

PR #42646 完整报告

vllm-project/vllm

[perf] Add gemma RMS AR fusion

合并时间: 2026-06-04 16:33

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/42646>

执行摘要

- 一句话: 集成 Flashinfer Gemma RMSNorm AR 融合, 优化 Qwen3.5 推理吞吐
- 推荐动作: 值得精读, 特别是模式匹配的注册技巧和 `extra_check` 的使用, 以及如何通过 `weight_bias` 抽象 Gemma 的特异性。展示了在 vLLM 编译 `passes` 中扩展新融合模式的标准流程。

功能与动机

利用 Flashinfer 新支持的 Gemma RMS AR 融合 (upstream PR#3322) 来优化 Gemma 系列模型的推理性能。通过将 `allreduce` 和 `RMSNorm` 合并为一个 `kernel`, 减少显存访问和 `kernel launch` 开销。PR body 给出了 Qwen3.5-397B-A17B-FP8 在 TP=4 上的详细 benchmark, 融合后吞吐量在所有并发度下均有提升, TTFT 在多数配置下降低。

实现拆解

1. 修改 `GemmaRMSNorm.forward_native` (`vllm/model_executor/layers/layernorm.py`): 简化实现, 移除手动 `dtype` 转换和上转型, 直接调用 `ir.ops.rms_norm` 或 `ir.ops.fused_add_rms_norm`, 权重偏移 +1.0 在输入 `dtype` 下计算。这使得 FX 子图模式统一, 便于后续模式匹配。
2. 新增 `dtype` 匹配检查 (`allreduce_rms_fusion.py`): 引入 `_norm_input_weight_dtype_match` 函数, 遍历匹配节点, 检查 `rms_norm` 或 `fused_add_rms_norm` 的输入与权重的 `dtype` 是否一致, 防止通用融合错误匹配 Gemma 的特殊情况。
3. 新增 Gemma 专属融合模式: `AllReduceGemmaRMSNormPattern` (无残差) 和 `AllReduceFusedAddGemmaRMSNormPattern` (有残差)。匹配 `allreduce -> ir.ops.rms_norm` 模式, 替换为 Flashinfer 的 `fused` 算子并传递 `weight_bias=1.0`。注册时使用 `extra_check=_norm_input_weight_dtype_match` 保护。
4. 扩展 `call_trtllm_fused_allreduce_norm` 签名: 增加 `weight_bias` 参数 (默认 0.0), 传递给 Flashinfer 底层融合 API。
5. 添加测试 (`test_fusion_all_reduce.py`): 新增 `TestAllReduceGemmaRMSNormModel`, 使用 4 个 `GemmaRMSNorm` 层和矩阵乘法构建模型, 验证融合后图中节点的 `weight_bias` 均为 1.0。在 ROCm 平台跳过。

关键文件:

- `vllm/compilation/passes/fusion/allreduce_rms_fusion.py` (模块 编译优化; 类别 `source`; 类型 `core-logic`; 符号 `_norm_input_weight_dtype_match`, `AllReduceGemmaRMSNormPattern`, `init`, `get_inputs`) : 核心实现: 新增 `dtype` 匹配检查函数、两个 Gemma 专属融合模式, 扩展 `call_trtllm_fused_allreduce_norm` 增加 `weight_bias` 参数。
- `tests/compile/passes/distributed/test_fusion_all_reduce.py` (模块 测试; 类别 `test`; 类型 `test-coverage`; 符号 `TestAllReduceGemmaRMSNormModel`, `init`, `forward`, `ops_in_model_before`) : 新增 `TestAllReduceGemmaRMSNormModel` 测试用例, 验证融合后 `weight_bias` 参数为 1.0。
- `vllm/model_executor/layers/layernorm.py` (模块 模型层; 类别 `source`; 类型 `data-contract`) : 重构 `GemmaRMSNorm.forward_native`, 简化精度处理并直接调用 `ir.ops`, 使得生成子图符合融合模式要求。

关键符号: `_norm_input_weight_dtype_match`, `AllReduceGemmaRMSNormPattern.init`, `AllReduceGemmaRMSNormPattern.register`, `AllReduceGemmaRMSNormPattern.pattern`, `AllReduceGemmaRMSNormPattern.replacement`, `AllReduceFusedAddGemmaRMSNormPattern`, `call_trtllm_fused_allreduce_norm`, `GemmaRMSNorm.forward_native`

