

| No. | Dataset name | Source(s) | Dataset description | Dataset URL | Prep/Type | Dataset location | Dataset release type | Dataset license | Dataset metadata | Dataset content | Dataset format | Dataset size | Number of instances per sample | Time span (years) | Validation split | Test split | No split data | Data size (MB) | Data size (KB) | Data size (bytes) | Dataset source | Dataset publisher | Dataset year | Dataset version | Dataset status | Dataset access | Dataset type | Dataset content | Dataset access | Dataset type | Dataset content | Dataset access | Dataset type | Dataset content | Dataset access | Dataset type | Dataset content | Dataset access |
|-----|--------------|-----------|---------------------------------|---|-----------|------------------|----------------------|-----------------|------------------|-----------------|----------------|--------------|--------------------------------|-------------------|------------------|------------|---------------|----------------|----------------|-------------------|----------------|-------------------|--------------|-----------------|----------------|----------------|--------------|-----------------|----------------|--------------|-----------------|----------------|--------------|-----------------|----------------|--------------|-----------------|----------------|
| 1 | Amharic | Amharic | A set of Amharic text documents | https://www.kaggle.com/datasets/Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic | Amharic |

| No. | Dataset name | Subject | Dataset description | Dataset URL | Project/Task URL | Dataset version | Dataset collection region | Dataset language | Dataset modality | Dataset domain | Dataset format | Dataset generation/collector type | Dataset generation/validation type | Number of records per sample | Sampling agreement | Task type | Validation split data size | Test split data size | No. split data size | Task size unit | Total data size | Dataset provider | Dataset year | Dataset year published | Publication venue | Dataset access | Dataset license | Does the data have any sensitive information? | Does the data have any personally identifiable information? | Submission date | Approval Status | Dataset owner | GitHub issue URL |
|-----|---------------------------|---|---|-------------|------------------|-----------------|---------------------------|------------------------------|------------------|----------------|------------------------------------|-----------------------------------|------------------------------------|------------------------------|--------------------|-----------|----------------------------|----------------------|---------------------|----------------|-----------------|---|------------------------------|------------------------|-------------------|------------------|-----------------|---|---|------------------|---------------------|--------------------------------------|--------------------------------------|
| 38 | Speech-captioning | | This is a subset of ASR dataset | | | | NA | Image-to-Text | Language Video | Audio | Creative Commons Expert generated | None | | | | | | | | Images | 2100 | NA | International Speech Library | 2022 | EMBLP | Free open access | CC-BY | No | No | 2023-10-02 08:00 | Approved | Shoun_Lightbulb | Link to GitHub issue |
| 40 | Speech-MT | | This version of ASR dataset | | | | NA | Language-Machine | Language Video | Audio | Other (other) | Clustering | None | | | | | | | Instances | 1200 | NA | International Speech Library | 2022 | EMBLP | Free open access | CC-BY | No | No | 2023-10-02 08:00 | Approved | Shoun_L | Link to GitHub issue |
| 41 | Speech-search | | This version of ASR dataset | | | | NA | Automatic Speech Recognition | Language Speech | Audio | Other (other) | Clustering | Manual (full) | | | | | | | Instances | 2000 | NA | International Speech Library | 2022 | EMBLP | Free open access | CC-BY | No | No | 2023-10-02 08:00 | Approved | Shoun_L | Link to GitHub issue |
| 43 | Speech | | ASR dataset | | | | Unknown | Text-to-Speech | Language Video | Audio | Creative Commons Expert generated | None | | | | | | | | Instances | 1477 | NA | International Speech Library | 2022 | EMBLP | Free open access | CC-BY | No | No | 2023-10-02 08:00 | Approved | Shoun_L | Link to GitHub issue |
