Хайстекс ОптСкейл

FinOps и MLOps платформа с открытым исходным кодом

Запускайте ML/AI или любой тип рабочей нагрузки с оптимальной производительностью и стоимостью инфраструктуры



Хайстекс



Компания основана в 2016 году, разработка и поддержка в РФ



Среди наших клиентов -Сбербанк, Яндекс.Облако, X5 Retail Group, Burger King, M. Видео, Спортмастер

Сценарии использования ОптСкейл

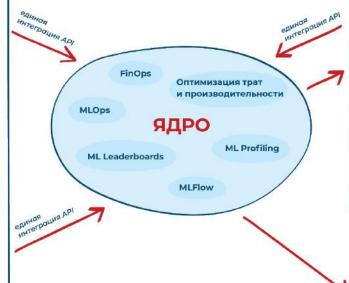




Источники данных



Рабочая схема ОптСкейл



Рекомендации

- S3 дедупликация
- VM Rightsizing
- Неиспользуемые ресурсы
- RDS, SageMaker instance family selection and rightsizing
- Reserved Instance/ Saving Plans and Spot Instance usage
- Power management
- Рекомендации по безопасности

Любые кастомные рекомендации

Дашборды

- ML leaderboard
- Model run details
- Cloud resource view
- · Pools
- Anomalies
- Cost explorer

* - в процессе разработки

FinOps и оптимизация расходов на облако

FinOps и оптимизация расходов на облако

- Прогнозирование и мониторинг затрат на ИТ-инфраструктуру
- Выявление потерь и оптимизация расходов
- Обеспечение прозрачности ресурсов/приложений/сервисов
- Управление ИТ-ресурсами
- Установка TTL и ограничений бюджета
- Внедрение долгосрочного процесса FinOps путем вовлечения команд инженеров

Хайстекс ОптСкейл vs cloud-native решения

- Прозрачность облачных ресурсов, включая различные облака, учетные записи и регионы
- Десятки сценариев оптимизации, не поддерживаемых облачными провайдерами, включая один из лучших механизмов rightsizing (перераспределение ресурсов инфраструктуры в зависимости от текущих потребностей)
- Распределение затрат не только по тегам, но и по другим свойствам
- Карта географического распределения и сетевого трафика
- Правила TTL и ограничений бюджета
- FinOps: ОптСкейл создан для инженеров, ответственных за свои облачные ресурсы

Оптимизация расходов на облако vs FinOps

Оптимизация расходов:

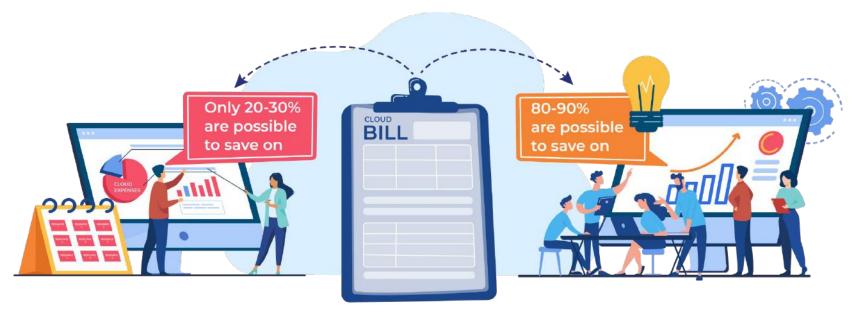
- Фокус на ИТ сотруднике, которому необходимо отслеживать R&D команды для тегирования, перераспределения, удаления неиспользуемых ресурсов
- Предоставление отчета, для решения проблем в краткосрочной перспективе, через несколько месяцев проблемы возвращаются
- R&D команда отключена от процесса экономии затрат и не несет ответственности

FinOps:

- Фокус на всей команде FinOps, включая инженеров, которые генерируют основную часть затрат
- Создание долгосрочного процесса экономии расходов путем вовлечения и обучения команды
- ИТ-специалисты несут ответственность за процесс управления; инженеры - за свои ресурсы и ТТL; ОптСкейл - за обучение команд и предоставление эффективных практик

