# Spacerini: Plug-and-play Search Engines with Pyserini and Hugging Face

Christopher Akiki
Leipzig University and ScaDS.AI

Odunayo Ogundepo
University of Waterloo

Aleksandra Piktus
Hugging Face and Sapienza University

Xinyu Zhang
University of Waterloo

Akintunde Oladipo
University of Waterloo

Jimmy Lin
University of Waterloo

Martin Potthast
Leipzig University and ScaDS.AI

## ABSTRACT

We present Spacerini, a modular framework for seamless building and deployment of interactive search applications, designed to facilitate the qualitative analysis of large scale research datasets. Spacerini integrates features from both the Pyserini toolkit and the Hugging Face ecosystem to ease the indexing text collections and deploy them as search engines for ad-hoc exploration and to make the retrieval of relevant data points quick and efficient. The user-friendly interface enables searching through massive datasets in a no-code fashion, making Spacerini broadly accessible to anyone looking to qualitatively audit their text collections. This is useful both to IR researchers aiming to demonstrate the capabilities of their indexes in a simple and interactive way, and to NLP researchers looking to better understand and audit the failure modes of large language models. The framework is open source and available on GitHub: https://github.com/castorini/hf-spacerini, and includes utilities to load, pre-process, index, and deploy local and web search applications. A portfolio of applications created with Spacerini for a multitude of use cases can be found by visiting https://hf.co/spacerini.

## CCS CONCEPTS

• **Information systems → Search interfaces**.

## KEYWORDS

information retrieval, search interface, text auditing

## 1 INTRODUCTION

The commodification of data has had a transformative effect on science in general [11] and machine learning in particular [21]. The race to train ever-larger language models is so contingent upon having access to immeasurable quantities of text [12] that datasets have become—as Bender et al. [5] contend—"too large to document" both technically [5] and methodologically [14]. This often leads to a "train first, ask questions later" approach to model training [2], itself a concrete instance of convenience experimentation [16], an approach to research that is modulated by the availability of a resource and ease of a method, rather than the suitability thereof to the problem at hand. In this sense, Beaulieu and Leonelli [4] deem it important to differentiate between the availability of data and its suitability, especially in light of the misconception that Web data is

representative of all human experiences and immune to the ever-widening *digital divide* [18]. They point out that the divide not only exists but actively limits this representativeness—both socio-economic and socio-cultural—of the Web, which in turn reinforces biases of the artifacts that leverage it [5].

Being unable to easily audit large datasets incentivizes researchers to release models trained on data they do not truly understand [21] leading to model behaviors that are hard to study, predict or trace [3, 21, 29]. This is especially problematic in light of the potential real-world harms that ensue [8, 9, 13, 31]. Being able to properly understand the limitations of our datasets is a necessary first step toward understanding the harmful biases and failure modes of the artifacts that build upon them. Understanding the training data is therefore a critical step in the process of auditing large language models [23].

It is from this vantage point that we initially developed Spacerini as an open-source tool for the quick indexing and deployment of shareable search engines, but have also since come to realize its potential in being useful for an even wider audience interested in making their text artifacts searchable. This includes IR students, Digital Humanists, Shared Tasks organizers, and digital investigative journalists. We cover these use cases in more detail in Section 4. Spacerini was designed as a modular framework that facilitates the indexing, creation, and free hosting of graphical search interfaces for text datasets. It helps streamline the process of auditing large datasets by allowing users to effortlessly index their text collections and deploy them as search systems that can be accessed by anyone. It does so by standing on the shoulders of battle-tested open-source libraries from the Castorini [20] and Hugging Face [1, 19] ecosystems and enhancing the interoperability between them to enable quick indexing and free deployment of search interfaces.

The key advantage of Spacerini is its ability to simplify the search process and retrieve relevant data points from massive datasets with ease, allowing researchers to conduct quick and efficient audits, while abstracting away all the minutiae of indexing data or hosting services. We believe that this provides an opportunity for collaboration and transparency in IR and NLP research. With the creation and sharing of search indexes publicly, practitioners, researchers and the general public can work together to pinpoint problematic content, find duplicates, and identify biases in datasets. Specific to the IR community, Spacerini allows a convenient and universal approach to index sharing, which greatly facilitates the demonstration of retrieval systems and the reproduction from other researcher teams.

