

Open Datasheets: Machine-readable Documentation for Open Datasets and Responsible AI Assessments

Anthony Cintron Roman^{1,*}, Jennifer Wortman Vaughan², Valerie See³, Steph Ballard⁴, Kevin Xu^{5,+}, Jehu Torres^{1,+}, Caleb Robinson^{1,+}, and Juan M. Lavista Ferres¹

¹Microsoft, AI For Good Research Lab, USA

²Microsoft, Research, USA

³Microsoft, Office of Responsible AI, USA

⁴Microsoft, Open Innovation, USA

⁵GitHub, USA

*anthony.cintron@microsoft.com

+these authors contributed equally to this work

ABSTRACT

This paper introduces a no-code, machine-readable documentation framework for open datasets, with a focus on responsible AI (RAI) considerations. The framework aims to improve comprehensibility, and usability of open datasets, facilitating easier discovery and use, better understanding of content and context, and evaluation of dataset quality and accuracy. The proposed framework is designed to streamline the evaluation of datasets, helping researchers, data scientists, and other open data users quickly identify datasets that meet their needs and organizational policies or regulations. The paper also discusses the implementation of the framework and provides recommendations to maximize its potential. The framework is expected to enhance the quality and reliability of data used in research and decision-making, fostering the development of more responsible and trustworthy AI systems.

Introduction

Machine-readable documentation, which refers to documentation that is structured in a format easily readable and processed by computers, plays a vital role in improving the accessibility, comprehensibility, and usability of open datasets. It has the potential to enable easier discovery and use of datasets, facilitate a better understanding of the content and context of a dataset, simplify the integration of datasets from various sources, and help with the evaluation of the quality and accuracy of datasets¹. With the exponential growth of open data, machine-readable metadata is becoming increasingly important for making this data easily discoverable and useable².

As open data is progressively employed in AI applications^{3,4}, it is imperative for researchers to consider the potential implications of these open datasets from a responsible AI (RAI) perspective¹. It is widely recognized that AI systems can inherit biases and other flaws from the data they are trained on^{1,5}, and with AI continuing to permeate various aspects of our daily lives, the transparency, reliability, trustworthiness, and fairness of AI systems becomes paramount. This begins with the quality of the data they are trained on⁶. Understanding the data being used in an AI system is crucial to uphold responsible AI implementations, including its source, the methodologies used in its collection, and other relevant factors discussed in this paper.

Evaluating a dataset involves multiple distinct goals or aspects. One is to determine if a dataset complies with organizational policies and/or regulations, such as ensuring appropriate consent was obtained. The other is to ascertain if a dataset is suitable for the task at hand, such as whether the labels are valid proxies for the real-world phenomena a model is being designed to predict. While researchers can often evaluate an open dataset and make determinations about some of its characteristics, such as demographic distribution, as part of their exploratory data analysis (EDA), there are properties that cannot be learned through data exploration, such as exactly how demographic information was obtained from data subjects or whether and how their consent was obtained. This makes searching for and evaluating datasets an overwhelming, if not impossible, task. It can lead to investing a significant amount of time in evaluating a dataset only to realize later that it was not collected properly and/or does not have a proper representation from an RAI perspective to meet their organizational policies, among other considerations. Such mistakes can be costly, with potential harm to individuals, and significant losses for organizations⁷⁻¹³.

Consider a researcher working on a project related to speech recognition technology. They need a dataset of audio recordings to train their AI model. Without proper documentation, they start searching the web for open datasets that might be suitable for

their project. They find a dataset that claims to contain a large collection of audio recordings. However, there is no detailed information about how the dataset was collected, who the speakers are, or any ethical considerations that were considered during the data collection process. The researcher is unsure if the dataset includes recordings of individuals without their consent or if it violates any privacy regulations. They also cannot easily determine whether the conditions under which the recordings were collected match the scenarios they are interested in modeling.

To ensure responsible AI practices, the researcher needs to evaluate the dataset thoroughly. They spend a significant amount of time manually examining the dataset, trying to find any clues about its quality, ethical considerations, and compliance with regulations. However, due to the lack of proper documentation, they cannot make an informed decision. After investing a considerable amount of time, the researcher realizes that the dataset does not meet their organizational policies and cannot be used for their project. They have to start the search process all over again, wasting valuable time and resources.

Therefore, the ability to filter out datasets that do not meet organizational policies or the specific objectives of the task at hand, can create significant time and cost savings. As a first step, automating and/or simplifying the evaluation of datasets through metadata can help researchers quickly identify datasets that do not meet their needs and/or organizational policies. However, it's important to note that this does not eliminate the need for EDA on a dataset that will be used to train production models, i.e. identifying properties of the dataset that can be learned through data exploration, like demographics. Instead, it aims to reduce the amount of time invested in filtering datasets that do not meet broader criteria.

