



CSCE 585: Machine Learning Systems

Lecture 6: Designing Machine Learning Systems

Pooyan Jamshidi

What is missing?

The gap between ML Research and Production

Chip Huyen @chipro · Jul 19, 2019
Replying to @chipro

Most candidates told me the hardest questions for them are the machine learning system design questions. They don't know what a good answer to these questions looks like. Interviewers: any tips?

18 11 132

Ravi Ganti @gmravi2003 · Jul 19, 2019

When I ask such questions, what I am looking for is the following. 1. Can the candidate break down the open ended problem into simple components (building blocks) 2. Can the candidate identify which blocks require ML and which do not.

9

What is missing?

The gap between ML Research and Production



Dmitry Kislyuk @dkislyuk · Jul 19, 2019

Replies to [@lishali88](#) and [@chipro](#)

Most candidates know the model classes (linear, decision trees, lstms, convnets) and memorize the relevant info, so for me the interesting bits in ML systems interviews are data cleaning, data prep, logging, eval metrics, scalable inference, feature stores (recommenders/rankers)



What is missing?

The gap between ML Research and Production



Illia Polosukhin @ilblackdragon · Jul 20, 2019

I think this is the most important question. Can person define problem, identify relevant metrics, ideate on data sources and possible important features, understands deeply what ML can do. ML methods change every year, solving problems stays the same.



In ML Systems, only a small fraction is comprised of actual ML code

Hidden Technical Debt in Machine Learning Systems

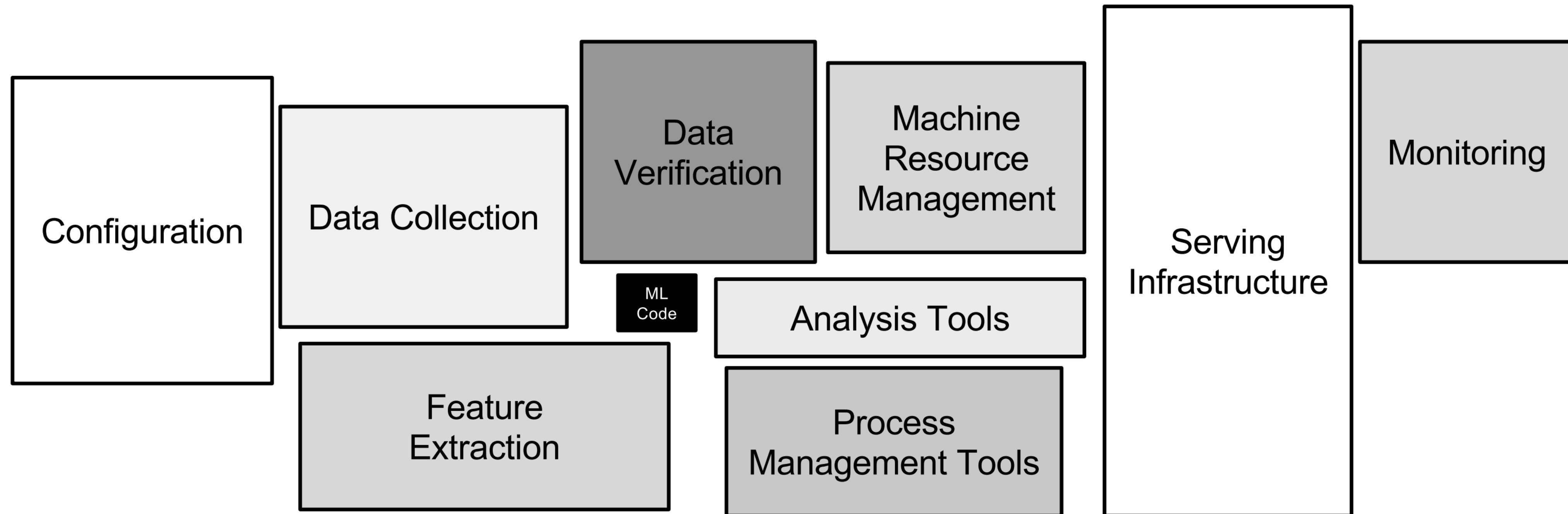
D. Sculley, Gary Holt, Daniel Golovin, Eugene Davydov, Todd Phillips
{dsculley, gholt, dg, edavydov, toddphillips}@google.com
Google, Inc.

Dietmar Ebner, Vinay Chaudhary, Michael Young, Jean-François Crespo, Dan Dennison
{ebner, vchaudhary, mwyoung, jfcrespo, dennison}@google.com
Google, Inc.

Abstract

Machine learning offers a fantastically powerful toolkit for building useful complex prediction systems quickly. This paper argues it is dangerous to think of these quick wins as coming for free. Using the software engineering framework of *technical debt*, we find it is common to incur massive ongoing maintenance costs in real-world ML systems. We explore several ML-specific risk factors to account for in system design. These include boundary erosion, entanglement, hidden feedback loops, undeclared consumers, data dependencies, configuration issues, changes in the external world, and a variety of system-level anti-patterns.

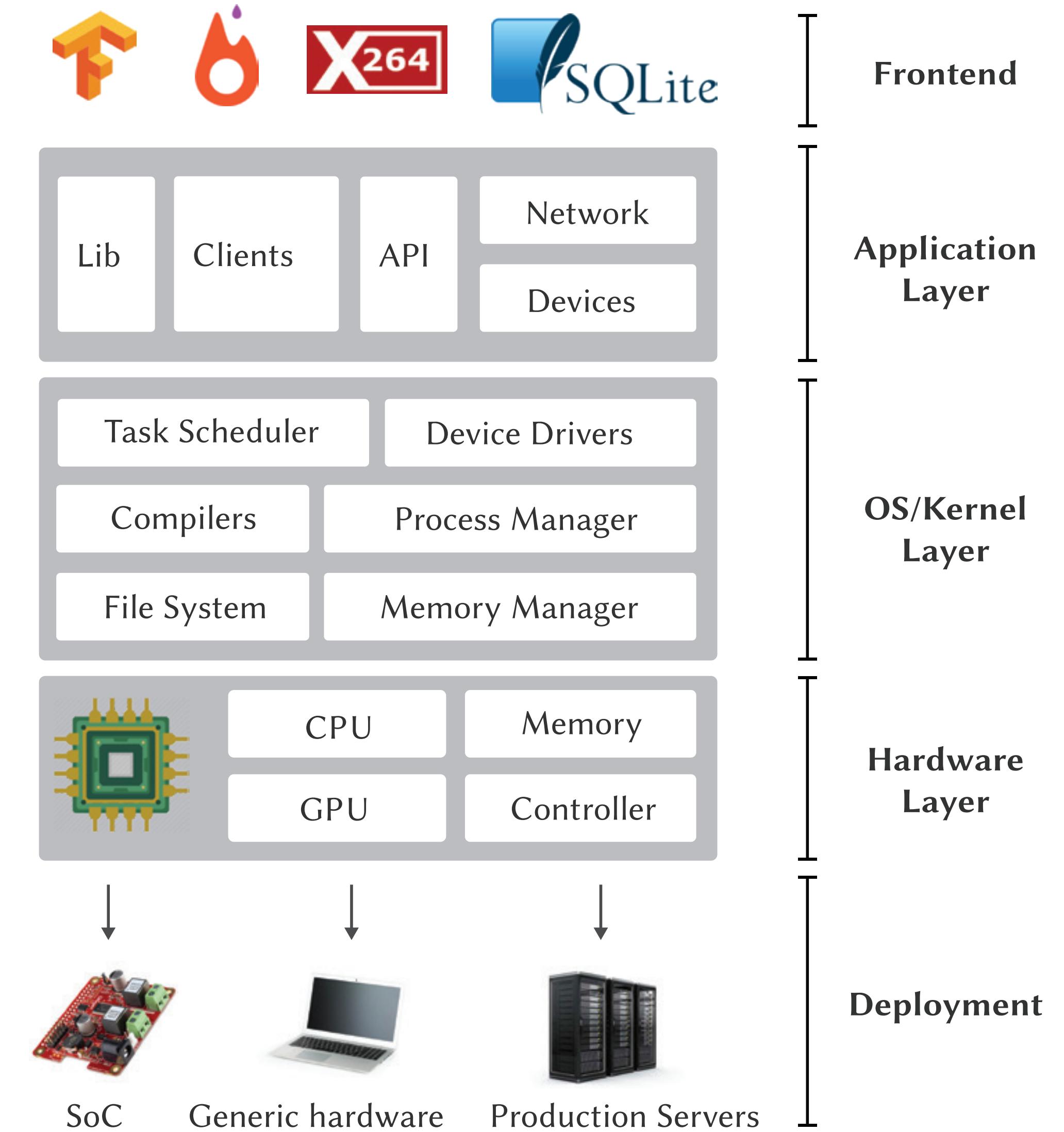
A vast array of surrounding infrastructure and processes is needed to support evolution of ML systems



Technical debt that can accumulate in ML systems

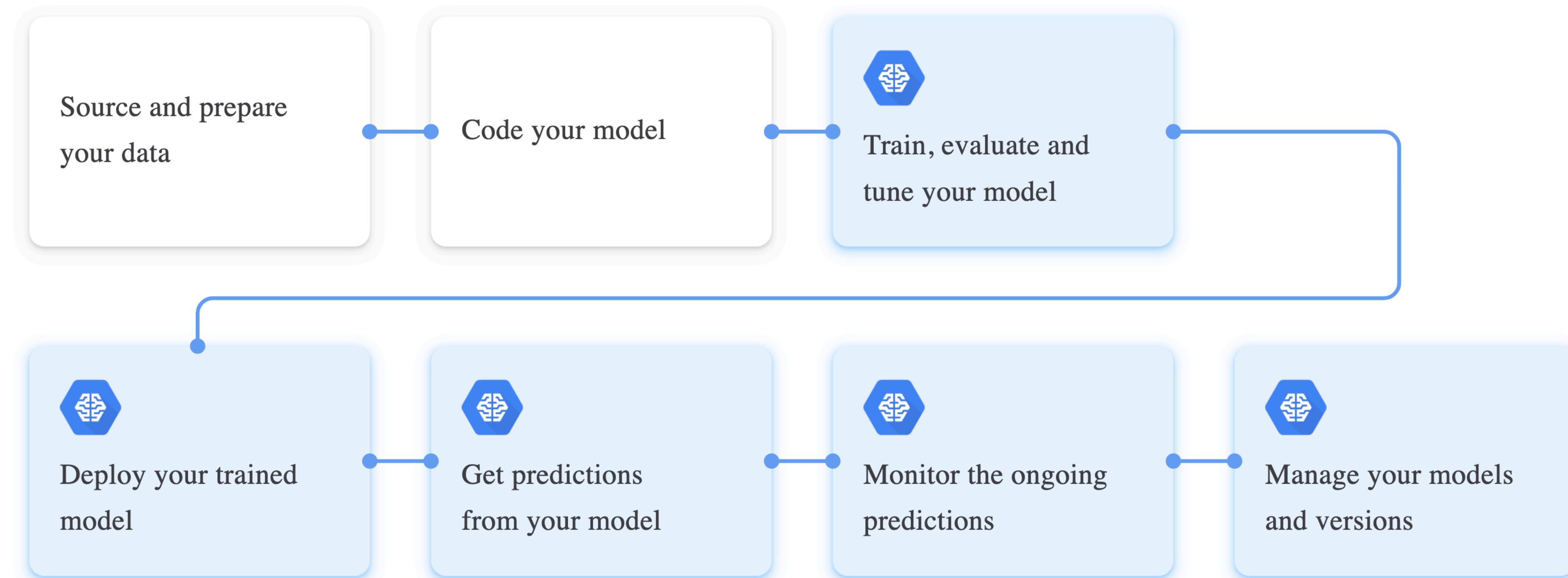
- Data dependencies
- Model complexity
- Reproducibility
- Testing
- Monitoring
- Configuration issues
- External changes

System = Software + Middleware + Hardware



The Building Process of ML Systems

Continuous Delivery for ML Systems



A Machine Learning System is more than just a model

The 3 axes of change in a ML System – data, model, and code – and a few reasons for them to change



Data

Schema

Sampling over Time

Volume



Model

Algorithms

More Training

Experiments



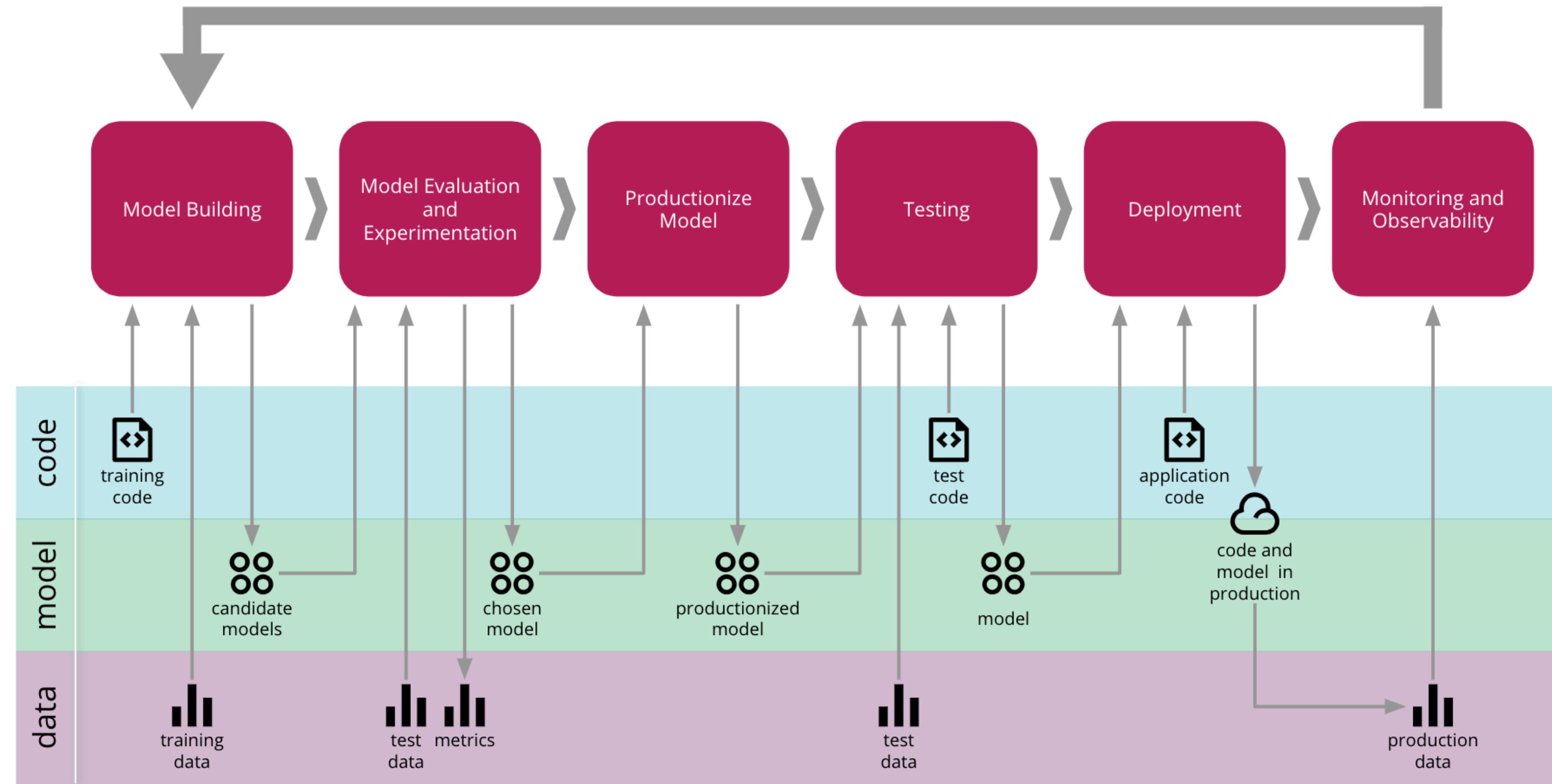
Code

Business Needs

Bug Fixes

Configuration

An End-to-End ML System Building Process



Continuous Delivery for Machine Learning Systems

Continuous Delivery for Machine Learning (CD4ML) is a software engineering approach in which a cross-functional team produces machine learning applications based on code, data, and models in small and safe increments that can be reproduced and reliably released at any time, in short adaptation cycles.

A Machine Learning Application for Sales Forecasting

Corporación Favorita Grocery Sales Forecasting

Can you accurately predict sales for a large grocery chain?



Overview Data Code Models Discussion Leaderboard Rules

Overview

Start

Oct 19, 2017

Close

Jan 15, 2018

Merger & Entry

Description

Brick-and-mortar grocery stores are always in a delicate dance with purchasing and sales forecasting. Predict a little over, and grocers are stuck with overstocked, perishable goods. Guess a little under, and popular items quickly sell out, leaving money on the table and customers fuming.

The problem becomes more complex as retailers add new locations with unique needs, new products, ever transitioning seasonal tastes, and unpredictable product marketing. [Corporación Favorita](#), a large Ecuadorian-based grocery retailer, knows this all too well. They operate hundreds of supermarkets, with over 200,000 different products on their shelves.

[Corporación Favorita](#) has challenged the Kaggle community to build a model that more accurately forecasts product sales. They currently rely on subjective forecasting methods with very little data to back them up and very little automation to execute plans. They're excited to see how machine learning could better ensure they please customers by having just enough of the right products at the right time.

Competition Host

Corporación Favorita



Prizes & Awards

\$30,000

Awards Points & Medals

Participation

9,832 Entrants

1,868 Participants

1,671 Teams

31,241 Submissions

Tags

Regression

Food

Tabular

Custom Metric

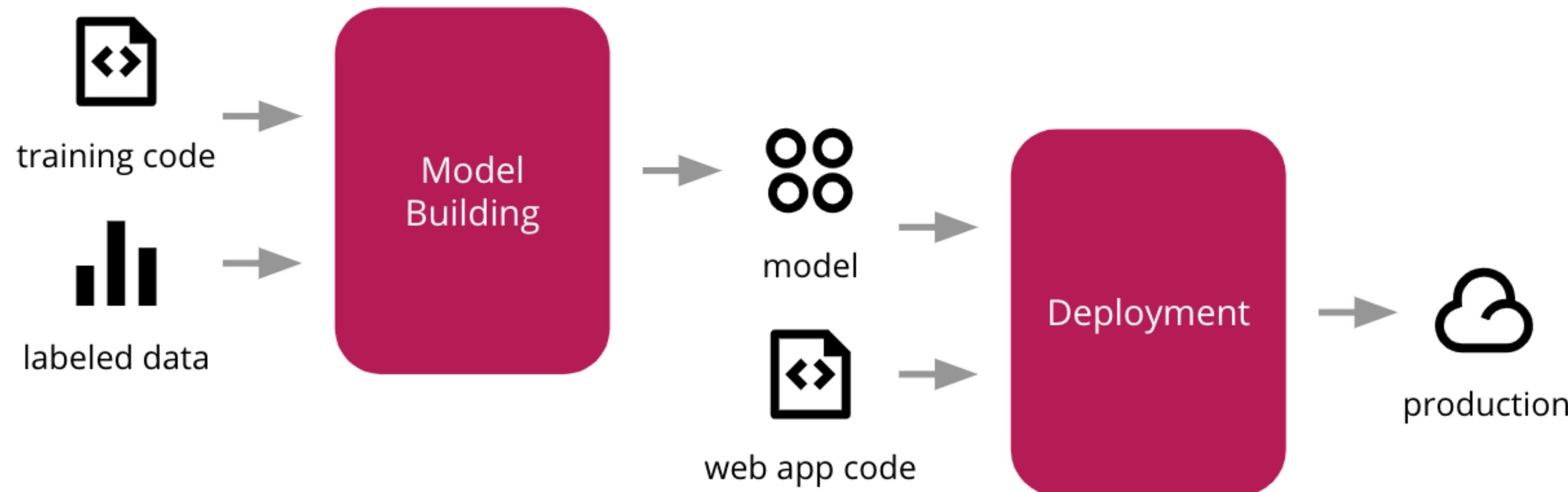
Table of Contents

Description

Evaluation

Prizes

Train ML model, integrate it with an application, and deploy into production



ML model behind a web application

A screenshot of a web browser window displaying a "Sales forecast" application. The address bar shows the URL as "localhost:5005". The main title is "Sales forecast". There are two input fields: one for "Date" (with placeholder "YYYY-MM-DD") and one for "Product" (set to "Milk", with a dropdown arrow). A blue "Submit" button is below the product field. At the bottom, there is a long, empty rectangular input field labeled "Prediction:".

← → ⌂ ⓘ localhost:5005

Sales forecast

Date

Product

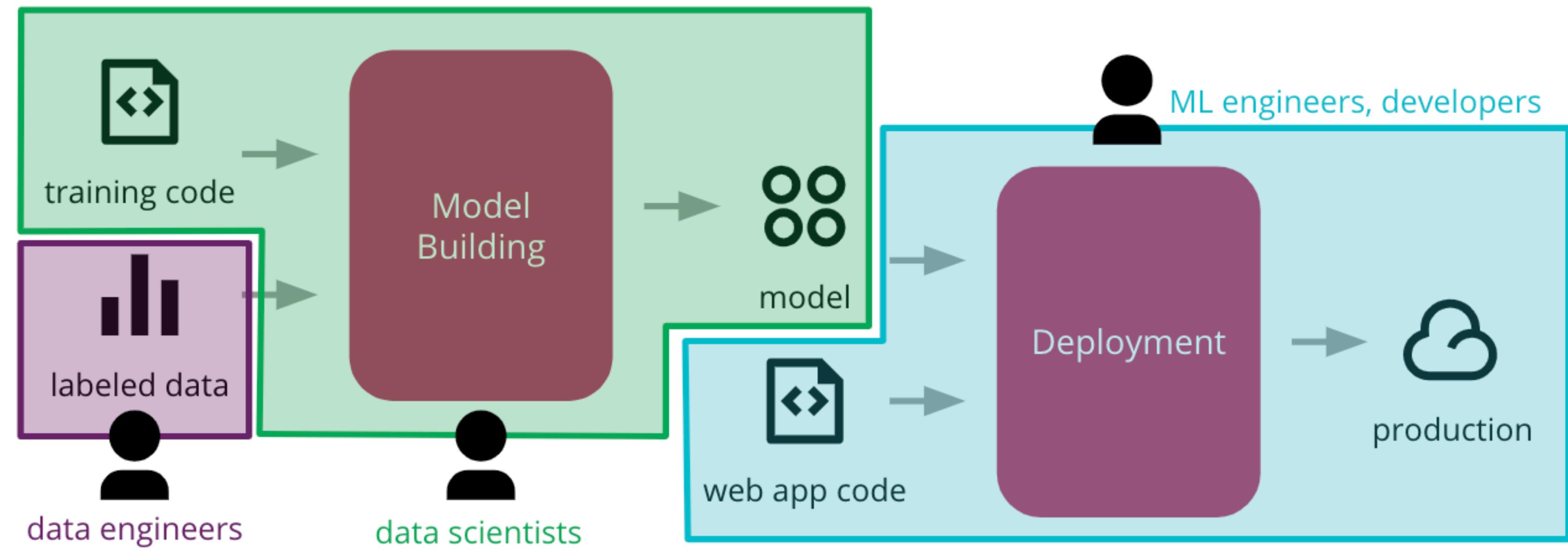
Submit

Prediction:

Common Challenges

1. **Organizational challenge:**

- Throw over the wall
- Models that only work in a lab environment
- Even if make it to production, they become stale and hard to update



2. **Technical challenge:** How to make the process

- Reproducible
- Auditable

Technical Components of CD for ML Systems

The first step was to understand the data set:

- **Products**, like their categorization and whether they are perishable or not
- **Stores**, like their location and how they are clustered together
- **Events**, like public holidays, seasonal events, or a 7.8 magnitude earthquake that struck Ecuador in 2016
- **Sales transactions**, including the number of units sold for a given product, date, and location

Discoverable and Accessible Data

- It is important that the data is easily discoverable and accessible.
- The harder it is for Data Scientists to find the data they need, the longer it will take for them to build useful models.
- We should also consider they will want to engineer new features on top of the input data that might help improve their model's performance.

