

Entity-driven Fact-aware Abstractive Summarization of Biomedical Literature

Amanuel Alambo*, Tanvi Banerjee*, Krishnaprasad Thirunarayan*, Michael Raymer*

*Wright State University

{alambo.2, tanvi.banerjee, t.k.prasad, michael.raymer}@wright.edu

Abstract—As part of the large number of scientific articles being published every year, the publication rate of biomedical literature has been increasing. Consequently, there has been considerable effort to harness and summarize the massive amount of biomedical research articles. While transformer-based encoder-decoder models in a vanilla source document-to-summary setting have been extensively studied for abstractive summarization in different domains, their major limitations continue to be entity hallucination (a phenomenon where generated summaries constitute entities not related to or present in source article(s)) and factual inconsistency. This problem is exacerbated in a biomedical setting where named entities and their semantics (which can be captured through a knowledge base) constitute the essence of an article. The use of named entities and facts mined from background knowledge bases pertaining to the named entities to guide abstractive summarization has not been studied in biomedical article summarization literature. In this paper, we propose an entity-driven fact-aware framework for training end-to-end transformer-based encoder-decoder models for abstractive summarization of biomedical articles. We call the proposed approach, whose building block is a transformer-based model, EFAS, Entity-driven Fact-aware Abstractive Summarization. We conduct a set of experiments using five state-of-the-art transformer-based encoder-decoder models (two of which are specifically designed for long document summarization) and demonstrate that injecting knowledge into the training/inference phase of these models enables the models to achieve significantly better performance than the standard source document-to-summary setting in terms of entity-level factual accuracy, N-gram novelty, and semantic equivalence while performing comparably on ROUGE metrics. The proposed approach is evaluated on ICD-11-Summ-1000, a dataset we build for abstractive summarization of biomedical literature, and PubMed-50k, a segment of a large-scale benchmark dataset for abstractive summarization of biomedical literature.

Index Terms—Transformers, Named Entity Recognition, Knowledge Bases, Abstractive Summarization, ICD-11, Knowledge Retrieval, Knowledge-enhanced Natural Language Generation

I. INTRODUCTION

Neural abstractive summarization is well explored for summarization of news articles [1]–[5], and scientific articles [6]–[8]. While there are some efforts to apply neural abstractive summarization techniques for biomedical literature [6], [9], the use of named entities and background knowledge bases to guide biomedical abstractive summary generation has not been explored. On the other hand, named entity recognition/understanding in biomedical literature has been extensively studied such that named entities are known to harbor significant semantics about a biomedical article [10]–[13]. Further, linking named entities to concept definitions in

background domain-specific knowledge bases boost semantic understanding of a biomedical article by improving comprehensiveness and contextualization as investigated in [14]–[16].

The use of named entities to guide abstractive summarization of news articles has been explored in recent studies by [17]–[20]. However, these studies do not leverage named entity-driven facts from background knowledge bases to guide abstractive summary generation. Inspired by recent advances in transformer-based encoder-decoder models [21]–[24] and knowledge and retrieval augmented natural language generation [25]–[30], we propose a technique for retrieval of entity-driven facts from biomedical knowledge bases and leveraging the retrieved facts as contextual signals for abstractive summarization of biomedical literature. Our proposed framework consists of two major components: 1) entity-driven knowledge retriever; and 2) knowledge-guided abstractive summarizer. The entity-driven knowledge retriever extracts facts from UMLS [31], ICD-10 [32], and SNOMED-CT [33] based on named entities in a biomedical article while the knowledge-guided abstractive summarizer is trained to generate summaries by attending to the source article to be summarized, the chain of named entities in the source article, and the retrieved facts from the knowledge bases.

We conduct experiments on two datasets: ICD-11-Summ-1000, and PubMed-50k. For curating ICD-11-Summ-1000, we conduct entity-driven clustering of PubMed abstracts collected per ICD-11 [34] chapter followed by entity-aware content selection to build an entity-aware pseudo extractive document. The pseudo extractive document is passed as an input source article to generate the abstractive summary during inference. The contributions of this study are summarized as follows:

- We introduce an approach for named entity-driven fact retrieval from biomedical knowledge bases using dense vector representations.
- We propose a framework based on transformer-based encoder-decoder models for knowledge-guided large-scale abstractive summarization of biomedical literature.
- We conduct extensive experiments and ablation studies to assess the efficacy of the proposed approach and show that injecting named entities and entity-aware facts mined from biomedical knowledge bases into abstractive summarization models can boost their performance in terms of entity-level factual consistency [19], [35]–[37], N-gram novelty [30], [38], and semantic equivalence

measured using BERTScore [39].

- We develop ICD-11-Summ-1000, a dataset consisting of clusters of biomedical articles collected from PubMed for the *special groups* chapters in the ICD-11 catalog ¹ and the abstractive summaries generated using different variations of our experimental design. Further, we plan to share our code, and data with other researchers ².

II. RELATED WORK

While abstractive summarization is well studied for summarization of news articles with success attributed to the availability of a massive amount of training data, their applicability to scholarly articles, particularly, in the biomedical domain is limited. Further, although named entities have been extensively studied to convey the semantics of an article (news, scientific, social media) and the saliency of individual sentences [18] within an article, they have not been widely used as part of modeling abstractive summarization. [17] performed entity-aware single-document abstractive summarization using reinforcement learning for training. Their pipeline-based approach consists of an entity-aware content selection module and abstract generation module. They evaluate their approach on the CNN/Daily Mail and NYT corpora. [18] perform entity-driven multi-document abstractive summarization of news articles (WikiSum, and Multi-News) using an encoder-decoder framework augmented with Graph Attention Network (GAT). [40] proposed EntityRank, an extension of the LexRank [41] graph-based algorithm, for entity-supported summarization of biomedical abstracts.

There have been a few recent efforts towards knowledge/fact-aware abstractive summarization in different domains. [42] introduced a Fact-aware Abstractive Summarization model called FaSum for improving the factual consistency of summaries in the domain of news articles. However, their approach does not leverage named entities for fact retrieval. [20] extended a transformer-based abstractive summarization model using entities disambiguated and linked to Wikidata knowledge graph and attending to the entities for summarization of news articles. Their approach, however, does not perform named-entity based fact retrieval from the knowledge base constrained by the article to be summarized and the named entities. [43] developed an unsupervised pipeline-based approach for knowledge-infused abstractive summarization for condensing patient-to-clinician diagnostic interviews based on Multi-Sentence Compression [44] and Integer Linear Programming [45]. Nevertheless, their approach uses domain-specific lexicons as knowledge source for filtering irrelevant utterances and for retrofitting language models [46] and, does not leverage named entities or facts as part of an end-to-end training of models. [47] proposed Biomed-Summarizer, a framework for extractive summarization of biomedical literature in a multi-document setting and evaluated on PubMed abstracts. [6] built a model for abstractive summarization of long documents

using a discourse-aware encoder-decoder framework and experimented on two large scale datasets including research articles collected from PubMed. To address the challenge associated with the scarcity of large-scale training data in the biomedical domain, [48] released MS2 (Multi-Document Summarization of Medical Studies). They experimented with BART [22] for abstractive summary generation on the dataset they introduced in a traditional multi-doc-to-summary setting.

