

Extractive Text Summarization with Latent Topics using Heterogeneous Graph Neural Network

Tuan-Anh Phan, Ngoc-Dung Nguyen, and Khac-Hoai Nam Bui *

Viettel Cyperspace Center, Viettel Group

Hanoi, Vietnam

{anhpt161, dungnn7, nambkh}@viettel.com.vn

Abstract

This paper presents a heterogeneous graph neural network (HeterGNN) model for extractive text summarization (ETS) by using latent topics to capture the important content of input documents. Specifically, topical information has been widely used as global information for sentence selection. However, most of the recent approaches use neural models, which lead the training models more complex and difficult for extensibility. In this regard, this study presents a novel graph-based ETS by adding a new node of latent topics into HeterGNN for summarization (TopicHeterGraphSum). Specifically, TopicHeterGraphSum includes three types of semantic nodes (i.e., topic-word-sentence) in order to enrich the cross-sentence relations for extractive summarization. Furthermore, an extended version of TopicHeterGraphSum for multi documents extraction is also taken into account to emphasize the advantage of the proposed method. Experiments on benchmark datasets such as CNN/DailyMail and Multi-News show the promising results of our method compared with state-of-the-art models.

1 Introduction

ETS is an important task of Natural Language Processing (NLP) in terms of extracting several relevant sentences from the original documents while keeping the main information. The traditional methods for ETS are TextRank (Mihalcea and Tarau, 2004)

and LexRank (Erkan and Radev, 2004), which focus on calculating the similarity between sentence scores. Sequentially, the rapid development of Deep Learning (DL) has brought breakthrough records by modeling a document as a sequence of sequences in order to deal with long-range inter-sentence relationships for the summarization (Cheng and Lapata, 2016; Cohan et al., 2018). However, cross-sentence relation is still a challenge in this research field (Liu and Lapata, 2019). Recent works focus on Graph Neural Networks (GNNs) (e.g., Graph Convolutional Network (GCN) (Kipf and Welling, 2017) or Graph Attention Network (GAT) (Velickovic et al., 2018)) to explore the cross-sentence relationships for the summarization task. The core idea is to represent inter-sentential graphs and use message passing to extract the complex relationship in the input documents. For instance, (Yasunaga et al., 2017) and (Xu et al., 2020) adopt discourse analysis to build document graphs. (Jia et al., 2020) and (Wang et al., 2020) built a bi-partite graph between words and sentences, which is regarded as a heterogeneous graph neural network. Moreover, modeling global information is also taken into account for sentence selection by using pre-trained models (Liu and Lapata, 2019; Zhang et al., 2019).

Sequentially, (Cui et al., 2020) utilized pre-trained BERT to learn contextual sentence representations and train jointly with latent topics using the neural topic model (NTM). (Nguyen et al., 2021) presents an extended version using NTM for abstractive text summarization indicating the capability of enriching the global information for the summarization.

Corresponding Author

Although the existing methods have provided remarkable results, there are several open research issues that need to take into account: i) the high performance mainly depends on pre-trained models for learning sentence representations, which is difficult for the extensibility, especially for low-resource languages; ii) the current external information (e.g, latent topics) are extracted by neural models, which requires more complex configurations of the training process. Furthermore, the model might be suffering because of the bias problem, especially in terms of small datasets; iii) multi-document summarization is still an open research issue, which requires a comprehensive summary for covering an event and avoiding redundancy. In this regard, this study proposes a new HeterGNN model for the ETS problem by adding latent topic nodes into a graph structure, in which the initialized topic features are extracted by well-known clustering methods such as K-means and Gaussian Mixture Models (GMM). The core idea is to investigate the impact of topical information on the EDS problem in terms of both single and multiple document extraction. To the best of our knowledge, this paper is the first study to adopt topical information for multi documents summarization by using the concept of the heterogeneous graph structure. More detail of the proposed model is described in the following sections.

