# Actionable Data Insights for Machine Learning

Nils Braun

# Agenda

Motivation

Implementation

Case Studies

Summary

# Motivation

Drive model performance with better data understanding

# Interview Learnings

Projects included: Image understanding, machine translation, text understanding, speech recognition

## Enable rapid data experiments

- High-level declarative interfaces
- User defined functions as first class citizen

## Flexibility

- Enable usage across different storage systems
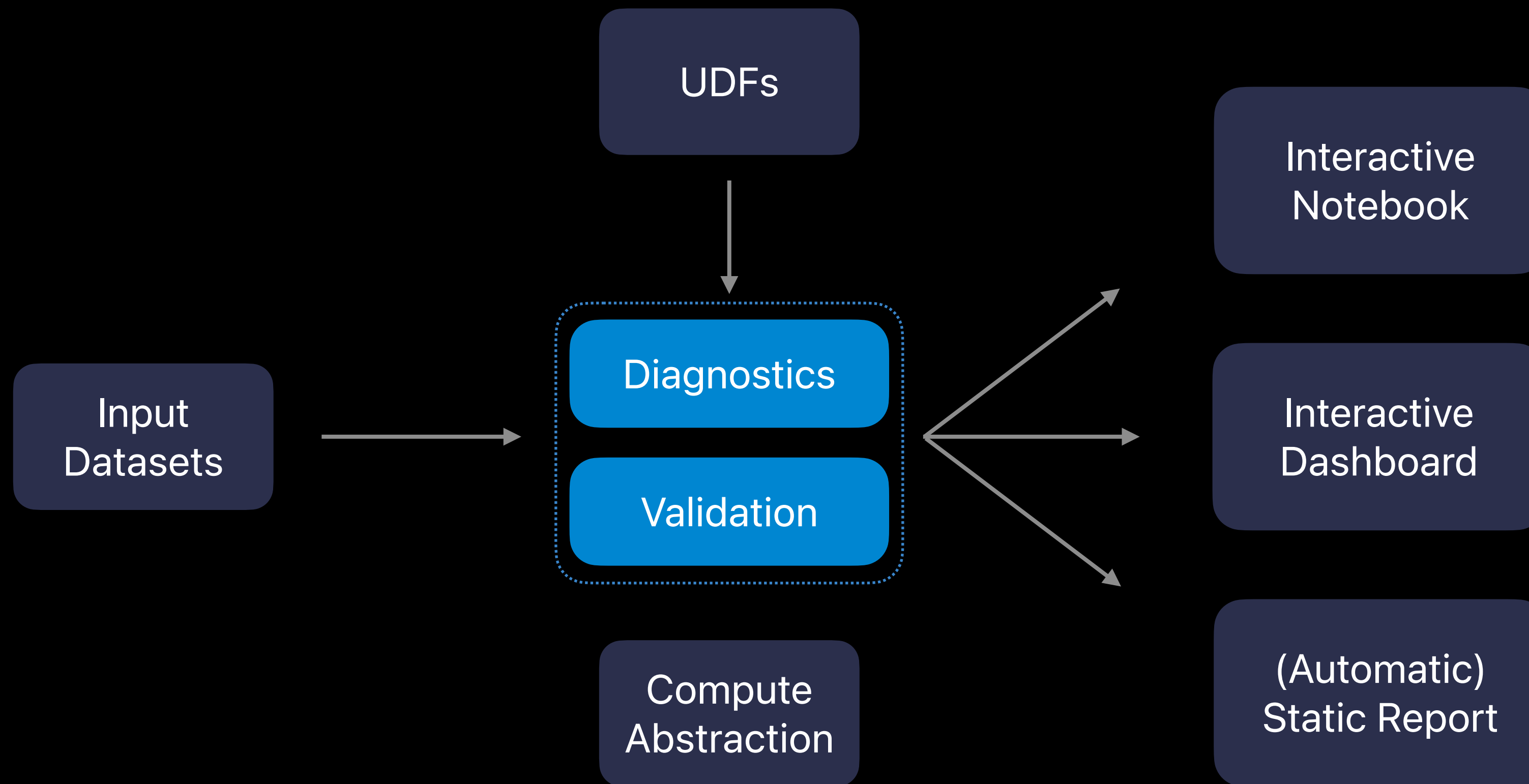- Enable portability from local machines to cloud to devices

