



# Vertex AI Forecast

Highest accuracy forecasts with  
state-of-the-art machine learning models

Google Cloud



Mark Tenenholtz 

@marktenenholtz

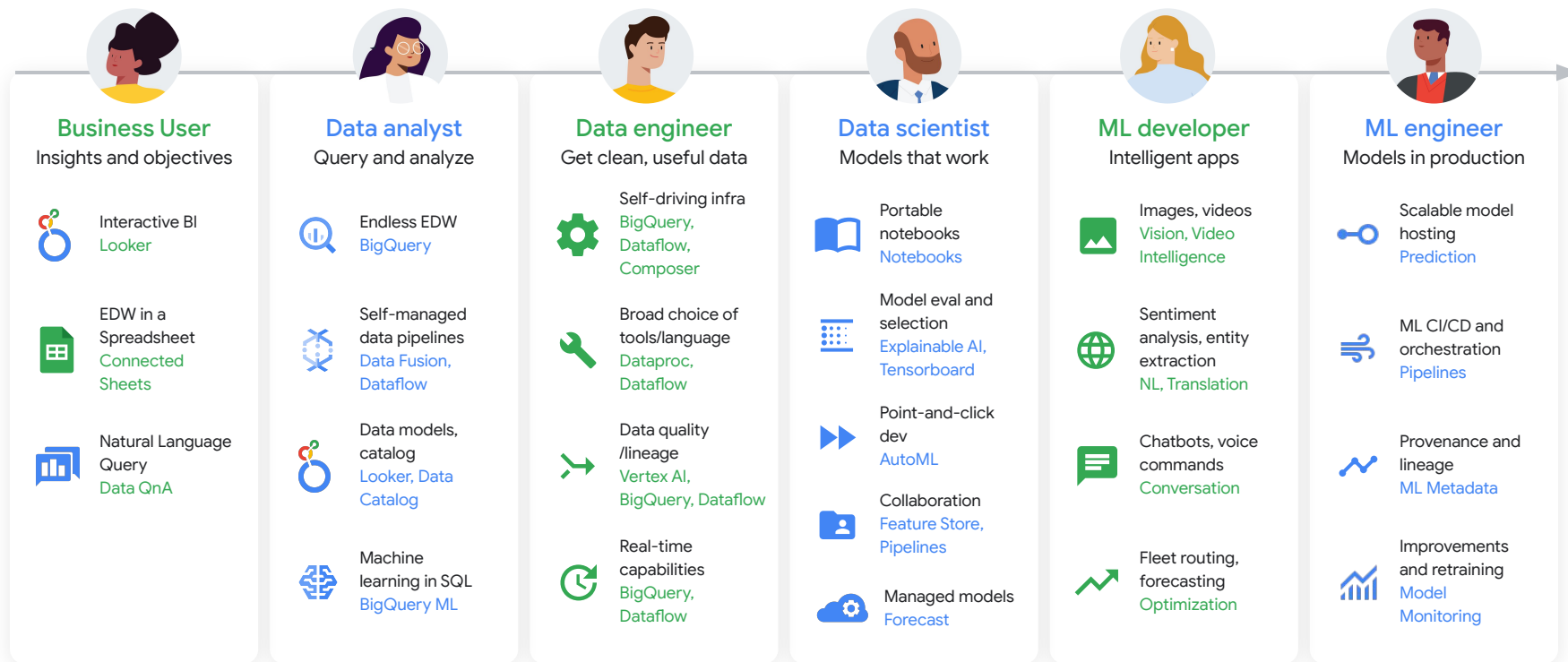
In NLP, it's really easy to consistently beat a traditional bag-of-words approach with a transformer.

In computer vision, it's really easy to consistently beat an SVM with a CNN.

In time-series forecasting, it's \$\*&%ing hard to consistently beat a seasonal rolling average.

5:00 AM · 09 Nov 22 · [Typefully](#)

# Google Cloud: Machine Learning tools for everyone





# Vertex AI

## Applications

Vision and Video

Conversation

Language

Structured Data

## Core

Workbench

Data Labeling

Deep Learning Env

Experiments

Metadata

AutoML

Training

Explainable AI

Feature Store

Vizier (Optimization)

Prediction

Continuous Monitoring

Pipelines


AI Accelerators

Hybrid AI

# Structured Data Problem Types

## Forecasting

How many products will be sold next month?

Sales	Date	Channel	SKU	...	Geo	Domain	Brand
\$1,200	Jun-10-2020	Website	12345	...	US	Shoes	Nike
\$1,350	Jun-11-2020	Email	54221	...	US	Shirts	Adidas
	Jan-10-2021	Website	22345	...	CA	Shoes	Asics

## Classification

Will this product sell in 7 days?

Sold	...	ID	Geo	Domain	Title	Description	Tags	Brand
YES	...	104	US	Shoes	"Dark red..."	"Try this soft..."	["A, B, ..."]	Nike
NO	...	204	US	Shoes	"Women's..."	"Medium-size..."	["A, B, ..."]	Adidas
	...	302	CA	Shoes	"Running..."	"All-terrain..."	["A, B, ..."]	Asics

## Regression

At which price will this product sell?

Price	...	ID	Geo	Domain	Title	Description	Tags	Brand
\$52	...	104	US	Shoes	"Dark red..."	"Try this soft..."	["A, B, ..."]	Nike
\$48	...	204	US	Shoes	"Women's..."	"Medium-size..."	["A, B, ..."]	Adidas
	...	302	CA	Shoes	"Running..."	"All-terrain..."	["A, B, ..."]	Asics

# Forecasting Methods

## Statistical

ARIMA, Exponential Smoothing, etc

- Theoretically grounded (>50 years of research)
- Mature tools and ecosystem
- Very popular
- Supported at Google via BQML ARIMA\_PLUS

## Neural Networks

CNN, RNNs, LSTMs

- New (~5 years)
- Active debate vs statistical methods
- Recently competed in prime competitions (M4, M5)
- Ecosystem is very raw
- Supported at Google via AutoML Forecast

## Statistical

## Neural Networks

### Data set

Univariate series



Multivariate series



Few features



Many features



### Patterns in Time Series

Repeated patterns



Feature driven patterns



### Other

Cold starts



Short life cycle products



Method Fit: great ok poor

## Tree-based model risks

- Can't extrapolate to values outside the training dataset.
- Doesn't capture the sequences in the dataset as well as DNN models.
- Returns good evaluation metrics (illusion risk) but not always business value. Can have issues with under forecasting critical SKU's.



