

# AI 驱动 软件研发 全面进入数字化时代

中国·深圳 11.24-25

AI+  
software  
Development  
Digital  
summit



## 阿里云人工智能平台PAI的MaaS实践

罗义云（一耘）

阿里云资深技术专家、PAI平台工程技术负责人

# 科技生态圈峰会 + 深度研习



—1000+ 技术团队的选择



K+全球软件研发行业创新峰会

会议时间: 2024.05.24-25



K+全球软件研发行业创新峰会

会议时间: 2024.09.20-21



AI+ 软件研发数字峰会

会议时间: 2023.11.24-25



AI+ 软件研发数字峰会

会议时间: 2024.07.19-20



AI+ 软件研发数字峰会

会议时间: 2024.11.15-16

# ▶ 演讲嘉宾



## 罗义云

阿里云资深技术专家

毕业于北京大学，曾任微软高级研发经理、旷视科技AI平台高级技术总监，现任阿里云资深技术专家、AI平台工程技术负责人。在机器学习、人工智能、大数据等方向有着深厚的技术积累和丰富的行业经验。

# 目录

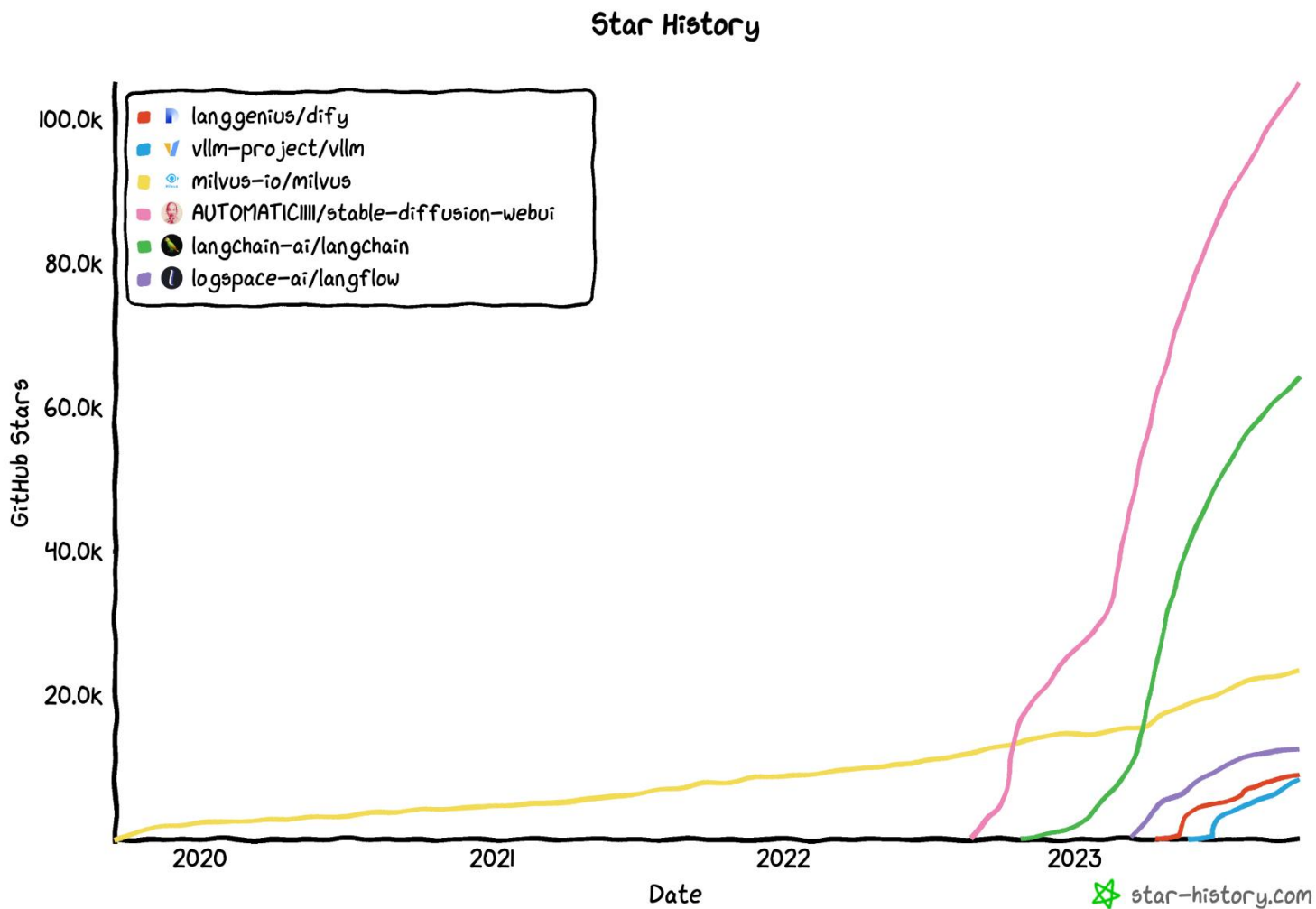
## CONTENTS

1. MaaS的起源和定义
2. PAI的MaaS实践
3. 未来展望

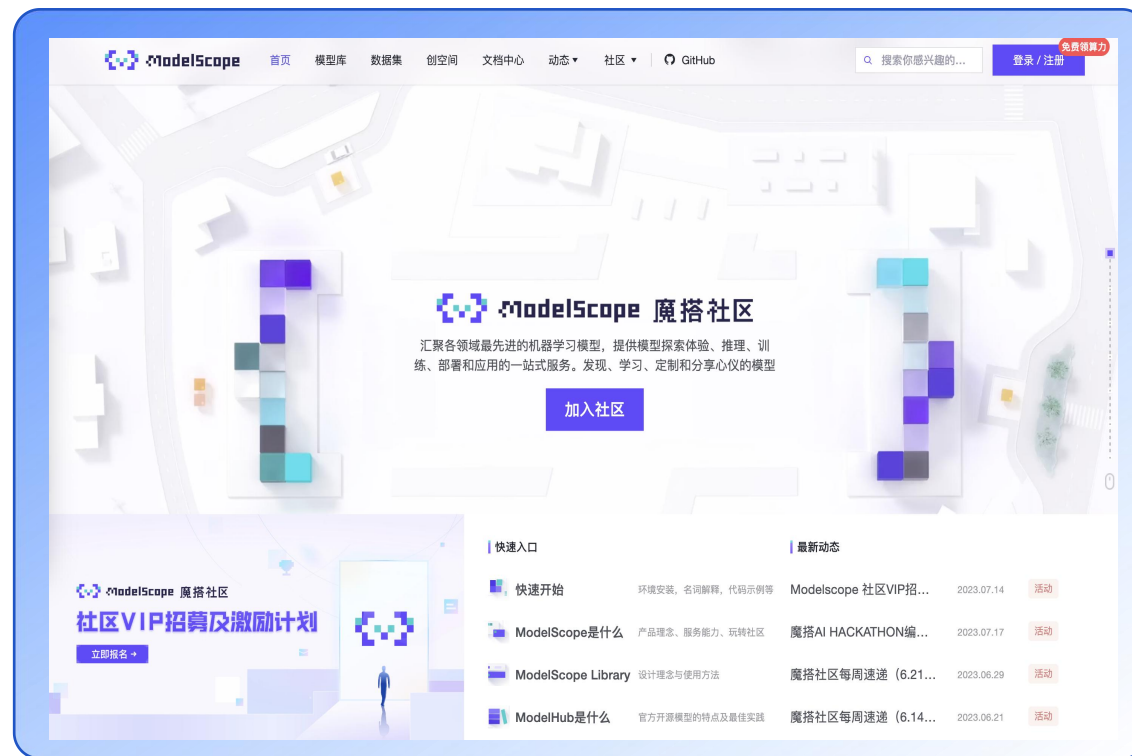
## **PART 01**

# **MaaS的起源和定义**

# ▶ 人工智能进入AIGC时代



# ▶ 开源模型社区



200万+

活跃用户

2300+

模型

1亿+

模型累计下载

20+

顶尖人工智能  
合作机构

AI驱动软件研发全面进入数字化时代

**AI+ 软件研发数字峰会**  
AI+ software Development Digital summit

# ▶ 什么是MaaS?





# ▶ 什么是MaaS?

👤 什么是MaaS (Model-as-a-Service) , 它的核心特点是什么?