Let's see what we need to do in our example scenario

Initial Exploratory Data Analysis

- To **de-normalize** multiple files into a single CSV file.
- To **cleanup** the data points that were not relevant or that could introduce unwanted noise into the model (like the negative sales).
- We then **store** the output in a cloud storage system like Amazon S3, Google Cloud Storage, or Azure Storage Account.
- Using this file to represent a snapshot of the input training data, we are able to devise a simple approach to **version our dataset** based on the folder structure and file naming conventions. Data versioning is an extensive topic, as it can change in two different axis:
 - Structural **changes** to its schema, and
 - The actual **sampling** of data over time.

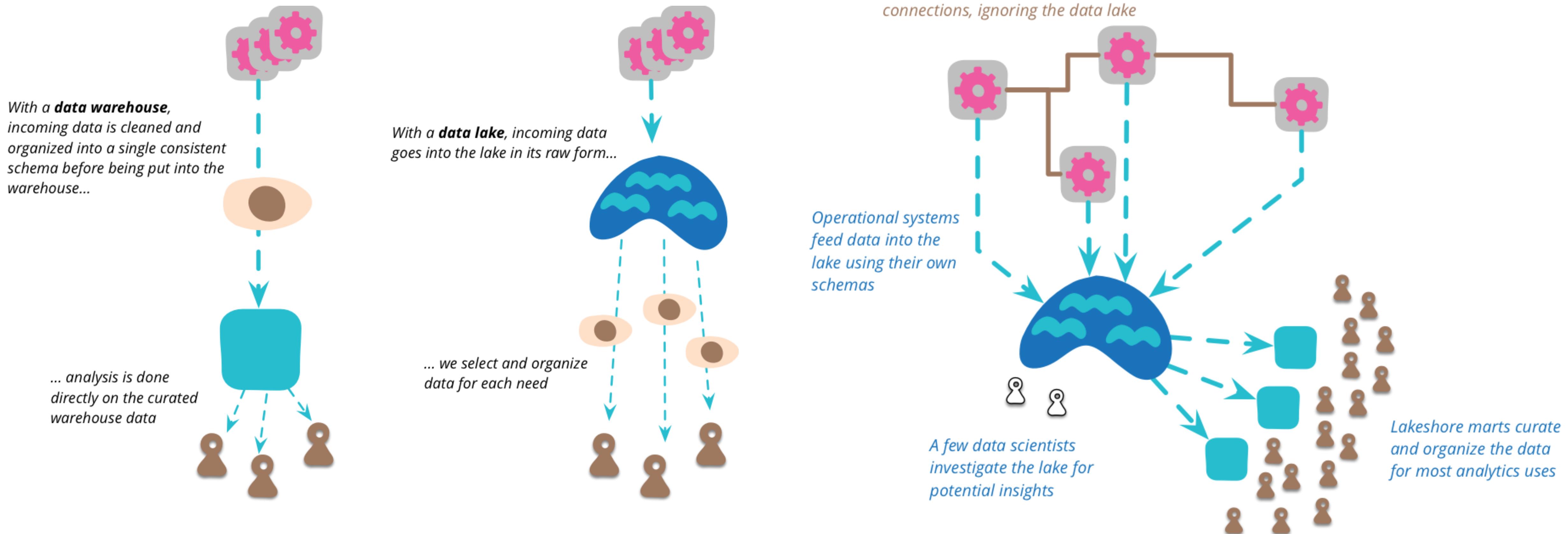
Discoverable and Accessible Data

In the real world, you will likely have more complex data pipelines to move data from multiple sources into a place where it can be accessed and used by Data Scientists.

- The most common **source of data** will be your core transactional systems.
- However, there is also value in bringing other data sources from **outside your organization**.
- A few common patterns for collecting and making data available:
 - **Data Lake** architecture: a traditional data warehouse
 - **Data Mesh** architecture: decentralized

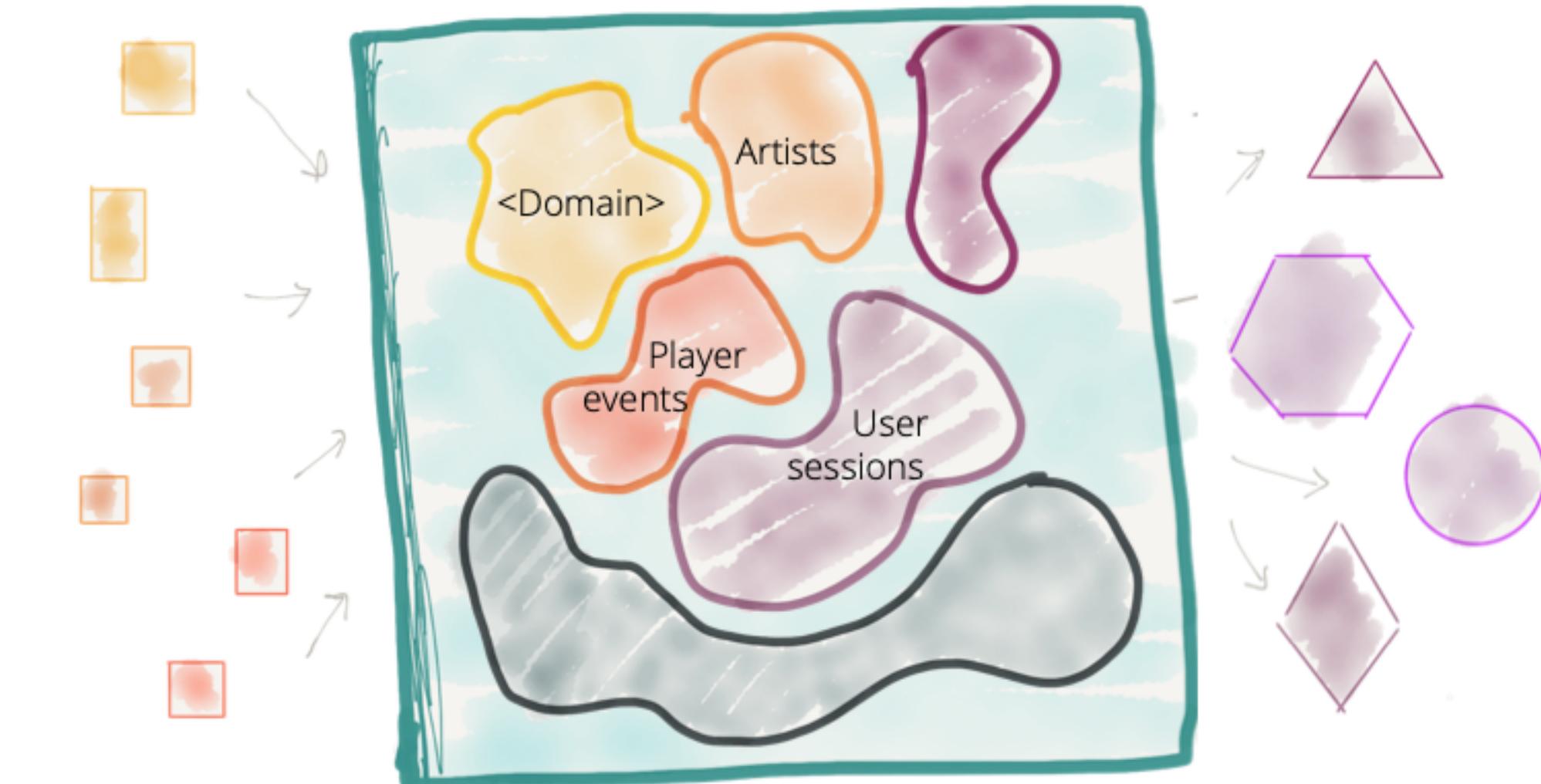
From Data Warehouse to Data Lake

The idea is to have a single store for all of the raw data that anyone in an organization might need to analyze.



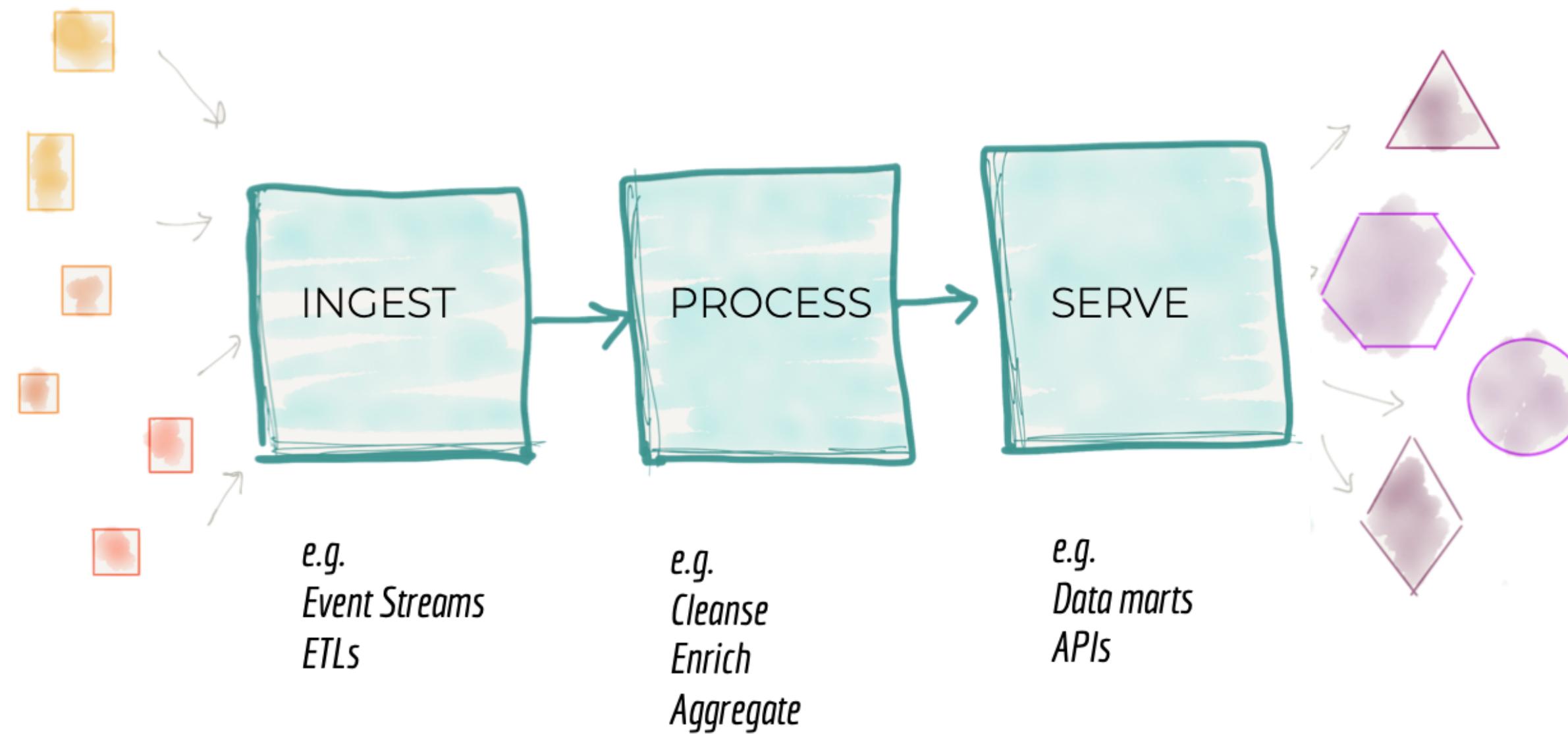
Architectural failure modes

Centralized data platform with no clear data domain boundaries and ownership of domain-oriented data

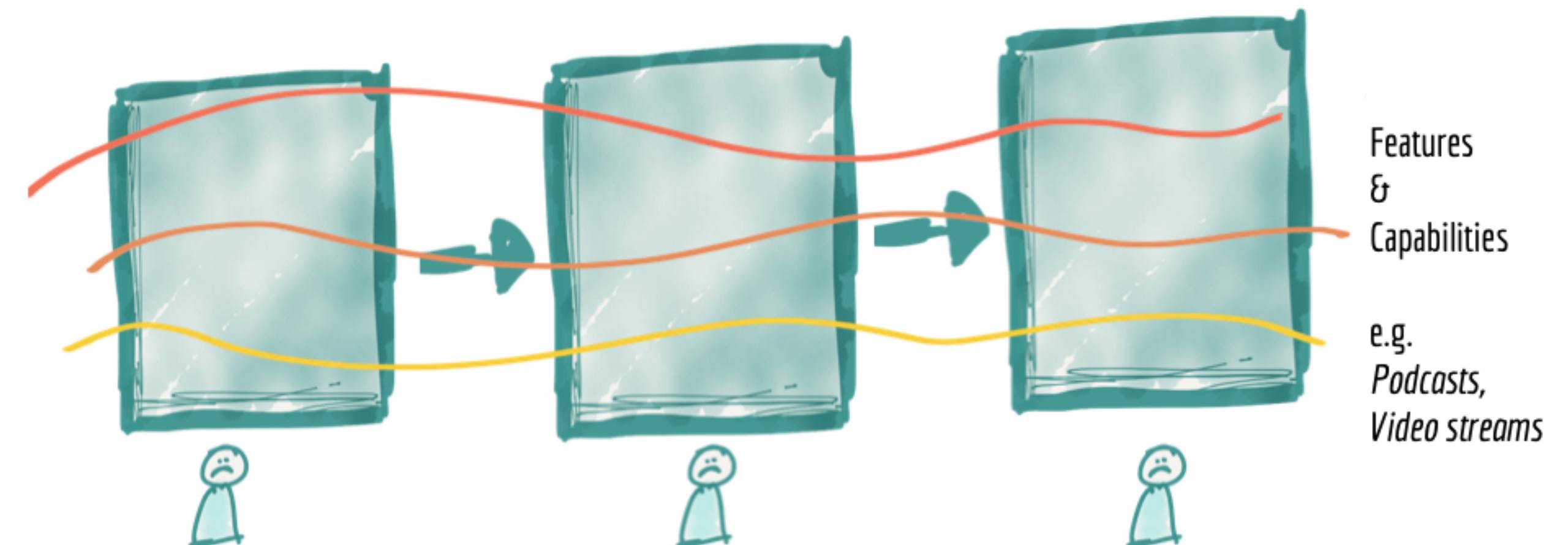


Architectural failure modes

Coupled pipeline decomposition



Architecture decomposition is orthogonal to the axis of change when introducing or enhancing features, leading to coupling and slower delivery



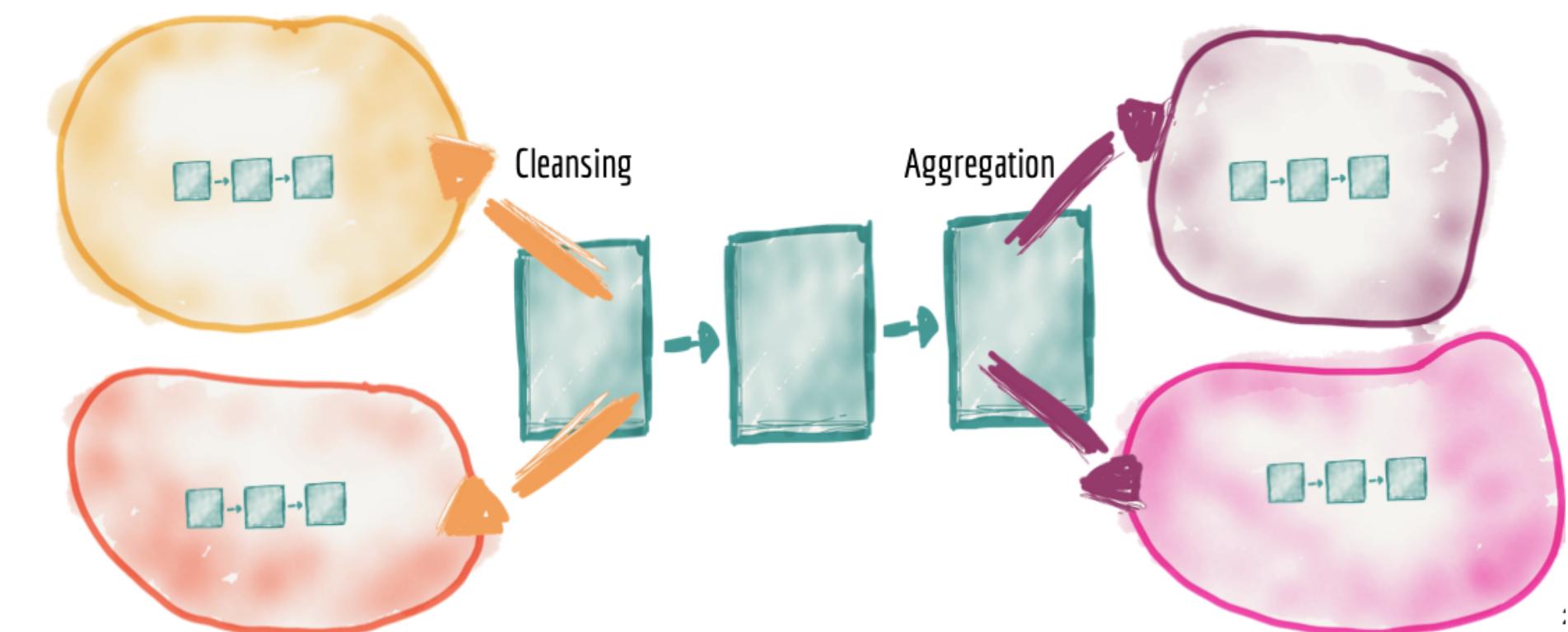
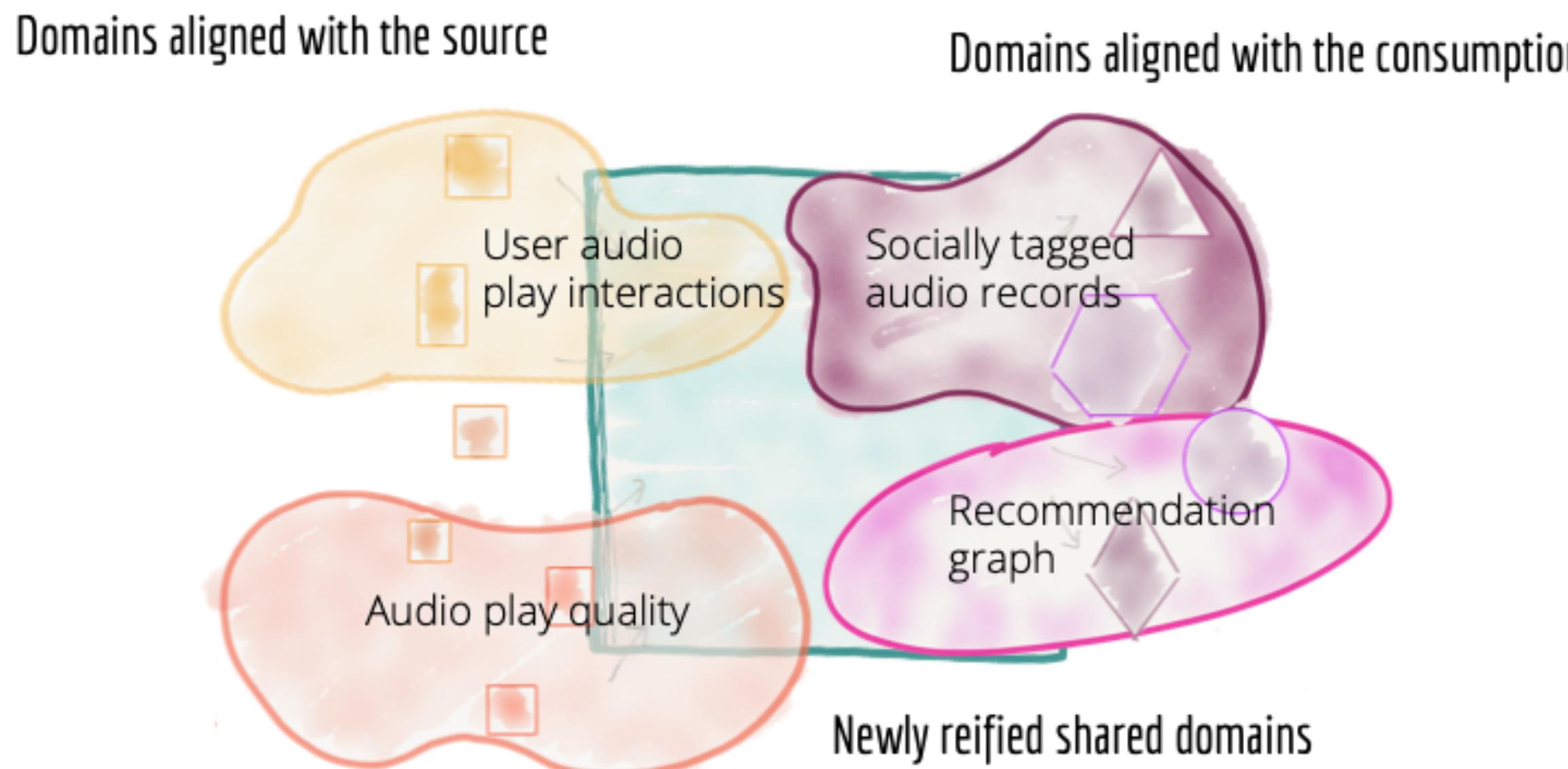
Architectural failure modes

Siloed and hyper-specialized ownership



Data mesh

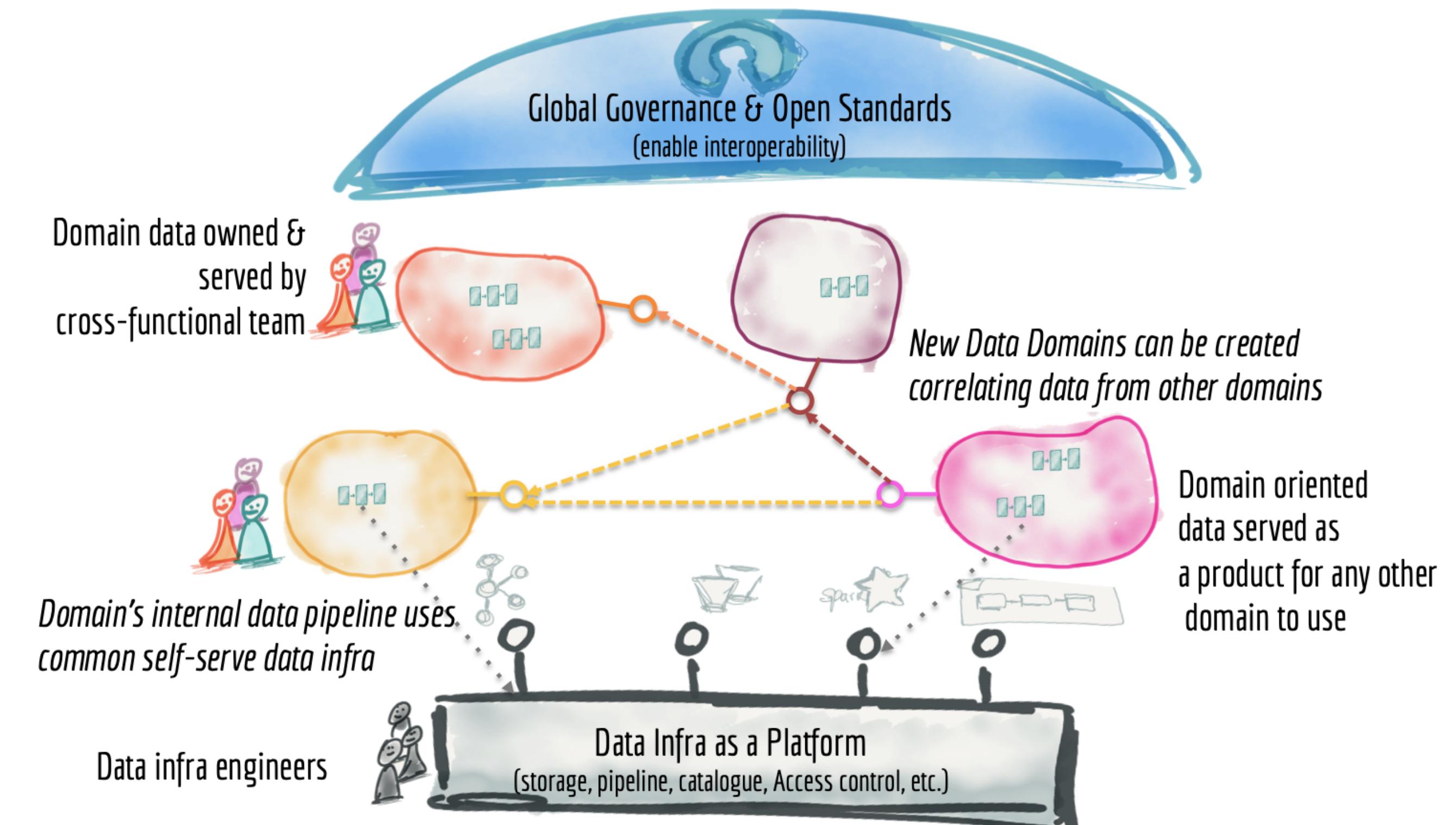
Decomposing the architecture and teams owning the data based on domains - source, consumer, and newly created shared domains



