

Идеальная «песочница» для ML- моделей:

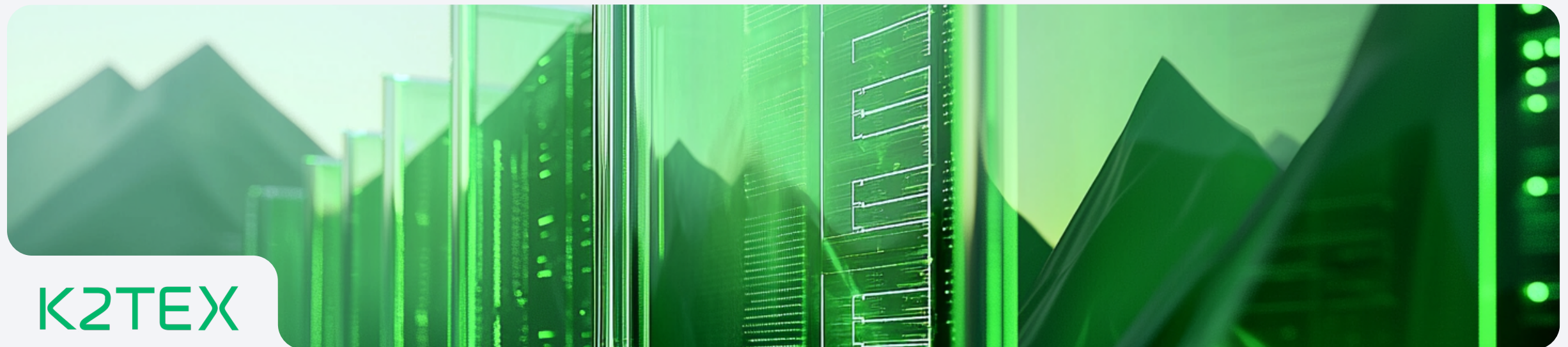
настраиваем контейнеризацию
без стресса



Даниил Салман

Tech Lead по практике
контейнеризации
K2 Tech

K2TEX



K2TEX

Эксперты цифровизации бизнеса

18

лет на рынке

1 500+

ИТ-профессионалов

250+

ведущих инженеров,
senior-разработчиков
и экспертов

1 200+

проектов в год
в сфере ИТ-услуг

10+

лет — средний стаж
работы специалистов

9

собственных
ИТ-продуктов

Помогаем клиентам достигать бизнес-целей, используя ИТ-инструменты. Внедряем технологические решения российских и зарубежных производителей, а также разрабатываем специализированные цифровые продукты.

Сегодня разберём:

С какими болями сталкиваются ML-команды каждый день

Почему «просто k8s-кластер» не рабочее решение

Какой стек нужен, чтобы собрать крепкий скелет ML-платформы

Практика: от запуска обучения через Airflow до деплоя и инференса модели в k8s

Сегодня разберём:

Кому будет интересно:

ML-инженеры, Data Scientists, архитекторы
и инфраструктурные команды

Что стоит знать заранее:

Базовое понимание Kubernetes, Docker
и Python-фреймворков

Что вы получите к концу доклада:

Разберётесь, как собрать свою ML-песочницу
без боли с драйверами и CUDA

Прогоните DAG в Airflow для обучения модели

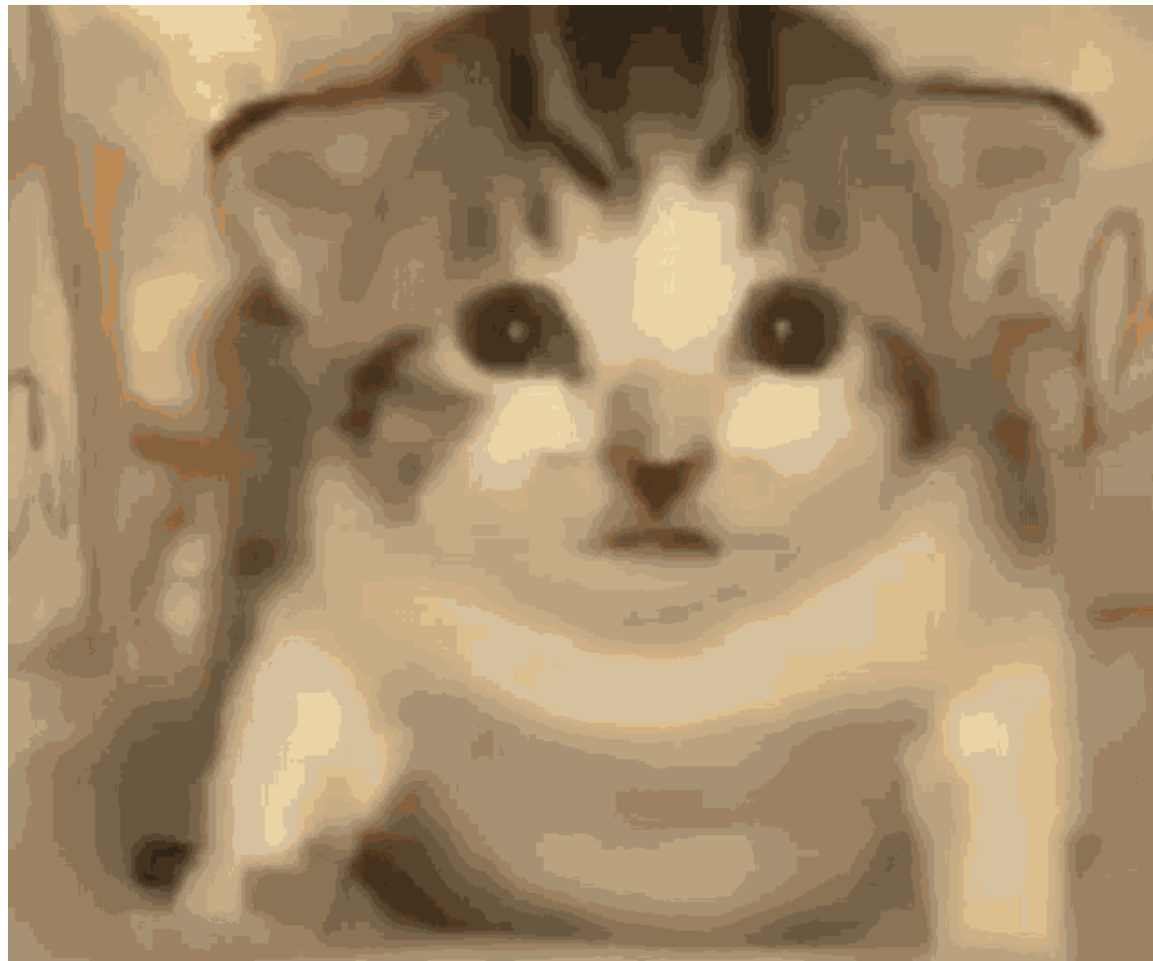
Научитесь централизованно логировать
метрики и артефакты моделей

Увидите, как модель превращается в inference-сервис

Немного истории Германа...

