

LipKey: A Large-Scale News Dataset with Abstractive Keyphrases and Their Benefits for Summarization

Anonymous ACL submission

Abstract

Summaries, keyphrases, and titles are different ways of concisely capturing the content of a document. While most previous work has addressed them separately, in this work, we jointly use the three elements via multi-task training and training as joint structured inputs, in the context of document summarization. We release LipKey, the largest news corpus with human-written summaries, titles, and keyphrases, as well as being the first large-scale Indonesian keyphrase dataset. We find that including keyphrases and titles as additional context to the source document improves transformer-based summarization models.¹

1 Introduction

Key content of an article can be presented in different ways, including summaries, keyphrases, and a title. While most previous research has addressed each element individually (e.g. summary generation (Zhang et al., 2020a; Lewis et al., 2020; Koto et al., 2020a) and keyphrase generation (Meng et al., 2017, 2021)), in this work we investigate the utility of keyphrases and titles in the context of single-document abstractive summarization.

The notion of incorporating keyphrases into summarization systems is not novel, and previous work has utilized *extractive* keyphrases obtained through unsupervised and supervised methods. For instance, traditional summarization models (Zhang et al., 2004; D’Avanzo and Magnini, 2005; Wan et al., 2007; Riedhammer et al., 2010; Qazvinian et al., 2010) and modern neural models (Müngen and Kaya, 2018; Nallapati et al., 2016; Liu et al., 2021) have been combined with the top- k frequent words, TF-IDF, and TextRank (Mihalcea and Tarau, 2004) to obtain keyphrases. Elsewhere, Gehrmann et al. (2018); Li et al. (2020) used words contained

in both the summary and article as keyphrases to improve summarization.

This paper aims to study how *abstractive* keyphrases (that are often not present in the article text) can be incorporated into summarization systems. Compared to *extractive* keyphrases used in previous work, *abstractive* keyphrases potentially better complement abstractive summarization methods. Previous work has been hindered by the unavailability of a large annotated dataset with gold-standard summaries and keyphrases, thus opting for *extractive* keyphrase extraction (Qazvinian et al., 2010; Liu et al., 2021).

We additionally study the utility of titles in summarization. The underlying hypothesis is that titles and keyphrases are concise, complementary representations of an article, and provide relevant clues for summarization. While previous summarization datasets such as CNNDM (Hermann et al., 2015), NYT (Sandhaus, 2008), and XSUM (Narayan et al., 2018) do not include keyphrases and titles, we present a novel large-scale dataset containing both.

Following Koto et al. (2020a), we crawl Liputan6² — an Indonesian news portal — to obtain 105K news articles with titles, abstractive summaries, and *abstractive* keyphrases, all authored by journalists. Note that the dataset of Koto et al. (2020a) is based on the time period 2000–2010, at which point Liputan6 did not include keyphrases, while our dataset is based on the time period 2019–2021.³ Furthermore, the fact that the dataset is in Indonesian contributes to language diversity in NLP (Joshi et al., 2020).

To summarize our contributions: (1) we release LipKey, the largest news corpus containing human-written summaries and keyphrases, as well as being the first large-scale Indonesian keyphrase dataset; (2) through extensive experimentation, we benchmark multi-task training and structured in-

¹Data and code used is available at: <https://anon.com>

²<https://www.liputan6.com>

³Koto et al. (2020a) also do not release the titles.

Dataset / Lang	Size	Includes Summ?	#Key per doc (%)	AbsKey (%)
LipKey (ours) / id	105,574	Yes	4.5	51.2
DUC-2001 / en	308	Yes	8.1	3.7
PT-BN-KP / en	110	No	23.7	2.5
KPCrowd / en	500	No	48.9	13.5
KPTimes / en	289,923	No	5.0	54.8
WikiNews / fr	100	No	11.8	5.0

Table 1: LipKey and other keyphrase datasets in the news domain. “AbsKey” is the percentage of “absent” keyphrases, relative to the source article.

put methods using keyphrases and titles for Indonesian text summarization over different pre-trained language models. We find that incorporating keyphrases and titles as structured inputs performs better than multi-task training, and consistently improves summary quality.

2 Related Work

Most keyphrase datasets are in the domain of English scientific publications (Hulth, 2003; Krapivin et al., 2009; Kim et al., 2010; Meng et al., 2021). In Table 1, we compare our corpus, LipKey, with other keyphrase datasets in the news domain. Most datasets such as DUC-2001 (Wan and Xiao, 2008), PT-BN-KP (Marujo et al., 2012), KPCrowd (Marujo et al., 2011), and WikiNews (Bougouin et al., 2013) are small in size and consist of *highly extractive* keyphrases, with KPTimes (Gallina et al., 2019) being the only exception. DUC-2001 is the only dataset with both keyphrases and summaries, but has only 308 documents. In comparison, LipKey is a large news corpus that includes human-written summaries and *abstractive* keyphrases, as well as being the first large-scale Indonesian keyphrase dataset.

Incorporating keyphrases into summarization has been explored in other languages such as Chinese (Jiang et al., 2018; Mihalcea and Tarau, 2004), but using *extractive* keyphrases. This is the first work to combine the two tasks in the Indonesian language, with previous work separately tackling: (1) keyphrase extraction, over Twitter (Mahfuzh et al., 2019), consumer-health questions (Saputra et al., 2018), or scientific articles (Asrori et al., 2020; Trisna and Nurwidyanoro, 2020) with limited data;⁴ or 2) document summarization in the news domain (Kurniawan and Louvan, 2018; Koto et al., 2020a).