Data mesh

The third and current generation data platforms are more or less similar to the previous generation

- A. Streaming for real-time data availability with architectures such as Kappa,
- B. Unifying the batch and stream processing for data transformation with frameworks such as Apache Beam, as well as
- C. Fully embracing cloud-based managed services for storage, data pipeline execution engines, and machine learning platforms



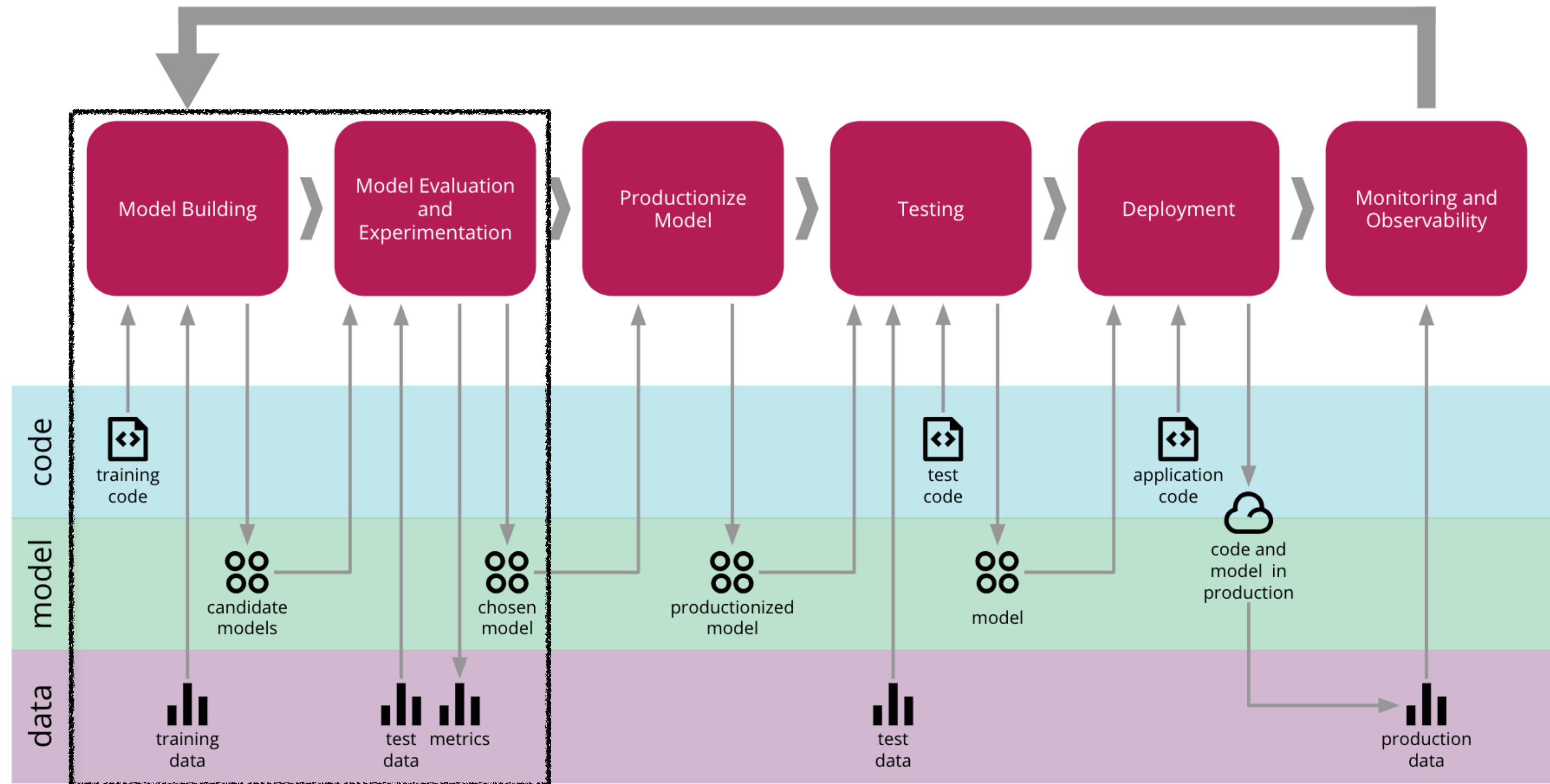
Discoverable Data

- A data product must be easily discoverable.
- A common implementation is to have a registry, a data catalogue, of all available data products with their meta information such as their owners, source of origin, lineage, sample datasets, etc.
- This centralized discoverability service allows data consumers, engineers and scientists in an organization, to find a dataset of their interest easily.
- Each domain data product must register itself with this centralized data catalogue for easy discoverability.

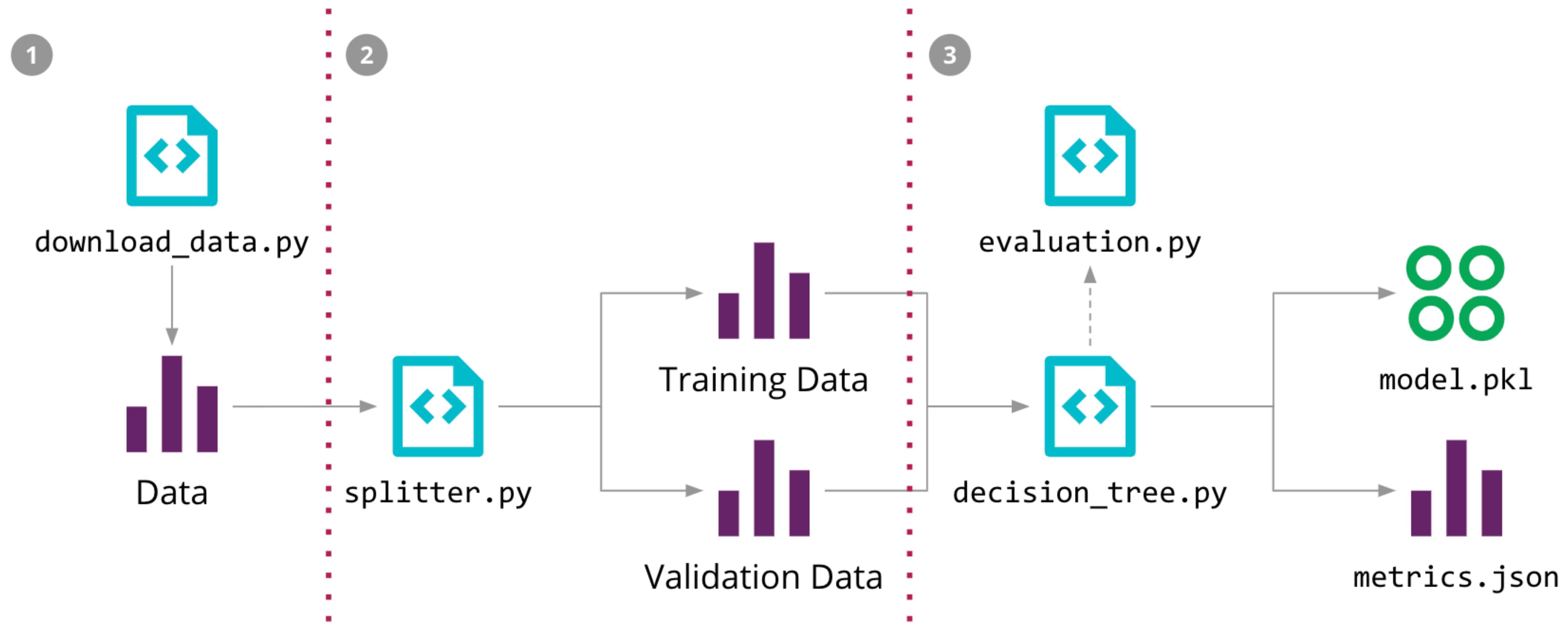
Addressable Data

- A data product, once discovered, should have a unique address following a global convention that helps its users to programmatically access it.
- Organizations may adopt different naming conventions for their data, depending on the underlying storage and format of the data.
- Considering the ease of use as an objective, in a decentralized architecture, it is necessary for common conventions to be developed.
- Different domains might store and serve their datasets in different formats, events might be stored and accessed through streams such as Kafka topics, columnar datasets might use CSV files, or AWS S3 buckets of serialized Parquet files.
- A standard for addressability of datasets in a polyglot environment removes friction when finding and accessing information.

An End-to-End ML System Building Process

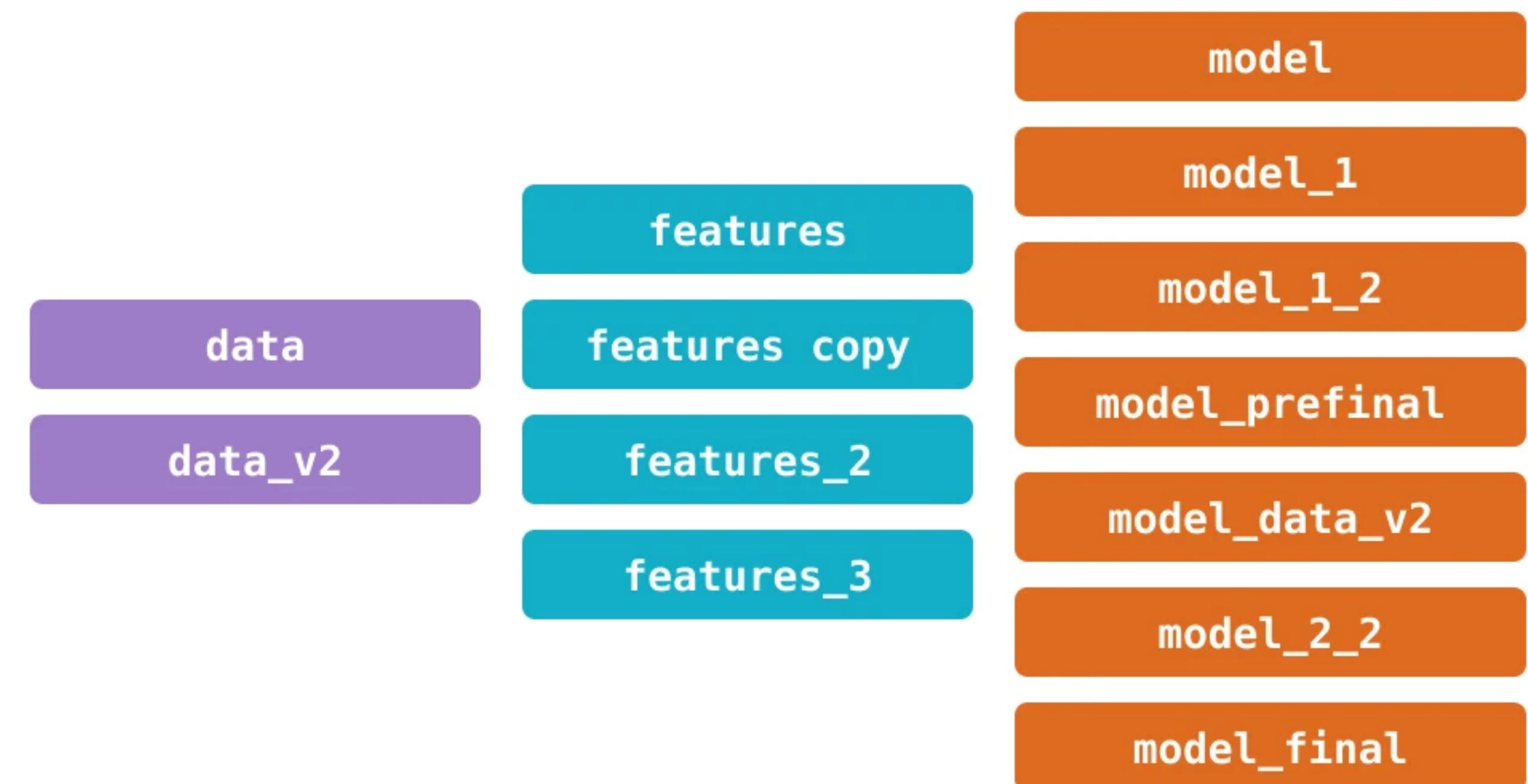


ML Model Training Pipeline



Versioning Data and Models

- How do we **keep track of changes** in data, source code, and ML models together?
- What's the best way to organize and store **variations** of these files and directories?

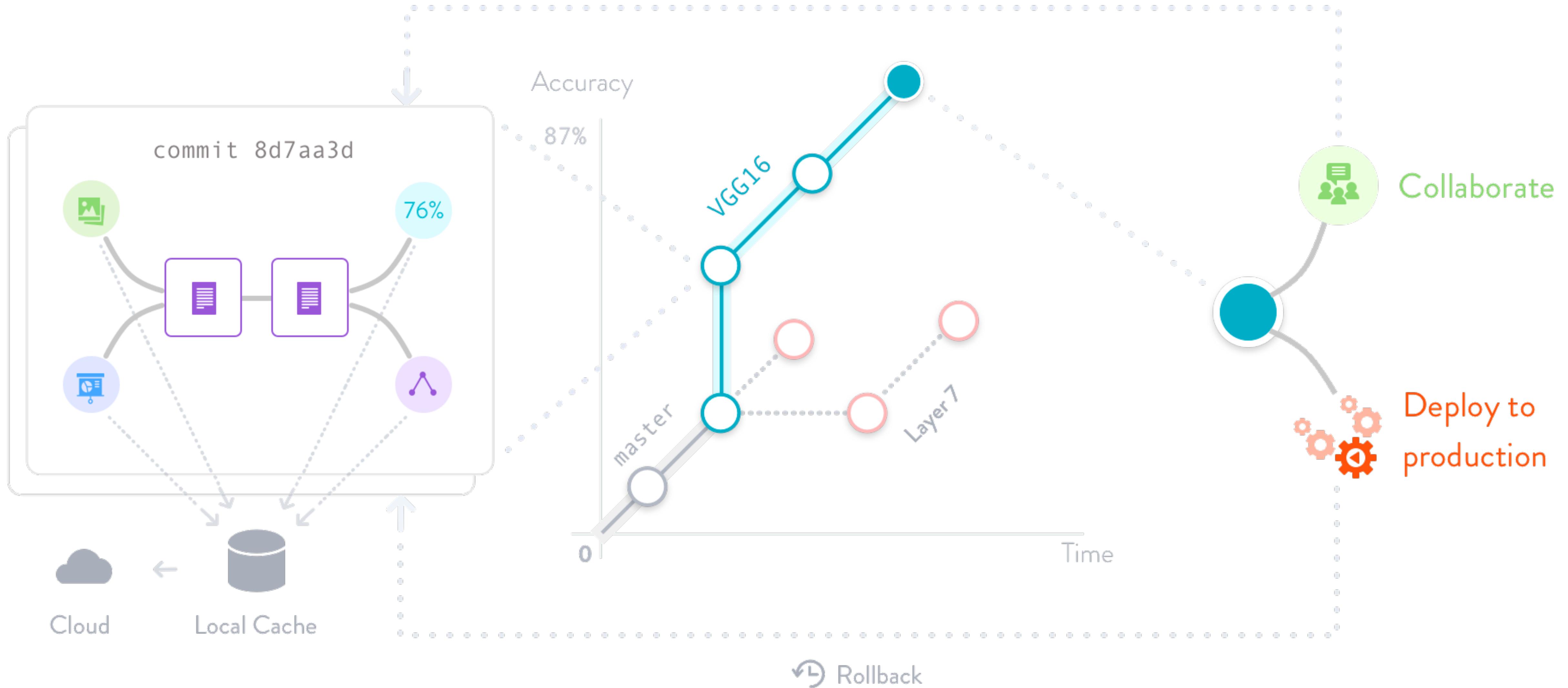


How about bookkeeping?

Another problem in the field has to do with bookkeeping

- being able to identify past data inputs and processes to understand their results, for knowledge sharing, or for debugging.

Configure ML pipeline: DVC tracks ML models and data sets



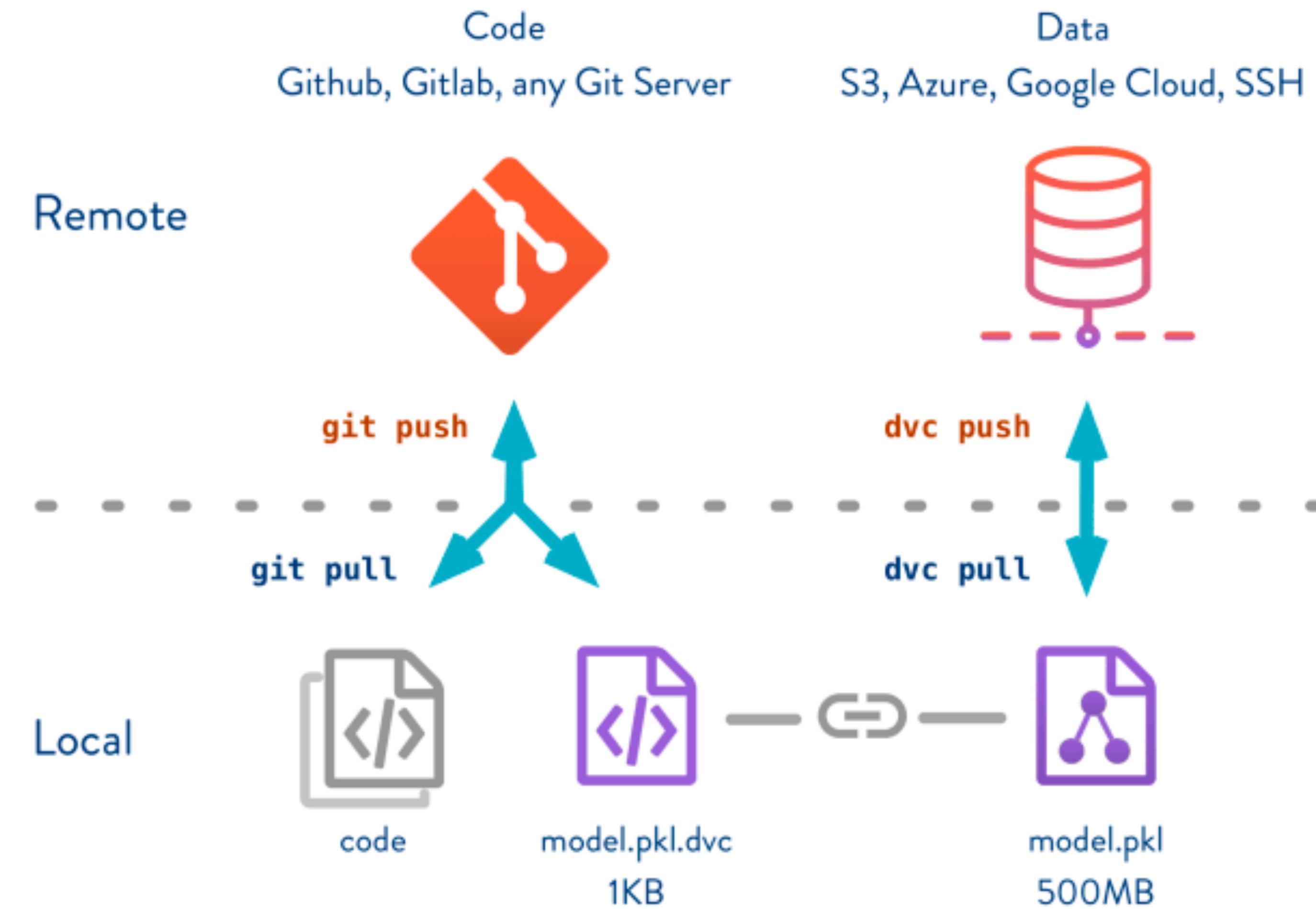
Configure ML pipeline: DVC tracks ML models and data sets

```
dvc run -f input.dvc \ ①
  -d src/download_data.py -o data/raw/store47-2016.csv python src/download_data.py
dvc run -f split.dvc \ ②
  -d data/raw/store47-2016.csv -d src/splitter.py \
  -o data/splitter/train.csv -o data/splitter/validation.csv python src/splitter.py
dvc run ③
  -d data/splitter/train.csv -d data/splitter/validation.csv -d src/decision_tree.py \
  -o data/decision_tree/model.pkl -M results/metrics.json python src/decision_tree.py
```

