

Discourse-Aware Unsupervised Summarization of Long Scientific Documents

Yue Dong*

MILA/McGill University
Montreal, QC, Canada
yue.dong2
@mail.mcgill.ca

Andrei Mircea*

MILA/McGill University
Montreal, QC, Canada
andrei.romascanu
@mail.mcgill.ca

Jackie C. K. Cheung

MILA/McGill University
Montreal, QC, Canada
jcheung
@cs.mcgill.ca

Abstract

We propose an unsupervised graph-based ranking model for extractive summarization of long scientific documents. Our method assumes a two-level hierarchical graph representation of the source document, and exploits asymmetrical positional cues to determine sentence importance. Results on the PubMed and arXiv datasets show that our approach¹ outperforms strong unsupervised baselines by wide margins in automatic metrics and human evaluation. In addition, it achieves performance comparable to many state-of-the-art supervised approaches which are trained on hundreds of thousands of examples. These results suggest that patterns in the discourse structure are a strong signal for determining importance in scientific articles.

1 Introduction

Single document summarization aims at shortening a text and preserving the most important ideas of the source document. While abstractive strategies generate summaries with novel words, extractive strategies select sentences from the source to form a summary (Nenkova et al., 2011). Despite recent advances in abstractive summarization, extractive models are still attractive in cases where faithfully preserving the original text is the priority. For example, legal arguments can hinge on the exact wording of a contract (Farzindar and Lapalme, 2004), and ensuring the factual correctness of a summary can be critical in the health or scientific domains, which is a known weakness of current abstractive methods (Kryściński et al., 2019).

Supervised neural-based models have been the dominant paradigm in recent extractive systems, at least for *short news summarization* (Nallapati et al.,

*Equal contribution.

¹Link to our code: <https://github.com/mirandrom/HipoRank>.

Introduction	anxiety affects quality of life in those living with parkinson’s disease (pd) more so than overall cognitive status, motor deficits, apathy, and depression.
Introduction	although anxiety and depression are often related and coexist in pd patients, recent research suggests that anxiety rather than depression is the most prominent and prevalent mood disorder in pd.
Related Work	furthermore, since previous work, albeit limited, has focused on the influence of symptom laterality on anxiety and cognition, we also explored this relationship .
Methodology	this study is the first to directly compare cognition between pd patients with and without anxiety.
Result	the findings confirmed our hypothesis that anxiety negatively influences attentional set-shifting and working memory in pd.
Result	moreover, anxiety has been suggested to play a key role in freezing of gait (fog), which is also related to attentional set-shifting.
Future work	s. future research should examine the link between anxiety, set-shifting, and fog, in order to determine whether treating anxiety might be a potential therapy for improving fog.

Table 1: Example of a PubMed article’s summary produced by our model HIPORANK. The hierarchical and directed graph combined with discourse-aware edge weighting allow HIPORANK to generate summaries that cover topics from different sections of the scientific article.

2017; Dong et al., 2018; Zhou et al., 2018; Liu and Lapata, 2019; Narayan et al., 2018b; Zhang et al., 2019b). These models usually employ the encoder-decoder structure and have achieved promising performance on news datasets such as CNN/DailyMail (Hermann et al., 2015), and NYT (Sandhaus, 2008).

However, these models cannot easily be adapted to out-of-domain data that have greater length and fewer training examples such as scientific article summarization (Xiao and Carenini, 2019) due to

two significant limitations. First, they require large domain-specific training pairs of source documents and gold-standard summaries, which are often not available or feasible to create (Zheng and Lapata, 2019). Second, the typical setup of using a token-level encoder-decoder with an attention mechanism does not scale well to longer documents (Shao et al., 2017), as the number of attention computations is quadratic with respect to the number of tokens in the input document.

We instead explore *unsupervised* approaches to address these challenges on long document summarization. We show that a simple unsupervised graph-based ranking model combined with proper sophisticated modelling of discourse information as an inductive bias can achieve unreasonable effectiveness in selecting important sentences from long scientific documents.

For the choice of unsupervised graph-based ranking model, we follow the paradigm of LexRank (Erkan and Radev, 2004) and PACSUM (Zheng and Lapata, 2019). In these methods, sentences are nodes and weighted edges represent the degree of similarity between sentences. Summary generation is formulated as a node selection problem, in which nodes (i.e., sentences) that are semantically similar to other nodes are chosen to be included in the final summary. In other words, they determine node importance by defining a notion of centrality in the graph.

In addition, we augment the document graph with directionality and hierarchy to reflect the rich discourse structure of long scientific documents. In particular, our method relies on two insights about the discourse structure of long scientific documents. The first is that important information typically occurs at the start and end of sections; i.e., they tend to appear near section boundaries (Baxendale, 1958; Lin and Hovy, 1997; Teufel, 1997). We implement this using an asymmetric edge weighting function in a *directed graph* which considers the distance of a sentence to a boundary. The second is that most sentences across section boundaries are unlikely to interact significantly with each other (Xiao and Carenini, 2019). We implement this insight by injecting *hierarchies* into our model, introducing section-level representations as graph nodes in addition to sentence nodes. By doing so, we convert a flat graph into a hierarchical non-fully-connected graph, which has two advantages: 1) reduced computational cost and 2) pruning of distracting weak

connections between sentences across different sections.

We call our approach **Hierarchical and Positional Ranking** model (HIPORANK) and evaluate it on summarizing long scientific articles from PubMed and arXiv (Cohan et al., 2018). Empirical results show that our method significantly improves performance over previous unsupervised models (Zheng and Lapata, 2019; Erkan and Radev, 2004) in both automatic and human evaluation. In addition, our simple unsupervised approach achieves performance comparable to many expensive state-of-the-art supervised neural models that are trained on hundreds of thousands of examples of long document pairs (Xiao and Carenini, 2019; Subramanian et al., 2019). This suggests that patterns in the discourse structure are highly useful for determining sentence importance in long scientific articles, and that explicitly building in biases inspired by this structure is a viable strategy for building summarization systems.