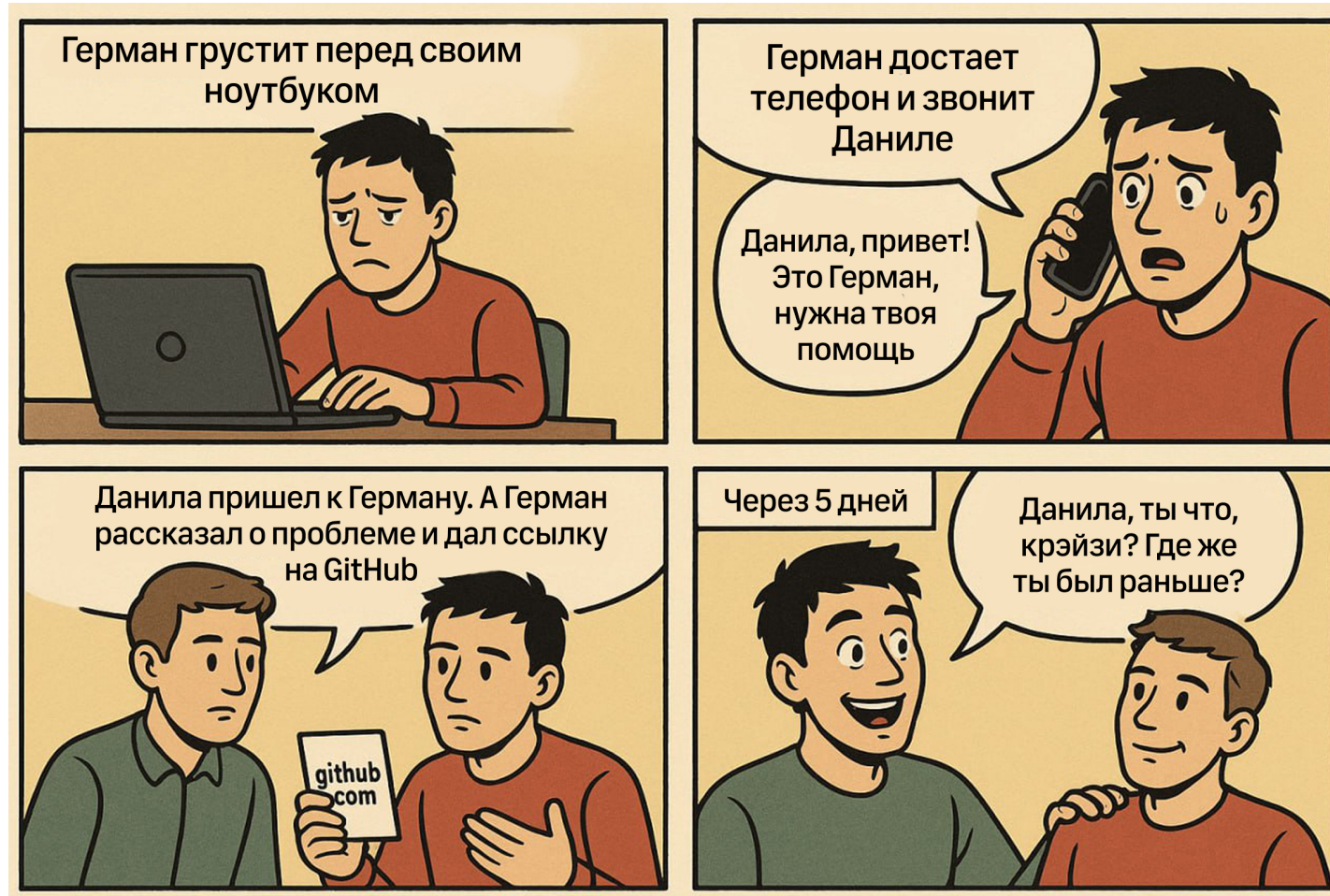


Герман:



source: <https://is.gd/sDxhKK>

Немного истории Германа...





source: Российское скетч-шоу «Даёшь молодёжь!»

С чем сталкиваются ML-инженеры

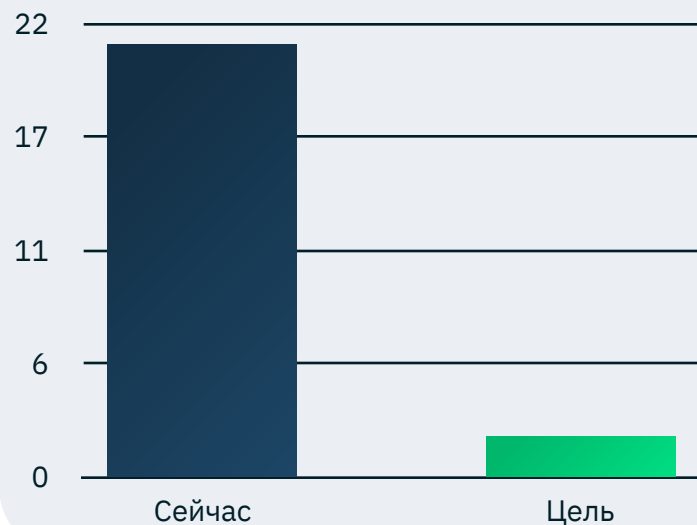
🕒 Time-to-first-training: 2–4 недели вместо 1–2 дней

🎮 GPU utilization: 30–40% вместо 70–80%

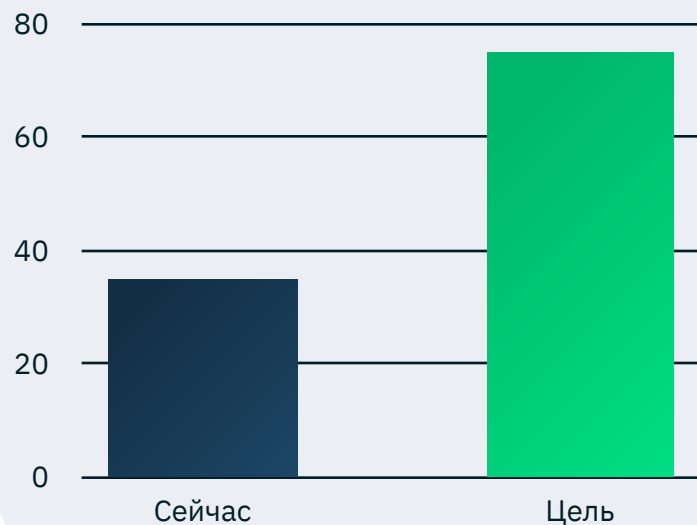
🔄 MTTR: дни вместо часов

⚠️ Каждая 3-я модель не воспроизводится в проде

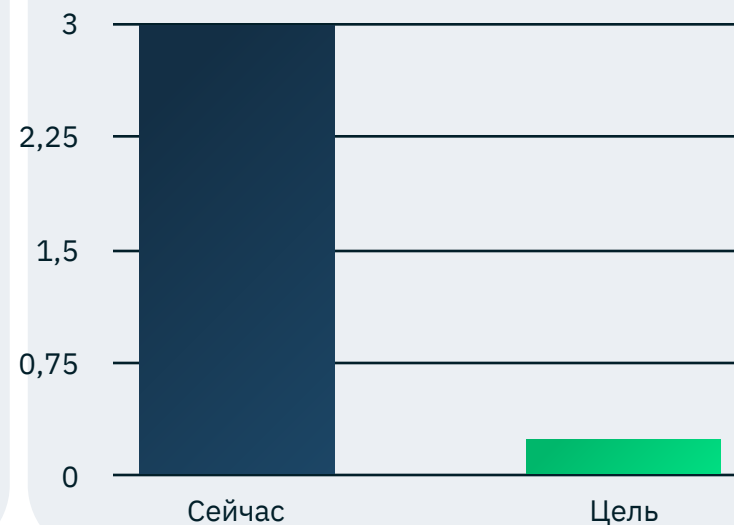
Time-to-first-training



GPU utilization



MTTR



Подходы к разворачиванию k8s

Способ	Для кого	Затраты	Контроль	Поддержка	Мониторинг	Плюсы	Минусы
kubeadm	Энтузиасты DevOps-гуру	Низкие	Максимальный	Сам	Свой	Максимальная гибкость, понимание изнутри	Много ручной работы
kuberspray	Рядовой DevOps- инженер	Средние	Высокий	Сообщество	Свой	Баланс гибкости и автоматизации, open source	Ansible
KaaS в облаке	Бизнес, готовый платить	Средние/ высокие	Средний	Облако	Облако	Удобно, быстро, поддержка	Нет всех сервисов “из коробки”, привязка к облаку
Вендорские решения	Бизнес, готовый платить	Высокие	Средний	Вендор	Встроенный	Интеграции, поддержка, все сервисы “с ходу”, быстрый старт	Стоимость, закрытость решений

Technology stack kubeadm

Observability	Security	Networking	Storage	Container Runtime
	RBAC PKI			
Service Discovery & Coordination	Service Proxy	Service Mesh	Backup	OS
CoreDNS etcd	kube-proxy			Any OS
Conf. Management	Infrastructure providers	UI	K8S services	Scaling
			kube-apiserver controller-manager scheduler etcd kubelet	