Furthermore, given the time-consuming and costly nature of finding the right dataset, an applied documentation framework becomes indispensable in making the process of discovering and evaluating open datasets more efficient and cost-effective. By implementing a comprehensive documentation framework, we can help address biases, enhance transparency, and promote responsible AI practices. Such a framework empowers researchers and practitioners to thoroughly understand the data they are using, enabling them to assess its suitability and identify any potential biases or limitations. Additionally, it streamlines the dataset discovery process by providing clear and standardized documentation, reducing the time and effort required to find and evaluate datasets.

By emphasizing the significance of a reliable documentation framework, we can elevate the quality and trustworthiness of AI systems, fostering responsible and ethical AI practices. In light of this, we propose a no-code, machine-readable documentation framework to assist in the evaluation of open datasets, improve usability, and consider responsible AI aspects. Our contributions include:

- Publication of the proposed JSON-based metadata framework for open datasets on GitHub; <https://github.com/microsoft/opendatasheets-framework>
- Implementation of a public no-code solution hosted on GitHub Pages to generate and evaluate this metadata; <https://microsoft.github.io/opendatasheets>
- Discussion of recommendations to maximize the potential of this framework, improving efficiency, and providing transparency for consumers of open datasets

Related Work

In recent years there have been calls for comprehensive documentation of the datasets used to train and evaluate AI systems^{1,14-18}. Several data documentation frameworks and tools have been proposed with the goal of encouraging thoughtful reflection on datasets and transparency about their makeup and creation process. Datasheets for datasets¹⁴ is a documentation framework designed to encourage dataset creators to reflect on the implicit and explicit choices behind their data and to enable those interested in using a dataset to train or evaluate their system to make more informed decisions about its appropriateness. Developed concurrently, data nutrition labels^{15,19} also include information about a dataset's provenance and makeup. The authors created a publicly available tool to create a nutrition label that highlights usage restrictions and potential harms that may arise from the dataset's use. Data statements for natural language processing¹⁶ are designed specifically for natural language datasets and include specialized information such as speaker demographics and language variety. More recently, researchers at Google released the Data Cards Playbook²⁰, while researchers at Microsoft released the Aether Data Documentation Framework, a variant of datasheets for datasets adapted to meet the needs of industry practitioners¹.

While research shows that well-designed data documentation can be effective at helping to identify ethical issues²¹, good documentation can be challenging to create. In a study of industry practitioners' data documentation perceptions, needs, challenges, and desiderata, Heger et al.¹ found that practitioners creating dataset documentation had trouble making connections between the questions they were asked to answer and their RAI implications and difficulty providing information that someone unfamiliar with their datasets would need to understand the data. Based on their findings, they derived seven design requirements that data documentation frameworks should satisfy. Briefly, these include 1) making the connection to RAI more explicit; 2) making data documentation frameworks practical; 3) adapting data documentation frameworks to different contexts; 4)

supporting simple tasks with automation without automating away responsibility; 5) clarifying the target audience for the documentation; 6) standardizing and centralizing data documentation; and 7) integrating data documentation frameworks into existing tools and workflows.

Inspired by and building on this line of work, our proposed framework, Open Datasheets, introduces a machine-readable metadata format for open datasets. Drawing on questions included in datasheets for datasets¹⁴ and the Aether data documentation template²², it includes detailed information about the datasets, including responsible AI considerations. Furthermore, it adheres to the seven design principles of Heger et al.¹, as discussed in the design and implementation sections. Unlike existing tools, the framework aims to integrate with existing open platforms and support conversion to other standard formats, such as JSON-LD, providing an applied and efficient open framework that is user-friendly for non-developers and other open data publishers and users.

Other documentation frameworks and formats, such as those utilized by HuggingFace^{23,24}, were also considered but not included here due to their specialization and association with proprietary platforms. The evaluation focused on flexible published frameworks that can be adapted, are open, and align closely with the 7 design principles proposed by Heger et al.¹.

Methods

The approach for designing and implementing the *Open Datasheets* framework involved a multistage process that incorporated feedback from various dataset producers. This approach also included researching existing data documentation frameworks^{14–16,19,20} and studying the challenges and opportunities identified in previous studies^{1,14–18} on data documentation.

Our primary research aim was to identify the areas that needed improvement in practical implementation for open data documentation. To achieve this, we collaborated with data scientists and researchers from the Microsoft AI for Good Research Lab²⁵, who leverage and publish open datasets for their efforts to solve societal challenges with AI. Through hands-on experience and collaboration, we aimed to fill the gaps in existing frameworks and documentation practices.

The second aim of our approach was to leverage an existing documentation standard to facilitate adoption, avoid reinvention, and focus on bridging the gaps identified in previous studies on responsible AI documentation. We selected a format and standard based on criteria such as simplicity, user usability, and practicality to support non-developers, researchers, and developers.