Оптимизация расходов на облако

FinOps



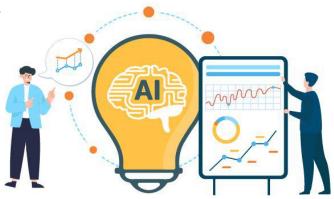
Облачные решения по управлению затратами созданы только для нескольких ИТ-сотрудников, ответственных за оптимизацию затрат, но у них ограничены управленческие функции и влияние на команды исследования и разработки

FinOps вовлекает в процессы оптимизации затрат руководителей компаний, финансовые и инженерные команды

MLOPS: ML/AI профайлинг и оптимизация

MLOps

- Запуск Runsets для автоматического масштабирования количества экспериментов
- Наблюдение за прогрессом команды и индивидуальных инженеров в области машинного обучения
- Профилирование задач ML/AI, выявление узких мест
- Рекомендации по оптимизации



Runsets

- Автоматический запуск нескольких экспериментов с настраиваемыми наборами данных, диапазонами гиперпараметров и версиями моделей
- Оптимальное использование оборудования с эффективным использованием Spot, Reserved Instances/Saving Plans
- Настроенные цели экспериментов и критерии успеха
- Различные условия завершения/прерывания выбор первого успешного, завершение всех
- Интегрированное профилирование для выявления узких мест



Runsets

Runset overview

AWS GPU Instances / #3_gentle_sky

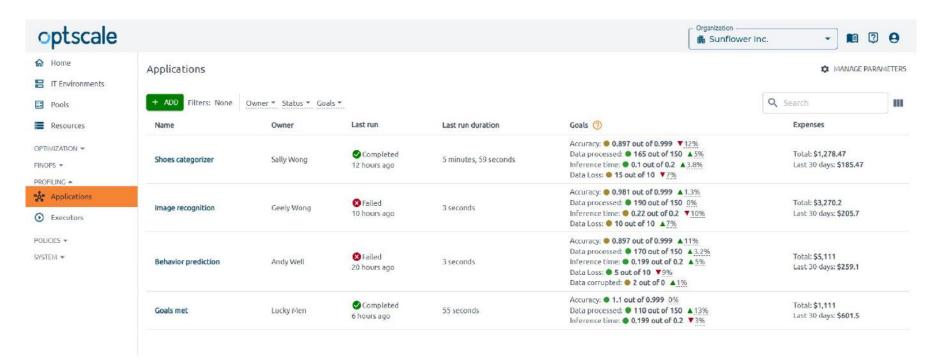


Прозрачность состояния ML R&D целей

- Список моделей с состоянием целей и активными рекомендациями
- Отслеживание количества и качества проведенных командой экспериментов
- Стоимость общей модели и отдельных экспериментов



Прозрачность состояния ML R&D целей



ML/AI профайлинг и оптимизация

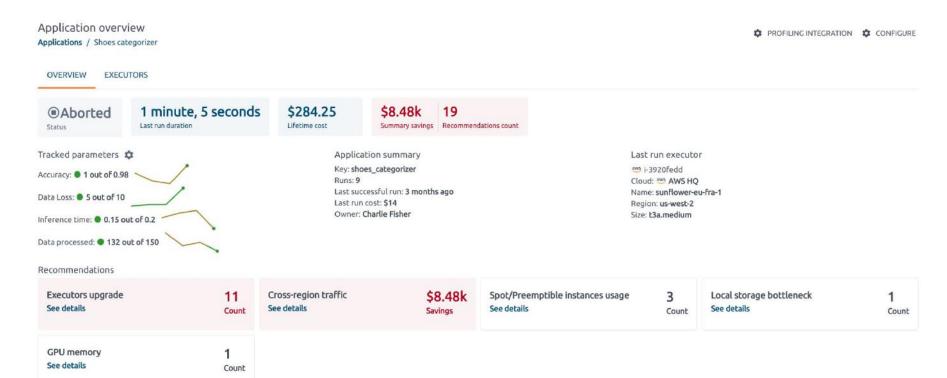
- Отслеживание и профилирование обучения моделей ML/AI, сбор метрик внутри и снаружи
- Отслеживание корреляции CPU/RAM/GPU/Disk IO
- Минимальные затраты на облачные ресурсы для экспериментов и разработки ML/Al за счет использования Reserved Instances/Saving Plans и десятков других сценариев оптимизации