⁴None of the datasets are publicly available.

	Vocab	#Word		#Sentence	
		mean	std	mean	std
Article	346,564	436.5	277.7	22	17.5
Summary	63,086	19	6.6	1.2	0.4
Title	58,113	10.1	2.2	1	0
Keyphrases	33,976	8.9	4.6	4.5	1.9

Table 2: Per-article summary statistics for LipKey. For keyphrases, #sentence indicates #keyphrases.

Dataset	Size	% of novel n -gram			
		1	2	3	4
IndoSum	18,764	3.1	10.8	16.2	20.3
Liputan6	215,827	12.9	41.6	57.6	66.9
LipKey (summary)	105,574	7.5	25.2	35.1	40.9
LipKey (title)	105,574	26.8	65.4	84.5	92.7

Table 3: Abtractiveness of summaries (and titles) in Indosum, Liputan6, and LipKey, compared to the article.

3 Data Construction

Liputan6 is one of the largest Indonesian news portals, containing news on topics such as politics, health, business, and popular culture.⁵ Koto et al. (2020a) found that Liputan6 summaries are highly abstractive, written by journalists, and suitable for Indonesian text summarization research. The summary and keyphrases are encapsulated in javascript variables `window.kmklabs.article` with the keys `shortDescription` and `keywords`, respectively.⁶

LipKey articles span the period December 2019 to March 2021, and each article is associated with a summary, title, and keyphrase(s). Given the time period, there is a prevalence of COVID-19 news in the data (see the Appendix for data examples). In Table 2 and Table 3, we show the overall data statistics of LipKey, and compare it with previous Indonesian summarization datasets: IndoSum (Kurniawan and Louvan, 2018) and Liputan6 (Koto et al., 2020a). We observe that summaries in LipKey are more abstractive than IndoSum in terms of novel n -grams (computed relatively to the article). Interestingly, we found that LipKey’s titles are even more abstractive than the summaries in all datasets. Note that the median summary length in LipKey is one sentence, and shorter than Liputan6 (Koto et al., 2020a) at two sentences, de-

⁵According to <https://www.alexacom>, Liputan6 was ranked 16th and 308th in Indonesia and worldwide, respectively, in November 2021 in terms of popularity.

⁶In 2012, Liputan6 added keyphrases for articles. These keyphrases are also assigned manually by the journalist.

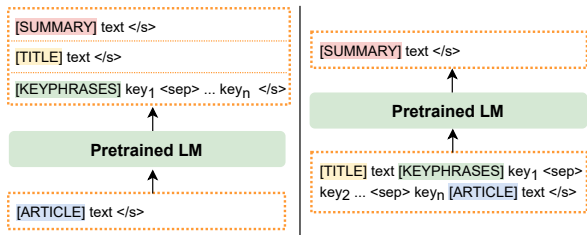


Figure 1: Experimental set-up. **Left**: multi-task training, **Right**: training with structured input.

143 spite both datasets being crawled from the same
 144 news portal.

145 In constructing LipKey, we discard instances
 146 where: (1) one of the keyphrases has more than
 147 6 words (which tends to be noise); (2) the arti-
 148 cle has less than 15 words; or (3) the summary
 149 has less than 5 words. This results in 105,574 in-
 150 stances that we split into 96,573/4,156/4,845 for
 151 train/development/test. In the final dataset, 33%
 152 and 43% of keyphrases consist of 1 and 2 words,
 153 respectively, with the remainder being 3–6 words
 154 (see Table 9 in the Appendix).

155 To better understand the abstractiveness of
 156 keyphrases in LipKey, we randomly sampled
 157 100 articles and manually examined “abstractive”
 158 keyphrases which didn’t occur in the article. We
 159 find that 80% of these partially matched the arti-
 160 cle or were word-order variants (see Table 8 in the
 161 Appendix). Moreover, 15%, 12% and 14% of ab-
 162 stractive keyphrases were acronyms, synonyms, or
 163 morphological variants.

164 4 Experiments

165 4.1 Set-Up

166 As described in Figure 1, we experiment in two
 167 settings: (1) multi-task training (title/keyphrases
 168 = output); and (2) training with structured input
 169 (title/keyphrases = input). For the first, we use sum-
 170 mary s , title t , and keyphrase(s) k as the separate
 171 target texts, and perform multi-task training with
 172 article a as the source text. The total loss \mathcal{L} for
 173 multi-task training is defined as $\mathcal{L}_s + \mathcal{L}_t + \mathcal{L}_k$. For
 174 the second, the goal is to learn $P(s|t, k, a)$ that is
 175 realized by concatenating title t , keyphrases k , and
 176 article a to form the source text, and use summary s
 177 as the target text. To distinguish the four text types
 178 and structure the input, we introduce the special to-
 179 kens of [SUMMARY], [TITLE], [KEYPHRASES], and
 180 [ARTICLE] for all pretrained language models. In
 181 the case of multiple keyphrases, we use <sep>
 182 as a separator. The maximum token length for