| 45 | ASR-PT3 | ASR-PT3 | This set is for ASR-PT3 | | | | Mexico | Communication | Language | Audio | Unknown | Unknown | Manual (full) | | | | | | | Instances | 40 | Unknown | NA | 2022 | NA | Free | CC-BY | No | No | 2024-03-14 08:00 | Approved | pt3 | Link to GitHub issue |
| 46 | ASR-C | The dataset for ASR-C | | | | | Unknown | Language | Language | Audio | Unknown | Unknown | Manual (partial) | | | | | | | Instances | 264700 | National Center for Speech and Language Acquisition | 2022 | CC-BY | Free | CC-BY | No | No | 2024-03-07 08:00 | Approved | ASR-C | Link to GitHub issue | |
| 48 | ASR-B | ASR-B is a set of ASR-B | | | | | Unknown | Language | Language | Audio | Unknown | Unknown | Manual (full) | | | | | | | Instances | 649100 | Yonsei University | 2022 | CC-BY | Free | CC-BY | No | No | 2024-03-14 08:00 | Approved | ASR-B | Link to GitHub issue | |
| 49 | ASR-T4 | ASR-T4 | The dataset for ASR-T4 | | | | Thailand | Communication | Language | Audio | Unknown | Unknown | Manual (full) | | | | | | | Instances | 0 | Surong University | 2022 | CC-BY | Free | CC-BY | No | No | 2023-11-09 08:00 | Approved | ASR-T4 | Link to GitHub issue | |
| 47 | CASA | CASA is a set of CASA | | | | | Unknown | Language | Language | Audio | Creative Commons Expert generated | Clustering | Manual (full) | 0 | | 0 | 0 | | | Instances | 1000 | Chulalongkorn University | 2018 | EMBLP | Free | CC-BY | No | No | 2023-11-09 08:00 | Approved | ASR-T4 | Link to GitHub issue | |
| 48 | CC-100 | | This corpus is for CC-100 | | | | NA | Language-Machine | Language | Text | MT (mt) | Clustering | None | | | | | | | Instances | 70000 | Meta | Unsupervised | 2020 | CC-BY | Free | CC-BY | No | No | 2023-11-09 08:00 | Approved | CC-100 | Link to GitHub issue |
| 49 | CC-Aligned | CC-Aligned | CC-Aligned | | | | NA | Machine-Translation | Language | Text | Unknown | Unknown | Manual (partial) | | | | | | | Instances | 6,274,000 | Meta | CC-Aligned | 2020 | EMBLP | Free | CC-BY | No | No | 2023-11-09 08:00 | Approved | CC-Aligned | Link to GitHub issue |
| 50 | CC-Aligned Sentence Pairs | | This dataset is for CC-Aligned Sentence Pairs | | | | NA | Machine-Translation | Language | Text | Unknown | Unknown | Manual (partial) | | | | | | | Instances | 2,000 | Meta | Low-Resource | 2020 | EMBLP | Free | CC-BY | No | No | 2023-11-09 08:00 | Approved | CC-Aligned | Link to GitHub issue |
| 51 | CC-100 | CC-100 | CC-100 | | | | Language-Machine | Language | Text | Other (other) | Clustering | None | | | | | | | | Instances | 30,000 | Meta | Unsupervised | 2020 | EMBLP | Free | CC-BY | No | No | 2023-11-09 08:00 | Approved | CC-100 | Link to GitHub issue |
| 52 | CC-100 | CC-100 | CC-100 | | | | NA | Image-to-Text | Language Video | Multi-modal | Creative Commons Machine generated | None | | 30,000 | | 0 | 0 | | | Instances | 100,000 | Google | Creative Commons | 2022 | EMBLP | Free | CC-BY | No | No | 2023-11-09 08:00 | Approved | CC-100 | Link to GitHub issue |
| 53 | CC-Media | The CC-Media | | | | | US | Language-Machine | Language | Text | BSD license | Clustering | Automatic | | | | | | | Instances | 0 | Meta | CC-Media | 2021 | CC-BY | Free | CC-BY | No | No | 2023-10-21 08:00 | Approved with notes | CC-Media | Link to GitHub issue |
| 54 | CC-Media | The CC-Media | | | | | NA | Image-to-Text | Language | Text | Creative Commons | Clustering | Manual (full) | 0 | 0 | 0 | 0 | | | Instances | 4,000 | National Center for Speech and Language Acquisition | 2022 | EMBLP | Free | CC-BY | No | No | 2023-11-09 08:00 | Approved | CC-Media | Link to GitHub issue | |
| 55 | CC-Media | CC-Media | | | | | US | Communication | Language | Text | CC-BY license | Clustering | Automatic | | | | | | | Instances | 300 | Meta | CC-Media | 2018 | EMBLP | Free | CC-BY | No | No | 2023-11-09 08:00 | Approved | CC-Media | Link to GitHub issue |
| 56 | CC-Media | CC-Media | | | | | US | Image-to-Text | Language | Text | Creative Commons | Clustering | Automatic | | | | | | | Instances | 100,000 | Meta | CC-Media | 2018 | EMBLP | Free | CC-BY | No | No | 2023-11-09 08:00 | Approved | CC-Media | Link to GitHub issue |
