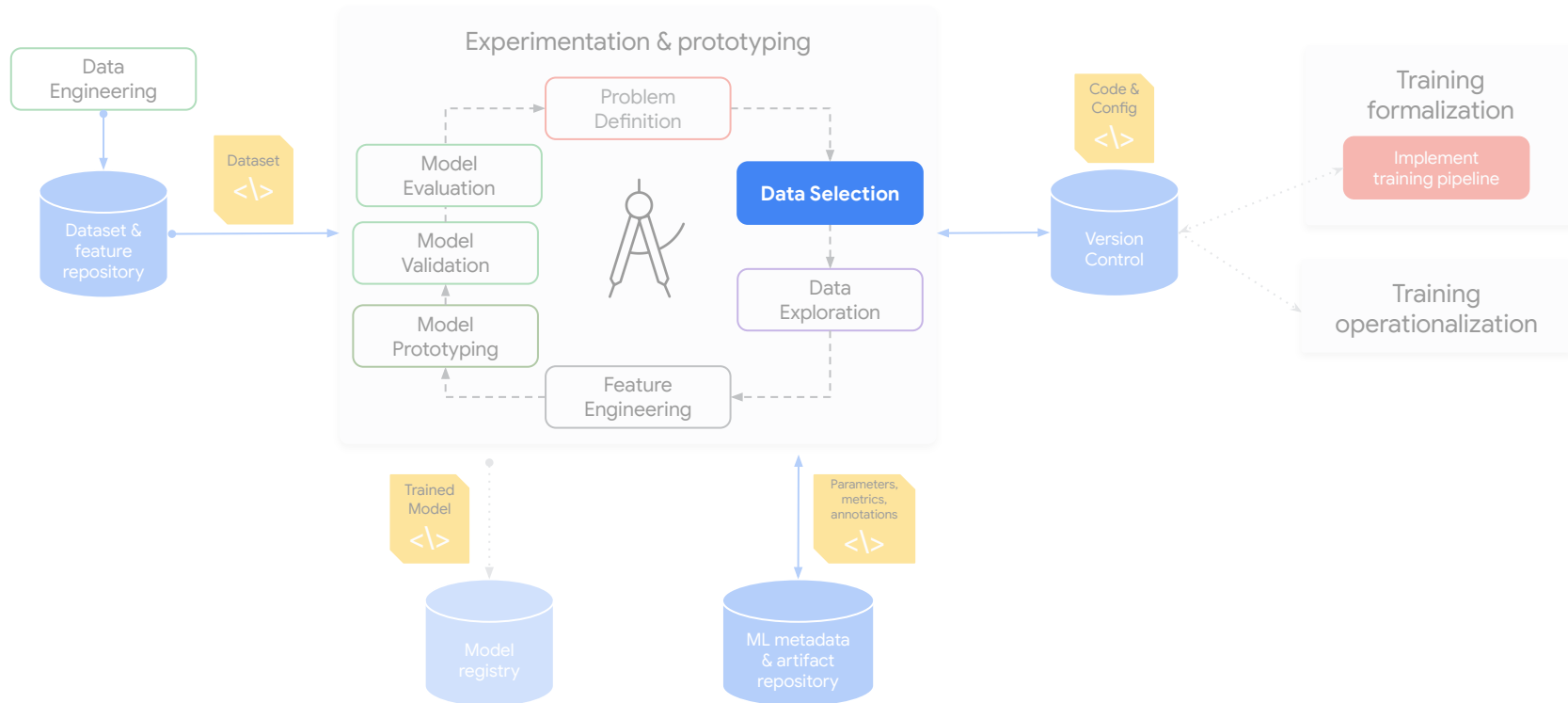# ML Development: Data Selection

# **MLOps:** ML Development

# The MLOps **Personas**

ML
Engineer

ML
Researcher

**Data
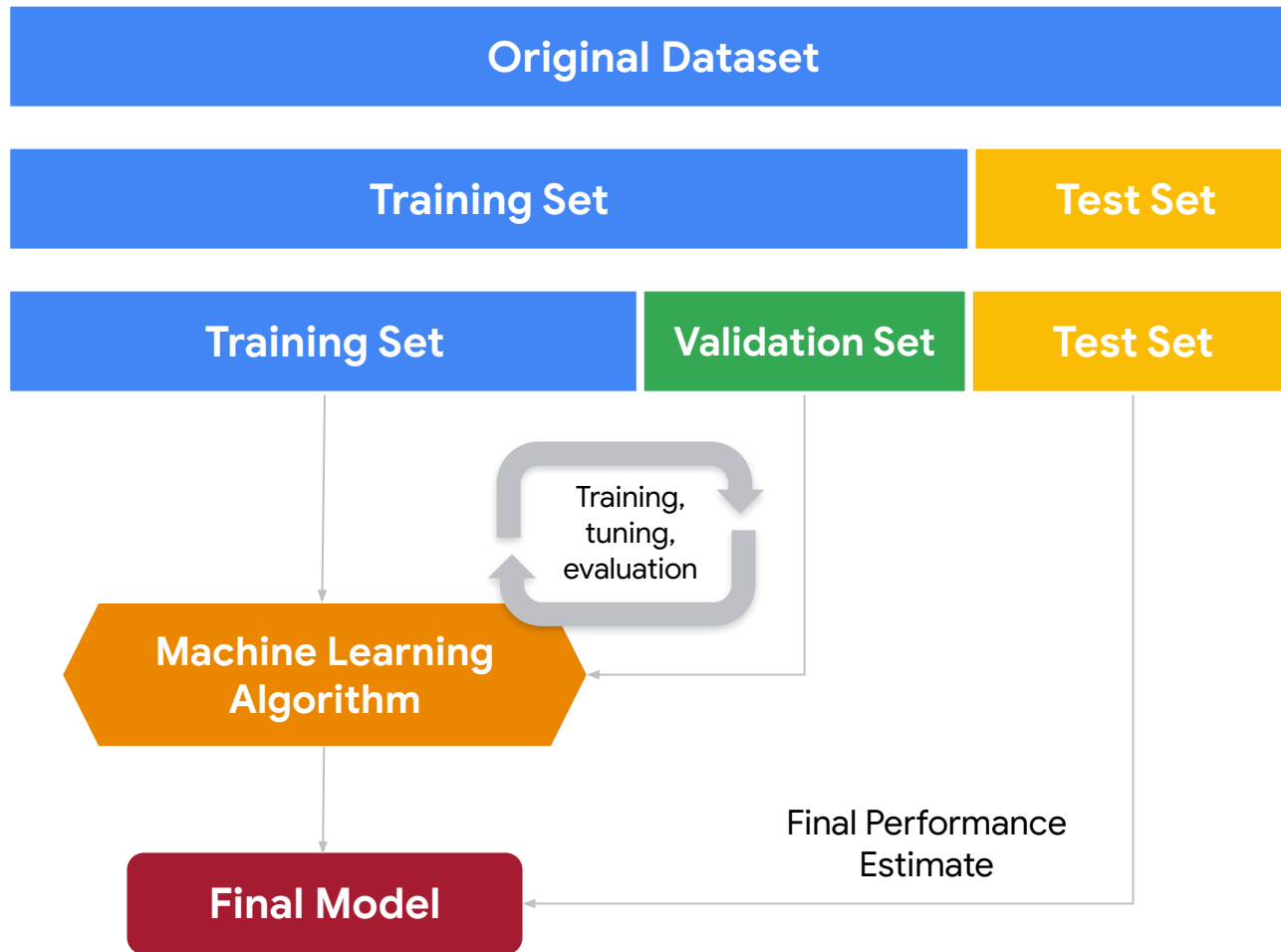Scientist**

