# Hystax OptScale

## MLOps open source platform

Multiply a number of ML/AI experiments with minimal cloud costs

# Hystax

Founded in 2016,
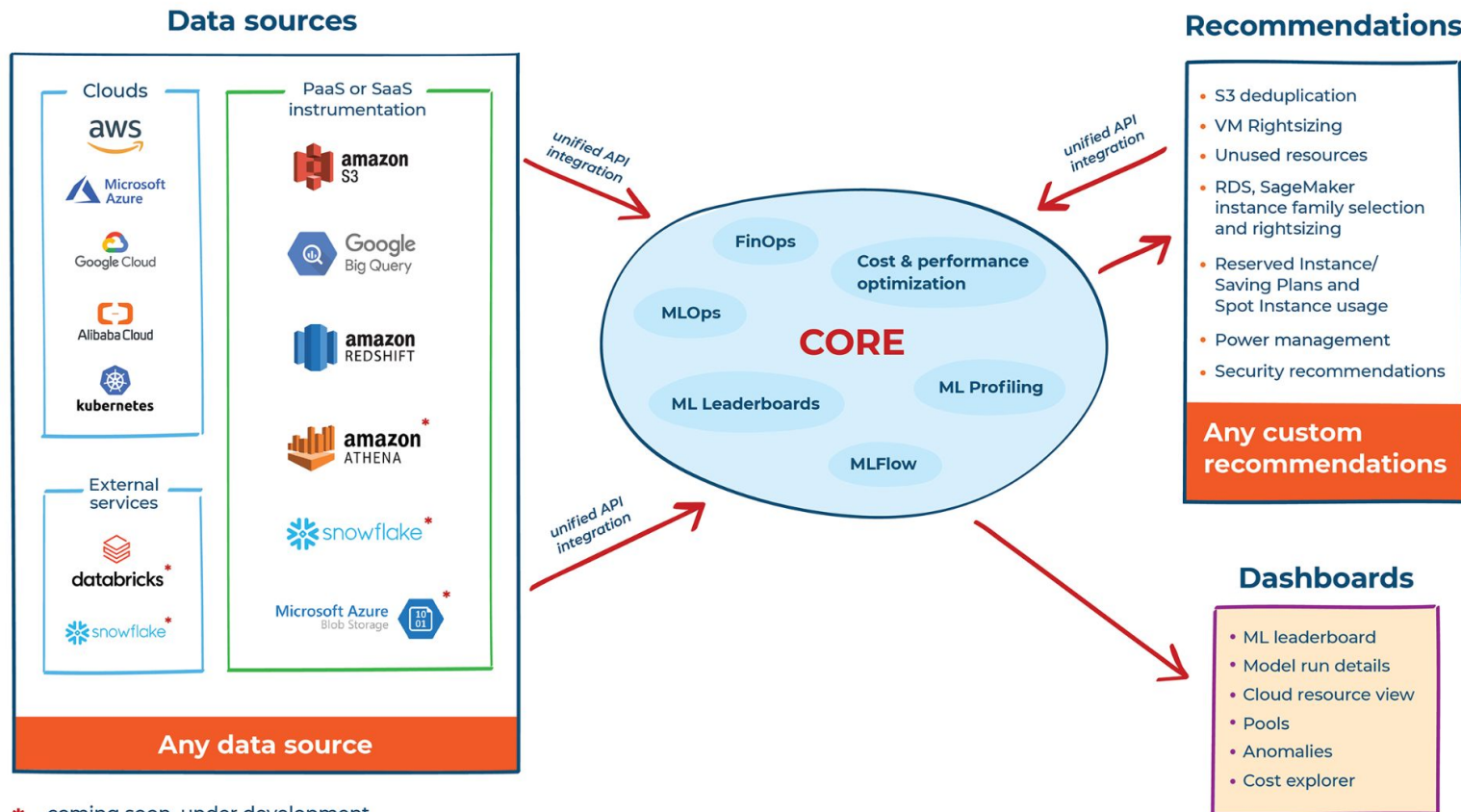customers in 48 countries