## 2 BACKGROUND

Large scale, predominantly web-mined text datasets have been proliferating in NLP recently, giving rise to publications [10, 17, 25, 27] which often contain interesting analyses of the specific datasets being presented, however, usually lack any comparison to existing resources beyond basic metrics such as sizes of the datasets or languages they contain.

As discussed in Section 1, in the face of an increased scrutiny of the models trained on datasets in question, the topic of data understanding and governance has been gaining more traction, being accepted as an important part of research. Efforts such as those of Mitchell et al. [21] contribute frameworks for more standardised and reproducible metrics and measurements of datasets, and we position ourselves as a complementary continuation of their work, focusing on a more curatorial and qualitative assessment that might not readily fit under the umbrella of "measurements". We therefore aim to fill the gap in the evaluation landscape by facilitating qualitative, rather than quantitative analysis of large scale data collections.

Similarly to the authors of Gradio [1], a Python package for fast development of Machine Learning demos, we believe that the accessibility of data and model analysis tools is crucial to building both the understanding of and the trust in the underlying resources. The potential of relevance-based interfaces to massive textual corpora, the creation of which can be facilitated by leveraging toolkits such as Pyserini [20], has previously been tapped into by the researchers at the Allen Institute of AI who propose a C4 [27] search engine[1]. Similar interfaces have also been found useful in more specialised domains, e.g. in COVID-related datasets [33], news quotes [30], or medical literature [24]. However, while these solutions are undeniably useful, they remain very contextual: dataset-specific, and project-specific. We believe Spacerini to be the first generalizable tool which proposes an end-to-end pipeline automating the route from raw text to qualitative analysis.

## 3 SPACERINI

Spacerini is a modular framework that integrates Pyserini with the Hugging Face ecosystem to streamline the process of going from any Hugging Face text dataset—either local or hosted on the Hugging Face Hub—to a search interface driven by a Pyserini index that can be deployed for free on the Hugging Face Hub.
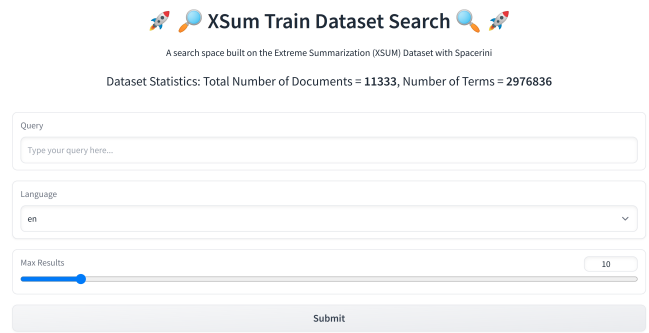
In what follows, we deconstruct an example script[2] to showcase the different features enabled by Spacerini. When run end-to-end, the script pulls a dataset from the Hugging Face hub, pre-processes it, indexes it, creates a gradio-based search interface and pushes that to the Web as a Hugging Face Spaces demo. This is only meant as a feature-complete demo, and we don't expect most people to want to integrate every step into their workflows, but rather to cherry-pick and decide what best to use depending on context.

### 3.1 Loading Data

All our workflows are backed by the Hugging Face `datasets` library [19], itself based on the extremely efficient Apache Arrow

---

[1]https://c4-search.apps.allenai.org/

[2]gradio-demo.py

---



**Figure 1: Search interface for the Extreme Summarization (xsum) dataset deployed on Hugging Face `Spaces` using Spacerini.**

---

format. `Datasets` is a mature library which provides a standardized interface to any tabular dataset, in particular, to tens of thousands of community datasets hosted on the Hugging Face Hub[3]. The `datasets` library gives fine-grained control over the lifecycle of tabular datasets, which we choose to abstract away through a set of opinionated data loading functions that cover the use cases we deem relevant to information retrieval. We also add new functionality, such as the ability to load any document dataset from the `ir_datasets` library using for example the following one-liner to load MS MARCO as a Hugging Face `datasets.Dataset` object using a function from the `data` subpackage:

```python
from spacerini.data import load_ir_dataset
hf_dset = load_ir_dataset("msmarco-passage")
```

We include wrappers to load database tables, pandas DataFrames [26, 32], and text datasets on disk, as well as the ability to load any dataset either in memory-mapped mode or in streaming mode: the former makes it possible to handle larger-than-memory datasets, and the latter larger-than-disk datasets that can be streamed from a remote location such as the Hugging Face Hub.