The third aim was to provide a framework that includes no-code tooling, resources, and guidance that facilitate the understanding of the importance and application of documentation for open datasets.

We engaged with a diverse cohort of open dataset producers from Microsoft, including data scientists, researchers, analysts, and program managers involved in dataset publication and documentation. Their varying perspectives provided valuable insights.

Initial Iteration

The initial evaluation phase occurred during the early development of the *Open Datasheets* framework web application. We employed the following methodologies:

- **Case Study Analysis:** We applied the *Open Datasheets* framework to three distinct open data publications hosted on GitHub^{26–28}. These projects were selected to represent a range of producers, including non-developer researchers and data scientists.
- **Qualitative Feedback:** We collected insights from publishers through observations and conversations. Their feedback helped us understand the importance of guidance for responsible AI documentation and determine what information and features should be included.
- **Usability Testing:** We conducted usability testing to identify challenges in the absence of automated features and inline guidance.

Challenges

Identified Publishers experienced difficulties in several areas:

- Comprehending the value of responsible AI documentation.
- Deciding on the required information for adequate documentation.
- Navigating the framework without inline guidance or automation tools.

These observations were consistent with findings from Heger et al.'s research¹, which highlighted similar challenges in documentation tools.

Iterative Improvement

In response to the feedback from the initial iteration, we refined the framework to provide a better balance of features:

- **Enhanced Inline Guidance:** We improved the inline help components to instruct users on creating effective responsible AI documentation, providing clear examples and definitions.
- **Automation Features:** We introduced automation tools to streamline the documentation process. These tools can infer and prefill certain parts of the datasheets based on the data itself, reducing manual input.
- **User Documentation:** We provided comprehensive user documentation, including a step-by-step guide and resources, to help dataset producers navigate the framework and understand best practices for responsible AI documentation.

Design and Implementation Insights

The refined design and implementation were tailored to address the key challenges identified, resulting in a more intuitive and resourceful framework. Automation and guidance components were fine-tuned to ensure dataset producers receive support throughout the documentation process.

The methods employed serve to refine the *Open Datasheets* framework into a tool that accommodates varied user needs while emphasizing the importance of robust documentation for open datasets. Each step of the process has been scrutinized and enhanced to align with feedback and prior research outcomes.

Design

The *Open Datasheets* framework is designed to streamline the documentation of open datasets. It aims to foster the inclusion of crucial information that assists users in comprehending potential biases, privacy concerns, and other elements of responsible AI. This is achieved by striking a balance between automation and thoughtful deliberation in the documentation process¹. The ultimate objective is to enhance data documentation for open datasets, improve their discoverability, and promote their reliability, fairness, and transparency when used for AI models.

The *Open Datasheets* format specification is rooted in the Datapackage standard, a product of the Frictionless Data project²⁹. This open-source initiative offers tools and standards for data management. The Datapackage specification, accessible on GitHub³⁰, is a user-friendly JSON based format that acts as a container for describing a set of data. Also, it supports different data types, including complex ones such as tabular and geographic data, allowing users to document the basic information of a dataset in a way that is tailored to their specific needs. This standard has been widely adopted on open data platforms such as <https://data.world>.

The *Open Datasheets* framework utilizes the Datapackage format specifications to consistently document basic dataset details. This approach enables users to easily comprehend the information on the datasheets by avoiding specialized terminology and adhering to the standard language and usage of the JSON format. This format is machine-readable and seamlessly integrates into data workflows, thereby enhancing data discovery, access, and usage. The adoption of these specifications fosters compatibility and uniformity in the realm of data documentation.

Listing 1. Datapackage Sample

```
{
  # general "metadata"
  "name" : "a-unique-human-readable-and-valid-url-identifier",
  "title" : "A descriptive title",
  "licenses" : [...],
  "sources" : [...],
  # list of the data resources
  "resources": [
    {
      ... resource information ...
    }
  ]
}
```

Although the Datapackage specification is beneficial for general data documentation, it may not suffice for organizations that need to comply with regulatory standards and ensure responsible AI practices. To cater to this need, the *Open Datasheets* framework incorporates concepts from “Datasheets for Datasets”¹⁴. These concepts, such as detailed descriptions of a dataset’s

privacy implications to encourage careful reflection, are integrated into a machine-readable format. This additional information is crucial for organizations to make informed decisions about dataset usage and ensure compliance with their policies.