**Customers:** Airbus, Nutanix,
Orange, Nokia, DHL, Burger King

# OptScale use cases

**MLOps**

FinOps & cloud
cost optimization

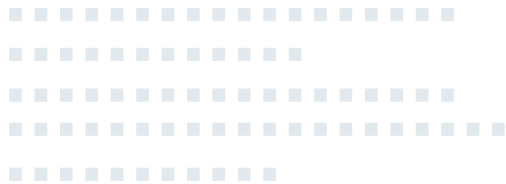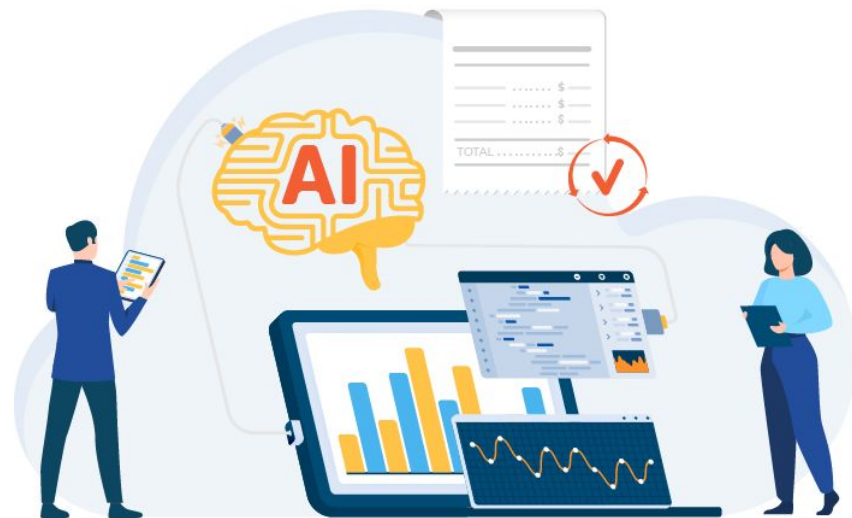**Open source under Apache 2.0:** <u>https://github.com/hystax/optscale</u>

# MLOps

- Team and individual ML engineer progress observability

- ML/AI task profiling, bottleneck identification

- PaaS or any external service instrumentation

- Optimization recommendations

- Runsets to automatically scale a number of experiments

# ML R&D status observability

- List of models with goals status and active recommendations

- Tracking the number and quality of experiments run by a team

- Cost of an overall model and individual experiments

# ML R&D status observability

# ML/AI profiling & optimization

- ML/AI model training tracking and profiling, inside and outside metrics collection

- CPU/RAM/GPU/Disk IO correlation tracking

- Minimal cloud cost for ML/AI experiments and development by utilizing Reserved Instances/Savings Plans and dozens of optimization scenarios

**Supported platforms:**

# ML/AI optimization recommendations

- Utilizing Reserved/Spot Instances and Savings Plans

- Rightsizing and instance family migration

- Detecting CPU, GPU, RAM, and IO bottlenecks

- Cross-regional traffic

- Experiment/run comparison

# ML/AI profiling & optimization

hx hystax

## Application overview

Applications / Shoes categorizer

PROFILING INTEGRATION    CONFIGURE

**OVERVIEW**    EXECUTORS

---

⊘ **Aborted**
Status

**1 minute, 5 seconds**
Last run duration

**$284.25**
Lifetime cost

**$8.48k**    **19**
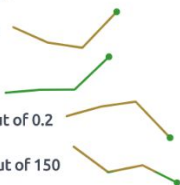Summary savings    Recommendations count

### Tracked parameters ⚙

Accuracy: ● 1 out of 0.98

Data Loss: ● 5 out of 10

Inference time: ● 0.15 out of 0.2

Data processed: ● 132 out of 150

### Application summary

Key: **shoes_categorizer**
Runs: **9**
Last successful run: **3 months ago**
Last run cost: **$14**
Owner: **Charlie Fisher**

### Last run executor

aws **i-3920fedd**
Cloud: aws **AWS HQ**
Name: **sunflower-eu-fra-1**
Region: **us-west-2**
Size: **t3a.medium**

## Recommendations

| Executors upgrade | Cross-region traffic | Spot/Preemptible instances usage | Local storage bottleneck |
|---|---|---|---|
| See details **11** Count | See details **$8.48k** Savings | See details **3** Count | See details **1** Count |

| GPU memory | | | |
|---|---|---|---|
| See details **1** Count | | | |

# ML/AI profiling & optimization

# PaaS or any external service instrumentation

- Cost, performance, and output details of any API call to PaaS or an external service