Technology stack kubespray

Observability	Security	Networking	Storage	Container Runtime
	RBAC PKI	<ul style="list-style-type: none">• Calico• Cilium• Flannel		<ul style="list-style-type: none">• Containerd (default)• Cri-o• Docker
Service Discovery & Coordination	Service Proxy	Service Mesh	Backup	OS
CoreDNS etcd	kube-proxy			<ul style="list-style-type: none">• Из коробки нет поддержки российских ОС
Conf. Management	Infrastructure providers	UI	K8S services	Scaling
			kube-apiserver controller-manager scheduler etcd kubelet	

Technology stack NOVA

Observability	Security	Networking	Storage	Container Runtime
<ul style="list-style-type: none"> • Prometheus • Grafana • Alertmanager • Thanos Query • Logging Operator • Opensearch • Fluentd 	<ul style="list-style-type: none"> • RBAC • Neuvector • Secrets Webhook • CertManager • Secrets Manager (CSI, PKI, OAuth) 	Calico Cilium	<ul style="list-style-type: none"> • Longhorn 	Containerd
Service Discovery & Coordination	Service Proxy	Service Mesh	Backup	OS
CoreDNS etcd	kube-proxy		<ul style="list-style-type: none"> • Velero 	<ul style="list-style-type: none"> • RED OS • AlmaLinux • Astra Linux • RHEL-based
Conf. Management	Infrastructure providers	UI	K8S services	Scaling
<ul style="list-style-type: none"> • Gitea • FluxCD • Authomation tools 	<ul style="list-style-type: none"> • Bare Metal • oVirt • zVirt • VMware vSphere 	<ul style="list-style-type: none"> • Nova Console 	kube-apiserver controller-manager scheduler etcd kubelet	<ul style="list-style-type: none"> • VPA • HPA

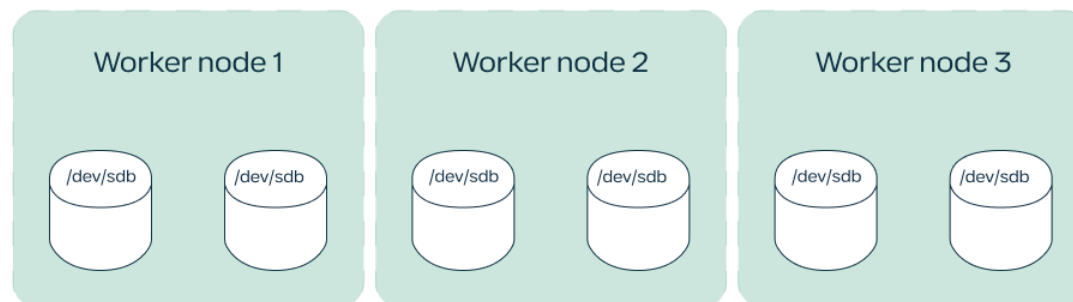


А что же есть в нашей песочнице?

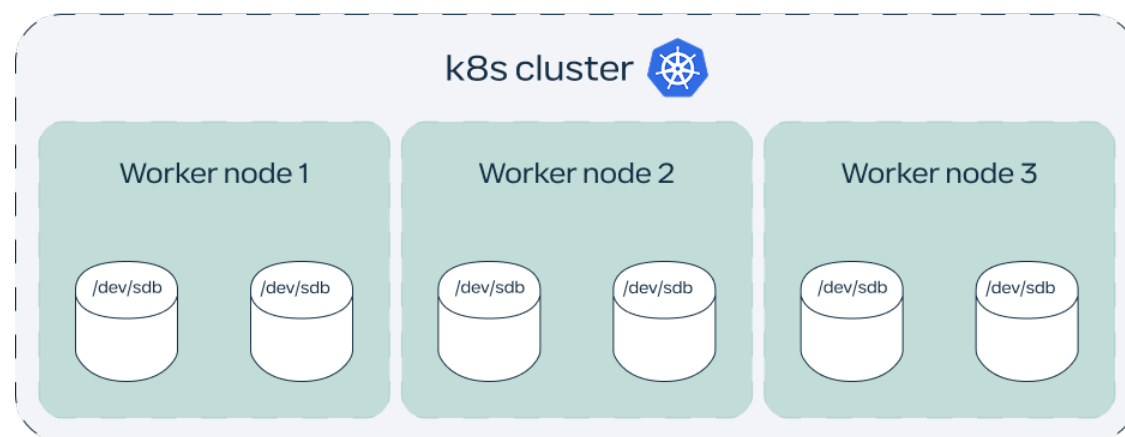
Аппаратный уровень

- NVIDIA GPU Operator (MIG / Time-slices)
- NVIDIA Network Operator (опционально)
- Longhorn

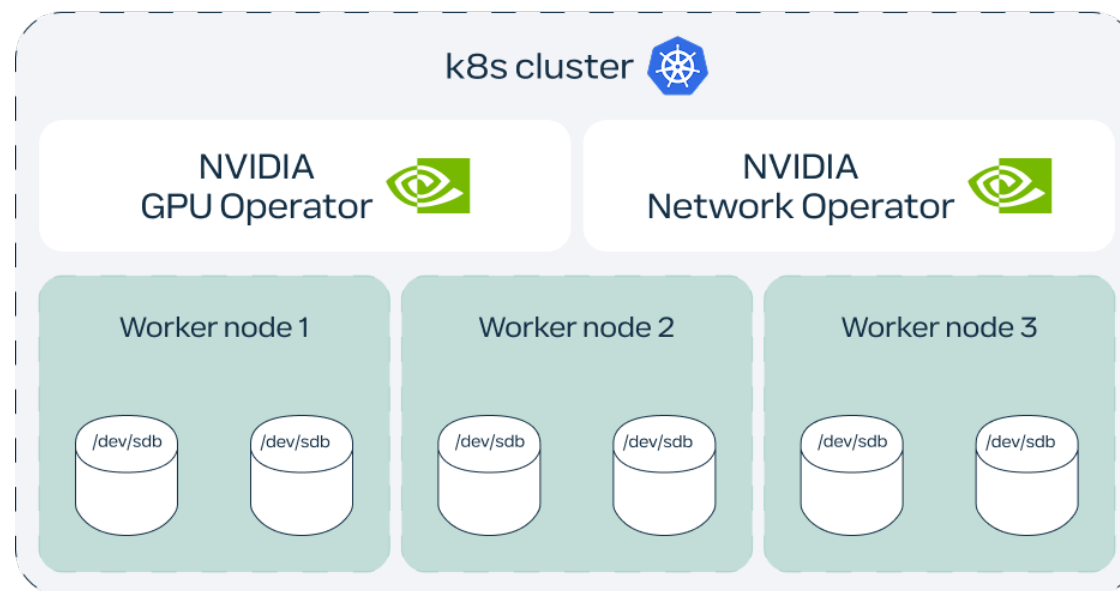
А что же есть в нашей песочнице?