2 Related Work

2.1 Extractive Summarization

Traditional extractive summarization methods are mostly unsupervised (Radev et al., 2000; Lin and Hovy, 2002; Wan, 2008; Wan and Yang, 2008; Hirao et al., 2013; Parveen et al., 2015; Yin and Pei, 2015; Li et al., 2017; Zheng and Lapata, 2019), utilizing a notion of sentence importance based on n-gram overlap with other sentences and frequency information (Nenkova and Vanderwende, 2005), relying on graph-based methods for sentence ranking (Erkan and Radev, 2004; Mihalcea and Tarau, 2004), or performing keyword extraction combined with submodular maximization (Tixier et al., 2017; Shang et al., 2018).

With the development of large-scale summarization datasets such as CNN/DailyMail (Hermann et al., 2015), NYT (Sandhaus, 2008), Newsroom (Grusky et al., 2018) and XSum (Narayan et al., 2018a), along with advancements in deep neural-based architectures, modern supervised neural network-based methods that employ encoder-decoder framework have become increasingly popular. These models have been proposed with extractive strategies (Cheng and Lapata, 2016; Nallapati et al., 2017; Wu and Hu, 2018; Dong et al., 2018; Zhou et al., 2018; Narayan et al., 2018b); abstractive strategies (See et al., 2017; Chen and Bansal,

2018; Gehrmann et al., 2018; Dong et al., 2019; Zhang et al., 2019a; Lewis et al., 2019); and hybrid strategies (Hsu et al., 2018; Bae et al., 2019; Moroshko et al., 2019).

More recently, extractive approaches leveraging transformer architectures (Vaswani et al., 2017) and their pretrained counterparts (Devlin et al., 2019; Lewis et al., 2019; Zhang et al., 2019a; Dong et al., 2019) have achieved state-of-the-art performances on the CNN/DailyMail news benchmark dataset (Zhang et al., 2019b; Liu and Lapata, 2019; Zhong et al., 2019). Furthermore, pretrained transformer models also provide better sentence representations for unsupervised summarization methods. For instance, PACSUM (Zheng and Lapata, 2019), a directed graph-based unsupervised model that utilizes BERT-based sentence representations, achieved comparable performance to supervised models on the CNN/DailyMail and NYT datasets.

2.2 Extractive Summarization of Long Scientific Papers

Despite the success of deep neural-based models on news summarization, these approaches typically face challenges when applied to long documents such as scientific articles. Furthermore, these approaches are often blind to the topical information resulting from the structured sections in scientific articles (Xiao and Carenini, 2019). Two recent neural supervised models address these issues. Subramanian et al. (2019) used the introduction section as a proxy for the whole document, while Xiao and Carenini (2019) divided articles into sections and used non-auto-regressive approaches to model global and local information.

Besides neural approaches, most previous scientific article summarization systems employ traditional supervised machine learning algorithms with surface features as input (Xiao and Carenini, 2019). Surface features such as sentence position, sentence and document length, keyphrase score, and fine-grain rhetorical categories are often combined with Naive Bayes (Teufel and Moens, 2002), CRFs and SVMs (Liakata et al., 2013), LSTM and MLP (Collins et al., 2017) for extractive summarization over long scientific articles. To the best of our knowledge, the only unsupervised extractive summarization model for long scientific documents relies on citation networks (Qazvinian and Radev, 2008; Cohan and Goharian, 2015), by extracting citation-contexts from citing articles and ranking

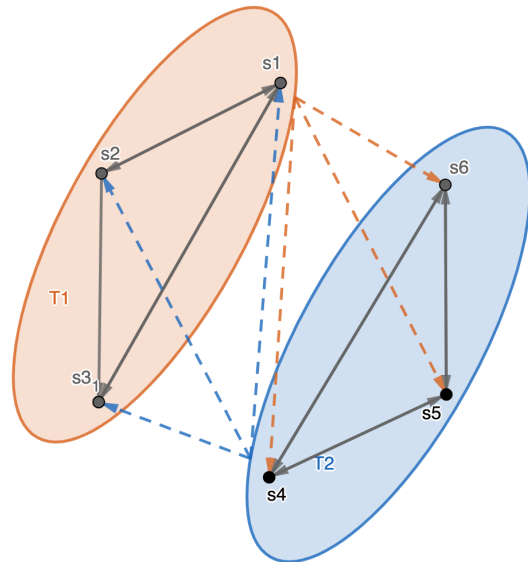


Figure 1: Example of a hierarchical document graph constructed by our approach on a toy document that contains two sections $\{T_1, T_2\}$, each containing three sentences for a total of six sentences $\{s_1, \dots, s_6\}$. Each double-headed arrow represents two edges with opposite directions. The solid and dashed arrows indicate intra-section and inter-section connections respectively. When compared to the flat fully-connected graph of traditional methods, our use of hierarchy effectively reduces the number of edges from 60 to 24 in this example.

these sentences to form the final summary. Our proposed method is different from their settings, where we perform single document summarization based on the long source article.

3 Method

Our proposed method combines simple graph-based ranking algorithms with a two-level hierarchical model of the rich discourse structures of long scientific documents (Teufel, 1997; Xiao and Carenini, 2019). We incorporate this discourse information into the graph as inductive biases through the construction of a *directed hierarchical graph* for document representation (Figure 1 and Section 3.2) and through the asymmetric edge weighting of edges with boundary functions (Section 3.3).

3.1 Graph-based Ranking Algorithm

Graph-based ranking algorithms for summarization represent a document as a graph $G = (V, E)$, where V is the set of vertices that represent sentences or other textual units in the document, and

E is the set of edges that represent interactions between sentences. The directed edge e_{ij} from node v_i to node v_j is typically weighted by $w_{ij} = f(\text{sim}(v_i, v_j))$, where sim is a measure of similarity between two nodes (e.g. cosine distance between their distributed representations), and f can be an additional weighting function. These algorithms select the most salient sentences from V based on the assumption that sentences that are similar to a greater number of other sentences capture more important content and therefore are more informative.

3.2 Hierarchical Document Graph Creation

To create a hierarchical document graph, we first split a document into its sections, then into sentences². To create the hierarchy, we allow two levels of connections in our hierarchical graph: intra-sectional connections and inter-sectional connections as shown in Figure 1.

Intra-sectional connections aim to model the *local* importance of a sentence within its section. It implements the idea that a sentence that is similar to a greater number of other sentences in the same topic/section should be more important. This is realized in our fully-connected subgraph for an arbitrary section I , where we allow *sentence-sentence* edges for all sentence nodes within the same section.

