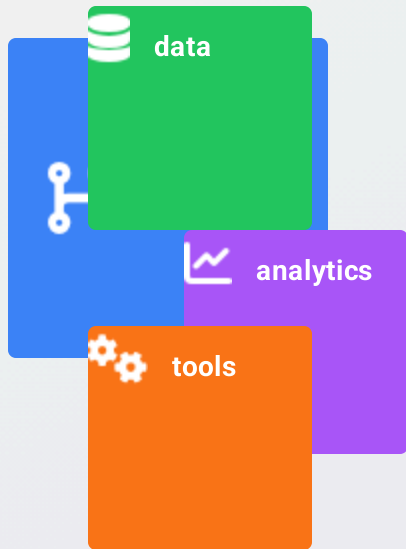


```
git clone github.com/organization/monorepo.git
```



Monorepo Strategy for Data Team

GitHub Implementation with UV Library

Streamlining code management and collaboration for our data team

```
git clone github.com/organization/monorepo
```

→

```
uv install -r requirements.txt
```

什么是 Monorepo 以及为什么考虑它？

🏗️ Monorepo 定义

Monorepo（单一代码库）是一种版本控制策略，将**多个项目、服务和库的代码存储在单一的中央仓库中**。

这与多仓库（Polyrepo）方法形成对比，后者是我们当前的方法，每个项目或服务维护自己的独立仓库。

⚠️ 当前挑战

🧩 代码分散在多个仓库中，难以保持一致性

🔗 依赖管理复杂，容易出现"依赖地狱"

👥 团队协作和代码共享受限

🔗 跨项目重构困难，难以保证一致性

多仓库 (Polyrepo)



仓库 A 仓库 B 仓库 C

🔗 服务间依赖难以管理

单一代码库 (Monorepo)



单一仓库

✓ 集中化管理，简化依赖

Monorepo 核心优势

- ✓ 简化依赖管理，共享库与服务就近存放
- ✓ 增强团队协作与代码共享
- ✓ 统一版本控制与 CI/CD 流水线

Monorepo 核心优势



简化依赖管理

共享库和服务就近存放，便于管理跨项目依赖。对共享组件的更新可以立即针对所有依赖项目进行测试，减少集成问题。

- ✔ 减少依赖版本冲突



增强协作与代码共享

打破团队之间的壁垒，促进对各种组件的更轻松访问、审阅和贡献。常见模型、共享库和辅助代码可以轻松共享，减少重复并促进统一开发文化。

- ✔ 减少代码重复



统一版本控制与 CI/CD

实现所有项目的统一版本控制，确保兼容性并简化更改跟踪。允许在所有项目中使用单一一致的 CI/CD 流水线，简化部署流程并减少开销。

- ✔ 简化发布流程



原子提交与重构

执行应用范围的重构变得更加简单，因为所有代码修改都可以在单个拉取请求下提交，确保受影响库之间的一致性。

- ✔ 确保代码一致性

潜在挑战与缓解策略



可扩展性问题

随着仓库规模增长，可能出现构建时间变慢、管理复杂度增加等问题。

影响：构建时间延长，开发效率下降

缓解策略

使用智能构建工具，如 Turborepo 或 Nx，实现增量构建和缓存。



访问控制复杂性

在单一仓库中管理不同团队的访问权限可能变得复杂。

影响：权限管理困难，安全隐患

缓解策略

利用 GitHub 的 CODEOWNERS 文件，为特定目录设置代码所有者和审查流程。



紧密耦合风险

如果不遵循最佳实践，可能导致不相关的项目之间出现紧密耦合。

影响：增加重构难度，阻碍独立扩展

缓解策略

建立清晰的项目边界和API约定，使用工作区 (workspaces) 管理依赖关系。



关键洞察

这些挑战可以通过适当的工具选择和明确的实践规范得到有效缓解。Monorepo 的优势通常超过这些挑战，特别是对于数据团队这样协作紧密的团队。

GitHub Monorepo 实现：目录结构

建议的目录结构

monorepo/

apps/

dashboard-app/

api-service/

reporting-tool/

packages/

data-connectors/

utility-functions/

ml-models/

docs/

project-overview.md

setup-guide.md

config/

dev.json

prod.json

目录说明

/apps

包含所有可部署的应用程序，如交互式数据仪表盘、API服务和报告工具。

/packages

包含共享库和可重用组件，如数据连接器、实用函数和机器学习模型。

/docs

专门用于整个monorepo的文档，包括项目概述、设置指南和API规范。

/config

集中化的配置文件，可能在项目间共享或模板化。

结构优势

- ✓ 促进模块化和可发现性
- ✓ 简化团队成员定位和贡献代码的过程

访问控制和 CI/CD 策略

访问控制 (CODEOWNERS)

GitHub 提供了 **CODEOWNERS** 文件来管理访问权限，放置在仓库根目录。

```
# CODEOWNERS 文件示例
# api-service 目录由 API 团队管理
/apps/api-service/ @api-team
# data-connectors 包由数据工程团队管理
/packages/data-connectors/ @data-engineering-team
```

CODEOWNERS 的主要优势：

- ✔ **责任明确**：为每个组件建立明确的所有权
- ✔ **质量控制**：由最了解代码的团队进行审查
- ✔ **安全保障**：防止对关键部分的未经授权修改