# Forecasting Methods

**Neural Network** builds a global model for all series, “cross-learns” from series.

It works well when

- there is lots of data (wide and long)
  - large number (100+) of time series (especially short series w/o data on repeated seasons)
  - Variable length time series, including short histories
  - lots of features
  - data has unstructured features (like product description text fields)
- series are strongly driven by features
- cold-starts

**Statistical (autoARIMA)** builds a separate model for each individual series. Typically retrained often.

It works well when

- series are strongly driven by seasonality and/or trends
- no limitation on dataset sizes, works equally well with
  - small number of series or univariate data
  - Long histories - enough to capture 2 or more seasonal
  - few or no features

# Why use Deep Learning Models for Forecasting?

- Creates one “global” model for all time series
- Learns patterns **across** time series

## Can use **large number of drivers**

- ✓ **Rich metadata**  
(eg: product attributes, location attributes)
- ✓ **Historical factors**  
(eg. inventory, weather)
- ✓ **Factors known in the future**  
(eg: planned promotions/events, holidays)

## Can model **complex scenarios**

- ✓ Cold start / new items
- ✓ Short product life cycles
- ✓ Burstiness, sparsity
- ✓ Unstructured data such as text descriptions



# Vertex AI Forecast

1. Prepare Dataset

2. Train Model(s)

3. Forecast

