# Agenda

Background

Problem statement

Proposed solution

Demo

Learnings

QnA

# Technology @ Intuit

Intuit is leading the way in building an AI-native development platform using cloud native open source technology. We're committed to building tools that scale and giving back to the open source community.

INTUIT  ✓ turbotax   ck creditkarma   qb quickbooks   mailchimp

**~100M**
customers

**107B**
consumer tax refunds per year

**$2T+**
invoices managed on our platform per year

**18M**
total US workers paid via QB payroll

# AI-native development platform

**AI-powered app experiences**

**4M**

Models running in production per day

**AI-assisted development: coding, testing, debugging**

**8x**

Developer velocity increase in past four years

**AI-powered app centric runtime**

**60B**

Machine learning predictions per day

**Smart operations using AIOps**

**40M+**

AIOps inferences/day

# Problem statement

# Problem statement

**1**

IT'S GONNA BE CHALLENGING

## App + k8s can not scale up fast enough

- High pod startup time
- Node scaling up take ~5m
- Image pulling
- 5xx errors, the app can't handle the surging traffic
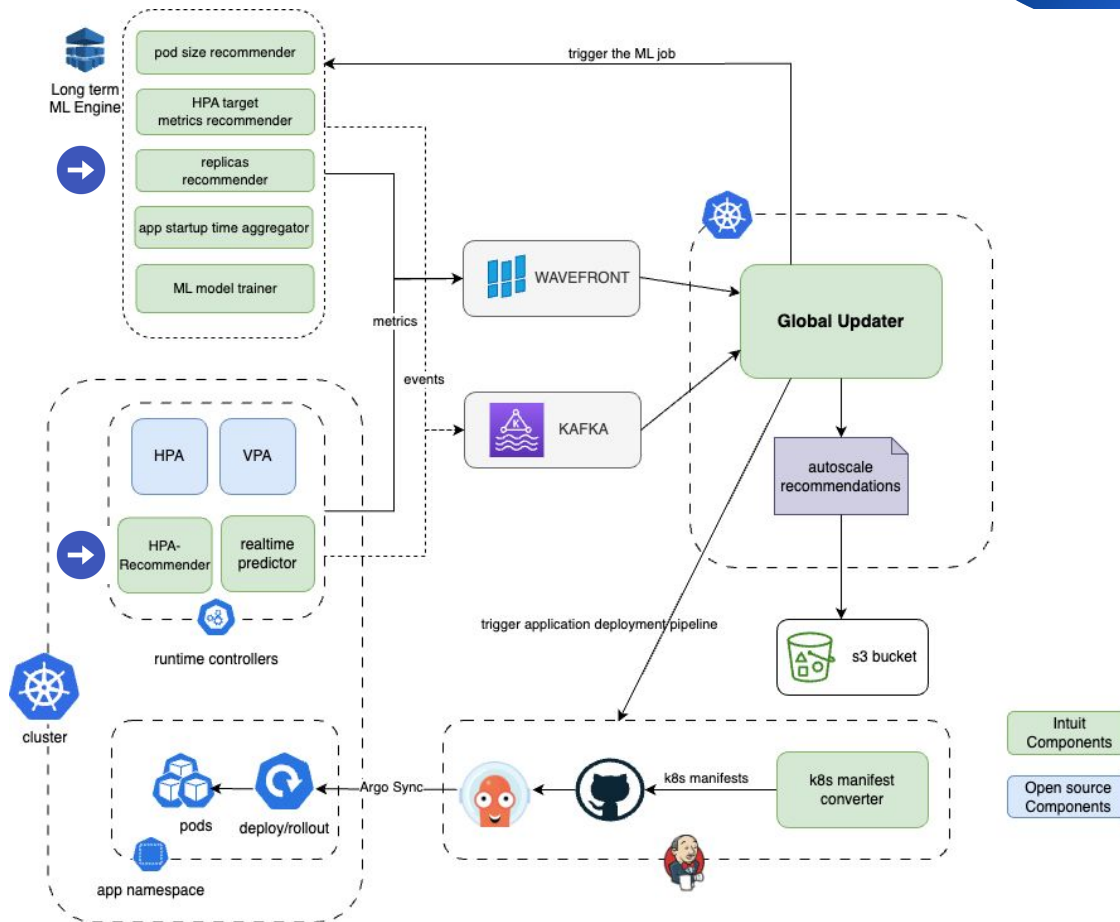


**2**

WHY DOES COST SO MUCH?

