

Sri Lanka Document Datasets: A Large-Scale, Multilingual Resource for Law, News, and Policy

Nuwan I. Senaratna
nuwans@alumni.stanford.edu

October 5, 2025

Abstract

We present a collection of open, machine-readable document datasets covering parliamentary proceedings, legal judgments, government publications, news, and tourism statistics from Sri Lanka. The collection currently comprises 215,552 documents (59.6 GB) across 13 datasets in Sinhala, Tamil, and English, updated daily and mirrored on GitHub and Hugging Face. These datasets aim to support research in computational linguistics, legal analytics, socio-political studies, and multilingual natural language processing. We describe the sources, collection pipeline, formats, and potential use cases, while discussing licensing and ethical considerations.

aims to bridge this gap by collecting, cleaning, and organizing key public documents into standardized, machine-readable formats. It unifies diverse materials—from Hansards and court judgements to news articles and tourism reports—under a common data framework. All datasets are openly licensed and continuously updated to ensure reproducibility and public transparency.

This effort is particularly significant for data-driven research in low-resource contexts. By providing structured data in Sinhala, Tamil, and English, the project supports the development of natural language processing models, cross-lingual studies, and digital humanities research. The datasets also enable policy analysis, legal precedent tracking, and media monitoring in a transparent, open science environment.

1 Introduction

Sri Lanka’s digital record of law, policy, and media is fragmented across numerous government and private sources. Much of this information exists as PDFs or web pages, often lacking machine-readable structure or public archival consistency. This fragmentation limits access for citizens, journalists, and researchers interested in the island’s governance, history, and socio-economic trends.

In this paper, we describe the scope and structure of these datasets, outline the scraping and curation processes, and highlight their potential applications in AI, governance, and public knowledge. Our goal is to create a living data archive that strengthens civic engagement and academic research through open, verifiable information.

The Sri Lanka Document Datasets initiative

2 Related Work

The study of open datasets has been central to the development of natural language processing (NLP) and computational social science. Large corpora such as Common Crawl[1], Wikipedia Dumps[2], and OpenWebText[3] have powered models that generalize across domains. However, these resources are dominated by data from high-resource languages and global institutions.

Regional initiatives have sought to address this imbalance by creating domain-specific collections. Examples include the Indian Kanoon legal corpus[4], the OpenSubtitles multilingual dataset[5], and the African News Corpus[6]. Such datasets have improved representation for under-resourced languages and enabled comparative linguistic research.

In South Asia, efforts remain scattered and often focus on individual media outlets or institutions. Sri Lanka, in particular, lacks consolidated and machine-readable documentation of its public records. Prior datasets were either limited in size, language coverage, or temporal continuity[7][8].

The Sri Lanka Document Datasets aim to fill this gap by aggregating diverse sources—parliamentary debates, court judgments, gazettes, press releases, and news—into a unified, open, and multilingual repository. This complements global initiatives by providing a structured view of a unique national information ecosystem.

3 Datasets

3.1 Hansard

A Hansard is the official verbatim record of parliamentary debates, preserving lawmakers' words

and decisions for history, law, and public accountability.

- Time Updated: 2025-10-04 18:36:13
- Number of Documents: 1,665
- Date Range: 2006-02-01 to 2025-09-24
- Dataset Size: 17.9 GB

3.2 Appeal Court Judgements

A Court of Appeal judgment is a higher court ruling that reviews decisions of lower courts, shaping legal precedent and protecting citizens' rights.

- Time Updated: 2025-10-04 18:44:41
- Number of Documents: 10,146
- Date Range: 2012-04-23 to 2025-10-03
- Dataset Size: 10.3 GB

3.3 Supreme Court Judgements

A Supreme Court judgment is a binding legal decision that interprets the Constitution and laws, shaping justice, governance, and citizens' rights.

