

Attention-based Clinical Note Summarization

Neel Kanwal
University of Stavanger
Norway, Stavanger
neel.kanwal@uis.no

Giuseppe Rizzo
Links Foundation
Torino, Italy
giuseppe.rizzo@linkfoundation.com

ABSTRACT

The trend of deploying digital systems in numerous industries has induced a hike in recording digital information. The health sector has observed an extensive adoption of digital devices and systems that generate large volumes of personal medical records. Electronic health records contain valuable information for retrospective and prospective analysis that is often not entirely exploited because of the dense information storage. The crude purpose of condensing health records is to select the information that holds most characteristics of the original documents based on reported disease. These summaries may boost diagnosis and extend a doctor's time with the patient during a high workload situation like the COVID-19 pandemic. In this paper, we propose applying a multi-head attention-based mechanism to perform extractive summarization of meaningful phrases in clinical notes. This method finds major sentences for a summary by correlating tokens, segments, and positional embeddings. The model outputs attention scores that are statistically transformed to extract key phrases and can be used to projection on the heat-mapping tool for visual and human use.

KEYWORDS

Natural Language Processing, Electronic Medical Records, Electronic Health Records, Text Summarization, Multi-head Attention, Transformer Models, Deep Learning, ICD-9, MIMIC-III, Clinical Notes

1 INTRODUCTION

Presenting text in a shorter form has been practiced in human history long before the birth of computers. A summary is defined as a document that conveys valuable information with significantly less text than usual [40]. Summarization can be sensitive in the medical domain owing to the nature of medical abbreviations and technicalities. The task of summarization can be categorized into two categories from a linguistic perspective. *Extractive* summarization is an indicative approach where phrases are scored based on similarity weights and chosen to produce verbatim. Contrarily, *abstractive* summarization is an informative approach that requires understanding a topic and generating a new text using fusion and compression. It relies on novel phrases, lexicon, and parsing for language generation [8].

Natural language processing (NLP) has been valuable to clinicians in saturated work environments. Health information systems and chatbots have reduced the workload of doctors during the COVID-19 pandemic. The notion of presenting a condensed version of a document using computers and algorithms became part of significant research in the late 1950s. Among the first ideas of machine summarization, researchers [32] presented a method for the

creation of literature abstracts in exploratory research of IBM Journal. This statistical approach is based on finding word frequency and scoring them based on higher significance. Baxendale et al. [7] came up with the idea based on grammatical position in the text. Other researchers like Lin [29] and Strzalkowski [44] proposed query-based summarization that is a direct descendant of information retrieval technique. The method was similar to multi-query vector in multi-head attention to performs an effective mapping using a query and a key-value pair to an output.

Some fundamental questions that arise while summarization are i) which content is essential to select and ii) how to create a shorter version of it [41]. In this research work, we propose a transformer-based method for selecting meaningful phrases from clinical discharge summaries and extracting them by preserving the sense of clinical notes based on the identified disease. Transformer [48] is being used as a new tool for text analysis tasks in NLP [6]. Inspired by [3], we have further tuned a BERT model utilizing discharge notes classified by diseases of the MIMIC-III dataset, Which captures important syntactic information based on International Classification of Diseases (ICD-9) labels. We extract a discrete attention distribution from the last layer. This probabilistic distribution is later translated using power transforms [18] to create monotonic data over bell curve [12]. Finally, the summary is comprised of sentences that have higher attention scores than the mean attention scores.

This paper organizes sections as follows: Section 2 illustrates various types of extractive summarization approaches and their usescases on medical documents. In Section 3, we described the methodology for intended task. Section 4 presents various evaluation methods and states a suitable method for this task. We have displayed results in Section 5. Finally, we conclude the discussion in Section 6 with limitations of interpretation in the medical domain in Section 7.

2 RELATED WORK

The early works on summarization are based on many different surface-level approaches for the intermediate representation of text documents. These methods focus on selecting top sentences based on some greedy algorithms and aim for maximizing coherence and minimizing redundancy [1] [54]. These techniques can be further generalized into:

- **Corpus-based Approach:** It is a frequency-driven approach based on common words that are often repeated and do not carry salient information. It relies on an information retrieval paradigm in which common words are considered as query-words. SumBasic [47] is a similar centroid-based method that uses word probability as a factor of importance for sentences. Words in each centroid with higher probabilistic values are selected for a summary.

- **Cohesion-based Approach:** Some techniques fail while extraction when the text is bound to anaphoric expressions.¹ Lexical chains are used to find a relation between two sentences. Brin et al. [9] proposed a co-reference system that uses cohesion for this purpose in web search. In clinical notes, anaphoric expressions are used frequently, but they refer to the same subjects meaning that we have a single relation to the patient.
- **Rhetoric-based Approach:** This method relies on forming text organization in a tree-like representation [19] [26]. Text units are extracted based on their position close to the nucleus. For clinical summaries, we often may have multiple nuclei for different diseases.
- **Graph theoretic Approach:** A few popular algorithms like HITS [24] and Google’s PageRank [33] instigated base for graph-based summarization. It helps to visualize intra-topic similarity where nodes present a number of topics and vertices show their cosine similarity with sentences [20]. It makes visual representation easy with different tools like MDIGESTS [10].
- **Machine Learning Approach:** ML models are outperforming for nearly all kinds of tasks, including text summarization. A recent trend of analyzing text using Bayesian models has gained popularity [4]. Neural networks better exploit hidden features from the text. Attention mechanism coupled with convolution layers helps to select importing phrases based on their position in the document. ML models treat summarization as a classification task. Transformer-based models combined with clustering methods are used in an unsupervised fashion. Miller et al. [34] used BERT to make text encoding and applied K-means to find sentences close to the centroid to learn health-informatics lecture summaries. This model offered a weakness for large documents since the extraction ratio is fixed for K sentences. Liu et al. [31] trained an extractive BERT model from abstractive summaries using a greedy method to generate an oracle summary for maximizing the ROGUE score. BERTSUM [31] used a trigram blocking method to extract candidate sentences based on golden abstractive summaries in CNN/daily mail dataset.

