

Topic-controllable Abstractive Summarization

Anonymous ACL submission

Abstract

Existing approaches for topic-controllable summarization either incorporate topic embeddings or modify the attention mechanism. The incorporation of such approaches in a particular summarization model requires the adaptation of its codebase, a process that can be complex and time-consuming. Instead, we propose a model-agnostic topic-controllable summarization method employing a simple tagging-based formulation that can effortlessly work with any summarization model. In addition, we propose a new topic-oriented evaluation measure to quantitatively evaluate the generated summaries based on the topic affinity between the generated summary and the desired topic. Experimental results show that the proposed tagging-based formulation can achieve similar or even better performance compared to the embedding-based approach, while being at the same time significantly faster.

1 Introduction

The exponential rise in the volume of textual data available through various sources, ranging from social media to financial reports, makes it virtually impossible for humans to digest all the important information for their needs, without spending an enormous amount of effort. Automatic summarization methods can mitigate this problem, by shortening texts to a more concise form (Nallapati et al., 2016; Celikyilmaz et al., 2018; Liu and Lapata, 2020; Song et al., 2019).

Even though early methods had limited success on this task, mainly focusing on *extractive summarization* (Fang et al., 2017; Mao et al., 2019), the advent of deep learning led to much more powerful neural *abstractive summarization* (See et al., 2017; Song et al., 2019; Dong et al., 2019; Lewis et al., 2020; Zhang et al., 2020) methods. These methods go beyond extracting unaltered sentences from the input, allowing for generating the summary using

novel words and phrases that are not necessarily part of the input text.

Despite the success of deep learning models, there is often the need to go beyond delivering a generic summary of the document, and instead produce a summary that focuses on a specific topic that pertains to the user’s interests. For example, a newswire article may discuss two topics, such as sports and politics, yet the user may be interested only in the sports aspect. Existing topic-controllable summarization models address this need either by incorporating topic embeddings into the model’s architecture (Krishna and Srinivasan, 2018) or by modifying the attention mechanism (Bahrainian et al., 2021). However, they are restricted to very specific neural architectures and it is not straightforward to use them with any summarization model.

At the same time, there is no clear way to evaluate such approaches, since there is no evaluation measure designed specifically for topic-controllable summarization. Indeed, existing methods just use the typical ROUGE score (Lin, 2004) for measuring the summarization accuracy and then employ user studies to qualitatively evaluate whether the topic of the generated summaries indeed matches the users’ needs (Krishna and Srinivasan, 2018; Bahrainian et al., 2021).

Based on the aforementioned observations, we propose a model-agnostic topic-controllable summarization method that can be effortlessly combined with any neural architecture. Given a topic labeled collection, the proposed method works by first extracting keywords that are semantically related to the topic the user requested and employing special tokens to tag them before feeding the document to the summarization model. Experimental results show that this can be an effective and efficient way to influence summarization models towards the users’ needs.

Furthermore, we propose a topic-aware evalu-

082 ation measure for quantitatively evaluating topic-
083 controllable summarization methods in an objec-
084 tive way without involving expensive and time-
085 consuming user studies. In particular, we propose
086 calculating prototype term weighting representa-
087 tions, namely tf-idf, of different topics, and then
088 calculating the cosine similarity between the gener-
089 ated summaries and the prototype topic vectors.

090 The contributions of this paper can be summa-
091 rized as follows:

- 092 • We propose a simple, yet effective and effi-
093 cient model-agnostic way to perform topic-
094 controllable summarization.
- 095 • We adapt an existing topic-controllable
096 method to work with Transformer-based archi-
097 tectures, scaling up from existing RNN-
098 based formulations, establishing a strong, yet
099 computationally demanding baseline for topic-
100 oriented summarization.
- 101 • We propose a topic-oriented measure to quan-
102 titatively evaluate the generated summaries
103 without the need for resorting to human stud-
104 ies.
- 105 • We provide an extensive empirical evaluation
106 as well as a zero-shot experimental evalua-
107 tion, demonstrating both the generality of the
108 proposed method, as well as its effectiveness.

109 The rest of the paper is organized as follows.
110 In Section 2 we review the existing topic-oriented
111 summarization related literature. In Section 3 we
112 introduce the proposed methods while in Section 4
113 we provide the experimental results. Finally, con-
114 clusions are drawn and interesting feature research
115 directions are discussed in Section 5.

116 2 Topic-oriented Summarization

117 Methods for topic-oriented summarization belong
118 to two broader categories: a) methods that em-
119 ploy topical information to enhance the quality of
120 the generated summaries and b) topic-controllable
121 methods that use topical information to control the
122 output of the generated summaries.

123 2.1 Improving summarization using topical 124 information

125 The integration of topic modeling into summariza-
126 tion models has been initially used in the literature
127 to improve the quality of existing state-of-the-art

128 models (Ailem et al., 2019; Wang et al., 2020; Liu
129 and Yang, 2021). Statistical topic models such
130 as Latent Dirichlet Allocation (LDA) (Blei et al.,
131 2003) or Poisson Factor Analysis (PFA) (Zhou
132 et al., 2012) are used to supply summarization mod-
133 els with global topic semantics, allowing the gener-
134 ation of more coherent and consistent summaries.

