

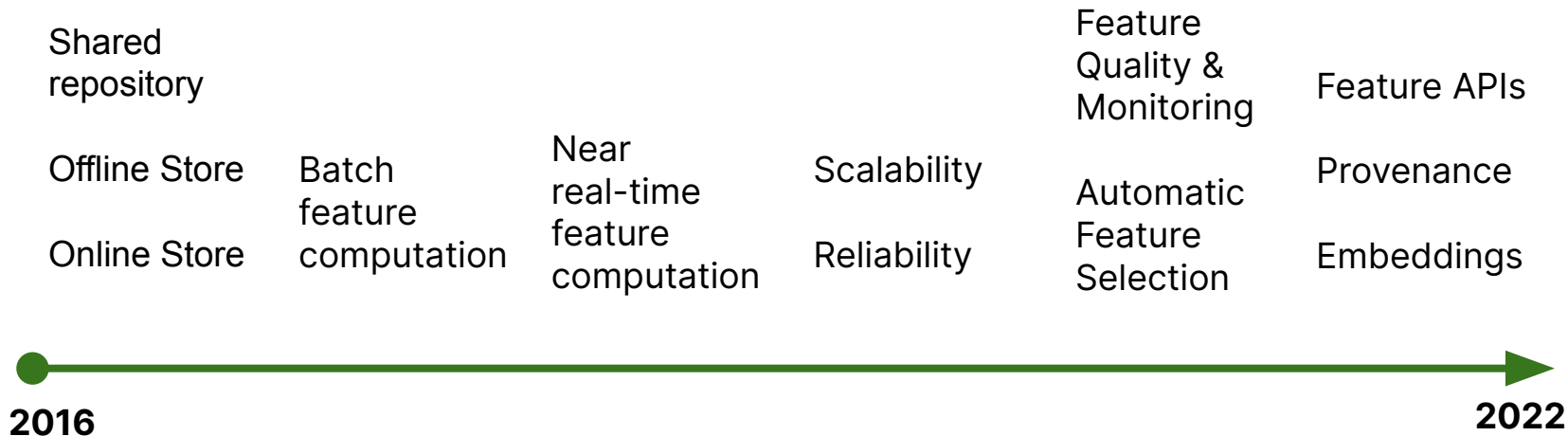
# Accelerating ML at Uber with Michelangelo Palette

[Amit Nene](#), Tech Lead Manager

Uber Technologies Inc.