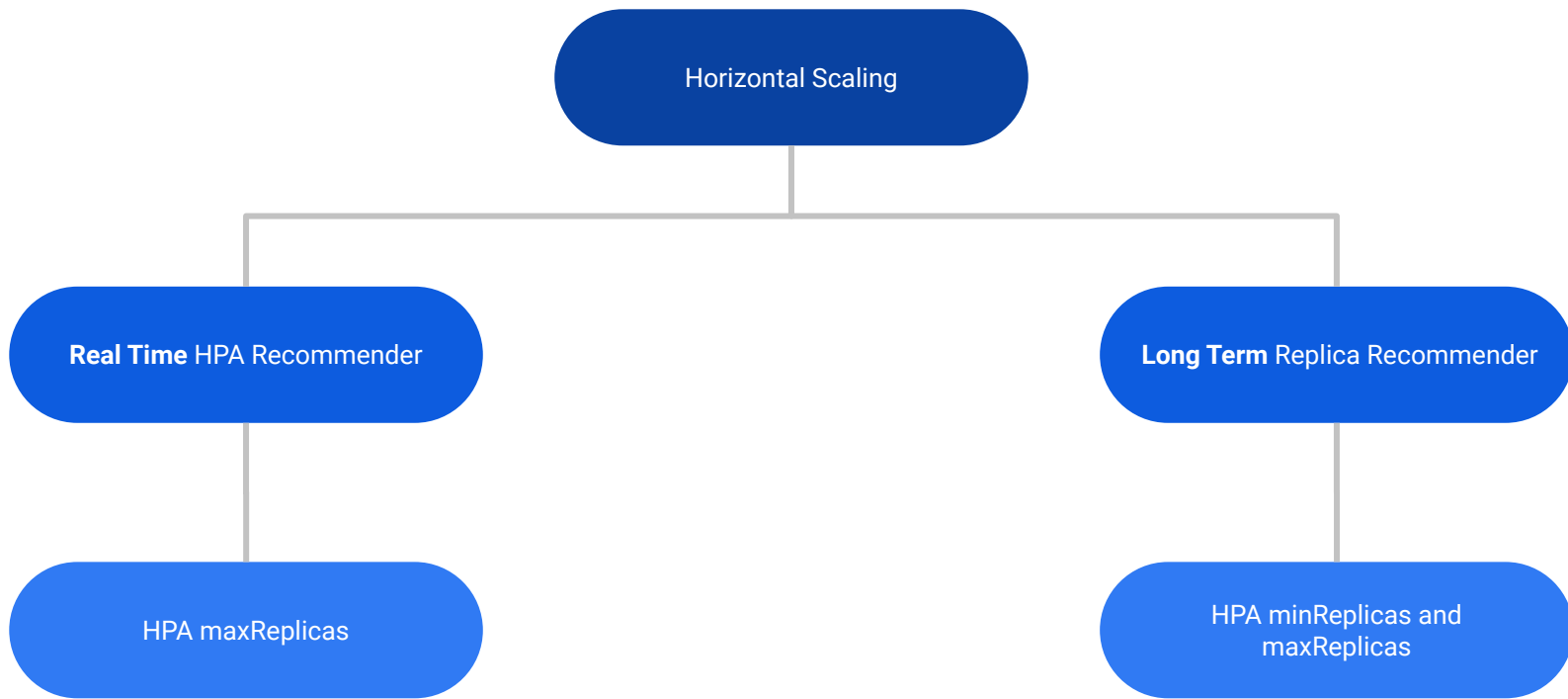## Developers have difficulty in right sizing the application and tune HPA settings