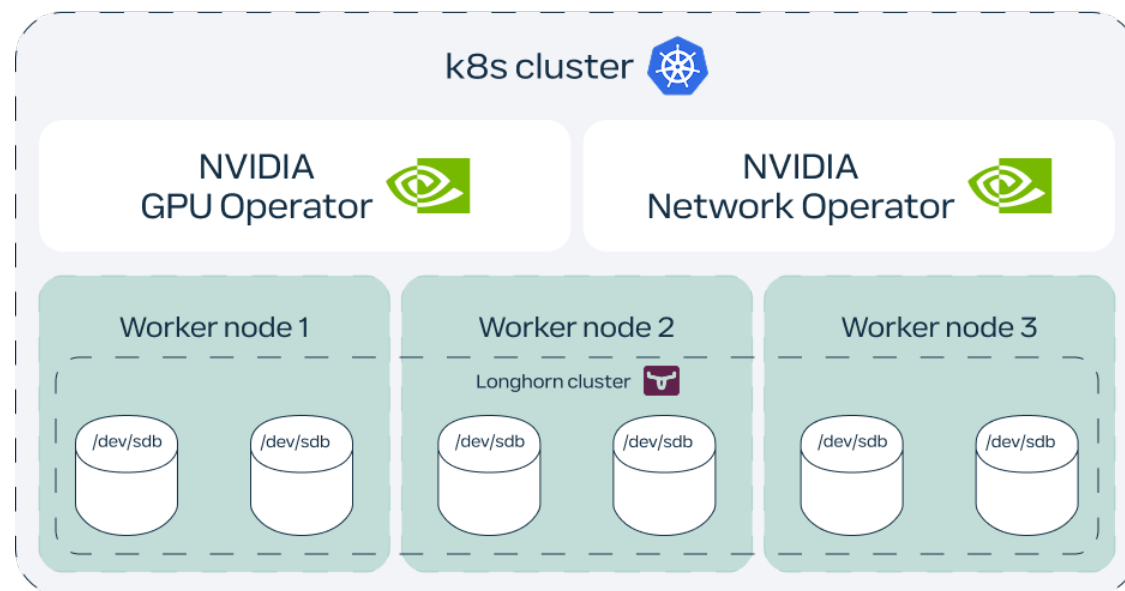
А что же есть в нашей песочнице?



А что же есть в нашей песочнице?



А что же есть в нашей песочнице?



А что же есть в нашей песочнице?

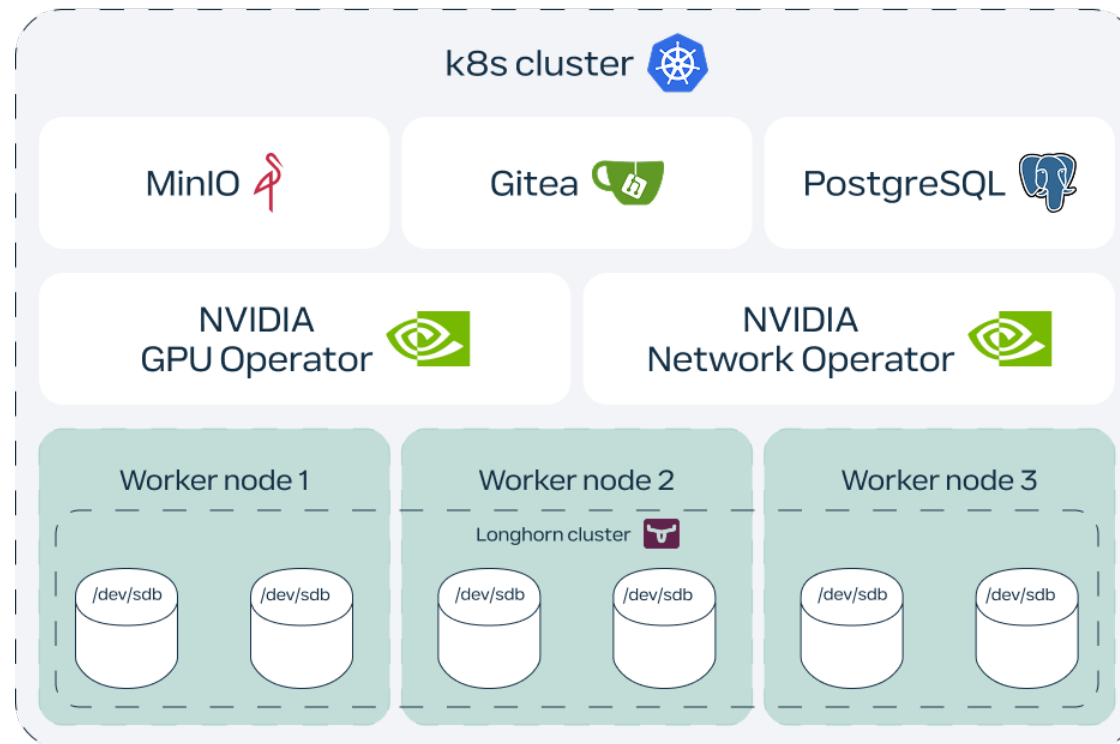
Аппаратный уровень

- NVIDIA GPU Operator (MIG / Time-slices)
- NVIDIA Network Operator (опционально)
- Longhorn

Хранилища

- MinIO
- PostgreSQL
- Gitea

А что же есть в нашей песочнице?



А что же есть в нашей песочнице?

Аппаратный уровень

- NVIDIA GPU Operator (MIG / Time-slices)
- NVIDIA Network Operator (опционально)
- Longhorn

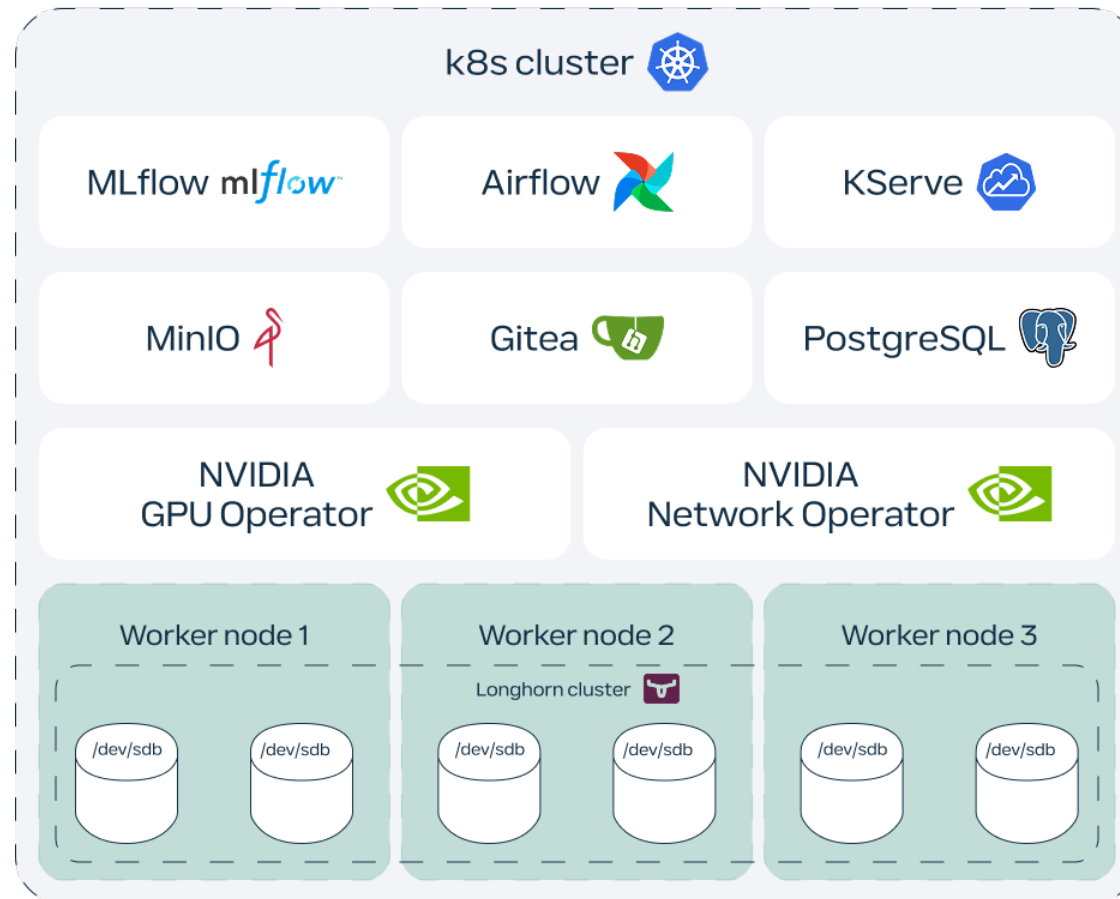
Хранилища

- ML-инженеры
- Data Scientists
- Архитекторы и инфраструктурные команды

Сервисный слой

- Mlflow
- Airflow
- KServe

А что же есть в нашей песочнице?



А что же есть в нашей песочнице?

