# Outline

1. **Metadata Analysis**

2. **Time Series Forecasting**

3. **Model Comparison**

INSERT DESIGNATOR, IF NEEDED

# Metadata Analysis

# Anatomy of a Metric

# Ex: Docker Latency

# Metric: kubelet docker operation latency microseconds



Instance Number (y-axis)

Timestamp (x-axis)

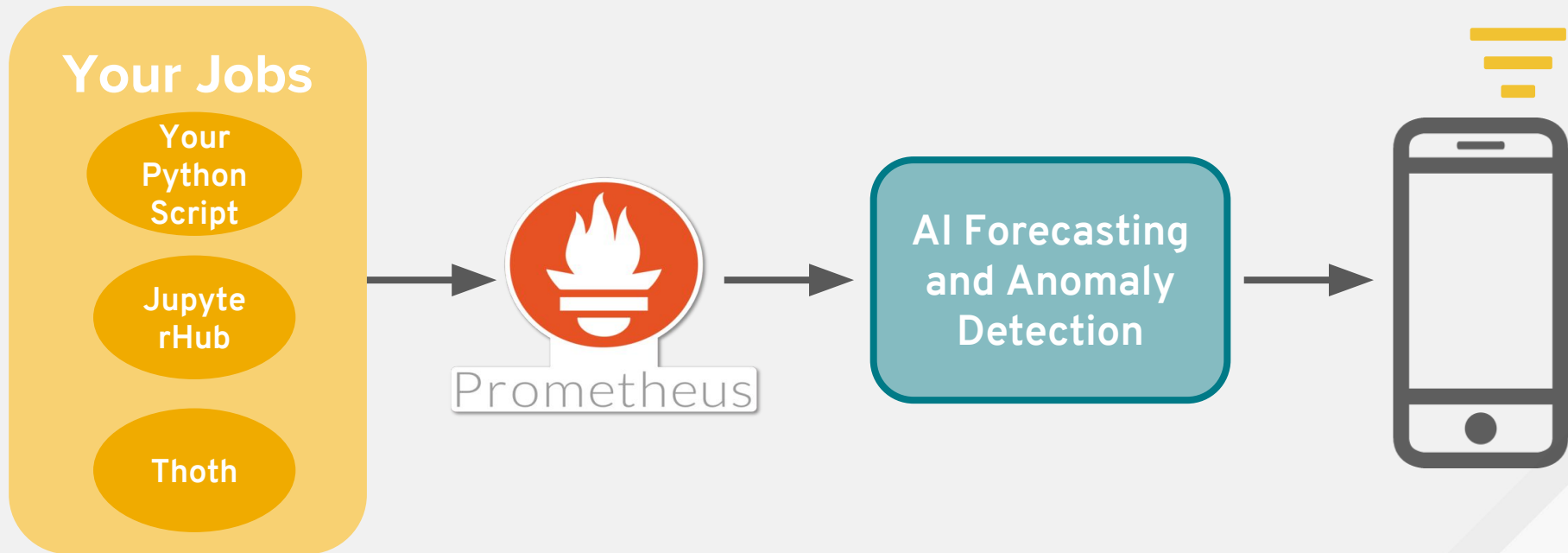| 0 | ip-172-31-78-254.us-east-2.compute.internal | 4 | ip-172-31-76-218.us-east-2.compute.internal |
| 1 | ip-172-31-73-251.us-east-2.compute.internal | 5 | ip-172-31-75-193.us-east-2.compute.internal |
| 2 | ip-172-31-65-74.us-east-2.compute.internal | 6 | ip-172-31-71-195.us-east-2.compute.internal |
| 3 | ip-172-31-74-247.us-east-2.compute.internal | 7 | ip-172-31-69-53.us-east-2.compute.internal |

redhat.

