

# PR #42083 完整报告

vllm-project/vllm

[Feat] Add support for per GPU worker RDMA NIC selection

合并时间: 2026-05-29 03:45

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/42083>

## 执行摘要

- 一句话: 支持 per-GPU worker RDMA NIC 选择
- 推荐动作: 值得精读, 尤其 PCI BDF 规范化和 sysfs 遍历的实现可供其他 RDMA 相关特性参考。设计决策 (仅 NVML、去除 prefetch) 体现了简化优先的务实思路。

## 功能与动机

在 multi-NIC multi-GPU 节点上, RDMA 性能取决于 GPU-NIC PCIe 亲和性, 但某些拓扑 (如 Azure ND96isr\_v5 的 flat PCIe) 下网络库无法自动发现最优配对, 导致所有 worker 竞争同一 NIC。PR body 指出: 'In multi-NIC, multi-GPU nodes, RDMA performance depends on GPU-NIC PCIe affinity. When the topology is flat or not properly discoverable by networking libraries, all workers may default to the same NIC or incorrect NICs, leading to suboptimal RDMA performance.'

## 实现拆解

1. 环境变量注册(vllm/envs.py): 添加 VLLM\_GPU\_NIC\_PCIE\_MAPPING 和 VLLM\_NIC\_SELECTION\_VARS, 必须同时设置。
2. 平台抽象扩展(vllm/platforms/interface.py, vllm/platforms/cuda.py): 在 Platform 接口中添加 get\_all\_gpu\_pci\_bus\_ids() 抽象方法, NvmlCudaPlatform 通过 NVML 实现; 非 NVML 平台抛出 NotImplementedError。
3. 核心映射与设备发现模块(vllm/v1/executor/vllm\_net\_devices.py): 实现 normalize\_pci() 规范化 PCI BDF 地址, parse\_gpu\_nic\_mapping() 解析映射字符串, rdma\_name\_for\_nic\_pci() 通过 /sys/class/infiniband 反向查找 RDMA 设备名, set\_worker\_net\_device() 串联上述步骤为当前 worker 设置环境变量。
4. 执行器集成(multiproc\_executor.py, uniproc\_executor.py): 在 worker 进程初始化前调用 set\_worker\_net\_device(local\_rank, vllm\_config)。
5. 单元测试(tests/v1/executor/test\_vllm\_net\_devices.py): 覆盖 normalize\_pci 的多种格式和异常路径。

关键文件:

- vllm/v1/executor/vllm\_net\_devices.py (模块 执行器; 类别 source; 类型 core-logic; 符号 normalize\_pci, parse\_gpu\_nic\_mapping, rdma\_name\_for\_nic\_pci, parse\_nic\_selection\_vars) : 核心新增模块, 包含所有 GPU-NIC 映射逻辑, 是整个 PR 的

技术基础。

- tests/v1/executor/test\_vllm\_net\_devices.py (模块 测试; 类别 test; 类型 test-coverage ; 符号 test\_normalize\_pci\_full\_domain, test\_normalize\_pci\_short\_form, test\_normalize\_pci\_case\_insensitive, test\_normalize\_pci\_strips\_whitespace) : 新增完整的 normalize\_pci 单元测试, 确保解析逻辑正确性。
- vllm/platforms/cuda.py (模块 平台层; 类别 source; 类型 core-logic; 符号 get\_all\_gpu\_pci\_bus\_ids) : 实现 NVML 路径的 GPU PCI Bus ID 查询, 是 NIC 映射的数据来源。
- vllm/platforms/interface.py (模块 平台层; 类别 source; 类型 core-logic; 符号 get\_all\_gpu\_pci\_bus\_ids) : 在 Platform 基类中添加抽象方法, 确保跨平台兼容性。
- vllm/envs.py (模块 配置; 类别 source; 类型 configuration) : 注册新环境变量, 统一管理配置入口。
- vllm/v1/executor/multiproc\_executor.py (模块 执行器; 类别 source; 类型 dependency-wiring) : 在 MultiprocExecutor 的 worker\_main 中调用 set\_worker\_net\_device, 是生产环境的主要集成点。
- vllm/v1/executor/uniproc\_executor.py (模块 执行器; 类别 source; 类型 dependency-wiring) : 在 UniProcExecutor 中集成 set\_worker\_net\_device, 覆盖单进程 TP=1 场景。

关键符号: normalize\_pci, parse\_gpu\_nic\_mapping, rdma\_name\_for\_nic\_pci, parse\_nic\_selection\_vars, set\_worker\_net\_device, get\_all\_gpu\_pci\_bus\_ids

