







黄敏杰 Daocloud云原生高级研发工程师

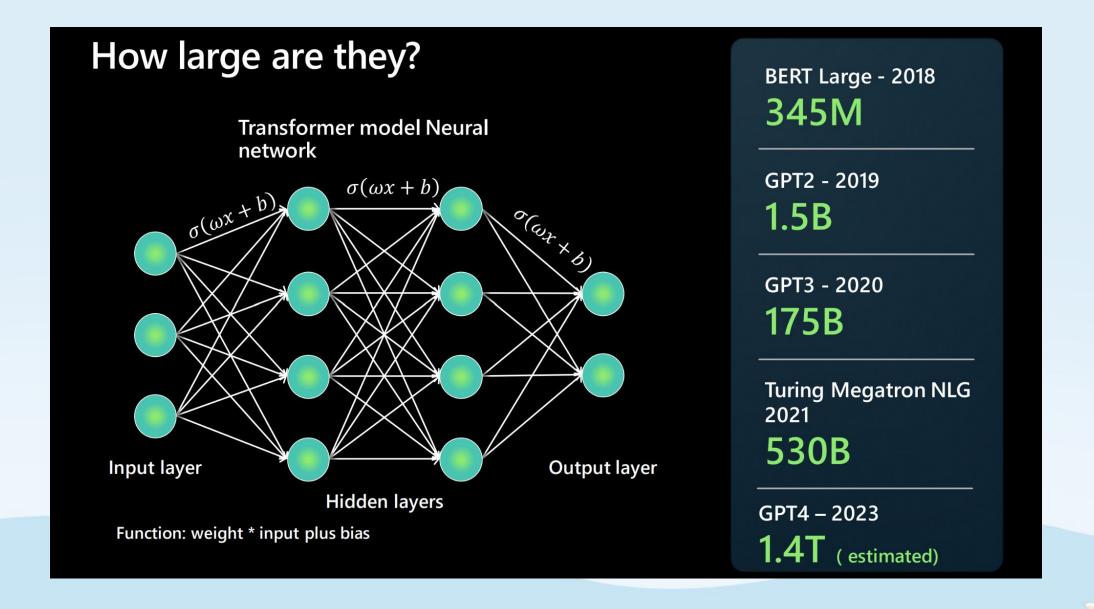


肃伦 Daocloud云原生高级研发工程师



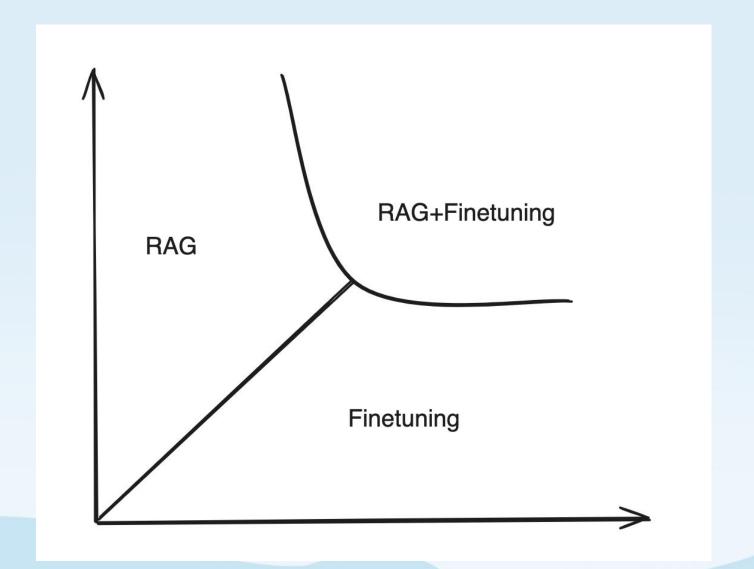
大模型LLM





大模型落地





大模型微调的挑战



1 高质量的数据集

微调需要大量与目标领域或任务相关的高质量数据,获取和构建这样的数据集是一个挑战。

2 超参数微调

微调的过程非常复杂,涉及控制优化过程的各种超参数,如学习率,批大小等。

3 模型评估

传统的评估指标并不能完全反应LLM的能力,如何精确评估LLM能力会是一个挑战。

DataTunerX



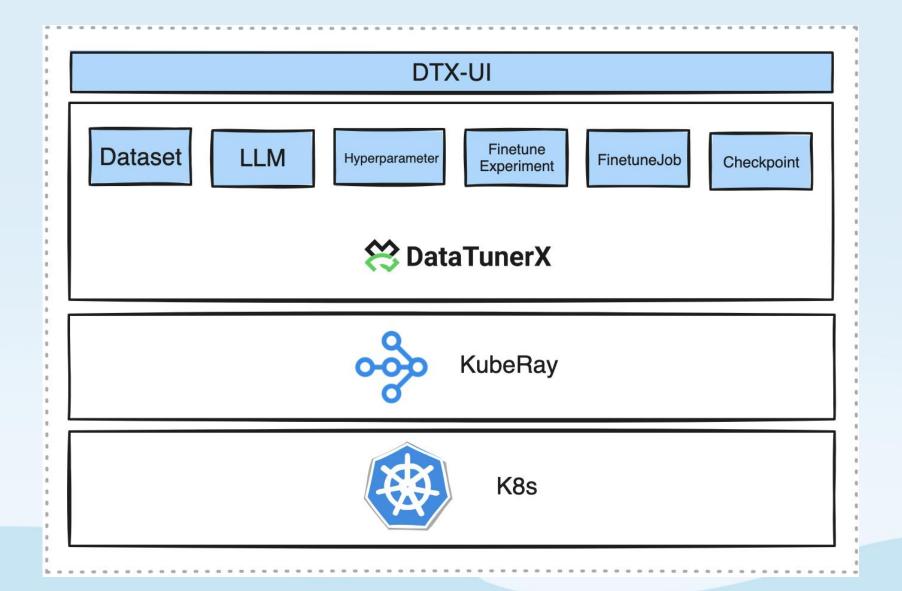


- 由DaoCloud和Futurewei 联合发起的一个基于云原生的大模型高效微调的项目
- 特点:
 - 易于使用,高效,快速一键微调