Model	R1	R2	R3	Foc.	Cov.
Lead-1	36.6	26.1	34.1	58.5	71.8
Oracle	69.2	58.9	66.9	76.4	87.2
IndoBERT (base) + raw decoder with 153M parameters					
summary	41.8	30.1	39.3	66.6	73.8
<i>multi-task training</i>					
summary, keyphrase	41.8	30.1	39.3	66.6	73.9
summary, title	42.9	31.1	40.4	67.0	74.4
summary, keyphrase, title	42.6	31.1	40.2	66.8	74.5
<i>training with additional context</i>					
+ keyphrase	43.4	31.8	41.0	67.2	<u>74.6</u>
+ title	43.2	31.5	40.7	<u>67.4</u>	74.3
+ keyphrase + title	<u>43.7</u>	<u>31.9</u>	<u>41.2</u>	<u>67.4</u>	<u>74.6</u>
mBART (large) with 600M parameters					
summary	43.1	31	40.5	67.6	73.9
<i>multi-task training</i>					
summary, keyphrase	43.6	31.3	41.0	68.1	74.0
summary, title	42.2	30.0	39.5	67.3	73.4
summary, keyphrase, title	43.5	31.6	40.8	67.8	74.1
<i>training with additional context</i>					
+ keyphrase	43.5	31.2	40.9	68.1	73.8
+ title	43.1	30.8	40.4	67.7	73.8
+ keyphrase + title	<u>44.8</u>	<u>32.3</u>	<u>42.0</u>	68.8	<u>74.6</u>
mT5 (base) with 580M parameters					
summary	45.2	33.7	42.7	67.5	76.2
<i>multi-task training</i>					
summary, keyphrase	44.7	33.2	42.1	66.9	76.3
summary, title	44.6	33.1	42.0	66.6	76.4
summary, keyphrase, title	43.7	32.0	41.0	66.1	76.0
<i>training with additional context</i>					
+ keyphrase	46.4	34.8	43.8	68.2	76.6
+ title	45.4	33.8	42.9	67.5	76.4
+ keyphrase + title	46.7	35.1	44.2	<u>68.4</u>	76.9

Table 4: Summarization results on LipKey. “Foc” and “Cov” are Focus and Coverage, resp., of FFCI. Entries in bold and underline refer to the best overall score and the best score for each model, respectively. “Oracle” is obtained by greedily selecting the subset of sentences in the article that maximizes the ROUGE score based on the reference summary.

183 the article is 512, and for the summary, title, and
 184 keyphrases it is 100.

185 We use the huggingface PyTorch framework
 186 (Wolf et al., 2020) for our experiments with three
 187 pretrained language models: IndoBERT⁷ (Koto
 188 et al., 2020b), mT5 (base)⁸ (Xue et al., 2021), and
 189 mBART (large)⁹ (Liu et al., 2020). For the monolin-
 190 gual IndoBERT, we follow Liu and Lapata (2019)
 191 in adding a raw transformer decoder (layers = 6,
 192 hidden size = 768, feed-forward = 2,048, and heads
 193 = 8) on top of IndoBERT, and train it on 4 × V100
 194 16GB GPUs for 200K steps. For the multilingual

⁷indolem/indobert-base-uncased

⁸google/mt5-base

⁹facebook/mbart-large-50

Model	R1	Foc.	Cov.	$F_1@5$	$F_1@O$	$F_1@M$
RAKE	7.8	40.0	58.7	1.0	1.0	1.0
IndoBERT	58.8	74.2	79.4	45.5	45.2	46.5
mT5 (base)	62.0	75.7	81.7	53.3	52.9	54.5
mBART (large)	63.4	76.4	81.9	54.5	54.4	56.0

Table 5: Keyphrase generation results on LipKey.

mT5 and mBART, we train them on $4 \times V100$ 32GB GPUs for 60 epochs (around 20K steps) with an initial learning rate of $1e-4$ (Adam optimizer). We pick the best checkpoint based on ROUGE scores (Lin, 2004) on the development set (See the Appendix for more details of hyper-parameters).

Additionally, we train keyphrase generation (KPG) models (Seq2Seq) with the same architectures and configurations as the summarization models. We compare the generated keyphrases with: (a) human-written keyphrases; and (b) keyphrases from RAKE, an unsupervised language-independent keyphrase extraction method (Rose et al., 2010).

For evaluating the summarization models, we use F1 of ROUGE scores (R1, R2, and R3), and Focus and Coverage from the FFCI framework (Koto et al., 2020), computed based on Precision and Recall of BERTSCORE (Zhang et al., 2020b) using mBERT uncased.¹⁰ For evaluating KPG, we use macro-averaged $F_1@5$, $F_1@O$, and $F_1@M$, following Meng et al. (2021), and additionally report R1, Focus, and Coverage. Detailed definitions of the metrics are provided in the Appendix.

4.2 Results

In Table 4, we show the full experimental results on the test set. First, we observe that vanilla models (trained only using the article) substantially outperform Lead-1 for all models.¹¹ We find that the vanilla model of mT5 performs better than IndoBERT and mBART, with an improvement of +3.4 and +2.1 R1, respectively.

Training with additional context as structured input consistently improves over multi-task training, with the best results generally being obtained with both keyphrases and title, and mT5 being the best model. When incorporating each element separately, keyphrases are generally better than titles, improving over the vanilla model, with IndoBERT (with multi-task training) being the notable excep-

¹⁰For details of BERT layer selection, see Koto et al. (2021).

