

# Semi-Supervised Low-Resource Style Transfer of Indonesian Informal to Formal Language with Iterative Forward-Translation

Haryo Akbarianto Wibowo  
Kata AI Research Team  
Kata.ai  
Jakarta, Indonesia  
haryo@kata.ai

Tatag Aziz Prawiro  
Faculty of Computer Science  
Universitas Indonesia  
Depok, Indonesia  
tatag.aziz@ui.ac.id

Muhammad Ihsan  
School of Computer Science  
Bina Nusantara University  
Jakarta, Indonesia  
muhammad.ihsan004@binus.ac.id

Alham Fikri Aji  
Kata AI Research Team  
Kata.ai  
Jakarta, Indonesia  
aji@kata.ai

Radityo Eko Prasoj  
Kata AI Research Team  
Kata.ai  
Jakarta, Indonesia  
ridho@kata.ai

Rahmad Mahendra  
Faculty of Computer Science  
Universitas Indonesia  
Depok, Indonesia  
rahmad.mahendra@cs.ui.ac.id

Suci Fitriany  
Kata AI Linguist Team  
Kata.ai  
Jakarta, Indonesia  
suci@kata.ai

**Abstract**—In its daily use, the Indonesian language is riddled with informality, that is, deviations from the standard in terms of vocabulary, spelling, and word order. On the other hand, current available Indonesian NLP models are typically developed with the standard Indonesian in mind. In this work, we address a style-transfer from informal to formal Indonesian as a low-resource machine translation problem. We build a new dataset of parallel sentences of informal Indonesian and its formal counterpart. We benchmark several strategies to perform style transfer from informal to formal Indonesian. We also explore augmenting the training set with artificial forward-translated data. Since we are dealing with an extremely low-resource setting, we find that a phrase-based machine translation approach outperforms the Transformer-based approach. Alternatively, a pre-trained GPT-2 finetuned to this task performed equally well but costs more computational resource. Our findings show a promising step towards leveraging machine translation models for style transfer. Our code and data are available in <https://github.com/haryoa/stif-indonesia>.

**Index Terms**—style-transfer, Indonesian, machine translation, colloquial, semi-supervised, natural language processing

## I. INTRODUCTION

Irregular or informal text has often been a problem for the conversational or social media domain in Indonesia. Most people write without minding the choice of words or structure, so typos and slangs become common. Furthermore, modern conversational Indonesians are heavily influenced by loan-words, both foreign or traditional. It is very common to see Indonesian tweets mixed with English, Javanese, Sundanese, (romanized) Arabic, or (romanized) Korean. On the other hand, currently available Indonesian NLP models mostly handle formal Indonesian, and therefore perform comparatively worse when exposed to informal texts. Examples of such a phenomenon is found in machine translation [4]

One of the main challenges of expanding the model capability to handle informal Indonesian is the availability of labelled informal training data for each particular NLP task, or lack thereof. Therefore, an independent style-transfer model to be utilized as a preprocessing step can potentially improve the downstream models without any reconfiguration or retraining necessary. Previous work has resulted in a treebank [11] or a dictionary of informal words in Indonesian [16], which by nature is limited to word-level standardization without carrying sentence-level context, and therefore further exploration of standardization in sequence level is still an open problem.

In this work, we focus on investigating sequence-level style transfer from informal Indonesian to its formal counterpart. However, since no dataset for this purpose currently exists. Therefore, we build a new dataset of parallel sentences in informal and formal Indonesian. Besides, we also explore adding an artificial dataset to leverage our training resource. Our main contributions are as follow.

- 1) We build a new dataset of parallel sentences in informal and formal Indonesian.
- 2) We benchmark several style-transfer strategies, starting from a dictionary-based approach as a baseline and some statistical and neural machine translation approaches.
- 3) We explore the use of synthetic dataset for informal to formal Indonesian style-transfer.

The remainder of the paper is structured as follow. We discuss our informal Indonesian text data in Section 2, including how we collect them, and their statistics. Section 3 describes the style-transfer approaches that we explore. Section 4 describes our experiment result and analysis.

