GILEAD
Creating Possible

# GnA DNA Intern Project
*Forecasting SPARC Request Volumes*

# Meeting Agenda

## Things I Learned

Regression, SQL, Starburst, Time Series Models

## Project Overview

Forecasting SPARC Service Request Volume

## Data Analysis and Model Selection

EDA (Exploratory Data Analysis), Regression, Fourier series modeling & final model selection

## Model Tuning

Hyperparameter Tuning

## Results

Accuracy and performance on historicals and 90-day forecast windows

# Things I learned

*Regression, SQL, Starburst, Time Series Models*

- *Time Series Models*
  - *Trend, Seasonality, Confidence Intervals*
- *Data Analysis*
  - *Periodogram, Seasonality, Trends, Volumes*
- *Querying Data via SQL*
- *Data Cleaning*
- *Hyperparameters and Impacts Thereof,*
- *Evaluating Model performance*
  - *MAPE, Percent Error, Accuracy across windows (90-days, etc.)*

# High Level Project Overview

## Impetus
- A value proposition to demonstrate the potential benefits of a demand forecast for SPARC request volumes
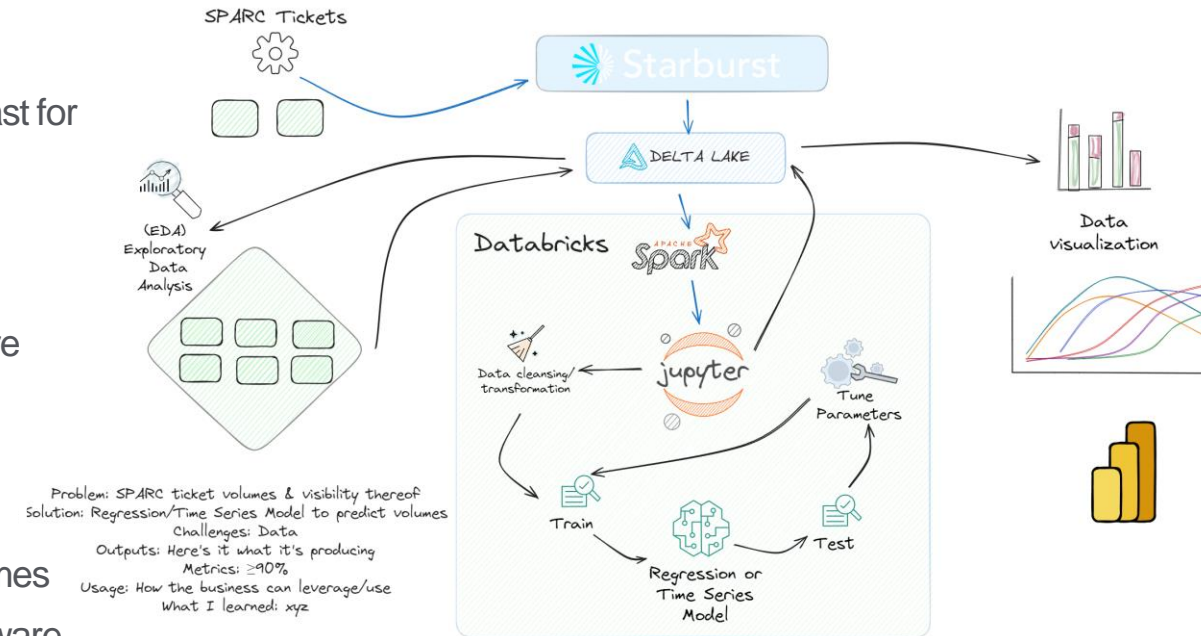
## Problem
- SPARC ticket volumes can fluctuate based on system outages, new software offerings, organization changes and hiring patterns

## Proposed Solution
- Build a demand forecast model to help inform the business of potential volumes
- Provide a solution that can incorporate external variables such as user hardware updates, software launches, new country rollouts, shifts in work patterns (*remote vs office*), etc.