### 3.2 Pre-processing

Spacerini also provices a `preprocess` subpackage which offers a range of customizable pre-processing options for preparing datasets. This module includes a sharding utility that enables the partitioning of large datasets into smaller, more manageable chunks for efficient parallel processing.

```python
from spacerini.preprocesss import shard_dataset
shard_dataset(
    hf_dataset=hf_dset,
    shard_size="1GB",
    column_to_index="text",
    shards_paths="msmarco-shards",
)
```

---

[3]https://hf.co/datasets

## 3.3 Indexing

Spacerini's `index` subpackage leverages Pyserini to provide very efficient Lucene indexing and allow users to easily and quickly index large datasets, either sharded in the pre-processing step, or any text format accepted by Pyserini, and streaming text datasets, such as those returned by Spacerini's `data` subpackage. This subpackage also exposes several tokenization options using existing language-specific analyzers [4] as well as Hugging Face subword tokenizers [22].

```python
from spacerini.index import index_json_shards
index_json_shards(
    shards_path="msmarco-shards",
    index_path="app/index",
)
```

## 3.4 Template-based Search Interfaces

Having indexed a collection, one can easily spin up a frontend using the `frontend` subpackage and one of many provided templates[5]. These are built using `cookiecutter`,[6] a Python templating library for software projects. We provide a few batteries-included frontend templates based both the Gradio[7] and Streamlit [8] demo app frameworks, both of which are natively supported by Hugging Face `Spaces`. Figure 1 showcases a search engine built using one of our Streamlit templates.

```python
from spacerini.frontend import create_app
cookiecutter_vars = {"dset_text_field": "text,
                "metadata_field": "docid",
                "space_title": "MS MARCO Search",
                "local_app": "app"}
create_app(
    template="gradio-vanilla",
    extra_context_dict=cookiecutter_vars,
    output_dir=".",
)
```

## 3.5 Deployment to Hugging Face Spaces

The local apps developed in the previous subsection can then be pushed to Hugging Face Spaces and hosted there for free. One can then further customize the running app from the browser, for example to add functionality not provided by the chosen template.

```python
from spacerini.frontend import create_space_from_local
create_space_from_local(
    space_slug="msmarco-passage-search,
    organization="spacerini,
    space_sdk="gradio",
    local_dir=LOCAL_APP,
    delete_after_push=True,
)
```

## 3.6 Sharing Indexes as Hugging Face Datasets

Orthogonal to the workflow presented so far, is the ability to upload Lucene indexes to the Hugging Face Hub using shareable dataset repositories and enabling reproducible retrieval experiments.

```python
from spacerini.index import push_index_to_hub

push_index_to_hub(
    dataset_slug="lucene-english-analyzer-msmarco",
    index_path="index",
)
```

Any index can then just as easily be downloaded for local use:

```python
from spacerini.index import load_index_from_hub

index_path = load_index_from_hub("lucene-fr-analyzer-")
```

## 3.7 Search and Pagination

Search features are provided by the `search` subpackage and leverage the memory-mapping feature of Arrow tables to load the entire table of results—no matter how big—only materializing the specific shard that corresponds to the requested result page.

```python
from spacerini.search import result_indices, result_page

ix = result_indices(
    "Lorem Ipsum",
    num_results=1_000,
    INDEX_PATH,
)

last_results_page = result_page(
    hf_dset,
    ix,
    page=-1,
    results_per_page=20,
).to_pandas()
```

# 4 USE CASES AND EXAMPLES

Spacerini is designed to enable qualitative analysis of large-scale textual corpora without the need for extensive engineering work. The tool can be used in dataset auditing campaigns, such as those carried out by Kreutzer et al. [15] or in data annotation efforts. It may be applied to inspect failures of large scale language model predictions and find potential sources of memorized generations. AI Ethics researchers can employ the tool to find evidence supporting their hypotheses about the content of the models' training datasets.

