

# PR #38016 完整报告

vllm-project/vllm

[gRPC] Add standard gRPC health checking (grpc.health.v1) for Kubernetes native probes

合并时间: 2026-04-23 05:31

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/38016>

## 执行摘要

- 一句话: 为 vLLM gRPC 服务器添加标准 gRPC 健康检查服务, 支持 Kubernetes 原生探针。
- 推荐动作: 建议技术管理者和工程师精读 `vllm/entrypoints/grpc_server.py` 中的健康服务集成部分, 关注关机处理和异常捕获设计; 同时查看测试文件以理解健康检查的各种场景。对于使用 gRPC 部署的用户, 此 PR 提供了重要的运维增强功能。

## 功能与动机

根据 PR body, 目的是“Add standard `grpc.health.v1.Health` service to the vLLM gRPC server for native Kubernetes gRPC health probes (`livenessProbe.grpc / readinessProbe.grpc`, GA since K8s 1.27)”, 以支持 Kubernetes 部署中的原生 gRPC 健康检查, 提升容器化环境的运维便利性和系统可靠性。

## 实现拆解

1. 导入和注册健康服务: 在 `vllm/entrypoints/grpc_server.py` 中, 新增导入 `grpc_health.v1.health_pb2_grpc` 和 `VllmHealthServicer`, 创建健康服务实例并注册到 gRPC 服务器, 同时更新反射服务列表以包含 `grpc.health.v1.Health`。
2. 优雅关机处理: 在服务器关机的 `finally` 块中, 调用 `health_servicer.set_not_serving()` 以设置健康状态为 `NOT_SERVING`, 并用 `try/except` 捕获异常避免影响关机流程。
3. 依赖版本升级: 在 `setup.py` 中, 将可选依赖 `grpc` 的 `smg-grpc-servicer[vllm]` 版本从 `>=0.5.0` 升级到 `>=0.5.2`, 以兼容新增的健康服务。
4. 测试配套: 新增 `tests/entrypoints/test_grpc_health.py` 文件, 包含 10 个单元测试, 覆盖健康检查的 `Check` 和 `Watch` 方法, 验证服务状态、错误处理和协议一致性。
5. 文档更新: 在 `docs/deployment/k8s.md` 中添加“Serving with gRPC”章节, 提供 gRPC 健康探针的 Kubernetes 配置示例和验证命令。

关键文件:

- `tests/entrypoints/test_grpc_health.py` (模块 健康检查; 类别 `test`; 类型 `test-coverage`; 符号 `async_llm, context, servicer, request_msg`): 新增的单元测试文件, 全面覆盖健康检查的 `Check` 和 `Watch` 方法, 验证服务状态、错误处理和协议一致性, 是确保功能正确性的关键。
- `vllm/entrypoints/grpc_server.py` (模块 入口点; 类别 `source`; 类型 `dependency-wiring`): gRPC 服务器入口文件, 关键变更包括导入健康服务、注册到服务器、更新反射列表和添

加优雅关机处理，直接影响服务器功能。

- setup.py (模块 依赖管理; 类别 infra; 类型 configuration) : 更新依赖版本文件, 将 smg-grpc-servicer[vllm] 从  $\geq 0.5.0$  升级到  $\geq 0.5.2$ , 确保兼容新增的健康服务功能。
- docs/deployment/k8s.md (模块 部署文档; 类别 docs; 类型 documentation) : 文档更新文件, 添加 gRPC 健康检查的配置说明, 帮助用户在实际部署中集成 Kubernetes 探针。

关键符号: VllmHealthServicer.Check, VllmHealthServicer.Watch,  
test\_check\_serving\_overall, test\_check\_serving\_vllm\_service,  
test\_check\_not\_serving\_engine\_errored, test\_check\_not\_serving\_shutting\_down

## 关键源码片段

### tests/entrypoints/test\_grpc\_health.py

新增的单元测试文件, 全面覆盖健康检查的 Check 和 Watch 方法, 验证服务状态、错误处理和协议一致性, 是确保功能正确性的关键。

```
# 测试健康检查 Check 方法的基本场景
@pytest.mark.asyncio
async def test_check_serving_overall(servicer, request_msg, context, async_llm):
    request_msg.service = "" # 设置服务名为空, 表示整体健康检查
    response = await servicer.Check(request_msg, context) # 调用健康检查方法
    assert response.status == SERVING # 验证返回状态为 SERVING
    async_llm.check_health.assert_awaited_once() # 确保委托给 AsyncLLM 的 check_health 方法

# 测试健康检查 Watch 方法的基本场景
@pytest.mark.asyncio
async def test_watch_yields_serving(servicer, request_msg, context, async_llm):
    request_msg.service = ""
    watch_iter = servicer.Watch(request_msg, context) # 调用 Watch 方法获取迭代器
    first = await anext(watch_iter.__aiter__()) # 获取第一个响应
    assert first.status == SERVING # 验证状态为 SERVING
```

