

PR #17883 完整报告

sgl-project/sglang

[NPU] Support GGUF quantization for Ascend NPU (dense + MoE)

合并时间: 2026-04-25 22:16

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/17883>

执行摘要

- 一句话: 在 Ascend NPU 上支持 GGUF 量化模型推理
- 推荐动作: 该 PR 功能完整, 设计上采用预去量化策略简化 NPU 推理路径, 是一个合理的权衡。建议关注以下设计决策: 预去量化的时机选择 (加载时而非推理时) 和 MoE 权重在 FusedMoE 中的延迟材料化。对于 NPU 平台的用户, 这是一个关键的功能补充。建议团队在未来的量化重构中对 NPU 变体进行统一抽象, 以降低维护成本。

功能与动机

在 Ascend NPU 上运行 GGUF 格式的量化 LLM 模型 (如 Q4_K_M、Q8_0、Q5_K_M 等), 实现原生 NPU 支持并优化性能。PR body 说明: "Enable GGUF quantized models to run on Ascend NPU hardware. GGUF is a popular format for quantized LLM models, and this PR adds native NPU support with optimized performance."

实现拆解

实现拆解如下:

1. 新增 NPU 专用量化方法: 在 `python/sglang/srt/layers/quantization/gguf.py` 中新增三个 Ascend 特化类: `GGUFLinearAscendMethod` (线性层)、`GGUFMoEAscendMethod` (MoE 层) 和 `GGUFEmbeddingAscendMethod` (Embedding 层)。它们继承自 `LinearMethodBase`、`FusedMoEMethodBase` 等基类, 并在 `process_weights_after_loading` 中执行预去量化 (pre-dequantization), 将量化参数还原为全精度 `half` 或 `bfloat16` 张量, 加载到 NPU 后推理时不再需要量化计算。新增 `ggml_dequantize_ascend` 函数封装 `gguf.dequantize` 的 CPU 参考实现。
2. MoE 权重加载适配: 在 `python/sglang/srt/layers/moe/fused_moe_triton/layer.py` 中新增 `_load_gguf_weight` 方法, 在 `_weight_loader_impl` 中优先处理 GGUF 类型的参数: 将权重数据暂存在 `param.data_container` 和 `param.expert_data_map` 中, 并处理 TP 分片。新增 `materialize_gguf_weights` 方法, 在所有权重加载完成后组装 `w13` (gate+up 合并) 和 `w2` 专家权重。
3. 权重加载器扩展: 在 `python/sglang/srt/model_loader/weight_utils.py` 中重写了 `gguf_quant_weights_iterator`, 使其能够识别 MoE 专家权重的打包格式 (`blk.{layer_id}.ffn_{gate/up/down}_exps.weight`), 将其拆分为每个专家的独立权重并生成相应的 `qweight` 和 `qweight_type` 张量名称。

4. 模型代码微调：修改 `python/sglang/srt/models/qwen3_moe.py`，在 GGUF 场景下将 `norm_topk_prob` 设置为 `False`，避免与预去量化权重不兼容；修改 `python/sglang/srt/models/qwen2_moe.py` 和 `python/sglang/srt/layers/linear.py` 以传递 `quant_config` 参数。
5. 测试与文档：新增两个集成测试文件 `test_npu_gguf.py`（密集模型）和 `test_npu_gguf_moe.py`（MoE 模型），基于 gsm8k 数据集验证精度，准确率分别达到 0.80 和 0.85 以上。文档更新了 Ascend NPU 量化指南和通用量化页面。