Inter-sectional connections aim to model the *global* importance of a sentence with respect to other topics/sections in the document, as a sentence that is similar to a greater number of other topics is deemed more important. However, calculating sentence-sentence connections across different sections is computationally expensive and may also suffer from performance degradation due to weak edges between sentences that are unrelated as a result of being from different sections (Mihalcea and Tarau, 2004). To address these issues, We introduce section nodes on top of sentence nodes to form a hierarchical graph. For inter-section connections, we only allow *section-sentence* edges for modeling the global information. This choice makes our approach more computationally efficient while greatly limiting the number of irrelevant inter-section edges that arise from the fact that sections in scientific documents typically have independent

²Our approach is agnostic to the sentence/section splitting method. In our experiments, articles in the datasets are already split into sections and sentences.

topics (Xiao and Carenini, 2019). In contrast, traditional graph-based ranking algorithms have a flat fully-connected graph document with no sections.

3.3 Asymmetric Edge Weighting by Boundary functions

To calculate the weight of an edge, we first measure similarity between a sentence-sentence pair $\text{sim}(v_j^I, v_i^I)$ and a section-sentence pair $\text{sim}(v_j^I, v_i^I)$. While our method is agnostic to the measure of similarity, we use cosine similarity with different vector representations in our experiments, averaging a section’s sentence representations to obtain its own.

While the similarities of two graph nodes are symmetric, one may be more salient than the other when considering their discourse structures (Baxendale, 1958; Teufel, 1997). Based on these discourse hypotheses of long scientific documents, we capture this asymmetry by making our hierarchical graph *directed* and inject *asymmetric* edge weighting over intra-section and inter-section connections.

Asymmetric edge weighting over sentences

Our asymmetric edge weighting is based on the hypothesis that important sentences are near the boundaries (start or end) of a text (Baxendale, 1958). We reflect this hypothesis by defining a *sentence boundary function* d_b over sentences v_i^I in section I such that sentences closer to the section’s boundaries are more important:

$$d_b(v_i^I) = \min(x_i^I, \alpha(n^I - x_i^I)), \quad (1)$$

where n^I is the number of sentences in section I and x_i^I represents sentence i ’s position in the section I . $\alpha \in \mathbb{R}^+$ is a hyper-parameter that controls the relative importance of the start or end of a section or document.

The sentence boundary function allow us to incorporate directionality in our edges, and weight edges differently depending on if they are incident to a more important or less important sentence in the same section. Concretely, we define the weight w_{ji}^I for intra-section edges (incoming edges for i) as:

$$w_{ji}^I = \begin{cases} \lambda_1 * \text{sim}(v_j^I, v_i^I), & \text{if } d_b(v_i^I) \geq d_b(v_j^I), \\ \lambda_2 * \text{sim}(v_j^I, v_i^I), & \text{if } d_b(v_i^I) < d_b(v_j^I) \end{cases} \quad (2)$$

where $\lambda_1 < \lambda_2$ such that an edge e_{ji} incident to i is weighted more if i is closer to the text bound-

ary than j . Edges with a weight below a certain threshold β can be pruned (i.e., set to 0).

Asymmetric edge weighting over sections

Similarly, to reflect the hierarchy hypothesis over long scientific documents proposed by Teufel (1997), we also define a *section boundary function* d_b to reflect that sections near a document’s boundaries are more important:

$$d_b(v^I) = \min(x^I, \alpha(N - x^I)), \quad (3)$$

where N is the number of sections in the document and x^I represents section I ’s position in the document.

This section boundary function allows us to inject asymmetric edge weighting w_i^{JI} to inter-section edges:

$$w_i^{JI} = \begin{cases} \lambda_1 * \text{sim}(v^J, v_i^I), & \text{if } d_b(v^I) \geq d_b(v^J). \\ \lambda_2 * \text{sim}(v^J, v_i^I), & \text{if } d_b(v^I) < d_b(v^J) \end{cases} \quad (4)$$

where $\lambda_1 < \lambda_2$ such that an edge e_i^{JI} incident to $i \in I$ is weighted more if section I is closer to the text boundary than section J .

3.4 Importance Calculation

We compute the overall importance of sentence v_i^I as the weighted sum of its inter-section and intra-section centrality scores:

$$c(v_i^I) = \mu_1 \cdot c_{\text{inter}}(v_i^I) + c_{\text{intra}}(v_i^I) \quad (5)$$

$$c_{\text{intra}}(v_i^I) = \sum_{v_j^I \in I} \frac{w_{ji}^I}{|I|} \quad (6)$$

$$c_{\text{inter}}(v_i^I) = \sum_{v^J \in D} \frac{w_i^{JI}}{|D|},$$

where I is the set of sentences neighbouring v_i^I and D is the set of neighbouring sections in the hierarchical document graph; μ_1 is a weighting factor for inter-section centrality.

3.5 Summary Generation

Lastly, we generate a summary by greedily extracting sentences with the highest importance scores until a predefined word-limit L is passed. Most graph-based ranking algorithms recompute importance after each sentence is extracted in order to prevent content overlap. However, we find that the

Dataset	# docs	avg. doc. len.	avg. summ. len.
CNN	92K	656	43
Daily Mail	219K	693	52
NYT	655K	530	38
PubMed	133K	3,016	203
arXiv	215K	4,938	220

Table 2: Dataset statistics on news articles (CNN, DailyMail, and NYT) and long scientific documents (PubMed and arXiv).

asymmetric edge scoring functions in (2) and (4) naturally prevent redundancy, because similar sentences have different boundary positional scores. Our method thus successfully extracts diverse sentences without recomputing importance.

4 Experimental Setup

This section describes the datasets, the hyperparameter choices, the baseline models, and the evaluation metrics used in the experiments.

4.1 Datasets

Our experiments are conducted on PubMed and arXiv (Cohan et al., 2018), two large-scale datasets of long and structured scientific articles with abstracts as summaries. The average source article length is four to seven times longer than popular news benchmarks (Table 2), making them ideal candidates to test our method.

