

Cross-Document Distillation via Graph-based Summarization of Extracted Essential Knowledge

Luca Ragazzi¹, Gianluca Moro², Lorenzo Valgimigli³, Riccardo Fiorani

Abstract—Abstractive multi-document summarization aims to generate a comprehensive summary that encapsulates crucial content derived from multiple input documents. Despite the proficiency exhibited by language models in text summarization, challenges persist in capturing and aggregating salient information dispersed across a cluster of lengthy sources. To accommodate more input, existing solutions prioritize sparse attention mechanisms, relying on sequence truncation without incorporating graph-based modeling of multiple semantic units to locate essential facets. Furthermore, the limited availability of training examples adversely impacts performance, thereby compromising summarization quality in real-world few-shot scenarios. In this paper, we present G-SEEK-2, a graph-enhanced approach designed to distill multiple topic-related documents by pinpointing and processing solely the pertinent information. We use a heterogeneous graph to model the input cluster, interconnecting various encoded entities via informative semantic edges. Then, a graph neural network locates the most salient sentences that are provided to a language model to generate the summary. We extensively evaluate G-SEEK-2 across seven datasets spanning various domains—including news articles, lawsuits, government reports, and scientific texts—under few-shot settings with a limited training sample size of only 100 examples. The experimental findings demonstrate that our model consistently outperforms advanced summarization baselines, achieving improvements as measured by syntactic and semantic metrics.

Index Terms—Multi-Document Summarization, Graph-Enhanced Transformer, Few-Shot Learning

I. INTRODUCTION

With a constant influx of new digital information, we are witnessing an exponential proliferation of textual data. Documentation plays a crucial role in grasping useful insights in various domains, including healthcare, law, and science journalism. In fact, workers invest considerable time in summarizing multiple topic-related documents into a unified text, whether it involves compiling outcomes from various lawsuits [1] or detecting key events from collections of news articles [2]. As a consequence, the proliferation of unstructured information has expanded the documentation workload, directly contributing to increased stress and burnout [3]. Even for attorneys with a high level of expertise, this intricate task naturally demands hours to accomplish, posing challenges for timely production [1]. Hence, there is a need to develop automated tools to accelerate human productivity.

In light of recent advances in natural language processing (NLP), there has been a surge in interest in abstractive summarization, which surpasses traditional extractive methods by adeptly paraphrasing the most significant details of documents.

Automatic text summarization tools play a vital role in helping people access the information they need, including lay summarization to increase readability and comprehension for non-experts [4]. In this context, a particularly challenging and practical task is the processing, identification, and synthesis of key information from a multitude of related sources, known as multi-document summarization (MDS) [5]. These assistive tools have received widespread attention across needs, ranging from query-focused summarization [6] to opinion summarization [7]. However, the complexities arising from the large volume of information and the inherent nature of documents—which often complement, overlap, or even contradict each other [8], [9]—contribute to the strong attention that the research community devotes to the advancement of MDS.

State-of-the-art MDS solutions are predominantly based on transformers [10], characterized by a structural constraint that links memory usage directly with input size, making them excessively resource-intensive when processing long texts. This limitation restricts the models to read only a fixed number of tokens,¹ introducing complications in MDS as it results in the truncation of any surplus information. In fact, unlike single document summarization, MDS methods must handle multiple texts, whose concatenation can form extremely lengthy inputs (e.g., 119,072.6 average words in MULTILEXSUM [1]). Therefore, standard sequence-to-sequence models, such as BART [12], are inadequate for MDS as they inevitably truncate long inputs, causing information loss and model degradation [13].² These problems are further exacerbated in real-world low-resource scenarios [15] characterized by a shortage of labeled instances available for model training supervision [16]. First, within small and medium organizations, creating gold-standard summaries from multiple lengthy documents can be costly, time-consuming, and may necessitate the expertise of domain specialists. Second, poorly correlated input-output training pairs³ can hinder effective learning [17]. Therefore, few-shot MDS emerges as a significant research area that deserves more attention from the NLP community [18]. Addressing these limitations presents an opportunity to accelerate real-world processes, thus also alleviating the burden associated with documentation.

In this work, we draw inspiration from the conventional approach humans adopt when engaging in text summarization

¹Tokens are subwords yielded by a subword tokenizer [11].

²Following [14], the input is the concatenation of documents in the cluster, each truncated according to the input length limit divided by the total number of sources, ensuring that each one is represented in the input (see Figure 1).

³Training examples composed of `<truncated_document, summary>` may inevitably lack syntactic and semantic correlation.

The authors contributed equally to this work.

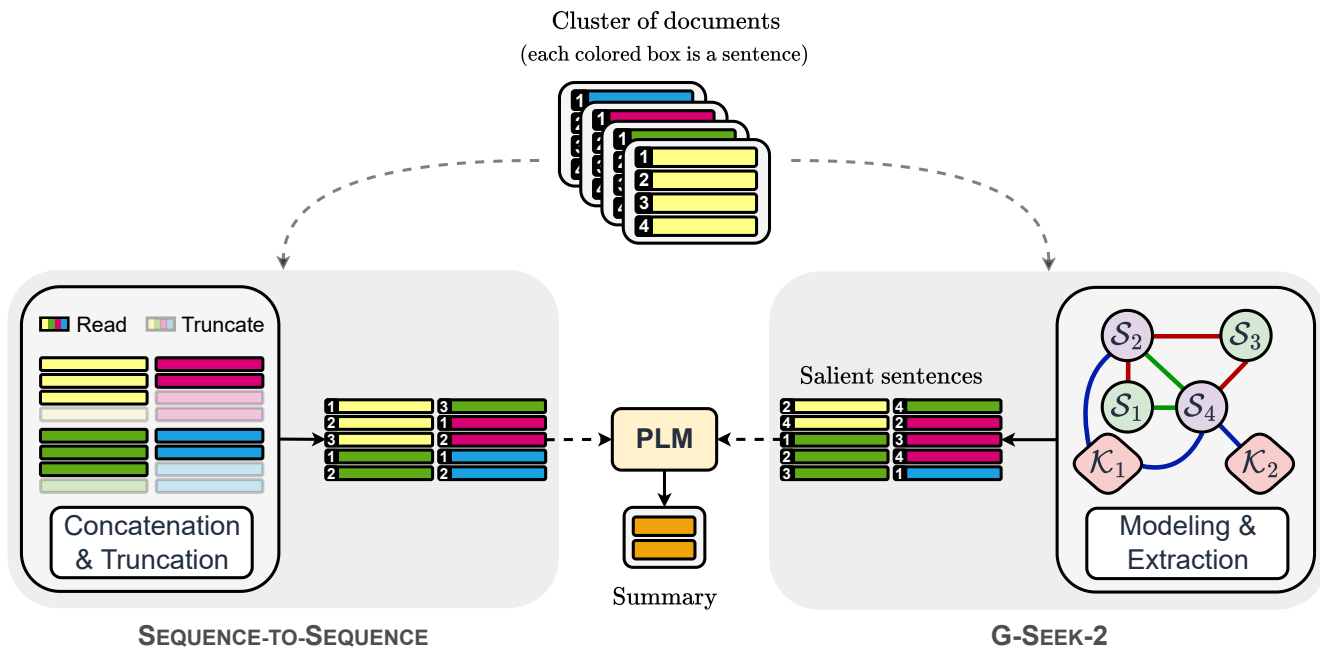


Fig. 1: Overview of our approach (right). Unlike standard sequence-to-sequence solutions (left), we first convert the documents into a heterogeneous graph: the pink diamonds represent the keyword nodes \mathcal{K} , the violet circles denote the sentence nodes \mathcal{S} containing keywords, the green circles indicate the context of \mathcal{S} , and the different segments between nodes symbolize edges. The salient sentences are then extracted and given to a generative PLM to produce the summary.

tasks. Specifically, humans read documents and highlight sentences that are deemed to be of greater importance. Subsequently, they review the underscored text and generate a summary based on it. This natural approach enables them to successfully synthesize extensive articles. A promising strategy to emulate this process involves using a two-stage pipeline approach that first extracts salient snippets and then summarizes them [19]. We face this task by using a semantic graph to represent the documents. This approach intuitively aggregates all source information, facilitating the identification and extraction of summary-worthy sentences. Previous contributions have used graph representations for text summarization [20]. Yet, they exhibit the following limitations: (i) they are mainly proposed for extractive summarization [21], [22], [23]; (ii) they are tailored for short texts [24], [25], which diverge highly from MDS settings; (iii) they do not leverage state-of-the-art generative pre-trained language models (PLMs) [26].

In light of this, we present G-SEEK-2 (Figure 1),⁴ a graph-based summarization of extracted essential knowledge. Our approach selects the most relevant sentences from a cluster of related documents and feeds them to a PLM to generate the summary. Technically, we model documents with a heterogeneous graph composed of multiple semantic edges and nodes of different granularities (i.e., keywords and sentences). Then, a graph neural network (GNN) is trained to select the salient sentence nodes—soft labeled with a heuristic—which are thereby provided to a PLM to produce the summary. This approach bypasses the limitation of feeding models only until

their maximum input size, which otherwise prevents them from fully processing the entire source, leading to performance drop [27]. Furthermore, by modeling a heterogeneous graph, we effectively capture cross-document relationships, a crucial aspect in handling multiple inputs [7]. Experimentally, we benchmark our solution in a realistic low-resource setting where a limited number of labeled training instances are available. This scenario is motivated by two key factors: (i) PLMs exhibit enhanced performance in generating summaries when trained with highly correlated source–target samples [17]; (ii) The limited number of trainable parameters of our learnable module (4M) allows G-SEEK-2 to avoid overfitting with a small number of examples. The experimental results yielded by quantitative and qualitative analyses register improved performance of summarization baselines equipped with G-SEEK-2 across various established evaluation metrics and datasets.

A conference version of this paper was presented in the Main Track of the 26th European Conference on Artificial Intelligence (ECAI 2023) [28]. In this manuscript, we enhance our work by refining model design, conducting additional in-depth comparative experiments, and expanding applications. The main differences are summarized as follows.

- To dissect the effectiveness of accurately capturing sentence node saliency in the graph, we conduct experiments with six different GNNs, such as Graph Convolutional Networks (GCN) [29], GraphSAGE [30], Graph Attention Networks (GAT) [31], Graph Isomorphism Networks (GIN) [32], Deep-GCN [33], and EdgeGCN [34].
- We evaluate our solution using three additional MDS datasets: MULTI-NEWS [35], WCEP [2], and WIKICAT-

⁴The code will be publicly released in case of acceptance.

SUM-ANIMAL [36], expanding our testbed to a total of seven corpora from various domains—news articles, lawsuits, government reports, and scientific texts—thoroughly exploring the generalizability of our approach.

