# Enhancing Abstractive Text Summarization
# with an Upfront Extractive Layer

**Prakash Krishnan**

## Abstract

Abstractive text summarization of long-form, unstructured documents has been a challenge for researchers. Encoder-only and encoder-decoder transformer models have changed the paradigm and offer many alternatives for an effective summarization. In prior research work, extractive and abstractive summarization were viewed as independent language model tasks. Here in this research paper, we explore a two-step process of abstractive summarization for long-form unstructured documents. As a first step, we conduct an extractive summarization and feed the extractive summarization through an abstractive layer. For the extractive layer, we test our hypothesis by applying different encoders (SBERT and BERT) and classification layers (Clustering, Query Retrieval, and Cosine Similarity). For the abstractive layer, we explore the use of a pre-trained transformer (PEGASUS, Checkpoint-CNN / Daily Mail).

## 1 Introduction

Abstractive document summarization as a language model task is popular given its application to several use cases such as news summaries, social media reviews, financial summaries, scientific articles, and legal documents. Abstractive summarization is challenging as the input text data can be presented in several formats such as long or short-form text, organized or unorganized, and narrative or conversational style. A single methodology to generate abstractive summarization across all input data types may not work optimally.

Our inspiration for this paper came from the research published by Yang Liu et al., 2019, where a novel document-level encoder based on BERT expressed semantics and generate representations of its sentences. This represented a breakthrough in text summarization by leveraging pre-trained language models as encoders for sentence and paragraph-level documents.

The BERT model suggested by Yang Liu et. al had several challenges for sentence classification such as:

- BERT model was optimized for word or token-level embeddings and not sentence-level embeddings.
- Application of BERT to text summarization is tricky as BERT is trained as a masked-language model and the output vectors are tied to tokens and not sentences.
- To compute accurate sentence similarity with BERT a cross-encoder structure is required where two sentences are passed through a BERT model with a classification head at the top of BERT to derive a similarity score. This is not scalable and causes computational overhead.

Nils Reimers et al., 2019 developed a modified BERT network using siamese network architecture (Sentence-BERT) to derive semantically meaningful sentence embeddings. This opened new tasks such as large-scale semantic comparison, information retrieval, and clustering. These were not possible in the original BERT architecture. See Figure 1 for a Sentence BERT architecture.
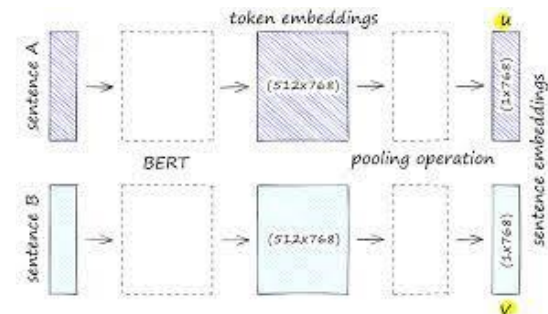


Figure 1: Sentence BERT architecture. Example an input length of 512, an embedding dimension of 768, and a pooling function.

In this research paper, we explore the following novel research questions:

1. For long-form texts such as news articles, would a two-step summarization process with an upfront extractive summarization generated using a sentence transformer and followed by an abstractive summarization lead to superior abstractive summaries as evaluated by ROUGE scores?
2. How can we improve extractive summaries of long documents by leveraging the capabilities of sentence transformers and query retrieval?
3. How do the models perform for different types of datasets such as conversational style MediaSum or more organized CNN/Daily Mail?

## 2 Background

The novelty of our approach is that prior work mostly treated extractive and abstractive summarizations as separate language model tasks.

Here in this paper, for conversational style text, we explore whether an intermediate enhanced extractive layer with a following abstractive model can lead to superior abstractive summaries.

Models from prior research include the following:

**Pre-Transformer Models -**
Before transformer models, sentence classification was derived by leveraging topic representation techniques such as Latent Semantic Analysis (Harendra Bhandari et al., 2008), Frequency Driven such as word probability, TF-IDF, and Topic Word approach such as Luhan's technique.

A Graph-based ranking approach for sentence extraction was developed by Rada Mihalcea (2005, Dept of CS, University of Texas).

Finally on the Machine Learning side, Sequence-to-Sequence RNNs have been proposed by Ramesh Nallapati et al., 2016 for text summarization.