4. Visualize, Integrate



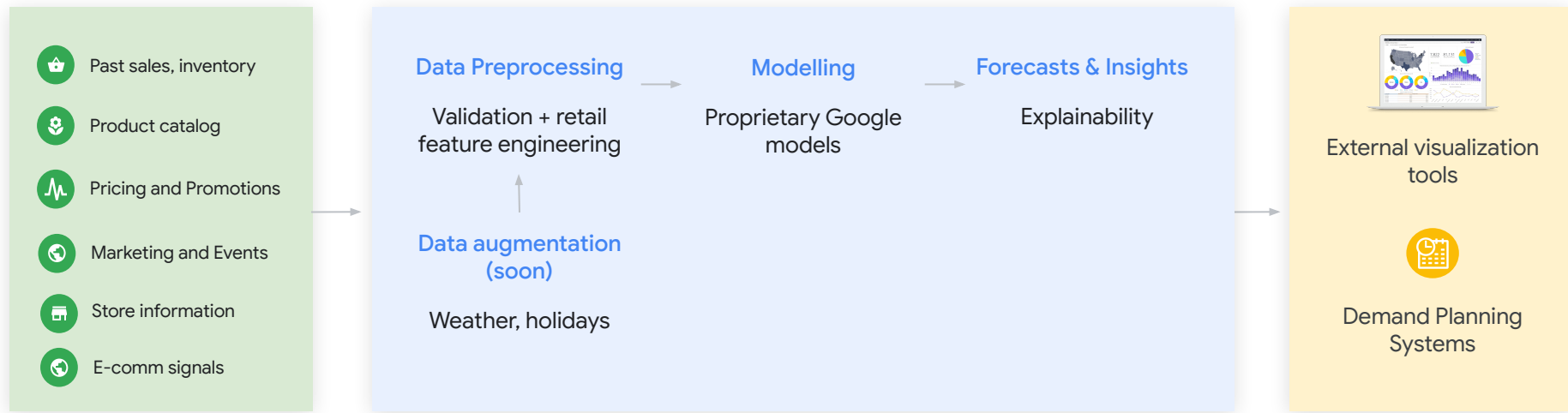
## Retailer Data



## Vertex AI Forecast

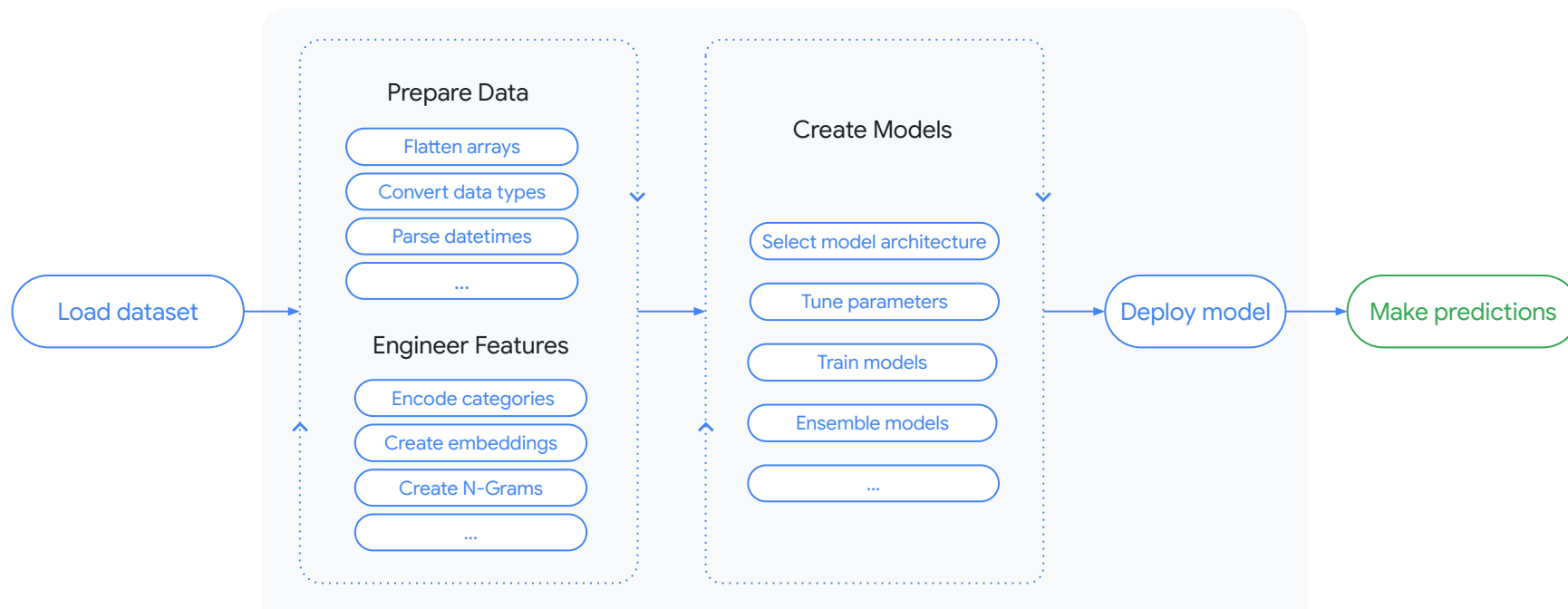


## Integration



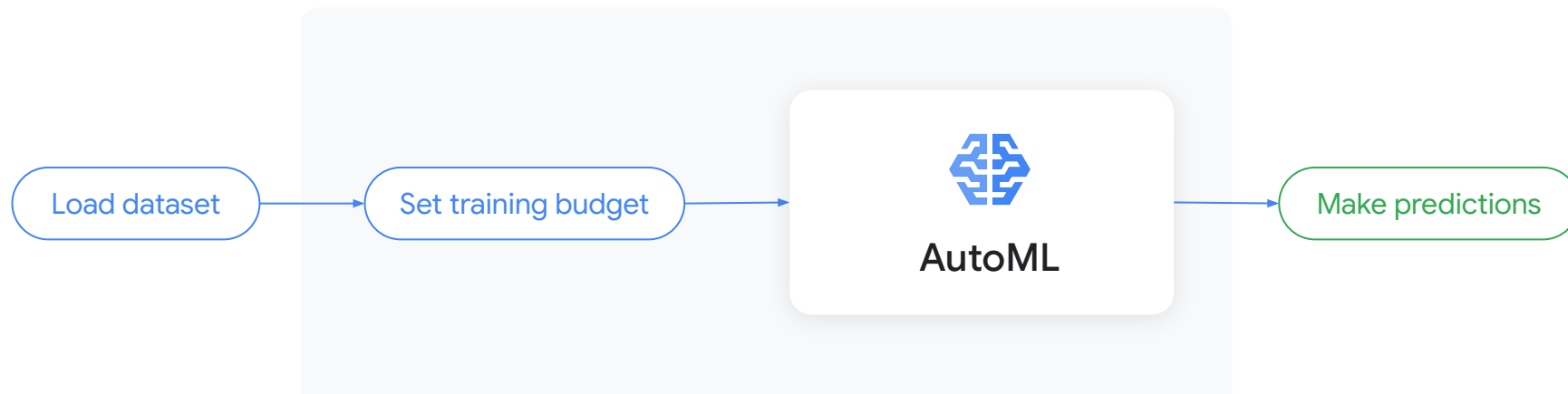
# AutoML - Fastest path from data to value

## Traditional Machine Learning Workflow

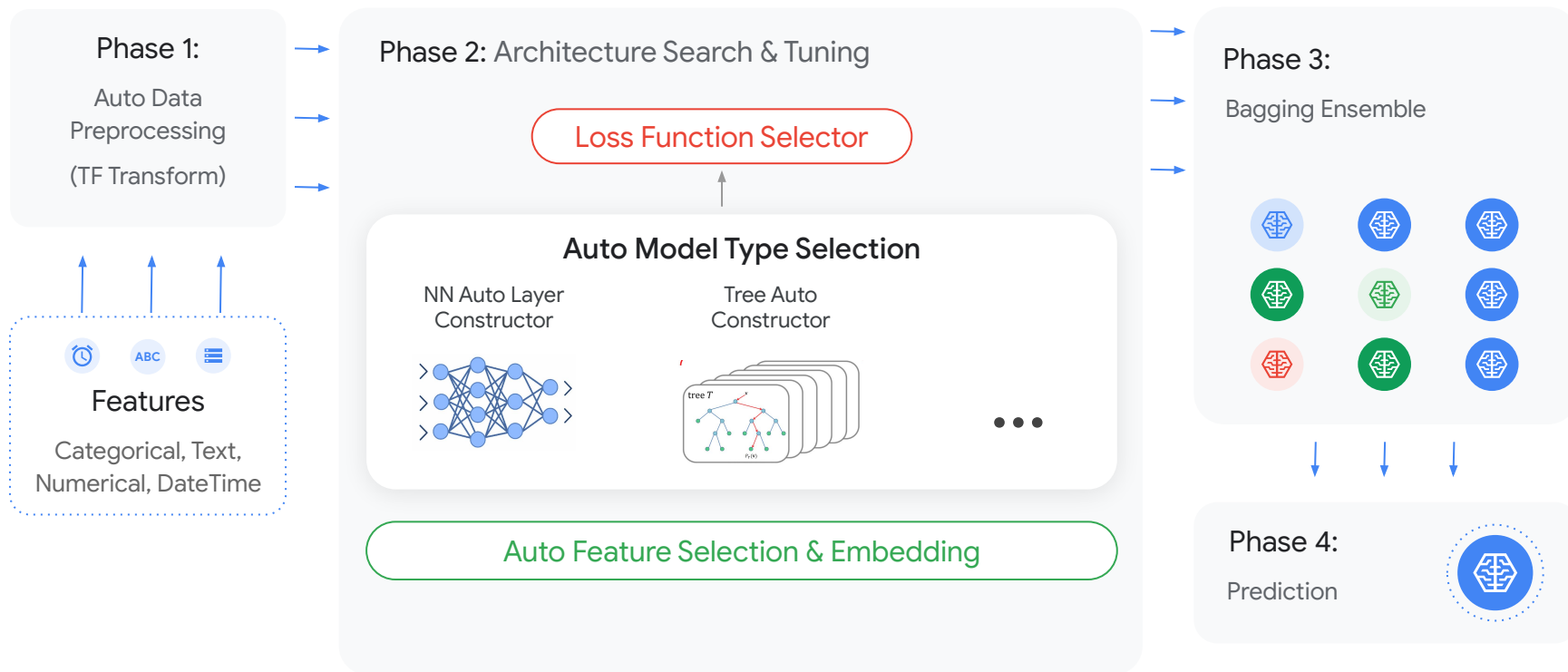


# AutoML - Fastest path from data to value

## AutoML Workflow



# Powered by latest research from Google Brain



# Automated feature engineering

## Best practice transformations for all data types



**Numbers:** generate quantiles, log, z\_score transforms



**Datetime:** extract year, month, day, weekday, categorize



**Text:** tokenize, generate n-grams, create embeddings



**Arrays of categories:** convert to lookup index, generate embeddings



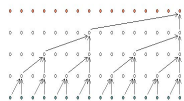
**Categories:** one-hot encoding, grouping, embeddings



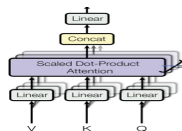
**Nested fields:** flatten, apply type transformations

## Automated model architecture search

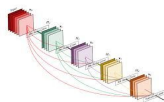
## Evaluating Google's best model architectures for forecasting



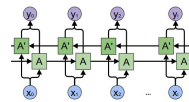
## Convolutions



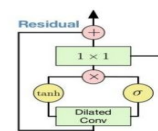
## Attention



## Skip Connections



## LSTM



## Gating



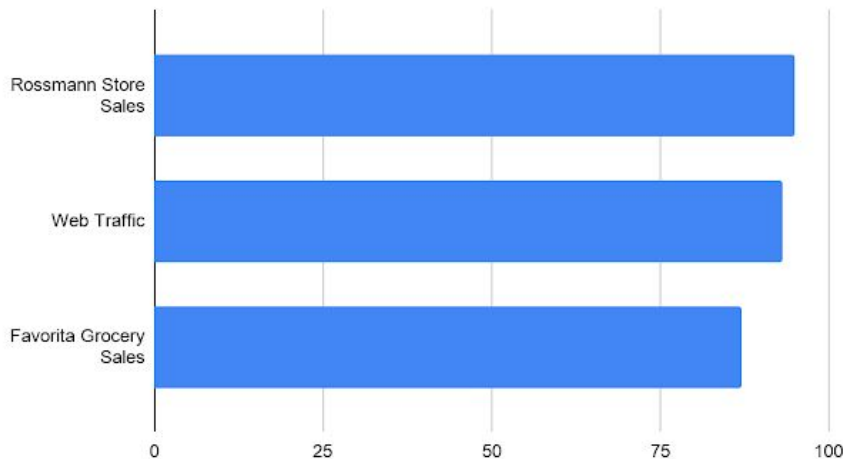
## Benchmarks on Kaggle Datasets

Finished in **top 2.5%** (138 out of 5558 teams) in World's Top Forecasting competition.