Though all the aforementioned studies conduct abstractive summarization of biomedical literature or the use of facts mined from knowledge bases for a different domain, they follow the well-established paradigm of source-document vs summary pairing during training/inference of models. Our approach is different in that we augment the state-of-the-art abstractive summarization models with additional contextual signals during training/inference and apply them to the biomedical domain.

III. DATA PREPARATION

In addition to the benchmark PubMed-50k [6] dataset which we use to train our models, we curate ICD-11-Summ-1000 dataset using a data preparation pipeline which follows two stages: 1) ICD-11 disease lexicon curation for querying PubMed for abstracts; and 2) Entity-aware pseudo-document generation for a collection of semantically related PubMed abstracts in an ICD-11 chapter. Thus, for querying PubMed for abstracts for an ICD-11 chapter, we first query the biomedical knowledge bases for “disease” keywords using the names of each ICD-11 chapter and build a lexicon of diseases corresponding to each chapter. Figure-1 shows what a lexicon build-up for an ICD-11 chapter looks like. Once the disease related keywords are identified for an ICD-11 chapter, we use these keywords to query PubMed via the Bio Entrez parser ³ to capture the first 1000 abstracts (PMIDs) spanning a period of last 90 days from the moment we initiated the query. We do this for each of the eight *special groups* ICD-11 chapters.

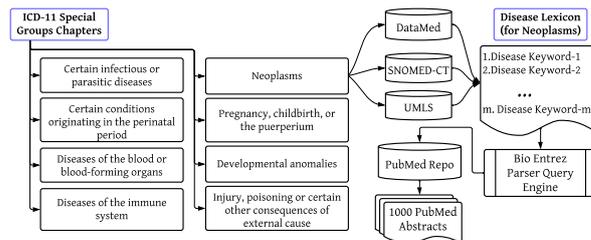


Fig. 1: ICD-11 based lexicon construction and querying for abstracts from PubMed using Bio Entrez parser. For illustration purpose, we show the pipeline for ICD-11 chapter *Neoplasms*

Figure-2 shows our ICD-11-Summ-1000 dataset preparation pipeline. Once we have queried the 1000 PubMed abstracts for an ICD-11 chapter, we conduct named entity recognition (NER) on each of the abstracts within a chapter using the SciSpacy NER model trained using the BC5CDR corpus [49]. Since we are interested in entity-level clustering of PubMed

¹<https://bit.ly/3GEFWvM>

²https://github.com/AmanuelF/biomed_abstractive_summarization

³<https://biopython.org/docs/1.75/api/Bio.Entrez.html>

abstracts within an ICD-11 chapter, we first conduct clustering of the named entities using agglomerative clustering as used in [18]. We use BioBERT [50] for named entity representation followed by agglomerative clustering. Once the named entities pertaining to an ICD-11 chapter are clustered into different bins, our next task is to cluster the PubMed abstracts into a bin based on how the named entities within the abstracts are related to the entities within a cluster. We use cosine similarity between named entities identified in a PubMed abstract and entities characterizing a cluster to determine an entity-aware cluster the abstract belongs to. Next, for each cluster, we perform named entity-aware salient content selection to produce an extractive pseudo-document for each cluster. This paradigm of reducing a multi-document corpus (i.e., a cluster consisting of PubMed abstracts grouped based on entity-relatedness) into an extractive pseudo-document is explored for different tasks in previous studies [51]–[53]. As part of the NER task, we use coreference resolution [54], [55] after named entities are extracted using SciSpacy to cluster the biomedical named entities and their coreferenced mentions spanning the multiple abstracts within an ICD-11 chapter.

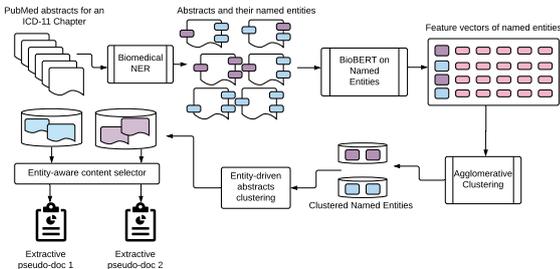


Fig. 2: Entity-aware content selection to produce extractive pseudo-documents. The light blue and light lavender colored documents in the final bins represent abstracts whose named entities are semantically similar to one another.

During entity-aware content selection to produce an extractive pseudo-doc for a cluster of PubMed abstracts that are clustered based on named entity relatedness, we preserve the positioning of sentences within an abstract. We also use the following heuristics while constructing the extractive pseudo-doc: 1) a sentence shall have at least one named entity identified using SciSpacy-BC5CDR NER model; and 2) the selected sentences from an abstract are placed in the same order as they appear in the abstract. Further, we also take into account abstracts’ relative importance scores where abstracts with higher document importance scores [56] have their sentences precede sentences from abstracts with less document importance scores while generating the extractive pseudo-document. Document importance D_{imp} of target abstract d_i is determined using pairwise cosine similarity between the BioBERT [50] embedding of d_i and other abstracts within the same cluster \mathcal{C} . Formally, document importance is defined as

$$D_{imp} = \frac{\sum_{d_i, d_j \in \mathcal{C}} \text{cossim}(d_i, d_j)}{|\mathcal{C}| - 1}, (i \neq j) \quad (1)$$

For all tasks throughout this paper involving initializing of networks or for representation learning, we use BioBERT [50].

IV. PROPOSED METHOD

Our proposed approach is a two-stage framework consisting of 1) an entity-driven knowledge retriever, and 2) a knowledge-guided abstractive summarizer. In this section, we discuss both modules in detail.

A. Entity-driven Knowledge Retriever

For each extractive pseudo-document generated in the data preparation stage for ICD-11 or input article for PubMed-50k designated by \mathcal{D} , we identify the named entities in the input document. The identified named entities are then used to retrieve facts from biomedical knowledge bases (UMLS, ICD-10, and SNOMED-CT). We use PyMedTermino [57] to work with the entire dump of UMLS [31] available at ⁴. For m named entities (and their coreferenced mentions), we have a set of pairs of entities $\{(e_i, e_j) \mid 0 \leq i < j < m\}$ extracted from \mathcal{D} , where each pair (e_i, e_j) is used to query for c candidate facts $F_1, F_2, F_3, \dots, F_{|c|}$ denoted collectively by $F_{\mathcal{D}}^{i,j}$ from the background knowledge bases \mathcal{K} using full text search. The complete set of facts retrieved for all pairs of named entities in source document \mathcal{D} is denoted by $\mathcal{F}_{\mathcal{D}}$.