Given its tight integration with Pyserini, Spacerini can also be leveraged by IR researchers to experiment with modifications of their retrieval pipelines in user studies (e.g. by exposing retrieval options such as BM25 and RM3 parameters in the the user interface) or to deploy demos of their working prototypes. Reproducibility

---

[4] https://lucene.apache.org/core/9_5_0/analysis/common/index.html
[5] https://github.com/castorini/hf-spacerini/tree/main/templates
[6] https://github.com/cookiecutter/cookiecutter
[7] https://gradio.app/
[8] https://streamlit.io/

for IR experiments is further enhanced thanks to the index sharing abilities introduced in Section 3.6.

Spacerini can also be leveraged by Digital Humanists, Archivists and Librarians looking to index their collections. Indeed, GLAM (Galleries, libraries, archives, and museums) collections are increasingly being made available as datasets. Furthermore, there is a growing interest in the digital humanities in training and using languages models, as demonstrated by the success of projects such as the *BERT for Humanists Project*[9]. This makes it especially pressing to have easy access ways to critically examine the data that goes into such models. In the context of the digital and computational humanities, indexing data relevant to these efforts is often not an easy task, and often project-based and contingent upon precarious funding arrangements. Having a project-agnostic tool like Spacerini could prove valuable to this community and a useful addition to toolkits such as the GLAM Workbench [28].

Given its low engineering barrier of entry, Spacerini can be a good addition to IR courses with a practical component, where students are tasked with developing search engines, by providing an easy-to-deploy interface for their developed retrieval systems.

Spacerini can also be leveraged by organizers of shared tasks such as MIRACL [34] and Touché [6], who want to help participants explore the datasets without forcing them to download large volumes of data. It can also be used as a platform for participants to deploy working prototypes with a unified interface.

Spacerini can help data journalists and digital investigative journalists index, explore, and understand open data, in a similar vein to the functionality provided by the Aleph suite.[10] Providing technical tools to data journalists is a crucial in uncovering matters of public interest, as was evident by role played by the collaborative use of the Neo4j graph database in unraveling the corrupt network surrounding tax havens [7].

Finally, three features of Hugging Face Spaces make them especially attractive for users: (1) they can leverage private datasets, meaning that one can provide search access to a dataset without sharing the underlying data,(2) they can be seamlessly embedded into HTML, specifically Gradio-based Spaces which can be embedded as *Web Components*[11] so that users can easily integrate a Spacerini-based search feature into their own sites[12], and (3) Gradio-based Spaces expose a FastAPI[13] endpoint that can be queried to access the functionality of the space, making deployed search engines accessible through HTTP calls.

## 5 LIMITATIONS AND FUTURE PLANS

The main limitation of the off-the-shelf variant of Spacerini is the disk space limit imposed by Hugging Face Spaces, which is currently set to 50 GB.[14] While not enough to accommodate entire corpora such as ROOTS or The Pile, such large datasets are most often amalgamations of constitutents datasets which can each be studied independently. This limit has no bearing on Spacerini search apps deployed locally. Should users still want to get more disk space for

---

[9] https://www.bertforhumanists.org/
[10] https://docs.alephdata.org/
[11] https://developer.mozilla.org/en-US/docs/Web/Web_Components
[12] For example: https://cakiki.github.io/search-engine/
[13] https://github.com/tiangolo/fastapi
[14] https://huggingface.co/docs/hub/spaces-overview#hardware-resources

their Spaces-hosted indexes, they are welcome to apply for community grants offered by Hugging Face.

Planned improvements to the library include automating the creation of dataset cards (Or rather "index cards") when pushing an index to the Hugging Face Hub, more support for dense indexing features provided by Pyserini, as well as more fine-grained support of document tokenization. We also look forward to community contributions both to the codebase and to the frontend templates.