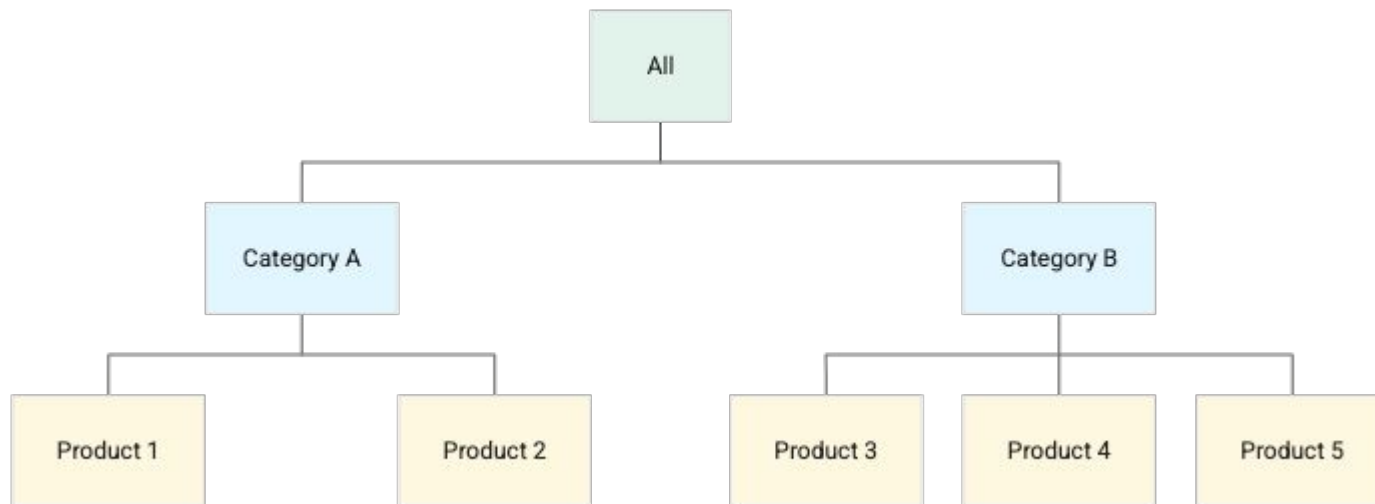
[M5: Estimate the sales of Walmart retail goods.](#)

Consistently ranks in top 20% across a variety of datasets in different industries.

Better than % of Kaggle Participants

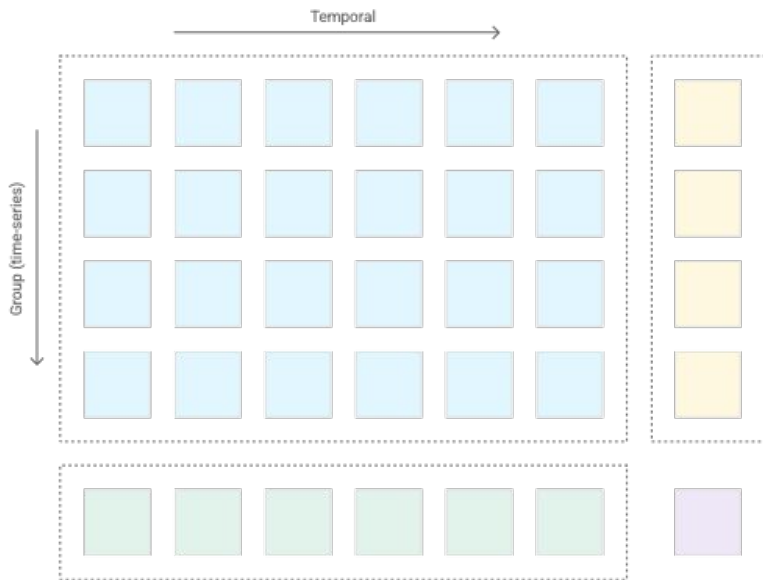


# Hierarchical Forecasting





# Hierarchical Forecasting

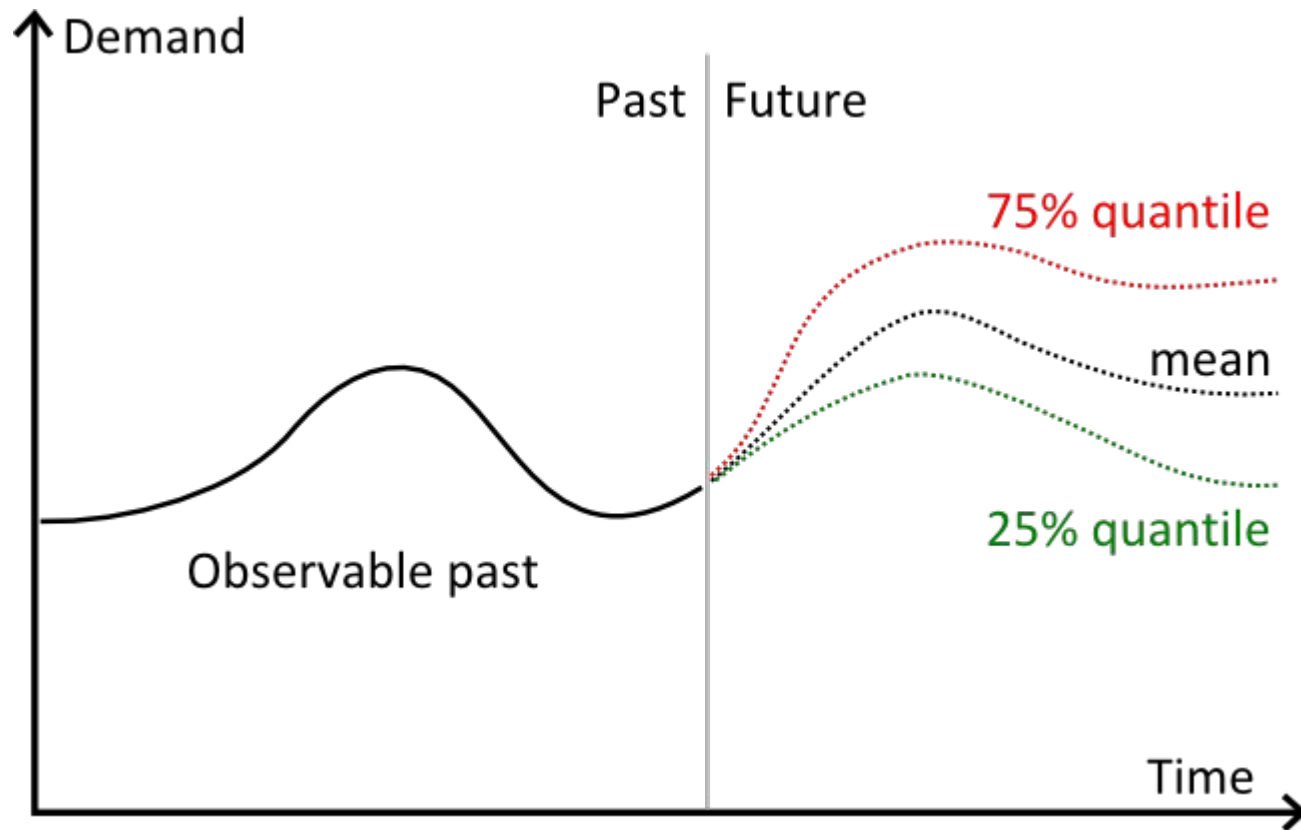


Hierarchical loss = 1 x loss  
+ temporal total weight x temporal total loss  
+ group total weight x group total loss  
+ group temporal total weight x group temporal total loss

