze the top predicted quotes from MTR. Here we randomly select 50 examples each from the Gandhi and Quotus data combined, and collect human judgements from two undergraduate students about the ‘context relevance’ and ‘quotability’. Given the context and predicted quote, we ask them to select ‘Yes/No’ for the ‘context relevance’ and ‘quotability’. Overall, the two annotators respectively mark ‘Yes’ in 46, 44 cases for ‘context relevance’, and 47, 41 cases for ‘quotability’.

11 Conclusion

In this work, we proposed a method to recommend quotes from large texts given a context. We employed a novel multi-task framework for quote prediction, which can in parallel predict the span of text and identify the quotable phrases. We constructed three datasets of different genres and experimented on them. We believe that our methodology and datasets will be beneficial for future research.

Acknowledgments

We sincerely thank Saketh Konda and Agnik Saha for their invaluable contributions to this work. Their assistance in conducting qualitative analyses and running baseline experiments played an integral role in shaping this research. We also extend our gratitude to the anonymous reviewers for their thoughtful feedback and constructive suggestions, which significantly enhanced the quality of this paper.

Limitations

In this section we will discuss the limitations of our study. While it is evident that the quotes are available in different regional languages, all of our experiments are conducted for the English version of the datasets. Few of the pre-processing steps might not be suitable for the languages with different morphosyntactic structures. Further the base models will also need to be changed.

Ethics statement

We used three datasets for our experiments. The QuoteR dataset was released publicly by the authors of (Qi et al., 2022). Besides, we extracted all the paragraphs from open corpora, including free public domain e-books. The quotes of Gandhi were collected from the free quote repository and the context were extracted from the publicly available portal. Both the quotes and the contexts for the Quotus data were collected from the open corpora. The annotators voluntarily annotated the predictions for our analysis, and we did not retain any of their private information.

References

- Sayantan Adak, Atharva Vyas, Animesh Mukherjee, Heer Ambavi, Pritam Kadasi, Mayank Singh, and Shivam Patel. 2020. [Gandhipedia: A one-stop ai-enabled portal for browsing gandhian literature, life-events and his social network](#). In *Proceedings of the ACM/IEEE Joint Conference on Digital Libraries in 2020, JCDL '20*, page 539–540, New York, NY, USA. Association for Computing Machinery.
- Michael Bendersky and David Smith. 2012. [A dictionary of wisdom and wit: Learning to extract quotable phrases](#). In *Proceedings of the NAACL-HLT 2012 Workshop on Computational Linguistics for Literature*, pages 69–77, Montréal, Canada. Association for Computational Linguistics.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. [Reading Wikipedia to answer open-domain questions](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1870–1879, Vancouver, Canada. Association for Computational Linguistics.
- Zhuyun Dai and Jamie Callan. 2019. [Deeper text understanding for ir with contextual neural language modeling](#). In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR'19*, page 985–988, New York, NY, USA. Association for Computing Machinery.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Aaron Gokaslan and Vanya Cohen. 2019. Open-webtext corpus. <http://Skylion007.github.io/OpenWebTextCorpus>.
- Miguel Grinberg. 2018. *Flask web development: developing web applications with python*. " O'Reilly Media, Inc."
- Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2021. Unsupervised dense information retrieval with contrastive learning. *arXiv preprint arXiv:2112.09118*.
- Gautier Izacard and Édouard Grave. 2021. Leveraging passage retrieval with generative models for open domain question answering. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 874–880.
- Hanbit Lee, Yeonchan Ahn, Haejun Lee, Seungdo Ha, and Sang-goo Lee. 2016. [Quote recommendation in dialogue using deep neural network](#). In *Proceedings of the 39th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '16*, page 957–960, New York, NY, USA. Association for Computing Machinery.
- Jinhyuk Lee, Seongjun Yun, Hyunjae Kim, Miyoung Ko, and Jaewoo Kang. 2018. [Ranking paragraphs for improving answer recall in open-domain question answering](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 565–569, Brussels, Belgium. Association for Computational Linguistics.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. [Retrieval-augmented generation for knowledge-intensive nlp tasks](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 9459–9474. Curran Associates, Inc.