2.1 Implementations for Medical Documents

Clinicians heavily rely on text information to analyze the condition of the patient. Vleck et al. [46] followed a cognitive walk-through methodology by identifying certain relevant phrases to medical understanding. They developed a new extrinsic evaluation for the results. Laxmisan et al. [27] formed a clinical summary screen to integrate with existing personal health record systems. The core purpose was to avail more interaction time for a clinician.

Felblowitz et al. [17] proposed a five-stage architecture to facilitate clinical summarization tasks. It was more based on producing short laboratory reports. The AORTIS model was evaluated using cohesion index Kappa. The abbreviation of AORTIS described distinct phases of the framework, namely Aggregation, Organization, Reduction, Interpretation and Synthesis (AORTIS) [50]. Alsentzer

¹Anaphoric expressions are words that relate the sentences using pronouns such as he, himself, that.

et al. [2] did a similar job using Bayesian modeling. They utilized heterogeneous sampling and topic modeling stated in the research paper [39]. They materialized a Concept Unique Identifier CUI-Upper bound to choose a phrase that has a high probability of being classified as a disease.

Thomas et al. [23] presented a semi-supervised graph-based method for summarization using neural networks and node classification. The model was limited to datasets other than the clinical domain. Other researchers followed a similar kind of method, G-FLOW for Multi-document Summarization [11] [51]. Azadani et al. [5] carried out these ideas to biomedical summarization. They formulated a model based on graph clustering that forms a minimum spanning tree using Unified Medical Language System (UMLS).

An ontology-oriented graphical representation was proposed in this article [53]. It described a clustering method that uses data centrality and mutual refinement to sample to limit compression. Mis-Classification index (MI), a new evaluation metric to verify cluster purity was used as primary evaluation. These implementations were done on different datasets and the organization of clinical notes. Generally, graphical methods concluded better results in most of the cases.

3 METHODOLOGY

We propose a method to extract important phrases from a clinical note. This approach uses a base BERT-model fined-tuned on ICD-9 labeled MIMIC-III discharge notes. The model was trained mainly to identify ICD-9 labels based on described symptoms and diagnosis information. This summarization approach works effectively when reference or human-made summaries are not available. The model outputs attention scores for all sentences from discharge notes. We are extracting sentences whose attention scores are higher than the mean value of all other sentences in the original note. Our model is compared against three baselines [16] [51] [34] using divergence methods of word probability distributions for quantitative analysis. Table 2 represents qualitative analysis against chosen baseline approaches.

3.1 Dataset

MIMIC-III is an open-access publicly available database of unstructured health records [22]. This dataset is available online.² It includes raw notes of 36,998 patients for each hospital stay. Each discharge note is tagged with a unique label for the identified disease. In total, we have 47,724 clinical discharge notes that comprise several details from radiology, nursing, and prescription. These notes can be equated with a multi-topic document based on the nature of multiple labels for each medical note. A reduced sample of randomly selected 100 notes of this dataset is selected to quantitatively and qualitatively assess the performance of the model. The dataset was only published for labeling diagnostics and does not contain any kind of reference summary.

3.2 Neural Architecture

We have utilized neural architecture that is built on top of transformer [48]. Bidirectional Encoder Representation from Transformer (BERT) is multi-layer neural architecture. It has two major variants.