🔗 MaaS (Model-as-a-Service) 是一种云计算服务模式, 它以云产品为核心, 为企业提供全方位、全周期的云服务。相比于传统的IaaS、PaaS和SaaS等云服务模式, MaaS具有更强的灵活性和适应性。

MaaS的核心特点可以概括为以下几点:

1. 模型即服务: MaaS将机器学习模型和数据处理能力作为一种服务提供给企业, 企业可以根据自己的需求选择不同的模型进行使用, 无需自己开发和维护模型。
2. 集成化: MaaS将数据处理和模型服务集成到企业的业务流程中, 提供一站式解决方案, 可以大大降低企业的开发和维护成本。
3. 灵活性和适应性: MaaS提供了多种不同的模型和数据处理能力, 企业可以根据自己的需求进行选择和组合, 满足不同的业务场景和需求。
4. 持续优化: MaaS通过不断优化和更新模型, 可以提供更准确、更高效的服务, 帮助企业提高业务效率和效果。

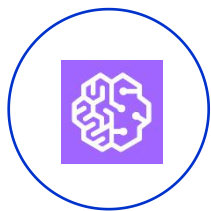
总之, MaaS提供了一种全新的、更灵活、更便捷的机器学习和数据处理服务模式, 可以帮助企业降低开发和维护成本, 提高业务效率和效果。

👍 | 🗣️ | 📄 | 🔄 重新生成

# ▶ MaaS的组成要件



## ▶ 主流云厂商是怎么做的？



SageMaker

Foundation models  
Jumpstart



AzureML

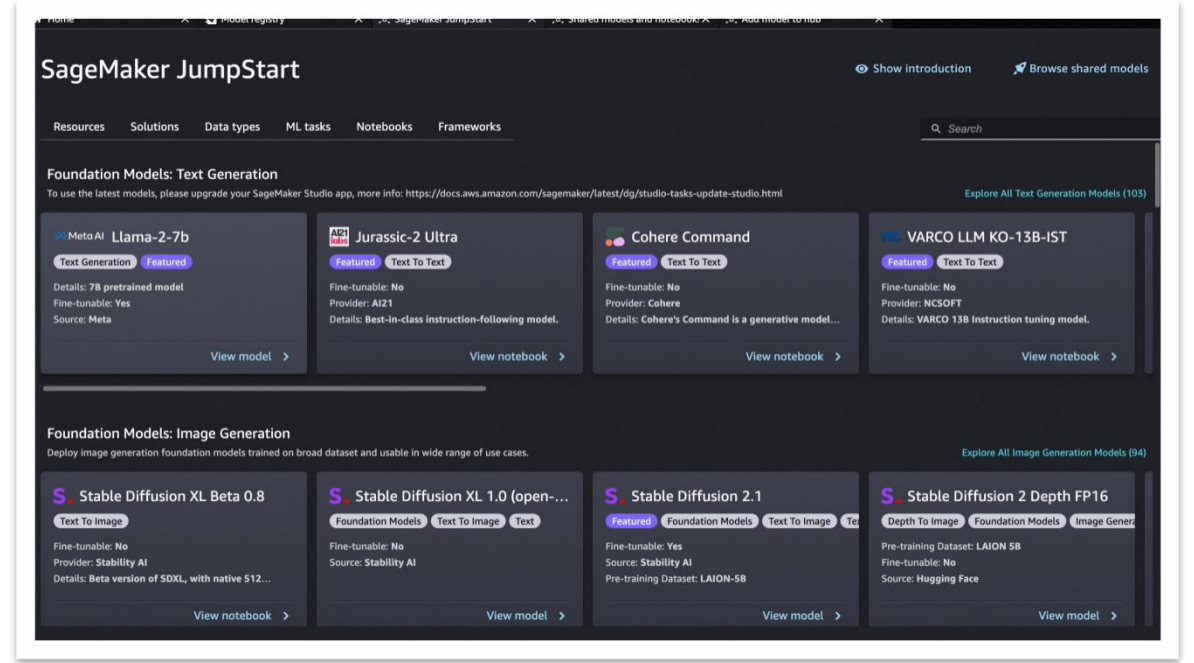
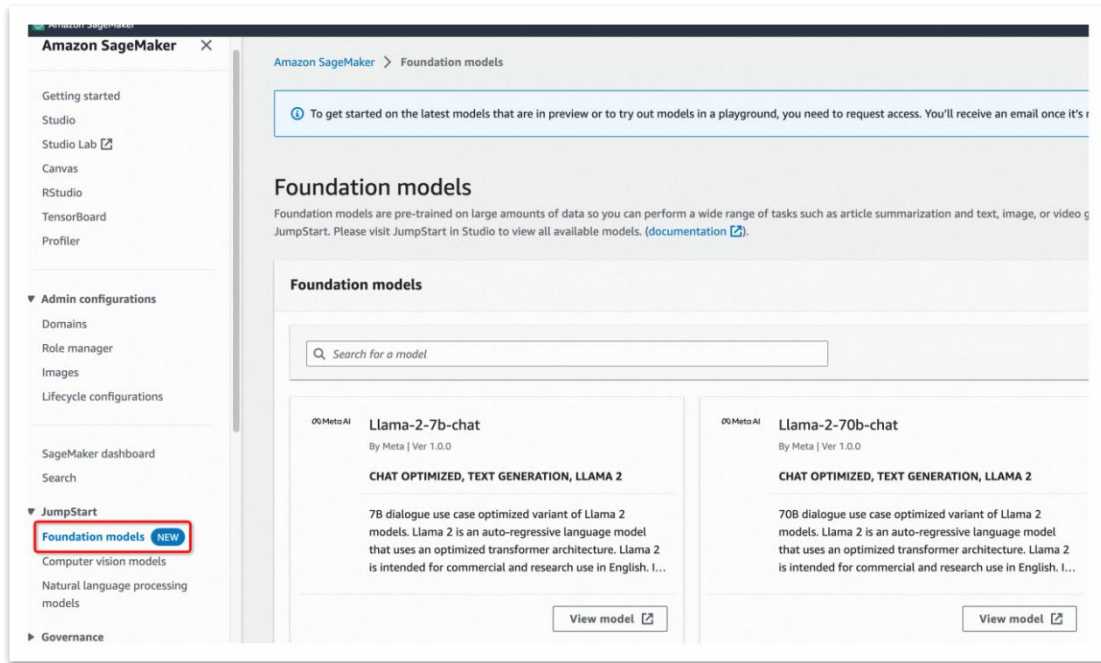
Model catalog

