

vllm-project/vllm 2026 年第 21 周周报 (05-18 至 05-24)

vllm-project/vllm

周期: 2026-05-18 至 2026-05-24

来源 PR: 224 · 重点 PR: 24 · 自动生成

原文链接: <http://prhub.com.cn/vllm-project/vllm/reports/2026-05-18-to-2026-05-24>

执行摘要

本周 2026-05-18 至 2026-05-24, 仓库共合并 224 个 PR (其中重点 24 个), 代码变更活跃。核心脉络是 DeepSeek V4 模型体系化迁移与优化、MoE 量化后端统一模块化、KV 连接器能力增强, 以及 Rust 前端正式并入主仓库。同时, 性能优化、跨平台适配 (CPU/XPU/ROCm) 和工具调用基础设施也取得显著进展。整体上, 本周是一次“向内重构”与“向外扩展”并重的密集交付周。

本周重点变化

1. DeepSeek V4 模型栈全面重构: WoosukKwon 和 zhyongye 主导的系列 PR (#43004、#43039、#43073、#43077、#43149) 将 DeepSeek V4 模型代码从 `model_executor/models` 迁移至 `models/deepseek_v4/`, 采用 `nvidia/` 和 `amd/` 硬件隔离目录结构, 为多后端模型架构奠定基础。同时, 稀疏 MLA 实现被提取并引入平台选择函数, ROCm 平台获得独立的 MTP 支持 (#43385), NVFP4 量化集成 (#42209) 和混合注意力缓存支持 (#42828) 进一步丰富了 DSV4 的生态。
2. MoE 量化后端模块化迁移: bedeks 等人的系列 PR (#42680、#42483) 将 W4A8 和 AWQ Marlin 等量化方案从直接调用特定内核迁移到统一的 `oracle` 后端选择框架。该模式已被推广至 NVFP4 (#40082) 和 CPU MXFP4 (#41922), 显著减少了重复代码, 并为未来新增量化后端提供了清晰的扩展路径。
3. Rust 前端正式纳入主仓库: BugenZhao 和 njhill 通过 #43283 和 #40848 将之前独立开发的 vLLMRust 前端 (`vllm-frontend-rs`) 全部源码迁入 `rust/` 目录, 并完成构建系统集成。默认关闭, 可通过环境变量启用。此里程碑标志着 vLLM 开始探索高性能 HTTP 前端的 Rust 原生实现。
4. KV 连接器 (Disaggregated Inference) 生态成熟: MooncakeStore 连接器新增 RPC 操作 Prometheus 指标 (#43392) 和混合注意力模型缓存支持 (#42828), 使其可用于 DeepSeek V4 等多注意力类型场景。offloading 连接器重构为类注册工厂 (#42529), 修复了请求挂起等关键 bug。示例隐藏状态连接器通过异步写入实现 1.45x 加速 (#37374)。
5. 性能优化多领域开花: Mamba 混合模型状态后处理融合内核 (#40172) 消除 GPU 气泡, 延迟降低 17%; Step3VL 多模态模型启用 Encoder CUDA Graph (#42224), 首字延迟降低 22%; FP8 线性层 padding 预计算 (#42651) 提升 TTFT 13.5%; Gemma4 视觉编码批量化 (#43169) 吞吐提升最高 3.8x; 多个 `zeros`→`empty` 替换 (#42988) 消除额外清零开销。

- 跨平台适配加速：CPU 后端实验性启用 Triton 与 Model Runner V2 (#43225)，新增 AMX 融合 GDN 算子 (#42707)；XPU 支持稀疏 attention 统一 custom op (#37888)、GPTQ int4 量化 (#37844)、FP8 block-scaled 量化 (#42952)；ROCm 修复多项 CI 稳定性并新增 DSV4 MTP 支持，同时修复 GDN 导入崩溃 (#43486)。
- 工具调用基础设施统一：sfeng33 等人通过 #43006、#43025、#43140 将多个 tool parser 中重复的 schema 类型转换逻辑提取到共享 `utils.py`，降低了维护成本。同时修复了流式推理内容丢失 (#42691) 和结构化标签受限生成 (#42452) 等边界 bug。

模块与主题趋势

从标签分布看，bugfix (88 个) 仍占主导，但 feature (36) 和 refactor (39) 也相当活跃。v1 标签持续高企 (62 个)，表明 v1 版本仍是主要开发线。性能优化 (37 个) 和内核开发 (25 个) 表明团队持续关注推理效率。

热点文件集中于 DeepSeek V4 相关 (`flashmla.py`、`rocm.py`、`model.py`) 和 MoE oracle 路径，验证了上述重点变化。