# Michelangelo Palette

Evolution into a Feature Engineering Platform



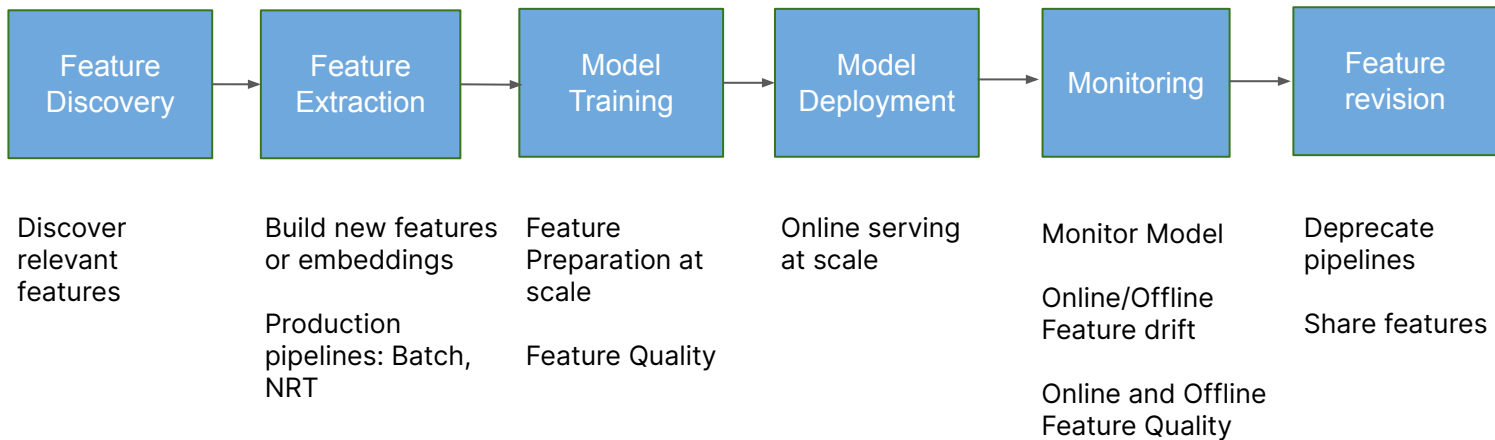


## Feature Engineering: User journey

*Number of trips  
taken over last N  
days*

*Number of orders  
from a particular  
restaurant*

*Average rating of  
user*



# Feature Discovery

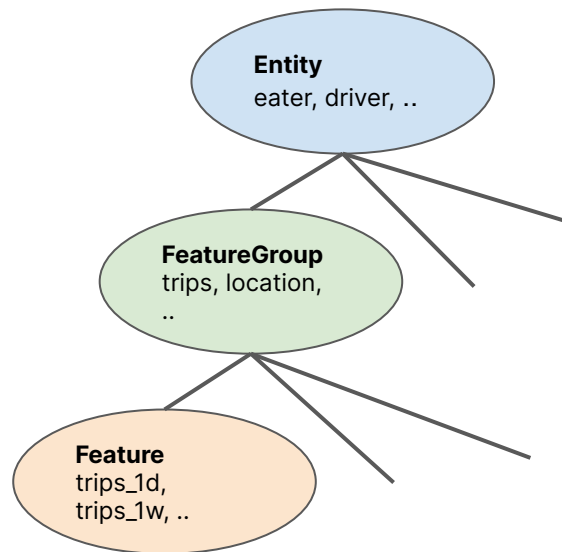
Browsing through Git repo

Entity: Primary entity associated with a Feature Group (table)

Feature Group: A logical table of individual features

Feature: A logical column within a Feature Group

Metadata: Type information, Join/Lookup Key, Feature extraction ETL, Online dispersal info



# Feature Discovery

## Keyword Search

**Databook** Explore Ownership Create ? Q eater

**Results for "eater"** x

All Dataset Metric Report Dashboard **ML Features** Schema

**Tier** — 58 results

- ☒ Tier 1
- ☒ Tier 2
- ☐ Tier 3
- ☐ Tier 4
- ☐ Tier 5
- ☐ Tier missing

**Team owner** —

Q uOwn Group

**eater\_source** Tier 1 ML Feature

[URL, Source] [https://featuregraph.uber.internal.com/uber-internal/data/michalangelu-features/data/michalangelu-features/eater\_source/michalangelu\_eater\_source.sql]

Entity: Feature Group: Join Key: Owner:

eater\_source user\_user Data Discovery Data Rotation

**eater\_source\_prices\_v1** Tier 1 ML Feature

[URL, Source] [https://featuregraph.uber.internal.com/uber-internal/data/michalangelu-features/data/michalangelu-features/eater\_source\_prices\_v1/michalangelu\_eater\_source\_prices\_v1.sql]

Entity: Feature Group: Join Key: Owner:

eater\_source\_prices\_v1 user\_user/michalangelu\_eater Data Discovery Data Rotation