 - 并行,分布式,充分利用算力资源
 - 模型一键部署, 推理
 - 可自定义数据,评估插件

项目地址: https://github.com/DataTunerX/datatunerx

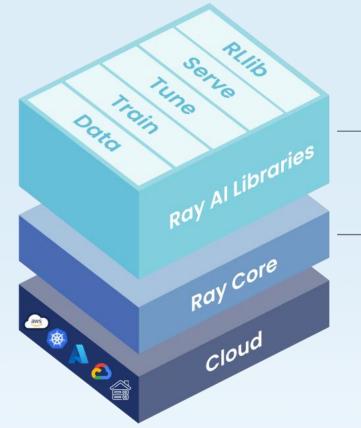
DTX架构





What is KubeRay?





high-level libraries that enable simple scaling of Al workloads

a low-level distributed computing framework with a concise core and Python-first API

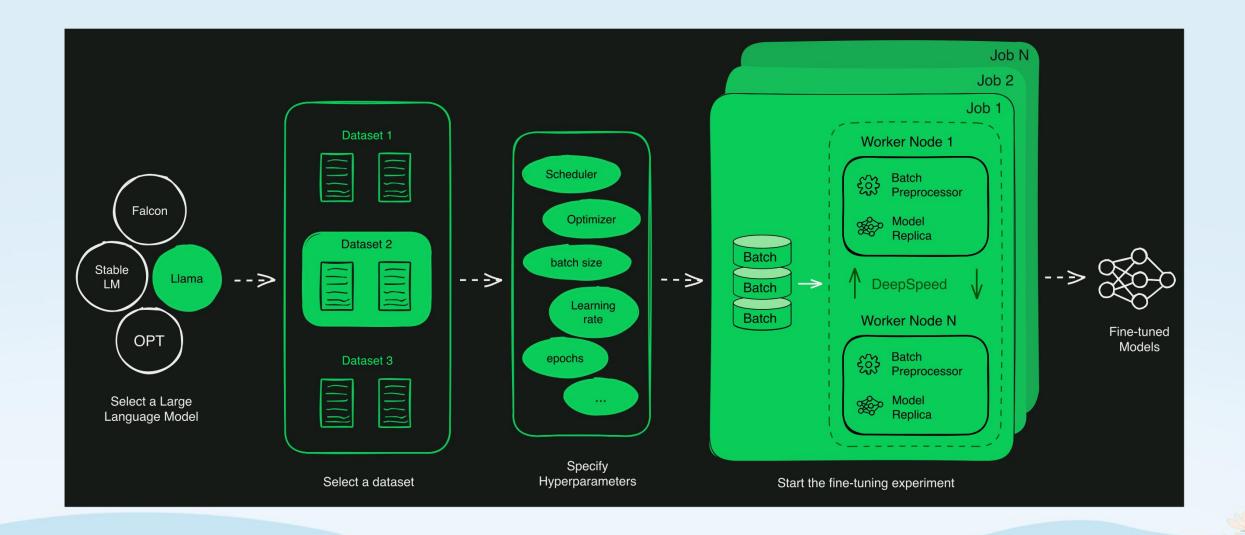
- Ray 是一个用于扩展 AI 和 Python 的 开源统一框架。
- Ray为并行处理提供了计算层,因此无 需成为分布式系统专家。
- 主要组件: Ray AIR, Ray Core, 集成和实用工具。

