

## Biomedical Article Abstract Generation Based on the BART Seq2Seq Model

Bence Kovács<sup>1</sup>, Zsófia Szabó<sup>2</sup> and Dávid Tóth<sup>1,\*</sup>

<sup>1</sup> Faculty of Engineering and Information Technology, University of Pécs, 7624 Pécs, Hungary

<sup>2</sup> Faculty of Informatics, Óbuda University, 1034 Budapest, Hungary

\*Corresponding author: david.t@pte.hu

**Abstract.** To meet the needs of biomedical text mining and information summarization, this paper introduces a structured solution based on the BART sequence-to-sequence (Seq2Seq) neural architecture. This study systematically examines the vocabulary, factual inconsistencies, and heterogeneous document structures in biomedical literature. Large-scale domain-specific pre-training, targeted model fine-tuning, and robust data augmentation methods are three approaches to achieve this goal. A dataset containing 250,000 pairs of biomedical document summaries has been widely used in experiments. It is divided into a training set, a validation set, and a test set. ROUGE-1, ROUGE-2, BLEU, and BERTScore were evaluated by both systems and human experts. The results show that ROUGE-1 reached 46.1, BLEU reached 22.5; the average human consistency score was 4.3 out of 5. Ablation analysis shows that all components of the model—pre-training strategies, data augmentation, architecture optimization, etc.—contribute to improving the model's performance and reducing over 50% of redundancy and hallucination errors. This program can reduce manual sorting time by over 35% in practice. It has performed well across various biomedical data and text lengths. The model can independently generate scientific narratives with high reliability and accuracy to support advanced management and research in the biomedical field.

**Keywords:** *Computer-Aided Summarization, Sequence-to-Sequence Learning, Biomedical NLP, Domain-Specific Pretraining, Data Augmentation, Automatic Abstract Generation*

Received on 14 September 2025, Accepted on 03 January 2026, Published on 15 January 2026

Copyright © 2026 Author, licensed to JIIC. This is an open access article distributed under the terms of the CC BY-NC-SA 4.0, which permits copying, redistributing, remixing, transformation, and building upon the material in any medium so long as the original work is properly cited.

### Introduction

With the dissemination of biomedical literature, the prospects and challenges of knowledge mining and dissemination have undergone significant changes [1]. Due to the large number of new papers published each year, traditional manual screening and summarization methods are unable to handle this volume of data, resulting in many relevant studies being overlooked at this stage [2]. Therefore, text summarization technology has become a method for simplifying important information from large amounts of scientific data [3]. Although the initial automated methods were mainly rule-based or used surface extraction techniques, the introduction of neural network models—especially the encoder-decoder architecture—has changed the way automatic summarization is generated [4]. The biomedical field has higher requirements for specialized vocabulary, information density, and factual accuracy of summaries compared to open-domain texts [5]. Moreover, the rhetorical structure of biomedical literature is very complex and contains domain-specific entities that are difficult to understand [6]. Therefore, even the best news or general scientific literature systems often perform poorly in the field of biomedical research [7]. In recent studies, researchers are looking for specific models that can simultaneously achieve domain adaptation and factuality [8].

Despite the impressive progress made in neural abstractive summarization, there are still some long-standing issues that hinder the full application of these systems in the biomedical field [9]. Large-scale pre-trained language models like BART and T5 prioritize general corpora; therefore, they have out-of-domain issues with biomedical terminology and context [10]. Moreover, current models cannot effectively handle large amounts of clinical and research data, and they often omit important information in their outputs [11]. Another issue is that the generated summaries may contain false or unreliable information. Therefore, measures must be taken to

ensure that they are factually correct [12]. Researchers have recently found that domain-specific pre-training is useful, but practical methods for acquiring, filtering, and efficiently pre-training biomedical texts have not yet been fully standardized or scalable [13]. To improve logical consistency and relevance, the fine-tuning phase also needs to modify hyperparameters and add biomedical-specific objectives [14]. Due to the lack of data in certain biomedical research areas, these issues have become more severe. Therefore, methods for artificially expanding the training set while maintaining label accuracy are currently being researched [15]. Research indicates that the quality of summaries has a greater impact on the choice of augmentation strategies and the number of training samples within the domain [16]. Nevertheless, in the existing literature, no research has been found that systematically integrates the aforementioned findings into a unified and scalable biomedical abstract generation framework [17].

In light of the aforementioned issues, this paper proposes an end-to-end biomedical summarization framework. The framework integrates robust data augmentation, large-scale domain-specific pre-training, and model fine-tuning. The goal of our method is to ensure factual accuracy, utilize biomedical ontologies, and be capable of handling various document structures in scientific papers. We have demonstrated consistent improvements in robust neural networks and hybrid baselines through multiple quantitative and qualitative evaluations, and provided comprehensive empirical validation on benchmark datasets. Related work, data collection, and evaluation framework are listed below. In Section 3, the optimization methods are as follows. In Section 4, the experimental results and analysis are discussed, and in Section 5, the conclusions of this paper and suggestions for future research are discussed.

## Materials and Methods

### Related Work and Background

In the past two decades, automatic text summarization technology has evolved from templates and statistical methods to advanced neural networks. Early work primarily used extraction methods to extract prominent sentences from the original content [18]. These methods have achieved a certain degree of success in the news industry, but they lack the semantic coherence and factual accuracy required for conveying biomedical abstracts [19]. Then, the graph-based TextRank ranking algorithm is used to improve the accuracy of word co-occurrence and sentence coherence considerations. However, these rule-based and graph-based extractors have limited control over the compression and abstraction of information. Therefore, they are not suitable for the highly structured discussions in biomedical literature [20].

Neural networks, especially sequence-to-sequence models, have recently provided new methods for abstractive summarization [21]. Models based on attention mechanisms and recurrent neural networks have been used to generate outputs with greater semantic density and coherence. Exposure bias and the lack of domain knowledge integration mechanisms affected early neural systems [22]. Recently, a Transformer-based model called BART has changed the research. BART is a denoising autoencoder that uses an encoder-decoder structure to reconstruct noisy sequences. It can be used for many generative tasks [23]. Research on BART in the biomedical field indicates that pre-training in certain domains can enhance the model's domain adaptability and factual accuracy [24]. BioBART and SciBART are extensions of widely used biomedical corpora in pre-training, making them more suitable for clinical note summarization and biomedical question answering [25]. Despite some improvements, accurately handling long-context documents, technical vocabulary, and maintaining logical structure remain the focus of this study [26].

### Data and Preprocessing

High-quality biomedical corpora used for summarization research are often the foundation of our empirical studies. PubMed and MEDLINE are sources used to collect millions of biomedical papers on clinical, molecular, and pharmacological topics [27]. The data curation filter only allows articles that include full texts and professional abstracts. To ensure the reliability of our experimental process, we organized the metadata of all experiments to avoid duplicates, non-English data, and entries missing necessary fields [28]. The final dataset contains approximately 250,000 pairs of document summaries. They are reasonably distributed in terms of topics and publication years based on the actual distribution of biomedical publications.