## Collaboration & Reproducibility

- Documentation for data and models
- Share findings

# Implementation

# Data Tooling
High-level overview

UDFs

Input Datasets

Diagnostics

Validation

Compute Abstraction

Interactive Notebook

Interactive Dashboard

(Automatic) Static Report

# Implementation
Required components for data insights

## Data Diagnostics

- Interactive roll-up and drill-down
- Raw data preview

## Data Validation
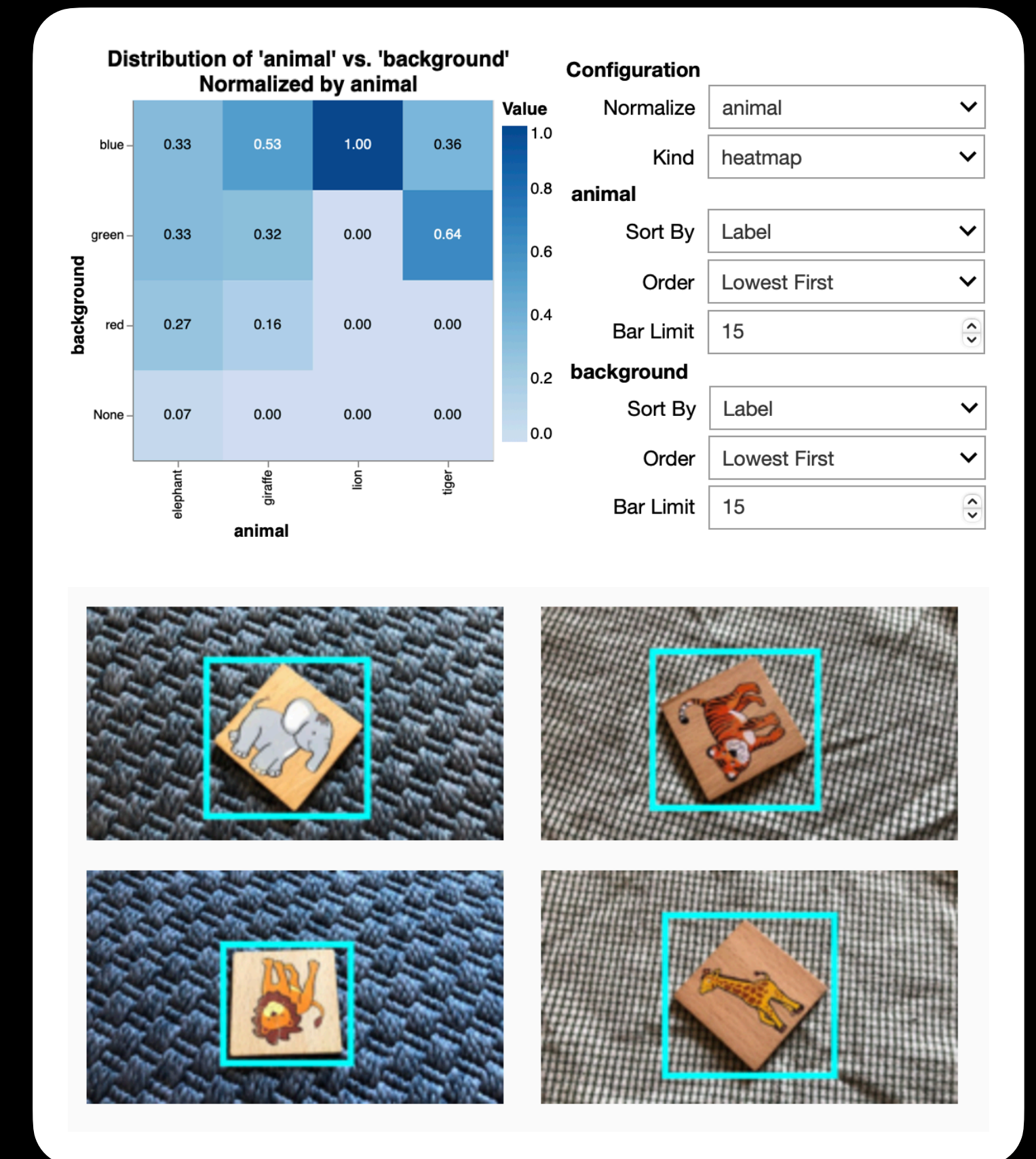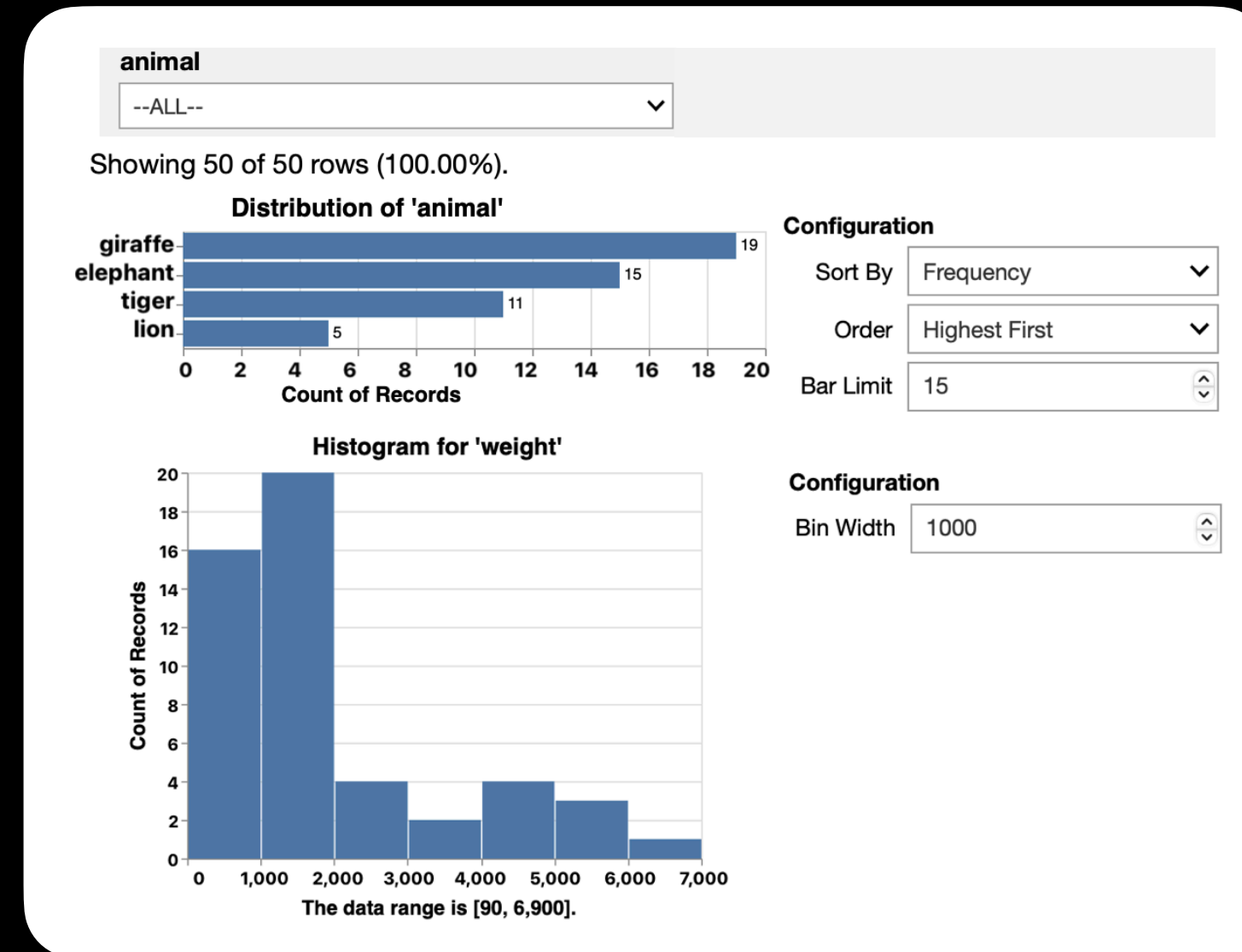
- Define constraints on sample or aggregate level