KubeRay 是一个功能强大的开源 Kubernetes Operator, 可简化 Ray 应用程序在 Kubernetes 上的部署和管理。

特点:自动化,高可用,弹性伸缩,资源调度,监控等。

微调流程





微调实验

插件n

──生成数据

Hyperparam-0

启动微调→



FinetuneExperiment 微调实验 数据集 数据/模型/参数组 微调 Checkpoint 推理服务 模型评估 Images Data-0 插件1 ──生成数据--Ray Services v1.0 一启动微调→ -构建镜像--▶ 一创建推理服务→ -微调-Hyperparam-0 Data-1 ─构建镜像── 微调任务 -启动微调→ 一微调一 ─-创建推理服务 ➤ Ray Services v1.1 Hyperparam-1 插件2 ─启动实验→ Data-2 **→**模型评估-插件2 Ray Services v1.2 ──生成数据 → 启动微调→ 一微调一 构建镜像→▶ 一创建推理服务 → Hyperparam-2 Data-3 构建镜像→▶ Ray Services v1.3 -启动微调→ 一微调──► 一创建推理服务→ Hyperparam-0 插件n Data-4

-微调-

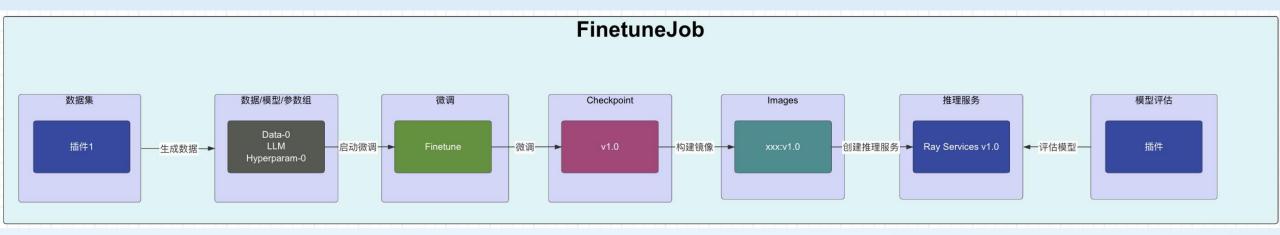
构建镜像─►

Ray Services v1.4

一创建推理服务 →

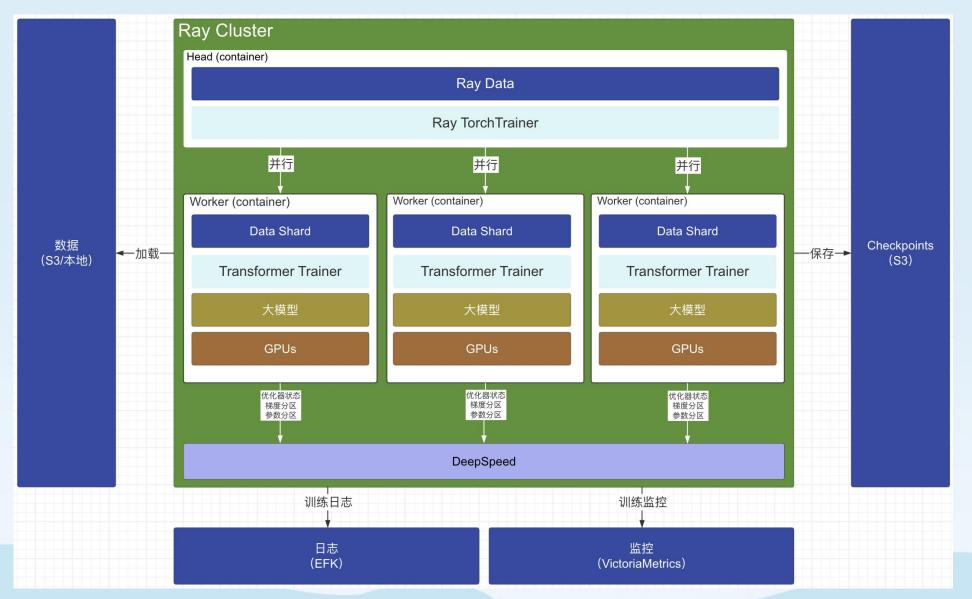
微调任务





finetune





Dataset

```
apiVersion: extension.datatunerx.io/v1beta1
kind: Dataset
metadata:
 name: example-dataset
spec:
  datasetCard:
    datasetCardRef: configmap-dataset-readme
  datasetFiles:
    source: www.test.com
  datasetMetadata:
    datasetInfo:
      features:
      - dataType: string
        mapTo: Content
        name: instruction
      - dataType: string
        mapTo: Result
        name: response
      subsets:
      - name: Default
        splits:
          test:
            file: s3://test.com
          train:
            file: s3://test.com
          validate:
           file: s3://test.com
    languages:
    - zh-CN
    license: CC BY-SA
    plugin:
      loadPlugin: true
      name: demo
      parameters: '{}'
    size: n < 1K
      name: Text Generation
```



```
plugin:
   loadPlugin: true
   name: demo
   parameters: '{}'
```



LLM



```
apiVersion: core.datatunerx.io/v1beta1
kind: LLM
metadata:
 name: llama2-7b
spec:
  llmMetdata:
    computeInfrastructure:
     hardware:
        cpu: 4
        memory: 28Gi
       vRam: 6Gi
    domain:
    - NLP
    frameworks:
    - PyTorch
    languages:
    - English
    license:
    - https://ai.meta.com/resources/models-and-libraries/llama-downloads/
    llmImageConfig:
      image: harobor.custome.com/llama2-7b:v1
      path: /data/llama2-7b
    name: Llama-2-7b
    tasks:
    - Text Generation
```

Hyperparameter



```
apiVersion: core.datatunerx.io/v1beta1
kind: Hyperparameter
metadata:
  name: demo
spec:
  objective:
   type: SFT
  parameters:
    FP16: false
   PEFT: true
    batchSize: 32
    blockSize: 512
    epochs: 10123156464
    gradAccSteps: 1
    int4: false
    int8: false
    learningRate: "0.001"
    loRA_Alpha: "32.0"
    loRA_Dropout: "0.1"
    loRA_R: "4"
    optimizer: AdamW
    scheduler: Cosine
    trainerType: Standard
    warmupRatio: "0.1"
    weightDecay: "0.0001"
```