ML/AI оптимизационные рекомендации

- Использование Reserved/Spot instances и Saving Plans
- Rightsizing и миграция инстансов
- Выявление узких мест в CPU, GPU, RAM и IO
- Межрегиональный трафик
- Сравнение экспериментов/запусков



ML/AI профайлинг и оптимизация



ML/AI профайлинг и оптимизация



Управление ИТ-окружениями

Управление ИТ-окружениями

- Управление списком ИТ-окружений, их состоянием и доступностью
- Бронирование ИТ-окружений и организация общего использования

Отслеживание истории развертывания, просмотр версий программного обеспечения

- Планирование ресурсов через Jira, Slack или ОптСкейл UI
- Эффективное управление и оптимизация затрат
- Мониторинг производительности ИТ-окружений











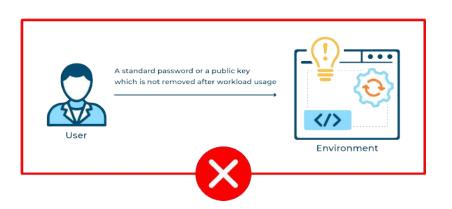




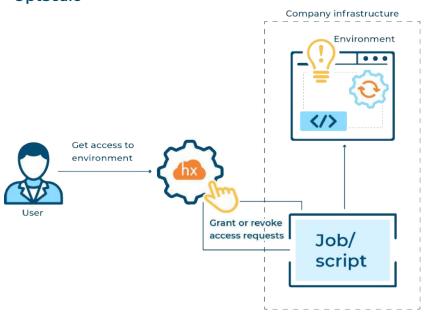


Управление доступом к ИТ-окружениям

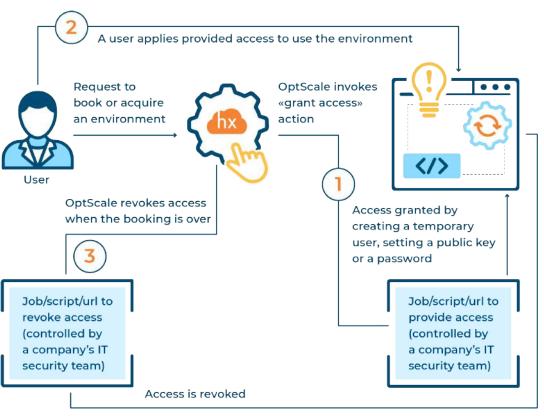
Традиционная система управления доступом к ИТ-окружениям



Система управления доступом к ИТ-окружениям с OptScale



Временный и отзываемый доступ



- ОптСкейл контролирует доступ к внутренним рабочим нагрузкам
- Скрипт или hook вызывается, когда пользователь запрашивает доступ. Сценарий предоставления временного доступа принадлежит самой компании
- Когда пользователь завершает работу с нагрузками, вызывается другой hook для отзыва доступа
- Доступны журналы аудита
- Доступны образцы скриптов для быстрой настройки

Внедрение ОптСкейл

UI и API

UI для управления настройками и просмотра отчетов, API для интеграции с джобами и пайплайнами

• Простота в использовании. Интеграция с R&D инструментами

Вашей команде не нужно изучать новый инструмент. 90% функционала доступно через Jira и Slack

- SaaS или частное развертывание Две опции доступны в продукте
- **5 минут для настройки** Нет долгой конфигурации и развертывания



Совместное использование ресурсов и управление их жизненным циклом

• Группировка ресурсов и установление их владельцев Представление кластеров, стеков, задач, а не только отдельных ресурсов. Получение, освобождение и планирование общего использования