### vllm/entrypoints/grpc\_server.py

gRPC 服务器入口文件, 关键变更包括导入健康服务、注册到服务器、更新反射列表和添加优雅关机处理, 直接影响服务器功能。

```
# 导入健康检查相关模块
try:
    import grpc
    from grpc_health.v1 import health_pb2_grpc # 导入 gRPC 健康检查协议
    from grpc_reflection.v1alpha import reflection
    from smg_grpc_proto import vllm_engine_pb2, vllm_engine_pb2_grpc
    from smg_grpc_servicer.vllm.health_servicer import VllmHealthServicer # 导入健康服务实现
    from smg_grpc_servicer.vllm.servicer import VllmEngineServicer
except ImportError as e:
    raise ImportError("gRPC模式需要smg-grpc-servicer, 请安装vllm[grpc]") from e

# 在服务器设置中添加健康服务
```

```
health_servicer = VllmHealthServicer(async_llm) # 创建健康服务实例, 委托给 AsyncLLM
health_pb2_grpc.add_HealthServicer_to_server(health_servicer, server) # 注册健康服务到 gRPC
服务器

# 更新反射服务列表以包含健康服务
service_names = (
    vllm_engine_pb2.DESRIPTOR.services_by_name["VllmEngine"].full_name,
    "grpc.health.v1.Health", # 添加健康服务到反射, 便于工具发现
    reflection.SERVICE_NAME,
)
reflection.enable_server_reflection(service_names, server)

# 在关机流程中设置健康状态为 NOT_SERVING
try:
    health_servicer.set_not_serving() # 通知健康服务服务器正在关闭
except Exception: # 宽泛异常捕获, 确保不影响关机流程
    logger.warning("Failed to set health status to NOT_SERVING", exc_info=True)
```

## 评论区精华

review 中主要讨论了两个问题:

- 日志记录: `gemini-code-assist[bot]` 指出健康检查异常时未记录日志, 可能影响调试; 作者 V2arK 回应已在 `VllmHealthServicer` 的 `except` 块中添加 `logger.exception` 调用。
- Watch 方法协议一致性: `gemini-code-assist[bot]` 提到当前 `Watch` 实现只发送一次状态, 不符合 gRPC 健康检查协议的连续流要求; 作者解释这为简化实现, 与 `SGLang` 保持一致, 且 Kubernetes 探针仅使用 `Check` 方法, 后续可改进。结论: 日志问题已修复, `Watch` 方法作为待优化项保留。
- 健康检查异常日志记录 (`correctness`): 问题已修复, 在 `VllmHealthServicer` 的 `except` 块中添加了日志记录。
- Watch 方法协议一致性 (`design`): 作为待优化项保留, 后续可改进以实现连续流。

## 风险与影响

- 风险: 技术风险包括:
  - 依赖风险: `smg-grpc-servicer` 版本升级至 0.5.2 可能引入不兼容或新 bug, 影响 gRPC 服务器稳定性。
  - 协议不完整: `Watch` 方法未实现连续状态更新, 可能不符合某些客户端对 gRPC 健康检查协议的期望。
  - 异常处理宽泛: 关机流程中捕获宽泛 `Exception` 可能隐藏底层错误, 但设计上为确保关机不中断。
  - 测试覆盖局限: 单元测试虽全面, 但未覆盖真实 Kubernetes 环境下的集成场景。
  - 影响: 对用户影响: Kubernetes 部署者可直接使用原生 gRPC 探针配置健康检查, 简化运维配置, 提升部署可靠性。对系统影响: 增加健康检查端点, 轻微增加请求处理开销, 但通过委托 `AsyncLLM.check_health()` 复用现有健康逻辑。对团队影响: 需维护外部依

赖 smg-grpc-servicer, 并关注后续协议改进。

- 风险标记: 依赖版本升级, Watch 方法不完整, 异常处理宽泛

## 关联脉络

- 暂无明显关联 PR