2 Background

The proposed model is based on the concept of a HeterGNN model, which is proposed by (Wang et al., 2020), for enriching the relationships between sentences by adding nodes with semantic features. Specifically, Fig. 1 illustrates the HeterGNN for the ETS problem. Particularly, the model includes three main components, such as initialized graph structure, graph layer, and sentence selection module. Graph structure is initialized by the set of word nodes, which is encoded using Glove (Pennington et al., 2014) as the addition node, and sentence features, which are calculated by combining CNN for extracting the local n-gram feature of each sentence and bidirectional Long Short-Term Memory (BiLSTM) for extracting the sentence-level feature, respectively. In this regard, the feature of the sentence

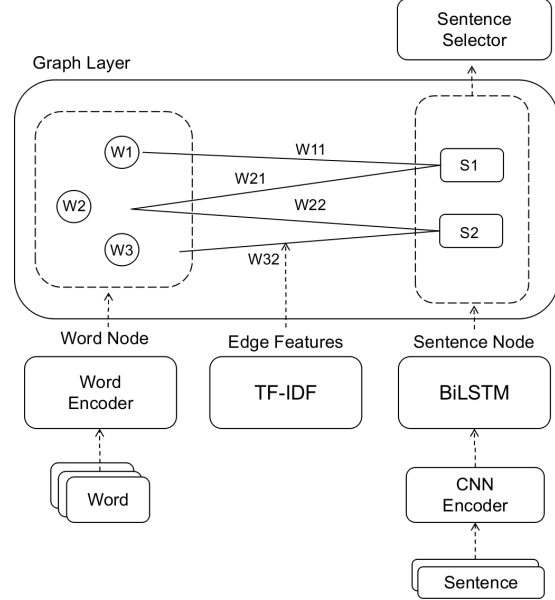


Figure 1: Model overview of HeterSumGraph

s_j can be obtained as follows:

$$X_{s_j} = CNN(x_{1:p}) \oplus BiLSTM(x_{1:p}) \quad (1)$$

where p denotes the number of words in the sentence. Moreover, tf-idf is adopted for further approval of information between words and sentences. Sequentially, the graph layer is updated using GAT(Velickovic et al., 2018), with a modification for the heterogeneous graph. Specifically, the updated node representation with modified GAT can be formulated as follows:

$$z_{ij} = LeakyReLU(W_a[W_q h_i; W_e h_j; \bar{e}_{ij}]) \quad (2)$$

where \bar{e}_{ij} denotes the multi-dimensional embedding space ($\bar{e}_{ij} \in R^{d_e}$), which is mapped from edge weight e_{ij} . Thereby, the sentences with their neighbor word nodes are updated via modified-GAT and Position-Wise Feed-Forward (FFN) layer, which can be sequentially formulated as follows:

$$\begin{aligned} U_{s \leftarrow w}^1 &= GAT(H_s^0, H_w^0, H_w^0) \\ H_s^1 &= FFN(U_{s \leftarrow w}^1 + H_s^0) \end{aligned} \quad (3)$$

where H_w^0 and H_s^0 are the node features of word X_w ($X_w \in R^{m \times d_w}$) and sentences X_s ($X_s \in R^{n \times d_s}$), respectively. Therefore, the new representations of word nodes can be obtained using the updated sentence nodes and further updated sentences

or query nodes, iteratively. Each iteration contains a sentence-to-word and a word-to-sentence update process, which can be demonstrated as follows:

$$\begin{aligned}
U_{w \leftarrow s}^{t+1} &= GAT(H_w^t, H_s^t, H_s^t) \\
H_w^{t+1} &= FFN(U_{w \leftarrow s}^{t+1} + H_w^t) \\
U_{s \leftarrow w}^{t+1} &= GAT(H_s^t, H_w^{t+1}, H_w^{t+1}) \\
H_s^{t+1} &= FFN(U_{s \leftarrow w}^{t+1} + H_s^t)
\end{aligned} \tag{4}$$

The output of the new sentence representation is input into a sentence classifier, which uses cross-entropy loss, for ranking the classification.

