

Compliance Rating Scheme: A Data Provenance Framework for Generative AI Datasets

Matyas Bohacek*
maty@stanford.edu
Stanford University

Ignacio Vilanova Echavarri*
i.vilanova21@imperial.ac.uk
Imperial College London

Abstract

Generative Artificial Intelligence (GAI) has experienced exponential growth in recent years, partly facilitated by the abundance of large-scale open-source datasets. These datasets are often built using unrestricted and opaque data collection practices. While most literature focuses on the development and applications of GAI models, the ethical and legal considerations surrounding the creation of these datasets are often neglected. In addition, as datasets are shared, edited, and further reproduced online, information about their origin, legitimacy, and safety often gets lost. To address this gap, we introduce the Compliance Rating Scheme (CRS), a framework designed to evaluate dataset compliance with critical transparency, accountability, and security principles. We also release an open-source Python library built around data provenance technology to implement this framework, allowing for seamless integration into existing dataset-processing and AI training pipelines. The library is simultaneously reactive and proactive, as in addition to evaluating the CRS of existing datasets, it equally informs responsible scraping and construction of new datasets.

CCS Concepts

• **Information systems** → *Data provenance; Data extraction and integration; Multimedia and multimodal retrieval*; • **Computing methodologies** → *Artificial intelligence; Machine learning*; • **Security and privacy** → *Human and societal aspects of security and privacy*.

Keywords

Datasets, Provenance, Generative AI, Ethics, Transparency

1 Introduction

As Generative Artificial Intelligence (GAI) applications become increasingly intuitive and their results more realistic, their adoption is becoming widespread [69]. This exponential growth in performance and adoption in recent years is partly facilitated by the abundance of large-scale open-source datasets, which are often created through unrestricted and opaque data collection practices [16, 46]. Datasets play a crucial role in the AI ecosystem [54] as they are the primary source of training for most AI systems. While much of the literature focuses on the development and applications of GAI models, many ethical and legal considerations surrounding dataset creation currently remain unaddressed [16, 46].

Many AI researchers and practitioners obtain training data on the internet [36, 75]. With thousands of publicly available datasets [17, 97], platforms like Hugging Face [59] and GitHub [32] have become the backbone of today’s AI infrastructure. This model of dataset sharing emerged organically [34] in the 2000s with pioneering datasets such as ImageNet [33] and has been present for several years with no oversight or formal framing [52, 83]. When the recent boom of GAI erupted in 2023, the same model of dataset sharing practices prevailed—and with it, the same legal and ethical challenges [65].

The democratization of GAI has equally caused a surge in malicious activity [16, 81], notably through impersonations, copyright infringement, and deepfake pornographic footage [16]. Yet, the complexity and opacity of the structure of advanced GAI models result in a lack of traceability and accountability of the datasets that fuel them. To illustrate this challenge, we pose the following example: When a researcher or a practitioner wants to use a publicly available dataset for AI training, standard practice suggests the use of a publicly available image dataset. While the dataset’s license is valid, and the use case permitted, much of the dataset is often scraped without the consent of its content creators. Consequently, while the researcher accepts the terms of use outlined by the dataset’s authors, and be under the impression that the use of this dataset is legitimate, this might not always be the case. Indeed, for large datasets with millions (and more often than not, billions) of data points are required to train advanced GAI models, and thus, manually inspecting each data point becomes virtually impossible. Yet this common practice might result in a series of grave legal and ethical consequences, varying from copyright infringement, to using illegal material such as Child Sexual Abuse Material (CSAM) — unknown to the researcher.

This scenario was the case of the LAION-5B dataset [82], powering popular AI image generators such as Stable Diffusion [76]. This dataset ultimately had to be removed from distribution as a result of the two issues described above [89].

There are two critical moments that have led to this undesirable outcome. The first pertains the researcher’s review of the license and acceptance of the dataset’s terms of use. In this largely unregulated landscape [10, 22], where copyright infringement remains a contentious issue [27, 58], licenses and terms of use present some of the few recognized legal standards [11, 52, 56]. However, as the licenses are often written directly by the dataset authors, users can become overwhelmed and misled by the licensing requirements [44]. Moreover, a recent study revealed that nearly half of popular AI

*Both authors contributed equally to this research.

training datasets exhibit similar issues: they include data whose creators were not asked or informed about its inclusion, potentially violating copyright, while the license makes it look like the use of the data is permissible [63]. As ethical and legal frameworks for AI datasets are still in their infancy [57, 92], it is not clear who is responsible for such misconduct. Currently, the responsibility is often placed on the authors of the datasets [60]. Nevertheless, the dataset license agreed upon by the researcher could put the liability on them as well [106]. Depending on the jurisdiction and context, this could render the researcher liable instead of the dataset author.

The second critical moment involves the researcher’s inability to verify the dataset authors’ claims of ethical and legal data sourcing. The researcher has no alternative but to trust the dataset authors on this (i.e., a trust-based system). From an ethical and legal perspective, these two key moments beg the following questions: first, on the side of the dataset authors, can the practice of including unauthorized data points during scraping be prevented? Second, on the side of the researcher, can the accuracy of the presented license and dataset policies be verified to prevent misuse?

To the best of our knowledge, current literature and practice shows that no. The only way for the dataset author to verify that all scraped data points satisfy their chosen criteria is to manually inspect each one. The same limitation applies to the inverse case of a user testing whether the dataset’s self-reported license and policies match reality.

To address these critical questions, we propose a set of four practical principles for accountable, license-compliant datasets. These principles are informed by existing dataset sharing practices and latest data provenance technologies. We then conceptualize the Compliance Rating Scheme (CRS) as a trustless tool to evaluate a given dataset’s compliance with these principles. Finally, we develop a Python library, *DatasetSentinel*, which allows dataset authors to integrate these principles into their scraping pipelines and enables users to verify the CRS of datasets they are considering. We open-source the library at [anonymized].

The rest of the paper is outlined as follows. We first contextualize GAI within the broader context of datasets for AI training and provide an overview of current concerns. We then position these concerns within the existing ethical and legal frameworks, which we synthesize into the practical principles that guide our solution. This discussion transitions into a description of the CRS. Next, we introduce the development, structure, and evaluation of *DatasetSentinel*, the Python library we open-source. Finally, we discuss the implications, limitations, and future directions of this research.

2 Background and Related Work

In this section, we provide an overview of the GAI landscape, focusing on training datasets and their misuse. We identify a gap in existing work, into which we later position our contribution.

2.1 Generative AI

Generative AI (GAI) [13] refers to AI systems that can synthesize novel text [105], image [15], video [62, 86], audio [103], and other modalities. These systems have recently undergone a substantial leap in generation quality and ease of use [66]. We first saw this in the domain of text generation with the advent of large language

models (LLMs) [21, 90]. Based on a simple premise of completing the next words in sentences, LLMs have improved to the point where they evince emergent capabilities [91, 95], such as text analysis and question answering. In fact, the quality of AI-generated text today is such that humans, in some cases, cannot distinguish AI-generated texts from human-written texts [18, 23].

Similar leaps are being made on the front of image, audio, and video generation [101]. It was not long ago when text-to-image generation was restricted to Generative adversarial networks (GANs) [39], which could only generate domain-specific images [49]. Today, diffusion-based models [102] work across a wide range of subject domains and create text-conditioned images of high quality. Similar to text, the latest methods for image generation got to a point where, in some contexts, humans cannot distinguish AI-generated images from real photos [68]. Methods for video [50, 85, 94, 99] and audio [48, 67, 100] generation are a more recent addition to the scope of modalities generated by AI, but we can expect them to follow a similar trajectory to the text and image modalities.

2.2 Datasets

The scale and quality of AI training datasets have been essential to the recent leap in GAI systems [40, 66]. As such, datasets play an essential role in the AI ecosystem [54] because they are the primary source of training for most AI systems. That is because—despite advances in reinforcement learning and other modes of AI—supervised training of AI models, in which a model is trained once before put to use, still dominates. Beyond training, datasets also allow teams to compare the performance of their solutions against a standardized benchmark. More broadly, datasets steer the focus and work within the community [80].

Albeit significant progress has been made in the areas of foundation models, fine-tuning, and knowledge transfer, most AI systems require task-specific datasets to achieve good performance [40, 93]. Therefore, new datasets are constantly being released by research institutions [72], companies [74], and laypeople [26] alike as new AI tasks and contexts emerge. Notably, the data acquisition and annotation of a large dataset is financially demanding [31] and, while some companies can afford to undergo such a project with manual curators and annotators, many entities resort to a less expensive data acquisition mode through internet scraping [53]. Many of today’s datasets for AI training are thus indiscriminately scraped from the internet [28, 73], often with little or no consideration of the ethical and legal implications of such practice.

2.3 Dataset Life Cycle and Stakeholders

We consider the distinction between the dataset author and AI practitioner in the dataset life cycle essential, as both parties approach it with different incentives and risks [43]. As such, we frame the life cycle of a dataset by its creation (performed by the dataset author) and by its use (performed by the AI practitioner). The author of the dataset (e.g., a research institution, a company, or an individual) first identifies the scope of the dataset: they select which task(s) and context(s) the dataset will address and create fundamental policies about the ingested data [46, 79]. Next, they set up the data ingestion sources (e.g., custom capture, data purchase, or internet data scraping) and annotation mechanisms (e.g., hiring annotators or

employing automated solutions) and, finally, proceed with the construction itself. Once constructed, the dataset’s license and terms of use are packaged with the dataset and shared. Over time, the owner may decide to make changes to the dataset. Such modifications may include simple error fixes or, if the changes are significant enough to justify so, constitute a new version of the dataset. As *Hutchinson et al.* argue, datasets powering AI are often used, shared, and reused with little visibility into the processes of deliberation that lead to their creation [46].