¹¹We choose Lead-1 because the average #sentence of the summary is 1.2 in Table 2.

Model	R1	R2	R3	Foc.	Cov.
Vanilla	45.2	33.7	42.7	67.5	76.2
+ keyphrases (RAKE)	44.8	33.3	42.3	66.5	75.6
+ keyphrases (Seq2Seq*)	46.0	34.4	43.5	68.1	76.4
+ keyphrases (Human)	46.4	34.8	43.8	68.2	76.6
Vanilla + title	45.4	33.8	42.9	67.5	76.4
+ keyphrases (RAKE)	43.7	32.1	41.2	67.3	76.1
+ keyphrases (Seq2Seq*)	45.9	34.1	43.3	67.9	76.4
+ keyphrases (Human)	46.7	35.1	44.2	68.4	76.9

Table 6: Ablation study of mT5 (base) over different keyphrases on test set. * denotes using mBART (large).

tion. We also observe that mBART (large) and mT5 (base) are similar in parameter size (600M), but mT5 is substantially better. The FFCI framework shows that both models have similar Focus (= precision), but mT5 has higher Coverage (= recall).

Next, in Table 5, we present results for keyphrase generation on the LipKey test set, and observe that mBART (large) achieves the best performance across all metrics. Interestingly, RAKE performs very poorly,¹² in part emphasizing the limitations of the extractive RAKE method (vs. the highly *abstractive* keyphrases in LipKey).

Lastly, we perform an ablation study over the best summarization model, mT5, using keyphrases sourced through three different methods: (1) RAKE, (2) Seq2Seq, and (3) human-assigned. As seen in Table 6, adding RAKE keyphrases hurts summarization results, but when using Seq2Seq keyphrases (generated by mBART), the performance consistently improves across all metrics, close to the performance of human-assigned keyphrases. Considering this finding, it would be interesting to explore the transferability of keyphrase generation models to other languages, to see if the result can be reproduced.

5 Conclusion

In this paper, we release LipKey, the largest news corpus with human-written summaries, titles, and keyphrases, which is also the first-large scale Indonesian keyphrase dataset. We experimented with incorporating keyphrases (and titles) into summarization training via multi-task training or as structured inputs, and found that the latter works better. Our results also show that abstractive keyphrases benefit summarization systems more than extractive ones.

¹²For each article, we pick the top-5 keyphrases based on RAKE scoring.

6 Ethical Considerations

According to Indonesian Copyright Law number 28 year 2014 article 44, the use, retrieval, reproduction, and/or change of works and/or related rights products in whole or substantial part are not regarded as a copyright infringement if the source is mentioned or cited in full for the purpose of education and research.¹³

References

Riris Bayu Asrori, Robert Setyawan, and Muljono Muljono. 2020. Performance analysis graph-based keyphrase extraction in indonesia scientific paper. In *2020 International Seminar on Application for Technology of Information and Communication (iSemantic)*.

Adrien Bougouin, Florian Boudin, and Béatrice Daille. 2013. TopicRank: Graph-based topic ranking for keyphrase extraction. In *Proceedings of the Sixth International Joint Conference on Natural Language Processing*, pages 543–551, Nagoya, Japan. Asian Federation of Natural Language Processing.

Ernesto D’Avanzo and Bernado Magnini. 2005. A keyphrase-based approach to summarization: the lake system at duc-2005. In *Proceedings of DUC*.

Ygor Gallina, Florian Boudin, and Beatrice Daille. 2019. KPTimes: A large-scale dataset for keyphrase generation on news documents. In *Proceedings of the 12th International Conference on Natural Language Generation*, pages 130–135, Tokyo, Japan. Association for Computational Linguistics.

Sebastian Gehrmann, Yuntian Deng, and Alexander Rush. 2018. Bottom-up abstractive summarization. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4098–4109, Brussels, Belgium. Association for Computational Linguistics.

Karl Moritz Hermann, Tomáš Kočiský, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend. In *NIPS’15 Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 1*, volume 28, pages 1693–1701.

Anette Hulth. 2003. Improved automatic keyword extraction given more linguistic knowledge. In *Proceedings of the 2003 Conference on Empirical Methods in Natural Language Processing*, pages 216–223.

Xiaoping Jiang, Po Hu, Liwei Hou, and Xia Wang. 2018. Improving pointer-generator network with

keywords information for chinese abstractive summarization. In *CCF International Conference on Natural Language Processing and Chinese Computing*, pages 464–474.

Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. 2020. The state and fate of linguistic diversity and inclusion in the NLP world. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6282–6293, Online. Association for Computational Linguistics.

Su Nam Kim, Olena Medelyan, Min-Yen Kan, and Timothy Baldwin. 2010. Semeval-2010 task 5 : Automatic keyphrase extraction from scientific articles. In *Proceedings of the 5th International Workshop on Semantic Evaluation*, pages 21–26.

Fajri Koto, Timothy Baldwin, and Jey Han Lau. 2020. Ffci: A framework for interpretable automatic evaluation of summarization. *arXiv preprint arXiv:2011.13662*.

Fajri Koto, Jey Han Lau, and Timothy Baldwin. 2020a. Liputan6: A large-scale Indonesian dataset for text summarization. In *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing*, pages 598–608, Suzhou, China. Association for Computational Linguistics.