3 Methodology

In this study, our model is proposed for single document summarization (SDS), however, it can be extended for multi documents (MDS) with minor modifications. The methods for the two aforementioned problems are described in the following sections.

3.1 Single Document Summarization

Given an arbitrary document $d = \{s_1, \dots, s_n\}$, which includes n sentences, the objective of EDS for single document problem is to predict a set of binary label $\{y_1, \dots, y_n\}$ ($y_j \in [0, 1]$), which determine that the sentence in the summary or not. Figure 2 illustrates the structure of the proposed HeterGNN model.

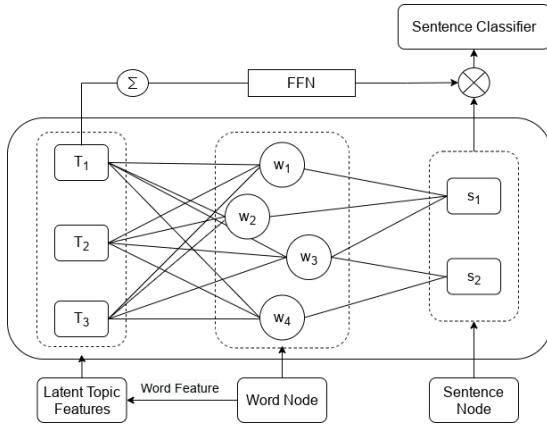


Figure 2: Overview of TopicHeterGraphSum for single document summarization. The initialized word node and sentence node features are processed following the work in (Wang et al., 2020). Furthermore, we provide latent topics as addition nodes into the Hetergraph.

Specifically, compared with previous works, the main idea of the proposed model is to enrich global

information. Accordingly, instead of using neural models for generating latent topics, we first extract the initialized topic feature of each document using simple clustering methods (e.g., K-mean and GMM) of pre-trained word embeddings (Sia et al., 2020). In particular, the initialized topic feature is calculated as follows:

$$X_T = \underset{\sum}{\operatorname{argmin}} \begin{cases} \|c^{(i)} - x_j\|, Kmean \\ \theta_i f(x_j | c^{(i)}, \Sigma_i), GMM \end{cases} \tag{5}$$

where θ_i denotes topic proportions. $c^{(i)}$ and x_j represent the cluster center and word vector, respectively. Sequentially, the latent topics are put into a graph layer for extracting semantic information. Similar to sentence representation calculation in Eq. 3, the topic representation can be updated via modified GAT as follows:

$$\begin{aligned}
U_{T \leftarrow w}^1 &= GAT(H_T^0, H_w^0, H_w^0) \\
H_T^1 &= FFN(U_{T \leftarrow w}^1 + H_T^0)
\end{aligned} \tag{6}$$

Each iteration contains word-to-sentence, sentence-to-word, and word-to-topic for the update process, which can be formulated as follows:

$$\begin{aligned}
U_{w \leftarrow s}^{t+1} &= GAT(H_w^t, H_s^t, H_s^t) \\
U_{w \leftarrow T}^{t+1} &= GAT(H_w^t, H_T^t, H_T^t) \\
U_{w \leftarrow s, T}^{t+1} &= \sigma(U_{w \leftarrow s}^{t+1} + U_{w \leftarrow T}^{t+1}) \\
H_w^{t+1} &= FFN(U_{w \leftarrow s, T}^{t+1} + H_w^t) \\
U_{s \leftarrow w}^{t+1} &= GAT(H_s^t, H_w^{t+1}, H_w^{t+1}) \\
H_s^{t+1} &= FFN(U_{s \leftarrow w}^{t+1} + H_s^t) \\
U_{T \leftarrow w}^{t+1}, A_{T \leftarrow w}^{t+1} &= GAT(H_T^t, H_w^{t+1}, H_w^{t+1}) \\
H_T^{t+1} &= FFN(U_{T \leftarrow w}^{t+1} + H_T^t)
\end{aligned} \tag{7}$$

where $A_{T \leftarrow w}$ denotes the attention matrix from the word node to the topic node. Subsequently, the topic representation of the input document is calculated by combining all topic features, which are learned using GAT as follows:

$$\begin{aligned}
\alpha_i &= \frac{\sum_{n=1}^{N_d} c(w_n) * A_{i,n}}{\sum_{j=1}^K \sum_{n=1}^{N_d} c(w_n) * A_{j,n}} \\
H_{T_d} &= \sum_{i=1}^K \alpha_i * H_{T_i}
\end{aligned} \tag{8}$$

where $A_{i,j}$ indicates the amount of information word j contributes to topic i . $c(w_n)$ is the frequency of w_n

in the document, K is the number of topics and α_i refers to the level dominant of topic i th to the total document-topic. Sequentially, each sentence hidden state is integrated with the above topic vector to capture sentence-topic representation as follows:

$$H_{s_i, T_d} = FFN(H_{T_d}) \oplus H_{s_i} \quad (9)$$

Finally, the output sentence-topic representation is used for sentences classification by using cross-entropy loss as the training objective:

$$\mathcal{L} = \sum_{i=1}^n y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i) \quad (10)$$

3.2 Multi Documents Summarization

Currently, there are not many studies for multi-document summarization. The main challenge of MDS is that the input documents may differ in terms of main focus and point of view (Fabbri et al., 2019). Intuitively, enriching global information is able to improve the performance of the MDS problem in which latent topics, extracted from word nodes, are considered for whole sentences in the multi documents. Therefore, in this paper, we take MDS into account by extending our proposed HeterGNN model. Fig. 3 demonstrates the modification of our model for the multi documents. In particular, compared with the original model for SDS, there are several minor modifications. Firstly, latent topics are generated for covering the topics of whole relevant documents. In this regard, instead of combining all topic features for the topic representation, we keep each topic feature representation separately to maintain the information. Secondly, the word node and sentence node are generated by a set of relevant documents, which include a list of sentences and a set of unique words from multiple documents instead of a single document in the SDS problem. Specifically, supporting $D = \{d_1, d_2, \dots, d_n\}$ denotes the set of each input multi documents, the output sentence-topic representation s_i is re-calculated as follows:

$$\begin{aligned} \bar{H}_{T_D} &= FFN(\parallel_{k=1}^K H_{T_k}) \\ H_{s_i, T_D} &= \sigma(FFN(\bar{H}_{T_D} \oplus H_{s_i})) \end{aligned} \quad (11)$$

where K denotes the number of topics for the multiple documents and \parallel represents the concatenation

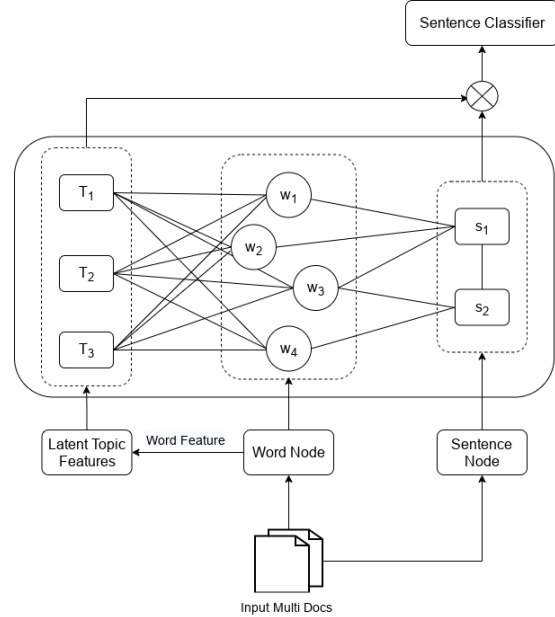


Figure 3: Overview of TopicHeterGraphSum for multi documents summarization.

operation. Sequentially, the output matrix is transformed into the vector by a flattened layer for the final classification.

4 Experiment

4.1 Experimental Setting

Datasets: Two benchmark datasets are considered for the evaluation such as CNN/DailyMail (Nalapaty et al., 2016) (single document dataset) and Multi-News (Fabbri et al., 2019) (multi documents dataset). For the data processing, we use the same split as the work in (Wang et al., 2020). Specifically, the statistics of two benchmark datasets are illustrated in Tab. 1.