## Potential Benefits
- Headcount projections (*required personnel to handle volume*) and business cost calculations thereof
- Accommodate what-if scenarios by correlating ticket inflow with external factors and applying them to projections
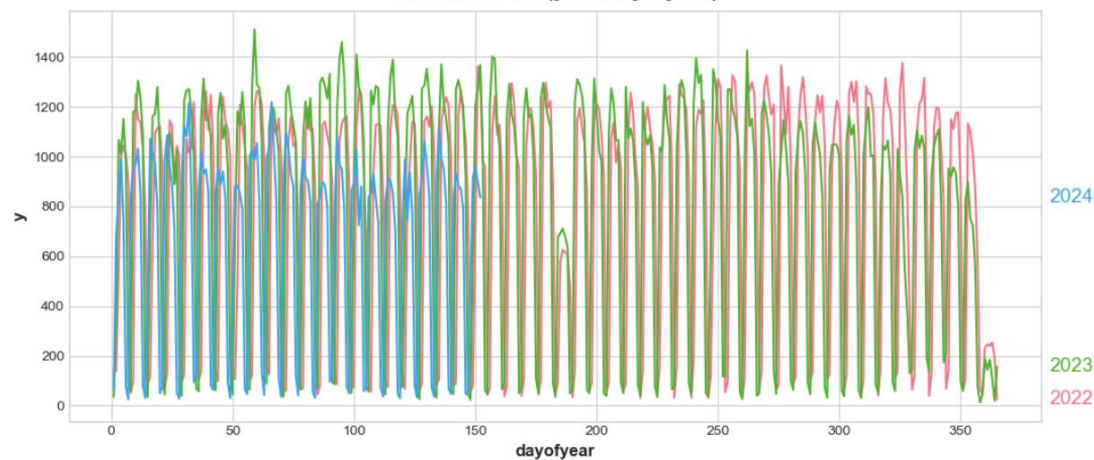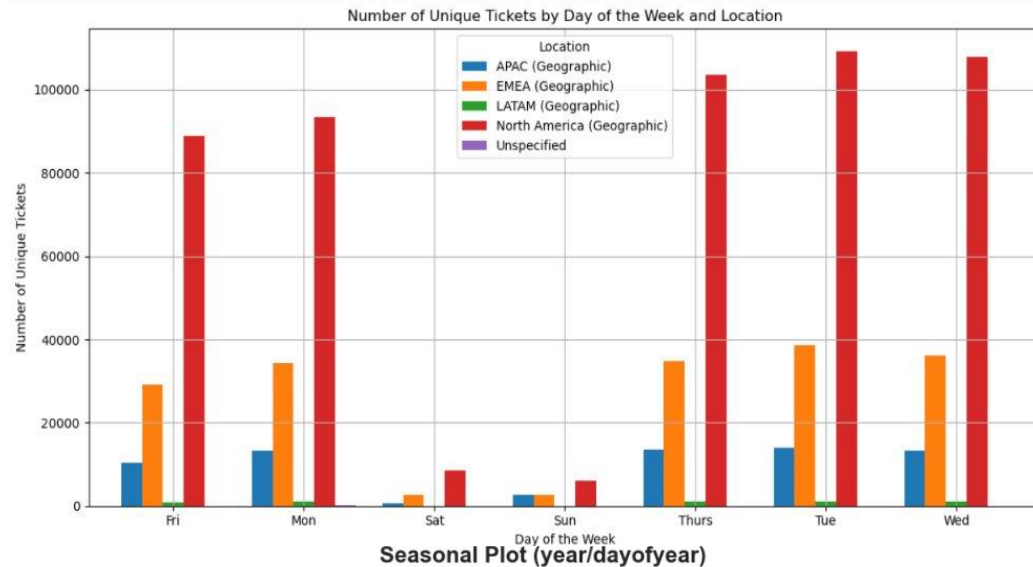