## 6 CONCLUSION

We presented Spacerini, a modular framework that enables the quick and free deployment and serving of template-based search indexes as interactive applications. The need for such a tool is especially pressing as large language models have come to consume inordinate amounts of text data, reinforcing the need for a qualitative exploration and understanding of datasets to assess them in a way that is impenetrable to quantitative analyses alone.

Spacerini leverages features from both the `Pyserini` toolkit and the Hugging Face ecosystem to facilitate the creation and hosting of user-friendly search systems for text datasets. Users can easily index their collections and deploy them as ad-hoc search interfaces, making the retrieval of relevant data points a quick and efficient process. The user-friendly interface enables non-technical users to effectively search massive datasets, making Spacerini a valuable tool for anyone looking to audit their text collections qualitatively. The framework is open-source and available on GitHub: https://github.com/castorini/hf-spacerini.

Finally, we emphasize that Spacerini is a first step in the direction of systematic dataset auditing, and more work is still needed to create standardized structures that leverage tools such as Spacerini to properly document the different axes of interest that are appropriate for a given usage context.

## REFERENCES

[1] Abubakar Abid, Ali Abdalla, Ali Abid, Dawood Khan, Abdulrahman Alfozan, and James Zou. 2019. Gradio: Hassle-Free Sharing and Testing of ML Models in the Wild. *arXiv preprint arXiv:1906.02569* (2019).

[2] Christopher Akiki, Giada Pistilli, Margot Mieskes, Matthias Gallé, Thomas Wolf, Suzana Ilic, and Yacine Jernite. 2022. BigScience: A Case Study in the Social Construction of a Multilingual Large Language Model. *CoRR* abs/2212.04960 (2022). https://doi.org/10.48550/arXiv.2212.04960 arXiv:2212.04960

[3] Ekin Akyurek, Tolga Bolukbasi, Frederick Liu, Binbin Xiong, Ian Tenney, Jacob Andreas, and Kelvin Guu. 2022. Towards Tracing Knowledge in Language Models Back to the Training Data. In *Findings of the Association for Computational Linguistics: EMNLP 2022*. Association for Computational Linguistics, Abu Dhabi, United Arab Emirates, 2429–2446. https://aclanthology.org/2022.findings-emnlp.180

[4] Anne Beaulieu and Sabina Leonelli. 2021. *Data and Society: A Critical Introduction*. Sage.

[5] Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency* (Virtual Event, Canada) *(FAccT '21)*. Association for Computing Machinery, New York, NY, USA, 610–623. https://doi.org/10.1145/3442188.3445922

[6] Alexander Bondarenko, Maik Fröbe, Johannes Kiesel, Shahbaz Syed, Timon Gurcke, Meriem Beloucif, Alexander Panchenko, Chris Biemann, Benno Stein, Henning Wachsmuth, Martin Potthast, and Matthias Hagen. 2022. Overview of Touché 2022: Argument Retrieval. In *Experimental IR Meets Multilinguality, Multimodality, and Interaction - 13th International Conference of the CLEF Association, CLEF 2022, Bologna, Italy, September 5-8, 2022, Proceedings (Lecture Notes in Computer Science, Vol. 13390)*, Alberto Barrón-Cedeño, Giovanni Da San Martino, Mirko Degli Esposti, Fabrizio Sebastiani, Craig Macdonald, Gabriella Pasi, Allan Hanbury, Martin Potthast, Guglielmo Faggioli, and Nicola Ferro (Eds.). Springer, 311–336. https://doi.org/10.1007/978-3-031-13643-6_21

[7] Emilia Díaz-Struck and Mar Cabra. 2018. *Uncovering International Stories with Data and Collaboration*. Springer International Publishing, Cham, 55–65. https://doi.org/10.1007/978-3-319-97283-1_6

[8] Ethan Fast, Tina Vachovsky, and Michael S. Bernstein. 2016. Shirtless and Dangerous: Quantifying Linguistic Signals of Gender Bias in an Online Fiction Writing Community. *CoRR* abs/1603.08832 (2016). arXiv:1603.08832 http://arxiv.org/abs/1603.08832