**Data
Engineer**

Software
Engineer

DevOps

Business
Analyst

**Original Dataset**

**Training Set** | **Test Set**

**Training Set** | **Validation Set** | **Test Set**

Training, tuning, evaluation

**Machine Learning Algorithm**

**Final Model**

Final Performance Estimate

**Original Dataset**

**Training Set** | **Test Set**

**Training Set** | **Validation Set** | **Test Set**

Training, tuning, evaluation

**Machine Learning Algorithm**

**Final Model**

Final Performance Estimate
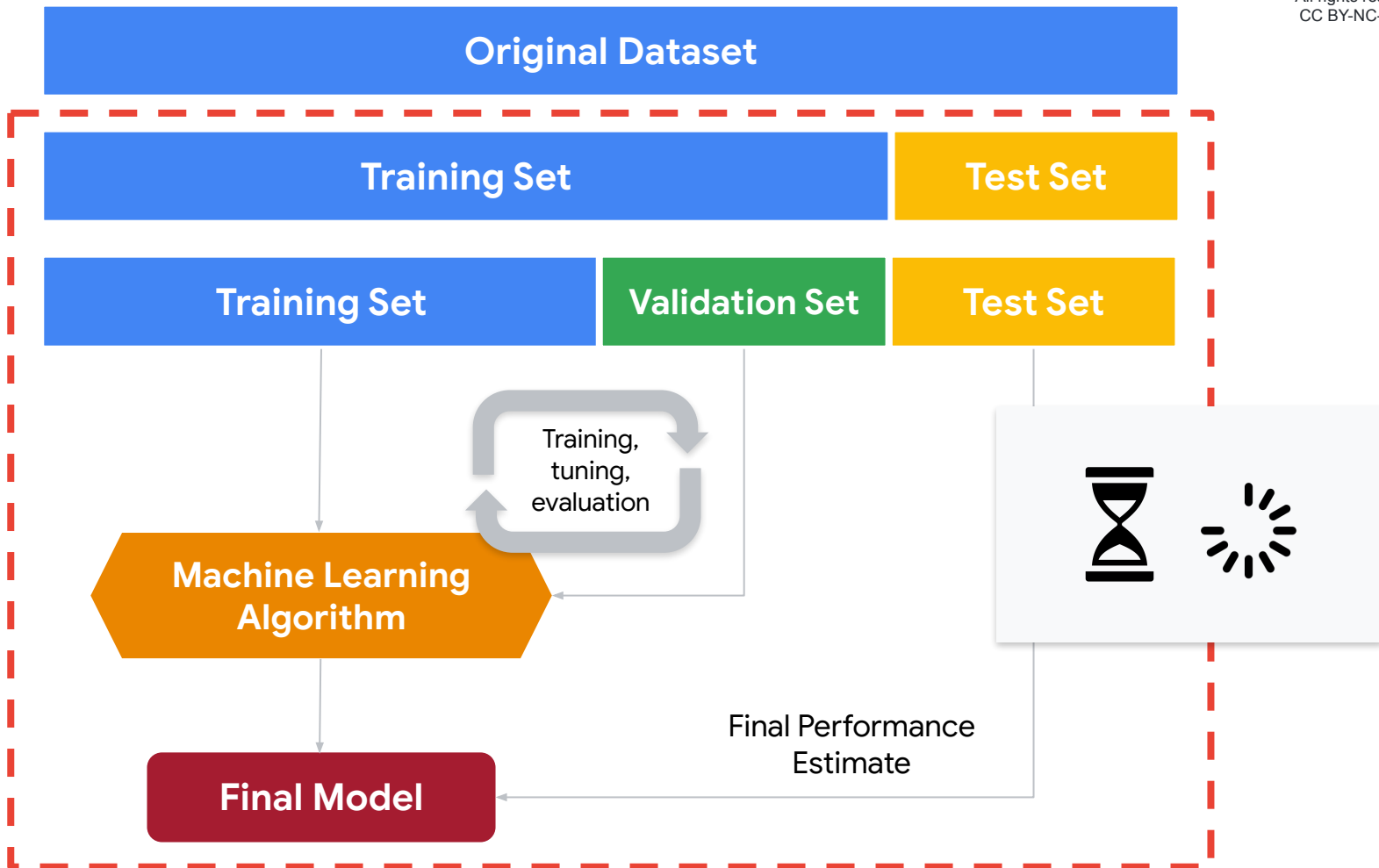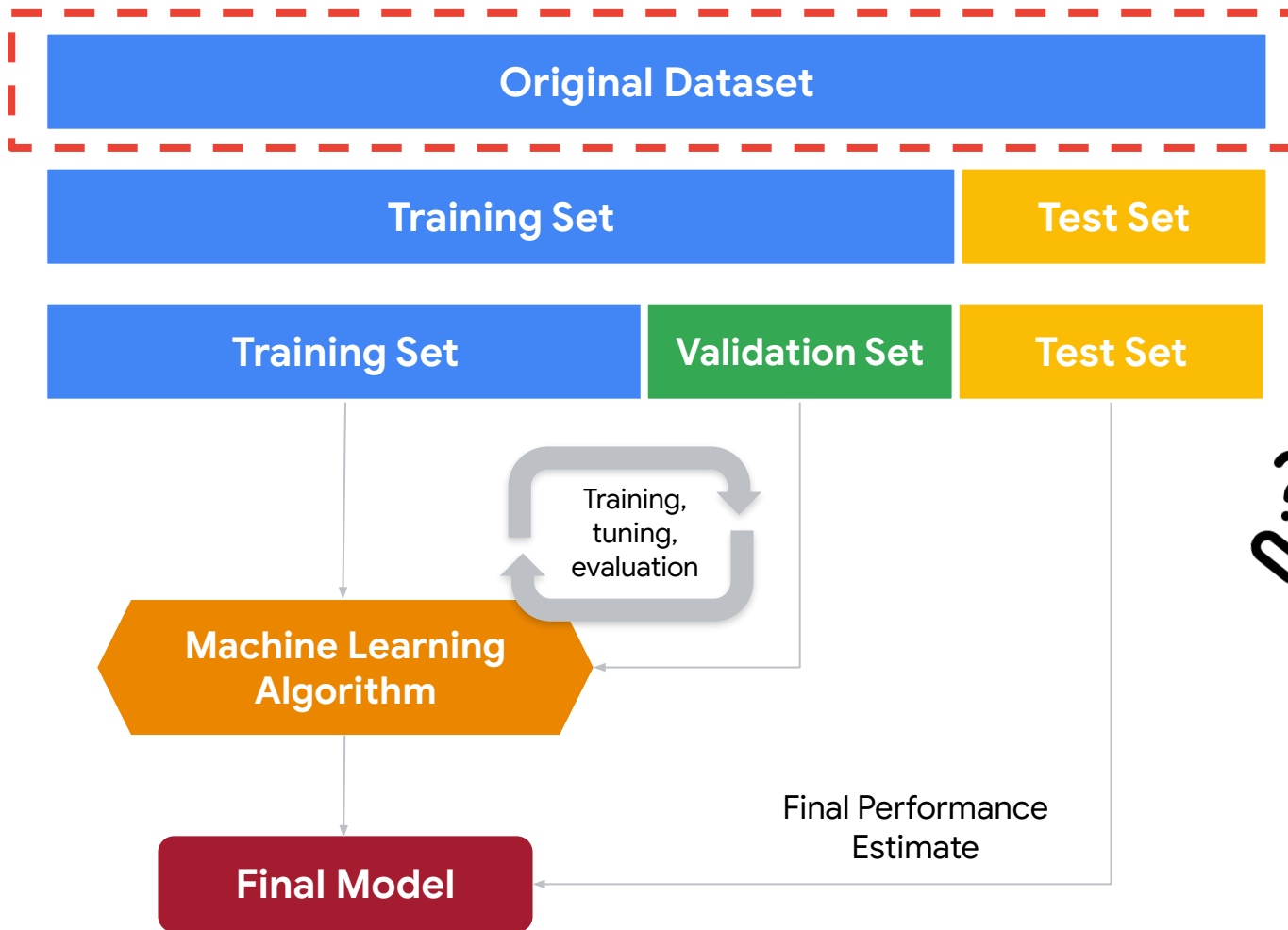
# Dataset Sources