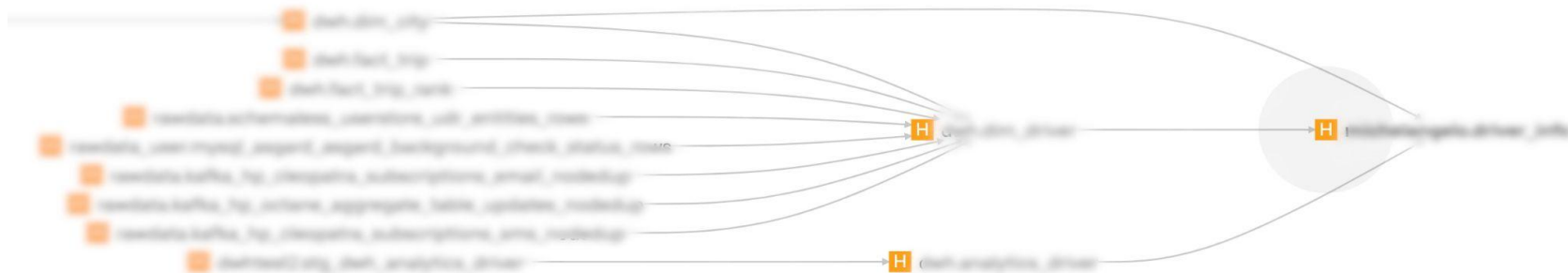
# Feature Discovery

## Data Lineage

**Warehouse: dim, fact, business tables, event logs**

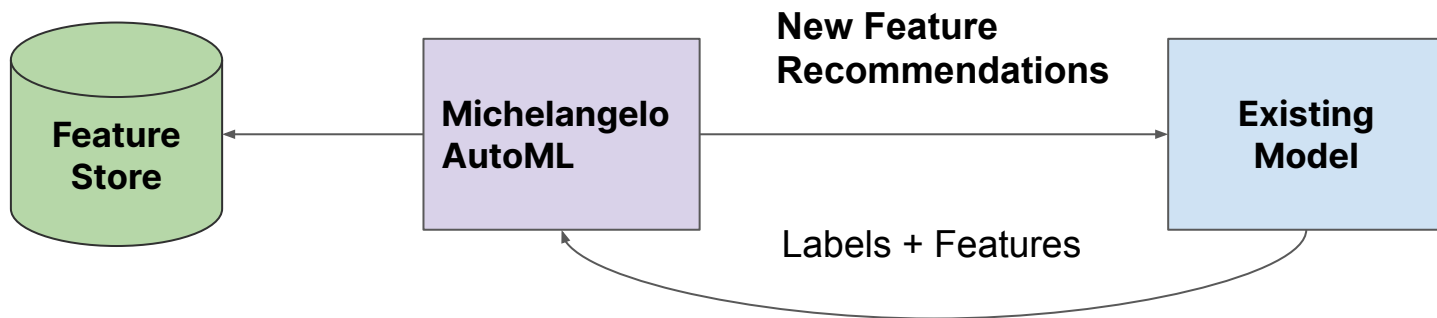
**Staging tables: joins, aggregations**

**Feature Store tables (Point-in-time snapshots)**



## Feature Discovery

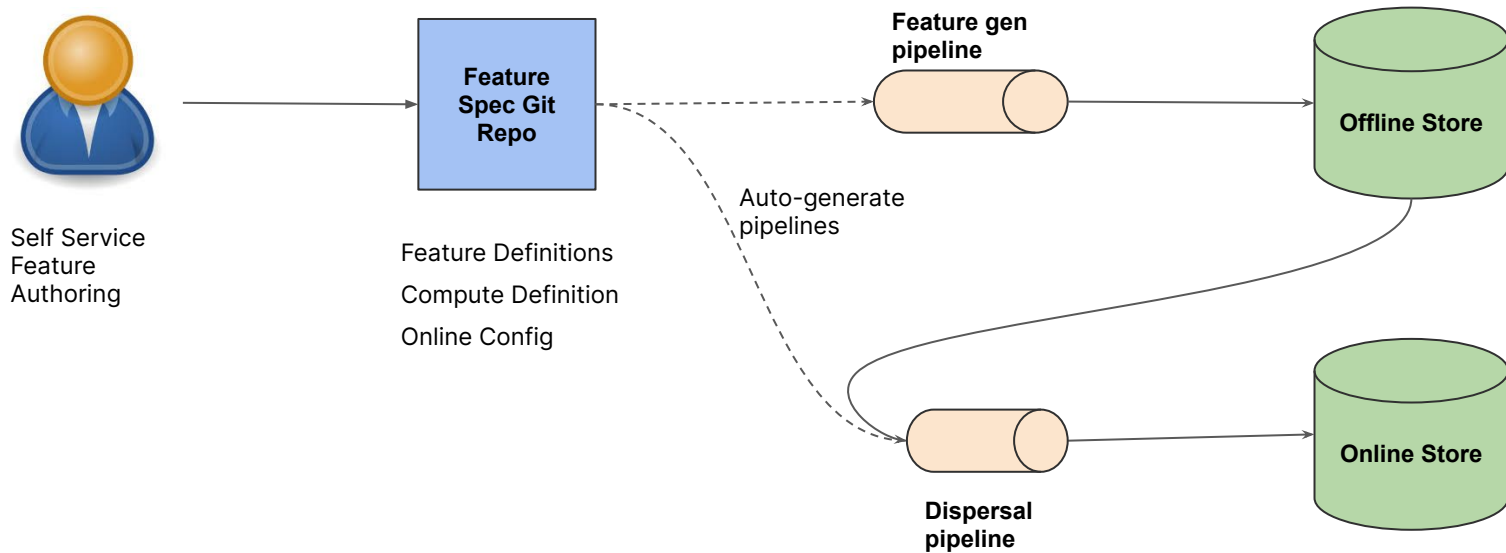
