

# Pre-training LongT5 for Vietnamese Mass-Media Multi-document Summarization Task

Nicolay Rusnachenko<sup>1</sup>[0000–0002–9750–5499], The Anh Le<sup>2,3</sup>[0000–0003–0740–6380],  
and Ngoc Diep Nguyen<sup>3</sup>

<sup>1</sup> [rusnicolay@gmail.com](mailto:rusnicolay@gmail.com)

<sup>2</sup> Vietnam Maritime University, Hai Phong, Viet Nam  
[anhlt@vamaru.edu.vn](mailto:anhlt@vamaru.edu.vn)

<sup>3</sup> Cyber Intellect, Ha Noi, Viet Nam  
[diepnn83@gmail.com](mailto:diepnn83@gmail.com)

**Abstract** Multi-document Summarization task aimed to extract the most salient information from the set of input documents. One of the main challenges that face this task is a long-term dependency problem. When we deal with texts written in Vietnamese it is also accompanied by the specific syllable-based text representation, and lack of labeled datasets. The recent advances in machine translation problem results in a significant impact on the related architecture, dubbed as *transformers*. Being pre-trained on large amounts of raw texts, transformers allows providing a deep knowledge of the texts. In this paper, we survey the findings of the language model applications for text summarization problems, including the remarkable Vietnamese text summarization models. According to the latter, we select LongT5 to pre-train and then fine-tune it for the Vietnamese Multi-document text summarization problem from scratch. We provide a result model analysis and experiments with Multi-document Vietnamese datasets, including ViMs, VMDS, and VLSP2022. We conclude that using a transformer-based model pre-trained on a large amount of unlabeled Vietnamese texts allows us to achieve promising results, with further enhancement via fine-tuning within the small amount of manually summarized texts. The pre-trained model utilized in the experiment section is published<sup>4</sup>.

**Keywords:** Vietnamese Multi-document Summarization · Text Summarization · Transformers · Language Models

## 1 Introduction

At present, the drastically huge growth of news and event recordings becomes one of the main reasons why most of the mass-media platforms become saturated with mass-media information. Such factor becomes a crucial for manual daily news reading making the related approach unfeasible. As a task *text summarization* [12] aims to create a short version of the original texts by keeping the most

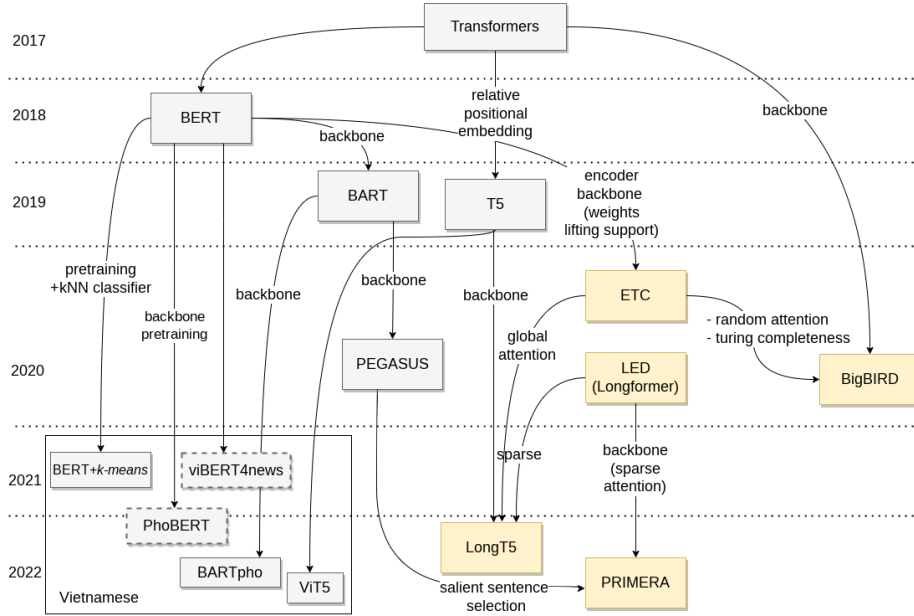
---

<sup>4</sup> <https://github.com/nicolay-r/ViLongT5-newscorpora>

concise, coherent, and salient information. Shortening the long documents by keeping the most meaningful information represents a quite consumptive task for manual execution involving the analysis and content understanding. To the best of our knowledge, such factors necessities studies in automatic text summarization approaches and systems built upon them. In terms of the result summary, such systems might be categorized as *extractive*, or *abstractive*. Summarization systems of extractive type [7] aim to rank sentences in the given text by relying on their meaning and importance, with further extraction of the high-ranked one. In turn, the abstractive type systems are focused on generative result in essay format for a given text [22,9,20].

The appearance of an attention mechanism that addresses the problem of capturing distant information in long input sequences in the Machine Translation (MT) task [2] cause a significant impact on further studies and attention implementations [29,35]. The attention mechanism represents a module in the neural network which aims to assess the importance of the given information by assigning *weights* to its components. A significant amount of investigations were directed to experiments with attention implementations as well as the integration of such modules into target-oriented machine learning models aside from the MT, including the text summarization domain. The further appearance of the *self-attention* mechanism [29] as an internal component of encoder-decoder architecture, results in a *transformer*. Transformer-based models cause a significant breakthrough in MT, resulting in further modifications [34]. The transition towards texts of a single language for transformers results in the appearance of *language models* that become recently both popular and standardized solutions in other natural language processing (NLP) domains including text summarization [9,33,15]. This paper focuses on the analysis of the recent advances of language-models to choose the promising solution for the Vietnamese Multi-document Summarization problem [28,17] of mass-media documents. It is worth noting that Multi-document Summarization faces the problem of the long contents where the importance of information might be spread in different per each document. To the best of our knowledge, we are the first who pretrain and fine-tune the Vietnamese LARGE-sized LongT5 model for Multi-document text summarization from scratch.

The remaining part of this paper is organized as follows. Section 2 provides an overview of the recent advances of transformer-based models along with their architectural updates and training techniques, with Vietnamese texts oriented models in Section 2.1 and sparse attention-based models in Section 2.2. Section 3 lists the resources that were adopted in LongT5 model pre-training and fine-tuning. The detailed description of the model pre-training process as well as the further experiments are covered in Sections 4 and 5 respectively, including comparison with the other extractive and abstractive text summarization baselines for VLSP2022<sub>test</sub> dataset. We accompany our results and relative analysis with text summarization examples during the preliminary stage of the training as well as for the published LongT5 model (ViLongT5<sub>NewsCorpus</sub>) in Appendix A.