The remainder of this paper is organized as follows. Section II offers an overview of related work on MDS. We present our method in Section IV and the experimental setup in Section V. Our results are reported and discussed in Section VI. Lastly, Section VII provides final remarks, highlighting limitations and future directions.

II. RELATED WORK

In this section, we provide a brief overview of related work on MDS, focusing on commonly used approaches and the utilization of graphs to model cross-document relationships.

A. Sequence-to-Sequence Pre-trained Language Models

Generative PLMs have demonstrated robust performance and adaptability to MDS [5], also when enhanced with reinforcement learning [37], and long-input summarization tasks in general. These sequence-to-sequence models are built on the transformer encoder–decoder architecture, denoted by stacks of self-attention layers. Vanilla transformers, such as BART [12] and PEGASUS [38], are limited to process up to 1024 tokens due to their quadratic memory and time complexity w.r.t. the input size. Consequently, they are not suitable for long sequences—including MDS—made up of tens of thousands of tokens. To address this obstacle, linear PLMs such as LED [39] and PRIMERA [14] feature a sparse attention mechanism that replaces the full quadratic self-attention by allowing models to scale linearly w.r.t. the input length. As a result, these models are capable of reading longer texts (e.g., up to 4096 tokens for PRIMERA), making them more suitable for MDS tasks. Despite the advantage of linear transformers in handling extensive information, they still rely on input truncation, like their quadratic counterparts. This constraint entails processing the source only up to the model's maximum input size (see Figure 1), thus overlooking potentially relevant details that merit inclusion in the summary.

B. Approaches for Multi-Document Summarization

MDS methods that do not adhere to a standard transformer-based sequence-to-sequence approach for handling multiple documents fall primarily into three categories.

a) Two-Phase Solutions: These methods, often referred to as “extract-then-abstract,” operate by selecting sentences deemed suitable for summarization, which are subsequently condensed by a PLM. Beyond unsupervised alternatives for sentence extraction [19] and abstractive summarization from clustered documents [40], most of the contributions embrace supervised techniques. Previous strategies relying on TF-IDF [41] have been replaced by ROUGE-based extractors, which aim to label sentences as relevant by measuring their similarity w.r.t. the summary [42]. This approach has also been addressed by jointly training the extractor with the summarization module [43]. In contrast to prior works, we

explore different metrics for sentence soft-labeling, covering BLEU [44] and all ROUGE [45] variants, such as ROUGE-1, ROUGE-2, and ROUGE-L, with precision and F1 assessments.

b) Aggregation-based Methods: These techniques involve combining hidden-states to aggregate information sourced from various snippets. Fusion-in-decoder [46] generates a unified hidden-state by concatenating multiple representations before decoding. On the other hand, marginalization-based models employ logit likelihood summation during decoding across inputs to weigh the probability of the next token [47], [48], [49]. Nevertheless, these solutions face difficulties when dealing with a limited number of labeled samples, primarily due to the initial cold-start phase required to train the model to effectively exploit this aggregated representation.

c) Hierarchical Models: These solutions aim to capture intricate interactions among documents to attain semantic-rich representations. Efforts have focused on enhancing transformers through graph-based techniques [50], [51], [52], [53], multi-head grouping and inter-paragraph attention [54], [55], maximal marginal relevance [35], and the inclusion of global and local attention [56], [57], [58]. Although hierarchical solutions have shown promising results, they struggle to effectively accommodate and leverage state-of-the-art PLMs [14].

C. Graph-based Summarization

Graphs and GNNs have emerged as integral components in MDS [20], offering enhanced scalability [59] and improved domain modeling [60] to mitigate transformer flaws. Several contributions leverage GNNs as standalone solutions [61], where the summary is generated by composing sentences extracted from input documents [62], [63]. Conversely, GNNs can be embedded with abstractive summarization models to improve performance [52]. Along this thread, BASS [26] introduces a unified semantic graph to represent the collection of texts and modifies the transformer architecture to interact with the graph. SKGSUM [24] exploits nodes at various levels to guide the summary generation process. To capture cross-document relationships, many studies have sought to construct different types of homogeneous graphs (e.g., topic graphs, discourse graphs) [64], [52], [65], [66]. Several studies have also delved into the use of heterogeneous graphs [67], [68], [69], [70]. However, these graphs accommodate only different types of nodes (e.g., word and sentence nodes) without considering diverse meaningful edges. On the other hand, HGSUM [71] extends a linear transformer with a heterogeneous graph of multiple semantic nodes and edges, but requires joint training that comes with the drawback of a costly pipeline.

III. RESEARCH OBJECTIVE

The goal of this research is to develop an automatic MDS approach that can address the following real-world challenges.

- **Long-input processing:** Design an extract-to-abstract framework to adeptly manage large information prevalent in MDS, bypassing input truncation drawbacks encountered in conventional sequence-to-sequence solutions.
- **Few labeled examples:** Integrate a powerful PLM to work successfully under few-shot conditions.

Algorithm 1 Sentence Soft-Labeling

Input:
 $\mathcal{X} = \{x_1, \dots, x_x\}$ {Input sentences}
 $\mathcal{Y} = \{y_1, \dots, y_y\}$ {Output sentences}
Parameters: \mathcal{M} {Similarity metric}
Output: \mathcal{S} {Set of scores}

```

0:  $\mathcal{S} \leftarrow \emptyset$ 
0: for  $x_i \in \mathcal{X}$  do
0:    $s \leftarrow \emptyset$ 
0:   for  $y_i \in \mathcal{Y}$  do
0:      $s.append(\mathcal{M}(x_i, y_i))$ 
0:   end for
0:    $\mathcal{S}.append(max(s))$ 
0: end for
0: return  $\mathcal{S} = 0$ 

```

TABLE I: MDS results on MULTI-LEXSUM with PRIMERA using various sentence soft-labeling methods as similarity functions. Notably, ROUGE-2-F1 results the most effective metric for the sentence labeling task.

Metric	R-1 _{f1}	R-2 _{f1}	R-L _{f1}
BLEU	43.62	19.18	28.58
R1-F1	44.06	19.33	29.32
R1-P	40.97	16.96	26.51
R2-F1	45.29	20.20	30.19
R2-P	43.32	19.20	29.14
RL-F1	43.18	19.34	28.30
RL-P	43.90	19.14	29.15

- **Cross-document relationships:** Leverage a heterogeneous graph to accurately encode document interrelations, and employ state-of-the-art GNNs to analyze their role in discerning significant patterns.
- **Generalizability:** Analyze system effectiveness across a set of multiple corpora from different domains through both syntactic and semantic evaluation metrics.

IV. METHOD

We present G-SEEK-2, a *graph-based summarization of extracted essential knowledge* (see Figure 1). Section IV-A delineates the preliminary procedures integral to our method, including the required sentence labeling and the compilation of passages essential for graph construction from extended inputs, using semantic and structural data. Section IV-B delves into the array of GNNs used to discern the interrelations among graph nodes, facilitating the identification of key sentences. Finally, Section IV-C provides the summarization pipeline.

A. Preliminaries

1) *Sentence Soft-Labeling:* Two-stage approaches have underscored the need to select suitable summary-worthy sentences from the source documents [72], [43]. Conceptually, each sentence can be categorized as relevant or irrelevant to the intended summary, enabling the training of a model to discern such noteworthy sentences. However, in the absence of ground-truth relevance labels—which reflects a characteristic of real datasets—we need to heuristically mark the salience of the input sentences w.r.t. the gold summary, namely performing a soft-labeling strategy. Formally, let $\mathcal{X} = \{x_1, \dots, x_{|\mathcal{X}|}\}$

and $\mathcal{Y} = \{y_1, \dots, y_{|\mathcal{Y}|}\}$ be the long input (i.e., the concatenated documents of the cluster) and the corresponding summary, respectively, where each $x_i \in \mathcal{X}$ and $y_i \in \mathcal{Y}$ is a sentence. We perform a greedy algorithm (Algorithm 1). Initially, for each instance x_i , we yield a list of relevance scores $\in [0, 1]$ through pairwise similarity computations among the sentences in the corresponding instance y_i . Subsequently, the highest scoring value is selected and assigned to x_i . For the assessment of sentence similarity, we examine various evaluation metrics, such as BLEU and ROUGE- $\{1,2,L\}$, considering both precision and F1 scores. Evaluation is carried out using the MULTI-LEXSUM-SHORT dataset [1] as testbed, using the first 100 samples from both training and validation sets. To appraise the summarization quality, we employ PRIMERA [14] as the backbone model, which is a transformer with linear complexity in the input length pre-trained with an MDS-specific objective. Practically, after assigning a summary-relevance score to each $x_i \in \mathcal{X}$, PRIMERA is furnished only with sentences (sequenced as per the source text) possessing the highest scores, up to the model’s maximum input size, which is 4096 tokens. The results, detailed in Table I, highlight ROUGE-2 F1 as the most effective metric for soft labeling. Consequently, this metric is adopted for the relevance labeling of sentences across all experiments. It is pertinent to note that Table I encapsulates the results derived from also extending soft labeling to the validation set. Pointedly, during inference, an oracle simulation is performed by accessing the ground-truth target summaries, thus facilitating the examination of upper-bound performance.

2) *Heterogeneous Graph:* We delineate the construction of our heterogeneous graph through the following steps, as illustrated in Figure 2.

a) *Keyword Extraction:* First, we eliminate English stopwords and general domain-specific terms appearing in more than 40% of the cluster. Then, we employ KEYBERT [73], [74] to identify the top- k keywords, which is a lightweight method in contrast to more resource-intensive alternatives [75]. Operationally, we extract k keywords for each document in the cluster and combine these keyword lists into a unique set, removing any duplicates.

b) *Sentence Filtering:* We partition the documents into sentences and select those containing at least one keyword. Additionally, we capture the surrounding context by retrieving the n sentences preceding and succeeding the selected one. Formally, we define $\{[x_1^1, \dots, x_{|\mathcal{X}^1|}^1], \dots, [x_1^z, \dots, x_{|\mathcal{X}^z|}^z]\}$ as the cluster of sources z . Let x_b^1 and x_e^1 be two sentences with keywords. We then select $[x_{b-n}^1, \dots, x_{b+n}^1, x_{e-n}^1, \dots, x_{e+n}^1]$, where $\langle x_{b-n}^1, \dots, x_{b+n}^1 \rangle - \{x_b\}$ is the context of x_b .

c) *Sentence & Keyword Embedding:* We use a frozen pre-trained DISTILROBERTA model [76] to generate the embeddings of three types of texts with different semantics: (1) keywords, (2) sentences containing keywords, and (3) neighboring sentences (i.e., the surrounding context) of the latter. DISTILROBERTA is characterized by a relatively small parameter count (82M), ensuring efficiency in terms of GPU memory utilization and computation. Notably, the model is already pre-trained to create sentence embeddings through a