4.2 Implementation Details

Our model’s hyperparameters for testing are chosen from the ablation studies on the validation sets. The test results are reported with the following hyperparameter settings: $\lambda_1 = 0.0$, $\lambda_2 = 1.0$, $\alpha = 1.0$, with $\mu_1 = 0.5$ for PubMed and $\mu_1 = 1.0$ for arXiv. We fix λ_2 to 1 and the choices of $\lambda_1 \in \{-0.2, 0, 0.5\}$. represent whether the edge between a less boundary-important sentence and a more boundary-important sentence is 1) negatively weighted, 2) pruned, or 3) down-weighted. $\lambda_1 < \lambda_2$ such that an edge e_{ji} incident to i is weighted more if i is closer to the text boundary than j . $\alpha \in \{0, 0.5, 0.8, 1.0, 1.2\}$ controls the relative importance of the start or end of a section or document. $\mu_1 \in \{0.5, 1.0, 1.5\}$ controls how much we weigh intra-section sentence importance vs. inter-section sectional importance.

For each dataset, we experimented with different pretrained distributional sentence representation models. The dimension of sentence representations is model-dependent (details in Section

6.2). We used the publicly released BERT model³ (Devlin et al., 2019), PACSUM BERT model⁴ (Zheng and Lapata, 2019), SentBERT and SentRoBERTa⁵ (Reimers and Gurevych, 2019), and BioMed word2vec representations⁶ (Moen and Ananiadou, 2013). A section’s representation is calculated as the average of its sentences’ representations. The similarity between sentences or sections is defined to be the cosine similarity between the distributed representations.

4.3 Baselines

We compare our approach with previous unsupervised and supervised models in extractive summarization. In addition, we also compare it with recent neural abstractive approaches for completeness.

For unsupervised extractive summarization models, we compare with SumBasic (Vanderwende et al., 2007), LSA (Steinberger and Jezek, 2004), LexRank (Erkan and Radev, 2004) and PACSUM (Zheng and Lapata, 2019). For supervised neural extractive summarization models, we compare with a vanilla encoder-decoder model (Cheng and Lapata, 2016), SummaRuNNer (Nallapati et al., 2017), GlobalLocalCont (Xiao and Carenini, 2019), Sent-CLF and Sent-PTR (Subramanian et al., 2019). We also compare with neural abstractive summarization models as reported in Xiao and Carenini (2019): Attn-Seq2Seq (Nallapati et al., 2016), Pntr-Gen-Seq2Seq (See et al., 2017) and Discourse-aware (Cohan et al., 2018). In addition, we report the lead baseline that selects the first k tokens as a summary ($k = 203, = 220$ for PubMed and arXiv respectively). Lastly, we report baselines for an Oracle summarizer (Nallapati et al., 2017).

4.4 Evaluation Methods

We evaluate our method with automatic evaluation metrics - ROUGE F1 scores (Lin, 2004). ROUGE-1 and ROUGE-2 compute unigram and bigram overlaps between system summaries and reference summaries, while ROUGE-L computes the longest common sub-sequence of the two.

In addition, we design a human evaluation experiment (details in Section 5.2) to compare our model with the best unsupervised summarization model - PACSUM (Zheng and Lapata, 2019). As far as we know, we are the first to perform human evaluation

³<https://github.com/huggingface/transformers>

⁴<https://github.com/mswellhao/PACSUM>

⁵<https://github.com/UKPLab/sentence-transformers>

⁶<http://bio.nlplab.org/word-vectors>

Model	ROUGE-1	ROUGE-2	ROUGE-L
Lead	35.63	12.28	25.17
Oracle (ROUGE-2, F1)	55.05	27.48	38.66
Supervised Abstractive			
Attn-Seq2Seq (2016)	31.55	8.52	27.38
Pntr-Gen-Seq2Seq (2017)	35.86	10.22	29.69
Discourse-aware (2018)	38.93	15.37	35.21
Supervised Extractive			
Cheng & Lapata (2016)	43.89	18.53	30.17
SummaRuNNer (2017)	43.89	18.78	30.36
GlobalLocalCont (2019)	44.85	19.70	31.43
Sent-CLF (2019)	45.01	19.91	41.16
Sent-PTR (2019)	43.30	17.92	39.47
Unsupervised Extractive			
SumBasic (2007)	37.15	11.36	33.43
LSA (2004)	33.89	9.93	29.70
LexRank (2004)	39.19	13.89	34.59
PACSUM (2019)	39.79	14.00	36.09
HIPORANK (ours)	43.58	17.00	39.31

Table 3: Test set results on PubMed (ROUGE F1).

on the 2018 PubMed and arXiv datasets (Cohan et al., 2018). Human evaluation over long scientific articles require annotators to comprehend a full domain-specific long article and compare multiple summaries for quality evaluation. Due to the challenging nature of the task, previous papers choose to skip it and purely rely on automatic evaluations to judge the system performance.

5 Results

5.1 Automatic Evaluation Results

Tables 3 and 4 summarize our automatic evaluation results on the PubMed and arXiv test sets with the best hyperparameters, as described in Section 4.2.

The first blocks in Table 3,4 include the lead and the oracle baselines. The second and the third blocks in the tables present the results of supervised abstractive models, and of supervised extractive models. ROUGE-2 oracle summaries are used as gold standard summaries for training supervised extractive models, which likely contributes to their better ROUGE-2 scores.

The last blocks compare previous unsupervised models with our approach. Our model outperforms all other unsupervised approaches by wide margins in terms of ROUGE-1,2,L F1 scores on both PubMed and arXiv datasets. We also show that PACSUM is biased towards selecting sentences that

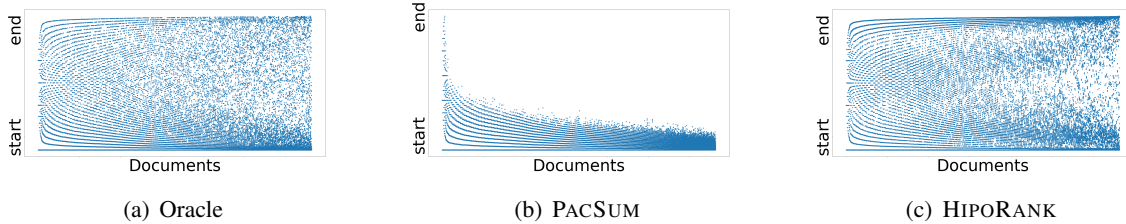


Figure 2: Sentence positions in source document for extractive summaries generated by different models on the PubMed validation set. Documents on the x-axis are ordered by increasing article length from shortest to longest. We also see a similar trend on arXiv (the plots with more details can be found in the appendix).