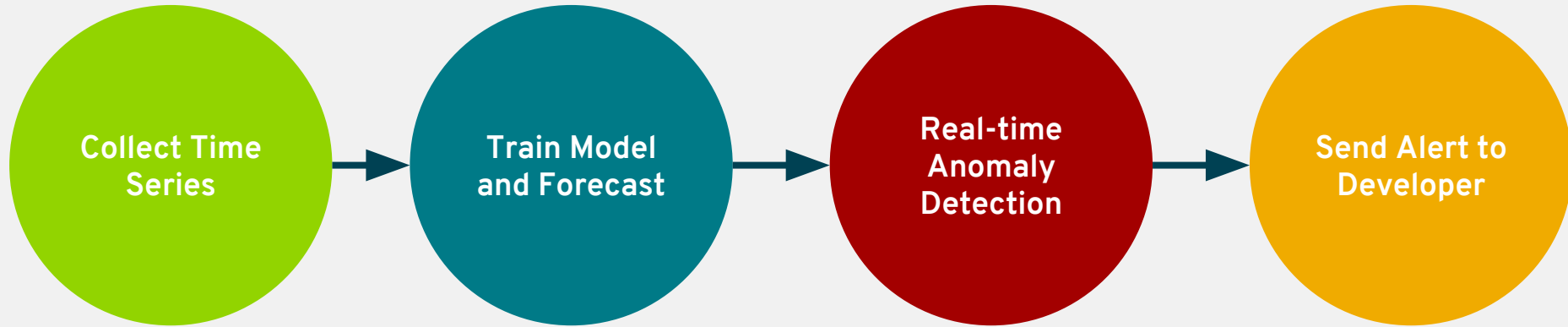# T-SNE Embedding of Metric Metadata



Natasha Frumkin

# Time Series Forecasting and Anomaly Detection

# Goal: to provide automatic alerting to developers when there are anomalies in Prometheus data



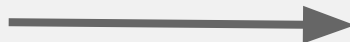Natasha Frumkin

# The Data Transfer Pipeline



Collect Time Series → Train Model and Forecast → Real-time Anomaly Detection → Send Alert to Developer

Natasha Frumkin

redhat.