²<https://physionet.org/content/mimiciii-demo/1.4/>

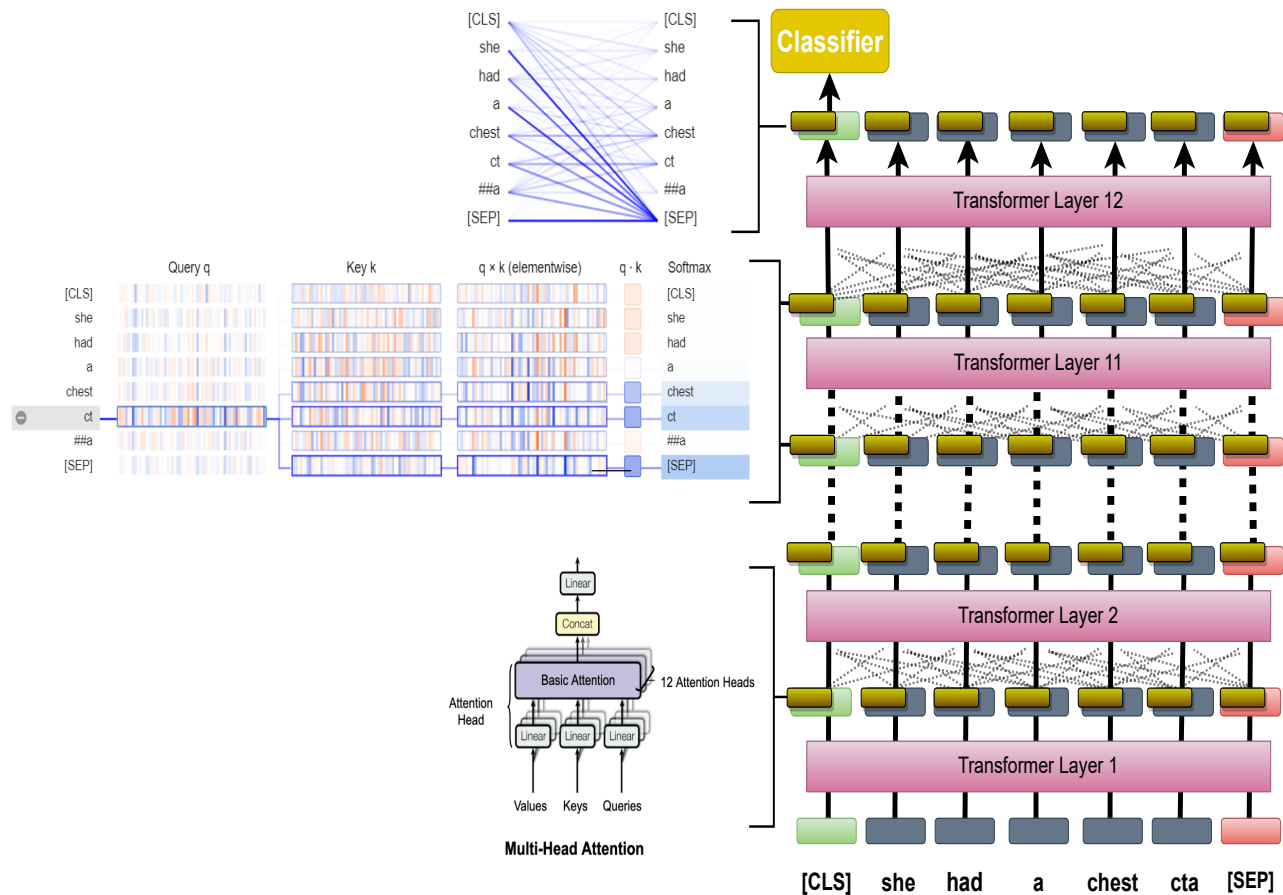


Figure 1: BERT base architecture with 12 transformer layers. Every layer carries embeddings (blue) and attention scores (green) for corresponding tokens. This multi-head attention correlates to every other word as seen in the images at left. The classifier token from the last layer is used for every sentence are described in the methodology.

We have used a base variant with 12 layers (transformer blocks), 768 hidden sizes, 12 attention heads, and 110 million parameters. The language model has shown significant improvements in various language processing tasks with fine-tuning [21] [13]. The BERT model usually creates embeddings in both directions for the representation of the inputs [15]. We have found that attention heads corresponding to delimiter tokens are remarkably effective for semantic understanding. This approach is trained using the classes available in the MIMIC-III dataset that map clinical notes to pathologies.

3.3 Implementation Details

The work presented in this paper employs fine-grained understanding of the notes to signify sentences that are relevant for the classification and brings more information to a summary. Finally, we have demonstrated attention scores for sentences using a highlighting tool (see Sec. 3.4) to inspect the output results. Figure 1 shows the model along with attention flow.

Token Representation: A sentence flows downstream as a sequence of tokens accompanied by two special tokens. An input representation for any token is formed by combining token, position, and segment embedding. [CLS] is the first token that classifies a sentence and appears in the beginning. [SEP] is a separator token used to identify the end of the stream. At the output [CLS], representation can be fed to the classifier for various tasks. We have only used the corresponding attention head for the task of summarization.

Pre-processing: Clinical documents contain many irregular abbreviations and periods for their particular formatting. Some fragments of notes are in grammatical order, whereas other parts are written as review keywords. We have used the custom tokenizer presented in [35] to formulate data in a listed manner. It removes tokens that contain no alphabetic character and are used as a percentage of drug prescriptions.

Fine-Tuning: Fine-tuning has a huge effect on performance for supervised tasks [25]. We have fine-tuned BERT model using maximum sequence length, batch size 8, learning rate 3e-5, ADAM

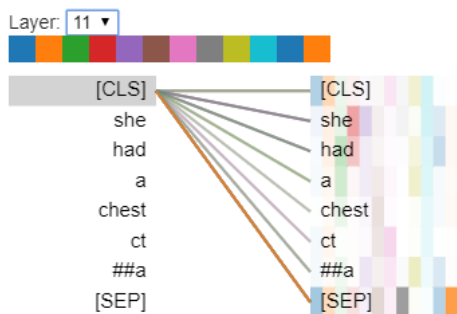


Figure 2: Attention Visualization shows how every attention-head in BERT architecture finds word useful to other words in a sentence. We have a different color for every head identifying its position in the last layer. This demonstration is performed using BertViz tool [49].

optimizer with epsilon 1e-8, and keeping other hyper-parameters the same as that of pre-training. The fine-tuned model can be used to classify [CLS] token for maximum-likelihood of ICD-9 label.

Attention Extraction: Fine-tuning helps to encode semantic knowledge in self-attention patterns [25]. Multi-head attention mechanism embeds every sentence of clinical note with a special token and feeds it to the first layer. Since the last layer of BERT model is considered vital to task-specific [45], we capture the first attention head of the last layer that is important for cross-sentence relation as observed with BertViz. This attention head focuses on a special [CLS] token. Attention scores from the first head of the last layer correspond to the whole sentence and are used as a measure of significance in a sentence. Equation 1 and Equation 2 show how the dot-product attention is calculated in layers.

$$a_i = \text{Softmax}(f(Q, K_i)) = \frac{\exp(f(Q, K_i))}{\sum_i \exp(f(Q, K_i))} \quad (1)$$

$$\text{Attention}(Q, K, V) = \sum_i a_i * V_i \quad (2)$$

A set of pre-processed clinical notes in the form of a list of sentences are fed to the BERT encoder, which creates embedding and attention scores at each layer. The attention score corresponding to [CLS] decides whether a sentence is a good candidate for the summary. Figure 2 illustrates the positional relevance of the [CLS] token. We have later selected sentences whose attention scores of the constituting words are above the average attention value of the original note. This extraction incentives dynamic selection, unlike a fixed sentence summary ratio. For example, a sentence with an attention score of 0.14 is chosen if it is greater than the average of attentions of all sentences in the document.

