



Machine Learning Design Patterns for ML Ops

Google Cloud

Valliappa (Lak) Lakshmanan
Global Head of Analytics & AI Solutions,
Google Cloud
[Twitter: @lak_gcp](#)
lak@google.com



ginablaber

@ginablaber

Follow



The story of enterprise Machine Learning: “It took me 3 weeks to develop the model. It’s been >11 months, and it’s still not deployed.”

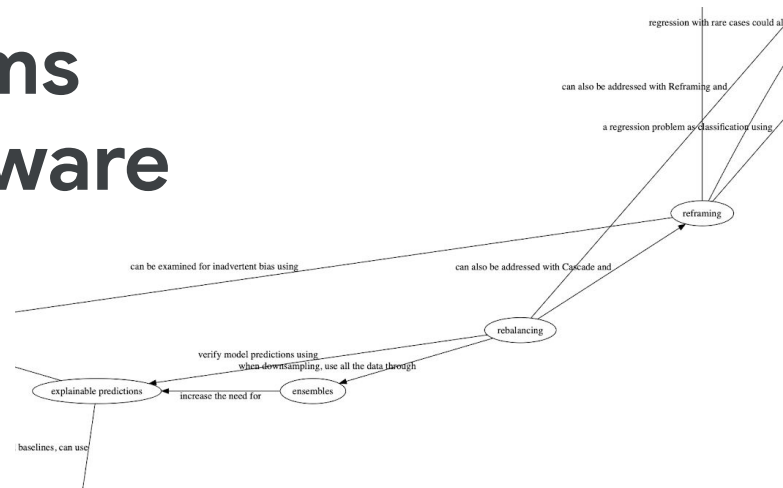
[@DineshNirmalIBM](#) [#StrataData](#) [#strataconf](#)

10:19 AM - 7 Mar 2018



Google Cloud

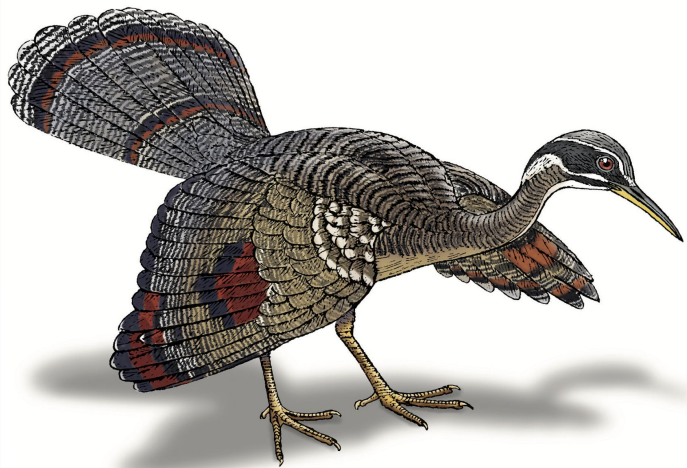
Design patterns are formalized best practices to solve common problems when designing a software system.



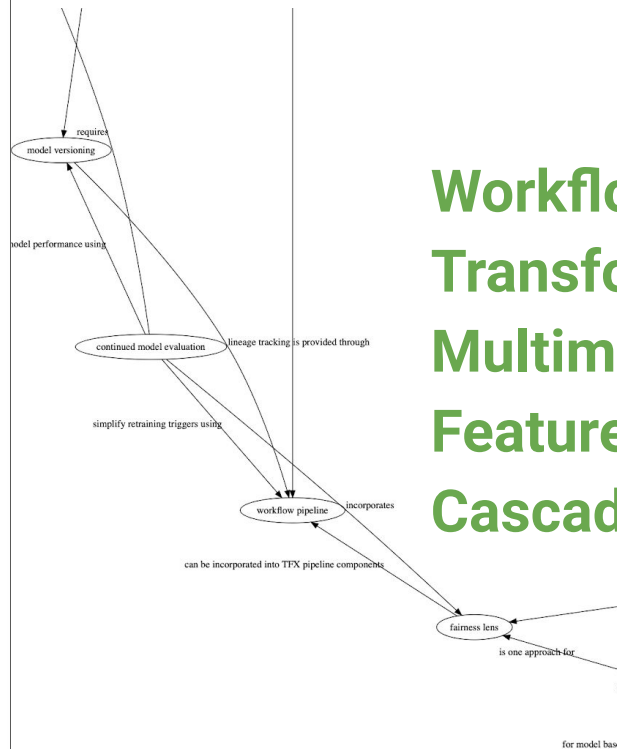
O'REILLY®

Machine Learning Design Patterns

Solutions to Common Challenges in Data Preparation, Model Building, and MLOps



Valliappa Lakshmanan,
Sara Robinson & Michael Munn



Workflow Pipeline
Transform
Multimodal Input
Feature Store
Cascade

<https://github.com/GoogleCloudPlatform/ml-design-patterns>

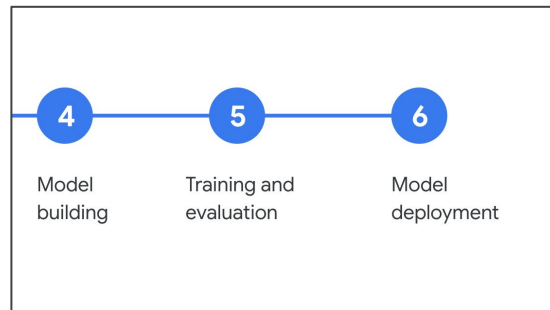
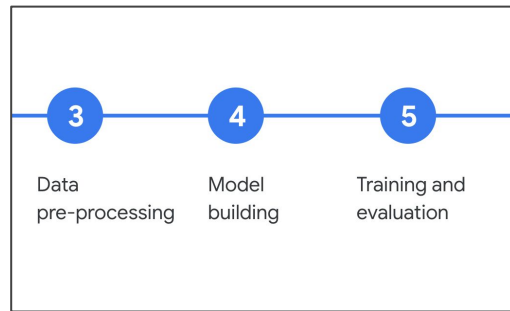
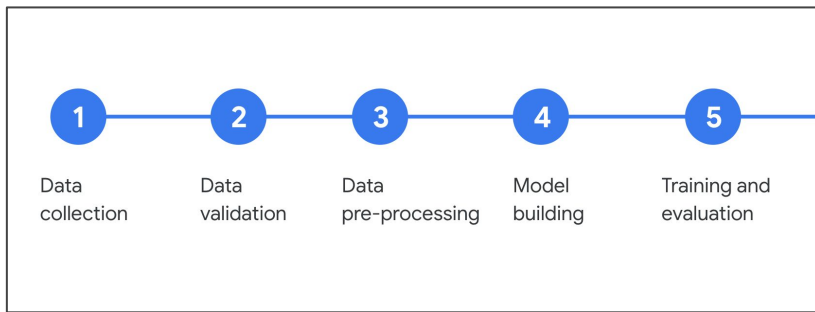
Workflow Pipeline



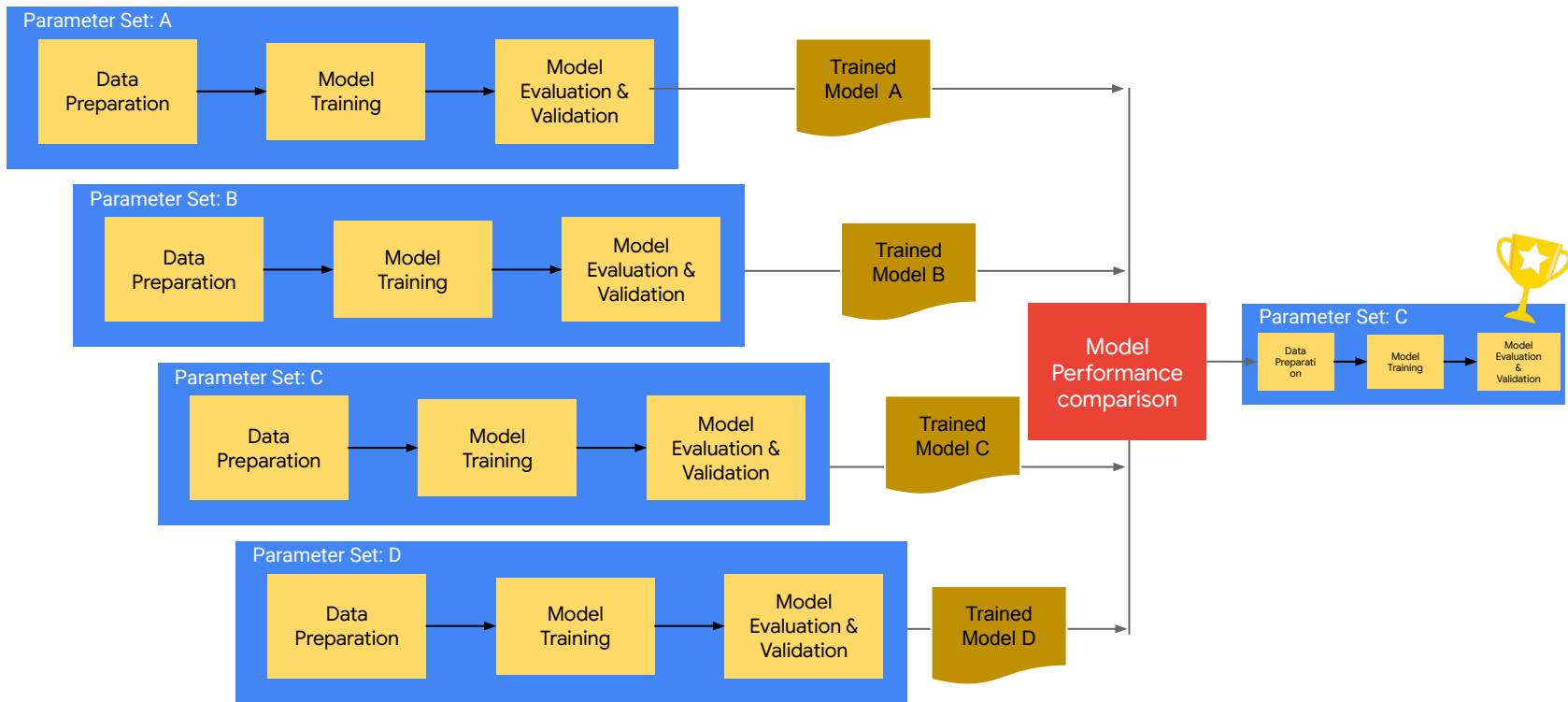
Typical development workflow is monolithic