# Data Processing

**Bz2 File**
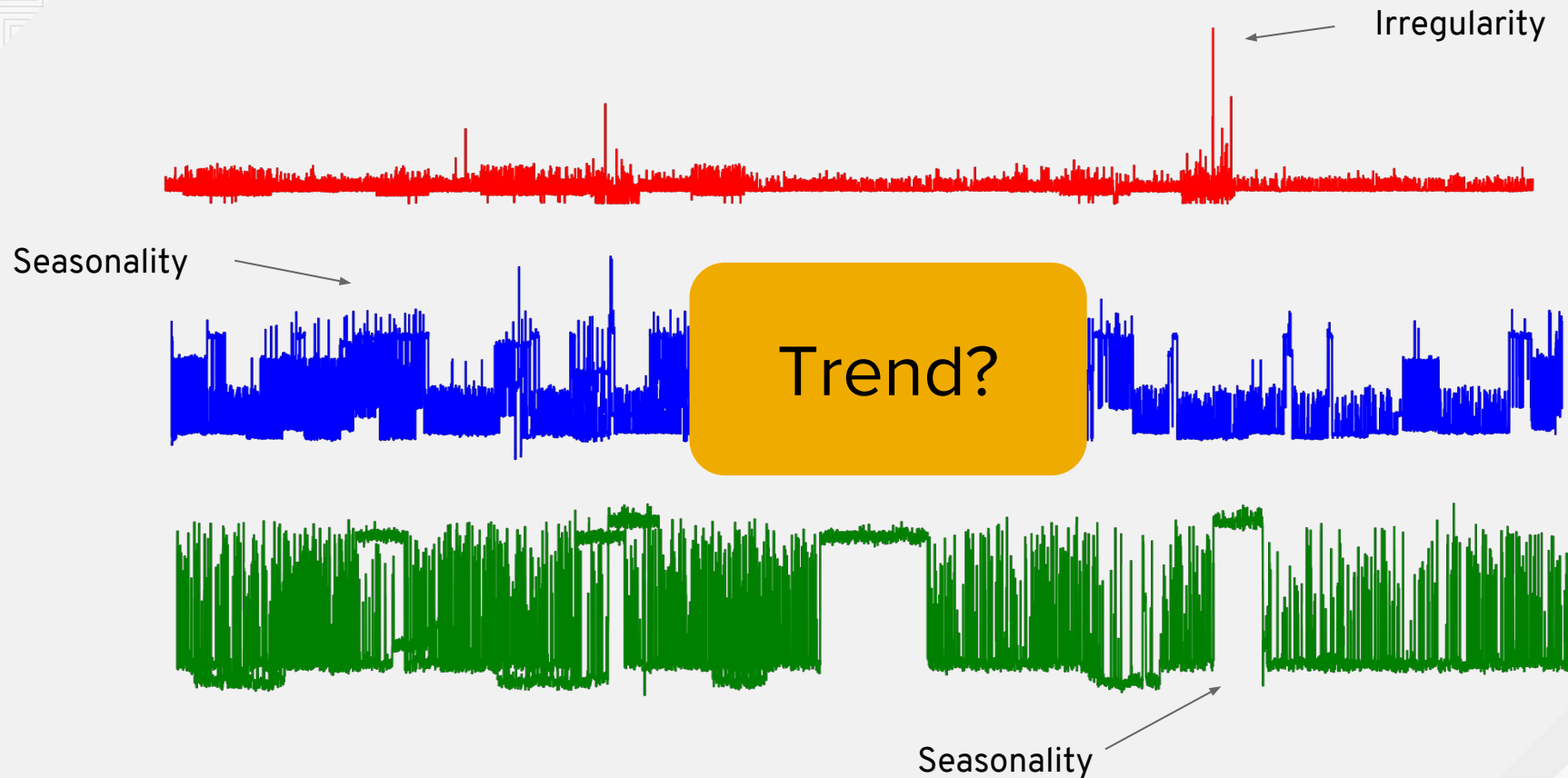Text file of Json packets

**Pickle File**
Dict of DataFrames
Key = str(metadata)
Value = Pandas DataFrame

{"metric": {"__name__": "kubelet_docker_operations_latency_microseconds",
"beta_kubernetes_io_arch": "amd64", "beta_kubernetes_io_fluentd_ds_ready":
"true", "beta_kubernetes_io_instance_type": "m4.xlarge", "beta_kubernetes_io_os":
"linux", "clam_controller_enabled": "True",
"failure_domain_beta_kubernetes_io_region": "us-east-2"},
"values": [[1518307199, "12844"], [1518308638, "13212"], [1518310077, "13830"],
[1518311516, "13395"], [1518312955, "16546"], [1518314394, "15174"], [1518315833,
"14455"], [1518317272, "12949"], [1518318711, "13439"], [1518320150, "14386"],
[1518321589, "12447"], [1518323028, "15947"], [1518324467, "14893"], [1518325906,
"14096"], [1518327345, "14735"], [1518328784, "12969"], [1518330223, "14067"],
[1518331662, "16286"], [1518333101, "14008"], [1518334540, "12923"], [1518335979,
"11888"], [1518337418, "12263"], [1518338857, "11751"], [1518340296, "13534"],
[1518341735, "15522"], [1518343174, "14912"], [1518390661, "12235"], [1518392100,
"14209"], [1518393539, "15757"]]}

| Timestamps | Values |
|---|---|
| 03-03-18 12:30:15 | 32193 |
| 03-03-18 12:31:15 | 33210 |
| 03-03-18 12:32:15 | 32184 |

redhat.

# Metric: http request duration



Irregularity

Seasonality

Trend?

Seasonality

redhat.