- Metrics tracking and visualization

- Performance and cost optimization of API calls

- Cross-regional traffic

- S3, Redshift, BigQuery - ready, unified way to add more services

# Runsets

- Automated run of a number of experiments with configurable datasets, hyperparameter ranges and model versions

- Optimal hardware with cost-efficient usage of Spot, Reserved Instances / Savings Plans

- Configurable experiment goals and success criteria

- Various complete/abort conditions - take first successful, complete all

- Integrated profiling to identify bottlenecks

# Runsets

## Runset overview

**AWS GPU Instances** / #3_gentle_sky

| **6** | **1** | **$73.2** |
|---|---|---|
| Configurations tried | Runs met goals | Total expenses |

Application: **Shoes categorizer**

### Parameters
Data source: aws **AWS HQ**
Region: aws **us-east-1**
Instance type: aws **t3a**
Maximum parallel runs: **14**

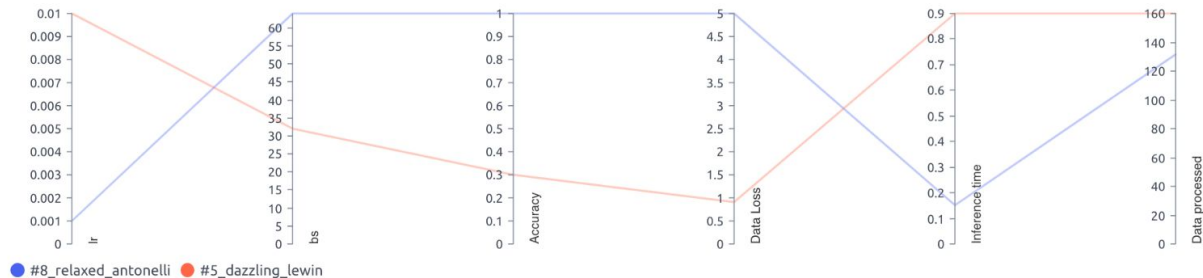### Hyperparameters count
Model path: **2** ⓘ
Dataset path: **3** ⓘ
Learning rate: **4** ⓘ

### Stop conditions
Stop runset when projected expenses exceed 20$
Stop individual run if its duration exceeds 3 minutes

### Correlations ▥

● #8_relaxed_antonelli    ● #5_dazzling_lewin

---

**RUNS**    **EXECUTORS**

Filters: None  Status ▾  Goals ▾

🔍 Search

| # ↓ | Status ❓ | Hyperparameters | Goals ⚙ | Started at | Duration | Executors |
|---|---|---|---|---|---|---|
| **#9_whispering_fog** | ◉ Aborted<br>Reached plateau for the **Accuracy** goal with the value of **0.87** | lr: **0.01**<br>bs: **32** | Accuracy: -<br>Data Loss: -<br>Inference time: -<br>Data processed: - | 01/26/2023 05:55 AM | 1 minute, 5 seconds | aws i-3920fedd<br>Size: **t3a.medium**<br>Expenses: **$14** |

# Roadmap



- Cost plugin for MLflow, WanDB, and neptune.ai

- Integration with Optuna to optimize Reserved Instance and other hardware parameter usage

- Model versioning

- Better hardware selection recommendations based on usage patterns and algorithms

# FINOPS & CLOUD COST OPTIMIZATION

# FinOps and cost management

- Forecast and monitor an IT infrastructure cost

- Identify wastage and optimize IT expenses

- Bring resource / application / service observability

- IT asset management

- Set TTL and budget constraints

- Establish a long-term FinOps process by engaging engineering teams

**Supported platforms:**

# OptScale vs cloud-native cost explorer

- Cloud resource visibility and filtering across all the clouds, accounts and regions

- Dozens of optimization scenarios not supported by clouds incl. one of the best rightsizing engines

- Cost allocation not just by tags but other properties

- Geo and network traffic map

- TTL rules and budget constraints

- FinOps: OptScale is built for engineers to be responsible for their cloud resources

# Contacts

📞  +1 628 251 1280

🌐  hystax.com

✉️  info@hystax.com

📍  1250 Borregas Avenue, Sunnyvale, CA, 94089