

《YOLO 代码复现框架》

项目概要设计

V1.0

版本历史

版本/ 状态	作者	参与者	日期	备注

目 录

第一部分 引言	5
一、编写目的	5
二、读者对象	5
三、术语与缩写解释	5
1.YOLO 目标检测	5
2.置信度	5
3.边界框	5
4.非极大值抑制	6
5.推理速度	6
6.平均精度均值	6
7.网格划分	6
8.交并比	6
四、参考资料	6
第二部分 项目概述	7
一、项目描述	7
二、项目功能描述	7
1、数据处理与增强	7
2、模型训练与微调	8
3、模型评估与可视化	8
4、模型推理与部署	9
5、系统管理与运维	9
第三部分 设计约束	9
一、需求约束	10
1. 本系统应当遵循的技术标准	10
2. 软、硬件环境标准	10
3. 接口/协议标准	10
4. 用户界面	11
5. 软件质量	11
二、隐含约束	11
第四部分 功能结构设计	12
一、YOLO 代码复现框架相关业务流程	12
1、模型训练与评估流程	12
2、模型微调流程	13
3、Web 调用模型流程	14
二、业务功能概要结构	15
1、数据处理与增强模块	15
1.1、数据源接入与格式兼容	15
1.2、数据预处理与数据增强	15
2、模型训练与微调模块	16
2.1、训练任务配置与提交	16
2.2、训练过程监控	16

2.3、训练任务管理	17
3、模型评估与可视化模块	17
3.1、标准化模型评估	17
3.2、模型性能分析	17
3.3、检测结果数据分析	18
3.4、生成评估结果报告	18
4、模型推理与部署模块	18
4.1、实时推理	18
4.2、模型格式转换	19
4.3 模型优化与加速	19
4.4、标准化部署接口	19
三、模块定义	20
第五部分 E-R 实体设计	21
一、E-R 实体结构图	21
二、实体描述	21
1、datasets 实体描述	21
2、annotations 实体描述	22
3、Experiments 实体描述	22
4、models 实体描述	23
5、nodes 实体描述	23
6、audit_logs 实体描述	24
第六部分 用户界面设计	24
一、UI 展示	25
第七部分 总体设计	25
一、模型架构	25
1、训练器	25
2、模型优化板块	25
3、模型评估	25
二.数据流转框架	26
1、数据预处理	26
2、数据加载	26
第八部分 运行环境和部署	错误！未定义书签。
一、运行环境	26
1、客户机器环境	26
2、开发环境要求	26
二、系统性能要求	27

第一部分 引言

一、编写目的

编写本文的主要目的是把需求分析得到的用例模型转换为软件结构和数据结构。设计软件结构的具体任务是：将一个复杂系统按功能进行模块划分、建立模块的层次结构及调用关系、确定模块间的接口及人机界面等。数据结构设计包括数据特征的描述、确定数据的结构特性，以及数据库的设计。

本设计是指导详细设计和项目实施的重要指导性文件，也是进行系统集成测试的重要依据。

二、读者对象

该文档的读者为用户代表、软件分析人员、开发管理人员和测试人员。

三、术语与缩写解释

1.YOLO 目标检测

YOLO(You Only Look Once) 是一种端到端的实时目标检测算法。YOLO 将目标检测任务转化为单一回归问题，直接从图像像素到边界框和类别概率的预测，具有速度快、精度高的特点，广泛应用于自动驾驶、安防监控等领域。

2.置信度

置信度 (confidence) 是目标检测模型对检测结果可信程度的评分，表示模型认为某个检测框中确实存在目标的概率。

3.边界框

边界框 (Bounding Box) 是目标检测算法用于标记目标在图像中的具体位置的矩形框，通常用坐标表示。

4.非极大值抑制

非极大值抑制（NMS）是一种后处理算法，用于去除目标检测结果中的重复边界框，保留置信度最高的检测结果。

5.推理速度

推理速度（Inference Speed）指目标检测模型在实际应用中处理图像的速度，通常以帧每秒（FPS）为单位，影响系统的实时性。

6.平均精度均值

平均精度均值（mAP）是目标检测模型性能评估的重要指标，衡量模型在不同类别上的检测准确率，mAP 越高表示模型检测效果越好。

7.网格划分

网格划分（Grid Division）是 YOLO 算法的核心思想，将输入图像划分为 $S \times S$ 个网格单元，每个网格单元负责检测其中心点落在该单元内的目标。YOLOv1 中通常使用 7×7 的网格划分。

8.交并比

交并比（IoU, Intersection over Union）是衡量两个边界框重叠程度的指标，计算公式为重叠面积除以并集面积。IoU 广泛用于训练损失计算、NMS 阈值判断和 mAP 评估。

四、参考资料

《项目需求分析》

第二部分 项目概述

一、项目描述

随着视频采集设备与边缘计算能力的大范围普及,以及行业对实时目标检测与批量模型管理需求的快速增长(如安防监控、工业质检、农业监测等),企业需要一套统一、可复用的目标检测开发与部署平台。为满足这些通用需求,本项目开发了“YOLO 复现与部署平台”,通过统一的数据接入与模型流水线,实现了数据标注与预处理、模型训练与版本管理、模型优化与导出、在线/离线推理接口以及可视化示例演示等功能。例如,在一个厂区人员计数项目中,通过应用 YOLO 模型,实现了对厂区内人员的实时计数以及佩戴安全帽的自动识别,利用数学算法实现通道人员进出实时计数,从而在不依赖人工干预的情况下,实时监控厂区内员工的安全状况。该平台支持多版本 YOLO 的对比分析、快速微调与一键部署(适用于边缘计算和服务端),并集成了监控、日志和权限管理功能,使得团队能够在安防监控、自动驾驶、工业自动化等多个业务场景中高效复用和迭代优化模型。

二、项目功能描述

本项目旨在实现 YOLO 系列目标检测算法的代码复现、训练优化与模型部署,并提供一个可交互的演示系统。项目功能主要包括数据处理、模型训练、模型评估、推理部署和系统管理五大模块。

1.数据处理与增强

本模块主要负责输入数据的标准化处理和多样性增强,为后续模型训练提供高质量的数据基础。系统支持从多种数据源获取数据,并通过一系列预处理和增强技术提升模型的泛化能力。

数据处理与增强模块主要包括:

- 1) 数据源接入:支持处理本地图片、视频文件以及实时视频流(如 RTSP)。
- 2) 格式兼容:兼容 COCO、VOC 等主流数据集格式,并支持用户导入自定义格式的数据集。
- 3) 数据预处理:提供图像尺寸归一化、色彩空间转换(如 RGB 转灰度图)等基础