3.4 Sentence Attention Visualization

The attention distribution for tokens on the last layer has an irregular pattern. In order to perform a heat-mapping, we have utilized a tool namely Neat-Vision.³

³<https://github.com/cbaziotis/neat-vision>

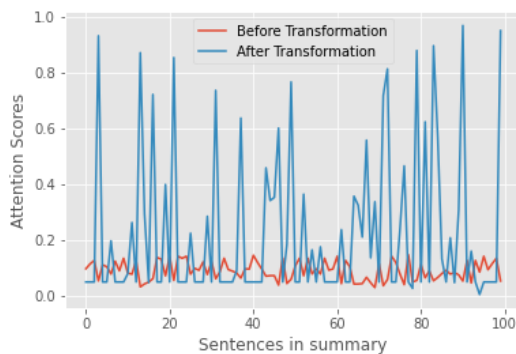


Figure 3: Effect of quantile transform on attention scores of all sentence in the document. The red line shows obtained attention distribution with a low variance where as the blue line shows higher variance that can be efficiently utilized for heat-mapping.

This tool requires a fixed type of input format and outputs a text heat-map. It demands input data be organized in a particular structure for vibrant coloring. We have stratified the distribution obtained from the neural architecture to the Gaussian distribution using the *Quantile Transformation* [18]. This turns a sentence with great attention more fragrant in visibility and vice versa. In other words, it makes sentences with higher attention scores rosier than the other ones. Figure 3 shows the impact of transformation on attention series data. Here x-axis presents the number of sentences, and the y-axis accounts for the value of the attention score for the corresponding sentence. This demonstration will have a huge impact on clinician practice by alienating time spent while reading long health records. Figure 4 exhibits the usefulness of heat-mapping concepts for health systems.

4 EVALUATION

Evaluation in summarization has been a critical issue, mainly due to the absence of a gold standard. Many competitions such as DUC⁴, TREC⁵, SUMMAC⁶ and MUC⁷ propose different metrics. The interpretation of these metrics is not so simple, mainly because the same summary receives different scores under different measures. Automatic evaluation for the quality of the summary is an ambitious task and can be performed by making a comparison with a human-generated summary. For this reason, evaluation is normally limited to domain-specific and opinion-oriented areas [41].

Formally, these evaluation methods can be divided into two areas. In *extrinsic* evaluation, it is manually analyzed how useful the summary is for the supplied document. For instance, this can be done by a clinician in our case and may result in different opinions based on his understanding. Miller et al. [34] leveraged this manual clinical evaluation to compare the performance of their model. In *intrinsic* evaluation, the extracted summary is directly matched with the ideal summary created by humans. The latter can be divided

⁴<http://duc.nist.gov/>

⁵<http://trec.nist.gov/>

⁶https://www-nlpir.nist.gov/related_projects/tipster_summac/

⁷<http://www.itl.nist.gov/iad/894.02/relatedprojects/muc/proceedings/muc7toc.html>

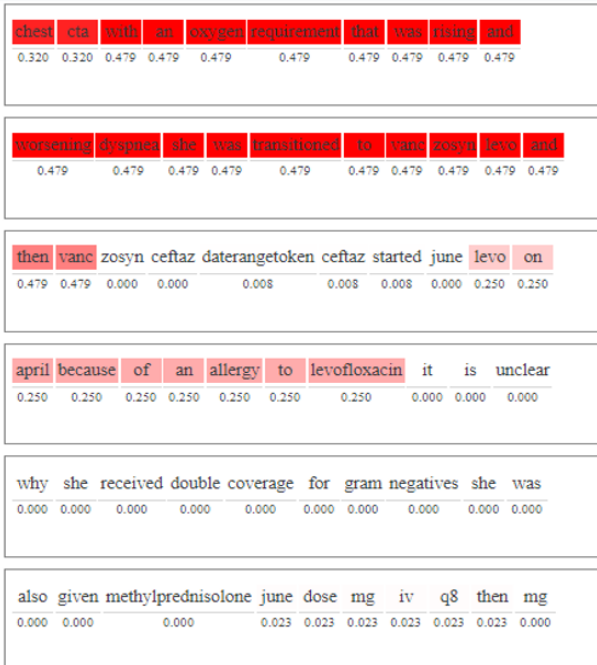


Figure 4: An excerpt of heat-mapping on Neat-Vision tool with transformed attention values to highlight importance of sentence with red color.

into two classes, primarily because it is hard to establish an ideal reference summary by a human.

- *Text quality Evaluation:* It is more related to linguistic check that examines grammatical and referential clarity. This assessment is not complete predominately owing to the fact that medical summaries are unstructured documents with lots of abbreviations and clinical jargon.
- *Content-based Evaluation:* It rates summaries based on provided reference summaries [43]. Some of the common approaches are ROGUE (Recall-Oriented Understudy for Gisting Evaluation) [28], Cosine Similarity and Pyramid Method [36]. Liu et al. [31] studied unigram and bigram ROGUE overlaps for different components of BERTSUM for their single-document summaries.

Sripada et al. [42] in their work presented that a summary can be considered effective if it has a similar probability distribution as that of the original document. The hypothesis was compared in other research works [37] [52] where this light-weight and less complex method demonstrated finer results. This criterion is a more suitable evaluation for our methodology since we do not have the reference summaries in MIMIC dataset. We will compare the distribution of words in the original and summary document to identify their effectiveness. We will use two tests for evaluating the goodness of our synopsis, namely Kullback–Leibler divergence (KLD) [30] and Jensen–Shannon divergence (JSD) [38].

KL Divergence: It is a measure of the difference between two distributions. This measure is asymmetric, and the minimum KLD

value shows better relative interference for distributions; for discrete probability distributions P and Q mapped on probability space ζ , KLD from Q to P is defined in Equation 3 [14];

$$KLD(P||Q) = \sum_{x \in \zeta} P(x) \log \frac{P(x)}{Q(x)} \quad (3)$$

JSD Divergence: It is an extension of KL divergence that quantifies the difference in a slightly modified way. It is a smoothed and normalized form assuring symmetry among inputs as reported in Equation 4 [38].

$$JSD(P||Q) = \frac{1}{2}KLD(P||M) + \frac{1}{2}KLD(Q||M). \quad (4)$$

where,

$$M = \frac{1}{2}(P + Q) \quad (5)$$

5 RESULTS

We have performed a comparison of our model with three baselines for extractive summarization methods. First, Part-of-Speech-based sentence tagging is designed on the empirical frequency selection method. The second one centers around the graphical method, and the third one uses BERT combined with K-means to find top k sentences. The results can be reproduced, and source code for replication is available on our Github repository.