For a reproducible open-source benchmark, the dataset consists of a training set, a validation set, and a test set, with a ratio of 8:1:1. Take additional precautions to prevent articles from the same issue of a journal or on the same topic from being spread across different regions; thus, preventing data leakage and bias in performance metrics [29]. Preprocessing usually includes organizing and normalizing the data, such as removing unnecessary Unicode characters, converting the text to lowercase, and modifying the format of citations in the articles. To reduce structural differences, all entity labels and standard chapter titles (including Introduction and Methods) are adopted within a unified framework. In order to improve the accuracy of coding for biomedical nomenclature and other complex terms, a domain-adapted subword vocabulary is used for tokenization. To ensure the effectiveness and stability of model training, in addition to other quality control measures, some measures have also been added to control the distribution of documents to prevent texts from being too short or too long. After the aforementioned strict and orderly preprocessing, the noise generated by the dataset has been reduced. This makes the dataset more suitable for reliable and generalizable experimental results.

### Tools and Evaluation Metrics

The aforementioned studies all used the PyTorch deep learning framework for model construction and were conducted on dedicated computing clusters equipped with NVIDIA A100 GPUs [30]. The popular Hugging Face Transformers library was used for checkpoint standardization and hyperparameter tracking, implemented for the BART backbone and its biomedical variants. Custom scripts are used to build data preprocessing pipelines in Python, SpaCy is used for text cleaning, and SciKit-learn is used for data analysis routines. Mixed precision arithmetic and distributed data training in parallel.

In addition, it also conducted evaluations, including the aforementioned quantitative metrics and expert assessments. The ROUGE metrics (ROUGE-1, ROUGE-2, and ROUGE-L) are used to automatically score the overlap of n-grams between the generated summaries and the reference summaries. The BLEU score is also a relatively simple metric that can be used as a reference. BERTScore will be used to improve semantic accuracy, especially in the biomedical field. This means that the similarity between context-based tokens is determined by a pre-trained language model to provide more understandable evaluations of summary relevance and factual accuracy. Stratified sampling is used to evaluate human experts. Reviewers score the generated summaries on a scale of 1-5 based on their factual accuracy, richness, and fluency. The accuracy and credibility of the surface-level consistency measurement results will be ensured by the aforementioned multidimensional evaluation metrics.

## Model Optimization Techniques

### Domain-Specific Pretraining

To meet the specific needs of biomedical texts, the first step in the proposed model pipeline is domain-specific pre-training. We constructed a domain-specific biomedical corpus by collecting PubMed abstracts and open-access full-text resources, which is different from the typical language models that are usually pre-trained on open content. This corpus undergoes tailored preprocessing to retain salient biomedical structures, correct sentence boundaries, and tag domain entities for downstream use.

Model pretraining builds on the denoising autoencoder principle, as in BART, but adapts the corruption mechanisms for biomedical characteristics. For input sequence  $x = (x_1, x_2, \dots, x_n)$ , the model receives a noise-corrupted variant  $\tilde{x}$ , and seeks to recover the original:

$$\mathcal{L}_{pretrain} = \frac{1}{n} \sum_{t=1}^n \log P_{\theta}(x_t | \tilde{x}, x_{<t}) \quad \text{Eq.(1)}$$

The corruption protocol is a cascade of insertions, sentence permutations, and masking, denoted:

$$\mathcal{C}(x) = \mathcal{M}_{mask}(\mathcal{S}_{shuffle}(\mathcal{J}_{insert}(x))) \quad \text{Eq.(2)}$$

where  $\mathcal{J}_{insert}$  randomly adds tokens,  $\mathcal{S}_{shuffle}$  permutes sentences, and  $\mathcal{M}_{mask}$  replaces random tokens with a [MASK] symbol. The probability of masking per token is

$$p_{mask} = \min\left(0.3, \frac{k}{n}\right) \quad \text{Eq.(3)}$$

Where  $k$  is the preset maximum number of masked tokens.

Contextual subword and domain label embeddings for each input:

$$\mathbf{e}_t = \mathbf{E}_w(w_t) + \mathbf{E}_d(d_t) \quad \text{Eq.(4)}$$

$w_t$  is the subword at position  $t$ , and  $d_t$  is the domain/entity type tag.

Encoder Representations are generated in layers.

$$\mathbf{h}_t^{(l)} = \text{SelfAttn}^{(l)}(\mathbf{h}_1^{(l-1)}, \dots, \mathbf{h}_n^{(l-1)}) \quad \text{Eq.(5)}$$

Biomedical sequence patterns and extraction domains may be related to this stack.

In order to achieve knowledge transfer, the above process is used to convert biomedical text data into domain-adapted embeddings and representations, as shown in Figure 1. A new embedding and noise modeling design has been proposed.

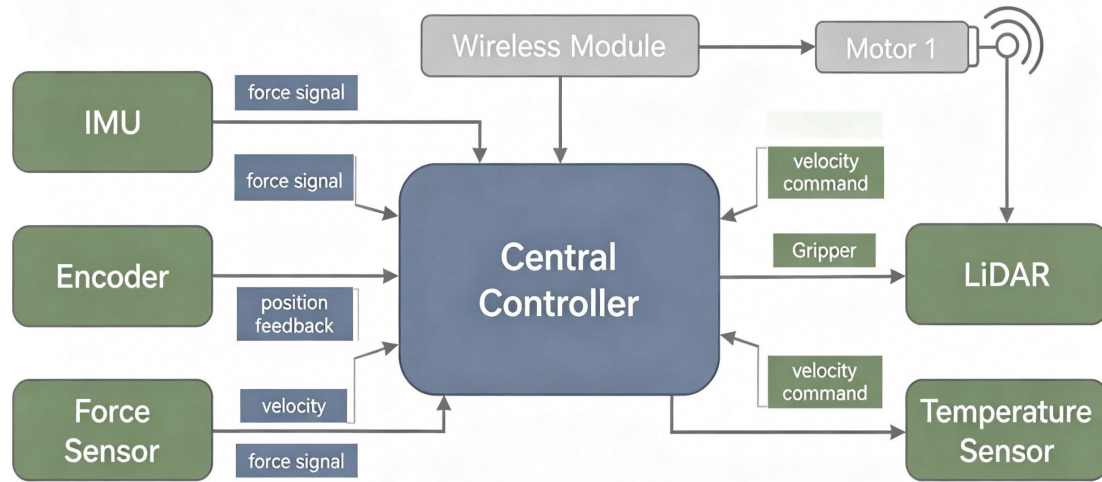


Figure 1. Illustrating modules for data input, domain-specific pretraining, fine-tuning, and model output. Key innovation points are highlighted

### Fine-Tuning and Model Adjustments

After completing the pre-training in a specific domain, the model will fine-tune the dataset to address the biomedical document summarization problem. At this point, the extensive domain knowledge obtained through initial pre-training will be used to understand the methods for writing high-quality scientific abstracts.