- Time Updated: 2025-10-04 18:44:12
- Number of Documents: 1,575
- Date Range: 2016-07-22 to 2025-10-03
- Dataset Size: 1.3 GB

3.4 Police Press Releases

A police press release is an official update from law enforcement on crimes, arrests, safety alerts, or public notices, ensuring transparency and public awareness.

- Time Updated: 2025-10-04 18:28:49
- Number of Documents: 736
- Date Range: 2025-05-01 to 2025-10-03
- Dataset Size: 250.5 MB

3.5 Acts

A legal act is a law passed by Parliament that governs rights, duties, economy, and society, shaping daily life and national policy.

- Time Updated: 2025-10-04 18:42:02
- Number of Documents: 3,928
- Date Range: 1981-01-22 to 2025-09-22
- Dataset Size: 6.8 GB

3.6 Bills

A Bill is a draft law proposed in Parliament. It becomes binding once passed and enacted, shaping governance, rights, and daily life in the country.

- Time Updated: 2025-10-04 18:37:26
- Number of Documents: 4,077
- Date Range: 2010-05-10 to 2025-10-26
- Dataset Size: 1.8 GB

3.7 Extraordinary Gazettes

An Extraordinary Gazette is an official government publication used to announce urgent laws, regulations, or public notices with immediate effect.

- Time Updated: 2025-10-04 18:39:51
- Number of Documents: 101,532
- Date Range: 2010-01-01 to 2025-10-03
- Dataset Size: 19.4 GB

3.8 Cabinet Decisions

A Sri Lanka Cabinet Decision is an official policy or action agreed by the Cabinet of Ministers, shaping governance, law, and national development in the country.

- Time Updated: 2025-10-04 18:08:09

- Number of Documents: 10,369
- Date Range: 2010-09-27 to 2025-09-22
- Dataset Size: 125.3 MB

3.9 Treasury Press Releases

A Sri Lanka Treasury press release shares key govt financial updates—on budgets, debt, or policy—vital for transparency, guiding investors, citizens, and markets on the nation’s economic direction.

- Time Updated: 2025-10-04 18:18:09
- Number of Documents: 133
- Date Range: 2015-09-08 to 2025-07-30
- Dataset Size: 142.9 MB

3.10 Pmd Press Releases

A Sri Lanka Presidential Media Division press release shares official updates on national decisions, policies, or events. It’s vital as the authoritative source ensuring transparency and public awareness.

- Time Updated: 2025-09-26 08:23:47
- Number of Documents: 2,182
- Date Range: 2024-09-23 to 2025-09-24
- Dataset Size: 55.9 MB

3.11 News

A collection of news documents.

- Time Updated: 2025-10-04 18:41:30
- Number of Documents: 79,049
- Date Range: 2021-09-12 to 2025-10-04
- Dataset Size: 1.2 GB

3.12 Tourism Weekly Reports

Report on Weekly Tourist Arrivals to Sri Lanka.

- Time Updated: 2025-10-04 18:36:40
- Number of Documents: 33
- Date Range: 2023-01-01 to 2025-09-01
- Dataset Size: 91.5 MB

3.13 Tourism Monthly Reports

Report on Monthly Tourist Arrivals to Sri Lanka.

- Time Updated: 2025-10-04 18:36:59
- Number of Documents: 127
- Date Range: 2015-01-01 to 2025-08-01
- Dataset Size: 308.9 MB

4 Data Collection Pipeline

Data Collection Pipeline Our pipeline is automated, reproducible, and resilient. It continuously discovers, ingests, parses, validates, and versions documents from official Sri Lankan sources.[9]

GitHub Actions orchestrates the workflow. Cron jobs run several times per day, per dataset. A matrix strategy isolates each source, allowing independent retries without blocking others. Secrets manage tokens; caches speed I/O.[10]

Each run is idempotent and incremental. Before crawling, we load a manifest of known items. New or changed items are detected by stable keys (URL + date) and content hashes. Only deltas are committed to the repository.[11]