Moreover, the *Open Datasheets* Framework adheres to the 7 design principles outlined in section 5.2 of Heger et al research¹. It explicitly connects data documentation with responsible AI through inline guidance based in part on Microsoft's Aether data documentation template²². The framework is practical, integrating with GitHub and providing a user-friendly interface. It adapts to different contexts, allowing customization. Automation is balanced with responsibility, automating foundational metadata extraction while guiding users on responsible AI metadata. It clarifies the target audience to data publishers, focusing on potential users of open datasets. Standardization is achieved through a machine-readable format, and integration with GitHub promotes collaboration. The framework seamlessly integrates into existing tools and workflows, associating documentation with the data publication and development lifecycle on GitHub.

Furthermore, with the emergence of Large Language Models (LLMs), machine-readable interpretation of complex documentation has greatly advanced. The *Open Datasheets* framework enables organizations to leverage these advancements by documenting datasets with more descriptive information, allowing for more effective analysis and interpretation. This includes automated interpretation as a preliminary step to exclude datasets that do not align with organizational policies.

Listing 2. RAI documentation Sample

```
"privacy": [{
  "sensitivity": {
    "description": "sensitivity types description",
    "types": [
      {
        "name": "political opinions",
        "description": "description of the content related to this type"
      }
    ]
  },
  "confidentiality": {
    "path": "https://microsoft.github.io/opendatasheets/confidentiality",
    "description": "description of the process to ensure the confidentiality
of the data subjects"
  }
}],
"procedures": {
  "collection": [{
    "description": "Procedure description",
    "path": "",
    "contributors": [],
    "methods": [
      {
        "name": "focus group",
        "description": "focus group description",
        "path": "/focusgroup.txt"
      }
    ]
  },
  "consent": [
    {
      "title": "consent form",
      "description": "consent form description",
      "path": "/consentform.pdf"
    }
  ]
}]
}}
```

Implementation

To achieve a no-code solution, the framework implements a user-friendly web application on GitHub Pages (Figure 1), a managed service by the GitHub Platform to host static websites on GitHub repositories. This implementation ensures the longevity and support of the framework. The web application features a wizard-style interface that assists in standardizing the documentation according to the framework's metadata format. It includes metadata parsers for common data file formats on the GitHub platform³¹, such as CSV, TSV, JSON, and others. These parsers extract metadata related to the file structure, including field names, types, and sample values. General metadata about the data file, such as filename, encoding, and size, is also extracted. This approach allows data publishers to focus solely on filling in the responsible AI metadata.

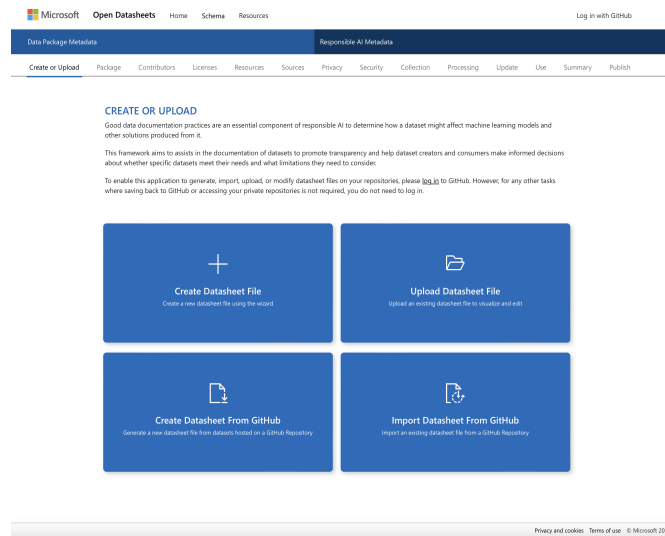


Figure 1. Open Datasheets Web Application (<https://microsoft.github.io/opendatasheets>)

The integration with the GitHub Platform, one of the largest open data platforms³¹, promotes openness and community. Currently, the web application enables seamless documentation of datasets on GitHub, enhancing discoverability and transparency for datasets hosted on the platform. However, the *Open Datasheets* framework is not limited to documenting datasets exclusively on GitHub. It provides the flexibility to create documentation for datasets and download the metadata file for publishing on any preferred platform. Moreover, the framework is openly available on GitHub, ensuring easy accessibility for data publishers.

The *Open Datasheets* framework simplifies the data documentation process for data publishers and promotes the inclusion of responsible AI elements. By automating the extraction of foundational metadata from data files and providing visual guidance for responsible AI considerations, the framework reduces the time and effort required for extensive documentation. This addresses the reluctance of data publishers to write lengthy documentation, as highlighted in Heger et al research¹.

For data users, the *Open Datasheets* framework offers detailed and machine-readable documentation, enabling informed decisions regarding dataset selection and use. The framework encourages data publishers to document potential biases and privacy implications associated with the dataset. This documentation empowers data users to evaluate the reliability and fairness of the dataset. The comprehensive documentation encouraged by the framework enhances transparency and trustworthiness, allowing data users to assess the dataset's quality, understand its limitations, and ensure it aligns with their ethical standards and requirements.