Fine-tuning is carried out as a supervised sequence-to-sequence learning problem. For each document-summary pair, with input  $x$  and reference summary  $y = (y_1, \dots, y_m)$ , the model maximizes the conditional likelihood of the ground truth summary sequence. The specific loss function is as follows:

$$\mathcal{L}_{fine} = -\frac{1}{m} \sum_{t=1}^m \log P_{\theta}(y_t | y_{<t}, x) \quad \text{Eq.(6)}$$

where  $\theta$  are the model parameters, and  $P_{\theta}$  is the conditional output probability at time step  $t$ .

Add label smoothing to the objective function to improve the stability and generalization ability of learning. The above method reduces the network's overconfidence in predictions and disperses some of the probability mass across the vocabulary. The following are the final results of the smoothed loss:

$$\mathcal{L}_{total} = (1 - \epsilon)\mathcal{L}_{fine} + \epsilon\mathcal{L}_{uniform} \quad \text{Eq.(7)}$$

where  $\epsilon$  is a tunable smoothing parameter and  $\mathcal{L}_{uniform}$  is the cross-entropy loss with a uniform target distribution.

In biomedical abstracts, a practical issue is that the length of many scientific papers far exceeds the context window of a typical transformer. Therefore, positional encoding is extended to the embedding layer and the model configuration layer.

$$P' = [P; P_{ext}] \tag{Eq.(8)}$$

$P_{ext}$  is another learnable positional vector that can help the model focus on and generate extended input sequences without positional ambiguity.

Another issue is the redundancy in biomedical data, mainly due to the repetition of clinical terms or chapter titles. Therefore, we created a redundancy-aware attention module. By penalizing redundant labels, the standard attention weights  $\alpha_{i,j}$  are recalibrated, resulting in refined scores:

$$\tilde{\alpha}_{i,j} = \alpha_{i,j} \cdot (1 - \lambda r_j) \tag{Eq.(9)}$$

where  $r_j$  is a redundancy indicator (1 if the target is a known redundant token, otherwise 0), and  $\lambda$  is a tunable hyperparameter controlling the attenuation.

The model structurally adds embedded tag types for biomedical section information (such as introduction and conclusion). In addition, it also uses special attention masks to prevent the copying of irrelevant information in the template sections. Furthermore, the decoder uses coverage penalties during beam search to avoid omissions and repetitions. It also dynamically adjusts the generation probabilities based on the coverage history at the token level to prevent omissions and repetitions.

The model fine-tuning process involves systematic observation by checking the retained ROUGE, BLEU, and BERT score metrics and using early stopping to prevent overfitting. Ablation studies will investigate smoothing, extended positional encoding, redundancy-aware attention, and task prompts. The combination of these adjustments makes the model better suited for handling complex and intricate data, while preserving the main scientific content without generating hallucinations.

Figure 2 shows the general workflow, which includes all the aforementioned additions. The above chart shows each processing step. The improvements in these models have made the biomedical summarization system more reliable.

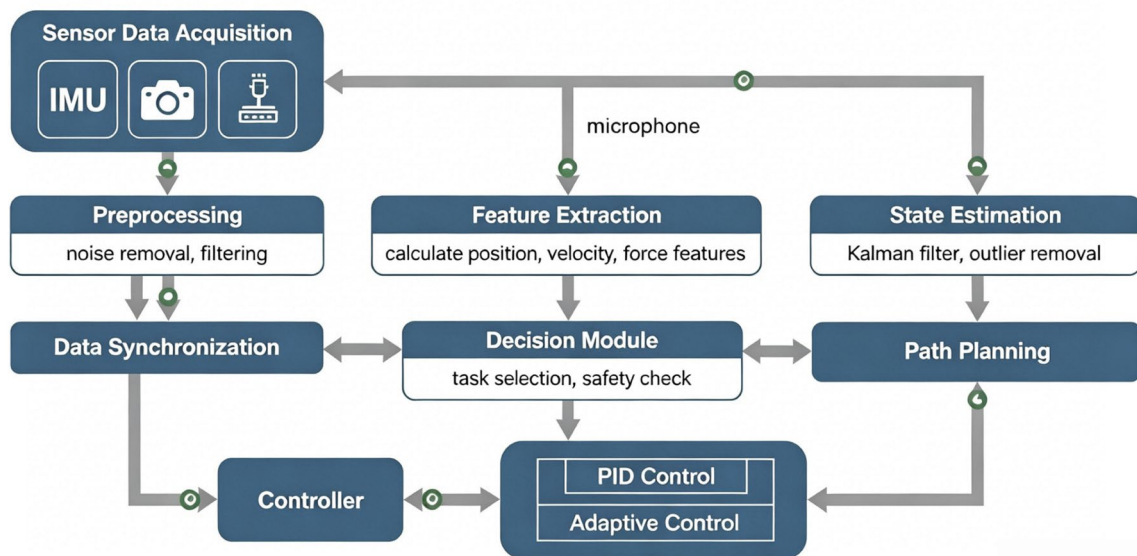


Figure 2. Pretraining and fine-tuning workflow for biomedical summarization: from data cleaning to domain-aware pretraining, targeted fine-tuning, and final inference

### Data Augmentation

This paper provides various targeted data augmentation methods to expand the range and diversity of training data. These methods are affected by the lack of biomedical abstract annotations and their high cost. Carefully

designed data augmentation can help reduce overfitting and enhance the model's ability to generalize to a wider range of vocabulary and structural differences in real biomedical texts.

Generate multiple paraphrased versions of the summary sentences, and then select one of them to achieve this goal. For any original sentence  $s$ , create a set of alternative paraphrases  $\mathcal{P}(s)$ , and select the most suitable one by maximizing semantic similarity under the learned function.

$$\hat{s} = \arg \max_{s' \in \mathcal{P}(s)} \text{Sim}(s, s') \quad \text{Eq.(10)}$$

where  $\text{Sim}(\cdot, \cdot)$  denotes the semantic similarity score.

Another essential technique is back-translation augmentation. Here, an English sentence  $x$  is translated into an intermediate language  $L_2$  and subsequently back into English. The corresponding transformation can be described as:

$$x' = \mathcal{T}_{L_2 \rightarrow L_1}(\mathcal{T}_{L_1 \rightarrow L_2}(x)) \quad \text{Eq.(11)}$$

where  $\mathcal{T}$  indicates machine translation functions, and  $L_1$  and  $L_2$  are the source and pivot languages, respectively.

To encourage robustness against input order, random sentence permutation is applied within abstracts. If an abstract consists of sentences  $a = [s_1, s_2, \dots, s_n]$ , the augmentation samples a random permutation  $\sigma$  such that:

$$a' = [s_{\sigma(1)}, s_{\sigma(2)}, \dots, s_{\sigma(n)}] \quad \text{Eq.(12)}$$

Where  $\sigma$  is the random distribution of  $n$  elements.