FinetuneExperiment

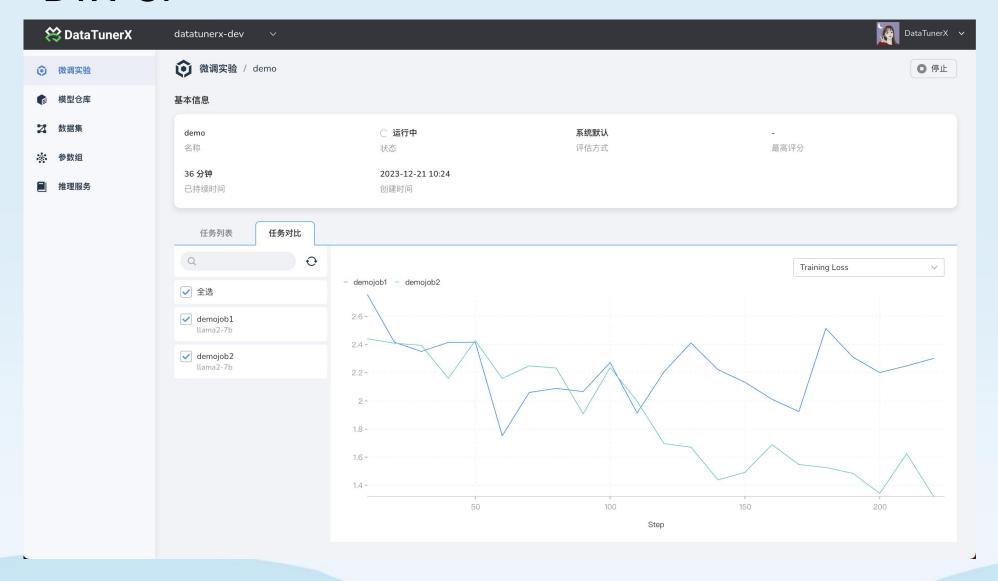
```
apiVersion: finetune.datatunerx.io/v1beta1
kind: FinetuneExperiment
metadata:
 name: test
spec:
  finetuneJobs:
 - name: test
    spec:
      finetune:
        finetuneSpec:
          dataset: example-dataset-liantiao
          hyperparameter:
           hyperparameterRef: happu
           overrides:
              batchSize: 1
              blockSize: 10
              epochs: 1
          image:
            imagePullPolicy: IfNotPresent
          llm: llama2-7b
          resource:
            limits: -
            requests: -
      scoringConfig:
       name: BuildInScoringPlugin
       parameters: '{}'
      serveConfig:
       nodeSelector:
          nvidia.com/gpu: present
```



```
finetuneSpec:
    dataset: example-dataset-liantiao
    hyperparameter:
    hyperparameterRef: happu
    overrides:
        batchSize: 1
        blockSize: 10
        epochs: 1
    image:
    imagePullPolicy: IfNotPresent
    llm: llama2-7b
```