The proposed architecture shows significant improvement compared with baseline approaches. Divergence scores significantly show how estimating differences in distributions can help in anticipating the word distribution of both documents. Table 1 exhibits that our extracted summaries are more informative than others based on lower average of KL and JS divergence scores. Frequency-based approach outcomes highest divergence among others. JSD and KLD scores for the graph-based method show a relative amelioration compared to frequency-based methods. There is a little quantitative difference of values with the centroid-based approach due to their nature of calculating sentence embeddings in a similar way. Our method is dynamic in choosing the length of the summary, which overcomes the weakness of fixed K sentences described in the paper [34]. Overall, the attention mechanism poses great abstraction power for the summarization task.

Models	KLD↓	JSD↓
Frequency-Based Approach	0.892	0.426
Graph-Based Approach	0.827	0.408
Centroid-Based K-means Approach	0.80	0.41
Our Proposed Architecture	0.795	0.405

Table 1: Experimental results on a reduced sample-set of 100 random clinical notes from MIMIC-III dataset compared with Frequency-Based Approach [16], Graph-Based Approach [51] and Centroid based K-means Approach [34] using KLD and JSD Values. A lower value pertains to a better-correlated summary.

Centroid-Based K-means Summary: daily disp tablet delayed release e.c. lastnametoken on february at 15pm cardiologist dr. lastnametoken on february at 30am wound check on thurs january at am with cardiac surgery on hospitaltoken please call to schedule appointments with your primary care dr. lastnametoken in march weeks please call cardiac surgery office with any questions or concerns telephonenumber token answering service will contact on call person during off hours completed by january.

Frequency-Based Summary: disp tablet refills ranitidine hcl mg tablet sig one tablet daily. refills tramadol tablet two tablet q6h hours as needed for pain. tablet senna mg tablet One tablet daily, disp tablet refills furosemide mg tablet for Mitral valve repair coronary artery bypass. graft x left internal mammary artery to left anterior descending history of present illness year old female who was told she had mvp since age currently quite active but has noticed some dyspnea on exertion when walking up hills most recent echo revealed severe mvp and moderate to severe Daily daily disp tablet delayed release e.c. s refills docusate sodium mg capsule sig one capsule po bid times Disp tablet er particles crystals s refills discharge disposition home with service facility hospitaltoken vna discharge diagnosis mitral regurgitation coronary artery disease.

Graph-Based Summary: refills docusate sodium mg capsule one capsule a day magnesium hydroxide suspension thirty ml at bedtime as needed for constipation atorvastatin tablet one tablet daily. disp tablet s refills furosemide tablet once a day for days disp tablet refills ranitidine hcl mg tablet daily. please shower daily including washing incisions gently with mild soap no baths or swimming until cleared by surgeon look at your incisions daily for redness or drainage. please no lotions cream powder or ointments to incisions each morning you should weigh yourself and then in the evening take your temperature these should be written down on the chart no driving for approximately one month and while taking narcotics will be discussed at follow up appointment with surgeon when you will be able to drive no lifting more than pounds for weeks please call with any questions or concerns telephonenumber token females please wear bra to reduce pulling on incision avoid rubbing on lower edge. please call cardiac surgery office with any questions or concerns telephonenumber token answering service will contact on call person during off hours followup instructions you are scheduled for the following appointments surgeon dr. lastnametoken on february at 15pm cardiologist dr. lastnametoken on february at 30am wound check on thurs january at am with cardiac surgery on hospitaltoken.

Our Proposed Approach: old female who was told she had mvp currently quite active but has noticed some dyspnea on exertion when walking up hills. she presents for surgical consultation past medical history mitral regurgitation copd secondary to asbestos exposure as a child arthritis cataracts headaches lactose intolerance r wrist and elbow surgery. widowed occupation retired disabled nurse tobacco quit smoking in father died suddenly at cause unknown physical exam. no spontaneous echo contrast is seen in the left atrial appendage there is a small pfo with left to right flow overall left ventricular systolic function is normal lvef in the face of mr there is normal free wall contractility there are simple atheroma in the descending thoracic aorta the aortic valve leaflets are mildly thickened trace aortic regurgitation is seen the posterior leaflet is very degenerate and there is moderate to severe mitral regurgitation. there is no pericardial effusion the tip of the sgc is seen at the pa bifurcation post cpb the patient is av paced on no inotropes the pfo is closed normal biventricular systolic fxn there is a mitral ring prosthesis which is well seated trace mr residual mean gradient with an area of no ai aorta intact. mrs. lastnametoken was a same day admit after undergoing all pre operative work. she was tolerating a full oral diet her incisions were healing well and she was ambulating in the halls without difficulty it was felt that she was safe for discharge home at this time with vna services all appropriate follow up appointments were arranged.

Table 2: Qualitative Analysis For Different Approaches

As noted in Figure 5 and 6, there are some summaries where distributional similarity does not outperform in cases where clinical notes are shorter in length. The curve presents that attention-based extraction is more impactful than other counterparts. JS divergence metrics show less fluctuation than other the metric because of its averaging symmetry mechanism. Summaries from each method are placed in Table 2 for qualitative analysis. It can be observed that summaries generated by our proposed architecture have more coherence and make it easier to adapt clinical understanding. On the other hand, baselines approaches provide short and incoherent sentences for the selected note. Shorter summaries are more likely to lose discriminatory information and affect the degree of understanding; thus, evaluating the usefulness of a summary in terms

of sentences may not be optimal. As described in Section 4, it may be hard for a non-specialist to understand the relative usefulness of each summary. This method shows the applicative benefits of dynamic summarization in healthcare systems. Furthermore, it is more helpful for a physician to grab the essence of diagnosis via highlighting tools as displayed in figure 4.

