

---

# The Data Cards Playbook

A toolkit for purposeful and people-centric dataset documentation for transparency in AI systems.

<https://pair-code.github.io/datacardsplaybook/>

#datacardsplaybook



THE DATA CARDS PLAYBOOK

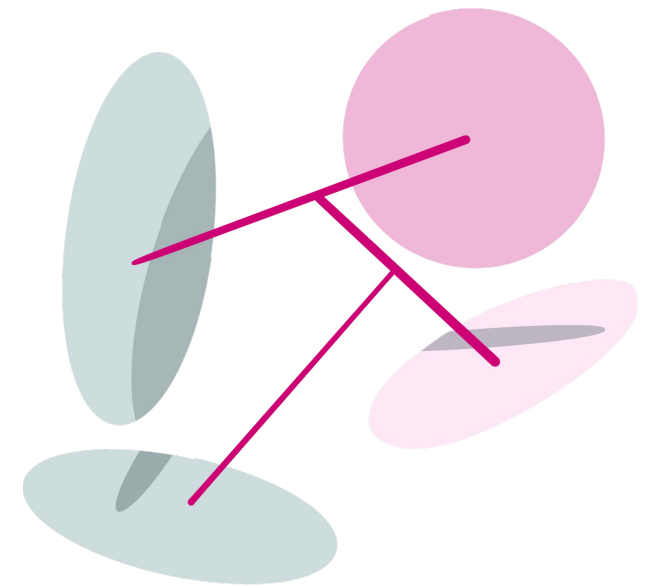
# Introduction

**01 Ask**

**02 Inspect**

**03 Answer**

**04 Audit**



THE DATA CARDS PLAYBOOK

# Your Data Card Brief

IN THIS SECTION

Define the scope, values, and vision for your Data Card(s) effort.

## INSTRUCTIONS

Through a series of critical reflections and discussions, decide how Data Card(s) might help your datasets' and transparency goals.

## OUTCOMES

A brief that describes why you're creating Data Card(s), what success looks like, and what's out of scope.

## ACTIVITY LEVEL

Basic

# Consider the following types of dataset documentation:

## Observable

Shape and size of data, pipelines, access, licenses. Context and information that can be easily acquired from the data.

UTILITARIAN DOCUMENTATION

## Explainable

Documentation ++ Processes, rules, rationales that shape the data. Context and information that cannot be learned from the data.

TRANSPARENT DOCUMENTATION

## Understandable

Transparent documentation **intentionally written for humans** and required for making responsible decisions about dataset use.

RESPONSIBLE DOCUMENTATION

Data Cards summarize critical information about datasets that help people to make informed decisions about how data is used in ML systems for product, policy, and research.

Translate your definition of transparency into the scope and utility for Data Cards.