作者方面，haosdent、yewentao256、njhill 贡献最多，且 haosdent 大量集中在 CI 和 bugfix，显示基础设施和稳定性投入。

风险观察

- 测试覆盖缺口：29 个 PR 标注“缺少测试覆盖”，尤其在 DeepSeek V4 新代码、MoE 重构、Rust 前端中。建议对关键路径（如稀疏 MLA、MTP 推测解码、Mooncake HMA）补充单元和集成测试。
- 核心路径频繁变更：26 个 PR 影响调度器、模型加载、KVCache 等核心路径。建议建立更细粒度的性能基准，避免回归。
- Rust 前端遗留问题：CR 指出的异步 I/O 阻塞、递归栈溢出、安全漏洞等未在合并前修复，需尽快跟进。
- 依赖版本倒挂：为修复崩溃而降级 `nvidia-cutlass-dsl` 和 `triton_kernels`，需在后续版本中回归正常版本，避免技术债务。
- MooncakeStore HMA 未解决高优先级评论：合并时仍存在潜在 `ZeroDivisionError` 和段注册错误，虽可能不影响常见路径，但建议优先验证。

重点 PR 速览

PR	标题	要点	风险
#43385	[ROCm] [DSv4] [Perf] Support DeepSeek v4 MTP	在 ROCm 上实现 DSV4 多令牌预测，新增独立 AMD 模型文件	高并发性能退化、缺少测试

PR	标题	要点	风险
#43 392	[Mooncake] Add metrics for MooncakeStoreConnector	添加 RPC 耗时 / 键数 / 字节数等 Prometheus 指标	指标开销低
#42 680	[MoE] Migrate W4A8 CT to oracle kernel setup	将 W4A8 迁移到 oracle 框架, 新增早期 dtype 验证	缺少测试
#40 881	elastic_ep: stage/commit MoE quant method	弹性 EP 重新配置时 MoE 量化方法阶段 / 提交机制	缺少测试、NIXL 依赖
#37 374	[Perf] Optimize hidden state extraction	异步化 DtoH 拷贝与磁盘写入, 1.45x 加速	路径遍历、safetensors 版本
#43 149	[Refactor] Extract DeepSeek V4 sparse MLA	稀疏 MLA 从 backends 迁移至 model 目录, 硬件隔离	CUDA-ROCM 耦合
#43 405	[Rust] Extract UtilityCallId newtype	新类型保持 MessagePack 整数表示	负值假设
#41 234	[Multimodal] Simplify ViT CUDA graph interfaces	合并三个接口为一个, 降低维护负担	AssertionError
#41 126	[Attention] Mamba attention module refactor	提取基类, 拆分模型专有实现	缺少测试、conv1d 形状风险
#43 225	[CPU] Experimentally enable Triton and MRV2	CPU 平台实验性启用 Triton 与 V2 运行器	实验性、构建脆弱
#37 888	[XPU] Enable multiple key kernels for sparse attention	XPU 稀疏 attention 统一 custom op	依赖外部版本

PR	标题	要点	风险
#40 841	[Frontend] DP Supervisor	节点级多端口健康聚合, 支持 Kubernetes 探针	进程管理复杂度
#43 283	[Rust] Move code from vllm-frontend-rs	将 Rust 前端完整历史迁入主仓库	大量新代码缺乏测试
#43 105	[Core] Add native ModelExpress load format	新增 modelexpress 加载格式, 动态委托	可选依赖
#40 172	[Perf][Hybrid] Fused Triton kernel for Mamba state	融合内核消除 GPU 气泡, 延迟降低 17%	平台绑定
#40 848	[Frontend][RFC] Rust front-end integration	集成 Rust 前端作为实验性替代	Rust 工具链依赖
#40 082	Integrate flashinfer b12x MoE and FP4 GEMM	为 SM12x GPU 添加 FlashInfer 后端	依赖特定版本
#42 654	[Model] Openvlla support	新增 OpenVLA 机器人模型支持	timm 差异
#42 828	[KVConnector][DSV4] HMA support for Mooncake	MooncakeStore 支持混合注意力缓存	ZeroDivisionError 等
#43 004	[Model Refactoring] Migrate DSV4 to vllm/models/	DSV4 迁移至硬件隔离目录	模型加载路径变更
#42 529	Tier offload followup	重构 offloading 工厂模式, 修复关键 bug	接口不兼容
#42 483	Refactor AWQ Marlin MoE onto modular WNA16 oracle	AWQ Marlin 迁移至 oracle 框架	AWQ 权重路径变更

PR	标题	要点	风险
#41 922	[CPU] Add MXFP4 W4A16 MoE support	CPU 新增 MXFP4 量化 MoE 内核	CPU AMX 依赖
#42 224	[MM][CG] Enable encoder Cudagraph for Step3VL	为 Step3VL 启用 Encoder CUDA Graph , TTFT -22%	核心路径变更

后续建议

1. 优先补充测试：对新迁移的 DeepSeek V4 代码、MoE oracle 框架、Rust 前端进程管理、Mooncake HMA 等模块补充自动化测试，尤其是集成和回归测试。
2. 跟踪 Rust 前端遗留问题：尽快解决 review 中的性能和安全问题，确保在默认关闭期间修复完毕。
3. 建立 MoE 性能基准：随着多个后端整合到 oracle，建议设置端到端性能基准，监控不同量化方案和硬件组合的吞吐和延迟变化。
4. 关注依赖版本：降级的 cutlass-dsl 和 triton_kernels 应尽快回归最新兼容版本，并验证无回归。
5. 加强跨平台 CI 稳定性：CPU/XPU/ROCm 的 PR 数量增多，需确保 CI 覆盖率和稳定性，避免阻塞其它开发者。