 Aggregate over time-series in group  Aggregate over time  Aggregate over both

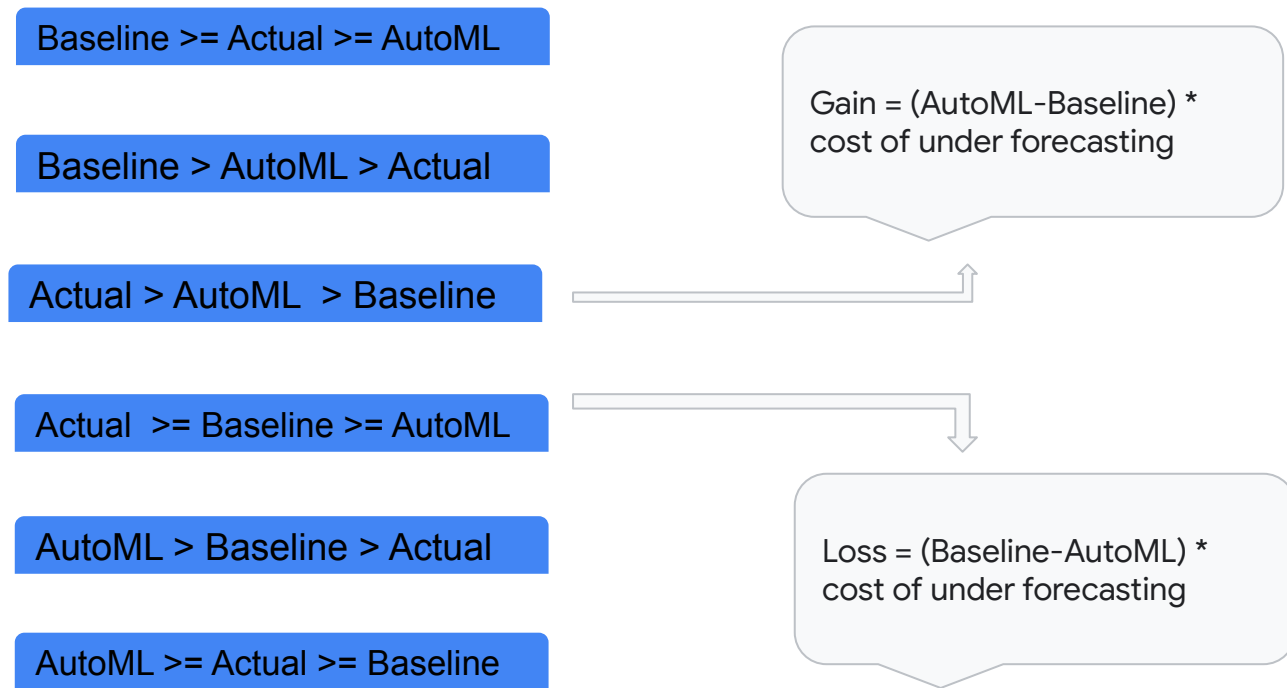
- **Reduce overall bias** to improve metrics over all time series (total sales).
- **Reduce temporal bias** to improve metrics over the horizon (season sales).
- **Reduce group level bias** to improve metrics over a group of time series (item sales)

# Quantiles

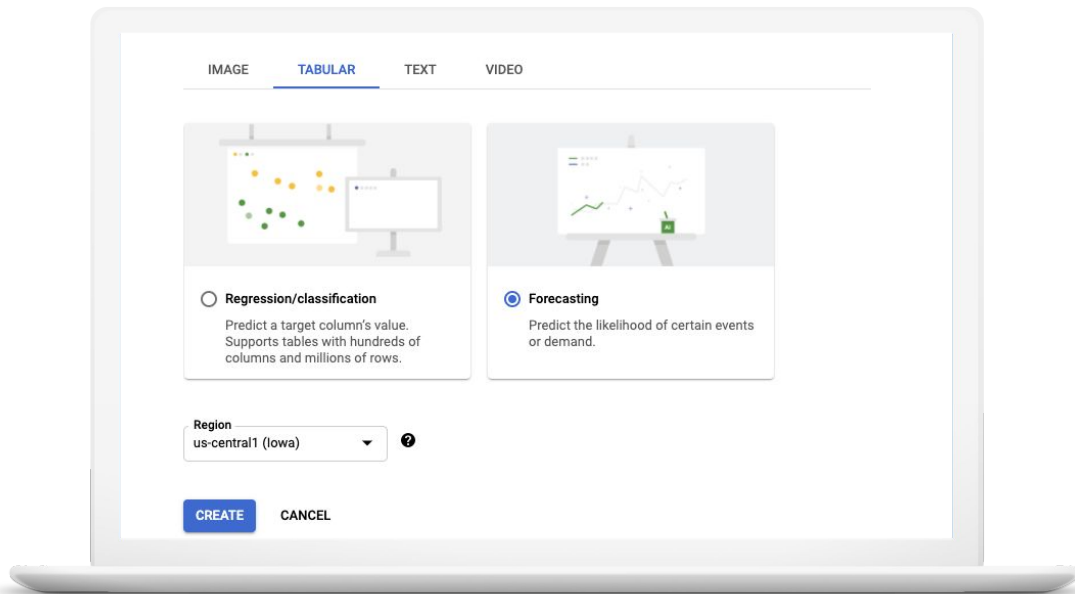


# Business and Financial Evaluation

Given you have the baseline, AutoML and actuals you can run a financial analysis as traditional ML metrics fail to observe the financial impact of over forecasting vs under forecasting.