6 CONCLUSION

The immense increase in digital text information has certainly emphasized the need for universal summarization frameworks. Abstractive summarization has been an area of research debatable for certain scenarios, e.g., medical, because of the risk of generating summaries that deliver different meanings of the original notes

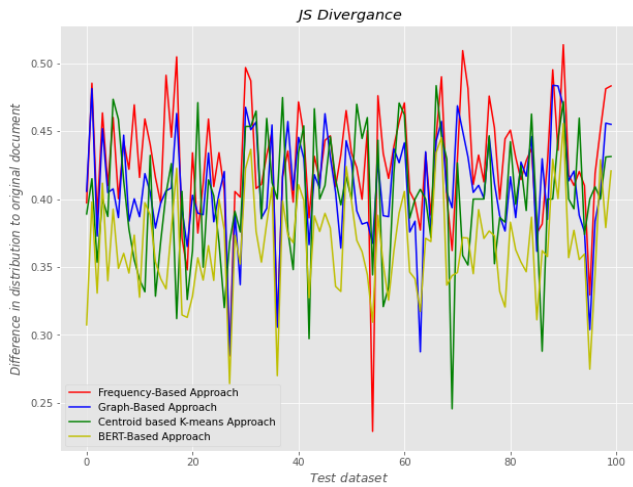


Figure 5: Line chart for JSD values for experimented models over sampled-set of clinical notes from MIMIC-III dataset.

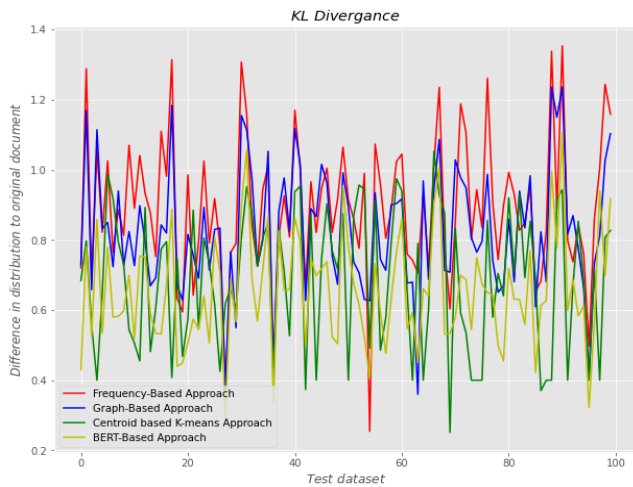


Figure 6: Line chart for KLD values for experimented models over sampled-set of clinical notes from MIMIC-III dataset.

reported by physicians. Extractive summarization techniques are relatively reasonable in the clinical domain. Research interest in creating synopsis has been re-surfed after the advent of machine learning techniques. Evaluation in medical summarization is at the toughest degree compared to other domains. We have utilized statistical analysis methods to understand the magnitude of relevance between summary and original clinical notes. The evaluation criteria of finding divergence among distributions are suitable when ideal summaries are not present.

In this paper, we have elucidated a neural architecture for extracting summaries based on multi-head attentions. The proposed model is domain-specific and outperforms other methods debated in the literature. The architecture achieves better results on a set of MIMIC-III clinical notes, outperforming frequency, graph-oriented,

and centroid-based approaches. Furthermore, our proposed model can be integrated into a decision-support system to provide a better interpretation of clinical information by highlighting diagnostically related phrases.

7 LIMITATIONS AND FUTURE WORK

Medical summarization is a special and delicate task. It is quite hard to evaluate whether the obtained summary is a general well-condensed representation of the document, and it would require manual labor and a trial with physicians for a clinical assessment. Moreover, performing a qualitative evaluation is subjective and may highly depend on the physician’s personal experience of dealing with similar diseases.

The utilized clinical attention-based model is fine-tuned on the MIMIC-III dataset. Therefore, it may not perform well on a different kinds of clinical notes, written in a different structure and mapped onto a different set of diseases. ICD-9 offers wide coverage and accurate cataloging; however, we are considering ICD-10 in current research activities. As future work, a concoction of abstractive and extractive summarization using a neural network language generation model may be more bankable. A universal medical summarizer may omit limitations arising from diverse writing style and reduce computational complexity.

REFERENCES

- [1] Mehdi Allahyari, Seyed Amin Pouriyeh, Mehdi Assefi, Saeid Safaei, Elizabeth D. Trippe, Juan B. Gutierrez, and Krys Kochut. 2017. Text Summarization Techniques: A Brief Survey. *CoRR* abs/1707.02268 (2017). arXiv:1707.02268 <http://arxiv.org/abs/1707.02268>
- [2] Emily Alsentzer and Anne Kim. 2018. Extractive Summarization of EHR Discharge Notes. *CoRR* abs/1810.12085 (2018). arXiv:1810.12085 <http://arxiv.org/abs/1810.12085>
- [3] Emily Alsentzer, John R. Murphy, Willie Boag, Wei-Hung Weng, Di Jin, Tristan Naumann, and Matthew B. A. McDermott. 2019. Publicly Available Clinical BERT Embeddings. *CoRR* abs/1904.03323 (2019). arXiv:1904.03323 <http://arxiv.org/abs/1904.03323>
- [4] T. Ayodele, R. Khusainov, and D. Ndzi. 2007. Email classification and summarization: A machine learning approach. In *2007 IET Conference on Wireless, Mobile and Sensor Networks (CCWMSN07)*. 805–808.
- [5] Mozghan Azadani, Nasser Ghadiri, and Ensieh Davoodijam. 2018. Graph-based biomedical text summarization: An itemset mining and sentence clustering approach. *Journal of Biomedical Informatics* 84 (06 2018). <https://doi.org/10.1016/j.jbi.2018.06.005>
- [6] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473* (2014).
- [7] P. B. Baxendale. 1958. Machine-Made Index for Technical Literature—An Experiment. *IBM Journal of Research and Development* 2, 4 (1958), 354–361.
- [8] H. Borko and Ch Bernier. 1978. Indexing Concepts and Methods. (01 1978).
- [9] Sergey Brin and Lawrence Page. 1998. The Anatomy of a Large-Scale Hypertextual Web Search Engine. *Comput. Netw. ISDN Syst.* 30, 1–7 (April 1998), 107–117. [https://doi.org/10.1016/S0169-7552\(98\)00110-X](https://doi.org/10.1016/S0169-7552(98)00110-X)
- [10] Rocio Chongtay, Mark Last, Mathias Verbeke, and Bettina Berendt. 2013. Summarize to learn: summarization and visualization of text for ubiquitous learning.
- [11] Janara Christensen, Mausam, Stephen Soderland, and Oren Etzioni. 2013. Towards Coherent Multi-Document Summarization. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, Atlanta, Georgia, 1163–1173. <https://www.aclweb.org/anthology/N13-1136>
- [12] Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. 2019. What Does BERT Look At? An Analysis of BERT’s Attention. *CoRR* abs/1906.04341 (2019). arXiv:1906.04341 <http://arxiv.org/abs/1906.04341>
- [13] Andrew M Dai and Quoc V Le. 2015. Semi-supervised Sequence Learning. In *Advances in Neural Information Processing Systems* 28, C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett (Eds.). Curran Associates, Inc., 3079–3087. <http://papers.nips.cc/paper/5949-semi-supervised-sequence-learning.pdf>
- [14] Dipanjan Das and André Martins. 2007. A survey on automatic text summarization. (12 2007).