The paths of the dataset author and the AI practitioner intersect at the dataset distribution platform. Often, the datasets are shared on Hugging Face Datasets¹, Kaggle², GitHub³, or custom websites. At this point, the AI practitioner is considering which datasets would best fit their use case [104]. Once they decide which dataset(s) to use, they obtain the data from the platform and proceed to the model training, evaluation, and potential deployment (inference).

2.4 Challenge to Address

Manually inspecting every data point included in a dataset is virtually impossible as it is not uncommon for a single dataset to contain millions to billions of such data points. This makes it challenging for dataset authors to filter incoming data and for AI practitioners to verify whether the contents of a dataset match its description. There needs to be a systematic trustless approach to infer the provenance of a single data point that would enable the dataset author to filter incoming data points effectively and AI practitioners to assess a considered dataset. Existing metadata (e.g., EXIF) sometimes includes relevant information about the license, author, AI opt-out, etc., but this information is often missing, and even when present, it is inconsistent in formatting and terminology. There needs to be a systematic approach to address data provenance, in which license, AI training consent, and other preferences would be automatically embedded.

Moreover, the abundance of large-scale open-source datasets derives from the legal vacuum of online data collection and use practices. This state of affairs has led to calls for a more responsible, transparent, accountable, and human-centered approach to AI dataset practices [16, 46]. Consequently, we argue that a new framework is needed to better these unrestricted and unaccountable practices. We identify the modalities of image, video, and audio datasets as the most pressing to address. These modalities constitute some of the most prominent kinds of data points in datasets for GAI training and, as mentioned above, can, in some contexts, pose imminent privacy concerns for individuals.

2.5 Dataset Principles

The literature on dataset ownership encompasses a wide range of sub-themes, such as privacy, security, stewardship and governance, and transparency [12]. Yet, there is no standardized definition of what such ownership entails [12]. This begs the question of whether users’ data should be considered some sort of property, and whether this status would change when aggregated to a dataset protected under Intellectual Property (IP) laws. Legal theories suggest that

data (information) cannot be owned [45]. Common law does not recognize property in facts or information and considers data as such. Continental (Civil) law follows a similar approach but presents data rights as an “extension or subset of fundamental human rights” [45], which is also unsuitable for proprietization and commercialization. For this reason, most legal scholars focus on the protection of data instead [29]. However, as we have discussed, individuals have little to no practical rights on how their personal data is used in the context of AI training datasets.

In the case of AI, applications involve primarily two categories of tort: dignitary and property. The distinction derives from the nature of the harm caused. Dignitary torts typically encompass harm to a person’s reputation, honor, or dignity, such as unauthorized use of personal images – like pornographic deepfakes. Conversely, property torts in AI contexts often relate to interference with one’s property rights – such as copyright infringement with artists’ work used to train GAI models. While legal categorization is relatively straightforward, successful prosecution and liability remain incredibly complex.

We break down this complexity into three main areas. First, data collection and dataset practices: AI models often use publicly available data (such as photos, videos, and voice samples) scraped from the internet to generate context-specific outcomes. These data collection practices are often unrestricted and opaque, and rarely have explicit consent for specific uses [8, 35, 87, 88]. Second, overlapping jurisdiction and areas of law: Part of the legal complexity and ambiguity derives from a series of overlapping legal areas, such as artistic freedom, freedom of expression, the right to information, the right to privacy, and personality rights, among others. Furthermore, different jurisdictions might interpret these rights differently. Third, there is a lack of liability and accountability. The anonymity of the internet makes it difficult to determine who created or distributed the AI application’s output (e.g., a deepfake).

While the data coming out of AI systems has been an active area of study (from deepfake detection [19, 20] to watermarking to tracing detailed provenance information of AI-generated content [77, 78]), the data coming into these systems during training, which leads to AI and dataset misuse, has not [9, 14]. Modern data protection laws such as the EU’s General Data Protection Regulation (GDPR) (2018) [6] and the California Consumer Privacy Act (2018) [5] are built on The Fair Information Practice Principles (FIPPS) published by the Organization for Economic Cooperation and Development (OECD) in 1980 [1, 4], and has ever since been the guiding model for data protection. These principles have been adopted through various institutions and improved through frameworks such as the EU Data Protection Directive Principles (1992) [2], the Federal Trade Commission Privacy Principles (1998) [3], and the Asia-Pacific Economic Cooperation Privacy Framework (2004) [4]. Yet, as discussed, the Fair Information Practice Principles and most of the laws derived from them have failed in practice [24], as the data protection regimes built on them come short in providing a high standard of effective and efficient data protection and use [24]. As such, data protection is not an end in itself, but rather a tool for enhancing individual and societal welfare. We aim to pursue this goal by proposing an initial set of four practical principles to consider for dataset compliance in the context of AI. Inspired by prior work and data protection laws, these four principles are designed to

¹<https://huggingface.co/datasets/>

²<https://www.kaggle.com>

³<https://github.com>

be technologically implementable, and to provide actionable measures for prosecution in the eventuality of misuse. These principles consist of:

- (1) Responsibility and Liability
- (2) Effective and Efficient Enforcement
- (3) Prevention of Harm
- (4) Transparency and Fair Use

2.6 Data Provenance

Data provenance [70] refers to the records about the origin, ownership, and evolution of a file. It is concerned with any relevant information from the moment the file was created—be it as an authentic recording or as a synthetic digital product—to its present form. This information includes details about the entities, software, and specific changes, if applicable, that have in any way manipulated the file from its inception [64, 96]. Moreover, the data provenance may capture additional information about its author’s decisions for sharing it with third parties, including the license under which it is shared, whether or not it may be included in AI training, etc.

Establishing data provenance for files disseminated over the internet may be challenging [25], especially as they may be stripped of their basic metadata or additional attachments. Therefore, the literature has studied cryptographic methods for provenance [37], which allow for verifying any assertions made about a file in its provenance metadata. While there are many contexts in which establishing data provenance may be essential, this technology has gained most of its recognition recently amidst a wave of fake AI-synthesized images on the internet [84].

The Coalition for Content Provenance and Authenticity (C2PA) data provenance specification [77] created a standardized framework for data provenance metadata. As of yet, this is the largest effort striving to establish a standardized approach to deployable data provenance on the internet, and it has received traction from industry and academia alike. Content Authenticity Initiative (CAI) then materialized this standard into a functional, cryptography-based library and metadata scheme [78].

3 Compliance Rating Scheme

Our contribution comprises two parts, the first is the Compliance Rating Scheme (CRS). It is a set of criteria and a summarizing score that together serve as an intuitive indicator of a given dataset’s compliance with the principles outlined above. The CRS score is evaluated based on the following six criteria:

- (1) The sourcing, filtering, and pre-processing employed during data acquisition and annotation of the dataset is transparent. The code for these processes is either fully open-sourced or is described at a level of detail that would enable full reproduction of the dataset.
- (2) The dataset complies with the license and allowed use described in the provenance metadata of each included data point. This means that the licenses and allowed use of each individual data point fall within the scope and allowed use of the dataset as a whole.
- (3) The dataset flags any data points where compliance with the provenance metadata is inconclusive.

- (4) The dataset has an opting-out mechanism, allowing authors of the included data points to request their removal from the dataset if they had not previously given consent.
- (5) Any changes made to the content of the dataset—both to the data points themselves and their annotations—are traceable. There is a designated trace log that includes dated records of changes, listing which data points were impacted and how.
- (6) The dataset adds the dataset source and the retention period into the provenance metadata of each included data point.

The CRS score summarizes the dataset’s compliance with these criteria into a letter on the scale from "A" (the highest, most compliant score) to "G" (the lowest, least compliant score). Starting at "G", each satisfied criterion moves the CRS of the evaluated dataset up by one letter grade. This means that, if a dataset does not meet any of these criteria, it receives a CRS of "G". Contrarily, if the dataset meets all criteria, it receives a CRS of "A".

While there are many contexts in which this assessment could be desired, it is primarily targeted at AI practitioners when they are deciding which dataset(s) to use for training in their AI project. Returning to our example of an AI researcher from Section 1, the researcher can benefit from the CRS score to determine which datasets out of the ones she was considering satisfy the legal and ethical standards she desired. Even if a dataset’s description claimed so, she could verify that through the CRS score, and thus remove the element of trust in the dataset’s creator good faith from the equation.

4 Library

The second part of our contribution is *DatasetSentinel*, an open-source Python library implementing the CRS. The library, available at [anonymized], is written in Python, the most popular programming language for AI research and development [38]. It can be easily integrated into existing dataset and AI pipelines as it is compatible with PyTorch [71], TensorFlow [7], MLX [42], HuggingFace [98], Kaggle, and custom dataset-sharing platforms, requiring minimal changes to existing code structures.

