# NusaCrowd: A Call for Open and Reproducible NLP Research in Indonesian Languages

**Samuel Cahyawijaya**\*, **Alham Fikri Aji**\*, **Holy Lovenia**\*,
**Genta Indra Winata**\*, **Bryan Wilie**\*, **Rahmad Mahendra**, **Fajri Koto**,
**David Moeljadi**, **Karissa Vincentio**, **Ade Romadhony**, **Ayu Purwarianti**
IndoNLP

## Abstract

At the center of the underlying issues that halt Indonesian natural language processing (NLP) research advancement, we find data scarcity. Resources in Indonesian languages, especially the local ones, are extremely scarce and underrepresented. Many Indonesian researchers refrain from publishing and/or releasing their dataset. Furthermore, the few public datasets that we have are scattered across different platforms, thus makes performing reproducible and data-centric research in Indonesian NLP even more arduous. Rising to this challenge, we initiate the first Indonesian NLP crowdsourcing effort, NusaCrowd. NusaCrowd strives to provide the largest datasheet aggregation with standardized data loading for NLP tasks in all Indonesian languages. By enabling open and centralized access to Indonesian NLP resources, we hope NusaCrowd can tackle the data scarcity problem hindering NLP progress in Indonesia and bring NLP practitioners to move towards collaboration.

## 1   What is 𝒩 NusaCrowd?

Natural language processing (NLP) resources in Indonesian languages, especially the local language ones, are extremely scarce and underrepresented in the research community. This introduces bottlenecks to Indonesian NLP research, restraining it from opportunities and hindering its progress. In response to this issue, several Indonesian and NLP communities have sourced various types of datasets to also be available in Indonesian languages (Winata et al., 2022; Koto et al., 2021a; Cahyawijaya et al., 2021; Koto et al., 2021b; Mahendra et al., 2021; Aji et al., 2021; Wilie et al., 2020; Abdul-Mageed et al., 2020; Ilmania et al., 2018; Pan et al., 2017). However, a significant mass of the local resources is scattered across different platforms (Etsa et al., 2018; Apriani et al.,
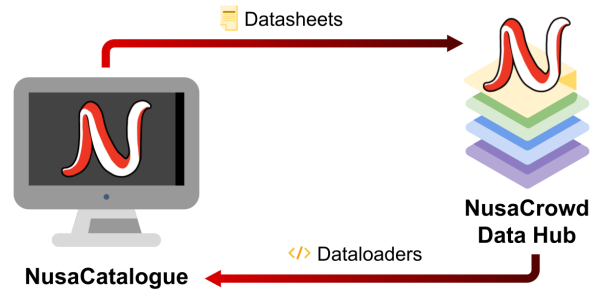


Figure 1: Open access to the datasheets collected is provided through **NusaCatalogue**, and the dataloader scripts to retrieve the resources are implemented in **NusaCrowd Data Hub**.

2016; Dewi et al., 2020), and a lack of access to public datasets still persists (Aji et al., 2022).

To address this vital problem, inspired by other open collaboration projects (Orife et al., 2020; Nekoto et al., 2020; Dhole et al., 2021; Alyafeai et al., 2021; McMillan-Major et al., 2022; Srivastava et al., 2022; Fries et al., 2022; Gehrmann et al., 2022), we take a step and initiate NusaCrowd, a joint movement to collect and centralize NLP datasets in Indonesian and various Indonesia's local languages, and engage the linguistics community in collaboration.

Powered by the collective effort of our contributors, NusaCrowd aims to increase the accessibility of these datasets and promote reproducible research on Indonesian languages through three fundamental facets: 1) Curated public corpora datasheet sourcing[1], 2) Open-access centralized data hub[2], and 3) Promoting private-to-public data access. We maintain the quality of the contributions, both the consolidation efforts and the programmatic means, by enforcing a quality control with a mix of automatic and manual evaluation schemes. NusaCrowd is currently open for contribution, the movement

---

\* Equal contribution.