处理功能。

- 4) 数据增强：通过实施多种图像增强技术，如随机裁剪、水平翻转、色彩抖动和亮度调整，可以有效提升模型的性能和泛化能力，同时减少过拟合现象。
- 5) 数据集管理：支持对导入的数据集进行版本控制、标注查看与管理。

其中，数据增强策略可根据具体任务需求进行组合配置，所有处理步骤均可通过配置文件灵活调用。

2.模型训练与微调

本模块为框架的核心，旨在为用户提供一个灵活、高效、可复现的模型训练环境。用户可根据需求选择不同版本的 YOLO 算法，进行从零训练或基于预训练模型的微调。

模型训练与微调模块主要包括：

- 1) 多版本支持：支持多个主流 YOLO 系列算法的复现和训练。
- 2) 训练模式：支持从零开始训练（train from scratch）和基于官方预训练权重的微调（fine-tuning）。
- 3) 参数配置：支持自定义配置学习率策略、优化器（如 SGD, Adam）、批次大小（Batch Size）、训练轮次（Epochs）等超参数。
- 4) 训练监控：实时显示训练过程中的损失函数（Loss）变化曲线、学习率变化等关键指标。
- 5) 断点续训：自动保存训练检查点（Checkpoint），支持从中断处恢复训练，保障长周期任务的稳定性。

用户可根据不同的硬件配置和任务需求，灵活调整训练参数，以达到最佳的训练效果和效率。

3.模型评估与可视化

本模块主要用于对训练完成后的模型进行客观、定量的性能评估。系统提供标准化的评估指标和丰富的可视化工具，帮助用户直观地分析模型性能优劣。

模型评估与可视化模块主要包括：

- 1) 性能指标计算：自动在验证集或测试集上计算模型的精确率（Precision）、召回率（Recall）、mAP（mean Average Precision）等核心指标。
 - 2) P-R 曲线绘制：生成 Precision-Recall 曲线，直观展示模型在不同置信度阈值下的性能表现。
 - 3) 损失曲线可视化：绘制训练过程中的各类损失（如定位损失、分类损失）曲线图，便于分析模型的收敛情况。
 - 4) 检测结果可视化：在样本图片或视频上绘制模型的预测框，直观展示检测效果。
 - 5) 评估报告生成：一键生成包含各项核心指标和可视化图表的综合评估报告。
- 所有评估结果支持导出和保存，便于不同模型或不同训练策略间的横向对比分析。

4.模型推理与部署

本模块旨在将训练好的模型应用于实际场景,提供高效的推理接口和便捷的部署方案。系统支持对模型进行优化和格式转换,以适应不同硬件平台的部署要求。

模型推理与部署模块主要包括:

- 1) 实时推理: 支持用户上传单张图片、视频文件或输入实时视频流,进行目标检测并实时返回结果。
- 2) 模型格式转换: 提供将 PyTorch 等框架训练出的模型转换为 ONNX、TensorRT 等标准化格式的工具,便于跨平台部署。
- 3) 模型优化: 支持模型剪枝、量化等压缩与加速技术,减小模型体积,提升在边缘设备或移动端上的推理速度。
- 4) 部署接口: 提供标准化的 API 接口(如 RESTful API),方便将检测功能集成到其他业务系统中。

通过本模块,用户可以快速完成从算法验证到产品化部署的全过程。

5.系统管理与运维

本模块为整个框架的稳定运行提供支撑,提供友好的用户交互界面和完善的后台管理功能,确保系统的易用性、安全性和可维护性。

系统管理与运维模块主要包括:

- 1) Web 用户界面: 提供图形化的 Web UI,用户可通过浏览器完成数据集管理、模型训练、结果查看等所有操作。
- 2) 任务管理与监控: 集中管理和监控所有训练、评估和推理任务的运行状态,实时查看资源占用情况(如 CPU、GPU)。
- 3) 日志管理: 详细记录系统运行日志、训练过程日志和错误信息,支持按级别、按时间查询,便于问题排查与追溯。
- 4) 模型版本管理: 对训练生成的不同版本模型进行统一存储和管理,记录每个模型的版本号、训练参数和性能指标。
- 5) 系统配置: 支持对系统级参数进行配置,如默认数据存储路径、可用 GPU 设备等。

第三部分 设计约束

一、需求约束

1. 本系统应当遵循的技术标准

1) 代码风格与可重现性

Python 代码遵循 PEP8/Black 格式化规范;PyTorch 模型实现遵循官方最佳实践和 PyTorch Lightning (如采用)

实验可复现性: 训练脚本必须可设置随机种子、记录依赖版本 (Python、PyTorch、CUDA、相关包) 并导出完整训练日志与配置文件 (YAML/JSON)。

2) 数据与标注规范

支持标准数据格式: COCO JSON、Pascal VOC XML、YOLO TXT 格式; 自定义数据需提供转换工具并遵循统一命名规则 (见下)。

3) 模块与导出标准

模型应支持导出为常用中间格式 (TorchScript、ONNX) 并保证在主流推理后端 (ONNX Runtime、TensorRT、OpenVINO) 上可运行。

4) 接口与协议规范

服务端推理与管理 API 使用 RESTful + JSON 为主, 必要时提供 gRPC 接口; 模型文件与配置通过安全的文件传输或 CI/CD 管道分发。

2. 软、硬件环境标准

1) 开发与训练环境

推荐操作系统: Ubuntu 20.04 LTS (或等效 Linux); Python 3.8+; PyTorch 与 CUDA 版本需在项目 README 中固定并可重建。

训练硬件: NVIDIA GPU (支持 CUDA 11.x) 或云端等效算力

2) 推理/部署环境

支持容器化部署 (Docker), 并提供 Dockerfile、docker-compose 和 Kubernetes manifests。

在边缘计算场景下, 部署需支持轻量级后端框架如 NCNN、TensorRT、TFLite 等, 并结合模型压缩和量化技术, 以优化模型性能并适应资源受限的边缘设备。

3) 依赖管理

使用虚拟环境或 Conda 管理依赖, 并提供锁定依赖文件 (requirements.txt / environment.yml)

3. 接口/协议标准

1) 数据接入与标注接口

提供数据导入 API / CLI, 支持批量导入、验证与预处理流水线。

2) 训练与推理接口

训练任务通过 CLI/REST 提交, 返回任务 ID 和日志访问地址; 推理提供同步 (REST)