The library leverages the Content Authenticity Initiative’s (CAI) library [78] and the Coalition for Content Provenance and Authenticity’s (C2PA) data provenance standard [77]. CAI’s library is the official implementation of C2PA, the leading data provenance standard widely adopted across social media platforms and hardware products. Note, however, that the CRS is not dependent on C2PA and CAI; we simply found it to be the most suitable and adopted data provenance framework to date.

4.1 Features

The library has two features: (1) determining whether a single data point considered for inclusion in a dataset would be compliant with the CRS and (2) calculating the overall CRS score of a dataset. We expect feature 1 to be used during the creation of a new dataset, as the dataset author is deciding which data points to include. On the other hand, we expect feature 2 to be used primarily by AI practitioners as they are deciding whether to use a dataset.

4.1.1 Feature 1. Feature 1, determining whether a single data point considered for inclusion in a dataset would be compliant with the

CRS, requires that data point-level criteria (C2, C3, and C6) be evaluated. The dataset author can pass a considered data point (e.g., an image, video, or audio file) to *DatasetSentinel*. The library will return a boolean indicating whether the data point is compliant. If not, it will list which criteria are violated and provide a description of the reasoning. The schematic overview of this feature is shown in Figure 3 (Appendix B). Put into practice, if the dataset author wants their dataset to remain CRS-compliant, they would call this function for every considered data point and drop those for which the assessment is negative.

4.1.2 Feature 2. Feature 2, calculating the overall CRS score of a dataset, requires that both dataset- (C1, C4, and C5) and data point-level (C2, C3, and C6) criteria be evaluated. The AI practitioner can provide the full dataset for consideration to *DatasetSentinel*. This dataset can be stored locally or on a dataset sharing platform. The library will return a final CRS score, along with the reasoning: for each criterion, it indicates whether the dataset is compliant, and lists data points that are in violation, if applicable. The schematic overview of this feature is shown in Figure 4 (Appendix B). Put into practice, if the AI practitioner wants to ensure the legal and ethical standing of a considered dataset, they would call this function on the dataset, review the assessment, and decide whether it is appropriate to move forward with it.

4.1.3 Dataset-level Criteria. Criteria C1, C4, and C5 concern features of the dataset that are determined by the means of distribution. The compliance of a given dataset with these criteria can thus be determined by the inspection of the dataset’s page on the distribution platform. For datasets hosted on Hugging Face and Kaggle, *DatasetSentinel* can infer much of this information from the standardized metadata on the dataset’s page. For GitHub and custom-hosted datasets, however, there is no standardized way of representing these features, and so *DatasetSentinel* uses an LLM to scan the content of the dataset repository and decide the compliance. To prevent false positive or false negative hits in such cases, *DatasetSentinel* presents the compliance with C1, C4, and C5 for the user to review. The user has the ability to manually override the library’s inference.

4.1.4 Data Point-level Criteria. Criteria C2, C3, and C6 concern features of data points included in the dataset. A given dataset is compliant with these criteria only when all data points satisfy the criterion. *DatasetSentinel* thus inspects each data point individually and verifies its compliance, which can be derived based on the provenance metadata of the data point (extracted using C2PA) and a set of conditions comparing the provenance information (including the license, whether the content creator opted out of AI training, etc.) to the dataset setting. Unlike dataset-level criteria, these criteria can be clearly determined without user confirmation.

5 Evaluation

In this section, we describe two modes of evaluation we employed for *DatasetSentinel* and CRS: an automated code quality assessment and a preliminary user study, surveying 5 recruited AI experts through a purposive (non-probability) sampling method.

#	Question
1	How easy is it to navigate the documentation?
2	How understandable is the documentation?
3	How understandable are the tutorials and examples?
4	How easily does the library design integrate into your development workflow?
5	How similar is the structure of the library interface to other libraries you have used before?
6	How likely are you to use the library in your workflow while working on a ML project?

Table 1: Survey questions used as a part of *DatasetSentinel* library usability evaluation

#	G	Nat.	Q1	Q2	Q3	Q4	Q5	Q6
P1	M	USA	7	7	7	4	5	4
P2	F	SWE	5	5	5	7	5	4
P3	M	IND	5	6	6	3	3	2
P4	F	USA	7	6	6	6	5	7
P5	M	USA	4	5	6	7	3	5
P6	M	IND	6	6	5	4	7	5
P7	M	USA	5	7	6	7	6	6
P8	M	SWE	7	7	6	7	6	7
P9	M	BGD	6	5	6	6	7	6
P10	M	NGA	5	6	5	5	5	6
P11	M	USA	7	7	7	7	7	6
P12	F	USA	7	7	7	7	7	7
P13	M	USA	3	6	6	4	4	4
P14	M	AUT	6	5	3	3	5	3
P15	M	HKG	4	3	4	5	6	7
Avg.			5.6	5.9	5.7	5.5	5.4	5.3

Table 2: Results of the library usability evaluation

5.1 Methodology

5.1.1 Code Quality. We used the Wily maintainability score⁴ on the scale from 0 to 100, with a higher score reflecting a better evaluation of the complexity, readability, and in-code documentation of the *DatasetSentinel* library. This suite of metrics is based on the Halstead complexity measures [41], which have been shown to increase code readability and minimize down-stream fault rates of the evaluated codebase [30, 51].

5.1.2 Library Usability. In addition, to better understand how this prototype would perform in real-life applications, we recruited 14 participants through a purposive sampling technique to evaluate its usability and robustness. Participants were recruited based on their expertise in the field of AI, and half of them are from the United States (0.5), the rest being from Sweden (0.14), India (0.14), Nigeria (0.07), Bangladesh (0.07) and Austria (0.07). The majority are male (0.85). Participants were asked to implement our script into a database and answer 6 evaluative questions (presented in

⁴<https://github.com/tonybaloney/wily>

Dataset	Source	Modality	C1	C2	C3	C4	C5	C6	CRS Score
SOD4SB	GitHub	Images	✓	✓	✓	✓	✗	✗	C
MS COCO	Custom website	Images	✓	✗	✗	✗	✗	✗	F
RANDOM People	Hugging Face	Videos	✓	✓	✓	✓	✓	✗	B
TikTok Dataset	Kaggle	Videos	✗	✗	✗	✗	✗	✗	G

Table 3: Results of the CRS case studies on four publicly available datasets. For each dataset, we report whether it satisfies CRS criteria C1 through C6, and to which CRS score this translates.

Table 1) on a 7-point Likert scale, where 1 very difficult and 7 very easy, with a high score reflecting greater usability. Participants were given the opportunity to add comments on their experience for a simple qualitative evaluation. The participants were recruited based on their technical expertise in AI and ML. Participants were not recorded; only their written answers were collected and fully anonymized.

5.2 Results

5.2.1 Code Quality. *DatasetSentinel*'s codebase obtained a mean Wily maintainability score of over 85, indicating an overall good code quality. The files that were indicated as lower-ranking mostly included connections between our framework and the provenance metadata flags of the C2PA library; we thus make it our priority to keep improving the library in this regard, mainly by adding in-code documentation.

5.2.2 Library Usability. The participants' answers are presented in Table 3. Quantitatively, we observe an overall positive response to our prototype as all questions are, on average, rated positively (≥ 5.6). While these results are preliminary and further work needs to be conducted to further improve the *DatasetSentinel* Library, they are nevertheless encouraging and optimistic. The majority of participants seem to agree that the script is easy to navigate (5.6/7), understand (5.9/7), and well documented with tutorials and examples (5.7/7) (questions #1, #2, and #3). The results relating to ease of integration (question #4, #5, and #6) were slightly less positive (5.5, 5.4, and 5.3 respectively), but encouraging nevertheless.

We would like to highlight that the responses to questions #5 and #6 depend on the type of AI project into which the user is integrating *DatasetSentinel*. For instance, one participant stated that "the primary reason I am unlikely to use this library in my projects is that I almost exclusively work with tabular data" (P3). This is a valid point, although in its current state, our library is designed to address the concerns resulting from image, audio, and video data files. Similarly, another participant stated that "my projects aren't really about ethics, which is the only reason I put only a 4" (P1). We find this statement to be a good reflection on the general dissociation found among practitioners between AI applications and ethics. Another participant stated that the CRS score's function was unclear, as they could not find any information online regarding this tool and asked for clarifications: "It is unclear whether the CRS score is something you invented or an agreed-upon standard. Searching for the CRS score take me to the Canadian government site..." (P2). This confusion was caused by the fact that we could not reveal the manuscript where we introduced and explained the CRS score to maintain high discretion and total anonymity. Similarly,

another participant stated that "I wanted to learn more about C2PA – a brief explanation and link would be great" (P13). We agree with this comment, as we believe it is crucial not to assume that every AI practitioner might be familiarised with C2PA, and how does the CRS score differ from it: "What's the difference between [DatasetSentinel] library and the C2PA Python library? Is C2PA more low-level and this one provides nicer abstractions, or is there something functionally different?" (P13). We have addressed these comments and updated the documentation. Other participants seem to appreciate the value of this work, as one mentioned that "I think that this library is very well organized, thoughtful and important for today's modern tech world" (P4) and another that "this project looks incredibly useful and helpful" (P5).

6 Case Studies

To put the CRS framework and *DatasetSentinel* library to practice, we applied them to four open-source datasets from different modes of distribution (GitHub, Hugging Face, Kaggle, and custom website). Next, we briefly describe these datasets and present their CRS assessment. The results are summarized in Table 3.