| 57 | CC-Media | CC-Media | | | | | NA | Text-to-Speech | Language | Text | Unknown | Unknown | Manual (partial) | | | | | | | Instances | 0 | Meta | CC-Media | 2022 | EMBLP | Free | CC-BY | No | No | 2024-03-02 08:00 | Approved with notes | CC-Media | Link to GitHub issue |

| No. | Dataset name | Subject | Dataset description | Dataset URL | Project/Year | Dataset collection region | Dataset language | Dataset modality | Dataset domain | Dataset format | Dataset acquisition collection style | Dataset annotation style | Number of instances per sample | Image resolution | Track length (s) | Video length (s) | Text length (words) | No. of pages | Task size (words) | Task size (images) | Dataset provider | Dataset project URL | Dataset of original paper URL | Publication year | Dataset license | Derived from | Does the data have sensitive information? | Does the data have personally identifiable information (PII)? | Submission date | Approval status | Dataset owner | GitHub repo URL |
|-----|--------------|------------------|---|--------------------------------------|--------------|---------------------------|------------------|------------------|----------------|---|--------------------------------------|--------------------------|--------------------------------|------------------|------------------|------------------|---------------------|------------------------------------|--|--------------------|------------------|---------------------|-------------------------------|------------------|-----------------|------------------|---|---|--|-----------------|---------------|-----------------|
| 68 | OOO Captions | Image Captioning | The dataset is a collection of image-caption pairs from the COCO dataset. | https://coco dataset | 2017 | NA | English | Image | General | Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License | Machine-generated | None | 113,267 | 8,000 | 8,000 | words | 120,000 | Microsoft, Ray A. Rasmussen et al. | https://github.com/matterport/mc2017 | 2018 | CC-BY-NC-SA | 2018 | CC-BY-NC-SA | No | No | 2023-10-26 16:12 | Approved | hny_sml | https://github.com/hny_sml | | | |
| 69 | OOO-ML | Image Captioning | The dataset is a collection of image-caption pairs from the COCO dataset. | https://coco dataset | 2017 | NA | English | Image | General | Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License | Machine-generated | None | 113,267 | 8,000 | 8,000 | words | 120,000 | Microsoft, Ray A. Rasmussen et al. | https://github.com/matterport/mc2017 | 2018 | CC-BY-NC-SA | 2018 | CC-BY-NC-SA | No | No | 2023-11-16 21:05 | Approved | hny_sml | https://github.com/hny_sml | | | |
| 70 | OOO-ML | Image Captioning | The dataset is a collection of image-caption pairs from the COCO dataset. | https://coco dataset | 2017 | NA | English | Image | General | Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License | Machine-generated | None | 113,267 | 8,000 | 8,000 | words | 120,000 | Microsoft, Ray A. Rasmussen et al. | https://github.com/matterport/mc2017 | 2018 | CC-BY-NC-SA | 2018 | CC-BY-NC-SA | No | No | 2023-11-16 21:05 | Approved | hny_sml | https://github.com/hny_sml | | | |
| 71 | OOO-ML | Image Captioning | The dataset is a collection of image-caption pairs from the COCO dataset. | https://coco dataset | 2017 | NA | English | Image | General | Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License | Machine-generated | None | 113,267 | 8,000 | 8,000 | words | 120,000 | Microsoft, Ray A. Rasmussen et al. | https://github.com/matterport/mc2017 | 2018 | CC-BY-NC-SA | 2018 | CC-BY-NC-SA | No | No | 2023-11-16 21:05 | Approved | hny_sml | https://github.com/hny_sml | | | |
| 72 | OOO-ML | Image Captioning | The dataset is a collection of image-caption pairs from the COCO dataset. | https://coco dataset | 2017 | NA | English | Image | General | Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License | Machine-generated | None | 113,267 | 8,000 | 8,000 | words | 120,000 | Microsoft, Ray A. Rasmussen et al. | https://github.com/matterport/mc2017 | 2018 | CC-BY-NC-SA | 2018 | CC-BY-NC-SA | No | No | 2023-11-16 21:05 | Approved | hny_sml | https://github.com/hny_sml | | | |