- Yuanchao Liu, Bo Pang, and Bingquan Liu. 2019. [Neural-based Chinese idiom recommendation for enhancing elegance in essay writing](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5522–5526, Florence, Italy. Association for Computational Linguistics.
- Ansel MacLaughlin, Tao Chen, Burcu Karagol Ayan, and Dan Roth. 2021. [Context-based quotation recommendation](#). *Proceedings of the International AAAI Conference on Web and Social Media*, 15(1):397–408.
- Ansel MacLaughlin and David Smith. 2021. [Content-based models of quotation](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2296–2314, Online. Association for Computational Linguistics.
- Henry B Mann and Donald R Whitney. 1947. On a test of whether one of two random variables is stochastically larger than the other. *The annals of mathematical statistics*, pages 50–60.
- Vlad Niculae, Caroline Suen, Justine Zhang, Cristian Danescu-Niculescu-Mizil, and Jure Leskovec. 2015. Quotus: The structure of political media coverage as revealed by quoting patterns. In *Proceedings of the 24th International Conference on World Wide Web*, pages 798–808.
- Rodrigo Nogueira, Zhiying Jiang, Ronak Pradeep, and Jimmy Lin. 2020. [Document ranking with a pre-trained sequence-to-sequence model](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 708–718, Online. Association for Computational Linguistics.
- Rodrigo Nogueira, Wei Yang, Kyunghyun Cho, and Jimmy Lin. 2019. Multi-stage document ranking with bert. *arXiv preprint arXiv:1910.14424*.
- Silvia Pareti, Tim O’Keefe, Ioannis Konstas, James R. Curran, and Irena Koprinska. 2013. [Automatically detecting and attributing indirect quotations](#). In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 989–999, Seattle, Washington, USA. Association for Computational Linguistics.
- Beatrice Portelli, Daniele Passabì, Edoardo Lenzi, Giuseppe Serra, Enrico Santus, and Emmanuele Chersoni. 2021. Improving adverse drug event extraction with spanbert on different text typologies. In *International Workshop on Health Intelligence*, pages 87–99. Springer.
- Fanchao Qi, Yanhui Yang, Jing Yi, Zhili Cheng, Zhiyuan Liu, and Maosong Sun. 2022. Quoter: A benchmark of quote recommendation for writing. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 336–348.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21(1).
- Christian Scheible, Roman Klinger, and Sebastian Padó. 2016. [Model architectures for quotation detection](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1736–1745, Berlin, Germany. Association for Computational Linguistics.
- Jiwei Tan, Xiaojun Wan, Hui Liu, and Jianguo Xiao. 2018. Quoterec: Toward quote recommendation for writing. *ACM Transactions on Information Systems (TOIS)*, 36(3):1–36.
- Jiwei Tan, Xiaojun Wan, and Jianguo Xiao. 2015. [Learning to recommend quotes for writing](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 29(1).
- Jiwei Tan, Xiaojun Wan, and Jianguo Xiao. 2016. [A neural network approach to quote recommendation in writings](#). In *Proceedings of the 25th ACM International Conference on Information and Knowledge Management, CIKM ’16*, page 65–74, New York, NY, USA. Association for Computing Machinery.
- Nikos Voskarides, Edgar Meij, Sabrina Sauer, and Maarten de Rijke. 2021. News article retrieval in context for event-centric narrative creation. In *Proceedings of the 2021 ACM SIGIR International Conference on Theory of Information Retrieval*, pages 103–112.
- Lingzhi Wang, Xingshan Zeng, and Kam-Fai Wong. 2021. Quotation recommendation and interpretation based on transformation from queries to quotations. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 754–758.
- Zhiguo Wang, Patrick Ng, Xiaofei Ma, Ramesh Nalapat, and Bing Xiang. 2019. Multi-passage bert: A globally normalized bert model for open-domain 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)*, pages 5878–5882.
- Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. 2016. Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*.
- Zeynep Akkalyoncu Yilmaz, Shengjin Wang, Wei Yang, Haotian Zhang, and Jimmy Lin. 2019. Applying bert to document retrieval with birch. 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): System Demonstrations, pages 19–24.

Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *Proceedings of the IEEE international conference on computer vision*, pages 19–27.

A Ablation for the Multi-task reader

In this section we ablate the sequence tagger module of multi-task reader to observe the performance variations of the reader. We use a simple SpanBERT module in the multi-task framework to train on the QuoteR training set. Then we use the QuoteR test set for running the inference. The result is shown in the Table 9. We observe that the SpanBERT (span prediction) along with SpanBERT-CRF (for sequence tagging) is working the best for the multi-task framework.

Approach	EM	BoW F1
SpanBERT (span prediction)		
⊕ SpanBERT-CRF (sequence tagging)	77±0.4	86±0.2
⊕ SpanBERT (sequence tagging)	74.1±0.3	82.3±0.4

Table 9: Ablation for the Multi-task reader. ⊕ denotes multi-task framework.

B Results using different LLMs as reader

Extending the Table 3, we demonstrate the results while using different other LLMs.

C Deployment status

We have deployed the RA-MTR framework in a flask (Grinberg, 2018) based web application (link will be made public upon acceptance). We plan to integrate this system with the publicly available and fully searchable historical encyclopedia (e.g., Gandhipedia¹⁹). We present an example page of our demo system in Figure 8. The figure shows the result when a user searches for the query “Find the famous quotes that Mahatma Gandhi had made about *health*”. The system extracts the most relevant quotes from the entire 100 volumes of the Collected Works of Mahatma Gandhi and highlights them in yellow.

D Data preprocessing details

QuoteR data : The *Project Gutenberg corpus* comprises more than 73000 e-books in textual form. We assign each of these books a unique bookID and divide each book into fixed-length (i.e., 200 word length) paragraphs, and assign each such paragraph a unique paragraphID. The distribution of the number of paragraphs per book and the number of tokens per quote is presented in Figure 9 and 10. We construct a TF-IDF weighted word-doc sparse matrix (Chen et al., 2017) from

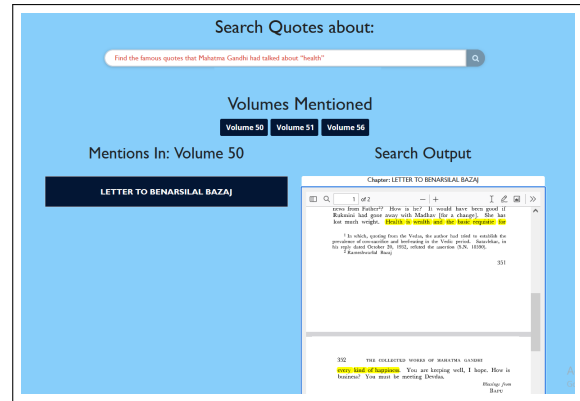


Figure 8: Example of a real time quote extraction from the Collected Works of Mahatma Gandhi. The output quote is highlighted in yellow in the pdf.

all the documents, index, and store this content in the sqlite db. For each of the quotes present in the QuoteR dataset we recursively search for the appearance of the quote in each of these books. Once a search gets a hit, we link the bookID with that particular quote (to be used for training the paragraph retrieval model). Since the authors in (Qi et al., 2022) stored the context from different sources and the correct mapping to the books is not present, we consider the 40 words preceding and following it as its left and right contexts, respectively. Similar to (Qi et al., 2022) the concatenation of the left and right contexts forms a complete context. We then store the context, quote, and positive paragraph (to be used for training the quote prediction model). Out of the 6108 unique quotes, we are able to find the occurrences of 1708 quotes from the *Project Gutenberg* and we finally construct 4889 quote-context-paragraph (one quote may contain multiple contexts) triples as examples for training and evaluating. The algorithm for generating the quote-context-paragraph triples is presented in Algorithm 1.