- Too conservative HPA minReplica and pod size
- Tune HPA
- High cost

# How the new auto scaling system works

# Horizontal scaling

# DeveloperPortal view



**Events Log**

**Asset ID:** 291625669146692353
**Service Name:** Attribute Generator

| Operation | Workspace | Environment | Start Date | End Date | |
|---|---|---|---|---|---|
| All | All | All | MM/DD/YYYY | MM/DD/YYYY | Clear Filters |

| Timestamp ⇅ | Workspace | Environment | Operation | Actions |
|---|---|---|---|---|
| Sunday, November 3rd, 12:08:23 am | * | * | PodSizeChange | Details |
| Thursday, October 31st, 4:11:53 pm | capital-loan-agsattributegenerator-prd | * | HpaMinMaxChange | Details |
| Thursday, October 31st, 4:11:53 pm | * | * | HpaMetricsChange | Details |
| Wednesday, October 30th, 2:59:16 pm | capital-loan-agsattributegenerator-prd | * | HpaMinMaxChange | Details |
| Wednesday, October 30th, 2:59:16 pm | * | * | HpaMetricsChange | Details |
| Wednesday, October 30th, 2:59:15 pm | * | * | PodSizeChange | Details |
| Tuesday, October 29th, 8:06:59 am | * | * | PodSizeChange | Details |

Close

# Events Log

✕

**Asset ID:** 291625669146692353

**Service Name:** Attribute Generator

| Operation | Workspace | Environment | Start Date | End Date | |
|---|---|---|---|---|---|
| All ▾ | All ▾ | All ▾ | MM/DD/YYYY 📅 | MM/DD/YYYY 📅 | **Clear Filters** |

| Timestamp ⇕ | Workspace | Environment | Operation | Actions |
|---|---|---|---|---|
| Sunday, November 3rd, 12:08:23 am | * | * | PodSizeChange | **Details** |
| Thursday, October 31st, 4:11:53 pm | capital-loan-agsattributegenerator-prd | * | HpaMinMaxChange | **Details** |
| Thursday, October 31st, 4:11:53 pm | * | * | HpaMetricsChange | **Details** |

## Details

✕

Old

New

| | 📄 data `CHANGED` | | | |
|---|---|---|---|---|
| | @@ -1,4 +1,4 @@ | | | |
| 1 | { | 1 | | { |
| 2 | −   "hpaMaxReplicas": "11", | 2 | + | "hpaMaxReplicas": "10", |
| 3 | −   "hpaMinReplicas": "4" | 3 | + | "hpaMinReplicas": "3" |
| 4 | } | 4 | | } |

# Recommendation

2024-11-03T07:08:23Z new autoscale recommendation with ID 1730617661 synced to s3 {"components":[{"name":"ags-attribute-generator","vertical":[{"containerName":"app","new":
{"memMin":"1860000000","memMax":"1875000Ki"},"old":{"memMin":"2281906618","memMax":"2281906618"}}]}]}
2024-10-31T23:11:53Z new autoscale recommendation with ID 1730377693 synced to s3 {"components":[{"name":"ags-attribute-generator","horizontal":[{"environment":"prd","new":
{"min":"3","max":"10"},"old":{"min":"4","max":"11"}}]}]}
2024-10-30T21:59:15Z new autoscale recommendation with ID 1730325003 synced to s3 {"components":[{"name":"ags-attribute-generator","vertical":[{"containerName":"app","new":
{"memMin":"2281906618","memMax":"2281906618"},"old":{"memMin":"2162292018","memMax":"2162292018"}}],"horizontal":[{"environment":"prd","new":{"min":"4","max":"11"},"old":
{"min":"3","max":"10"}}]}]}