SPARC Tickets

Starburst

DELTA LAKE

(EDA) Exploratory Data Analysis

Databricks Spark

Data cleansing/ transformation

jupyter

Tune Parameters

Train

Test

Regression or Time Series Model

Data visualization

Problem: SPARC ticket volumes & visibility thereof
Solution: Regression/Time Series Model to predict volumes
Challenges: Data
Outputs: Here's it what it's producing
Metrics: ≥90%
Usage: How the business can leverage/use
What I learned: xyz

# Exploratory Data Analysis

- *Pulled data down from Starburst via SQL*
- *Removed duplicates*
- *Broke down data into regions*
- *Analyzed Gilead vs Kite*
- *Cleaned data by removing spikes/ outliers*
- *The top 5 countries by volume represented 86% of the overall volume*
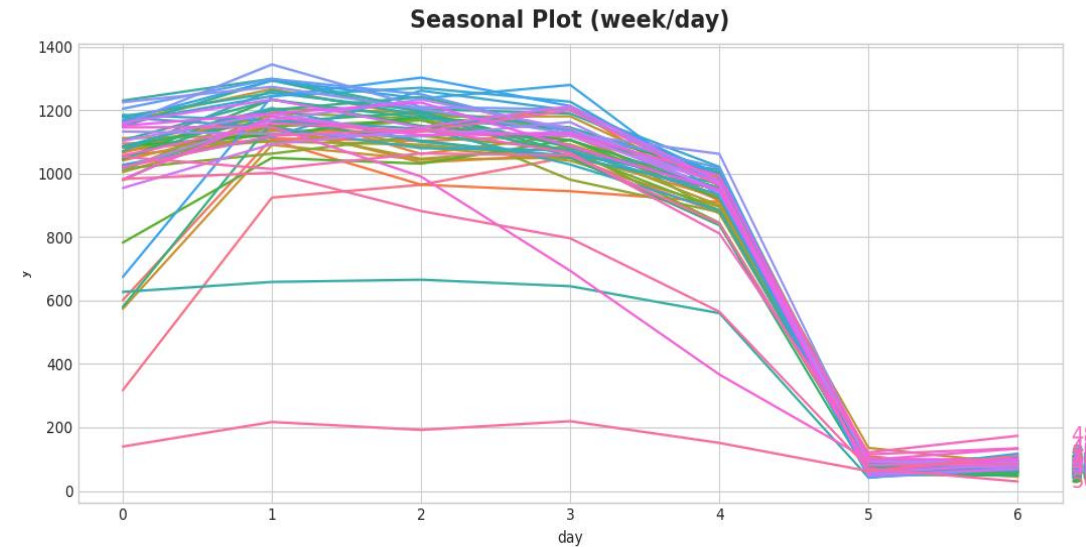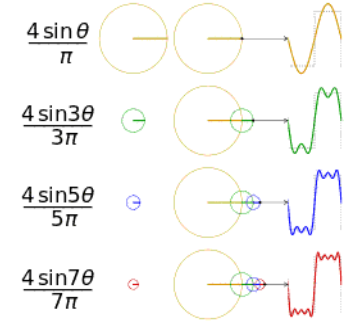
# Exploratory Data Analysis, cont.
## *Seasonality (Month, Week, Day)*