Configure ML pipeline: DVC tracks ML models and data sets

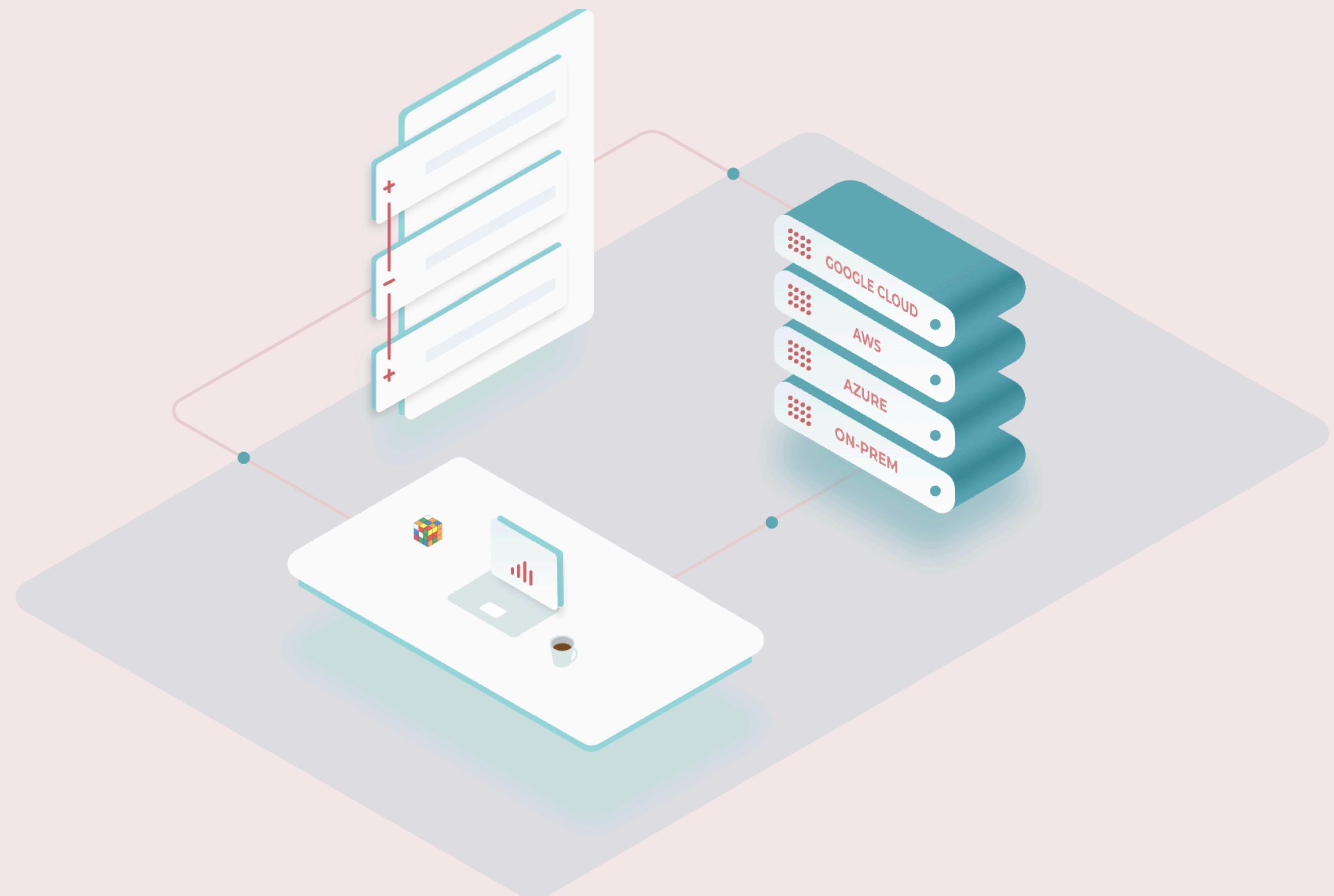
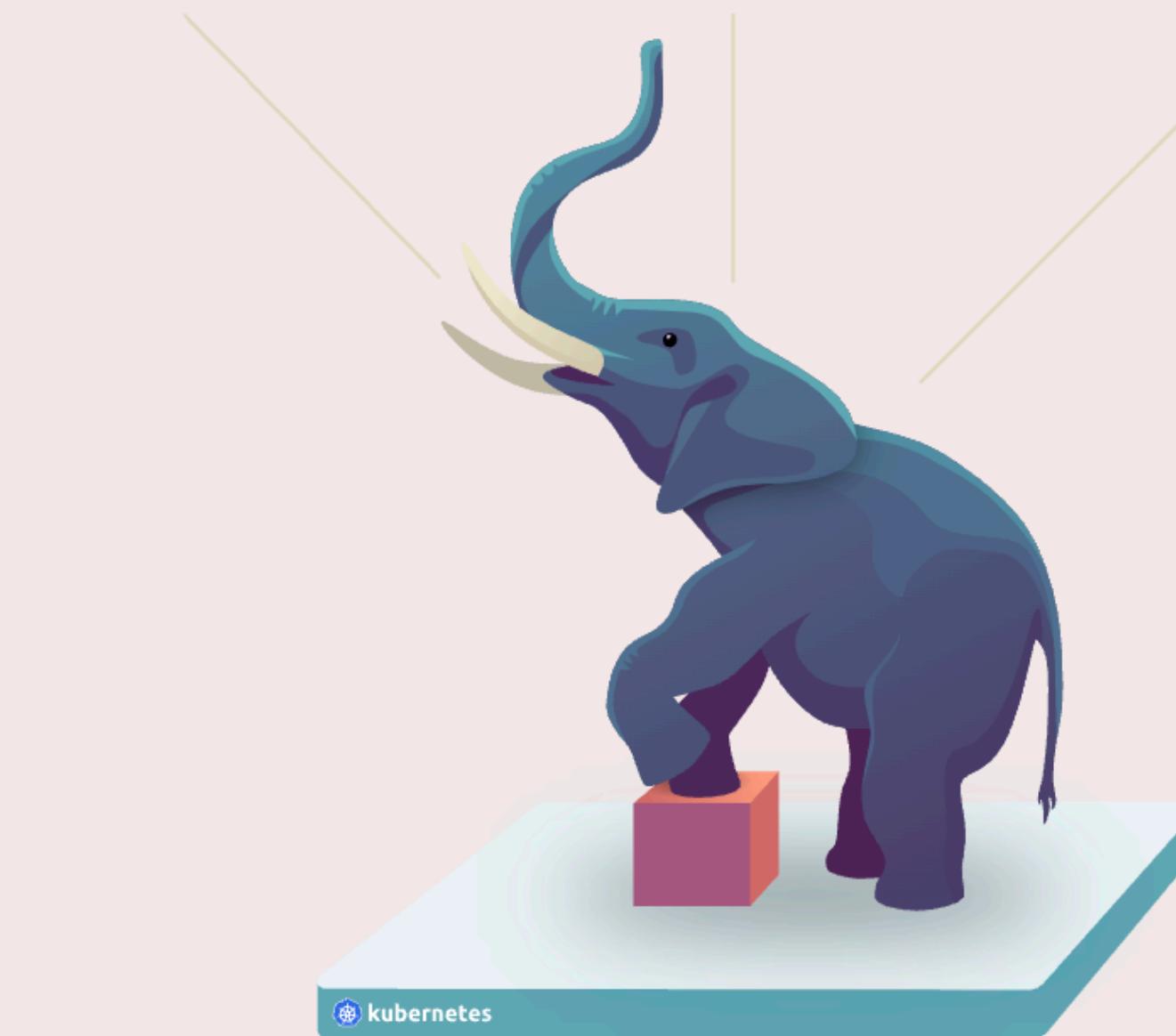
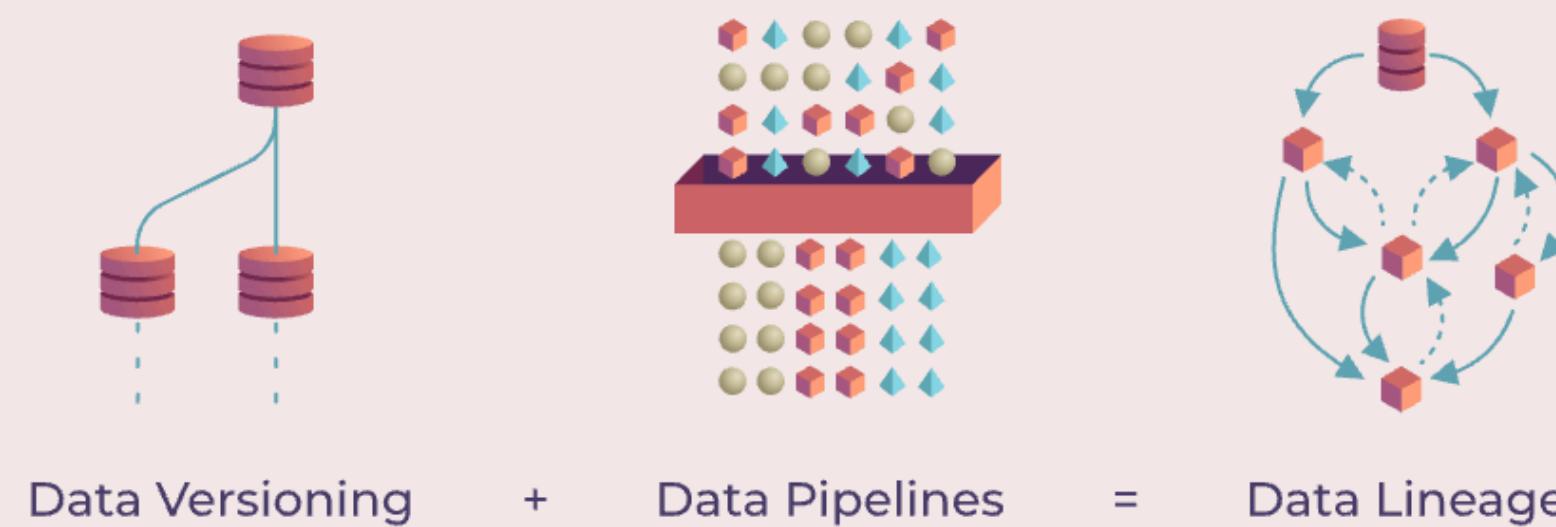
- Each run will create a file, that can be committed to version control
- DVC allows other people to reproduce the entire ML pipeline, by executing the *dvc repro* command.
- Once we find a suitable model, we will treat it as an artifact that needs to be *versioned* and *deployed* to production.
- With DVC, we can use the *dvc push* and *dvc pull* commands to publish and fetch it from *external storage*.

Configure ML pipeline: DVC tracks ML models and data sets

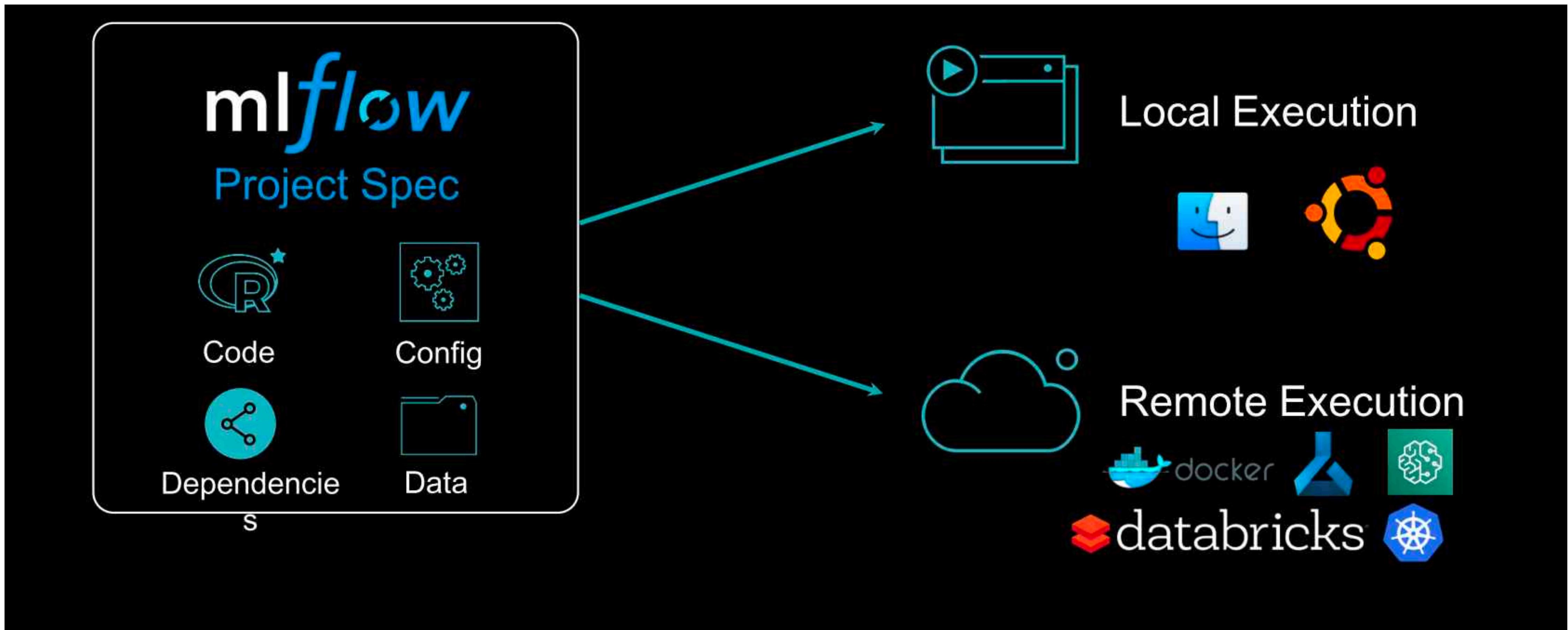


There are other open source tools for versioning

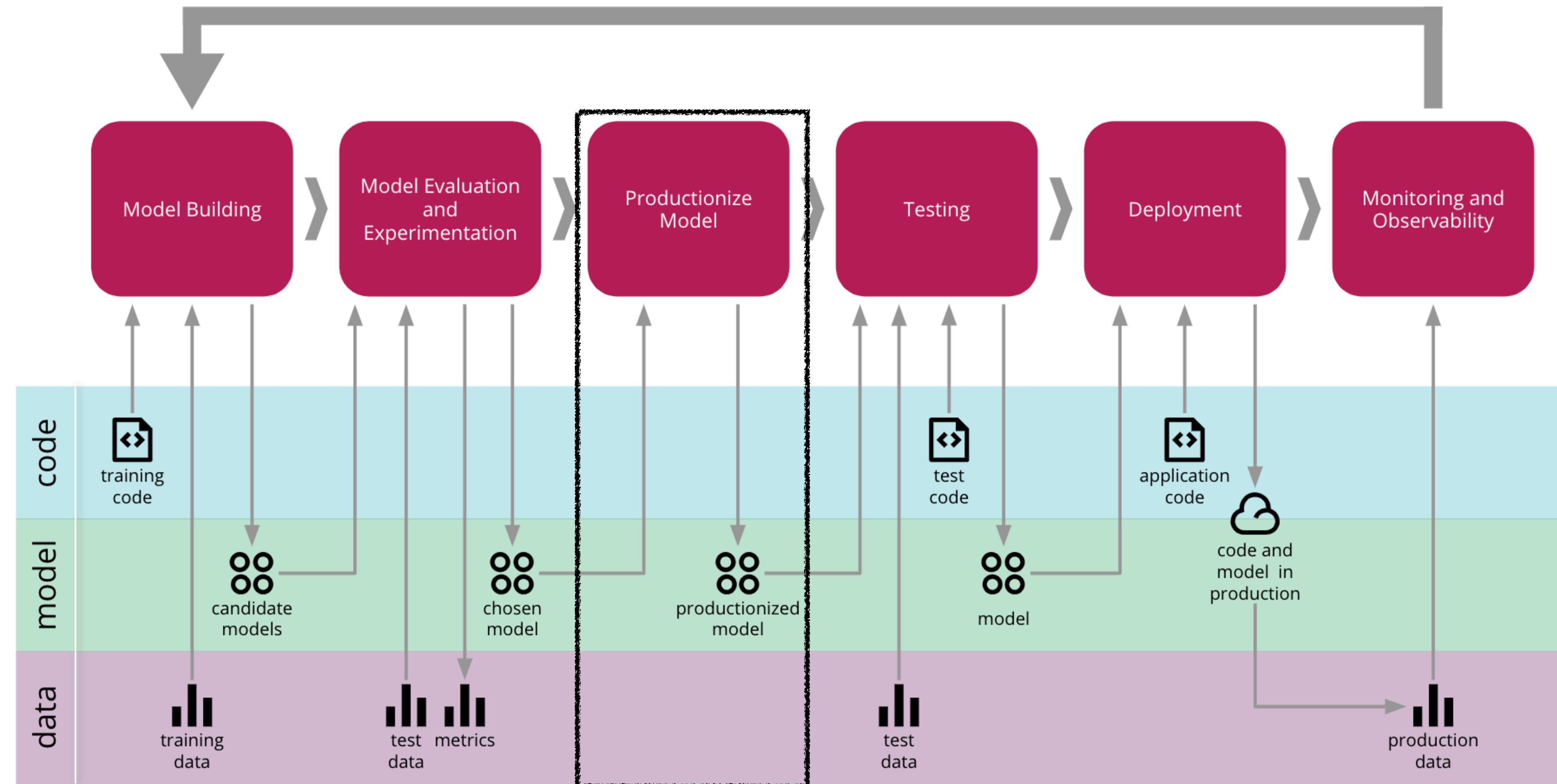
Pachyderm



There are other open source tools for versioning MLflow



An End-to-End ML System Building Process

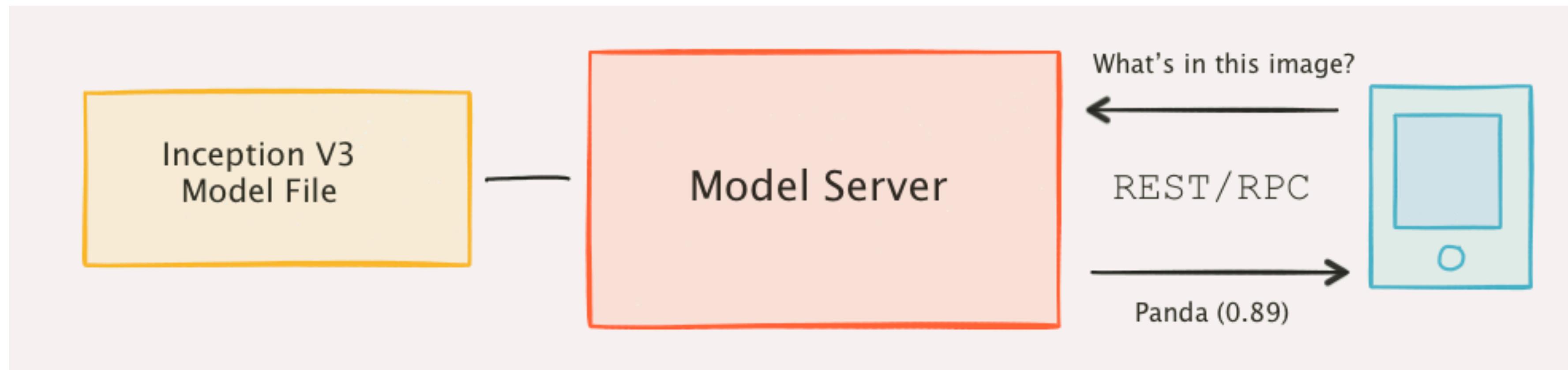


Model serving patterns

- **Embedded model:** treat the model artifact as a dependency that is built and packaged within the consuming application.
- **Model deployed as a separate service:** The model is wrapped in a service that can be deployed independently of the consuming applications.
- **Model published as data:** The model is also treated and published independently, but the consuming application will ingest it as data at runtime.