[1] https://indonlp.github.io/nusa-catalogue/
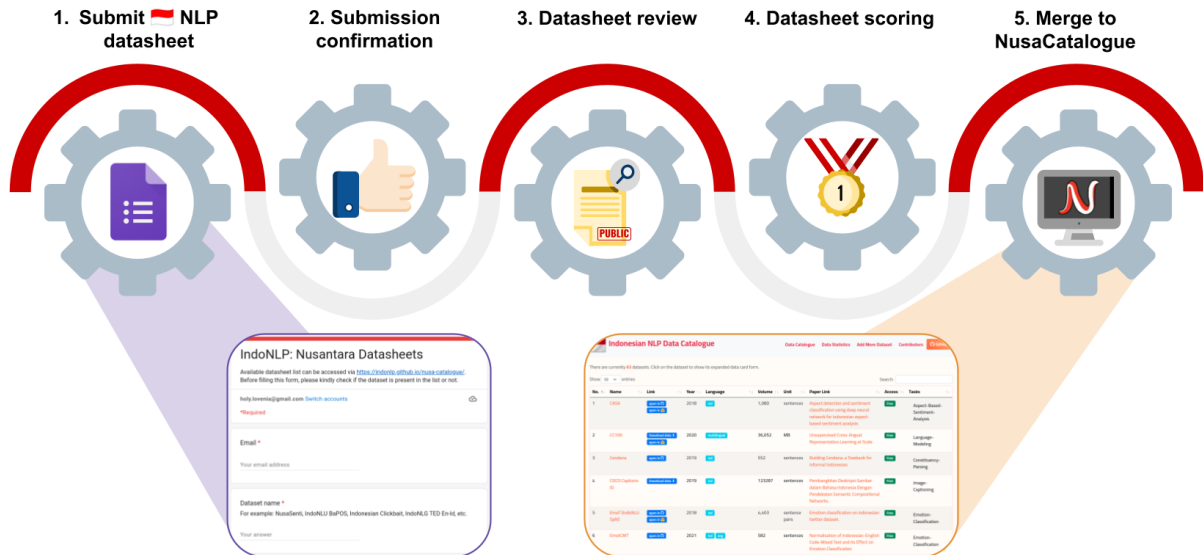[2] https://github.com/IndoNLP/nusa-crowd

Figure 2: The outline of public Indonesian NLP datasheet submission to NusaCatalogue.

is held from 25 June 2022 to 18 November 2022. Let's bring Indonesian NLP research one step forward together.

## 2 Contributing in 𝒩 NusaCrowd

Together, contributors in NusaCrowd drive Indonesian NLP forward by developing a multi-faceted solution to improve data accessibility and research reproducibility. To assist our widespread open collaboration and collective progression, we formulate three main ways of contributing in NusaCrowd in the following sections, each corresponds to a fundamental aspect in NusaCrowd's main objective.

### 2.1 Submit public 🇮🇩 NLP datasheet

We encourage contributors to register the datasheets of public datasets on **NusaCatalogue**[3] by submitting them through an online form at `https://forms.gle/31dMGZik25DPFYFd6`. NusaCatalogue is a public datasheet catalogue website, inspired by Alyafeai et al. (2021), in which we list all datasets collected in NusaCrowd. We build NusaCatalogue to improve the discoverability of Indonesian NLP datasets and to assist users in searching and locating Indonesian NLP datasets based on their metadata.

A datasheet is a dataset metadata which contains various information about the dataset, including but not limited to: dataset name, original resource URL, relevant publication, supported tasks, and dataset licence. The datasheet will be reviewed and scored within a week or two. The contribution point calculation for the public Indonesian NLP corpora datasheet is based on three criteria: 1) whether the relevant dataset is previously public or not, 2) dataset quality, and 3) language rarity. Details on the scoring mechanism will be explained in §3.1.

Once the datasheet passes the review, we will notify the responsible contributor and list the approved dataset's datasheet on NusaCatalogue and also reported on NusaCrowd data hub task list[4] so its dataloader could be implemented. The complete flow of how to submit the public Indonesian NLP corpora datasheet is shown in Figure 2.

### 2.2 Implement dataloader(s) for 𝒩 NusaCrowd data hub

A large-scale centralized data hub has to be equipped with the capability of a simple and standardized programmatic data access that spans across diverse resources, regardless of their separate hosting locations, distinct data structures or formats, and different configurations. For this purpose, building NusaCrowd data hub requires a few key elements: datasheet documentation (§2.1), task schema standardization to support common NLP tasks, and dataloader implementation. While the datasheets become the backbone of NusaCrowd and standardized task schemas compose the skeleton, the dataloader implementation is the heart of
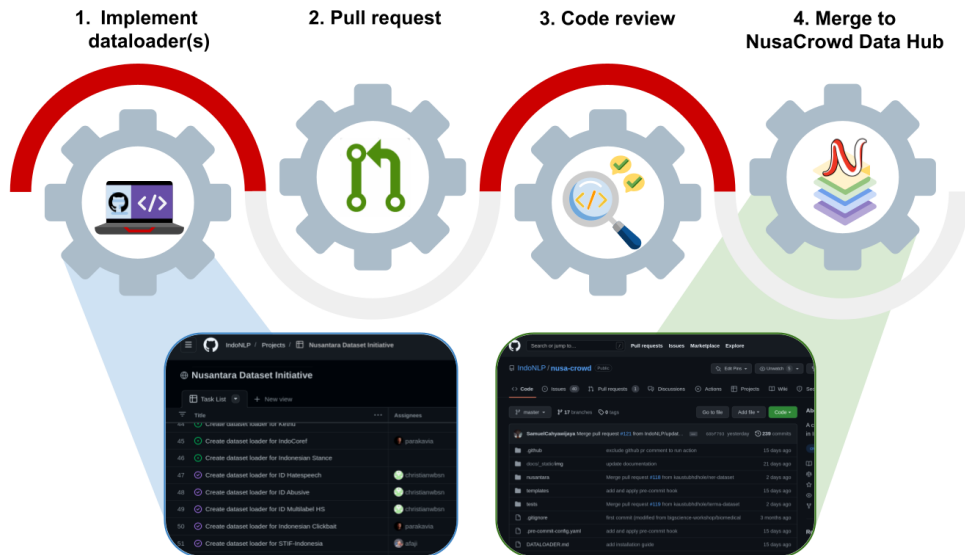
---