Experiments are ad-hoc and not repeatable



Hyperparameter tuning has wasteful repetition

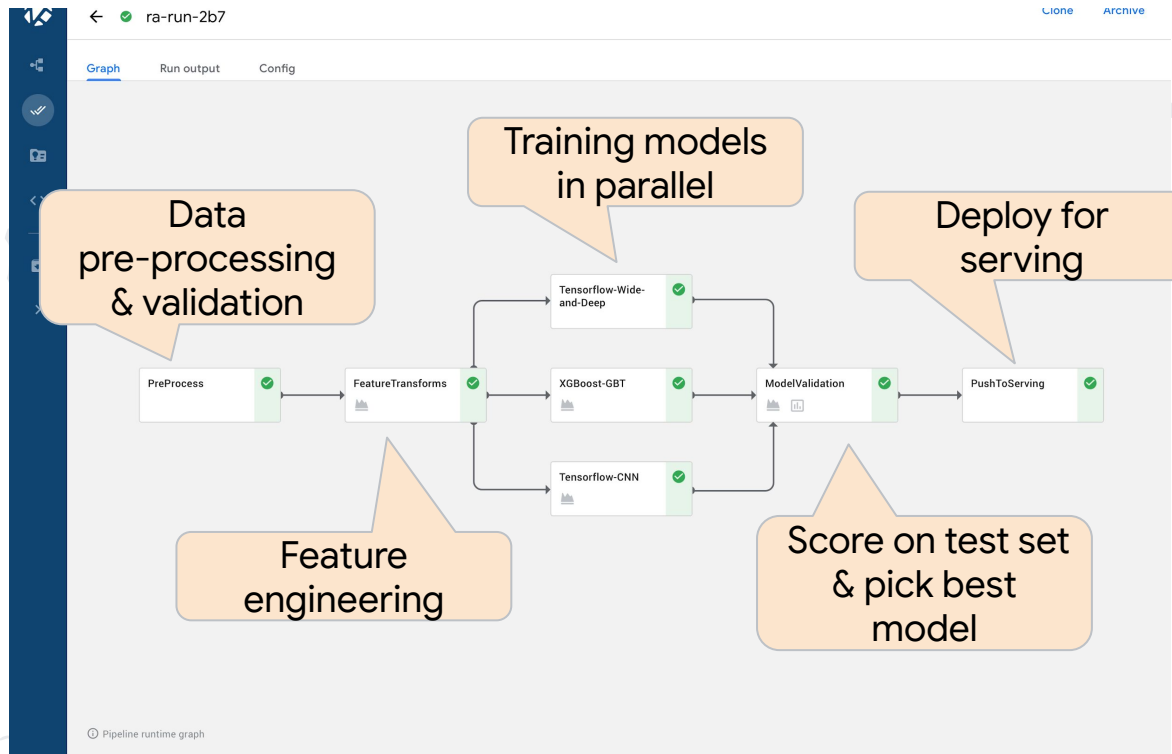


**MLOps has some DevOps concepts, and adds
data validation & continuous evaluation/training**

But MLOps differs from DevOps in important ways

DevOps	MLOps
1 Test and validate code and components	Also test and validate data, data schemas, and models
2 Focus on a single software package or service	Also consider the whole system, the ML training pipeline
3 Deploy code and move to the next task	Constantly monitor, retrain and serve the model