Gandhi data : Similarly, for the Gandhi quotes, we search for the quotes in the CWMG and find their appearance in a particular chapter of a book in the CWMG. We utilize the *Gandhipedia* (Adak et al., 2020) engine, which uses an elasticsearch based search engine to locate the quotes. We consider the 40 preceding and following words from each quote in the particular chapter as its context. In addition, we find that out of all the Gandhi quotes, three quotes are already present in the QuoteR set. We, therefore, remove them from the

¹⁹<https://gandhipedia150.in/en/>

Method	QuoteR (test)		Dataset Gandhi		Quotus	
	Top-1 F1	Top-3 F1	Top-1 F1	Top-3 F1	Top-1 F1	Top-3 F1
Vector-store based retriever + LLM reader (FLAN-T5-large)	14.5	19.23	18.2	22.5	25.27	31.2
Vector-store based retriever + LLM reader (FLAN-T5-xl)	13.32	19.4	16.54	24.33	35.0	38.25
Vector-store based retriever + LLM reader (bloomz-3b)	10.33	12.12	9.19	13.48	16.43	21.31
Vector-store based retriever + LLM reader (Falcon-7b)	10.08	13.34	17.05	22.35	29.73	36.73
Vector-store based retriever + LLM reader (Llama-3-8b-instruct)	14.81	18.76	17.53	23.28	31.38	39.77

Table 10: The result for the quote extraction using the different LLMs as reader.

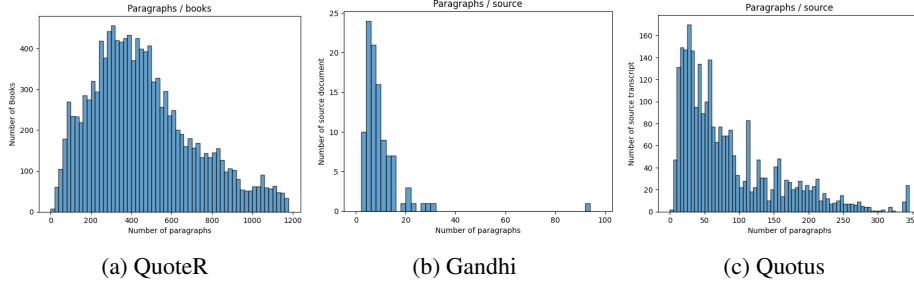


Figure 9: # of paragraphs per source documents.

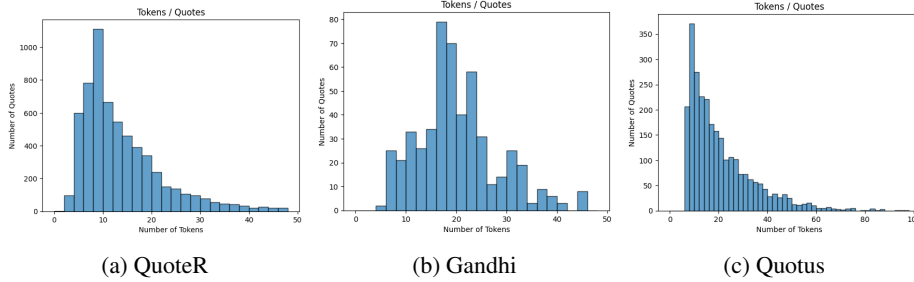


Figure 10: Distributions over source document, paragraphs, and quote lengths.

Gandhi data. Finally, we obtain 737 quote-context-paragraph triples.

Generating target context : Unlike the Quotus data, we do not have explicit target documents (i.e., where the quote needs to be recommended from source) for the QuoteR and Gandhi dataset. We synthesize the target context by paraphrasing the original context in the corresponding dataset. This is performed to reduce the overlapping words of the target context and the source document and to effectively evaluate the robustness of the methodology. We use ChatGPT²⁰ API to paraphrase the context. The examples of such paraphrased context are provided in Appendix G. The prompt used for paraphrasing:

As a paraphrasing expert can you rephrase the following input text? Ensure the rephrased text incorporates a different range of vocabulary compared to the

original text.

