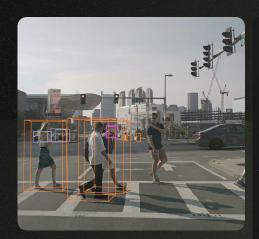


Scaling Computer Vision in the Enterprise







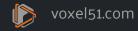
Pest Infestation



Why Data-Centric Computer Vision?

Agenda

- > Why data-centric computer vision?
- Dataset management at scale
- Data annotation and auto-labeling
- > Out-of-the-box workflows
- Collaborating on datasets
- Safety and security
- Next steps



3



DATA EATS MODELS FOR LUNCH

Unlocking the potential of your data is key to competitive advantage



"I do believe the models are getting commoditized...models by themselves are not sufficient, but having a full system stack and great successful products, those are the two places" where companies need to focus now."

Satya Nadella, Microsoft CEO



Al projects fail 85% of the time



Amazon's "Just Walk Out" store vision system could not handle the full diversity of real shopper behavior and product interactions. Human reviewers had to verify about 70% of transactions, because the Al wasn't accurate enough on its own.

GUARDIAN



NHTSA has reports of 16 crashes, including seven injury incidents and one death, involving Tesla vehicles in Autopilot that had struck stationary first-responder and road maintenance vehicles.

REUTERS

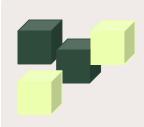


GE's Caption AI miscalculated key heart metrics during ultrasound exams due to a data-handling bug that included incorrect video frames. The flaw triggered a Class II FDA recall.

FDA



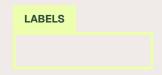
30% of model errors are due to bad data



Unstructured data is difficult to search, slice, and debug



Complex integration across modalities, labels, and metadata



Hard to annotate and validate at scale



ML teams spend 39% of their time wrangling data

It is simply *impossible* to touch each data sample manually to ensure quality. Yet, even top teams have no good alternatives

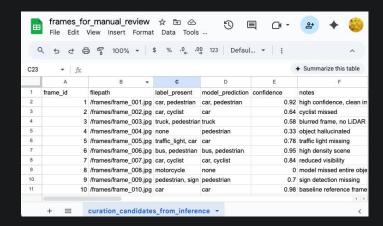
This is what it looks like when your data stack wasn't built for visual data.

It's time-consuming, painful, and impossible to scale.