**Figure 1.** Tree diagram of the transformer-based models [29], placed in order of their appearance from top to bottom; *arrows* illustrate the most significant findings were found in successor models; *blocks* illustrate models with: original self-attention (gray), sparse self-attention mechanism (yellow); the highlight Vietnamese-targeted models for text summarization problems are bordered, with general application BERT-based models are dotted.

## 2 Background

Since the text summarization problem is commonly treated as an extractive or abstractive tasks, both encoder and decoder components of the transformers could be used as *backbones*. Considering BERT architecture as a backbone, it finds its application in extractive-based text summarization problems. Due to the architecture specifics, which are considered to encode information bidirectionally, BERT could not be easily adopted for generative task format [5]. As for extractive task format, we may consider BERT as a sentence encoder, complemented with a clustering-type algorithm.

However, to address the generative limitations, in [11] the authors propose the BART framework, which represents a BERT (bidirectional transformer) complemented by an autoregressive decoder (GPT). BART proposed a denoising sequence-to-sequence framework, in which the pre-training stage includes: (1) corrupted text restoration and (2) original text reconstruction, i.e. translation. Architecturally, BART is a standard Transformer-based neural machine translation architecture [29] with the potential for customization of encoder and decoder transformer parts, including pre-training schemes modification. Being

particularly effective for text generation tasks, including text summarization, at the time of the model announcement, the authors mention a significant improvement of LARGE-sized BART over previous works on the XSum [14] dataset (Table 1).

BART has become a fundamental architecture for a variety set of text summarization oriented frameworks, such as follows. In [33], authors propose PEGASUS framework, in which the sentence-based masking strategy is based on the invented *salient sentence selection* algorithm. With the latter, authors proposed a sentence assessment metric with a limited selection of the top  $k$ -scored sentences. According to the extensive experiments on XSum and CNN/DailyMail [13] collections with LARGE-sized models (Table 1), authors illustrate that the result PEGASUS model [21] outperforms the other transformer-based solutions, such as BART, T5 [22]. Text-To-Text Transfer Transformer (T5) [22] is based on the original transformer [29] complemented by the following modifications toward layer normalization techniques and token positioning [24]. Analyzing the results of LARGE-based versions, the T5 model with the *principle sentences generation* strategy [33] in pretraining significantly outperforms the rest of the models discussed above on several common datasets (see Table 1).

## 2.1 Vietnamese Multi-document Summarization Models

One of the main traits of Vietnamese texts is *syllable-based* sentence segmentation – the atomic part of sentences are *syllables*. To the best of our knowledge, recent advances in Vietnamese text processing for Multi-document Summarization problems are limited by the original self-attention-based transformers application. Figure 1 illustrates the recent advances in transformer-based models for Vietnamese (bordered bottom-left corner). In this section, we overview the recent advances in *extractive* and *abstractive* text summarization approaches.

For extractive text summarization, several studies of non-transformer-based approaches application [18] address such training techniques examination as: *distant supervision*, *supervised learning*. The adaptation of BERT towards the downstream tasks for texts in Vietnamese results in the appearance of PhoBERT [16], and viBERT4news<sup>5</sup>. In [25] authors combines Vietnamese-oriented BERT-based pretrained states and *k-means*, and the result BERT+*k-means* illustrated top results on the VMDS dataset<sup>6</sup> compared with prior methods.

In the case of abstractive summarization, BARTpho [26] represents an initial study with BART-based [11] architecture application towards the domain of Vietnamese texts. The authors mentioned the importance of vocabulary by gluing syllables into complete words. Due to the latter and in terms of BERT-based approaches, the word-based model performed better than the default, syllable-based representation, and tokenization. Recently, authors of ViT5 [20] experimented with a transformer-based encoder-decoder model for the Vietnamese language, based on the T5 self-supervised pretraining. The latter illustrates the

<sup>5</sup> <https://huggingface.co/NlpHUST/vibert4news-base-cased>

<sup>6</sup> Dataset details in Section 3

**Table 1.** LARGE-sized transformer-based model performances in text summarization problems; models are grouped by self-attention mechanism into original self-attention [29] (512 tokens input limit) and sparsed version (4K+ token input limit); dataset names with best results are bolded; best and second best results are highlighted in gray; according to the results, models with sparse attention tend to perform better due to the longer input sequences

| Model                | Architectural Features  | Dataset              | R-1   | R-2   | R-L   |
|----------------------|---|----------------------|-------|-------|-------|
| BART [11]            | Bidirectional encoder + autoregressive decoder  | XSum                 | 45.14 | 22.27 | 37.25 |
| PEGASUS [33]         | Transformer + Gap-Sentence Selection  | <b>CNN/DailyMail</b> | 44.17 | 21.47 | 41.11 |
|                      |   | Multi-News           | 47.52 | 18.72 | 24.91 |
|                      |   | arXiv                | 44.21 | 16.95 | 38.83 |
|                      |   | CNN/DailyMail        | 43.41 | 20.99 | 40.77 |
| T5 [22]              | Transformer + relative token positioning<br>+ layer norm bias and normalization changes<br>PEGASUS pretraining strategy | Multi-News           | 47.48 | 18.60 | 24.31 |
|                      |   | <b>BigPatent</b>     | 67.05 | 52.24 | 58.70 |
|                      |   | arXiv                | 45.86 | 18.40 | 41.62 |
|                      |   | PubMed               | 48.94 | 22.92 | 45.40 |
| LED (16K) [3]        | Transformer with windowed attention<br>LED + sparse attention (encoder side)  | arXiv                | 46.63 | 19.62 | 41.83 |
|                      |   | arXiv                | 46.63 | 19.02 | 41.77 |
| BigBird-PEGASUS [32] | + random attention mask<br>PEGASUS (PSG) pretraining strategy   | PubMed               | 46.32 | 20.65 | 42.33 |
|                      |   | BigPatent            | 60.64 | 42.46 | 50.01 |
| PRIMERA [31]         | Longformer, Entity Pyramid Strategy   | arXiv                | 47.60 | 20.80 | 42.60 |
|                      |   | <b>Multi-News</b>    | 49.90 | 21.10 | 25.90 |
|                      |   | CNN/DailyMail        | 42.49 | 20.51 | 40.18 |
| LongT5 (4K) [9]      | T5 + global-local attention from LED  | <b>BigPatent</b>     | 70.38 | 56.81 | 62.73 |
|                      |   | <b>arXiv</b>         | 48.28 | 21.63 | 44.11 |
|                      |   | <b>PubMed</b>        | 49.98 | 24.69 | 46.46 |

recent advances in *abstractive text summarization* and *named entity recognition* (NER) problems [20].