Domain	Twitter IDs
Telecommunication	@telkomsel, @indosatcare, @smartfrenicare, @myXLCare, @3CareIndonesia
Banking	@kontakBRI, @mandiricare, @HaloBCA, @BNICustomerCare, @CIMBNiaga
E-Wallet	@ovo_id, @danawallet, @linkaja, @jeniushelp
E-Commerce	@TokopediaCare, @ShopeeCare, @BukaBantuan, @LazadaIDCare, @BlibliCare, @csjd_id, @ZaloraID
Logistics	@JNECare, @IdTiki, @jntexpressid, @PosIndonesia
Ride-Hailing	@gojekindonesia, @GrabID

TABLE I: The list of scraped tweet IDs

## II. INFORMAL-FORMAL INDONESIAN PARALLEL DATA

### A. Data Collection

We access Twitter API using tweepy<sup>1</sup>, focusing on the customer service domain. We identify Twitter’s customer service account in Indonesia across different business areas. The full list of scraped Twitter ID is shown in Table I.

### B. Data Filtering

We remove the hashtags and deduplicate the tweets. We filter out any instances which contain less than 5 or more than 25 tokens. We include code-mixed tweets since the tweets containing mixed English and Indonesian words are commonly used in the customer service domain. As we want to ensure that the data does not contain predominantly foreign words, we ignore any tweets which have 60% or more English words. We collect 52.5k informal Indonesian tweets. We then sample 2500 tweets to be annotated into formal Indonesian. Both the 50k raw informal tweets and 2.5k annotated informal-formal parallel tweets are publicly accessible<sup>2</sup>.

### C. Informal to Formal Data Annotation

We annotate the informal data into a high form of *Bahasa Indonesia* [12]. The high form *Bahasa Indonesia* is the standard literary Indonesian omitting the regional varieties and the colloquial words, but including some loan words, online words, and some forms of common conversational Indonesian words. A common theme of these words is English words occurring within a computer or internet-related context (like “mouse”, “keyboard”, “voucher”, “tweet”, etc.), which became common with the rise of the Internet before they are standardized into formal Indonesian. For simplicity, We refer to the high form *Bahasa Indonesia* data as formal data and the opposite as informal data.

We annotate the data by rewriting the given informal text to its formal form. Here are several things that we consider when we annotate the data:

- **Punctuation** (e.g.: ‘Saya bisa admin’ → ‘Saya bisa, admin’) (‘I can, admin’)
- **Capitalization** (e.g.: ‘nama saya haryo’ → ‘Nama saya Haryo’) (‘my name is Haryo’)
- **Word order** (e.g.: ‘Admin, bisa seperti itu kenapa?’ → ‘Admin, kenapa bisa seperti itu?’) (‘Hey Admin, how is that possible?’)

<sup>1</sup><https://www.tweepy.org/>

<sup>2</sup><https://github.com/haryoa/stif-indonesia>

- **Colloquial/Shorten Word** (e.g.: ‘gw knp tdk bs’ → ‘Saya kenapa tidak bisa?’) (‘Why Can’t I do this?’)
- **Affixes or Suffixes** (e.g.: ‘saya mesenin taxi’ → ‘Saya memesan taxi’) (‘I ordered a taxi’)

To ensure the quality of our annotation, we maintain a strict annotation guideline and perform cross-review between annotators. For the experiment purposes, we split the data into train, development, and test. The data is divided into 1922, 214, and 364 for the train, dev, and test data respectively.

### D. Data Preprocessing

We use different pre-processing on each informal and formal data. For informal data, we lowercase the data since the text is irregular, which may hinder the learning process of the model. Then we apply text tranformation when finding more than 2 consecutive characters (e.g.: ‘makannn’ becomes ‘makann’). Finally, we mask the number, account, date, percentage token

### E. Data Analysis

Stat	Informal data	Formal data
Number of unique tokens	5381	4036
Number of tokens	37040	39467

TABLE II: Descriptive statistic of informal and formal data

In Table II, we show descriptive statistics of the informal and formal data. Where the statistics also include token punctuation and mask token.

The number of unique tokens in informal data is more than formal data. This shows that informal tokens are more varied than formal data. This is because a formal word can be represented with multiple informal tokens. For example, an informal token of ‘tidak’ (‘no’) could be ‘gak’, ‘ga’, ‘nggak’, or ‘engga’.