Ailem et al. (2019) use LDA to influence the
135 model to generate summaries based on both the
136 input text and the underlying document topics and
137 as a result to improve the quality of the gener-
138 ated summary. To achieve this, the decoder of
139 a pointer generator network is enhanced with the
140 information of the latent topics that are derived
141 from an LDA model. Thus, the integration of
142 topic modeling can capture hidden semantic struc-
143 tures based on word co-occurrences, allowing the
144 model to generate better summaries conditioned
145 on a more global context. Similar methods have
146 been applied by Wang et al. (2020) using PFA with
147 a plug-and-play architecture that can be adapted
148 to any Transformer-based model. This architec-
149 ture consists of 3 independent modules: Semantic-
150 informed attention (SIA), Topic Embedding with
151 Masked Attention (TEMA), and Document-related
152 modulation (DRM). SIA is embedded as an addi-
153 tional head into the multi-head attention mecha-
154 nism. This added head is extracted from a fixed
155 semantic-similarity attention matrix for each topic.
156 TEMA uses topic embeddings as an additional de-
157 coder input based on the top- n topics from the input
158 document. Since a topic can be represented as a
159 distribution over all the tokens from the vocabulary,
160 topic embeddings can be derived from a mixture of
161 all the corresponding token embeddings. Finally,
162 DRM is used to modulate a hidden layer for each
163 decoder adding a topic feature bias vector. 164

Liu and Yang (2021) propose to enhance summa-
165 rization models using an Extreme Multi-Label Text
166 Classification (XMTC) model to improve the con-
167 sistency between the underlying topics of the input
168 document and the summary, leading to summaries
169 of higher quality. 170

Even though Wang et al. (2020) refers to the
171 potential of controlling the output conditioned on
172 a specific topic using GPT-2 (Radford et al., 2019)
173 with TEMA, all the aforementioned approaches
174 are focused on improving the accuracy of existing
175 summarization models. 176

2.2 Topic-controllable summarization methods

Some steps towards controlling the output of a summarization model conditioned on a thematic category have been made by Krishna and Srinivasan (2018) proposing a controllable summarization setting that builds upon the pointer generator network (See et al., 2017). The topical information is integrated into the model as a topic vector, which is then concatenated with each of the word embeddings of the input text. Each topic vector is computed as a Bag of Words (BoW) representation that is derived from Vox Dataset (Vox Media, 2017), a news dataset that contains articles from 185 different news topics.

Krishna and Srinivasan (2018) created a topic-oriented training dataset that builds upon CNN/DailyMail as follows. First, the dot-product between the BoW representation of the summary and all the BoW topic representations is computed. The topic with the highest similarity is assigned to the corresponding article while articles with more than one dominant topic are discarded. All the topic-assigned articles are used to compile a temporary intermediate dataset. To create the topic-oriented dataset, two articles a_1 and a_2 with different topics, are randomly selected from the intermediate dataset. A new article a' is created by sequentially selecting sentences from both articles. The new article a' is assigned with the summary from one out of two selected articles and the same process is repeated to create a new article a'' which is now assigned with the remaining summary. Then, the initially selected articles a_1 and a_2 are discarded from the intermediate dataset. This process is continued until there are no articles in the intermediate dataset or all the remaining articles belong to the same topic. Finally, the new topic-oriented dataset consists of super-articles that discuss two distinct topics but are assigned each time with one of the corresponding summaries so the model learns to distinguish the most important sentences for the corresponding topic during training.

Recently, Bahrainian et al. (2021) propose to incorporate the topical information from each document to modify the attention mechanism of the pointer generator network (See et al., 2017). The modification of the attention mechanism is introduced as topical attention generated by an LDA model. More specifically, each word is represented

as a topic vector that is derived from LDA and then is combined with the original attention weights of the model to compute the final attention weights. It is important to note that even though the model is trained with the topical attention mechanism during training, no topical information is used during inference. Thus, the aforementioned method allows for controlling the topic of the generated summary only from the perspective of the restriction of unwanted topics during training, contrary to the proposed method, which allows for guiding the generation towards a topic, during inference.

3 Contributions

In this section, we present the main contributions of this paper. More specifically, we introduce two different topic-controllable methods to guide the summary generation towards a specific topic: a) tagging-based formulation and b) embedding-based formulation. We also present the proposed topic-oriented similarity measure which is used for evaluating the topic affinity between the desired topics and the generated summaries.

3.1 Tagging-based formulation

The proposed tagging-based method employs a trivial, yet effective mechanism to shift the summary generation towards the desired topic, assuming the existence of a set of representative terms for each thematic category. More specifically, after lemmatization, the most representative words for the desired topic are tagged with special tag tokens before feeding to the summarization model. As demonstrated in Section 4, this can be an effective way to intuitively guide the model towards the tagging words during both training and inference.