## 2.2 Sparse Self-Attention

The main task solved by attention is the connectivity of the particular token with respect to the other mentioned in the text. However, the crucial payment of this solution lies in its computational ineffectiveness. The computational complexity of full self-attention for an input size of  $n$  is  $O(n^2)$  [29]. Besides the BERT [5], such mentioned models as BART [11] and T5 [22] nested the self-attention mechanisms, and hence in practice input sequences tend to be limited by 512 tokens [32].

To address the shortcomings of the self-attention application towards longer input sequences, the series of independent works were accomplished to its *sparse* variations [3,1,32]. To manage attention behavior on that matter, in [1] authors propose Extended Transformer Construction (ETC). Alongside the other works, authors invent *relative token positioning* [3,32] as a preliminary step for attention sparsification. To distribute attention between distant tokens, authors introduce a *global-local* attention mechanism by expanding the original (local) input with *global tokens* under the following restriction: the length of the global token sequence ( $n_g$ ) is expected to be significantly less than the original input sequence length ( $n_l$ ). Considering the latter, authors split attention calculation into parts

and prove the resulting complexity of  $O(n_g^2 + n_g \cdot n_l)$  remains linearly dependent on the original input length  $n_l$ . The relative token positioning encoding as long as sparse attention implementation [32] allows to train ETC with longer input sequences and hence caused a significant impact on the result performance in question answering (QA) tasks [1]. In [3], authors propose Longformer and the related encoder-decoder (LED), which represents a modification of the original transformer with *windowed attention* variations. Similar to implementation in ETC, the latter denotes that, for a particular token only  $r$  (radius parameter) left and right neighbored tokens are considered to attend. Figure 1 illustrates models with sparse-attention mechanisms (yellow colored). Authors experiment with LARGE sized models towards arXiv summarization dataset [4] and illustrate a better performance of LED (447M params) over PEGASUS (4K) and equal to BigBird (4K) once input size has been increased from 4K up to 16K tokens. [3] LED architecture caused a significant effect on models that appear further, such as PRIMERA [31] with salient sentences masking approach, LongT5 [9] described further in this section.

Alongside the findings of the ETC application, in [32] authors treat the computational problem of self-attentive mechanism connection as a *graph sparsification*. Complementing sliding window and global attention mechanisms [1] with Erdős-Rényi model [6] of independently choosing edge with fixed probability, authors aim to prove Turing Completeness of the sparse attention mechanism behind the proposed BigBird model, which is computationally linear in the number of tokens. The outcome of the latter is as follows: the more sparse the graph, the more layers are required to reach completeness. For text summarization, authors experiment with sparse attention at encoder side<sup>7</sup>, using pre-trained schemes from PEGASUS [33] for LARGE-sized models. The resulting model is dubbed BigBird-PEGASUS [32].

LongT5 represents a modified version of the T5 [22] which adopts the sparse attentive mechanism variations, proposed with the ETC model, including windowed attention and global-local variation, dubbed as TGlobal [9]. The latter introduces local sparsity in the attention mechanism, which allows the reduction of the quadratic cost when scaling to long inputs. Unlike T5, the modified LongT5 can handle longer input sequences before reaching the out-of-memory exceptions. It is worth noting LongT5 (4K input) reaching top results across the variety of text-generative on almost every text summarization dataset: arXiv summarization dataset, PubMed, BigPatent [23], and MediaSum [9]. As for PRIMERA (447M) model, the latter illustrates the best results in MultiNews across other models listed in Table 1 due to the specifics and news-related information utilized at the pretraining stage. Analyzing the results across multiple datasets, LongT5 illustrates the best performance across the other models discussed above. The cost of the LongT5 architectural traits is the leveraged amount of the hidden parameters. The LARGE-sized version of LongT5 [21] with a 4K input token size results in  $\approx 780$ M parameters, which is almost two times larger than PRIMERA (447M) and comparable with the size of LED with 16K token input size.

<sup>7</sup> Since the output is relatively short compared with the size of input

### 3 Resources

To the best of our knowledge, there are few Vietnamese single-document summarization datasets and only three Vietnamese multi-document summarization datasets. All of them are abstractive datasets. The details of these datasets are described below, with the brief statistics described in Table 2.

NewsCorpus<sup>8</sup> represents a relatively large collection of 14.9M documents with unlabeled summaries crawled from about 143 Vietnamese news websites. This can be treated as a single-document summarization dataset, in which each document yields the title and sampled content.

VNDS<sup>9</sup> created by Nguyet et al. [19], is a single document summarization dataset, including articles collected from three sources tuoitre.vn, vnexpress.net, and nguoiduatin.vn and classified into *the world*, *news*, *law*, and *business* categories. The body of each article is used as the source document, and its samples treated as a summary. Many documents consist of only one sentence in their summary. The dataset contains  $\approx 151K$  samples, divided into three sets with a ratio of 70%, 15%, and 15% for training, validation, and test sets.

VMDS<sup>10</sup> is a multi-document dataset collected from a Vietnamese online news provider baomoi.com. This dataset contains 600 documents categorized into 200 topics.

ViMs<sup>11</sup> represents a multi-document dataset released by Nghiem et al. [27]. This corpus was collected from different Google News domains. In total, the authors collect 1945 documents from popular news websites in Vietnam.

