

BIG DATA ANALYTICS



Final presentation
25.02.2019

TU Berlin

... Agenda

- 1** Use case presentation
- 2** System architecture
- 3** Game plan & Challenges
- 4** Data preprocessing
- 5** Anomaly detection models explanation and comparison
- 6** Prediction models explanation and comparison
- 7** Notification functions
- 8** Recap and Recommandations

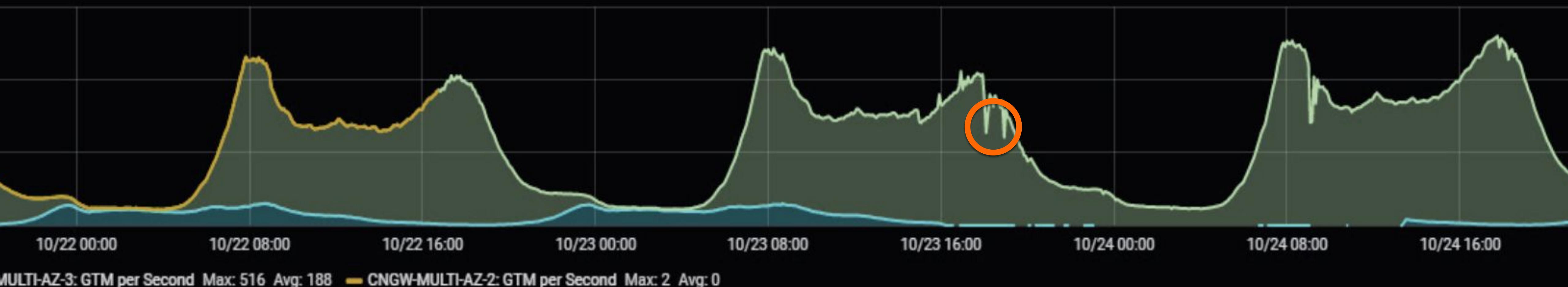
8 Mio.

ConnectedDrive Cars

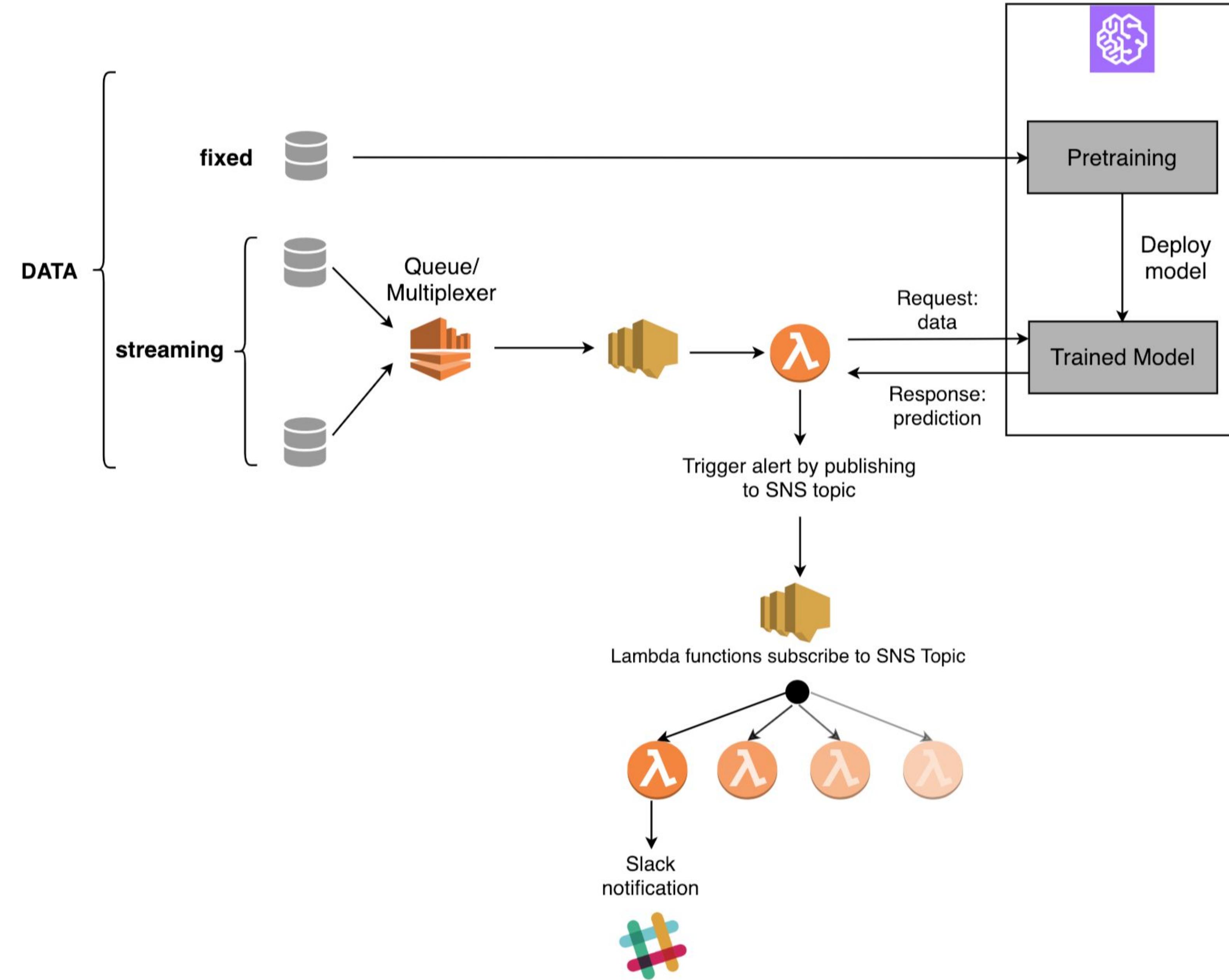
Use Case Presentation

...
.