- [15] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *CoRR abs/1810.04805* (2018). arXiv:1810.04805 <http://arxiv.org/abs/1810.04805>
- [16] H. P. Edmundson. 1969. New Methods in Automatic Extracting. *J. ACM* 16, 2 (April 1969), 264–285. <https://doi.org/10.1145/321510.321519>
- [17] Joshua C. Feblowitz, Adam Wright, Hardeep Singh, Lipika Samal, and Dean F. Sittig. 2011. Summarization of clinical information: A conceptual model. *Journal of Biomedical Informatics* 44, 4 (2011), 688 – 699. <https://doi.org/10.1016/j.jbi.2011.03.008>
- [18] W.G. Gilchrist. 2000. *Statistical modelling with quantile functions*. 1–325 pages.
- [19] Jade Goldstein, Mark Kantrowitz, Vibhu Mittal, and Jaime Carbonell. 1999. Summarizing Text Documents: Sentence Selection and Evaluation Metrics. In *Proceedings of the 22nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '99)*. Association for Computing Machinery, New York, NY, USA, 121–128. <https://doi.org/10.1145/312624.312665>
- [20] Vishal Gupta and Gurpreet Singh Lehal. 2010. A survey of text summarization extractive techniques. *Journal of emerging technologies in web intelligence* 2, 3 (2010), 258–268.
- [21] Jeremy Howard and Sebastian Ruder. 2018. Fine-tuned Language Models for Text Classification. *CoRR abs/1801.06146* (2018). arXiv:1801.06146 <http://arxiv.org/abs/1801.06146>
- [22] Alistair Johnson, Tom Pollard, Lu Shen, Li-wei Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Celi, and Roger Mark. 2016. MIMIC-III, a freely accessible critical care database. *Scientific Data* 3 (05 2016), 160035. <https://doi.org/10.1038/sdata.2016.35>
- [23] Thomas N. Kipf and Max Welling. 2016. Semi-Supervised Classification with Graph Convolutional Networks. *CoRR abs/1609.02907* (2016). arXiv:1609.02907 <http://arxiv.org/abs/1609.02907>
- [24] Kevin Knight and Daniel Marcu. 2000. Statistics-Based Summarization - Step One: Sentence Compression. In *Proceedings of the Seventeenth National Conference on Artificial Intelligence and Twelfth Conference on Innovative Applications of Artificial Intelligence*. AAAI Press, 703–710.
- [25] Olga Kovaleva, Alexey Romanov, Anna Rogers, and Anna Rumshisky. 2019. Revealing the Dark Secrets of BERT. 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)*. Association for Computational Linguistics, Hong Kong, China, 4365–4374. <https://doi.org/10.18653/v1/D19-1445>
- [26] Julian Kupiec, Jan Pedersen, and Francine Chen. 1995. A Trainable Document Summarizer. In *Proceedings of the 18th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '95)*. Association for Computing Machinery, New York, NY, USA, 68–73. <https://doi.org/10.1145/215206.215333>
- [27] Archana Laxmisan, Allison McCoy, Adam Wright, and Dean Sittig. 2012. Clinical Summarization Capabilities of Commercially-available and Internally-developed Electronic Health Records. *Applied clinical informatics* 3 (02 2012), 80–93. <https://doi.org/10.4338/ACI-2011-11-RA-0066>
- [28] Chin-Yew Lin. 2004. ROUGE: A Package For Automatic Evaluation Of Summaries. In *ACL 2004*.
- [29] Chin-Yew Lin and Eduard Hovy. 2000. The Automated Acquisition of Topic Signatures for Text Summarization. In *Proceedings of the 18th Conference on Computational Linguistics - Volume 1 (COLING '00)*. Association for Computational Linguistics, USA, 495–501. <https://doi.org/10.3115/990820.990892>
- [30] D. V. Lindley. 1959. *Information Theory and Statistics*. Solomon Kullback. New York: John Wiley and Sons, Inc.; London: Chapman and Hall, Ltd.; 1959. Pp. xvii, 395. \$12.50. *J. Amer. Statist. Assoc.* 54, 288 (1959), 825–827. <https://doi.org/10.1080/01621459.1959.11691207> arXiv:<https://doi.org/10.1080/01621459.1959.11691207>
- [31] Yang Liu. 2019. Fine-tune BERT for extractive summarization. *arXiv preprint arXiv:1903.10318* (2019).
- [32] H. P. Luhn. 1958. The Automatic Creation of Literature Abstracts. *IBM Journal of Research and Development* 2, 2 (1958), 159–165.
- [33] Rada Mihalcea and Paul Tarau. 2004. TextRank: Bringing Order into Text. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Barcelona, Spain, 404–411. <https://www.aclweb.org/anthology/W04-3252>
- [34] Derek Miller. 2019. Leveraging BERT for extractive text summarization on lectures. *arXiv preprint arXiv:1906.04165* (2019).
- [35] James Mullenbach, Sarah Wiegrefe, Jon Duke, Jimeng Sun, and Jacob Eisenstein. 2018. Explainable Prediction of Medical Codes from Clinical Text. *CoRR abs/1802.05695* (2018). arXiv:1802.05695 <http://arxiv.org/abs/1802.05695>