VLSP2022<sup>12</sup> is a dataset is provided in a competition hosted by the Association for Vietnamese Language and Speech Processing. The provided data is Vietnamese news on various topics, including the economy, society, culture, science, and technology. Every document includes: title, anchor text and body text of single documents, summary, category tag. It is divided into train (VLSP2022<sub>train</sub>) validation (VLSP2022<sub>valid</sub>) and test datasets (VLSP2022<sub>test</sub>). The datasets contain several document clusters. Each cluster has 3-5 documents that illustrate the same topic. There are only 300 samples in training and validation sets in total (VLSP2022<sub>train+valid</sub>). The compression ratio of the summaries provided per every split of the dataset represents 9%.

### 4 Experiential Setup

We experiment with case insensitive LongT5 version with Transient Global Attention mechanism with 2048/512 input/output tokens respectively, and size of the original T5<sub>LARGE</sub><sup>13</sup> [22]. We refer to this model as LongT5<sub>LARGE</sub>-TGlobal

<sup>8</sup> <https://github.com/binhvq/news-corpus>

<sup>9</sup> <https://github.com/ThanhChinhBK/vietnews>

<sup>10</sup> <https://github.com/lupanh/VietnameseMDS>

<sup>11</sup> <https://github.com/CLC-HCMUS/ViMs-Dataset>

<sup>12</sup> <https://vlsp.org.vn/vlsp2022/eval/abmusu>

<sup>13</sup> GIN scripts reference to <https://github.com/nicolay-r/ViLongT5-newscorpora>

**Table 2.** Statistics of the Vietnamese datasets utilized for the model training and evaluation; NewsCorpus dataset represents only raw clustered documents without summaries.

| Dataset                         | #doc       | #samples | #docs<br>per cluster | #words<br>per document | #words<br>per summary |
|---------------------------------|------------|----------|----------------------|------------------------|-----------------------|
| NewsCorpus                      | 14 896 998 | –        | –                    | –                      | –                     |
| VNDS                            | 150 704    | –        | –                    | 41 337                 | 2 848                 |
| VMDS                            | 628        | 300      | 3.0                  | 1 308                  | 153                   |
| ViMs                            | 1 945      | 300      | 6.5                  | 2 208                  | 192                   |
| VLSP2022 <sub>train</sub>       | 621        | 200      | 3.11                 | 1925.75                | 168.48                |
| VLSP2022 <sub>valid</sub>       | 304        | 100      | 3.04                 | 1815.41                | 167.68                |
| VLSP2022 <sub>train+valid</sub> | 925        | 300      | 3.0                  | 1 853                  | 162                   |
| VLSP2022 <sub>test</sub>        | 914        | 300      | 3.05                 | 1762.40                | 153.05                |

(2K/512) in further. Next, we provide the details of input data preparation and organization of the pre-training using Vietnamese datasets described in Section 3.

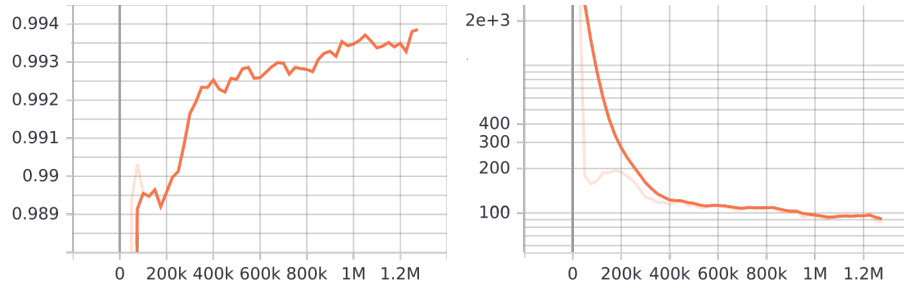
We consider NewsCorpus dataset for the LongT5<sub>LARGE</sub>-TGlobal (2K/512) pretraining. Precisely speaking, we select the first  $10^6$  documents from the whole NewsCorpus. Due to the specifics of this dataset, which consists of raw documents only (Table 2), additional post-processing was applied toward document clustering and summary generation for the composed clusters. We perform the artificial transformations of the documents to the multi-document by interpreting every document as a *cluster* – a list of paragraphs, where every paragraph is considered as a subdocument of the original content. For the preliminary document summarization, we consider the *principle sentence generation* strategy from PEGASUS [33] by relying on the results of the extensive experiments [9]. For each document we select the five most salient sentences by **pyramid-rouge** [31] score. To emphasize a separation between documents in a cluster, we consider an auxiliary document separation token  $\langle doc-sep \rangle$ . To emphasize the end of each sentence and the whole input sequence, we adopt  $\langle sent-sep \rangle$  and  $\langle eos \rangle$  auxiliary tokens respectively.

By default, the core LongT5 [9] is designed for the “Sentence Piece” based tokenization model [10]. To meet these requirements, we then compose case-insensitive Vietnamese language-oriented **SentencePiece** model<sup>14</sup>. To compose this model, we consider original documents from all datasets mentioned in Section 3, with NewsCorpus limited by the first  $10^6$  documents. Due to the specifics of the Vietnamese texts, all syllables were merged into words with the auxiliary “\_” (underscore) character. We apply the stemming and lowercasing. In terms of stemming operation, all the syllables of the related word were concatenated with the underscore character. For such operation, the **VnCoreNLP** [30] library<sup>16</sup> was considered. The size of the result vocabulary was established as 32K tokens.

<sup>14</sup> We adopt the native Google SentencePiece library<sup>15</sup>

<sup>16</sup> We **wseg** annotation type





**Figure 2.** *Accuracy* (left) and *Loss* (right) parameter dynamics during the LongT5<sub>LARGE</sub>-TGlobal (2K/512) pretraining stage over NewsCorpus dataset documents (details in Section 4); Y-axis corresponds to logarithmic-scaled values; X-axis represents the number of steps passed from 0 to 1.4M, where each step involves feed-forward and back-propagation over a single batch.

For the LongT5<sub>LARGE</sub>-TGlobal (2K/512) pretraining, the default configuration and hyperparameters setup was considered [22]. The whole process lasts 12 days and is performed on  $2 \times$  NVIDIA A100 GPUs (40GB each). For such parameters, the maximum possible batch size of the model was considered to be set as 8.