- Examined seasonality across all significant intervals
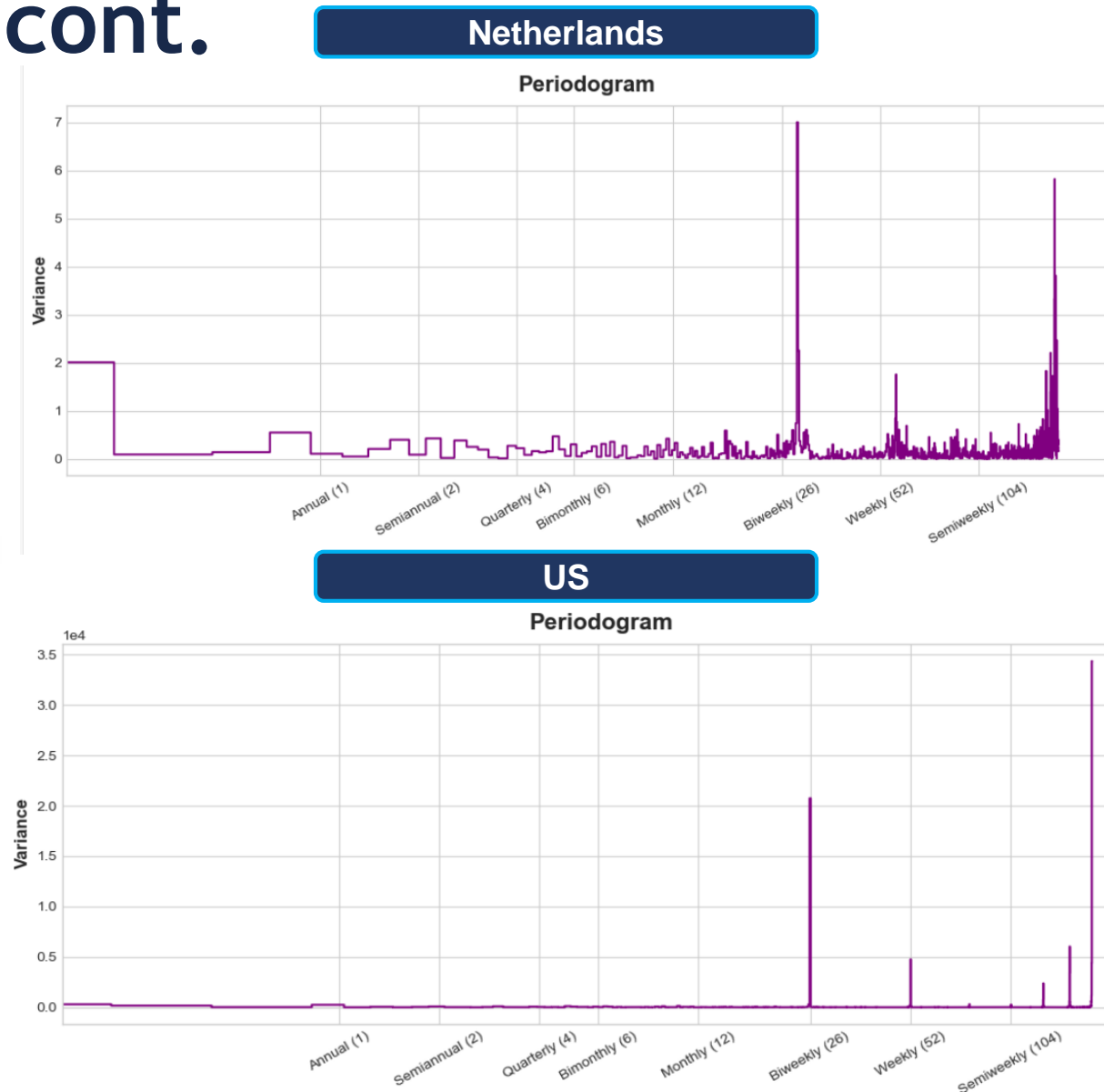  - Quarters, Months, Weeks
- As well as across the top 5 countries and 4 regions