Model	ROUGE-1	ROUGE-2	ROUGE-L
Lead	33.66	8.94	22.19
Oracle (ROUGE-2, F1)	53.88	23.05	34.90
Supervised Abstractive			
Attn-Seq2Seq (2016)	29.30	6.00	25.56
Pntr-Gen-Seq2Seq (2017)	32.06	9.04	25.16
Discourse-aware (2018)	35.80	11.05	31.80
Supervised Extractive			
Cheng&Lapata (2016)	42.24	15.97	27.88
SummaRuNNer (2017)	42.81	16.52	28.23
GlobalLocalCont (2019)	43.62	17.36	29.14
Sent-CLF (2019)	34.01	8.71	30.41
Sent-PTR (2019)	42.32	15.63	38.06
Unsupervised Extractive			
SumBasic (2007)	29.47	6.95	26.30
LSA (2004)	29.91	7.42	25.67
LexRank (2004)	33.85	10.73	28.99
PACSUM (2019)	38.57	10.93	34.33
HIPORank (ours)	39.34	12.56	34.89

Table 4: Test set results on arXiv (ROUGE F1).

appear at the beginning of a document while our method selects sentences in every section and near the article boundaries, similar to the oracle (Figure 2). This overlap with gold standard summaries suggests our use of discourse structure and hierarchy plays a significant role in our method’s performance.

Interestingly, despite limited access to only the validation set for hyperparameter tuning, our method achieves performance scores that are competitive with supervised models that require hundreds of thousands of training examples, outperforming almost all abstractive and extractive models on ROUGE-L. This suggests that our discourse-aware unsupervised model is surprisingly effective at selecting salient sentences in long scientific document and perhaps should be used as a strong

Model	Content-coverage	Importance
PACSUM	30.52	48.70
HIPORank (ours)	42.13	59.06

Table 5: Human evaluation results on 20 sampled reference summaries with 281 system summary sentences from PubMed. Each reference summary-sentence pair is annotated by two annotators with an average annotator agreement of 73.24%. The results are averaged across 127 sentences from HipoRank and 154 sentences from state-of-the-art unsupervised extractive summarization system PACSUM (Zheng and Lapata, 2019)..

baseline to accessing the merits of supervised approaches for learning content beyond discourse.

5.2 Human Evaluation

We asked the human judges⁷ to read the reference summary⁸ (abstract) and present extracted sentences from different summarization systems in a random and anonymized order. The judges are asked to evaluate the system summary sentence according to two criteria: 1) *content coverage* (whether the presented sentence contains content from the abstract); and 2) *importance* (whether the presented sentence is important for a goal-oriented reader even if isn’t in the abstract (Lin and Hovy, 1997)).

Table 5 presents the human evaluation results. HIPORank is shown to be significantly better than PACSUM in both content coverage and importance ($p = 0.002$ and $p = 0.007$ with Mann-Whitney U tests, respectively). We also measure inter-rater reliability using Fleiss’ κ (46.56 for *content-coverage* and 41.37 for *importance*). These results help sup-

⁷All judges are native English speakers with at least a bachelor’s degree and experience in scientific research. We compensated the judges at an hourly rate of \$20.

⁸We made the decision to not present the whole article, which would create a large cognitive burden on judges and incentivize them to take shortcuts.

Model	ROUGE-1	ROUGE-2	ROUGE-L
HIPORANK + Different Positional Functions			
lead	37.43	12.13	33.68
undirected	40.66	13.41	36.55
boundary-distance (ours)	43.42	16.76	39.23
HIPORANK + Different Hierarchical Functions			
w.o. hierarchy	41.88	15.39	37.91
w. hierarchy (ours)	43.42	16.76	39.23

Table 6: Results on the PubMed validation set with different positional function or hierarchical information.

port that our method’s use of hierarchy and discourse structure improves summarization quality.

6 Ablation Studies

6.1 Component-wise Analysis

Table 6 presents the ablation study to assess the relative contributions of the boundary function and the hierarchical information. We keep all the hyperparameters unchanged with respect to the best setting in Section 4.2 and either vary the positional function or the hierarchical structures. We also found that the improvement of each components are stable across all the hyperparameters we tested (more details in the appendix).

The first block of Table 6 reports the ablation results with different positional functions: no positional function (Erkan and Radev, 2004; Mihalcea and Tarau, 2004), lead bias function (Zheng and Lapata, 2019), and our proposed boundary function. We can see that using the wrong positional function hurts the model’s performance when comparing no positional function with lead bias function. Our boundary positional function outperforms the lead or no positional functions significantly.

The second block of Table 6 reports the results with or without the hierarchical structure. We observe that adding the hierarchical information results in a huge performance improvement.

6.2 Effect of Embeddings

To disentangle the effect of sentence representation, we show PubMed test set results of our best model with different sentence embeddings in Table 7. While pretrained transformer models finetuned on sentence similarity improve performance, HIPORANK still consistently outperforms previous state-of-the-art unsupervised models (Table 3) even with random embeddings. These results once

Model	ROUGE-1	ROUGE-2	ROUGE-L
Lead	35.63	12.28	25.17
Oracle	55.05	27.48	38.66
HIPORANK with Different Embeddings			
Random Embedding (d=200)	43.05	16.69	38.63
Biomed-w2v (d=200)	43.70	17.06	39.19
BERT (d=768)	42.91	16.27	38.52
PACSUM-BERT (d=768)	43.58	17.00	39.31
SentBERT (d=768)	43.59	17.08	39.07
SentRoBERTa (d=1024)	43.55	17.06	39.07

Table 7: PubMed test set results with HIPORANK framework and different pretrained sentence and section embeddings.

again suggest that our method’s improvement can indeed be attributed to the use of hierarchy and discourse structure, rather than to the the choice of representations.

6.3 Stability of Hyperparameters

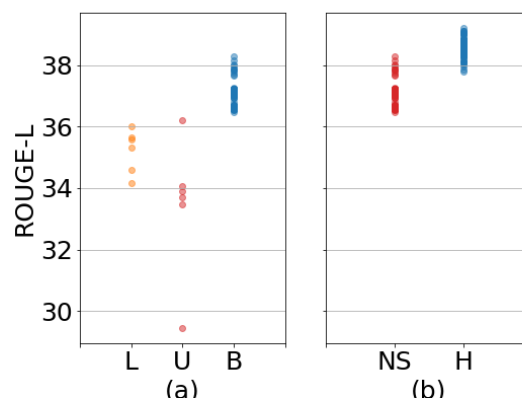


Figure 3: ROUGE-L scores for (a) different positional functions (L=lead, U=undirected, B=boundary) and (b) different graph hierarchies (NS=no.section, H=hierarchical). Each point corresponds to one configuration of the hyperparameter gridsearch described in Section 4.2.