6.1 Use Case: SOD4SB

The SOD4SB dataset [55], released as a part of the MVA2023 Spotting Birds challenge, contains 39,070 images annotated with bounding boxes of birds. These images were taken by the dataset's authors. The dataset is distributed through GitHub. As with the previous dataset, it is not compliant with criterion C6. Additionally, it is not compliant with criterion C5, as there is no trace log of changes. This results in the CRS Score "C".

6.2 Use Case: MS COCO

The MS COCO dataset [61] contains over 300,000 images with annotations for object detection, segmentation, captioning, and keypoint detection. The images were gathered from Flickr. The dataset is distributed through a custom website. As with the previous dataset, it is not compliant with the criteria C5 and C6. Additionally, it is not compliant with criterion C4, as there is no opting-out mechanism; C3, as the data points with inconclusive provenance metadata are not flagged; and C2, as some data points are used against their license. This results in the CRS Score "F".

6.3 Use Case: RANDOM People

The RANDOM People dataset⁵ contains videos with human protagonists performing actions around the house, generated using a

⁵<https://anonymous.4open.science/r/random-people-dataset-D70F/>

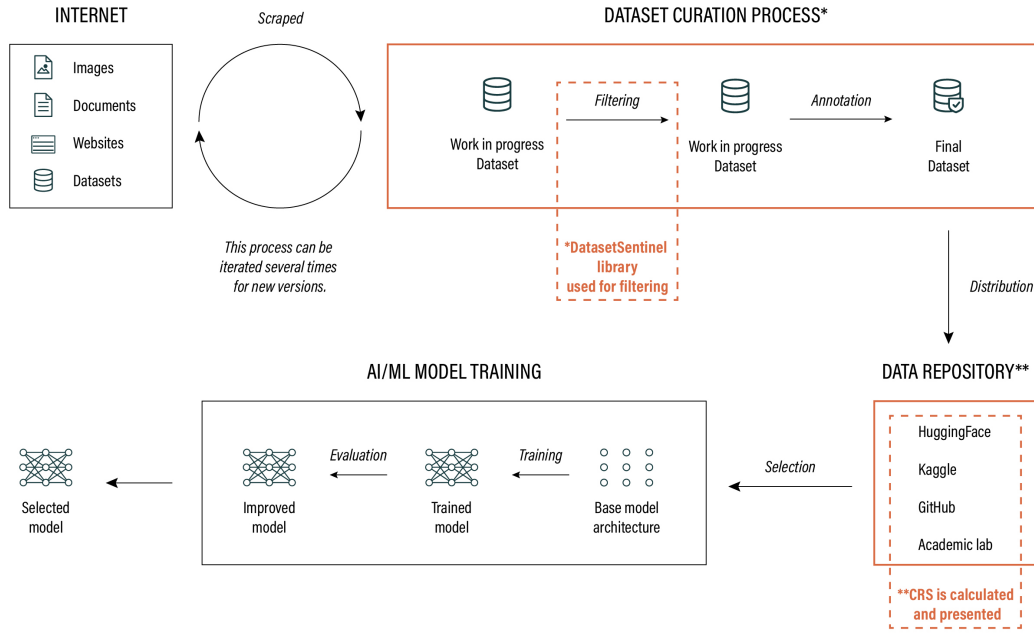


Figure 1: A schematic overview of the AI ecosystem workflow with the main stages of dataset and model development

pose-transfer AI model, along with the annotations of these actions. The identities used as a reference for pose transfer were consenting individuals gathered by the dataset authors, and the driving videos were from an open-source database whose creator had permission from all depicted participants. The dataset is distributed on Hugging Face. It is not compliant with criterion C6, as the dataset source and the retention period are not added to the provenance metadata of data points, resulting in the CRS Score "B".

6.4 Use Case: The TikTok Dataset

The TikTok Dataset [47] contains 300 dance videos, 10 to 15 seconds in length, sourced from TikTok. Additional 3D representations are also provided. The dataset is distributed through Kaggle. As with the previous dataset, it is not compliant with criteria C2, C3, C4, C5, and C6. Additionally, it is not compliant with criterion C1, as the sourcing, filtering, and pre-processing are not detailed at a level that would enable reproducing the dataset. This results in the CRS Score "G".

7 Discussion

By proposing a set of four practical principles to consider for dataset compliance in the context of AI, we aim to provide a framework that raises the discussion on the legality and ethics of AI applications. However, similar to most principles, these can be interpreted as highly conceptual and disconnected from current practices, often making them either irrelevant or challenging to implement. Precisely for this reason, we attempted to move away from a purely descriptive contribution to the literature, and provide a tangible

and prescriptive approach through our CRS tool and *DatasetSentinel* library. We highlight the specific points of the AI workflow at which we target our contribution, aiming to reduce the misuse of personal data for GAI training models and applications by introducing traceability and accountability of the datasets used for harmful purposes.

To this end, the first line of defense is with the *DatasetSentinel* library, which can be used by practitioners to filter the collected data. Using provenance metadata, the tool ensures that the data is compliant with the purpose of the dataset. The second line of defense is the CRS score, which calculates and informs the practitioners about the dataset's compliance with the practical principles embedded in its structure. These two intervention points in the life cycle of a dataset are illustrated in Figure 1.

We believe that the benefits of implementing this tool are twofold. In the long term, it benefits the AI field and, more broadly, society as a whole. Over time, poorly rated datasets (E and below) would stop being used as much and eventually become less impactful. We ground this belief in studies about consumers' quality standards expectations, showing that 92% of consumers tend to purchase products with at least a 4-star rating⁶. We believe the field of AI is no different. To induce this effect in AI practitioners while choosing datasets, we propose accompanying visuals for the CRS scores shown in Appendix 9. In the short term, it benefits the individual user as it removes the heavy lifting of manually conducting this type of analysis and helps protect themselves from any liability of data misuse.

⁶<https://explodingtopics.com/blog/online-review-stats>

We intend to render it more challenging for defendants accused of malicious activity through AI applications to plead ignorance about the nature or compliance of any given dataset. In the eventuality of a legal demand, the CRS score enables developers and regulators to gauge and verify the transparency, accountability, and security of any given dataset, with the ultimate objective of providing traceability and accountability. By doing so, we hope to help reduce the gap between digital technological innovation and ethics by providing a framework to responsibility, liability, and legal enforcement of data malpractices in the context of AI.

In the future, dataset-sharing platforms may adopt this tool on their end, which would remove the heavy lifting (of running this analysis) from individual users. Shown in Figures 7, 5, and 6 (Appendix C) are mockups that fictitiously contextualizes the CRS score in an online repository, as practitioners would perceive it. As observed in these mock-ups, the CRS score seemingly integrates with the rest of the dataset’s information, while providing a clear reading .

Regarding the adoption of CRS and the *DatasetSentinel* library, we do not expect them to be a mandatory requirement but rather a tool to support the AI community. By providing an overview of the compliance of any given dataset, both dataset owners and users can better reflect on their responsibility and liability towards the AI community, and make a more informed decision on the resources they use in their projects.

We are witnessing a growing interest among software and hardware companies in tracing the provenance of media in an attempt to fight misinformation and other malicious content. This trend is manifesting itself, for example, by an uptick of organizations joining coalitions such as the Coalition for Content Provenance and Authenticity (C2PA) [77]. It seems that there is a growing trend towards data traceability and immutability within the digital sphere. We therefore reiterate our belief in this project and its potential positive impact within the field of AI.

8 Limitations

As our prototype is in its infancy, we acknowledge its limitations and that there is still much research to be conducted until this framework can become a standard for ethical GAI use. For instance, as data provenance technologies are just rolling out, the majority of digital media available online still lacks provenance metadata. Nonetheless, many technological companies – both in software and hardware – are starting to deploy or announce the integration of data provenance technologies into their products. There are indications that this trend is becoming more and more common, as users express concerns over the use of their images, artwork and intellectual property; and companies are attempting to solve this. For instance, we do not discard the possibility of smartphone operating systems introducing an "opt-out" feature for all (or only selected images and videos) taken on the smartphone for AI training. As such, we expect that, within a few years, the vast majority of new digital media distributed on the internet will have provenance metadata. Another limitation is that the library is dependent on the existing data provenance protocols. To that end, our library can only analyze data types that are supported by these protocols and other dependencies. This should not pose a problem for most

current use cases, as the protocols support the most common data types for image, video, audio, and 3D objects. Still, moving forward, this dependency could introduce a delay in introducing support for new data types.

9 Conclusion

We call for a larger discussion confronting the unsustainable dataset practices in the AI community. While we recognize that the dataset sharing platforms have substantial power to influence the practical rules and guidelines, we argue that a value shift is also needed. Specifically, a broader awareness and appreciation of ethical and legal considerations surrounding datasets must be established for the rules and guidelines of dataset sharing platforms to have a meaningful impact. Our framework and tangible outputs can serve as a springboard for piloting and implementing these values into existing workflows.