To apply this mechanism, a topic-oriented training set is required. However, this is not a straightforward process due to the lack of appropriate topic-oriented summarization datasets. Indeed, there are no existing datasets for summarization that contain multiple summaries for each input document, according to the different topical aspects of the text (Krishna and Srinivasan, 2018). Thus, we adopt the same approach with (Krishna and Srinivasan, 2018) to create a topic-oriented dataset that builds upon the CNN/DailyMail (Hermann et al., 2015). We apply the tagging mechanism to each document of the topic-oriented dataset according to the assigned topic of the corresponding summary. More specifically, for each document of the com-

277 piled dataset, we tag the terms that belong to the
 278 intersection of words between the lemmatized doc-
 279 ument and the top- N most representative terms for
 280 the corresponding topic.

281 The most representative words can be extracted
 282 either by simple prototype term weighting represen-
 283 tations such as BoW or tf-idf, statistical topic mod-
 284 eling algorithms such as Labeled LDA (Ramage
 285 et al., 2009) or even more sophisticated keyword ex-
 286 traction models (Ding and Luo, 2021; Liang et al.,
 287 2021). In this work, we use tf-idf to demonstrate
 288 the efficacy of our method, even when a simple
 289 mechanism is employed.

290 More specifically, we use tf-idf to extract docu-
 291 ment representations and then calculate the topical
 292 vectors. Given a corpus \mathcal{D} , we can represent a docu-
 293 ment d as a vector \mathbf{x}_d which contains the tf-idf
 294 scores for each term of the document. The tf-idf
 295 score for each term t of a document d , belonging
 296 to a corpus \mathcal{D} , is computed as:

$$297 \quad x_{dt} = tf(t, d, \mathcal{D}) \times idf(t, \mathcal{D}), \quad (1)$$

298 where $tf(t, d, \mathcal{D})$ indicates the number of times
 299 that term t appears in document d , while $idf(t, \mathcal{D})$
 300 indicates the inverse document frequency of term t
 301 in corpus \mathcal{D} which is computed as follows:

$$302 \quad idf(t, \mathcal{D}) = \log \frac{|\mathcal{D}| + 1}{df(t, \mathcal{D}) + 1} + 1, \quad (2)$$

303 where $df(t, \mathcal{D})$ is the frequency of term t in \mathcal{D} .
 304 Note that the length of each tf-idf vector is equal to
 305 the size of the vocabulary V of the corpus \mathcal{D} , i.e.,
 306 $\mathbf{x}_d \in \mathbb{R}^{|\mathcal{V}|}$, where $|\mathcal{V}|$ denote the cardinality of the
 307 vocabulary \mathcal{V} . Finally, we normalized the extracted
 308 vectors to have unit length as:

$$309 \quad \mathbf{x}_d^{(n)} = \frac{\mathbf{x}_d}{\|\mathbf{x}_d\|_2}, \quad (3)$$

310 where $\|\mathbf{x}\|_2$ is the l_2 norm of the vector \mathbf{x}_d .

311 Then, given a topic-assigned collection of docu-
 312 ments \mathcal{X} , we can follow the aforementioned pro-
 313 cedure to extract a topical vector representation \mathbf{y}_c
 314 for each topic c , by grouping together documents
 315 of the same topic and averaging their tf-idf repre-
 316 sentation as follows:

$$317 \quad \mathbf{y}_c = \frac{1}{|\mathcal{X}_c|} \sum_{\mathbf{x} \in \mathcal{X}_c} \mathbf{x} \quad (4)$$

318 The topical vector extraction is summarized in
 319 Figure 1.

Table 1: Representative terms for topics from 2017
 KDD Data Science+Journalism Workshop (Vox Media,
 2017)

Topic	Terms
Politics	policy, president, state, political, vote, law, country, election
Sports	game, sport, team, football, fifa, nfl, player, play, soccer, league
Health Care	patient, uninsured, insurer, plan, coverage, care, insurance, health
Education	student, college, school, educa- tion, test, score, loan, teacher
Movies	film, season, episode, show, movie, character, series, story
Space	earth, asteroid, mars, comet, nasa, space, mission, planet, astronaut

320 Finally, we extract the top- N most important
 321 terms for each topic according to the top tf-idf
 322 scores of each topical vector. An example of some
 323 indicative representative words for a number of
 324 topics in a topical corpus is shown in Table 1.

325 Given the set of representative words for each
 326 topic, a document, and the desired topic, the tag-
 327 ging mechanism works as follows:

- 328 1. All the words of the input document are lem-
 329 matized to their roots.
- 330 2. We identify the common words between the
 331 existing lemmatized tokens and the represen-
 332 tative words for the desired topic.
- 333 3. Finally, we tag each token of the input docu-
 334 ment with a special token, i.e., [TAG], only if
 335 the lemmatized form of this token is contained
 336 in the set of the most representative words for
 337 the corresponding topic.

338 For example, suppose that we pre-process the
 339 sentence below, as a part of an input document,
 340 from which we aim to guide the generation towards
 341 the topic “*Business & Finance*”.

342 “By one estimate, American individuals
 343 and **businesses** together spend 6.1 **bil-**
 344 **lion** hours complying with the **tax** code
 345 every year.”