Anomaly detection on time series data
(requests per second)



System architecture



Deployment

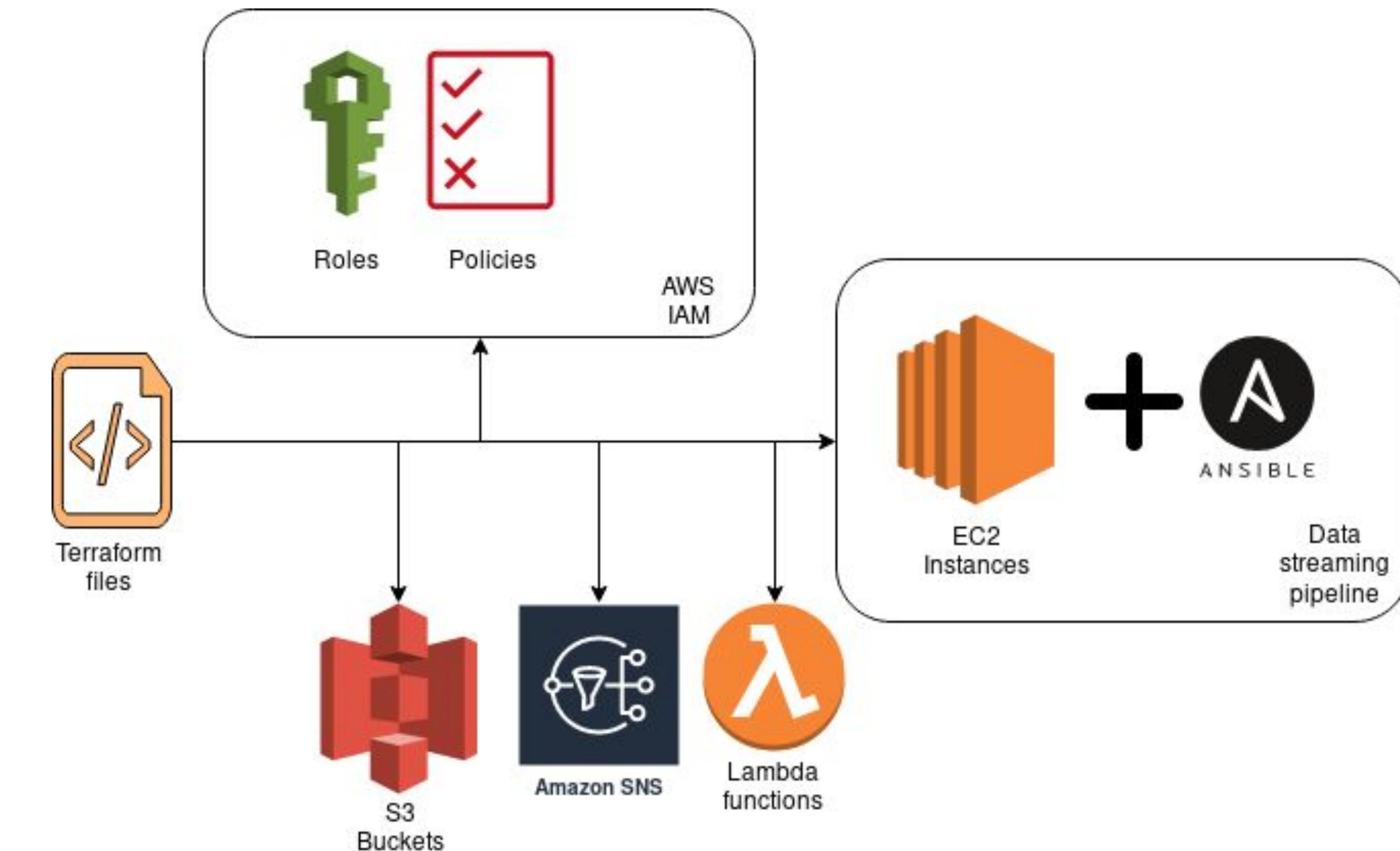
...

Infrastructure managed with
Terraform

Provisioning with Ansible

Goal: get rid of most of the setup
complexity

BUT some features are not
supported yet

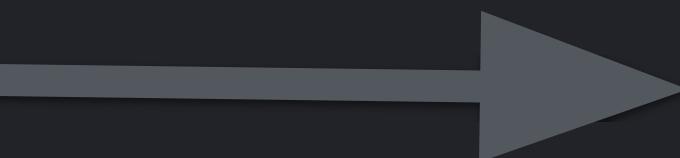


... Game plan

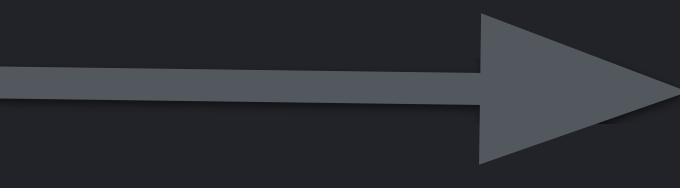


Challenges

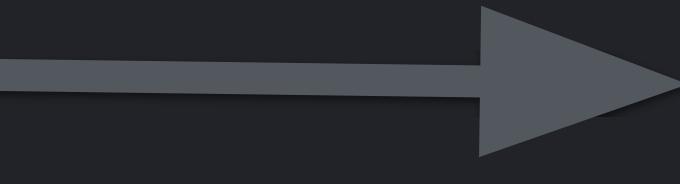
1. Understanding and preprocessing data



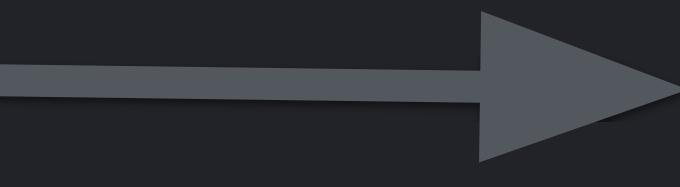
2. Faking a streaming data pipeline



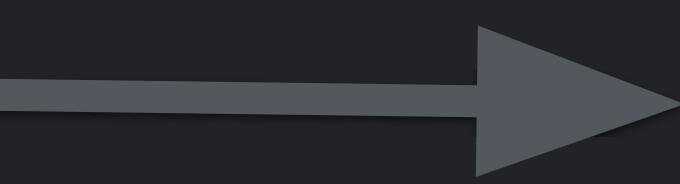
3. Lack of extensive AWS knowledge



4. Exploring Machine Intelligence solutions



5. Small amount of data



Solutions

1. Create 1 min and 5 min granularity buckets

2. Feeding portions of data once every 15 min

3. Organizing knowledge sharing sessions

4. Trying and comparing a couple algorithms

5. Generating new data with our prediction models

•
○
○
○
○

... Data processing

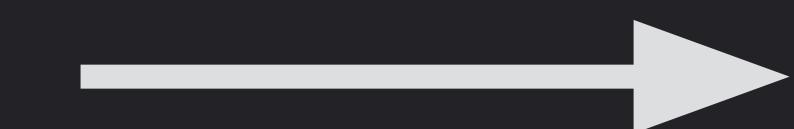
```
version account-id interface-id srcaddr dstaddr srcport dstport protocol  
packets bytes start-time end-time log-status  
2 unknown eth0:sub2 740:fcff:0:0:0:0:0:0 10:24:23 1545404465 - NODATA
```