References

- [1] 1980. OECD Privacy Principles. <http://oecdprivacy.org/>
- [2] 1992. EUR-Lex - 32016R0679 - EN - EUR-Lex. <https://eur-lex.europa.eu/eli/reg/2016/679/oj> Doc ID: 32016R0679 Doc Sector: 3 Doc Title: Regulation (EU) 2016/679 of the European Parliament and of the Council of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation) (Text with EEA relevance) Doc Type: R Usr_lan: en.
- [3] 1998. Privacy Online: A Report to Congress 7. <https://www.ftc.gov/reports/privacy-online-report-congress>
- [4] 2004. APEC Privacy Framework. <https://www.apec.org/publications/2005/12/apec-privacy-framework>
- [5] 2018. *California Consumer Privacy Act 2018*. <https://oag.ca.gov/privacy/ccpa>
- [6] 2018. General Data Protection Regulation (GDPR) – Official Legal Text. <https://gdpr-info.eu/>
- [7] Martin Abadi, Paul Barham, (...), and Xiaoqiang Zhang. 2016. TensorFlow: A system for large-scale machine learning. In *USENIX Symposium on Operating Systems Design and Implementation*. <https://api.semanticscholar.org/CorpusID:6287870>
- [8] Idris Adjerid, Alessandro Acquisti, and George Loewenstein. [n. d.]. Choice Architecture, Framing, and Cascaded Privacy Choices. 65, 5 ([n. d.]), 2267–2290. <https://pubsonline.informs.org/doi/pdf/10.1287/mnsc.2018.3028>
- [9] Leah Ajmani, Logan Stapleton, Mo Houtti, and Stevie Chancellor. 2024. Data Agency Theory: A Precise Theory of Justice for AI Applications. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*. 631–641.
- [10] Adam J Andreotta, Nin Kirkham, and Marco Rizzi. 2022. AI, big data, and the future of consent. *AI & Society* 37, 4 (2022), 1715–1728.
- [11] Anmol Arora, Joseph E Alderman, Joanne Palmer, Shaswath Ganapathi, Elinor Laws, Melissa D McCradden, Lauren Oakden-Rayner, Stephen R Pföhl, Marzyeh Ghassemi, Francis McKay, et al. 2023. The value of standards for health datasets in artificial intelligence-based applications. *Nature Medicine* 29, 11 (2023), 2929–2938.
- [12] Jad Asswad and Jorge Marx Gómez. [n. d.]. Data Ownership: A Survey. 12, 465 ([n. d.]). <https://www.mdpi.com/2078-2489/12/11/465>
- [13] Leonardo Banh and Gero Strobel. 2023. Generative artificial intelligence. *Electronic Markets* 33 (2023), 1–17. <https://api.semanticscholar.org/CorpusID:265675536>
- [14] Eshta Bhardwaj, Harshit Gujral, Siyi Wu, Ciara Zogheib, Tegan Maharaj, and Christoph Becker. 2024. Machine learning data practices through a data curation lens: An evaluation framework. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*. 1055–1067.
- [15] Fengxiang Bie, Yibo Yang, Zhongzhu Zhou, Adam Ghanem, Minjia Zhang, Zhewei Yao, Xiaoxia Wu, Connor Holmes, Pareesa Ameneh Golnari, David A. Clifton, Yuxiong He, Dacheng Tao, and Shuaiwen Leon Song. 2023. RenAIssance: A Survey into AI Text-to-Image Generation in the Era of Large Model. *ArXiv abs/2309.00810* (2023). <https://api.semanticscholar.org/CorpusID:265821110>
- [16] Abeba Birhane, Vinay Uday Prabhu, and Emmanuel Kahembwe. 2021. Multimodal datasets: misogyny, pornography, and malignant stereotypes. *arXiv preprint arXiv:2110.01963* (2021).
- [17] Abeba Birhane, Vinay Uday Prabhu, and Emmanuel Kahembwe. 2021. Multimodal datasets: misogyny, pornography, and malignant stereotypes. *arXiv preprint arXiv:2110.01963* (2021).
- [18] Matyas Bohacek. 2023. The Unseen A+ Student: Navigating the Impact of Large Language Models in the Classroom. <https://api.semanticscholar.org/CorpusID:262939886>
- [19] Matyas Bohacek and Hany Farid. 2022. Protecting President Zelenskyy against deep fakes. *arXiv preprint arXiv:2206.12043* (2022).
- [20] Matyas Bohacek and Hany Farid. 2022. Protecting world leaders against deep fakes using facial, gestural, and vocal mannerisms. *Proceedings of the National Academy of Sciences* 119, 48 (2022), e2216035119.
- [21] Tom B. Brown, Benjamin Mann, (...), and Dario Amodei. 2020. Language Models are Few-Shot Learners. *ArXiv abs/2005.14165* (2020). <https://api.semanticscholar.org/CorpusID:218971783>
- [22] Denise Carter. 2020. Regulation and ethics in artificial intelligence and machine learning technologies: Where are we now? Who is responsible? Can the information professional play a role? *Business Information Review* 37, 2 (2020), 60–68.
- [23] J. Elliott Casal and Matthew Kessler. 2023. Can linguists distinguish between ChatGPT/AI and human writing?: A study of research ethics and academic publishing. *Research Methods in Applied Linguistics* (2023). <https://api.semanticscholar.org/CorpusID:260713371>
- [24] Fred H. Cate. 2006. The Failure of Fair Information Practice Principles. <https://papers.ssrn.com/abstract=1156972>
- [25] Ang Chen, Yang Wu, Andreas Haeberlen, Boon Thau Loo, and Wenchao Zhou. 2017. Data Provenance at Internet Scale: Architecture, Experiences, and the Road Ahead. In *Conference on Innovative Data Systems Research*. <https://api.semanticscholar.org/CorpusID:559852>
- [26] Chun-Wei Chiang and Ming Yin. 2022. Exploring the Effects of Machine Learning Literacy Interventions on Laypeople’s Reliance on Machine Learning Models. *27th International Conference on Intelligent User Interfaces* (2022). <https://api.semanticscholar.org/CorpusID:247585115>
- [27] Timothy Chu, Zhao Song, and Chiwun Yang. 2024. How to Protect Copyright Data in Optimization of Large Language Models?. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 17871–17879.
- [28] Ilya V. Chugunkov, Dmitry V. Kabak, Viktor N. Vyunnikov, and Roman E. Aslanov. 2018. Creation of datasets from open sources. *2018 IEEE Conference of Russian Young Researchers in Electrical and Electronic Engineering (EIConRus)* (2018), 295–297. <https://api.semanticscholar.org/CorpusID:3899704>
- [29] Ignacio Cofone. 2021. Beyond data ownership. *Cardozo L. Rev.* 43 (2021), 501.
- [30] Rodrigo Tavares Coimbra, Antônio Resende, and Ricardo Terra. 2018. A correlation analysis between halstead complexity measures and other software measures. In *2018 XLIV Latin American Computer Conference (CLEI)*. IEEE, 31–39.
- [31] Zicun Cong, Xuan Luo, Jian Pei, Feida Zhu, and Yong Zhang. 2021. Data pricing in machine learning pipelines. *Knowledge and Information Systems* 64 (2021), 1417 – 1455. <https://api.semanticscholar.org/CorpusID:237194666>
- [32] Valerio Cosentino, Javier Luis, and Jordi Cabot. 2016. Findings from GitHub: methods, datasets and limitations. In *Proceedings of the 13th International Conference on Mining Software Repositories*. 137–141.
- [33] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. ImageNet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*. IEEE, 248–255.
- [34] Emily Denton, Alex Hanna, Razvan Amironesei, Andrew Smart, and Hilary Nicole. 2021. On the genealogy of machine learning datasets: A critical history of ImageNet. *Big Data & Society* 8, 2 (2021), 20539517211035955.
- [35] Nora A Draper and Joseph Turow. [n. d.]. The corporate cultivation of digital resignation. 21, 8 ([n. d.]), 1824–1839. <https://doi.org/10.1177/1461444819833331> Publisher: SAGE Publications.
- [36] Yanai Elazar, Akshita Bhagia, Ian Magnusson, Abhilasha Ravchander, Dustin Schwenk, Alane Suhr, Pete Walsh, Dirk Groeneveld, Luca Soldaini, Sameer Singh, et al. 2023. What’s In My Big Data? *arXiv preprint arXiv:2310.20707* (2023).
- [37] Shamaria Ingram, Tyler Kaczmarek, Alice Lee, and David Bigelow. 2021. Proactive Provenance Policies for Automatic Cryptographic Data Centric Security. In *International Provenance and Annotation Workshop*. <https://api.semanticscholar.org/CorpusID:235266153>
- [38] Danielle Gonzalez, Thomas Zimmermann, and Nachiappan Nagappan. 2020. The State of the ML-universe: 10 Years of Artificial Intelligence & Machine Learning Software Development on GitHub. *2020 IEEE/ACM 17th International Conference on Mining Software Repositories (MSR)* (2020), 431–442. <https://api.semanticscholar.org/CorpusID:220117963>
- [39] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. 2014. Generative adversarial networks. *Commun. ACM* 63 (2014), 139 – 144. <https://api.semanticscholar.org/CorpusID:1033682>
- [40] Alon Y. Halevy, Peter Norvig, and Fernando C Pereira. 2009. The Unreasonable Effectiveness of Data. *IEEE Intelligent Systems* 24 (2009), 8–12. <https://api.semanticscholar.org/CorpusID:14300215>
- [41] Maurice H Halstead. 1977. *Elements of Software Science (Operating and programming systems series)*. Elsevier Science Inc.
- [42] Awni Hannun, Jagrit Digani, Angelos Katharopoulos, and Ronan Collobert. 2023. *MLX: Efficient and flexible machine learning on Apple silicon*. <https://github.com/ml-explore>
- [43] Amy Kathleen Heger, Elizabeth B. Marquis, Mihaela Vorvoreanu, Hanna M. Wallach, and Jenn Wortman Vaughan. 2022. Understanding Machine Learning Practitioners’ Data Documentation Perceptions, Needs, Challenges, and Desiderata. *Proceedings of the ACM on Human-Computer Interaction* 6 (2022), 1 – 29. <https://api.semanticscholar.org/CorpusID:249431472>
- [44] Amy K Heger, Liz B Marquis, Mihaela Vorvoreanu, Hanna Wallach, and Jennifer Wortman Vaughan. 2022. Understanding machine learning practitioners’ data documentation perceptions, needs, challenges, and desiderata. *Proceedings of the ACM on Human-Computer Interaction* 6, CSCW2 (2022), 1–29.
- [45] Patrik Hummel, Matthias Braun, and Peter Dabrock. [n. d.]. Own Data? Ethical Reflections on Data Ownership. ([n. d.]). <https://doi.org/10.1007/s13347-020-00404-9>
- [46] Ben Hutchinson, Andrew Smart, A. Hanna, Emily L. Denton, Christina Greer, Oddur Kjartansson, Parker Barnes, and Margaret Mitchell. 2020. Towards Accountability for Machine Learning Datasets: Practices from Software Engineering and Infrastructure. *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency* (2020). <https://api.semanticscholar.org/CorpusID:225067460>
- [47] Yasamin Jafarian and Hyun Soo Park. 2021. Learning high fidelity depths of dressed humans by watching social media dance videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12753–12762.