Crawling is implemented in Python with Selenium in headless Chromium. We wait for dynamic content via explicit conditions (e.g., presence of selectors), handle pagination, and

capture canonical URLs.[12]

Politeness is enforced. We respect robots.txt, throttle requests, randomize delays, and apply exponential backoff on transient failures. Errors are logged and surfaced in the Actions summary for rapid triage.[13]

Raw artifacts are preserved alongside parsed representations. For each document we store the fetched HTML or PDF, plus normalized JSON with metadata (title, date, source, language, hashes) to enable downstream reproducibility.[14]

PDF parsing uses PyMuPDF (fitz). For each PDF, we extract text, metadata, and layout blocks, retain page boundaries, and normalize Unicode. When images contain embedded text, PyMuPDF's text extraction captures vector text regions.[15]

The parser records coordinates for blocks, allowing approximate structure recovery (sections, headings, tables) where present. Heuristics join hyphenated lines and preserve numbering and legal citations.[16]

Quality gates run in CI. Schemas are validated, required fields are enforced, and checksums guard against corruption. Unit tests cover fetching, parsing, and serialization, and fail the job on regressions.[17]

Historical coverage was built via a back-population pipeline. We iterate over archival indexes and date ranges, enqueue jobs in batches, checkpoint progress, and resume safely after interruptions.[18]

Transparency is prioritized. Run metadata, document counts, and last-updated timestamps are published to README badges. Commit messages summarize deltas, enabling clear, auditable provenance across releases.[19]

6 Conclusion and Future Work

5 Licensing and Access

All datasets and code are openly available to the public. The repositories are hosted on GitHub under the MIT License, which permits reuse, modification, and redistribution with attribution to the original author.[20]

This permissive model encourages transparency and collaboration. Researchers, developers, and institutions can integrate the datasets into their pipelines without restrictive terms or commercial barriers.[19]

Each dataset repository includes structured metadata, versioned releases, and README files with descriptive statistics and provenance. All assets are mirrored on Hugging Face to ensure redundancy and faster global access.[21]

Public accessibility is a design principle. Automated GitHub Actions update metadata badges and commit summaries whenever new data are ingested. Users can clone, fork, or download any subset directly without authentication.[10]

Licensing notices are embedded in every dataset directory, clarifying reuse rights and responsibilities. The datasets intentionally avoid any personally identifiable information or restricted content to uphold ethical data-sharing standards.[22]

The open license facilitates reproducible science and supports downstream applications in natural language processing, digital governance, and policy research. By ensuring public access, the project aligns with FAIR principles—Findable, Accessible, Interoperable, and Reusable.[23]

This project establishes an open, reproducible, and scalable foundation for Sri Lankan document datasets, spanning legal, governmental, and media sources. The pipeline integrates crawling, parsing, and versioning into a unified ecosystem for data-driven research.[19]

The datasets already serve as a valuable resource for natural language processing, computational law, and policy analysis. They enable quantitative and qualitative insights into governance, lawmaking, and civic communication over time.[23]

Future development focuses on three priorities:

First, expanding coverage by adding new datasets from additional government agencies, media sources, and historical archives.[24]

Second, improving the linguistic accuracy of Sinhala and Tamil parsing, particularly for complex sentence structures and transliterated terms. Enhancements in tokenization, font handling, and multilingual embeddings are planned.[25]

Third, integrating OCR parsing for PDFs with unstructured or scanned content. We are experimenting with deep-learning-based OCR pipelines that combine layout recognition and language modeling to recover high-quality text from low-quality sources.[26]

Together, these directions will further improve coverage, data quality, and usability, ensuring that the Sri Lanka Datasets initiative remains a sustainable open infrastructure for the region's digital and academic ecosystem.[27]