VPC flowlogs

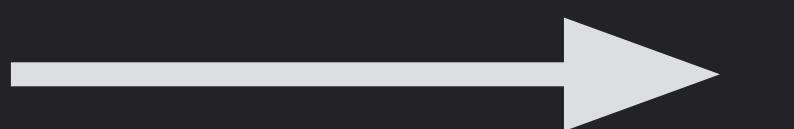
```
version account-id interface-id srcaddr dstaddr srcport dstport protocol  
packets bytes start-time end-time log-status  
2 unknown eth0:sub2 740:fcff:0:0:0:0:0:0 10:24:23 1545404465 - NODATA
```



```
{'RequestID': '8b2cdeb5-260b-11e9-9198-c961bc291f20', 'HTTPStatusCode': 200, 'HTTPHeaders': {'x-amzn-requestid': '8b2cdeb5-260b-11e9-9198-c961bc291f20', 'content-type': 'text/xml', 'content-length': '345', 'date': 'Fri, 01 Feb 2019 10:24:23 GMT'}, 'RetryAttempts': 0}, 'response-code-4xx': {'Label': 'response-code-4xx', 'Datapoints': []}, 'response-code-5xx': {'Label': 'response-code-5xx', 'Datapoints': []}, 'ResponseMetadata': {'RequestId': '8b2f0199-260b-11e9-9198-c961bc291f20', 'HTTPStatusCode': 200, 'HTTPHeaders': {'x-amzn-requestid': '8b2f0199-260b-11e9-9198-c961bc291f20', 'content-type': 'text/xml', 'content-length': '340', 'date': 'Fri, 01 Feb 2019 10:24:23 GMT'}, 'RetryAttempts': 0}, 'response-code-200': {'Label': 'response-code-200', 'Datapoints': []}, 'ResponseMetadata': {'RequestId': '8b303a1b-260b-11e9-9198-c961bc291f20', 'HTTPStatusCode': 200, 'HTTPHeaders': {'x-amzn-requestid': '8b303a1b-260b-11e9-9198-c961bc291f20', 'content-type': 'text/xml', 'content-length': '340', 'date': 'Fri, 01 Feb 2019 10:24:23 GMT'}, 'RetryAttempts': 0}, 'response-code-200': {'Label': 'response-code-200', 'Datapoints': []}, 'ResponseMetadata': {'RequestId': '8b3172c-260b-11e9-9198-c961bc291f20', 'HTTPStatusCode': 200, 'HTTPHeaders': {'x-amzn-requestid': '8b3172c-260b-11e9-9198-c961bc291f20', 'content-type': 'text/xml', 'content-length': '340', 'date': 'Fri, 01 Feb 2019 10:24:23 GMT'}, 'RetryAttempts': 0}, 'response-code-200': {'Label': 'response-code-200', 'Datapoints': []}, 'ResponseMetadata': {'RequestId': '8b32ab1f-260b-11e9-9198-c961bc291f20', 'HTTPStatusCode': 200, 'HTTPHeaders': {'x-amzn-requestid': '8b32ab1f-260b-11e9-9198-c961bc291f20', 'content-type': 'text/xml', 'content-length': '345', 'date': 'Fri, 01 Feb 2019 10:24:23 GMT'}, 'RetryAttempts': 0}, 'response-code-200': {'Label': 'Count', 'Datapoints': []}, 'ResponseMetadata': {'RequestId': '8b33e3a4-260b-11e9-9198-c961bc291f20', 'HTTPStatusCode': 200, 'HTTPHeaders': {'x-amzn-requestid': '8b33e3a4-260b-11e9-9198-c961bc291f20', 'content-type': 'text/xml', 'content-length': '328', 'date': 'Fri, 01 Feb 2019 10:24:23 GMT'}, 'RetryAttempts': 0}, 'api-gw-IntegrationLatency': {'Label': 'IntegrationLatency', 'Datapoints': []}, 'ResponseMetadata': {'RequestId': '8b354337-260b-11e9-9198-c961bc291f20', 'HTTPStatusCode': 200, 'HTTPHeaders': {'x-amzn-requestid': '8b354337-260b-11e9-9198-c961bc291f20', 'content-type': 'text/xml', 'content-length': '341', 'date': 'Fri, 01 Feb 2019 10:24:23 GMT'}, 'RetryAttempts': 0}, 'api-gw-4XXError': {'Label': '4XXError', 'Datapoints': []}, 'ResponseMetadata': {}}
```