Fajri Koto, Jey Han Lau, and Timothy Baldwin. 2021. Evaluating the efficacy of summarization evaluation across languages. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 801–812, Online. Association for Computational Linguistics.

Fajri Koto, Afshin Rahimi, Jey Han Lau, and Timothy Baldwin. 2020b. IndoLEM and IndoBERT: A benchmark dataset and pre-trained language model for Indonesian NLP. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 757–770, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Mikalai Krapivin, Aliaksandr Autaeu, and Maurizio Marchese. 2009. Large dataset for keyphrases extraction.

Kemal Kurniawan and Samuel Louvan. 2018. Indosum: A new benchmark dataset for indonesian text summarization. In *2018 International Conference on Asian Language Processing (IALP)*, pages 215–220.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. 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*

¹³<https://wipolex-res.wipo.int/edocs/lexdocs/laws/en/id/id064en.pdf>

- 487 I Nyoman Prayana Trisna and Arif Nurwidyantoro.
488 2020. Single document keywords extraction in ba-
489 hasa indonesia using phrase chunking. *TELKOM-*
490 *NIKA Telecommunication Computing Electronics*
491 *and Control*, 18(4):1917–1925.
- 492 Xiaojun Wan and Jianguo Xiao. 2008. Single doc-
493 ument keyphrase extraction using neighborhood
494 knowledge. In *AAAI’08 Proceedings of the 23rd na-*
495 *tional conference on Artificial intelligence - Volume*
496 *2*, pages 855–860.
- 497 Xiaojun Wan, Jianwu Yang, and Jianguo Xiao. 2007.
498 [Towards an iterative reinforcement approach for si-](#)
499 [multaneous document summarization and keyword](#)
500 [extraction](#). In *Proceedings of the 45th Annual Meet-*
501 *ing of the Association of Computational Linguistics*,
502 pages 552–559, Prague, Czech Republic. Associa-
503 tion for Computational Linguistics.
- 504 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien
505 Chaumond, Clement Delangue, Anthony Moi, Pier-
506 ric Cistac, Tim Rault, Remi Louf, Morgan Funtow-
507 icz, Joe Davison, Sam Shleifer, Patrick von Platen,
508 Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu,
509 Teven Le Scao, Sylvain Gugger, Mariama Drame,
510 Quentin Lhoest, and Alexander Rush. 2020. [Trans-](#)
511 [formers: State-of-the-art natural language process-](#)
512 [ing](#). In *Proceedings of the 2020 Conference on Em-*
513 *pirical Methods in Natural Language Processing:*
514 *System Demonstrations*, pages 38–45, Online. Asso-
515 ciation for Computational Linguistics.
- 516 Linting Xue, Noah Constant, Adam Roberts, Mi-
517 hir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya
518 Barua, and Colin Raffel. 2021. [mT5: A massively](#)
519 [multilingual pre-trained text-to-text transformer](#). In
520 *Proceedings of the 2021 Conference of the North*
521 *American Chapter of the Association for Computa-*
522 *tional Linguistics: Human Language Technologies*,
523 pages 483–498, Online. Association for Computa-
524 tional Linguistics.
- 525 Jingqing Zhang, Yao Zhao, Mohammad Saleh, and
526 Peter Liu. 2020a. Pegasus: Pre-training with ex-
527 tracted gap-sentences for abstractive summarization.
528 In *ICML 2020: 37th International Conference on*
529 *Machine Learning*, volume 1, pages 11328–11339.
- 530 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q.
531 Weinberger, and Yoav Artzi. 2020b. Bertscore:
532 Evaluating text generation with bert. In *ICLR 2020*
533 *: Eighth International Conference on Learning Rep-*
534 *resentations*.
- 535 Yongzheng Zhang, Evangelos Milios, and Nur Zincir-
536 Heywood. 2004. A comparison of keyword-and
537 keyterm-based methods for automatic web site sum-
538 marization. In *AAAI04 Workshop on Adaptive Text*
539 *Extraction and Mining*, pages 15–20.

A Overview and Analysis of Keyphrases

Text	% Avg. of present keyphrases	
	Keyphrases level	Word level
Article	66.8	86.9
Summary	39.8	55.3
Title	48.2	62.9

Table 7: Proportion of keyphrases which match article, summary, and title

Category	%	Examples (keyphrase vs. article)
Acronym	15	<i>manchester united</i> vs. <i>man utd</i>
Synonym	12	<i>kepribadian</i> “characteristic” vs. <i>sifat</i> “characteristic”
Morphology	14	<i>pemotor</i> “motorcyclist” vs. <i>motor</i> “motorcycle”
Different order or partial	80	<i>Virus Corona di Aceh</i> “coronavirus in Aceh” vs. <i>Virus Corona</i> “coronavirus”
Not a synonym but related	44	N/A
Found in title (not in article)	24	N/A

Table 8: Analysis of keyphrases from 100 random samples.

#Word	Freq	Example
1	167,203	<i>COVID-19; Netflix</i>
2	214,672	<i>New Normal; Diego Michels</i>
3	84,162	<i>Klasemen Liga Inggris</i> “Premier League”
4	22,780	<i>Ganjil Genap Kota Bogor</i> “odd-even policy in Bogor”
5	6,366	<i>Kru KM Lambelu Positif Covid-19</i> “KM Lambelu Crew Positive Covid-19”
6	1,623	<i>Cara Menulis Daftar Pustaka dari Internet</i> “ways to write a bibliography from Internet”

Table 9: Frequency of keyphrases in LipKey based on #Word.

