

PR #26914 完整报告

sgl-project/sglang

[AMD] Remove BF16-to-FP32 elementwise cast from compressor GEMM on HIP

合并时间: 2026-06-04 15:58

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/26914>

执行摘要

- 一句话: 移除 AMD 上 compressor GEMM 的 BF16→FP32 类型转换
- 推荐动作: 该 PR 值得精读, 特别是对于在 AMD 平台上部署 DeepSeek-V4 模型的团队。核心设计决策 (在 HIP 路径绕过昂贵的类型转换, 同时在 Triton kernel 中添加显式类型处理) 展示了平台特定优化的典型方法。性能数据详实, aiter 库的使用也值得关注。

功能与动机

在 AMD MI355/HIP 平台上, `Compressor.forward()` 调用的 `linear_bf16_fp32()` 函数内部执行 `tgemm.mm(x, y, otype=x.dtype).float()`, 其中 `.float()` 会触发独立的 `bfloat16tofloat32_copy_kernel_cuda` 元素级 kernel (每次约 4.2 μ s)。对于 DeepSeek-V4 模型, 每 iteration 共有 90 次这样的转换 (C 型层 30x2 + B 型层 30x1), 累计约 384 μ s, 占 ITL 的 1.87%。B200 没有此问题, 因其 GEMM 后端原生支持 FP32 输出。

实现拆解

1. 在 `compressor.py` 中引入 aiter 条件导入和 GEMM 路径: 新增 `SGLANG_USE_AITER` 环境变量检查 (需同时满足 HIP 平台), 当启用时从 `aiter.tuned_gemm` 导入 `tgemm`, 并在 `compute_kv_score()` 方法中优先调用 `_tgemm.mm(x, self.wkv_gate.weight, otype=x.dtype)` 直接返回 BF16 结果, 而非调用 `linear_bf16_fp32()` (后者会额外执行 `.float()` 转换)。
2. 在 `fused_compress_triton.py` 中调整 Triton kernel 的类型处理: 放宽 `_check_common()` 函数中对 `kv_score_input` 和 `out` 的 `dtype` 断言, 允许接受 `torch.bfloat16` (原仅允许 `torch.float32`)。同时, 在所有加载 `kv_in_ptr` 的 Triton kernel (`_c4_decode_kernel`、`_c128_decode_kernel`、`_c4_prefill_compress_kernel`、`_c4_prefill_write_kernel`、`_c128_prefill_compress_kernel`) 中, 对加载的值显式调用 `.to(torch.float32)` 以确保后续 FP32 算术的正确性; 在最终存储输出时, 将结果显式转换回 `out_ptr.dtype.element_ty` 以匹配目标张量的实际类型。

关键文件:

- `python/sglang/srt/layers/attention/dsv4/compressor.py` (模块 压缩器; 类别 source; 类型 core-logic; 符号 `compute_kv_score`, `_use_aiter`, `_tgemm`): 引入 aiter 条件导入和 GEMM 调用路径, 是移除类型转换的核心决策点。

- `python/sglang/srt/layers/attention/dsv4/fused_compress_triton.py` (模块 压缩核; 类别 source; 类型 core-logic; 符号 `_check_common`, `_c4_decode_kernel`, `_c128_decode_kernel`, `_c4_prefill_compress_kernel`) : 调整 Triton kernel 以正确处理 BF16 输入输出, 确保类型安全。

关键符号: `compute_kv_score`, `_check_common`, `_c4_decode_kernel`, `_c128_decode_kernel`, `_c4_prefill_compress_kernel`, `_c4_prefill_write_kernel`, `_c128_prefill_compress_kernel`

关键源码片段

`python/sglang/srt/layers/attention/dsv4/compressor.py`

引入 `aiter` 条件导入和 GEMM 调用路径, 是移除类型转换的核心决策点。

```
# file: python/sglang/srt/layers/attention/dsv4/compressor.py
# 模块顶部添加条件导入
from sglang.srt.utils import add_prefix, get_bool_env_var

# 仅在 HIP 且启用 SGLANG_USE_AITER 时导入 aiter 的 tgemm
_use_aiter = get_bool_env_var("SGLANG_USE_AITER") and _is_hip
_tgemm = None
if _use_aiter:
    from aiter.tuned_gemm import tgemm
    _tgemm = tgemm

class Compressor(nn.Module):
    ...
    def compute_kv_score(self, x: torch.Tensor, forward_batch: ForwardBatch):
        # HIP/aiter 路径: 直接调用 tgemm.mm, 返回 BF16 输出
        # 避免 linear_bf16_fp32 中的 .float() 转换
        if _tgemm is not None:
            kv_score = _tgemm.mm(x, self.wkv_gate.weight, otype=x.dtype)
        else:
            # 非 HIP 或未启用 aiter 时使用原有逻辑
            kv_score = linear_bf16_fp32(x, self.wkv_gate.weight)

        # CUDA path: delegate to backend
        if dsa_use_prefill_cp(forward_batch):
            kv_score = cp_all_gather_rerange_output(
                kv_score, get_attention_cp_size(), forward_batch, torch.cuda.current_stream()
            )
        return kv_score
```