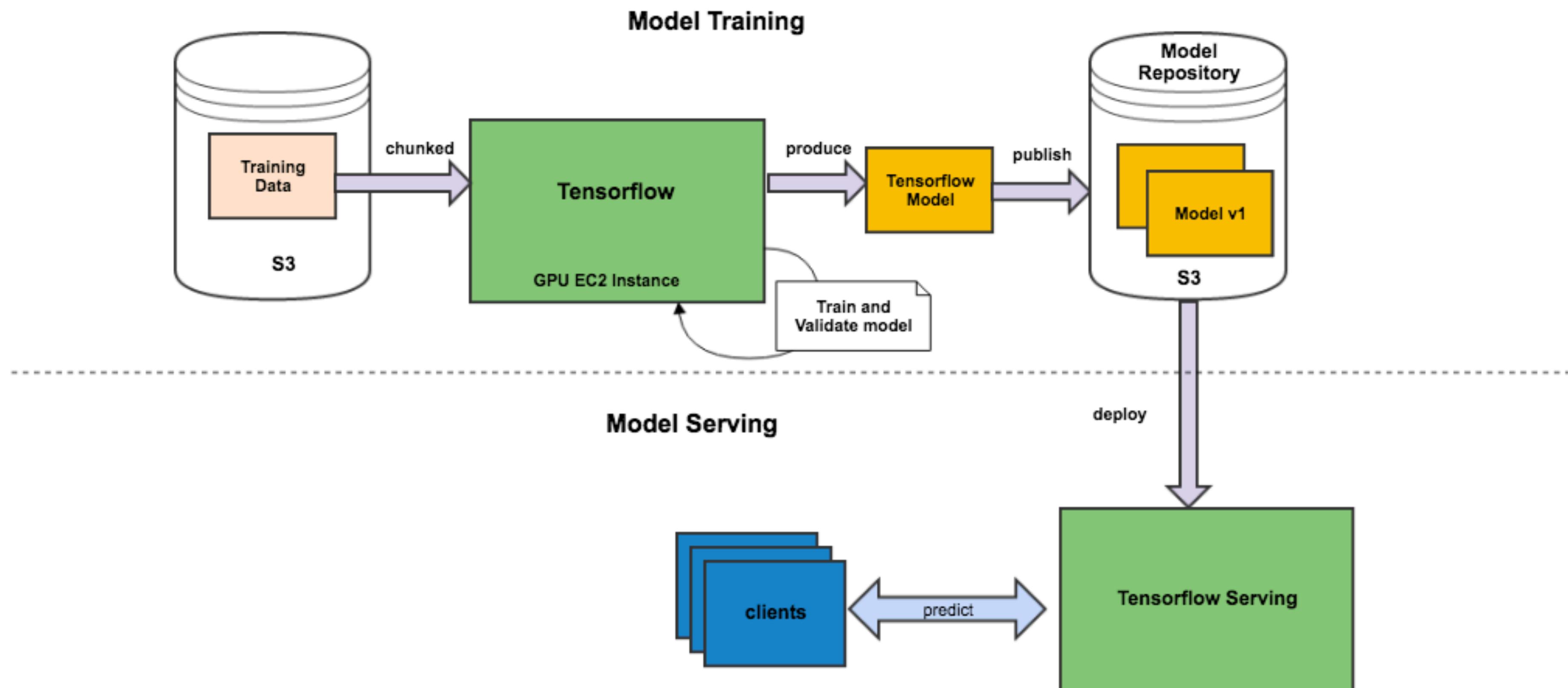
Model Serving

Abstract level



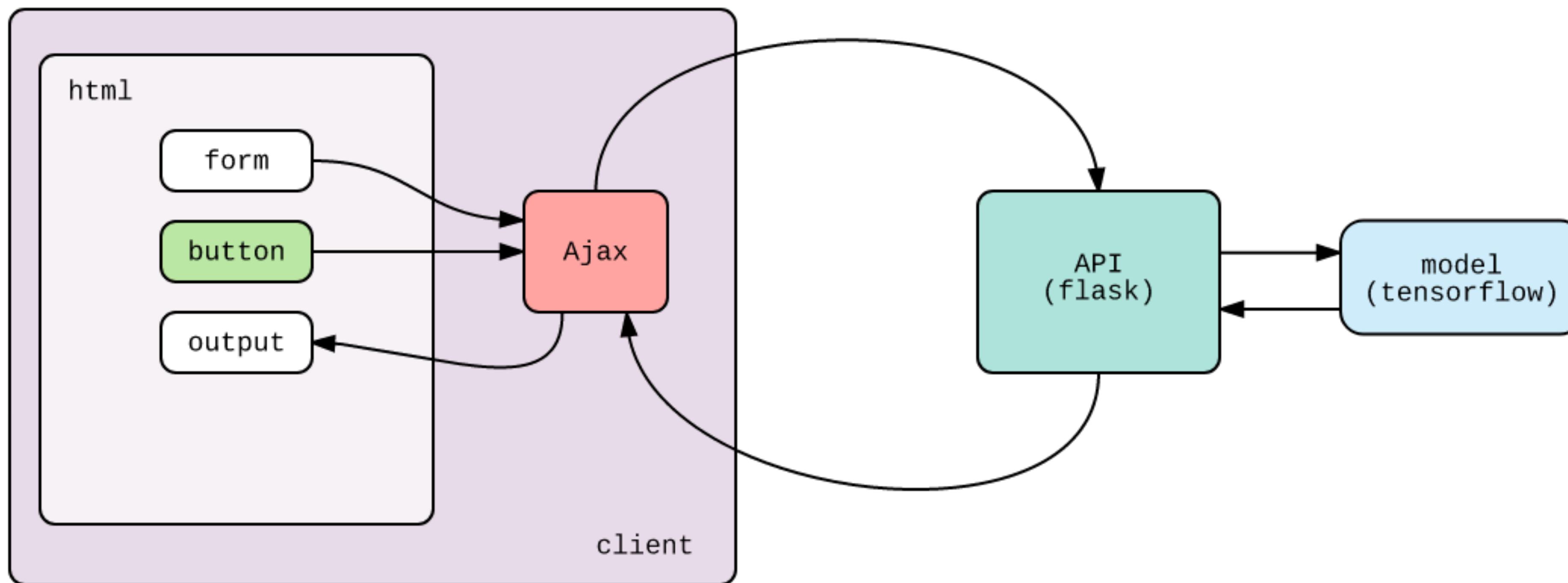
Model Serving

TF Serving



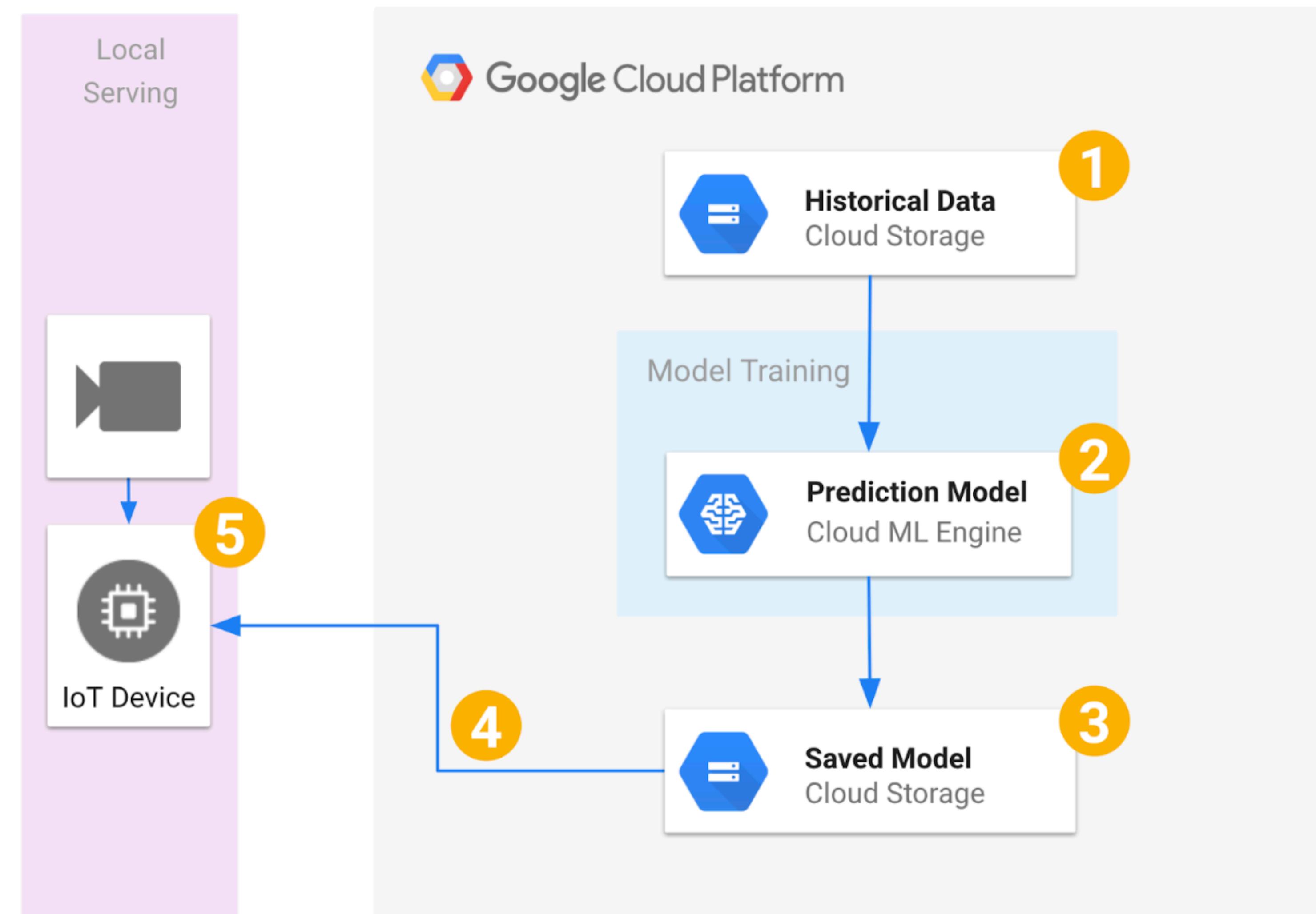
Model Serving

Web app

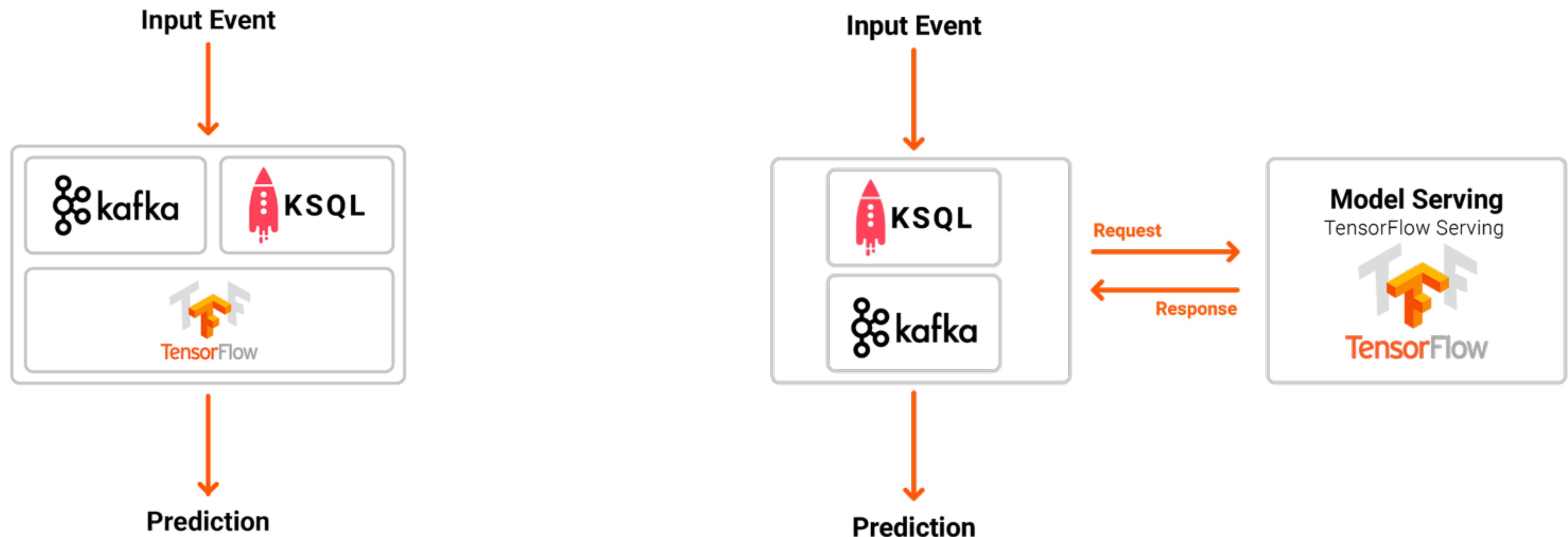


Model Serving

Internet of Thing



Model Serving Stream Processing System



Model Serving

Embedded model

- Simple approach
- You treat the model artifact as a dependency that is built and packaged within the consuming application.
- You can treat the application artifact and version as being a combination of the application code and the chosen model.

Model Serving

Model deployed as a separate service

- The model is wrapped in a service that can be deployed independently of the consuming applications.
- This allows updates to the model to be released independently, but it can also introduce latency at inference time
- There will be some sort of remote invocation required for each prediction.

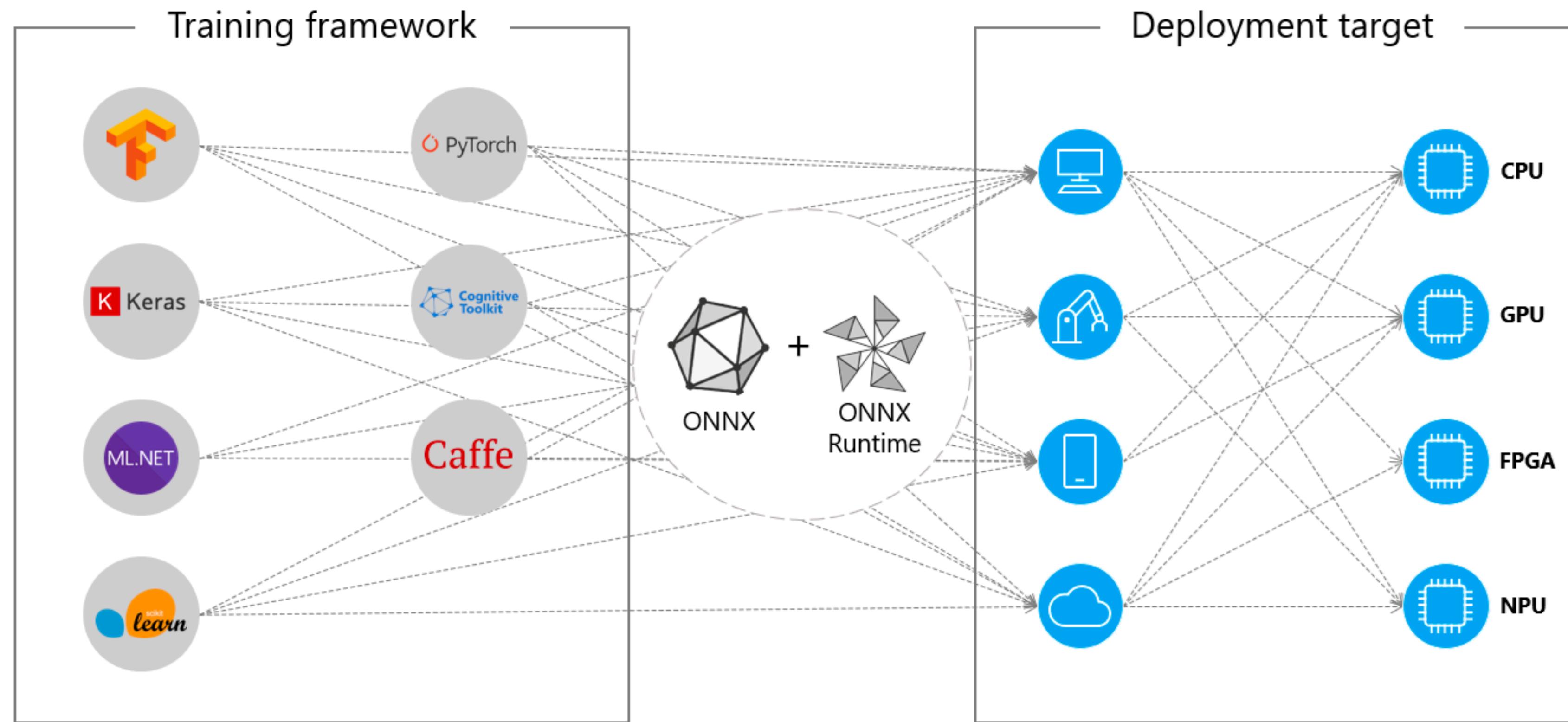
Model Serving

Model published as data

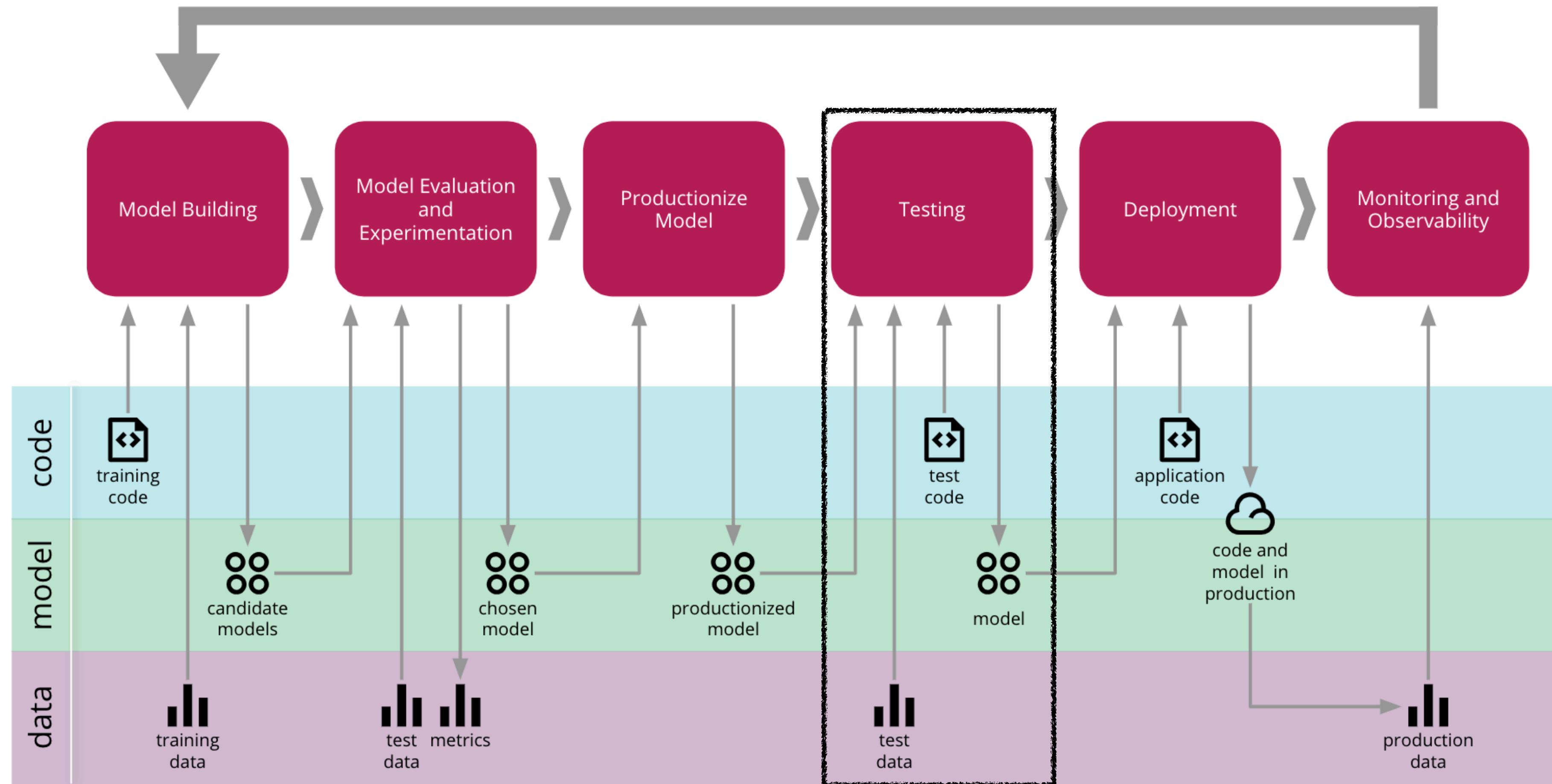
- The model is also treated and published independently,
- But the consuming application will ingest it as data at runtime.
- We have seen this used in streaming/real-time scenarios where the application can subscribe to events that are published whenever a new model version is released, and ingest them into memory while continuing to predict using the previous version.
- Software release patterns such as Canary Releases can also be applied in this scenario.

Export ML models to production environment

Open Neural Network Exchange



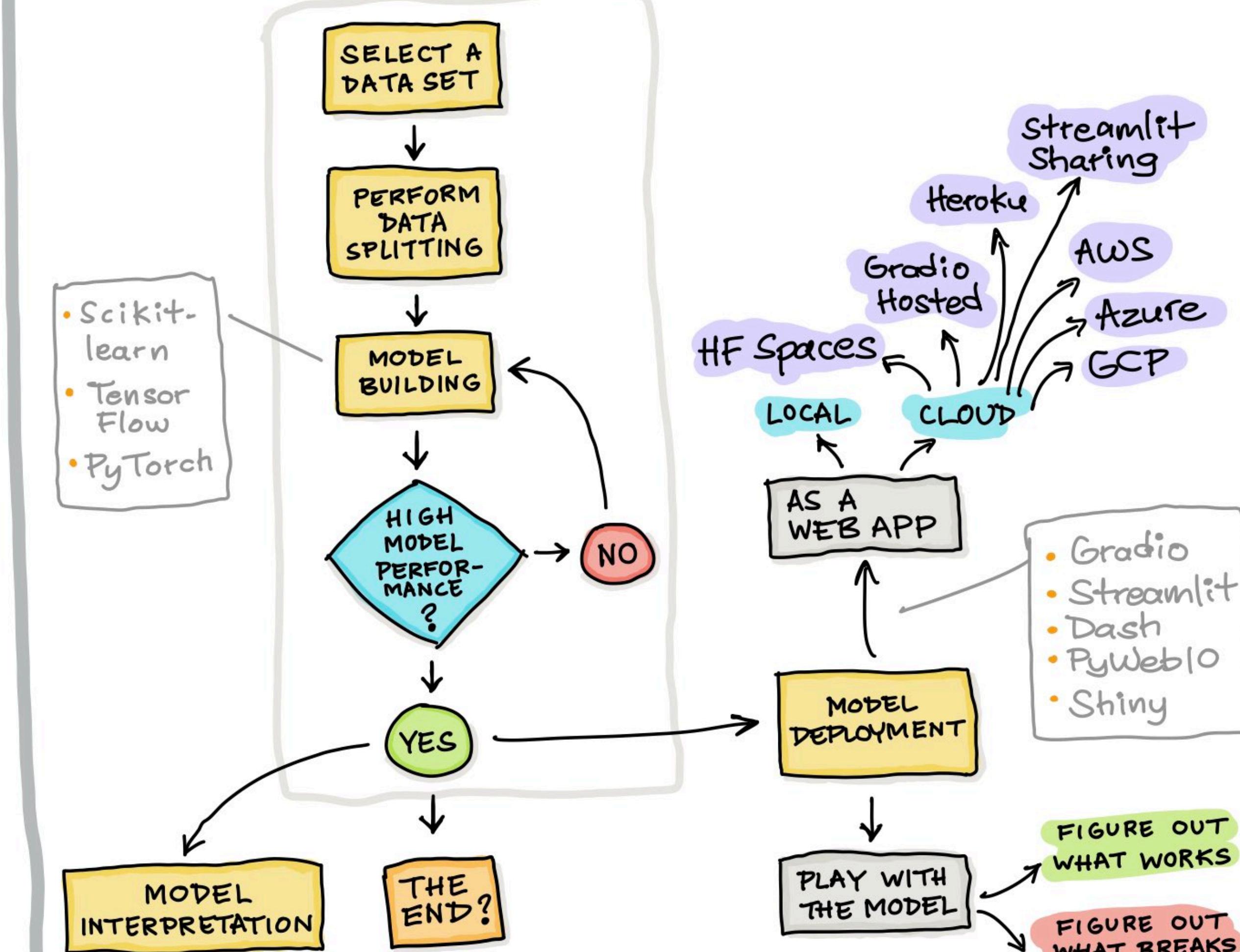
An End-to-End ML System Building Process



Testing and Quality in Machine Learning

- Regardless of which pattern you decide to use, there is always an implicit contract between the model and its consumers.
- The model will usually expect input data in a certain shape, and if Data Scientists change that contract to require new input or add new features, you can cause integration issues and break the applications using it.
- So testing becomes important.

QUICKLY DEPLOY MACHINE LEARNING MODELS

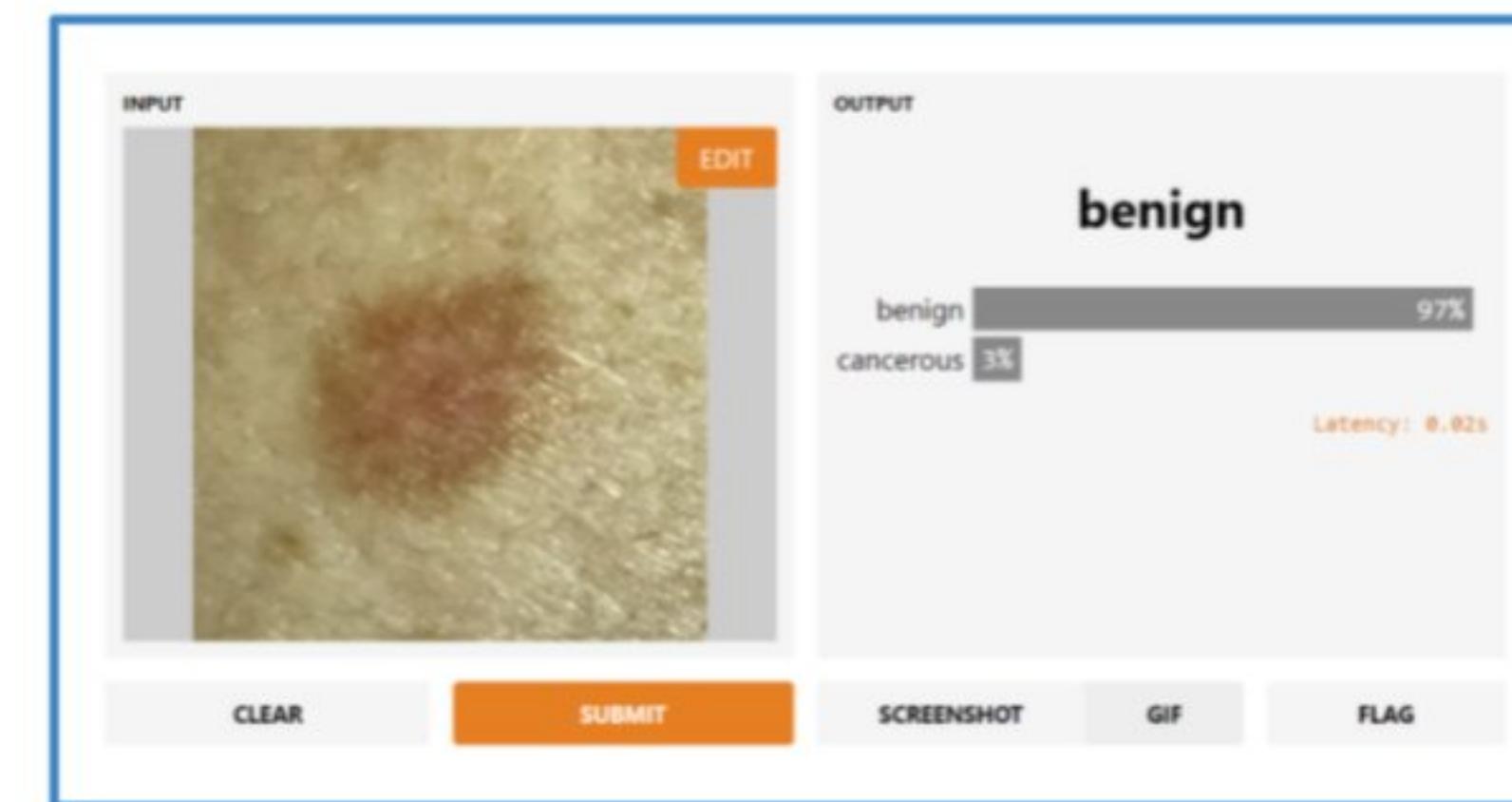


DRAWN BY CHANIN NANTASENAMAT
DATA PROFESSOR <http://youtube.com/dataprofessor>

Gradio powering clinical trials of machine learning models



1 Interdisciplinary team build/improve skin cancer classification model



2 They deploy it easily using Gradio



4 Model's mistakes collected



3 Model undergoes clinical validation at multiple sites

Testing Machine Learning Systems

Validating data

- Tests to **validate input data** against the expected schema, or to validate our **assumptions** about its valid values:
 - Values fall within expected ranges
 - Values are not null
- Unit tests to check **features** are calculated correctly:
 - Numeric features are scaled or normalized,
 - One-hot encoded vectors contain all zeroes and a single 1
 - Missing values are replaced appropriately

Testing Machine Learning Systems

Validating component integration

- Test the **integration** between different services:
 - Contract Tests to validate that the expected model interface is compatible with the consuming application.
- Test that the **exported model** still produces the same results:
 - Running the original and the productionized models against the same validation dataset, and comparing the results are the same.

Testing Machine Learning Systems

Validating the model quality

- ML **model performance** is non-deterministic.
- Collect and monitor **metrics** to evaluate a model's performance,
 - Error rates, accuracy
 - Precision, recall
 - AUC, ROC, confusion matrix
- **Threshold Tests** in our pipeline, to ensure that new models don't degrade against a known performance baseline.

Testing Machine Learning Systems

Validating model bias and fairness

- Check how the model performs against **baselines** for specific **data slices**:
 - Inherent bias in the training data where there are many more data points for a given value of a feature (e.g. race, gender, or region) compared to the actual distribution in the real world.
 - A tool like **Facets** can help you visualize those slices and the distribution of values across the features in your datasets.

Testing Machine Learning Systems

Integration Test

- When models are **distributed or exported** to be used by a different application, the engineered features are **calculated differently** between training and serving time.
- Distribute a **holdout dataset** along with the model artifact, and allow the consuming application team to reassess the model's performance against the holdout dataset after it is integrated.
- This would be the equivalent of a broad **Integration Test** in traditional software development.

Governance process for ML Systems

Experiments Tracking

- To capture and display information that will allow humans to decide if and which model should be promoted to production.
- It is common that you will have multiple experiments being tried in parallel, and many of them might not ever make it to production.
- The code for many of these experiments will be thrown away, and only a few of them will be deemed worthy of making it to production.
- Different Git branches to track the different experiments in source control.
- Tools such as DVC can fetch and display metrics from experiments running in different branches or tags, making it easy to navigate between them.