A pipeline is an executable DAG of ML steps



Each step of the pipeline is a container

Pipeline runs can be grouped into Experiments

Logs are associated with each step of each run

An artifact repository stores the parameters & data passed between steps



Experiments > Default

← ✓ train_and_deploy 2020-02-02 20:40-19

Retry Clone run Terminate Archive

Graph Run output Config

preprocess ✓

hypertrain ✓

traintuned ✓

deploycmle ✓

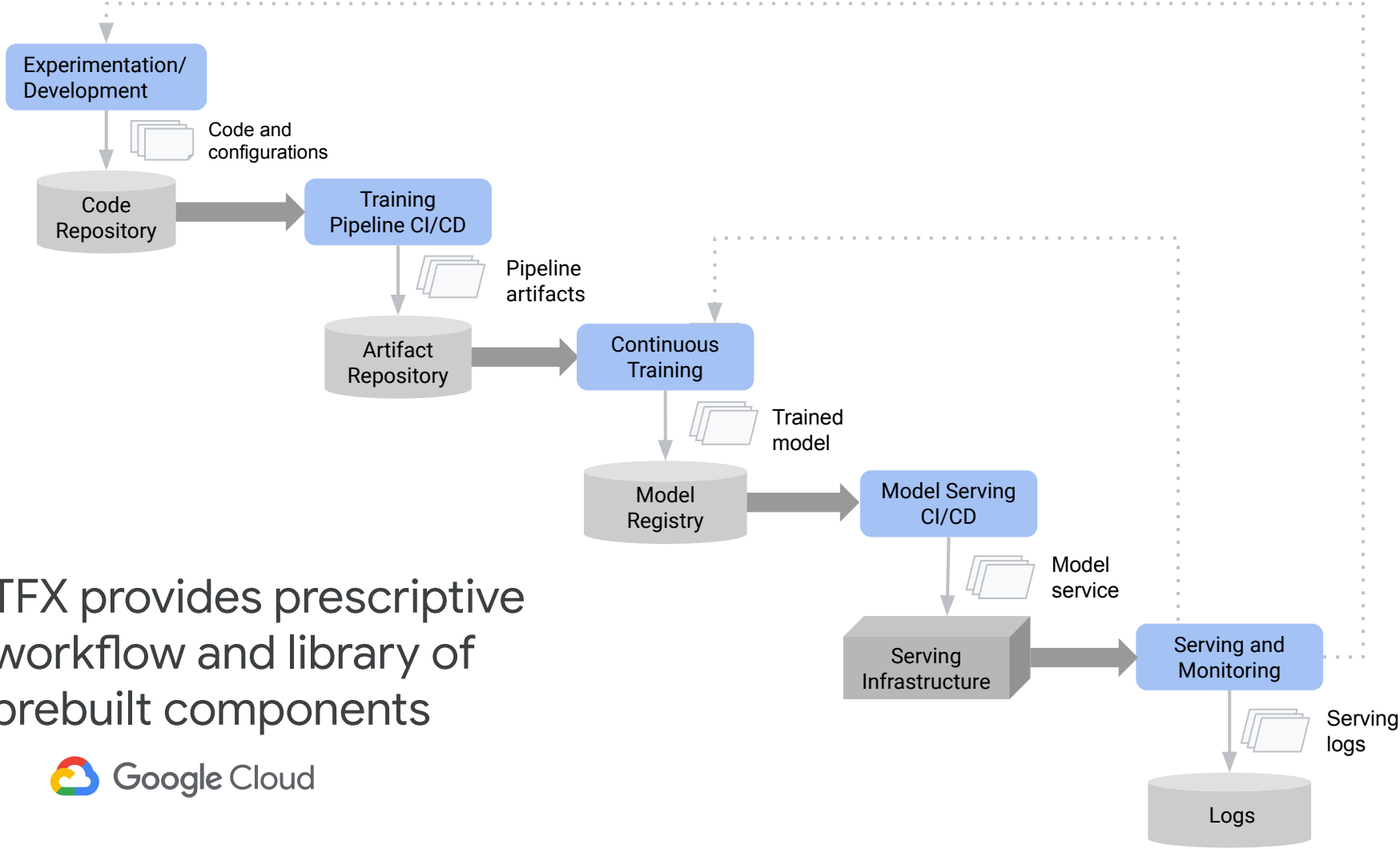
deployapp ✓

Log viewer for babyweight-74ggh-3448946795

	Artifacts	Input/Output	Volumes	Manifest	Logs
19092	INFO	2020-02-02 22:37:44	+0000	master-replica-0	14
19093	INFO	2020-02-02 22:37:44	+0000	master-replica-0	14
19094	INFO	2020-02-02 22:37:44	+0000	master-replica-0	14
19095	INFO	2020-02-02 22:37:44	+0000	worker-replica-0	14
19096	INFO	2020-02-02 22:37:44	+0000	master-replica-0	14
19097	INFO	2020-02-02 22:37:44	+0000	worker-replica-3	14
19098	INFO	2020-02-02 22:37:44	+0000	worker-replica-2	14
19099	INFO	2020-02-02 22:37:44	+0000	master-replica-0	14
19100	INFO	2020-02-02 22:37:45	+0000	master-replica-0	14
19101	INFO	2020-02-02 22:37:45	+0000	worker-replica-1	14
19102	INFO	2020-02-02 22:37:45	+0000	master-replica-0	14
19103	INFO	2020-02-02 22:37:45	+0000	master-replica-0	14
19104	INFO	2020-02-02 22:37:45	+0000	master-replica-0	14
19105	INFO	2020-02-02 22:37:45	+0000	worker-replica-0	14
19106	INFO	2020-02-02 22:37:45	+0000	worker-replica-3	14
19107	INFO	2020-02-02 22:37:45	+0000	master-replica-0	14
19108	INFO	2020-02-02 22:37:45	+0000	worker-replica-1	14
19109	INFO	2020-02-02 22:37:45	+0000	worker-replica-2	14
19110	INFO	2020-02-02 22:37:45	+0000	master-replica-0	14
19111	INFO	2020-02-02 22:37:45	+0000	master-replica-0	14
19112	INFO	2020-02-02 22:37:45	+0000	master-replica-0	14
19113	INFO	2020-02-02 22:37:45	+0000	master-replica-0	14
19114	INFO	2020-02-02 22:37:45	+0000	worker-replica-0	14
19115	INFO	2020-02-02 22:37:45	+0000	master-replica-0	14
19116	INFO	2020-02-02 22:37:46	+0000	worker-replica-3	14
19117	INFO	2020-02-02 22:37:46	+0000	worker-replica-1	14
19118	INFO	2020-02-02 22:37:46	+0000	master-replica-0	14
19119	INFO	2020-02-02 22:37:46	+0000	worker-replica-2	14
19120	INFO	2020-02-02 22:37:46	+0000	master-replica-0	14
19121	INFO	2020-02-02 22:37:46	+0000	master-replica-0	14
19122	INFO	2020-02-02 22:37:46	+0000	master-replica-0	14
19123	INFO	2020-02-02 22:37:46	+0000	master-replica-0	14
19124	INFO	2020-02-02 22:37:46	+0000	worker-replica-0	14
19125	INFO	2020-02-02 22:37:46	+0000	worker-replica-3	14
19126	INFO	2020-02-02 22:37:46	+0000	master-replica-0	14

Log viewer content (partial):

```
global_step/sec: 789.983
global_step/sec: 868.457
global_step/sec: 861.257
loss = 11.854738, step = 2495198 (0.577)
global_step/sec: 844.096
loss = 3.0272508, step = 2495236 (0.589)
loss = 14.983984, step = 2495331 (0.642)
global_step/sec: 868.489
loss = 12.244221, step = 2495351 (0.627)
loss = 5.5081087, step = 2495366 (0.532)
global_step/sec: 857.556
global_step/sec: 857.131
global_step/sec: 837.97
loss = 3.9065936, step = 2495691 (0.582)
loss = 5.2942133, step = 2495718 (0.561)
global_step/sec: 862.088
loss = 6.4086, step = 2495840 (0.552)
loss = 6.0803595, step = 2495851 (0.607)
global_step/sec: 883.29
loss = 24.808087, step = 2495886 (0.628)
global_step/sec: 878.07
global_step/sec: 801.749
loss = 4.4279103, step = 2496170 (0.563)
global_step/sec: 863.605
loss = 4.549688, step = 2496206 (0.571)
loss = 10.025039, step = 2496297 (0.538)
global_step/sec: 848.600
loss = 4.4648595, step = 2496378 (0.624)
global_step/sec: 812.315
loss = 10.613207, step = 2496444 (0.668)
global_step/sec: 803.133
global_step/sec: 818.936
loss = 13.342169, step = 2496653 (0.582)
loss = 2.6249943, step = 2496697 (0.593)
global_step/sec: 858.356
```



TFX provides prescriptive workflow and library of prebuilt components

Workflow Pipeline

Capture ML workflows in a DAG

The pipeline DAG:

- Executable as a whole or in parts
- Can be triggered by events
- Logging and monitoring for each step, run, experiment
- Artifacts stored in repository

Prebuilt components, which:

- Understand training vs. inference
- Avoid rerunning if output artifact already up-to-date in repository
- Run as containers on the pipelines platform

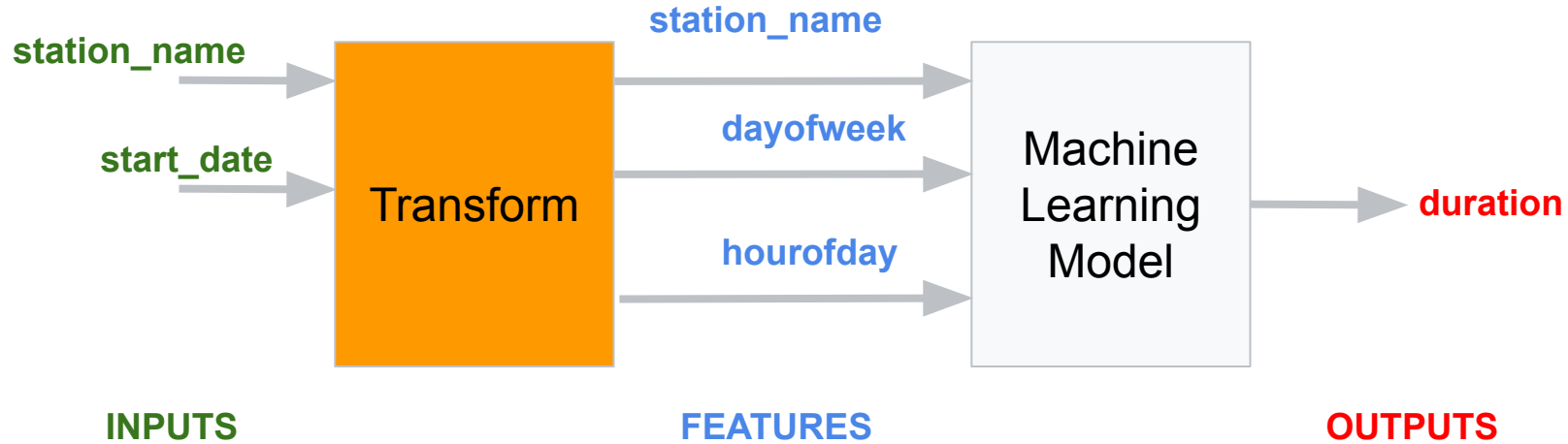
```
pipeline.Pipeline(  
    pipeline_name='hurricane_prediction',  
    pipeline_root='path/to/pipeline/code',  
    components=[  
        bigquery_gen, statistics_gen, schema_gen, train, model_pusher  
    ]  
)
```

```
bigquery_gen = BigQueryExampleGen(query=query)  
  
statistics_gen = StatisticsGen(examples=example_gen.outputs['examples'])  
  
...  
  
model_pusher = Pusher(  
    model=trainer.outputs['model'],  
    model_blessing=evaluator.outputs['blessing'],  
    push_destination=pusher_pb2.PushDestination(  
        filesystem=pusher_pb2.PushDestination.Filesystem(  
            base_directory=serving_model_dir)))
```