# Variety of Interfaces



Intuitive Web UI

## API + SDK

```
In [ ]: # The number of hours to train the model.
model_train_hours = 1 ##param {type:'integer'}

create_model_response = tables_client.create_model(
    model_display_name=MODEL_DISPLAY_NAME,
    dataset=dataset,
    train_budget_milli_node_hours=model_train_hours*1000,,
    exclude_column_spec_names=['fnlwgt','income'],
)

operation_id = create_model_response.operation.name

print('Create model operation: {}'.format(create_model_response.operation))
```

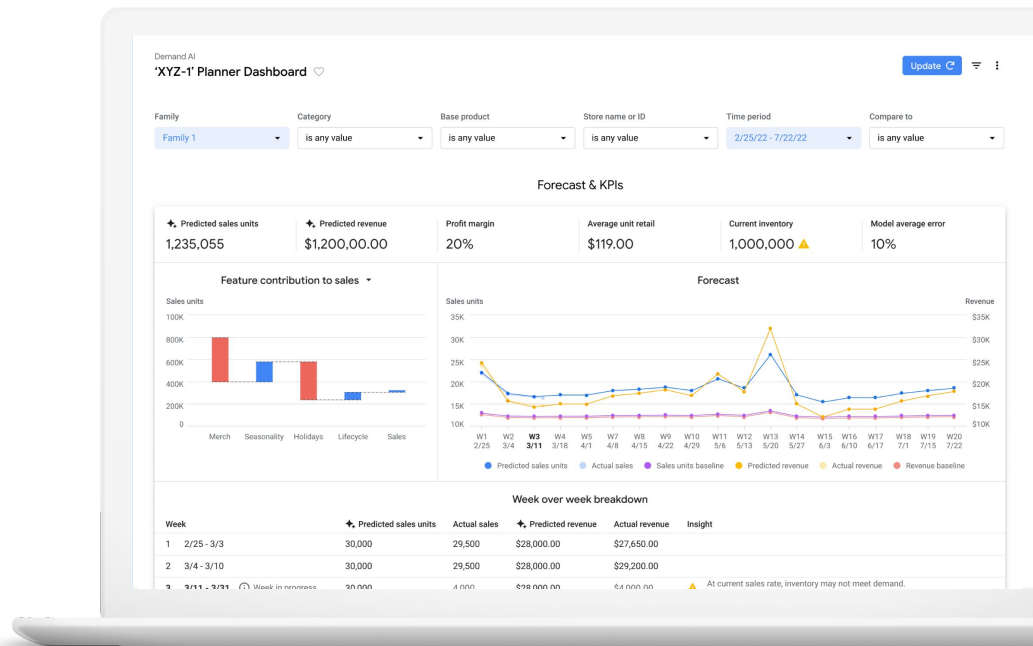
## BigQuery SQL\*

```
CREATE OR REPLACE MODEL project_id.mydataset.mymodel
  OPTIONS(model_type='AUTOML_REGRESSOR',
    input_label_cols=['fare_amount'],
    budget_hours=1.0)

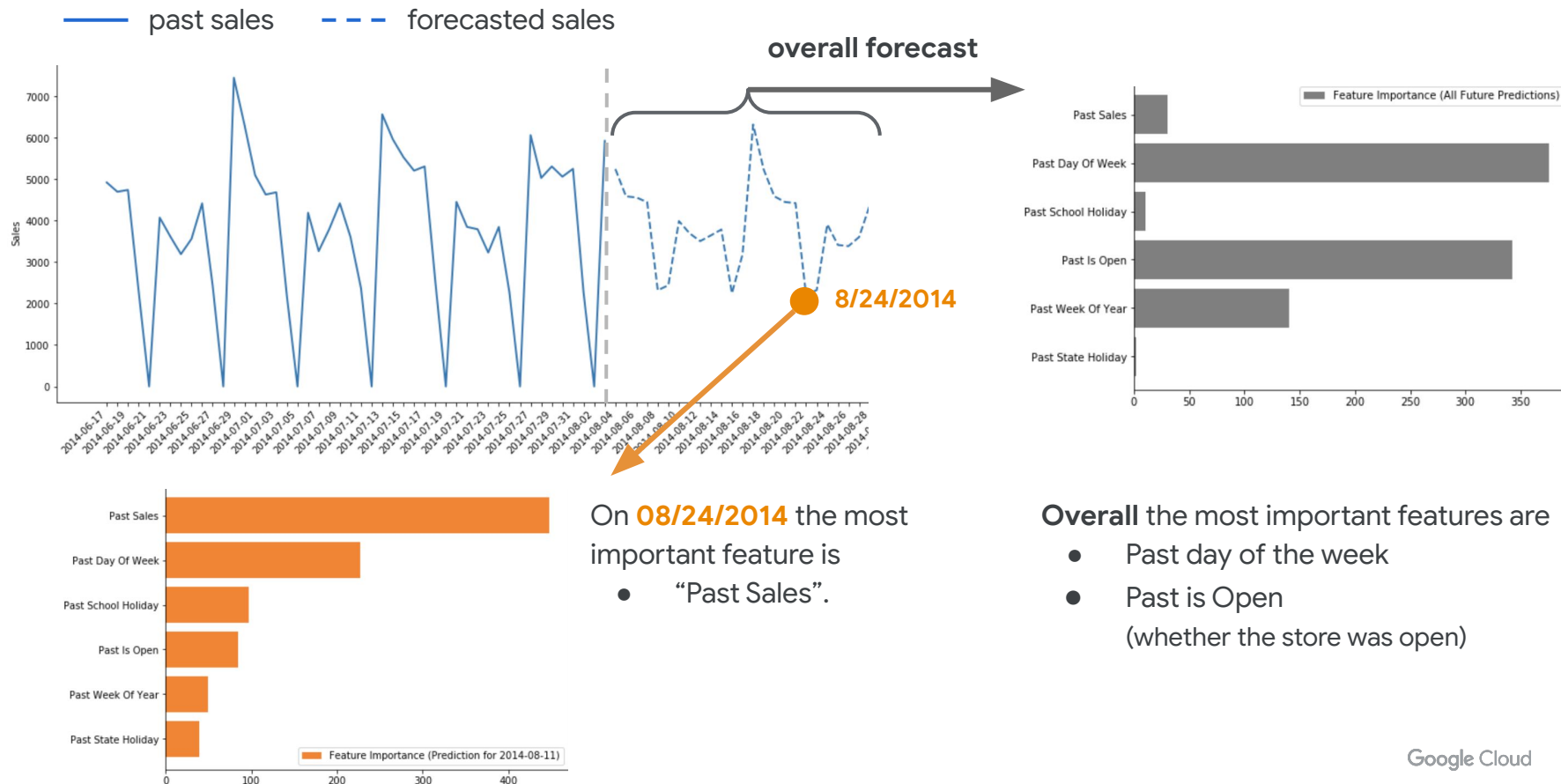
AS SELECT
  (tolls_amount + fare_amount) AS fare_amount,
  pickup_longitude,
  pickup_latitude,
  dropoff_longitude,
  dropoff_latitude,
  passenger_count
FROM `nyc-tlc.yellow.trips`
WHERE ABS(MOD(FARM_FINGERPRINT(CAST(pickup_datetime AS STRING)), 100000)) = 1
```

# XAI for time-series / forecasting

- Feature attributions for time-series models based on Sampled Shapley method
- **Forecasting models differ from normal tabular models** in the form of their inputs and the way we construct baselines.
  - 3-dimensional data (num\_samples X num\_historical\_instances X num\_features)
  - Attributions are generated per-feature, *per-historical* instance
  - Baselines are generated taking into account the time-ordering of samples
- **Aggregations.** These can be aggregated over historical instances to obtain aggregate feature-attributions.
  - E.g. for a retail sales forecasting use-case, feature-attribution aggregations can be by *product, location or time-slice*



# Applying feature attributions to time-series





# Case Studies

Google Cloud

# Large US-based Apparel Retailer

## Challenges

For this **high-quality, on-trend apparel retailer**, demand planning is key to business success. With disparate systems, the company recognized that it required a new approach to address challenges involving big data, granularity of predictions, and real-time demand signals.