**Post-Transformer Models-**
In the original BERT-based encoder-decoder model (Yang Liu et al., 2019) the extractive model is built on top of the BERT encoder by stacking inter-sentence transformer layers to capture document-level features. For the abstractive model, an encoder and decoder model with the encoder being the extractive model as described above and the decoder being a 6-layered transformer model.

Sentence-BERT (Nils Reimers et al., 2019) is a novel approach to sentence-level embeddings and text classification using sentence transformers. This is done by mapping each sentence to a vector space and pre-computing sentence vectors that can be stored and then used whenever required. This gave rise to two approaches as illustrated in the paper Sentence-BERT published by Nils Reimers et al., 2019. The first approach was to leverage the standard BERT model and build a sentence embedding by averaging the values across all token embedding outputs or using the first [CLS] token. The second approach is the siamese network architecture (Sentence Level BERT (SBERT)) that derives fixed-size vectors for input sentences. Using a similarity measure like Cosine Similarity, Manhattan, or Euclidean Distance, similar sentences that are semantically close can be ascertained. Both approaches are presented in Section 6 Experimental Results.

## 3 Dataset

MediaSum data set on Hugging Face is used for model evaluation. This large-scale media interview dataset contains 463.6K transcripts with abstractive summaries, collected from interview transcripts and overview/topic descriptions from NPR and CNN.

For our study, we restricted the dataset size to 1,000 examples for extractive summarization and 100 for abstractive summarization due to computational limitations.

CNN / Daily Mail dataset was also used for comparative purposes given its narrative non-conversational structured nature.

To split the original text into sentences, pre-processing was performed to remove punctuations. (Example: /n, '?</s>, \, _, |, /, Dr., !, ?)

## 4 Methodology

In this section, we describe our methodology and modeling approach. See Figure 2 for the End-To-End Pipeline.
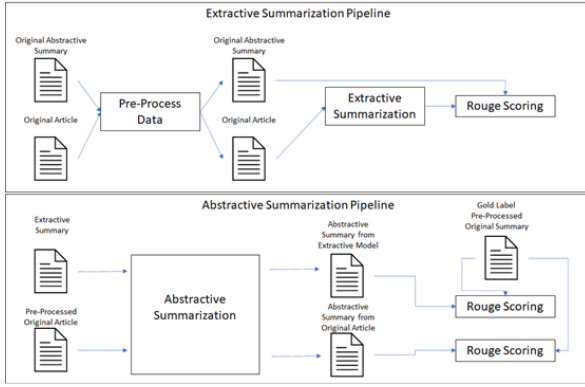
Figure 2: End-To-End Pipeline. First, sentence summarization is generated using a Classification Layer. Sentence extractions are fed into an encoder-decoder PEGASUS

As a first step, several extractive models were developed using pre-trained encoders BERT and SBERT (Sentence BERT), and a classification layer on top of the encoder. A naive baseline model using TF-IDF was also developed to provide a baseline metric.

Classification techniques explored for the extractive model were:

● Centroid-Sentence Based Summarization where a centroid vector of embedding representation of the input document is derived and the cosine similarity of each sentence vector from the centroid is computed. The Top 'N' sentences with the highest similarity scores were selected for extraction. Top 'N' sentences are computed using a compression ratio that is applied to the original document. The compression ratio is a hyperparameter and user selectable. Trigram blocking is applied to the extraction to ensure words are not repeated in the extractive summary.

● Section-Clustering where the sentences in the input document are clustered using KMeans and the abstractive label is used as a query vector to retrieve sentences from each cluster nearest to the label. The compression ratio is a hyperparameter and user selectable. Trigram blocking is applied to the extraction to ensure words are not repeated in the extractive summary.

The extractive models were run as both unsupervised and supervised models. See Table 1 for a list of extractive models evaluated.

For the extractive summarization, we did not have a gold label extractive summary to assess summarization quality. We used the original abstractive summary as a label to evaluate the different models and classification techniques. We recognize this is not an ideal comparison but offers an intermediate evaluation of the extractive models.

For the supervised extractive model, the input document was first clustered using KMeans, then a query embedding vector was generated using the original article abstractive summary. The distance of the query vector to each cluster centroid was computed. Top N Clusters were derived and then applied KMeans Nearest Neighbor to generate the extractive text summarization. The following hyper-parameters were used - Number of clusters = 5, number of nearest neighbors = 10, and a compression ratio of 0.3.

For the abstractive summarization, we feed the extractive summaries generated into the abstractive model to generate the model summaries. Summaries were generated both from the original full article and also from the extractive summary to determine whether an enhanced extraction provided a superior abstraction. Pre-Trained PEGASUS was used as the abstractive model.