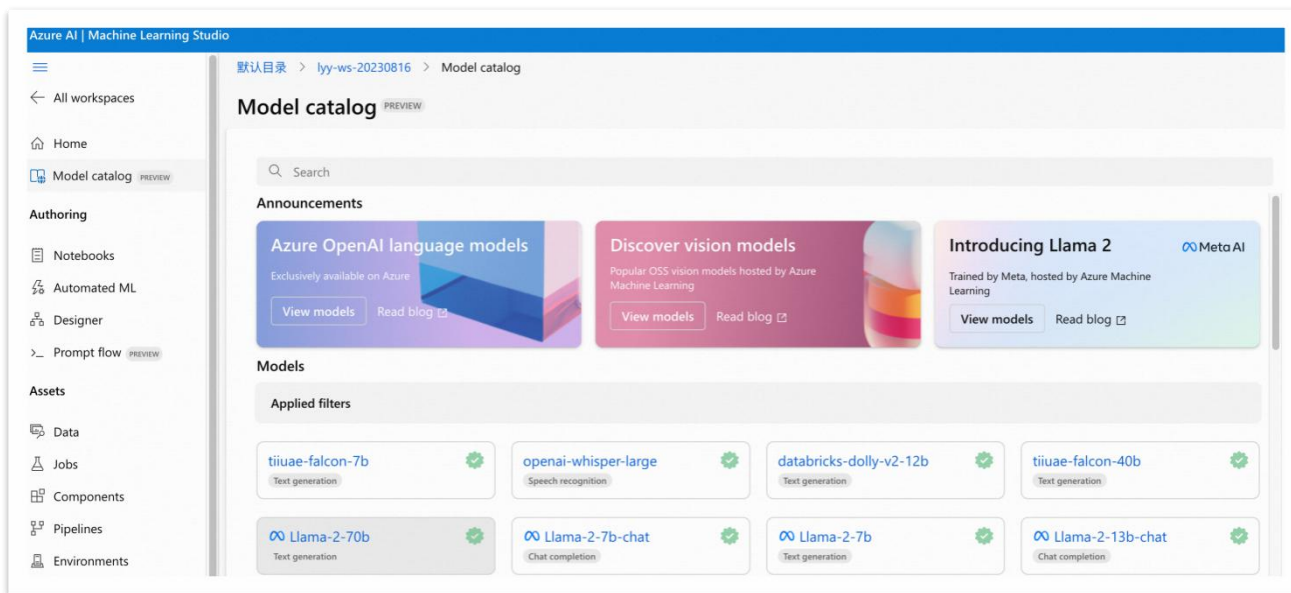


火山引擎

方舟大模型平台

# ▶ SageMaker





### tiuae-falcon-40b

Overview Versions Artifacts

Task: Text generation Finetuning task: text-classification Languages: en, de, es, fr License: apache-2.0

Refresh Evaluate Finetune Deploy View license

#### Description

Falcon-40B is a large language model (LLM) developed by the Technology Innovation Institute (TII) with 40 billion parameters. It is a causal decoder-only model trained on 1 trillion tokens from the RefinedWeb dataset, enhanced with curated corpora. Falcon-40B supports English, German, Spanish, and French languages, with limited capabilities in Italian, Portuguese, Polish, Dutch, Romanian, Czech, and Swedish. It is available under the Apache 2.0 license.

Falcon-40B is considered the best open-source model currently available, optimized for inference with features such as FlashAttention and multiquery. However, it is recommended to fine-tune the model for specific use cases.

The training of Falcon-40B involved using 384 A100 40GB GPUs and took two months. The model carries biases and stereotypes encountered online and requires appropriate precautions for production use. It is suggested to finetune the model for specific tasks and consider guardrails. The technical specifications, training details, and evaluation results are provided in the summary.

The above summary was generated using ChatGPT. Review the [original model card](#) to understand the data used to train the model, evaluation metrics, license, intended uses, limitations and bias before using the model.

## 功能全景

### 模型全景

火山方舟提供模型训练、推理、评测、精调等全方位功能与服务，并重点支撑大模型生态。

01. 模型广场
02. 模型体验
03. 模型训练推理
04. 模型应用

产品咨询



## PART 02

# PAI的MaaS实践

# PAI产品架构

AI应用  
模型服务 (MaaS)  
灵骏智算服务 & 机器学习框架 (PaaS)  
计算资源 & 基础设施 (IaaS)