Аппаратный уровень

- NVIDIA GPU Operator (MIG / Time-slices)
- NVIDIA Network Operator (опционально)
- Longhorn

Хранилища

- ML-инженеры
- Data Scientists
- Архитекторы и инфраструктурные команды

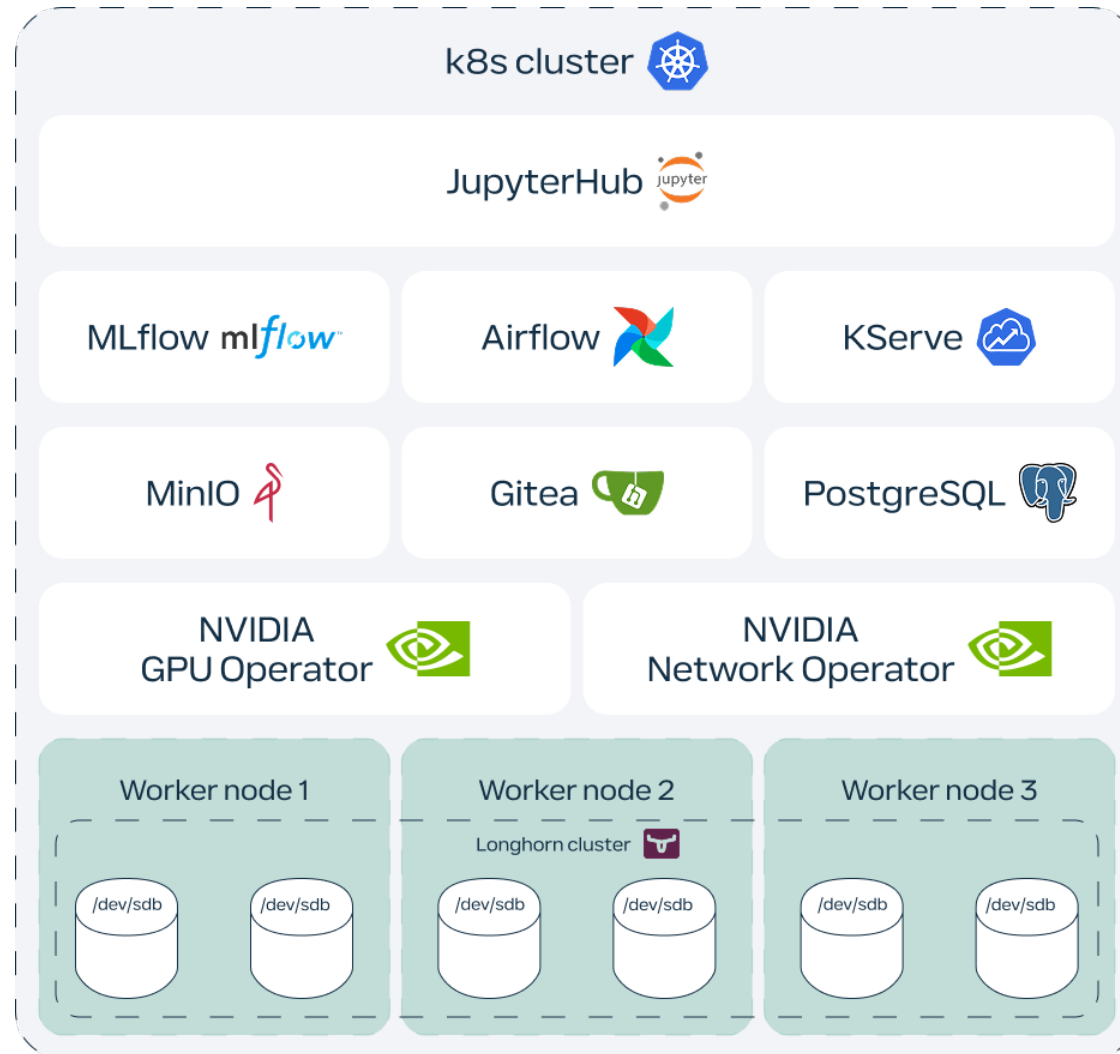
Сервисный слой

- Mlflow
- Airflow
- KServe

Интерфейс

- JupyterHub

А что же есть в нашей песочнице?



А что же есть в нашей песочнице?

Аппаратный уровень

- NVIDIA GPU Operator (MIG / Time-slices)
- NVIDIA Network Operator (опционально)
- Longhorn

Хранилища

- ML-инженеры
- Data Scientists
- Архитекторы и инфраструктурные команды

Сервисный слой

- Mlflow
- Airflow
- KServe

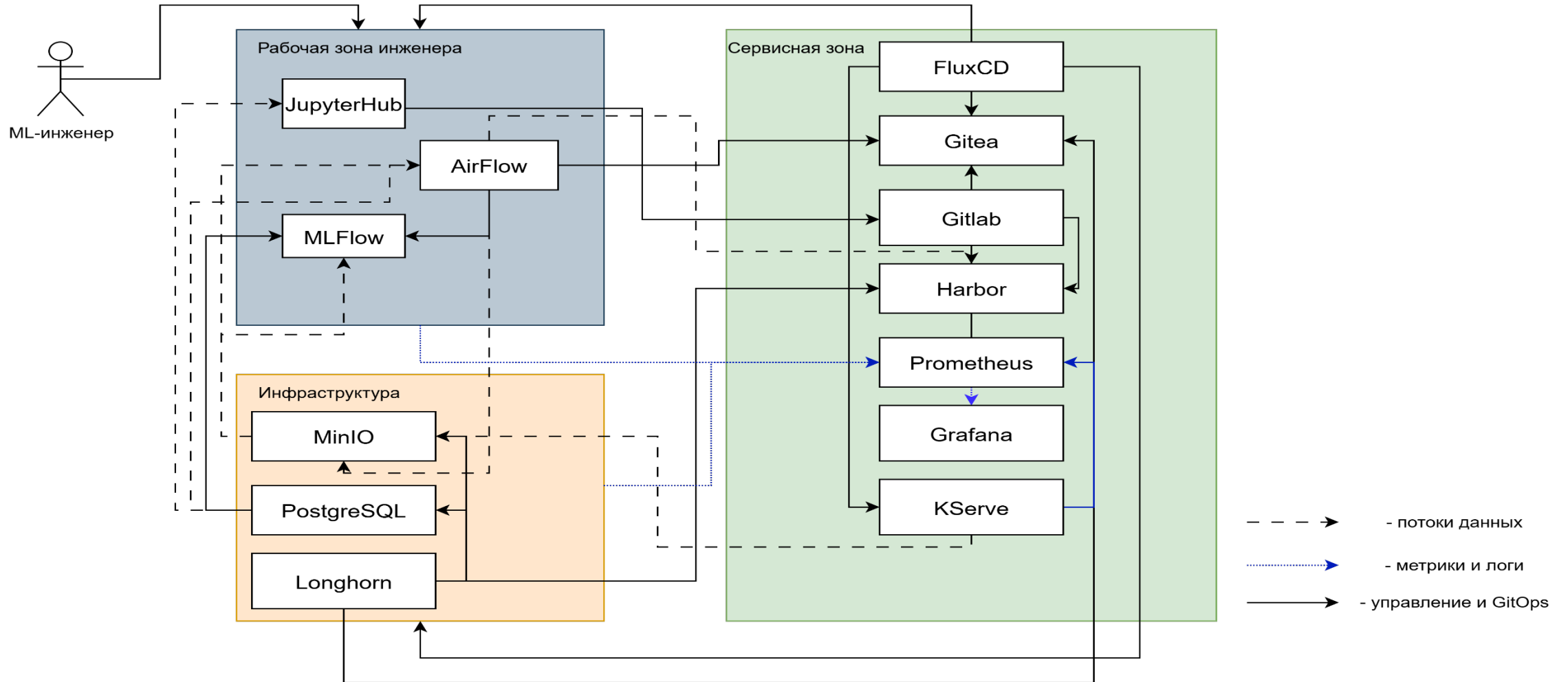
Интерфейс

- JupyterHub

Интеграции и автоматизация

- Harbor/Nexus
- GitLab

Архитектурная схема ML-платформы



Юху, практика!