$$\frac{4\sin\theta}{\pi}$$

$$\frac{4\sin 3\theta}{3\pi}$$

$$\frac{4\sin 5\theta}{5\pi}$$

$$\frac{4\sin 7\theta}{7\pi}$$

# Exploratory Data Analysis, cont.
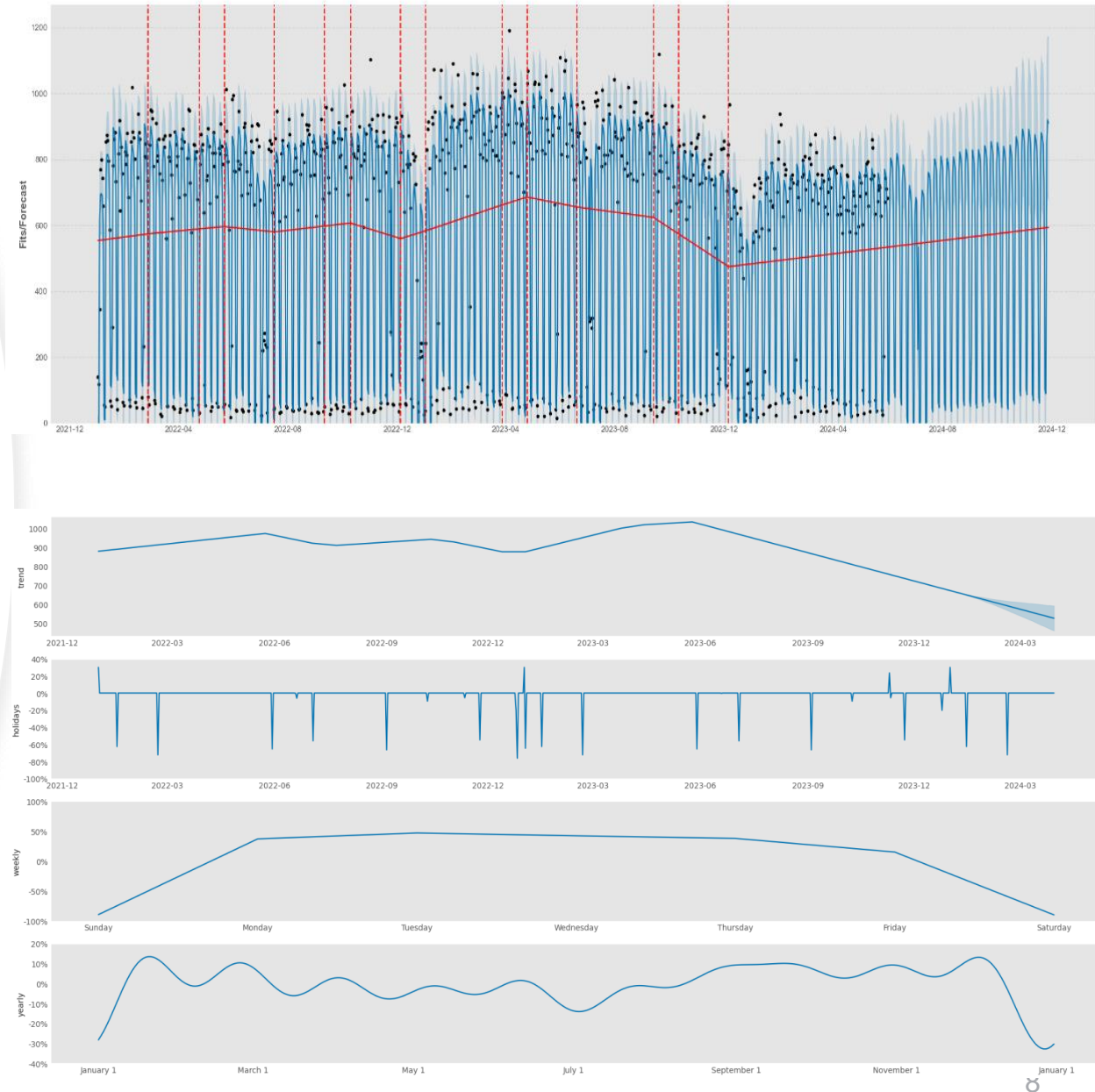## *Periodicity*

- *The variance in the data by country/region indicated a model would be required that could accommodate more than simple trend and growth factors for SPARC tickets*

- *Examining the unique seasonality patterns within the data's subsets helped identify that segregating the data by the top 5 countries, and remaining countries by region would compensate for the unique seasonality patterns within regions that would impact demand*
  - *Some areas had unique periodicity, i.e Netherlands*
  - *Netherlands and US pictured*