Input text: {<Input text>}

Rephrased text:

To analyse the hardness of the generated target contexts, we measure the word overlap between the original context and the rephrased context. We observe that the average word overlap ratio between the original and the rephrased contexts are - 0.19 and 0.18 for QuoteR and Gandhi data respectively. This indicates that the rephrased target context has significantly different words thus making the task of the paragraph retriever harder. We also measure whether the meaning of the rephrased contexts get significantly deviated from the original context. We use *GPT-4* to calculate the widely used *Faithfulness* of the rephrased contexts with respect to the original context. We observe an average *Faithfulness* of 0.84 and 0.88 for the QuoteR and Gandhi data respectively, which ensures that the rephrased contexts preserves the meaning of original contexts.

²⁰<https://openai.com/blog/chatgpt>

Quotus data : For the Quotus dataset, we utilize the Quoting POTUS website²¹ to collect a set of examples for our experiments. They release the transcripts and the collection of aligned quotes, containing the text of the quote in the news article, its aligned position within the source transcript, and the corresponding news article metadata (title, url, timestamp). We crawl the provided news article URLs and extract the body content of each news article using BeautifulSoup²². We are able to extract 10,114 news articles in this way (some of the links were not working and could not be crawled). To locate the quotes within the news articles, we utilize regular expressions and identify the appearance of 2,698 quotes. We then consider the 40 preceding and following words from each quote in the news article as its context. In the released dataset, the source transcript is already divided into several paragraphs, and the alignment of the quotes to the positive paragraph is also provided. As a result, we did not need to explicitly create the quote-paragraph alignment. This yields a total of 2,698 quote-context-paragraph triples, which we use for our experiments.

The Algorithm 1 shows the step-by-step procedure to prepare the dataset for our experiments. The auxiliary functions (i.e., Algorithms 3, 4 and 2) used in the algorithms are depicted in the subsequent algorithms.

E Baseline methods

Baselines: Both the retriever and the reader can have many variants which serve as ideal baselines. In the retriever part we use vanilla BM25 as a first baseline. Apart from the simple BM25 retriever, we employ BERT and T5 based re-ranking to improve paragraph selection. For input to BERT we tokenize the contexts and source document paragraphs into wordpieces (Wu et al., 2016) and cap them at predetermined lengths chosen as hyperparameters. BERT uses a special token [SEP] to separate paragraph from the context. So the final wordpiece input to the BERT is:

[CLS] context [SEP] paragraph [SEP]

Following (Wang et al., 2019), we fine-tune BERT-base using the pairwise loss. Thus, a single training example for paragraph BERT consists of $n + 1$ instances, i.e., one positive instance plus n negative

instances. Each of the $n + 1$ packed input sequences are fed to BERT independently. We use the final hidden vector $\mathbf{C} \in \mathbb{R}^h$ corresponding to the first input token [CLS] as the representation for each of the $n + 1$ sequences, where h is the size of the final hidden layer. In addition, we also fine-tune encoder-decoder based (T5, FLAN-T5) and decoder-only (Llama-3) re-ranking models in the same way as discussed in section 5.1.

For the reader part, as primitive baselines, we consider using the first and last sentences of each paragraph as potential quotes. To further explore, we also fine-tune the BERT and the SpanBERT pre-trained models on the BERT question answering architecture. We keep the same hyperparameter settings as the multi-task framework. Again, we fine-tune on 80% of the QuoteR data, and use 10% for validation before testing on the remaining 10%. In addition, we conduct similar few-shot experiments with the Gandhi and the Quotus dataset.