To further inspect our model’s stability across different hyperparameter choices, we conducted fine-grained analysis across all different hyperparameter settings as below.

Stability w.r.t. Discourse Structure To evaluate the impact and the stability of discourse structure informed edge weighting (Section 3.3), we first compared our *boundary* positional function (Eqn. 1,3) to PACSUM’s *lead* positional function, as well as the standard *undirected* approach over different

hyperparameter settings. Figure 3 (a) shows that our method consistently performed better on the PubMed validation set, across *different hyperparameters and embedding models* outlined in Section 4.2.

Stability w.r.t. Hierarchy We then evaluated the effect of adding hierarchy (Section 3.2) on top of our boundary positional function. In addition to decreasing the computational cost, Figure 3 (b) shows that incorporating hierarchy further improved ROUGE-L consistently across *different hyperparameters and embedding models* we tested.

Application to other genres While our work here is focused on long scientific document summarization, we believe that our approach is promising for other genres of text, provided that the right discourse-aware biases are given to the model. Indeed, one version of our model with our proposed boundary function can be seen as a generalization of PACSUM, which achieves state-of-the-art performance on unsupervised summarization of news by exploiting the well known lead bias of news text (Zheng and Lapata, 2019; Grenander et al., 2019). We leave such explorations of adapting HIPORANK to other genres to future work.

7 Conclusion

We presented an unsupervised graph-based model for long scientific document summarization. The proposed approach augments the measure of sentence centrality by inserting directionality and hierarchy in the graph with boundary positional functions and hierarchical topic information grounded in discourse structure. Our simple unsupervised approach with rich discourse modelling outperforms previous unsupervised graph-based summarization models by wide margins and achieves comparable performance to state-of-the-art supervised neural models. This makes our model a lightweight but strong baseline for assessing the performance of expensive supervised approaches for long scientific document summarization.

Acknowledgments

This work is supported by the Natural Sciences and Engineering Research Council of Canada, Compute Canada, and the CIFAR Canada AI Chair program. We would like to thank Hao Zheng, Wen Xiao, and Sandeep Subramanian for useful discussions.

References

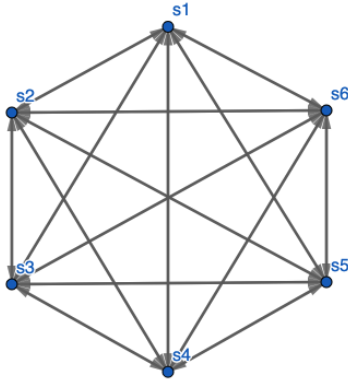
- Sanghwan Bae, Taeuk Kim, Jihoon Kim, and Sang-goo Lee. 2019. Summary level training of sentence rewriting for abstractive summarization. In *Proceedings of the 2nd Workshop on New Frontiers in Summarization*, pages 10–20.
- Phyllis B Baxendale. 1958. Machine-made index for technical literature—an experiment. *IBM Journal of research and development*, 2(4):354–361.
- Yen-Chun Chen and Mohit Bansal. 2018. Fast abstractive summarization with reinforce-selected sentence rewriting. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 675–686.
- 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 (Volume 1: Long Papers)*, pages 484–494, Berlin, Germany. Association for Computational 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 *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 615–621.
- Arman Cohan and Nazli Goharian. 2015. [Scientific article summarization using citation-context and article’s discourse structure](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 390–400, Lisbon, Portugal. Association for Computational Linguistics.
- Edward Collins, Isabelle Augenstein, and Sebastian Riedel. 2017. A supervised approach to extractive summarisation of scientific papers. In *Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017)*, pages 195–205.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186.
- Li Dong, Nan Yang, Wenhui Wang, Furu Wei, Xiaodong Liu, Yu Wang, Jianfeng Gao, Ming Zhou, and Hsiao-Wuen Hon. 2019. Unified language model pre-training for natural language understanding and generation. In *Advances in Neural Information Processing Systems*, pages 13042–13054.
- Yue Dong, Yikang Shen, Eric Crawford, Herke van Hoof, and Jackie Chi Kit Cheung. 2018. Bandit-sum: Extractive summarization as a contextual ban-

