

PR #38065 完整报告

vllm-project/vllm

[Perf] FP8 FlashInfer Attn for ViT

合并时间: 2026-04-27 13:44

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/38065>

执行摘要

- 一句话: ViT 编码器注意力 FP8 量化加速
- 推荐动作: 建议精读。特别是 FP8 缩放策略的设计 (动态 vs 静态)、自动保存机制以及 Triton kernel 的 stride-aware 实现。对于需要优化多模态流水线的工程师具有直接参考价值。

功能与动机

在视觉理解负载中, 大图像和短文本的 ViT 编码器注意力成为显著瓶颈, 尤其是文本模型已被量化 (如 nvfp4)。本 PR 为此场景提供可选的 FP8 量化加速, 属于 NVIDIA MLPerf Inference V6.0 的一部分。

实现拆解

1. 配置层扩展: 在 `vllm/config/multimodal.py` 新增 `mm_encoder_attn_dtype` 等字段, 替代原有的环境变量方式; 同步更新 `vllm/engine/arg_utils.py` 和 `vllm/config/model.py` 以支持 CLI 参数和模型配置解析。
2. FP8 量化 Kernel: 在 `vllm/kernels/triton/qkv_padded_fp8_quant.py` 实现一个 stride-aware 的 Triton 量化 kernel, 能将 Q/K/V 量化到 FP8 并将 `head_dim` 填充到 16 的倍数。该 kernel 通过三维步长直接读取非连续的视图, 避免了额外的数据搬移。
3. MMEncoderAttention 集成: 在 `mm_encoder_attention.py` 中添加 `_init_fp8_state`、`_load_fp8_scales_file`、`_record_amax_and_update_scales` 和 `_maybe_save_fp8_scales` 等方法。支持两种模式: 动态缩放 (通过循环 `amax` 历史记录更新 `scale`) 和静态缩放 (从 JSON 文件读取 per-layer 的 Q/K/V `scale`)。
`process_weights_after_loading` 方法负责初始化 FP8 量化器并加载缩放文件。
4. cu_seqlens 兼容: FP8 量化后 Q/K/V 变为独立的连续张量, 不再需要 V 的特殊偏移。在 `vllm/model_executor/models/vision.py` 中提供 `get_fp8_padded_hidden_size` 辅助函数, 并在 `qwen3_vl.py` 的视觉编码器初始化时计算并传递给 `cu_seqlens` 生成逻辑。
5. 平台门控与动态保存: 在 `vllm/utils/flashinfer.py` 添加 `is_flashinfer_cudnn_fp8_prefill_attn_supported`, 检查 cuDNN 版本和 GPU 计算能力 (需 $\geq 9.17.1$ 且 Hopper+)。自动保存功能通过模块级全局变量, 在某个层的 `amax` 缓冲回绕时一次性将所有层的最新 `scale` 写入 JSON。

6. 测试与基准：新增三个测试文件覆盖量化 kernel 的正确性、缩放逻辑（动态 / 静态 / 保存）和端到端 FP8 注意力路径。新增基准脚本 `benchmark_vit_fp8_attn.py`，支持 CUDA Graph 和 PyTorch profiler 两种计时方式。

关键文件：

- `vllm/model_executor/layers/attention/mm_encoder_attention.py`（模块 编码器注意力；类别 source；类型 core-logic；符号 `_load_fp8_scales_file`, `_maybe_save_fp8_scales`, `_init_fp8_state`, `process_weights_after_loading`）：核心集成文件，添加了 FP8 缩放加载、动态 / 静态缩放、自动保存和 FP8 注意力前向路径。
- `vllm/kernels/triton/qkv_padded_fp8_quant.py`（模块 FP8 量化；类别 source；类型 core-logic；符号 `_quantize_pad_fp8_kernel`, `_get_fp8_pad_quant_config`, `quantize_fp8_pad_head_dim_triton`, `quantize_fp8_maybe_pad_head_dim`）：新增的 Triton kernel 文件，实现 stride-aware 的 FP8 量化并填充 `head_dim` 到 cuDNN 需要的倍数。
- `vllm/utils/flashinfer.py`（模块 FlashInfer 工具；类别 source；类型 dependency-wiring；符号 `is_flashinfer_cudnn_fp8_prefill_attn_supported`）：添加了 FP8 注意力支持的门控函数，检查 cuDNN 版本和 GPU 能力。
- `vllm/config/multimodal.py`（模块 多模态配置；类别 source；类型 configuration）：定义了新的配置字段 `mm_encoder_attn_dtype`、`mm_encoder_fp8_scale_path` 等，是功能入口配置。
- `tests/kernels/core/test_vit_fp8_attn.py`（模块 FP8 注意力测试；类别 test；类型 test-coverage；符号 `_has_flashinfer_cudnn`, `_fp8_attention`, `_build_cu_seqlens_and_meta`, `test_fp8_attn_output_shape`）：端到端 FP8 注意力测试，验证量化 +cuDNN+ 反填充的全路径正确性。
- `benchmarks/kernels/benchmark_vit_fp8_attn.py`（模块 性能测试；类别 source；类型 benchmark；符号 `_setup_fp8_attention`, `_build_meta`, `run_benchmark`, `bf16_fn`）：FP8 与 BF16 注意力的性能对比基准，包含 CUDA Graph 和 Profiler 模式。