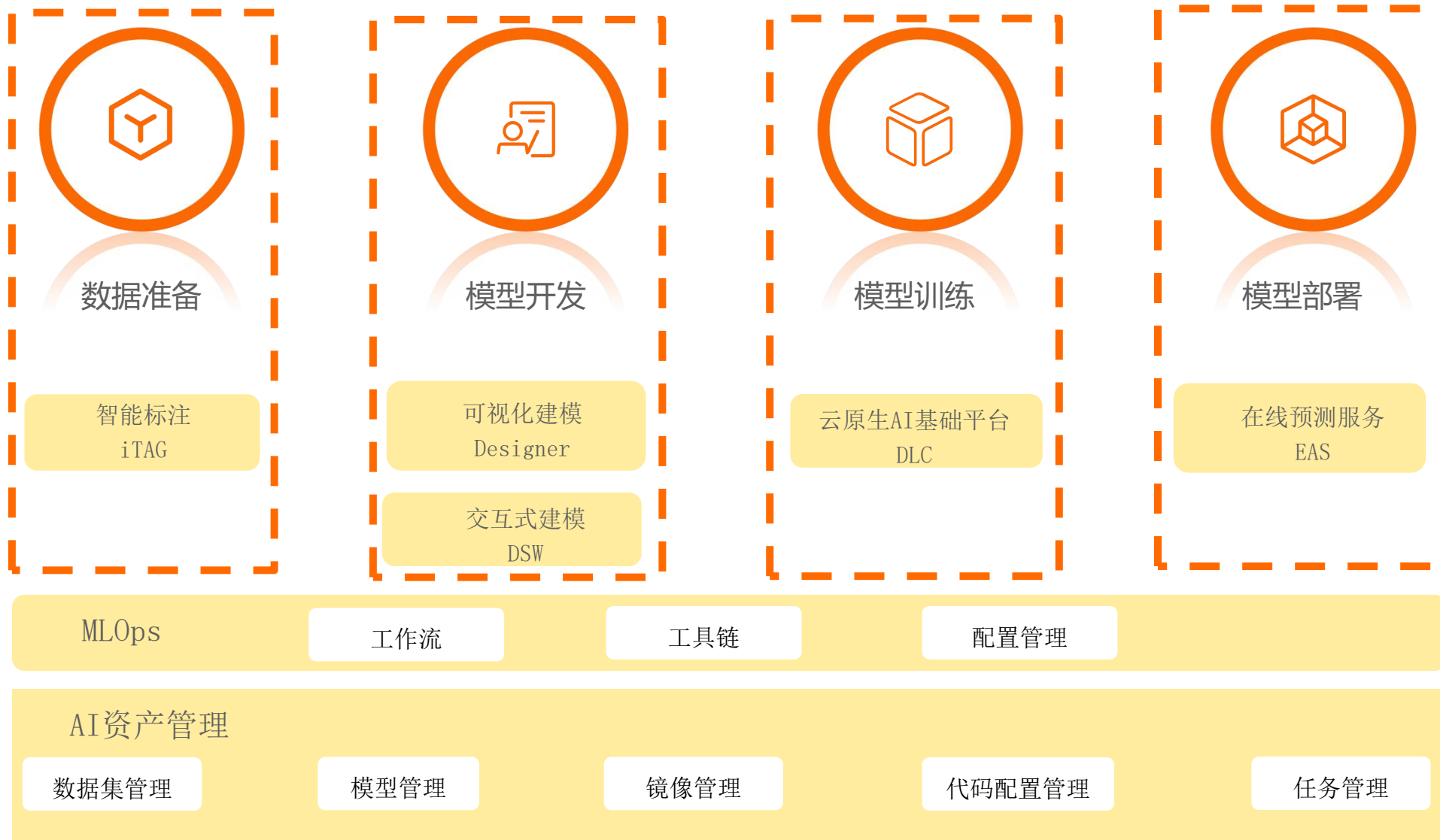
AI驱动软件研发全面进入数字化时代

**AI+ 软件研发数字峰会**  
AI\* software Development Digital summit



# ▶ PAI 提供一站式AI研发生命周期管理

AI工作空间



# ▶ PAI-iTAG 智能标注

阿里云智能标注，提供全场景、高质量、智能化的标注平台服务和人力标注服务



## 标注能力

图像、文本、视频、语音、PDF、多模态、自定义等全方位的标注能力和场景

## 智能标注

打通PAI-EAS部署的模型，让模型服务进行主动预标注

## 预标工具

预置OCR、ASR等预标工具，大幅提升标注员的效率和准确率

## 数据安全

阿里云最新的数字安全传输技术，保护标注数据不泄露

## 任务分发与验收

标注+质检+验收的任务分发机制，保证标注数据的高质量交付

## 人员与权限管理

管理员+标注组长+标注员的人员与权限管理，保证各角色的权限隔离和数据安全

## 人力标注服务：专业、全托管的数据服务

公益模式标注基地，专业培训上岗，全托管的项目管理，极大降低标注人力成本

## 云原生交互式建模

### 实例环境持久化

- 支持保存环境为镜像

### 灵活的环境支持

- 支持预置镜像和自定义镜像

### 多数据集支持

- 支持同时挂载多个数据集

### 实例生命周期管理

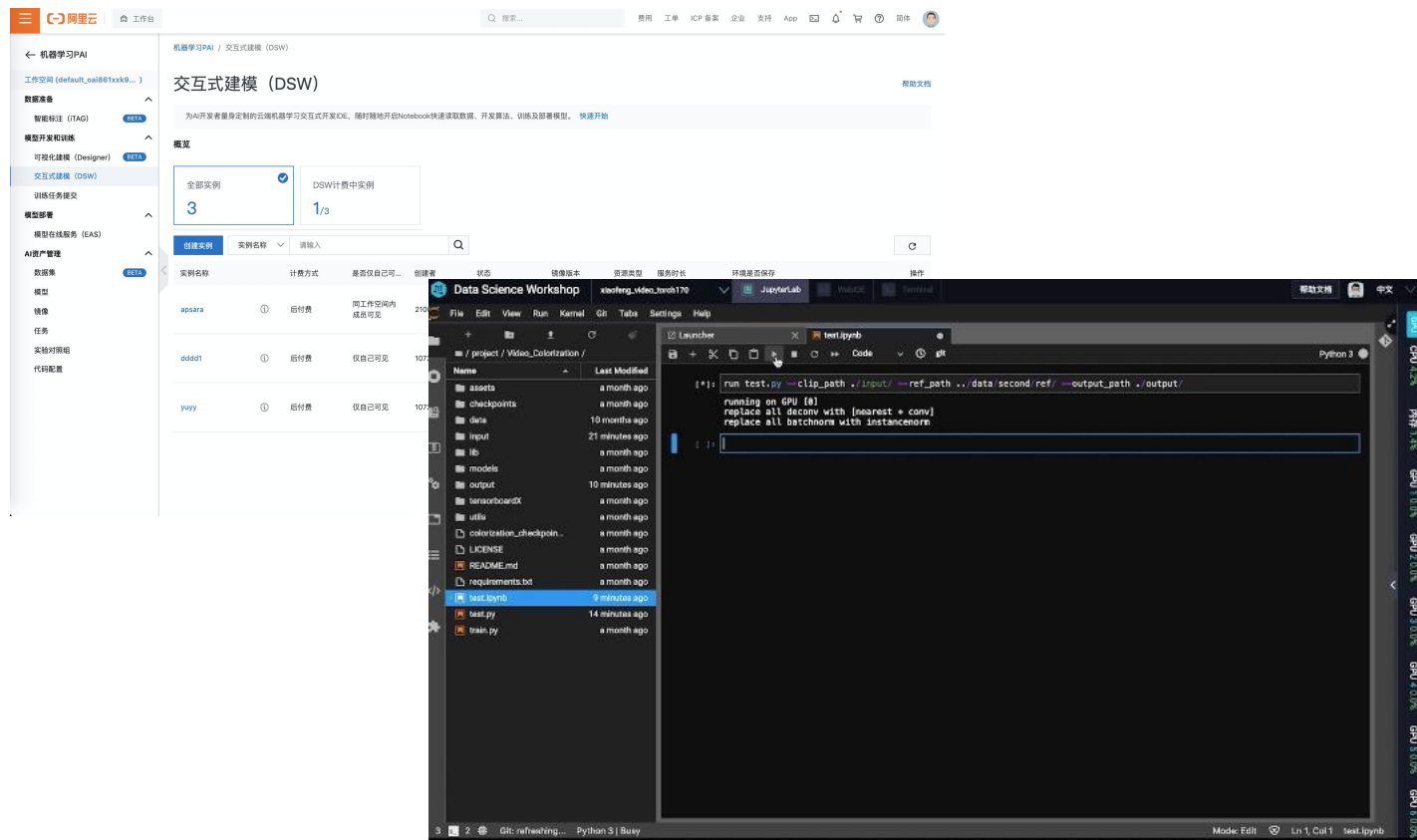
- 支持定时停止

### 实例权限控制

- 支持同一工作空间内成员间分隔

### 开放被集成

- 开放OpenAPI



### 多资源组支持

- 支持公共资源与专有资源

### 灵活的环境支持

- 支持预置镜像和自定义镜像

### 多数据集支持

- 支持同时挂载多个数据集

### 任务生命周期管理

- 任务全流程把控

### 实例权限控制

- 同一工作空间内成员间分隔

### 开放被集成

- 开放OpenAPI

The screenshot displays the PAI-DLC console interface. The top navigation bar includes the Alibaba Cloud logo, a search bar, and various utility icons. The main content area shows the details of a task named 'sample-job' which has completed successfully. A progress bar at the top of the task details indicates the stages: '开始创建' (00:03:21), '资源准备' (00:00:55), and '开始运行' (00:00:55). Below the progress bar, the '基本任务配置' (Basic Task Configuration) section shows the execution command: 'echo 'Hello World''. The '任务资源配置' (Task Resource Configuration) section shows the worker configuration: 'pytorch-training:1.7.1-gpu-py37-cu110-ubuntu18.04' with 1 node. At the bottom, a table lists the task instances.