Table 1: Extractive Model Details

| Model | Supervised | Selection Criteria |
|---|---|---|
| TFIDF | No | Top N Sentences |
| Pre-trained sentence transformer ('all-MiniLM-L6-v2') with sentence-level embeddings. | No | Agglomerative Clustering |
| Pre-trained sentence transformer ('all-MiniLM-L6-v2') with sentence-level embeddings. | No | Cosine similarity and degree centrality. |

| | | | | | | |
|---|---|---|---|---|---|---|
| Sentence level built from scratch using the original Bert Model. | No | Cosine similarity and degree centrality. | | | | |
| Pre-trained sentence transformer ('all-MiniLM-L6-v2') with sentence-level embeddings. | Yes | KMeans and Nearest Neighbors Clustering<br><br>Abstractive Gold Label used as a query | | | | |

We selected PEGASUS as the abstractive model as it was pre-trained for abstractive summarization. PEGASUS uses an encoder-decoder model for sequence-to-sequence learning. PEGASUS is unique in its pre-training as it factors important sentences in its output closely resembling a human abstractive process. CNN Dailymail checkpoint was used for PEGASUS and no fine-tuning was performed.

## 5 Evaluation Metrics

For the extractive summarization, we looked at both RL-P and RL-F1 to compare model performances. For this intermediate layer, RL-P was important to ensure the extractive summary captured the key semantics of the abstractive gold summary. However, RL-F1 was a better measure for abstractive summarization.

For the abstractive layer, we compare the model abstractions versus the original abstractive summaries as gold labels.

Additionally, a subjective (manual) evaluation of a random example was used to evaluate for coherence, fluidity, and faithfulness.

## 6 Experimental Results

The experimental results on MediaSum and CNN/DailyMail datasets are shown in Table 2 and Table 3.

Table 2 presents the experimental results from the extractive models. TF-IDF is the baseline model to compare against.

Table 3 presents the experimental results from the abstractive models. The baseline model is no extractive summarization is performed before abstractive text summarization.

Table 2: Extractive Summarization Results

| ID | Model | Data | R1 | R2 | RL-P | RL-F1 |
|---|---|---|---|---|---|---|
| | | | Baseline Model | | | |
| 1.a | TF-IDF | Media Sum | 13.71 | 1.56 | 15.88 | 9.21 |
| 1.b | TF-IDF | CNN/DM | 15.38 | 2.96 | 11.86 | 11.27 |
| | | | Unsupervised Models | | | |
| 2.a | SBert Summ | Media Sum | 12.69 | 4.56 | 46.04 | 8.82 |
| 2.b | SBert Summ | CNN/DM | 23.45 | 9.00 | 40.34 | 16.38 |
| 3.a | SBert Trans | Media Sum | 12.30 | 4.49 | 46.66 | 8.80 |
| 3.b | SBert Trans | CNN/DM | 22.24 | 8.98 | 39.78 | 15.22 |
| 4.a | Bert Trans | Media Sum | 10.16 | 3.45 | 47.04 | 7.32 |
| 4.b | Bert Trans | CNN/DM | 21.03 | 8.52 | 41.24 | 14.17 |
| | | | Supervised Models | | | |
| 5.a | Supv-SBert | Media Sum | **13.93** | **6.51** | **54.73** | **10.39** |
| 5.b | Supv-SBert | CNN/DM | **26.58** | **12.92** | **48.70** | **18.65** |

Table 3: Abstractive Summarization Results

| ID | Model | Data | Step 1 Ext Layer RL | Step 2 Abs Layer RL | Delta |
|----|-------|------|---------------------|---------------------|-------|
| 2.a | SBert Summarizer + PEGASUS | MediaSum | 8.82 | 17.93 22.14** | 9.11 |
| 2.b | SBert Summarizer + PEGASUS | CNN/DM | 16.38 | 23.10 41.52** | 6.72 |
| 3.a | SBert Transformer + PEGASUS | MediaSum | 8.80 | 20.35 22.14** | 11.55 |
| 3.b | SBert Transformer + PEGASUS | CNN/DM | 15.22 | 21.06 41.52** | 5.84 |
| 4.a | Bert Transformer + PEGASUS | MediaSum | 7.32 | 19.10 22.14** | 11.78 |
| 4.b | Bert Transformer + PEGASUS | CNN/DM | 14.17 | 23.39 41.52** | 9.22 |
| 5.a | SBert Transformer + PEGASUS (Supervised) | MediaSum | **10.39** | **23.66** 22.14** | **13.27** |
| 5.b | SBert Transformer + PEGASUS (Supervised) | CNN/DM | **18.65** | **29.17** 41.52** | **10.52** |