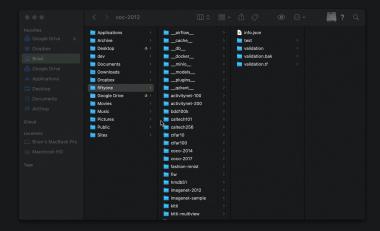
Python scripts and Jupyter notebooks



Spreadsheets



Manual debugging loops



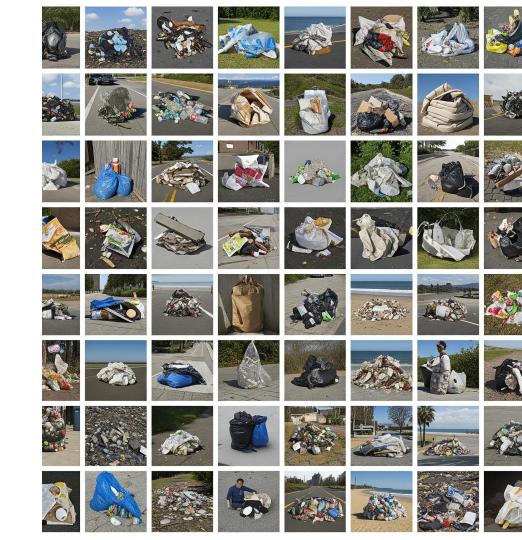


DATA CURATION

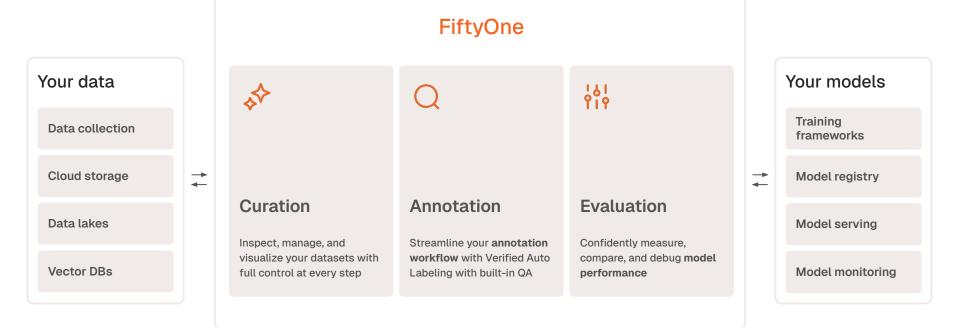
Garbage in, garbage out

Poor dataset curation leads to:

- Misleading metrics and poor generalization in production
- Missed edge cases that compromise safety and robustness
- Incorrect labels that silently degrade model performance



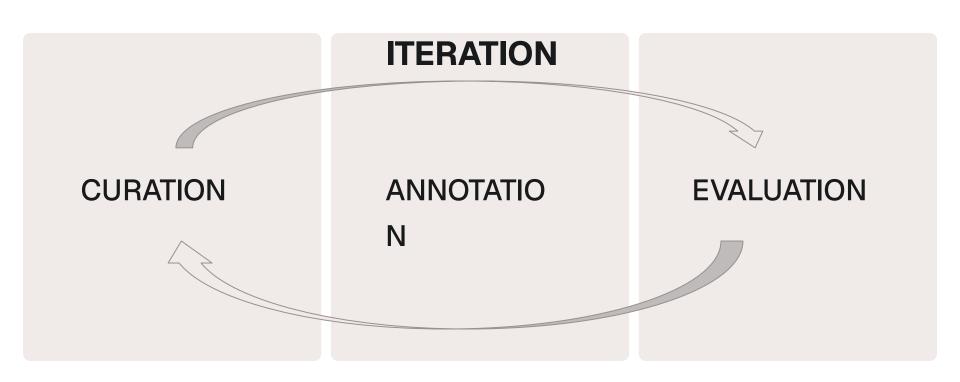




Your compute



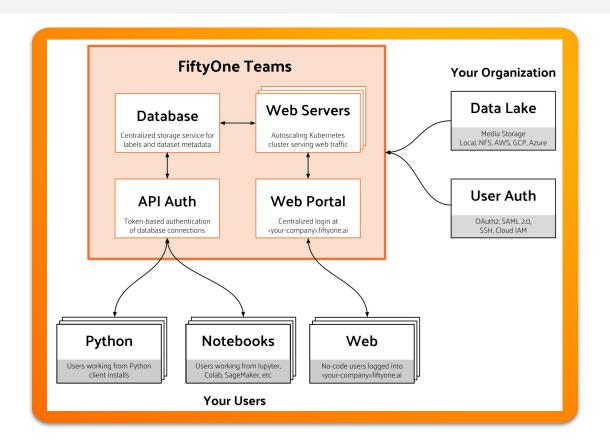
Building a Successful Pipeline





Dataset Management at Scale

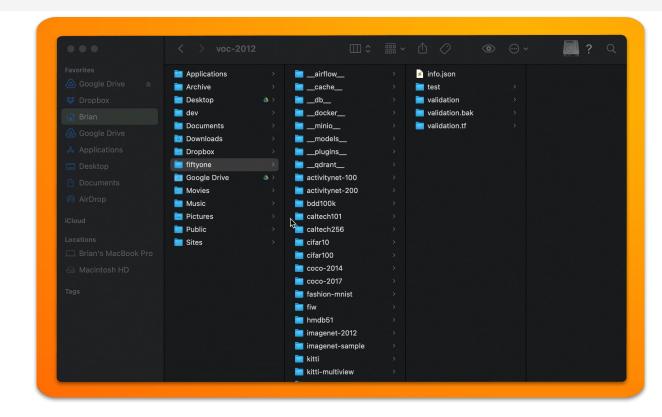
Platform Architecture





Where Do Datasets Live?

For DIY and some open source tools, it's (generally) on your computer's local filesystem.