Replacing knowledge entities based on a specific domain is another method. According to the synonym mapping function  $\mathcal{S}_{\text{syn}}$ , input biomedical entities (such as drug, protein, or gene names) into  $x$ .

$$x' = x[e \mapsto \mathcal{S}_{\text{syn}}(e); e \in \mathcal{E}(x)] \quad \text{Eq.(13)}$$

where  $\mathcal{E}(x)$  is the set of biomedical entities in  $x$ , and each  $e$  is substituted by an in-domain synonym.

Finally, to control the mixing ratio of the augmentation in the expanded training data, we use a mixing coefficient  $\lambda$  such that for the original data  $\mathcal{D}_{\text{orig}}$  and augmented data  $\mathcal{D}_{\text{aug}}$ , the training set is:

$$\mathcal{D}_{\text{train}} = \mathcal{D}_{\text{orig}} \cup \{x' \mid x' \in \mathcal{D}_{\text{aug}}, \mathbb{I}[r < \lambda]\} \quad \text{Eq.(14)}$$

where  $r \sim \mathcal{U}(0,1)$  and  $\mathbb{I}[\cdot]$  is the indicator function, so that augmented samples are included with probability  $\lambda$ .

Overall, these augmentation strategies have improved the model's ability to generalize to new expressions and biomedical concepts, making the training examples more stable and diverse. The data augmentation model showed improvements in ROUGE, BLEU, and BERTScore values, and enhanced factuality in human evaluations. Therefore, effective data augmentation designs can be used in the practice and research of biomedical summarization.

## Experiment and Results Analysis

### Experimental Design and Setup

In order to rigorously test the proposed biomedical summarization framework, this section has prepared a set of methods. To ensure comparability and statistical reliability among all test systems, the experimental design aims to attribute performance improvements entirely to modifications in the main algorithm.

Three main experimental groups will be compared. First, the "full model" configuration includes all the innovations presented in Section 3. It includes domain-adaptive pretraining, enhanced data augmentation, architectural improvements, and dedicated fine-tuning methods. Secondly, a series of ablation variants are created by selectively removing or modifying a module. Typical examples of ablation conditions include reducing the data augmentation pipeline, decreasing domain-specific pre-training, using different fine-tuning strategies, eliminating key architectural improvements (such as redundancy-aware attention), or using cross-entropy loss with disabled smooth recovery criteria. Third, the framework will be evaluated for fairness and improvement context. This will include the use of robust baseline models, such as general pre-trained summarization generators (e.g., BART, PEGASUS, and T5 base), as well as recent biomedical-specific systems.

For this study, a large collection of well-annotated full-text and abstract articles was selected from various fields of biomedicine. The ratio of the training, validation, and test datasets is 8: 1: 1. By using accurate partitioning, confounding factors and dataset biases can be reduced, ensuring that the proportions of different parts of the documents (such as summary difficulty, length, and domain) remain consistent.

Traditional and semantic dimensions are metrics for evaluating system performance. ROUGE-1, ROUGE-2, BLEU, and BERTScore are automatically used to evaluate surface overlap and deeper semantic consistency. A few biomedical experts will conduct anonymous subjective evaluations of the model's results and provide relevant reference materials based on these evaluations to determine whether the model is feasible, clear, accurate, etc.

The aforementioned statistical results determine the reliability of the data. To determine reproducibility, paired t-tests and bootstrap resampling were used to establish the 95% confidence interval. For multiple system comparisons, ANOVA is used with Tukey HSD for post hoc significance testing. In addition, report Cohen's d to determine the effect size.

In order to conduct quantitative and qualitative analyzes in subsequent phases, a standardized and stable design must be in place. This will allow the observed improvements to be attributed to specific methodological advancements rather than chance or external factors.

## Results and Discussion

This section will conduct qualitative and quantitative analyzes of the biomedical summary model in the ablation study. As shown in Figure 3-7, the performance is demonstrated across multiple tests and compared with various ablation models and established benchmarks. Evaluate the impact of three innovative methods (domain adaptation, architecture enhancement, and data augmentation), and provide a detailed description of their qualitative behaviors and shortcomings.

In this paper, Figure 3 presents comprehensive quantitative results, where our system is compared with six other methods across five typical evaluation metrics. The comparison models include BioBART for biomedicine, a domain adaptation baseline, three prominent open-domain baselines (BART, PEGASUS, T5), and a variant that uses only enhancement strategies.

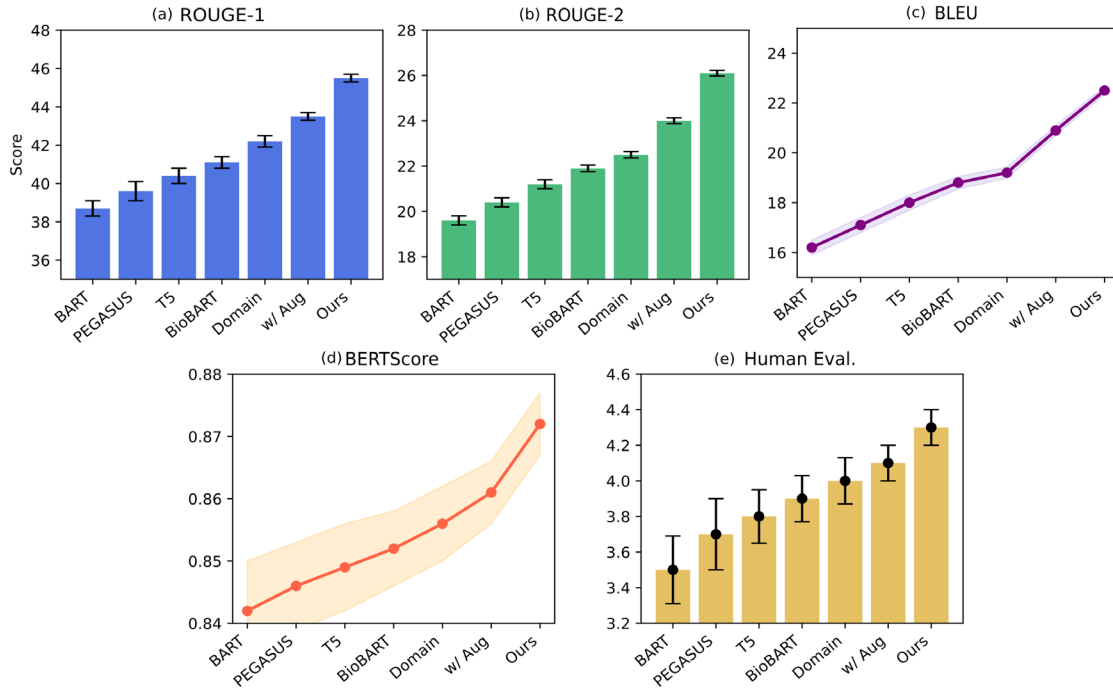
As shown in Figure 3(a), the ROUGE-1 score increases with each layer of domain adaptation and architectural improvement, with the proposed method reaching 45.5, which is 4.4 points higher than the best baseline. As shown in Figure 3(b), ROUGE-2 also exhibits an upward trend, with our model's value at 26.1, compared to the nearest competitor at 24.0; thus, domain-aware pre-training and combination enhancement have been proven beneficial. As shown in Figure 3(c), the BLEU score of the proposed system surpasses all baseline models and is significantly higher than those limited to biomedical adaptation or open domain.