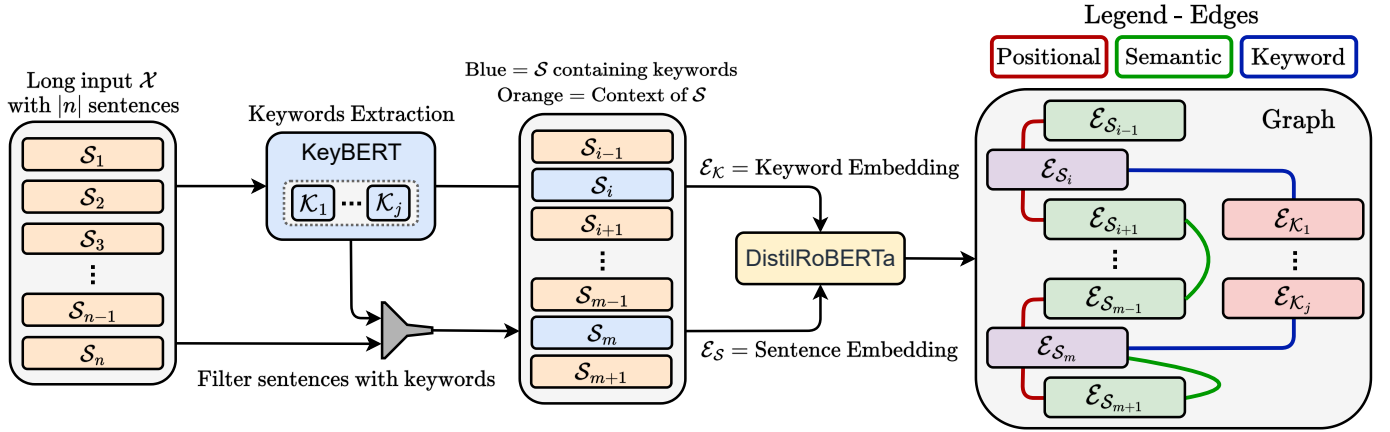


Fig. 2: The adopted pipeline for generating a semantic heterogeneous graph from multiple documents begins with providing the textual information to KEYBERT, which generates a set of unique keywords. Subsequently, the pipeline proceeds as follows: (i) sentences containing at least one keyword, (ii) their surrounding context (i.e., sentences immediately preceding or following), and (iii) the identified keywords, are transformed into embeddings using DISTILROBERTA, serving as the new graph nodes interconnected by various meaningful edges, i.e., positional, semantic, and keyword edge.



Fig. 3: Overview of our learnable module to model cross-document relationships over the heterogeneous graph. Notably, “GNN Layer” is compatible with various GNN architectures, ensuring our solution’s adaptability and flexibility.

self-supervised contrastive learning objective.⁵ About creating sentence embeddings, the model produces a representation for each token within a sentence. Then, following [77], we employ mean pooling to aggregate token embeddings, yielding a final single vector representation denoted as e_i^x .

d) *Graph Creation*: All keyword and sentence embeddings, denoted as KE and SE, respectively, are represented as nodes within our graph. Inspired by [26], we designate KEs as supernodes, indicating that every sentence containing a keyword establishes bidirectional *Keyword Edges* (cf. the blue lines in Figure 2) with the corresponding keyword node, rather than forming connections solely among themselves. Then, we introduce bidirectional *Positional Edges* (cf. the red lines in Figure 2) between two SEs if they appear consecutively in the source text. Finally, in alignment with [52] and [23], we incorporate *Semantic Edges* (cf. the green lines in Figure 2) between e_i and e_j if their cosine similarity is greater than a threshold t . The ultimate graph has as many nodes as the combined number of sentences and keywords.

B. Cross-Document Modeling

To capture cross-document relationships, we employ a learnable module featured by a GNN on our heterogeneous graph, discerning the significance of sentences by assigning an

unbounded positive relevance score to each sentence node. We harness both nodes (semantic information) and edges (structural information), essential for propagating information across nodes. This enables the GNN to attain a deeper comprehension of the context and meaning of each sentence. Our module comprises the following layers (see Figure 3):

- **Reprojection Layer** comprises two linear feed-forward layers (FFL) tasked with learning the transformation of the node embeddings x within the vector space, resulting in reprojected embeddings denoted as x' . Functionally, it augments the dimension d_x of the input embeddings n ($\mathbb{R}^{n \times 768}$) by a factor termed the *Boom Factor* (BF), drawing inspiration from the transformer architecture [10].

$$x' = \text{FFL}_\sigma(\text{FFL}_\gamma(x, d_x \cdot \text{BF}), d_x) \quad (1)$$

where σ and γ represent learnable parameters associated with distinct linear layers.

- **GNN Layer** operates on graph-structured data by exploiting structural information to improve the semantic representation of nodes and capture relational dependencies and patterns. We test multiple GNNs.⁶
 - GCN (Graph Convolutional Network) [29] uses convolutional operations to propagate information between nodes in the graph by leveraging a localized aggregation of neighboring node features.
 - GRAPHSAGE [30] is an inductive algorithm that learns a function to generate node embeddings by sampling and aggregating features from the local neighborhood of a target node.
 - GAT (Graph Attention Network) [31] assigns attention weights to each node, allowing the network to focus on the most influential neighbors throughout the information propagation process.

⁵<https://huggingface.co/sentence-transformers/all-distilroberta-v1>.

⁶We direct the reader to consult original papers for technical information.

- GIN (Graph Isomorphism Network) [32] aims to generalize the Weisfeiler-Lehman graph isomorphism test to achieve the maximum discriminative power.
- DEEP-GCN [33] leverages residual and dense connections and dilated convolutions into GCNs.
- EDGEGCN [34] uses multi-dimensional edge features for explicit relationship modeling.
- **Scoring Layer** is composed of two linear layers aimed at reducing the dimensionality of each node (\hat{x}) to a singular real number (s), serving as the relevance score for the sentence linked to the node (\mathbb{R}^n).

$$s = \text{FFL}_\beta(\text{FFL}_\theta(\hat{x}, d_{\hat{x}} \cdot \text{BF}), 1) \quad (2)$$

C. Summarization Pipeline

After converting the multi-document input into a graph and determining the relevance scores for each sentence, we employ a PLM to generate the output summary from the most pertinent sentences. Technically, based on their relevance scores, we select the most salient sentences and construct an input text for the model with fewer tokens than its maximum input size (i.e., without text overflow that would require truncation). Consequently, the new input comprises only important sentences arranged in the order of their occurrence in the original source.

We train the summarizer using the standard cross-entropy loss, wherein the model is tasked with predicting the next token w_i of the target sequence \mathcal{Y} given the input \mathcal{X} and the previous target tokens $w_{1:i-1}$, formulated as follows:

$$\mathcal{L}_{\text{ce}} = - \sum_{i=1}^{|\mathcal{Y}|} \log p_\tau(w_i | w_{1:i-1}, \mathcal{X}) \quad (3)$$

where τ denotes the model parameters and p represents the predicted probability distribution over the vocabulary. To create a more efficient training process, we do not jointly train G-SEEK-2 with the summarization model. Accordingly, our overall solution pipeline is specifically designed to work effectively with a limited number of labeled examples.

V. EXPERIMENTAL SETUP

In this section, we describe the experimental datasets, evaluation metrics, implementation details, and baselines used for comparison. Our research focuses primarily on the task of few-shot summarization, characterized by a real-world scenario marked by limited data availability for model supervision, mainly due to the high cost associated with annotation. In line with prior studies [78], [17], we select the first 100/10/100 samples from the training, validation, and test sets of all datasets without engaging in further data pre-processing.

A. Datasets

We perform experiments using multiple datasets from various domains that are publicly available in HuggingFace, serving as widely recognized benchmarks for MDS tasks. **MULTI-LEXSUM**⁷ [1] comprises real-world federal civil rights lawsuits accompanied by summaries authored by experts. The

⁷https://huggingface.co/datasets/allenai/multi_lexsum

primary challenge encountered in **MULTI-LEXSUM** lies in the extended length of the source documents and the varying granularity of the summaries, including tiny, short, and long versions (see Table II). Owing to this multifaceted nature, we conduct experiments using three distinct dataset renditions as testbed. **MULTI-NEWS**⁸ [35] is a large-scale dataset in which each instance comprises multiple news articles gathered from various sources, accompanied by a summary crafted by professional editors. **WCEP**⁹ [2] is constructed based on news events sourced from the Wikipedia Current Events Portal. It comprises clusters of news events along with human-authored summaries. We focus on the **WCEP-10** dataset, which features 10 related articles per instance. **WIKICAT-SUM-ANIMAL**¹⁰ [36] is a collection of news articles related to the domain of animals. Finally, **GOVREPORT**¹¹ [80] consists of lengthy reports from the US government. Although it is used in long document summarization benchmarks, we regard **GOVREPORT** as an intriguing testbed owing to its notably large input size, similar to [71] with the **ARXIV** dataset.

The key statistics of the datasets are presented in Table II. Specifically, we provide the number of samples in the corpus, the average number of documents per cluster, and the average number of words and sentences in both the source documents and target summaries—computed using the **NLTK** library [81]. Additionally, we report the average coverage, density, and compression ratio of extractive fragments, as defined by [79]. Technically, coverage is defined as the average fraction of token spans that can be identified in both the source and target. For instance, a coverage of 0.94 indicates that 94% of the summary words appear in extractive source fragments. Density, on the other hand, represents the average length of the extractive fragments. Finally, the compression ratio quantifies the extent to which a text is condensed to produce its summary.

Pointedly, our evaluation testbed encompasses several challenges: (i) processing very long legal reports (i.e., $> 100K$ words) with an extremely high compression ratio [1]; (ii) generating extremely short summaries [2], [1]; (iii) handling highly abstractive syntheses, where the targets contain few source-related tokens, as indicated by the low coverage [36]; and (iv) managing highly dense summary phrases [35], [80].

B. Metrics

We use standard **ROUGE**-{1,2,L}¹² F1 [45] and **BERTScore** F1 (BS) [83] to quantify the syntactic and semantic correspondence, respectively, between the generated summaries and the ground-truth. Furthermore, we compute $\mathcal{R} = \frac{\text{avg}(r_1, r_2, r_L)}{1 + \sigma_r^2}$ [82] to aggregate the ROUGE evaluation, where σ_r^2 represents the variance of the average ROUGE scores, penalizing discrepancies in performance across different dimensions. Note that all metrics lie within the range [0, 1], with higher scores indicating better performance. Table III reports additional details.

⁸https://huggingface.co/datasets/multi_news

⁹<https://huggingface.co/datasets/ccdv/WCEP-10>

¹⁰https://huggingface.co/datasets/GEM/wiki_cat_sum

¹¹<https://huggingface.co/datasets/ccdv/govreport-summarization>

¹²For ROUGE-L, we utilize summary-level computation where each summary is segmented into sentences.