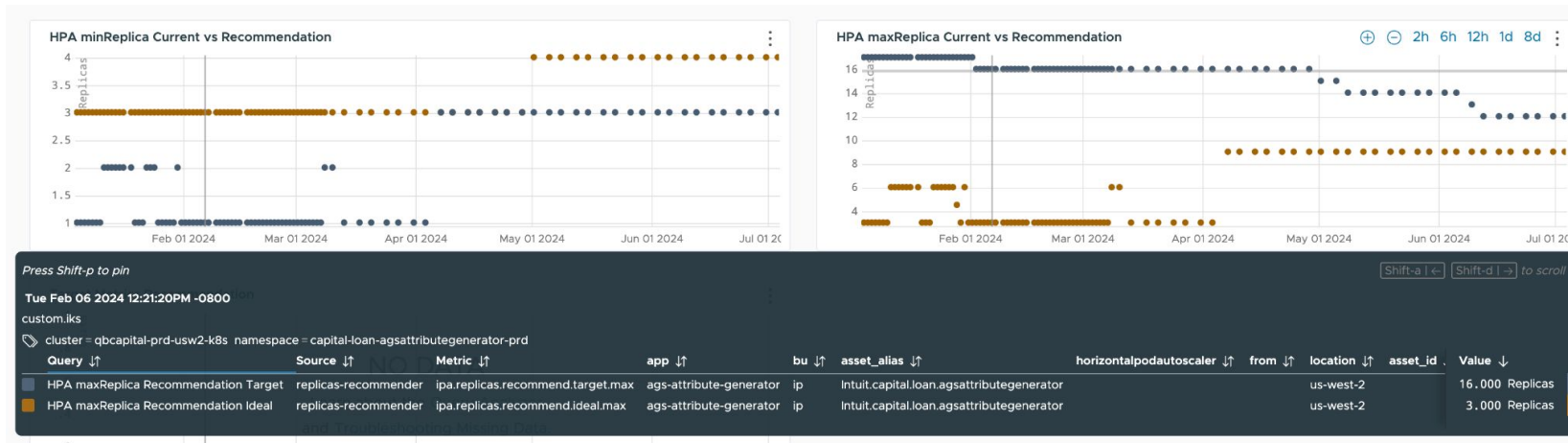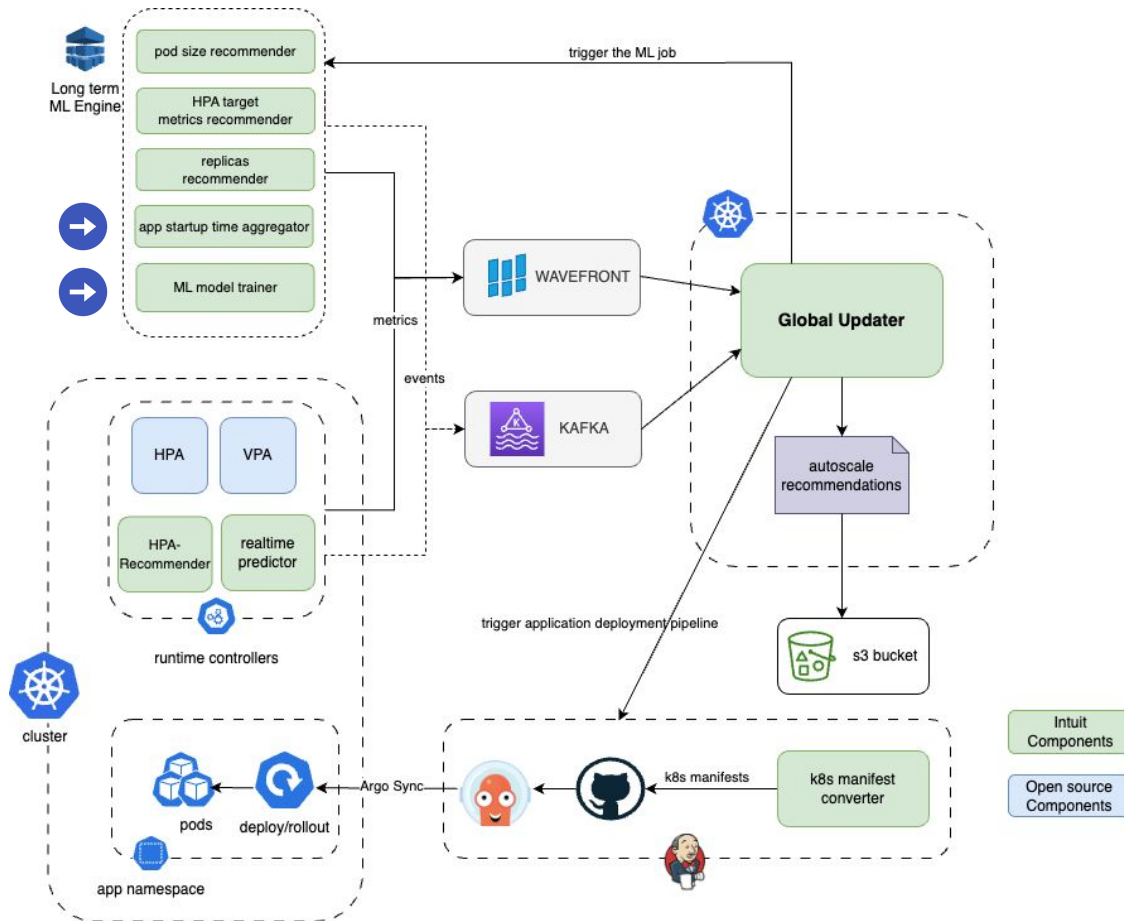| | | | |
|---|---|---|---|
| 🔲 **svc-express** new autoscale recommendation synced to s3 | ✕ fb2a9e5  yesterday | 🕐 **236** commits |
| 📄 .gitpod.yml | CWS-191: Automated Changes, to Gitpod Java Template by Cloud W… | last year |
| 📄 .iks-express.log | new autoscale recommendation synced to s3 | yesterday |
| 📄 Jenkinsfile | QBF-28520 add east2 deployment approval (#20) | 4 months ago |
| 📄 Jenkinsfile.pci | QBF-27325-Add/Update Environments & Sample Application Yaml (#… | 6 months ago |
| 📄 README.md | Initial commit | last year |
| 📄 iks-express.yaml | [Changed] - infrastructure deployment: [env=all, jenkins_build=jenki… | 4 months ago |
| 📄 msaas-deployment-config.yaml | QBF-27325-Add/Update Environments & Sample Application Yaml (#… | 6 months ago |