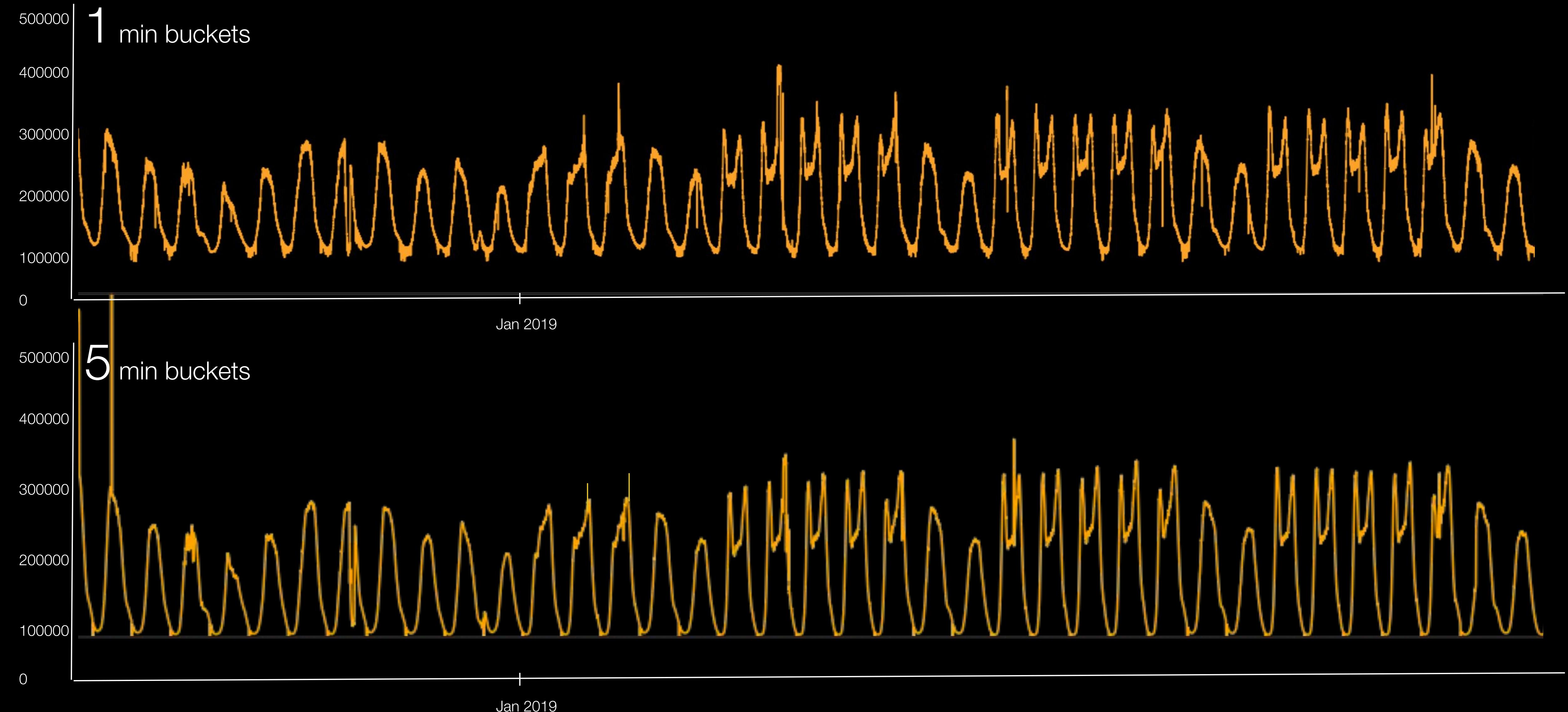
1) Aggregate logs into
1 min buckets



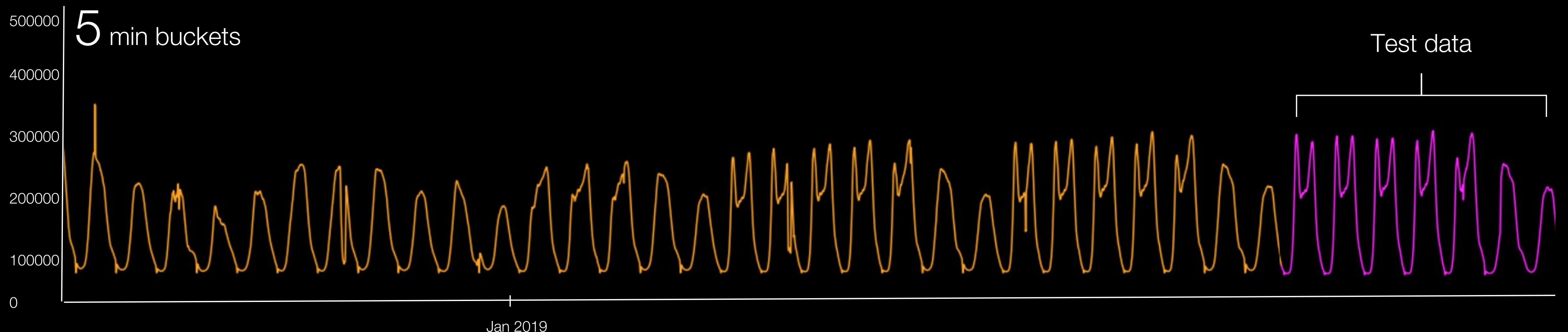
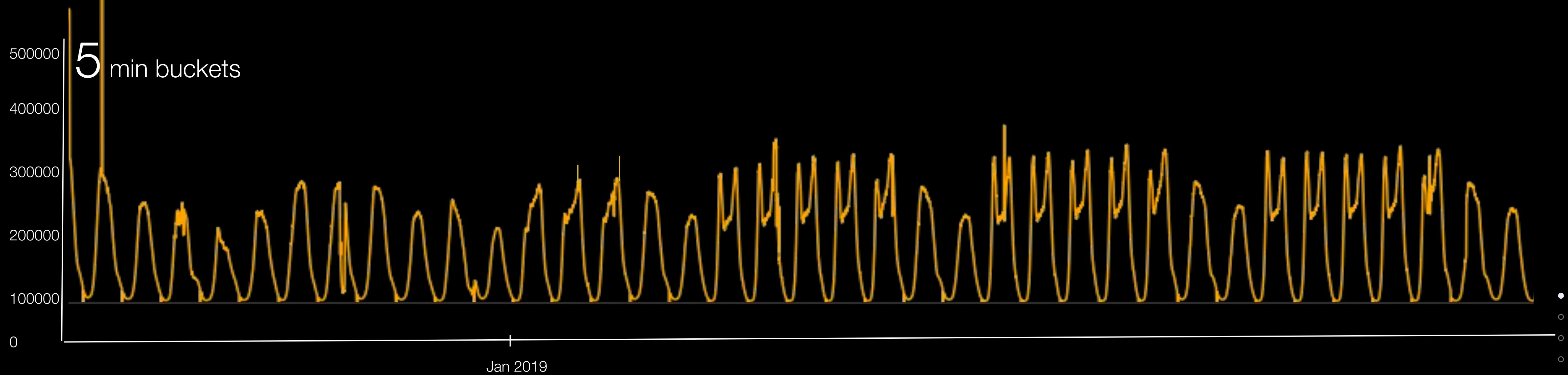
1) Extract the
Datetime & Sum