- What **relevant datasets** are available?

*Audio*
*Image*
*Image Classification*
*Object Detection*
*Question Answering*
*Structured*
*Summarization*
*Text*
*Translate*
*Video*

**TensorFlow**
Datasets Catalog

# Dataset Sources

- What **relevant datasets** are available?

*Audio*
*Image*
*Image Classification*
*Object Detection*
*Question Answering*
*Structured*
*Summarization*
*Text*
*Translate*
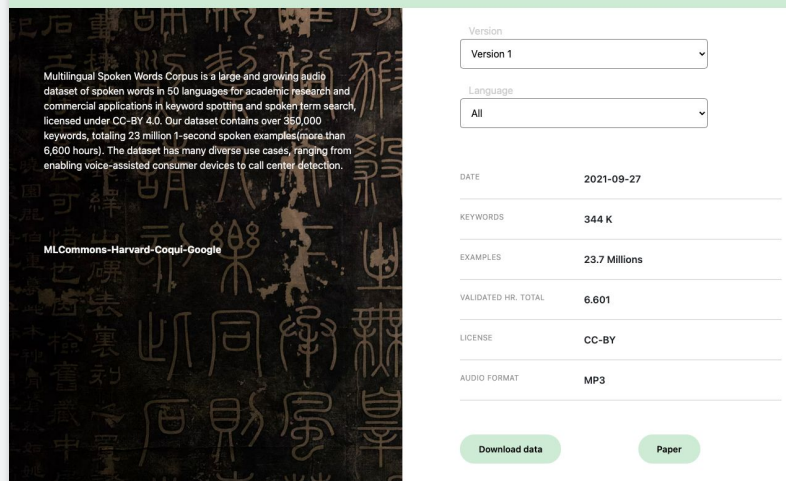*Video*

## Mozilla
## Common Voice

# Dataset Sources

- What **relevant datasets** are available?

*Audio*
*Image*
*Image Classification*
*Object Detection*
*Question Answering*
*Structured*
*Summarization*
*Text*
*Translate*
*Video*



ML
•Commons

**Multilingual Spoken Words Corpus**

Multilingual Spoken Words Corpus is a large and growing audio dataset of spoken words in 50 languages for academic research and commercial applications in keyword spotting and spoken term search, licensed under CC-BY 4.0. Our dataset contains over 350,000 keywords, totaling 23 million 1-second spoken examples(more than 6,600 hours). The dataset has many diverse use cases, ranging from enabling voice-assisted consumer devices to call center detection.

**MLCommons-Harvard-Coqui-Google**

| Version | |
| --- | --- |
| Version 1 | |

| Language | |
| --- | --- |
| All | |

| DATE | 2021-09-27 |
| --- | --- |
| KEYWORDS | 344 K |
| EXAMPLES | 23.7 Millions |
| VALIDATED HR. TOTAL | 6.601 |
| LICENSE | CC-BY |
| AUDIO FORMAT | MP3 |

Download data          Paper

# Dataset Sources

- What **relevant datasets** are available?
- Is this data **accurate and reliable**?

What **problem** are you trying to *solve*?

Both *quantity* and *quality* will influence your model's performance

# Dataset Sources

- What **relevant datasets** are available?

- Is this data **accurate and reliable**?

  ➔ **2,618** volunteers
    - ◆ Variety of accents
  ➔ 1,000+ examples for **each** word
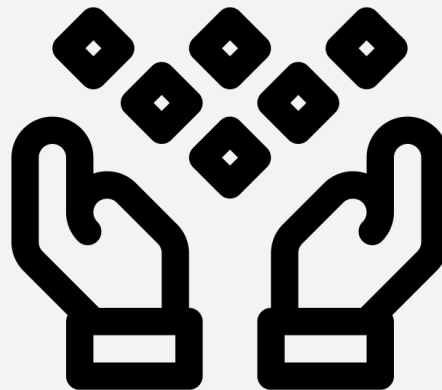  ➔ All 105,829 remaining utterances **manually reviewed** through crowdsourcing