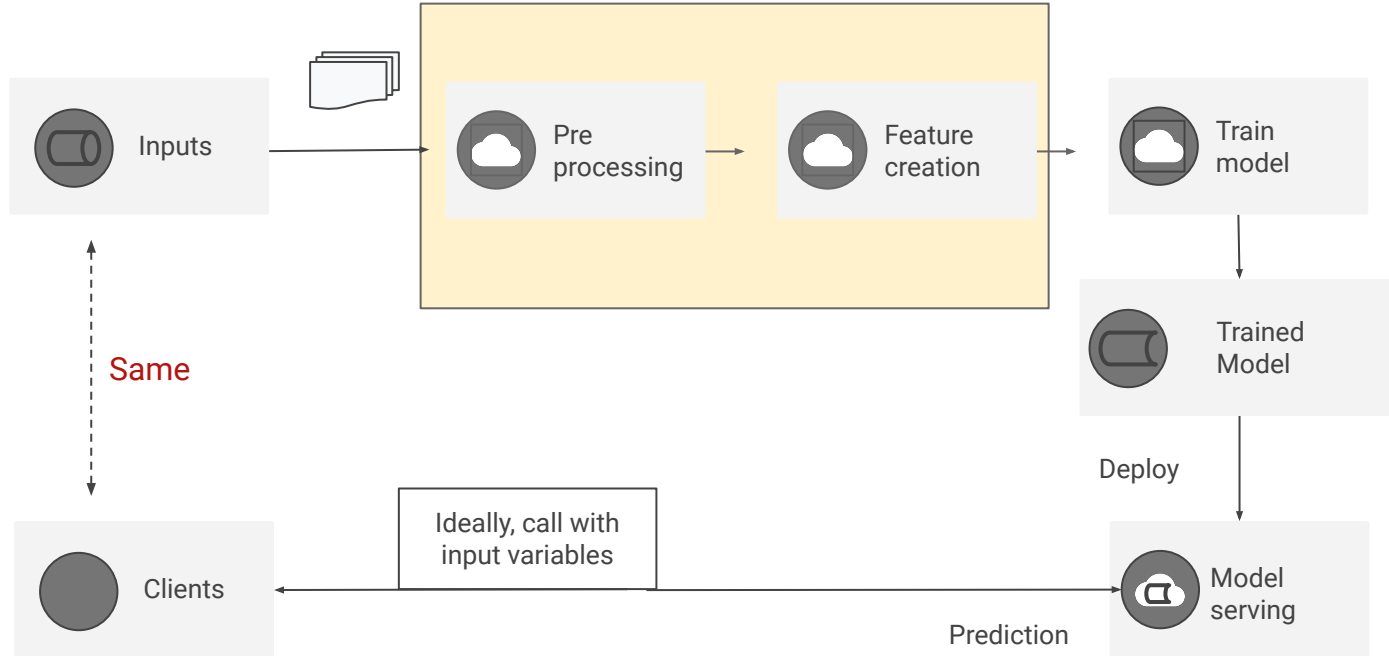
Transform



Imagine a model to predict the length of rentals



Who does transformation during prediction?



Ideally, client code does not have to know about all the transformations that were carried out

```
CREATE OR REPLACE MODEL ch09edu.bicycle_model
OPTIONS(input_label_cols=['duration'],
        model_type='linear_reg')
AS

SELECT
  duration
  , start_station_name
  , CAST(EXTRACT(dayofweek from start_date) AS STRING)
    as dayofweek
  , CAST(EXTRACT(hour from start_date) AS STRING)
    as hourofday
FROM
  `bigquery-public-data.london_bicycles.cycle_hire`
```



Leading cause of
training-serving skew

```
SELECT * FROM ML.PREDICT(MODEL ch09edu.bicycle_model,(
  350 AS duration
  , 'Kings Cross' AS start_station_name
  , '3' as dayofweek
  , '18' as hourofday
))
```

TRANSFORM ensures transformations are automatically applied during ML.PREDICT

```
CREATE OR REPLACE MODEL ch09edu.bicycle_model
OPTIONS(input_label_cols=['duration'],
        model_type='linear_reg')
AS
SELECT
  duration
  , start_station_name
  , CAST(EXTRACT(dayofweek from start_date) AS STRING)
    as dayofweek
  , CAST(EXTRACT(hour from start_date) AS STRING)
    as hourofday
FROM
  `bigquery-public-data.london_bicycles.cycle_hire`
```

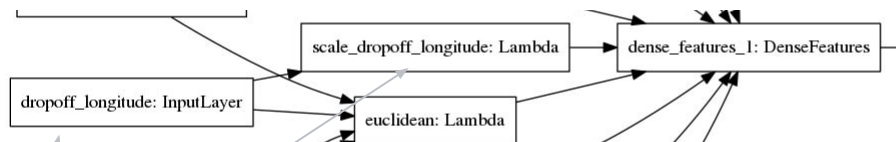


```
SELECT * FROM ML.PREDICT(MODEL ch09edu.bicycle_model,(
  350 AS duration
  , 'Kings Cross' AS start_station_name
  , '3' as dayofweek
  , '18' as hourofday
))
```

```
CREATE OR REPLACE MODEL ch09edu.bicycle_model
OPTIONS(input_label_cols=['duration'],
        model_type='linear_reg')
TRANSFORM(
  SELECT * EXCEPT(start_date)
    , CAST(EXTRACT(dayofweek from start_date) AS STRING)
      as dayofweek
    , CAST(EXTRACT(hour from start_date) AS STRING)
      as hourofday
)
AS
SELECT
  duration, start_station_name, start_date
FROM
  `bigquery-public-data.london_bicycles.cycle_hire`
```

```
SELECT * FROM ML.PREDICT(MODEL ch09edu.bicycle_model,(
  350 AS duration
  , 'Kings Cross' AS start_station_name
  , CURRENT_TIMESTAMP() as start_date
))
```

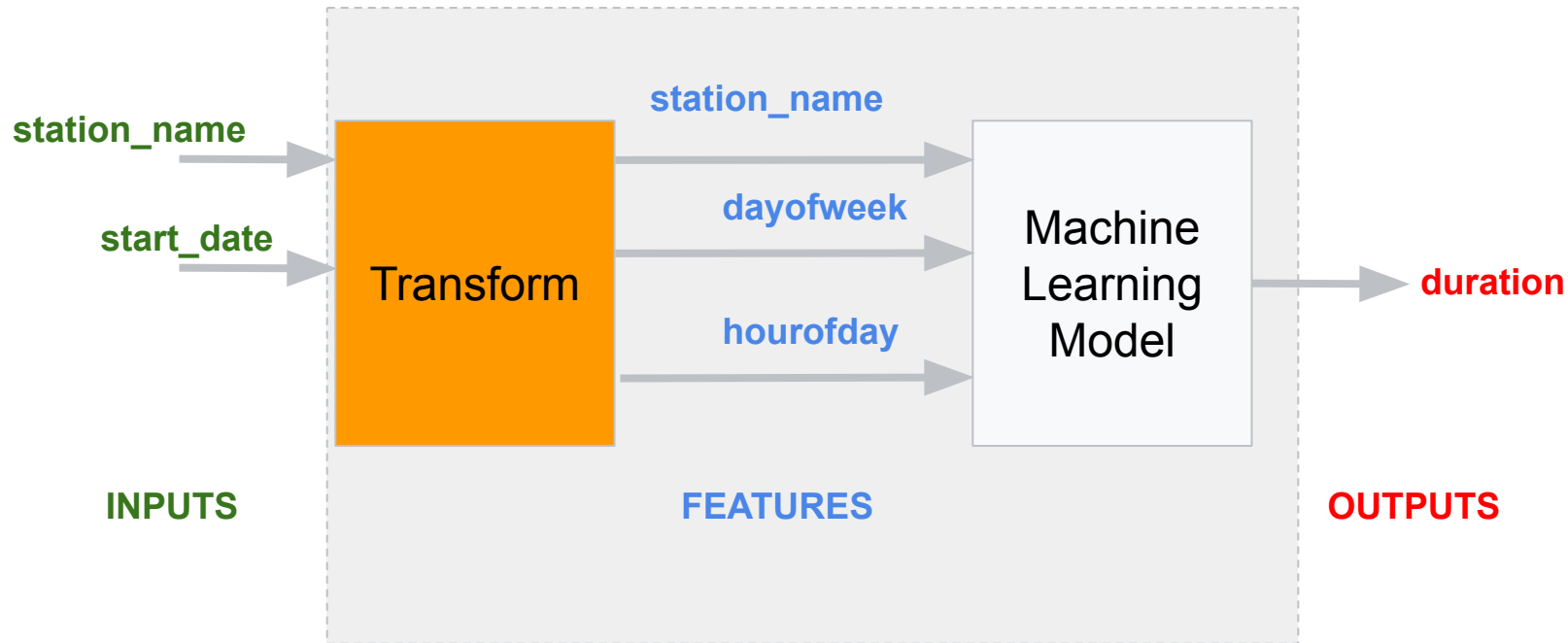
In TensorFlow/Keras, do transformations in Lambda Layers so that they are part of the model graph



```
for lon_col in ['pickup_longitude', 'dropoff_longitude']: # in range -70 to -78
    transformed[lon_col] = tf.keras.layers.Lambda(
        lambda x: (x+78)/8.0,
        name='scale_{}'.format(lon_col)
    )(inputs[lon_col])
```

Moving an ML model to production is much easier if you keep inputs, features, and transforms separate

Transform pattern: the model graph should include the transformations

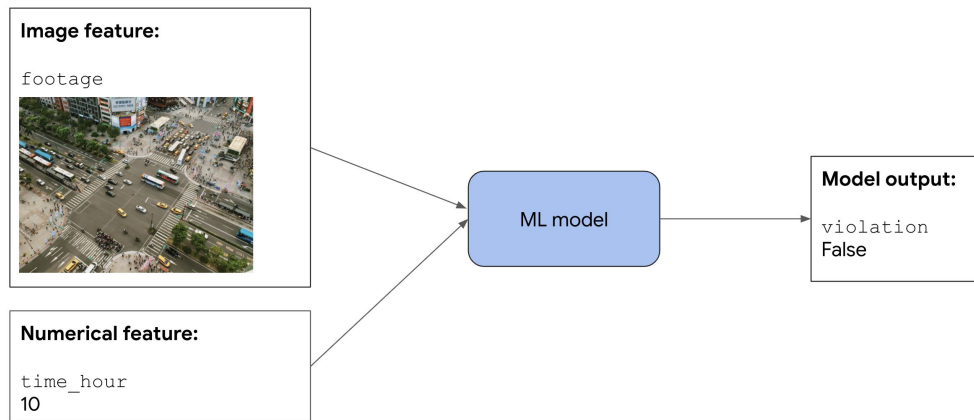


Multimodal input



Is this an image classification problem?