# Fourier Extrapolation



Training Set

Incoming Data Points
Forecast

Natasha Frumkin

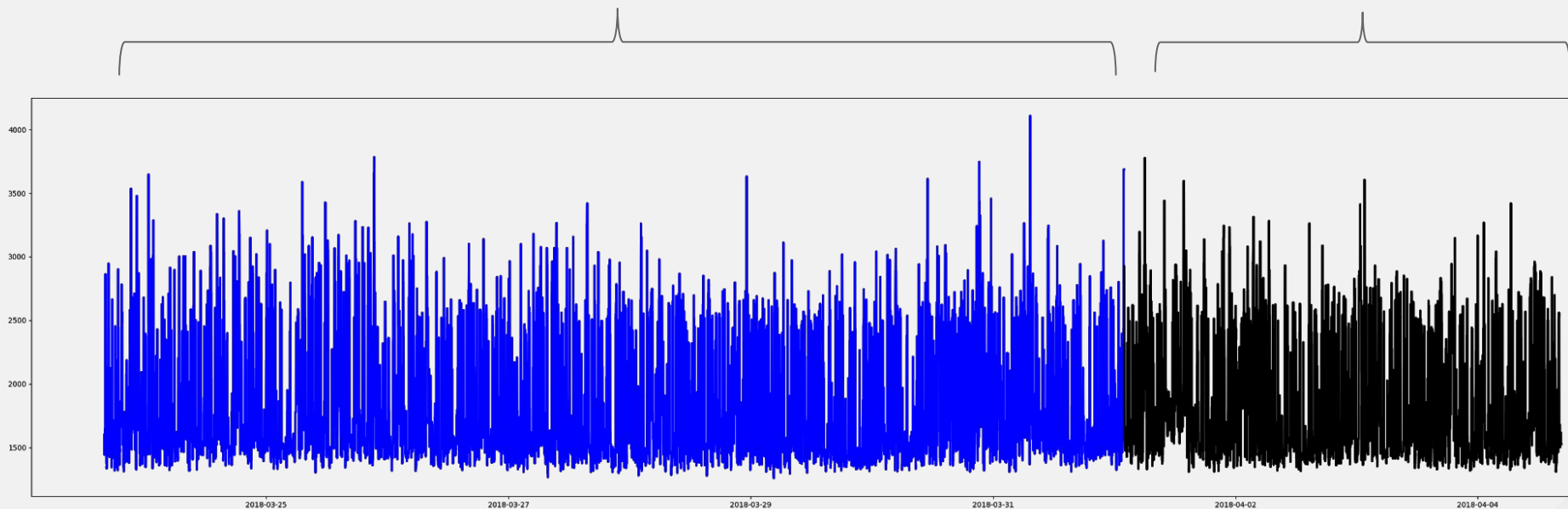# Problem: What is an anomaly?