	CNN/Daily Mail	Multi-News
Train	287,227	44,972
Val	13,368	5,622
Test	11,490	5,622
Vocab Size	717,951	666,515

Table 1: Statistics of the evaluated datasets.

Hyperparameter Setting: Regarding the word node generation, the vocabulary is limited to 50,000. The tokens are initialized with 100 dimensions using Glove embeddings (Pennington et al., 2014). The

multi-head of the GAT layer for word-to-sentence and word-to-topic is set to 4 and 1, respectively. The maximum number of sentences in each document is set to 100. The initialized dimensions of sentence embedding and topic embedding are set to 128 and 100, respectively. The dimension of the final output representation of all models is set to 64. Regarding the decoder process, we select top-3 for CNN/DailyMail and top-11 sentences for Multi-News following the performance of the validation set. Furthermore, n-gram Blocking (Liu and Lapata, 2019) is also taken into account to improve the performance. Specifically, we vary the values of n-gram from 3 to 6 in order to determine the best results. The number of latent topics is set to 5 both single and multi documents, respectively.

Baseline: For the SDS problem, several state-of-the-art non-pretrained models, which have recently introduced, are taken into account such as BAND-DITSUM (Dong et al., 2018), JECS (Xu and Durrett, 2019), HER (Luo et al., 2019), Topic GraphSum (non-pretrained version) (Cui et al., 2020), HSG (Wang et al., 2020), and Multi GraS (Jing et al., 2021). Regarding the MDS problem, the most recent state-of-the-art methods using pre-trained models are proposed for abstractive summarization (Xiao et al., 2021). Consequently, we follow the reports in (Wang et al., 2020) to take the comparison. The proposed model, TopicHeterGraphSum (THGS) is executed with two versions, by adopting two clustering algorithms for initialized latent topic features, such as K-Mean (THGS-KMean) and GMM (THGS-GMM).

4.2 Main Results

Single Document Summurization: Table 2 shows the results of our evaluation on the CNN/DailyMail dataset. As result, our model outperforms the state-of-the-art models in this research field. Specifically, the results show that initialized features of latent topics by using GMM achieves better results than K-Mean.

Multi Document Summurization: Table 3 shows the results on the Multi-News dataset for the MDS problem. Specifically, the results indicate that enriching global information by using latent topics is able to improve the performance of the MDS problem.

Model	R-1	R-2	R-L
BANDITSUM	41.50	18.70	37.60
JECS	41.70	18.50	37.90
HER	42.30	18.90	37.90
Topic-GraphSum	41.93	19.15	38.22
HSG	42.95	19.76	39.23
Multi-GraS	43.16	20.14	39.49
THGS-Kmean (ours)	43.25	20.20	39.62
THGS-GMM (ours)	43.28	20.31	39.67

Table 2: Results on CNN/DailyMail dataset. Report results are obtained from respective papers. Bold texts indicate the best results in each column.

Model	R-1	R-2	R-L
TextRank	41.95	13.86	38.07
LexRank	41.77	13.81	37.87
PG-BRNN	45.27	15.32	41.38
Hi-MAP	45.21	16.29	41.39
HDSG	46.05	16.35	42.08
THGS-Kmean (ours)	46.60	16.81	42.63
THGS-GMM (ours)	46.66	16.90	42.73

Table 3: Results on Multi-News dataset. Reported results are obtained from (Wang et al., 2020). Bold texts indicate the best results in each column.

4.3 Results with Varying Hyperparameters

We execute experiments to evaluate the impact of important hyperparameters on the performance of the proposed model. Due to the limitation of our resources, we mainly focus on the Multi-News datasets.