Figure 3: The outline of dataloader implementation for NusaCrowd data hub.

NusaCrowd data hub.

Each dataset requires a specific dataloader script tailored to its source task type, structure, and configuration to enable easy loading and enforce interoperability. To centralize all of the Indonesian NLP resources, NusaCrowd needs a large number of proper dataloaders to be implemented. Therefore, we invite all collaborators to contribute through creating these dataloaders via NusaCrowd's GitHub[5].

Firstly, a contributor can view the task list in NusaCrowd Github project, then choose the dataset they want to implement by assigning themself to the related issue. Afterwards, the contributor can start setting up the environment needed for development. To help with the dataloader implementation, we provide a template script specifying all the parts the contributor needs to complete, and task schemas for common NLP tasks, e.g., knowledge base, question answering, text classification, text-to-text, text pairs, question answering, and more. We also equip NusaCrowd's repository with several working dataloader scripts that the contributor can check for examples. To fill in the details of the dataset in the dataloader, such as its source URL or its publication, the contributor can refer to the corresponding datasheet recorded in NusaCatalogue.

The contributor can ensure that their dataloader is implemented correctly through a manual inspection, a direct attempt of execution, and a unit test provided in the repository. We also encourage the contributor to tidy up their code accordingly with our formatter before they make a pull request to

submit their changes. The implemented dataloader then will go through a code review process by NusaCrowd maintainers. If an adjustment in the code is required, the maintainer will request some changes and provide their feedback on a comment on the respective pull request, so the contributor will be able to improve the dataloader accordingly. Once two maintainers give their approvals, the dataloader will be merged to NusaCrowd repository. The flow overview of the dataloader implementation is depicted in Figure 3. A comprehensive guide for contributing through dataloader implementation can be accessed here[6].

## 2.3 Provide information on private 🇮🇩 NLP datasets

Many studies in Indonesian NLP use private datasets, which in turn hampers the research reproducibility. Therefore, the last method to contribute is to list research papers of Indonesian NLP in which the data is not shared publicly. We then contact the author to participate in open data access and ask for their approval to include the dataset in NusaCrowd. Contribution points for the authors that release their data to public will follow the scoring defined in §3.3. The paper lister will also be awarded with a contribution point.

As far as we know, there is no data or analysis yet on why local researchers prefer not to share their dataset publicly. Some reasons that we are aware of include: 1) not accustomed to open data
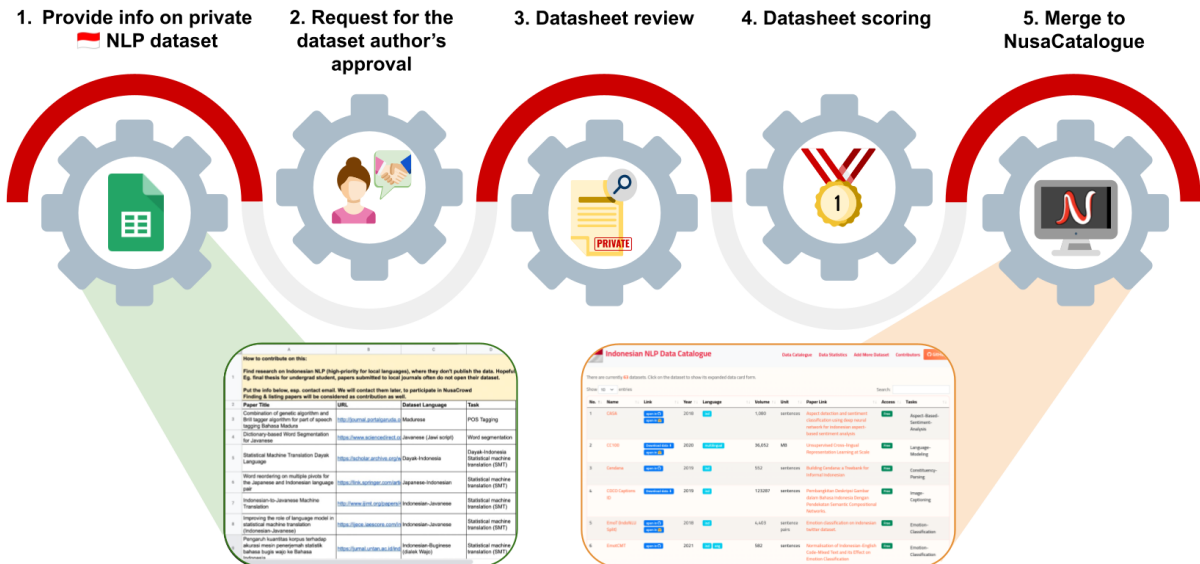
Figure 4: The outline of private Indonesian NLP dataset listing in NusaCatalogue.