Punctuation	Informal data	Formal data
.	1386	3488
,	1289	1821
?	1021	1346

TABLE III: Punctuation number of informal and formal data

We found an interesting finding that the number of tokens on the informal data is lower than the formal data, which is the opposite of the unique token number statistics. It shows that there is an increase in the number of punctuation tokens from informal to formal. Table III shows the changes of the punctuation ‘.’, ‘,’, And ‘?’. This result makes sense, as people usually pay less attention to punctuation when writing informal text. For example, “akuu dari awal tahunn gila ga

sih” (“I was crazy from the beginning of the year”) which becomes ‘aku sejak awal tahun, gila bukan?’ in the formal form. Furthermore, there are informal sentences that do not use ‘.’ at the end of the sentence. We observe that there are 2148 instances in our data where there are no ‘.’ at the end of the sentence.

Token in Informal Data	Freq.	Token in Formal Data	Freq.
saya (‘I’)	763	saya (‘I’)	1058
min (informal of ‘admin’)	585	tidak (‘no’)	913
ya (‘yes’)	554	admin	635
di (‘in’)	510	sudah (‘done’)	545
ini (‘this’)	425	bisa (‘can’)	510
bisa (‘can’)	421	ini (‘this’)	482
ada (‘available’)	324	di (‘in’)	410
ga (informal of ‘tidak’; ‘no’)	272	yang (‘that’)	398
mau (‘want’)	268	ada (‘available’)	354
dm (direct message)	253	mengapa (‘why’)	283
yg (informal of ‘yang’; ‘that’)	232	terima (‘thank’)	277
kok (‘why’ informal)	229	kasih (‘you’)	277
tolong (‘help’)	226	tolong (‘help’)	265
ke (‘to’)	215	dari (‘from’)	253
sudah (‘done’)	213	mau (‘want’)	240

TABLE IV: Top 15 tokens found in informal and formal data. Punctuation and masked token are excluded

Table IV shows the top 15 formal and informal data, which both rank differently according to their word frequency. Tokens ‘saya’, ‘di’, ‘ini’, ‘sudah’, ‘bisa’, ‘ada’, ‘mau’, and ‘tolong’ appear frequently in both informal and formal data. While, several pairs of words in different style are found in the Table, i.e. ‘min’ - ‘admin’, ‘ga’ - ‘tidak’, ‘kok’ - ‘mengapa’, and ‘yang’ - ‘yg’.

### III. METHOD

We address our informal to formal style transfer as a sequence to sequence problem. In this section, we describe our approaches.

#### A. Dictionary-based Translation

One naive idea is to simply transform the sequence on a word-level basis with an aid of a dictionary. As a baseline, we developed a simple dictionary-based translator system. We make use of the word-level Indonesian formal-informal dictionary<sup>3</sup> that has been used in multiple works. This system simply translates an informal word into its formal form if it appears in the dictionary.

#### B. Phrase-Based Statistical Machine Translation

Phrase-Based machine translation (PBSMT) has shown to perform well when the dataset is scarce [9], as in our case. Hence, we explore using PBSMT as one of the baselines. We employ Moses [8] to develop our PBSMT system. We use MGIZA [3] on aligning the phrases.

<sup>3</sup><https://github.com/ialfina/ID-Kamus-Typo>

#### C. Neural Machine Translation

The Transformer architecture [20] has been state-of-the-art for neural machine translation. Therefore, we explore this architecture as one of our baselines. We employ the standard transformer architecture consisting of 6 layers encoder and decoder. Our model is trained with Marian toolkit [7].

#### D. Pretrained Language Modeling with GPT-2

In pre-trained language models, initially, a model (usually Transformer-based) is trained to learn general language modelling. Then, the model will be fine-tuned to the downstream task. Pretrained LM has been shown to adapt to the downstream task without requiring a large dataset, which is suitable for our case. We chose GPT-2 in particular due to its flexibility in fine-tuning downstream generative tasks [14].

We first train a GPT-2 based Indonesian language model, which is designed for task generation. Our GPT-2 model is trained on the OSCAR corpus [19]. Then, we fine-tune our GPT-2 for style-transfer by modifying our parallel corpus into a Informal\_sentence <STIF> formal\_sentence format and train with the modified dataset. (For example, cabs kuy <STIF> ayo pergi). To translate informal sentences, we simply feed the GPT-2 with the informal part and <STIF> tag and let the model complete the sentence.

#### E. Synthetic Data Generation

Commonly, back-translation is used as a synthetic dataset in addition to the parallel corpus. In that case, we need to gather formal Indonesian corpus in customer services domain to be translated to informal text. Unfortunately, formal customer services domain text is uncommon. Therefore, we opt to use a forward-translated synthetic dataset instead. forward-translated data has shown to be beneficial towards to model’s performance [2].