与异步（消息队列）两种方式。

3) 运维与监控接口

提供 Prometheus 指标暴露端点与日志收集（ELK/EFK），并支持 Webhook 告警。

4. 用户界面

1) 管理与可视化

提供基于 Web 的管理界面（React/Vue），实现任务提交、训练监控、模型版本管理、评估报告与可视化（检测叠加、混淆矩阵），利用机器学习技术优化 Web 管理后台的识别与安全。

2) 可访问性

界面采用响应式设计，关键操作须有清晰提示与日志导出功能。

5. 软件质量

1) 正确性

训练与评估流程应保证可重复得到相同或近似结果；关键模块须有单元/集成测试覆盖。

2) 健壮性

作业调度与断点续训支持；训练/推理异常（OOM、节点失联）能触发重试或回退策略。

3) 效率性

提供性能基准：在指定硬件上给出训练速率、单帧推理延迟与吞吐量（例如边缘设备目标延迟 $\leq X$ ms，服务端吞吐量 $\geq Y$ FPS）。例如，NVIDIA A100 GPU 在推理测试中，与 3080 相比，可以提供 2.5 倍的推理吞吐量。

4) 易用性

提供完整文档、示例脚本和一键复现说明，API 与配置项尽量自解释并有合理默认值。

5) 安全性

模型与数据传输需使用 TLS；模型文件与关键配置使用签名或校验；对敏感数据（含原始视频）应有访问控制与脱敏策略。

6) 可扩展性

模块化设计（数据 \rightarrow 训练 \rightarrow 模型 \rightarrow 导出 \rightarrow 推理），支持插件式新增模型架构、后处理策略与评估指标。

7) 可维护性

代码须具备清晰模块边界、文档与 CI 流程，保证多人协作和长期维护。

二、隐含约束

1) 用户假设

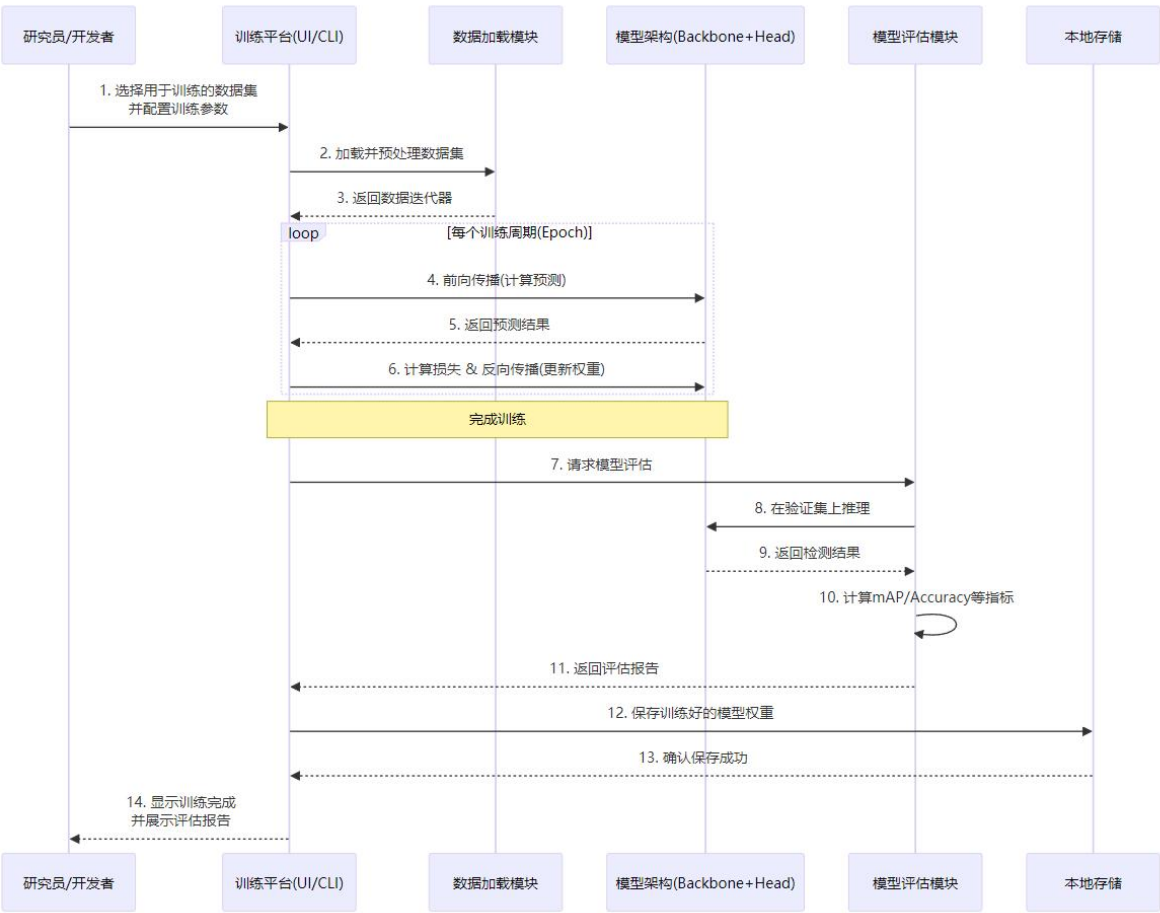
假设使用者具备基础机器学习/深度学习知识与常用命令行操作能力；运维人员熟悉容器与 Kubernetes。

- 2) 可重现性与可配置性
训练参数、数据路径、依赖版本应放入配置文件并纳入版本控制；实验结果可通过导出的配置 + 随机种子复现。
- 3) 合规与隐私
项目需遵守数据使用规范，对敏感图像/视频实施最小化保存与访问审计。