Natasha Frumkin

# Forecast Comparison

## Prophet vs. Fourier

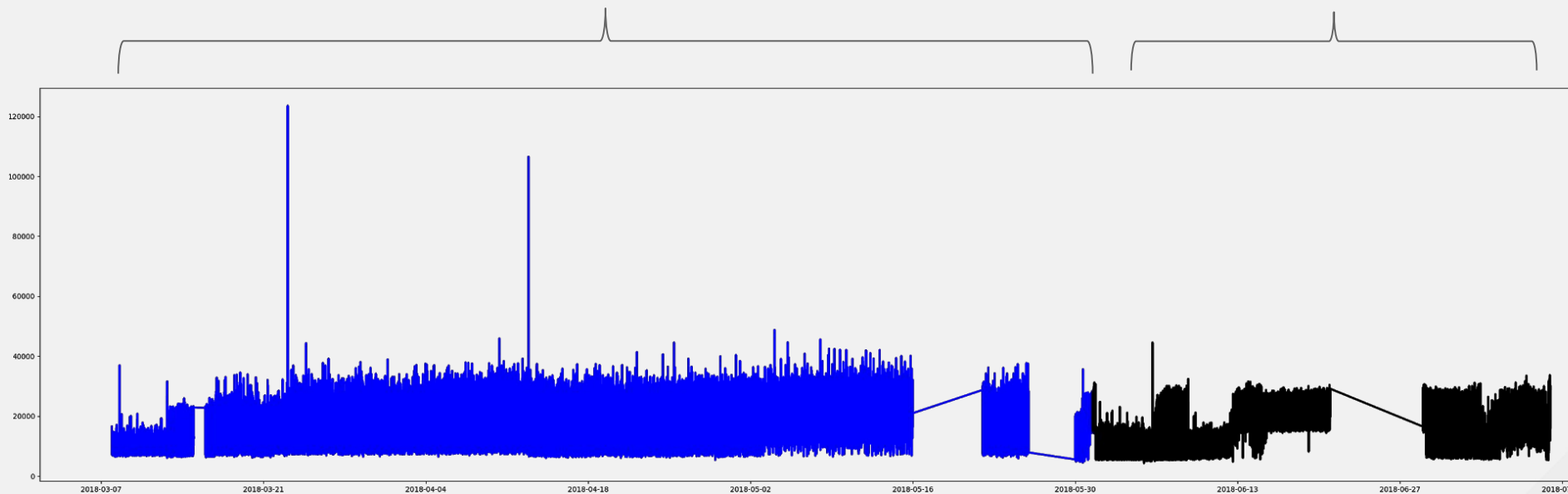# Metric: http request duration



Training Set

Incoming Data Points

Natasha Frumkin

redhat.

Prophet

Fourier

19 Natasha Frumkin