The reason we use a pair of named entities to perform lexical query from \mathcal{K} is to capture the relationship between a pair of named entities as it appears in a knowledge base since relationships among named entities harbor semantics in addition to the entities themselves and help with disambiguating relevant facts from irrelevant facts. After the candidate facts $\mathcal{F}_{\mathcal{D}}$ are retrieved from the knowledge bases \mathcal{K} , we embed the candidate facts using BioBERT. Then, we perform efficient vector similarity search using Maximum Inner Product Search (MIPS) [58] implemented in the FAISS library ⁵ to query for the top-k facts among the candidate facts ($\mathcal{F}_{\mathcal{D}}$) using the input document \mathcal{D} as the query. Formally, we define the similarity between fact $F_i \in \mathcal{F}_{\mathcal{D}}$ and document \mathcal{D} as

$$\text{sim}(F_i, \mathcal{D}) = \vec{V}(F_i)^T \vec{V}(\mathcal{D}) \quad (2)$$

where $\vec{V}(F_i)$ - Vector representation of Fact F_i ; $\vec{V}(\mathcal{D})$ - Vector representation of document \mathcal{D}

Thus, after the knowledge retrieval task, we have 1) the input document \mathcal{D} which is obtained during the data preparation phase for ICD-11 and readily available for PubMed-50k; 2) the named entity chain (i.e., chain of named entities extracted from the pseudo-doc) \mathcal{E} [19]; and 3) top-k facts $F_1, F_2, F_3, \dots, F_{|K|}$ retrieved from the background knowledge bases collectively represented as $F_K \subseteq \mathcal{F}_{\mathcal{D}}$. We set the value of K to 3 following the study by [30]. We experiment with different values of K as detailed in the ablation studies section. The combination of these contextual signals will be used to guide the summarization model at training/inference time. The rationale for using maximum inner product search for knowledge retrieval

⁴<https://bit.ly/3E0zrll>

⁵<https://github.com/facebookresearch/faiss>

is inspired by the works of [26]–[28], [59], [60], albeit they used it mainly for open domain question answering [61], [62]. [30] use a similar approach for exemplar retrieval in their RetrievalSum model which is based on contrastive learning [63] using a Siamese network [64] to learn representations for an input document and the exemplars and guide their summary generation. Our problem of retrieving the most relevant facts from the background KB, however, is framed as a dense passage retrieval problem. Named entities from the input document are extracted using the SciSpacy NER model trained on the BC5CDR corpus [49]. Table-I shows sample facts, as they appear and retrieved from UMLS KB for an input article with a given pair of named entities identified.

Named Entity Pair	Entity-driven Facts from UMLS KB
(iron, anemia)	Iron deficiency anemia secondary to inadequate dietary iron intake. Iron deficiency anemia in mother complicating childbirth.
(dementia, depression)	Primary degenerative dementia of the Alzheimer type, presenile onset, with depression. Arteriosclerotic dementia with depression.
(diabetes, hypertension)	Hypertension in chronic kidney disease due to type 1 diabetes mellitus. Hypertension concurrent and due to end stage renal disease on dialysis due to type 2 diabetes mellitus.

TABLE I: Pairs of named entities and sample facts mined from UMLS for each pair.

B. Knowledge-guided Abstractive Summarizer

The backbone component of our knowledge-guided abstractive summarizer, which is a transformer encoder-decoder model, is based on the work by [65]. Figure-3 shows the proposed end-to-end model architecture. We use this architectural setup for all the models we experiment with. We designate a model augmented with one of the knowledge signals as model-EFAS. We train the models on the 50k samples obtained from PubMed abstractive scientific summarization dataset [6] using different combinations of signals (with and without named entities and facts). The top-k facts retrieved by the biomedical knowledge retriever, corresponding to each pair of named entities in an input extractive pseudo-doc or input article, are separated by a special token [SEP]. The input article is passed as one input document prepended with [CLS] and appended with [SEP] token. The named entity chain is passed as one segment prepended with [CLS] and appended with [SEP] token. There have been different approaches to combining different signals such as concatenating the different pieces to prime the generation component such as the one proposed in Fusion-in Decoder [59] and [60]. The top-k retrieved facts are initialized using BioBERT and the concatenated encoding is then passed through a sequence of transformer layers to be projected onto a 768- dimension vector to later be attended to by the autoregressive decoder. Similarly, the named entity chain is initialized with BioBERT and passed through a sequence of transformer encoders. Each transformer encoder layer is composed of self-attention and feed forward sub-layers. At training time, a batch of input-output pairings is passed to the encoder and decoder respectively in the form $\langle x, y \rangle$. The encoder undergoes the following transformations to the input sequence x which in this formalism is used to represent the first hidden layer h^0 of the stacked sequence of l transformer encoder layers.

$$\begin{aligned}\tilde{h}_x^l &:= \text{LayerNorm}(h_x^{l-1} + \text{MHAtt}(h_x^{l-1})) \\ h_x^l &:= \text{LayerNorm}(\tilde{h}_x^l + \text{FFN}(\tilde{h}_x^l))\end{aligned}$$

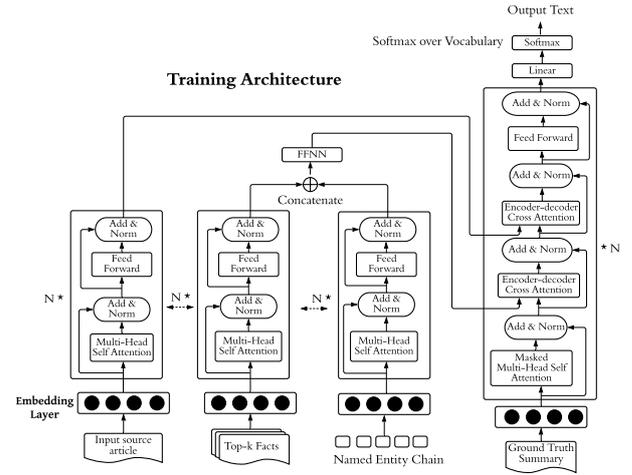


Fig. 3: The Proposed Framework. The encoder networks have their parameters shared. The two cross attention sub-layers in the decoder attend to the input source article, and a linear transformed projection of encodings of facts, and the chain of named entities. This architecture best represents the three traditional transformer models. For BigBird, and LED, the full self attention layer gets replaced with sparse attention.