We generate our forward translation model by translating a sample of 5000 tweets with our best-performing model according to fully-supervised learning. We then add this synthetic dataset on top of our original parallel corpus. Finally, we re-train our model with this new dataset composition.

In this experiment, we create an artificial corpus by forward-translating informal text. Hoang et al. [5] has shown that the quality of the synthetic dataset matters. They suggest that iteratively re-construct the synthetic back-translation can improve translation quality. We adopt their finding and attempt an iterative forward-translation: We train our model across different iterations. At the  $i$ -th iteration, our model is trained with the parallel corpus and additional synthetic forward-translated corpus generated by the model from the  $i - 1$ -th iteration.

## IV. EXPERIMENT AND RESULT

### A. Model Benchmark

We first answer the question of which approach is the best for low-resource informal-formal Indonesian style transfer. For this purpose, we trained our models with our parallel corpus and evaluate their performance with sacreBLEU [13].

Method	Examples
Input	ku coba resto lain juga sama. jd gmn sih sistemnya?
Ground-truth	Saya mencoba <b>restoran</b> lain juga sama. Jadi bagaimana ini sistemnya? (I tried other restaurants as well. So how is this the system?)
Dict-based	ku coba resto lain juga sama. jadi bagaimana sih sistemnya?
PBSMT	kucoba resto lain juga sama. jadi bagaimana sih sistemnya?
Transformer	Aku coba <b>telpon</b> lain juga sama. Jadi bagaimana? Apa sistemnya? (I tried other phones, it is the same. So how? What is the system?)
GPT-2	aku coba resto lain juga sama. jadi bagaimana sistemnya?
Input	kenapa <b>pas</b> mau login <b>kaya</b> gini terus ya?
Ground-truth	mengapa ketika mau login seperti begini terus? (why it is always like this whenever I want to log in?)
Dict-based	kenapa <b>pas</b> mau login <b>kaya</b> gini terus ya?
PBSMT	kenapa saat mau login seperti ini terus?
Transformer	kenapa saat mau login seperti ini terus?
GPT-2	kenapa saat mau login seperti ini terus?

TABLE V: some examples of formal text generated from our system.

Method	BLEU
No Modification	35.32
Dictionary-Based	42.11
PBSMT	<b>49.39</b>
Transformer	27.50
GPT-2 Pretrain	<b>49.28</b>

TABLE VI: Style-transfer benchmark with different approaches.

As shown in Table VI, we present our result of dictionary-based, PBSMT, Transformer, and GPT-2 based approaches. As an additional baseline, we also include *No Modification*, where the input text is unchanged at all.

Generally, informal Indonesian consists of using colloquial terms of certain words. Therefore, our dictionary-based model performed well with 42.11 BLEU score. One flaw of the dictionary-based model is that it cannot translate words that already exist in the formal dictionary, but have different informal meaning. For example, in the Table V, *kaya* (rich, wealthy) is a formal Indonesian word. However, it is also used informally with a different meaning (*kaya* in colloquial Indonesian means ‘similar’, or ‘look like’). The word *pas* also has the same issue (it means ‘precise’ in formal Indonesian, but means ‘whenever’ informally).

Besides, informal Indonesian is also more flexible in the sentence structure. The drawback of this approach is that the model does not consider any word removal or addition. Furthermore, it does not consider words re-alignment. Therefore, it does not perform well when word removal, addition or swap is required. For example, the dictionary-based model cannot remove the suffix ‘ya’ in the Table V.

The result in Table VI shows that the Transformer performed the worst. Consistent with prior work [9], we find that the Transformer is incapable of performing well under the extreme low-resource setting. The Transformer performance is worse compared to not modifying the informal text at all. From manual evaluation, we see that the Transformer model often generates output with a different meaning from the source. For example, in the Table V, the Transformer output changes the sentence meaning significantly, from complaining about

BLEU vs Iteration

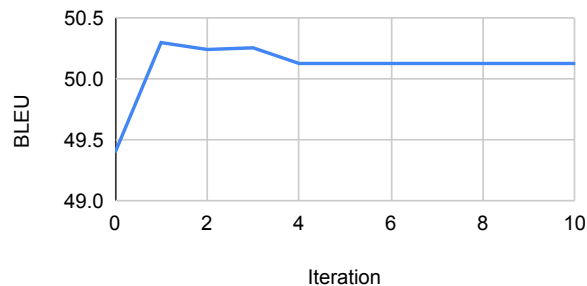


Fig. 1: The performance of style-transfer with iterative forward-translation performance in terms of BLEU.

the restaurant to phones. It also alters the sentence structure significantly. We usually find this problem on longer and more complicated sentences.