The standardized, machine-readable documentation also saves time and effort for data users, enabling them to programmatically evaluate the dataset's suitability for their needs and seamlessly integrate it into their organizational workflow for automated initial evaluation of usability and potential issues. Data users can also assess the documentation using the web application, which provides a user-friendly graphical representation of the metadata.

The *Open Datasheets* framework metadata format comprises two main sections: the dataset foundational metadata and the responsible AI metadata. The dataset foundational metadata extends the Datapackage standard and includes essential information such as the dataset's name, title, licenses, sources, and a comprehensive list of the data resources within the package. On the other hand, the responsible AI metadata builds on the concepts from "Datasheets for datasets"¹⁴ and Microsoft's Aether Data Documentation Template²². It provides information about the data's origin, processing methods, privacy implications, potential biases, and other relevant aspects. These components are further discussed in the following sections, specifically the *Datapackage* and *Responsible AI* sections.

Datapackage

The Datapackage module describes the foundational metadata of the framework. This module includes the package description, licensing information, contributors, resources, and sources associated with the dataset. This helps in understanding the composition of a dataset and the structure of the data assets contained within the datasets, which is crucial. The datapackage provides essential information for this purpose.

Package includes descriptors like name and description, version and creation date.

Licensing provides the terms under which the data is shared, an important factor in determining appropriate use cases for the given dataset.

Contributors provides information about the people or organizations that created and/or collected, and are sharing the dataset, which assists in identifying the provenance of the dataset.

Resources describe the list of data files included in the datapackage, including their data fields and value types in the case of tabular data.

Sources describe the list of data sources utilized for the creation and/or collection of the dataset, which also assist in identifying the provenance of the dataset.

Responsible AI

The Responsible AI module of the *Open Datasheets* framework focuses on documenting aspects that can impact user privacy, security, and biases that can undermine the trustworthiness of AI systems when used for training AI models. Evaluating datasets is crucial for creating responsible AI systems because the quality of the data directly affects the accuracy and fairness of the resulting system. Biases in the data can lead to biased outcomes, which can have serious consequences for individuals and society³²⁻³⁴. For instance, if an AI system used for hiring decisions is trained on biased data, it may perpetuate discrimination against certain groups.

To address this issue, organizations must evaluate their datasets for potential biases and take steps to mitigate them. The *Open Datasheets* framework provides a standardized approach to document and share dataset information, including metadata that aids in the evaluation of responsible AI. This metadata includes details about the data collection process, preprocessing or cleaning steps, and any identified biases.

By incorporating this metadata into their responsible AI workflows, organizations can automate much of the initial evaluation process and ensure that open datasets align with their responsible AI policies. This not only saves time and resources but also contributes to promoting fairness, reducing bias, and increasing transparency in the resulting AI system for all users.

However, it is important to note that human review will still be necessary to make decisions based on organizational policies. The *Open Datasheets* framework is not a substitute for human judgment but rather a tool to assist in the evaluation of responsible AI. By combining human expertise with machine-readable metadata, organizations can create more responsible and trustworthy AI systems.

With this objective in mind, the framework encompasses the following responsible AI areas:

- Privacy
- Data Access
- Collection Procedures
- Processing Procedures
- Update Procedures
- Use Cases

Privacy and Use Terms

To strike a balance between the advantages of open data and safeguarding individuals' privacy rights while promoting ethical data practices, it is crucial to evaluate the privacy implications of open datasets. The *Open Datasheets* framework promotes thoughtful consideration of a dataset's privacy implications by integrating metadata on confidentiality, data sensitivity, and usage terms for the data.

Confidentiality is a key aspect of privacy assessments. It refers to the protection of personal information from unauthorized access, use, or disclosure³⁵. Therefore, having this incorporated into the framework helps the users of the dataset understand the measures that were taken, if any, by the publishers to avoid the disclosure of personal information in their dataset.

Data Sensitivity in the context of privacy, refers to the level of sensitivity or potential harm associated with certain types of personal information included in a dataset. This includes sensitive attributes such as race, sex, religion, sexual orientation, health information, financial data, and other personally identifiable information (PII)^{36,37}. These sensitive attributes have the potential to cause harm or discrimination if mishandled, accessed, or disclosed without proper consent or safeguards.

Use Terms define the terms and conditions under which the data can be accessed, used, and shared. These terms help ensure that the data is used in a responsible and ethical manner, and that the privacy and confidentiality of individuals' personal information is protected. Use Terms can include restrictions on the types of analyses that can be performed on the data, limitations on the sharing or redistribution of the data, and requirements for obtaining consent or anonymizing the data before use. By defining these terms, data providers can help prevent the misuse or unauthorized access of sensitive data and ensure that the data is used for its intended purpose.