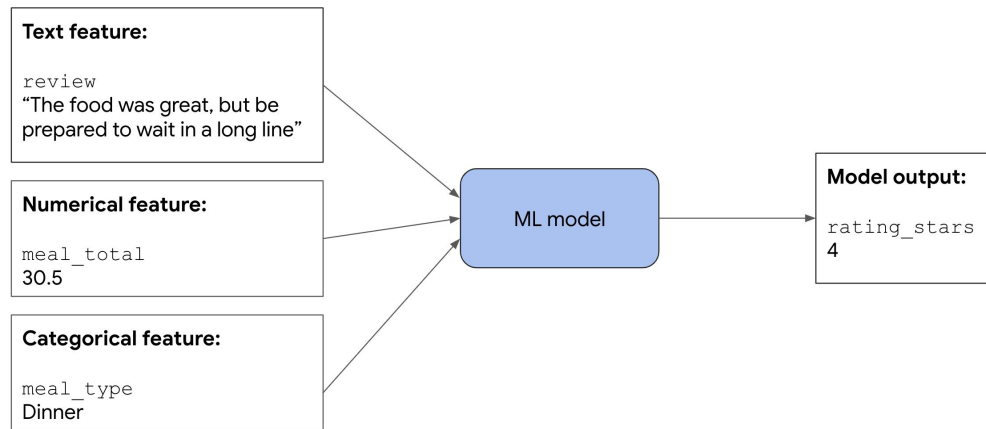
The input is multimodal: an image and some structured data



Much clearer in the case of free form text in tabular data

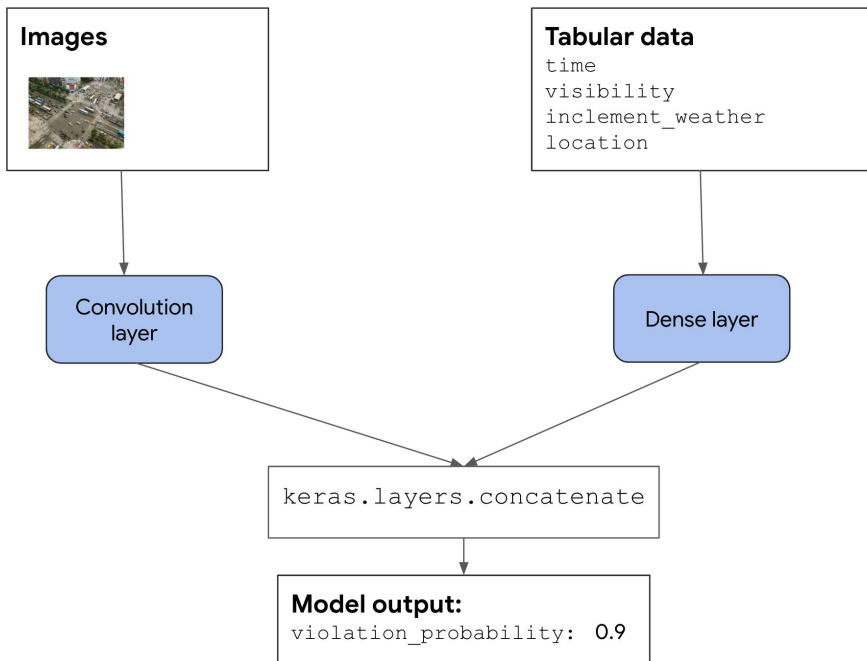
There are multiple ways of representing the text feature (length, sentiment, language, etc.)

Plus, most models will need a mix of structured and unstructured data



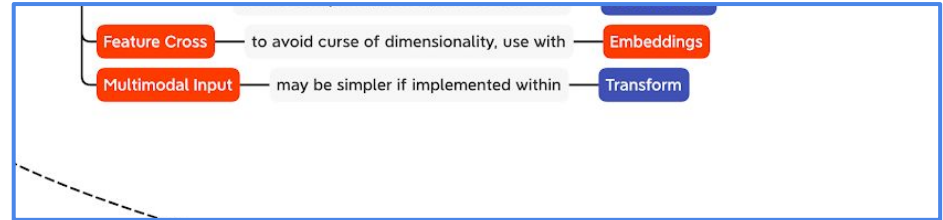
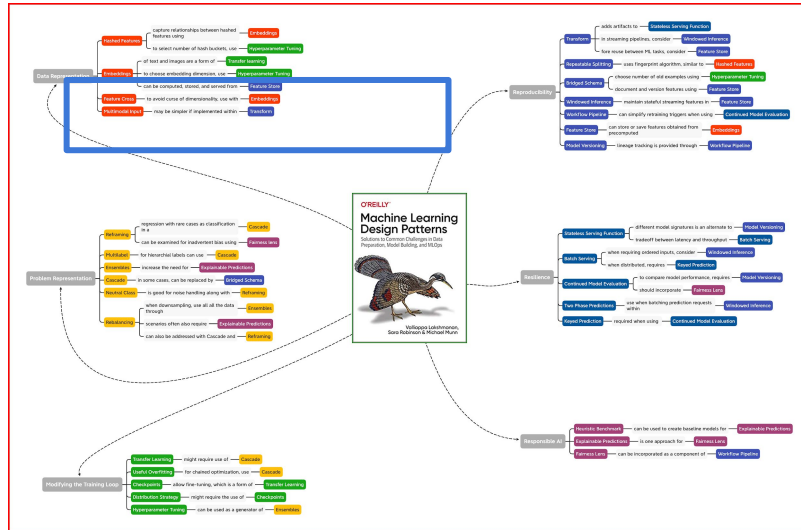
Concatenate the multimodal representations

Concatenate the inputs as a layer in the model, so that it is part of the model graph (like the Transform pattern)



Chapter	Design pattern	Problem solved	Solution
Data Representation	Multimodal Input	Choose between several potential data representations	Concatenate all the available data representations
Problem Representation	Cascade	Simpler if used within	
Reproducibility	Transform	The inputs to a model must be transformed to create the features the model expects and that process must be consistent between training and serving	Explicitly capture and store the transformations applied to convert the model inputs into features
	Workflow Pipeline	When scaling the ML workflow, you need a way to run trials independently and track performance for each step of the pipeline.	Is part of Make each step of the ML workflow a separate, containerized service which can be chained together to make a pipeline that can be run with a single REST API call
	Feature Store		