2) Aggregate results into
5 min buckets

... Data visualisation



••• Data processing



Anomaly detection models

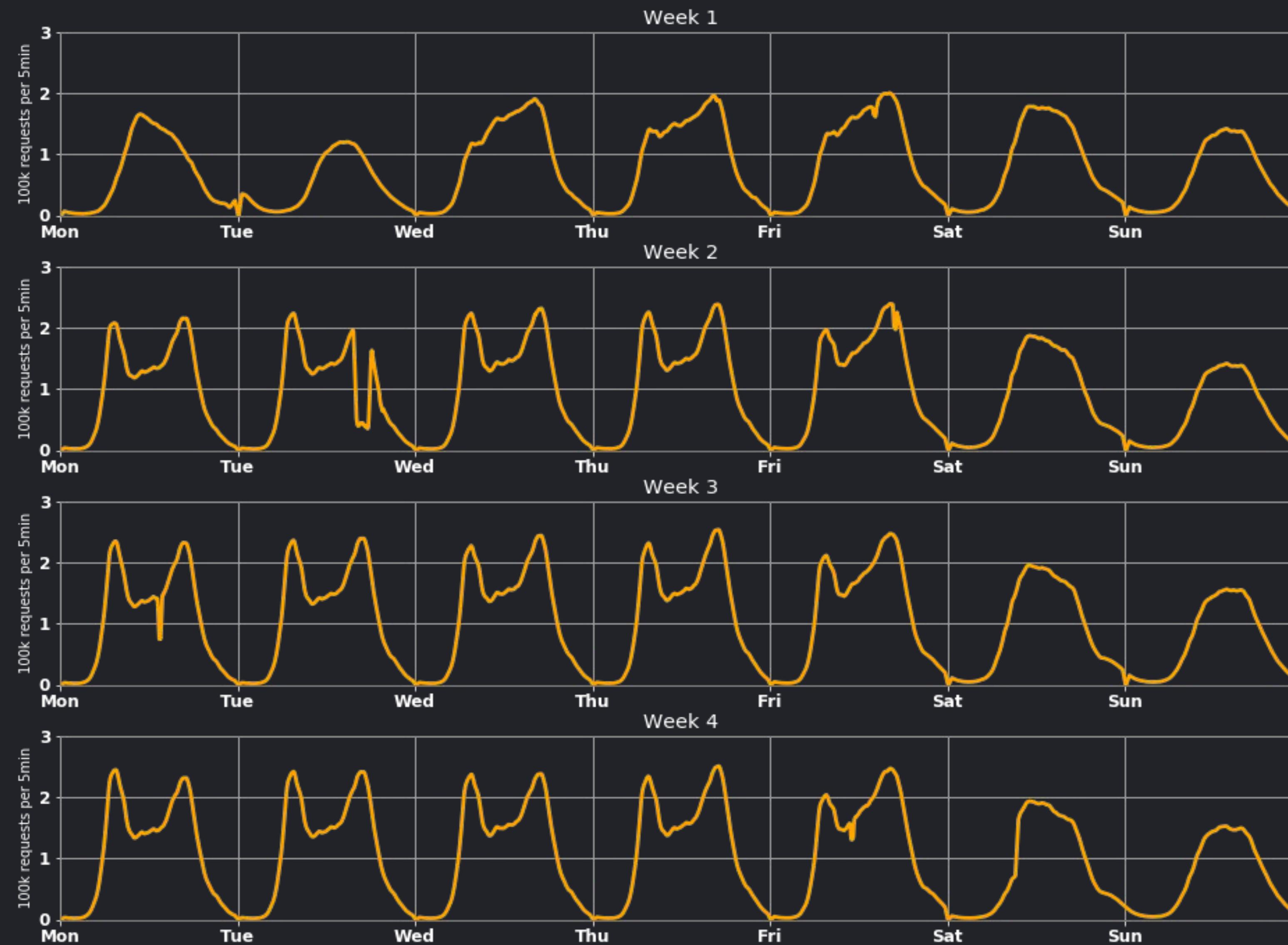


1 - Mean Predictor

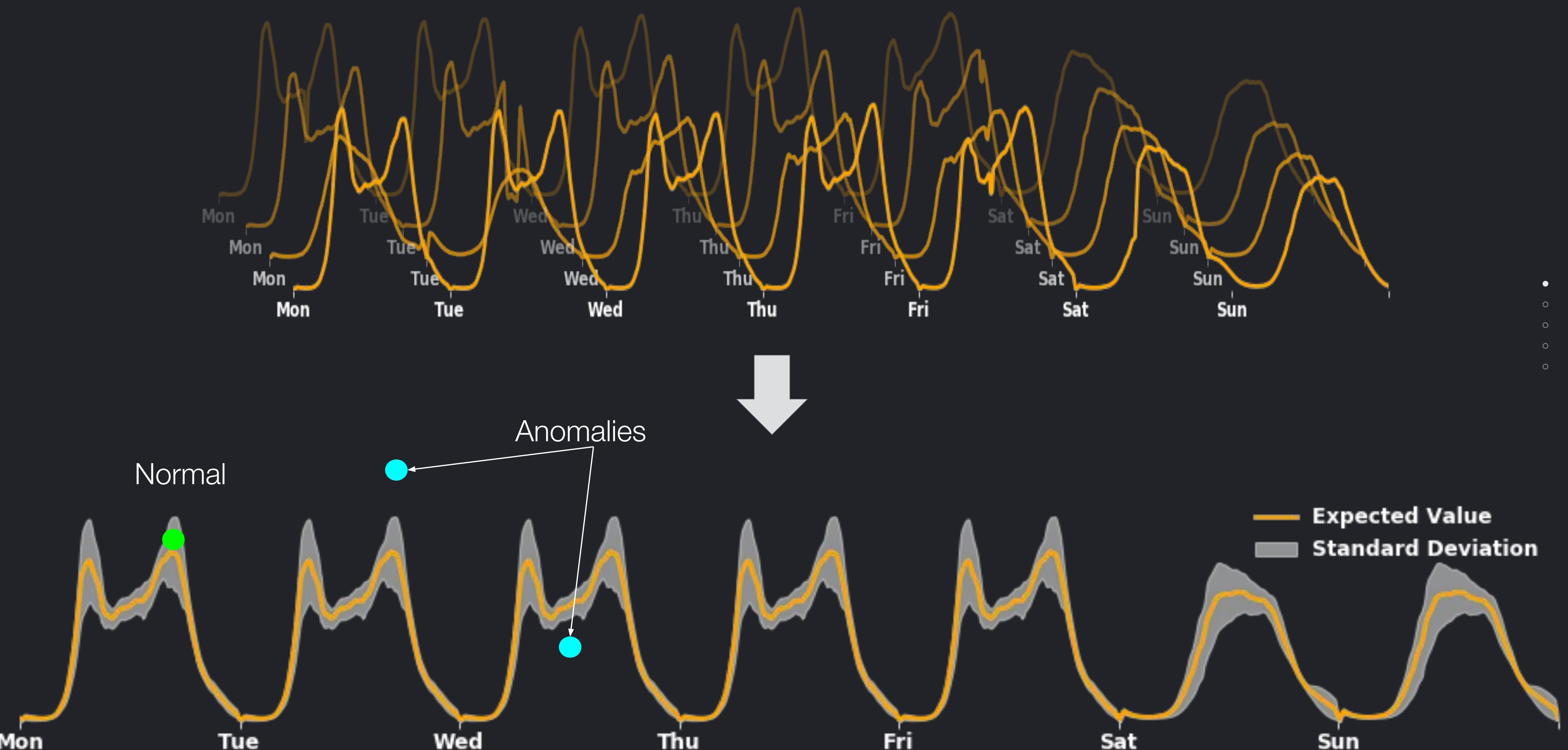
•
○
○
○
○
○

o

Data visualization



... Mean predictor



PROS

- Requires very little data to train
- Trains very fast
- Requires very little computation power
- Robust against outliers
- Can be used both for forecasting as well as anomaly detection

Can only take into account a single source of information

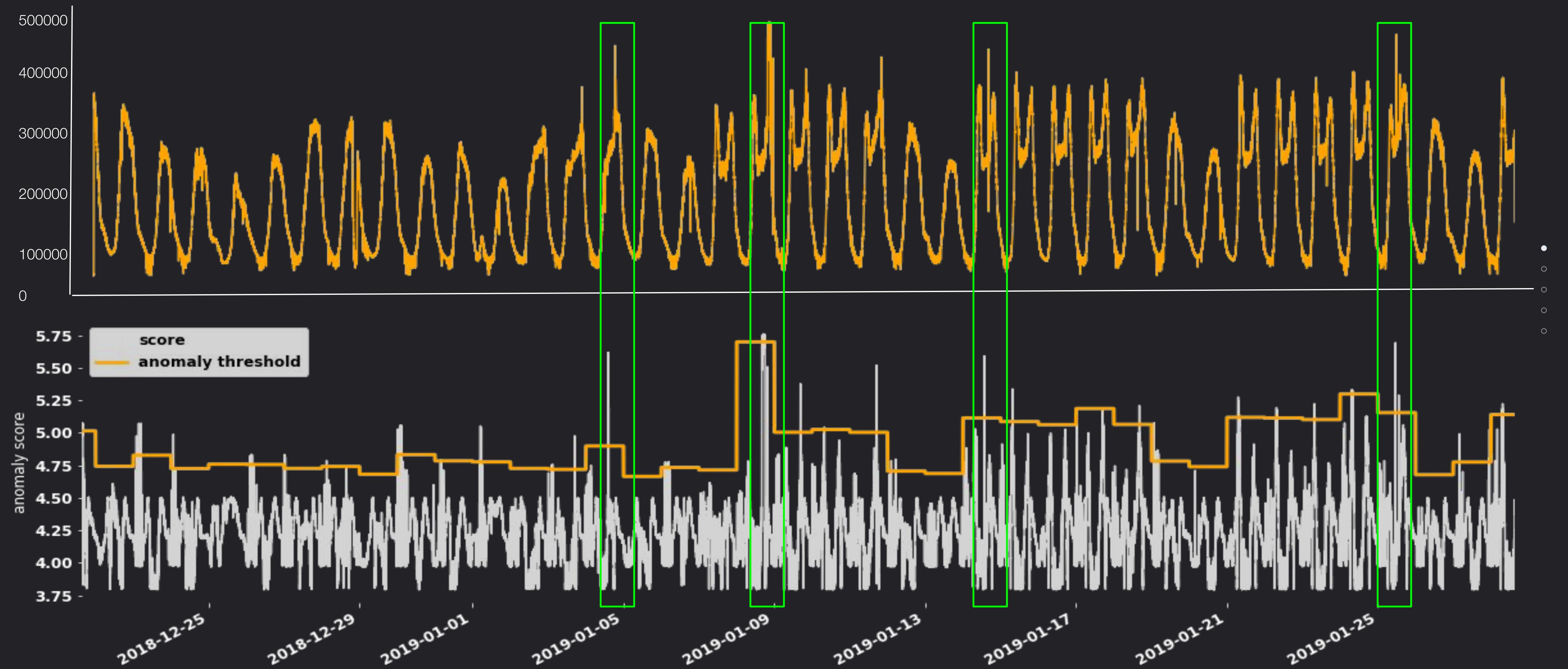
Not an Out-Of-The-Box solution from AWS