第四部分 功能结构设计

一、YOLO 代码复现框架相关业务流程

1.模型训练与评估流程



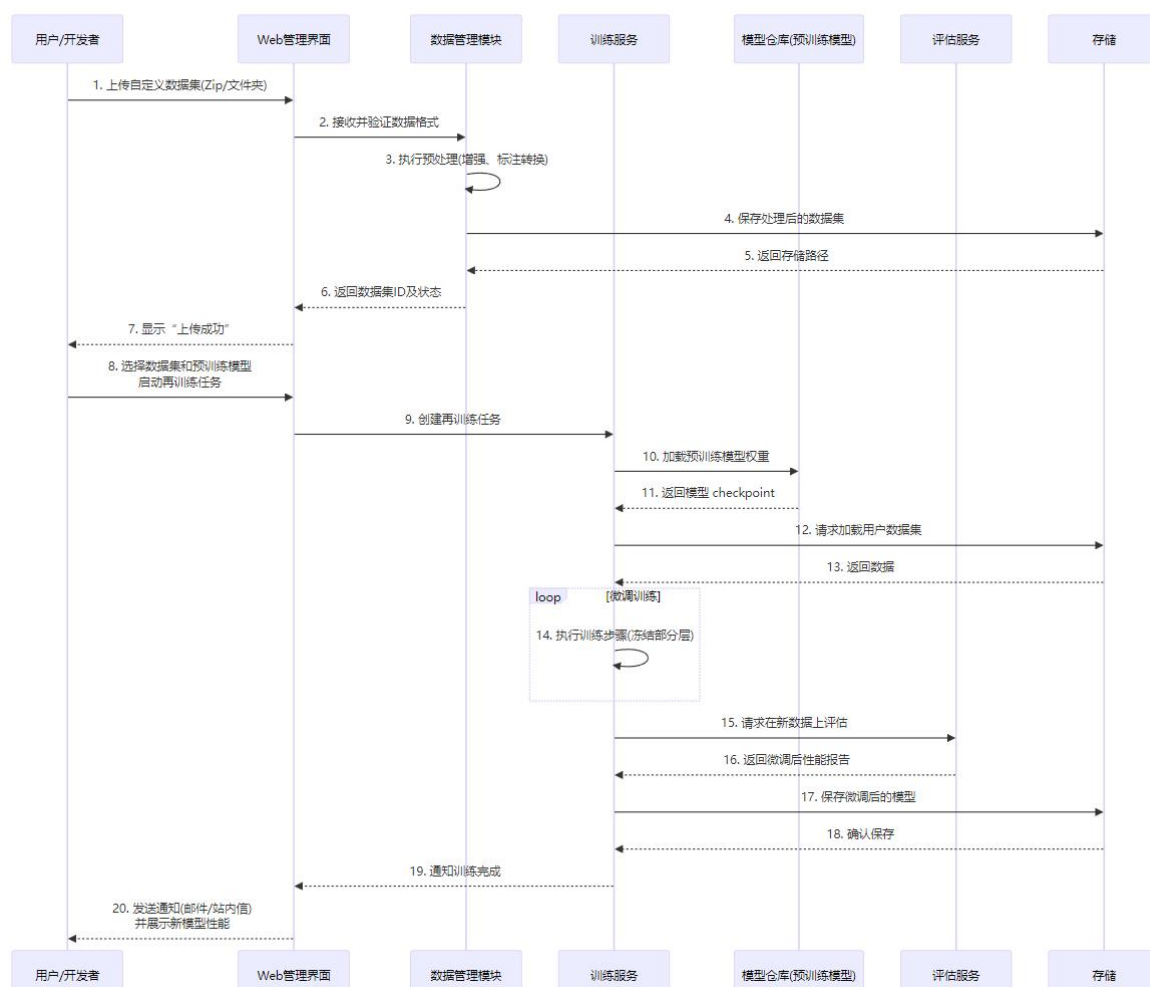
说明：

模型训练与评估流程包括：

- ✓ 开发者选择用于训练的数据集并配置相关的训练参数

- ✓ 对数据集中的图片进行数据增强与预处理
- ✓ 对预处理后的图片进行卷积等方式提取特征信息
- ✓ 模型训练并通过反向传播更新模型参数
- ✓ 模型使用测试集评估现有模型的准确率
- ✓ 保存模型到指定的路径

2. 模型微调流程



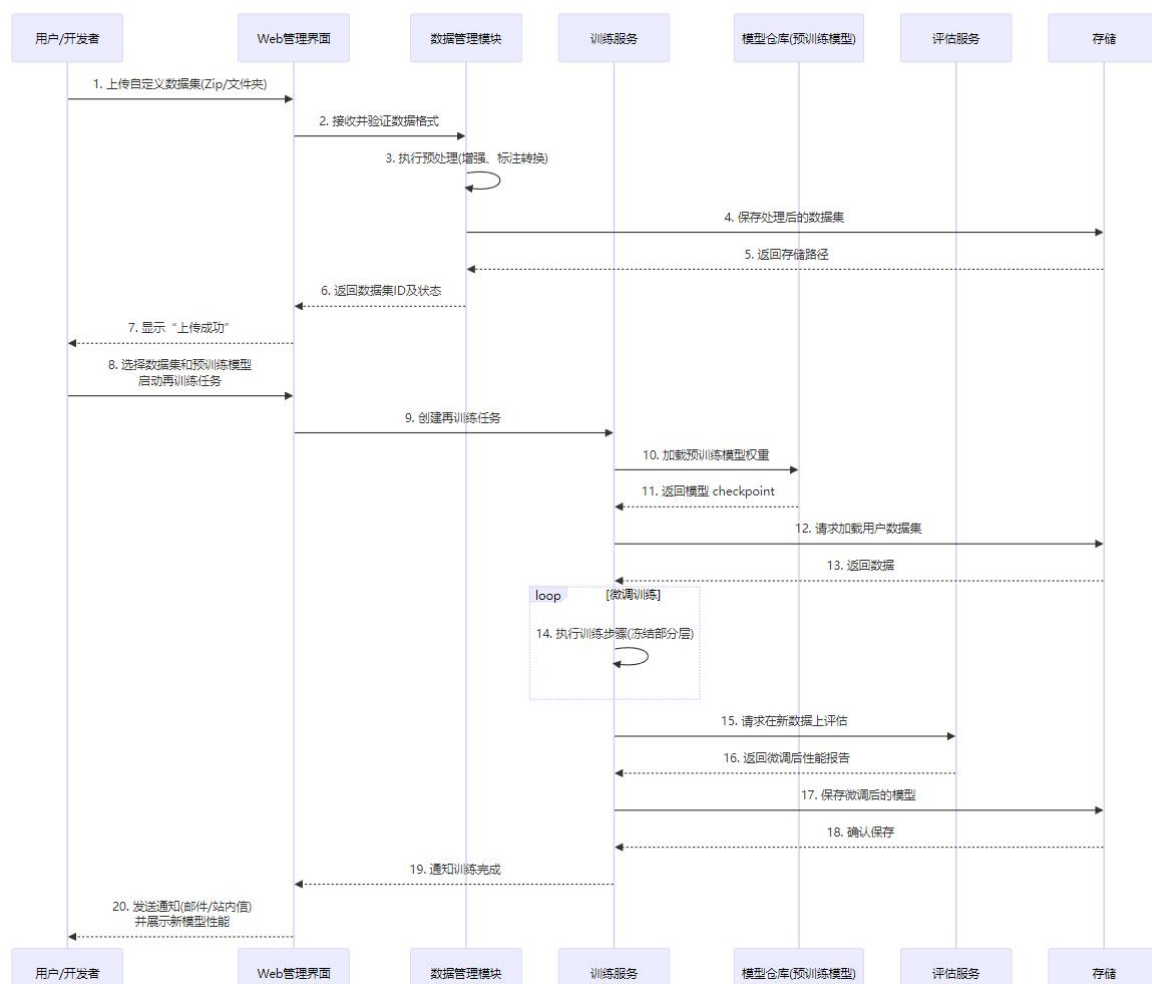
说明:

模型微调的主要业务流程包括:

- ✓ 上传自定义数据集;
- ✓ 接收并验证数据格式, 对数据进行预处理, 保存处理后的数据集;
- ✓ 返回处理好的数据集, 显示“上传成功”;
- ✓ 选择数据集预训练的模型, 加载数据集和预训练模型的权重, 启动二次训练任务;
- ✓ 训练结束后, 在新数据上评估模型性能, 并返回性能评估报告。

- ✓ 保存模型权重，发送通知。

3.Web 调用模型流程



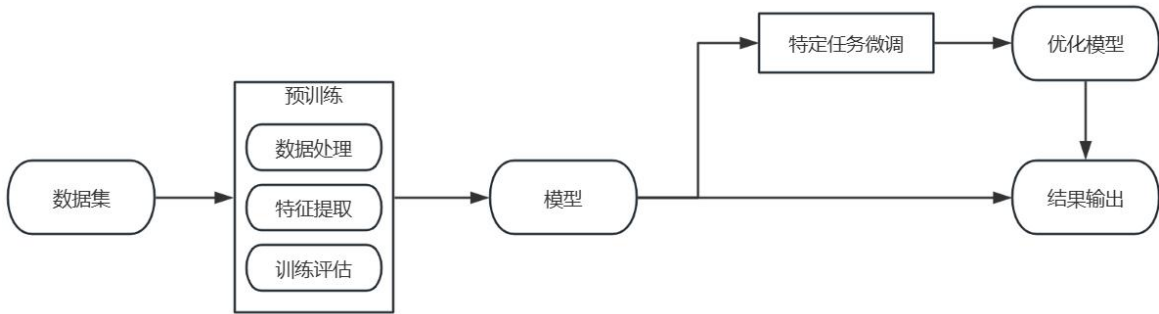
说明：

此流程用于将本地训练好的模型发布到云端模型仓库，便于部署、版本管理和团队协作。

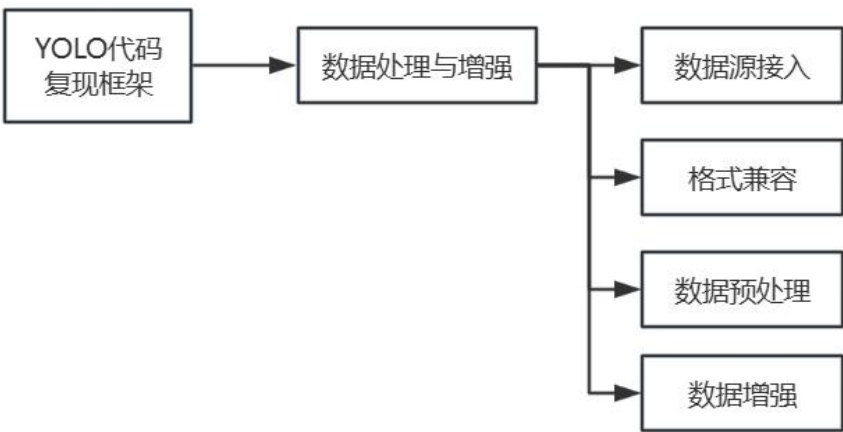
- ✓ 用户需手动选择本地模型文件以进行信息访问
- ✓ 平台自动解析模型并生成元数据
- ✓ 系统管理员手工补全模型元信息
- ✓ 平台自动校验模型兼容性与安全性
- ✓ 校验成功，模型发布至云端仓库
- ✓ 校验失败，返回错误信息

二、业务功能概要结构

逻辑结构图如下：



1.数据处理与增强模块



1.1、数据源接入与格式兼容

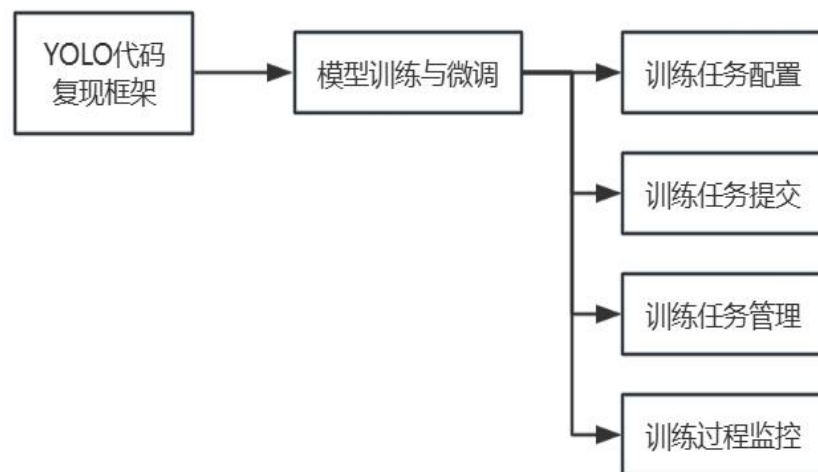
- 1.在系统初始化的时候，首先判定接入的数据类型，如图片或视频源等。如果是图片，则正常训练，如果是视频源，需要将视频转化为逐帧的图片。如果传入的格式不符合现有格式要求，需提示“你传入的文件类型有误，请检查”。
- 2.数据接入完成后，系统输出提示：“你的数据接入已完成”。

1.2、数据预处理与数据增强

- 1.数据接入完成后。

- 2.对已有数据进行清洗，去除过于模糊的图片和空图。
- 3.对输入图像进行尺寸归一化等基础处理。
- 4.对图像进行反转、裁剪等数据增强策略

2.模型训练与微调模块



2.1、训练任务配置与提交

- 1.用户新建训练任务，进入训练任务配置流程。
- 2.选择模型版本，选择一个 YOLO 算法版本。
- 3.选择训练模式，选择“从零开始训练”（预训练）或是“基于预训练权重微调”。
- 4.配置参数：设置学习率、优化器（如 SGD, Adam）、批次大小（Batch Size）、训练轮次（Epochs）等。
- 5.对参数进行校验，校验通过后将训练任务加入任务队列，并根据系统资源进行调度。