and open research, 2) restricted by the university or funding policy, or 3) keep the data private as property. By listing research with private datasets, we can additionally ask the authors about their consideration to improve our understanding on this matter. Steps to provide information on private Indonesian NLP datasets are illustrated in Figure 4.

## 3 Contribution Point

To support fairness and transparency for all of our contributors, we establish a scoring system of which co-authorship eligibility will be decided from. To be eligible as a co-author in the upcoming NusaCrowd publication, a contributor needs to earn at least 10 contribution points. The score is aggregated from all the contributions made by the contributors. In order to earn contribution points, we introduce three different methods to contribute (for the method details, see §2): 1) submitting public Indonesian NLP corpora datasheet, 2) implementing dataloader(s) for NusaCrowd data hub, and 3) provide information on non-public Indonesian NLP datasets. The point for each type of contribution is described in the following paragraphs.

### 3.1 Public 🇮🇩 NLP corpora datasheet

A contributor can help to register public NLP corpora in NusaCrowd. For any datasheet listed, the contributor is eligible for +2 contribution point as a referrer. To support the development of local language datasets, we provide additional contribution points according to the rarity of the dataset

language. Specifically, a contributor of any Sundanese (sun), Javanese (jav), or Minangkabau (min) dataset, will receive +2 contribution points, while a contributor of any other local language dataset will be granted +3 contribution points.

In addition, to encourage more diverse NLP corpora, we provide additional +2 contribution points for tasks that are considered rare. Based on our observation, we find that the common NLP tasks in Indonesian languages include: machine translation (MT), language modeling (LM), sentiment analysis (SA), and named entity recognition (NER). All other NLP tasks are considered rare and are eligible for the +2 contribution points. Lastly, we also notice that publicly available Indonesian NLP corpora involving another modality (e.g., speech or image) are very scarce, for instance: speech-to-text or automatic speech recognition (ASR), text-to-speech (TTS) or speech synthesis, image-to-text (e.g., image captioning), text-to-image (e.g., controllable image generation), etc. To encourage more coverage over these data, we will give additional +2 contribution points for the relevant datasheets submitted.

We understand that dataset quality can vary a lot. To support fairness in scoring datasets with different qualities, for any dataset that does not pass a certain minimum standard, 50% penalty will be applied. This penalty affects any dataset that is collected with: 1) crawling without manual validation process, 2) machine or heuristic-rule labelled dataset without manual validation, and 3) machine-
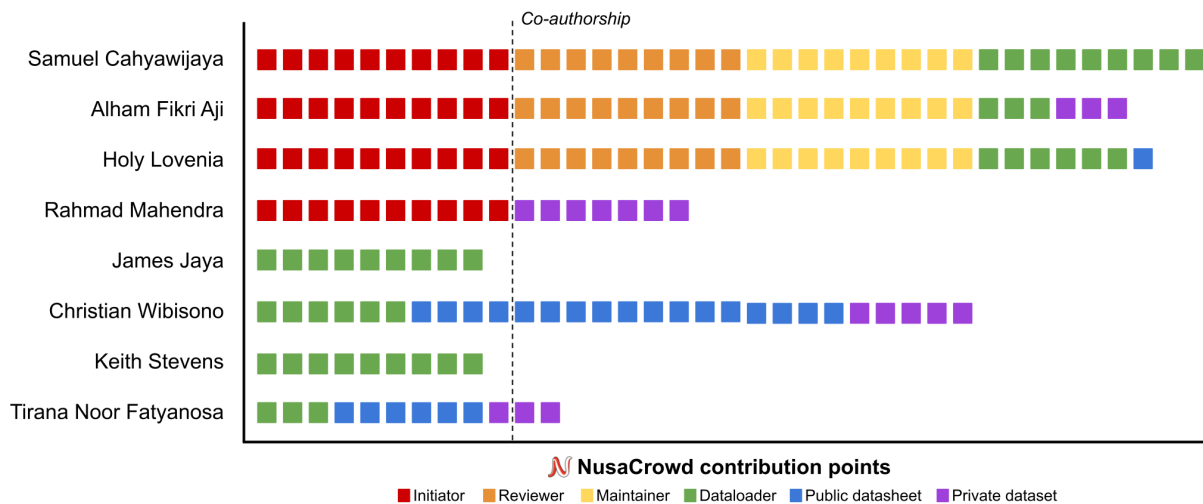
Figure 5: A glance at the contribution matrix used for recapitulation. The contributor list is clipped for simplicity.