| No | Dataset name | Dataset URL | Dataset paper URL | Submission Date |
|----|---|---|---|------------------|
| 1 | generated_reviews_enth | https://github.com | https://arxiv.org/p | 2023-11-17 14:29 |
| 2 | Dengue Filipino | https://huggingfa | https://ieeexplore | 2023-11-20 01:58 |
| 3 | Leipzig Corpora Collection | https://wortschat | http://www.lrec-c | 2023-12-15 15:11 |
| 4 | GlobalVoices | https://opus.nlpl.e | http://www.lrec-c | 2023-12-19 18:38 |
| 5 | DeepLontar | https://doi.org/10 | https://www.natu | 2023-12-23 17:47 |
| 6 | Corpus Crawler | https://github.com | | 2023-12-23 20:50 |
| 7 | OpenMSD | https://github.com | https://arxiv.org/p | 2023-12-26 08:11 |
| 8 | QED | https://opus.nlpl.e | https://aclantholo | 2023-12-31 13:21 |
| 9 | GNOME | https://opus.nlpl.e | https://aclantholo | 2023-12-31 16:00 |
| 10 | CCMatrix | https://opus.nlpl.e | https://aclantholo | 2023-12-31 16:30 |
| 11 | Lio and the Central Flores languages | https://archive.m | https://studentthe | 2024-01-07 15:00 |
| 12 | Thai Gov v2 Corpus | https://github.com | | 2024-01-13 22:38 |
| 13 | ParaNames | https://github.com | https://aclantholo | 2024-01-16 16:21 |
| 14 | CLIRMatrix | https://github.com | https://aclantholo | 2024-02-02 14:38 |
| 15 | MM-Sum | https://drive.goo | https://aclantholo | 2024-02-17 10:38 |
| 16 | BEST | https://github.com | http://pioneer.chu | 2024-03-04 00:00 |
| 17 | LEXiTRON | https://opend-por | https://www.sem | 2024-03-10 16:11 |
| 18 | multilingual-NLI-26lang-2mil7 | https://huggingfa | https://www.cam | 2024-03-16 01:47 |
| 19 | Onto4All | https://huggingfa | https://huggingfa | 2024-03-16 11:38 |
| 20 | Amanatun wordlist | https://catalog.pa | https://www.ocse | 2024-03-22 13:11 |
| 21 | ProSub | https://github.com | https://www.anlp | 2024-03-22 17:00 |
| 22 | InterBEST-2009 | http://thailang.ne | https://ieeexplore | 2024-03-23 13:21 |
| 23 | BKD-Prosub | https://babyai-hu | https://www.virac | 2024-03-29 11:12 |
| 24 | AlloVera | https://github.com | https://aclantholo | 2024-04-01 11:28 |
| 25 | OpenSpeech Dataset V1 by Wang | https://www.wang | | 2024-04-04 15:00 |
| 26 | Thai Handwritten Free Datasets by Wang: Data Market | https://www.wang | | 2024-04-04 15:10 |
| 27 | SEACrowd Instruct Multi-task Collection | https://github.com | | 2024-06-19 00:00 |
| 28 | Philippine Language Database (PLD) | Gitlab link not yet | https://aclantholo | 2024-07-03 14:11 |
| 29 | VLUE | https://uitnlprou | https://arxiv.org/p | 2024-07-03 14:21 |

| No | Dataset name | Dataset URL | Dataset paper URL | Submission Date | Notes |
|----|--|---|---|------------------|---|
| 1 | Belebele | https://github.com | https://github.com | 2023-11-01 10:00 | Duplicated #113 |
| 2 | Dakshina | https://github.com | https://aclanthology.org | 2023-11-21 00:21 | This Tamil is not from SEA |
| 3 | SLR65 | https://www.open | https://aclanthology.org | 2023-11-30 00:41 | this Tamil dataset is not specifically collected from SEA region |
| 4 | SLR80 | https://www.open | https://aclanthology.org | 2023-11-30 00:51 | Based on discussion with Holy, openSLR is considered as one single dataset. So, submissions ID 272 and 273 are put together as ID 278 |
| 5 | SLR42 | https://www.open | https://www.isca | 2023-11-30 00:51 | Based on discussion with Holy, openSLR is considered as one single dataset. So, submissions ID 272 and 273 are put together as ID 278 |
| 6 | BLOOM LM | https://huggingface | https://huggingface | 2023-12-03 21:11 | Duplicated with #277 |
| 7 | SEA Wikisource | https://huggingface | | 2023-12-19 02:11 | It is another Wikipedia dump, with no peer-reviewed publication. We had already SEA Wiki, a Wikipedia dump dataset (submission #127) Another submission about Wikipedia dump is submission #308 |
| 8 | VISTEC-TP-TH-21 | https://github.com | https://aclanthology.org | 2023-12-21 13:01 | duplicate to submission #237 |
| 9 | Thai-Alpaca | https://huggingface | | 2023-12-22 16:41 | The dataset uses gibberish Thai language—possibly due to bad results from Google Cloud Translation, which is used to translate the data from English to Thai. Two Thai researchers confirmed the data's unusability. |