关键源码片段

`vllm/compilation/passes/fusion/allreduce_rms_fusion.py`

核心实现: 新增 `dtype` 匹配检查函数、两个 Gemma 专属融合模式, 扩展 `call_trtllm_fused_allreduce_norm` 增加 `weight_bias` 参数。

以下代码展示了新增的 `dtype` 匹配检查函数和 `call_trtllm_fused_allreduce_norm` 中 `weight_bias` 参数的传递。

```
def _norm_input_weight_dtype_match(match: pm.Match) -> bool: """Prevent fusion when the norm input and weight dtypes differ (e.g. a Gemma fp32 weight.float()+1 gamma), covering rms_norm and fused_add_rms_norm.""" # 遍历匹配到的节点, 分别处理 rms_norm 和 fused_add_rms_norm
for node in match.nodes:
    if node.target == _IR_RMS_NORM_OP: # rms_norm 的参数顺序: x, weight, eps
        x, weight = node.args[0], node.args[1]
    elif node.target == _IR_FUSED_ADD_RMS_NORM_OP: # fused_add_rms_norm 的参数顺序: x, residual, weight, eps
        x, weight = node.args[0], node.args[2]
    else:
        continue
    if isinstance(x, fx.Node) and isinstance(weight, fx.Node):
        # 检查输入和权重的实际 dtype 是否一致
        return x.meta["val"].dtype == weight.meta["val"].dtype
    return True # 未找到相关节点则通过 #
call_trtllm_fused_allreduce_norm 新增 weight_bias 参数并传递给 Flashinfer 工作区
@torch.library.impl(triton_ops._flashinfer_trtllm_fused_allreduce_norm_op.default, "CUDA")
defcall_trtllm_fused_allreduce_norm(
    allreduce_in: torch.Tensor, residual: torch.Tensor, rms_gamma: torch.Tensor, rms_eps: float, world_size: int, launch_with_pdl: bool, fp32_acc: bool, max_token_num: int, pattern_code: int, norm_out: torch.Tensor | None = None, quant_out: torch.Tensor | None = None, scale_out: torch.Tensor | None = None, scale_factor: torch.Tensor | None = None,
```

```

weight_bias: float = 0.0, # 新增参数, Gemma 场景设为 1.0 ) -> None: # ... ( 参数
验证和内存检查 ) ... # 将 weight_bias 传递给 Flashinfer 的融合 API workspace.run(
    allreduce_in=allreduce_in, residual=residual,
rms_gamma=rms_gamma, rms_eps=rms_eps, world_size=world_size,
launch_with_pdl=launch_with_pdl, fp32_acc=fp32_acc,
max_token_num=max_token_num, pattern_code=pattern_code,
norm_out=norm_out, quant_out=quant_out, scale_out=scale_out,
scale_factor=scale_factor, layout_code=layout_code,
use_oneshot=use_oneshot, weight_bias=weight_bias, # 传递 weight_bias
trigger_completion_at_end=num_tokens > PDL_ADVANCE_LAUNCH_TOKENS, )

```

[tests/compile/passes/distributed/test_fusion_all_reduce.py](#)

新增 TestAllReduceGemmaRMSNormModel 测试用例, 验证融合后 weight_bias 参数为 1.0。

```

class TestAllReduceGemmaRMSNormModel(torch.nn.Module):
    """Test model for Gemma-style RMSNorm + AllReduce fusion."""
    def __init__(self, hidden_size=16, token_num=16, eps=1e-6, dtype: torch.dtype = torch.
float16):
        super().__init__()
        self.hidden_size = hidden_size
        self.eps = eps
        # 使用 4 个 GemmaRMSNorm 层, 权重初始化为小随机数
        self.norm = [GemmaRMSNorm(hidden_size, eps) for _ in range(4)]
        for n in self.norm:
            n.weight.data.normal_(mean=0.0, std=0.1)
        self.w = [torch.rand(hidden_size, hidden_size) for _ in range(3)]

    def forward(self, x):
        z = torch.relu(x)
        x = resid = tensor_model_parallel_all_reduce(z)
        y = self.norm[0](x)
        z2 = torch.mm(y, self.w[0])
        x2 = tensor_model_parallel_all_reduce(z2)
        y2, resid = self.norm[1](x2, resid)
        z3 = torch.mm(y2, self.w[1])
        x3 = tensor_model_parallel_all_reduce(z3)
        y3, resid = self.norm[2](x3, resid)
        z4 = torch.mm(y3, self.w[2])
        x4 = tensor_model_parallel_all_reduce(z4)
        y4, resid = self.norm[3](x4, resid)
        return y4

    def ops_in_model_before(self):
        return [torch.ops.vllm.all_reduce.default]

    def ops_in_model_after(self):
        # 融合后期望只出现 flashinfer 的 fused 算子
        return [torch.ops.vllm.flashinfer_trtllm_fused_allreduce_norm.default]