The decoder component, which is trained using teacher forcing [66] at training time, consists of two cross-attention sub-layers to attend to: 1) the input source article; and 2) the affine transformed concatenation of facts and named entity chain’s encodings. The following formulations show the transformations in the decoder component where the ground truth output sequence y is passed to the sequence of transformer decoder layers and is used to initialize the first hidden layer h^0 of the decoder network. Note that we have a Masked Multi-head Self-Attention in the decoder network denoted by $MMHAtt$.

$$\begin{aligned}\tilde{h}_y^{l-2} &:= \text{LayerNorm}(h_y^{l-3} + \text{MMHAtt}(h_y^{l-3})) \\ \tilde{h}_y^{l-1} &:= \text{LayerNorm}(\tilde{h}_y^{l-2} + \text{CrossAtt}(\tilde{h}_y^{l-2}, F_K, \mathcal{E})) \\ \tilde{h}_y^l &:= \text{LayerNorm}(\tilde{h}_y^{l-1} + \text{CrossAtt}(\tilde{h}_y^{l-1}, x)) \\ h_y^l &:= \text{LayerNorm}(\tilde{h}_y^l + \text{FFN}(\tilde{h}_y^l))\end{aligned}$$

V. MODEL TRAINING

All models are trained with a cross-entropy loss using back-propagation:

$$\mathcal{L}_\theta = -\frac{1}{n} \sum_{k=1}^n \mathcal{P}(t_k | t_{<k}, X, \mathcal{E}, F_K; \theta) \quad (3)$$

Where X - the input sequence to be summarized; \mathcal{E} - the named entities chain in the input sequence X ; and F_K - top-k facts extracted from biomedical KB; θ - model parameters.

We train each model using cross entropy loss to generate the ground truth summaries for the PubMed-50k dataset. We use the following hyperparameters setting: number of epochs to 5, fixed learning rate of 5e-5 with adam optimizer [67], batch size to 8, beam size of 5 with a length penalty [68] α between the range of 0.6 and 1 [5] at inference time; to deal with long-document summarization using the traditional transformer encoder-decoder models, we split the source article into

chunks of a maximum of 512 tokens and independently encode each chunk, after which we concatenate and project back to 768 dimension using a linear layer. The approach of splitting the long input sequence into smaller chunks of 512 tokens and then embedding independently is motivated by the recent work by [48]. For Longformer-Encoder-Decoder (LED) [69] and BigBird [70], however, we set the maximum length of the input sequence to 8192 tokens since they can deal with long input sequences without having to truncate; maximum output sequence length is set to 210 tokens following the experiments by [6]. To mitigate redundancy in the generated summaries, we enable trigram blocking [71] during inference. For each backbone model, we use its base variant with 12 encoder and 12 decoder layers. The train/validation/test sizes for PubMed-50k is 50,000/5,000/5,000 and each model is trained using early stopping. A checkpoint of the model that performs the best (in terms of validation loss) on the validation set across different epochs is saved to the file system. All models are built and trained using PyTorch on NVIDIA Tesla T4 GPU. We perform model training experiments with different input guidance settings: input document only, input document + named entities chain, input document + named entities chain + knowledge facts. For our base summarization model, we experiment with five transformer-based encoder-decoder models and show that our entity-driven knowledge-aware approach enables us to achieve the best performance in entity-level factual consistency, N-gram novelty, and semantic equivalence while performing comparably on the commonly used ROUGE metrics. At inference time, we experiment with two settings (w/o named entities, and w/ named entities).

VI. EXPERIMENTS AND RESULTS

While all models are trained using the PubMed-50k corpus, they are evaluated using a hold-out test set from the original PubMed dataset as well as the ICD-11-Summ-1000 corpus we curate. The experimental results are shown in Table-II through Table-IX. Results of evaluation w.r.t source articles reported are average results for both the ICD-11-Summ-1000 and PubMed corpora since the ICD-11 pseudo-extractive documents do not have a ground truth summary. For lexical (ROUGE) evaluation, we report ROUGE F1 scores [72]. Similarly, here for evaluation conducted w.r.t source articles, the results reported are average results across the PubMed and ICD-11-Summ-1000 corpora. Entity-level factual accuracy [37] is measured in terms of precision, and recall w.r.t ground truth summary (for PubMed), and w.r.t source articles (for both PubMed and ICD-11-Summ-1000). Entity-level precision and recall w.r.t ground truth summaries are denoted with precision-target and recall-target; similarly, entity-level precision, and recall w.r.t the source article are designated with precision-source, and recall-source, respectively. The F1 score is the harmonic mean of the precision and recall for either case. For measuring semantic equivalence between generated summaries and ground truth summaries, we leverage BERTScore as proposed by [39]; specifically, we use BioBERT for representing each token in a generated summary and the ground truth

Backbone Model	Training Config ($K=3$)	R-1	R-2	R-L
T5	T5 Vanilla (Baseline)	31.333	12.821	29.018
	T5 w/ named entities (Ours)	29.915	11.352	27.667
	T5 w/ named entities /w facts - EFAS (Ours)	28.643	11.286	26.591
BART	BART Vanilla (Baseline)	34.214	13.830	31.545
	BART w/ named entities (Ours)	32.377	11.733	29.910
	BART w/ named entities /w facts - EFAS (Ours)	31.283	10.528	28.174
Pegasus	Pegasus Vanilla (Baseline)	28.851	11.274	26.859
	Pegasus w/ named entities (Ours)	30.365	11.483	28.003
	Pegasus w/ named entities /w facts - EFAS (Ours)	30.872	12.031	28.263
BigBird	BigBird Vanilla (Baseline)	35.426	13.801	32.537
	BigBird w/ named entities (Ours)	33.491	12.362	30.184
	BigBird w/ named entities /w facts - EFAS (Ours)	31.936	13.162	28.730
LED	LED Vanilla (Baseline)	36.218	14.173	32.862
	LED w/ named entities (Ours)	33.734	13.825	30.614
	LED w/ named entities /w facts - EFAS (Ours)	33.283	13.582	29.038

TABLE II: Lexical (ROUGE) Evaluation w.r.t Ground Truth Summary (*vanilla input @ inference time*). The input in this experimental setting is the *raw input article to be summarized (i.e., w/o named entity chain)*

Backbone Model	Training Config ($K=3$)	Entity-level Factual Consistency		
		Precision-target	Recall-target	F1 score-target
T5	T5 Vanilla (Baseline)	27.008	21.175	23.738
	T5 w/ named entities (Ours)	27.564	19.246	22.666
	T5 w/ named entities /w facts - EFAS (Ours)	27.329	19.136	23.310
BART	BART Vanilla (Baseline)	28.315	20.404	23.718
	BART w/ named entities (Ours)	27.949	19.105	22.695
	BART w/ named entities /w facts - EFAS (Ours)	27.241	18.792	22.241
Pegasus	Pegasus Vanilla (Baseline)	17.911	20.212	18.992
	Pegasus w/ named entities (Ours)	22.950	21.335	22.113
	Pegasus w/ named entities /w facts - EFAS (Ours)	23.572	22.956	23.260
BigBird	BigBird Vanilla (Baseline)	16.523	19.384	17.840
	BigBird w/ named entities (Ours)	23.273	21.831	22.529
	BigBird w/ named entities /w facts - EFAS (Ours)	25.317	23.839	24.586
LED	LED Vanilla (Baseline)	17.830	20.173	18.929
	LED w/ named entities (Ours)	24.528	22.573	23.510
	LED w/ named entities /w facts - EFAS (Ours)	26.827	25.322	26.053

TABLE III: Entity-level Factual Consistency Evaluation w.r.t Ground Truth Summary (*vanilla input @ inference time*). The input in this experimental setting is the *raw input article to be summarized (i.e., w/o named entity chain)*

summary after which we perform pairwise cosine similarity as proposed in [39]. All experimental results are reported in percentages. The average full text length of input source articles in PubMed-50k is 3,224 words and the average abstract length is 218 words, while for ICD-11-Summ-1000, the average length of an extractive pseudo-doc (i.e., input source article) is 4816 words.