## Empowering business decisions with improved forecasting

Leveraging Vertex AI Forecast, the company saw increased accuracy for demand **and** labor forecasting while integrating new supply chain requirements into its forecasting capabilities.



**2-3 month**

reduction in model development time

**\$5M-\$10M per year**

savings from enhanced labor efficiency

**\$40M+**

estimated additional revenue from improved supply chain and product allocation



## Case Study 1: Pre-season forecast for Apparel

- **Both Pre- and in-season:** forecast up to more than a year, based on data up to *three months* before the season starts
- **Cold start items:** Products relationships have to be learned automatically.
- **Fine granularity:** breakdown demand forecast to product variants (eg. color x size) down to the day x SKU level.



?

# Retail Forecasting Primer: Item Demand

## PRODUCT METADATA



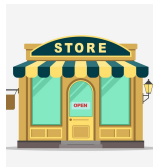
- Summer dress
- Lightweight. Pattern.
- Long-sleeve. V-neck.
- 55% Linen. 45% Cotton
- Machine wash

## DEMAND DRIVERS



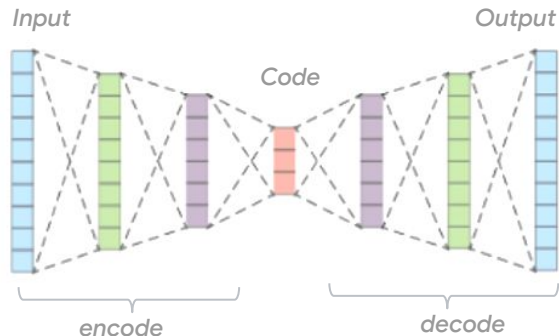
- Sales
- Price
- Competitor price
- Promotions/Events
- Holidays

## STORE DATA



- Store description: Large, small, specialty
- Store location
- Foot traffic

## AUTOML FORECAST



Machine learning models understand rich metadata, relationships between products and the joint effect of pricing, competition and product lifecycle.

## MEDIUM HORIZON 12-16 month

- e.g. pre-season planning
- Buy/order planning
- New & Cold start items

## SHORT HORIZON 0-8 weeks

- in-season planning
- Replenishment, inventory
- Pricing, allocation
- SKU-Level



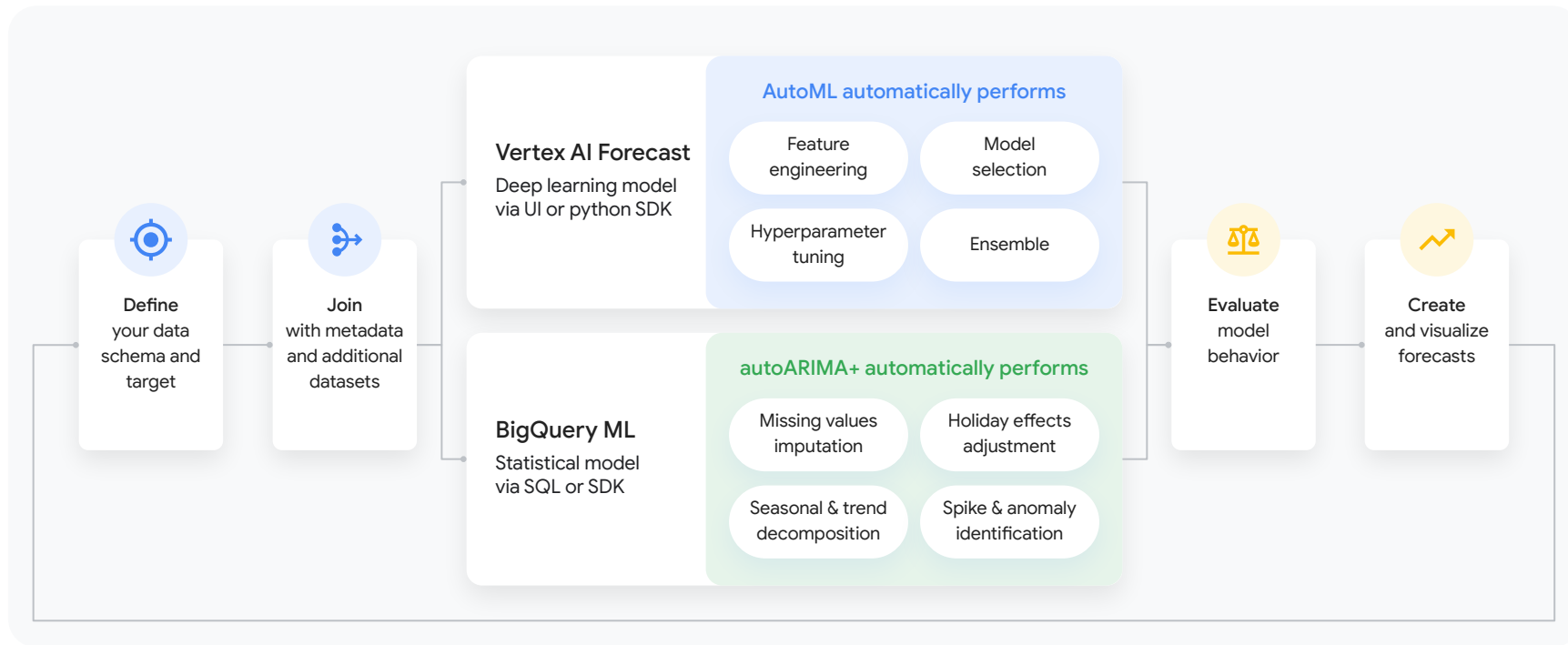
# Vertex AI Forecast: Workflow & Demo

Google Cloud

# Forecasting Workflows on Google

Proprietary + Confidential

Build high quality, scalable forecasting solutions with deep learning and statistical models from Google.