## 5 Result Analysis And Discussion

The pre-training statistics of such parameters as *accuracy* and *loss* are illustrated in Figure 2. We terminate the pre-training process once it reached 1.4M steps over NewsCorpus documents (details in Section 4), where each step includes a feedforward and back-propagation of a whole input batch.

According to the Figure 2 (left), once LongT5<sub>LARGE</sub>-TGlobal (2K/512) reached  $\approx 100K$  pre-training steps, it illustrates a relatively high training accuracy of 0.989, with further parameter value increment up to 0.994. In terms of the *loss* variation, it is possible to investigate a significant decrease within the first  $\approx 600K$  steps and reach the flat once getting closer to 1.4M which finally leads us to the termination of the pretraining process. During the pre-training stage, we examine the preliminary state once the model reached 925K steps. Since texts of the other collections differ from NewsCorpus and provide the manually written summaries, we additionally experiment with finetuning with 5K steps. Table 4 illustrates an example of the preliminary model summarization. It is possible to investigate the resulting model works close to the extractive, with the output summary representing the copies of the original sentence parts (see Table 4).

We use checkpoint of model pre-trained with 1.4M steps to continue finetuning with extra 10K steps on small Vietnamese multi-document summarization datasets, which we divide into the train, validation, and test sets with the proportion of 8:1:1. Appendix section also illustrates the behavior of the evaluated model. According to the results of Table 5 and 6 (Appendix sec-

**Table 3.** Results of the baseline models in comparison with ViLongT5<sub>NewsCorpus</sub> – case-insensitive LongT5<sub>LARGE-TGlobal</sub> (2K/512), pretrained on NewsCorpus with 1.4M steps and finetuned with extra 10K steps on VMDS, ViMs and VLSP2022<sub>train+valid</sub> datasets (training folds); «\*» corresponds to the preliminary state finetuned with 10K steps excluding VLSP2022<sub>valid</sub> dataset; results corresponding to the proposed model for VLSP2022 datasets are highlighted

| Model                                | Dataset  | Rouge Scores (F1) |       |       |        |
|--------------------------------------|--|-------------------|-------|-------|--------|
|                                      |  | R-1               | R-2   | R-L   | AVG. R |
| ViLongT5 <sub>NewsCorpus</sub>       | VLSP2022 <sub>train+valid</sub>                                  | 62.00             | 39.20 | 38.30 | 46.50  |
| ViLongT5 <sub>NewsCorpus</sub>       | VLSP2022 <sub>train+valid</sub> + ViMs + VMDS <sub>(test)</sub>  | 62.90             | 39.60 | 37.20 | 46.50  |
| ViLongT5 <sub>NewsCorpus</sub>       | VLSP2022 <sub>train+valid</sub> + ViMs + VMDS <sub>(valid)</sub> | 52.90             | 33.20 | 33.30 | 39.80  |
| hybrid <sub>the_coach</sub>          | VLSP2022 <sub>valid</sub>  | 51.68             | 31.50 | 48.93 | —      |
| LexRank + MMR <sub>baseline</sub>    | VLSP2022 <sub>valid</sub>  | 48.36             | 26.50 | 44.21 | —      |
| rule <sub>baseline</sub>             | VLSP2022 <sub>valid</sub>  | 46.40             | 25.82 | 42.84 | —      |
| ViLongT5 <sub>NewsCorpus</sub> *     | VLSP2022 <sub>valid</sub>  | 45.70             | 24.83 | 42.85 | —      |
| anchor <sub>baseline</sub>           | VLSP2022 <sub>valid</sub>  | 43.81             | 19.31 | 39.28 | —      |
| ViT5 <sub>abstractive-baseline</sub> | VLSP2022 <sub>valid</sub>  | 31.29             | 30.77 | 27.97 | —      |
| hybrid <sub>the_coach</sub>          | VLSP2022 <sub>test</sub>   | 49.62             | 29.37 | 47.01 | —      |
| LexRank + MMR <sub>baseline</sub>    | VLSP2022 <sub>test</sub>   | 47.72             | 26.25 | 43.39 | —      |
| rule <sub>baseline</sub>             | VLSP2022 <sub>test</sub>   | 46.27             | 26.11 | 42.73 | —      |
| ViLongT5 <sub>NewsCorpus</sub>       | VLSP2022 <sub>test</sub>   | 45.16             | 24.48 | 42.08 | —      |
| anchor <sub>baseline</sub>           | VLSP2022 <sub>test</sub>   | 43.21             | 18.86 | 38.69 | —      |
| ViT5 <sub>abstractive-baseline</sub> | VLSP2022 <sub>test</sub>   | 32.26             | 14.97 | 28.95 | —      |

tion) the specifics of the pretraining, which relies on NewsCorpus raw texts, results in extracting alike behavior of the LongT5<sub>LARGE-TGlobal</sub> (2K/512), with the repetitive manner in some cases (Table 6). Considering the results of behavioral aspects mentioned above, we provide post-processing involving output trimming by keeping only information until the first appeared  $\langle eos \rangle$  in output. Table 3 illustrates the obtained results for: 1) VLSP2022<sub>train+valid</sub>, 2) VLSP2022<sub>train+valid</sub> + ViMs + VMDS (test/validation), and 3) VLSP2022<sub>test</sub> according to the related competitions<sup>17</sup>.

In terms of the VLSP2022<sub>test</sub> assessment, the proposed ViLongT5<sub>NewsCorpus</sub> in 13th out of 20 participants on VLSP2022<sub>valid</sub><sup>18</sup> and 10th place on VLSP2022<sub>test</sub>. Models ranked by R2-F1 measure results. Table 3 lists the results of other baselines as well as the top submissions for comparison (hybrid<sub>the\_coach</sub>). First it is worth to mention that abstractive approaches with generative texts are tend to perform worse than generative in terms of the result assessment systems. Wit that idea in mind, considering the baseline results of the purely extractive and abstractive approaches it is possible to investigate the large gap in the obtained results and importance of the originally salient sentences in the result summary, especially with long-common-sequence assessment (R-L). The

<sup>17</sup> <https://aihub.ml/competitions/341>

<sup>18</sup> Preliminary finetuned version, in which VLSP2022<sub>valid</sub> dataset has been excluded

hybrid approach ( $\text{hybrid}_{\text{the\_coach}}$ ) text summarization approach illustrates the highest result. Application of the LexRank [7] + MMR [8] correspond to extractive baseline approach ranked by 6 and 8 in VLSP2022<sub>valid</sub> and VLSP2022<sub>test</sub> respectively. Results of the  $\text{rule}_{\text{baseline}}$  correspond to the case of selecting first and last sentence for each cluster of the documents, ranked with #7 and #10 in VLSP2022<sub>valid</sub> and VLSP2022<sub>test</sub> respectively. The  $\text{anchor}_{\text{baseline}}$  is a result of the input duplication, results in #19 rank. Model ViT5 has been adopted as a zero-shot abstractive baseline, which was ranked #20.