# Metric: http request duration



Training Set

Incoming Data Points

Natasha Frumkin

Prophet

Fourier

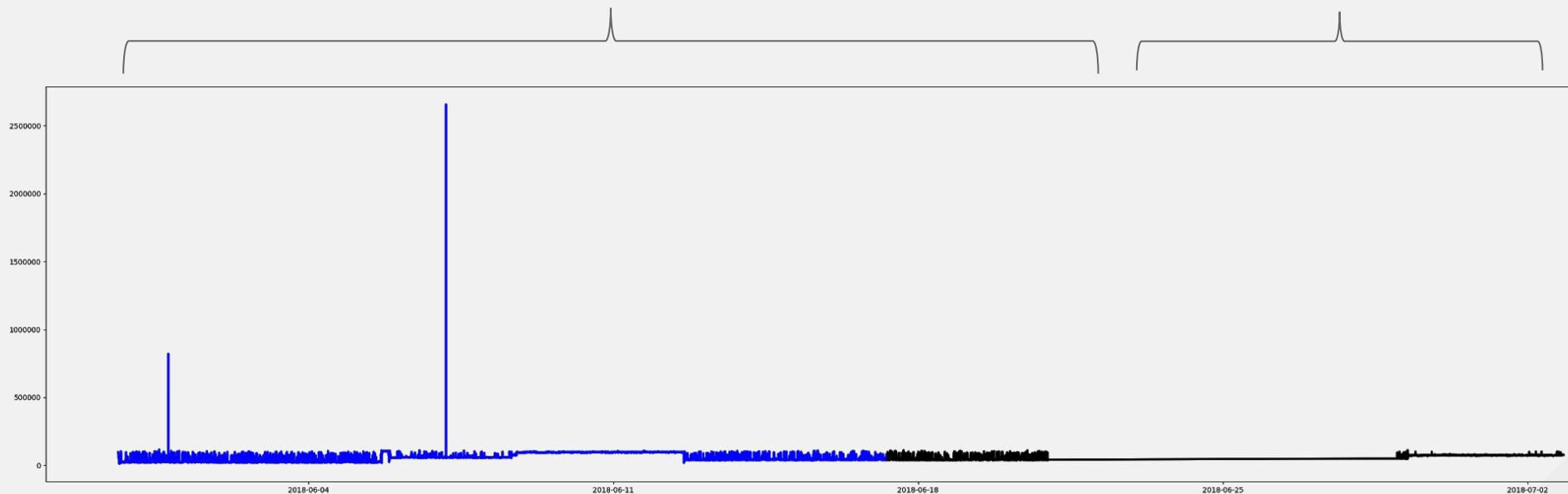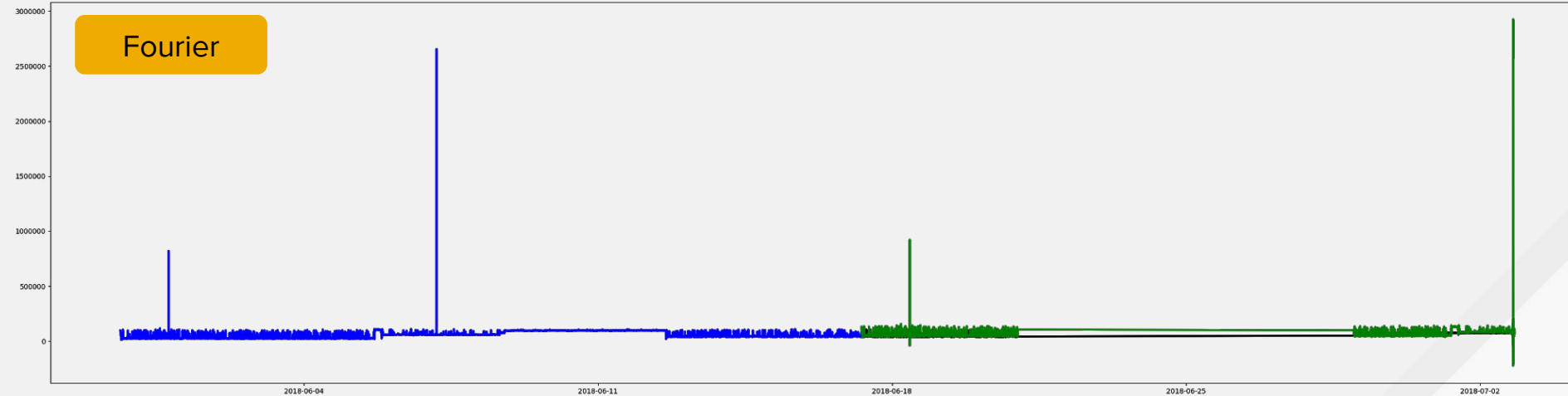# Metric: http request duration