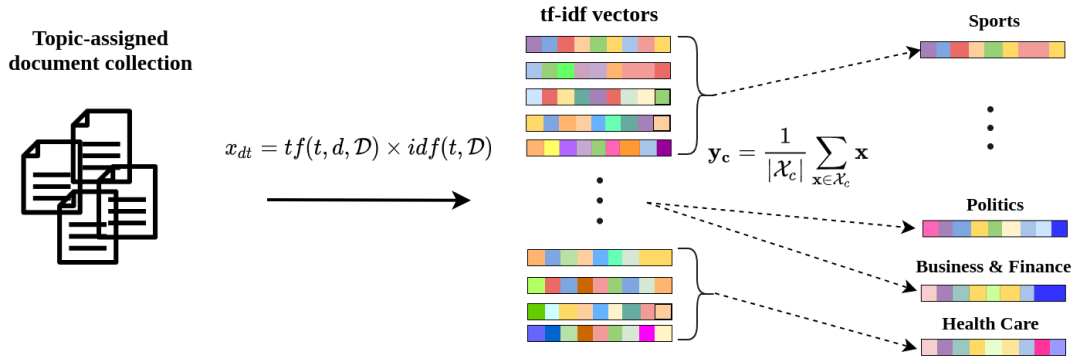


Figure 1: Topical vector extraction using tf-idf scores, given a topic-assigned document collection. First, we calculate tf-idf scores for each document. Then, documents of the same topic are grouped and their tf-idf representation is averaged.

Following the aforementioned procedure, we will enclose with special tokens, the words “businesses”, “billion” and “tax” since they belong to the set of the most representative words for the desired topic.

During training, the model learns to intuitively give more “attention” to the tagged words and as a result shift the generation towards the desired topic. The tagging mechanism can be used during inference to guide the summary generation towards the user-requested topic provided by any set of representative terms. Also, since this method does not affect the architecture of the summarization model, it can easily be applied to any model’s architecture.

3.2 Embedding-based formulation

To establish a strong baseline for comparing the tagging-based method with existing methods in the literature, we adapted the method proposed in Krishna and Srinivasan (2018) to work with Transformer-based architectures. As described in Section 2, Krishna and Srinivasan (2018) use a pointer generator network (See et al., 2017) to concatenate topic embeddings with token embeddings allowing for generating topic-oriented summaries. The topic embeddings are represented as one-hot encoding vectors with a size equal to the number of the total topics. During training, the model takes as inputs the corresponding topic embedding along with the input document.

However, this method cannot be directly applied to pre-trained Transformer-based models due to the different shapes of initialized weights of the word and position embeddings. Unlike RNNs, Transformer-based models are typically trained for general tasks and then fine-tuned with less data for

more specific tasks like Summarization. Thus, the architecture of a pre-trained model is already defined and cannot be altered easily to initialize the pre-trained model’s weights with the exact same shape of the concatenated word and topic embeddings. Another option would be to initialize the model from scratch with random weights with the appropriate shape of the concatenated word and topic embeddings but this would be very computationally demanding as it would require a large amount of data and time for training.

To this end, instead of concatenation, we propose to sum the topic embeddings following the same concept with positional encoding where token embeddings are summed with positional encoding representations to create an input representation that contains the position information. Instead of one-hot encoding embeddings, we use trainable embeddings allowing the model for optimizing them accordingly during training. The topic embeddings have the same dimensionality as the token embeddings.

To sum the trainable topic embeddings with token and positional embeddings, we modify the input representation as follows:

$$z_i = WE(x_i) + PE(i) + TE(i), \quad (5)$$

where WE, PE and TE are the word embeddings, positional encoding and topic embeddings respectively, for token x_i in position i .

Then, we use the same created topic-oriented dataset from Krishna and Srinivasan (2018) to fine-tune the summarization model for topic-oriented summarization allowing for establishing a strong comparison between the proposed tagging-based method and the more powerful embedding-based

one.

3.3 Topic-focused evaluation measure

As explained in Section 1, there is currently no structured way to evaluate the performance of topic-oriented summarization methodologies. To this end, we propose a new topic-oriented measure, Summarization Topic Affinity Score (STAS), to evaluate the generated summaries according to the semantic similarity between the vector representation of the desired topic and the generated summary. More specifically, we compute the similarity between the vector representation of the summary and the vector representation of the desired topic, divided by the maximum value of all the similarities between the vector representation of the summary and all the topic vector representations. Given the vector of the target topic \mathbf{x}_t and the vector representation of the predicted summary \mathbf{x}_s , STAS is computed as follows:

$$STAS(\mathbf{x}_s, \mathbf{x}_t) = \frac{s(\mathbf{x}_s, \mathbf{x}_t)}{\max\{s(\mathbf{x}_s, \mathbf{x}_{ti}) : i = 1 \dots N_t\}}, \quad (6)$$

where N_t is the number of topic and $s(\mathbf{x}_s, \mathbf{x}_t)$ indicates the cosine similarity between the two vectors \mathbf{x}_s and \mathbf{x}_t which is computed as follows:

$$s(\mathbf{x}_t, \mathbf{x}_s) = \frac{\mathbf{x}_t \mathbf{x}_s}{\|\mathbf{x}_t\| \|\mathbf{x}_s\|}. \quad (7)$$

Thus, summaries that are similar to the requested topic are rewarded while summaries that are dissimilar are penalized.