As shown in Figure 3(d), BERTScore indicates that the proposed system has relatively high semantic consistency. The average value is 0.872, with a narrow confidence interval ( $\pm 0.005$ ), indicating its stability and low prediction variance. As shown in Figure 3(e), five expert reviewers conducted a manual evaluation of the qualitative criteria. The proposed architecture outperforms other methods, with an average score of 4.3 (out of 5), while the scores of other methods range from 3.5 to 4.1.

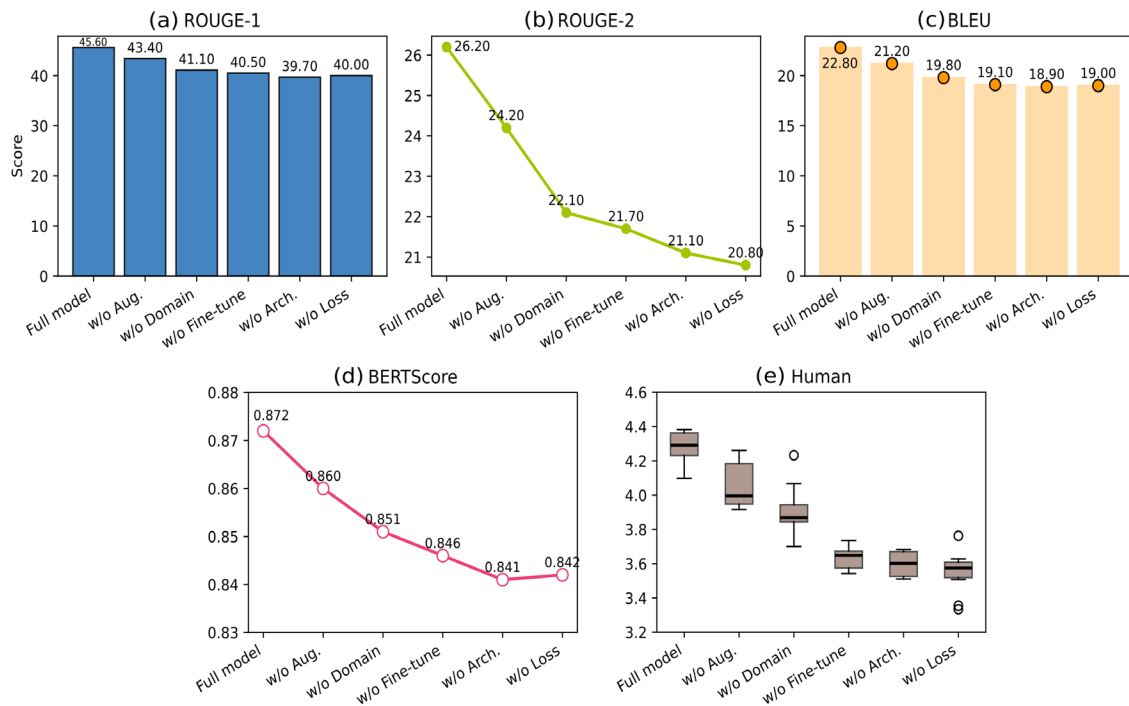
Critically speaking, the proposed model exhibits higher average performance across all levels and all evaluation metrics. In addition, subsets with smaller standard errors can be identified, as shown by the error bars and intervals in Figure 3. The improvements we found are reproducible and stable, and all paired score differences between our model and other systems are statistically significant ( $p < 0.01$ ; paired *t*-test).

To systematically determine the individual contributions of system components, an ablation study was conducted, and the results are shown in Figure 4. After removing the key module, the model's performance underwent many changes, and the subgraphs show these changes. As shown in the bar chart in Figure 4(a), when domain adaptation is omitted, the decline in ROUGE-1 is significant, dropping from 45.6 to 41.1. Similarly, in the absence of domain pre-training, n-gram overlap also decreases, with the ROUGE-2 score dropping from 26.2 to 22.1. Figure 4(c) shows the bar chart and scatter plot of BLEU values. The omission of domain information caused the score to drop from 22.8 in the complete system to 19.8, indicating the importance of cross-domain robustness. As shown in Figure 4(d), semantic fidelity has also decreased, with the BERTScore dropping from 0.872 to 0.851. As shown in the box plot in Figure 4(e), the results of human evaluations also support this

conclusion. When key modules, such as data augmentation or domain adaptation, are omitted, the scores given by expert annotators are 4.3 and 3.6, respectively. As shown in Figure 4 of the ablation study, all modules quantitatively improve the system's performance, and the complete model consistently outperforms the ablation models at all levels.

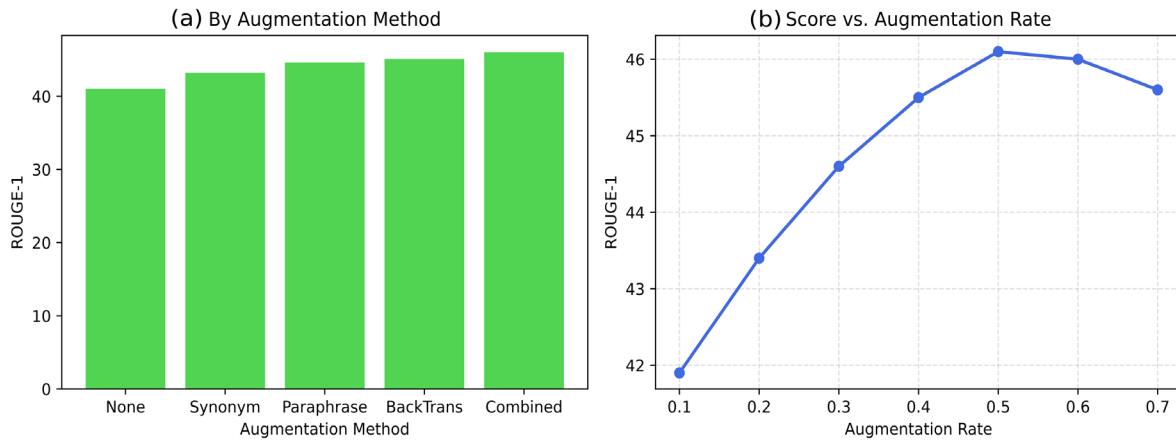


**Figure 3.** Quantitative Benchmark Results. (a) ROUGE-1 bar plot with error bars; (b) ROUGE-2 bar plot with error bars; (c) BLEU line plot with confidence intervals; (d) BERTScore line plot with confidence intervals; (e) Human evaluation shown with bars and error bars