## 关键源码片段

### tests/v1/executor/test\_vllm\_net\_devices.py

新增完整的 normalize\_pci 单元测试, 确保解析逻辑正确性。

```
import pytest

from vllm.v1.executor.vllm_net_devices import normalize_pci

# 测试完整的 domain:bus:dev.fn 格式
@pytest.mark.parametrize(
    "addr, expected",
    [
        ("0000:3f:00.0", (0, 0x3F, 0, 0)),
        ("0001:00:00.0", (1, 0, 0, 0)),
        ("00000001:00:00.0", (1, 0, 0, 0)),
        ("0000:0a:1f.7", (0, 0x0A, 0x1F, 7)),
    ],
)

def test_normalize_pci_full_domain(addr, expected):
    assert normalize_pci(addr) == expected

# 测试简短的 bus:dev.fn 格式 (domain 默认为 0)
@pytest.mark.parametrize(
```

```

    "addr, expected",
    [
        ("01:00.0", (0, 1, 0, 0)),
        ("3f:00.0", (0, 0x3F, 0, 0)),
        ("ff:1f.7", (0, 0xFF, 0x1F, 7)),
    ],
)
def test_normalize_pci_short_form(addr, expected):
    assert normalize_pci(addr) == expected

# 验证大小写不敏感
def test_normalize_pci_case_insensitive():
    assert normalize_pci("0A:1F.7") == normalize_pci("0a:1f.7")

# 验证去除首尾空白和 0x 前缀
def test_normalize_pci_strips_whitespace():
    assert normalize_pci(" 0001:00:00.0 ") == (1, 0, 0, 0)

def test_normalize_pci_strips_0x_prefix():
    assert normalize_pci("0x0001:00:00.0") == (1, 0, 0, 0)

# 测试各种异常输入
def test_normalize_pci_missing_function_raises():
    with pytest.raises(ValueError, match="missing function suffix"):
        normalize_pci("0001:00:00")

def test_normalize_pci_invalid_function_char_raises():
    with pytest.raises(ValueError, match="invalid PCI function"):
        normalize_pci("0001:00:00.z")

def test_normalize_pci_too_many_segments_raises():
    with pytest.raises(ValueError, match="invalid PCI BDF"):
        normalize_pci("a:b:c:d.0")

def test_normalize_pci_bus_out_of_range_raises():
    with pytest.raises(ValueError, match="out of range"):
        normalize_pci("0000:1ff:00.0")

def test_normalize_pci_device_out_of_range_raises():
    with pytest.raises(ValueError, match="out of range"):
        normalize_pci("0000:00:20.0")

def test_normalize_pci_empty_string_raises():
    with pytest.raises(ValueError):
        normalize_pci("")

```

## 评论区精华

主要讨论聚焦于实现简化:

- tlrnchlsmith 建议放弃非 NVML 路径和 nvidia-smi 子进程探测，仅保留 pynvml 实现，从而去除 prefetch 机制，使 PR 大幅简化 (commit 1843c28)。
- tlrnchlsmith 提议使用 envs.py 中注册的变量而非直接 os.environ.get，已被采用。
- tlrnchlsmith 要求为 normalize\_pci 添加单元测试，rajkiranjoshi 在 commit 7385778 中实现。
- 关于 local\_rank 传递的疑问，rajkiranjoshi 解释了 kwargs.get('local\_rank', 0) 仅是 fallback，实际值在 spawn 时已正确设置，tlrnchlsmith 表示理解并同意简化调用。
- 简化：仅支持 NVML 路径 (design): rajkiranjoshi 采纳建议，移除了非 NVML 支持和 prefetch 机制，使实现更加简洁。
- 使用 envs.py 注册的变量 (style): rajkiranjoshi 在后续 commit 中改为引用 vllm.envs 中的变量。
- worker 间 local\_rank 传递 (correctness): tlrnchlsmith 表示理解并同意保留当前调用方式。
- 添加 normalize\_pci 单元测试 (testing): rajkiranjoshi 在 commit 7385778 中新增了完整的单元测试套件，覆盖正常和异常路径。

## 风险与影响

- 风险：新模块 vllm\_net\_devices.py 依赖 /sys/class/infiniband，仅 Linux 有效。NVML 路径被设为唯一实现，非 NVML 平台（如 Jetson）或 AMD GPU 将抛出 NotImplementedError，需要后续扩展。环境变量配置错误（如 PCI BDF 格式不匹配）会抛出 ValueError 并终止 worker 初始化。在多 DeepEP 场景中已验证性能提升，但其他 RDMA 库（如 UCCL）的实际测试尚未覆盖。
- 影响：用户侧需要手动设置两个环境变量，但 PR body 提供了详细示例，降低了使用门槛。系统侧正确配置后可显著降低 RDMA 传输延迟，提升跨节点通信效率。团队侧新增一个独立模块，后续维护成本低，且库无关设计易于扩展支持新通信库。
- 风险标记：仅 NVML 平台支持，依赖 /sys/class/infiniband（仅 Linux），环境变量必须同时设置，缺少非 NVML 平台的实现在未来可能影响扩展，深层性能测试仅覆盖 UCX 和 NVSHMEM

## 关联脉络

- 暂无明显关联 PR