- Data Access abstraction enables loading data from many sources
- Compute abstraction enables portability

# Data Diagnostics

```python
def show_sample(row, **kwargs):
    # visualise image with bounding box

sample_view = SampleView(
    details [CustomDetail(show_sample)]
)

Explorer(
    filters ["animal"],
    views [
        "animal",
        "weight",
        ("animal", "background"),
        sample_view
    ]
).show(data)
```
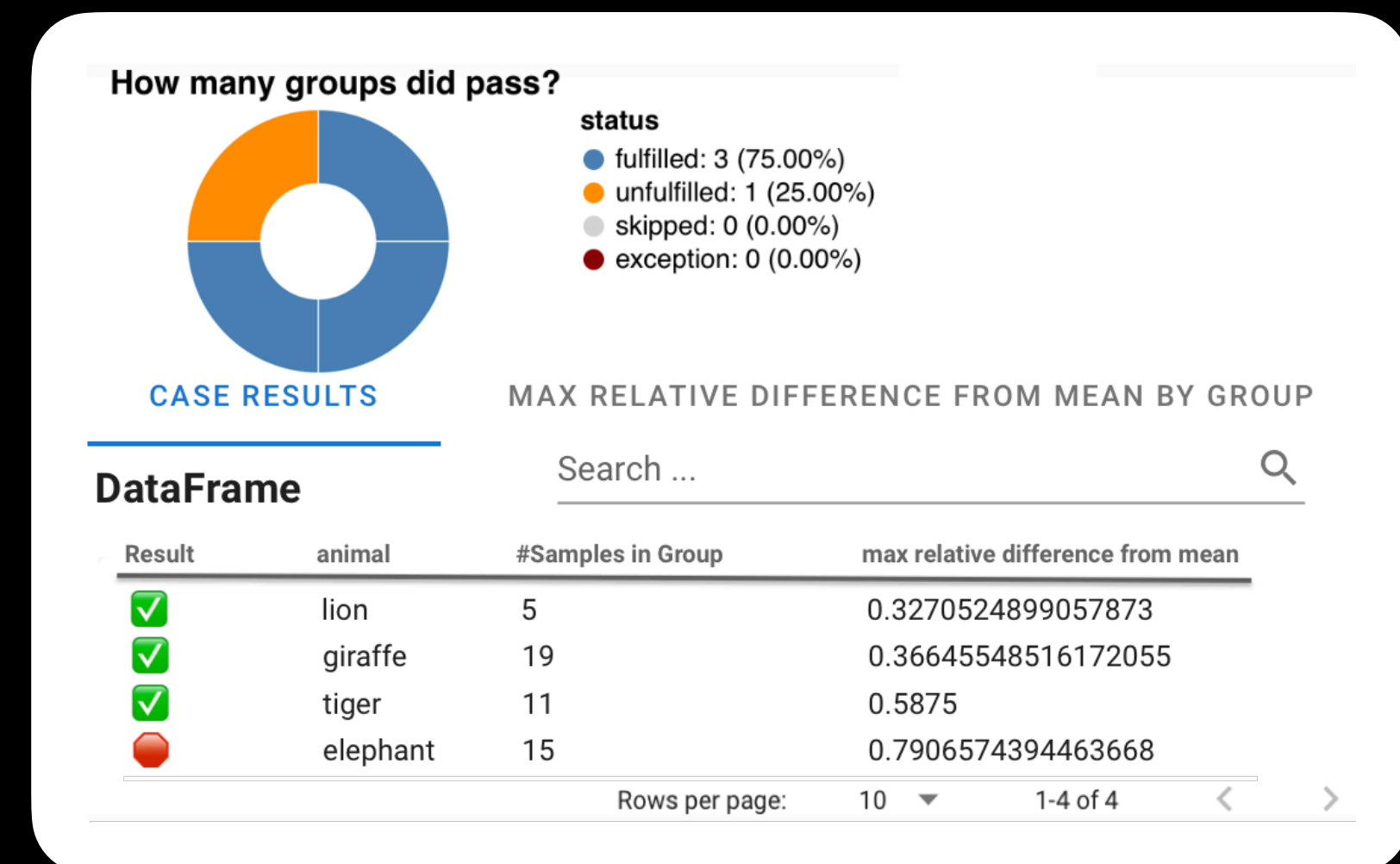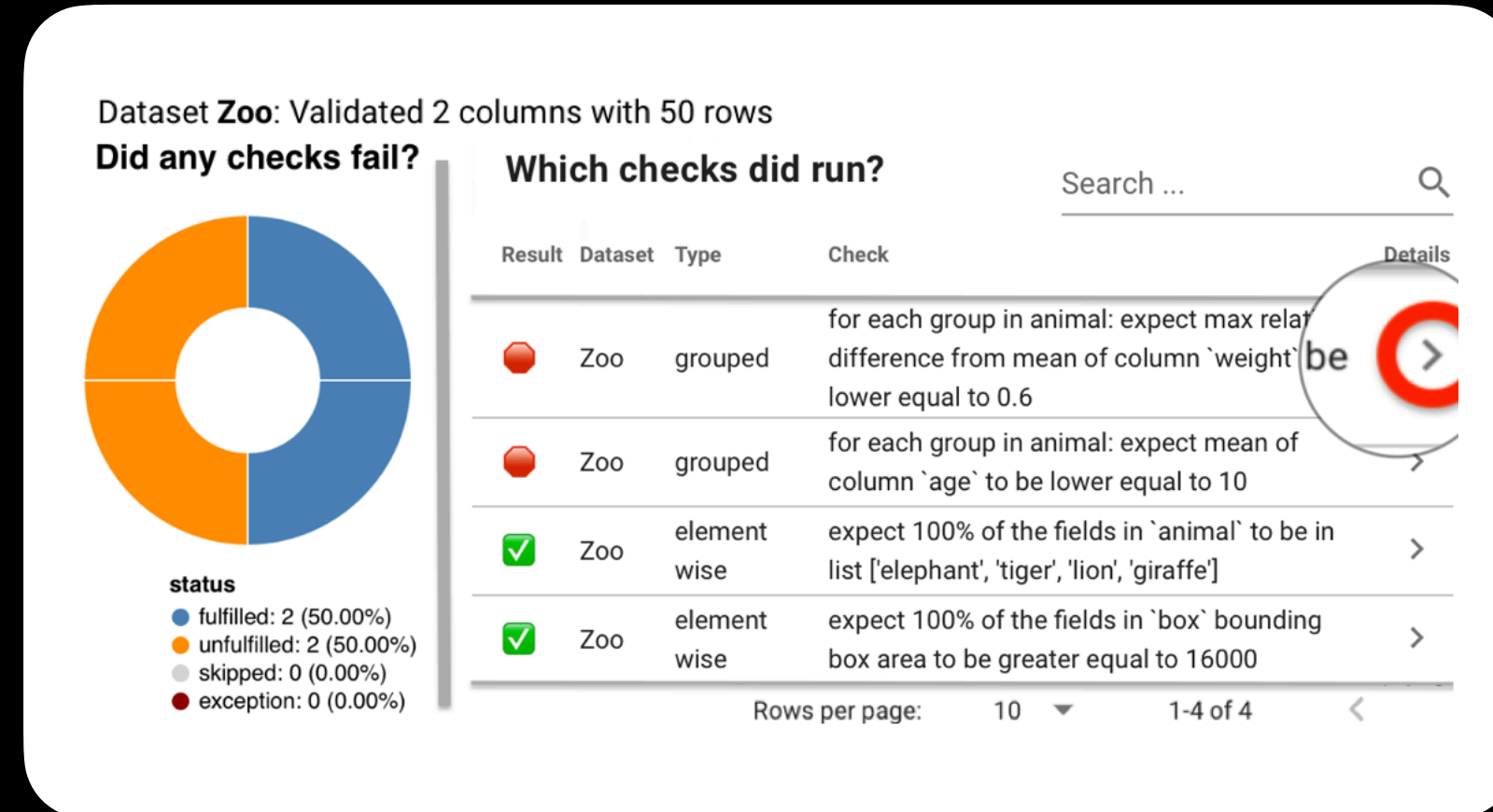
# Data Validation

```python
from udf import bbox_area, mean, max_diff_from_mean

checks = [
    Cell("animal") in ["elephant", "tiger", "lion", "giraffe"],
    Cell("box", bbox_area) >= 16000,
    Group("animal", Column("age", mean) < 10),
    Group("animal",
        Column("weight", max_diff_from_mean) < 0.6
    )
]

Validator(expectations=checks).run(data)
```