LLM based baselines: With the advancement of large language models (LLMs) such as T5 (Raffel et al., 2020), GPT-3 (Brown et al., 2020) it is important to observe their ability to perform the task of quote extraction. These models have proven to be highly valuable for contextual learning when provided with specific prompts in zero-shot scenarios. We replace the multi-task reader with different medium sized open-source LLMs such as FLAN-T5-large²³, FLAN-T5-XL²⁴, Bloomz-3b²⁵, Falcon-7b²⁶, Llama-3-8b²⁷ models to predict the most relevant quote given the paragraph and context. We use the below prompt:

You are an AI assistant in recommending a suitable 'quote' based on the context and your task is to extract a relevant quote from the given paragraph based on the context. Note that, the context and the paragraph may contain grammatical errors. DO NOT use any external information.

Context: "{context}"

Paragraph: "{paragraph}"

²³<https://huggingface.co/google/flan-t5-large>

²⁴<https://huggingface.co/google/flan-t5-xl>

²⁵<https://huggingface.co/bigscience/bloomz-3b>

²⁶<https://huggingface.co/tiiuae/falcon-7b>

²⁷<https://huggingface.co/meta-llama/>

Meta-Llama-3-8B-Instruct

²¹<http://snap.stanford.edu/quotus/vis/>

²²<https://www.crummy.com/software/BeautifulSoup/bs4/doc/>

Algorithm 1 Paragraph retrieval data generation

Require: *list_of_quotes*: list of selected quotes; *corpus_directory*: directory of the corpus (ex. Gutenberg)

```
1: quoteid_book_mapping ←  
   CREATE_QUOTE_TO_BOOK_MAPPING(list_of_quotes, corpus_directory)  
2: ctxid ← 0 ▷ Initialize Context Id  
3: ctxid_to_text ← {} ▷ Initialize Context Id to Context text mapping  
4: quoteid_to_ctxid ← {} ▷ Initialize Quote Id to Context Id mapping  
5: for all (quoteid, list_of_book_paths) in quoteid_book_mapping do  
6:   dataset ← [] ▷ Dataset to be used for training and testing paragraph retrieval  
7:   quoteid_to_ctxid[quoteid] ← []  
8:   for all book_path in list_of_book_paths do  
9:     paragraphs ← SEGMENT_BOOK(book, paragraph_length = 200) ▷ Segment the book  
     contents into several paragraphs  
10:    save(docid_to_text)  
11:    for all paragraph do  
12:      if quote in paragraph then  
13:        ctx ← CREATE_CONTEXT(quote, paragraph) ▷ Creating context for a quote  
14:        ctxid_to_text[ctxid] ← ctx  
15:        dataset.append([ctxid, [pos_para_id], [candidate_id]])  
16:        ctxid ← ctxid + 1  
17:        if quoteid in quoteid_to_ctxid.keys() then  
18:          quoteid_to_ctxid[quoteid].append(ctxid)  
19:        else  
20:          quoteid_to_ctxid[quoteid] ← [ctxid]  
21:        end if  
22:      end if  
23:    end for  
24:    save(dataset)  
25:  end for  
26: end for  
27: save(ctxid_to_text)  
28: save(quoteid_to_ctxid)
```

Algorithm 2 Create quote to book mapping

```
1: function CREATE_QUOTE_TO_BOOK_MAPPING(list_of_quotes, corpus_directory)
2:   Input: list_of_quotes, corpus_directory
3:   Output: quote_to_book_mapping
4:   quote_to_book_mapping  $\leftarrow$  {}
5:   for all quote in list_of_quotes do
6:     for all book_path in corpus_directory do
7:       if quote found in book_path then
8:         if quote in quote_to_book_mapping then
9:           quote_to_book_mapping[quote].append(book_path)
10:        else
11:          quote_to_book_mapping[quote]  $\leftarrow$  [book_path]
12:        end if
13:      end if
14:    end for
15:  end for
16:  return quote_to_book_mapping
17: end function
```