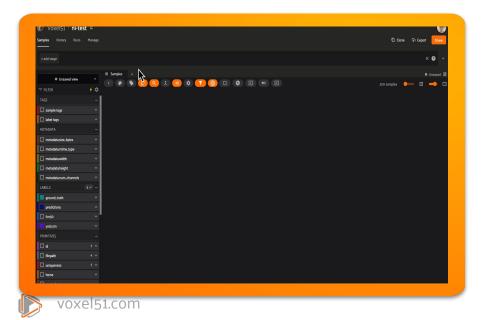




Where Do Datasets Live?

For enterprise vendors, cloud-backed media or data lake

Cloud-Backed Media



Data Lake



Dataset Storage

Cloud-Backed Media

- ✓ Any S3-compatible object store
- Cached locally when running workflows requiring pixel-level access
- ✓ The sample's file path will be the object storage path
 - E.g., gs://path/to/sample
- Admins can configure cloud credentials in the App

Data Lake Integration

- Connects to and imports from your data lake (e.g., Databricks, BigQuery)
- Makes use of operators for defining connection to data lake, query parameters, preview behavior
- Import selected number of samples into new or existing dataset



Data Lens: Powered by FiftyOne Operators

```
import json
import time
from typing import Generator

import fiftyone as fo
from databricks.sdk import WorkspaceClient
from databricks.sdk.service.sal import (
    StatementResponse, StatementState, StatementParameterListItem
)
from fiftyone import operators as foo
from fiftyone.operators import types
from fiftyone.operators.data lens import (
    DataLensOperator, DataLensSearchRequest, DataLensSearchResponse
)
```

```
class DatabricksConnector(DataLensOperator):
    """Data Lens operator which retrieves samples from Databricks."""
    @property
    def config(self) -> foo.OperatorConfig: --
    def resolve_input(self, ctx: foo.ExecutionContext): --
   def handle_lens_search_request(...
class DatabricksHandler:
    """Handler for interacting with Databricks tables."""
    def __init__(self): --
    def handle request( ...
   def _init_client(self, ctx: foo.ExecutionContext):--
    def _start_warehouse(self) -> str: --
    def _iter_data(self, request: DataLensSearchRequest) -> Generator[dict, None, None]: ...
    def _transform_sample(self, sample: dict) -> dict: --
    def _build_detections(self, sample: dict) -> fo.Detections: --
   def _response_to_dicts(self, response: StatementResponse) -> list[dict]: "
   def _check_for_error(self, response: StatementResponse): --
```





Auto Labeling

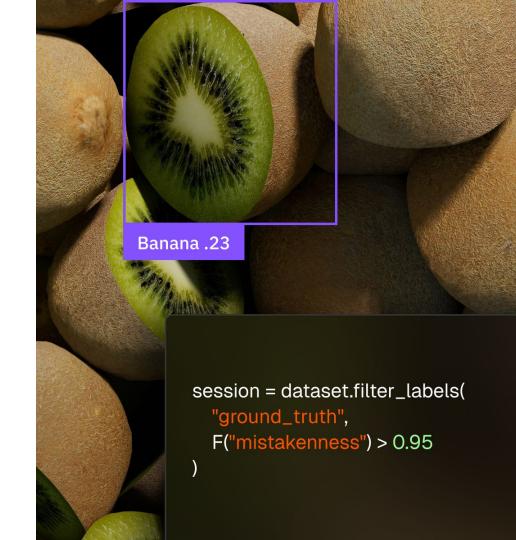


ANNOTATION

Annotation bottlenecks stall AI development

Manual annotation is slow, inconsistent, and expensive.