Netherlands
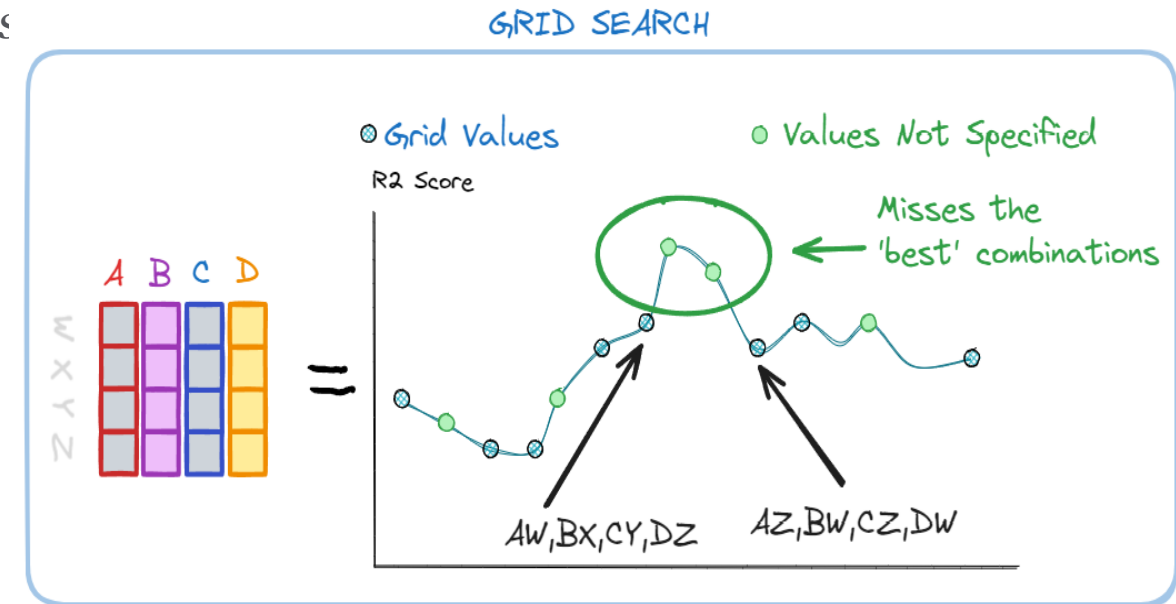


US

# Model Selection

- *Tested 9 Different Algorithms*
- *All model testing leveraged Fourier weights for seasonality (month, week and day of year) informed by the EDA using a deterministic process*
  - *AutoARIMA*
  - *Croston Classic*
  - *Dynamic Optimized Theta*
  - *Generalized Additive (GAM)*
  - *Historic Average*
  - *Holt Winters*
  - *Linear Regression*
  - *Random Forest Regression*
  - *Seasonal Naïve*
- *Highest Performance (MAPE) was with:*
  - *GAM model (prophet model from Meta)*
  - *Benefits:*
    - Robust seasonality
    - Accommodates exogeneous regressors (holidays, etc.)
    - Robust trend weights with changepoint inflection markers
    - Numerous hyper-parameters for tuning