2.2、训练过程监控

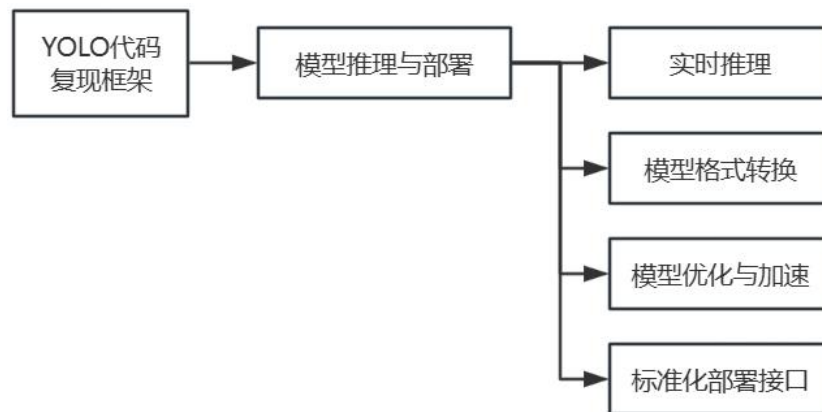
- 1.训练任务开始后，可查看指定训练任务的监控信息，主要包括：
 - 1) 损失曲线：以图表形式动态展示训练过程中损失函数变化。

2) 关键指标：实时显示学习率等关键训练指标。

2.3、训练任务管理

- 1.系统在训练过程中会定时自动保存训练检查点（Checkpoint），以支持断点续训。
- 2.可以随时执行“暂停”或“终止”操作。
- 3.若任务被暂停或意外中断，用户可选择从最新的检查点恢复训练，保障长周期任务的稳定性。

3.模型评估与可视化模块



3.1、标准化模型评估

- 1.对训练好的模型提供标准化行业评估指标。
- 2.加载训练好的模型，选择测试数据集，设置评估参数（NMS 阈值、置信度阈值）。
- 3.评估指标包括： $mAP@0.5$ （即 IoU 设为 0.5 时的平均精度均值）， $mAP@0.5:0.95$ （即在 IoU 从 0.5 到 0.95 不同阈值上的平均精度均值），以及精确率（Precision）、召回率（Recall）、F1 分数和交并比（IoU）分布。
- 4.利用特定场景对模型进行鲁棒性评估。

3.2、模型性能分析

- 1.测试模型在不同硬件条件下的性能

- 2.测试并计算运行时的资源消耗情况，如 CPU/GPU 利用率，内存/显存占用情况等。
- 3.测量模型推理速度（FPS）是通过计算模型每秒可以处理的帧数来实现的，例如，如果模型处理一张图片需要 50ms，则 FPS 为 $1000\text{ms}/50\text{ms}=20$ 。
- 4.测量延迟（Latency）。

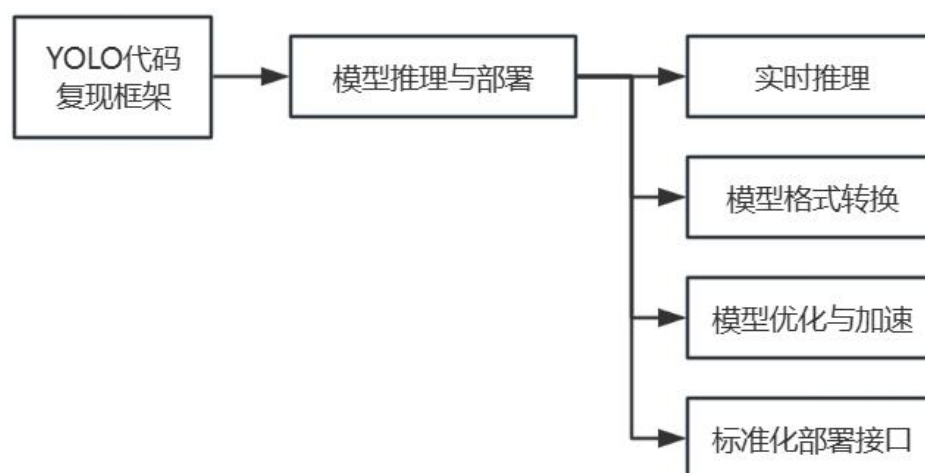
3.3、检测结果数据分析

- 1.汇总统计数据，包括识别目标数量、类别分布、目标尺寸分布。
- 2.对模型推理产生的数据进行分析，挖掘数据价值。

3.4、生成评估结果报告

- 1.整合模型相关信息、数据集信息。
- 2.整合模型测试结果、模型性能各类指标、训练曲线等信息。
- 3.绘制可视化统计图表，包括 mAP 曲线、PR 曲线等。
- 4.生成 PDF/HTML 格式的完整模型评估报告。

4.模型推理与部署模块



4.1、实时推理

- 1.提供 Web 界面或 API，允许用户上传单张图片（JPG/PNG）。

- 2.支持上传视频文件（MP4/AVI），系统进行逐帧分析。
- 3.支持接入 RTSP、HTTP-FLV 等主流网络视频流协议，或调用本地/USB 摄像头，进行实时帧-by-frame 分析，并低延迟地返回结果。
- 4.在图片/视频帧上绘制边界框、类别标签、置信度分数。
- 5.同步返回 JSON 格式的结构化结果，包含每个目标的坐标（xmin, ymin, xmax, ymax）、类别、置信度、时间戳（用于视频和流）。
- 6.对于视频流，提供 Web 端或客户端的实时检测结果预览界面。

4.2、模型格式转换

- 1.主要支持转换 PyTorch (.pt, .pth) 和 TensorFlow (.pb, .h5) 格式的训练模型。
- 2.能够将模型转换为标准的 ONNX 格式（.onnx），实现框架无关性。
- 3.针对 NVIDIA GPU 平台，进一步将 ONNX 模型优化并序列化为 TensorRT 引擎（.engine），最大化推理性能。
- 4.支持转换为 TFLite（用于移动端）、OpenVINO IR（用于 Intel 硬件）、CoreML（用于 Apple 设备）等格式。
- 5.提供工具验证转换后的模型与原始模型的预测输出是否一致，确保转换过程无损。