- [48] Nal Kalchbrenner, Erich Elsen, Karen Simonyan, Seb Noury, Norman Casagrande, Edward Lockhart, Florian Stimberg, Aäron van den Oord, Sander Dieleman, and Koray Kavukcuoglu. 2018. Efficient Neural Audio Synthesis. In *International Conference on Machine Learning*. <https://api.semanticscholar.org/CorpusID:3524525>
- [49] Tero Karras, Samuli Laine, and Timo Aila. 2018. A Style-Based Generator Architecture for Generative Adversarial Networks. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2018), 4396–4405. <https://api.semanticscholar.org/CorpusID:54482423>
- [50] Levon Khachatryan, Andranik Movsisyan, Vahram Tadevosyan, Roberto Henschel, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. 2023. Text2Video-Zero: Text-to-Image Diffusion Models are Zero-Shot Video Generators. *2023 IEEE/CVF International Conference on Computer Vision (ICCV)* (2023), 15908–15918. <https://api.semanticscholar.org/CorpusID:257687280>
- [51] Bilal Khan and Aamer Nadeem. 2023. Evaluating the effectiveness of decomposed Halstead Metrics in software fault prediction. *PeerJ Computer Science* 9 (2023), e1647.
- [52] Mehtab Khan and Alex Hanna. 2022. The subjects and stages of ai dataset development: A framework for dataset accountability. *Ohio St. Tech. LJ* 19 (2022), 171.
- [53] Moaiad Ahmad Khder. 2021. Web Scraping or Web Crawling: State of Art, Techniques, Approaches and Application. *International Journal of Advances in Soft Computing and its Applications* (2021). <https://api.semanticscholar.org/CorpusID:245584401>
- [54] Bernard Koch, Emily L. Denton, A. Hanna, and Jacob Gates Foster. 2021. Reduced, Reused and Recycled: The Life of a Dataset in Machine Learning Research. *ArXiv abs/2112.01716* (2021). <https://api.semanticscholar.org/CorpusID:244894836>
- [55] Yuki Kondo, Norimichi Ukita, Takayuki Yamaguchi, Hao-Yu Hou, Mu-Yi Shen, Chia-Chi Hsu, En-Ming Huang, Yu-Chen Huang, Yu-Cheng Xia, Chien-Yao Wang, et al. 2023. Mva2023 small object detection challenge for spotting birds: Dataset, methods, and results. In *2023 18th International Conference on Machine Vision and Applications (MVA)*. IEEE, 1–11.
- [56] Martyna Kusak. 2022. Quality of data sets that feed AI and big data applications for law enforcement. In *ERA Forum*, Vol. 23. Springer, 209–219.
- [57] Ignasi Labastida and Thomas Margoni. 2020. Licensing FAIR data for reuse. *Data Intelligence* 2, 1-2 (2020), 199–207.
- [58] Amanda Levendowski. 2018. How copyright law can fix artificial intelligence's implicit bias problem. *Wash. L. Rev.* 93 (2018), 579.
- [59] Quentin Lhoest, Albert Villanova del Moral, Yacine Jernite, Abhishek Thakur, Patrick von Platen, Suraj Patil, Julien Chaumond, Mariama Drame, Julien Plu, Lewis Tunstall, et al. 2021. Datasets: A community library for natural language processing. *arXiv preprint arXiv:2109.02846* (2021).
- [60] Hanlin Li. 2023. Data scraping makes AI systems possible, but at whose expense? <https://www.techpolicy.press/data-scraping-makes-ai-systems-possible-but-at-whose-expense/>
- [61] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*. Springer, 740–755.
- [62] Chang Liu and Han Yu. 2021. AI-Empowered Persuasive Video Generation: A Survey. *Comput. Surveys* 55 (2021), 1 – 31. <https://api.semanticscholar.org/CorpusID:245329411>
- [63] Shayne Longpre, Robert Mahari, Anthony Chen, Naana Obeng-Marnu, Damien Sileo, William Brannon, Niklas Muennighoff, Nathan Khazam, Jad Kabbara, Kartik Perisetla, et al. 2023. The data provenance initiative: A large scale audit of dataset licensing & attribution in ai. *arXiv preprint arXiv:2310.16787* (2023).
- [64] Barbara Magagna, Doron Goldfarb, Paul Martin, Malcolm Atkinson, Spiros Koulouzis, and Zhiming Zhao. 2020. Data provenance. In *Towards Interoperable Research Infrastructures for Environmental and Earth Sciences: A Reference Model Guided Approach for Common Challenges*. Springer, 208–225.
- [65] Gonzalo Martínez, Lauren Watson, Pedro Reviriego, José Alberto Hernández, Marc Juárez, and Rik Sarkar. 2023. Towards understanding the interplay of generative artificial intelligence and the Internet. In *International Workshop on Epistemic Uncertainty in Artificial Intelligence*. Springer, 59–73.
- [66] Nestor Maslej, Loredana Fattorini, Erik Brynjolfsson, John Etchemendy, Katrina Ligett, Terah Lyons, James Manyika, Helen Ngo, Juan Carlos Niebles, Vanessa Parli, Yoav Shoham, Russell Wald, Jack Clark, and Ray Perrault. 2023. Artificial Intelligence Index Report 2023. *ArXiv abs/2310.03715* (2023).
- [67] Anastasia Natsiou and Seán O'Leary. 2021. Audio representations for deep learning in sound synthesis: A review. *2021 IEEE/ACS 18th International Conference on Computer Systems and Applications (AICCSA)* (2021), 1–8. <https://api.semanticscholar.org/CorpusID:245827795>
- [68] Sophie J. Nightingale and Hany Farid. 2022. AI-synthesized faces are indistinguishable from real faces and more trustworthy. *Proceedings of the National Academy of Sciences of the United States of America* 119 (2022). <https://api.semanticscholar.org/CorpusID:246827447>
- [69] Jonas Oppenlaender, Aku Visuri, Ville Paananen, Rhema Linder, and Johanna Silvennoinen. 2023. Text-to-Image Generation: Perceptions and Realities. *arXiv preprint arXiv:2303.13530* (2023).
- [70] Bofeng Pan, Natalia Stakhanova, and Suprio Ray. 2023. Data Provenance in Security and Privacy. *Comput. Surveys* 55 (2023), 1 – 35. <https://api.semanticscholar.org/CorpusID:258259369>
- [71] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Neural Information Processing Systems*. <https://api.semanticscholar.org/CorpusID:202786778>
- [72] Amandalynne Paullada, Inioluwa Deborah Raji, Emily M. Bender, Emily L. Denton, and A. Hanna. 2020. Data and its (dis)contents: A survey of dataset development and use in machine learning research. *Patterns* 2 (2020). <https://api.semanticscholar.org/CorpusID:228084012>
- [73] Amandalynne Paullada, Inioluwa Deborah Raji, Emily M. Bender, Emily L. Denton, and A. Hanna. 2020. Data and its (dis)contents: A survey of dataset development and use in machine learning research. *Patterns* 2 (2020). <https://api.semanticscholar.org/CorpusID:228084012>
- [74] Neoklis Polyzotis, Sudip Roy, Steven Euijong Whang, and Martin Zinkevich. 2018. Data lifecycle challenges in production machine learning: a survey. *ACM SIGMOD Record* 47, 2 (2018), 17–28.
- [75] Gopi Krishnan Rajbahadur, Erika Tuck, Li Zi, Dayi Lin, Boyuan Chen, Zhen Ming, Daniel M German, et al. 2021. Can I use this publicly available dataset to build commercial AI software?—A Case Study on Publicly Available Image Datasets. *arXiv preprint arXiv:2111.02374* (2021).
- [76] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 10684–10695.
- [77] Leonard Rosenthal. 2022. C2PA: the world's first industry standard for content provenance. In *Applications of Digital Image Processing XLV*, Vol. 12226. SPIE, 122260P.
- [78] Leonard Rosenthal, Andy Parsons, Eric Scouten, Jatin Aythor, Bruce MacCormack, Paul England, Marc Levallee, Jonathan Dotan, Sherif Hanna, Hany Farid, et al. 2020. The content authenticity initiative: Setting the standard for digital content attribution. *Adobe Whitepaper* (2020).
- [79] Nithya Sambasivan, Shivani Kapania, Hannah Highfill, Diana Akrong, Praveen K. Paritosh, and Lora Aroyo. 2021. "Everyone wants to do the model work, not the data work": Data Cascades in High-Stakes AI. *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems* (2021). <https://api.semanticscholar.org/CorpusID:231829607>
- [80] Morgan Klaus Scheuerman, Emily L. Denton, and A. Hanna. 2021. Do Datasets Have Politics? Disciplinary Values in Computer Vision Dataset Development. *Proceedings of the ACM on Human-Computer Interaction* 5 (2021), 1 – 37. <https://api.semanticscholar.org/CorpusID:236965639>
- [81] Patrick Schramowski, Manuel Brack, Björn Deiseroth, and Kristian Kersting. 2023. Safe latent diffusion: Mitigating inappropriate degeneration in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 22522–22531.
- [82] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. 2022. LAION-5B: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems* 35 (2022), 25278–25294.
- [83] Mark L Shope. 2021. Lawyer and judicial competency in the era of artificial intelligence: Ethical requirements for documenting datasets and machine learning models. *Geo. J. Legal Ethics* 34 (2021), 191.
- [84] Emily Sidnam-Mauch, Bernat Ivancsics, Ayana Monroe, Evelien Bergrath Washington, Errol Francis, Kelly E. Caine, Joseph Bonneau, and Susan E. McGregor. 2022. Usable Cryptographic Provenance: A Proactive Complement to Fact-Checking for Mitigating Misinformation. In *ICWSM Workshops*. <https://api.semanticscholar.org/CorpusID:249668756>
- [85] Uriel Singer, Adam Polyak, Thomas Hayes, Xiaoyue Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, Devi Parikh, Sonal Gupta, and Yaniv Taigman. 2022. Make-A-Video: Text-to-Video Generation without Text-Video Data. *ArXiv abs/2209.14792* (2022). <https://api.semanticscholar.org/CorpusID:252595919>
- [86] Aditi Singh. 2023. A Survey of AI Text-to-Image and AI Text-to-Video Generators. *2023 4th International Conference on Artificial Intelligence, Robotics and Control (AIRC)* (2023), 32–36. <https://api.semanticscholar.org/CorpusID:264977095>
- [87] Daniel J. Solove. 2012. Introduction: Privacy Self-Management and the Consent Dilemma Symposium: Privacy and Technology. *Harvard Law Review* 126, 7 (2012), 1880–1903. <https://heinonline.org/HOL/P?h=hein.journals/hlr126&i=1910>