B Training configurations

Summarization and keyphrase generation use the same models and architecture. For IndoBERT, we follow the Liu and Lapata (2019) architecture by adding a raw transformer decoder (layers = 6, hidden size = 768, feed-forward = 2,048, and heads = 8) on top of IndoBERT, and train it on 4×V100 16GB GPUs for 200K steps with the Adam optimizer and learning rate $lr = 2e^{-3} \cdot \min(\text{step}^{-0.5}, \text{step} \cdot 20,000^{-1.5})$ and $0.1 \cdot \min(\text{step}^{-0.5}, \text{step} \cdot 10,000^{-1.5})$ for IndoBERT and the transformer decoder, respectively. We use a warmup of 20,000, a dropout of 0.2, a batch size total of 200 (10 x 4 GPUs x gradient accumulation of 5), and save checkpoints for every 10,000 steps.

We compute ROUGE scores (R1) to pick the best checkpoint based on the development set.

For mT5 and mBART, we train them on 4×V100 32GB GPUs for 60 epochs (around 20K steps) with an initial learning rate of 1e-4 (Adam optimizer). We use a total batch size of 400 (10 x 4 GPUs x gradient accumulation of 10), a warmup of 10% of total steps, and save checkpoints for every 1,000 steps. We also compute ROUGE scores (R1) to pick the best checkpoint based on the development set.

C Evaluation Metrics

For summarization, we use ROUGE scores (Lin, 2004), and Focus and Coverage from the FFCI framework (Koto et al., 2020). Following Koto et al. (2021), for non-English text, focus and coverage are computed based on the precision and recall of BERTSCORE (Zhang et al., 2020b) using mBERT uncased at layers 12 and 6, respectively. For Y and Y' as the reference and system summary, BERTSCORE is computed as follows:

$$\mathcal{P}_{\text{BERT}} = \frac{1}{|Y'|} \sum_{t_i \in Y'} \max_{s_j \in Y} t_i^T s_j$$

$$\mathcal{R}_{\text{BERT}} = \frac{1}{|Y|} \sum_{s_j \in Y} \max_{t_i \in Y'} t_i^T s_j$$

$$\mathcal{F}_{\text{BERT}} = 2 \frac{\mathcal{P}_{\text{BERT}} \cdot \mathcal{R}_{\text{BERT}}}{\mathcal{P}_{\text{BERT}} + \mathcal{R}_{\text{BERT}}}$$

where s_j and t_i are token embeddings of Y and Y' .

For evaluating the keyphrase generation model, we use macro-averaged $F_1@5$, $F_1@O$, and $F_1@M$, following Meng et al. (2021). Given gold-standard keyphrases \mathcal{Y} and the prediction $\hat{\mathcal{Y}} = \{y'_1, \dots, y'_n\}$, we truncate the prediction to $\hat{\mathcal{Y}} = \{y'_1, \dots, y'_{\min(k,m)}\}$ when only the top k predictions are used for evaluation. Precision, Recall, and F_1 are consequently conditioned on k , and computed as follows:

$$\text{P}@k = \frac{|\hat{\mathcal{Y}}_{:k} \cap Y|}{|\hat{\mathcal{Y}}_{:k}|}$$

$$\text{R}@k = \frac{|\hat{\mathcal{Y}}_{:k} \cap Y|}{|\mathcal{Y}|}$$

$$F_1@k = \frac{2 * \text{P}@k * \text{R}@k}{\text{P}@k + \text{R}@k}$$

Thus $F_1@5$ is $F_1@k$ when $k = 5$, $F_1@O$ is $F_1@k$ when k is the number of oracle (ground truth) keyphrases, and $F_1@M$ is when $k = |\hat{\mathcal{Y}}|$.