CONS

2- Random Cut Forest

•
○
○
○
○
○

○○○ Data visualization



PROS

Good performance for long-lived data streams

Short training time

Relatively easy to create and deploy the model (Out-Of-The-Box solution from Amazon)

Has a stream-friendly version

Supports streaming only on kinesis streams

Requires huge amounts of data

CONS

Technical decisions

We have:

- A lack of long-term data (FlowLogs collection started in December)
- A unique time series to analyze

....
o
o
o
o

We need:

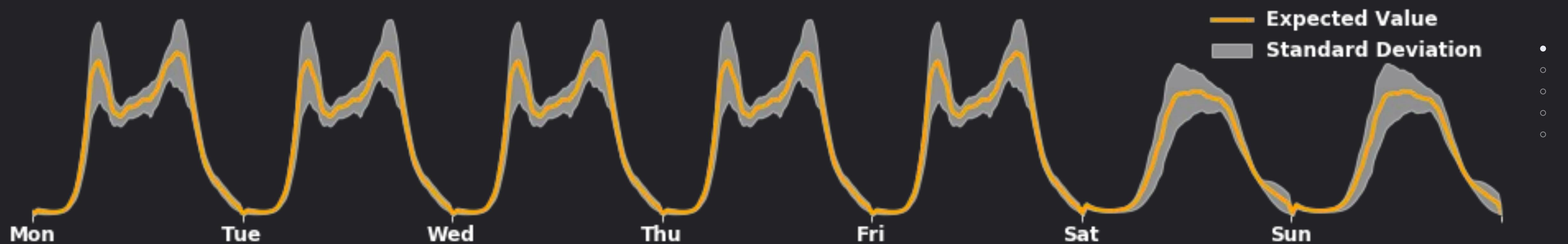
- A tweakable solution
- Short train times