- dit. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3739–3748.
- Günes Erkan and Dragomir R Radev. 2004. Lexrank: Graph-based lexical centrality as salience in text summarization. *Journal of artificial intelligence research*, 22:457–479.
- Atefeh Farzindar and Guy Lapalme. 2004. [Legal text summarization by exploration of the thematic structure and argumentative roles](#). In *Text Summarization Branches Out*, pages 27–34, Barcelona, Spain. Association for Computational Linguistics.
- Sebastian Gehrmann, Yuntian Deng, and Alexander M Rush. 2018. Bottom-up abstractive summarization. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4098–4109.
- Matt Grenander, Yue Dong, Jackie Chi Kit Cheung, and Annie Louis. 2019. [Countering the effects of lead bias in news summarization via multi-stage training and auxiliary losses](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6019–6024, Hong Kong, China. Association for Computational Linguistics.
- Max Grusky, Mor Naaman, and Yoav Artzi. 2018. Newsroom: A dataset of 1.3 million summaries with diverse extractive strategies. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 708–719.
- Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend. In *Advances in neural information processing systems*, pages 1693–1701.
- Tsutomu Hirao, Yasuhisa Yoshida, Masaaki Nishino, Norihito Yasuda, and Masaaki Nagata. 2013. [Single-document summarization as a tree knapsack problem](#). In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1515–1520, Seattle, Washington, USA. Association for Computational Linguistics.
- Wan-Ting Hsu, Chieh-Kai Lin, Ming-Ying Lee, Kerui Min, Jing Tang, and Min Sun. 2018. [A unified model for extractive and abstractive summarization using inconsistency loss](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 132–141, Melbourne, Australia. Association for Computational Linguistics.
- Wojciech Kryściński, Bryan McCann, Caiming Xiong, and Richard Socher. 2019. Evaluating the factual consistency of abstractive text summarization. *arXiv preprint arXiv:1910.12840*.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*.
- Piji Li, Zihao Wang, Wai Lam, Zhaochun Ren, and Lidong Bing. 2017. Saliency estimation via variational auto-encoders for multi-document summarization. In *Thirty-First AAAI Conference on Artificial Intelligence*.
- Maria Liakata, Simon Dobnik, Shyamasree Saha, Colin Batchelor, and Dietrich Rebholz-Schuhmann. 2013. [A discourse-driven content model for summarising scientific articles evaluated in a complex question answering task](#). In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 747–757, Seattle, Washington, USA. Association for Computational Linguistics.
- Chin-Yew Lin. 2004. [ROUGE: A package for automatic evaluation of summaries](#). In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Chin-Yew Lin and Eduard Hovy. 1997. Identifying topics by position. In *Fifth Conference on Applied Natural Language Processing*, pages 283–290.
- Chin-Yew Lin and Eduard Hovy. 2002. From single to multi-document summarization. In *Proceedings of the 40th annual meeting of the association for computational linguistics*, pages 457–464.
- Yang Liu and Mirella Lapata. 2019. Text summarization with pretrained encoders. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3721–3731.
- Rada Mihalcea and Paul Tarau. 2004. Textrank: Bringing order into text. In *Proceedings of the 2004 conference on empirical methods in natural language processing*, pages 404–411.
- SPFGH Moen and Tapio Salakoski2 Sophia Ananiadou. 2013. Distributional semantics resources for biomedical text processing. *Proceedings of LBM*, pages 39–44.
- Edward Moroshko, Guy Feigenblat, Haggai Roitman, and David Konopnicki. 2019. An editorial network for enhanced document summarization. In *Proceedings of the 2nd Workshop on New Frontiers in Summarization*, pages 57–63.
- Ramesh Nallapati, Feifei Zhai, and Bowen Zhou. 2017. Summarunner: a recurrent neural network based sequence model for extractive summarization of documents. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, pages 3075–3081.

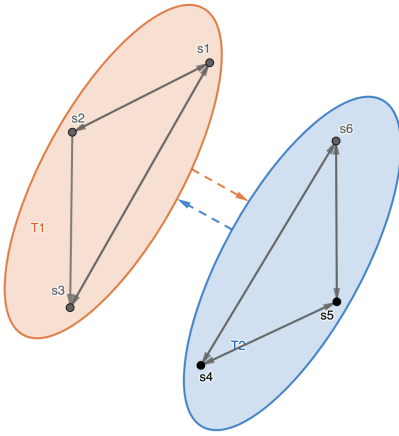
- Ramesh Nallapati, Bowen Zhou, Cicero dos Santos, Caglar Gulcehre, and Bing Xiang. 2016. Abstractive text summarization using sequence-to-sequence rnns and beyond. In *Proceedings of The 20th SIGNLL Conference on Computational Natural Language Learning*, pages 280–290.
- Shashi Narayan, Shay B Cohen, and Mirella Lapata. 2018a. Don’t give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1797–1807.
- Shashi Narayan, Shay B Cohen, and Mirella Lapata. 2018b. Ranking sentences for extractive summarization with reinforcement learning. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1747–1759.
- Ani Nenkova, Kathleen McKeown, et al. 2011. Automatic summarization. *Foundations and Trends® in Information Retrieval*, 5(2–3):103–233.
- Ani Nenkova and Lucy Vanderwende. 2005. The impact of frequency on summarization. *Microsoft Research, Redmond, Washington, Tech. Rep. MSR-TR-2005*, 101.
- Daraksha Parveen, Hans-Martin Ramsel, and Michael Strube. 2015. Topical coherence for graph-based extractive summarization. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1949–1954.
- Vahed Qazvinian and Dragomir R. Radev. 2008. [Scientific paper summarization using citation summary networks](#). *CoRR*, abs/0807.1560.
- Dragomir R. Radev, Hongyan Jing, and Malgorzata Budzikowska. 2000. [Centroid-based summarization of multiple documents: sentence extraction, utility-based evaluation, and user studies](#). In *NAACL-ANLP 2000 Workshop: Automatic Summarization*.
- Nils Reimers and Iryna Gurevych. 2019. [Sentence-BERT: Sentence embeddings using Siamese BERT-networks](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.
- Evan Sandhaus. 2008. The new york times annotated corpus. *Linguistic Data Consortium, Philadelphia*, 6(12):e26752.
- Abigail See, Peter J Liu, and Christopher D Manning. 2017. Get to the point: Summarization with pointer-generator networks. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1073–1083.
- Guokan Shang, Wensi Ding, Zekun Zhang, Antoine Tixier, Polykarpos Meladianos, Michalis Vazirgiannis, and Jean-Pierre Lorré. 2018. Unsupervised abstractive meeting summarization with multi-sentence compression and budgeted submodular maximization. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 664–674.
- Yuanlong Shao, Stephan Gouws, Denny Britz, Anna Goldie, Brian Strope, and Ray Kurzweil. 2017. Generating high-quality and informative conversation responses with sequence-to-sequence models. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2210–2219.
- Josef Steinberger and Karel Jezek. 2004. Using latent semantic analysis in text summarization and summary evaluation. *Proc. ISIM*, 4:93–100.
- Sandeep Subramanian, Raymond Li, Jonathan Pilaft, and Christopher Pal. 2019. On extractive and abstractive neural document summarization with transformer language models. *arXiv preprint arXiv:1909.03186*.
- Simone Teufel. 1997. Sentence extraction as a classification task. In *Intelligent Scalable Text Summarization*.
- Simone Teufel and Marc Moens. 2002. Summarizing scientific articles: experiments with relevance and rhetorical status. *Computational linguistics*, 28(4):409–445.
- Antoine Tixier, Polykarpos Meladianos, and Michalis Vazirgiannis. 2017. Combining graph degeneracy and submodularity for unsupervised extractive summarization. In *Proceedings of the workshop on new frontiers in summarization*, pages 48–58.
- Lucy Vanderwende, Hisami Suzuki, Chris Brockett, and Ani Nenkova. 2007. Beyond subbasic: Task-focused summarization with sentence simplification and lexical expansion. *Information Processing & Management*, 43(6):1606–1618.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.
- Xiaojun Wan. 2008. [An exploration of document impact on graph-based multi-document summarization](#). In *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing*, pages 755–762, Honolulu, Hawaii. Association for Computational Linguistics.
- Xiaojun Wan and Jianwu Yang. 2008. Multi-document summarization using cluster-based link analysis. In *Proceedings of the 31st annual international ACM SIGIR conference on Research and development in information retrieval*, pages 299–306.