4 Experimental Evaluation

In this section, we present the experimental results of the proposed method. First, we introduce the experimental setup used for the evaluation, including the dataset generation procedure, the evaluation metrics, and employed deep learning architectures. Then, we proceed by presenting and discussing the experimental evaluation using both the proposed tagging-based method, as well as the embedding-based method, appropriately adapted to work on Transformers.

4.1 Experimental setup

Datasets and Evaluation Metrics In order to create the topic-oriented dataset as described in Section 3, we use the Vox Dataset (Vox Media, 2017), which consists of 23,024 news articles of 185 different topical categories. We discarded topics with

relatively low frequency, i.e. lower than 20 articles, as well as articles assigned to general categories that do not discuss explicitly a topic, i.e. “*The Latest*”, “*Vox Articles*”, “*On Instagram*” and “*On Snapchat*”.

In the experiments, we investigate two different setups: a) fine-tuning without pre-processing the Vox dataset, keeping also noisy categories that do not discuss a particular topic, and b) fine-tuning after pre-processing the Vox dataset as described. All summaries of the created dataset are assigned with a topic according to the similarity between the derived topical vector representations and the vectorized summary. Thus, keeping noisy topics might lead to false topic assignments to the training summaries.

After pre-processing, we end up with 14,312 articles from 70 categories out of the 185 initial topical categories. Then, following the same procedure as Krishna and Srinivasan (2018), we create the topic-oriented dataset combining sentences from article-pairs from the CNN/DailyMail (Hermann et al., 2015). We use the anonymized version of CNN/Dailymail similar to See et al. (2017). The final topic-oriented dataset consists of 132,766, 5,248, and 6,242 articles for training, validation, and test, respectively. The average document and summary length of the created dataset is 1,544 and 56 tokens, respectively.

All the tags for the tagging-based method were applied to the dataset after lemmatization using NLTK (Bird, 2006) based on the top- $N=100$ most representative terms for each topic. We also use the Vox Dataset (Vox Media, 2017) to extract the tf-idf vector representations for each document in the corpus. To this end, we employed the tf-idf vectorizer provided by the Scikit-learn library (Pedregosa et al., 2011).

All methods were evaluated using both the well-known ROUGE (Lin, 2004) score, to measure the quality of the generated summary, as well as the proposed STAS measure.

Models and Training For all the conducted experiments we have employed a BART-large (Lewis et al., 2020) architecture, which is a transformer-based model with a bidirectional encoder and an auto-regressive decoder. BART-large consists of 12 layers for both encoder and decoder and 406M parameters. We used the implementation provided by Hugging Face for the BART-large architecture (Wolf et al., 2020).

We fine-tune all the models for 100,000 steps with a learning rate of 0.00003 and batch size 4 with early stopping on the validation set. We use the established parameters for BART-large architecture using label smoothed cross-entropy loss (Pereyra et al., 2017) with the label smoothing factor set to 0.1.

For all the experiments, we use PyTorch version 1.10 and Hugging Face version 4.11.0. All the models were trained using available GPUs in Google Colab¹, with approximate average training runtime 9.5 and 18 hours for the tagging-based and embedding-based method, respectively. Both data and code will be publicly available.

4.2 Results

The evaluation results on the generated dataset are shown in Table 2. We report results using five different methods. First, we employ both the generic Pointer Generation method (“PG”) (See et al., 2017), as well as the topic-oriented PG (“Topic-Oriented PG”) (Krishna and Srinivasan, 2018). We also use the generic BART (Lewis et al., 2020) model (“BART”) fine-tuned on the regular CNN/DailyMail dataset for summarization, as well as both the adapted embedding-based formulation (“BART_{emb}”) and the tagging-based formulation (“BART_{tag}”).

The experimental results reported in Table 2 for the two different pre-processing setups indicate that topic-oriented methods indeed perform significantly better compared to the baseline methods that do not take into account the topic requested by the user. Furthermore, the proposed BART-based formulation significantly outperforms the generic PG approach, regardless of the applied topic mechanism (BART_{emb} or BART_{tag}). Also, the proposed tag-based mechanism seems to be more robust to noise, leading to slightly better results when no pre-processing is applied. On the other hand, when the data are pre-processed, both the embedding and the topic tagging approach lead to quite similar results. However, as we further demonstrate later, the proposed tagging method is significantly faster than embedding-based approaches, leading to the overall best trade-off between accuracy and speed.

The results of the inference time for both methods are shown in Table 3. The inference time of the proposed method is significantly smaller, improving the performance of the model by almost one

Table 2: Experimental results on the created topic-oriented dataset based on CNN/DailyMail dataset. We report f-1 scores for ROUGE-1 (R-1), ROUGE-2 (R-2) and ROUGE-L (R-L).