Open Images Extended - Crowdsourced		Open Images Extended - Crowdsourced intends to capture global representation. This dataset comprises over 478,000 images and associated labels from otherwise under-represented populations. It can be used with Open Images V4.
<b>PUBLISHER(S)</b> Google LLC	<b>INDUSTRY TYPE</b> Corporate - Tech	<b>INTENDED USE CASE(S)</b> <ul style="list-style-type: none"> <li>Identify objects or context of photos visually (e.g., through Lens or Camera)</li> <li>Find objects, plants, animals, etc. through search in Photos or Image Search</li> </ul>
<b>PRIMARY DATA TYPE</b> Image Data	<b>KEY APPLICATION</b> Machine Learning, Object Recognition	<b>NATURE OF CONTENT</b> Labeled images of objects (household goods, commercial products), vehicles, plants, animals and people (faces blurred).
<b>DATASET FUNCTION(S)</b> Training Testing	<b>DATASET CHARACTERISTICS</b> (All numbers are approximate) Total Instances 478k+ Total Classes 6k+ Total Labels 1.27m+ Algorithmically Generated Labels 1.11m+ User Contributed Labels 505k+ Human Verified Labels All labels verified	<b>EXCLUDED DATA</b> All EXIF data including location has been removed  <b>PRIVACY</b> PII associated with human subjects removed
<b>LICENSE TYPE(S)</b> CC-BY-4.0	<b>LAST UPDATED</b> Oct 2018 <b>VERSION</b> 1.0 <b>STATUS</b> Actively Maintained	<b>SUMMARY OF LICENSE PERMISSIONS (CC-BY-4.0)</b> <ul style="list-style-type: none"> <li>You are free to share and adapt</li> <li>Attribution required</li> <li>You cannot apply any additional restrictions</li> </ul> <b>ACCESS COST</b> Open Access
<b>DATA COLLECTION METHOD(S)</b> Crowdsourced	<b>DATA SOURCE(S)</b> <ul style="list-style-type: none"> <li>Contributions by global users of the <a href="#">Crowdsourcing</a> app</li> <li>Vendor data collection efforts</li> </ul>	<b>DATA SELECTION</b> All images are opted-in for open-sourcing by Crowdsourcing app contributors
<b>SAMPLING METHOD(S)</b> Unsampled	<b>GEOGRAPHIC DISTRIBUTION</b> 83% India 2% Vietnam 2% Brazil 1% Israel 1% Nigeria 1% Thailand 1% Colombia 1% UAE 8% Others (each less than 1%)	<b>FILTERING CRITERIA</b> <ul style="list-style-type: none"> <li>PII: Name tags, Unblurred faces, etc.</li> <li>Inappropriate Content</li> <li>Unusable Imagery</li> </ul>
<b>LABELING METHOD(S)</b> Human Labels Algorithmic Labels	<b>LABEL TYPE(S)</b> Human Labels Free-form text labels Algorithmic Labels Additional labels  <b>LABEL SOURCE(S)</b> Human Labels Image owners Algorithmic Labels Google's internal image annotation algorithm	<b>LABELING PROCEDURE - HUMAN</b> Free-form labels are provided by users of the Crowdsourcing app. The user who has taken the picture provides the labels.  <b>LABELING PROCEDURE - ALGORITHMIC</b> Labels are resolved against known entity names from <a href="#">Knowledge Graph</a> . Additional labels are added based on Google's internal image annotation system.
<b>VALIDATION METHOD(S)</b> Human Validated	<b>VALIDATION TASK(S)</b> <ul style="list-style-type: none"> <li>Human validators verify labels</li> <li>Human validators flag PII</li> <li>Human validators filter data</li> </ul> <b>VALIDATOR DESCRIPTION(S)</b> Compensated workers based out of India	<b>VALIDATION POLICY SUMMARY</b> Algorithmic and user contributed labels are verified by human validators based out of India. There is a known overlap in algorithmic and user contributed labels. Validators flag any PII content.



# As a group, answer the following questions:

What are the top 3 goals of your dataset(s)?

What are the top 3 goals of your Data Card(s)?

3 circumstances that describe the experience of transparency that your Data Card(s) should deliver?

3 circumstances in which your Data Card(s) are a wasted investment?



# What are the **top 3 goals** of your dataset(s)?

**Consider:** What motivates you to create them? What does success of your dataset look like? Is it adoption or purchase? Improvement of production systems? Furthering a research agenda?

1. 📝

2. 📝

3. 📝



# What are the **top 3 goals** of your Data Card(s)?

**Consider:** How does the goal of your Data Card support the goals of your dataset(s)? How must a Data Card be used so that it contributes to the success of your dataset(s)?

1. 📝

2. 📝

3. 📝

## 3 circumstances that describe success for your Data Card(s)?

Think about your own experience of transparency, the experience of using dataset documentation, and successful outcomes of transparency in data.

1. 📝

2. 📝

3. 📝

# 3 circumstances in which your Data Card(s) are **a wasted investment?**

Consider obvious and mundane tasks, or easily preventable failures of Data Cards. What type of documentation included in a Data Card might make it worthwhile?

1. 📝

2. 📝

3. 📝

# Your Data Card Brief

## By using Data Cards, we hope to...

✍️ short statements that describe major use cases, objectives and value of Data Cards

## Successful Data Cards look like ...

✍️ explicit statements for measuring the success of implementing Data cards, telling us if they have met our objectives.

## Non-goals of Data Cards are...

✍️ statements about what has been purposefully excluded from your Data Card efforts, and things that should not be related.

# Checklist

YOU SHOULD NOW HAVE DEFINED

–

- The goals and objectives of your dataset(s)
- The goals and objectives of your Data Card(s)
- What's in scope and out of scope for your Data Card(s)
- When your Data Card(s) do not work



#datacardsplaybook



[The Data Cards Playbook](#) <sup>↗</sup> is an adaptable toolkit of participatory activities, conceptual frameworks, and guidance that support Responsible AI practices for transparency in dataset documentation.

If you've adapted, implemented, or have feedback for this guidance, we'd love to hear from you at <https://github.com/pair-code/datacardsplaybook> <sup>↗</sup>.

Find the complete playbook at  
<https://pair-code.github.io/datacardsplaybook> <sup>↗</sup>



The [Data Cards Playbook](#) <sup>↗</sup> by [the People + AI Research Initiative](#) <sup>↗</sup> at [Google Research](#) <sup>↗</sup> is licensed under a Creative Commons Attribution-ShareAlike 4.0 International License. You are free to share and adapt this work under the [appropriate license terms](#) <sup>↗</sup>.