translated dataset without manual validation.

## 3.2 Implementing 𝒩 dataloader

A contributor can help to implement a dataloader for any dataset listed on the task list in NusaCrowd GitHub project (see §3.2). As a rule of thumb, one dataloader implementation is generally worth 3 contributions points. However, there are some exceptions where a dataloader can be worth more. The contribution points will be counted once the respective pull request is merged to `master`.

## 3.3 Finding and opening private 🇮🇩 NLP datasets

Contributors can help to list research papers introducing non-public Indonesian NLP dataset. For every private dataset listed, the corresponding contributor will be eligible for +1 contribution point. While for the original author of the dataset, if the author agrees to make the dataset publicly available, the author will be eligible for +3 contribution points. Note that we might request the author to clean up their data until it is proper for public release (i.e. formatting, consistency, or additional filtering). In addition to the contribution points obtained from publicly releasing the dataset, the author is also eligible for additional points from the datasheet listing, as mentioned in §3.1, when the datasheet is submitted to NusaCatalogue.

## 3.4 Other ways to contribute in 𝒩

Other than the previously mentioned contributions methods, we also open for other forms of contribution, subject to NusaCrowd's open discussion. To get more details on the open discussion, please join our Slack and Whatsapp group (see §7).

## 3.5 Contribution point recapitulation

The total contribution point for all contributors will be recapped every week by a maintainer, and a contribution matrix will be published and updated on a weekly basis. The contribution matrix and more detailed information to the contribution point can be accessed on the following link[7]. The final score will be recapped in November and the finalized contribution matrix (see Figure 5 for an example) will be published along with the research paper.

## 4 Dataset Licensing and Ownership

NusaCrowd does not make a clone or copy the submitted dataset. The owner and copyright holder will remain to the original data owner. All data access policy will follow the original data licence without any modification from NusaCrowd. NusaCrowd simply downloads and reads the file from the original publicly available data source location when creating the dataloader, and, in addition, the NusaCatalogue datasheet also directly points to the original data site and publication.

## 5 Timeline

The current NusaCrowd movement is started from 25 June 2022 and will be closed on 18 November

---

[7]https://docs.google.com/spreadsheets/d/e/2PACX-1vS3Kbi9s3o_V-1yFRHeONOI7jFnMlUswqKj-D6cpgiSYOSxbijC4DIrjAstqxj-H-EI6I2lFhhyKe5s/pubhtml
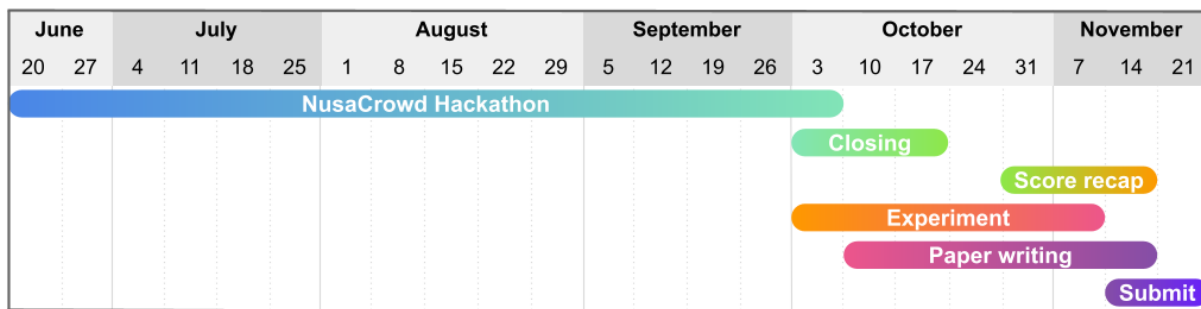
Figure 6: Timeline of the NusaCrowd movement. The datasheet collection and dataloader implementation start from 25 June 2022 to 2 October 2022. Then NusaCrowd will continue with experiments, paper writing, and point recapitulation. The paper will be submitted to a top-level computational linguistics conference, ACL 2023.