Iteration Numbers: In order to select the best number of iterations of GAT, we measure the performance of different numbers of iterations on the validation sets. Table 4 shows the results in which the number of iterations ranges from 1 to 3. Accord-

Iteration	R-1	R-2	R-L
1	46.16	16.57	42.05
2	46.66	16.90	42.73
3	46.53	16.87	42.67

Table 4: Results on Multi-News dataset with different number of iteration of GAT.

ingly, the larger number of iterations does not make a further substantial gain. Therefore, we select the number of iterations equal to 2 for the Multi-News

datasets. In the case of the CNN/DailyMail dataset, the number of iterations is set to 3, corresponding to the best performance.

Number of Sentences for Decode: Normally, the number of sentences for decoding is determined based on the average length of the human-written summaries. Accordingly, the average length of CNN/DailyMail and Multi-News are 3 and 9, respectively. However, we take this issue into account by varying the number of sentences. Table 5 shows the results of various numbers of sentences for decoding the Multi-News dataset. As result, we select

Num. of Sent.	R-1	R-2	R-L
9	46.16	16.57	42.05
10	46.55	16.80	42.53
11	46.66	16.90	42.73
12	46.53	16.90	42.71
13	46.53	16.89	42.62

Table 5: Results on Multi-News dataset with different number of sentence for decoding.

the top-11 sentences for Multi-News datasets following the performance of the validation set.

N-gram Blocking: Trigram blocking is adopted to reduce redundancy for the decode process (Liu and Lapata, 2019). In this study, we vary the values of the n-gram from 3 to 6 to determine the best value. Accordingly, we set the value of n equal to 5 for the

n-gram	R-1	R-2	R-L
3	46.04	16.07	42.04
4	46.61	16.73	42.65
5	46.66	16.90	42.73
6	46.60	16.93	42.68

Table 6: Results on Multi-News dataset with different numbers of n-gram blocking.

Multi-News dataset, which provides the best performance in terms of R-1 and R-L.

5 Conclusion and Future Work

We introduce a new method for the EDS problem by enriching global information using latent topics. Specifically, we first generate the latent topics using well-known clustering algorithms. The outputs are put into a HeterGNN as additional nodes for

enriching the feature representations of sentences. The experiment on two benchmark datasets such as CNN/DailyMail and Multi-News of both SDS and MDS indicates the promising results of the proposed method in this research field. A major drawback of this study is that we use the same latent topic aggregation method for both SDS and MDS problems. Specifically, latent topics are suitable for the MDS problem, which has been proved in this study. However, since the complex relationship between word node and sentence node in multiple documents, a further investigation on exploiting the relationship between two types of nodes across multiple documents is able to improve the performance. Therefore, further exploitation of topic aggregation for the MDS problem is considered as our future work regarding this study.

References

- Jianpeng Cheng and Mirella Lapata. 2016. Neural summarization by extracting sentences and words. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 1: Long Papers*. The Association for Computer Linguistics.
- Arman Cohan, Franck Dernoncourt, Doo Soon Kim, Trung Bui, Seokhwan Kim, Walter Chang, and Nazli Goharian. 2018. A discourse-aware attention model for abstractive summarization of long documents. In Marilyn A. Walker, Heng Ji, and Amanda Stent, editors, *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 2 (Short Papers)*, pages 615–621. Association for Computational Linguistics.
- Peng Cui, Le Hu, and Yuanchao Liu. 2020. Enhancing extractive text summarization with topic-aware graph neural networks. In Donia Scott, Núria Bel, and Chengqing Zong, editors, *Proceedings of the 28th International Conference on Computational Linguistics, COLING 2020, Barcelona, Spain (Online), December 8-13, 2020*, pages 5360–5371. International Committee on Computational Linguistics.
- Yue Dong, Yikang Shen, Eric Crawford, Herke van Hoof, and Jackie Chi Kit Cheung. 2018. Banditsum: Extractive summarization as a contextual bandit. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun’ichi Tsujii, editors, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*,