- Правила TTL
 Правила TTL для отдельных ресурсов, групп и бюджетов
- Политики тегов и автоматическое назначение ресурсов Установление и управление правилами тегов и автоматическое назначение ресурсов группам или бюджетам
- Простота использования
 Управление TTL и другими параметрами ресурсов через Slack



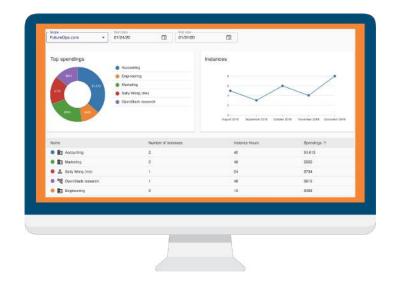
Прозрачность и оптимизация облачных затрат

• Бюджеты и автоматическое назначение ресурсов

Отслеживание затрат каждого бизнес-юнита, команды, пользователя или приложения

- Обнаружение аномалий в затратах Выявление пиков и мгновенные оповещения
- Глубокий анализ затрат
 История расходов и бюджетов на ресурсы
- Прогнозы бюджетов, пороговые значения и аналитика

Получение аналитики, прогнозов и рекомендаций по оптимизации



Внедрение FinOps процесса

- Продукт для построения FinOps
 Видимость, оптимизация, контроль и сотрудничество
- **Вовлечение инженеров**Члены команды несут ответственность за свои ресурсы, TTL и расходы в облаке
- Нет необходимости в внедрении новых инструментов. Используйте Slack Сценарии Уничтожения, Уведомления, Уведомить и Уничтожить



Текущие клиенты и <u>партнеры</u>

Telecoms

Integrators

Cloud service providers

Cloud vendors

End customers































vorbi















AIRBUS



















Контакты

□ sales@hystax-team.ru



() +7 495 204 28 77

9 123112, Москва, Пресненская наб., д.12, эт/пом/ком/45/1, 2/82

Приложение

Сложности в управлении затратами на облако. Почему FinOps?

• Инженеры не занимаются процессами оптимизации затрат.

Получение длинного списка сценариев оптимизации не помогает, так как несколько сотрудников DevOps или руководители ИТ-команд не могут решить все проблемы оптимизации без связи с владельцами ресурсов, у которых есть другие приоритеты. В результате реализуется только 20-30% рекомендаций.

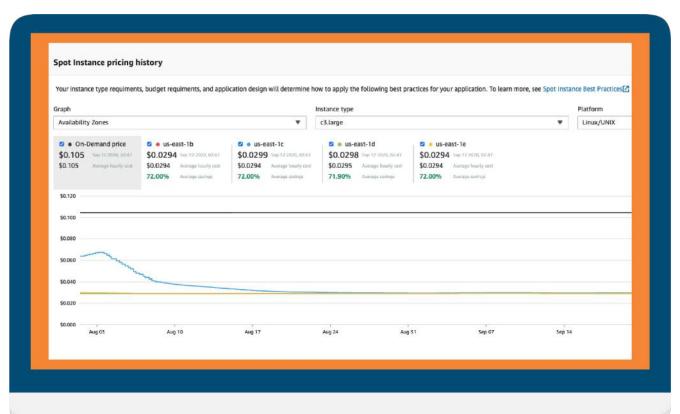
• Отсутствие управления жизненным циклом ресурсов

Инструменты управления затратами в облаке не предоставляют возможности управления жизненным циклом ресурсов.

• Отсутствие прозрачности и гибкости

Инструменты, разработанные облачными провайдерами, не обеспечивают достаточной детализации и прозрачности в рамках бюджетов, команд, кластеров и приложений.

Spot Instance - диапозон цен



Что нового? (Май 2023)

- Прозрачность и визуализация Reserved Instances/Saving Plans
- Обнаружение аномалий и ограничений
- Профилирование и оптимизация ML или любый других приложений
- MLOps
- Масштабируемость и повышение производительности

Roadmap (версии с открытым исходным кодом и SaaS)

- Улучшение рекомендаций по RI/SP/Spot
- Оптимизация с учетом RAM и GPU
- S3 дубликаты, тиринг, профайлинг