Figure 7: Join NusaCrowd's Slack (`https://join.slack.com/t/nusacrowd/shared_invite/zt-1b61t06zn-rf2rWw8WFZCpjVp3iLXf0g`), Whatsapp group (`https://chat.whatsapp.com/Jn4nM6l3kSn3p4kJVESTwv`), and Github (`https://github.com/IndoNLP/nusa-crowd`).

2022. The registration of datasheets and dataloaders will be completed on 2 October 2022. From October onwards, we will focus more on preparing extension and set of experiments to show the benefit of having NusaCrowd platform. In addition, the contribution point for each contributors (see Figure 5 for an example) and the research paper will also be finalized by early November, followed by the final submission of the paper to the Association of Computational Linguistics (ACL) 2023 conference. The detailed phase for the paper development is shown in Figure 6.

## 6 Summary

NLP resources in Indonesian languages, especially the local ones, are extremely low-resource and underrepresented in the research community. There are multitudes factors causing this limitation. Here we solve this problem by initiating the largest Indonesian NLP crowd sourcing efforts, NusaCrowd. In the spirit of fairness, openness, and transparency; NusaCrowd comes with various ways to contribute towards openness and standardization in Indonesian NLP, while at the same time, introducing a scoring mechanism that provides equal chance for all contributors to show the best out of their contributions. We hope that, NusaCrowd can bring a new perspective to all Indonesian NLP practitioners to focus more on openness and collaboration through code and data sharing, complete documentation, and community efforts.

## 7 Call for participation

We invite all Indonesian NLP enthusiasts to participate in NusaCrowd. For any inquiry and further information, contributors can join our community channel on Slack or Whatsapp Group (see Figure 7). Let's work together to advance the progress of Indonesian language NLP! 🇮🇩🇮🇩🇮🇩

## Acknowledgements

# References

Muhammad Abdul-Mageed, AbdelRahim Elmadany, El Moatez Billah Nagoudi, Dinesh Pabbi, Kunal Verma, and Rannie Lin. 2020. Mega-cov: A billion-scale dataset of 100+ languages for covid-19. *arXiv preprint arXiv:2005.06012*.

Alham Fikri Aji, Radityo Eko Prasojo Tirana Noor Fatyanosa, Philip Arthur, Suci Fitriany, Salma Qonitah, Nadhifa Zulfa, Tomi Santoso, and Mahendra Data. 2021. ParaCotta: Synthetic multilingual paraphrase corpora from the most diverse translation sample pair. In *Proceedings of the 35th Pacific Asia Conference on Language, Information and Computation*, pages 533–542, Shanghai, China. Association for Computational Lingustics.

Alham Fikri Aji, Genta Indra Winata, Fajri Koto, Samuel Cahyawijaya, Ade Romadhony, Rahmad Mahendra, Kemal Kurniawan, David Moeljadi, Radityo Eko Prasojo, Timothy Baldwin, Jey Han Lau, and Sebastian Ruder. 2022. One country, 700+ languages: NLP challenges for underrepresented languages and dialects in Indonesia. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7226–7249, Dublin, Ireland. Association for Computational Linguistics.

Zaid Alyafeai, Maraim Masoud, Mustafa Ghaleb, and Maged Saeed AlShaibani. 2021. Masader: Metadata sourcing for arabic text and speech data resources. *CoRR*, abs/2110.06744.

Tri Apriani, Herry Sujaini, and Novi Safriadi. 2016. Pengaruh kuantitas korpus terhadap akurasi mesin penerjemah statistik bahasa bugis wajo ke bahasa indonesia. *J. Sist. dan Teknol. Inf*, 1(1):1–6.

Samuel Cahyawijaya, Genta Indra Winata, Bryan Wilie, Karissa Vincentio, Xiaohong Li, Adhiguna Kuncoro, Sebastian Ruder, Zhi Yuan Lim, Syafri Bahar, Masayu Leylia Khodra, et al. 2021. Indonlg: Benchmark and resources for evaluating indonesian natural language generation. *arXiv preprint arXiv:2104.08200*.

Nindian Puspa Dewi, Joan Santoso, Ubaidi Ubaidi, and Eka Rahayu Setyaningsih. 2020. Combination of genetic algorithm and brill tagger algorithm for part of speech tagging bahasa madura. *Proceeding of the Electrical Engineering Computer Science and Informatics*, 7(2):38–42.

