

PR #38011 完整报告

vllm-project/vllm

Add ``/v1/chat/completions/batch`` endpoint for batched chat completions

合并时间: 2026-03-26 12:13

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/38011>

执行摘要

本 PR 为 vLLM 的 OpenAI 兼容 API 引入了新的 `/v1/chat/completions/batch` 端点，支持在单个请求中处理多个独立对话，旨在减少 HTTP 开销并提升应用程序效率。实现通过新增 `OpenAIServingChatBatch` 子类复用现有逻辑，确保响应按索引排序，并包含完整的测试和示例。该变更不影响现有单对话端点，是一个有意义的 API 扩展，适合需要批量处理场景的用户。

功能与动机

此变更是为了解决应用程序中需要同时处理多个对话的需求，例如从多个文档提取结构化数据。根据 PR body，批处理能“减少 HTTP 开销和简化结果处理，因为所有输出在一个响应中到达”。在 Issue 评论中，维护者 DarkLight1337 指出：“This is outside of OpenAI API spec so we will not support this for Chat Completions API to avoid bloating the existing functionality”，这促使作者改为创建一个独立的端点，以保持 API 规范兼容性并避免现有功能复杂化。

实现拆解

实现按模块拆解如下：

1. 路由与集成：

- 文件: `vllm/entrypoints/openai/chat_completion/api_router.py`
- 添加新端点 `/v1/chat/completions/batch`，使用 `create_batch_chat_completion` 函数处理请求。
- 在 `vllm/entrypoints/openai/generate/api_router.py` 中初始化 `openai_serving_chat_batch` 状态，确保批处理器与单处理器分离。

2. 协议模型：

- 文件: `vllm/entrypoints/openai/chat_completion/protocol.py`
- 定义 `BatchChatCompletionRequest` 类，扩展 `ChatCompletionRequest`，但 `messages` 字段改为 `list[list[...]]`。
- 添加 `model_validator` 强制约束：

```
python if data.get("use_beam_search"): raise ValueError("Batch chat completions do not support beam search.") if n is not None and n != 1: raise ValueError("Batch chat completions do not support `n > 1`.")
```

3. 批处理服务逻辑：

- 文件: `vllm/entrypoints/openai/chat_completion/batch_serving.py`
- `OpenAIServingChatBatch` 类继承 `OpenAIServingChat`, 核心方法 `render_batch_chat_request` 批量预处理对话, `create_batch_chat_completion` 使用 `asyncio.gather` 并发执行。
- 代码片段示例:

```
python async def create_batch_chat_completion(...): result = await handler.create_batch_chat_completion(request, raw_request) return JSONResponse(content=result.model_dump())
```

4. 测试与示例:

- 新增测试文件 `tests/entrypoints/openai/chat_completion/test_batched_chat_completions.py`, 覆盖基本批处理、JSON 架构和正则约束场景。
- 示例脚本 `examples/online_serving/batched_chat_completions.py` 提供端到端使用演示。

评论区精华

review 讨论中最有价值的交锋集中在设计正确性和效率上:

- 关于 `n` 参数验证: `gemini-code-assist[bot]` 强调: “For batched requests, `n` (the number of choices to generate for each prompt) should be restricted to 1. If a user provides `n > 1`, the response will contain choices with duplicate index values”。作者回应并修复, 添加验证器确保合规。
- 关于代码组织: `DarkLight1337` 提出: “Can we put these logics into a separate subclass of `OpenAIServingChat`? This class is long enough as it is”, 作者采纳建议, 创建独立子类提升模块化。
- 关于效率: `gemini-code-assist[bot]` 指出重复调用 `to_chat_completion_request` 可能低效, 但 `DarkLight1337` 认为“`No need`”优化, 展示了在性能与代码简洁性之间的权衡。

风险与影响

风险:

- 正确性风险: `echo=True` 逻辑在初始实现中有误, 但提交历史显示已修复; 验证器依赖 `Pydantic`, 如果客户端绕过可能引发未定义行为。
- 性能风险: 批处理并发可能增加内存和 CPU 使用, 尤其是在处理大量对话时; 建议监控生产环境负载。
- 兼容性风险: 新端点限制如不支持流式传输, 可能影响期望完整 `OpenAI` 功能的用户; 文档已说明, 但需确保用户知晓。

影响:

- 对用户: 提供更高效率的请求方式, 减少网络延迟, 适用于批量处理场景; 不影响现有单对话 API, 迁移成本低。
- 对系统: 扩展 API 层, 轻微增加维护复杂性, 但实现复用核心引擎, 整体架构稳定。
- 对团队: 引入新的代码模式, 如子类化处理批处理, 为未来 API 扩展提供参考。

关联脉络

从近期历史 PR 看，本 PR 与 vLLM 仓库中前端和 API 层的持续改进相关：

- PR #38029（工具解析器重构）和 #35182（输入模块重构）都涉及对 API 基础设施的优化，显示团队在标准化和扩展 API 功能上的努力。
- 本 PR 的独立端点设计呼应了维护者对保持 API 规范简洁性的重视，避免在现有端点上添加非标准功能，这与仓库中其他注重兼容性的变更（如 ROCm 性能优化 PR）一脉相承。
- 未来可能基于此批处理端点，进一步集成流式传输或工具支持，但当前限制文档清晰，为后续演进预留了空间。