## 6 Conclusion

The recent appearance of language models significantly addresses the long-range information memorizing is what becomes a result of the vast amount of further studies, focused on context length increment. In this work, we survey the transformers and their variations and evolution toward the internal self-attention mechanism implementations. The main highlights that overcome the main problem of self-attention with its computational complexity were shown. Considering highlights and the lack of their recent application findings for the Vietnamese language, we adopt and experiment with one of the promising models (LongT5) for the abstractive multi-document text summarization in mass-media texts. One of the largest and publicly available NewsCorpus of raw texts has been adopted for the initial pre-training. We experiment with the finetuned version and due to the pre-train specifics investigate with the summaries representing the retelling of the most salient sentences.

## References

1. Ainslie, J., Ontanon, S., Alberti, C., Cvicek, V., Fisher, Z., Pham, P., Ravula, A., Sanghai, S., Wang, Q., Yang, L.: ETC: Encoding long and structured inputs in transformers. In: Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP). pp. 268–284. Association for Computational Linguistics, Online (Nov 2020). <https://doi.org/10.18653/v1/2020.emnlp-main.19>, <https://aclanthology.org/2020.emnlp-main.19>
2. Bahdanau, D., Cho, K., Bengio, Y.: Neural machine translation by jointly learning to align and translate. arXiv preprint arXiv:1409.0473 (2014)
3. Beltagy, I., Peters, M.E., Cohan, A.: Longformer: The long-document transformer. ArXiv **abs/2004.05150** (2020)
4. Cohan, A., Derroncourt, F., Kim, D.S., Bui, T., Kim, S., Chang, W., Goharian, N.: A discourse-aware attention model for abstractive summarization of long documents. arXiv preprint arXiv:1804.05685 (2018)
5. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: 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). pp. 4171–4186. Association for Computational Linguistics, Minneapolis, Minnesota (Jun 2019). <https://doi.org/10.18653/v1/N19-1423>, <https://aclanthology.org/N19-1423>

6. Erdős, P., Rényi, A., et al.: On the evolution of random graphs. *Publ. Math. Inst. Hung. Acad. Sci* **5**(1), 17–60 (1960)
7. Erkan, G., Radev, D.R.: LexRank: Graph-based lexical centrality as salience in text summarization. *Journal of Artificial Intelligence Research* **22**, 457–479 (12 2004). <https://doi.org/10.1613/jair.1523>,
8. Goldstein, J., Carbonell, J.G.: Summarization:(1) using mmr for diversity-based reranking and (2) evaluating summaries. In: *TIPSTER TEXT PROGRAM PHASE III: Proceedings of a Workshop held at Baltimore, Maryland, October 13-15, 1998*. pp. 181–195 (1998)
9. Guo, M., Ainslie, J., Uthus, D., Ontanon, S., Ni, J., Sung, Y.H., Yang, Y.: LongT5: Efficient text-to-text transformer for long sequences. In: *Findings of the Association for Computational Linguistics: NAACL 2022*. pp. 724–736. Association for Computational Linguistics, Seattle, United States (Jul 2022). <https://doi.org/10.18653/v1/2022.findings-naacl.55>, <https://aclanthology.org/2022.findings-naacl.55>
10. Kudo, T., Richardson, J.: SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. pp. 66–71. Association for Computational Linguistics, Brussels, Belgium (Nov 2018). <https://doi.org/10.18653/v1/D18-2012>, <https://aclanthology.org/D18-2012>
11. Lewis, M., Liu, Y., Goyal, N., Ghazvininejad, M., Mohamed, A., Levy, O., Stoyanov, V., Zettlemoyer, L.: BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. pp. 7871–7880. Association for Computational Linguistics, Online (Jul 2020). <https://doi.org/10.18653/v1/2020.acl-main.703>, <https://aclanthology.org/2020.acl-main.703>
12. Luhn, H.P.: The automatic creation of literature abstracts. *IBM Journal of research and development* **2**(2), 159–165 (1958)
13. Nallapati, R., Zhou, B., dos Santos, C., Gulçehre, Ç., Xiang, B.: Abstractive text summarization using sequence-to-sequence RNNs and beyond. In: *Proceedings of the 20th SIGNLL Conference on Computational Natural Language Learning*. pp. 280–290. Association for Computational Linguistics, Berlin, Germany (Aug 2016). <https://doi.org/10.18653/v1/K16-1028>, <https://aclanthology.org/K16-1028>
14. Narayan, S., Cohen, S.B., Lapata, M.: Don’t give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization. *arXiv preprint arXiv:1808.08745* (2018)
15. Nenkova, A., Passonneau, R.: Evaluating content selection in summarization: The pyramid method. In: *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics: HLT-NAACL 2004*. pp. 145–152. Association for Computational Linguistics, Boston, Massachusetts, USA (5 2004), <https://aclanthology.org/N04-1019>
16. Nguyen, D.Q., Nguyen, A.T.: Phobert: Pre-trained language models for vietnamese. *CoRR abs/2003.00744* (2020), <https://arxiv.org/abs/2003.00744>
17. Nguyen, M.T., Nguyen, H.D., Nguyen, T.H.N., Nguyen, V.H.: Towards state-of-the-art baselines for vietnamese multi-document summarization. In: *2018 10th International Conference on Knowledge and Systems Engineering (KSE)*. pp. 85–90 (2018). <https://doi.org/10.1109/KSE.2018.8573420>