VII. ABLATION STUDIES

To assess the impact of facts mined on the quality of summaries generated, we conduct an ablation study where we experiment with different values of K in top-k for the backbone models. Figure-4 shows results of ablation to assess precision-source and recall-target. Since we want to minimize entity hallucination which is measured in terms of precision-source and want to maximize the number of entities in the ground truth summary that are retrieved in the generated summary as measured by recall-target, we report the impact of different values of K for these two metrics. As shown in the two plots, precision-source and recall-target consistently improve as we

Backbone Model	Training Config ($K=3$)	Entity-level Factual Consistency		
		Precision-source	Recall-source	F1 score-source
T5	T5 Vanilla (Baseline)	55.076	7.976	13.934
	T5 w/ named entities (Ours)	54.015	7.232	12.756
	T5 w/ named entities /w facts - EFAS (Ours)	53.284	6.275	11.238
BART	BART Vanilla (Baseline)	58.592	5.623	10.261
	BART w/ named entities (Ours)	60.422	5.361	9.848
	BART w/ named entities /w facts - EFAS (Ours)	61.593	4.739	8.801
Pegasus	Pegasus Vanilla (Baseline)	33.821	7.401	12.144
	Pegasus w/ named entities (Ours)	46.757	7.743	13.286
	Pegasus w/ named entities /w facts - EFAS (Ours)	48.387	8.263	14.116
BigBird	BigBird Vanilla (Baseline)	34.288	9.261	14.583
	BigBird w/ named entities (Ours)	48.283	8.625	14.636
	BigBird w/ named entities /w facts - EFAS (Ours)	48.572	9.583	16.008
LED	LED Vanilla (Baseline)	59.361	6.731	12.091
	LED w/ named entities (Ours)	62.479	6.382	11.581
	LED w/ named entities /w facts - EFAS (Ours)	63.731	6.821	12.323

TABLE IV: Entity-level Factual Consistency w.r.t *source article*. The input in this experimental setting is the *raw input article to be summarized @ inference time (i.e., w/o named entity chain)*.

Backbone Model	Training Config ($k=3$)	R-1	R-2	R-L
		T5	T5 Vanilla (Baseline)	29.837
	T5 w/ named entities (Ours)	32.183	13.725	29.398
	T5 w/ named entities /w facts - EFAS (Ours)	29.372	9.682	28.275
BART	BART Vanilla (Baseline)	34.762	12.592	29.387
	BART w/ named entities (Ours)	35.281	12.938	31.276
	BART w/ named entities /w facts - EFAS (Ours)	33.731	11.923	30.285
Pegasus	Pegasus Vanilla (Baseline)	26.592	10.052	24.386
	Pegasus w/ named entities (Ours)	32.562	13.864	30.174
	Pegasus w/ named entities /w facts - EFAS (Ours)	33.824	13.841	30.639
BigBird	BigBird Vanilla (Baseline)	28.174	11.371	25.692
	BigBird w/ named entities (Ours)	32.281	14.263	31.863
	BigBird w/ named entities /w facts - EFAS (Ours)	34.728	13.264	31.752
LED	LED Vanilla (Baseline)	34.265	10.826	26.173
	LED w/ named entities (Ours)	36.840	13.773	32.155
	LED w/ named entities /w facts - EFAS (Ours)	34.927	14.003	30.851

TABLE V: Lexical (ROUGE) Evaluation w.r.t Ground Truth Summary (*input article + named entity chain @ inference time*); i.e., the input in this experimental setting is the *raw input article to be summarized with the named entities (i.e., w/ named entity chain)*.

Backbone Model	Training Config ($K=3$)	Entity-level Factual Consistency		
		Precision-target	Recall-target	F1 score-target
T5	T5 Vanilla (Baseline)	26.194	19.759	22.526
	T5 w/ named entities (Ours)	29.826	22.952	25.941
	T5 w/ named entities /w facts - EFAS (Ours)	28.582	20.738	24.036
BART	BART Vanilla (Baseline)	26.581	18.381	21.733
	BART w/ named entities (Ours)	27.949	19.105	22.696
	BART w/ named entities /w facts - EFAS (Ours)	27.341	18.793	22.341
Pegasus	Pegasus Vanilla (Baseline)	15.386	19.382	17.154
	Pegasus w/ named entities (Ours)	24.638	23.329	24.071
	Pegasus w/ named entities /w facts - EFAS (Ours)	25.498	24.374	24.923
BigBird	BigBird Vanilla (Baseline)	15.942	19.873	17.692
	BigBird w/ named entities (Ours)	26.315	24.728	25.497
	BigBird w/ named entities /w facts - EFAS (Ours)	26.638	24.163	25.340
LED	LED Vanilla (Baseline)	17.284	20.692	18.835
	LED w/ named entities (Ours)	28.173	25.866	26.970
	LED w/ named entities /w facts - EFAS (Ours)	26.116	26.830	26.468

TABLE VI: Entity-level Factual Consistency w.r.t Ground Truth Summary. The input in this experimental setting is the *raw input article to be summarized with the named entities (i.e., w/ named entity chain) @ inference time*.

Backbone Model	Training Config ($K=3$)	Entity-level Factual Consistency		
		Precision-source	Recall-source	F1 score-source
T5	T5 Vanilla (Baseline)	52.183	5.792	10.427
	T5 w/ named entities (Ours)	56.803	10.816	18.172
	T5 w/ named entities /w facts - EFAS (Ours)	55.728	8.629	14.944
BART	BART Vanilla (Baseline)	56.611	8.051	9.241
	BART w/ named entities (Ours)	62.385	7.284	13.045
	BART w/ named entities /w facts - EFAS (Ours)	61.938	6.382	11.572
Pegasus	Pegasus Vanilla (Baseline)	31.492	6.792	11.174
	Pegasus w/ named entities (Ours)	48.389	8.396	14.309
	Pegasus w/ named entities /w facts - EFAS (Ours)	48.964	9.491	15.900
BigBird	BigBird Vanilla (Baseline)	31.882	8.177	13.016
	BigBird w/ named entities (Ours)	48.733	9.267	15.575
	BigBird w/ named entities /w facts - EFAS (Ours)	50.373	11.274	18.424
LED	LED Vanilla (Baseline)	58.316	6.472	11.651
	LED w/ named entities (Ours)	63.722	8.537	15.057
	LED w/ named entities /w facts - EFAS (Ours)	65.180	8.374	14.841

TABLE VII: Entity-level Factual Consistency w.r.t *input source article (input article + named entity chain @ inference time)*; i.e., the input in this experimental setting is the *raw input article to be summarized with the named entities (i.e., w/ named entity chain)*.