Advanced: Michelangelo AutoML tool



Boost Model Performance through  
Feature Discovery

# Feature Extraction

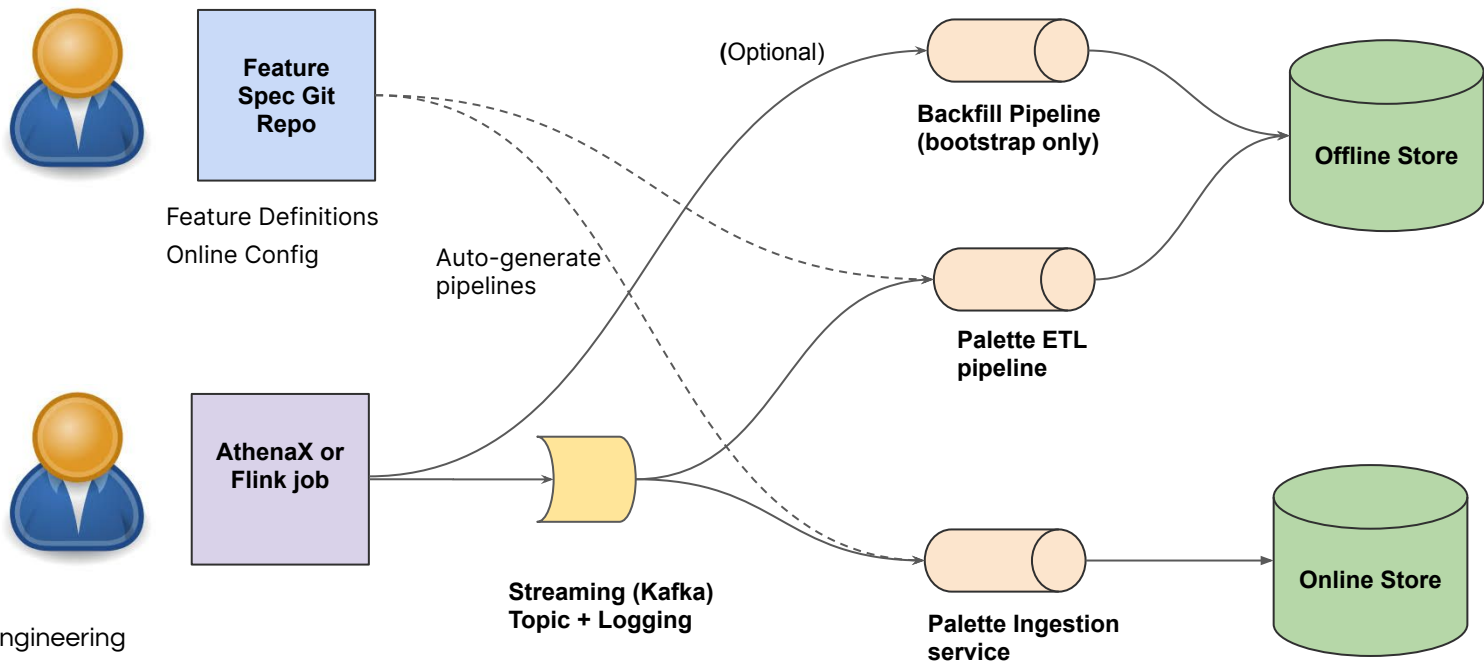
Batch computed features





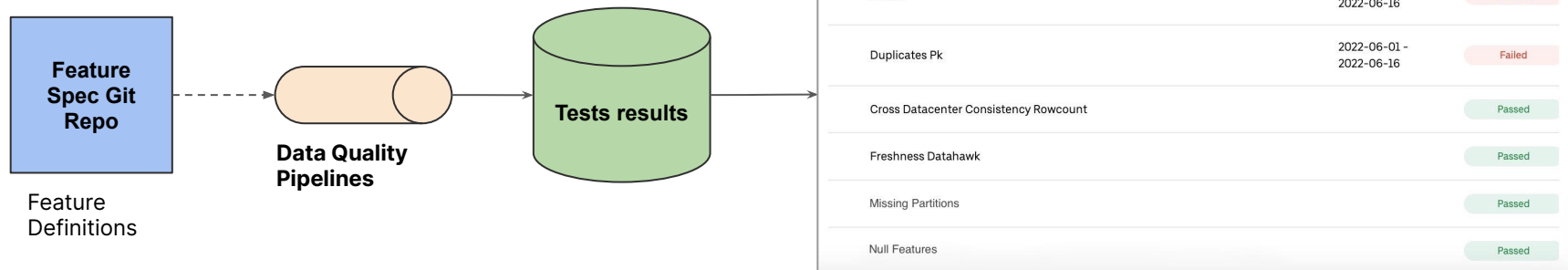
# Feature Extraction

Near real-time features: automation critical for iteration velocity!