# Recommendation over time



**Conservative, with a bias towards availability and correctness**

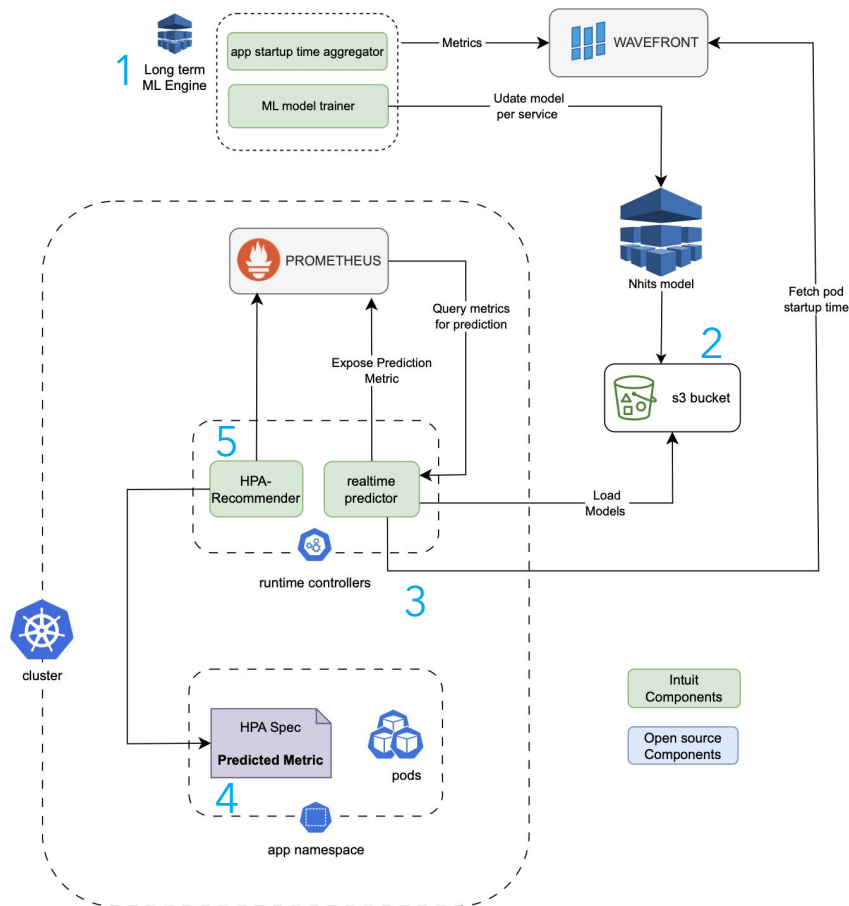# How the new auto scaling system works

# Proposed solution

Time series forecasting for real-time cpu usage prediction

- Integrate predicted metrics with Prometheus

- Use predicted metrics as custom metric in HPA

- Proactively increases desired replicas x minutes ahead of time, based on pod start up time