| 10 | Thai Databricks Dolly | https://huggingface | | 2023-12-22 16:51 | The dataset uses gibberish Thai language—possibly due to bad results from Google Cloud Translation, which is used to translate the data from English to Thai. Two Thai researchers confirmed the data's unusability. |
| 11 | Thai HH-RLHF | https://huggingface | | 2023-12-22 16:51 | The dataset uses gibberish Thai language—possibly due to bad results from Google Cloud Translation, which is used to translate the data from English to Thai. Two Thai researchers confirmed the data's unusability. |
| 12 | Thai GPTeacher | https://huggingface | | 2023-12-22 17:01 | The dataset uses gibberish Thai language—possibly due to bad results from Google Cloud Translation, which is used to translate the data from English to Thai. Two Thai researchers confirmed the data's unusability. |
| 13 | Singlish treebank | https://github.com | https://aclanthology.org | 2023-12-23 17:01 | This dataset is covered in the latest release (STB_EXT, submission #323) |
| 14 | Singlish treebank v2 (STB-EXT) | https://github.com | https://dl.acm.org | 2023-12-23 17:11 | duplicate to submission #258 |
| 15 | eBible | https://github.com | | 2023-12-23 20:41 | This doesn't seem like a bible corpus, probably need to contact the submitter > Update 10 Jan, it is a bible corpus, the script crawl it in some ways I guess, no publication, how to proceed? > Update 23 Jan by Holy: Joel (the contributor) contacted me and said he had submitted a better version of eBible datasheet (no. 345). I'll make that one in review and reject this one. |
| 16 | M3IT-80 | https://huggingface | https://arxiv.org/p | 2023-12-25 23:31 | duplicate with M3IT |
| 17 | JWSign | https://github.com | https://arxiv.org/a | 2023-12-26 07:31 | The github page mentioned that the data should have been released on around November - December 2023 Per 1 April 2024, data is not available (NOT free upon request). Notes: As complementary information, JW300 dataset was pulled out from public access due to license problem |
| 18 | Leipzig Corpora Collection | https://corpora.ur | https://link.spring | 2023-12-27 23:21 | duplicate with 289 |
| 19 | MuMIN | https://data.bris.a | https://dl.acm.org | 2023-12-29 16:21 | 31/05/2024: All data cannot be decoded after unzipping. To compile the dataset a Twitter API key is required, which has become more difficult these days. However, the submission includes a link to a website where a full raw version of the dataset can be downloaded Among the languages listed, it is not clear if these Tweets were generated in Southeast Asia The other issue is that the non-English data is machine translated, so this data is not generated in SEA - Added the total number of Tweets, although the SEA language tweets will only be a very small percentage of the total - License is unclear, the submitted license is the one stated in the associated publication and in the readme file that downloads with the data, but the dataset website link states another license |
| 20 | MuMIN | https://data.bris.a | https://dl.acm.org | 2023-12-29 16:31 | duplicate with submission #361 |
| 21 | TED2020 | https://opus.nlpl.€ | https://arxiv.org/a | 2024-01-01 18:01 | Invalid link |
| 22 | CLICK-ID | https://data.menc | https://www.scier | 2024-01-04 16:31 | update data size (only the headlines) holy: duplicate #47. removed its dataloader issue. |
| 23 | Multilingual Open Text (MOT) | https://github.com | https://aclanthology.org | 2024-01-16 13:11 | duplicate with no.400 |
| 24 | Mesolithica muftwp.gov.my scrape | https://huggingface | | 2024-01-19 09:21 | duplicate with no.411 |
| 25 | Malaysia-AI government websites scrape | https://huggingface | | 2024-01-21 05:01 | We decided that the crawled data from mesolithica https://huggingface.co/datasets/mesolithica/crawl-my-website . Since there are 3 datasets pointing to these URL, we will just accept one of the three |
| 26 | Malaysia AI GovDocs scrape | https://huggingface | | 2024-01-21 05:11 | We decided that the crawled data from mesolithica https://huggingface.co/datasets/mesolithica/crawl-my-website . Since there are 3 datasets pointing to these URL, we will just accept one of the three |