[9] Antigoni-Maria Founta, Constantinos Djouvas, Despoina Chatzakou, Ilias Leontiadis, Jeremy Blackburn, Gianluca Stringhini, Athena Vakali, Michael Sirivianos, and Nicolas Kourtellis. 2018. Large Scale Crowdsourcing and Characterization of Twitter Abusive Behavior. *CoRR* abs/1802.00393 (2018). arXiv:1802.00393 http://arxiv.org/abs/1802.00393

[10] Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. 2021. The Pile: An 800GB Dataset of Diverse Text for Language Modeling. https://doi.org/10.48550/ARXIV.2101.00027

[11] Tony Hey, Stewart Tansley, Kristin Tolle, and Jim Gray. 2009. *The Fourth Paradigm: Data-Intensive Scientific Discovery*. Microsoft Research. https://www.microsoft.com/en-us/research/publication/fourth-paradigm-data-intensive-scientific-discovery/

[12] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. 2022. Training Compute-Optimal Large Language Models. *CoRR* abs/2203.15556 (2022). https://doi.org/10.48550/arXiv.2203.15556 arXiv:2203.15556

[13] Ben Hutchinson, Vinodkumar Prabhakaran, Emily Denton, Kellie Webster, Yu Zhong, and Stephen Denuyl. 2020. Social Biases in NLP Models as Barriers for Persons with Disabilities. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 5491–5501. https://doi.org/10.18653/v1/2020.acl-main.487

[14] Eun Seo Jo and Timnit Gebru. 2020. Lessons from Archives: Strategies for Collecting Sociocultural Data in Machine Learning. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency* (Barcelona, Spain) *(FAT* '20)*. Association for Computing Machinery, New York, NY, USA, 306–316. https://doi.org/10.1145/3351095.3372829

[15] Julia Kreutzer, Isaac Caswell, Lisa Wang, Ahsan Wahab, Daan van Esch, Nasanbayar Ulzii-Orshikh, Allahsera Tapo, Nishant Subramani, Artem Sokolov, Claytone Sikasote, Monang Setyawan, Supheakmungkol Sarin, Sokhar Samb, Benoît Sagot, Clara Rivera, Annette Rios, Isabel Papadimitriou, Salomey Osei, Pedro Ortiz Suarez, Iroro Orife, Kelechi Ogueji, Andre Niyongabo Rubungo, Toan Q. Nguyen, Mathias Müller, André Müller, Shamsuddeen Hassan Muhammad, Nanda Muhammad, Ayanda Mnyakeni, Jamshidbek Mirzakhalov, Tapiwanashe Matangira, Colin Leong, Nze Lawson, Sneha Kudugunta, Yacine Jernite, Mathias Jenny, Orhan Firat, Bonaventure F. P. Dossou, Sakhile Dlamini, Nisansa de Silva, Sakine Çabuk Ballı, Stella Biderman, Alessia Battisti, Ahmed Baruwa, Ankur Bapna, Pallavi Baljekar, Israel Abebe Azime, Ayodele Awokoya, Duygu Ataman, Orevaoghene Ahia, Oghenefego Ahia, Sweta Agrawal, and Mofetoluwa Adeyemi. 2022. Quality at a Glance: An Audit of Web-Crawled Multilingual Datasets. *Transactions of the Association for Computational Linguistics* 10 (2022), 50–72. https://doi.org/10.1162/tacl_a_00447

[16] Ulrich Krohs. 2012. Convenience experimentation. *Studies in History and Philosophy of Science Part C: Studies in History and Philosophy of Biological and Biomedical Sciences* 43, 1 (2012), 52–57. https://doi.org/10.1016/j.shpsc.2011.10.005 Data-Driven Research in the Biological and Biomedical Sciences On Nature and Normativity: Normativity, Teleology, and Mechanism in Biological Explanation.