Indonesian	English (translation)
<p>Title: Ada Warga Positif Corona di KRL, Ini Kata Kemenhub</p> <p>Gold Keyphrases: krl, COVID-19, Corona</p> <p>Article: Liputan6 . com , Jakarta Kementerian Perhubungan (Kemenhub) memastikan pelaksanaan protokol di Kereta Rangkaian Listrik (KRL) Jabodetabek terus berjalan . Pernyataan ini dikeluarkan pasca adanya 3 penumpang asal Bogor yang dinyatakan positif corona pasca dilakukan test swab . Juru Bicara Kementerian Perhubungan Adita Irawati menyatakan , pihaknya telah mengeluarkan Permenub Nomor 18/2020 yang telah mengatur operasional moda transportasi di masa pandemi . Khususnya pula di daerah yang telah menjalankan Pembatasan Sosial Berskala Besar (PSBB) seperti di Jabodetabek . Perlu dipahami bahwa penularan Covid-19 bisa terjadi dimana saja , tidak hanya di di KRL , " kata Adita , Selasa (5/5/2020) . Adita mengatakan , Permenhub 18/2020 secara tegas telah menyatakan adanya beberapa syarat yang wajib dipenuhi penumpang moda transportasi publik seperti KRL . Pertama , penumpang wajib menggunakan masker . Kedua , sambungnya , petugas mengecek suhu tubuh penumpang .</p> <p>[254 words are abbreviated from here]</p> <p>Gold summaries: Kemenhub menyebutkan Permenhub 18/2020 secara tegas telah menyatakan adanya beberapa syarat yang wajib dipenuhi penumpang moda transportasi publik seperti KRL .</p> <p>IndoBERT: Kementerian perhubungan (Kemenhub) memastikan pelaksanaan protokol di kereta rangkaian listrik (KRL) Jabodetabek terus berjalan</p> <p>IndoBERT with additional contexts (+ keyphrases + titles): Kemenhub memastikan pelaksanaan protokol di kereta rangkaian listrik (KRL) jabodetabek terus berjalan .</p> <p>mBART: Permenhub 18/2020 secara tegas telah menyatakan adanya beberapa syarat yang wajib dipenuhi penumpang moda transportasi publik seperti KRL.</p> <p>mBART with additional contexts (+ keyphrases + titles): Adita mengatakan, Permenhub 18/2020 secara tegas telah menyatakan adanya beberapa syarat yang wajib dipenuhi penumpang moda transportasi publik seperti KRL.</p> <p>mT5: Kemenhub memastikan pelaksanaan protokol di Kereta Rangkaian Listrik (KRL) Jabodetabek terus berjalan.</p> <p>mT5 with additional contexts (+ keyphrases + titles): Juru Bicara Kementerian Perhubungan Adita Irawati menyatakan, pihaknya telah mengeluarkan Permenub Nomor 18/2020 yang telah mengatur operasional moda transportasi di masa pandemi.</p>	<p>Title: Corona positive passengers are detected on the KRL, this is what the Ministry of Transportation says</p> <p>Gold Keyphrases: krl, COVID-19, Corona</p> <p>Article: Liputan6 . com , Jakarta The Ministry of Transportation (Kemenhub) ensures that the implementation of the protocol on the Jabodetabek Electric Circuit Train (KRL) continues. This statement was issued after 3 passengers from Bogor were tested positive for corona after a swab test was carried out. Spokesman for the Ministry of Transportation, Adita Irawati, stated that her party had issued Permenub No. 18/2020 which had regulated the operation of transportation modes during the pandemic. This is particularly the case in areas that have implemented Large-Scale Social Restrictions (PSBB) such as in Jabodetabek. It should be understood that the transmission of Covid-19 can occur anywhere, not only in KRL," said Adita, Tuesday (5/5/2020). First, passengers are required to wear masks. Second, he continued, officers check passengers' body temperatures.</p> <p>[254 words are abbreviated from here]</p> <p>Gold summaries: The Ministry of Transportation stated that Permenhub 18/2020 has explicitly stated that there are several conditions that must be met by passengers of public transportation modes such as KRL.</p> <p>IndoBERT: The Ministry of Transportation (Kemenhub) ensures that the implementation of the protocol on the Jabodetabek Electric Circuit Train (KRL) continues</p> <p>IndoBERT with additional contexts (+ keyphrases + titles): The Ministry of Transportation ensures that the implementation of the protocol on the Jabodetabek electric circuit train (KRL) continues.</p> <p>mBART: Permenhub 18/2020 has explicitly stated that there are several conditions that must be met by passengers of public transportation modes such as KRL.</p> <p>mBART with additional contexts (+ keyphrases + titles): Adita said that Permenhub 18/2020 has explicitly stated that there are several conditions that must be met by passengers of public transportation modes such as KRL.</p> <p>mT5: The Ministry of Transportation ensures that the implementation of the protocol on the Jabodetabek Electric Circuit Train (KRL) continues.</p> <p>mT5 with additional contexts (+ keyphrases + titles): Spokesperson for the Ministry of Transportation, Adita Irawati, stated that her party had issued Permenub No. 18/2020 which regulates the operation of transportation modes during the pandemic.</p>

Figure 2: Example from the LipKey dataset, with gold-standard and generated summaries.

