

PR #42968 完整报告

vllm-project/vllm

[Feature] Add `--cpu-distributed-timeout-seconds` CLI Option for CPU Process Group Timeout

合并时间: 2026-05-22 06:42

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/42968>

执行摘要

- 一句话: 为 CPU 通信组新增独立超时配置项
- 推荐动作: 该 PR 设计清晰, 命名合理, 向后兼容, 建议合并。读者可关注其参数命名讨论和工具函数提取模式, 体现的“独立关注点分离”设计值得借鉴。

功能与动机

vLLM 初始化分设备 (NCCL) 和 CPU (gloo) 两种进程组用于分布式执行。现有的 `--distributed-timeout-seconds` 仅作用于设备通信 (见 PR#36047)。而 CPU 进程组仍依赖 PyTorch 硬编码的 1800 秒默认超时。通过引入 `--cpu-distributed-timeout-seconds`, 运维人员可为 CPU 侧协调配置更短的超时, 在不影响 NCCL 超时语义的前提下实现快速故障恢复。

实现拆解

1. 配置定义: 在 `vllm/config/parallel.py` 的 `ParallelConfig` 类中新增 `cpu_distributed_timeout_seconds: int | None = None` 字段, 并添加文档字符串说明默认使用 Gloo 的 1800 秒。
2. CLI 注册与值传递: 在 `vllm/engine/arg_utils.py` 的 `EngineArgs` 中添加同名属性, 并在 `add_cli_args` 中注册 `--cpu-distributed-timeout-seconds` 参数; 在 `create_engine_config` 中将该值传入 `ParallelConfig` 构造函数。
3. 工具函数提取: 在 `vllm/distributed/utils.py` 中新增 `get_cpu_distributed_timeout_or_none()` 函数, 通过当前 vLLM 配置读取超时秒数并返回 `timedelta | None`; 同时修改 `stateless_init_torch_distributed_process_group`, 在 backend 为 gloo 时调用此函数并覆盖该分支的默认 timeout。
4. 集成到分布式组初始化: 在 `vllm/distributed/parallel_state.py` 的 `GroupCoordinator.__init__` 中导入并使用同一工具函数, 将返回的超时值传递给 `torch.distributed.new_group(..., timeout=timeout)`, 从而控制所有 gloo 后端的进程组超时。

未添加自动化测试, 但通过手动注入 `sleep` 验证了超时生效。

关键文件:

- `vllm/config/parallel.py` (模块 并行配置; 类别 source; 类型 core-logic; 符号 `cpu_distributed_timeout_seconds`): 定义 `cpu_distributed_timeout_seconds` 字段, 是配置数据的源头。

- `vllm/engine/arg_utils.py` (模块 引擎参数; 类别 `source`; 类型 `core-logic`; 符号 `cpu_distributed_timeout_seconds`) : 注册 CLI 参数并将值传入 `ParallelConfig`, 是用户接口层。
- `vllm/distributed/parallel_state.py` (模块 分布式状态; 类别 `source`; 类型 `dependency-wiring`; 符号 `GroupCoordinator.init`) : 在 `GroupCoordinator` 初始化中实际应用超时到 `gloo` 组, 是运行时生效的关键点。
- `vllm/distributed/utils.py` (模块 分布式工具; 类别 `source`; 类型 `core-logic`; 符号 `get_cpu_distributed_timeout_or_none, stateless_init_torch_distributed_process_group`) : 新增核心工具函数 `get_cpu_distributed_timeout_or_none`, 并修改了 `stateless_init_torch_distributed_process_group` 中的 `gloo` 超时逻辑, 是复用关键。

关键符号: `get_cpu_distributed_timeout_or_none, stateless_init_torch_distributed_process_group, GroupCoordinator.init`