- Yuxiang Wu and Baotian Hu. 2018. Learning to extract coherent summary via deep reinforcement learning. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- Wen Xiao and Giuseppe Carenini. 2019. Extractive summarization of long documents by combining global and local context. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3002–3012.
- Wenpeng Yin and Yulong Pei. 2015. Optimizing sentence modeling and selection for document summarization. In *Twenty-Fourth International Joint Conference on Artificial Intelligence*.
- Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter J Liu. 2019a. Pegasus: Pre-training with extracted gap-sentences for abstractive summarization. *arXiv preprint arXiv:1912.08777*.
- Xingxing Zhang, Furu Wei, and Ming Zhou. 2019b. Hibert: Document level pre-training of hierarchical bidirectional transformers for document summarization. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5059–5069.
- Hao Zheng and Mirella Lapata. 2019. Sentence centrality revisited for unsupervised summarization. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6236–6247.
- Ming Zhong, Pengfei Liu, Danqing Wang, Xipeng Qiu, and Xuan-Jing Huang. 2019. Searching for effective neural extractive summarization: What works and what’s next. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1049–1058.
- Qingyu Zhou, Nan Yang, Furu Wei, Shaohan Huang, Ming Zhou, and Tiejun Zhao. 2018. Neural document summarization by jointly learning to score and select sentences. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 654–663.

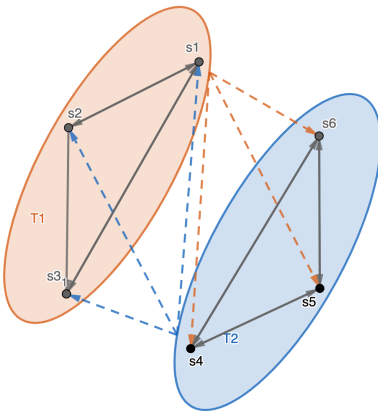
A Appendices



(a) Flat fully-connected graph



(b) Section-section hierarchical multiplication (hierarchy-multiply, ours)



(c) Section-sentence hierarchical addition (hierarchy-add, ours)

Figure 4: Comparison of the flat fully-connected graph used in Erkan and Radev (2004); Mihalcea and Tarau (2004); Zheng and Lapata (2019) to the hierarchical graph used in our models (b) and (c). Although the section-section multiplication reduces the edge computation proportionally to the number of sections, we found it oversimplifies the graph by assuming independence between sentences across different sections. Our final model loosens the assumption by including section-sentence connections as shown in sub-figure (c).

A.1 Different Hierarchical Structure

Besides our proposed hierarchical model (Figure 4 (c), hierarchy-add) in the paper, we also proposed and experimented with another novel hierarchical graph by introducing section-section connections (Figure 4 (b), hierarchy-multiply). In this hierarchical setting, we multiply a sentence’s sectional importance with its sentence importance (Eqn. (2)) to form the final centrality score:

$$c(v_j^I) = \mu_1 \cdot c_{\text{inter}}(v_j^I) \times c_{\text{intra}}(v_j^I). \quad (7)$$

Model	ROUGE-1	ROUGE-2	ROUGE-L
Various Hierarchical Centrality			
no-hierarchy	41.88	15.39	37.91
hierarchy-multiply (ours)	43.04	16.76	38.77
hierarchy-add (ours)	43.42	16.76	39.23

Table 8: Results on the PubMed validation set with different positional function or different hierarchical information.

Our empirical results indicate the hierarchy-multiply model always outperforms no-hierarchy models ((Figure 4 (a)) but under performs hierarchy-add. Nevertheless, Table 8 shows that adding any hierarchical structure results in performance improvement by wide margins when compared to the no-hierarchy model.

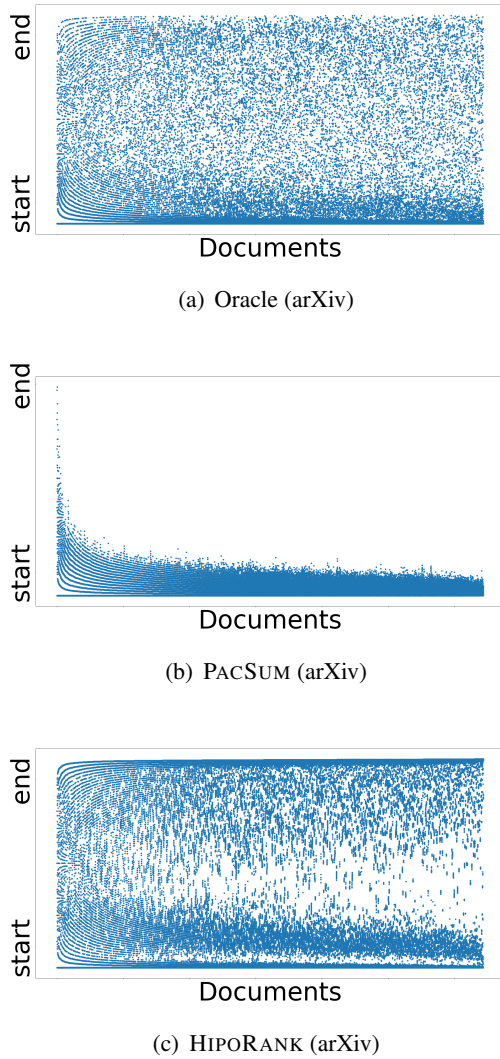


Figure 5: Sentence positions in source document for extractive summaries generated by different models on the arXiv validation set. Documents on the x-axis are ordered by increasing article length from shortest to longest.

A.2 Sentence Position Comparison

Figure 5 shows the sentence positions in source document for extractive summaries generated by different models on the arXiv validation set. We can again see that PACSUM is biased towards selecting sentences that appear at the beginning of a document while our method selects sentences in every section and near the article boundaries, similar to the oracle.