关键符号：`_load_fp8_scales_file`, `_maybe_save_fp8_scales`, `_init_fp8_state`, `process_weights_after_loading`, `_record_amax_and_update_scales`, `quantize_fp8_maybe_pad_head_dim`, `quantize_fp8_pad_head_dim_triton`, `is_flashinfer_cudnn_fp8_prefill_attn_supported`

关键源码片段

`vllm/model_executor/layers/attention/mm_encoder_attention.py`

核心集成文件，添加了 FP8 缩放加载、动态 / 静态缩放、自动保存和 FP8 注意力前向路径。

```
# 模块级常量：amax 历史缓冲区长度
_FP8_AMAX_HISTORY_LEN = 16
```

```
# 模块级状态：单个保存路径和尺度引用缓存
_fp8_scale_save_path: str | None = None
_fp8_saved_scale_refs: dict[str, tuple[torch.Tensor, ...]] = {}
```

```

@functools.cache
def _load_fp8_scales_file(path: str | None) -> dict[str, dict[str, float]]:
    """加载 per-layer FP8 Q/K/V 缩放系数 (JSON 格式) , 结果会被缓存。
    支持简单格式 ({"layer": {"q": ..., "k": ..., "v": ...}}) 和嵌套格式 ({"layers": ...}) 。
    """
    if path is None:
        return {}
    with open(path) as f:
        data = json.load(f)
    # 兼容嵌套的 "layers" 键
    if "layers" in data and isinstance(data["layers"], dict):
        data = data["layers"]
    scales = {}
    for layer_name, layer_scales in data.items():
        q = layer_scales.get("q", layer_scales.get("q_scale"))
        k = layer_scales.get("k", layer_scales.get("k_scale"))
        v = layer_scales.get("v", layer_scales.get("v_scale"))
        if all(s is not None for s in (q, k, v)):
            q_f, k_f, v_f = float(q), float(k), float(v)
            if q_f <= 0 or k_f <= 0 or v_f <= 0:
                raise ValueError(
                    f"FP8 scales must be positive, got q={q_f}, "
                    f"k={k_f}, v={v_f} for layer '{layer_name}'")
            scales[layer_name] = {"q": q_f, "k": k_f, "v": v_f}
    logger.info_once("Loaded FP8 scales from %s (%d layers)", path, len(scales))
    return scales

def _maybe_save_fp8_scales(layer_name, q_scale, k_scale, v_scale, buffer_wrapped):
    """在首个 amax 缓冲区回绕时, 将累计的缩放系数写入 JSON 文件。
    避免每个 step 都做 GPU→CPU 同步, 只在回绕时一次性 .item()。
    """
    global _fp8_scale_save_path
    if _fp8_scale_save_path is None:
        return
    # 保存张量引用, 不立即同步
    _fp8_saved_scale_refs[layer_name] = (q_scale, k_scale, v_scale)
    if not buffer_wrapped:
        return
    # 缓冲区回绕: 释放保存路径, 写入文件
    path, margin = _fp8_scale_save_path, _fp8_scale_save_margin
    scales = {
        name: {
            "q": q.item() * margin,
            "k": k.item() * margin,
            "v": v.item() * margin
        }
    }
    for name, (q, k, v) in _fp8_saved_scale_refs.items()
}