- [88] Daniel J. Solove. 2021. The Myth of the Privacy Paradox. *George Washington Law Review* 89, 1 (2021), 1–51. <https://heinonline.org/HOL/P?h=hein.journals/gwlr89&i=15>
- [89] David Thiel, Melissa Stroebel, and Rebecca Portnoff. 2023. Generative ML and CSAM: Implications and Mitigations.
- [90] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. LLaMA: Open and Efficient Foundation Language Models. *ArXiv abs/2302.13971* (2023). <https://api.semanticscholar.org/CorpusID:257219404>
- [91] Karthik Valmeekam, Sarath Sreedharan, Matthew Marquez, Alberto Olmo Hernandez, and Subbarao Kambhampati. 2023. On the Planning Abilities of Large Language Models (A Critical Investigation with a Proposed Benchmark). *ArXiv abs/2302.06706* (2023).
- [92] Stefaan G. Verhulst, Laura Sandor, and Julia Stamm. 2023. The Urgent Need to Reimagine Data Consent. (2023). <https://doi.org/10.48558/TDS9-6Y22>
- [93] Pablo Villalobos, Jaime Sevilla, Lennart Heim, Tamay Besiroglu, Marius Hobbhahn, and An Chang Ho. 2022. Will we run out of data? An analysis of the limits of scaling datasets in Machine Learning. *ArXiv abs/2211.04325* (2022). <https://api.semanticscholar.org/CorpusID:253397775>
- [94] Xiang Wang, Hangjie Yuan, Shiwei Zhang, Dayou Chen, Jiuniu Wang, Yingya Zhang, Yujun Shen, Deli Zhao, and Jingren Zhou. 2023. VideoComposer: Compositional Video Synthesis with Motion Controllability. *ArXiv abs/2306.02018* (2023). <https://api.semanticscholar.org/CorpusID:259075720>
- [95] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed Huai hsin Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022. Emergent Abilities of Large Language Models. *Trans. Mach. Learn. Res.* 2022 (2022).
- [96] Karl Werder, Balasubramaniam Ramesh, and Rongen Zhang. 2022. Establishing Data Provenance for Responsible Artificial Intelligence Systems. *ACM Transactions on Management Information Systems (TMIS)* 13 (2022), 1 – 23. <https://api.semanticscholar.org/CorpusID:247395133>
- [97] Wikipedia. 2024. List of datasets for machine-learning research — Wikipedia, The Free Encyclopedia. <http://en.wikipedia.org/w/index.php?title=List%20of%20datasets%20for%20machine-learning%20research&oldid=1221075088>. [Online; accessed 09-May-2024].
- [98] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, and Jamie Brew. 2019. HuggingFace’s Transformers: State-of-the-art Natural Language Processing. *ArXiv abs/1910.03771* (2019). <https://api.semanticscholar.org/CorpusID:208117506>
- [99] Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Weixian Lei, Yuchao Gu, Wynne Hsu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. 2022. Tune-A-Video: One-Shot Tuning of Image Diffusion Models for Text-to-Video Generation. *2023 IEEE/CVF International Conference on Computer Vision (ICCV)* (2022), 7589–7599. <https://api.semanticscholar.org/CorpusID:254974187>
- [100] Jinlong Xue, Yayue Deng, Yingming Gao, and Ya Li. 2024. Auffusion: Leveraging the Power of Diffusion and Large Language Models for Text-to-Audio Generation. *ArXiv abs/2401.01044* (2024). <https://api.semanticscholar.org/CorpusID:266725678>
- [101] Chenshuang Zhang, Chaoning Zhang, Mengchun Zhang, and In So Kweon. 2023. Text-to-image diffusion model in generative ai: A survey. *arXiv preprint arXiv:2303.07909* (2023).
- [102] Chenshuang Zhang, Chaoning Zhang, Mengchun Zhang, and In-So Kweon. 2023. Text-to-image Diffusion Models in Generative AI: A Survey. *ArXiv abs/2303.07909* (2023). <https://api.semanticscholar.org/CorpusID:257505012>
- [103] Chenshuang Zhang, Chaoning Zhang, Sheng Zheng, Mengchun Zhang, Maryam Qamar, Sung-Ho Bae, and In-So Kweon. 2023. A Survey on Audio Diffusion Models: Text To Speech Synthesis and Enhancement in Generative AI. *ArXiv abs/2303.13336* (2023). <https://api.semanticscholar.org/CorpusID:257913174>
- [104] Dawen Zhang, Boming Xia, Yue Liu, Xiwei Xu, Thong Hoang, Zhenchang Xing, Mark Staples, Qinghua Lu, and Liming Zhu. 2023. Navigating Privacy and Copyright Challenges Across the Data Lifecycle of Generative AI. *ArXiv abs/2311.18252* (2023). <https://api.semanticscholar.org/CorpusID:265506320>
- [105] Wayne Xin Zhao, Kun Zhou, Junyi Li (...), and Ji rong Wen. 2023. A Survey of Large Language Models. *ArXiv abs/2303.18223* (2023). <https://api.semanticscholar.org/CorpusID:257900969>
- [106] Christopher T Zirpoli. 2023. Generative artificial intelligence and copyright law. (2023).