Dataset **Zoo**: Validated 2 columns with 50 rows
**Did any checks fail?**      **Which checks did run?**      Search ...

| Result | Dataset | Type | Check | Details |
|--------|---------|------|-------|---------|
| 🔴 | Zoo | grouped | for each group in animal: expect max relative difference from mean of column `weight` be lower equal to 0.6 | |
| 🔴 | Zoo | grouped | for each group in animal: expect mean of column `age` to be lower equal to 10 | |
| ✅ | Zoo | element wise | expect 100% of the fields in `animal` to be in list ['elephant', 'tiger', 'lion', 'giraffe'] | > |
| ✅ | Zoo | element wise | expect 100% of the fields in `box` bounding box area to be greater equal to 16000 | > |

status
- fulfilled: 2 (50.00%)
- unfulfilled: 2 (50.00%)
- skipped: 0 (0.00%)
- exception: 0 (0.00%)

Rows per page: 10    1-4 of 4    <



**How many groups did pass?**

status
- fulfilled: 3 (75.00%)
- unfulfilled: 1 (25.00%)
- skipped: 0 (0.00%)
- exception: 0 (0.00%)

CASE RESULTS      MAX RELATIVE DIFFERENCE FROM MEAN BY GROUP

**DataFrame**      Search ...

| Result | animal | #Samples in Group | max relative difference from mean |
|--------|--------|-------------------|-----------------------------------|
| ✅ | lion | 5 | 0.3270524899057873 |
| ✅ | giraffe | 19 | 0.36645548516172055 |
| ✅ | tiger | 11 | 0.5875 |
| 🔴 | elephant | 15 | 0.7906574394463668 |

Rows per page: 10    1-4 of 4    <    >

# Case Studies

# Robust and Reproducible Dataset Documentation

- Capture context about datasets in documentation

- Enrich with summary statistics and raw data preview

- Transition seamlessly to interactive usage
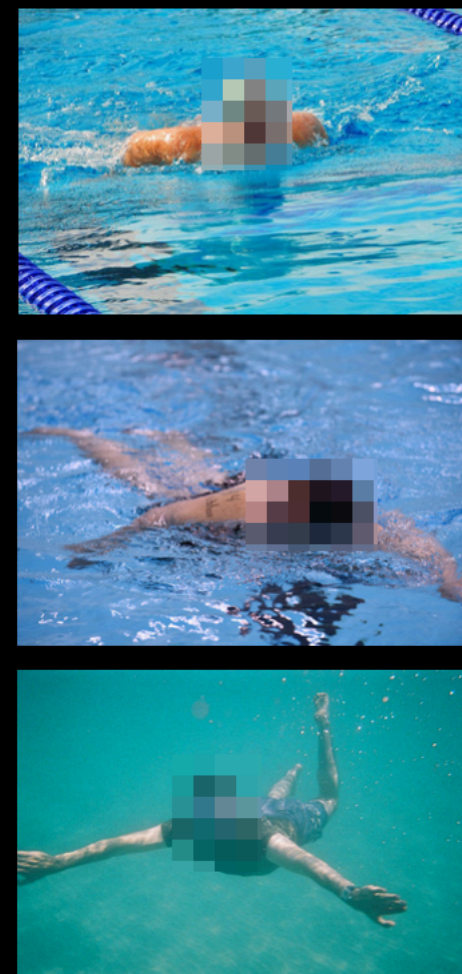
# Large Corpora Curation with Linguists in the Loop

- Empower domain experts to set filters while inspecting raw data

- Enable interactivity for 100M+ paragraphs

- Seamless collaboration between ML engineers and domain experts

# Actionable Failure Analysis of an Image Classifier

- When available metadata does not explain wrong predictions, visual inspection can reveal clues