4.3 模型优化与加速

- 1.提供自动化工具，移除模型中冗余的权重或通道，生成更稀疏、更小的模型。
- 2.对已有模型进行 INT8 量化，显著减小模型体积、提升速度，精度损失可控。
- 3.在训练阶段模拟量化操作，获得更高精度的 INT8 模型。
- 4.自动将多个层（如 Conv-BN-ReLU）融合为单一算子，减少计算开销。
- 5.针对特定硬件（如 NVIDIA GPU、Intel CPU、ARM NPU）进行底层优化，充分利用硬件加速能力。

4.4、标准化部署接口

- 1.提供 POST /v1/detect 接口，接收图片，返回 JSON 结果。

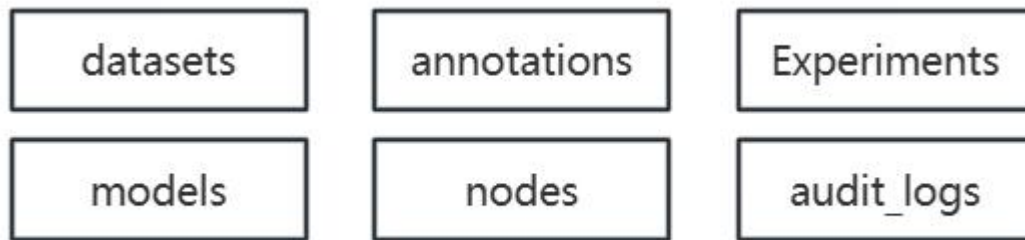
2. 提供 POST /v1/tasks 接口，用于提交视频文件等耗时任务，通过 GET /v1/tasks/{task_id} 查询任务状态和结果。
3. 提供交互式 API 文档，方便开发者测试和理解。
4. 为高吞吐、低延迟的内部服务间通信提供 gRPC 接口和 Protocol Buffers 格式的定义文件。
5. 针对视频流应用，提供 WebSocket 或 gRPC Stream 接口，建立双向通信通道，持续传输视频帧并接收检测结果。

三、模块定义

模块	组建	规格/型号
数据处理与增强	数据源接入与格式兼容	支持图片、视频源接入；视频自动逐帧转换；支持格式校验与提示。
	数据预处理与数据增强	支持模糊/空图片清洗；支持尺寸归一化；支持图像翻转、裁剪等增强策略。
模型训练与微调	训练任务配置与提交	支持 YOLO 不同版本选择；支持从零训练与微调模式；支持学习率、优化器、Batch Size 等超参数配置与校验。
	训练过程监控	实时图表展示损失函数曲线；实时显示学习率等关键指标。
	训练任务管理	自动保存 Checkpoint；支持任务暂停、终止、断点续训。
模型评估与可视化	标准化模型评估	计算 mAP, Precision, Recall, F1 Score 等指标；支持鲁棒性评估。
	模型性能分析	计算 CPU/GPU 利用率、内存/显存占用；测量推理速度（FPS）与延迟（Latency）。
	检测结果数据分析	统计识别目标数量、类别分布、尺寸分布；支持数据价值挖掘。
	生成评估结果报告	整合模型/数据集信息、指标、曲线图；生成 PDF/HTML 格式报告。
模型推理与部署模块	实时推理	支持图片（JPG/PNG）、视频（MP4/AVI）、网络流（RTSP）输入；同步返回 JSON 结果与可视化预览。
	模型格式转换	支持 PyTorch/TF 转 ONNX、TensorRT、TFLite 等；提供转换一致性验证
	模型优化与加速	支持剪枝、INT8 量化、算子融合及特定硬件优化。
	标准化部署接口	提供 RESTful API (POST /v1/detect)、gRPC、WebSocket 接口；提供交互式 API 文档。

第五部分 E-R 实体设计

一、E-R 实体结构图



二、实体描述

1.datasets 实体描述

编号	英文名	中文名	数据类型
1	id	数据集 ID	NSInteger
2	name	数据集名称	NSString
3	owner_id	所有者 id	NSInteger
4	format	数据格式	NSString
5	storage_path	存储路径	NSString
6	status	状态	NSString
7	num_images	图片数量	NSInteger
8	description	说明	NSString
9	created_at	创建时间	NSDate

2.annotations 实体描述

编号	英文名	中文名	数据类型
1	id	标注 ID	NSInteger
2	image_id	图片 ID	NSInteger
3	annotator_id	标注人 ID	NSInteger
4	boxes	标注框	NSMutableArray
5	status	状态	NSString
6	created_at	创建时间	NSDate
7	verified_at	审核时间	NSDate
8	notes	备注	NSString

3.Experiments 实体描述

编号	英文名	中文名	数据类型
1	id	实验 ID	NSInteger/NSNumber (PK)
2	name	实验名称	NSString
3	user_id	提交用户 ID	NSInteger/NSNumber (FK -> users.id)
4	dataset	训练数据集 ID	NSInteger/NSNumber (FK -> datasets.id)
5	config	训练配置	NSString/NSDictionary(JSON 存储)
6	status	状态	NSString(queued/running/success/failed)
7	metrics	训练/评估指标	NSString/NSDictionary ryNSString
8	checkpoint_path	checkpoint 路径	NSString

9	start_at	开始时间	NSDate
10	end_at	结束时间	NSDate
11	created_at	创建时间	NSDate

4.models 实体描述

编号	英文名	中文名	数据类型
1	id	模型 ID	NSInteger/NSNumber(PK)
2	experiment_id	来源实验 ID	NSInteger/NSNumber(PK)
3	version	版本号	NSString
4	format	格式	NSString(pt/onnx/tensorrt/tflite)
5	storage_path	存储路径	NSString
6	size_bytes	大小（字节）	NSInteger/NSNumber
7	metrics	评估指标	NSString/NSDictionary
8	checksum	校验和	NSString/NSDictionary
9	signed_by	签名用户	NSString
10	created_at	创建时间	NSDate
11	note	备注	NSString