Our PBSMT approach achieves the best performance of 49.39 BLEU. This result confirms that the PBSMT approach performed better under the extreme low-resource setting. Interestingly, our GPT-2 approach can fine-tune well even with only 2.5k sentence pairs and achieved comparably near-best performance of 49.28 BLEU.

The BLEU scores of GPT-2 and PBSMT approaches are comparable, but the PBSMT is more efficient in terms of computational resource. Therefore, our follow-up experiment will utilize the PBSMT model.

### B. Semi-supervised Style-transfer with Forward-Translation

We explore adding 2500 forward-translated synthetic sentences on top of the original parallel dataset. We train the model up to 10 iterations. For each iteration, the model is trained normally until convergence. After that, we use the current model to generate the synthetic dataset for the next iteration and then re-start the training again. Model at the iteration 0 is trained only with the original parallel corpus. Our experiment result is shown in the Figure 1. We use the PBSMT approach as discussed previously.

As shown in the Figure 1, adding the synthetic dataset improves the performance to 50.3 BLEU. However, adding more iteration does not seem to improve the performance. This result suggests that the synthetic dataset generated by models at different iterations to be varied enough to produce significantly different performances. The slight decrease in BLEU in later iterations also suggests that the synthetic data introduces noise into the overall training data. We leave the verification of these findings and investigation into their improvements as future work.

## V. RELATED WORK

Rao and Tetreault [15] formulated the problem of formality of text as a style transfer. They released the Grammarly’s Yahoo Answers Formality Corpus (GYAFC) dataset, which is parallel data of informal and formal data in the Entertainment & Music and Family Relationship domain in English. They also tried several sequence-to-sequence approaches to benchmark the data they released. Jhamtani et al. [6] also applied sequence-to-sequence approach on style transfer problem. They conducted a style transfer approach by using a modern paraphrase of Shakespeare’s play released by Xu et al. [21]. We are inspired by their work to apply it into Indonesian informal and formal style. Unfortunately, there are no parallel informal and formal Indonesia data that is ready to be used.

Barik et al. [1] aimed to tackle one of the informal Indonesia text problems, namely code-mixed, in which a sentence contains the words in more than one language. They proposed a pipeline solution, i.e. word tokenization, language identification, normalization, and translation. They identified the language of each token in the sentence, then translated the words into Indonesian.

Shang et al. [17] leveraged a semi-supervised approach of style changes by projecting latent space. Yang et al. [22] also utilized semi-supervised learning with a generator and discriminator approach to study representations of the latent space of each style.

On the other hand, several works proposed the unsupervised approach. Luo et al. [10] performed unsupervised style transfer with Dual Reinforcement Learning, where the model was trained adversarially. It used the results of the style classifier as the reward system. Shen et al. [18] utilized the Variational Auto Encoder (VAE) approach to form latent representation and studied mapping functions to map the latent representation of each style.

## VI. CONCLUSION

In this paper, we have explored Indonesian informal to formal style-transfer as a low-resource machine translation problem. We build a new parallel informal - formal Indonesian data by annotating tweets from the customer service domain. We conclude that the pretrained GPT-2 and PBSMT approaches achieve the best performance in terms of BLEU. Training with an additional synthetic dataset in the form of forward-translated informal Indonesian text improves the performance.

## ACKNOWLEDGEMENT

This research was supported by the research grant from Universitas Indonesia, namely Publikasi Terindeks Internasional (PUTI) Saintekkes year 2020 no NKB-2142/UN2.RST/HKP.05.00/2020