Connections between patterns

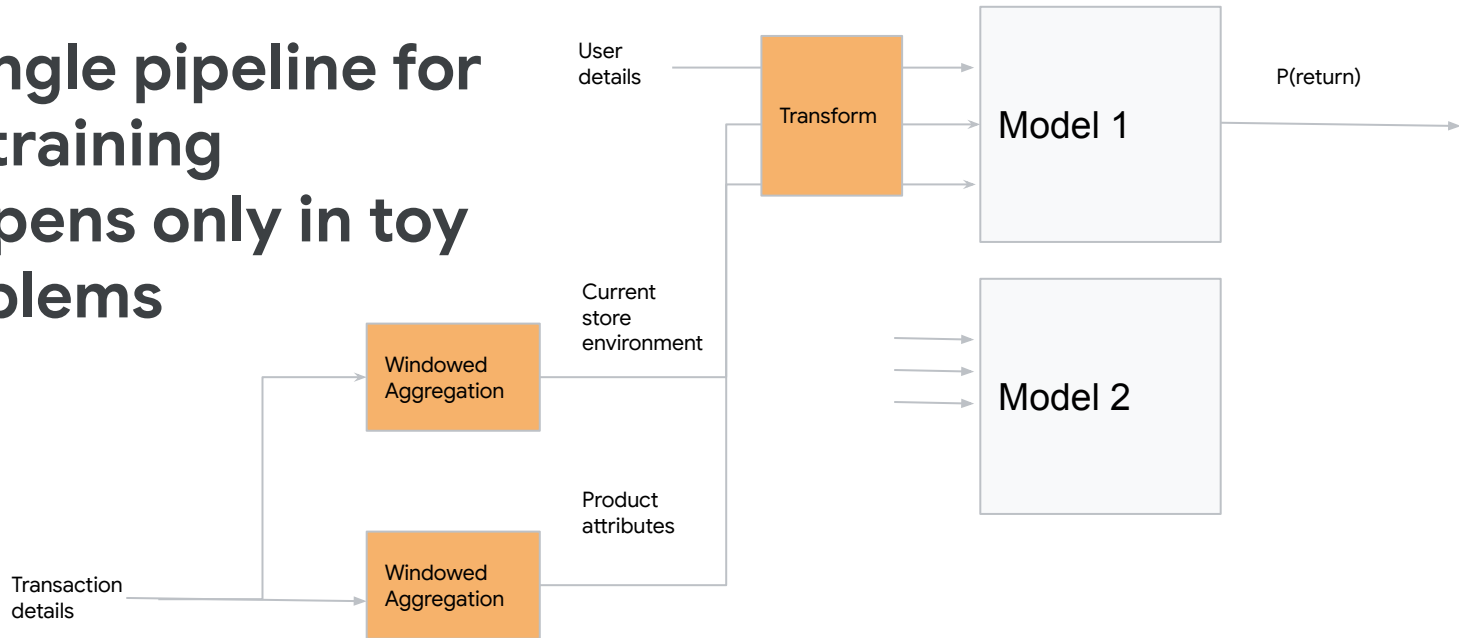


https://github.com/GoogleCloudPlatform/ml-design-patterns/blob/master/08_connected_patterns/machine-learning-design-patterns.png

Feature Store

D4

A single pipeline for ML training happens only in toy problems



Features get reused
Aggregations get reused
Transforms get reused
Some Transforms are real-time

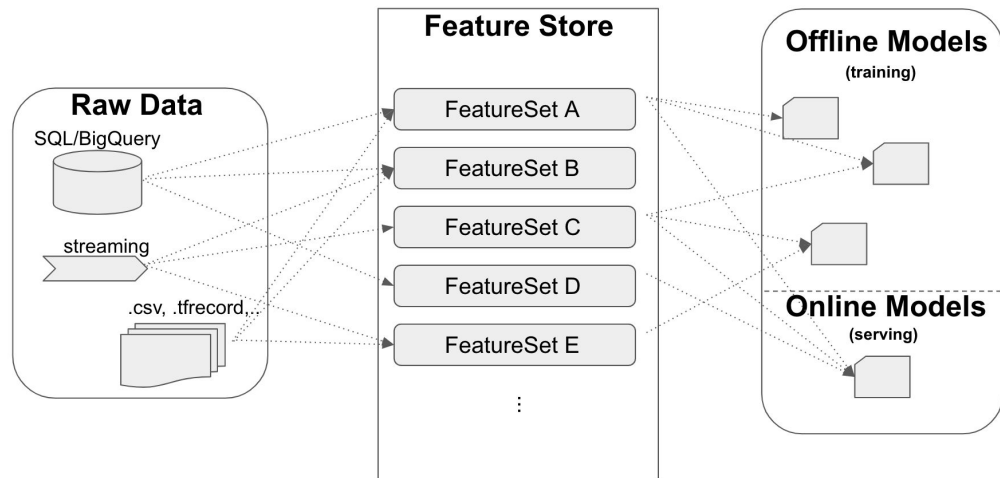
Feature Store:

Centralized location to store feature sets (data storage)

Lambda architecture:

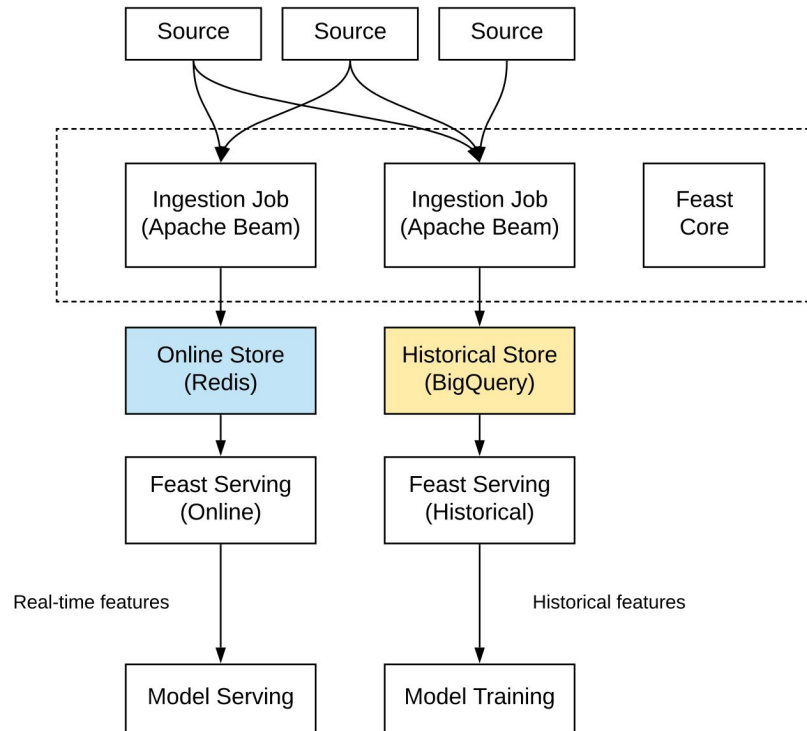
- Store/serve features for low-latency inference
- Store/serve features for large batch access in training
- Metadata layer for versioning of different feature sets

API to manage loading and retrieving feature data.



Feast

Developed by Google and GoJek

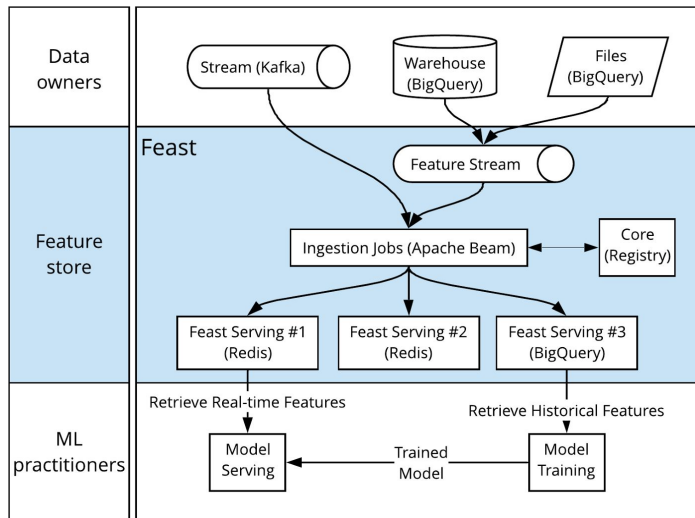
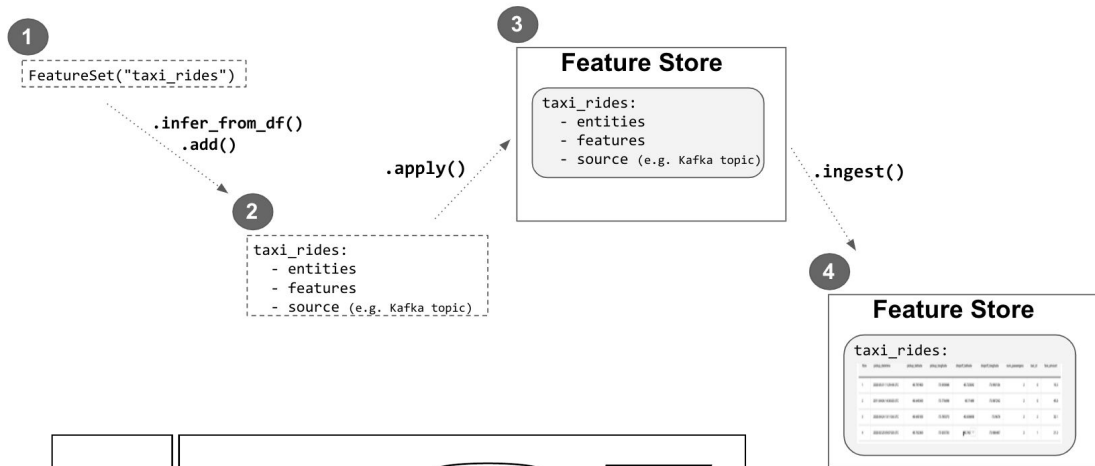


Ingesting into Feast

Four steps:

1. create a FeatureSet
2. add entities and features
3. register the FeatureSet (returns JSON schema)
4. ingest feature data into the FeatureSet.