- [36] Ani Nenkova, Rebecca Passonneau, and Kathleen McKeown. 2007. The pyramid method: Incorporating human content selection variation in summarization evaluation. *ACM Transactions on Speech and Language Processing (TSLP)* 4, 2 (2007), 4–es.
- [37] Ani Nenkova, Lucy Vanderwende, and Kathleen McKeown. 2006. A compositional context sensitive multi-document summarizer: exploring the factors that influence summarization. In *Proceedings of the 29th annual international ACM SIGIR conference on Research and development in information retrieval*. 573–580.
- [38] Frank Nielsen. 2019. On a generalization of the Jensen-Shannon divergence and the JS-symmetrization of distances relying on abstract means. *CoRR abs/1904.04017* (2019). arXiv:1904.04017 <http://arxiv.org/abs/1904.04017>
- [39] Rimma Pivovarov. 2015. Electronic Health Record Summarization Over Heterogenous and Irregularly Sampled Clinical Data. *Columbia University* (2015). <https://doi.org/10.7916/D89W0F6V>
- [40] Dragomir R. Radev, Eduard Hovy, and Kathleen McKeown. 2002. Introduction to the Special Issue on Summarization. *Comput. Linguist.* 28, 4 (Dec. 2002), 399–408. <https://doi.org/10.1162/089120102762671927>
- [41] Horacio Saggion and Thierry Poibeau. 2013. *Automatic Text Summarization: Past, Present and Future*. 3–21. https://doi.org/10.1007/978-3-642-28569-1_1
- [42] Sandeep Sripada and Jagadeesh Jagarlamudi. 2009. Summarization Approaches Based on Document Probability Distributions. In *Proceedings of the 23rd Pacific Asia Conference on Language, Information and Computation, Volume 2*. City University of Hong Kong, Hong Kong, 521–529. <https://www.aclweb.org/anthology/Y09-2010>
- [43] Josef Steinberger and Karel Ježek. 2012. Evaluation measures for text summarization. *Computing and Informatics* 28, 2 (2012), 251–275.
- [44] Tomek Strzalkowski and Jose Perez Carballo. 1996. Natural language information retrieval: TREC-4 report. In *Text REtrieval Conference*. 245–258.
- [45] Betty van Aken, Benjamin Winter, Alexander Löser, and Felix A Gers. 2019. How does bert answer questions? a layer-wise analysis of transformer representations. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*. 1823–1832.
- [46] Tielman T Van Vleck, Daniel M Stein, Peter D Stetson, and Stephen B Johnson. 2007. Assessing data relevance for automated generation of a clinical summary. *AMIA ... Annual Symposium proceedings: AMIA Symposium* (October 2007), 761–765. <https://europepmc.org/articles/PMC2655814>
- [47] Lucy Vanderwende, Hisami Suzuki, Chris Brockett, and Ani Nenkova. 2007. Beyond SumBasic: Task-Focused Summarization with Sentence Simplification and Lexical Expansion. *Inf. Process. Manage.* 43 (11 2007), 1606–1618. <https://doi.org/10.1016/j.ipm.2007.01.023>
- [48] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention Is All You Need. *CoRR abs/1706.03762* (2017). arXiv:1706.03762 <http://arxiv.org/abs/1706.03762>
- [49] Jesse Vig. 2019. A Multiscale Visualization of Attention in the Transformer Model. *CoRR abs/1906.05714* (2019). arXiv:1906.05714 <http://arxiv.org/abs/1906.05714>
- [50] Adam Wright, Dean F Sittig, Joan S Ash, Joshua Feblowitz, Seth Meltzer, Carmit McMullen, Ken Guappone, Jim Carpenter, Joshua Richardson, Linas Simonaitis, et al. 2011. Development and evaluation of a comprehensive clinical decision support taxonomy: comparison of front-end tools in commercial and internally developed electronic health record systems. *Journal of the American Medical Informatics Association* 18, 3 (2011), 232–242.
- [51] Michihiro Yasunaga, Rui Zhang, Kshitij Meelu, Ayush Pareek, Krishnan Srinivasan, and Dragomir R. Radev. 2017. Graph-based Neural Multi-Document Summarization. *CoRR abs/1706.06681* (2017). arXiv:1706.06681 <http://arxiv.org/abs/1706.06681>
- [52] Wen-tau Yih, Joshua Goodman, Lucy Vanderwende, and Hisami Suzuki. 2007. Multi-Document Summarization by Maximizing Informative Content-Words.. In *IJCAI*, Vol. 7. 1776–1782.
- [53] Illhoi Yoo, Xiaohua Hu, and Il-Yeol Song. 2007. A coherent graph-based semantic clustering and summarization approach for biomedical literature and a new summarization evaluation method. *BMC bioinformatics* 8 Suppl 9 (02 2007), S4. <https://doi.org/10.1186/1471-2105-8-S9-S4>
- [54] Ming Zhong, Pengfei Liu, Yiran Chen, Danqing Wang, Xipeng Qiu, and Xuanjing Huang. 2020. Extractive summarization as text matching. *arXiv preprint arXiv:2004.08795* (2020).