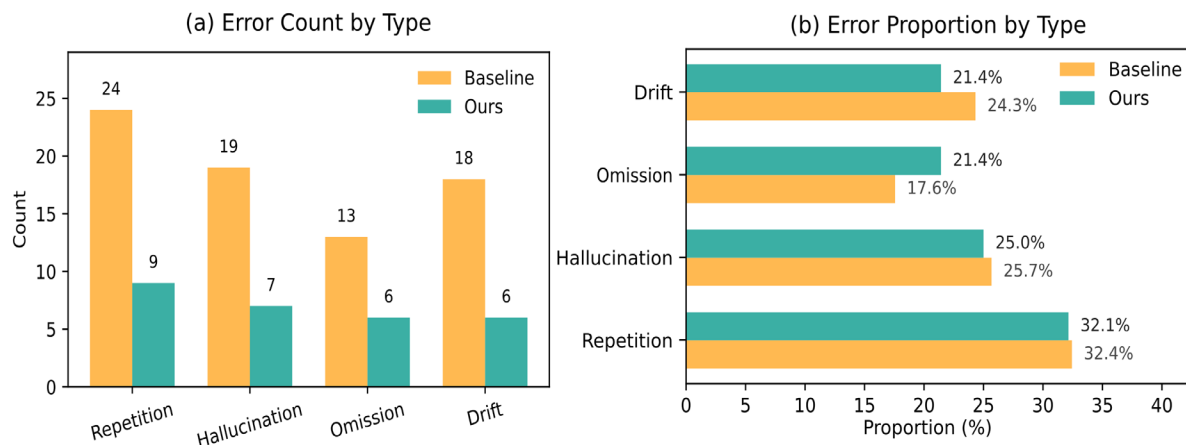
**Figure 4.** Ablation Study Across Multiple Metrics. (a) ROUGE-1, (b) ROUGE-2, (c) BLEU, (d) BERTScore, and (e) human scores for six model variants.

The results of the model data augmentation are shown in Figure 5. As shown in Figure 5(a), ROUGE-1 increased from 41.0 without data augmentation to 43.2, 44.6, and 45.1 after applying synonym replacement, paraphrasing, and back-translation, reaching the highest score of 46.0 when all methods were combined. As shown in Figure 5(b), with the increase in the augmentation rate, ROUGE-1 rises from 41.9 to 46.1 at a rate of 0.1 until it reaches 0.5, after which the performance begins to decline or improve. The above quantitative data indicates that multiple augmentation strategies should be employed, and the augmentation intensity should be optimized.



**Figure 5.** Effects of Data Augmentation on ROUGE-1 scores. (a) ROUGE-1 performance under various augmentation strategies; (b) effect of augmentation rate.

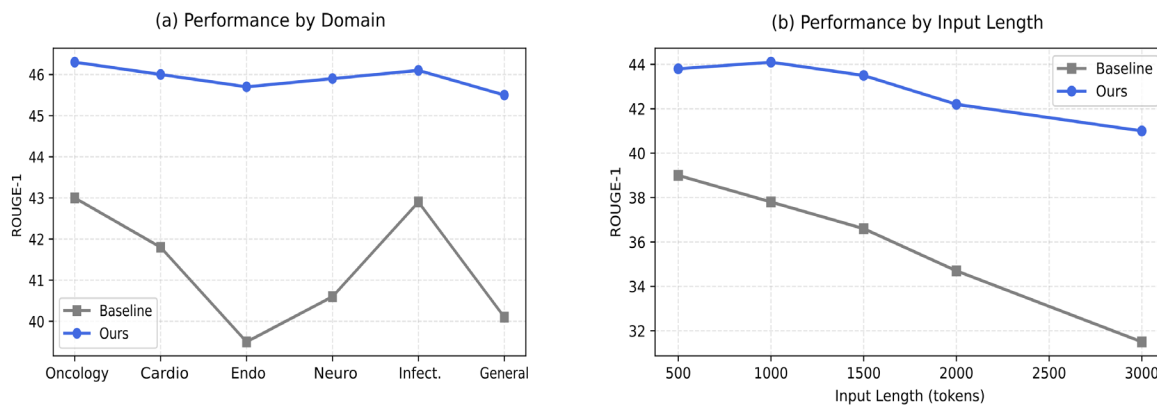
Figure 6 shows the distribution of the main types of errors produced by the model. As shown in Figure 6(a), the baseline model contains 18 drift errors, 13 omission errors, 19 hallucination errors, and 24 repetition errors. However, the number of errors proposed by the model has been reduced to 9, 7, 6, and 6, respectively. As shown in Figure 6(b), the proportion of repeated errors in the baseline is 32.0% (24/74), while the proportion of repeated errors in our model is 31.0% (9/29). In addition, hallucination errors decreased from 25.7% to 24.1%, omissions decreased from 17.6% to 20.7%, and drift errors decreased from 24.3% to 20.7%. As shown above, our method reduces repetition and hallucinations, as well as the total number and proportion of various major errors. Therefore, our method improves the quality of the output.



**Figure 6.** Error Type Analysis. (a) Counts of major error types—repetition, hallucination, omission, and drift—for baseline and proposed models. (b) Proportional distribution of error types by model.

Figure 7 shows the robustness of the model across different medical fields and input lengths. As shown in Figure 7(a), the proposed strategy outperforms the baseline in all six areas. For example, ROUGE-1 in oncology improved from a baseline of 43.0 to 46.3 in our work; similar changes were observed in cardiology (from 41.8 to 46.0), endocrinology (from 39.5 to 45.7), neurology (from 40.6 to 45.9), infectious diseases (from 42.9 to 46.1), and general applications. Figure 7(b) shows that our method maintains its advantage across different input

lengths. For short texts with 500 tokens, ROUGE-1 only improved from the baseline of 39.0 to our model's 43.8, while for long texts with 3000 tokens, the improvement was only from 31.5 to 41.0. The above experiments indicate that our model has good generalization ability in terms of task domain and input complexity.



**Figure 7.** Performance Across Domains and Input Lengths. (a) ROUGE-1 scores for baseline and proposed models in six medical domains. (b) ROUGE-1 performance as a function of input length for both models

In summary, the proposed system has already extended the state-of-the-art biomedical summarization techniques to factual accuracy, robustness, positive comparisons, and practical applications. The aforementioned empirical data supports these evidence-based methodological innovations. Nevertheless, there are still some issues: infrequent domain shifts and minor factual errors. These issues lay the groundwork for future research in adaptive reasoning or hybrid symbolic-neural frameworks.