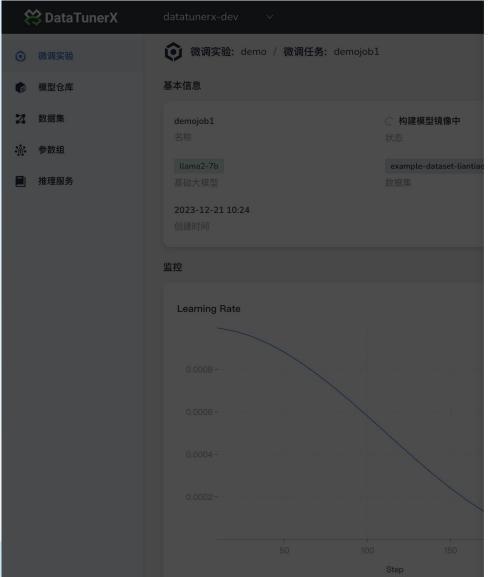






- 微调实验
- 模型仓库
- 数据集管理
- 参数组管理
- 推理服务



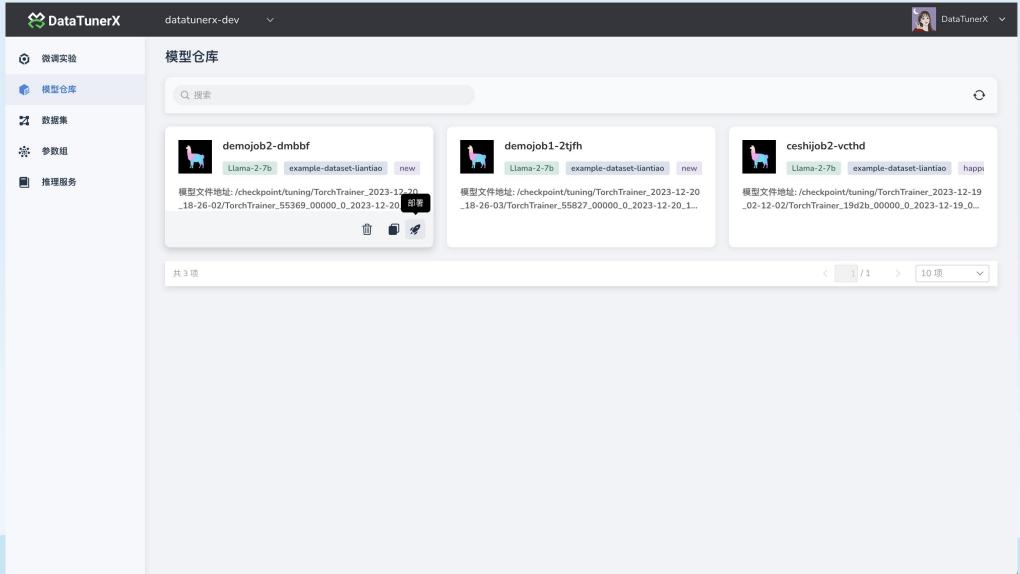


日志 × 耐 默认展示 100 行,查看更多日志信息或者下载日志请联系管理员 Using mask_token, but it is not set yet. Using sep_token, but it is not set yet. Using cls_token, but it is not set yet. Using mask_token, but it is not set yet. Training finished iteration 1 at 2023-12-20 18:33:43. Total running time: 7min 20s Training result checkpoint_000000 checkpoint_dir_name time_this_iter_s 384.59703 384.59703 time_total_s training_iteration epoch 2.2299 train_loss 169.5806 train runtime train samples per second 1.315 train_steps_per_second 1.315 Using sep_token, but it is not set yet. Using cls_token, but it is not set yet. Using mask_token, but it is not set yet. Using sep_token, but it is not set yet. Using cls_token, but it is not set yet. Using mask_token, but it is not set yet. Training saved a checkpoint for iteration 1 at: (s3)datatunerx/tuning/TorchTrainer_2023-12-20_18-2 6-03/TorchTrainer_55827_00000_0_2023-12-20_18-26-29/checkpoint_000000 Training completed after 1 iterations at 2023-12-20 18:33:50. Total running time: 7min 27s Using sep_token, but it is not set yet. Using cls_token, but it is not set yet. Using mask_token, but it is not set yet. Using sep_token, but it is not set yet. Using cls token, but it is not set yet. Using mask_token, but it is not set yet. 2023-12-20 18:34:59,544 WARNING experiment_state.py:355 -- Syncing the experiment checkpoint to cl oud took a long time with 68.74 seconds. This can be due to a large number of trials, large es, or throttling from the remote storage provider for too frequent syncs. If your `Checkpoi ig.num_to_keep` is a low number, this can trigger frequent syncing, in which case you should incre ase it. result path datatunerx/tuning/TorchTrainer_2023-12-20_18-26-03/TorchTrainer_55827_00000_0_2023-12-20 18-26-29/checkpoint 000000 2023-12-20 18:35:09,714 SUCC cli.py:60 --2023-12-20 18:35:09,714 SUCC cli.py:61 -- Job 'demojob1-42bjr' succeeded 2023-12-20 18:35:09,714 SUCC cli.py:62 -- --