Natasha Frumkin

Prophet

Fourier

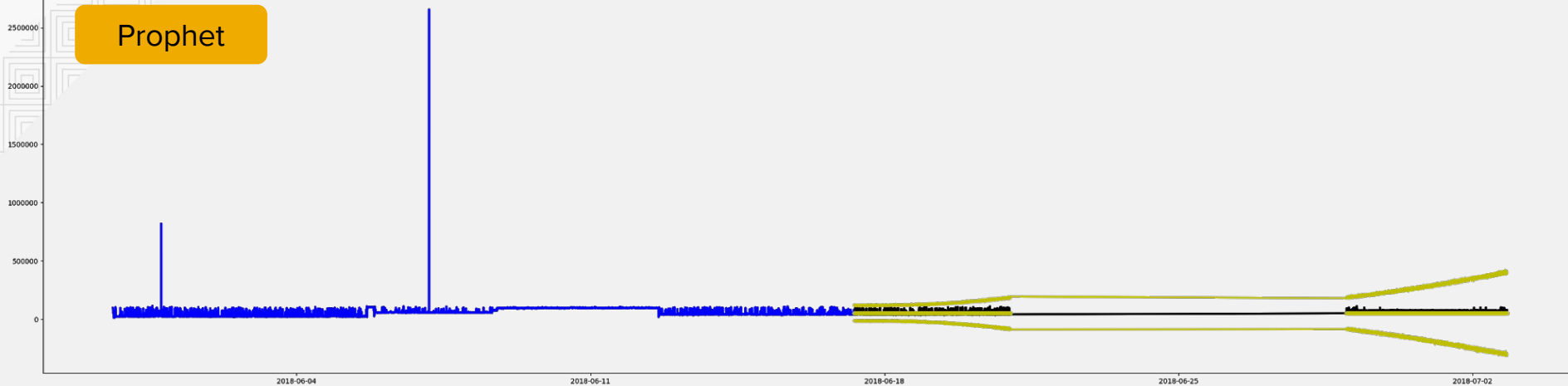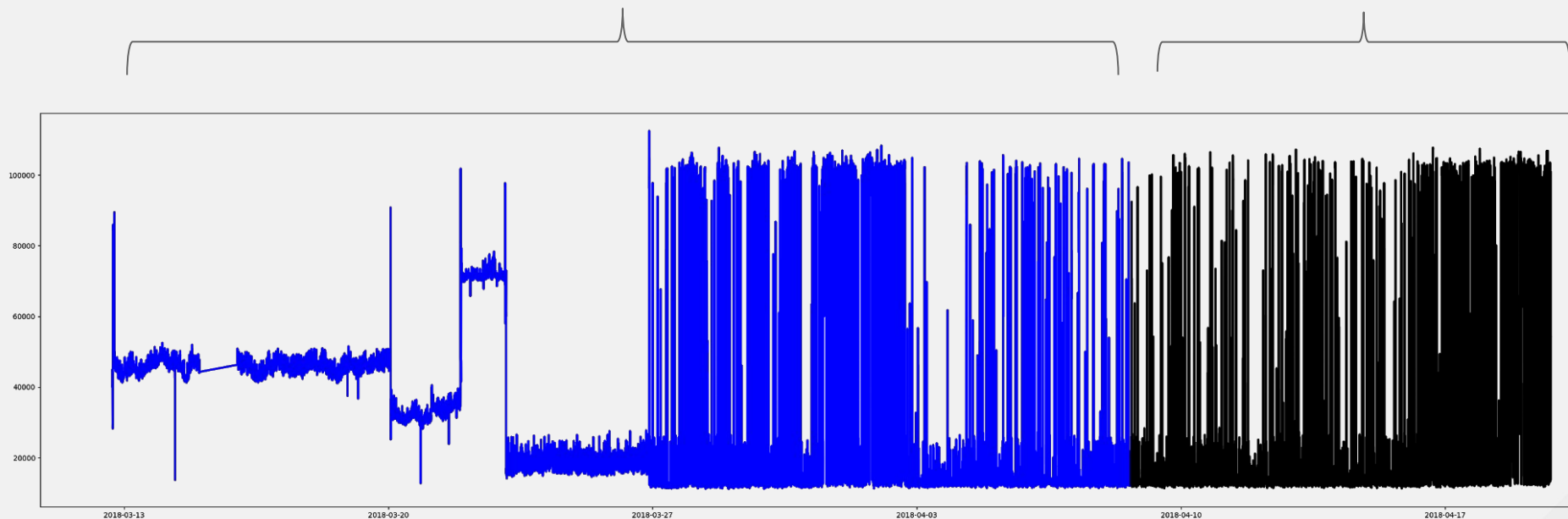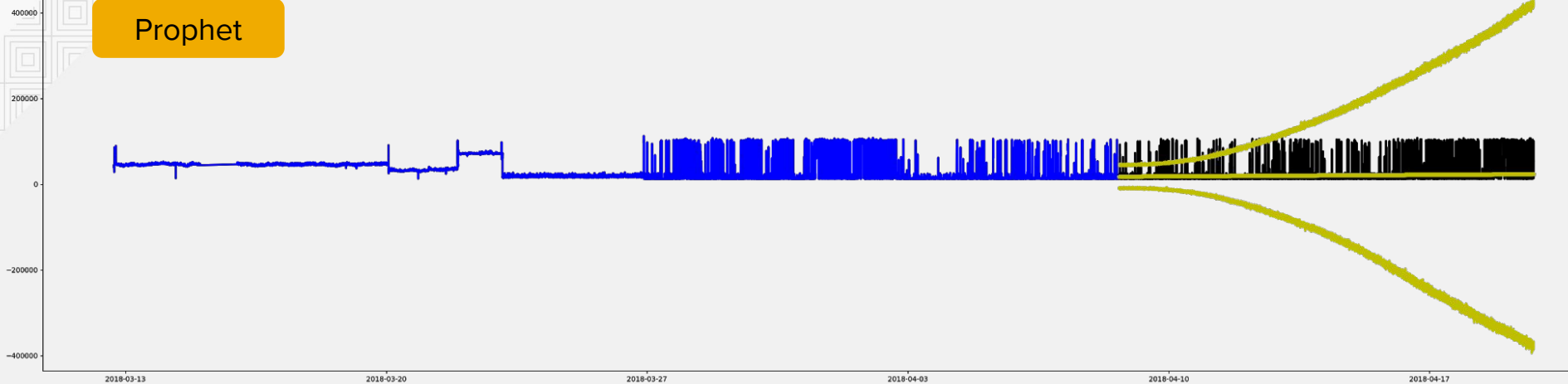# Metric: http request duration