关键源码片段

`vllm/distributed/parallel_state.py`

在 `GroupCoordinator` 初始化中实际应用超时到 `gloo` 组, 是运行时生效的关键点。

```
# 在 GroupCoordinator.__init__ 中, 创建 gloo 进程组时应用独立超时
from vllm.distributed.utils import get_cpu_distributed_timeout_or_none

timeout = get_cpu_distributed_timeout_or_none()

for ranks in group_ranks:
    device_group = torch.distributed.new_group(
        ranks, backend=torch_distributed_backend
    )
    with suppress_stdout():
        # 将超时参数传给 gloo 后端; None 表示使用 PyTorch 默认 (1800 秒)
        cpu_group = torch.distributed.new_group(
            ranks, backend="gloo", timeout=timeout
        )
```

`vllm/distributed/utils.py`

新增核心工具函数 `get_cpu_distributed_timeout_or_none`, 并修改了 `stateless_init_torch_distributed_process_group` 中的 `gloo` 超时逻辑, 是复用关键。

```
def get_cpu_distributed_timeout_or_none() -> timedelta | None:
    """
    从全局 vLLM 配置中读取 CPU 分布式超时秒数,
    若配置存在且设置了值则返回 timedelta, 否则返回 None。
    """
    # 延迟导入以避免循环依赖
    from vllm.config import get_current_vllm_config_or_none
    vllm_config = get_current_vllm_config_or_none()
    if vllm_config is None:
        return None
```

```
timeout_seconds = vllm_config.parallel_config.cpu_distributed_timeout_seconds
# 仅当显式设置时才返回 timedelta
return timedelta(seconds=timeout_seconds) if timeout_seconds is not None else None

# 在 stateless_init_torch_distributed_process_group 中使用
timeout = _get_default_timeout(backend)
# 对 gloo 后端尝试使用独立的 CPU 超时
if backend == "gloo":
    gloo_timeout = get_cpu_distributed_timeout_or_none()
    if gloo_timeout is not None:
        timeout = gloo_timeout
```

评论区精华

- 命名讨论 (tlrmchlsmth) : 建议将 `--gloo-timeout-seconds` 改为 `--cpu-distributed-timeout-seconds` 以保持通用性。hmellor 赞同，作者接受并修改。
- 重构建议 (SageMoore) : 建议消除两处重复的配置读取逻辑，提取为 `get_cpu_distributed_timeout_or_none` 工具函数。作者采纳并实现。
- 超时传递争议 (gemini-code-assist) : 指出直接传递 `timeout=None` 可能覆盖 PyTorch 默认行为。作者回应称 `timeout=None` 与省略参数等价，最终版本保留直接传递方式。
- 参数命名讨论: gloo vs cpu-distributed (design): 决定采用 `--cpu-distributed-timeout-seconds`，以保持与设备超时解耦并适应未来可能的后端。
- 提取工具函数消除重复 (design): 作者接受建议，提取出函数并替换两处调用。
- 传递 `timeout=None` 的影响 (correctness): 作者解释 `timeout=None` 与省略参数在 `new_group` 中行为一致，维持原设计。代码审查者未进一步反对。

风险与影响

- 风险:
 - 向后兼容: 低风险。新参数默认 `None`，不影响现有逻辑。
 - 误配置风险: 若设置过短超时（如 1 秒），在高延迟环境（如多节点）中可能引发非预期的超时，但这属于用户操作问题。
 - 维护一致性: 未来若引入更多通信后端（如 NCCL 的其他变体），需要确保与之类似的独立超时参数保持一致设计。
- 影响:
 - 用户: 进行多节点分布式部署、使用 CPU 通信进行协调的用户（如 DP 元数据交换）可获得更灵活的故障恢复配置。
 - 系统: 不影响正常的推理性能；仅在故障时影响超时检测速度。
 - 团队: 新增一个 CLI 选项和一个内部工具函数，维护负担较小。
 - 风险标记: 向后兼容保障，超时覆盖默认行为，新增配置项

关联脉络

- PR #36047 Distributed timeout option: Referenced in PR body as the original timeout that only applies to device communication; this PR is a complementary extension for CPU.