

PR #42710 完整报告

vllm-project/vllm

[MRV2][XPU] add Model Runner V2 log

合并时间: 2026-05-17 12:15

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/42710>

执行摘要

- 一句话: XPU Worker 新增 V2 Model Runner 日志
- 推荐动作: 值得合并。变更简洁、无风险, 为 XPU 后端调试提供便利。可作为 V2 Model Runner 迁移状态的一个轻量级 markers。

功能与动机

PR body 中提及 'add Model Runner V2 log in XPU worker', 目的是在 XPU (Intel GPU) 后端中增加日志, 以便在运行时明确知晓是否启用了 V2 Model Runner, 方便开发和调试。

实现拆解

在 `vllm/v1/worker/xpu_worker.py` 的 `init_device` 方法中, 于 `oneCCL` 预热之后、设置随机种子之前, 插入条件判断:

1. 在 `torch.distributed.all_reduce` 之后、`set_random_seed` 之前, 添加 `if self.use_v2_model_runner: logger.info_once("Using V2 Model Runner")`。
2. 使用 `info_once` 确保该日志只在进程生命周期内输出一次, 避免反复打印。

变更仅涉及 3 行新增, 无删除。

关键文件:

- `vllm/v1/worker/xpu_worker.py` (模块 XPU Worker; 类别 source; 类型 core-logic) : 唯一的变更文件, 新增一行日志记录 V2 Model Runner 启用状态。

关键符号: 未识别

关键源码片段

`vllm/v1/worker/xpu_worker.py`

唯一的变更文件, 新增一行日志记录 V2 Model Runner 启用状态。

```
# vllm/v1/worker/xpu_worker.py (片段)

# 全局 all_reduce 用于 overall onecccl warm up
if torch.distributed.is_xccl_available():
    torch.distributed.all_reduce(torch.zeros(1).xpu())
```

```
# 新增日志：当使用 V2 Model Runner 时输出一次
if self.use_v2_model_runner:
    logger.info_once("Using V2 Model Runner")

# 设置随机种子
set_random_seed(self.model_config.seed)
```

评论区精华

无实质 review 讨论。xinyu-intel 和 jikunshang 均直接批准，仅 bot 自动评论无进一步反馈。

- 暂无高价值评论线程

风险与影响

- 风险：风险极低。新增的日志语句不会影响任何控制流或性能，仅在配置使用 V2 Model Runner 时多打印一行信息，且使用 info_once 限制为单次输出。不会引入回归或安全风险。
- 影响：影响范围极小：仅影响 XPU（Intel GPU）后端的调试日志输出。对系统行为、兼容性、性能均无影响。开发者可借此快速确认模型运行器版本。
- 风险标记：暂无

关联脉络

- 暂无明显关联 PR