	R-1	R-2	R-L
PG (See et al., 2017)	26.8	9.2	24.5
BART (Lewis et al., 2020)	30.46	11.92	20.57
Topic-Oriented PG (Krishna and Srinivasan, 2018)	34.1	13.6	31.2
Proposed BART _{emb} (all topics)	37.64	16.94	26.20
Proposed BART _{tag} (all topics)	37.94	17.21	26.49
Proposed BART _{tag} (pre-processed topics)	39.30	18.06	27.49
Proposed BART _{emb} (pre-processed topics)	40.15	18.53	28.06

order of magnitude. Indeed, the proposed method can perform inference on 100 articles in less than 40 seconds, while the embedding-based formulation requires more than 300 seconds for the same task.

Table 3: Inference time for 100 articles. All numbers are reported in seconds.

	Tagging	Inference	Total time
BART _{emb}	-	303.0	303.0
BART _{tag}	7.1	32.0	39.1

In Table 4, we also provide an experimental evaluation using the proposed Summarization Topic Affinity Score (STAS) measure. The effectiveness of using topic-oriented approaches is further highlighted using the proposed method since the improvements acquired when applying the proposed method are much higher compared to the ROUGE score. Also, both the embedding and the tagging method lead to similar results (~68.5%) using STAS measure, even though the tagging approach is significantly faster and easier to apply. Note that when no pre-processing is used, the tagging-based approach is more robust to noise, leading to a better STAS score (49.65%) compared to the embedding-based approach (46.70%).

4.3 Zero-shot experimental evaluation

The tagging mechanism allows the model to intuitively guide the summary generation according to the tagged words of the desired topic which can also be an effective way to generalize to unseen topics. To demonstrate the efficacy of the tagging-based model on unseen topics, we fine-tune the BART model on the same training set of the created topic-oriented dataset but removing 5% of the

¹<https://research.google.com/colaboratory/>

Table 4: Evaluation based on the proposed Summarization Topic Affinity Score (STAS).

	STAS (%)
BART (Lewis et al., 2020) (all topics)	33.99
Proposed BART _{emb} (all topics)	46.70
Proposed BART _{tags} (all topics)	49.65
BART (Lewis et al., 2020) (pre-processed topics)	51.86
Proposed BART _{tags} (pre-processed topics)	68.42
Proposed BART _{emb} (pre-processed topics)	68.50

589 topics. More specifically, we randomly remove 3
590 topics out of the 70 topics (i.e., “Movies”, “Trans-
591 portation” and “Podcasts”) of the training set and
592 evaluate the model both on the test set of seen top-
593 ics and on the zero-shot test, which consists of 264
594 articles of unseen topics, as shown in Table 5.

Table 5: Experimental results on both test set with seen topics and zero-shot test set with unseen topics. We report STAS measure scores and f-1 scores for R-1, R-2 and R-L.

	R-1	R-2	R-L	STAS (%)
BART _{tag} (seen topics)	38.31	17.27	26.48	68.21
BART _{tag} (unseen topics)	37.52	16.99	26.71	74.80

595 Even though the model has not seen the zero-
596 shot topics during training, it can successfully gener-
597 ate topic-oriented summaries for these topics
598 achieving similar results in terms of ROUGE-1
599 score (~38% for both test sets) and even better
600 results in terms of STAS measure on the zero-shot
601 test (~68%) compared to the test set with the seen
602 topics (~74%). This finding confirms the capa-
603 bility of the tagging-based method to generalize
604 successfully to unseen topics, provided that a set
605 of representative terms is given.

606 4.4 Examples of generated summaries

607 We present some examples generated by the
608 tagging-based model on the created dataset for dif-
609 ferent topics as shown in Table 6. Indeed, the pro-
610 posed model can shift the generation towards the
611 desired topic of the super-article which contains
612 different topics. Furthermore, the generation of
613 the summary according to the corresponding topic
614 is not affected by the presence of the other topic
615 which is also discussed in the input article.

Table 6: Generated summaries of our proposed tagging-based model according to the two different topics of the super-article containing articles of these topics. Part of summaries is truncated due to size limitations.

Sports: Jenson Button and Fernando Alonso failed to finish the Malaysian Grand Prix ... Button lasted double the amount of time as his teammate.

Gun Violence: Adam Lanza killed his mother, Nancy, inside the home before killing 20 first-graders and six members of staff at Sandy Hook Elementary School in 2012. ...

Transportation: Ford unveiled two prototype electric bikes at Mobile World Congress in Barcelona. ... The bikes are part of an experiment by Ford called Angle on Mobility.

Neuroscience: Researchers from Bristol University measured biosonar bat calls to calculate what members of group perceived as they foraged for food ...

616 5 Conclusions and Future Work

617 We proposed a model-agnostic topic-controllable
618 method that can work with any summarization
619 model to influence the summary generation towards
620 the desired topic. The proposed method works
621 by employing special tokens to tag semantically-
622 related words for each topic and then guide the
623 generation towards this topic. To establish a
624 strong baseline, we also adapt an existing topic-
625 controllable embedding-based method to a more
626 powerful Transformer-based model, scaling up
627 from traditional RNNs. We also proposed STAS, a
628 structured way to evaluate the generated summaries
629 according to the affinity of the requested topic with
630 the topic of the generated summary. Experimental
631 results under two different pre-processing setups
632 demonstrate that the proposed method can achieve
633 similar or even better performance than the adapted
634 embedding-based mechanism, while being signifi-
635 cantly faster and easier to apply.