Инфраструктура:

Платформа

k8s - Nova Container Platform

Воркеры

Два сервера Dell с NVIDIA T4

Внешний GitLab сервер

Задача

Классифицировать изображение

Датасет

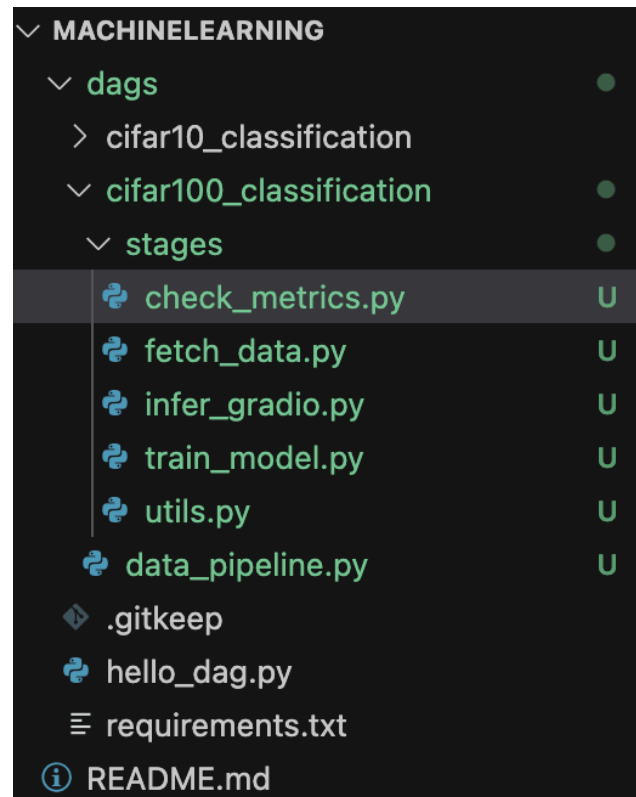
CIFAR100

Цель

Бесшовно пройти весь процесс ML, не думая о бэкенде, где всё это работает (на самом деле я хочу узнать, как классифицирует моего кота ИИ)

Юху, практика!

PUSH!





Юху, практика!






MachineLearning / Repository

 main ▼ machinelearning / dags / + ▼


Find file Edit ▼ Code ▼

 **cifar100 model added**
dsalman authored just now

2df40664  History

Name	Last commit	Last update
..		
 cifar100_classification	cifar100 model added	just now
 cifar10_classification	Fixed max val loss	2 weeks ago
 .gitkeep	Add new directory	3 weeks ago
 hello_dag.py	Add new file	3 weeks ago
 requirements.txt	Test	2 weeks ago

Юху, практика!

 Airflow

DAGs

Cluster Activity



Datasets

Security

Browse

Admin

Docs

 13:18 UTC  AU

Usage of a dynamic webserver secret key detected. We recommend a static webserver secret key instead. See the Helm Chart Production Guide for more details.

DAGs

All 3Active 0Paused 3

Running 0Failed 0

Filter DAGs by tag

Search DAGs

☒ Auto-refresh

<<

<


1

>


>>

Showing 1-3 of 3 DAGs

Юху, практика!



[DAGs](#) [Cluster Activity](#) [Datasets](#) [Security](#) [Browse](#) [Admin](#) [Docs](#)

 18:47 UTC AU

Usage of a dynamic webserver secret key detected. We recommend a static webserver secret key instead. See the Helm Chart Production Guide for more details.


DAGs










All 3 Active 1 Paused 2

Running 1 Failed 0

Filter DAGs by tag

Search DAGs

☒ Auto-refresh 

DAG	Owner	Runs	Schedule	Last Run	Next Run	Recent Tasks	Actions	Links
<input checked="" type="checkbox"/> cifar100_classification	airflow	<div><div></div><div></div><div>1</div><div></div></div>	None	2025-06-09, 14:25:24		<div><div>3</div><div></div><div></div><div></div><div></div><div>1</div><div>1</div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div></div>	<div></div> ...	
<input type="checkbox"/> cifar10_classification	airflow	<div><div></div><div></div><div></div><div></div></div>	None			<div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div></div>	<div></div> ...	
<input type="checkbox"/> hello_git_dag git	airflow	<div><div></div><div></div><div></div><div></div></div>	@daily		2025-06-08, 00:00:00	<div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div></div>	<div></div> ...	

«

<

1

>

»

Showing 1-3 of 3 DAGs

Юху, практика!

Namespace: airflow

Pods

Создать Pod

Фильтр

ИмяПоиск по имени...

Имя	Статус	Готовность	Перезапуски	Владелец	Память	CPU	Дата создания
airflow-scheduler-5c54648fd4-vrld5	Запущен	3/3	0	airflow-scheduler-5c54648fd4	3 204,0 MiB	0,133 ядра	26 мая 2025 г., 15:40
airflow-triggerer-0	Запущен	2/2	0	airflow-triggerer	2 888,9 MiB	0,061 ядра	26 мая 2025 г., 15:40
airflow-postgresql-0	Запущен	1/1	0	airflow-postgresql	76,9 MiB	0,019 ядра	26 мая 2025 г., 15:40
airflow-statsd-6bdddbb65b-7x64s	Запущен	1/1	0	airflow-statsd-6bdddbb65b	23,3 MiB	0,003 ядра	26 мая 2025 г., 15:40
airflow-webserver-c4cbb5766-mgx7s	Запущен	1/1	0	airflow-webserver-c4cbb5766	1 952,2 MiB	0,003 ядра	26 мая 2025 г., 15:40
cifar100-classification-data-preparation-7z2p2kzi	Запущен	1/1	0	Владелец не указан	-	-	9 июня 2025 г., 17:25

Юху, практика!

MINIO

OBJECT STORE

LICENSE

User

Object Browser

Access Keys

Documentation

Administrator

Buckets

Policies

Identity

Monitoring

Object Browser

Start typing to filter objects in the bucket

data

Created on: Thu, May 15 2025 21:43:07 (GMT+3) Access: PRIVATE 679.0 MiB - 14 Objects

Rewind

Refresh

Upload

< data

Create new path

<input type="checkbox"/>	Name	Last Modified	Size
<input type="checkbox"/>	cifar10		-
<input type="checkbox"/>	cifar100		-

Юху, практика!

Namespace: airflow ▾

Pods > Информация

P cifar100-classification-data-preparation-7z2p2kzi Запущен

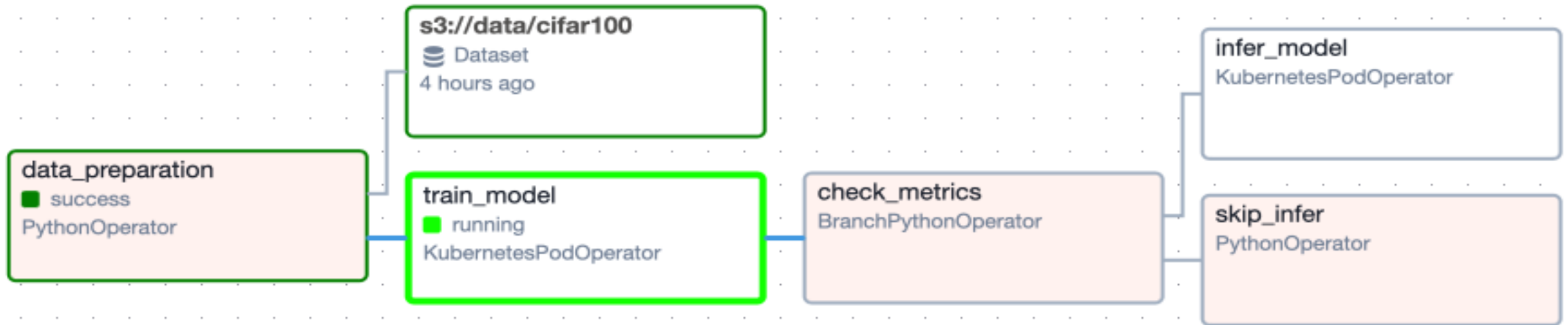
Подобности Метрики YAML Переменные окружения Логи События Терминал