Kaustubh D. Dhole, Varun Gangal, Sebastian Gehrmann, Aadesh Gupta, Zhenhao Li, Saad Mahamood, Abinaya Mahendiran, Simon Mille, Ashish Srivastava, Samson Tan, Tongshuang Wu, Jascha Sohl-Dickstein, Jinho D. Choi, Eduard Hovy, Ondrej Dusek, Sebastian Ruder, Sajant Anand, Nagender Aneja, Rabin Banjade, Lisa Barthe, Hanna Behnke, Ian Berlot-Attwell, Connor Boyle, Caroline Brun, Marco Antonio Sobrevilla Cabezudo, Samuel Cahyawijaya, Emile Chapuis, Wanxiang Che, Mukund Choudhary, Christian Clauss, Pierre Colombo, Filip Cornell, Gautier Dagan, Mayukh Das, Tanay Dixit, Thomas Dopierre, Paul-Alexis Dray, Suchitra Dubey, Tatiana Ekeinhor, Marco Di Giovanni, Rishabh Gupta, Rishabh Gupta, Louanes Hamla, Sang Han, Fabrice Harel-Canada, Antoine Honore, Ishan Jindal, Przemyslaw K. Joniak, Denis Kleyko, Venelin Kovatchev, Kalpesh Krishna, Ashutosh Kumar, Stefan Langer, Seungjae Ryan Lee, Corey James Levinson, Hualou Liang, Kaizhao Liang, Zhexiong Liu, Andrey Lukyanenko, Vukosi Marivate, Gerard de Melo, Simon Meoni, Maxime Meyer, Afnan Mir, Nafise Sadat Moosavi, Niklas Muennighoff, Timothy Sum Hon Mun, Kenton Murray, Marcin Namysl, Maria Obedkova, Priti Oli, Nivranshu Pasricha, Jan Pfister, Richard Plant, Vinay Prabhu, Vasile Pais, Libo Qin, Shahab Raji, Pawan Kumar Rajpoot, Vikas Raunak, Roy Rinberg, Nicolas Roberts, Juan Diego Rodriguez, Claude Roux, Vasconcellos P. H. S., Ananya B. Sai, Robin M. Schmidt, Thomas Scialom, Tshephisho Sefara, Saqib N. Shamsi, Xudong Shen, Haoyue Shi, Yiwen Shi, Anna Shvets, Nick Siegel, Damien Sileo, Jamie Simon, Chandan Singh, Roman Sitelew, Priyank Soni, Taylor Sorensen, William Soto, Aman Srivastava, KV Aditya Srivatsa, Tony Sun, Mukund Varma T, A Tabassum, Fiona Anting Tan, Ryan Teehan, Mo Tiwari, Marie Tolkiehn, Athena Wang, Zijian Wang, Gloria Wang, Zijie J. Wang, Fuxuan Wei, Bryan Wilie, Genta Indra Winata, Xinyi Wu, Witold Wydmański, Tianbao Xie, Usama Yaseen, M. Yee, Jing Zhang, and Yue Zhang. 2021. Nl-augmenter: A framework for task-sensitive natural language augmentation.

Muhammad Dwi Etsa, Herry Sujaini, and Novi Safriadi. 2018. Pengaruh metode dictionary lookup pada cleaning korpus terhadap akurasi mesin penerjemah statistik indonesia–melayu pontianak. *JEPIN (Jurnal Edukasi dan Penelitian Informatika)*, 4(1):49–53.

Jason Alan Fries, Leon Weber, Natasha Seelam, Gabriel Altay, Debajyoti Datta, Samuele Garda, Myungsun Kang, Ruisi Su, Wojciech Kusa, Samuel Cahyawijaya, Fabio Barth, Simon Ott, Matthias Samwald, Stephen Bach, Stella Biderman, Mario Sänger, Bo Wang, Alison Callahan, Daniel León Periñán, Théo Gigant, Patrick Haller, Jenny Chim, Jose David Posada, John Michael Giorgi, Karthik Rangasai Sivaraman, Marc Pàmies, Marianna Nezhurina, Robert Martin, Michael Cullan, Moritz Freidank, Nathan Dahlberg, Shubhanshu Mishra, Shamik Bose, Nicholas Michio Broad, Yanis Labrak, Shlok S Deshmukh, Sid Kiblawi, Ayush Singh, Minh Chien Vu, Trishala Neeraj, Jonas Golde, Albert Villanova del Moral, and Benjamin Beilharz. 2022. Bigbio: A framework for data-centric biomedical natural language processing.

Sebastian Gehrmann, Abhik Bhattacharjee, Abinaya Mahendiran, Alex Wang, Alexandros Papangelis, Aman Madaan, Angelina McMillan-Major, Anna Shvets, Ashish Upadhyay, Bingsheng Yao, et al.