# Autoscaling ahead of time



ML Model Training

Exposing Prometheus Metric

Use Custom metric in HPA

# Models under consideration

We evaluated 4 separate time series forecasting models

**Prophet** model

**TimesFM** forecasting model
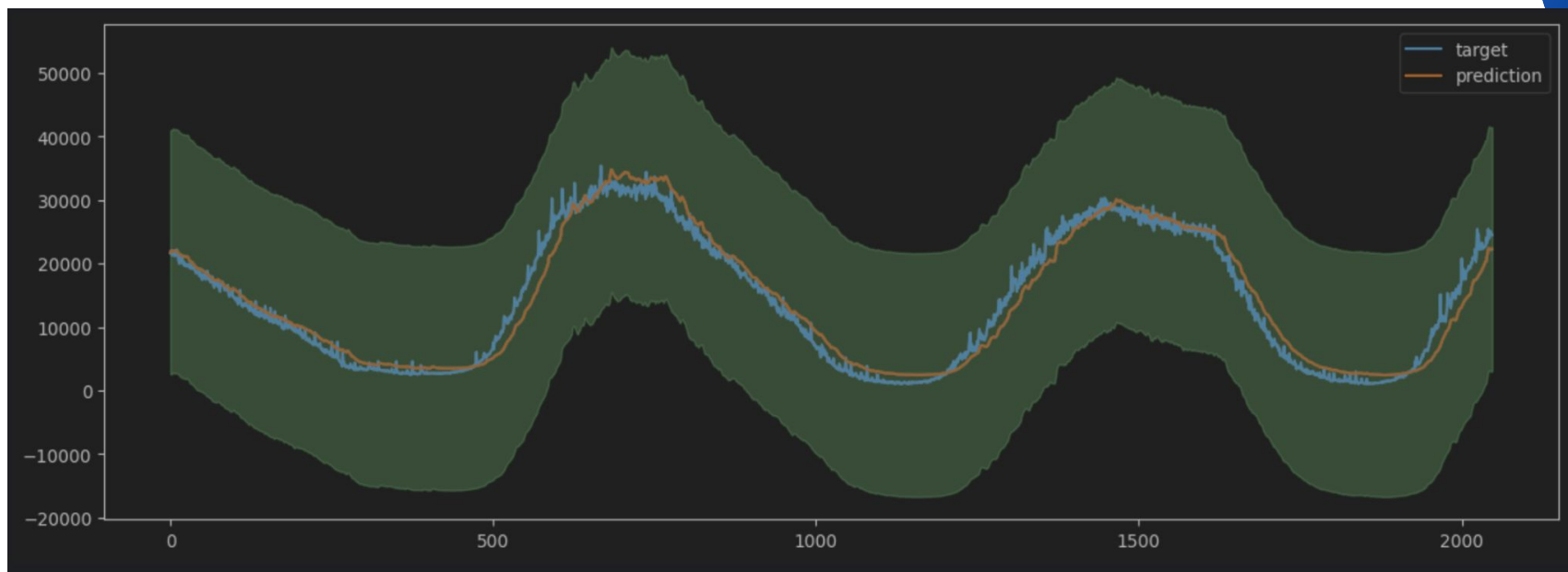
In-**House GRU** model

**Nhits** forecasting model

# Prophet model



Actual vs Predicted CPU Usage

# Google's TimesFM model

# In house RNN forecasting model

# Nhits model
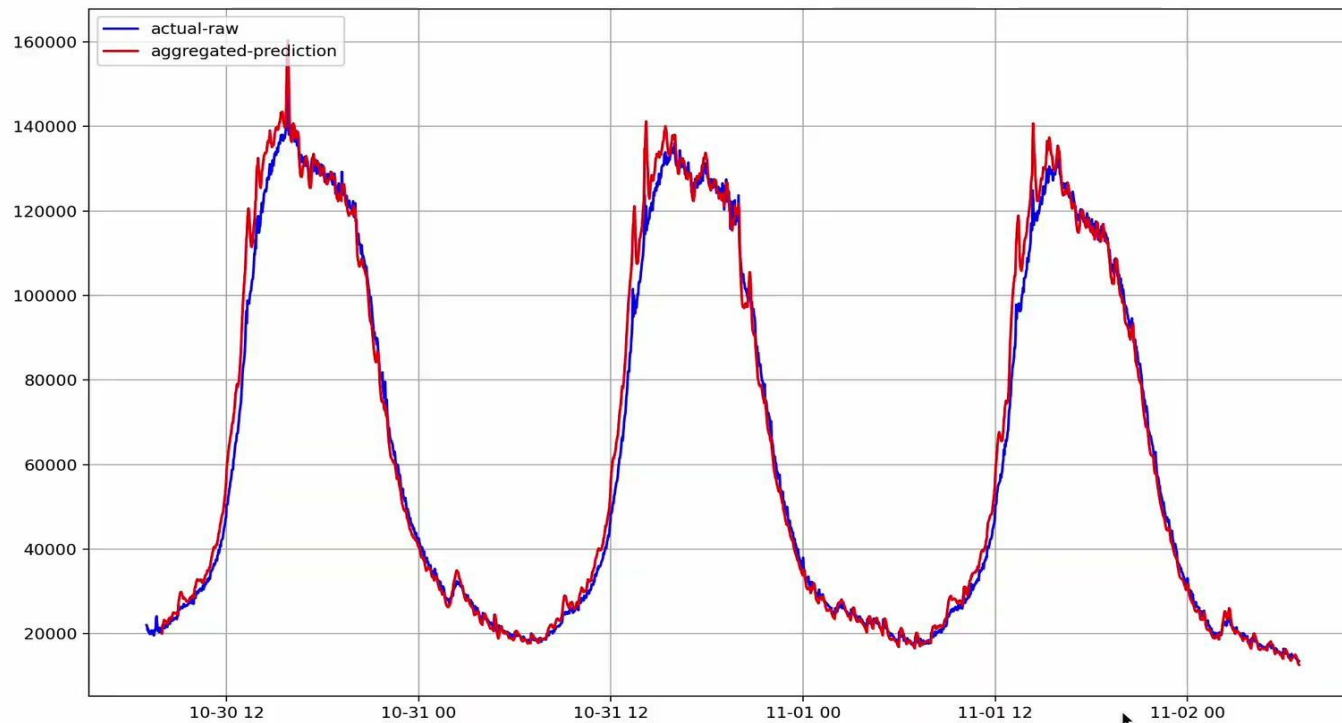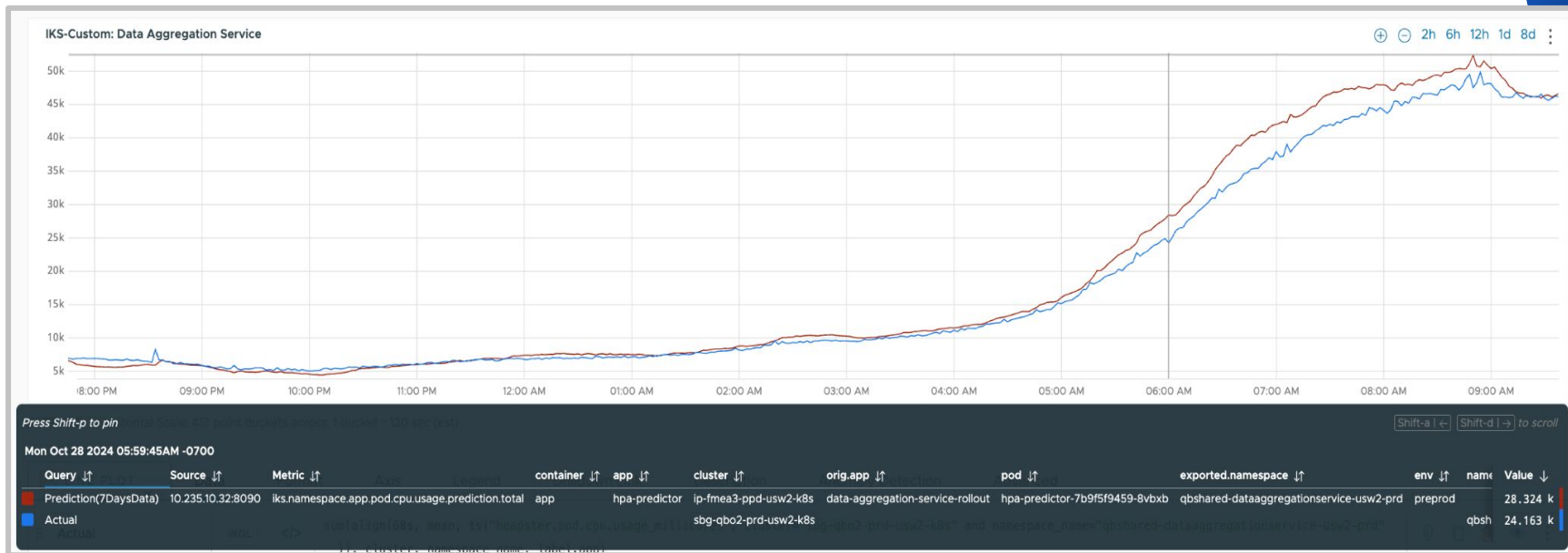
Demo

# 1. Train Nhits model

# 2. Expose metric

# 3. Use predicted metric in HPA

```
spec:
  maxReplicas: 20
  metrics:
  - object:
      describedObject:
        apiVersion: apps/v1
        kind: Deployment
        name: test-app
      metric:
        name: avg_cpu_utilization_metric
      target:
        type: Value
        value: "50"
      type: Object
  - object:
      describedObject:
        apiVersion: apps/v1