名称	类型	状态	创建时间	启动时间	结束时间	执行时长	操作
dlc1tc29kxwck4rd-master-0	master	执行成功	2021-11-18 08:54:12		2021-11-18 08:55:06	00:00:54	日志

# ▶ PAI-EAS

基于异构硬件(CPU/GPU/NPU/FPGA)提供机器学习深度学习模型快速部署的微服务平台

## 支持多种框架模型

提供通用机器学习模型 (PMML, TensorFlow, Pytorch等) 一键部署成在线服务。

## 自定义多语言processor

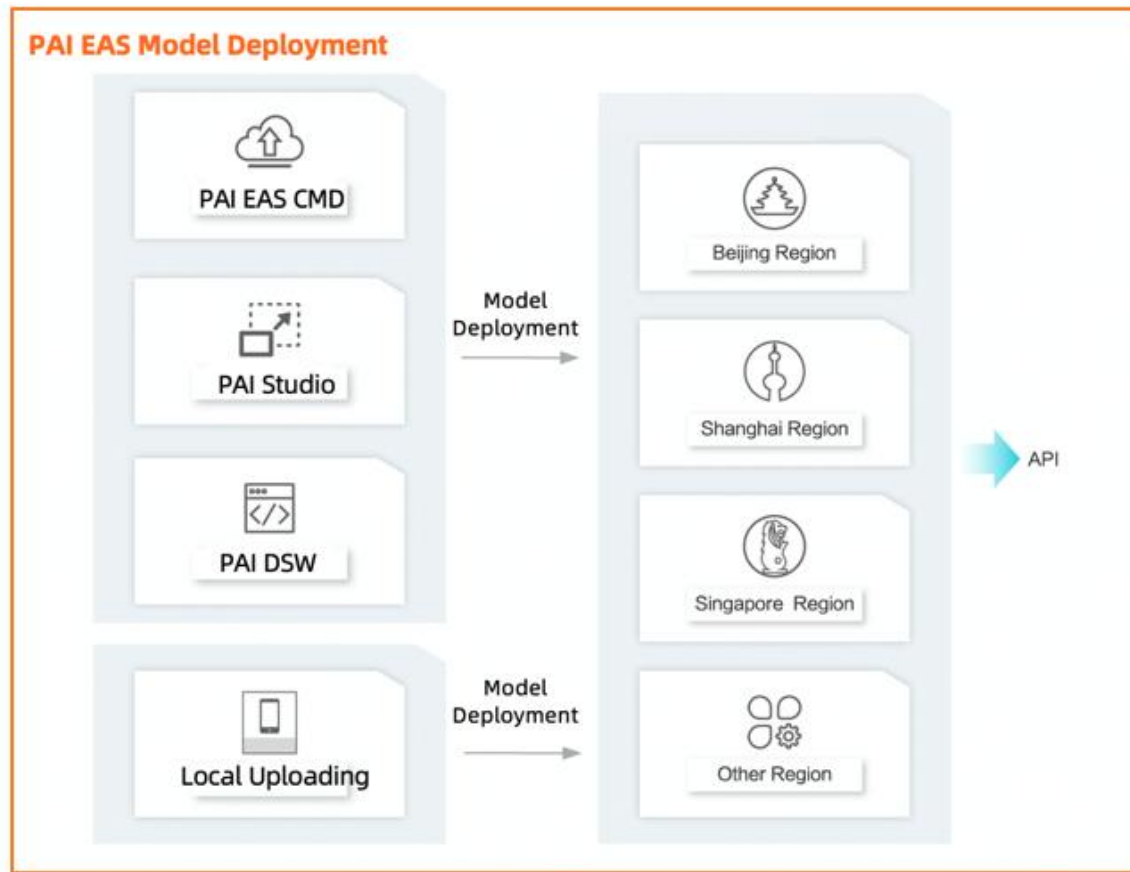
提供三种主流语言(C++/Python/Java)开发用户自定义的预测逻辑, 将模型包装成端对端的产品服务。

## 高性能

支持阿里集团数十个BU的模型预测服务及大量外部客户, 单服务QPS峰值30w, 大量QPS过万的图像类服务

## 开放被集成

开放OpenAPI



## ▶ PAI的MaaS理念

Model Centric

一站式

白盒化

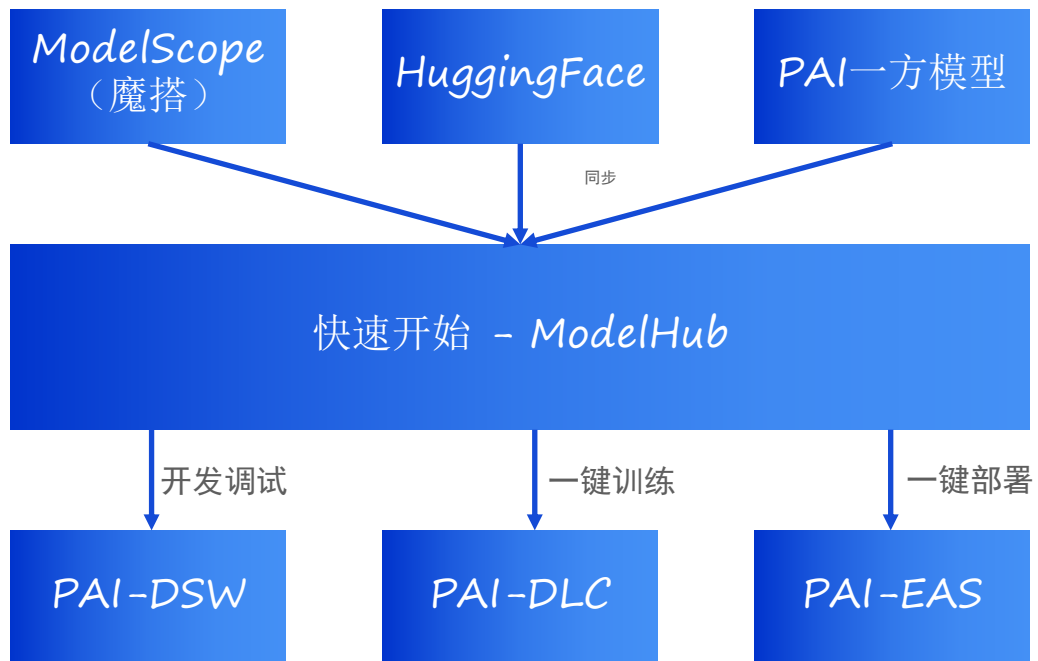
## ▶ PAI的MaaS实践

快速开始  
(ModelHub)

PAI SDK

PAI × ModelScope

## 快速开始 (ModelHub)



- 丰富的模型支持
- 一站式模型训练->部署的全链路体验
- 内置算法工程优化, 提升迭代效率
- 结合平台优化能力, 提供极致的性能和性价比



# ▶ PAI SDK对开源模型的支持

```
# Huggingface模型训练
from pai.huggingface.estimator import HuggingFaceEstimator

# 通过 git_config 来指定 transformers 的 git 地址和分支
git_config = {
    "repo": "https://github.com/huggingface/transformers.git",
    "branch": "v4.29.2",
}

# 通过 hyperparameters 来指定训练参数, 这里以 GLUE 的 MRPC 任务为例
hyperparameters = {
    "model_name_or_path": "bert-base-uncased",
    "output_dir": "/ml/output/model/",
    "task_name": "mrpc",
    "do_train": True,
    "do_eval": True,
    "max_seq_length": 128,
    "per_device_train_batch_size": 32,
    "learning_rate": 2e-5,
    "num_train_epochs": 3,
}

# 创建 HuggingFaceEstimator 对象
est = HuggingFaceEstimator(
    source_dir="./examples/pytorch/text-classification", # 训练脚本在 git repo 中的相对路径
    git_config=git_config, # 指定 transformers 的 git 地址和分支
    command="python3 run_glue.py $PAI_USER_ARGS", # 指定训练命令, 通过 $PAI_USER_ARGS 传入 hyperparameters
    instance_type="ecs.gn7i-c32g1.8xlarge", # 指定训练机型
    transformers_version="4.29.2", # 指定训练镜像所使用的 transformers 版本, 目前支持 4.29.2 或 latest
    hyperparameters=hyperparameters, # 指定训练参数
    base_job_name="huggingface-sdk-train", # 指定训练任务名称
)

# 进行训练
est.fit()
```