2022. Gemv2: Multilingual nlg benchmarking in a single line of code. *arXiv preprint arXiv:2206.11249*.

Arfinda Ilmania, Abdurrahman, Samuel Cahyawijaya, and Ayu Purwarianti. 2018. Aspect detection and sentiment classification using deep neural network for indonesian aspect-based sentiment analysis. In *2018 International Conference on Asian Language Processing (IALP)*, pages 62–67.

Fajri Koto, Jey Han Lau, and Timothy Baldwin. 2021a. 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, Jey Han Lau, and Timothy Baldwin. 2021b. IndoBERTweet: A pretrained language model for Indonesian Twitter with effective domain-specific vocabulary initialization. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 10660–10668, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Rahmad Mahendra, Alham Fikri Aji, Samuel Louvan, Fahrurrozi Rahman, and Clara Vania. 2021. IndoNLI: A natural language inference dataset for Indonesian. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 10511–10527, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Angelina McMillan-Major, Zaid Alyafeai, Stella Biderman, Kimbo Chen, Francesco De Toni, Gérard Dupont, Hady Elsahar, Chris Emezue, Alham Fikri Aji, Suzana Ilić, et al. 2022. Documenting geographically and contextually diverse data sources: The bigscience catalogue of language data and resources. *arXiv preprint arXiv:2201.10066*.

Wilhelmina Nekoto, Vukosi Marivate, Tshinondiwa Matsila, Timi Fasubaa, Taiwo Fagbohungbe, Solomon Oluwole Akinola, Shamsuddeen Muhammad, Salomon Kabongo Kabenamualu, Salomey Osei, Freshia Sackey, Rubungo Andre Niyongabo, Ricky Macharm, Perez Ogayo, Orevaoghene Ahia, Musie Meressa Berhe, Mofetoluwa Adeyemi, Masabata Mokgesi-Selinga, Lawrence Okegbemi, Laura Martinus, Kolawole Tajudeen, Kevin Degila, Kelechi Ogueji, Kathleen Siminyu, Julia Kreutzer, Jason Webster, Jamiil Toure Ali, Jade Abbott, Iroro Orife, Ignatius Ezeani, Idris Abdulkadir Dangana, Herman Kamper, Hady Elsahar, Goodness Duru, Ghollah Kioko, Murhabazi Espoir, Elan van Biljon, Daniel Whitenack, Christopher Onyefuluchi, Chris Chinenye Emezue, Bonaventure F. P. Dossou, Blessing Sibanda, Blessing Bassey, Ayodele Olabiyi, Arshath Ramkilowan, Alp Öktem, Adewale Akinfaderin, and Abdallah Bashir. 2020. Participatory research for low-resourced machine translation: A case study in African languages.

In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2144–2160, Online. Association for Computational Linguistics.

Iroro Orife, Julia Kreutzer, Blessing Sibanda, Daniel Whitenack, Kathleen Siminyu, Laura Martinus, Jamiil Toure Ali, Jade Z. Abbott, Vukosi N. Marivate, Salomon Kabongo, Musie Meressa, Espoir Murhabazi, Orevaoghene Ahia, Elan Van Biljon, Arshath Ramkilowan, Adewale Akinfaderin, Alp Öktem, Wole Akin, Ghollah Kioko, Kevin Degila, Herman Kamper, Bonaventure Dossou, Chris Emezue, Kelechi Ogueji, and Abdallah Bashir. 2020. Masakhane - machine translation for africa. *CoRR*, abs/2003.11529.

Xiaoman Pan, Boliang Zhang, Jonathan May, Joel Nothman, Kevin Knight, and Heng Ji. 2017. Cross-lingual name tagging and linking for 282 languages. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1946–1958.

Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. 2022. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint arXiv:2206.04615*.

Bryan Wilie, Karissa Vincentio, Genta Indra Winata, Samuel Cahyawijaya, Xiaohong Li, Zhi Yuan Lim, Sidik Soleman, Rahmad Mahendra, Pascale Fung, Syafri Bahar, et al. 2020. Indonlu: Benchmark and resources for evaluating indonesian natural language understanding. *arXiv preprint arXiv:2009.05387*.

Genta Indra Winata, Alham Fikri Aji, Samuel Cahyawijaya, Rahmad Mahendra, Fajri Koto, Ade Romadhony, Kemal Kurniawan, David Moeljadi, Radityo Eko Prasojo, Pascale Fung, et al. 2022. Nusax: Multilingual parallel sentiment dataset for 10 indonesian local languages. *arXiv preprint arXiv:2205.15960*.