Training Set

Incoming Data Points

Natasha Frumkin

redhat.

Prophet

Fourier

Natasha Frumkin

# Summary of our Techniques

Exponential Smoothing

ARIMA models

Fourier Analysis

Prophet models

RNNs

Thresholding

Gaussian Tail Probability

Accumulators

**training models**

**anomaly detection rules**
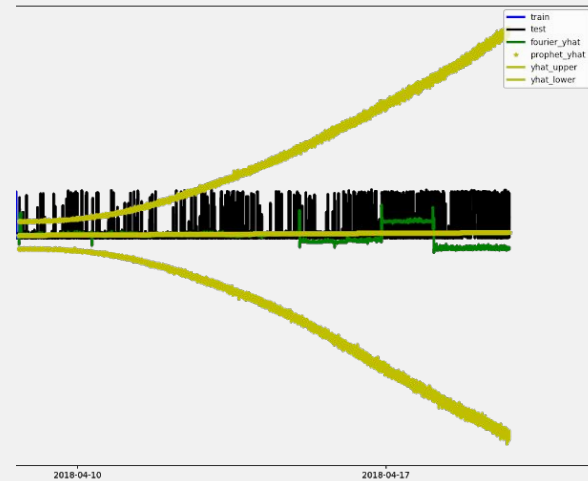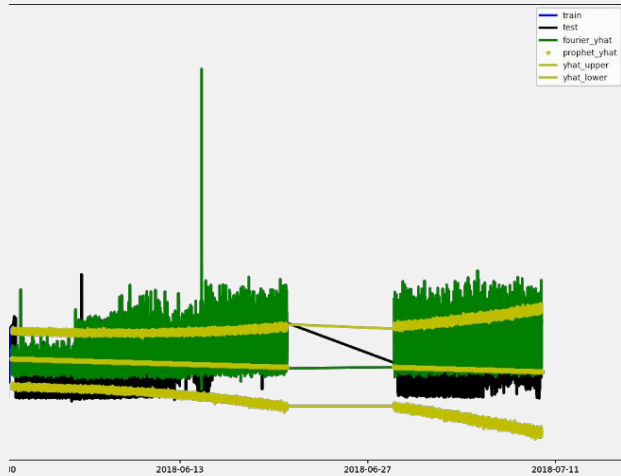
redhat.

# **Main Takeaways**

Metadata configurations are constantly changing.

Prophet vs. Fourier vs. RNNs    Which features do we care about?

Anomaly detection requires finesse.   Need to test parameters.

redhat.

# Next Steps

For which anomalies do we send alerts?     Threshold needed.

Dive deeper into more complex models.     Ensemble methods?

Scalability?     Which time series do we choose to monitor?

Natasha Frumkin

redhat.

# THANK YOU

Notebooks: Gitlab AICOE/jupyter-notebooks
Documentation and Scripts: github.com/nfrumkin/forecast-prometheus

# Challenges with Prometheus Dataset

- **Data comes from multiple sources**
  - Need to explore correct time series filtering
- **Data has holidays and season**
  - Leverage known smoothing and decomposition techniques
- **Wide range of metric types and behavior**
  - Possibly apply different AD techniques for different series
- **Training Data has hidden anomalies and dropouts**
  - Find a way to accurately prepare historical data for training

# The Data Transfer Pipeline

Natasha Frumkin