- 微调实验
- 模型仓库
- 数据集管理
- 参数组管理
- 推理服务

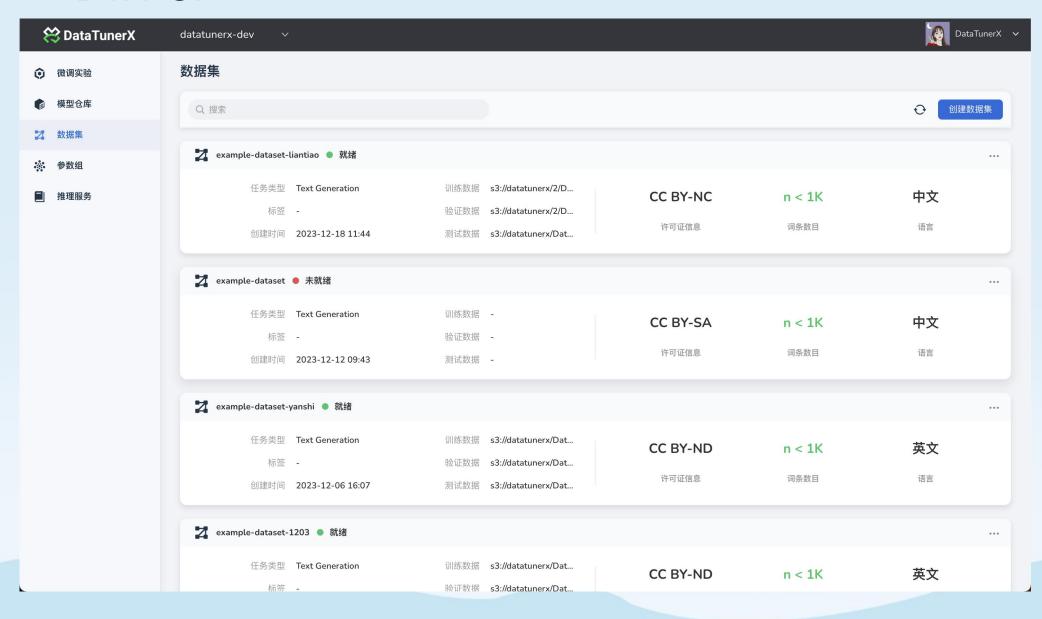






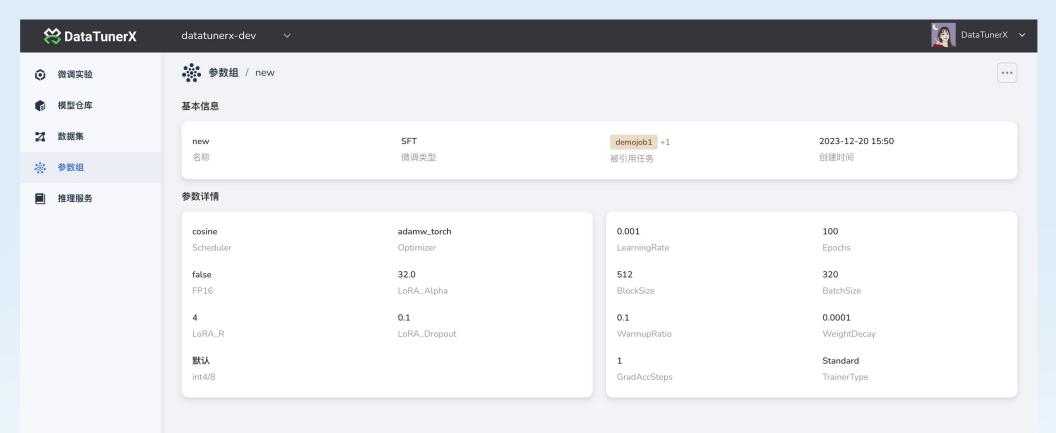
- 微调实验
- 模型仓库
- 数据集管理
- 参数组管理
- 推理服务





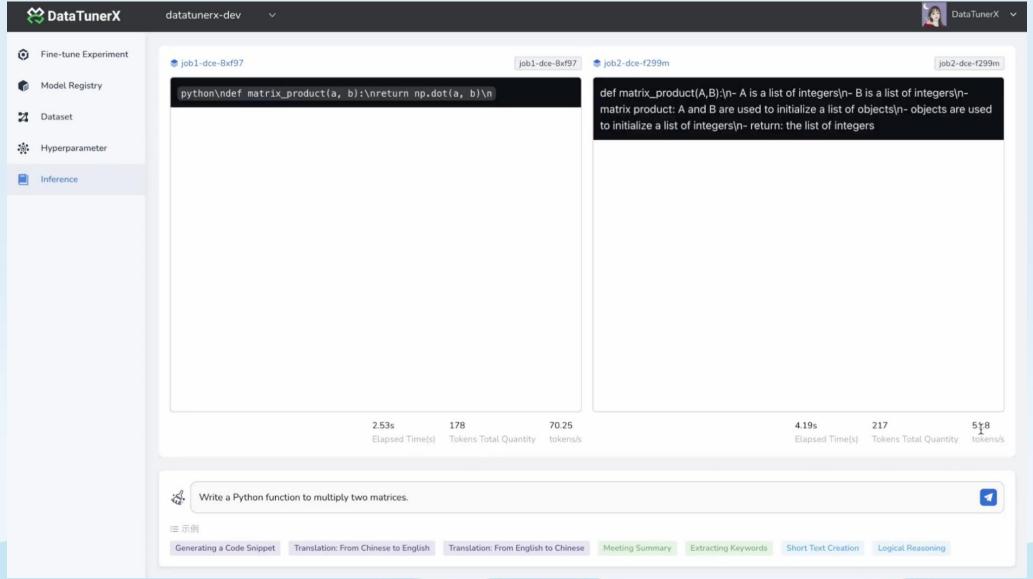
- 微调实验
- 模型仓库
- 数据集管理
- 参数组管理
- 推理服务





- 微调实验
- 模型仓库
- 数据集管理
- 参数组管理
- 推理服务





- 微调实验
- 模型仓库
- 数据集管理
- 参数组管理
- 推理服务

Demo 视频请查看:

https://www.bilibili.com/video/BV1eD421J76W

9 ®

usage step Tutorial