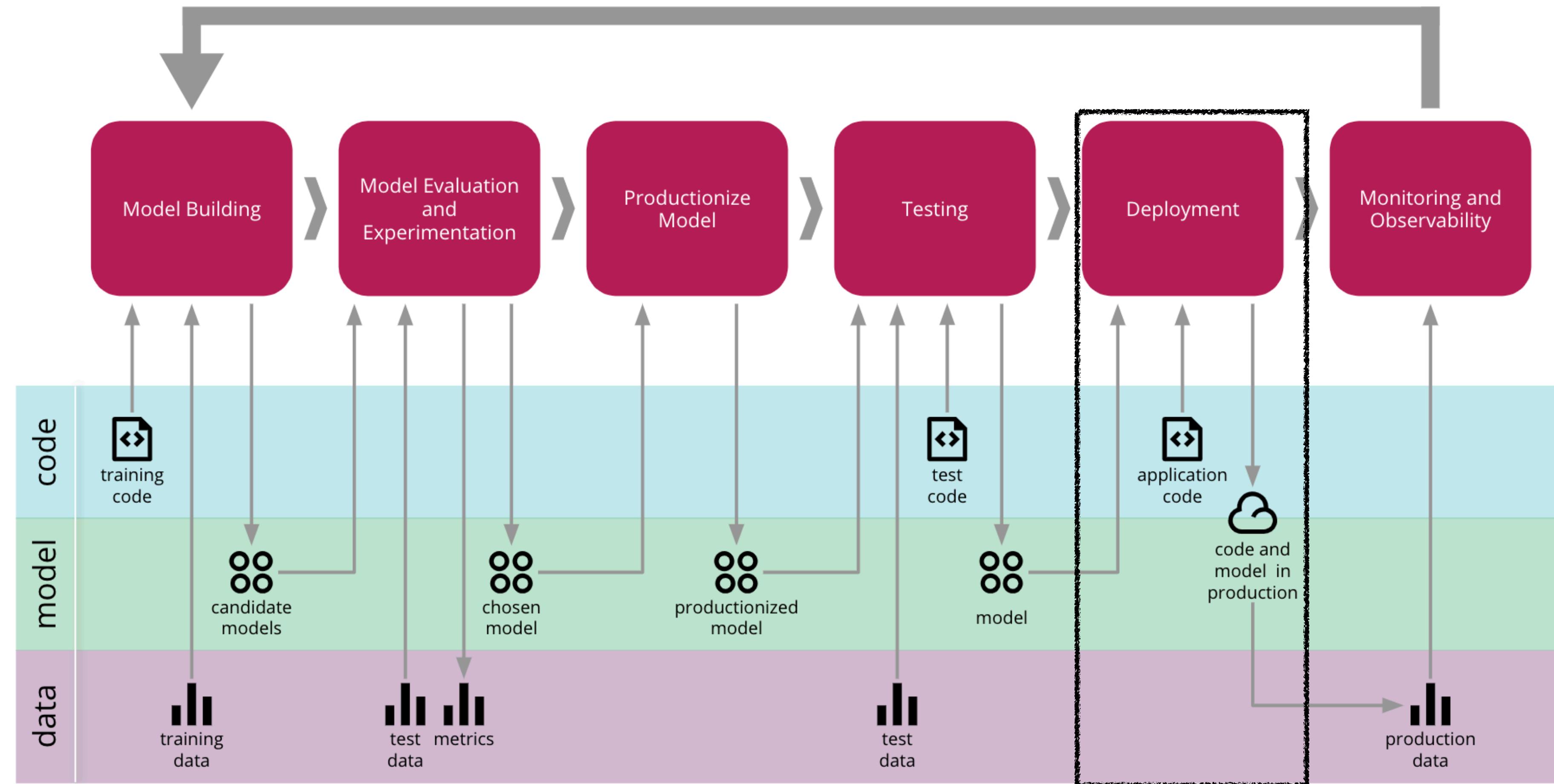
Governance process for ML Systems

MLflow Tracking web UI

The screenshot shows the MLflow Tracking web UI interface. At the top, there is a dark header bar with the 'mlflow' logo on the left and 'GitHub Docs' links on the right. Below the header, the main page title is 'Experiments' with a back arrow and the user identifier 'user1'. On the left side, there is a sidebar with user names: 'user2' (disabled), 'user1' (selected), and another 'user1' (disabled). The main content area displays experiment details for 'user1': 'Experiment ID: 1' and 'Artifact Location: gs://cd4ml-mlflow-tracking/1'. It includes search and filter controls: 'Search Runs: metrics.rmse < 1 and params.model = "tree"', 'State: Active', 'Search' button, 'Filter Params: alpha, lr', 'Filter Metrics: rmse, r2', and a 'Clear' button. Below these controls, there is a summary: '1 matching run' with buttons for 'Compare', 'Delete', 'Download CSV', and two icons. A table then lists the single matching run with columns: Date, User, Run Name, Source, Version, Parameters, and Metrics. The run details are: Date 2019-04-28 00:03:29, User go, Run Name 5, Source decision_tree.py, Version b24402, Parameters model: RANDOM_FOREST, n_estimators: 10, Metrics nwrmsle: 0.743, r2_score: 0.109.

Date	User	Run Name	Source	Version	Parameters	Metrics
2019-04-28 00:03:29	go	5	decision_tree.py	b24402	model: RANDOM_FOREST, n_estimators: 10	nwrmsle: 0.743, r2_score: 0.109

An End-to-End ML System Building Process



Model Serving vs. Deployment

What is the difference?

- **Model Serving Runtime:** Packaging a trained machine learning model into a container and setting up APIs so it can handle incoming requests. This allows the model to be used in a production environment, responding to data inputs with predictions (inference).
- **Model Serving Platform:** An environment designed to dynamically scale the number of model containers in response to incoming traffic. Tools like KServe are examples of serving platforms. They manage the infrastructure to deploy and scale models efficiently, responding to varying traffic without manual intervention.
- **Model Deployment:** The process of integrating a packaged model into a serving platform and connecting it to the broader infrastructure, such as databases and downstream services. This ensures the model can access necessary data, perform its intended functions, and deliver inference results to consumers.

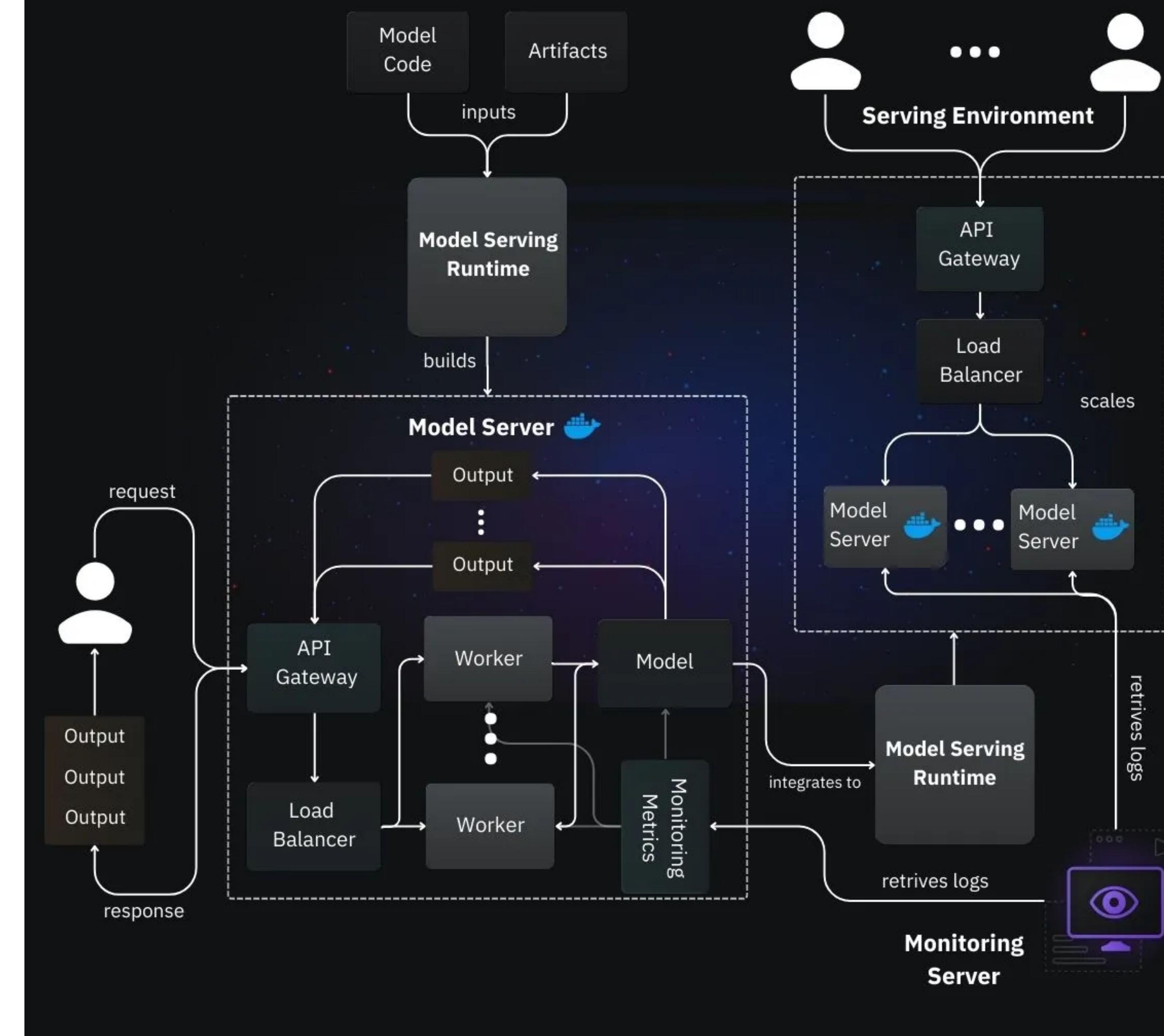
Model Serving vs. Deployment

Let's consider this typical ML model lifecycle

- **Train:** Suppose we trained an LLM on articles on USC website to assist students in taking courses each semester.
- **Package:** We use a model-serving runtime like BentoML to package our LLM into a Docker image, wrapped with a standardized, functional API.
- **Deploy:** We deploy our model packaged with BentoML to the KServe model-serving platform. This Kubernetes-based platform auto-scales the model containers according to incoming requests.
- **Integrate:** We connect the model with a chat widget in the sidebar of USC website so that students can type questions and receive the answers generated by our LLM. This requires us to integrate the necessary API requests into the website frontend code and ensure that our model's API can be reached publicly via the internet.
- **Bring value:** Finally, the model is ready to assist many students simultaneously.

Understanding model serving

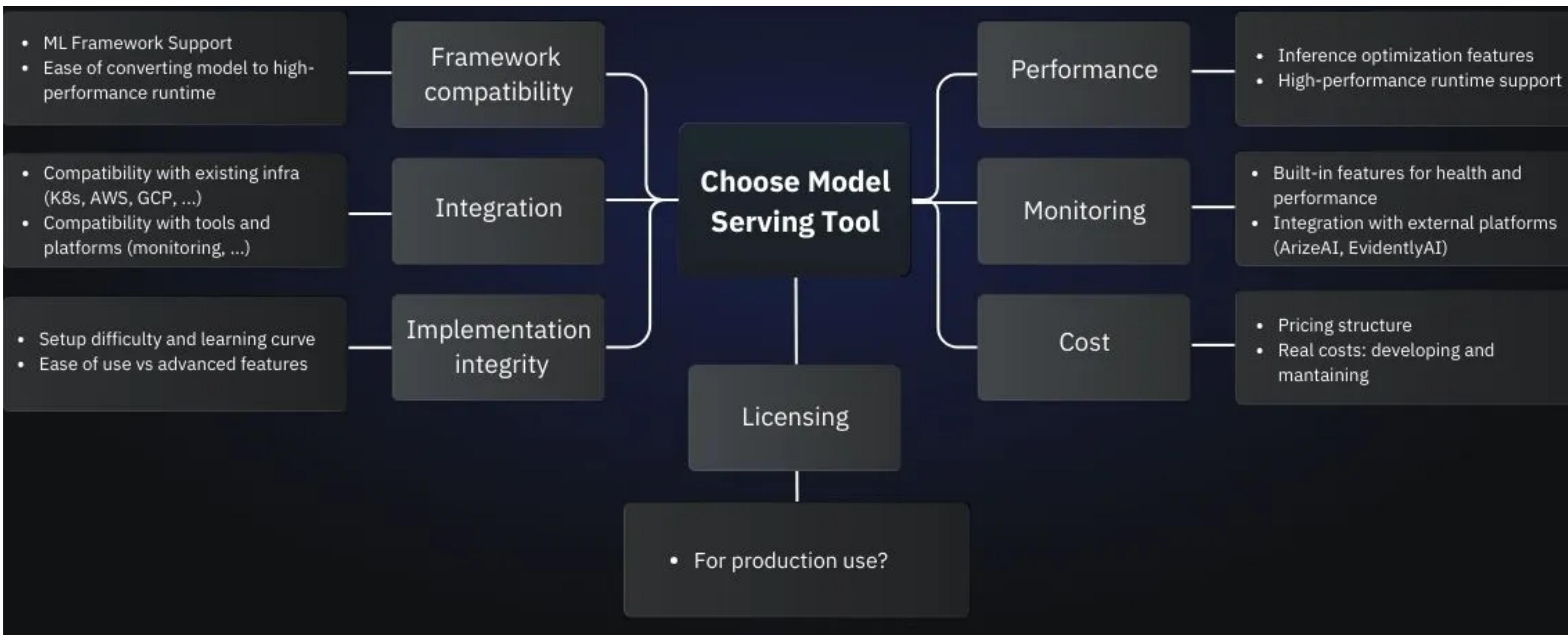
- Overview of the canonical model serving architecture and components.
- Given the model code and artifacts, we create a Docker image based on a model serving runtime.
- This Docker image contains the model server and is deployed to a model-serving platform that provides scalability and exposes the model to downstream users.



Three reasons to avoid using a model-serving runtime

- **Skill gap:** Some model-serving runtimes require significant software engineering skills on your team's part. If your team does not bring sufficient experience to the table, this can lead to challenges in setup, ongoing maintenance, and integration.
- **Batch processing:** When you don't need real-time inference, but all computations can be batch-processed, simpler solutions may be more straightforward and cost-effective than implementing a solution with a full-fledged serving runtime.
- **No scaling needs:** If your model does not need to be scaled because of low inference time or request volume, the benefits of using an ML-optimized container might not outweigh its engineering costs.

Criteria for selecting model-serving tools

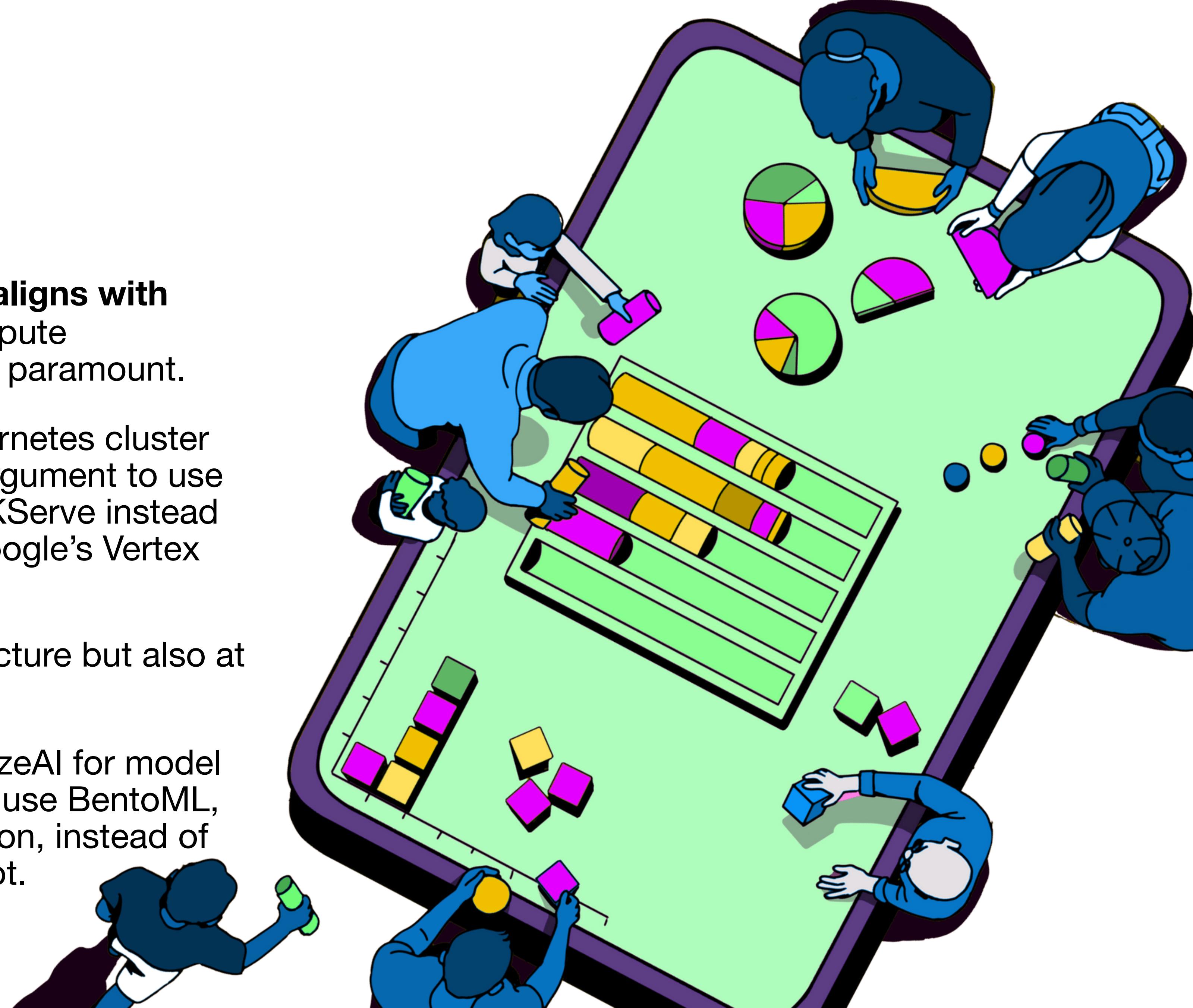


Framework compatibility

- Support for frameworks such as scikit-learn, TensorFlow, or PyTorch.
 - It would be unfortunate to select and begin setting up TorchServe only to discover later that it does not support the Keras model your colleague has trained.
- Whether the tool provides GPU support and works with the CUDA version you're on.
 - This is particularly important if you work with large deep-learning models.
- Support for distributed processing is crucial if you plan to scale your models across multiple machines to handle larger workloads.

Integration

- Assessing how a model-serving tool **aligns with your current MLOps stack** and compute infrastructure or cloud environment is paramount.
 - Suppose you already have a Kubernetes cluster running. That would be a strong argument to use a Kubernetes-native solution like KServe instead of a fully managed solution like Google's Vertex AI.
- This applies not only to your infrastructure but also at the **framework level**.
 - For example, if you plan to use ArizeAI for model observability, it would be better to use BentoML, which has out-of-the-box integration, instead of Tensorflow Serving, which does not.