[17] Hugo Laurençon, Lucile Saulnier, Thomas Wang, Christopher Akiki, Albert Villanova del Moral, Teven Le Scao, Leandro Von Werra, Chenghao Mou, Eduardo González Ponferrada, Huu Nguyen, Jörg Frohberg, Mario Šaško, Quentin Lhoest, Angelina McMillan-Major, Gérard Dupont, Stella Biderman, Anna Rogers, Loubna Ben allal, Francesco De Toni, Giada Pistilli, Olivier Nguyen, Somaieh Nikpoor, Maraim Masoud, Pierre Colombo, Javier de la Rosa, Paulo Villegas, Tristan Thrush, Shayne Longpre, Sebastian Nagel, Leon Weber, Manuel Romero Muñoz, Jian Zhu, Daniel Van Strien, Zaid Alyafeai, Khalid

Almubarak, Vu Minh Chien, Itziar Gonzalez-Dios, Aitor Soroa, Kyle Lo, Manan Dey, Pedro Ortiz Suarez, Aaron Gokaslan, Shamik Bose, David Ifeoluwa Adelani, Long Phan, Hieu Tran, Ian Yu, Suhas Pai, Jenny Chim, Violette Lepercq, Suzana Ilic, Margaret Mitchell, Sasha Luccioni, and Yacine Jernite. 2022. The BigScience ROOTS Corpus: A 1.6TB Composite Multilingual Dataset. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*. https://openreview.net/forum?id=UoEw6KigkUn

[18] Sabina Leonelli. 2020. Scientific Research and Big Data. In *The Stanford Encyclopedia of Philosophy* (Summer 2020 ed.), Edward N. Zalta (Ed.). Metaphysics Research Lab, Stanford University.

[19] Quentin Lhoest, Albert Villanova del Moral, Yacine Jernite, Abhishek Thakur, Patrick von Platen, Suraj Patil, Julien Chaumond, Mariama Drame, Julien Plu, Lewis Tunstall, Joe Davison, Mario Sasko, Gunjan Chhablani, Bhavitvya Malik, Simon Brandeis, Teven Le Scao, Victor Sanh, Canwen Xu, Nicolas Patry, Angelina McMillan-Major, Philipp Schmid, Sylvain Gugger, Clément Delangue, Théo Matussière, Lysandre Debut, Stas Bekman, Pierric Cistac, Thibault Goehringer, Victor Mustar, François Lagunas, Alexander M. Rush, and Thomas Wolf. 2021. Datasets: A Community Library for Natural Language Processing. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, EMNLP 2021, Online and Punta Cana, Dominican Republic, 7-11 November, 2021*, Heike Adel and Shuming Shi (Eds.). Association for Computational Linguistics, 175–184. https://doi.org/10.18653/v1/2021.emnlp-demo.21

[20] Jimmy Lin, Xueguang Ma, Sheng-Chieh Lin, Jheng-Hong Yang, Ronak Pradeep, and Rodrigo Nogueira. 2021. Pyserini: A Python Toolkit for Reproducible Information Retrieval Research with Sparse and Dense Representations. In *Proceedings of the 44th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2021)*. 2356–2362.

[21] Margaret Mitchell, Alexandra Sasha Luccioni, Nathan Lambert, Marissa Gerchick, Angelina McMillan-Major, Ezinwanne Ozoani, Nazneen Rajani, Tristan Thrush, Yacine Jernite, and Douwe Kiela. 2022. Measuring Data. *CoRR* abs/2212.05129 (2022). https://doi.org/10.48550/arXiv.2212.05129 arXiv:2212.05129

[22] Anthony MOI, Nicolas Patry, Pierric Cistac, Pete, Funtowicz Morgan, Sebastian Pütz, Mishig, Bjarte Johansen, Thomas Wolf, Sylvain Gugger, Clement, Julien Chaumond, Lysandre Debut, François Garillot, Luc Georges, dctelus, JC Louis, MarcusGrass, Taufiquzzaman Peyash, 0xflotus, Alan deLevie, Alexander Mamaev, Arthur, Cameron, Colin Clement, Dagmawi Moges, David Hewitt, Denis Zolotukhin, and Geoffrey Thomas. 2022. *huggingface/tokenizers: Rust 0.13.2*. https://doi.org/10.5281/zenodo.7298413