GAM (*prophet*)
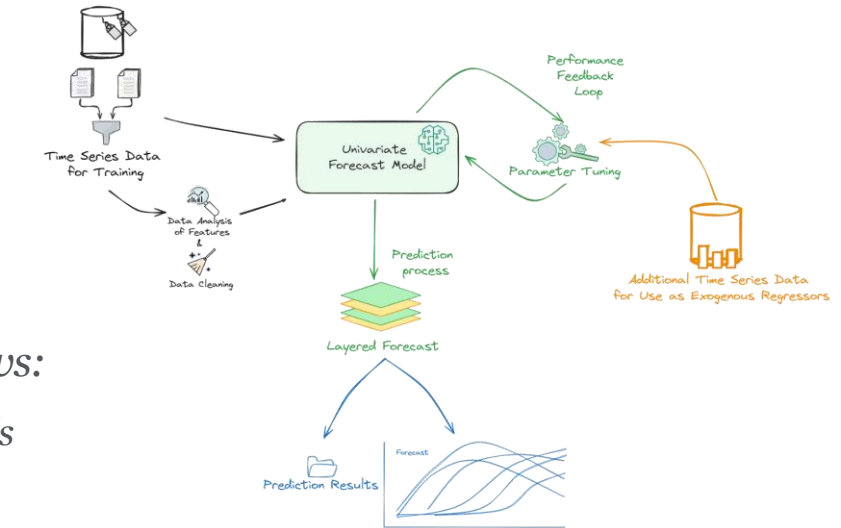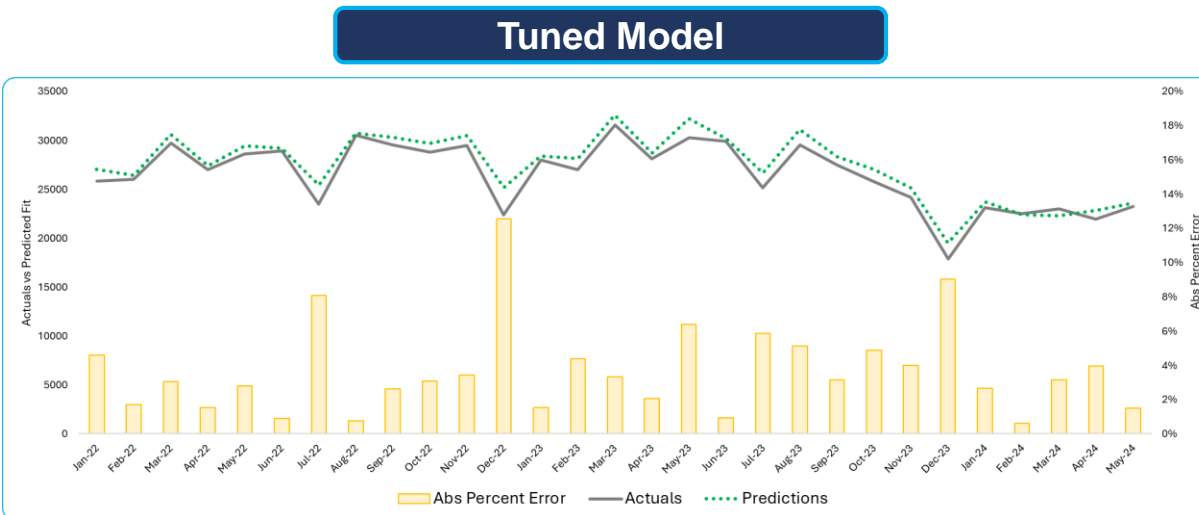
# Hyperparameter Optimization

- *Method: Bayesian Optimization*
    - *Faster model tuning as well as potentially improved accuracy for the model compared to traditional methods (e.g., random or grid searching)*
    - *Treats tuning like a regression problem by iterating through value combinations*
        - *Each run uses values incrementally close to the best previous or far removed to explore the range of values*
        - *Concentrates on trying combinations of 'good' values across the individual parameters*
    - *Other methods are limited by either:*
        - *A finite number of suggested values we enter – and the more values we suggest, the longer the run time*
        - *Random selection of values (i.e., not concentrating on combinations near previous good values)*
    - *Allows us to re-train the model quickly while still maximizing accuracy*

GRID SEARCH

Grid Values    Values Not Specified

R2 Score

Misses the 'best' combinations

A B C D
W X Y Z

AW,BX,CY,DZ    AZ,BW,CZ,DW

# Results and Evaluating Model Performance

*The layered (9 variable) forecast resulted in performance of:*

- *Tuned Model:*
  - *Fit w/ MAPE of 3.7%*
    - *St Dev of 2.7%*
- *Forecast results on a 90-day Horizon measured across chained windows:*
  - *E.g. when predicting a month's volume a quarter ahead from historicals (actuals through Dec, predict Mar)*
  - *Mean accuracy 92.5%, St Dev 3.2%*



**Chained Window Forecasts**

**Tuned Model**



Confidential - Internal Use Only