Algorithm 3 Segment book into paragraphs of fixed length

```
1: function SEGMENT_BOOK(text_document, paragraph_length)
2:   Input: text_document, paragraph_length
3:   Output: paragraphs
4:   paragraphs  $\leftarrow$  {}
5:   current_paragraph  $\leftarrow$  ""
6:   current_paragraph_id  $\leftarrow$  0
7:   for word in text_document.split() do
8:     current_paragraph  $\leftarrow$  current_paragraph + "" + word
9:     if len(current_paragraph)  $\geq$  paragraph_length then
10:      paragraphs[current_paragraph_id]  $\leftarrow$  current_paragraph.strip()
11:      current_paragraph  $\leftarrow$  ""
12:      current_paragraph_id  $\leftarrow$  current_paragraph_id + 1
13:    end if
14:  end for
15:  if len(current_paragraph)  $>$  0 then
16:    paragraphs[current_paragraph_id]  $\leftarrow$  current_paragraph.strip()
17:  end if
18:  return paragraphs
19: end function
```

Algorithm 4 Generate context for a quote in a paragraph

```
1: function CREATE_CONTEXT(quote, paragraph)
2:   Input: quote, paragraph
3:   Output: context
4:   context  $\leftarrow$  ""
5:   quote_position  $\leftarrow$  paragraph.find(quote)
6:   if quote_position  $\neq$  -1 then
7:     preceding_40  $\leftarrow$  paragraph[quote_position].split(" ")[-40:]
8:     following_40  $\leftarrow$  paragraph[quote_position + len(quote):].split(" ")[:40]
9:     context  $\leftarrow$  " ".join(preceding_40) " ".join(following_40)
10:  end if
11:  return context
12: end function
```

Just extract the relevant quote without any other sentence:

F Model implementation details

Retriever : For retriever we use *langchain API*²⁸, employ *recursive_text_splitter*²⁹ for splitting the document, *chromaDB* as vector store. For fine-tuning the reranking models we use *huggingface API*³⁰.

FLAN-T5: We fine-tune our T5 models (base³¹, large³²) and *FLAN-T5-large*³³ with a learning rate of $2e^{-5}$ and a weight decay of 0.01 for a maximum of 10 epochs with a batch size of 4. We use a maximum of 1024 input tokens and one output token. Training T5 base, large, and Flan-T5-large take approximately 2, 5, and 6 hours overall, respectively, on a single RTX 4090 GPU. We use greedy decoding during inference and used *output_logits=True* while generating text to retrieve unprocessed probabilities assigned to a token. We use same hyperparameter setting for *Llama-3-8b-instruct*

bert-base: For fine-tuning *bert-base*³⁴ for the paragraph retrieval task, we search over a batch-size $\in \{4, 8, 16\}$, and set the learning rate of $2e^{-5}$. We set the maximum number of epochs to 10. We also perform a search over $n \in \{3, 6, 9, 12\}$ sampled negative paragraphs per positive paragraph for our

paragraph ranking model. We select the best model using the dev set and the best paragraph model is trained with 9 negative examples and a batch size of 16. We used single NVIDIA Tesla P100 GPU for training the model.

Reader : For the span selection model (the multi-task and other transformers based baseline models), we cap the total length of the context and paragraph to 384 length wordpieces. In case the total length exceeds the maximum length (i.e., 384), we only truncate the paragraph. Similarly, for the quotable phrase identification model (i.e., the sequence tagger model in the multi-task setting) we select a maximum length of 384. We fine-tune the publicly available *spanbert-large*³⁵, by setting the batch-size $\in \{4, 8\}$, learning rate of $2e^{-5}$. We fine-tune the model over 10 epochs and use early stopping based on the dev set. Again we used single NVIDIA Tesla P100 GPU for training the model. For the multi-task framework, the training process took 3.5 hours to complete. For the LLM inference we use single NVIDIA Tesla P100 GPU. Additionally, we applied 4bit quantization while loading the larger LLMs as those models would not fit in our GPU.