→ Explains why we built our demo pipeline with Mean Predictor

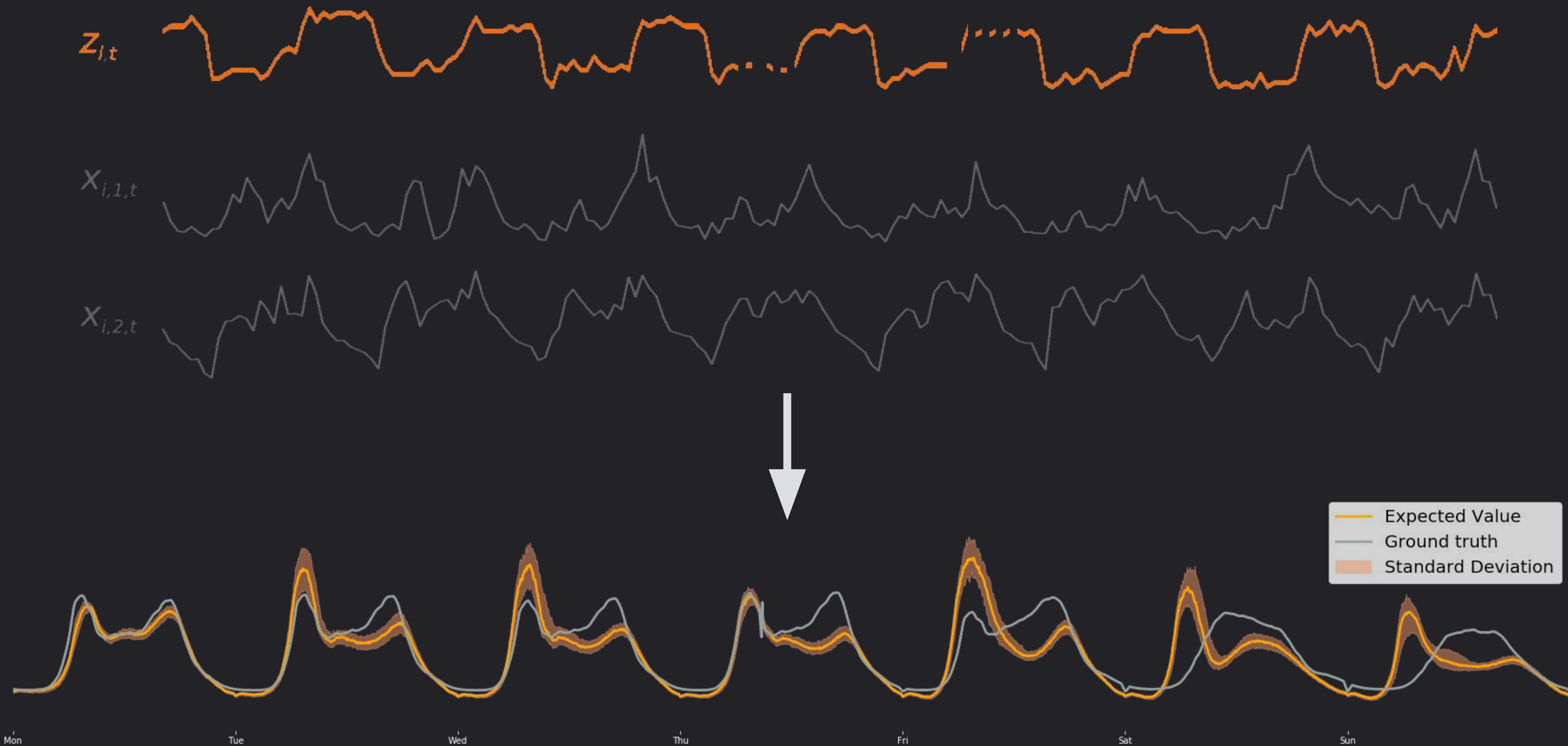
Prediction models



... Mean Predictor



DeepAR

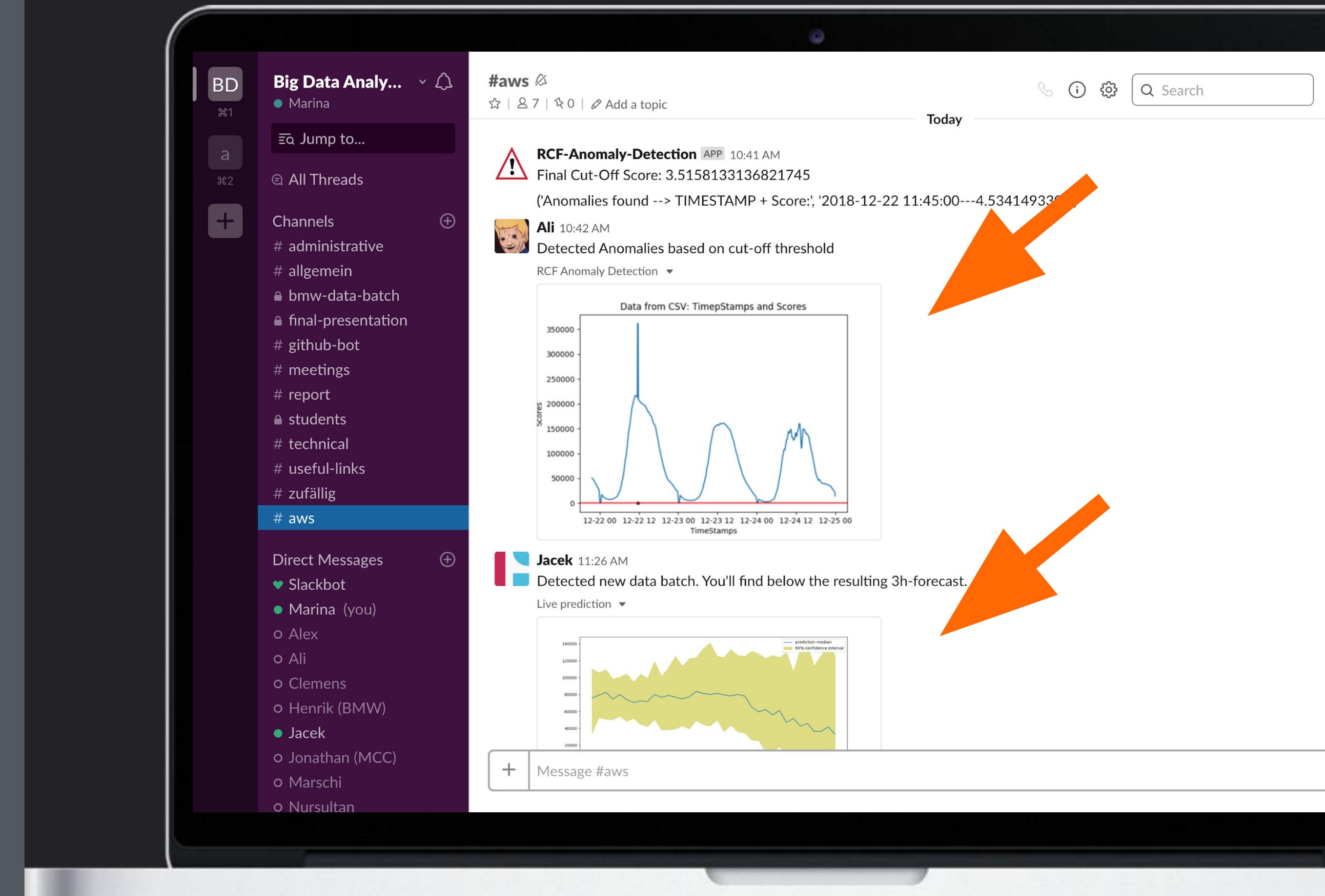


Notification functions



Anomaly detection notification

3 hour forecast notification



RECOMMENDATIONS

Implement the
Mean Predictor
Model

Switch to RCF
and DeepAR

Attach
geolocation to
data points

Gather 2-3
years of data

Look for
seasonal trends

Explore the
holidays periods



Thank you for your
attention

Slack File Edit View History Window Help

Big Data Analytics - C

Shows 21 · 5.7 · 9.0 · Add a topic

Today