```

```
_fp8_scale_save_path = None # 一击即止
with open(path, "w") as f:
    json.dump(scales, f, indent=2)
logger.info("Saved FP8 scales to %s (margin=%.3f)", path, margin)
```

vllm/kernels/triton/qkv_padded_fp8_quant.py

新增的 Triton kernel 文件，实现 stride-aware 的 FP8 量化并填充 head_dim 到 cuDNN 需要的倍数。

```
def quantize_fp8_maybe_pad_head_dim(
    tensor: torch.Tensor,
    scale: torch.Tensor,
    fp8_quant: QuantFP8,
    skip_scale: bool = False,
) -> torch.Tensor:
    """将输入张量量化到 FP8，如有必要将 head_dim 填充为 16 的倍数。

    优先使用 QuantFP8 CustomOp（当 head_dim 已对齐时），否则回退到 Triton kernel。
    支持 3D (S, H, D) 和 4D (B, S, H, D) 输入。
    """
    if tensor.dim() == 4:
        B, S, H, D = tensor.shape
        flat = tensor.reshape(B * S, H, D)
    else:
        flat = tensor
    padded = round_up(D, 16)
    if not skip_scale and D % 16 == 0:
        # head_dim 对齐：使用原生 QuantFP8（通常更快）
        return fp8_quant.quantize(flat, scale=scale, dtype=current_platform.fp8_dtype())
    else:
        # 未对齐：使用 Triton kernel 同时完成量化和填充
        result = quantize_fp8_pad_head_dim_triton(
            flat, scale, skip_scale=skip_scale
        )
    if tensor.dim() == 4:
        return result.view(B, S, H, -1)
    return result
```

评论区精华

Reviewer Isotr0py 提出将环境变量改为 `MultiModalConfig` 字段，以避免环境变量膨胀，作者采纳。Isotr0py 还建议简化缩放配置：如果提供了 `scales` 文件则使用静态缩放，否则默认动态缩放，作者重构。ProExpertProg 建议将 Triton kernel 独立到 `vllm/kernels/triton/` 目录下，以保持代码组织清晰。Isotr0py 建议将 `MMEncoderAttention` 加入模型加载器的 `process_weights_after_loading` 自动扫描，作者完成。Isotr0py 指出缺少 GPU 设备能力检查，作者添加了 `has_device_capability(90)` 门控。总体设计权衡聚焦于简单性与灵活性。

- 配置设计：环境变量 vs 配置字段 (design)：作者将所有环境变量替换为 `MultiModalConfig` 字段。

- 缩放策略简化 (design): 作者同意并简化: 只保留 `attn_dtype` 和 `scale_path`, 动态缩放作为缺省。
- 平台门控: GPU FP8 能力检查 (correctness): 作者添加了 `has_device_capability(90)` 门控, 确保仅 Hopper+ 启用。
- `process_weights_after_loading` 自动扫描 (design): 作者将 `MMEncoderAttention` 添加到自动扫描元组。

风险与影响

- 风险: 该功能依赖 FlashInfer cuDNN 后端和 `cuDNN ≥ 9.17.1`, 且仅支持 Hopper+ GPU (计算能力 9.0+), 限制了可部署性。动态缩放会引入额外的量化 kernel 开销, 在小序列长度下可能导致性能回退 (如 4096 序列长度 0.91x 加速比)。自动保存功能使用模块级全局状态, 在多 worker 环境中可能存在竞态, 目前仅适用于单进程 GPU。FP8 注意力仅在 FlashInfer 后端生效, 其他后端不受影响。新配置项增加了用户学习成本, 需文档辅助。
- 影响: 用户: 需要满足硬件和软件依赖才能获得加速; 对于大图像视觉任务 (如 ChartQA) 可降低 ViT 注意力延迟。系统: 新增编译级依赖 (Triton 和 cuDNN), 但功能默认不启用, 不影响现有部署。团队: 需维护新的 Triton kernel 和缩放逻辑, 但代码已遵循 vLLM 现有模式 (如 `process_weights_after_loading`) 。
- 风险标记: 硬件限制 (Hopper+ 和 `cuDNN ≥ 9.17.1`), 动态缩放小序列性能回退, 模块级全局状态竞态风险

关联脉络

- 暂无明显关联 PR