TABLE II: Statistics of the evaluation datasets including size, number of source documents per instance, number of total words in source and target texts, and source–target coverage, density, and compression ratio of words [79]. Except for the number of samples, all reported values are averaged across all instances. We observe a comprehensive range of input lengths, which challenges our model and provides an in-depth evaluation benchmark.

Dataset	Domain	Samples	Source			Target		Source → Target		
			Docs	Words [†]	Sents	Words	Sents	Coverage	Density	Compress
MULTI-LEXSUM-TINY [1]	Legal	1603	10.7	119072.6	5962.5	24.7	1.4	0.92	2.27	5449.6
MULTI-LEXSUM-SHORT [1]	Legal	3138	10.3	99378.2	5017.0	130.2	5.1	0.96	3.33	840.7
MULTI-LEXSUM-LONG [1]	Legal	4534	8.8	75543.2	3814.2	646.5	28.8	0.94	4.07	97.4
MULTI-NEWS [35]	News	56,206	2.8	2092.1	80.9	257.9	10.0	0.82	5.47	8.1
WCEP [2]	News	10200	10.0	4356.1	154.0	31.9	1.4	0.91	3.17	162.1
WIKICAT-SUM-ANIMAL [36]	Scientific	53,638	131.0	5091.8	288.5	92.0	4.6	0.78	3.07	81.4
GOVREPORT [80]	Legal	19,463	1.0	8765.0	298.7	556.3	18.1	0.94	9.08	17.9

[†] The total average number of words in the source cluster; we consider a single input by concatenating all the documents.

TABLE III: Hyperparameters initialization and description of the evaluation metrics used for the text summarization task.

Metric	Description	Bound [*]	Hyperparameters
ROUGE [45]	Unigrams (R-1), bigrams (R-2), and longest common sub-sequence (R-L) lexical overlaps (%).	[0, 1] ↑	<code>rouge_types=["rouge1", "rouge2", "rougeL"],</code> <code>use_aggregator=True,</code> <code>use_stemmer=True,</code> <code>metric_to_select="fmeasure"</code>
\mathcal{R} [82]	Aggregated ROUGE value penalizing results with discrepant R-1, R-2, R-L.	[0, 1] ↑	/
BERTScore [83]	n-gram hard-alignment via contextualized BERT embeddings.	[0, 1] ↑	<code>model_type="microsoft/deberta-xlarge-mnli",</code> <code>batch_size=32</code>

* ↑ = higher is better.

TABLE IV: The number of trainable parameters of generative PLMs and their maximum input sequence length. G-SEEK-2 uses the max input size of the downstream model but provides salient sentences instead of truncating the exceeding ones.

	URL	#Params	Input
Models			
BART	https://huggingface.co/facebook/bart-large	400M	1024
PEGASUS	https://huggingface.co/google/pegasus-large	568M	1024
LED	https://huggingface.co/allenai/led-large-16384	459M	4096
PRIMERA	https://huggingface.co/allenai/PRIMERA	447M	4096
G-SEEK-2	-	+4M	-

C. Baselines

While decoder-only architectures driven by large language models (LLMs) have gained popularity for news summarization [84], [85], their application in MDS remains unexplored, and recent research demonstrates that encode–decoder networks may still offer superior summarization performance [86]. Therefore, to assess the efficacy of our proposed method in filtering relevant information prior to inputting it into generative models, we select several widely recognized and leading MDS solutions notable for their distinct capabilities in handling various input sizes. We then compare their performance when enhanced with G-SEEK-2. **BART** [12] is a transformer-based model with quadratic memory and time complexity concerning input length. **PEGASUS** [38] is a quadratic transformer pre-trained specifically

for summarization tasks, employing an objective to predict gap sentences as pseudo summaries. **LED** [39] is a transformer model with linear memory complexity, attributed to a sparse attention mechanism. **PRIMERA** [14] is a linear transformer built upon the LED architecture but with a pre-training objective specifically tailored for MDS, generating pseudo summaries by automatically extracting text spans based on entity salience. Technically, we adopt the conventional approach of concatenating documents from the same cluster to form a single long input. Following [14], we introduce a special token `<doc-sep>` to separate individual documents. We use the large checkpoints for all models. Table IV presents the number of parameters and max input size of the models.¹³

D. Implementation Details

We fine-tune the models using the PyTorch [87] implementations provided by the HuggingFace library [88], ensuring reproducibility by setting the seed to 42. All experiments are conducted on an internal workstation equipped with a Nvidia RTX 3090 GPU with 24 GB of memory, 64 GB of RAM, and an Intel(R) Core(TM) i9-10900X CPU @ 3.70GHz processor. For the GNN module, training is carried out over 75 epochs with a learning rate of $5e^{-5}$, using AdamW as the optimizer with $\beta_1 = 0.9$ and $\beta_2 = 0.99$. It should be noted that we trained all GNNs once on MULTI-LEXSUM-SHORT

¹³The maximum input length is determined by the encoder architecture of the models, while the output size varies depending on the dataset.

TABLE V: Precision, recall, and f-measure results of different graph settings on the eval set of MULTI-LEXSUM-SHORT.

#	P	R	F1	#	P	R	F1
Keywords				Consecutive Sentences			
4	22.12	76.43	34.31	1	23.26	62.64	33.92
5	23.29	78.81	35.95	2	22.12	76.43	34.31
6	22.98	75.42	35.22	3	23.24	76.00	35.60
7	22.65	77.26	35.03	4	23.15	78.85	35.79
8	22.57	77.12	34.92	5	23.06	77.42	35.54

using soft labels as described in Section IV-A1. Regarding the summarization task, all models are trained for 5 epochs with a learning rate of $3e^{-5}$, using mixed precision and gradient checkpointing techniques to conserve memory. During decoding, we employ beam search with 5 beams and n-gram repetition blocking for $n > 5$.

VI. RESULTS AND DISCUSSION

In this section, we present comprehensive experiments and detailed analyses to demonstrate the advantages and effectiveness of our proposed G-SEEK-2 module when integrated with PLMs for MDS tasks.

A. Analysis of the Graph

We explore various configurations of graph settings, with particular emphasis on the *Sentence Filtering* module (see Section IV-A2), where we evaluate Precision and Recall of all labeled sentences among the selected salient ones. The results of these investigations are summarized in Table V, conducted over the validation set of the MULTI-LEXSUM-SHORT dataset and considering the following facets:

- **Keywords.** We analyze the maximum number of keywords extracted by KEYBERT for each document in the cluster. We observe that 5 keywords yields the most favorable outcomes.
- **Context.** We investigate varying numbers of consecutive neighboring sentences selected as the context of the salient ones (i.e., those containing keywords). Results suggest that 4 sequences result in improved performance.

In our hardware environment (see Section V-D), the average time to create the graph for a single long input containing approximately 100K words is ≈ 34 seconds. Regardless, it is pertinent to note that the current implementation does not incorporate any specific optimizations.

B. Analysis of the GNN-based Module

We investigate different settings of our learnable module related to the inner architecture and the hyperparameters.

1) *Architecture:* Table VI reports the results obtained through various GNN architectures. The performance of sentence classification is assessed using two commonly employed evaluation metrics, such as Precision and Recall. First, we used only MULTI-LEXSUM-SHORT as our benchmark dataset. The results reveal that DEEP-GCN outperforms other architectures, demonstrating superior performance across all metrics. This

suggests its efficacy in accurately identifying relevant sentences while minimizing false positives. Conversely, GCN and GIN yield comparatively poorer results, indicating that these architectures might not be well-suited for the dataset under consideration. Consequently, we conduct this experiment on all the evaluation datasets to further validate these observations, as can be seen from Table VI. In fact, the results reveal that no single GNN architecture outperforms the others in all datasets. Nevertheless, it is evident that DEEP-GCN emerges as the optimal selection for legal corpora, while GRAPH-SAGE is preferable for other types. These findings emphasize the importance of choosing the appropriate model based on the specific characteristics and requirements of the dataset. Upon closer examination, certain datasets such as MULTI-NEWS, WCEP, and WIKICAT-SUM-ANIMAL display identical values for precision and recall. These datasets are characterized by shorter average input lengths, as detailed in Table II. Therefore, this feature may contribute to a more balanced task environment for determining the relevance of sentences.

2) *Ablation Studies:* We conduct experiments with our trainable module using the validation set of MULTI-LEXSUM-SHORT to analyze the best setting for the Reprojection and Scoring layers. Specifically, we employ GAT as our GNN layer, performing 30 training epochs and assessing efficacy by selecting the top 100 sentences based on their assigned scores. Subsequently, we calculate precision, recall, and f-measure metrics for these 100 sentences.

Table VII presents the results of our experiments, where we evaluate the following items:

- **Boom Factor:** We examine the impact of the Boom Factor in the *Reprojection Layer*. We found that a value of 2 yields optimal performance.
- **Layers:** We test different numbers of layers and surprisingly find that having only 1 layer achieves the best results. This suggests that, due to the limited pool of training examples, a lightweight solution with fewer trainable parameters is preferable.
- **Cosine Similarity:** We analyze the threshold for creating *Semantic Edges* between nodes. We uncover that a threshold of 0.86 produces the best outcomes.

Table VII additionally displays the results of the optimal module configuration trained for 100 epochs. Subsequently, we identify the most favorable model checkpoint, which occurs after 75 epochs. The average duration for each epoch is approximately 60 seconds. It should be noted that, after extensive experimentation, we systematically varied each hyperparameter while maintaining the other two constants at their respective optimal values (marked in bold in Table VII).

C. Summarization Results

We train and evaluate all models on the evaluation datasets both with and without G-SEEK-2 to underscore our contribution. Additionally, we provide the results for G-SEEK, which represents the initial version of our previous contribution [28]. Regarding the complexity analysis, models equipped with G-SEEK-2 have the same memory growth w.r.t. the input size

TABLE VI: Sentence classification results with different GNNs on each evaluated dataset equipped with soft labels. For each corpus, the best results are bolded and the second-best results are underlined. Notably, DEEP-GCN consistently outperforms other models in legal datasets, marking it as the most effective choice for these corpora. Conversely, GRAPHSAGE proves to be the superior option for other types of datasets.