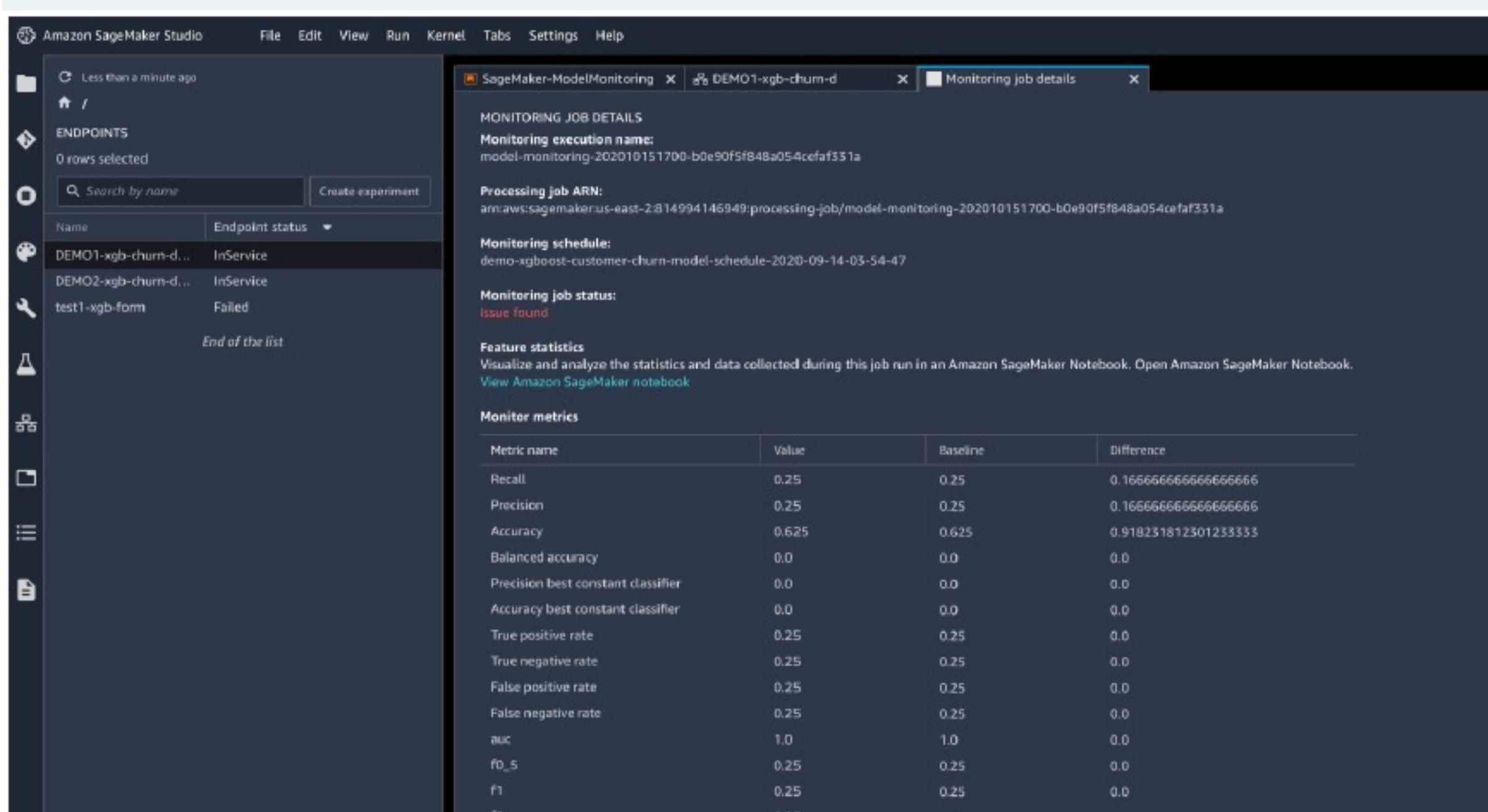
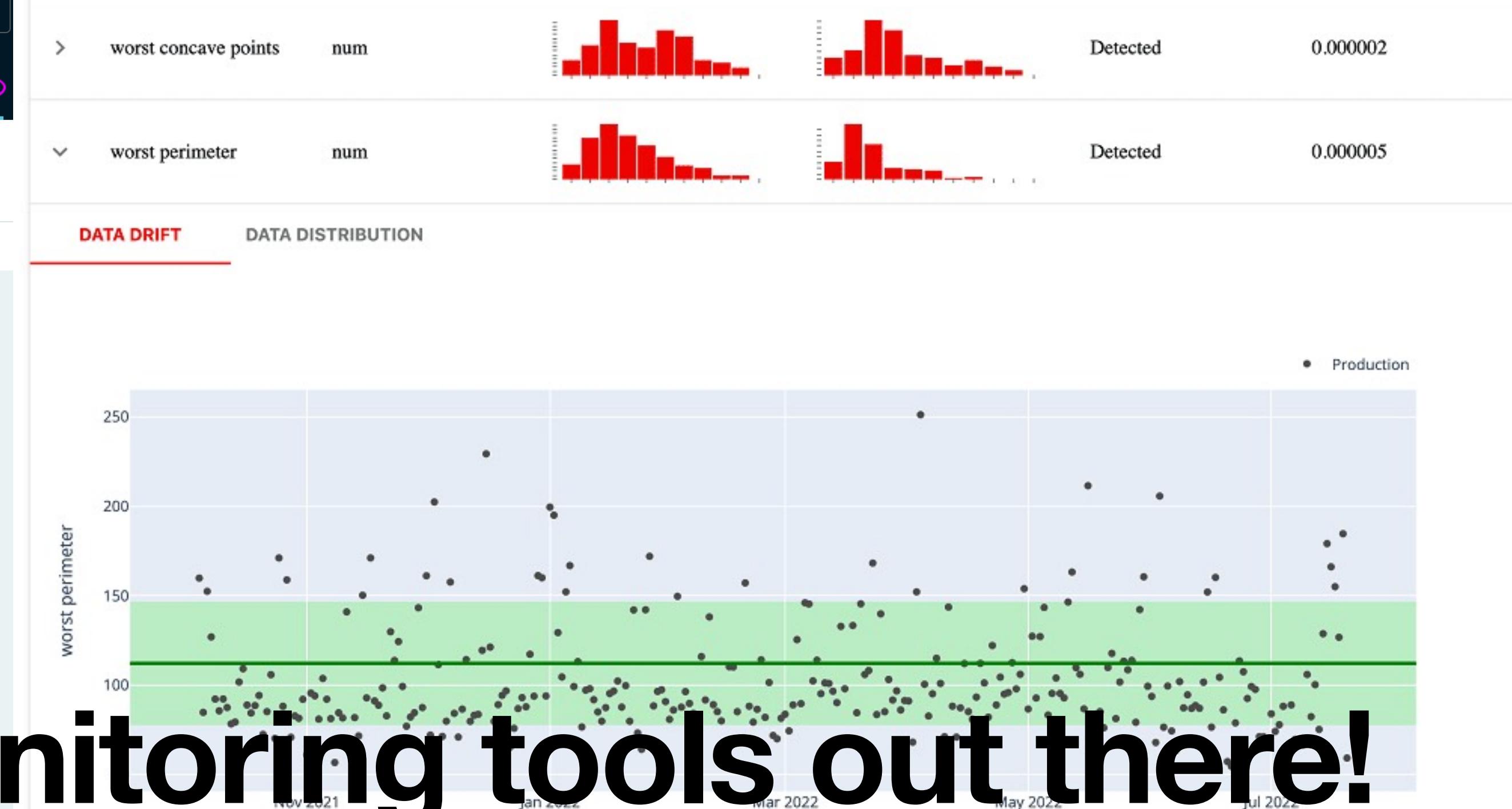
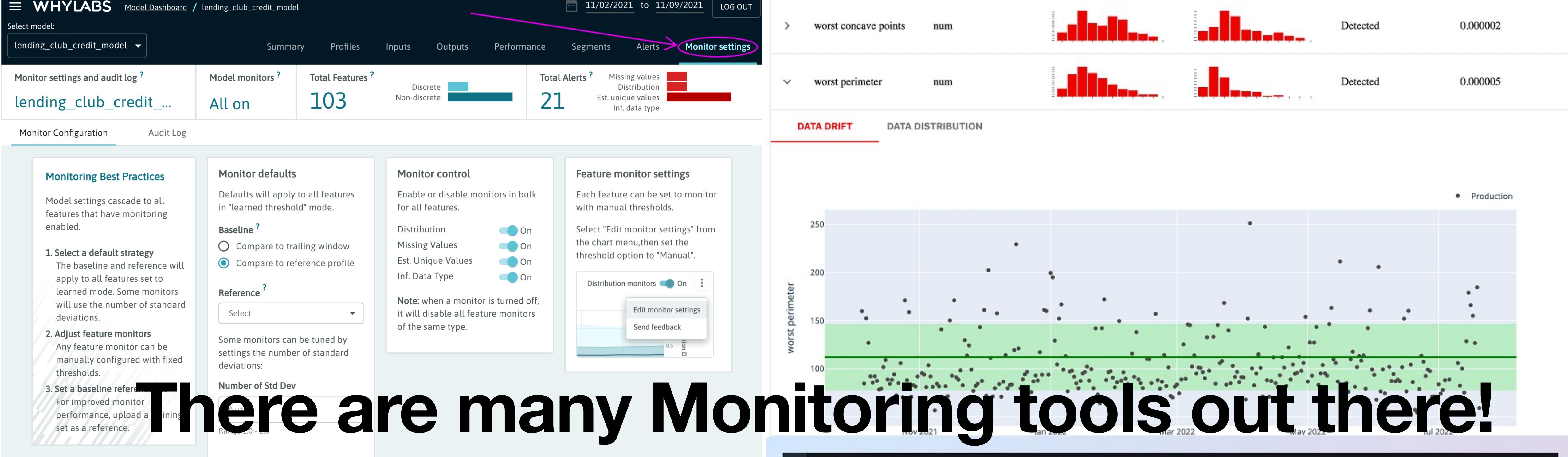


Implementation complexity

- Not every framework is **suitable for every team**,
 - The complexity of implementation and required background knowledge can vary significantly.
 - A tool that is difficult to use can slow down progress, especially if you are not familiar with the required technologies.
- Tools that provide **high flexibility** tend to be **more complex** and have a steeper learning curve.
 - This complexity arises because these tools offer more options and control to the user.
 - While this allows for better adaptation to specific needs, it also requires a deeper understanding.
- Choose the **simplest tool that meets your team's and project's needs**. This approach ensures you don't overcomplicate your setup with unnecessary features or struggle with a tool that's too limited for your requirements.

Performance

- **Concurrent model execution:** Spawns multiple instances of the same model to run simultaneously on a single hardware processor (GPU or CPU) and load balances the incoming requests across the instances.
 - This way, multiple smaller models can share one processor, saving costs.
- **Inference parallelization:** Distributes inference tasks across multiple hardware processors (GPU or CPU) to speed up processing.
- **Adaptive batching:** Allows the server to combine multiple inference requests into a single batch dynamically, optimizing throughput and latency.
- **High-performance runtime support:** Compute-intensive models benefit from conversion to a more efficient runtime such as TensorRT.
- **Asynchronous APIs:** Enable non-blocking requests, allowing the system to handle multiple requests at the same time.
 - This improves responsiveness as the system does not process the requests sequentially.
- **gRPC inference protocol:** Offers a more efficient alternative to traditional HTTP/REST for communication between services. In fact, the gRPC protocol has shown to be superior than REST in terms of response time.



There are many Monitoring tools out there!

Monitoring

Grafana + Prometheus



Prometheus

Logging

- Prometheus is a popular **open-source ML model monitoring** tool that collects multidimensional data and queries.
- The main advantages of Prometheus are tight integration with Kubernetes and many of the available exporters and client libraries, as well as a fast query language. Prometheus is also Docker-compatible and available on the Docker Hub.
- The Prometheus server has its own self-contained unit that does not depend on network storage or external services. So it doesn't require a lot of work to deploy additional infrastructure or software. Its main task is to store and monitor certain objects. An object can be anything: a Linux server, one of the processes, a database server, or any other component of the system. Each element that you want to monitor is called a metric.
- The Prometheus server **reads targets at an interval that you define to collect metrics** and stores them in a time series database. You set the targets and the time interval for reading the metrics. You query the Prometheus time series database for where metrics are stored using the PromQL query language.

Grafana Visualization

- Grafana allows you to visualize monitoring metrics. Grafana specializes in time series analytics. It can visualize the results of monitoring work in the form of line graphs, heat maps, and histograms.
- Instead of writing PromQL queries directly to the Prometheus server, you use Grafana GUI boards to request metrics from the Prometheus server and render them in the Grafana dashboard.

Model Deployment

Multiple models

- More than one model performing the same task.
 - Train a model to predict demand for each product.
 - Deploying the models as a separate service might be better for consuming applications to get predictions with a single API call.
 - You can later evolve how many models are needed behind that Published Interface.

Model Deployment

Shadow models

- Deploy the new model side-by-side with the current one, as a shadow model
- Send the same production traffic to gather data on how the shadow model performs before promoting it into the production.

Model Deployment

Competing models

- Multiple versions of the model in production – like an A/B test
 - Infrastructure and routing rules required to ensure the traffic is being redirected to the right models.
 - To gather enough data to make statistically significant decisions, which can take some time.
- Evaluating multiple competing models is Multi-Armed Bandits,
 - To define a way to calculate and monitor the reward associated with using each model.

Model Supervisor

No one
Likes dogs

French
Bulldog

French
Bulldog

French
Bulldog



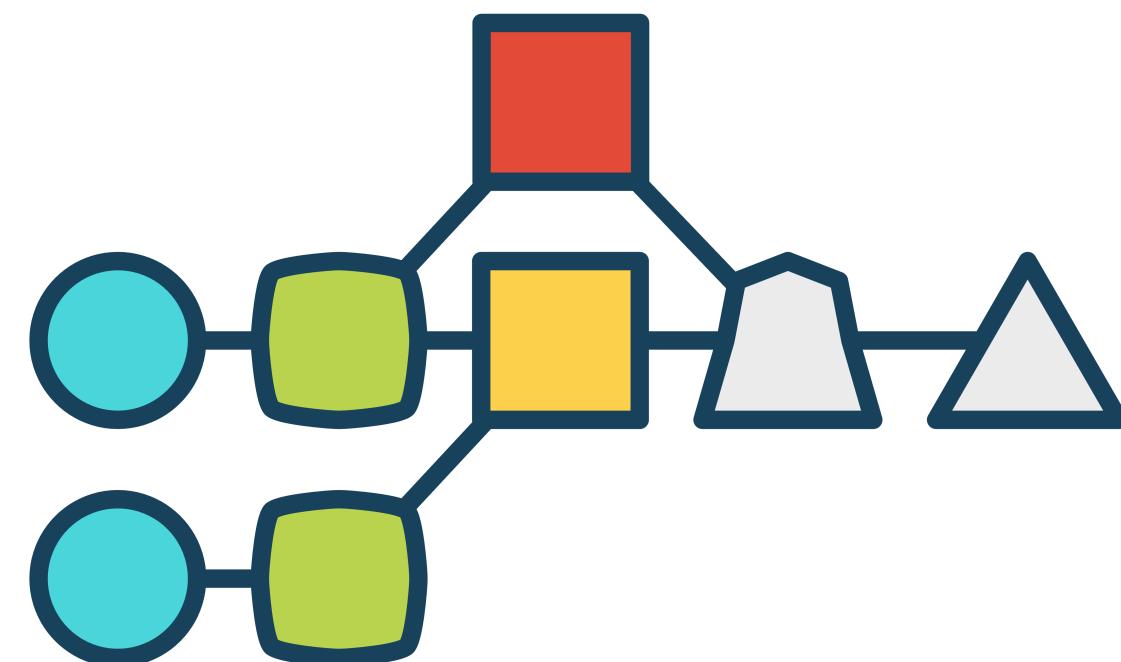
Model Deployment

Online learning models

- To use algorithms and techniques that can continuously improve its performance with the arrival of new data.
- Constantly learning in production.
- Extra complexities, as versioning the model as a static artifact won't yield the same results if it is not fed the same data.
- You will need to version not only the training data, but also the production data that will impact the model's performance.

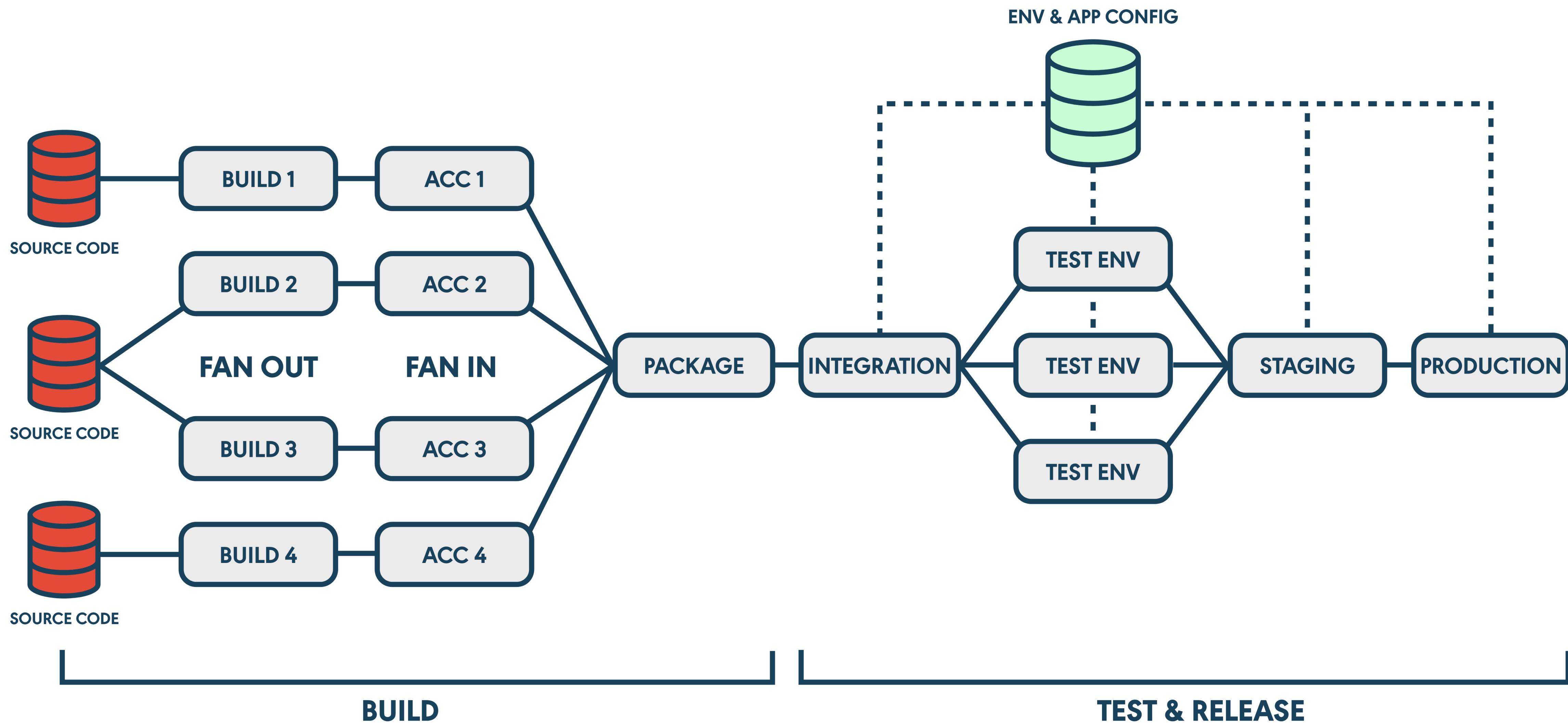
Orchestration in ML Pipelines

- Provisioning of infrastructure and the execution of the ML Pipelines to train and capture metrics from multiple model experiments
- Building, testing, and deploying Data Pipelines
- Testing and validation to decide which models to promote
- Provisioning of infrastructure and deployment of models to production



Continuous Integration and Delivery

GoCD



A Continuous Delivery Scenario for ML

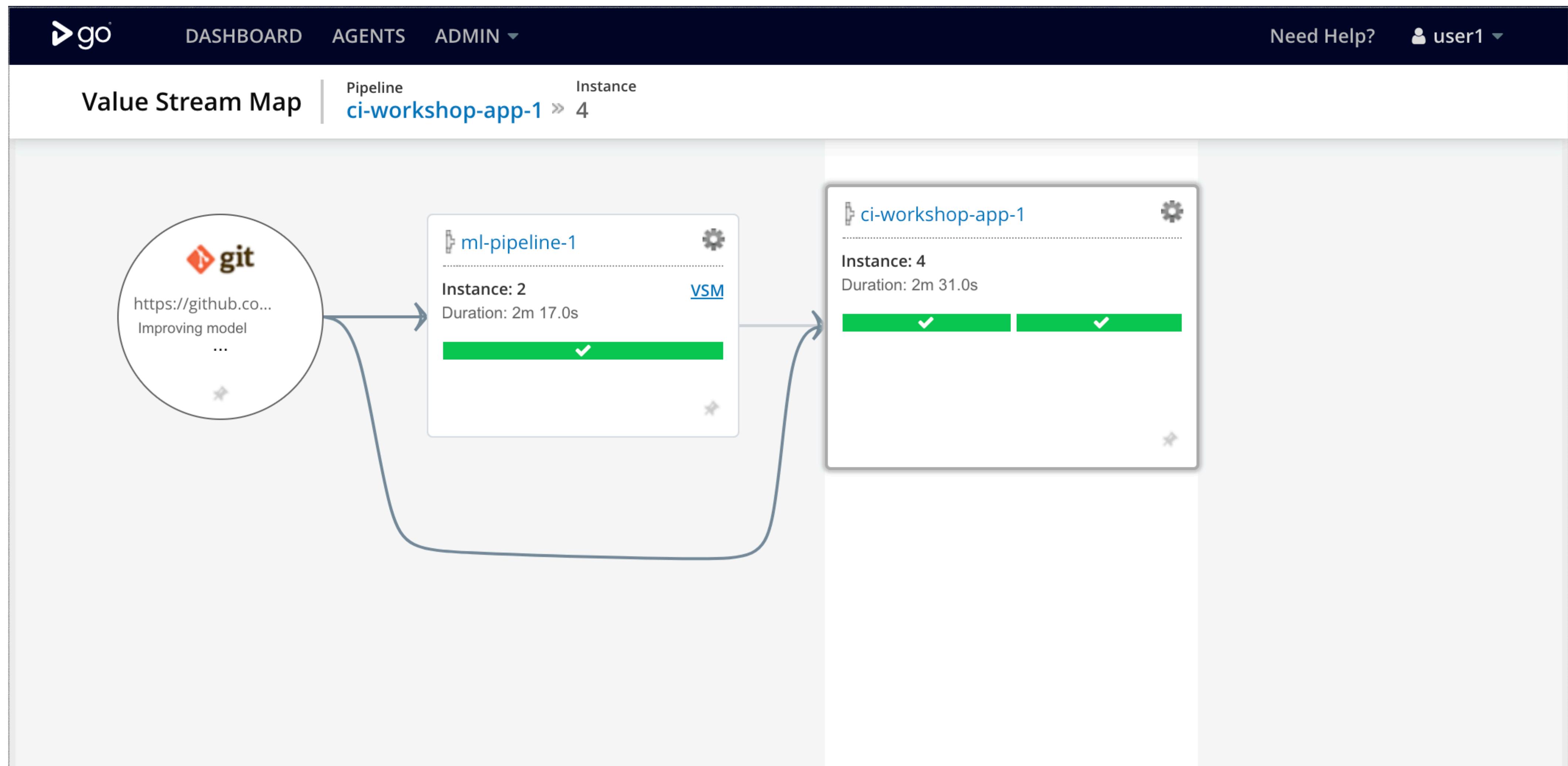
1. Machine Learning Pipeline:

- To train and evaluate ML models
- To execute threshold test to decide if the model can be promoted or not
- *dvc push* to publish it as an artifact

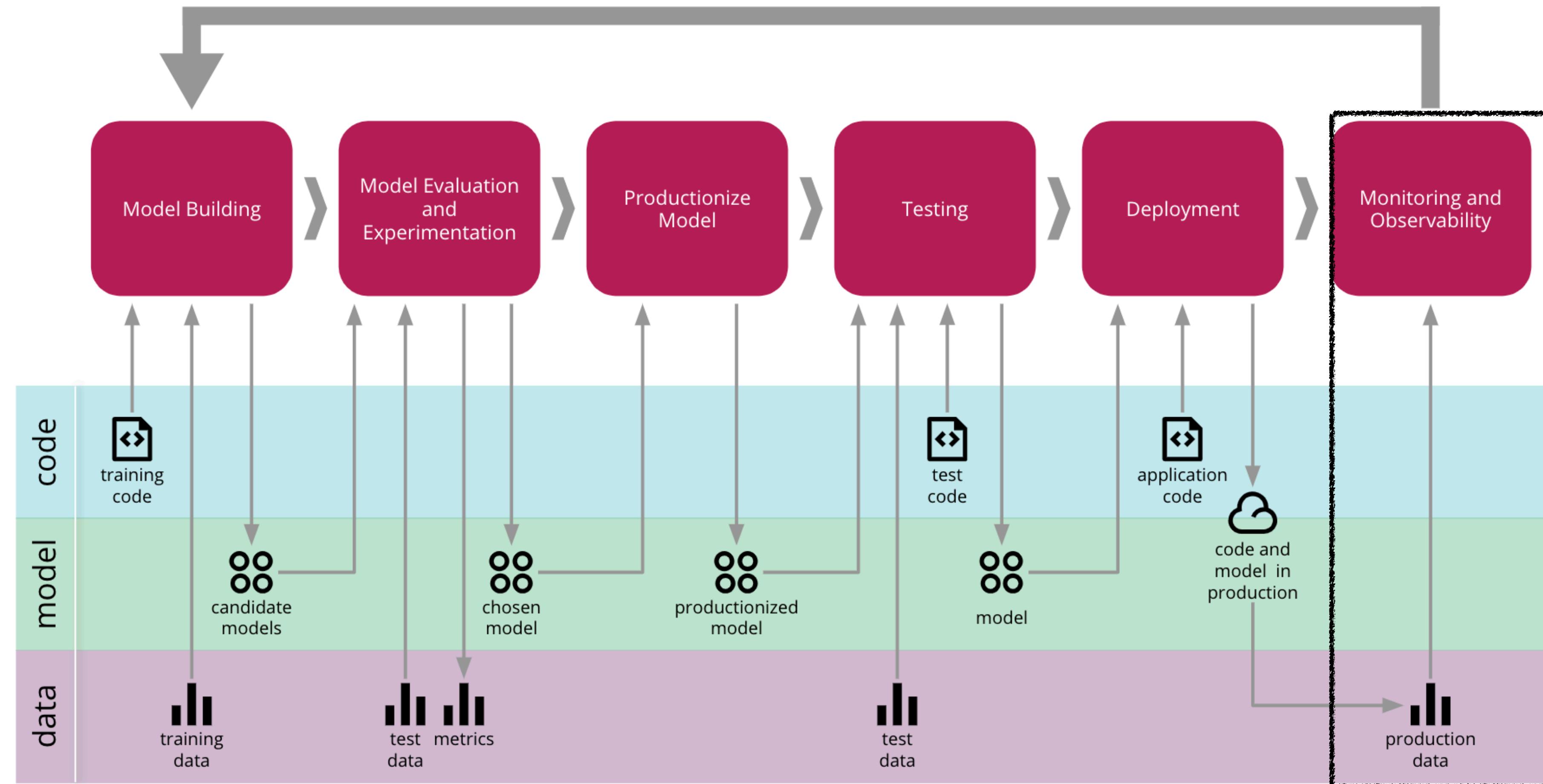
2. Application Deployment Pipeline:

- To build and test the application code
- To fetch the promoted model from the upstream pipeline using *dvc pull*
- To package a new combined artifact that contains the model and the application as a Docker image
- To deploy them to a production cluster

Combining Machine Learning Pipeline and Application Deployment Pipeline



An End-to-End ML System Building Process



ML Model Monitoring

How models perform in production and rollback mechanisms

- **Model inputs:**
 - What data is being fed to the models, identifying training-serving skew.
- **Model outputs:**
 - What predictions and recommendations are the models making from these inputs, to understand how the model is performing with real data.