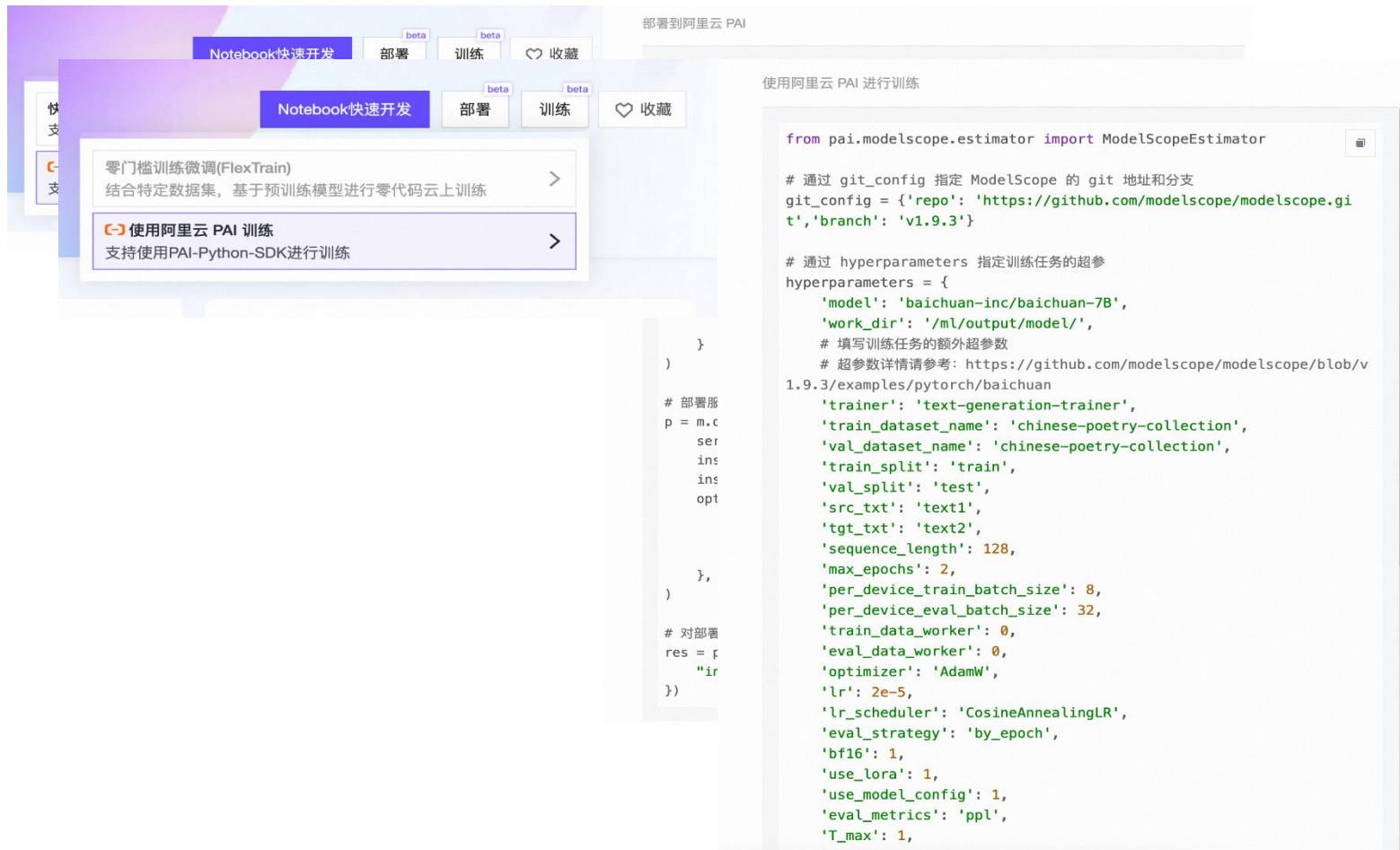
```
# Huggingface模型部署
from pai.common.utils import random_str
from pai.huggingface.model import HuggingFaceModel

# 创建 HuggingFaceModel 对象
m = HuggingFaceModel(
    command="python app.py", # 指定服务启动命令
    transformers_version="latest", # 指定服务镜像所使用的 transformers 版本, 目前支持 4.29.2 或 latest
    # 通过环境变量传入模型ID、任务类型和模型版本, 与 Huggingface Modelhub 保持一致
    environment_variables={
        "MODEL_ID": "distilbert-base-uncased-finetuned-sst-2-english",
        "TASK": "text-classification",
        "REVISION": "main",
    }
)

# 部署服务
p = m.deploy(
    service_name="hf_bert_serving_{}".format(random_str(6)),
    instance_count=1,
    instance_type="ecs.gn6i-c4g1.xlarge",
    options={
        "metadata.rpc.keepalive": 5000000,
        "features.eas.aliyun.com/extra-ephemeral-storage": "40Gi",
    },
)

# 对部署好的服务进行一些测试
res = p.predict(
    {
        "data": ["it's so easy!"]
    }
)
```