It is essential to evaluate a dataset for privacy implications because not only does it help protect the privacy of individuals whose data is being used, but also because many countries have laws and regulations that mandate organizations to safeguard individuals' privacy. By assessing a dataset for privacy implications, organizations can ensure that they comply with these laws and regulations.

Data Access

Including information on how to access a dataset is also essential for privacy and security considerations. One key reason is to clearly communicate whether the data can be accessed anonymously or if registration is required. This information helps potential users understand how to access the data and determine if they meet the necessary criteria to access it.

Furthermore, explaining the intended use of the dataset is an important aspect of data access documentation. By providing information on how the data can be accessed, dataset documentation allows users to assess whether the dataset aligns with their specific needs and purposes. For example, if the dataset is intended for academic research, requiring registration may be necessary to ensure that it is accessed only by qualified individuals or institutions.

Another important reason for specifying data access methods in the dataset documentation is to establish clear terms of use. This helps prevent non-qualified access or misuse of the data, preserving its integrity and privacy.

In line with these considerations, the *Open Datasheets* framework incorporates a module that enables data publishers to clearly specify how people can access the data. Additionally, it enables dataset providers to establish terms of use for the data and maintain data security. By providing this information, data providers can promote transparency, accountability, and responsible data usage.

Collection Procedures

Building upon the importance of evaluating datasets for privacy implications, the *Open Datasheets* framework also recognizes the significance of documenting data collection procedures as a crucial aspect of responsible AI considerations. By documenting these procedures, data users can gain insight into how and what data was collected. This promotes transparency and builds trust when using the data to train an AI system, reducing the potential for bias or discrimination in the system.

Additionally, documenting data collection procedures enables the replication of the process in the future, which is vital for researchers to validate the results of an AI system trained on that data. By replicating the study and following the documented collection procedures, researchers can confirm that the system is functioning as intended. If the replicated study yields different results, it suggests the need for further investigation to identify potential issues or biases within the original dataset or data collection procedures.

Furthermore, documenting data collection helps data users ensure that the process adheres to ethical and legal standards. The framework takes into consideration the methods, consent forms utilized for data subjects and the contributors of the collection procedures.

Methods describe the approach and instruments used to collect the data such as surveys, interviews, websites and others.

Consent is necessary to obtain when collecting personal data, especially sensitive data. Obtaining consent ensures compliance with privacy laws, respects individuals' rights, and promotes transparency and trust in data collection practices. For data users, understanding if consent was obtained provides insights into whether the data was collected ethically and in adherence to privacy laws and regulations.

Contributors describe who performed the data collection, if an individual, an organization or a third-party. This helps data users understand if the data publisher was not the entity that collected the data.

Processing Procedures

Documenting data pre-processing or processing procedures is important for reproducibility, transparency, accountability, and quality control, in addition to documenting collection procedures. This is because understanding how the data was processed makes it easier to identify errors or issues with the data. For example, let's say a dataset was pre-processed by removing all missing values without any imputation. If this pre-processing step is not documented, it may not be immediately apparent that some data points were removed, potentially leading to biased results or incorrect conclusions. Overall, documenting data preprocessing procedures is crucial for ensuring responsible and reliable AI systems.

As part of the data processing documentation, the *Open Datasheets* framework includes processing methods and contributors to this step.

Methods describe the approach and methodology used to process the data such as aggregation, anonymization, labeling and others.

Contributors describe who performed the data processing, if an individual, an organization or a third-party. Like the data collection contributors, data users can understand who processed the data.

Update Procedures

Datasets can be categorized as either static or periodically updated. A static dataset remains unchanged over time, while a periodically updated dataset undergoes regular updates to incorporate new data or revisions. When choosing between these options, researchers consider their specific objectives. If the research requires analyzing historical data or examining trends over a specific period, a static dataset may be suitable. Conversely, if the research involves analyzing real-time or evolving phenomena, a periodically updated dataset would be more appropriate. By documenting the procedures for updating the dataset, data users can gain insights into the timing and nature of the updates. This documentation promotes transparency and builds trust in the data, as users can verify its recency and reliability.

To facilitate this process, the *Open Datasheets* framework incorporates descriptors to identify whether a dataset is static or updated, the periodicity of updates, the method and versioning procedures, and the contributors for the update procedures.

Is Updated identifies whether a dataset is static or periodically updated.

Periodicity procedure describes the schedule or frequency of updates for the dataset.

Method describes the approach for updating the dataset, such as incrementally or a full refresh of the dataset.

Versioning procedure indicates how the data is versioned and how often.

Contributors describe who performs the data updates, whether it is an individual, an organization, or a third-party. Like the data collection and processing contributors, data users can understand who is responsible for updating the data.