GAT		GCN		GIN		DEEP-GCN		EDGE GCN		GRAPHSAGE	
P	R	P	R	P	R	P	R	P	R	P	R
MULTI-LEXSUM-TINY											
35.2	48.0	30.1	43.5	28.7	42.0	37.8	52.4	<u>36.9</u>	<u>51.5</u>	34.4	48.9
MULTI-LEXSUM-SHORT											
38.1	52.5	35.4	50.3	32.7	47.7	39.6	54.7	<u>39.1</u>	<u>54.2</u>	37.5	52.5
MULTI-LEXSUM-LONG											
37.0	51.0	33.0	47.5	31.5	45.8	38.5	53.2	<u>37.9</u>	<u>52.7</u>	36.2	50.6
MULTI-NEWS											
63.0	63.0	60.0	60.0	58.5	58.5	62.5	62.5	63.0	63.0	67.0	67.0
WCEP											
<u>78.3</u>	<u>78.3</u>	75.0	75.0	74.0	74.0	75.1	75.1	74.8	74.8	80.5	80.5
WIKICAT-SUM-ANIMAL											
<u>62.8</u>	<u>62.8</u>	60.5	60.5	59.0	59.0	60.2	60.2	61.2	61.2	63.4	63.4
GOVREPORT											
38.0	50.0	35.2	47.8	32.5	45.6	39.0	52.0	<u>38.5</u>	<u>51.5</u>	37.0	50.0

TABLE VII: The results of the learnable module on the validation set of MULTI-LEXSUM-SHORT under different settings with 30 training epochs. The “Final Module” is the best setting and checkpoint after 75 epochs.

Value	P	R	F1	Value	P	R	F1
Boom Factor							
1	31.34	38.43	34.52	Cosine Similarity			
2	32.12	39.17	35.29	0.80	32.02	39.01	35.17
3	31.69	38.80	34.88	0.82	32.10	39.18	35.28
4	31.52	38.54	34.68	0.84	32.10	39.17	35.28
GAT Layers				0.86	32.49	39.55	35.67
1	35.38	42.31	38.54	0.88	31.96	39.01	35.13
2	32.13	39.18	35.31	Final Module			
3	28.34	35.13	31.37	-	38.37	45.49	41.63
4	29.45	36.37	32.55				

of vanilla counterparts. More precisely, the space complexity to summarize the entire input document is $\mathcal{O}(L^2)$ for quadratic models (e.g., BART) and $\mathcal{O}(L)$ for linear ones (e.g., PRIMERA), where L is the minimum value between the source length and the model’s maximum input size.

Table VIII presents the performance of the systems on the datasets within the legal domain, while Table IX displays the results on the newly introduced corpora. Remarkably, our solution consistently enhances model performance across all datasets and metrics, underscoring its beneficial impact, which provides only salient information to generative PLMs. We positively highlight that G-SEEK-2 consistently surpasses on average our previous approach, underscoring the importance of selecting the appropriate GNN based on the data at hand.

To assess the effectiveness of input compression achieved

Model	Tokens	BS	Time (s)
BART	1024	76.59	8
w/ G-SEEK-2	1024	77.97	8
BART	512	75.33	6
w/ G-SEEK-2	512	77.87	6
BART	256	71.48	4
w/ G-SEEK-2	256	76.38	4

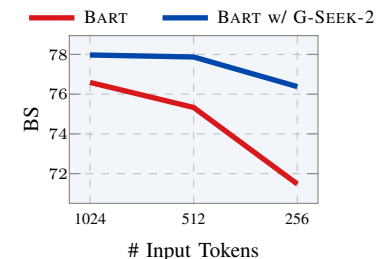


Fig. 4: Comparison of summarization quality with BERTScore (BS) by varying the number of input tokens on MULTI-LEXSUM-TINY. We report the time in seconds to compute the test set. On the right is furnished a graphical representation.

through graph-based processing, we conduct an experiment where we vary the number of input tokens provided to the generative PLM. Consequently, we choose different quantities of pertinent sentences to retain from the multi-document cluster. In Figure 4, we present the results based on semantic evaluation using BERTScore, employing BART as the summarizer and MULTI-LEXSUM-TINY as the dataset, and using three decreasing input sizes. Remarkably, as the input size decreases, the model equipped with G-SEEK-2 experience less impact compared to vanilla solutions that rely on reading the initial truncated tokens of the documents, which may include information unrelated to the summary. The drop in BERTScore as the number of input tokens decreases is due to important input information being truncated, leading to less relevant content in the summary and negatively impacting the score.

TABLE VIII: Evaluation F1 scores on MULTI-LEXSUM- $\{\text{TINY, SHORT, LONG}\}$ and GOVREPORT. The best intra-model score is bolded. † means statistically significant results of G-SEEK-2 (p-value < 0.05 with student t-test). Remarkably, G-SEEK-2 demonstrates superior performance compared to both standard solutions and those enhanced by its predecessor.

Model	MULTI-LEXSUM-TINY					MULTI-LEXSUM-SHORT					MULTI-LEXSUM-LONG					GOVREPORT				
	R-1	R-2	R-L	\mathcal{R}	BS	R-1	R-2	R-L	\mathcal{R}	BS	R-1	R-2	R-L	\mathcal{R}	BS	R-1	R-2	R-L	\mathcal{R}	BS
Quadratic																				
BART	22.37	7.91	19.74	16.61	76.17	41.45	18.74	35.81	31.70	79.89	41.41	16.47	38.98	31.88	79.13	48.46	14.03	44.83	34.94	80.82
w/ G-SEEK	24.46	7.70	20.34	17.41	77.07	41.57	16.72	35.78	31.01	79.97	43.92	17.43	41.08	33.67	80.19	51.46	17.12	48.05	37.97	81.78
w/ G-SEEK-2	27.72 †	11.04 †	22.64 †	20.49 †	79.58 †	42.48 †	15.98	37.57 †	31.85	80.77 †	48.08 †	22.06 †	45.57 †	38.48 †	82.15 †	52.30 †	17.39	48.65 †	39.45 †	81.77
PEGASUS	15.09	3.20	12.07	10.09	70.27	38.29	16.31	32.63	28.83	78.58	40.19	16.13	37.82	31.02	78.39	47.12	14.07	44.82	34.55	80.51
w/ G-SEEK	19.84	5.13	16.28	13.70	73.12	38.78	16.32	33.26	29.18	79.11	42.38	17.04	39.88	32.68	79.53	50.55	17.06	48.01	37.67	81.30
w/ G-SEEK-2	23.57 †	8.85 †	17.78 †	16.23 †	75.23 †	42.46 †	17.82 †	37.77 †	32.30 †	80.59 †	46.79 †	21.94 †	43.94 †	37.56 †	81.99 †	50.76	16.93	48.12	38.58 †	81.29
Linear																				
LED	22.86	7.98	18.86	16.50	76.20	40.09	17.50	35.15	30.63	79.51	45.26	20.01	42.66	35.52	81.31	53.86	19.53	49.28	39.96	82.65
w/ G-SEEK	24.39	7.96	20.55	17.55	77.19	40.95	16.28	35.31	30.51	80.56	45.42	18.87	42.93	35.23	81.03	55.63	21.08	50.67	41.49	82.77
w/ G-SEEK-2	27.74 †	12.99 †	22.38 †	20.94 †	78.43 †	46.98 †	21.08 †	41.42 †	36.49	82.70 †	50.94 †	25.75 †	47.58 †	41.48 †	83.34 †	58.29 †	22.37 †	54.04 †	44.87 †	83.01
PRIMERA	25.37	8.13	20.84	18.02	76.45	40.20	14.88	34.88	29.63	80.31	45.31	21.06	42.44	35.85	81.34	54.20	19.37	50.20	40.28	79.75
w/ G-SEEK	25.76	7.59	21.36	18.13	77.26	43.99	18.67	37.55	33.02	81.32	45.92	19.61	42.59	35.55	81.36	57.13	21.20	53.64	42.87	80.37
w/ G-SEEK-2	27.53 †	10.92 †	22.52 †	20.32 †	78.01 †	43.65	20.87 †	38.09	34.54 †	82.29 †	49.51 †	23.61 †	46.24 †	39.63 †	82.59 †	55.59	20.54	50.84	42.24	82.54 †

TABLE IX: Evaluation F1 scores on the new added datasets, such as WIKICAT-SUM-ANIMAL, MULTI-NEWS, and WCEP. The best intra-model score is bolded. † means statistically significant results of G-SEEK-2 (p-value < 0.05 with student t-test). G-SEEK-2 demonstrates superior performance compared to standard solutions and those enhanced by its predecessor.

Model	WIKICAT-SUM-ANIMAL					MULTI-NEWS					WCEP				
	R-1	R-2	R-L	\mathcal{R}	BS	R-1	R-2	R-L	\mathcal{R}	BS	R-1	R-2	R-L	\mathcal{R}	BS
Quadratic															
BART	37.04	17.76	27.24	30.10	78.00	43.94	13.42	20.08	32.34	79.10	41.33	24.89	33.72	33.87	83.00
w/ G-SEEK	45.86	14.49	24.04	27.62	78.29	38.25	10.54	18.83	29.20	78.36	36.51	15.79	26.49	26.93	80.56
w/ G-SEEK-2	39.64	14.28	25.38	29.74	79.20 †	40.68	9.97	17.78	28.37	78.21	41.35	18.09	29.28	29.75	84.39 †
PEGASUS	37.45	13.98	23.15	26.45	78.18	33.70	8.42	16.63	23.50	77.58	44.59	25.16	37.04	35.82	82.55
w/ G-SEEK	36.80	13.44	23.33	27.88	78.90	39.63	11.18	17.77	28.13	77.58	46.25	25.17	35.74	35.75	86.17
w/ G-SEEK-2	37.80	14.47 †	24.48 †	28.84 †	79.10 †	44.93 †	13.43 †	19.67 †	32.44 †	79.13 †	48.97 †	27.54 †	41.65 †	39.37 †	85.58
Linear															
LED	39.12	18.27	28.60	31.57	78.35	44.47	12.54	19.56	32.20	79.30	51.10	23.90	37.37	37.34	85.01
w/ G-SEEK	40.81	15.72	25.75	31.60	78.61	42.92	11.25	18.97	30.14	79.29	45.32	21.83	34.06	33.75	85.85
w/ G-SEEK-2	40.94	17.13	26.85	31.20	79.73 †	44.68	13.94 †	19.47	32.35	78.96	51.26	26.29 †	41.22 †	39.27 †	87.28 †
PRIMERA	44.18	19.49	28.83	33.16	80.00	39.90	9.79	18.67	28.88	78.30	46.24	26.33	38.07	37.15	84.49
w/ G-SEEK	40.71	16.69	27.58	31.92	79.66	42.49	10.34	18.00	29.38	78.41	43.71	25.82	35.33	35.45	84.63
w/ G-SEEK-2	44.32	21.06 †	29.33	35.70 †	80.71 †	40.20	9.09	17.02	28.18	78.89	48.13 †	25.42	38.11	36.91	86.06 †

VII. CONCLUSION

This study delves into the intricate domain of multi-document summarization, particularly simulating real-world scenarios where data availability is limited. We introduce G-SEEK-2, a graph-based method to distill essential insights from vast textual data, empowering abstractive summarization models to craft succinct and informative summaries. At the core of our approach lies the construction of a heterogeneous graph, representing a cluster of documents with various semantic units. This graph comprises distinct types of nodes and edges, meticulously designed to capture the nuanced relationships within the textual corpus. Through a tailored algorithm, we assign relevance scores to individual sentences, allowing us to pinpoint the most salient ones for inclusion in the summary. Experimental findings carried out in few-shot learning across multiple publicly available datasets demonstrate the remarkable performance enhancements achieved by G-SEEK-2. In particular, our approach significantly elevates both syntactic and semantic metrics reached by state-of-the-art summarization systems. Moreover, by fostering more coherent source-target pairs, we showcase how our solution facilitates faster learning for generative PLMs with limited labeled training instances. In future work, we will explore the development of lightweight end-to-end pipelines to jointly integrate our graph-based approach with generative PLMs to

enhance model interpretability [89]. Further, we aim to extend our methodology to take advantage of recent LLMs.