Backbone Model	Training Configuration ($K=3$)	N-gram Novelty	
		w/o named entities	w/ named entities
T5	T5 Vanilla (Baseline)	52.930	49.699
	T5 w/ named entities (Ours)	50.079	50.967
	T5 w/ named entities /w facts - EFAS (Ours)	53.817	52.841
BART	BART Vanilla (Baseline)	54.816	54.997
	BART w/ named entities (Ours)	54.959	57.811
	BART w/ named entities /w facts - EFAS (Ours)	57.360	61.370
Pegasus	Pegasus Vanilla (Baseline)	51.260	50.035
	Pegasus w/ named entities (Ours)	52.558	51.269
	Pegasus w/ named entities /w facts - EFAS (Ours)	54.621	52.702
BigBird	BigBird Vanilla (Baseline)	49.783	51.374
	BigBird w/ named entities (Ours)	52.729	54.836
	BigBird w/ named entities /w facts - EFAS (Ours)	53.661	53.827
LED	LED Vanilla (Baseline)	53.732	53.288
	LED w/ named entities (Ours)	55.826	58.637
	LED w/ named entities /w facts - EFAS (Ours)	59.283	61.482

TABLE VIII: N-gram Novelty w.r.t source articles w/o and w/ named entity chain during inference.

Backbone Model	Training Configuration ($K=3$)	BioBERTScore	
		w/o named entities	w/ named entities
T5	T5 Vanilla (Baseline)	52.269	51.682
	T5 w/ named entities (Ours)	51.868	52.739
	T5 w/ named entities /w facts - EFAS (Ours)	53.162	54.164
BART	BART Vanilla (Baseline)	51.799	50.283
	BART w/ named entities (Ours)	51.783	53.618
	BART w/ named entities /w facts - EFAS (Ours)	52.072	51.472
Pegasus	Pegasus Vanilla (Baseline)	53.168	51.381
	Pegasus w/ named entities (Ours)	53.401	55.761
	Pegasus w/ named entities /w facts - EFAS (Ours)	54.382	55.263
BigBird	BigBird Vanilla (Baseline)	55.271	53.620
	BigBird w/ named entities (Ours)	56.813	54.271
	BigBird w/ named entities /w facts - EFAS (Ours)	56.372	55.088
LED	LED Vanilla (Baseline)	53.732	52.427
	LED w/ named entities (Ours)	54.163	55.791
	LED w/ named entities /w facts - EFAS (Ours)	53.814	57.284

TABLE IX: Semantic Equivalence (BioBERTScore [39]) w.r.t ground truth summaries w/o and w/ named entity chain during inference. Since we are using BioBERT for representation learning, we refer to the metric as BioBERTScore, a variant of BERTScore.

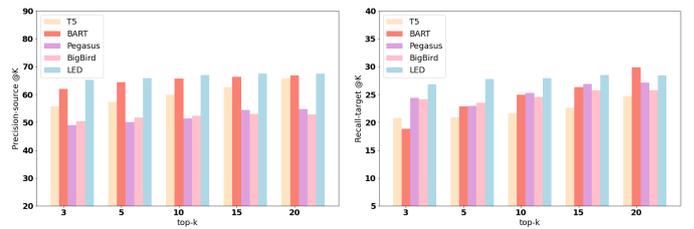


Fig. 4: Precision-source and Recall-target for different values of K . As can be seen, both metrics slightly increase as we increase the number of facts used to train a model. Note that this evaluation is done with source article and named entity chain passed to the trained models at inference time.

retrieve more relevant facts from the biomedical knowledge bases and train our models.

VIII. DISCUSSION OF RESULTS

From the results reported in the previous section, we generally see entity-level factual consistency (particularly, precision-source, and recall-target) improve when a model is trained with named entities and/or facts included as an additional signal in the training with the same objective of generating the ground truth summary using cross-entropy loss. The addition of more facts further improves entity-level factual consistency as shown in Figure-4. Further, we notice N-gram novelty improves with our proposed framework for the five backbone models. Semantic equivalence generally improves when named entities and/or facts are included during training for all models. Thus, the corresponding entries for the various models and training configurations show improvement in semantic based scores. The ROUGE scores, however, drop slightly from when there is no additional context at training or inference time. The drop in ROUGE is a result of augmenting the models with facts from background knowledge bases which in turn leads to higher N-gram novelty. Thus, the proposed framework enables us to achieve better abstractive scores in terms of entity-level factual consistency, paraphrasing and semantic equivalence.

IX. CONCLUSIONS AND FUTURE WORK

In this study, we proposed a framework to integrate named entities in a source article and facts extracted from biomedical knowledge bases pertaining to the named entities in transformer-based encoder-decoder models and applied to the task of abstractive summarization of biomedical literature. Through extensive experiments, we showed the proposed approach improves the semantics of generated summaries in terms of entity-level factual consistency and semantic equivalence while generating novel words. For future steps, we plan to jointly train the knowledge-retriever and the knowledge-guided abstractive summarizer in an end-to-end fashion. While our current architecture optimizes a single cost function given different input signals, we plan to augment the existing framework using multi-objective optimization to further enhance the factual accuracy and semantic equivalence of generated summaries using different cost functions during training.