# ▶ PAI × ModelScope



部署到阿里云 PAI

Notebook快速开发 部署 训练 收藏

零门槛训练微调(FlexTrain)  
结合特定数据集，基于预训练模型进行零代码云上训练

使用阿里云 PAI 训练  
支持使用PAI-Python-SDK进行训练

```
from pai.modelscope.estimator import ModelScopeEstimator

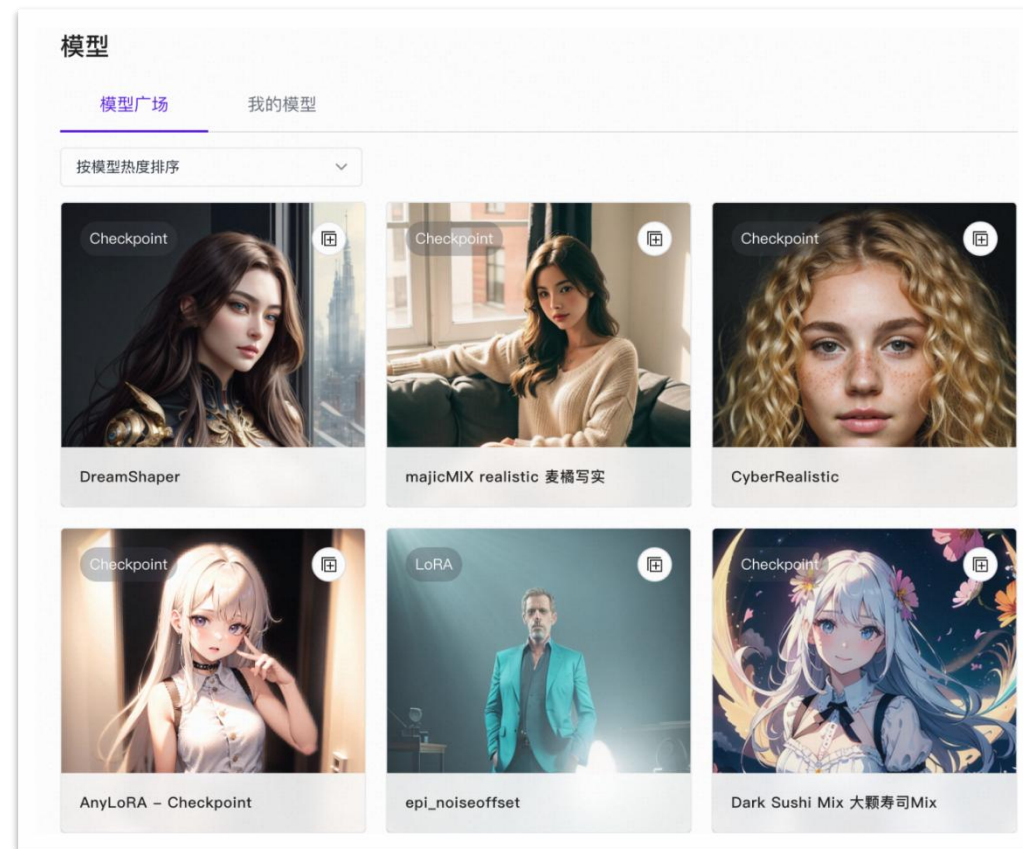
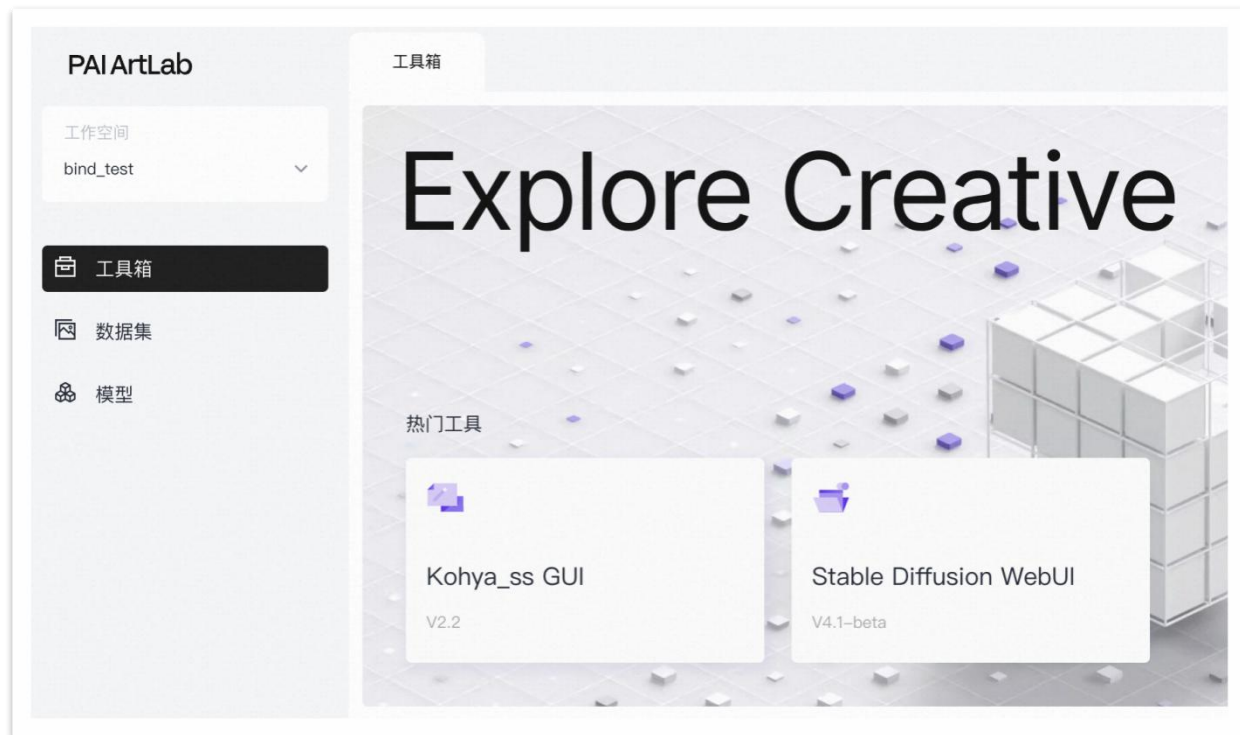
# 通过 git_config 指定 ModelScope 的 git 地址和分支
git_config = {'repo': 'https://github.com/modelscope/modelscope.git', 'branch': 'v1.9.3'}

# 通过 hyperparameters 指定训练任务的超参
hyperparameters = {
    'model': 'baichuan-inc/baichuan-7B',
    'work_dir': '/ml/output/model/',
    # 填写训练任务的额外超参数
    # 超参数详情请参考: https://github.com/modelscope/modelscope/blob/v
1.9.3/examples/pytorch/baichuan
    'trainer': 'text-generation-trainer',
    'train_dataset_name': 'chinese-poetry-collection',
    'val_dataset_name': 'chinese-poetry-collection',
    'train_split': 'train',
    'val_split': 'test',
    'src_txt': 'text1',
    'tgt_txt': 'text2',
    'sequence_length': 128,
    'max_epochs': 2,
    'per_device_train_batch_size': 8,
    'per_device_eval_batch_size': 32,
    'train_data_worker': 0,
    'eval_data_worker': 0,
    'optimizer': 'AdamW',
    'lr': 2e-5,
    'lr_scheduler': 'CosineAnnealingLR',
    'eval_strategy': 'by_epoch',
    'bf16': 1,
    'use_lora': 1,
    'use_model_config': 1,
    'eval_metrics': 'ppl',
    'T_max': 1,
}
```

```
}
)
# 部署服
p = m.c
ser
ins
ins
opt
},
)
# 对部署
res = {
    "ir
})
```



# ▶ Case study: PAI AI设计师专区



# ▶ Demo: 通过快速开始完成qwen的训练和部署

The screenshot displays the Alibaba Cloud PAI console interface. The main content area shows the 'Workspace Details' (工作空间详情) for workspace 'lyy\_358138\_20230801'. It includes a 'Quick Start' (快速开始) button and a 'Run Overview' (运行总览) table.

运行总览	DSW 实例	DLC 任务	EAS 服务	Designer 工作流
	0	0	5	0

AI资产				
数据集	模型	镜像	自定义组件	代码
0	0	0	0	0

Below the tables is a 'Cloud Native Development Scenario' (云原生开发场景) diagram illustrating the workflow: Data Collection (数据集) → Interactive Modeling (DSW) → Container Training (DLC) → Model Management (模型管理) → Model Online Service (EAS). The diagram includes a 'Start Using' (开始使用) button and a 'Mirror Management' (镜像管理) section.

The right sidebar contains workspace configuration details:

- Workspace Name: lyy\_358138\_20230801
- Workspace ID: 79567
- Current Account Role: Administrator (管理员)
- Storage Settings: Not configured (未设置)
- Compute Resources: public-cluster