BROADER IMPACT AND ETHICS STATEMENT

Our research introduces G-SEEK-2, a graph-based multi-document summarization method designed to capture and distill essential knowledge from multiple documents. The implementation of G-SEEK-2 holds significant potential to enhance the efficiency and accuracy of multi-document summarization across various fields.

By generating concise and informative summaries from vast amounts of information, our method can make knowledge more accessible to a broader audience, including researchers, professionals, and the general public. Automating the summarization process can save substantial time and effort for professionals who handle large volumes of text, such as legal experts [90], [91], researchers, and content creators, increasing productivity and focus on critical tasks. For example, in educational settings, G-SEEK-2 can assist students and educators by providing succinct summaries of academic resources, facilitating quicker understanding and knowledge acquisition. Further, by delivering high-quality summaries, our solution can help researchers stay updated on developments in their fields, enabling them to review more literature in less time and promoting faster advancements in research.

However, it is crucial to be aware of potential negative impacts. The development and deployment of G-SEEK-2 must be conducted with careful consideration of ethical implications. While our solution aims to enhance productivity, it is important to ensure its responsible use. The tool should complement, not replace, human expertise, especially in critical domains such as healthcare, law, and journalism, where nuanced understanding is essential. Additionally, the risk of disseminating biased or inaccurate summaries could misinform users. Therefore, continuous evaluation, transparency, and ethical use are imperative to maximize the benefits and minimize the drawbacks of G-SEEK-2 and all generative models in our era.

REFERENCES

- [1] Z. Shen, K. Lo, L. Yu, N. Dahlberg, M. Schlanger, and D. Downey, "Multi-lexsum: Real-world summaries of civil rights lawsuits at multiple granularities," in *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, Eds., 2022. [Online]. Available: http://papers.nips.cc/paper_files/paper/2022/hash/552ef803bef9368c29e53c167de34b55-Abstract-Datasets_and_Benchmarks.html
- [2] D. G. Ghalandari, C. Hokamp, N. T. Pham, J. Glover, and G. Ifrim, "A large-scale multi-document summarization dataset from the wikipedia current events portal," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, D. Jurafsky, J. Chai, N. Schluter, and J. R. Tetreault, Eds. Association for Computational Linguistics, 2020, pp. 1302–1308. [Online]. Available: <https://doi.org/10.18653/v1/2020.acl-main.120>
- [3] E. Gesner, P. Gazarian, and P. C. Dykes, "The burden and burnout in documenting patient care: An integrative literature review," in *MEDINFO 2019: Health and Wellbeing e-Networks for All - Proceedings of the 17th World Congress on Medical and Health Informatics, Lyon, France, 25-30 August 2019*, ser. Studies in Health Technology and Informatics, L. Ohno-Machado and B. Sérroussi, Eds., vol. 264. IOS Press, 2019, pp. 1194–1198. [Online]. Available: <https://doi.org/10.3233/SHTI190415>
- [4] L. Ragazzi, P. Italiani, G. Moro, and M. Panni, "What are you token about? differentiable perturbed top-k token selection for scientific document summarization," in *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, L. Ku, A. Martins, and V. Srikumar, Eds. Association for Computational Linguistics, 2024, pp. 9427–9440. [Online]. Available: <https://aclanthology.org/2024.findings-acl.561>
- [5] C. Ma, W. E. Zhang, M. Guo, H. Wang *et al.*, "Multi-document summarization via deep learning techniques: A survey," *ACM Comput. Surv.*, vol. 55, no. 5, pp. 102:1–102:37, 2023.
- [6] G. Moro, L. Ragazzi, L. Valgimigli, and D. Freddi, "Discriminative marginalized probabilistic neural method for multi-document summarization of medical literature," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, S. Muresan, P. Nakov, and A. Villavicencio, Eds. Association for Computational Linguistics, 2022, pp. 180–189. [Online]. Available: <https://doi.org/10.18653/v1/2022.acl-long.15>
- [7] M. Zhang, G. Zhou, N. Huang, P. He, W. Yu, and W. Liu, "Asu-osum: Aspect-augmented unsupervised opinion summarization," *Inf. Process. Manag.*, vol. 60, no. 1, p. 103138, 2023. [Online]. Available: <https://doi.org/10.1016/j.ipm.2022.103138>
- [8] S. Lamsiyah, A. E. Mahdaouy, B. Espinasse, and S. E. A. Ouatik, "An unsupervised method for extractive multi-document summarization based on centroid approach and sentence embeddings," *Expert Syst. Appl.*, vol. 167, p. 114152, 2021. [Online]. Available: <https://doi.org/10.1016/j.eswa.2020.114152>
- [9] C. Shen, L. Cheng, X. Nguyen, Y. You, and L. Bing, "A hierarchical encoding-decoding scheme for abstractive multi-document summarization," in *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, H. Bouamor, J. Pino, and K. Bali, Eds. Association for Computational Linguistics, 2023, pp. 5872–5887. [Online]. Available: <https://aclanthology.org/2023.findings-emnlp.391>
- [10] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit *et al.*, "Attention is all you need," in *NIPS*, 2017, pp. 5998–6008.
- [11] T. Kudo and J. Richardson, "Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing," in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, EMNLP 2018: System Demonstrations, Brussels, Belgium, October 31 - November 4, 2018*, E. Blanco and W. Lu, Eds. Association for Computational Linguistics, 2018, pp. 66–71. [Online]. Available: <https://doi.org/10.18653/v1/d18-2012>
- [12] M. Lewis, Y. Liu, N. Goyal, M. Ghazvininejad *et al.*, "BART: denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension," in *ACL 2020, Online, July 5-10, 2020*. ACL, 2020, pp. 7871–7880. [Online]. Available: <https://doi.org/10.18653/v1/2020.acl-main.703>
- [13] Y. Song, Y. Chen, and H. Shuai, "Improving multi-document summarization through referenced flexible extraction with credit-awareness," in *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL 2022, Seattle, WA, United States, July 10-15, 2022*, M. Carpuat, M. de Marneffe, and I. V. M. Ruíz, Eds. Association for Computational Linguistics, 2022, pp. 1667–1681. [Online]. Available: <https://doi.org/10.18653/v1/2022.naacl-main.120>
- [14] W. Xiao, I. Beltagy, G. Carenini, and A. Cohan, "PRIMERA: Pyramid-based masked sentence pre-training for multi-document summarization," in *ACL (Volume 1: Long Papers)*. Dublin, Ireland: Association for Computational Linguistics, May 2022, pp. 5245–5263. [Online]. Available: <https://aclanthology.org/2022.acl-long.360>
- [15] G. Moro, L. Ragazzi, L. Valgimigli, G. Frisoni, C. Sartori, and G. Marfia, "Efficient memory-enhanced transformer for long-document summarization in low-resource regimes," *Sensors*, vol. 23, no. 7, p. 3542, 2023. [Online]. Available: <https://doi.org/10.3390/s23073542>
- [16] T. Yu, Z. Liu, and P. Fung, "Adaptsum: Towards low-resource domain adaptation for abstractive summarization," in *NAACL-HLT*. Association for Computational Linguistics, 2021, pp. 5892–5904.
- [17] G. Moro and L. Ragazzi, "Semantic self-segmentation for abstractive summarization of long documents in low-resource regimes," in *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, Thirty-Fourth Conference on Innovative Applications of Artificial Intelligence, IAAI 2022, The Twelfth Symposium on Educational Advances in Artificial Intelligence, EAAI 2022 Virtual Event, February 22 - March 1, 2022*. AAAI Press, 2022, pp. 11 085–11 093. [Online]. Available: <https://doi.org/10.1609/aaai.v36i10.21357>
- [18] M. A. Hedderich, L. Lange, H. Adel, J. Strötgen *et al.*, "A survey on recent approaches for natural language processing in low-resource scenarios," in *NAACL*. Online: ACL, Jun. 2021, pp. 2545–2568. [Online]. Available: <https://aclanthology.org/2021.naacl-main.201>
- [19] S. Tu, J. Yu, F. Zhu, J. Li, L. Hou, and J. Nie, "UPER: boosting multi-document summarization with an unsupervised prompt-based extractor," in *Proceedings of the 29th International Conference on Computational Linguistics, COLING 2022, Gyeongju, Republic of Korea, October 12-17, 2022*, N. Calzolari, C. Huang, H. Kim, J. Pustejovsky, L. Wanner, K. Choi, P. Ryu, H. Chen, L. Donatelli, H. Ji, S. Kurohashi, P. Paggio, N. Xue, S. Kim, Y. Hahn, Z. He, T. K. Lee, E. Santus, F. Bond, and S. Na, Eds. International Committee on Computational Linguistics, 2022, pp. 6315–6326. [Online]. Available: <https://aclanthology.org/2022.coling-1.550>
- [20] M. F. Salchner and A. Jatowt, "A survey of automatic text summarization using graph neural networks," in *COLING*. International Committee on Computational Linguistics, 2022, pp. 6139–6150.
- [21] R. Jia, Y. Cao, H. Tang, F. Fang *et al.*, "Neural extractive summarization with hierarchical attentive heterogeneous graph network," in *EMNLP (1)*. Association for Computational Linguistics, 2020, pp. 3622–3631.
- [22] D. Wang, P. Liu, Y. Zheng, X. Qiu *et al.*, "Heterogeneous graph neural networks for extractive document summarization," in *ACL*. Association for Computational Linguistics, 2020, pp. 6209–6219.
- [23] M. Chen, W. Li, J. Liu, X. Xiao *et al.*, "Sgsum: Transforming multi-document summarization into sub-graph selection," in *EMNLP (1)*. ACL, 2021, pp. 4063–4074.
- [24] X. Ji and W. Zhao, "SKGSUM: abstractive document summarization with semantic knowledge graphs," in *IJCNN*. IEEE, 2021, pp. 1–8.
- [25] Y. Qiu and S. B. Cohen, "Abstractive summarization guided by latent hierarchical document structure," in *EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022*. ACL, 2022, pp. 5303–5317. [Online]. Available: <https://aclanthology.org/2022.emnlp-main.355>