636 Future research could examine other controllable
637 aspects, such as style (Fan et al., 2018) or enti-
638 ties (He et al., 2020). In addition, it would be very
639 interesting to extend the proposed method towards
640 working with any arbitrary topic, bypassing the re-
641 quirement of having a labeled document collection
642 of a topic to be able to guide the summary towards
643 this topic.

References

- 645 Melissa Ailem, Bowen Zhang, and Fei Sha. 2019. Topic
646 augmented generator for abstractive summarization.
647 *arXiv preprint arXiv:1908.07026*.
- 648 Seyed Ali Bahrainian, George Zerveas, Fabio Crestani,
649 and Carsten Eickhoff. 2021. **Cats: Customizable
650 abstractive topic-based summarization**. *ACM Trans.
651 Inf. Syst.*, 40(1).
- 652 Steven Bird. 2006. **NLTK: The Natural Language
653 Toolkit**. In *Proceedings of the COLING/ACL 2006
654 Interactive Presentation Sessions*, pages 69–72, Syd-
655 ney, Australia. Association for Computational Lin-
656 guistics.
- 657 David M Blei, Andrew Y Ng, and Michael I Jordan.
658 2003. Latent dirichlet allocation. *the Journal of
659 machine Learning research*, 3:993–1022.
- 660 Asli Celikyilmaz, Antoine Bosselut, Xiaodong He, and
661 Yejin Choi. 2018. **Deep communicating agents for
662 abstractive summarization**. In *Proceedings of the
663 2018 Conference of the North American Chapter of
664 the Association for Computational Linguistics: Hu-
665 man Language Technologies, Volume 1 (Long Pa-
666 pers)*, pages 1662–1675, New Orleans, Louisiana.
667 Association for Computational Linguistics.
- 668 Haoran Ding and Xiao Luo. 2021. **AttentionRank: Un-
669 supervised keyphrase extraction using self and cross
670 attentions**. In *Proceedings of the 2021 Conference
671 on Empirical Methods in Natural Language Process-
672 ing*, pages 1919–1928, Online and Punta Cana, Do-
673 minican Republic. Association for Computational
674 Linguistics.
- 675 Li Dong, Nan Yang, Wenhui Wang, Furu Wei, Xi-
676 aodong Liu, Yu Wang, Jianfeng Gao, Ming Zhou,
677 and Hsiao-Wuen Hon. 2019. Unified language model
678 pre-training for natural language understanding and
679 generation. In *Advances in Neural Information Pro-
680 cessing Systems*, pages 13042–13054.
- 681 Angela Fan, David Grangier, and Michael Auli. 2018.
682 **Controllable abstractive summarization**. In *Proce-
683 edings of the 2nd Workshop on Neural Machine Transla-
684 tion and Generation*, pages 45–54, Melbourne, Aus-
685 tralia. Association for Computational Linguistics.
- 686 Changjian Fang, Dejun Mu, Zhenghong Deng, and Zhi-
687 ang Wu. 2017. Word-sentence co-ranking for auto-
688 matic extractive text summarization. *Expert Systems
689 with Applications*, 72:189–195.
- 690 Junxian He, Wojciech Kryscinski, Bryan McCann,
691 Nazneen Rajani, and Caiming Xiong. 2020. Ctrl-
692 sum: Towards generic controllable text summariza-
693 tion. *ArXiv*, abs/2012.04281.
- 694 Karl Moritz Hermann, Tomas Kocisky, Edward Grefen-
695 stette, Lasse Espeholt, Will Kay, Mustafa Suleyman,
696 and Phil Blunsom. 2015. Teaching machines to read
697 and comprehend. In *Advances in neural information
698 processing systems*, pages 1693–1701.
- Kundan Krishna and Balaji Vasani Srinivasan. 2018. 699
Generating topic-oriented summaries using neural 700
attention. In *Proceedings of the 2018 Conference of 701
the North American Chapter of the Association for 702
Computational Linguistics: Human Language Tech- 703
nologies, Volume 1 (Long Papers)*, pages 1697–1705. 704
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan 705
Ghazvininejad, Abdelrahman Mohamed, Omer Levy, 706
Veselin Stoyanov, and Luke Zettlemoyer. 2020. 707
**BART: Denoising sequence-to-sequence pre-training 708
for natural language generation, translation, and com- 709
prehension**. In *Proceedings of the 58th Annual Meet- 710
ing of the Association for Computational Linguistics*, 711
pages 7871–7880, Online. Association for Computa- 712
tional Linguistics. 713
- Xinnian Liang, Shuangzhi Wu, Mu Li, and Zhoujun Li. 714
2021. **Unsupervised keyphrase extraction by jointly 715
modeling local and global context**. In *Proceedings of 716
the 2021 Conference on Empirical Methods in Nat- 717
ural Language Processing*, pages 155–164, Online 718
and Punta Cana, Dominican Republic. Association 719
for Computational Linguistics. 720
- Chin Yew Lin. 2004. **Rouge: A package for auto- 721
matic evaluation of summaries**. In *Proceedings of the 722
workshop on text summarization branches out (WAS 723
2004)*. 724
- Jingzhou Liu and Yiming Yang. 2021. Enhancing sum- 725
marization with text classification via topic consis- 726
tency. In *Joint European Conference on Machine 727
Learning and Knowledge Discovery in Databases*, 728
pages 661–676. Springer. 729
- Yang Liu and Mirella Lapata. 2020. **Text summariza- 730
tion with pretrained encoders**. In *EMNLP-IJCNLP 731
2019 - 2019 Conference on Empirical Methods in 732
Natural Language Processing and 9th International 733
Joint Conference on Natural Language Processing*, 734
Proceedings of the Conference. 735
- Xiangke Mao, Hui Yang, Shaobin Huang, Ye Liu, and 736
Rongsheng Li. 2019. Extractive summarization us- 737
ing supervised and unsupervised learning. *Expert 738
systems with applications*, 133:173–181. 739
- Ramesh Nallapati, Bowen Zhou, Cicero dos Santos, 740
Caglar Gulcehre, and Bing Xiang. 2016. **Abstrac- 741
tive Text Summarization using Sequence-to-sequence 742
RNNs and Beyond**. In *Proceedings of the 2016 743
SIGNLL Conference on Computational Natural Lan- 744
guage Learning*, pages 280–290, Stroudsburg, PA, 745
USA. Association for Computational Linguistics. 746
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, 747
B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, 748
R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, 749
D. Cournapeau, M. Brucher, M. Perrot, and E. Duch- 750
esnay. 2011. Scikit-learn: Machine learning in 751
Python. *Journal of Machine Learning Research*, 752
12:2825–2830. 753