关键文件：

- `python/sglang/srt/layers/quantization/gguf.py`（模块 量化层；类别 source；类型 core-logic；符号 `ggml_dequantize_ascend`, `GGUFLinearAscendMethod`, `GGUFMoEAscendMethod`, `GGUFEmbeddingAscendMethod`）：核心实现文件：添加了 NPU 专用的 GGUF 量化方法（线性层、MoE 层、Embedding 层）和去量化函数，以及在 `GGUFConfig` 中根据设备选择不同方法的调度逻辑。
- `python/sglang/srt/layers/moe/fused_moe_triton/layer.py`（模块 MoE 层；类别 source；类型 core-logic；符号 `_load_gguf_weight`, `materialize_gguf_weights`）：MoE 关键路径：新增 GGUF 权重加载钩子 `_load_gguf_weight` 和材料化方法 `materialize_gguf_weights`，实现 FusedMoE 对 GGUF 权重的支持。
- `python/sglang/srt/model_loader/weight_utils.py`（模块 权重加载；类别 source；类型 data-contract；符号 `gguf_quant_weights_iterator`）：权重加载器：重写 `gguf_quant_weights_iterator` 以支持 MoE 专家权重的打包格式，将其解析为每个专家的独立权重。
- `test/registered/ascend/basic_function/quant/test_npu_gguf_moe.py`（模块 NPU 测试；类别 test；类型 test-coverage；符号 `TestAscendGGUFMoE`, `setUpClass`, `test_a_gsm8k`）：MoE 模型的集成测试，验证 Qwen3-30B-A3B GGUF 模型在 NPU 上的推理准确率不低于 0.85。
- `test/registered/ascend/basic_function/quant/test_npu_gguf.py`（模块 NPU 测试；类别 test；类型 test-coverage；符号 `TestAscendGGUF`, `setUpClass`, `test_a_gsm8k`）：密集模型集成测试，验证 Qwen3-14B GGUF 模型在 NPU 上的推理准确率不低于 0.80。

关键符号：`ggml_dequantize_ascend`, `GGUFLinearAscendMethod.create_weights`, `GGUFLinearAscendMethod.process_weights_after_loading`, `GGUFMoEAscendMethod.apply`, `GGUFEmbeddingAscendMethod.apply`, `GGUFConfig.get_quant_method`, `_load_gguf_weight`, `materialize_gguf_weights`, `gguf_quant_weights_iterator`

关键源码片段

`python/sglang/srt/layers/quantization/gguf.py`

核心实现文件：添加了 NPU 专用的 GGUF 量化方法（线性层、MoE 层、Embedding 层）和去量化函数，以及在 `GGUFConfig` 中根据设备选择不同方法的调度逻辑。

```
# 文件：python/sglang/srt/layers/quantization/gguf.py
# 关键变更 1: 设备感知的量化方法分派
```

```

def get_quant_method(self, layer, prefix):
    """根据层类型和设备返回合适的量化方法"""
    from sglang.srt.layers.moe.fused_moe_triton import FusedMoE
    from sglang.srt.layers.vocab_parallel_embedding import VocabParallelEmbedding

    if isinstance(layer, LinearBase):
        if is_layer_skipped_gguf(prefix, self.modules_to_not_convert):
            return UnquantizedLinearMethod()
        # NPU 使用预去量化方案, CUDA/MUSA 使用原方案
        if _is_npu:
            return GGUFLinearAscendMethod(self)
        return GGUFLinearMethod(self)
    elif isinstance(layer, VocabParallelEmbedding):
        if _is_npu:
            return GGUFEmbeddingAscendMethod(self)
        return GGUFEmbeddingMethod(self)
    elif isinstance(layer, FusedMoE):
        if _is_npu:
            return GGUFMoEAscendMethod(self)
        return GGUFMoEMethod(self)
    return None

```

关键变更 2: NPU 专用的去量化函数 (在 CPU 上执行参考实现)

```

def ggml_dequantize_ascend(qweight, qweight_type, rows, cols, dtype):
    """

```

```

    使用 gguf 库的参考实现在 CPU 上进行去量化,
    支持所有 GGML 格式, 保证正确性。
    结果转移到 NPU 设备后供推理使用。
    """

```

```

    # 将量化权重移到 CPU 并通过 numpy 调用 gguf 库
    qweight_cpu = qweight.cpu().numpy()
    dequant_np = gguf_dequantize(qweight_cpu, qweight_type)
    # 转换回 torch 张量并移到 NPU
    result = torch.from_numpy(dequant_np).to(dtype=dtype, device=qweight.device)
    return result.reshape(rows, cols)