A A. CRS Scale Visuals

COMPLIANCE RATING SCHEME SCORE



COMPLIANCE RATING SCHEME SCORE



Figure 2: Proposed design interface for "A" and "C" score on the CRS scale

B B. Schematic Overviews of the *DatasetSentinel* Use Cases

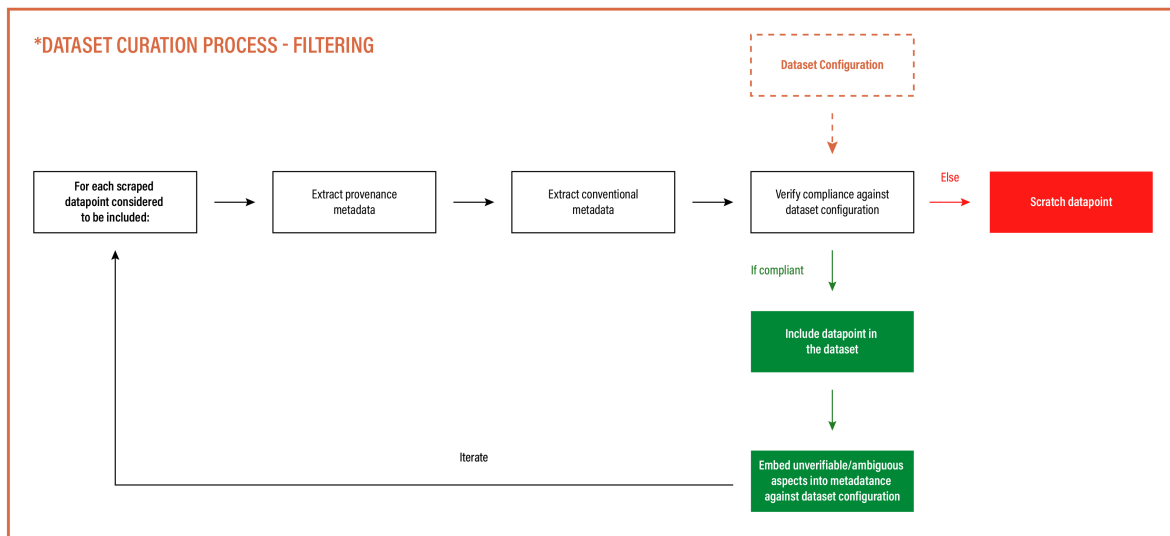


Figure 3: A schematic overview of DatasetSentinel's use case within the dataset curation stage of the dataset lifecycle.

Compliance Rating Scheme:
A Data Provenance Framework for Generative AI Datasets

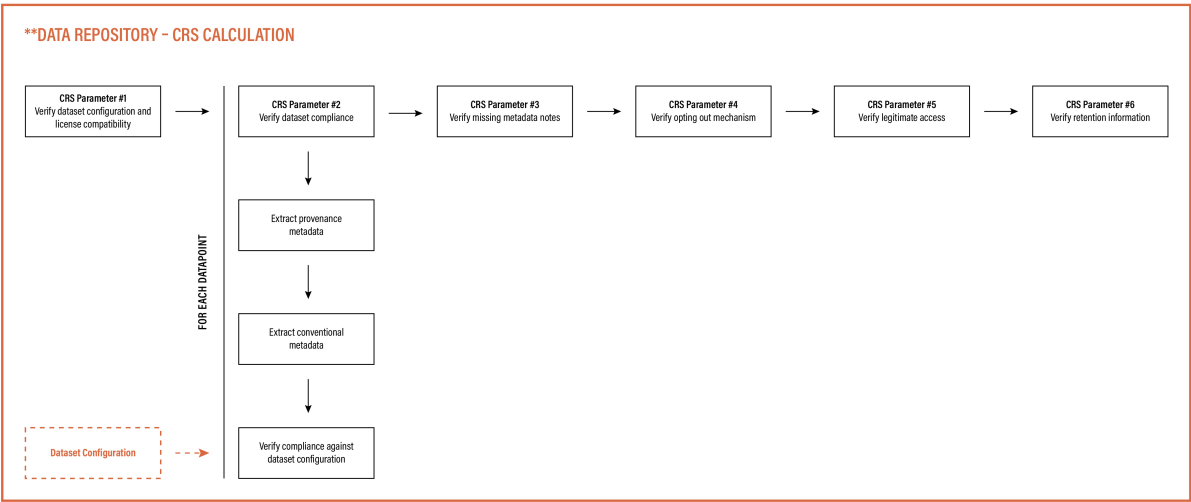


Figure 4: A schematic overview of CRS’ use case within the dataset repository stage of the dataset lifecycle.

C C. Additional CRS Mockups

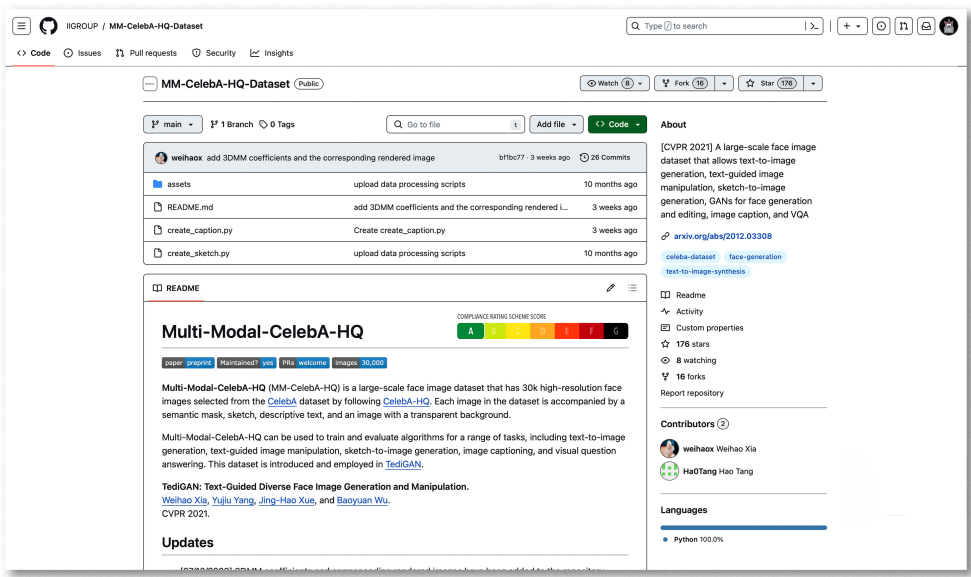


Figure 5: A fictitious "A" CRS score mock-up of a random GitHub dataset

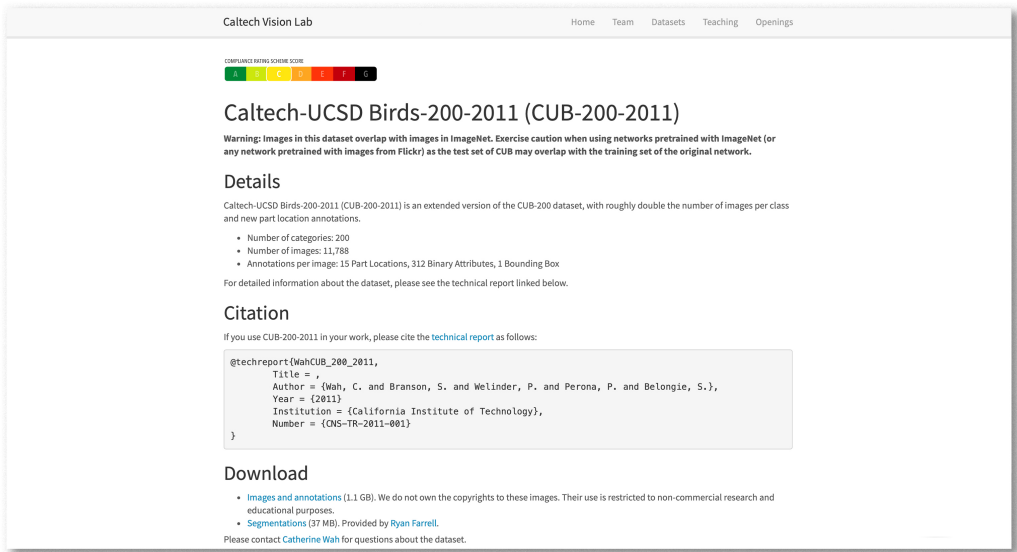


Figure 6: A fictitious CRS "C" score mock-up of a random academic repository dataset

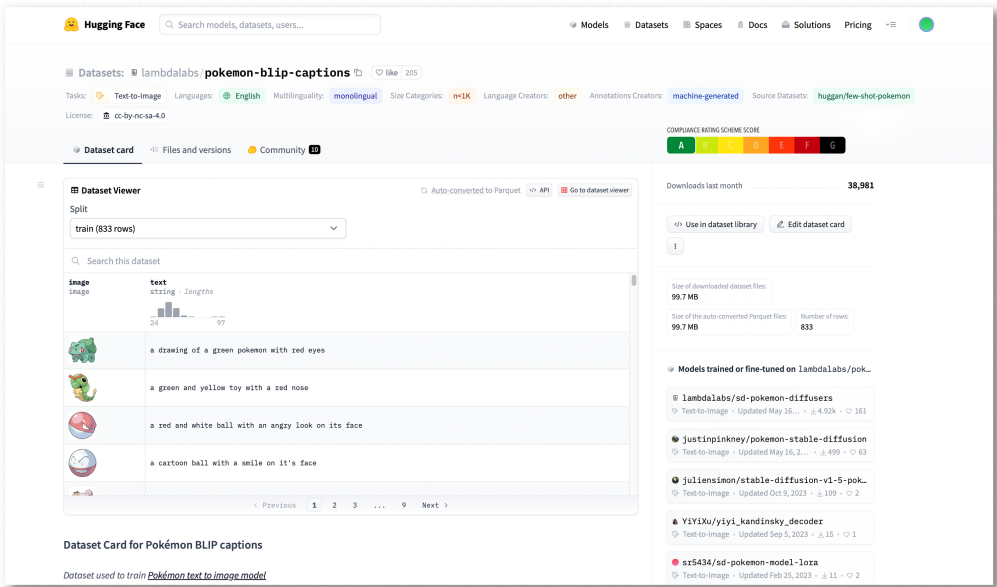


Figure 7: A fictitious CRS "A" score mock-up of a random Hugging Face dataset