Speech Commands: A Dataset for Limited-Vocabulary Speech Recognition

Pete Warden
Google Brain
Mountain View, California
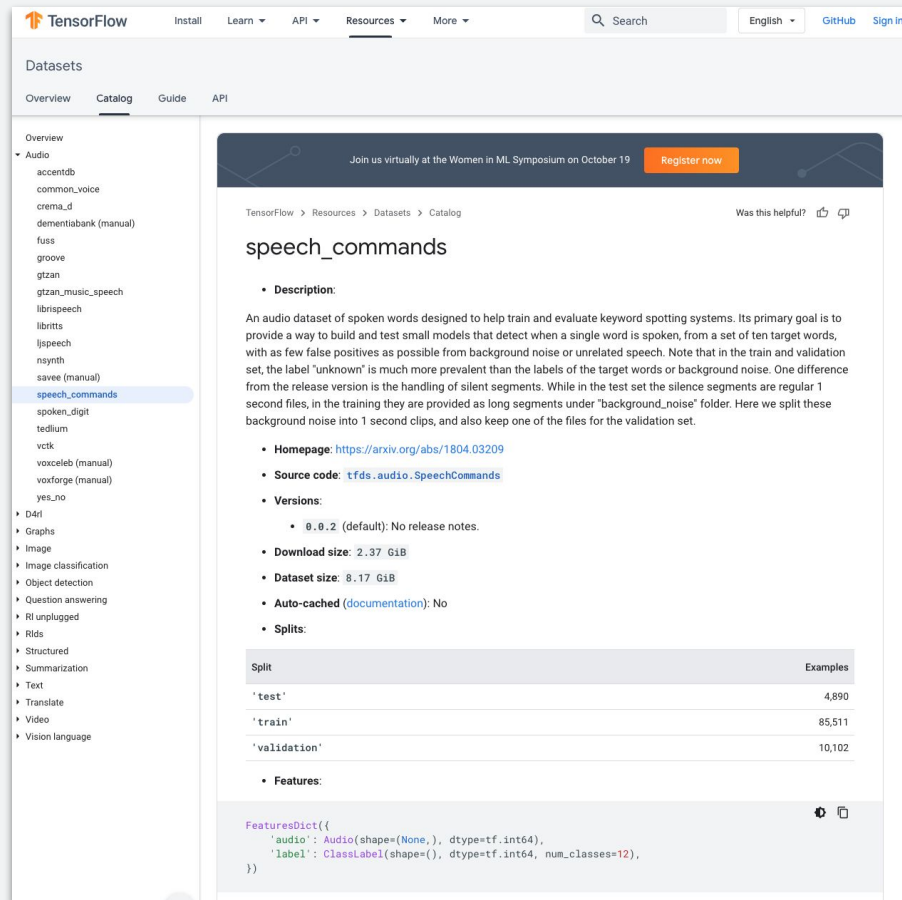petewarden@google.com

April 2018

# Dataset Sources

- What **relevant datasets** are available?
- Is this data **accurate and reliable**?
- Can stakeholders get **access to data**?

# Dataset Sources

- What **relevant datasets** are available?
- Is this data **accurate and reliable**?
- Can stakeholders get **access to data**?

# Dataset Sources

- What **relevant datasets** are available?
- Is this data **accurate and reliable**?
- Can stakeholders get **access to data**?
- What **properties** are in this dataset?

# Dataset Sources

- What **relevant datasets** are available?
- Is this data **accurate and reliable**?
- Can stakeholders get **access to data**?
- What **properties** are in this dataset?
- Does the data need to be **labeled**?

# Dataset Sources

- What **relevant datasets** are available?
- Is this data **accurate and reliable**?
- Can stakeholders get **access to data**?
- What **properties** are in this dataset?
- Does the data need to be **labeled**?
- How will it be **updated** post training?

# Dataset Sources

- What **relevant datasets** are available?
- Is this data **accurate and reliable**?
- Can stakeholders get **access to data**?
- What **properties** are in this dataset?
- Does the data need to be **labeled**?
- How will it be **updated** post training?
- What are the **licenses** and terms of use?

  ➜ Open?
  ➜ Copyrighted?
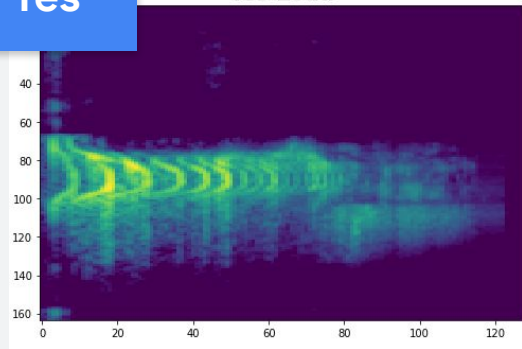  ➜ Licensed?
  ➜ Product users?