# Model Training

Feature Quality: avoid debugging training time failures



At Training kickoff..

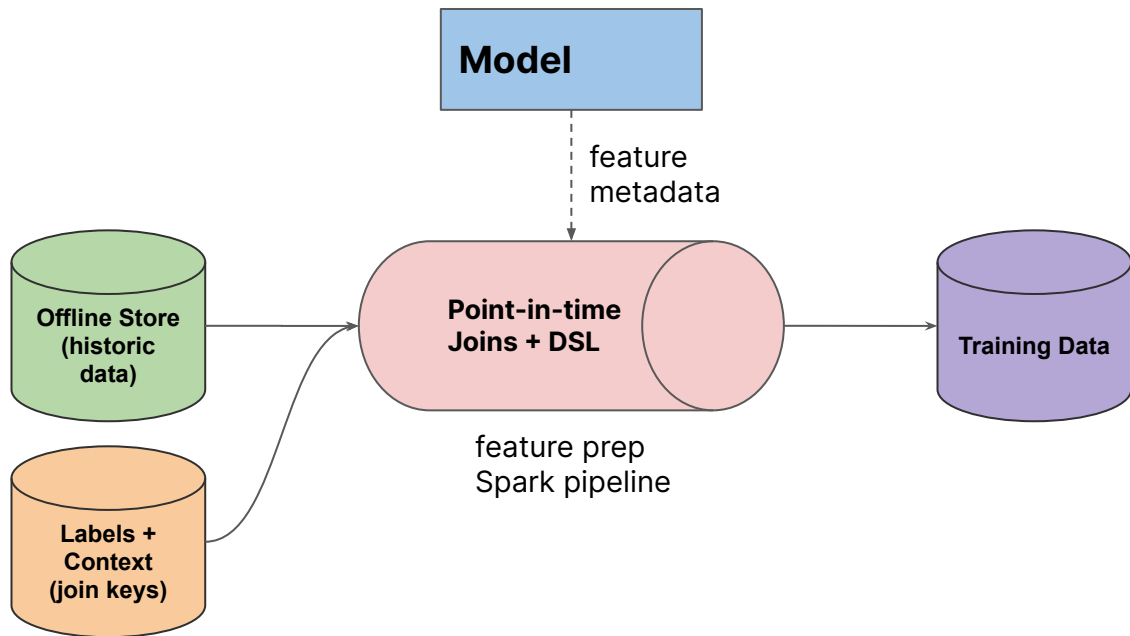
# Model Training

## Feature Preparation at Scale

Materialize at use: DSL for  
lightweight or  
model-specific xforms

Reducing shuffle overhead  
of joins is key

Incremental feature  
computation  
(Delta Store)