- Brussels, Belgium, October 31 - November 4, 2018, pages 3739–3748. Association for Computational Linguistics.
- Günes Erkan and Dragomir R. Radev. 2004. Lexrank: Graph-based lexical centrality as salience in text summarization. *J. Artif. Intell. Res.*, 22:457–479.
- Alexander R. Fabbri, Irene Li, Tianwei She, Suyi Li, and Dragomir R. Radev. 2019. Multi-news: A large-scale multi-document summarization dataset and abstractive hierarchical model. In Anna Korhonen, David R. Traum, and Lluís Màrquez, editors, *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 1074–1084. Association for Computational Linguistics.
- Ruipeng Jia, Yanan Cao, Hengzhu Tang, Fang Fang, Cong Cao, and Shi Wang. 2020. Neural extractive summarization with hierarchical attentive heterogeneous graph network. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 3622–3631. Association for Computational Linguistics.
- Baoyu Jing, Zeyu You, Tao Yang, Wei Fan, and Hanghang Tong. 2021. Multiplex graph neural network for extractive text summarization. *CoRR*, abs/2108.12870.
- Thomas N. Kipf and Max Welling. 2017. 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.
- Yang Liu and Mirella Lapata. 2019. Hierarchical transformers for multi-document summarization. In Anna Korhonen, David R. Traum, and Lluís Màrquez, editors, *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 5070–5081. Association for Computational Linguistics.
- Ling Luo, Xiang Ao, Yan Song, Feiyang Pan, Min Yang, and Qing He. 2019. Reading like HER: human reading inspired extractive summarization. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan, editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, pages 3031–3041. Association for Computational Linguistics.
- Rada Mihalcea and Paul Tarau. 2004. Textrank: Bringing order into text. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing, EMNLP 2004, A meeting of SIGDAT, a Special Interest Group of the ACL, held in conjunction with ACL 2004, 25-26 July 2004, Barcelona, Spain*, pages 404–411. ACL.
- Ramesh Nallapati, Bing Xiang, and Bowen Zhou. 2016. Sequence-to-sequence rnn for text summarization. *CoRR*, abs/1602.06023.
- Thong Nguyen, Anh Tuan Luu, Truc Lu, and Tho Quan. 2021. Enriching and controlling global semantics for text summarization. *CoRR*, abs/2109.10616.
- Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. Glove: Global vectors for word representation. In Alessandro Moschitti, Bo Pang, and Walter Daelemans, editors, *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar; A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 1532–1543. ACL.
- Suzanna Sia, Ayush Dalmaia, and Sabrina J. Mielke. 2020. Tired of topic models? clusters of pretrained word embeddings make for fast and good topics too! In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 1728–1736. Association for Computational Linguistics.
- Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. 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.
- Danqing Wang, Pengfei Liu, Yining Zheng, Xipeng Qiu, and Xuanjing Huang. 2020. Heterogeneous graph neural networks for extractive document summarization. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel R. Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 6209–6219. Association for Computational Linguistics.
- Wen Xiao, Iz Beltagy, Giuseppe Carenini, and Arman Cohan. 2021. PRIMER: pyramid-based masked sentence pre-training for multi-document summarization. *CoRR*, abs/2110.08499.
- Jiacheng Xu and Greg Durrett. 2019. Neural extractive text summarization with syntactic compression. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan, editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong*

- Kong, China, November 3-7, 2019, pages 3290–3301. Association for Computational Linguistics.
- Jiacheng Xu, Zhe Gan, Yu Cheng, and Jingjing Liu. 2020. Discourse-aware neural extractive text summarization. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel R. Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 5021–5031. Association for Computational Linguistics.
- Michihiro Yasunaga, Rui Zhang, Kshitijh Meelu, Ayush Pareek, Krishnan Srinivasan, and Dragomir R. Radev. 2017. Graph-based neural multi-document summarization. In Roger Levy and Lucia Specia, editors, *Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017), Vancouver, Canada, August 3-4, 2017*, pages 452–462. Association for Computational Linguistics.
- Xingxing Zhang, Furu Wei, and Ming Zhou. 2019. HIBERT: document level pre-training of hierarchical bidirectional transformers for document summarization. In Anna Korhonen, David R. Traum, and Lluís Màrquez, editors, *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 5059–5069. Association for Computational Linguistics.