Jacob 11:43 PM
Detected new data batch. You'll find below our 2h forecast.

Live 2h forecast +

Jacob 11:43 PM
Detected new data batch. You'll find below our 2h forecast.

Live 2h forecast +

Message Room

Search

Star

Share

File

Edit

View

History

Window

Help

Big Data Analytics

Shows 21 · 5.7 · 9.0 · Add a topic

Today

Jacob 11:43 PM
Detected new data batch. You'll find below our 2h forecast.

Live 2h forecast +

Jacob 11:43 PM
Detected new data batch. You'll find below our 2h forecast.

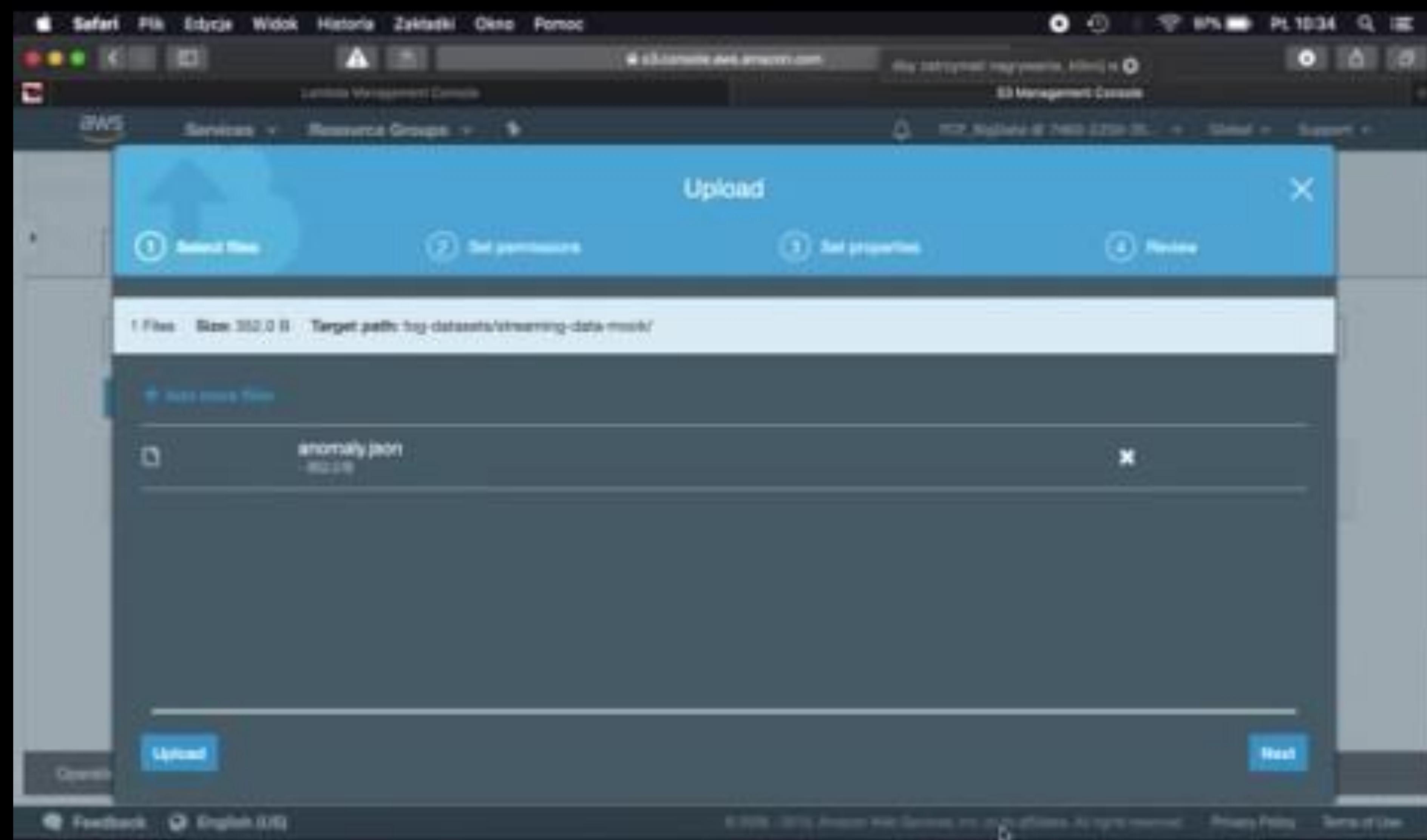
Live 2h forecast +

Message Room

Search

Star

Share



RCF Kinesis architecture