## REFERENCES

- [1] Anab Maulana Barik, Rahmad Mahendra, and Mirna Adriani. Normalization of Indonesian-English code-mixed Twitter data. In *Proceedings of the 5th Workshop on Noisy User-generated Text (W-NUT 2019)*, pages 417–424, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-5554.
- [2] Nikolay Bogoychev and Rico Sennrich. Domain, translationese and noise in synthetic data for neural machine translation. *arXiv preprint arXiv:1911.03362*, 2019.
- [3] Qin Gao and Stephan Vogel. Parallel implementations of word alignment tool. In *Software engineering, testing, and quality assurance for natural language processing*, pages 49–57, 2008.
- [4] Tri Wahyu Guntara, Alham Fikri Aji, and Radityo Eko Prasajo. Benchmarking multidomain english-indonesian machine translation. In *Proceedings of the 13th Workshop on Building and Using Comparable Corpora*, pages 35–43, 2020.
- [5] Vu Cong Duy Hoang, Philipp Koehn, Gholamreza Haffari, and Trevor Cohn. Iterative back-translation for neural machine translation. In *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation*, pages 18–24, 2018.
- [6] Harsh Jhamtani, Varun Gangal, Eduard Hovy, and Eric Nyberg. Shakespearizing modern language using copy-enriched sequence to sequence models. In *Proceedings of the Workshop on Stylistic Variation*, pages 10–19, Copenhagen, Denmark, September 2017. Association for Computational Linguistics. doi: 10.18653/v1/W17-4902.
- [7] Marcin Junczys-Dowmunt, Roman Grundkiewicz, Tomasz Dwojak, Hieu Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann, Alham Fikri Aji, Nikolay Bogoychev, André F. T. Martins, and Alexandra Birch. Marian: Fast neural machine translation in C++. In *Proceedings of ACL 2018, System Demonstrations*, pages 116–121, Melbourne, Australia, July 2018. Association for Computational Linguistics.
- [8] Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, et al. Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th annual meeting of the ACL on interactive poster and demonstration sessions*, pages 177–180. Association for Computational Linguistics, 2007.
- [9] Guillaume Lample, Myle Ott, Alexis Conneau, Ludovic Denoyer, and Marc’Aurelio Ranzato. Phrase-based & neural unsupervised machine translation. In *Proceedings*

- of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 5039–5049, 2018.
- [10] Fuli Luo, Peng Li, Jie Zhou, Pengcheng Yang, Baobao Chang, Zhifang Sui, and Xu Sun. A dual reinforcement learning framework for unsupervised text style transfer. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence, IJCAI 2019*, 2019.
- [11] David Moeljadi, Aditya Kurniawan, and Debaditya Goswami. Building cendana: a treebank for informal indonesian. 2019.
- [12] Scott Paauw. *The Malay contact varieties of Eastern Indonesia: A typological comparison*. State University of New York at Buffalo, 2009.
- [13] Matt Post. A call for clarity in reporting bleu scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, 2018.
- [14] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners.
- [15] Sudha Rao and Joel Tetreault. Dear sir or madam, may I introduce the GYAFC dataset: Corpus, benchmarks and metrics for formality style transfer. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 129–140, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-1012.
- [16] Nikmatun Aliyah Salsabila, Yosef Ardhito Winatmoko, Ali Akbar Septiandri, and Ade Jamal. Colloquial indonesian lexicon. In *2018 International Conference on Asian Language Processing (IALP)*, pages 226–229, 2018. doi: 10.1109/IALP.2018.8629151.
- [17] Mingyue Shang, Piji Li, Zhenxin Fu, Lidong Bing, Dongyan Zhao, Shuming Shi, and Rui Yan. Semi-supervised text style transfer: Cross projection in latent space. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4937–4946, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1499.
- [18] Tianxiao Shen, Tao Lei, Regina Barzilay, and Tommi Jaakkola. Style transfer from non-parallel text by cross-alignment. In *Advances in neural information processing systems*, pages 6830–6841, 2017.
- [19] Pedro Javier Ortiz Suárez, Benoît Sagot, and Laurent Romary. Asynchronous pipeline for processing huge corpora on medium to low resource infrastructures. In *7th Workshop on the Challenges in the Management of Large Corpora (CMLC-7)*. Leibniz-Institut für Deutsche Sprache, 2019.
- [20] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
- [21] Wei Xu, Alan Ritter, Bill Dolan, Ralph Grishman, and Colin Cherry. Paraphrasing for style. In *Proceedings of COLING 2012*, pages 2899–2914, Mumbai, India, December 2012. The COLING 2012 Organizing Committee.
- [22] Zichao Yang, Zhiting Hu, Chris Dyer, Eric P Xing, and Taylor Berg-Kirkpatrick. Unsupervised text style transfer using language models as discriminators. In *Advances in Neural Information Processing Systems*, pages 7287–7298, 2018.