18. Nguyen, M.T., Nguyen, H.D., Nguyen, T.H.N., Nguyen, V.H.: Towards state-of-the-art baselines for vietnamese multi-document summarization. In: 2018 10th International Conference on Knowledge and Systems Engineering (KSE). pp. 85–90 (2018). <https://doi.org/10.1109/KSE.2018.8573420>
19. Nguyen, V.H., Nguyen, T.C., Nguyen, M.T., Hoai, N.X.: Vnds: A vietnamese dataset for summarization. In: 2019 6th NAFOSTED Conference on Information and Computer Science. pp. 375–380 (2019). <https://doi.org/10.1109/NICS48868.2019.9023886>
20. Phan, L., Tran, H., Nguyen, H., Trinh, T.H.: ViT5: Pretrained text-to-text transformer for Vietnamese language generation. In: Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Student Research Workshop. pp. 136–142. Association for Computational Linguistics (2022), <https://aclanthology.org/2022.naacl-srw.18>
21. Phang, J., Zhao, Y., Liu, P.J.: Investigating efficiently extending transformers for long input summarization. arXiv preprint arXiv:2208.04347 (2022)
22. Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., Liu, P.J.: Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research* **21**(140), 1–67 (2020), <http://jmlr.org/papers/v21/20-074.html>
23. Sharma, E., Li, C., Wang, L.: BIGPATENT: A large-scale dataset for abstractive and coherent summarization. In: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics. pp. 2204–2213. Association for Computational Linguistics, Florence, Italy (Jul 2019). <https://doi.org/10.18653/v1/P19-1212>, <https://aclanthology.org/P19-1212>
24. Shaw, P., Uszkoreit, J., Vaswani, A.: Self-attention with relative position representations. In: Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers). pp. 464–468. Association for Computational Linguistics, New Orleans, Louisiana (Jun 2018). <https://doi.org/10.18653/v1/N18-2074>, <https://aclanthology.org/N18-2074>
25. To, H.Q., Nguyen, K.V., Nguyen, N.L.T., Nguyen, A.G.T.: Monolingual vs multilingual BERTology for Vietnamese extractive multi-document summarization. In: Proceedings of the 35th Pacific Asia Conference on Language, Information and Computation. pp. 692–699. Association for Computational Linguistics, Shanghai, China (11 2021), <https://aclanthology.org/2021.paclic-1.73>
26. Tran, N.L., Le, D.M., Nguyen, D.Q.: Bartpho: Pre-trained sequence-to-sequence models for vietnamese. In: Proceedings of the 23rd Annual Conference of the International Speech Communication Association (2022)
27. Tran, N.T., Nghiem, M.Q., Nguyen, N.T., Nguyen, N.L.T., Van Chi, N., Dinh, D.: Vims: a high-quality vietnamese dataset for abstractive multi-document summarization. *Language Resources and Evaluation* **54**(4), 893–920 (2020)
28. Ung, V.G., Luong, A.V., Tran, N.T., Nghiem, M.Q.: Combination of features for vietnamese news multi-document summarization. In: 2015 Seventh International Conference on Knowledge and Systems Engineering (KSE). pp. 186–191. IEEE (2015)
29. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. *Advances in neural information processing systems* **30** (2017)

30. Vu, T., Nguyen, D.Q., Nguyen, D.Q., Dras, M., Johnson, M.: VnCoreNLP: A Vietnamese natural language processing toolkit. In: Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Demonstrations. pp. 56–60. Association for Computational Linguistics, New Orleans, Louisiana (Jun 2018). <https://doi.org/10.18653/v1/N18-5012>, <https://aclanthology.org/N18-5012>
31. Xiao, W., Beltagy, I., Carenini, G., Cohan, A.: PRIMERA: Pyramid-based masked sentence pre-training for multi-document summarization. In: Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 5245–5263. Association for Computational Linguistics, Dublin, Ireland (May 2022). <https://doi.org/10.18653/v1/2022.acl-long.360>, <https://aclanthology.org/2022.acl-long.360>
32. Zaheer, M., Guruganesh, G., Dubey, K.A., Ainslie, J., Alberti, C., Ontanon, S., Pham, P., Ravula, A., Wang, Q., Yang, L., et al.: Big bird: Transformers for longer sequences. *Advances in Neural Information Processing Systems* **33** (2020)
33. Zhang, J., Zhao, Y., Saleh, M., Liu, P.J.: Pegasus: Pre-training with extracted gap-sentences for abstractive summarization. In: Proceedings of the 37th International Conference on Machine Learning. ICML’20, JMLR.org (2020)
34. Zheng, Z., Yue, X., Huang, S., Chen, J., Birch, A.: Towards making the most of context in neural machine translation. *arXiv preprint arXiv:2002.07982* (2020)
35. Zheng, Z., Yue, X., Huang, S., Chen, J., Birch, A.: Towards making the most of context in neural machine translation. In: Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence. IJCAI’20 (2021)

## A LongT5<sub>LARGE</sub>-TGlobal (2K/512) Examples

| Input   | Result Summary  |
|---|---|
| tại tuần lễ đổi mới sáng tạo win-win của huawei , giám đốc điều hành hội đồng quản trị kiêm chủ tịch hội đồng quản lý cơ sở hạ tầng ict của huawei , ông david wang , đã có bài phát biểu quan trọng mang tên "ổi mới sáng tạo và thắp sáng kỷ nguyên 5 . 5g " . trong bài phát biểu , ông wang nói nhiều về bước phát triển tiếp theo của công nghệ 5g mà công ty gọi là 5 . 5g , cũng như lộ trình đổi mới sáng tạo của ngành 5-10 năm tới . " tầm nhìn đến năm 2025 , yêu cầu dịch vụ mạng đa dạng và quy mô lớn sẽ tạo ra tiềm năng cực lớn cho thị trường mới "    | tại tuần lễ đổi mới sáng tạo win-win của huawei , giám đốc điều hành hội đồng quản trị kiêm chủ tịch hội đồng quản lý cơ sở hạ tầng ict của huawei , ông david wang , đã có bài phát biểu quan trọng mang tên "ổi mới sáng tạo và thắp sáng kỷ nguyên 5 . 5g <eos> " . trong bài phát biểu , ông wang nói nhiều về bước phát triển tiếp theo của công nghệ 5g mà công ty gọi là 5 . 5g , cũng như lộ trình đổi mới sáng tạo của ngành 5-10 năm tới . tầm nhìn đến năm 2025 , yêu cầu dịch vụ mạng đa dạng và quy mô lớn sẽ tạo ra tiềm năng cực lớn cho thị trường mới <eos>                                |
| <i>at huawei win-win innovation week, the executive director of the board and chairman of the ict infrastructure management board of huawei, mr. david wang, gave an important speech entitled "change" creative new and light up the 5. 5g era". during the keynote , mr. wang talked a lot about the next development of 5g technology which the company calls 5 . 5g , as well as the industry 's innovation roadmap for the next 5-10 years . "vision to 2025, diverse and large-scale network service requirements will create huge potential for new markets"</i> | <i>at huawei win-win innovation week, the executive director of the board and chairman of the ict infrastructure management board of huawei, david wang, gave a keynote speech entitled " Innovate and light up the era of 5. 5g &lt;eos&gt;". In his keynote , Mr. wang talked a lot about the next step of the development of 5g technology which the company calls 5. 5g, as well as the roadmap for the change innovation of the industry in the next 5-10 years with a vision to 2025, diverse and large-scale network service requirements will create huge potential for new markets &lt;eos&gt;</i> |