Note:
** RL-F1 Scores without Extractive Layer

Table 5: Manual Evaluation of Abstractive Summarization for MediaSum



Table 6: Manual Evaluation of Abstractive Summarization for CNN / Daily Mail



# 7 Discussion and Analysis

We evaluated the extractive and abstractive summarizations using ROUGE. We report unigram and bigram overlap (R1 and R2) and longest common subsequence (RL) for two datasets MediaSum and CNN/DailyMail.

From Table 1 Extractive Summarization Results, the following conclusions can be ascertained:

- TF-IDF baseline model outperformed the unsupervised transformer models for MediaSum on Rl-F1. As MediaSum dataset is conversational style question and answer format, the baseline model gives more importance to words that are frequently mentioned in the question and responses. These frequently mentioned words are also perhaps captured in the abstractive summary resulting in higher RL-F1 scores for the baseline model. (See 1.a, 2.a, 3.a, 4.a)

- We do not see this phenomenon with the CNN/Daily Mail dataset which is more structured, organized, and narrative style. All the

unsupervised models outperformed the baseline model on RL-F1 score. (See 1.b, 2.b, 3.b, 4.b)

- Pre-trained sentence transformer models (SBert Summarizer, SBert Transformer) performed better than the pre-trained word token-based transformers (Bert Transformer). (See 2.a, 2.b, 3.a, 3.b)
- The sentence transformers can capture better sentence-level context than the original BERT model. This was the case with both MediaSum and CNN / Daily Mail.
- A supervised model (SBert Transformer with Kmeans + Nearest Neighbor) leveraging the original abstractive summary as a query for sentence retrievals provided richer extraction as illustrated by the higher ROUGE Scores for both MediaSum and CNN / Daily Mail (See 5.a and 5.b).

From Table 2 Abstractive Summarization Results, the following conclusions can be ascertained:

- For unsupervised models (2.a, 2.b, 3.a, 3.b, 4.a, 4.b) there were improvements in the ROUGE-L scores for both MediaSum and CNN/Daily Mail.
- Extractive summarization seems to provide most benefit when the source document is large and unstructured. MediaSum models were most benefited via the two-step process. CNN/Daily_Mail Models, showed improvements but not as much as MediaSum.
- Extractive summarization generated via query retrieval yielded superior RL scores for the abstractive summarization (5.a, 5.b)
- For all the models in Table 2, stand-alone abstractive summarization yielded superior RL-F1 scores versus the two-step process.
- However, in looking at a random example in Table 5 and Table 6, it does appear that a prelim extractive layer can offer a more complete abstractive summarization as extractive layer can potentially determine underlying semantic similarities.

## 8    Conclusion and Next Steps

To conclude, there are benefits to performing a preliminary extractive summarization for large unstructured texts leveraging query retrieval where available. Query retrieval can help determine the sentence vectors closest to the query vector and capture semantic similarities. An abstractive layer

applied to the extractive summary can yield a coherent and fluent summarization of the original text.

For structured text, an extractive summarization can help but the benefits are lower. Finally, the nature of the input text plays a role in the summarization technique as a conversation-style article would need a dialog summarization model as opposed to narrative summarization. Fine-tuning the model can also be explored to improve ROUGE scores.

## 9    References

Samir Abdaljalil, and Houda Bouamor. 2021. An Exploration of Automatic Text Summarization of Financial Reports. *In Proceedings of the Third Workshop on Financial Technology and Natural Language Processing.*

Harendra Bhandari, Masashi Shimbo, Takahiko Ito, Yuji Matsumoto. 2008. Generic Text Summarization Using Probabilistic Latent Semantic Indexing. *Proceedings of the Third International Joint Conference on Natural Language Processing: Volume-I.*

Yang Liu and Mirella Lapata. 2019. Text Summarization with Pretrained Encoders, Yang Liu and Mirella Lapata, *In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing.*

Yang Liu. 2019. Fine-Tune BERT for Extractive Summarization, *arXiv*.

Jonathan Pilaut, Raymond Li, Sandeep Subramanian, Christopher Pal. 2020. On Extractive and Abstractive Neural Document Summarization with Transformer Language Models. *In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing.*

Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. *In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing.*