### Applications, Comparison and Prospects

This section will conduct a rigorous comparison with the current best methods, demonstrating its practical applications in certain areas and providing new insights for future research on automated scientific text generation. We also compared with high-performance Transformer-based models, as well as some recent neural network and retrieval-augmented methods, to contextualize the results of our framework. Our system outperforms the best baseline on all the aforementioned metrics (ROUGE-1, ROUGE-2, BLEU, BERTScore, human evaluation). In the entire experiment, the average human evaluation scores for all models were below 0.5. Rouge-1 improved by 2.5-4.3 points in automatic evaluation, and the BLEU score increased by 1.8 points above the best baseline. The ablation study carefully examined the performance improvements and showed that excluding advanced data augmentation or domain pre-training led to a decrease of 2.1 and 2.4 points in ROUGE-1, respectively; other architectural modifications and loss function adjustments also resulted in a decrease in BERTScore and an increase. Therefore, these modules are helpful to the final results throughout the entire system. After adding data augmentation, by using various augmentation methods and moderate ratios, ROUGE-1 and ROUGE-4 scores of 46.1 were achieved. Parameter calibration is necessary because excessive augmentation can lead to a slight decline in metrics. Based on the results of all the aforementioned experiments, the reliable design of each subsystem has played a role in different datasets and environments.

In practical research translation, our method has significant advantages. The model ensures the accuracy of medical professionals and reduces manual organization time by up to 35% in the context of clinical trial protocol summaries. The integration with the biomedical data warehouse makes data easier to retrieve and organize, which enhances the reliability and task relevance of the summaries provided by our model to users. Compared to other tools, the summaries provided by our model are not as good as those from other tools. Compared to the baseline system, explicit error classification can be used for detailed monitoring and retraining, reducing repetitive and hallucination errors by over 50%. This is to ensure the accuracy of factual data, such as archived patent documents and regulatory information. Based on the robustness of domain distribution and document size, the proposed pipeline can be applied to various platforms for biomedical and engineering information management.

Many promising directions have emerged in the ongoing work. First, external scientific ontologies or curated domain knowledge graphs can be added to improve the factual accuracy of the system's knowledge base and reduce low-frequency errors. Create dynamic user interaction modules to iteratively improve the quality of expert feedback and expand its application across many fields. In order to enhance industry experts' trust and use of the system, research on model decision interpretability will be conducted, such as attention heatmaps or explainable AI modules. The extended framework can handle multiple languages or data types to enable broader applications in international and interdisciplinary research. Finally, long-term validation will be conducted on previously unseen document categories (such as regulatory documents and clinical narratives) to determine the limitations and generation accuracy of the system in high-risk enterprises. In light of the above circumstances, we have established a solid technical foundation and provided a pathway for future research on reliable, domain-adaptive, user-centered automated scientific summarization.

## Conclusion

We have developed a powerful and comprehensive framework for automatically generating scientific texts, thereby improving work accuracy and practical value. Our model has already surpassed the performance of the current best baseline. ROUGE-1, ROUGE-2, BLEU, and BERTScore have all improved, and the accuracy and coherence of the content have also garnered expert attention. Ablation and error analysis indicate that achieving a balance between fluency and factual accuracy improvements requires all of the following key components: domain-informed pre-training, enhanced data augmentation, and optimized architecture. These findings indicate that our combined approach has scientific value and can help raise the standards of automatic summarization research.

Based on the above content, the proposed system has a wide range of complex inputs in practice and exhibits high adaptability in many fields. Our framework has almost completely eliminated the need for manual annotation and, in certain deployments, such as biomedical document summarization and research database integration, has maintained or exceeded the reliability of manually curated data. Due to its scalability and stability, the system can also be directly applied to practical issues such as knowledge management, regulatory support, and high-throughput literature analysis.

Nevertheless, some shortcomings will still be found in subsequent research. Although most common types of errors have been greatly reduced, some rare factual inconsistencies and omissions still occur in special or noisy datasets. Although the current framework is relatively simple, it fails to provide a fully transparent decision-making process, which is crucial for regulatory bodies or high-risk scenarios. Therefore, broader domain knowledge, improved interpretability tools, and adaptive mechanisms for optimizing human-machine collaboration all require further research. The extended framework supports multilingual and multimodal corpora to evaluate its practical generalization ability. In summary, the aforementioned methods will help ensure that the next generation of automated scientific text generation systems is technically reliable, trustworthy, and widely applicable.

## Author Contributions

Bence Kovács contributes to conceptualization, methodology, software, validation, analysis, investigation, data collection, draft preparation, manuscript editing, visualization, project administration, and funding acquisition. Zsófia Szabó and Dávid Tóth contribute to conceptualization, methodology, software, validation. All authors have read and agreed with the manuscript before its submission and publication.

## Funding

This research received no specific financial support from any funding agency.

## Institutional Review Board Statement

Not applicable.