Indonesian	English (translation)
<p>Gold Keyphrases: Relawan Uji Vaksin, Vaksin Sinovac</p> <p>Article: liputan6 . com , jakarta - manajer lapangan tim riset uji klinis vaksin covid-19 sinovac , dr eddy fadliyana menyebut sejauh ini sudah ada sekitar 1 . 020 calon relawan yang mendaftarkan diri untuk mengikuti uji vaksin dari tiongkok itu . dia mengatakan , pelaksanaan uji vaksin itu akan dilakukan selasa 11 agustus 2020 . pada hari pertama itu , uji vaksin bakal dilakukan di rumah sakit pendidikan (rsp) universitas padjadjaran , jalan eyckman , kota bandung . " sebetulnya sama saja , hanya pemeriksaan di rsp itu , tes usapnya (swab test) didahulukan . sama saja sih prosedurnya , tidak ada yang berbeda , besok rsp imunisasi , kalau di tempat lain baru tahap awal , " kata eddy di bandung , senin (10/8/2020) . dikutip dari antara , menurut eddy , semua tempat yang ditunjuk menjadi lokasi uji vaksin covid-19 ini dipastikan sudah siap . mulai dari sarana prasarananya , menurutnya sudah sesuai dengan protokol kesehatan yang berlaku . dia mengatakan , uji vaksin itu dilakukan di enam lokasi , di antaranya yakni rsp unpad , balai kesehatan unpad dipatiukur , puskesmas dago , puskesmas sukapakir , puskesmas garuda , dan puskesmas ciumbuleuit . dari seluruh calon relawan yang sudah mendaftar , menurutnya tak menutup kemungkinan sudah ada asn yang ikut mendaftar . karena , pendaftaran untuk menjadi relawan itu terbuka untuk umum . " dari asn mungkin ada , saya tidak melihat statusnya apa pokoknya masyarakat yang mau silakan saja , " katanya . meski terbuka untuk umum , menurutnya ada beberapa syarat yang perlu dipenuhi oleh calon relawan antara lain usia relawan dalam rentang 18 hingga 59 tahun , dan dalam keadaan sehat tanpa penyakit bawaan .</p>	<p>Gold Keyphrases: Vaccine Test Volunteers, Sinovac Vaccines</p> <p>Article: liputan6 . com , Jakarta - field manager of the Sinovac Covid-19 vaccine clinical trial research team, Dr. Eddy Fadliyana, said that so far there have been around 1.020 prospective volunteers who registered to take part in the vaccine test from China. He said the implementation of the vaccine test would be carried out on Tuesday, August 11, 2020. On that first day, the vaccine test will be conducted at the Teaching Hospital (RSP) at Padjadjaran University, Jalan Eyckman, Bandung City. "it's actually the same, only the examination at the rsp, the swab test takes precedence. the procedure is the same, nothing is different, tomorrow the immunization rsp, if it's in place others are only in the early stages," said Eddy in Bandung, Monday (10/8/2020). Quoted from Antara, according to Eddy, all the places designated to be the test locations for the COVID-19 vaccine are confirmed to be ready. starting from the infrastructure, according to him, it is in accordance with the applicable health protocol. he said the vaccine test was carried out in six locations, including the Unpad Hospital, Dipatiukur Health Center, Dago Health Center, Sukapakir Health Center, Garuda Health Center, and Ciumbuleuit Health Center. From all prospective volunteers who have registered, according to him, it is possible that there are already ASN who have registered. because , registration to become a volunteer is open to the public . " From the ASN there may be , I do not see what the status is , basically people who want to go ahead , " he said . although it is open to the public , according to him , there are several requirements that need to be fulfilled by prospective volunteers , including the age of volunteers in the range of 18 to 59 years , and in good health without any congenital disease .</p>
<p>Gold Keyphrases: buaya terkam warga, Sulbar</p> <p>Article: liputan6 . com , mamuju tengah - kejadian nahas menimpa h (40) warga desa barakkang , kecamatan budong-budong , mamuju tengah , sulawesi barat . ibu rumah tangga itu diterkam seekor buaya saat mandi dan buang air besar di sungai . kapolsek budong-budong akp suparman membenarkan peristiwa nahas itu , ia mengatakan , peristiwa terjadi pada selasa (4/8/2020) dini hari , sekitar pukul 05 . 30 wita . korban yang tengah buang air besar itu tiba-tiba diterkam buaya yang memiliki panjang kurang lebih 7 meter . " menurut saksi andi (38) yang merupakan adik korban , buaya itu tiba-tiba menerkam korban dari belakang , " kata suparman kepada liputan6 . com , petani labuhan batu utara diterkam buaya di depan anak istri suparman menambahkan , saksi juga sempat mendengarkan teriakan korban dan berusaha untuk menolong . namun , belum sempat menolong , buaya tersebut sudah terlebih dahulu menarik korban ke dalam air . " beberapa saat kemudian korban dan buaya muncul di permukaan air namun hanya sesaat lalu kemudian tenggelam lagi ke dalam air , " jelas suparman . hingga saat ini korban belum juga ditemukan , warga bersama pihak kepolisian sempat melakukan pencarian dengan peralatan seadanya . pihak bpbd mamuju tengah dan basarnas mamuju pun sudah dihubungi . " saat ini bpbd dan masyarakat serta basarnas sudah ada di tkp melakukan pencarian , " tutup suparman .</p>	<p>Gold Keyphrases: Crocodile devours residents, Sulbar</p> <p>Article: liputan6 . com , Mamuju - an unfortunate incident happened to H (40) a resident of Barakkang Village, Budong-Budong District, Central Mamuju, West Sulawesi. The housewife was attacked by a crocodile while bathing and defecating in the river. The head of the Budong-Budong Police, AK Suparman, confirmed the unfortunate incident, saying that the incident occurred on Tuesday (4/8/2020) early in the morning, around 05 am. 30 pm. The victim who was defecating was suddenly attacked by a crocodile which has a length of approximately 7 meters. " According to witness Andi (38) who is the victim 's younger brother , the crocodile suddenly pounced on the victim from behind , " said Suparman to liputan6 . com . The farmer in North Batu Harbor was attacked by a crocodile in front of his wife and children. Suparman added that the witness had also heard the victim's screams and tried to help. however , before they could help , the crocodile had already pulled the victim into the water . " a few moments later the victim and the crocodile appeared on the surface of the water , but only a moment later then sank again into the water , " explained Suparman . Until now the victim has not been found , residents together with the police had conducted a search with makeshift equipment . The Central Mamuju BPBD and Mamuju Basarnas have also been contacted. " Currently , BPBD and the community as well as the National Basis are already at the scene conducting a search , " concluded Suparman .</p>

Figure 3: Example of articles and keyphrases in the LipKey dataset. We highlight words in the article that match its abstractive keyphrases with different colours. Yellow means partial match, green means acronym, and blue means morphology variants. English translation is for illustration purposes.