        kind: Deployment
        name: test-app
      metric:
       name: avg_cpu_utilization_predictive_metric
      target:
        type: Value
        value: "50"
    type: Object
    minReplicas: 9
```

# Learnings

# Learnings

1. Performance test: small ramp up / ramp down time
   a. Prediction was having delay recognizing pattern in first 10 minutes
   b. Solution: HPA uses max of both actual and predicted metrics

2. Difficult to predict for spiky cpu usage
   a. Creates a spiky prediction metric
   b. Solution: Smooth the training data, prediction metric

3. Not enough data
   a. 14 days or less of data

# Future enhancements

## Apps with seasonal traffic pattern

- Tax peak
- Super Bowl event

## Use other metrics

- Custom metrics
- Jvm metrics
- Tps

## Multidimensional prediction

- Prediction metric not based on single metric
- Prevent from influencing real cpu usage data

# Stay in the loop

**FOLLOW**

# Intuit Open Source

Don't miss on exciting OSS events, activities & news

Scan or visit
**bit.ly/intuit-oss**

## Visit our Booth

Get some exciting OSS swag - while supplies last



## Check out our past presentations

**Platform Eng Day**
https://www.youtube.com/watch?v=z6ItgXM4RxE
**Autoscaling**
https://www.youtube.com/watch?v=h2zmITPG3GM
**Debuggability**
https://www.youtube.com/watch?v=bPa1PjY-Hg4

Proposed solution

# Proposed solution

**Time series forecasting for real-time cpu usage prediction**

- Integrate predicted metrics with Prometheus
- Use predicted metrics as custom metric in HPA
- Proactively increases desired replicas x minutes ahead of time, based on pod start up time

**Cost saving**

- Recommend an optimal minReplica value

# Benefits



### Improve service availability

Proactively increases desired replicas ~x minutes ahead of time



### Promising result for weekly/daily traffic pattern app

The ML Model is able to predict ahead of time



### Cost saving

Reduce number of minReplicas