Поток запущен... base ▾ Текущий лог ▾ ☐ Перенос строк [Открыть в i](#)

85 строка

```
57 Downloading nvidia_nvtx_cu12-12.6.77-py3-none-manylinux2014_x86_64.manylinux_2_17_x86_64.whl.metadata (1.6 kB)
58 Collecting nvidia_nvjitlink-cu12==12.6.85 (from torch==2.7.1->torchvision)
59 Downloading nvidia_nvjitlink_cu12-12.6.85-py3-none-manylinux2010_x86_64.manylinux_2_12_x86_64.whl.metadata (1.5 kB)
60 Collecting nvidia-cuffile-cu12==1.11.1.6 (from torch==2.7.1->torchvision)
61 Downloading nvidia_cuffile_cu12-1.11.1.6-py3-none-manylinux2014_x86_64.manylinux_2_17_x86_64.whl.metadata (1.5 kB)
62 Collecting triton==3.3.1 (from torch==2.7.1->torchvision)
63 Downloading triton-3.3.1-cp312-cp312-manylinux_2_27_x86_64.manylinux_2_28_x86_64.whl.metadata (1.5 kB)
64 Requirement already satisfied: botocore<1.37.0,>=1.36.3 in /home/airflow/.local/lib/python3.12/site-packages (from boto3) (1.36.3)
65 Requirement already satisfied: jmespath<2.0.0,>=0.7.1 in /home/airflow/.local/lib/python3.12/site-packages (from boto3) (0.10.0)
66 Requirement already satisfied: s3transfer<0.12.0,>=0.11.0 in /home/airflow/.local/lib/python3.12/site-packages (from boto3) (0.11.2)
67 Requirement already satisfied: python-dateutil<3.0.0,>=2.1 in /home/airflow/.local/lib/python3.12/site-packages (from botocore<1.37.0,>=1.36.3->boto3) (2.9.0.post0)
68 Requirement already satisfied: urllib3<2.2.0,>=1.25.4 in /home/airflow/.local/lib/python3.12/site-packages (from botocore<1.37.0,>=1.36.3->boto3) (2.3.0)
69 Requirement already satisfied: six>=1.5 in /home/airflow/.local/lib/python3.12/site-packages (from python-dateutil<3.0.0,>=2.1->botocore<1.37.0,>=1.36.3->boto3) (1.17.0)
70 Collecting mpmath<1.4,>=1.1.0 (from sympy==1.13.3->torch==2.7.1->torchvision)
71 Downloading mpmath-1.3.0-py3-none-any.whl.metadata (8.6 kB)
72 Requirement already satisfied: MarkupSafe>=2.0 in /home/airflow/.local/lib/python3.12/site-packages (from Jinja2->torch==2.7.1->torchvision) (3.0.2)
73 Downloading torchvision-0.22.1-cp312-cp312-manylinux_2_28_x86_64.whl (7.5 MB)
74 _____ 7.5/7.5 MB 69.8 MB/s eta 0:00:00
75 Downloading torch-2.7.1-cp312-cp312-manylinux_2_28_x86_64.whl (821.0 MB)
76 _____ 821.0/821.0 MB 36.3 MB/s eta 0:00:00
77 Downloading nvidia_cublas_cu12-12.6.4.1-py3-none-manylinux2014_x86_64.manylinux_2_17_x86_64.whl (393.1 MB)
78 _____ 393.1/393.1 MB 84.4 MB/s eta 0:00:00
79 Downloading nvidia_cuda_cupti_cu12-12.6.80-py3-none-manylinux2014_x86_64.manylinux_2_17_x86_64.whl (8.9 MB)
80 _____ 8.9/8.9 MB 125.3 MB/s eta 0:00:00
81 Downloading nvidia_cuda_nvrtc_cu12-12.6.77-py3-none-manylinux2014_x86_64.whl (23.7 MB)
82 _____ 23.7/23.7 MB 141.0 MB/s eta 0:00:00
83 Downloading nvidia_cuda_runtime_cu12-12.6.77-py3-none-manylinux2014_x86_64.manylinux_2_17_x86_64.whl (897 kB)
84 _____ 897.7/897.7 kB 75.5 MB/s eta 0:00:00
85 Downloading nvidia_cudnn_cu12-9.5.1.17-py3-none-manylinux_2_28_x86_64.whl (571.0 MB)
```

Юху, практика!



Юху, практика!

Namespace: airflow

Pods

Создать Pod

ФильтрИмяПоиск по имени...

Имя	Статус	Готовность	Перезапуски	Владелец	Память	CPU	Дата создания
airflow-postgresql-0	Запущен	1/1	0	airflow-postgresql	67,6 MiB	0,016 ядра	11 июня 2025 г., 18:41
airflow-scheduler-66dfc7bbdf-64jz8	Запущен	3/3	0	airflow-scheduler-66dfc7bbdf	1 683,2 MiB	0,126 ядра	11 июня 2025 г., 21:35
airflow-statsd-6bdddbb65b-g5qns	Запущен	1/1	0	airflow-statsd-6bdddbb65b	24,2 MiB	0,003 ядра	11 июня 2025 г., 18:41
airflow-triggerer-0	Запущен	2/2	0	airflow-triggerer	1 342,9 MiB	0,055 ядра	11 июня 2025 г., 21:35
airflow-webserver-66b7f5d584-hl8nt	Запущен	1/1	0	airflow-webserver-66b7f5d584	2 130,7 MiB	0,002 ядра	11 июня 2025 г., 21:35
cifar100-classification-train-model-992gx6mk	Запущен	1/1	0	Владелец не указан	-	-	16 июня 2025 г., 12:16
train-model-vrpwr2cm	Запущен	1/1	0	Владелец не указан	-	-	16 июня 2025 г., 12:17

Юху, практика!

Namespace: airflow

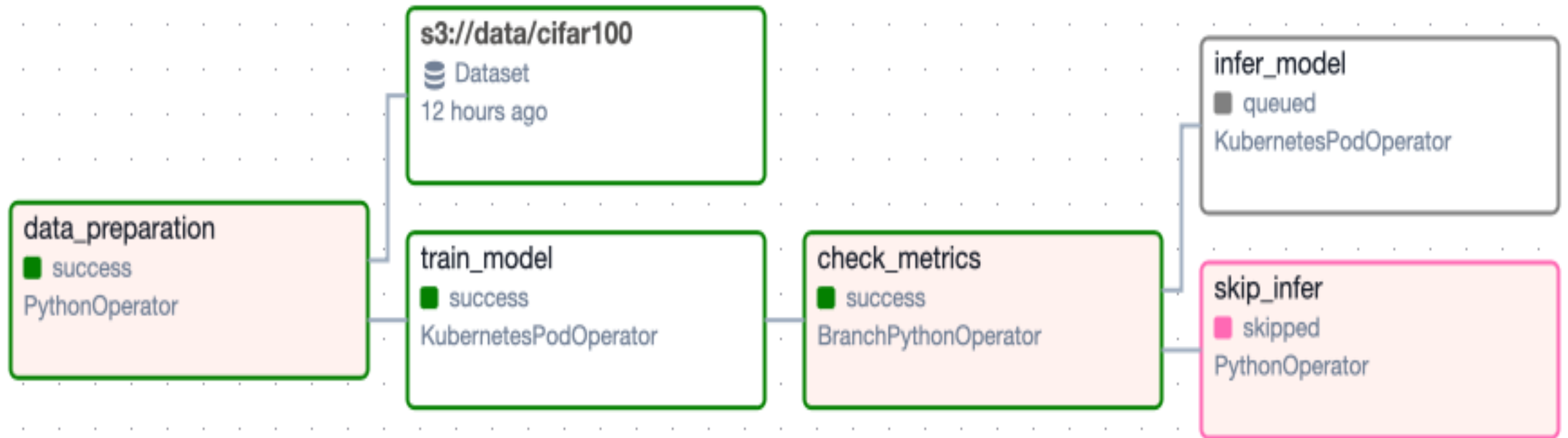
Pods

Создать Pod

ФильтрИмяПоиск по имени...