`python/sglang/srt/layers/attention/dsv4/fused_compress_triton.py`

调整 Triton kernel 以正确处理 BF16 输入输出, 确保类型安全。

```
# file: python/sglang/srt/layers/attention/dsv4/fused_compress_triton.py
# _check_common 函数: 放宽断言, 允许 bf16 输入
@staticmethod
```

```

def _check_common(kv_score_input, kv_score_buffer, coff, head_dim):
    assert kv_score_input.dim() == 2 and kv_score_input.dtype in (
        torch.float32, torch.bfloat16, # 新增 bfloat16 支持
    )
    ...

# 在 _c4_decode_kernel 中，加载 kv_in_ptr 时显式转换为 float32
@triton.jit
def _c4_decode_kernel(...):
    # 写入 buffer 时转为 float32
    val = tl.load(kv_in_ptr + in_base + ch_off + d_offs, mask=d_mask, other=0.0).to(tl.float32)
    ...
    # 处理输入行时转为 float32
    kv = tl.load(kv_in_ptr + in_base + kv_off + d_offs, mask=d_mask & valid, other=0.0).to(tl.float32)
    score = tl.load(kv_in_ptr + in_base + score_off + d_offs, mask=d_mask & valid, other=NEG_BIG).to(tl.float32)
    # 存储时转换为输出张量的 dtype
    tl.store(out_ptr + bid.to(tl.int64) * out_row_stride + d_offs,
            (weighted / running_sum).to(out_ptr.dtype.element_ty), mask=d_mask)

```

评论区精华

Review 中 [gemini-code-assist\[bot\]](#) 指出：仅放宽 `_check_common` 中的类型断言而不更新下游 Triton kernel 会导致 HIP 上的编译失败，因为 kernel 内部没有正确处理混合精度。具体提出了三点：

- kv_in_ptr 加载的值需要显式转换为 FP32；
- out_ptr 存储的值需要显式转换为输出张量的 dtype；
- buffer_ptr 相关操作需要 `.to(tl.float32)` 转换。

作者在最终提交中采纳了这些建议，在多个 kernel 中增加了 `to(tl.float32)` 和 `to(out_ptr.dtype.element_ty)`。

- Triton kernel 类型不匹配导致编译失败 (correctness)：作者在所有相关 kernel 中添加了 `.to(tl.float32)` 和 `.to(out_ptr.dtype.element_ty)` 转换。

风险与影响

- 风险：
 1. 回归风险：仅在 HIP 且启用了 SGLANG_USE_AITER 时生效，其他平台（NVidia B200、CUDA）路径不变，因此不影响已有功能。
 2. 类型安全：Triton kernel 中显式添加了 `.to(tl.float32)` 和 `.to(out_ptr.dtype.element_ty)`，避免了隐式类型转换的不确定性。
 3. 性能稳定性：未发现性能退化风险，实测显示吞吐量提升 2.1%，ITL 微小下降。
 4. 测试覆盖：PR 未新增单元测试，但提供了大规模 GSM8K 准确度测试（0.95 分，阈值 0.94）和端到端基准测试数据，验证了正确性和性能。- 影响：影响范围：仅限于 AMD

MI355 等 HIP 平台且使用 DeepSeek-V4 模型的用户。性能提升: output throughput 提升 +2.1%, median ITL 从 20.59 ms 降至 20.46 ms (-0.6%) , TPOT 改善 -2.6%。
兼容性: 完全向后兼容, 通过 SGLANG_USE_AITER 环境变量控制, 默认不启用。

- 风险标记: 缺少单元测试覆盖

关联脉络

- PR #26272 [AMD] Remove BF16-to-FP32 elementwise cast from compressor GEMM on HIP (amd/deepseek_v4 branch): 本 PR 是 26272 的 rebase 版本, 在 amd/deepseek_v4 分支上已被合并, 现提交到 main 分支。