# **PART 03**

# **未来展望**



# ▶▶ 未来展望



## ▶▶ Qwen训练的一些细节

~100

0

实验次数

2天

7B模型全参数微  
调训练时间

10mi

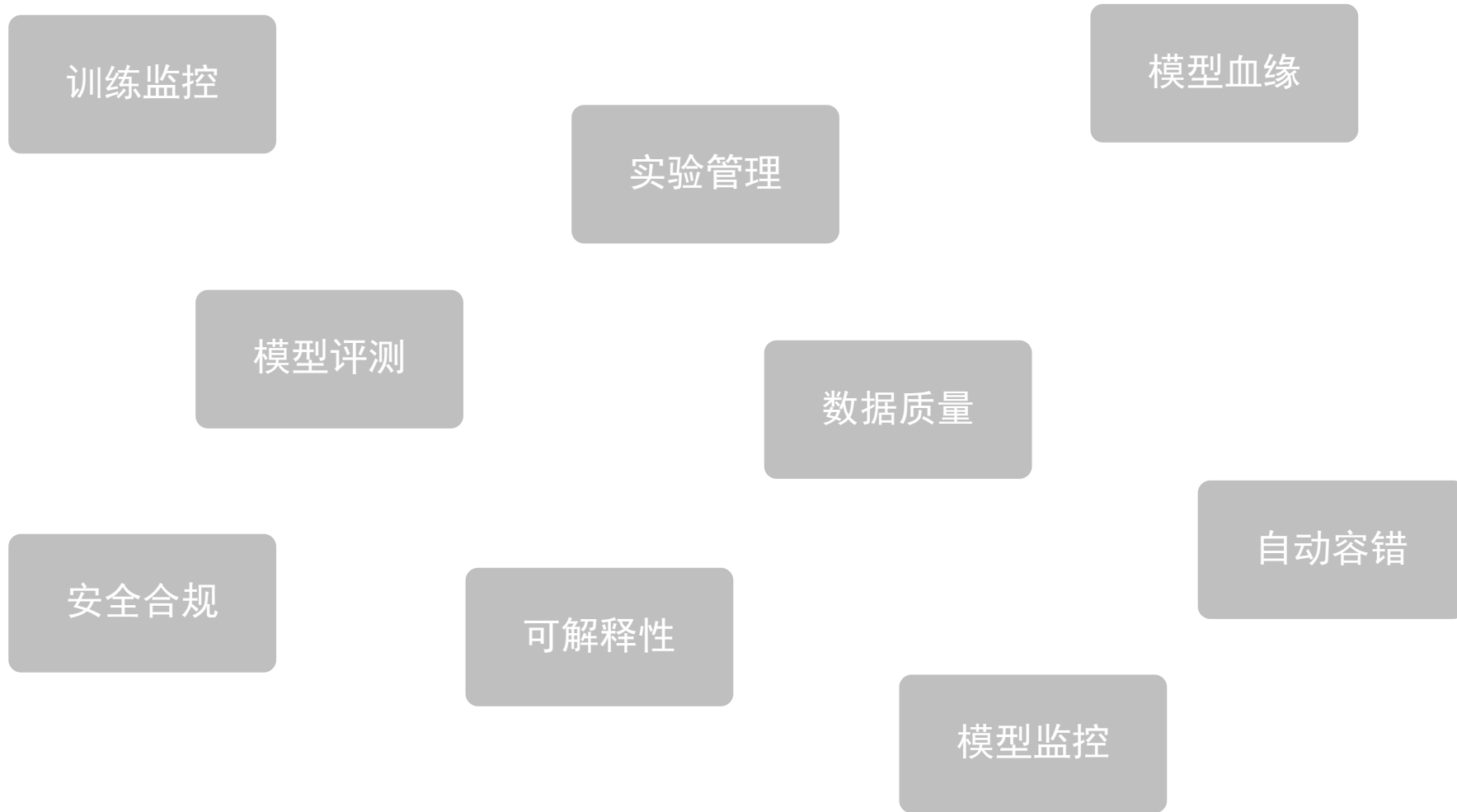
n

关注loss的时  
间间隔

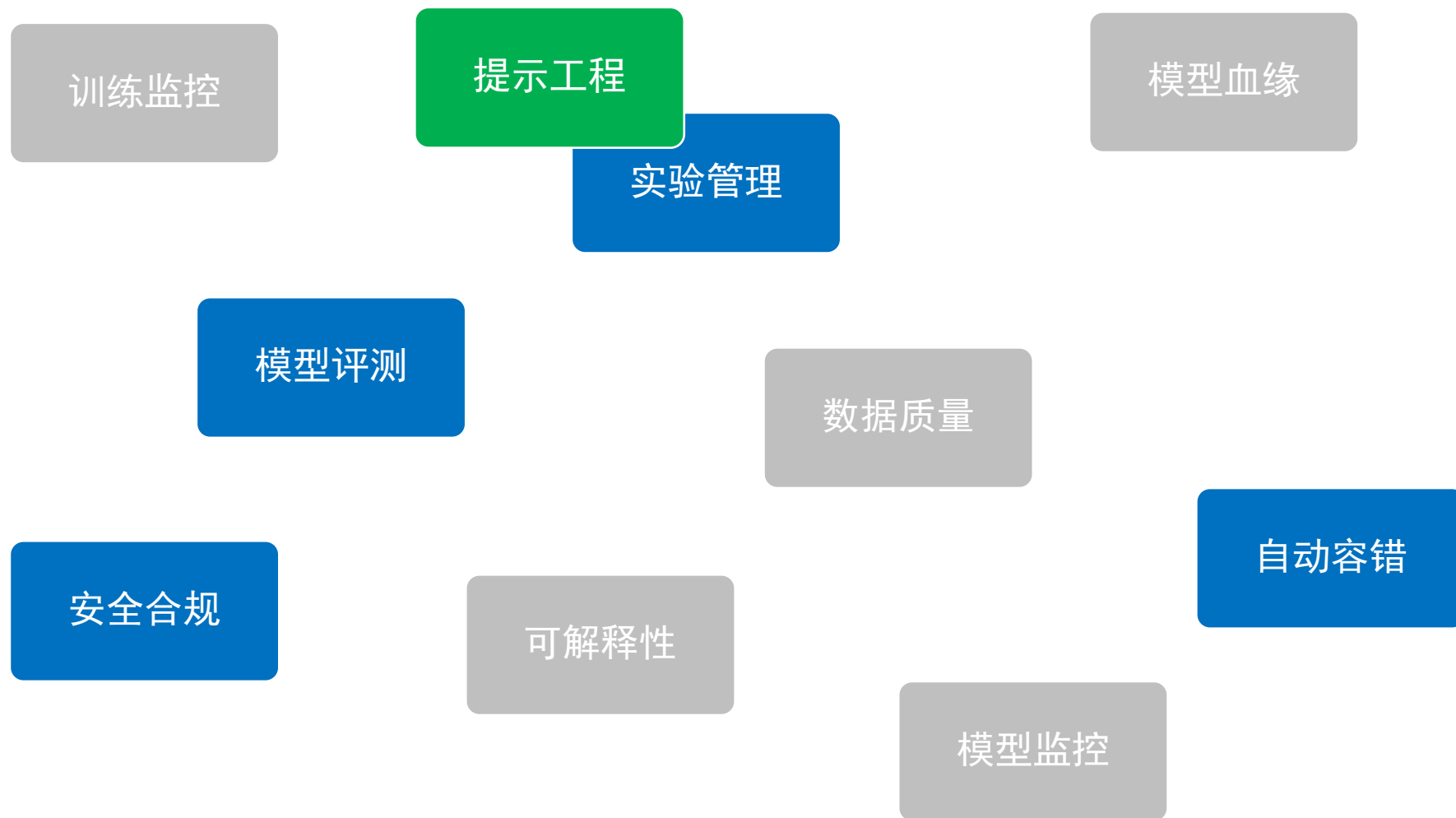
# MLOps: Efficiency is All You Need



# ▶ MLOps



# ▶ MLOps -> LLMOps



## ▶ 深度结合平台优化能力

PAI-ACC AI加速服务：提供AI workflow完整的加速能力



数据集加速器

*DatasetAccelerator*



分布式训练加速

*TorchAcc*



推理加速

*Blade BladeLLM*

## ▶ BladeLLM: LLM高性能推理引擎



1.7~3.8倍

服务吞吐提升



8.7~13.8倍

首包延迟降低

# THANKS