# Dataset Sources

- What **relevant datasets** are available?
- Is this data **accurate and reliable**?
- Can stakeholders get **access to data**?
- What **properties** are in this dataset?
- Does the data need to be **labeled**?
- How will it be **updated** post training?
- What are the **licenses** and terms of use?
- Is there **personally identifiable information** (PII) in the dataset?



Security & Privacy

# Dataset Sources

- What **relevant datasets** are available?
- Is this data **accurate and reliable**?
- Can stakeholders get **access to data**?
- What **properties** are in this dataset?
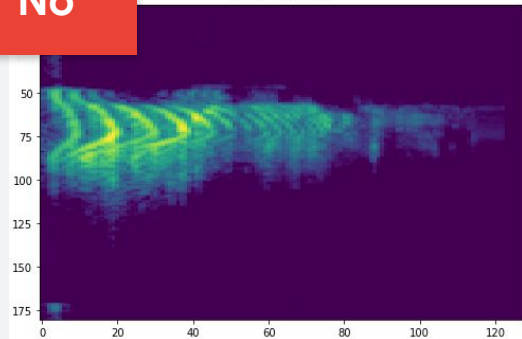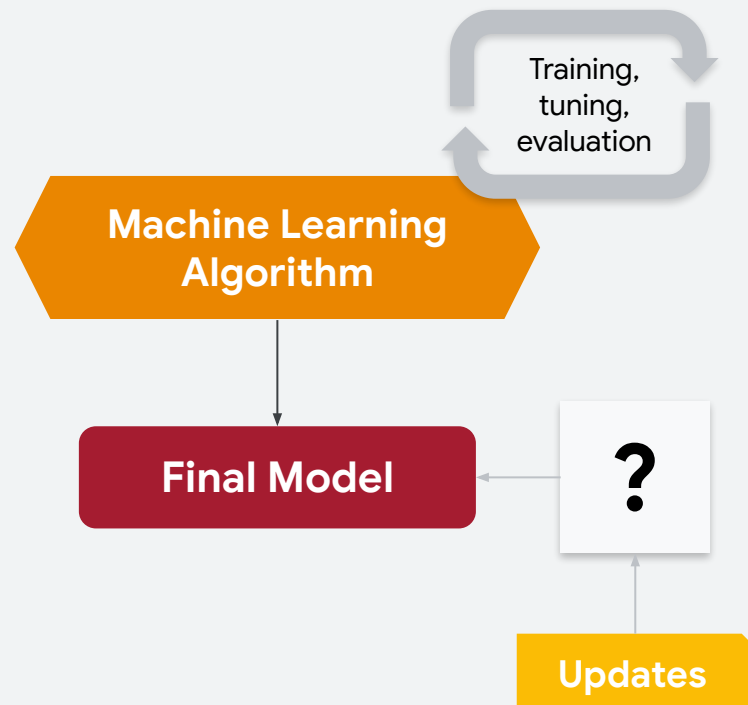- Does the data need to be **labeled**?
- How will it be **updated** post training?
- What are the **licenses** and terms of use?
- Is there **personally identifiable information** (PII) in the dataset?
- Are there **restricted features** that cannot be used for commercial applications?
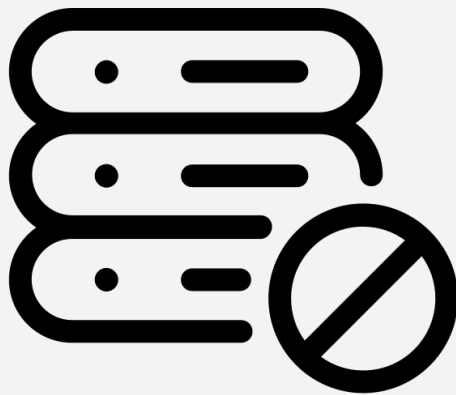
# Dataset Sources

- What **relevant datasets** are available?
- Is this data **accurate and reliable**?
- Can stakeholders get **access to data**?
- What **properties** are in this dataset?
- Does the data need to be **labeled**?
- How will it be **updated** post training?
- What are the **licenses** and terms of use?
- Is there **personally identifiable information** (PII) in the dataset?
- Are there **restricted features** that cannot be used for commercial applications?