Feast, like TFX, uses Beam for feature creation

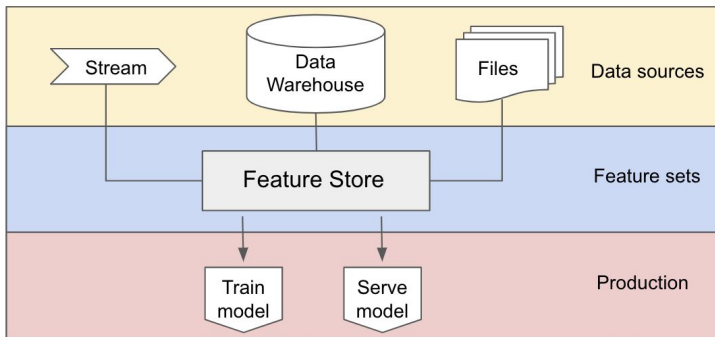
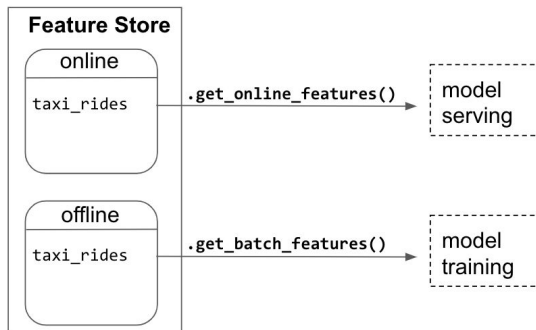


Retrieving from Feast

Feature data can be retrieved either offline, using historical features for model training, or online, for serving.

Two separate REST endpoints

Feast uses Redis and BigQuery for storage for online and batch respectively

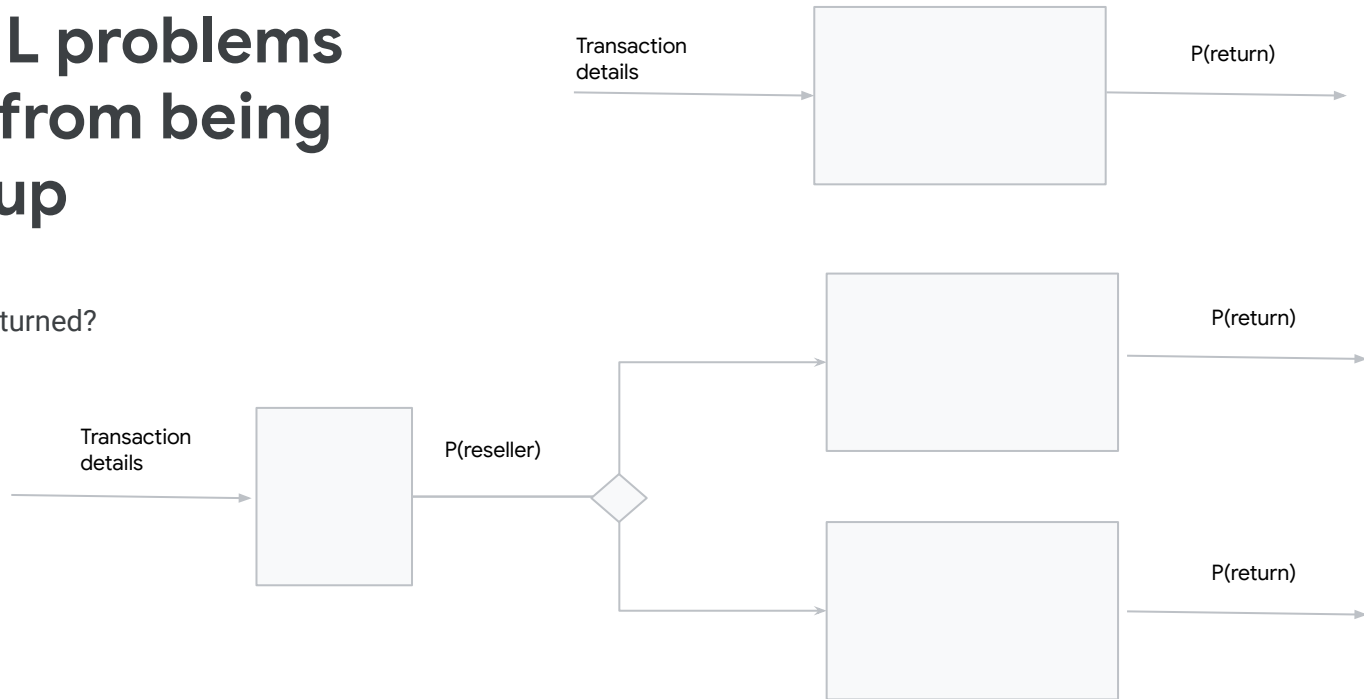


Cascade



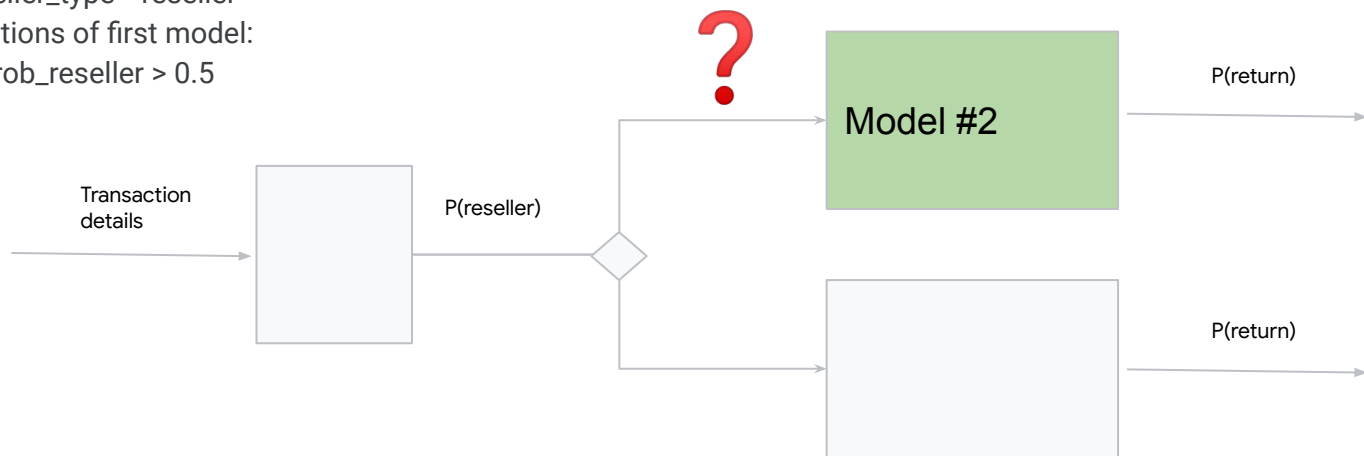
Some ML problems benefit from being broken up

Will this item be returned?



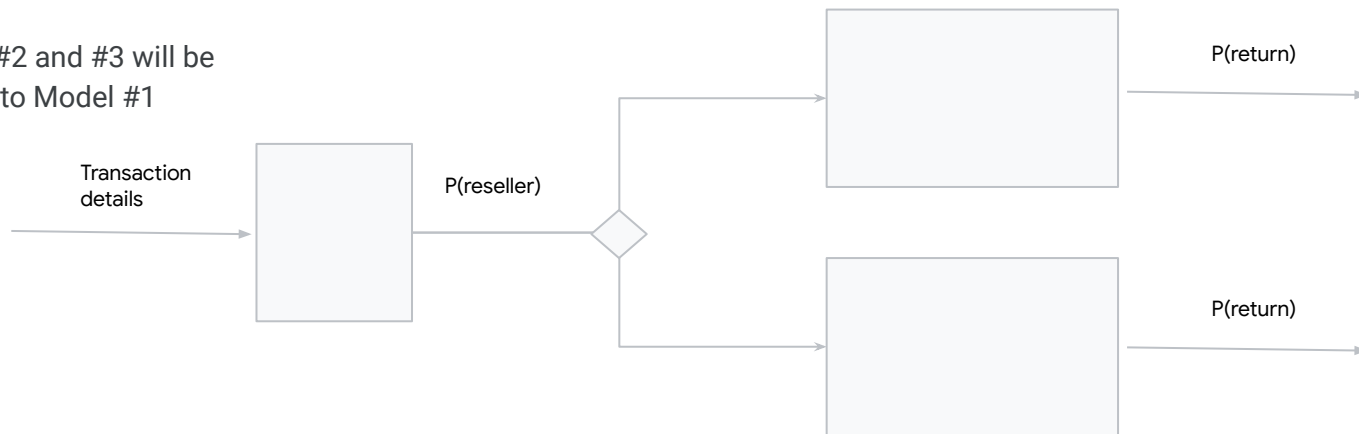
What is the training data for model #2?