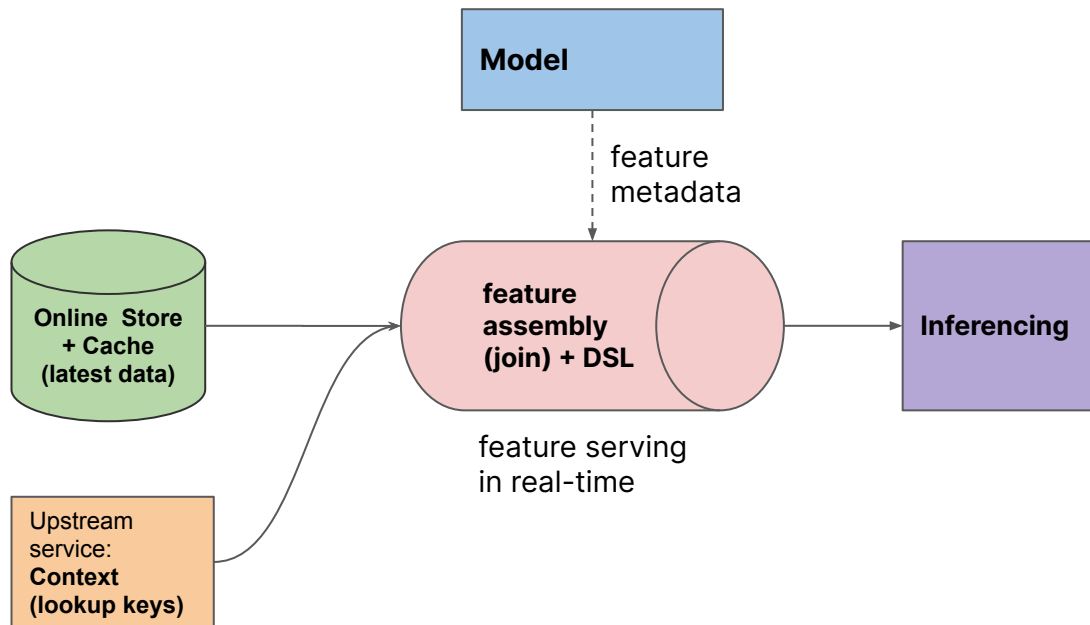
# Online Model Deployment

## Provisioning for Scale

Materialize at use: DSL +  
feature join has semantic  
parity to offline

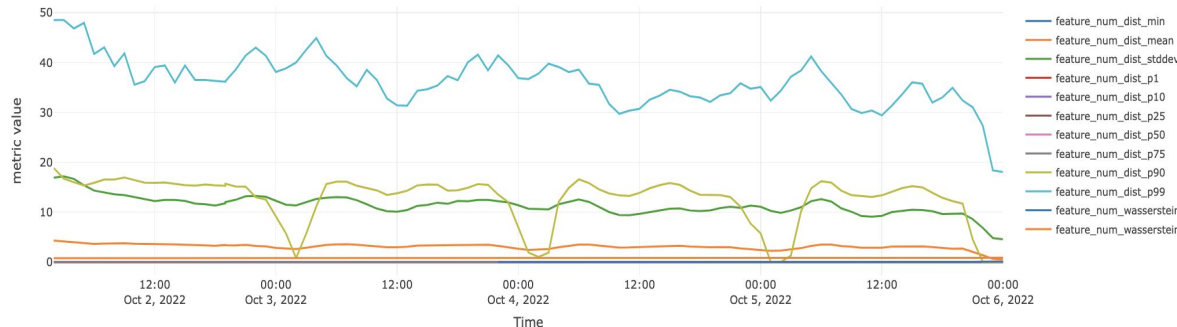
Process to estimate  
Redis + Cassandra capacity  
from access patterns and  
feature payload

Shadow Testing through  
cloning Production traffic



# Model Monitoring

Near real-time monitoring:  
use prediction logs to detect  
feature distribution shifts  
wrt to training/historic  
baseline



Feature Store monitoring:  
offline only (pre-training,  
pre-dispersal)

metric_name	min	max	avg	training_baseline	feature_importance
filter data...					
feature_num_dist_min	0	0	0	0	0.766606
feature_num_dist_mean	0.530115	4.285806	3.078492	25.7552	0.766606
feature_num_dist_stddev	4.556955	17.183177	11.366911	33.591095	0.766606
feature_num_dist_p1	0	0	0	0	0.766606
feature_num_dist_p10	0	0	0	0	0.766606
feature_num_dist_p25	0	0	0	19.375021	0.766606
feature_num_dist_p50	0	0	0	37.451324	0.766606

# Feature Revision

## Feature Sharing, Deprecation

Which Michelangelo models (including dev,deployed,online,offline) are using a Palette Feature?

```
# deployed_only: if set to true, only returns the online and offline models that are deployed
# palette_expression: A full or partial palette expression,|
# Return projects and all the models ordered by reverse timestamp (if show_models is set to true) that are using the palette_expression
# Note: 1. The model id prefix stands for tmYYYYMMDD-HHMMSS.
# 2. The output for the models that match the palette expression is limited to the most recent 10000

palette_expression = 'palette:reater:demo_feature'

get_models_for_features(
    deployed_only = False,
    palette_expression = palette_expression,
    show_models = False
)
```

59 Projects use this Palette Feature:

```
demo_12_31_2021 077 models
demo_12_31_2021 077 models
demo_12_31_2021 077 models
demo_12_31_2021 077 models
demo_12_31_2021 077 models
demo_12_31_2021 077 models
demo_12_31_2021 077 models
```

**Deprecate feature**

No new Models  
Existing Models allowed

**Remove Feature**

No Model references

Feature → Model Lineage via metadata



## Accelerating end-to-end ML

### Summing it up

- Feature Discovery through search and automatic feature selection
- Feature Extraction via automated pipeline generation from spec
- Feature Preparation through data validation and scalable point-in-time joins
- Feature Serving in real-time through a process for capacity estimation, and scalable online store
- Near real-time feature monitoring via prediction logs, daily monitoring of features
- Feature sharing and deprecation via Model-Feature lineage tracking

**Thank you!**