CI/CD 策略 (GitHub Actions)

使用路径过滤器配置 GitHub Actions，仅在特定文件或目录修改时运行工作流：

```
# .github/workflows/ci.yml 文件示例
name: CI
on:
  pull_request:
  paths:
    - 'apps/dashboard-app/**'
    - 'packages/utility-functions/**'
```



高效 CI/CD 的优势：

-  **更快的反馈循环**：开发者获得更快的更改反馈
-  **减少资源消耗**：避免不必要的构建和测试
-  **优化部署**：仅更改的应用程序或库被构建和部署

引入 UV: 下一代 Python 包管理器

UV 是一个用 Rust 编写的极快的 Python 包安装器和解析器，旨在通过提供传统工具的高性能替代方案来简化 Python 打包流程。



极速性能

比传统 pip 快 10-100 倍，特别是在使用热缓存时表现卓越。



统一工具链

作为单一工具，旨在替代 pip、pip-tools、pipx、virtualenv，甚至可能替代 poetry 和 pyenv。



高效磁盘使用

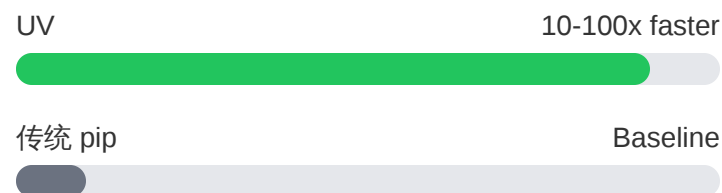
使用全局缓存进行依赖去重，最小化存储需求。



全面项目管理

支持锁文件、工作区和脚本内的内联依赖元数据。

性能对比



i 基于 warm cache 的测试结果



使用示例

```
$ uv pip install -r requirements.txt
$ uv sync
$ uv run python script.py
```


✔ 与现有 pip 生态系统兼容

UV 与现有工具对比


选择合适的 Python 包管理器对提高开发效率至关重要。以下是 UV 与常见工具的性能和功能对比：

特性	 uv	 pip + virtualenv	 Conda	 Poetry
实现语言	Rust	Python	Python	Python
速度	10-100x faster than pip	基准	Slower than pip	Faster than pip
内存使用	Very efficient	Higher	High	Moderate
环境管理	Built-in	需要单独工具	Built-in	Built-in
依赖解析	Fast, modern resolver	Basic	Comprehensive	Modern resolver
非 Python 包	No	No	Yes	No
锁定文件	Yes	No (basic `requirements.txt`)	Yes	Yes
项目结构	Yes	No	No	Yes

💡 关键优势

 **极速安装**：比传统 pip 快 10-100 倍

 **统一工具链**：替代 pip、pip-tools、pipx 等

 **高效磁盘使用**：全局依赖缓存 deduplication

必要的 Monorepo 管理工具

除了 Python 特定的包管理工具外，还有多种工具可以管理 monorepo 中固有的复杂性，特别是在跨多个项目时的构建、依赖和任务执行。



Nx

强大的构建系统，主要针对 TypeScript monorepo，但适用于多种语言。

核心功能：

- 依赖图可视化
- 受影响的项目检测
- 代码生成
- 计算缓存

适用于大型 TypeScript/JavaScript 项目



Turborepo

高性能构建系统，针对 JavaScript 和 TypeScript 代码库，注重快速任务缓存和并行执行。

核心功能：

- 任务并行执行
- 高效任务缓存
- 最小配置需求
- 增量构建

适用于需要快速构建执行的大型 JS/TS 项目



Lerna

专门用于管理具有多个包的 JavaScript 项目的工具。帮助处理版本管理和发布 workflow。

核心功能：

- 版本管理
- 发布 workflow
- 工作区支持
- 依赖更新

适用于管理前端组件或数据应用中的共享 JavaScript 库



mise

提供实验性的 monorepo 任务支持，允许统一的任务命名空间、智能工具和环境继承。

核心功能：

- 统一任务命名空间
- 智能工具和环境继承
- 强大的通配符模式
- 语言无关性

适用于需要平衡简洁与功能的项目

结论与下一步

✔ Monorepo 核心优势总结

- ✔ 简化依赖管理：共享库与服务就近存放，更新可立即测试
- ✔ 增强协作：打破团队隔阂，促进代码共享与贡献
- ✔ 统一版本控制：确保兼容性，简化发布流程
- ✔ 原子提交：支持跨项目重构，确保一致性

💡 为什么现在是时候行动？

Monorepo 提供了我们当前代码管理方法所缺乏的**一致性、协作和效率**。结合 UV 等现代工具，我们可以显著提升数据团队的开发体验。

👉 准备迎接变革

▶▶ 实施计划

试点项目：服务迁移

- 1 选择试点服务**
选择 2-3 个现有服务进行迁移，作为试点项目

dashboard-app api-service reporting-tool
- 2 设计目录结构**
根据之前讨论的结构创建目录组织

```
/apps  
/packages  
/docs  
/config
```
- 3 配置工具与 workflow**
设置 UV 依赖管理与 CI/CD 流水线
- 4 评估与调整**
收集反馈，优化流程，为全面采用做准备