**Table 4.** Example for LongT5<sub>LARGE</sub>-TGlobal (2K/512), pre-trained with 925K steps with 5K fine-tuning); *top*: input text and the related summary output *bottom*: translation of the input and summary in English (bottom); the model tends to copy entire sentences instead of generating some novel ones.

| Input   | Result Summary   |
|---|--|
| <p>chuyến bay của hãng viva_aerobus mang số_hiệu vb518 , có sức chứa 186 hành_khách , khởi_hành từ guadalajara , mexico , tối 24/8 ( giờ_địa_phương ) . chuyến bay dự_kiến kéo_dài 3 tiếng , tuy_nhiên , khoảng 10 phút sau khi cất_cánh , hành_khách phát_hiện tia lửa bắn ra từ động_cơ bên phải của máy_bay . một_số hành_khách cho biết đã nghe thấy “ âm_thanh của một_vụ_nổ ” . tình_huống này đã khiến nhiều hàng khách hoảng_loạn , một_số người than_khóc , la_hét và cầu_nguyện . phi_hành_đoàn nhận được cảnh_báo khi máy_bay đang ở độ cao tương_đương 4.000 m và đã quay trở_lại hạ_cánh tại sân_bay guadalajara , 45 phút sau khi cất_cánh . rất may_mắn , vụ_việc này không biến thành một_thảm_hoạ hàng_không . viva_aerobus cho biết tia lửa điện trong động_cơ là do “ ma_sát của kim_loại ...</p> <p><i>viva_aerobus airline’s flight number vb518, with a capacity of 186 passengers, departs from guadalajara, mexico, on the evening of august 24 (local time). the flight is expected to last 3 hours, however, about 10 minutes after takeoff, passengers emit sparks from the right engine of the plane. one passengers reported hearing “the sound of an explosion”. this situation caused many customers to panic, some people lamented, cried, shouted and prayed. the flight crew received a warning when the plane was at an altitude equivalent to 4,000 m and turned back and landed at the airport at guadalajara, 45 minutes after takeoff. very fortunately, this incident did not turn into a disaster. viva_aerobus indicates that sparks in the engine_motor are caused by “friction of needle_type ...</i></p> | <p>chuyến bay của hãng viva_aerobus mang số_hiệu vb518, có sức chứa 186 hành_khách, khởi_hành từ guadalajara, mexico, tối 24/8 (giờ_địa_phương). chuyến bay dự_kiến kéo_dài 3 tiếng, tuy_nhiên, khoảng 10 phút sau khi cất_cánh, hành_khách phát_hiện tia lửa bắn ra từ động_cơ bên phải của máy_bay. một_số hành_khách cho biết đã nghe thấy <i>(eos)</i></p> <p><i>flight of viva-aerobus airline with number vb518, line with number vb518, with a capacity of 186 passengers, departing from Guadalajara, mexico, on the evening of August 24 (local time). the flight was scheduled to last 3 hours, however, about 10 minutes after take-off, the passenger detected sparks coming from the plane’s right engine. some_passengers reported hearing</i></p> |

**Table 5.** *top*: input text and the related summary output by LongT5<sub>LARGE</sub>-TGlobal (2K/512) (pre-trained 1.4M steps with 10K fine-tuning); *bottom*: translation of the input and summary in English (bottom); the model tends to copy entire sentences instead of generate some novel ones; model summary was trimmed by keeping all the tokens including the first appearance of the auxiliary end-of-sequence token *(eos)*.



---

Output Summary

---

brordinary\_attorney\_woo được coi như một trong số dự án phim thành công nhất trong năm 2022 . theo điều trị ung thư dạ dày . ” . ” đó là vì con số này cũng bao gồm các trường hợp phát hiện bệnh sớm . với ung thư dạ dày giai đoạn 3 như anh , tỷ lệ sống được 5 năm sau phẫu thuật là trên 70% . ” . ” . ” . ” đó là vì con số này cũng bao gồm các trường hợp phát hiện bệnh sớm . với ung thư dạ dày giai đoạn 3 như anh , tỷ lệ sống được 5 năm sau phẫu thuật chỉ là 30-40% ” . ” . ”

*ordinary attorney woo is considered one of the most successful movie projects in 2022 . under the treatment of gastric cancer. ” . ” that’s because this number also includes cases with early detection. With stage 3 stomach cancer like him, the 5-year survival rate after surgery is over 70%. ” . ” . ” . ” that’s because this number also includes cases with early detection. with stage 3 stomach cancer like you, the 5-year survival rate after surgery is only 30-40% ” . ” . ”*

---

một chuyến bay khác đến los angeles vào sáng 24/8 . *<eos> <eos> <eos> . <eos> <eos> <eos> . <eos> <eos> <eos> . <eos> <eos> <eos>*

*another flight to los angeles on the morning of 24/8 . <eos> <eos> <eos> . <eos> <eos> <eos> <eos> . <eos> <eos> <eos> . <eos> <eos> <eos>*

---

**Table 6.** Example of the LongT5<sub>LARGE</sub>-TGlobal (2K/512) summaries, which represents a repetitive responses.