ML Model Monitoring

How models perform in production and rollback mechanisms

- **Model interpretability outputs:**

- Metrics such as model coefficients, ELI5, or LIME outputs that allow further investigation to understand how the models are making predictions to identify potential overfit or bias that was not found during training.



hi there, i am here looking for some help. my friend is a interie
graphics software on pc. any suggestion on which software to
sophisticated software(the more features it has,the better)

y=0 (probability 0.000) top features			y=1 (probability 0.100) top features			y=2 (probability 0.900) top features		
Contribution?	Feature	Value	Contribution?	Feature	Value	Contribution?	Feature	Value
+0.301	<BIAS>	1.000	+0.427	<BIAS>	1.000	+0.289	hue	0.670
+0.064	color_intensity	8.500	+0.033	proline	630.000	+0.272	<BIAS>	1.000
+0.004	malic_acid	4.600	+0.022	od280/od315_of_diluted_wines	1.920	+0.095	color_intensity	8.500
-0.018	alcalinity_of_ash	25.000	+0.009	alcalinity_of_ash	25.000	+0.083	flavanoids	0.960
-0.044	total_phenols	1.980	+0.006	total_phenols	1.980	+0.067	proline	630.000
-0.055	flavanoids	0.960	-0.003	proanthocyanins	1.110	+0.056	malic_acid	4.600
-0.100	proline	630.000	-0.010	alcohol	13.400	+0.038	total_phenols	1.980
-0.153	hue	0.670	-0.028	flavanoids	0.960	+0.010	alcohol	13.400
			-0.060	malic_acid	4.600	+0.009	alcalinity_of_ash	25.000
			-0.137	hue	0.670	+0.003	proanthocyanins	1.110
			-0.160	color_intensity	8.500	-0.022	od280/od315_of_diluted_wines	1.920

“Why Should I Trust You?”

Explaining the Predictions of Any Classifier

Example #3 of 6 True Class:  Atheism Instructions Previous Next

Algorithm 1

Words that A1 considers important:

GOD	
mean	
anyone	
this	
Koresh	
through	

Predicted:  Atheism

Prediction correct: 

Document

From: pauld@verdix.com (Paul Durbin)
Subject: Re: DAVID CORESH IS! **GOD!**
Nntp-Posting-Host: sarge.hq.verdix.com
Organization: Verdix Corp
Lines: 8

Algorithm 2

Words that A2 considers important:

Posting	
Host	
Re	
by	
in	
Nntp	

Predicted:  Atheism

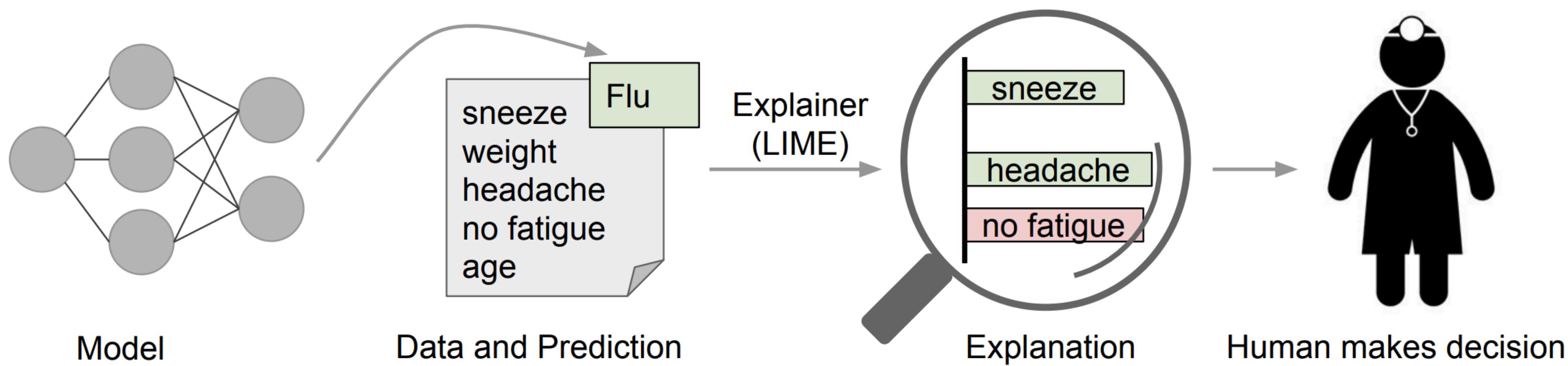
Prediction correct: 

Document

From: pauld@verdix.com (Paul Durbin)
Subject: **Re:** DAVID CORESH IS! GOD!
Nntp-Posting-Host: sarge.hq.verdix.com
Organization: Verdix Corp
Lines: 8

Explaining individual predictions

A model predicts that a patient has the flu, and LIME highlights the symptoms in the patient's history that led to the prediction



ML Model Monitoring

How models perform in production and rollback mechanisms

- **Model outputs and decisions:**

- What predictions our models are making given the production input data, and also which decisions are being made with those predictions.
- Sometimes the application might choose to ignore the model and make a decision based on pre-defined rules (or to avoid future bias).

ML Model Monitoring

How models perform in production and rollback mechanisms

- **User action and rewards:**
 - Based on further user action, we can capture reward metrics to understand if the model is having the desired effect.
 - For example, if we display product recommendations, we can track when the user decides to purchase the recommended product as a reward.

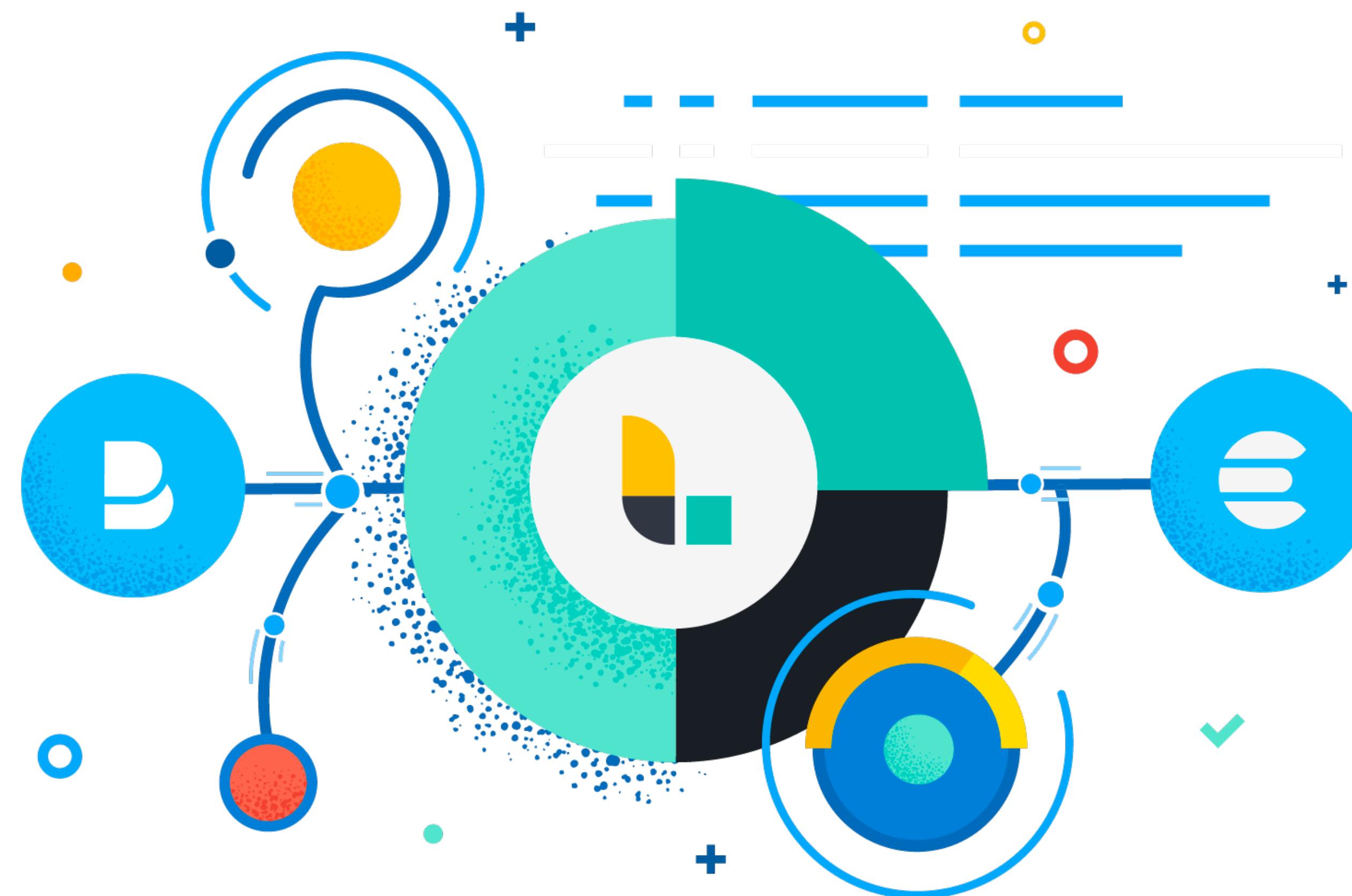
A pipeline for model monitoring

ELK

- **Elasticsearch**: an open source *search* engine.
- **Logstash**: an open source data collector for unified *logging* layer.
- **Kibana**: an open source web UI that makes it easy to explore and *visualize* the data indexed by Elasticsearch.

A pipeline for model monitoring

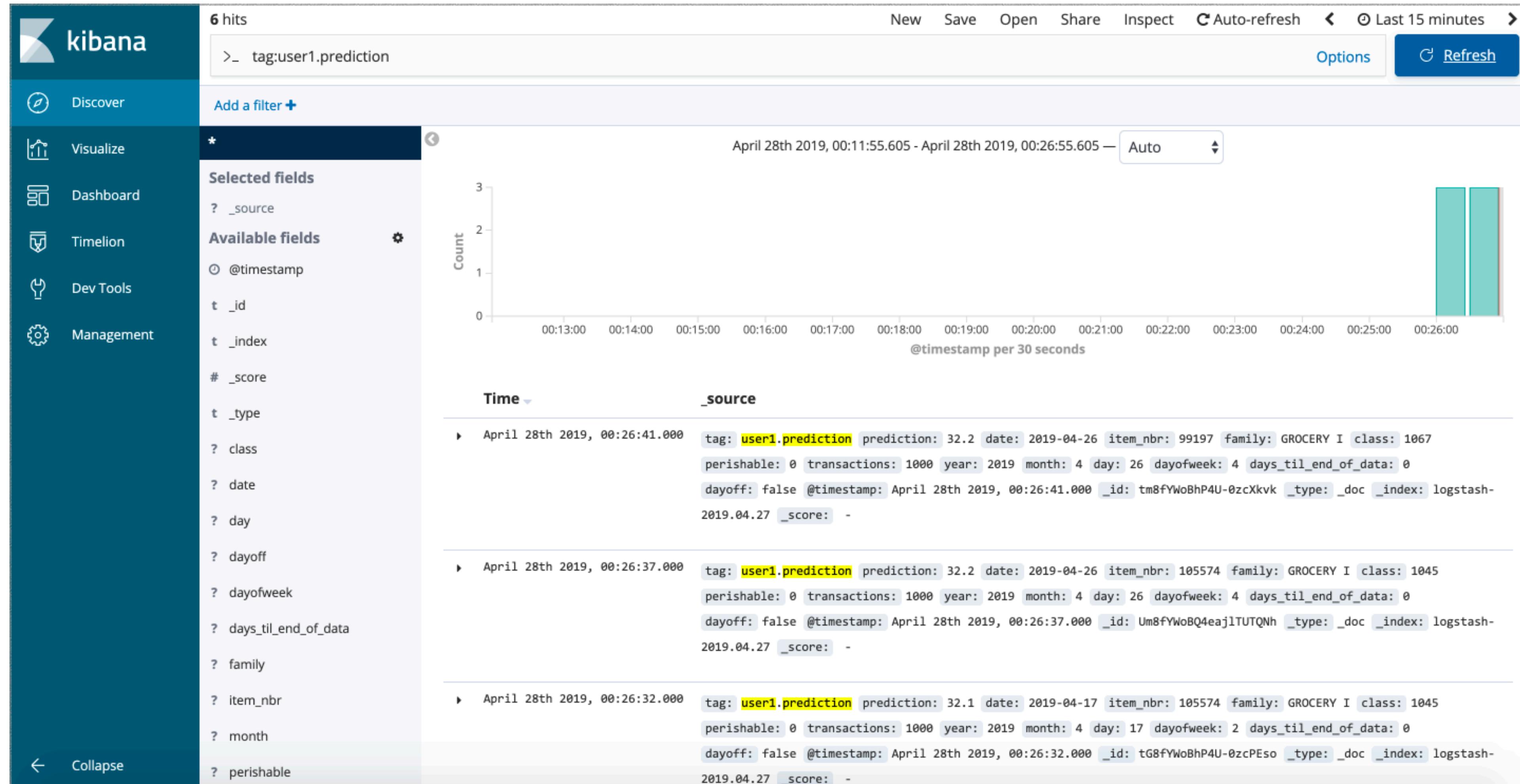
ELK

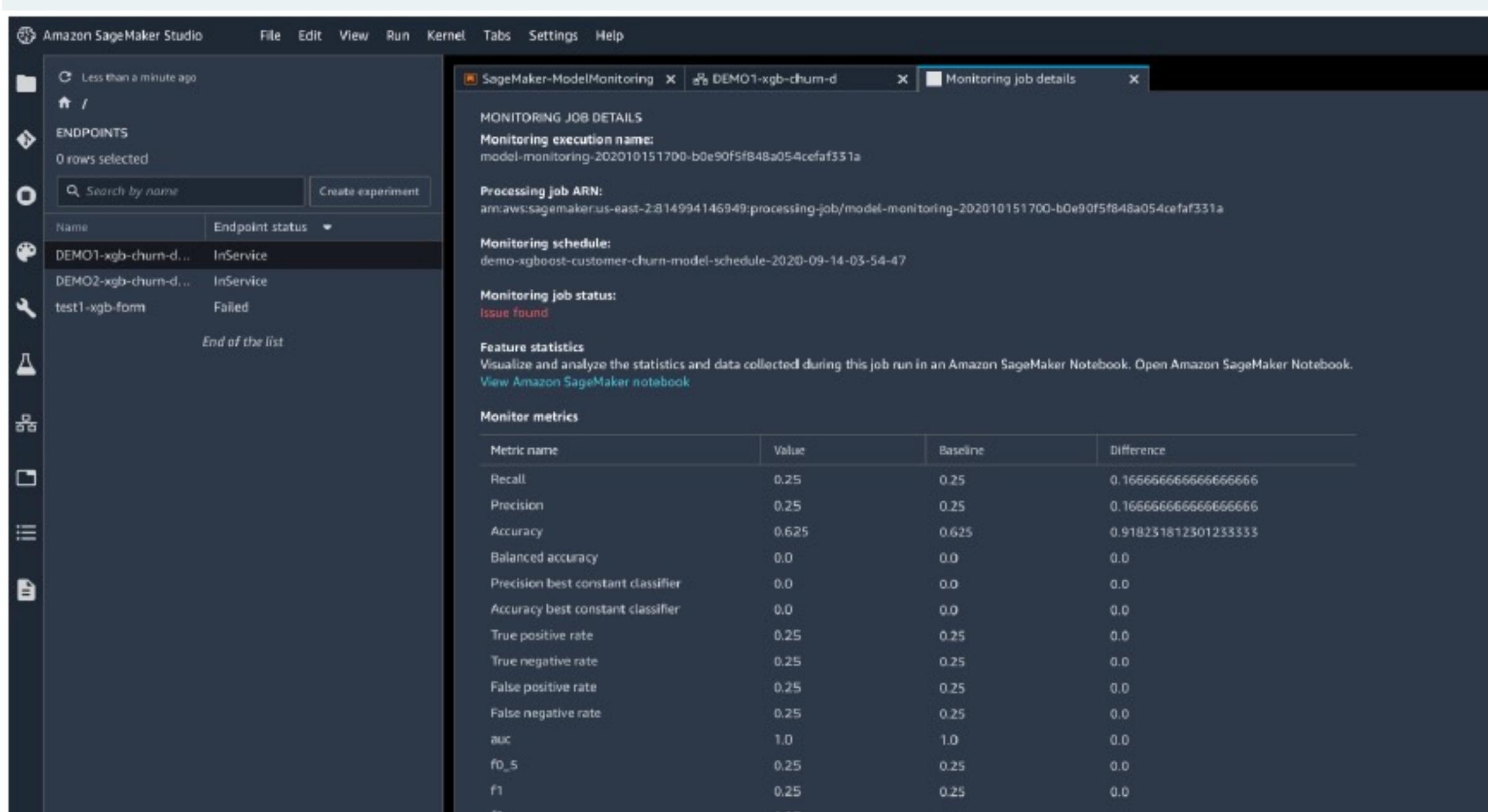
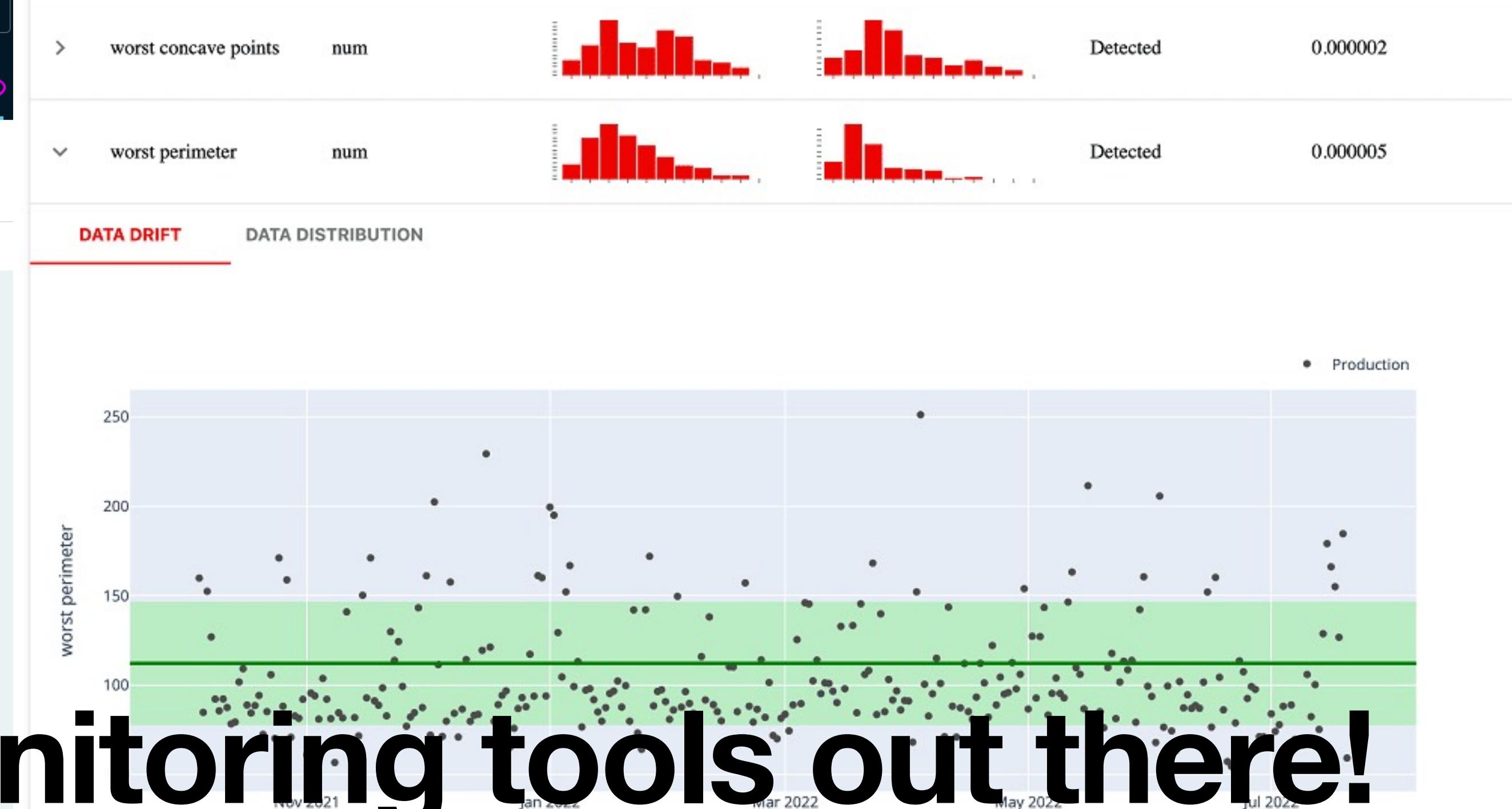
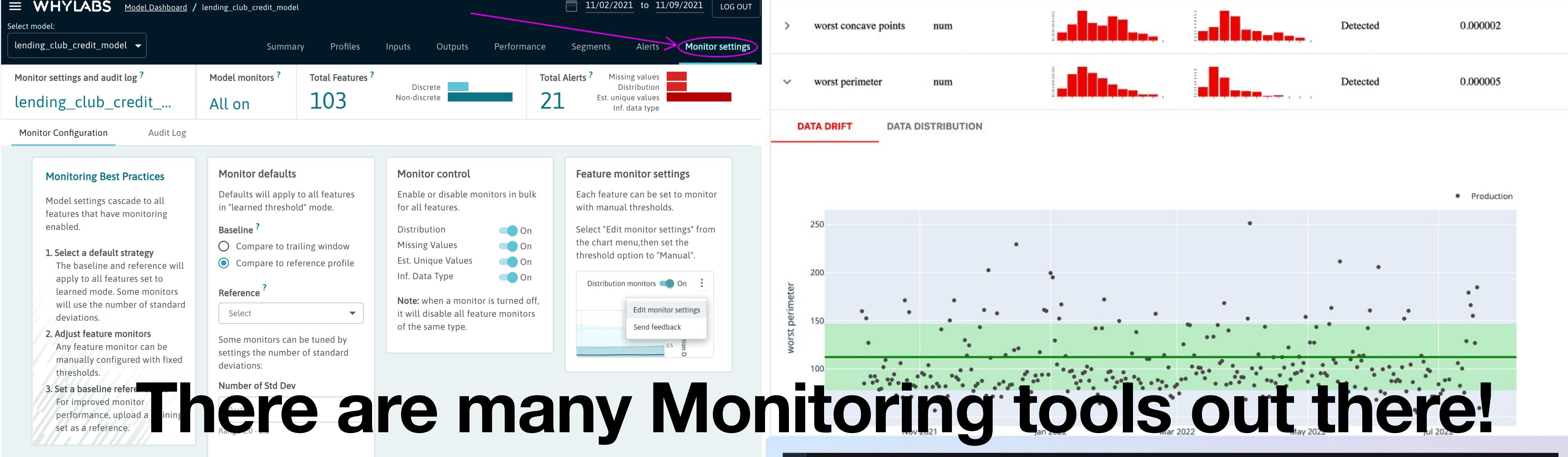


Logging

```
predict_with_logging.py...
df = pd.DataFrame(data=data, index=['row1'])
df = decision_tree.encode_categorical_columns(df)
pred = model.predict(df)
logger = sender.FluentSender(TENANT, host=FLUENTD_HOST, port=int(FLUENTD_PORT))
log_payload = {'prediction': pred[0], **data}
logger.emit('prediction', log_payload)
```

A pipeline for model monitoring





There are many Monitoring tools out there!

Monitoring

Grafana + Prometheus



Prometheus

Logging

- Prometheus is a popular **open-source ML model monitoring** tool that collects multidimensional data and queries.
- The main advantages of Prometheus are tight integration with Kubernetes and many of the available exporters and client libraries, as well as a fast query language. Prometheus is also Docker-compatible and available on the Docker Hub.
- The Prometheus server has its own self-contained unit that does not depend on network storage or external services. So it doesn't require a lot of work to deploy additional infrastructure or software. Its main task is to store and monitor certain objects. An object can be anything: a Linux server, one of the processes, a database server, or any other component of the system. Each element that you want to monitor is called a metric.
- The Prometheus server **reads targets at an interval that you define to collect metrics** and stores them in a time series database. You set the targets and the time interval for reading the metrics. You query the Prometheus time series database for where metrics are stored using the PromQL query language.

Grafana Visualization

- Grafana allows you to visualize monitoring metrics. Grafana specializes in time series analytics. It can visualize the results of monitoring work in the form of line graphs, heat maps, and histograms.
- Instead of writing PromQL queries directly to the Prometheus server, you use Grafana GUI boards to request metrics from the Prometheus server and render them in the Grafana dashboard.

An End-to-End ML System Building Process