| 27 | Aishell | https://www.aishell | https://arxiv.org/a | 2024-02-11 11:31 | It is not indigenous language in SEA, while the data was not collected in SEA region |
| 28 | THCHS-30 | http://www.opens | https://arxiv.org/a | 2024-02-11 17:51 | It is not indigenous language in SEA, while the data was not collected in SEA region |
| 29 | PFSA-ID-MED | https://github.com | https://doi.org/10 | 2024-03-07 09:21 | Merge this on into #462 |
| 30 | PFSA-ID-TEST | https://github.com | https://doi.org/10 | 2024-03-07 09:31 | Merge this on into #462 |
| 31 | Thai Wikipedia Dump | https://github.com | https://arxiv.org/p | 2024-03-10 16:21 | There is no wikipedia dump data provided in the repo, nonetheless there is the cleaned pretraining data of WangchanBERT: https://github.com/vistec-AI/thai2transformers/releases/tag/att-v1.0 |
| 32 | tamilmixsentiment | https://huggingface | https://aclanthology.org | 2024-03-15 12:01 | This Tamil is not from SEA region. |
| 33 | tamil-alpaca | https://huggingface | https://arxiv.org/a | 2024-03-15 12:11 | This Tamil is not from SEA region. |
| 34 | ULCA-ASR | https://github.com | https://bhashini.g | 2024-03-15 12:31 | This Tamil is not from SEA region. |
| 35 | IISc-MILE Tamil ASR Corpus | https://openslr.on | https://arxiv.org/a | 2024-03-15 12:31 | This Tamil is not from SEA region. |
| 36 | HopeEDI | https://github.com | https://aclanthology.org | 2024-03-15 12:51 | This Tamil is not from SEA region. |
| 37 | MinangNLP MT - Bilingual dictionary | https://github.com | https://aclanthology.org | 2024-03-15 20:21 | MinangNLP MT has been carried-over from NusaCrowd |
| 38 | CI-AVSR | https://github.com | https://aclanthology.org | 2024-03-15 21:31 | This is Cantonese HK dataset. |
| 39 | UIT-VNewsQA | https://sites.goog | https://dl.acm.org | 2024-03-15 22:01 | no clear information in the webpage how to obtain the dataset. consider to submit this data to private datasheet submission |
| 40 | UIT-VSFC | https://drive.goog | https://ieeexplore | 2024-03-15 22:11 | duplicate with #304 |
| 41 | OASST2 | https://github.com | https://drive.goog | 2024-03-16 11:11 | duplicate #448? please check |

| No | Dataset name | Dataset URL | Dataset paper URL | Submission Date | Notes |
|----|--|---|---|------------------|--|
| 42 | multi-figqa | https://github.com | | 2024-03-18 18:11 | duplicate with #225 |
| 43 | Burmese speech dataset | https://www.open | https://aclantholo | 2024-03-20 15:31 | duplicate with #278 |
| 44 | thai-policy-statements | https://github.com | | 2024-03-24 22:21 | <p>This is NOT data for a summarization task. It is unannotated corpus. No clear metadata explaining the data. No peer-reviewed publication supports whether the data is collected using academic methodology.</p> <p>License is incorrect. Should be CC-0 instead of Apache 2.0</p> |
| 45 | Thai WIKI QA | https://aiforthai.in | | 2024-03-24 23:01 | <p>The link does not provide dataset. It is a service / an API. If you thought that the service can provide the data, please submit into private datasheet.</p> |
| 46 | UIT-ViNewsQA | https://sites.goog | https://arxiv.org/p | 2024-04-01 14:41 | duplicate with #492 |
| 47 | Multilingual Lexical Simplification - Filipino | https://github.com | https://sites.goog | 2024-04-02 15:21 | The author of this dataset not allows anyone to re-release this dataset to any open-webs or any sites. |
| 48 | id-hatespeech-detection | https://github.com | https://ieeexplore | 2024-04-02 17:31 | Duplicate from #28 |
| 49 | final_project | https://github.com | https://arxiv.org/a | 2024-04-03 16:11 | <p>Duplicate of #128, submitter did not check for existing datasheet There are 2 datasets on sentiment analysis of AiryRooms review already This submission includes the link to the paper author's Github page rather than the IndoNLP page</p> |

Here are potential SEA datasets to be included in SEACrowd, but not yet submitted.

If a dataset has been submitted through the form, it will be automatically removed from this list if the dataset name is identical.

| No | Dataset Name | Dataset URL | Notes |
|----|--------------|---|---|
| 1 | CanVEC | https://github.com | https://aclanthology.org/2020.lrec-1.507.pdf |