5.nodes 实体描述

编号	英文名	中文名	数据类型
1	id	节点 ID	NSInteger

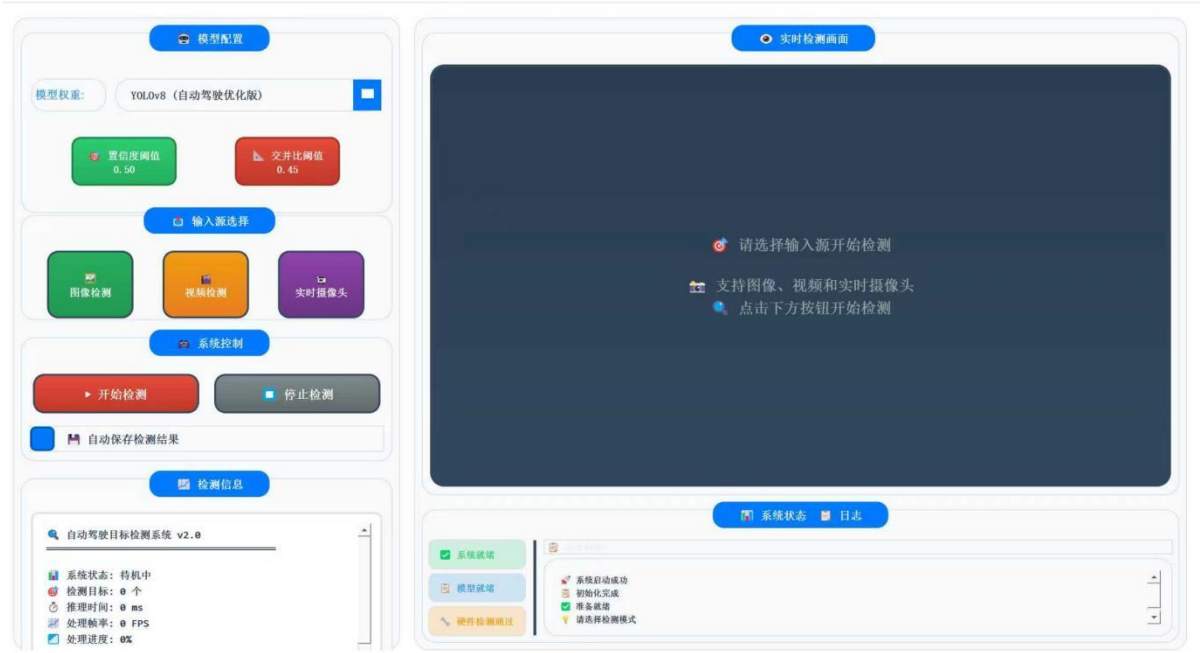
2	node_key	节点唯一标识	NSString
3	hostname	主机名	NSString
4	ip	IP 地址	NSString
5	capabilities	能力	NSDictionary
6	public_key	公钥 / 证书信息	NSString
7	status	状态	NSString
8	last_time	最后使用时间	NSDate
9	registered_by	注册用户 ID	NSInteger
10	created_at	注册时间	NSDate

6.audit_logs 实体描述

编号	英文名	中文名	数据类型
1	id	ID	NSInteger/ NSNumber (PK)
2	user_id	操作用户 ID	NSInteger/ NSNumber (FK -> users.id, 可空)
3	action	操作类型	NSString
4	target_table	目标表名	NSString
5	target_id	目标记录 ID	NSInteger/ NSNumber
6	details	操作详情	NSString/ NSDictionary (JSON)
7	ip	操作 ip	NSString
8	created_at	发生时间	NSDate

第六部分 用户界面设计

一、UI 展示



第七部分 总体设计

一、模型架构

1.训练器

负责模型训练的调度、参数管理、损失计算、优化器更新和梯度传播，保存模型权值文件和训练参数配置文件。使用 YAML 或 JSON 文件管理模型参数，管理日志、早停、模型训练点检查。

2.模型优化板块

集成 PyTorch 的 torch.quantization, torch.nn.utils.prune 等，或第三方库。

3.模型评估

负责模型推理能力的测试，对目标识别的结果进行分析，挖掘数据价值，生成模

型性能评估报告。

4. 部署与本地化

通过 Python API 提供直接训练好的模型进行推理。

二、数据流转框架

1.数据预处理

该部分负责通过定义 **Dataset** 类，接收各类型的训练数据并裁剪为同一大小的张量，通过色彩混合、翻转等数据增强手段对数据进行预处理。

2.数据加载

采用多进程 GPU 加速技术，借助 **Dataloader** 将训练数据输入模型。

第八部分：运行环境与部署

一、运行环境

1.客户机器环境

PC, Windows 7/10/11, Chrome/Firefox/Edge 浏览器

2.开发环境要求

项目	名称	版本
开发平台	Mac OS	10.7+
	Windows	10
开发工具	VS Code	-
	Pycharm	-
后端服务	Python	3.9+
	PyTorch	2.0+
服务器	Ubuntu	20.04 LTS

二、系统性能要求

#	项目	模块	级别	技术参数
1	YOLO 模型复现与训练	YOLO 模型训练速度	A	单卡 RTX3090: ≤ 24 小时/epoch。
2		GPU 利用率	A	训练过程中 $\geq 80\%$ 。
3		CPU 利用率	A	训练过程中 $\leq 50\%$ 。
4		训练稳定性	B	训练过程中损失曲线平滑，无异常震荡或爆炸。
5		配置加载/保存时间	A	加载/保存训练配置时间 $\leq 100\text{ms}$ 。
6	数据管理与预处理	数据加载吞吐量 (CPU)	A	单进程未经 GPU 加速时 $\geq 500\text{img/s}$
7		数据加载吞吐量 (GPU)	A	多进程经 GPU 加速时 $\geq 1500\text{img/s}$
8		数据增强算法效率	A	Mosaic/Mixup 等复杂处理单张时间 $\leq 10\text{ms}$
9		内存占用	B	数据加载和预处理过程内存占用 $\leq 1\text{GB}$ (batch size = 64)
11	模型评估与分析	mAP 计算速度	A	评估 1000 张图片 ≤ 5 分钟 (GPU)
12		推理速度	A	服务器端 TensorRT 优化 (GTX 1080 Ti, 640x640) $\geq 80\text{FPS}$, 移动端 Core ML (A14 Bionic, 640x640) $\geq 15\text{FPS}$ 。
13		报告生成时间	B	生成完整的评估报告 ≤ 1 分钟。
14		数据分析图表生成	C	生成 10000 张图片识别结果图表 ≤ 30 秒。
15	模型优化与部署	模型导出速度	A	ONNX/TFLite/Core ML 导出 ≤ 5 分钟
16		办理服务吞吐量	B	$\geq 100\text{QPS}$ (每秒查询数)。
17		API 服务启动时间	C	推理服务部署/启动 ≤ 1 分钟。
18	集成部署环境	服务器	A	Ubuntu 20.04 LTS
19		PC 电脑	A	Microsoft Windows 10 及以上

说明：级别（A：表示非常重要必须达到的技术性能要求，B：表示重要推荐达到的技术性能要求，C：表示非重要可以弱化的技术性能要求。）