G Examples of paraphrased context

Table 11 shows one paraphrased example from QuoteR and Gandi dataset which were used as the target context. Quotus dataset having a separate target article, we did not require paraphrasing the context.

²⁸https://python.langchain.com/docs/modules/data_connection/

²⁹https://python.langchain.com/docs/modules/data_connection/document_transformers/recursive_text_splitter

³⁰<https://huggingface.co/>

³¹<https://huggingface.co/t5-base>

³²<https://huggingface.co/t5-large>

³³<https://huggingface.co/google/flan-t5-large>

³⁴<https://huggingface.co/bert-base-uncased>

³⁵<https://huggingface.co/SpanBERT/spanbert-large-cased>

Dataset	Actual Context	Paraphrased Context
QuoteR	and for the great Peasant Revolt of 1381. John Ball's famous rhyme condensed the scorn for the nobles, the longing for just rule, and the resentment at oppression, of the peasants of that time and of all times:– " A hundred years after the Black Death the wages of a common English laborer—we have the highest authority for the statement—commanded twice the amount of the necessities of life which could have been obtained for the wages paid under	For the significant Peasant Revolt of 1381, John Ball's renowned rhyme encapsulated disdain for the nobles, the yearning for fair governance, and resentment towards oppression. A century after the Black Death, the wages of an ordinary English laborer, as verified by the highest authority, were double the necessities obtainable with previous wages.
Gandhi	For, highest perfection is unattainable without highest restraint. Suffering is thus the badge of the human tribe. The goal ever recedes from us. The greater the progress, the greater the recognition of our unworthiness. Full effort is full victory. Therefore, though I realize more than ever how far I am from that goal, for me the Law of complete Love is the law of my being. Each time I fail, my effort shall be	The pinnacle of perfection requires the utmost restraint, and suffering becomes the emblem of the human experience. The goal remains elusive, and progress accentuates our sense of unworthiness. Full effort equates to complete victory. Despite realizing the vast distance from the goal, the Law of complete Love governs my existence. Each failure only strengthens my resolve.

Table 11: Examples of paraphrased contexts for QuoteR and Gandhi datasets.

H Examples of LLM generated quotes

In Table 12 we provide examples of quotes extracted by different LLMs used in our experiments for a specific context and paragraph. We observe that, larger models (such as FLAN-T5-XL, Llama-3-8b) generate better quotes compared to the smaller models. However, Llama-3 merges some part of the context (“Sweet dewdrops”) in the predicted quote. This is one of the precise reasons why standalone LLMs cannot be reliably used in the quote extraction task, as it may augment the original text which could affect the sanctity of the quote.

Context	Paragraph	Actual quote	Used LLM	Extracted quote
Tonight, the heavens and the earth will mourn your passing, as sweet dewdrops fall from the sky. Your life is fleeting, and your time is coming to an end. The rose, with its	earth and sky, Sweet dewds shall weep thy fall tonight, For thou must die. Sweet rose, whose hue, angry and brave, Bids the rash gazer wipe his eye, Thy root is ever in its grave, And thou must die. Sweet spring, full of sweet days and roses, a box where sweets compacted lie , My music shows you have your closes, And all must die. Only a sweet and virtuous soul, Like seasoned timber, never gives But when the whole world turns to coal, Then chiefly lives. George Herbert. THE LIE. Go, Soul,	Sweet spring, full of sweet days and roses, a box where sweets compacted lie	Bloomz-3b	earth and sky, Sweet dewds shall weep thy
			FLAN-T5-large	rage, is ever in its grave
			FLAN-T5-XL	Sweet dewds shall weep thy fall tonight, For thou must die. Sweet rose, whose hue, angry and brave, Bids the rash
			Falcon-7b	Sweet dewdrops fall from the sky
			Llama-3-8b	Sweet dewdrops fall from the sky. For thou must die. Sweet rose, whose hue, angry and brave, Bids the rash gazer wipe his eye, Thy root is ever in

Table 12: Quotes extracted by different LLMs used for a specific context and paragraph