[23] Jakob Mökander, Jonas Schuett, Hannah Rose Kirk, and Luciano Floridi. 2023. Auditing large language models: a three-layered approach. https://doi.org/10.48550/ARXIV.2302.08500

[24] Danna Niezni, Hillel Taub-Tabib, Yuval Harris, Hagit Sason-Bauer, Yakir Amrusi, Dana Azagury, Maytal Avrashami, Shaked Launer-Wachs, Jon Borchardt, M Kusold, Aryeh Tiktinsky, Tom Hope, Yoav Goldberg, and Yosi Shamay. 2022. Extending the Boundaries of Cancer Therapeutic Complexity with Literature Data Mining. *bioRxiv* (2022). https://doi.org/10.1101/2022.05.03.490286 arXiv:https://www.biorxiv.org/content/early/2022/10/30/2022.05.03.490286.full.pdf

[25] Pedro Javier Ortiz Suárez, Benoît Sagot, and Laurent Romary. 2019. 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)*, Piotr Bański, Adrien Barbaresi, Hanno Biber, Evelyn Breiteneder, Simon Clematide, Marc Kupietz, Harald Lüngen, and Caroline Iliadi (Eds.). Leibniz-Institut für Deutsche Sprache, Cardiff, United Kingdom. https://doi.org/10.14618/IDS-PUB-9021

[26] The pandas development team. 2020. *pandas-dev/pandas: Pandas*. https://doi.org/10.5281/zenodo.3509134

[27] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2022. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *J. Mach. Learn. Res.* 21, 1, Article 140 (jun 2022), 67 pages.

[28] Tim Sherratt. 2021. *GLAM Workbench*. https://doi.org/10.5281/zenodo.5603060

[29] Shoaib Ahmed Siddiqui, Nitarshan Rajkumar, Tegan Maharaj, David Krueger, and Sara Hooker. 2022. Metadata Archaeology: Unearthing Data Subsets by Leveraging Training Dynamics. *CoRR* abs/2209.10015 (2022). https://doi.org/10.48550/ARXIV.2209.10015 arXiv:2209.10015

[30] Vuk Vuković, Akhil Arora, Huan-Cheng Chang, Andreas Spitz, and Robert West. 2022. Quote Erat Demonstrandum: A Web Interface for Exploring the Quotebank Corpus. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM. https://doi.org/10.1145/3477495.3531696

[31] Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, Zac Kenton, Sasha Brown, Will Hawkins, Tom Stepleton, Courtney Biles, Abeba

Birhane, Julia Haas, Laura Rimell, Lisa Anne Hendricks, William S. Isaac, Sean Legassick, Geoffrey Irving, and Iason Gabriel. 2021. Ethical and social risks of harm from Language Models. *CoRR* abs/2112.04359 (2021). arXiv:2112.04359 https://arxiv.org/abs/2112.04359

[32] Wes McKinney. 2010. Data Structures for Statistical Computing in Python. In *Proceedings of the 9th Python in Science Conference*, Stéfan van der Walt and Jarrod Millman (Eds.). 56 – 61. https://doi.org/10.25080/Majora-92bf1922-00a

[33] Edwin Zhang, Nikhil Gupta, Raphael Tang, Xiao Han, Ronak Pradeep, Kuang Lu, Yue Zhang, Rodrigo Nogueira, Kyunghyun Cho, Hui Fang, and Jimmy Lin. 2020. Covidex: Neural Ranking Models and Keyword Search Infrastructure for the COVID-19 Open Research Dataset. In *Proceedings of the First Workshop on Scholarly Document Processing*. Association for Computational Linguistics, Online, 31–41. https://doi.org/10.18653/v1/2020.sdp-1.5

[34] Xinyu Zhang, Nandan Thakur, Odunayo Ogundepo, Ehsan Kamalloo, David Alfonso-Hermelo, Xiaoguang Li, Qun Liu, Mehdi Rezagholizadeh, and Jimmy Lin. 2022. Making a MIRACL: Multilingual Information Retrieval Across a Continuum of Languages. https://doi.org/10.48550/ARXIV.2210.09984