- [26] W. Wu, W. Li, X. Xiao, J. Liu *et al.*, "BASS: boosting abstractive summarization with unified semantic graph," in *ACL/IJCNLP (1)*. Association for Computational Linguistics, 2021, pp. 6052–6067.
- [27] G. Moro and L. Ragazzi, "Align-then-abstract representation learning for low-resource summarization," *Neurocomputing*, vol. 548, p. 126356, 2023. [Online]. Available: <https://doi.org/10.1016/j.neucom.2023.126356>
- [28] G. Moro, L. Ragazzi, and L. Valgimigli, "Graph-based abstractive summarization of extracted essential knowledge for low-resource scenarios," in *ECAI 2023 - 26th European Conference on Artificial Intelligence, September 30 - October 4, 2023, Kraków, Poland - Including 12th Conference on Prestigious Applications of Intelligent Systems (PAIS 2023)*, ser. Frontiers in Artificial Intelligence and Applications, K. Gal, A. Nowé, G. J. Nalepa, R. Fairstein, and R. Radulescu, Eds., vol. 372. IOS Press, 2023, pp. 1747–1754. [Online]. Available: <https://doi.org/10.3233/FAIA230460>
- [29] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. [Online]. Available: <https://openreview.net/forum?id=SJU4ayYgl>
- [30] W. L. Hamilton, Z. Ying, and J. Leskovec, "Inductive representation learning on large graphs," in *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, I. Guyon, U. von Luxburg, S. Bengio, H. M. Wallach, R. Fergus, S. V. N. Vishwanathan, and R. Garnett, Eds., 2017, pp. 1024–1034. [Online]. Available: <https://proceedings.neurips.cc/paper/2017/hash/5dd9db5e033da9c6f5ba83c7a7e9bea9-Abstract.html>
- [31] P. Velickovic, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio, "Graph attention networks," in *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. [Online]. Available: <https://openreview.net/forum?id=rjXmpikCZ>
- [32] K. Xu, W. Hu, J. Leskovec, and S. Jegelka, "How powerful are graph neural networks?" in *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. [Online]. Available: <https://openreview.net/forum?id=ryGs6iA5Km>
- [33] G. Li, M. Müller, A. K. Thabet, and B. Ghanem, "Deepgcn: Can gcn be as deep as cnns?" in *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*. IEEE, 2019, pp. 9266–9275. [Online]. Available: <https://doi.org/10.1109/ICCV.2019.00936>
- [34] C. Zhang, J. Yu, Y. Song, and W. Cai, "Exploiting edge-oriented reasoning for 3d point-based scene graph analysis," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*. Computer Vision Foundation / IEEE, 2021, pp. 9705–9715. [Online]. Available: https://openaccess.thecvf.com/content/CVPR2021/html/Zhang_Exploiting_Edge-Oriented_Reasoning_for_3D_Point-Based_Scene_Graph_Analysis_CVPR_2021_paper.html
- [35] A. R. Fabbri, I. Li, T. She, S. Li *et al.*, "Multi-news: A large-scale multi-document summarization dataset and abstractive hierarchical model," in *ACL, Florence, Italy, July 28- August 2 2019*. ACL, 2019, pp. 1074–1084. [Online]. Available: <https://doi.org/10.18653/v1/p19-1102>
- [36] L. Perez-Beltrachini, Y. Liu, and M. Lapata, "Generating summaries with topic templates and structured convolutional decoders," in *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, A. Korhonen, D. R. Traum, and L. Márquez, Eds. Association for Computational Linguistics, 2019, pp. 5107–5116. [Online]. Available: <https://doi.org/10.18653/v1/p19-1504>
- [37] J. Parnell, I. J. Unanue, and M. Piccardi, "A multi-document coverage reward for relaxed multi-document summarization," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, S. Muresan, P. Nakov, and A. Villavicencio, Eds. Association for Computational Linguistics, 2022, pp. 5112–5128. [Online]. Available: <https://doi.org/10.18653/v1/2022.acl-long.351>
- [38] J. Zhang, Y. Zhao, M. Saleh, and P. J. Liu, "PEGASUS: pre-training with extracted gap-sentences for abstractive summarization," in *ICML 2020, 13-18 July 2020, Virtual Event*, ser. Proceedings of Machine Learning Research, vol. 119. PMLR, 2020, pp. 11328–11339. [Online]. Available: <http://proceedings.mlr.press/v119/zhang20ae.html>
- [39] I. Beltagy, M. E. Peters, and A. Cohan, "Longformer: The long-document transformer," *CoRR*, vol. abs/2004.05150, 2020. [Online]. Available: <https://arxiv.org/abs/2004.05150>
- [40] M. Trabelsi and H. Uzunalioglu, "Absformer: Transformer-based model for unsupervised multi-document abstractive summarization," in *Document Analysis and Recognition - ICDAR 2023 Workshops - San José, CA, USA, August 24-26, 2023, Proceedings, Part II*, ser. Lecture Notes in Computer Science, M. Coustaty and A. Fornés, Eds., vol. 14194. Springer, 2023, pp. 151–166. [Online]. Available: https://doi.org/10.1007/978-3-031-41501-2_11
- [41] P. J. Liu, M. Saleh, E. Pot, B. Goodrich *et al.*, "Generating wikipedia by summarizing long sequences," in *ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018*. OpenReview.net, 2018. [Online]. Available: <https://openreview.net/forum?id=Hyg0vbWC->
- [42] Y. Zhang, A. Ni, Z. Mao, C. H. Wu, C. Zhu, B. Deb, A. H. Awadallah, D. R. Radev, and R. Zhang, "Summ^sn^s: A multi-stage summarization framework for long input dialogues and documents," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, S. Muresan, P. Nakov, and A. Villavicencio, Eds. Association for Computational Linguistics, 2022, pp. 1592–1604. [Online]. Available: <https://doi.org/10.18653/v1/2022.acl-long.112>
- [43] Z. Mao, C. H. Wu, A. Ni, Y. Zhang *et al.*, "DYLE: dynamic latent extraction for abstractive long-input summarization," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*. ACL, 2022, pp. 1687–1698. [Online]. Available: <https://doi.org/10.18653/v1/2022.acl-long.118>
- [44] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, "Bleu: a method for automatic evaluation of machine translation," in *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*. Philadelphia, Pennsylvania, USA: Association for Computational Linguistics, Jul. 2002, pp. 311–318. [Online]. Available: <https://aclanthology.org/P02-1040>
- [45] C.-Y. Lin, "ROUGE: A package for automatic evaluation of summaries," in *Text Summarization Branches Out*. Barcelona, Spain: Association for Computational Linguistics, Jul. 2004, pp. 74–81. [Online]. Available: <https://aclanthology.org/W04-1013>
- [46] G. Izacard and E. Grave, "Leveraging passage retrieval with generative models for open domain question answering," in *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, EACL 2021, Online, April 19 - 23, 2021*, P. Merlo, J. Tiedemann, and R. Tsarfaty, Eds. Association for Computational Linguistics, 2021, pp. 874–880. [Online]. Available: <https://doi.org/10.18653/v1/2021.eacl-main.74>
- [47] X. Gu, Y. Mao, J. Han, J. Liu *et al.*, "Generating representative headlines for news stories," in *WWW 2020: Taipei, Taiwan, April 20-24, 2020*. ACM / IW3C2, 2020, pp. 1773–1784. [Online]. Available: <https://doi.org/10.1145/3366423.3380247>
- [48] C. Hokamp, D. G. Ghalandari, N. T. Pham, and J. Glover, "Dyne: Dynamic ensemble decoding for multi-document summarization," *CoRR*, vol. abs/2006.08748, 2020. [Online]. Available: <https://arxiv.org/abs/2006.08748>
- [49] G. Moro, L. Ragazzi, L. Valgimigli, and L. Molfetta, "Retrieve-and-rank end-to-end summarization of biomedical studies," in *Similarity Search and Applications - 16th International Conference, SISAP 2023, A Coruña, Spain, October 9-11, 2023, Proceedings*, ser. Lecture Notes in Computer Science, O. Pedreira and V. Estivill-Castro, Eds., vol. 14289. Springer, 2023, pp. 64–78. [Online]. Available: https://doi.org/10.1007/978-3-031-46994-7_6
- [50] M. T. Nayeem, T. A. Fuad, and Y. Chali, "Abstractive unsupervised multi-document summarization using paraphrastic sentence fusion," in *COLING, Santa Fe, New Mexico, USA, August 20-26 2018*. ACL, 2018, pp. 1191–1204. [Online]. Available: <https://aclanthology.org/C18-1102/>
- [51] D. Antognini and B. Faltings, "Learning to create sentence semantic relation graphs for multi-document summarization," *CoRR*, vol. abs/1909.12231, 2019. [Online]. Available: <http://arxiv.org/abs/1909.12231>
- [52] W. Li, X. Xiao, J. Liu, H. Wu *et al.*, "Leveraging graph to improve abstractive multi-document summarization," in *ACL*. Association for Computational Linguistics, 2020, pp. 6232–6243.
- [53] R. K. Amplayo and M. Lapata, "Informative and controllable opinion summarization," in *EACL, Online, April 19-23 2021*. ACL, 2021, pp. 2662–2672. [Online]. Available: <https://aclanthology.org/2021.eacl-main.229/>
- [54] Y. Liu and M. Lapata, "Hierarchical transformers for multi-document summarization," in *ACL (1)*. Association for Computational Linguistics, 2019, pp. 5070–5081.
- [55] H. Jin, T. Wang, and X. Wan, "Multi-granularity interaction network for extractive and abstractive multi-document summarization," in *ACL*,