# Easy to Get Started

## Data

Load Dataset

Supported sources:

- CSV
- BigQuery Table
- BigQuery View

Define Columns

Required inputs:

- Target
- Time Series Identifier
- Time

Define Features

Required inputs:

- Type Transformations
- Time Dependence
- Availability at Prediction Time

Generate Statistics

Run Validations

## Model

Setup Model

Required inputs:

- Forecast Horizon
- Context Window
- Data Split
- Optimization Objective

Train Model

Required inputs:

- Training Budget

Automated Feature Engineering

Auto Model Architecture Search

Evaluate Model

## Prediction

Create Request

Required inputs

- Time
- *Feature Values*

Create Forecast

Output to

- CSV
- BigQuery

# Demo: Create Dataset

Google Cloud Platform
Cloud AutoML Tables
Search products and resources

AI Platform (Unified)
Create dataset

Dashboard
Datasets
Labeling tasks
Notebooks
Pipelines
Training
Models
Endpoints
Batch predictions

Dataset name \*  
my\_dataset  
Can use up to 128 characters.

Select a data type and objective  
First select the type of data your dataset will contain. Then select an objective, which is the outcome that you want to predict.

IMAGE
TABULAR
TEXT
VIDEO

☐ Regression/classification

Predict a target column's value. Supports tables with hundreds of columns and millions of rows.

☒ Forecasting

Predict the likelihood of certain events or demand.

Region  
us-central1 (Iowa)

CREATE CANCEL

my\_dataset
SOURCE ANALYZE

Add data to your dataset  
Before you begin, read the [data guide](#) to learn how to prepare your data. Then choose a data source.

Select a data source

- CSV file: Can be uploaded from your computer or on Cloud Storage. [Learn more](#)
- BigQuery: Select a table or view from BigQuery. [Learn more](#)

☒ Upload CSV files from your computer
☐ Select CSV files from Cloud Storage
☐ Select a table or view from BigQuery

Upload CSV files from your computer  
Add up to 500 CSV files per upload. The files will be stored in a new Cloud Storage bucket ([charges apply](#)). Data from multiple files will be referenced as one dataset.

SELECT FILES

# Demo: Training

Cloud Platform

(Unified)

SOU

Data

Cre

Data

Data

test.

Field

confir

count

count

date

## Train new model

- ☒ Choose training method
- 2 Define your model**
- 3 Choose training options
- 4 Compute and pricing

**START TRAINING** CANCEL

Model name \*

my\_Model

?

Target column \*

confirmed\_cases

▼

Time series identifier column \*

county\_fips\_code

▼

?

Time column \*

date

▼

?

### Time series regularity

☒ Regular time frequency ☐ Irregular time frequency

Observations in the time series occur at equal time periods (for example, 1 hour between each observation)

### Forecasting configuration

Data frequency \*

Days

▼

The temporal distance between observations in your time series. The selection should match the time frequency of your input time series data. Sets the time frequency of the historical window and forecast horizon.

# Demo: Training

Cloud Platform

(Unified)

SOU

Data

Cre

Data

Data

test.

Field

confi

count

count

date

## Train new model

- ✓ Choose training method
- 2 Define your model**
- 3 Choose training options
- 4 Compute and pricing

**START TRAINING** CANCEL

## Forecasting configuration

Data frequency \*

Days

The temporal distance between observations in your time series. The selection should match the time frequency of your input time series data. Sets the time frequency of the historical window and forecast horizon.

Historical window

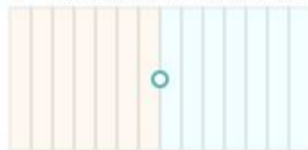
7

Forecast horizon \*

7

The historical window is the period of previous observations to be used for model training. The longer the window, the more historical data will be used. The forecast horizon is the length of time into the future for which predictions will be generated.

Historical context window (7 days) Forecast horizon (7 days)



Forecast starts right after the Historical context window

☐ Export test dataset to BigQuery



# Demo: Training

Cloud Platform

m (Unified)

SOU

Data

Cre

Data

Data

test

Field

confi

count

count

date

death

state

## Train new model

- ✓ Choose training method
- ✓ Define your model
- 3 Choose training options**
- 4 Compute and pricing

**START TRAINING** CANCEL

GENERATING STATISTICS...

Enter property name or value

	Field Name ↑	Transformation	Time dependence ?	Missing % (count) ?	Distinct values ?	
<input type="checkbox"/>	county	Auto ▾	Time dependent ▾ Unavailable at prediction time (historical data only)	-	-	⊖
<input type="checkbox"/>	county	Available at prediction time Unavailable at prediction time (historical data only)	Time independent Time dependent ▸	-	-	⊖
<input type="checkbox"/>	date Time column	Auto ▾	Time dependent ▾ Unavailable at prediction time (historical data only)	-	-	⊖
<input type="checkbox"/>	deaths	Auto ▾	Time dependent ▾ Unavailable at prediction time (historical data only)	-	-	⊖
<input type="checkbox"/>	state_name	Auto ▾	Time dependent ▾ Unavailable at prediction time (historical data only)	-	-	⊖

ADVANCED OPTIONS

CONTINUE

# Demo: Training

**Train new model**

- ✓ Choose training method
- ✓ Define your model
- ✓ Choose training options
- 4 Compute and pricing**

**START TRAINING** **CANCEL**

Enter the **maximum** number of node hours you want to spend training your model.

You can train for as little as 1 node hour. You may also be eligible to train with free node hours. [Pricing guide](#)

Budget \*  Maximum node hours ?

**Estimated completion date:** Mar 5, 2021 8 PM GMT-8

# Demo: Evaluate Model

Google Cloud Platform

Cloud AutoML Tables

Search products and resources



AI Platform (Unified)

[← demand\\_forecast\\_20212216593](#)[VIEW DATASET](#)

Dashboard

Datasets

Labeling tasks

Notebooks

Pipelines

Training

Models

Endpoints

EVALUATE

BATCH PREDICTIONS

MODEL PROPERTIES

Target column  
unit\_sales  
numeric

MAE ?  
5.915

MAPE  
 $\infty$

RMSE ?  
14.454

RMSLE ?  
-

R^2 ?  
0.002

# Demo: Create Prediction

Google Cloud Platform

Cloud AutoML Tables

Search

AI Platform (Unified)

← demand\_forecast\_2021221659

Dashboard

Datasets

Labeling tasks

Notebooks

Pipelines

Training

Models

Endpoints

Batch predictions

EVALUATE

BATCH PREDICTIONS

**Batch predictions**

Batch prediction intakes a group of prediction request location. Use batch prediction when you don't require process accumulated data with a single request. Batch [models](#) and [custom-trained models](#).

CREATE BATCH PREDICTION

Filter batch predictions

Batch prediction

No rows to display

New batch prediction

Batch prediction name \*

Model name  
demand\_forecast\_20212216593

Select source

☒ BigQuery table

☐ File on Cloud Storage (CSV)

Use a table from BigQuery

Google Cloud project ID \*

BigQuery dataset ID \*

BigQuery table or view ID \*

Select a Cloud Storage location

Prediction results will be stored in the selected Cloud Storage bucket

Output format  
BigQuery

Google Cloud project ID \*

CREATE

CANCEL



# Questions

Google Cloud