- High cost with little scalability
- Painfully slow for large datasets
- Quality varies across annotators and vendors
- Delays model iteration and validation



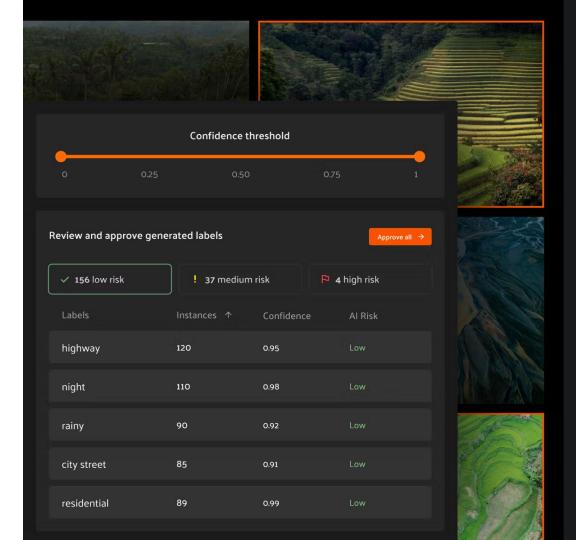


ANNOTATION

Verified Auto Labeling

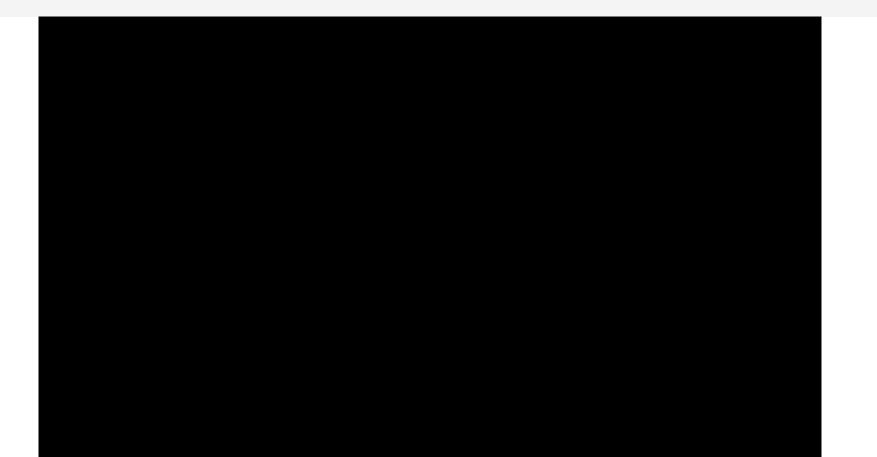
VAL uses foundation models to automatically generate labels — and adds confidence scoring to prioritize the ones that require human review.

- Reduce QA and annotation costs while maintaining near-human accuracy
- Bring your own foundational models
- Ensure data sovereignty by keeping more of your proprietary training data in-house



<u>Demo</u>

Verified Auto Labeling with FiftyOne







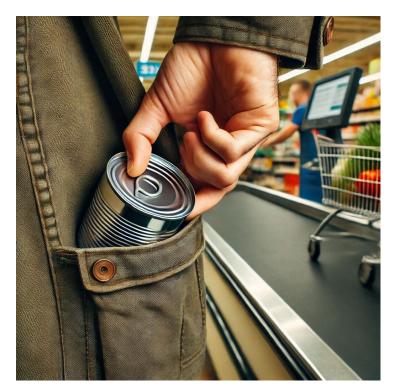
Model Evaluation

My model is broken! A mystery...





My model is broken! A mystery...



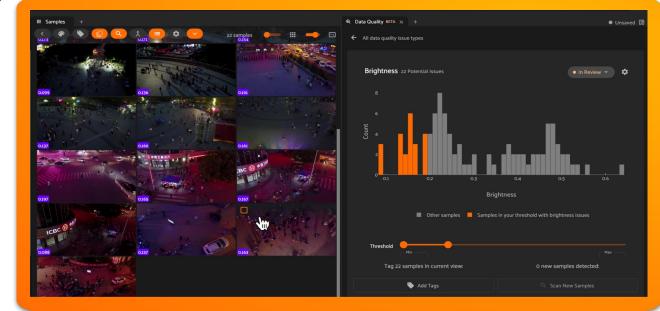




Data Quality

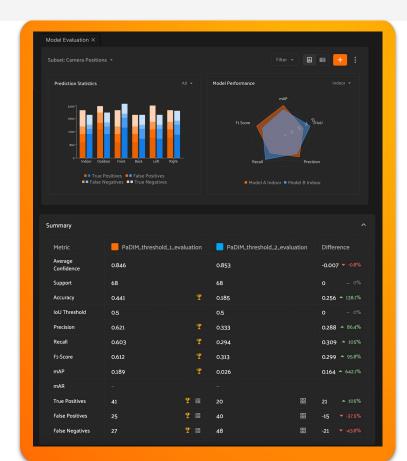
Scan for and analyze several different types of data quality issues