REFERENCES

- [1] A. M. Rush, S. Chopra, and J. Weston, "A neural attention model for abstractive sentence summarization," *arXiv preprint arXiv:1509.00685*, 2015.
- [2] S. Chopra, M. Auli, and A. M. Rush, "Abstractive sentence summarization with attentive recurrent neural networks," in *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2016, pp. 93–98.
- [3] A. See, P. J. Liu, and C. D. Manning, "Get to the point: Summarization with pointer-generator networks," *arXiv preprint arXiv:1704.04368*, 2017.
- [4] A. R. Fabbri, I. Li, T. She, S. Li, and D. R. Radev, "Multi-news: A large-scale multi-document summarization dataset and abstractive hierarchical model," *arXiv preprint arXiv:1906.01749*, 2019.
- [5] Y. Liu and M. Lapata, "Text summarization with pretrained encoders," *arXiv preprint arXiv:1908.08345*, 2019.
- [6] A. Cohan, F. Dernoncourt, D. S. Kim, T. Bui, S. Kim, W. Chang, and N. Goharian, "A discourse-aware attention model for abstractive summarization of long documents," *arXiv preprint arXiv:1804.05685*, 2018.
- [7] E. Sharma, C. Li, and L. Wang, "Bigpatent: A large-scale dataset for abstractive and coherent summarization," *arXiv preprint arXiv:1906.03741*, 2019.
- [8] I. Cachola, K. Lo, A. Cohan, and D. S. Weld, "Tldr: Extreme summarization of scientific documents," *arXiv preprint arXiv:2004.15011*, 2020.
- [9] A. Esteva, A. Kale, R. Paulus, K. Hashimoto, W. Yin, D. Radev, and R. Socher, "Co-search: Covid-19 information retrieval with semantic search, question answering, and abstractive summarization," *arXiv preprint arXiv:2006.09595*, 2020.
- [10] R. I. Doğan, R. Leaman, and Z. Lu, "Ncbi disease corpus: a resource for disease name recognition and concept normalization," *Journal of biomedical informatics*, vol. 47, pp. 1–10, 2014.
- [11] J. Li, Y. Sun, R. J. Johnson, D. Sciaky, C.-H. Wei, R. Leaman, A. P. Davis, C. J. Mattingly, T. C. Wieggers, and Z. Lu, "Biocreative v cdr task corpus: a resource for chemical disease relation extraction," *Database*, vol. 2016, 2016.
- [12] Z. Zhao, Z. Yang, L. Luo, L. Wang, Y. Zhang, H. Lin, and J. Wang, "Disease named entity recognition from biomedical literature using a novel convolutional neural network," *BMC medical genomics*, vol. 10, no. 5, pp. 75–83, 2017.
- [13] V. Kocaman and D. Talby, "Biomedical named entity recognition at scale," in *International Conference on Pattern Recognition*. Springer, 2021, pp. 635–646.
- [14] W. Shen, J. Wang, and J. Han, "Entity linking with a knowledge base: Issues, techniques, and solutions," *IEEE Transactions on Knowledge and Data Engineering*, vol. 27, no. 2, pp. 443–460, 2014.
- [15] J. G. Zheng, D. Howsmon, B. Zhang, J. Hahn, D. McGuinness, J. Hendler, and H. Ji, "Entity linking for biomedical literature," *BMC medical informatics and decision making*, vol. 15, no. 1, pp. 1–9, 2015.
- [16] M. Zhu, B. Celikkaya, P. Bhatia, and C. K. Reddy, "Latte: Latent type modeling for biomedical entity linking," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 05, 2020, pp. 9757–9764.
- [17] E. Sharma, L. Huang, Z. Hu, and L. Wang, "An entity-driven framework for abstractive summarization," *arXiv preprint arXiv:1909.02059*, 2019.
- [18] H. Zhou, W. Ren, G. Liu, B. Su, and W. Lu, "Entity-aware abstractive multi-document summarization," in *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, 2021, pp. 351–362.
- [19] S. Narayan, Y. Zhao, J. Maynez, G. Simoes, V. Nikolaev, and R. McDonald, "Planning with learned entity prompts for abstractive summarization," *arXiv preprint arXiv:2104.07606*, 2021.
- [20] B. Gunel, C. Zhu, M. Zeng, and X. Huang, "Mind the facts: Knowledge-boosted coherent abstractive text summarization," *arXiv preprint arXiv:2006.15435*, 2020.
- [21] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in neural information processing systems*, 2017, pp. 5998–6008.
- [22] M. Lewis, Y. Liu, N. Goyal, M. Ghazvininejad, A. Mohamed, O. Levy, V. Stoyanov, and L. Zettlemoyer, "Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension," *arXiv preprint arXiv:1910.13461*, 2019.
- [23] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, "Exploring the limits of transfer learning with a unified text-to-text transformer," *Journal of Machine Learning Research*, vol. 21, no. 140, pp. 1–67, 2020.
- [24] J. Zhang, Y. Zhao, M. Saleh, and P. Liu, "Pegasus: Pre-training with extracted gap-sentences for abstractive summarization," in *International Conference on Machine Learning*. PMLR, 2020, pp. 11 328–11 339.
- [25] Z. Zhang, X. Han, Z. Liu, X. Jiang, M. Sun, and Q. Liu, "Ernie: Enhanced language representation with informative entities," *arXiv preprint arXiv:1905.07129*, 2019.
- [26] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W.-t. Yih, T. Rocktäschel *et al.*, "Retrieval-augmented generation for knowledge-intensive nlp tasks," *arXiv preprint arXiv:2005.11401*, 2020.
- [27] K. Guu, K. Lee, Z. Tung, P. Pasupat, and M.-W. Chang, "Realm: Retrieval-augmented language model pre-training," *arXiv preprint arXiv:2002.08909*, 2020.
- [28] V. Karpukhin, B. Oğuz, S. Min, P. Lewis, L. Wu, S. Edunov, D. Chen, and W.-t. Yih, "Dense passage retrieval for open-domain question answering," *arXiv preprint arXiv:2004.04906*, 2020.
- [29] F. Petroni, A. Piktus, A. Fan, P. Lewis, M. Yazdani, N. De Cao, J. Thorne, Y. Jernite, V. Karpukhin, J. Maillard *et al.*, "Kilt: a benchmark for knowledge intensive language tasks," *arXiv preprint arXiv:2009.02252*, 2020.
- [30] C. An, M. Zhong, Z. Geng, J. Yang, and X. Qiu, "Retrievalsum: A retrieval enhanced framework for abstractive summarization," *arXiv preprint arXiv:2109.07943*, 2021.
- [31] O. Bodenreider, "The unified medical language system (umls): integrating biomedical terminology," *Nucleic acids research*, vol. 32, no. suppl_1, pp. D267–D270, 2004.
- [32] J. Hirsch, G. Nicola, G. McGinty, R. Liu, R. Barr, M. Chittle, and L. Manchikanti, "Icd-10: history and context," *American Journal of Neuroradiology*, vol. 37, no. 4, pp. 596–599, 2016.
- [33] K. Donnelly *et al.*, "Snomed-ct: The advanced terminology and coding system for ehealth," *Studies in health technology and informatics*, vol. 121, p. 279, 2006.
- [34] J. E. Harrison, S. Weber, R. Jakob, and C. G. Chute, "Icd-11: an international classification of diseases for the twenty-first century," *BMC medical informatics and decision making*, vol. 21, no. 6, pp. 1–10, 2021.
- [35] B. Goodrich, V. Rao, P. J. Liu, and M. Saleh, "Assessing the factual accuracy of generated text," in *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2019, pp. 166–175.
- [36] W. Kryściński, B. McCann, C. Xiong, and R. Socher, "Evaluating the factual consistency of abstractive text summarization," *arXiv preprint arXiv:1910.12840*, 2019.
- [37] F. Nan, R. Nallapati, Z. Wang, C. N. d. Santos, H. Zhu, D. Zhang, K. McKeown, and B. Xiang, "Entity-level factual consistency of abstractive text summarization," *arXiv preprint arXiv:2102.09130*, 2021.
- [38] M. Zhong, P. Liu, Y. Chen, D. Wang, X. Qiu, and X. Huang, "Extractive summarization as text matching," *arXiv preprint arXiv:2004.08795*, 2020.
- [39] T. Zhang, V. Kishore, F. Wu, K. Q. Weinberger, and Y. Artzi, "Bertscore: Evaluating text generation with bert," *arXiv preprint arXiv:1904.09675*, 2019.
- [40] F. Schulze and M. Neves, "Entity-supported summarization of biomedical abstracts," in *Proceedings of the Fifth Workshop on Building and Evaluating Resources for Biomedical Text Mining (BioTxtM2016)*, 2016, pp. 40–49.