- (a) Same as original training dataset:
WHERE seller_type="reseller"
- (b) Use predictions of first model:
WHERE prob_reseller > 0.5

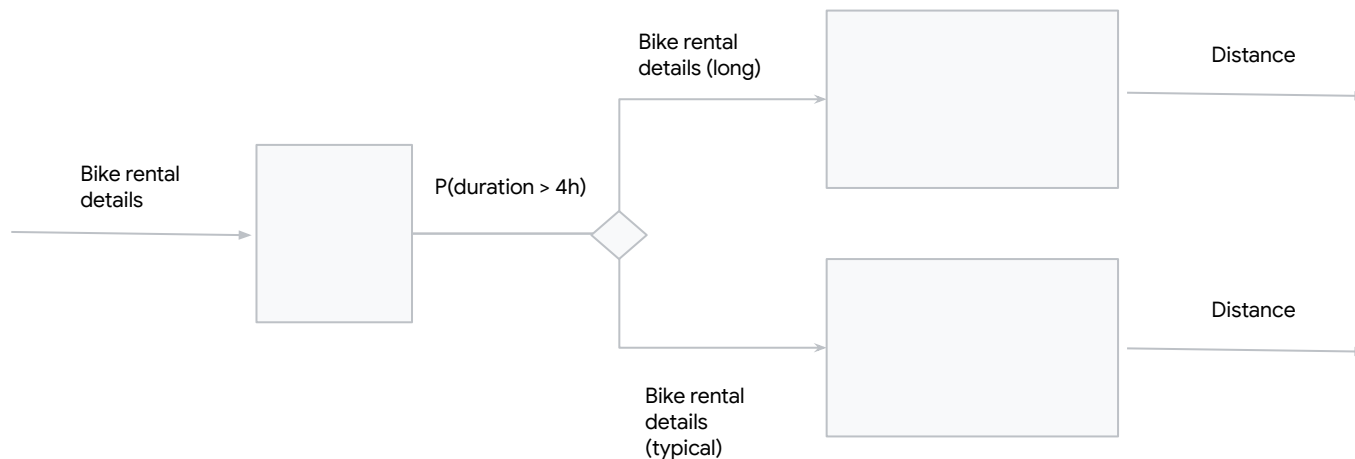


Training this model requires training a Cascade

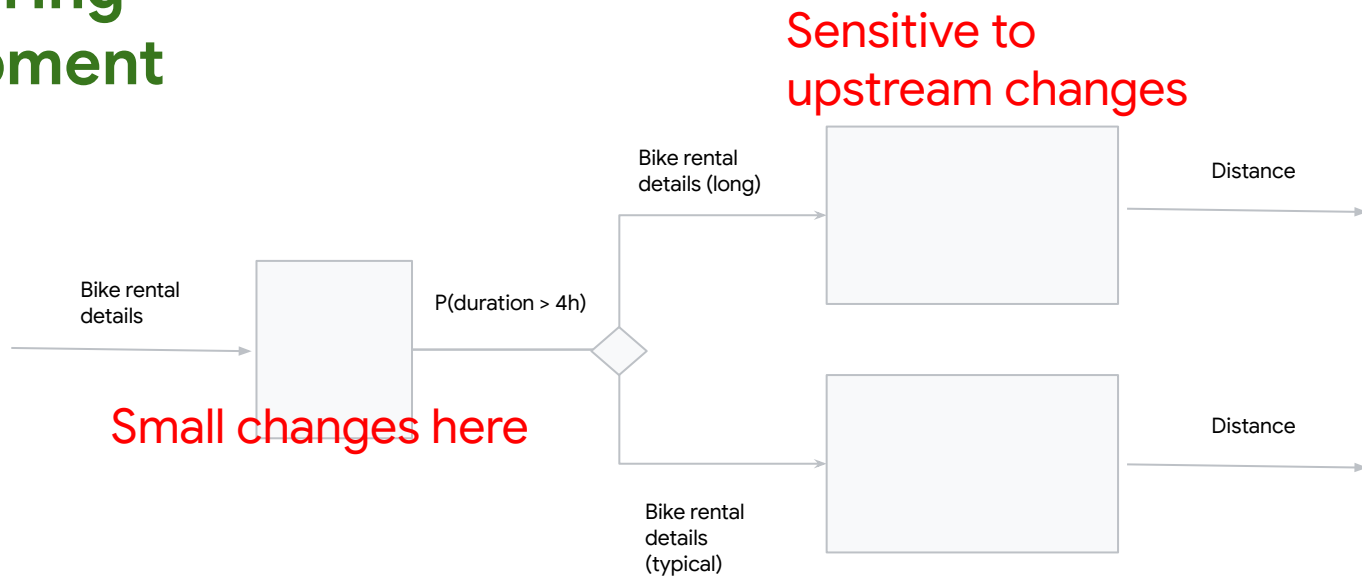
Otherwise, Model #2 and #3 will be
fragile to changes to Model #1



Use Cascade to handle rare scenarios with less data & simpler models



Use Workflow Pipeline to automate Cascade even during development



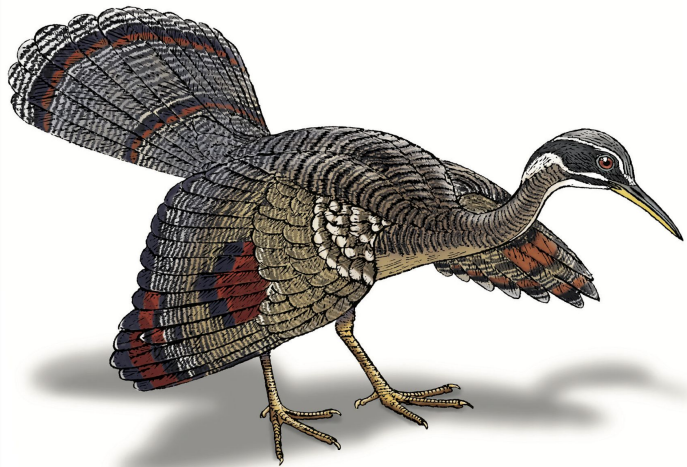
Summary

06

O'REILLY®

Machine Learning Design Patterns

Solutions to Common Challenges in Data Preparation, Model Building, and MLOps



Valliappa Lakshmanan,
Sara Robinson & Michael Munn

Workflow Pipeline
Transform
Multimodal Input
Feature Store
Cascade

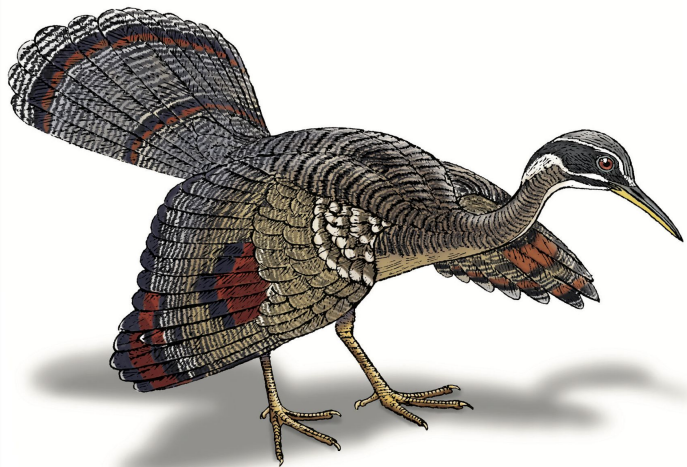
<https://github.com/GoogleCloudPlatform/ml-design-patterns>

Chapter	Design pattern	Problem solved	Solution
Data Representation	Multimodal Input	Choose between several potential data representations	Concatenate all the available data representations
Problem Representation	Cascade	Maintainability or drift issues when a machine learning problem is broken into a series of ML problems	Treat your ML system as a unified workflow for the purposes of training, evaluation, and prediction.
Reproducibility	Transform	The inputs to a model must be transformed to create the features the model expects and that process must be consistent between training and serving	Explicitly capture and store the transformations applied to convert the model inputs into features
	Workflow Pipeline	When scaling the ML workflow, you need a way to run trials independently and track performance for each step of the pipeline.	Make each step of the ML workflow a separate, containerized service which can be chained together to make a pipeline that can be run with a single REST API call
	Feature Store	The ad-hoc approach to feature engineering slows model development and leads to duplicated effort between teams as well as work stream inefficiency.	Create a Feature Store, a centralized location to store and document feature datasets that will be used in building machine learning models and can be shared across projects and teams

O'REILLY®

Machine Learning Design Patterns

Solutions to Common Challenges in Data
Preparation, Model Building, and MLOps



Valliappa Lakshmanan,
Sara Robinson & Michael Munn

Read the book (Nov 2020)

<https://bit.ly/ml-design-patterns>

Follow me on Twitter:

[@lak_gcp](https://twitter.com/lak_gcp)

Read my blog:

<https://medium.com/@lakshmanok>

Check out implementations:

<https://github.com/GoogleCloudPlatform/ml-design-patterns>