Use Cases

The framework also incorporates the concept of use cases into the documentation. This involves documenting how a dataset can be used and what it cannot be used for, which is crucial for responsible AI. By understanding the potential uses and limitations of a dataset, it becomes easier to comprehend its boundaries. This documentation plays a vital role in preventing unintended consequences that may arise from utilizing the data in unintended ways. Additionally, it helps identify ethical considerations that must be taken into account when using the data in an AI system.

Overall, documenting the permissible and impermissible uses of a dataset is a crucial step in ensuring responsible AI systems. It provides guidelines for utilizing the data in a responsible and ethical manner, while also highlighting clearly ill-advised uses.

As a result, the *Open Datasheets* framework considers including use cases or examples of use for the dataset.

Conclusion

In this publication, we have proposed and explored the integration of Responsible AI (RAI) documentation into a machine-readable metadata framework for open datasets. Our contribution includes a no-code web application designed to simplify the documentation process. The *Open Datasheets* framework we have developed provides a robust foundation for dataset documentation, offering simplicity and flexibility that caters to data publishers working with diverse data types and requiring a customizable data management approach.

The framework covers core areas crucial for evaluating RAI considerations in a machine-readable format. These areas include privacy, collection and processing procedures, limitations of use, and other relevant aspects. We believe that this framework can be integrated into automated workflows to determine whether a dataset satisfies an organization's compliance criteria and warrants human review, or if it should be automatically discarded.

By incorporating the concepts of "Datasheets for Datasets"¹⁴ and the design principles of Heger et al.¹ into our framework, we have developed a more comprehensive and transparent approach to sharing and documenting data. This can lead to improved decision-making and the development of more trustworthy AI systems.

Looking ahead, we propose several areas for future work on the *Open Datasheets* framework. These include extending metadata extraction automation to incorporate more data types, evaluating integration with data governance frameworks, fostering collaboration and community building around the data documentation framework, and automating the validation of the richness and quality of the free-form text for RAI evaluations.

In conclusion, the *Open Datasheets* framework is a significant contribution to the field of data science. It has the potential to enhance the quality and reliability of data used in research and decision-making, thereby fostering the development of more responsible and trustworthy AI systems.

References

1. Heger, A. K., Marquis, L. B., Vorvoreanu, M., Wallach, H. & Wortman Vaughan, J. Understanding machine learning practitioners' data documentation perceptions, needs, challenges, and desiderata. *Proc. ACM on Human-Computer Interact.* **6**, 1–29 (2022).
2. Contaxis, N. *et al.* Ten simple rules for improving research data discovery (2022).
3. Manyika, J. *et al.* Open data: Unlocking innovation and performance with Liquid Information. <https://www.mckinsey.com/capabilities/mckinsey-digital/our-insights/open-data-unlocking-innovation-and-performance-with-liquid-information> (2013).
4. Brinkhaus, H. O., Rajan, K., Schaub, J., Zielesny, A. & Steinbeck, C. Open data and algorithms for open science in ai-driven molecular informatics. *Curr. Opin. Struct. Biol.* **79**, 102542, DOI: <https://doi.org/10.1016/j.sbi.2023.102542> (2023).
5. Buolamwini, J. & Gebru, T. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*, 77–91 (PMLR, 2018).
6. Paullada, A., Raji, I. D., Bender, E. M., Denton, E. & Hanna, A. Data and its (dis) contents: A survey of dataset development and use in machine learning research. *Patterns* **2** (2021).
7. Fazlioglu, M. Training AI on personal data scraped from the web. <https://iapp.org/news/a/training-ai-on-personal-data-scraped-from-the-web> (2019). Accessed: 2023-12-08.
8. Price, W. N. & Cohen, I. G. Privacy in the age of medical big data. *Nat. medicine* **25**, 37–43 (2019).
9. Trivedi, A., Mukherjee, S., Tse, E., Ewing, A. & Ferres, J. L. Risks of using non-verified open data: A case study on using machine learning techniques for predicting pregnancy outcomes in india. *arXiv preprint arXiv:1910.02136* (2019).
10. Bruno, B. The True Cost Of Bad Data And How It Can Hinder The Benefits Of AI. <https://www.forbes.com/sites/forbestechcouncil/2023/09/01/the-true-cost-of-bad-data-and-how-it-can-hinder-the-benefits-of-ai> (2023).
11. Davie, M. Why Bad Data Could Cost Entrepreneurs Millions. <https://www.entrepreneur.com/en-au/growth-strategies/why-bad-data-could-cost-entrepreneurs-millions/332238> (2019).
12. Cote, C. WHAT IS DATA INTEGRITY AND WHY DOES IT MATTER? <https://online.hbs.edu/blog/post/what-is-data-integrity> (2021).
13. Sakpal, M. How to Improve Your Data Quality. <https://www.gartner.com/smarterwithgartner/how-to-improve-your-data-quality> (2021).
14. Gebru, T. *et al.* Datasheets for datasets. *Commun. ACM* **64**, 86–92 (2021).
15. Holland, S., Hosny, A., Newman, S., Joseph, J. & Chmielinski, K. The dataset nutrition label. *Data Prot. Priv.* **12**, 1 (2020).
16. Bender, E. M. & Friedman, B. Data statements for natural language processing: Toward mitigating system bias and enabling better science. *Transactions Assoc. for Comput. Linguist.* **6**, 587–604 (2018).