- Image-text encoder models can be used to source new data in the same interactive session
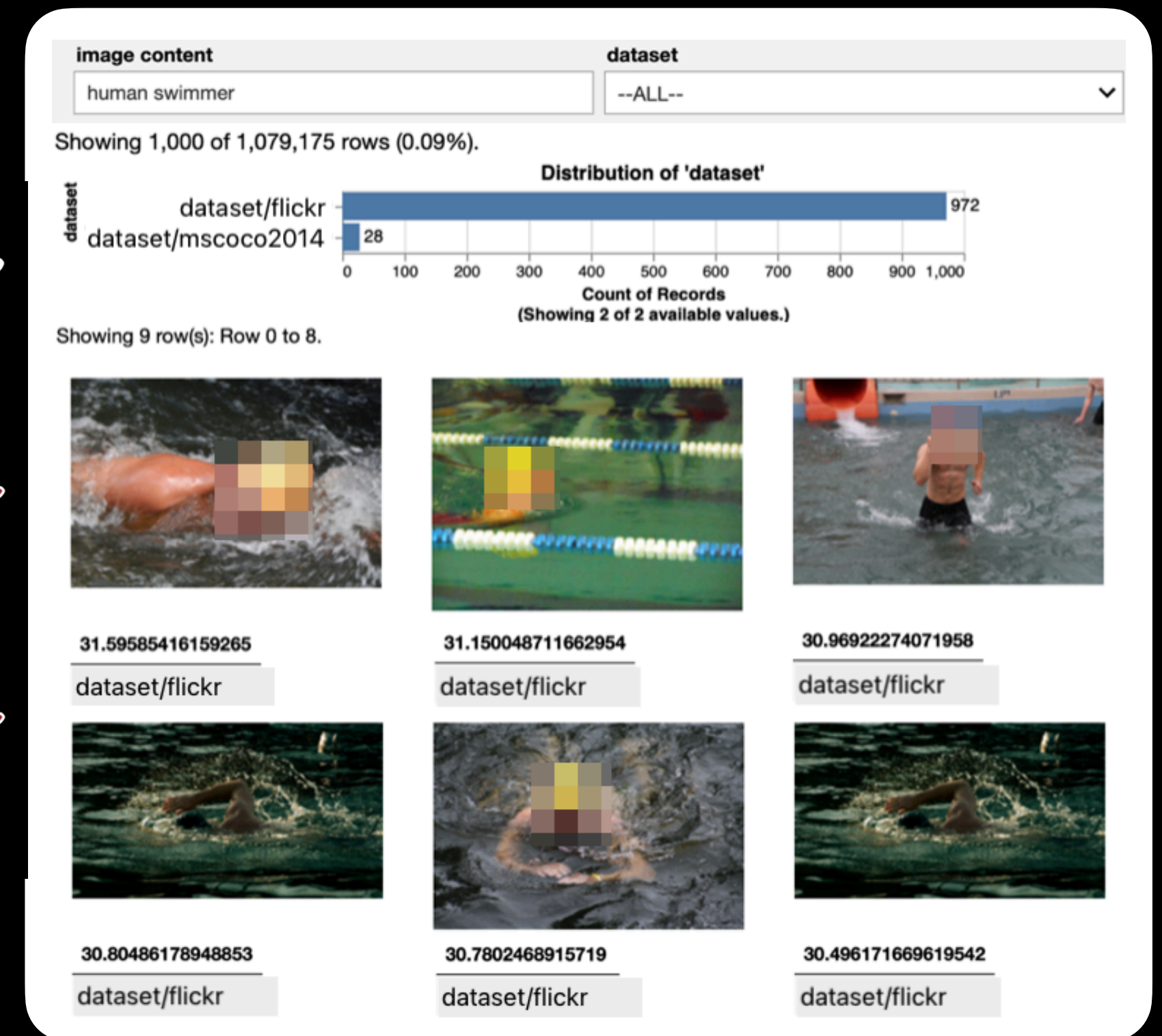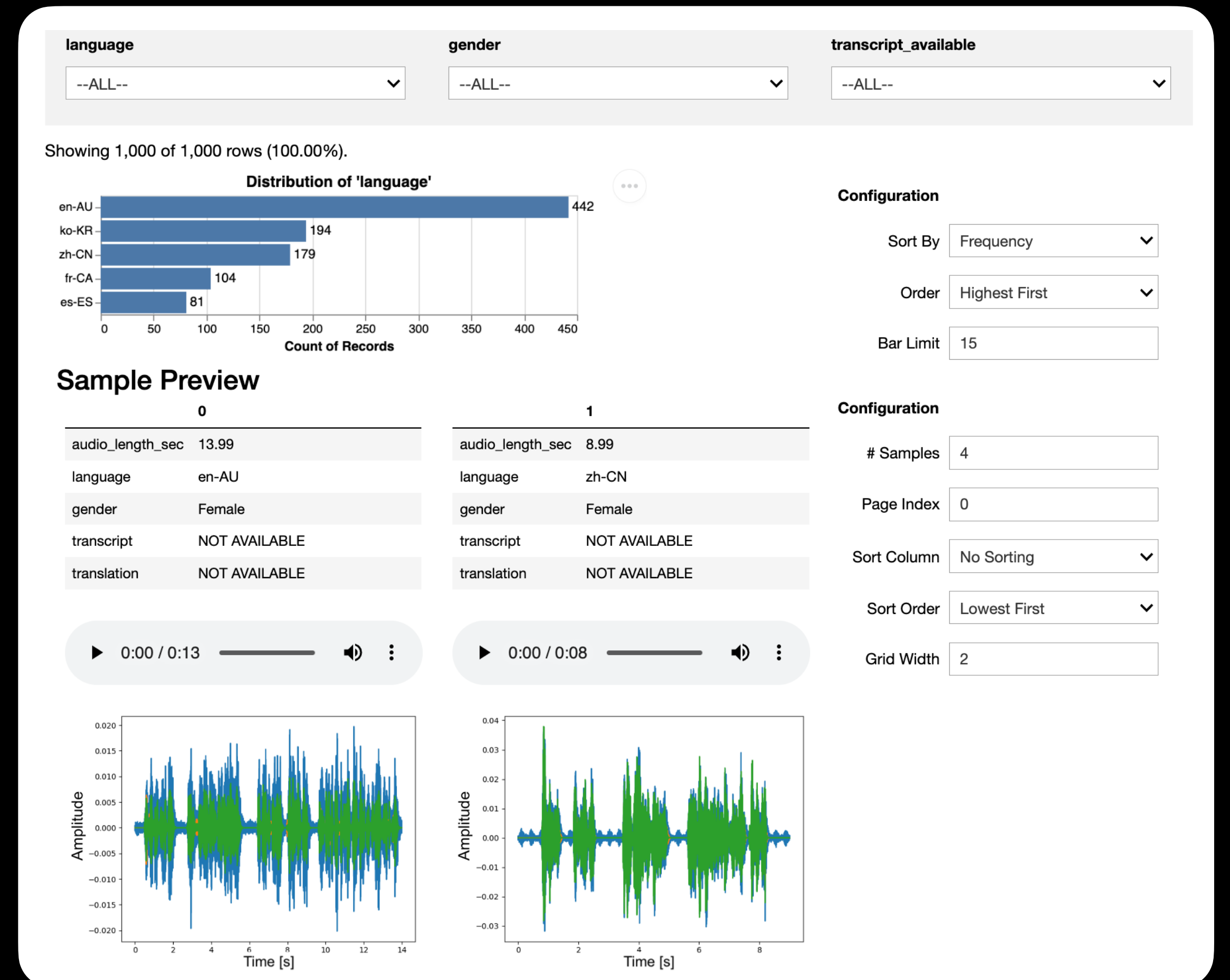
# Rapid Iteration on Multi-Modal Datasets
from heterogeneous data sources

- Multi-modal machine translation project

- Necessity to blend different data sources and modalities in a unified preview

- Learnings from exploration could be persisted in validation checks

# Conclusions

- We presented ADIML, a toolset to democratise data technology throughout the ML lifecycle and to enable the data-centric ML approach.

- The design of ADIML is based on the set of challenges and pain points we collected and validated from a wide range of ML teams.

- The case studies showing how easily ADIML can enable ML teams to focus and to improve data quality at scale are testimonies of its values.