- Online, July 5-10 2020. ACL, 2020, pp. 6244–6254. [Online]. Available: <https://doi.org/10.18653/v1/2020.acl-main.556>
- [56] C. Zhu, R. Xu, M. Zeng, and X. Huang, “A hierarchical network for abstractive meeting summarization with cross-domain pretraining,” in *Findings of the Association for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020*, ser. Findings of ACL, T. Cohn, Y. He, and Y. Liu, Eds., vol. EMNLP 2020. Association for Computational Linguistics, 2020, pp. 194–203. [Online]. Available: <https://doi.org/10.18653/v1/2020.findings-emnlp.19>
- [57] S. Li and J. Xu, “Hiermds: a hierarchical multi-document summarization model with global-local document dependencies,” *Neural Comput. Appl.*, vol. 35, no. 25, pp. 18 553–18 570, 2023. [Online]. Available: <https://doi.org/10.1007/s00521-023-08680-0>
- [58] V. Nguyen, S. T. Mai, and M. Nguyen, “Learning to summarize multi-documents with local and global information,” *Prog. Artif. Intell.*, vol. 12, no. 3, pp. 275–286, 2023. [Online]. Available: <https://doi.org/10.1007/s13748-023-00302-z>
- [59] Z. Jia, S. Lin, M. Gao, M. Zaharia *et al.*, “Improving the accuracy, scalability, and performance of graph neural networks with roc,” in *MLSys*. mlsys.org, 2020.
- [60] A. Hogan, E. Blomqvist, M. Cochez, C. d’Amato *et al.*, “Knowledge graphs,” *ACM Computing Surveys (CSUR)*, vol. 54, no. 4, pp. 1–37, 2021.
- [61] D. Muratore, M. Hagenbuchner, F. Scarselli, and A. C. Tsoi, “Sentence extraction by graph neural networks,” in *ICANN (3)*, ser. Lecture Notes in Computer Science, vol. 6354. Springer, 2010, pp. 237–246.
- [62] M. Yasunaga, R. Zhang, K. Meelu, A. Pareek, K. Srinivasan, and D. R. Radev, “Graph-based neural multi-document summarization,” in *Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017), Vancouver, Canada, August 3-4, 2017*, R. Levy and L. Specia, Eds. Association for Computational Linguistics, 2017, pp. 452–462. [Online]. Available: <https://doi.org/10.18653/v1/K17-1045>
- [63] X. Doan, L. Nguyen, and K. N. Bui, “Multi graph neural network for extractive long document summarization,” in *COLING*. International Committee on Computational Linguistics, 2022, pp. 5870–5875.
- [64] Y. Yin, L. Song, J. Su, J. Zeng, C. Zhou, and J. Luo, “Graph-based neural sentence ordering,” in *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI 2019, Macao, China, August 10-16, 2019*, S. Kraus, Ed. ijcai.org, 2019, pp. 5387–5393. [Online]. Available: <https://doi.org/10.24963/ijcai.2019/748>
- [65] R. Pasunuru, M. Liu, M. Bansal, S. Ravi, and M. Dreyer, “Efficiently summarizing text and graph encodings of multi-document clusters,” in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, K. Toutanova, A. Rumshisky, L. Zettlemoyer, D. Hakkani-Tür, I. Beltagy, S. Bethard, R. Cotterell, T. Chakraborty, and Y. Zhou, Eds. Association for Computational Linguistics, 2021, pp. 4768–4779. [Online]. Available: <https://doi.org/10.18653/v1/2021.naacl-main.380>
- [66] Z. Zhang, H. Elfordy, M. Dreyer, K. Small, H. Ji, and M. Bansal, “Enhancing multi-document summarization with cross-document graph-based information extraction,” in *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2023, Dubrovnik, Croatia, May 2-6, 2023*, A. Vlachos and I. Augenstein, Eds. Association for Computational Linguistics, 2023, pp. 1688–1699. [Online]. Available: <https://doi.org/10.18653/v1/2023.eacl-main.124>
- [67] H. Wang, J. Chang, and J. Huang, “User intention-based document summarization on heterogeneous sentence networks,” in *Database Systems for Advanced Applications - 24th International Conference, DASFAA 2019, Chiang Mai, Thailand, April 22-25, 2019, Proceedings, Part II*, ser. Lecture Notes in Computer Science, G. Li, J. Yang, J. Gama, J. Natwichai, and Y. Tong, Eds., vol. 11447. Springer, 2019, pp. 572–587. [Online]. Available: https://doi.org/10.1007/978-3-030-18579-4_34
- [68] P. Cui and L. Hu, “Topic-guided abstractive multi-document summarization,” in *Findings of the Association for Computational Linguistics: EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 16-20 November, 2021*, M. Moens, X. Huang, L. Specia, and S. W. Yih, Eds. Association for Computational Linguistics, 2021, pp. 1463–1472. [Online]. Available: <https://doi.org/10.18653/v1/2021.findings-emnlp.126>
- [69] M. Lu, L. Liang, and G. Liu, “Parallel relationship graph to improve multi-document summarization,” in *Artificial Neural Networks and Machine Learning - ICANN 2022 - 31st International Conference on Artificial Neural Networks, Bristol, UK, September 6-9, 2022, Proceedings, Part II*, ser. Lecture Notes in Computer Science, E. Pimenidis, P. P. Angelov, C. Jayne, A. Papaleonidas, and M. Aydin, Eds., vol. 13530. Springer, 2022, pp. 630–642. [Online]. Available: https://doi.org/10.1007/978-3-031-15931-2_52
- [70] Y. Zhao, L. Wang, C. Wang, H. Du, S. Wei, H. Feng, Z. Yu, and Q. Li, “Multi-granularity heterogeneous graph attention networks for extractive document summarization,” *Neural Networks*, vol. 155, pp. 340–347, 2022. [Online]. Available: <https://doi.org/10.1016/j.neunet.2022.08.021>
- [71] M. Li, J. Qi, and J. H. Lau, “Compressed heterogeneous graph for abstractive multi-document summarization,” pp. 13 085–13 093, 2023. [Online]. Available: <https://doi.org/10.1609/aaai.v37i11.26537>
- [72] A. Bajaj, P. Dangati, K. Krishna, P. Ashok Kumar *et al.*, “Long document summarization in a low resource setting using pretrained language models,” in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: Student Research Workshop*. Online: ACL, Aug. 2021, pp. 71–80. [Online]. Available: <https://aclanthology.org/2021.acl-srw.7>
- [73] M. Grootendorst, “Keybert: Minimal keyword extraction with bert,” 2020. [Online]. Available: <https://doi.org/10.5281/zenodo.4461265>
- [74] A. Mathur and M. Suchithra, “Application of abstractive summarization in multiple choice question generation,” in *2022 International Conference on Computational Intelligence and Sustainable Engineering Solutions (CISES)*, 2022, pp. 409–413.
- [75] N. Nikzad-Khasmakhi, M. Feizi-Derakhshi, M. Asgari-Chenaghlu, M. A. Balafar *et al.*, “Phraseformer: Multimodal key-phrase extraction using transformer and graph embedding,” *CoRR*, vol. abs/2106.04939, 2021. [Online]. Available: <https://arxiv.org/abs/2106.04939>
- [76] V. Sanh, L. Debut, J. Chaumond, and T. Wolf, “Distilbert, a distilled version of BERT: smaller, faster, cheaper and lighter,” *CoRR*, vol. abs/1910.01108, 2019.
- [77] N. Reimers and I. Gurevych, “Sentence-bert: Sentence embeddings using siamese bert-networks,” in *EMNLP/IJCNLP (1)*. Association for Computational Linguistics, 2019, pp. 3980–3990.
- [78] Y. Chen and H. Shuai, “Meta-transfer learning for low-resource abstractive summarization,” in *AAAI 2021, Virtual Event, February 2-9, 2021*. AAAI Press, 2021, pp. 12 692–12 700. [Online]. Available: <https://ojs.aaai.org/index.php/AAAI/article/view/17503>
- [79] M. Grusky, M. Naaman, and Y. Artzi, “Newsroom: A dataset of 1.3 million summaries with diverse extractive strategies,” in *NAACL, Volume 1 (Long Papers)*. New Orleans, Louisiana: Association for Computational Linguistics, Jun. 2018, pp. 708–719. [Online]. Available: <https://aclanthology.org/N18-1065>
- [80] L. Huang, S. Cao, N. N. Parulian, H. Ji *et al.*, “Efficient attentions for long document summarization,” in *NAACL-HLT, Online, June 6-11, 2021*. ACL, 2021, pp. 1419–1436. [Online]. Available: <https://doi.org/10.18653/v1/2021.naacl-main.112>
- [81] S. Bird, “NLTK: the natural language toolkit,” in *ACL*. The Association for Computer Linguistics, 2006. [Online]. Available: <https://aclanthology.org/P06-4018/>
- [82] G. Moro, L. Ragazzi, and L. Valgimigli, “Carburacy: Summarization models tuning and comparison in eco-sustainable regimes with a novel carbon-aware accuracy,” in *Thirty-Seventh AAAI Conference on Artificial Intelligence, AAAI 2023, Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence, IAAI 2023, Thirteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2023, Washington, DC, USA, February 7-14, 2023*, B. Williams, Y. Chen, and J. Neville, Eds. AAAI Press, 2023, pp. 14 417–14 425. [Online]. Available: <https://doi.org/10.1609/aaai.v37i12.26686>
- [83] T. Zhang, V. Kishore, F. Wu, K. Q. Weinberger *et al.*, “Bertscore: Evaluating text generation with BERT,” in *ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. [Online]. Available: <https://openreview.net/forum?id=SkeHuCVFDr>
- [84] T. Goyal, J. J. Li, and G. Durrett, “News summarization and evaluation in the era of GPT-3,” *CoRR*, vol. abs/2209.12356, 2022. [Online]. Available: <https://doi.org/10.48550/arXiv.2209.12356>
- [85] T. Zhang, F. Ladhak, E. Durmus, P. Liang, K. R. McKeown, and T. B. Hashimoto, “Benchmarking large language models for news summarization,” *Trans. Assoc. Comput. Linguistics*, vol. 12, pp. 39–57, 2024. [Online]. Available: https://doi.org/10.1162/tacl_a_00632
- [86] Z. Fu, W. Lam, Q. Yu, A. M. So, S. Hu, Z. Liu, and N. Collier, “Decoder-only or encoder-decoder? interpreting language model as a regularized encoder-decoder,” *CoRR*, vol. abs/2304.04052, 2023. [Online]. Available: <https://doi.org/10.48550/arXiv.2304.04052>
- [87] A. Paszke, S. Gross, F. Massa, A. Lerer *et al.*, “Pytorch: An imperative style, high-performance deep learning library,” in *NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, 2019, pp. 8024–8035. [Online]. Available: <https://proceedings.neurips.cc/paper/2019/hash/bdbca288fee7f92f2bfa9f7012727740-Abstract.html>

- [88] T. Wolf, L. Debut, V. Sanh, J. Chaumond *et al.*, “Huggingface’s transformers: State-of-the-art natural language processing,” *CoRR*, vol. abs/1910.03771, 2019. [Online]. Available: <http://arxiv.org/abs/1910.03771>
- [89] G. Moro, L. Ragazzi, L. Valgimigli, F. Vincenzi, and D. Freddi, “Revelio: Interpretable long-form question answering,” in *The Second Tiny Papers Track at ICLR 2024, Tiny Papers @ ICLR 2024, Vienna, Austria, May 11, 2024*. OpenReview.net, 2024. [Online]. Available: <https://openreview.net/forum?id=fyvEJXsaQf>
- [90] G. Moro, N. Piscaglia, L. Ragazzi, and P. Italiani, “Multi-language transfer learning for low-resource legal case summarization,” *Artificial Intelligence and Law*, pp. 1–29, 2023. [Online]. Available: <https://doi.org/10.1007/s10506-023-09373-8>
- [91] L. Ragazzi, G. Moro, S. Guidi, and G. Frisoni, “Lawsuit: a large expert-written summarization dataset of italian constitutional court verdicts,” *Artificial Intelligence and Law*, 2024. [Online]. Available: <https://api.semanticscholar.org/CorpusID:272530787>