754 Gabriel Pereyra, George Tucker, Jan Chorowski, Łukasz
755 Kaiser, and Geoffrey Hinton. 2017. Regularizing
756 neural networks by penalizing confident output dis-
757 tributions. *arXiv preprint arXiv:1701.06548*.

758 Alec Radford, Jeffrey Wu, Rewon Child, David Luan,
759 Dario Amodei, Ilya Sutskever, et al. 2019. Language
760 models are unsupervised multitask learners. *OpenAI*
761 *blog*, 1(8):9.

762 Daniel Ramage, David Hall, Ramesh Nallapati, and
763 Christopher D Manning. 2009. Labeled lda: A su-
764 pervised topic model for credit attribution in multi-
765 labeled corpora. In *Proceedings of the 2009 con-*
766 *ference on empirical methods in natural language*
767 *processing*, pages 248–256.

768 Abigail See, Peter J Liu, and Christopher D Manning.
769 2017. [Get To The Point: Summarization with Pointer-](#)
770 [Generator Networks](#). In *Proceedings of the 2017*
771 *Annual Meeting of the Association for Computational*
772 *Linguistics*, pages 1073–1083.

773 Shengli Song, Haitao Huang, and Tongxiao Ruan. 2019.
774 [Abstractive text summarization using LSTM-CNN](#)
775 [based deep learning](#). *Multimedia Tools and Applica-*
776 *tions*, 78.

777 Vox Media. 2017. [Vox Dataset \(DS+J Workshop\)](#).

778 Zhengjue Wang, Zhibin Duan, Hao Zhang, Chaojie
779 Wang, Long Tian, Bo Chen, and Mingyuan Zhou.
780 2020. [Friendly topic assistant for transformer based](#)
781 [abstractive summarization](#). In *Proceedings of the*
782 *2020 Conference on Empirical Methods in Natural*
783 *Language Processing (EMNLP)*, pages 485–497, On-
784 line. Association for Computational Linguistics.

785 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien
786 Chaumond, Clement Delangue, Anthony Moi, Pier-
787 ric Cistac, Tim Rault, Remi Louf, Morgan Funtow-
788 icz, Joe Davison, Sam Shleifer, Patrick von Platen,
789 Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu,
790 Teven Le Scao, Sylvain Gugger, Mariama Drame,
791 Quentin Lhoest, and Alexander Rush. 2020. [Trans-](#)
792 [formers: State-of-the-art natural language processing](#).
793 In *Proceedings of the 2020 Conference on Empirical*
794 *Methods in Natural Language Processing: System*
795 *Demonstrations*, pages 38–45, Online. Association
796 for Computational Linguistics.

797 Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter
798 Liu. 2020. [PEGASUS: Pre-training with extracted](#)
799 [gap-sentences for abstractive summarization](#). In *Pro-*
800 *ceedings of the 37th International Conference on*
801 *Machine Learning*, volume 119 of *Proceedings of*
802 *Machine Learning Research*, pages 11328–11339.
803 PMLR.

804 Mingyuan Zhou, Lauren Hannah, David Dunson, and
805 Lawrence Carin. 2012. Beta-negative binomial pro-
806 cess and poisson factor analysis. In *Artificial Intelli-*
807 *gence and Statistics*, pages 1462–1471. PMLR.