

# PR #43925 完整报告

vllm-project/vllm

[CI] Enable prefix caching in BFCL benchmark

合并时间: 2026-05-29 07:36

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/43925>

## 执行摘要

- 一句话: 启用 BFCL benchmark 的 prefix caching
- 推荐动作: 简单有效的小优化, 无需精读。

## 功能与动机

BFCL 是多轮 benchmark, prefix caching 能大幅提高吞吐量。PR body 中提到 'Because it's multiturn benchmark, prefix caching will greatly increase up benchmark throughput.'

## 实现拆解

1. 启用 prefix caching: 在 `.buildkite/scripts/tool_call/run-bfcl-eval.sh` 中, 将 vLLM 服务启动参数从 `--no-enable-prefix-caching` 改为 `--enable-prefix-caching`。
2. 修复安装命令: 将 `pip install` 替换为 `uv pip install`, 以使用更快的 UV 包管理器安装 `bfcl-eval`。

关键文件:

- `.buildkite/scripts/tool_call/run-bfcl-eval.sh` (模块 CI 脚本; 类别 other; 类型 core-logic)  
: 唯一变更文件, 包含两处修改: 启用 prefix caching 和改用 uv pip。

关键符号: 未识别

## 关键源码片段

`.buildkite/scripts/tool_call/run-bfcl-eval.sh`

唯一变更文件, 包含两处修改: 启用 prefix caching 和改用 uv pip。

```
# 安装 bfcl-eval 包: 从 pip 切换到 uv pip 以获得更快安装
uv pip install "bfcl-eval>=2025.10.20.1,<2026"
```

```
# 启动 vLLM 服务时启用 prefix caching, 加速多轮对话 benchmark
SERVE_ARGS=(
  --tensor-parallel-size "$TP_SIZE"
  --max-model-len "$MAX_MODEL_LEN"
  --enforce-eager
  --enable-prefix-caching # 之前为 --no-enable-prefix-caching
```

)

## 评论区精华

无讨论。

- 暂无高价值评论线程

## 风险与影响

- 风险：风险极低。变更仅涉及 CI 脚本中的两行参数，不影响核心代码。prefix caching 是 vLLM 已支持的特性，且测试表明准确率在运行间方差内，未出现回归。
- 影响：影响范围仅限于 BFCL 基准测试的 CI 运行。时间缩短约 50%，可提升 CI 效率。
- 风险标记：暂无

## 关联脉络

- 暂无明显关联 PR