```

```

# 在 pytest parametrize 中添加该模型，并跳过 ROCm
pytest.param(
    TestAllReduceGemmaRMSNormModel, False, False,
    marks=pytest.mark.skipif(current_platform.is_rocm(), reason="Not supported on ROCm
platform"),
),

# 验证融合后每个 fused_op 节点的 weight_bias 均为 1.0
if test_model_cls is TestAllReduceGemmaRMSNormModel:
    fused_op = torch.ops.vllm.flashinfer_trtllm_fused_allreduce_norm.default
    fused_nodes = list(find_op_nodes(fused_op, backend.graph_post_pass))
    assert fused_nodes
    assert all(n.kwargs.get("weight_bias") == 1.0 for n in fused_nodes)

```

评论区精华

- gemini-code-assist 自动审查指出新 pattern 缺少 `extra_check=_norm_input_weight_dtype_match`，建议添加 dtype 匹配检查以避免运行时错误。
- ZJY0516 质疑 `weight_bias` 参数的必要性，认为可直接在 `forward` 中计算 `weight + 1.0` 以复用标准接口。作者回复 `weight_bias` 更灵活，但可简化。
- mgoin 关注精度变化，作者回应改用 Flashinfer `weight_bias` 方式以保持单 GPU/unfused 路径一致。
- ZJY0516 指出测试类中 `dtype` 参数未使用。
- `AllReduceGemmaRMSNormPattern` 缺少 dtype 匹配检查 (`correctness`): 已添加 `extra_check` 参数，确保 fusion 仅在输入和权重 dtype 一致时触发。
- `weight_bias` 参数必要性 (design): 保留了 `weight_bias` 参数，但仅用于 Gemma 模式，标准模式不受影响。
- 精度变化风险 (`correctness`): 作者改为使用 Flashinfer `weight_bias` 方式，单 GPU/unfused 路径保持不变，融合路径由 Flashinfer 内部处理精度。
- 测试类中未使用的 dtype 参数 (other): 该参数确实多余，但无实际影响，未修改。

风险与影响

- 风险:
 - 依赖 Flashinfer 版本: 需 $\geq 0.6.12$ (参见关联 PR#44036)，未升级时融合不会生效。
 - 精度风险: 修改了 `GemmaRMSNorm.forward_native` 的精度处理方式 (移除 fp32 upcast 手动控制)，虽然 `ir.ops` 内部可能自动处理，但边界行为可能有差异。
 - 仅特定 GPU 架构和 world size 生效: 融合 kernel one-shot 大小限制依赖于 `device capability` 和 `world size`，大 batch 可能 fallback。
 - 测试覆盖有限: 仅单模型结构，缺乏多并发端到端集成测试。
- 影响:

- 用户：使用 Gemma、Qwen3-next、Qwen3.5 的用户直接受益于性能提升（吞吐量 +2~8%，TTFT 下降）。
- 系统：编译 pass 增加两个模式匹配，编译开销稍增但运行时无影响。
- 团队：开辟了通过 `weight_bias` 抽象模型特异性的模式，后续可扩展至其他类似模型。
- 风险标记：依赖外部库版本，精度敏感，仅有限 GPU 配置生效

关联脉络

- PR #44036 Update flashinfer to 0.6.12: 此 PR 依赖 Flashinfer 0.6.12 版本中新增的 Gemma RMS AR 融合支持，需先合并更新。
- PR #38983 [PR discussion] Related context for `weight_bias`: ZJY0516 引用了该 PR 中的讨论上下文，涉及 Gemma RMSNorm 的实现策略。
- PR #3322 Flashinfer upstream PR for Gemma RMS AR fusion: 本 PR 集成的 Flashinfer 特性来自 upstream PR#3322。