References

- [1] Common crawl: Open web corpus. <https://commoncrawl.org>, 2020.
- [2] Wikipedia dumps. <https://dumps.wikimedia.org>, 2018.
- [3] Aaron Gokaslan and Vanya Cohen. Open-webtext corpus. <https://skylion007.github.io/OpenWebTextCorpus>, 2019.
- [4] Indian kanoon legal corpus. <https://indiankanoon.org>, 2018.
- [5] Opensubtitles: Multilingual parallel corpus. <https://opus.nlpl.eu/OpenSubtitles.php>, 2016.
- [6] African news corpus. <https://data.africanlp.org>, 2021.
- [7] Sltalk: Sinhala and tamil social media dataset. <https://github.com/sltalk>, 2023.
- [8] Srilankanlp: Sinhala and tamil language resources. <https://github.com/SriLankaNLP>, 2022.
- [9] Omar Ashraf et al. A survey on mlops: Building reliable machine learning systems. *Journal of Systems and Software*, 2022. doi: 10.1016/j.jss.2022.111223.
- [10] GitHub, Inc. Github actions documentation. <https://docs.github.com/actions>, 2023.
- [11] Victoria Stodden, Jennifer Seiler, and Zhaokun Ma. An empirical analysis of journal policy for reproducibility. *Proceedings of the National Academy of Sciences*, 2017. doi: 10.1073/pnas.1708290114.
- [12] Selenium Project. Selenium webdriver documentation. <https://www.selenium.dev/documentation/>, 2023.
- [13] Gautam Pant and Padmini Srinivasan. Web crawling best practices: Ethics, efficiency, and effectiveness. In *Proceedings of the Web Conference*, 2021.
- [14] Nathan Blaustein et al. Data versioning in practice: Challenges and solutions. *Data Science Journal*, 2020. doi: 10.5334/dsj-2020-012.
- [15] Artifex Software. Pymupdf documentation. <https://pymupdf.readthedocs.io>, 2024.
- [16] Chen Li et al. Document layout analysis via geometric and semantic segmentation. *Pattern Recognition Letters*, 2021. doi: 10.1016/j.patrec.2021.07.014.
- [17] Leo Pipino et al. Data quality assessment frameworks in data engineering. *Information Systems Frontiers*, 2022. doi: 10.1007/s10796-022-10314-2.
- [18] Eytan Ben-David et al. Archival web data collection and preservation at scale. In *Proceedings of the Web Archiving Conference*, 2019.
- [19] Marijn Janssen et al. Open data practices and governance: A comprehensive framework. *Government Information Quarterly*, 2020. doi: 10.1016/j.giq.2020.101458.
- [20] Massachusetts Institute of Technology. The mit license (mit). <https://opensource.org/licenses/MIT>, 1988.
- [21] Quentin Lhoest et al. Datasets: A community library for natural language processing. *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, 2021. doi: 10.18653/v1/2021.emnlp-demo.21.

- [22] Luciano Floridi and Mariarosaria Taddeo. What is data ethics? *Philosophical Transactions of the Royal Society A*, 2019. doi: 10.1098/rsta.2018.0083.
- [23] Mark D. Wilkinson et al. The fair guiding principles for scientific data management and stewardship. *Scientific Data*, 2016. doi: 10.1038/sdata.2016.18.
- [24] Arjun Mishra et al. Scaling open data: Strategies for continuous expansion and maintenance. *Data Intelligence*, 2023. doi: 10.1162/dint_a_00123.
- [25] Alexis Conneau et al. Unsupervised cross-lingual representation learning at scale. *Transactions of the Association for Computational Linguistics*, 2022. doi: 10.1162/tacl_a_00424.
- [26] Benjamin Davis et al. Deep ocr: Neural architectures for document text recognition and layout parsing. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. doi: 10.1109/TPAMI.2021.3056728.
- [27] Anneke Zuiderwijk and Marijn Janssen. Sustainability of open data and data infrastructures: A research agenda. *Government Information Quarterly*, 2020. doi: 10.1016/j.giq.2020.101491.