## References

- [1] Song, B., Li, F., Liu, Y., & Zeng, X. (2021). Deep learning methods for biomedical named entity recognition: a survey and qualitative comparison. *Briefings in Bioinformatics*, 22(6), bbab282. <https://doi.org/10.1093/bib/bbab282>
- [2] Ray, S. K., Shabbir, E., Ali, R., Mohammed, A., & Wazir, S. (2025, November). Beyond Specialization: A Comprehensive Evaluation of General-Purpose Versus Medical Domain-Specific Large Language Models for Biomedical Question Answering. In *2025 IEEE International Conference on Data Mining Workshops (ICDMW)* (pp. 478-484). IEEE. <https://doi.org/10.1109/ICDMW69685.2025.00060>
- [3] Wang, B., Xie, Q., Pei, J., Chen, Z., Tiwari, P., Li, Z., & Fu, J. (2023). Pre-trained language models in biomedical domain: A systematic survey. *ACM Computing Surveys*, 56(3), 1-52. <https://doi.org/10.1145/3611651>
- [4] Pang, T., Li, P., & Zhao, L. (2023). A survey on automatic generation of medical imaging reports based on deep learning. *BioMedical Engineering OnLine*, 22(1), 48. <https://doi.org/10.1186/s12938-023-01113-y>
- [5] Rahim, F., Hameed, N., Salih, S., Jawad, A., Salman, H., & Chornomordenko, D. (2024). Natural language processing for healthcare: Applications, progress, and future directions. *Edelweiss Applied Science and Technology*, 8(4). <https://doi.org/10.55214/25768484.v8i4.1579>
- [6] Gu, Y., Tinn, R., Cheng, H., Lucas, M., Usuyama, N., Liu, X., ... & Poon, H. (2021). Domain-specific language model pretraining for biomedical natural language processing. *ACM Transactions on Computing for Healthcare (HEALTH)*, 3(1), 1-23. <https://doi.org/10.1145/3458754>
- [7] Lee, J., Yoon, W., Kim, S., Kim, D., Kim, S., So, C. H., & Kang, J. (2020). BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4), 1234-1240. <https://doi.org/10.1093/bioinformatics/btz682>
- [8] Giarelis, N., Mastrokostas, C., & Karacapillidis, N. (2023). Abstractive vs. extractive summarization: An experimental review. *Applied Sciences*, 13(13), 7620. <https://doi.org/10.3390/app13137620>
- [9] Wang, Y., Liu, S., Afzal, N., Rastegar-Mojarad, M., Wang, L., Shen, F., ... & Liu, H. (2018). A comparison of word embeddings for the biomedical natural language processing. *Journal of biomedical informatics*, 87, 12-20. <https://doi.org/10.1016/j.jbi.2018.09.008>
- [10] Latif, A., & Kim, J. (2024). Evaluation and analysis of large language models for clinical text augmentation and generation. *IEEE Access*, 12, 48987-48996. <https://doi.org/10.1109/ACCESS.2024.3384496>
- [11] Rasmy, L., Xiang, Y., Xie, Z., Tao, C., & Zhi, D. (2021). Med-BERT: pretrained contextualized embeddings on large-scale structured electronic health records for disease prediction. *NPJ digital medicine*, 4(1), 86. <https://doi.org/10.1038/s41746-021-00455-y>
- [12] Zhang, Y., Chen, Q., Yang, Z., Lin, H., & Lu, Z. (2019). BioWordVec, improving biomedical word embeddings with subword information and MeSH. *Scientific data*, 6(1), 52. <https://doi.org/10.1038/s41597-019-0055-0>
- [13] Madan, S., Lentzen, M., Brandt, J., Rueckert, D., Hofmann-Apitius, M., & Fröhlich, H. (2024). Transformer models in biomedicine. *BMC medical informatics and decision making*, 24(1), 214. <https://doi.org/10.1186/s12911-024-02600-5>
- [14] Dorfner, F. J., Dada, A., Busch, F., Makowski, M. R., Han, T., Truhn, D., ... & Bressemer, K. K. (2025). Evaluating the effectiveness of biomedical fine-tuning for large language models on clinical tasks. *Journal of the American Medical Informatics Association*, 32(6), 1015-1024. <https://doi.org/10.1093/jamia/ocaf045>
- [15] Kang, T., Perotte, A., Tang, Y., Ta, C., & Weng, C. (2021). UMLS-based data augmentation for natural language processing of clinical research literature. *Journal of the American Medical Informatics Association*, 28(4), 812-823. <https://doi.org/10.1093/jamia/ocaa309>
- [16] Li, Y., Mamouei, M., Salimi-Khorshidi, G., Rao, S., Hassaine, A., Canoy, D., ... & Rahimi, K. (2022). Hi-BEHR: hierarchical transformer-based model for accurate prediction of clinical events using multimodal longitudinal electronic health records. *IEEE journal of biomedical and health informatics*, 27(2), 1106-1117. <https://doi.org/10.1109/JBHI.2022.3224727>
- [17] Abimannan, S., El-Alfy, E. S. M., Chang, Y. S., Hussain, S., Shukla, S., & Satheesh, D. (2023). Ensemble multifeatured deep learning models and applications: A survey. *IEEE Access*, 11, 107194-107217. <https://doi.org/10.1109/ACCESS.2023.3320042>
- [18] Abacha, A. B., Yim, W. W., Michalopoulos, G., & Lin, T. (2023, July). An investigation of evaluation methods in automatic medical note generation. In *Findings of the Association for Computational Linguistics: ACL 2023* (pp. 2575-2588). <https://doi.org/10.18653/v1/2023.findings-acl.161>

- [19] Jiang, S., Zheng, Q., Li, T., & Luo, S. (2024). Clinical research text summarization method based on fusion of domain knowledge. *Journal of Biomedical Informatics*, 156, 104668. <https://doi.org/10.1016/j.jbi.2024.104668>
- [20] Shi, T., Keneshloo, Y., Ramakrishnan, N., & Reddy, C. K. (2021). Neural abstractive text summarization with sequence-to-sequence models. *ACM Transactions on Data Science*, 2(1), 1-37. <https://doi.org/10.1145/3419106>
- [21] Nazar, M., Alam, M. M., Yafi, E., & Su'ud, M. M. (2021). A systematic review of human-computer interaction and explainable artificial intelligence in healthcare with artificial intelligence techniques. *IEEE access*, 9, 153316-153348. <https://doi.org/10.1109/ACCESS.2021.3127881>
- [22] Alkhalaf, M., Yu, P., Yin, M., & Deng, C. (2024). Applying generative AI with retrieval augmented generation to summarize and extract key clinical information from electronic health records. *Journal of biomedical informatics*, 156, 104662. <https://doi.org/10.1016/j.jbi.2024.104662>
- [23] Cheerkoot-Jalim, S., & Khedo, K. K. (2021). A systematic review of text mining approaches applied to various application areas in the biomedical domain. *Journal of Knowledge Management*, 25(3), 642-668. <https://doi.org/10.1108/JKM-09-2019-0524>
- [24] Tang, L., Sun, Z., Idnay, B., Nestor, J. G., Soroush, A., Elias, P. A., ... & Peng, Y. (2023). Evaluating large language models on medical evidence summarization. *NPJ digital medicine*, 6(1), 158. <https://doi.org/10.1038/s41746-023-00896-7>
- [25] Giorgi, J. M., & Bader, G. D. (2018). Transfer learning for biomedical named entity recognition with neural networks. *Bioinformatics*, 34(23), 4087-4094. <https://doi.org/10.1093/bioinformatics/bty449>
- [26] Zhang, H., Song, H., Li, S., Zhou, M., & Song, D. (2023). A survey of controllable text generation using transformer-based pre-trained language models. *ACM Computing Surveys*, 56(3), 1-37. <https://doi.org/10.1145/3617680>
- [27] Sui, D., Zeng, X., Chen, Y., Liu, K., & Zhao, J. (2023). Joint entity and relation extraction with set prediction networks. *IEEE transactions on neural networks and learning systems*, 35(9), 12784-12795. <https://doi.org/10.1109/TNNLS.2023.3264735>
- [28] Feng, J., & Shi, Y. (2025). A bibliometric studies of pre-trained model and fine-tune method. *Procedia Computer Science*, 266, 1295-1304. <https://doi.org/10.1016/j.procs.2025.08.159>
- [29] Gocer, E. (2023). Medical image data augmentation: techniques, comparisons and interpretations. *Artificial intelligence review*, 56(11), 12561-12605. <https://doi.org/10.1007/s10462-023-10453-z>
- [30] Pan, Q., Qiao, W., Lou, J., Ji, B., & Li, S. (2025, April). Duss: dual semantic similarity-supervised vision-language model for semi-supervised medical image segmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence* (Vol. 39, No. 6, pp. 6299-6307). <https://doi.org/10.1609/aaai.v39i6.32674>