- [41] G. Erkan and D. R. Radev, "Lexrank: Graph-based lexical centrality as salience in text summarization," *Journal of artificial intelligence research*, vol. 22, pp. 457–479, 2004.
- [42] C. Zhu, W. Hinthorn, R. Xu, Q. Zeng, M. Zeng, X. Huang, and M. Jiang, "Enhancing factual consistency of abstractive summarization," *arXiv preprint arXiv:2003.08612*, 2020.
- [43] G. Manas, V. Aribandi, U. Kursuncu, A. Alambo, V. L. Shalin, K. Thirunarayan, J. Beich, M. Narasimhan, A. Sheth *et al.*, "Knowledge-infused abstractive summarization of clinical diagnostic interviews: Framework development study," *JMIR Mental Health*, vol. 8, no. 5, p. e20865, 2021.
- [44] K. Filippova, "Multi-sentence compression: Finding shortest paths in word graphs," in *Proceedings of the 23rd international conference on computational linguistics (Coling 2010)*, 2010, pp. 322–330.
- [45] A. Schrijver, *Theory of linear and integer programming*. John Wiley & Sons, 1998.
- [46] M. Faruqui, J. Dodge, S. K. Jauhar, C. Dyer, E. Hovy, and N. A. Smith, "Retrofitting word vectors to semantic lexicons," *arXiv preprint arXiv:1411.4166*, 2014.
- [47] M. Afzal, F. Alam, K. M. Malik, and G. M. Malik, "Clinical context-aware biomedical text summarization using deep neural network: Model development and validation," *Journal of medical Internet research*, vol. 22, no. 10, p. e19810, 2020.
- [48] J. DeYoung, I. Beltagy, M. van Zuylen, B. Kuehl, and L. L. Wang, "Ms2: Multi-document summarization of medical studies," *arXiv preprint arXiv:2104.06486*, 2021.
- [49] M. Neumann, D. King, I. Beltagy, and W. Ammar, "Scispacy: Fast and robust models for biomedical natural language processing," *arXiv preprint arXiv:1902.07669*, 2019.
- [50] J. Lee, W. Yoon, S. Kim, D. Kim, S. Kim, C. H. So, and J. Kang, "Biobert: a pre-trained biomedical language representation model for biomedical text mining," *Bioinformatics*, vol. 36, no. 4, pp. 1234–1240, 2020.
- [51] L. Lebanoff, K. Song, and F. Liu, "Adapting the neural encoder-decoder framework from single to multi-document summarization," *arXiv preprint arXiv:1808.06218*, 2018.
- [52] J. Zhang, J. Tan, and X. Wan, "Towards a neural network approach to abstractive multi-document summarization," *arXiv preprint arXiv:1804.09010*, 2018.
- [53] Y. Gao, W. Zhao, and S. Eger, "Supert: Towards new frontiers in unsupervised evaluation metrics for multi-document summarization," *arXiv preprint arXiv:2005.03724*, 2020.
- [54] K. Lee, L. He, M. Lewis, and L. Zettlemoyer, "End-to-end neural coreference resolution," *arXiv preprint arXiv:1707.07045*, 2017.
- [55] M. Gardner, J. Grus, M. Neumann, O. Tafjord, P. Dasigi, N. Liu, M. Peters, M. Schmitz, and L. Zettlemoyer, "Allennlp: A deep semantic natural language processing platform," *arXiv preprint arXiv:1803.07640*, 2018.
- [56] S. Banerjee, P. Mitra, and K. Sugiyama, "Multi-document abstractive summarization using ilp based multi-sentence compression," in *Twenty-Fourth International Joint Conference on Artificial Intelligence*, 2015.
- [57] L. Jean-Baptiste, "Using medical terminologies with pymedtermino and umls," in *Ontologies with Python*. Springer, 2021, pp. 207–239.
- [58] A. Shrivastava and P. Li, "Asymmetric lsh (alsh) for sublinear time maximum inner product search (mips)," *arXiv preprint arXiv:1405.5869*, 2014.
- [59] G. Izacard and E. Grave, "Leveraging passage retrieval with generative models for open domain question answering," *arXiv preprint arXiv:2007.01282*, 2020.
- [60] D. Singh, S. Reddy, W. Hamilton, C. Dyer, and D. Yogatama, "End-to-end training of multi-document reader and retriever for open-domain question answering," *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [61] J. M. Prager, "Open-domain question-answering," *Found. Trends Inf. Retr.*, vol. 1, no. 2, pp. 91–231, 2006.
- [62] P. Lewis, P. Stenetorp, and S. Riedel, "Question and answer test-train overlap in open-domain question answering datasets," *arXiv preprint arXiv:2008.02637*, 2020.
- [63] S. Chopra, R. Hadsell, and Y. LeCun, "Learning a similarity metric discriminatively, with application to face verification," in *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, vol. 1. IEEE, 2005, pp. 539–546.
- [64] G. Koch, R. Zemel, R. Salakhutdinov *et al.*, "Siamese neural networks for one-shot image recognition," in *ICML deep learning workshop*, vol. 2. Lille, 2015.
- [65] Z.-Y. Dou, P. Liu, H. Hayashi, Z. Jiang, and G. Neubig, "Gsum: A general framework for guided neural abstractive summarization," *arXiv preprint arXiv:2010.08014*, 2020.
- [66] R. J. Williams and D. Zipser, "A learning algorithm for continually running fully recurrent neural networks," *Neural computation*, vol. 1, no. 2, pp. 270–280, 1989.
- [67] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [68] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey *et al.*, "Google's neural machine translation system: Bridging the gap between human and machine translation," *arXiv preprint arXiv:1609.08144*, 2016.
- [69] I. Beltagy, M. E. Peters, and A. Cohan, "Longformer: The long-document transformer," *arXiv preprint arXiv:2004.05150*, 2020.
- [70] M. Zaheer, G. Guruganesh, K. A. Dubey, J. Ainslie, C. Alberti, S. Ontanon, P. Pham, A. Ravula, Q. Wang, L. Yang *et al.*, "Big bird: Transformers for longer sequences," in *NeurIPS*, 2020.
- [71] R. Paulus, C. Xiong, and R. Socher, "A deep reinforced model for abstractive summarization," *arXiv preprint arXiv:1705.04304*, 2017.
- [72] C.-Y. Lin and E. Hovy, "Automatic evaluation of summaries using n-gram co-occurrence statistics," in *Proceedings of the 2003 Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics*, 2003, pp. 150–157.