Имя	Статус	Готовность	Перезапуски	Владелец	Память	CPU	Дата создания
airflow-postgresql-0	Запущен	1/1	0	airflow-postgresql	71,7 MiB	0,018 ядра	11 июня 2025 г., 18:41
airflow-scheduler-66dfc7bbdf-64jz8	Запущен	3/3	0	airflow-scheduler-66dfc7bbdf	1 680,8 MiB	0,131 ядра	11 июня 2025 г., 21:35
airflow-statsd-6bdddbb65b-g5qns	Запущен	1/1	0	airflow-statsd-6bdddbb65b	23,6 MiB	0,003 ядра	11 июня 2025 г., 18:41
airflow-triggerer-0	Запущен	2/2	0	airflow-triggerer	1 340,3 MiB	0,053 ядра	11 июня 2025 г., 21:35
airflow-webserver-66b7f5d584-hl8nt	Запущен	1/1	0	airflow-webserver-66b7f5d584	1 717,8 MiB	0,003 ядра	11 июня 2025 г., 21:35
cifar100-classification-check-metrics-q8k4vcov	Запущен	1/1	0	Владелец не указан	-	-	16 июня 2025 г., 11:53

Юху, практика!



Юху, практика!

Манифест InferenceService

Дополнительно нужно создать две сущности:

- secret с s3-credentials
- service account с привязкой к этому секрету

```
1  apiVersion: serving.kserve.io/v1beta1
2  kind: InferenceService
3  metadata:
4    name: resnet18
5    namespace: ml-inference
6    annotations:
7      serving.kserve.io/s3-endpoint: "minio.svc.cluster.local"
8      serving.kserve.io/s3-credentials-secret-name: "s3-cred"
9  spec:
10   predictor:
11     serviceAccountName: resnet18-sa
12     triton:
13       runtimeVersion: "23.10-py3"
14       storageUri: "s3://data/models"
15       resources:
16         requests:
17           cpu: "500m"
18           memory: "1Gi"
19           nvidia.com/gpu: 1
20         limits:
21           cpu: "1"
22           memory: "2Gi"
23           nvidia.com/gpu: 1
24       nodeSelector:
25         nvidia.com/gpu.present: "true"
26       tolerations:
27         - key: "nvidia.com/gpu"
28           operator: "Exists"
29           effect: "NoSchedule"
30
```


Юху, практика!

Namespace: airflow

Pods

Создать Pod

Фильтр

ИмяПоиск по имени...

Имя	Статус	Готовность	Перезапуски	Владелец	Память	CPU	Дата создания
airflow-postgresql-0	Запущен	1/1	0	airflow-postgresql	72,0 MiB	0,018 ядра	11 июня 2025 г., 18:41
airflow-scheduler-66dfc7bbdf-64jz8	Запущен	3/3	0	airflow-scheduler-66dfc7bbdf	1 680,7 MiB	0,139 ядра	11 июня 2025 г., 21:35
airflow-statsd-6bddd5b65b-g5qns	Запущен	1/1	0	airflow-statsd-6bddd5b65b	23,7 MiB	0,003 ядра	11 июня 2025 г., 18:41
airflow-triggerer-0	Запущен	2/2	0	airflow-triggerer	1 340,4 MiB	0,053 ядра	11 июня 2025 г., 21:35
airflow-webserver-66b7f5d584-hl8nt	Запущен	1/1	0	airflow-webserver-66b7f5d584	1 719,6 MiB	0,003 ядра	11 июня 2025 г., 21:35
cifar100-classification-infer-model-pjs1hdbi	Запущен	1/1	0	Владелец не указан	-	-	16 июня 2025 г., 11:54
gradio-ui-mw0jfasw	Запущен	1/1	0	Владелец не указан	-	-	16 июня 2025 г., 11:55

Юху, практика!

CIFAR-100 Classifier

Upload an image and get top-5 prediction from a ResNet18 model trained on CIFAR-100

image



output

Flag

Clear


Submit

Юху, практика!

CIFAR-100 Classifier

Upload an image and get top-5 prediction from a ResNet18 model trained on CIFAR-100

image



output

flatfish

flatfish	81%
rocket	2%
house	1%
can	1%
elephant	1%

Flag

Clear

Submit

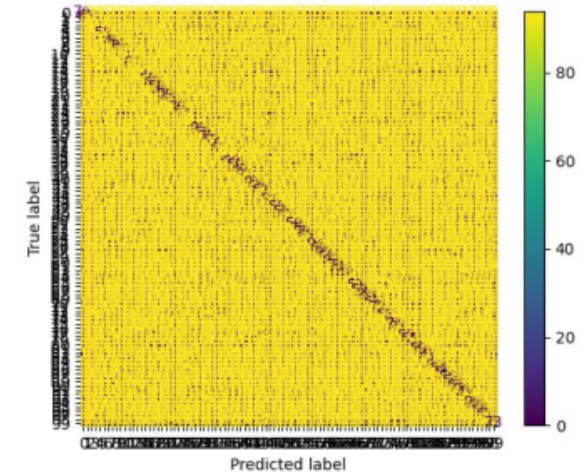
Юху, практика!

Metrics (6)

Search metrics	
Metric	Value
val_loss	85.10773599147797
val_accuracy	0.5678
train_loss	96.06210631132126
val_precision	0.5675951905553892
val_recall	0.5672150327612843
val_f1_score	0.5646417947926783

Parameters (6)

Search parameters	
Parameter	Value
dataset	cifar100
epochs	30
lr	0.001
model	ResNet18
batch_size	256
early_stopping_patience	20



✅ Обучение модели завершено!

🧠 Model: ResNet18
📅 Timestamp: 2025-06-16 19:52:16
🔄 Epochs: 30 | Batch: 256 | LR: 0.001
💻 Device: cuda

📊 Accuracy: 0.5678
📊 F1-score: 0.5646
🔄 Δ Accuracy: +0.0898
🔄 Δ F1-score: +0.0932

🔗 Открыть в MLflow
Run ID: cc80c2871f644d5fb051dc613faabd4b 23:03

Почему песочница «идеальная»?

Time-to-first-demo ↓ в 5–7 раз
(с недель до дней)

GPU-utilization ↑ 30–40% за счёт MIG/
Time-slicing и автоскейлинга

Time-to-production ↓ в 2–3 раза

Время обнаружения инцидента ↓
с часов до минут

Онбординг инженера ↓ с месяцев
до пары дней

Простой GPU ↓ ~50% → меньше
неиспользуемых карточек →
экономия бюджета

99% экспериментов воспроизводимы
через MLflow

Reuse пайплайнов и экспериментов
↑ 60–70%

А какие еще «игрушки» можно затащить в песочницу?

Schedulers

KAI, Volcano

Инференс / деплой моделей

Triton, Dynamo, KServe

Визуализация и анализ данных

Grafana, Superset, Evidently AI

Операторы для работы с видеокартами других вендоров

HAMi, ROCm

Безопасность

SSO

Подпись образов (Cosign / Notary)

SBOM (Syft)

Сканирование уязвимостей (Trivy)

Управление секретами (Vault)

Правила запуска контейнеров
(OPA Gatekeeper / Kyverno)

Настройка SSO для MLflow, Jupyterhub и Airflow



habr.com/ru/companies/orion_soft/articles/940592/

«Данила, план гениальный — как всегда!»
(с) Герман



Выводы

Платформа является кросс-доменной,
подходит для любой отрасли

Эффективное использование рабочего
времени и уменьшение ТРЗ на решение
смежных проблем

Контроль ресурсов, в том числе GPU —
больше не узкое место

Полноценный конвейер обкатки
ML-модели больше не головная боль

Мониторинг встроен в процесс и помогает отслеживать качество
и стабильность обкатывания моделей

K2TEX

Даниил Салман

Tech Lead по практике
контейнеризации K2 Tech

tg: @DV_Salman

Будьте в курсе новостей
K2 TEX

k2.tech
info@k2.tech
+7 (495) 797-85-84