- ✓ Brightness
- ✓ Blurriness
- ✓ Aspect Ratio
- ✓ Entropy
- ✓ Near Duplicates
- ✓ Exact Duplicates





Model Evaluation and Scenario Analysis



- Evaluate predictions visually and numerically
- Slice performance by conditions
- Compare multiple models or thresholds side-by-side
- Detect hidden failure modes and edge cases at a glance



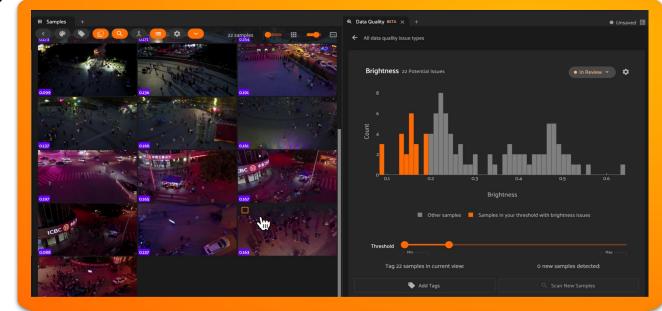


Iteration

Data Quality

Scan for and analyze several different types of data quality issues

- ✓ Brightness
- ✓ Blurriness
- ✓ Aspect Ratio
- ✓ Entropy
- ✓ Near Duplicates
- ✓ Exact Duplicates



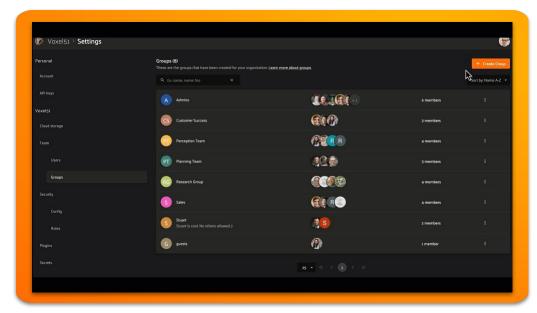




Collaborate Securely on Datasets

Collaborate With Your Team

- Many free and open source tools are (generally) single-user
- Enterprise platforms generally multi-user by design
 - Hooks into your org's auth provider (AD / LDAP / SAML)
- Role-based access control for users and datasets
 - Limit datasets to particular users/groups
 - Users/group actions limited by their role





Dataset Versioning

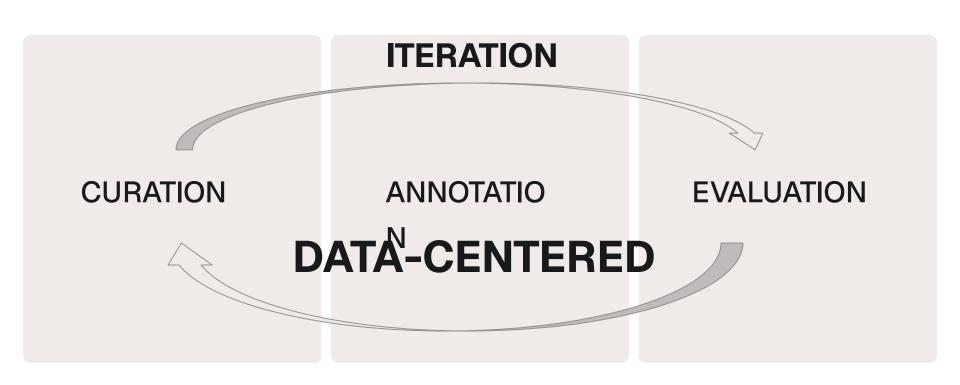


- Clone and export datasets
 - o DVC
- Create, share, and roll back dataset snapshots
- Compare diffs between snapshots





Building a Successful Pipeline





Thank you! Questions?