17. Hutchinson, B. *et al.* Towards accountability for machine learning datasets: Practices from software engineering and infrastructure. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, 560–575 (2021).
18. Raji, I. D. & Yang, J. About ml: Annotation and benchmarking on understanding and transparency of machine learning lifecycles. *arXiv preprint arXiv:1912.06166* (2019).
19. Chmielinski, K. S. *et al.* The dataset nutrition label (2nd gen): Leveraging context to mitigate harms in artificial intelligence. *arXiv preprint arXiv:2201.03954* (2022).
20. Pushkarna, M., Zaldivar, A. & Kjartansson, O. Data cards: Purposeful and transparent dataset documentation for responsible ai. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, 1776–1826 (2022).
21. Boyd, K. L. Datasheets for datasets help ml engineers notice and understand ethical issues in training data. *Proc. ACM on Human-Computer Interact.* **5**, 1–27 (2021).
22. Microsoft. AETHER DATA DOCUMENTATION TEMPLATE. <https://www.microsoft.com/en-us/research/uploads/prod/2022/07/aether-datadoc-082522.pdf>.
23. McMillan-Major, A. *et al.* Reusable templates and guides for documenting datasets and models for natural language processing and generation: A case study of the huggingface and gem data and model cards. *arXiv preprint arXiv:2108.07374* (2021).
24. Lhoest, Q. *et al.* Datasets: A community library for natural language processing. *arXiv preprint arXiv:2109.02846* (2021).
25. Microsoft AI For Good Lab. <https://www.microsoft.com/en-us/research/group/ai-for-good-research-lab/> (2024). Accessed: 2024-02-24.
26. clinical visit note summarization corpus. https://github.com/microsoft/clinical_visit_note_summarization_corpus (2023). Accessed: 2023-12-08.
27. clandestino. <https://github.com/microsoft/Clandestino/tree/main> (2023). Accessed: 2023-12-08.
28. RTP-LX. <https://github.com/microsoft/RTP-LX> (2023). Accessed: 2023-12-08.
29. Frictionless Data Team. Frictionless Data. <https://frictionlessdata.io> (2023). Accessed: 2023-12-08.
30. Team, F. D. Frictionless Data Specs. <https://github.com/frictionlessdata/specs> (2023). Accessed: 2023-12-08.
31. Roman, A. C. *et al.* Open data on github: Unlocking the potential of ai (2023). 2306.06191.
32. Obermeyer, Z., Powers, B., Vogeli, C. & Mullainathan, S. Dissecting racial bias in an algorithm used to manage the health of populations. *Science* **366**, 447–453 (2019).
33. Chen, Y., Clayton, E. W., Novak, L. L., Anders, S. & Malin, B. Human-centered design to address biases in artificial intelligence. *J. Med. Internet Res.* **25**, e43251 (2023).
34. Silberg, J. & Manyika, J. Tackling bias in artificial intelligence (and in humans). <https://www.mckinsey.com/featured-insights/artificial-intelligence/tackling-bias-in-artificial-intelligence-and-in-humans> (2019).
35. CSRC, N. confidentiality - Glossary: CSRC. <https://csrc.nist.gov/glossary/term/confidentiality>. Accessed: 2023-12-08.
36. Ringrose, K. New categories, new rights: The CPRA’s opt-out provision for sensitive data. <https://iapp.org/news/a/new-categories-new-rights-the-cpras-opt-out-provision-for-sensitive-data> (2021). Accessed: 2023-12-08.
37. ICO. What is special category data? <https://ico.org.uk/for-organisations/uk-gdpr-guidance-and-resources/lawful-basis/special-category-data/what-is-special-category-data>. Accessed: 2023-12-08.

Author contributions statement

A.C.R. and J.L.F. conceived the concept. A.C.R. implemented the framework. A.C.R. and V.S. conducted the evaluation. A.C.R. wrote the draft manuscript. S.B. contributed to the design of the responsible AI components. J.W.V., K.X., C.B. and J.T. contributed to the content of all sections of the manuscript. V.S. contributed to the legal and privacy sections. All authors reviewed the manuscript.

Additional information

Competing interests statement

The authors declare no competing interests.