```

python/sglang/srt/layers/moe/fused_moe_triton/layer.py

MoE 关键路径: 新增 GGUF 权重加载钩子 `_load_gguf_weight` 和材料化方法 `materialize_gguf_weights`, 实现 FusedMoE 对 GGUF 权重的支持。

文件 : python/sglang/srt/layers/moe/fused_moe_triton/layer.py
 # 关键方法 : 在权重加载时劫持 GGUF 参数, 暂存到 `data_container` 中

```

def _load_gguf_weight(self, param, loaded_weight, shard_id, expert_id, tp_rank):
    """尝试将 loaded_weight 作为 GGUF 权重处理。
    如果处理成功返回 True, 否则返回 False (由普通加载逻辑继续) 。
    """
    is_gguf_weight = getattr(param, 'is_gguf_weight', False)
    is_gguf_weight_type = getattr(param, 'is_gguf_weight_type', False)

```

```

if is_gguf_weight_type:
    # 存储该专家的量化类型（如 Q4_K_M），用于后续去量化
    param.weight_type = loaded_weight.item()
    return True

if is_gguf_weight:
    output_dim = getattr(param, 'output_dim', None)
    if self.moe_tp_size > 1:
        # 如果模型并行，对第一个维度进行分片
        if shard_id in ['w1', 'w3', 'w2'] and output_dim == 0:
            shard_size = loaded_weight.size(0) // self.moe_tp_size
            start_idx = tp_rank * shard_size
            loaded_weight = loaded_weight.narrow(0, start_idx, shard_size).clone()

        # 暂存到 expert_data_map 和 data_container 中
        if not hasattr(param, 'expert_data_map'):
            param.expert_data_map = {}
        key = (expert_id, shard_id)
        param.expert_data_map[key] = loaded_weight
        param.data_container.append(loaded_weight)
        return True

return False

```

python/sglang/srt/model_loader/weight_utils.py

权重加载器：重写 gguf_quant_weights_iterator 以支持 MoE 专家权重的打包格式，将其解析为每个专家的独立权重。

文件：python/sglang/srt/model_loader/weight_utils.py
关键变更：gguf_quant_weights_iterator 新增 MoE 专家权重解析

```

def gguf_quant_weights_iterator(gguf_file, gguf_to_hf_name_map):
    reader = gguf.GGUFReader(gguf_file)
    # MoE 权重在 GGUF 中以打包格式存储：blk.{layer_id}.ffn_gate_exps.weight
    MOE_WEIGHT_PATTERNS = {
        'ffn_gate_exps': 'gate_proj',
        'ffn_up_exps': 'up_proj',
        'ffn_down_exps': 'down_proj',
    }

    # 第一遍：输出权重类型 (qweight_type)
    for tensor in reader.tensors:
        # 判断是否为 MoE 专家权重
        if any(pattern in tensor.name for pattern in MOE_WEIGHT_PATTERNS):
            import re
            match = re.match(r'blk\.(?P<layer_id>)\.(?P<pattern>)\.weight', tensor.name)
            if match:
                layer_id = int(match.group(1))
                pattern = match.group(2)

```

```

hf_name = MOE_WEIGHT_PATTERNS.get(pattern)
if hf_name and tensor.tensor_type.name != 'F32':
    # 对打包中的每个专家输出一个 qweight_type
    num_experts = tensor.data.shape[0]
    for expert_id in range(num_experts):
        hf_type_name = f'model.layers.{layer_id}.mlp.experts.{expert_id}.{hf_name}.
        qweight_type'
        yield hf_type_name, torch.tensor(tensor.tensor_type)
elif tensor.name in gguf_to_hf_name_map:
    # 常规权重保持不变
    name = gguf_to_hf_name_map[tensor.name]
    if tensor.tensor_type.name != 'F32':
        yield name.replace('weight', 'qweight_type'), torch.tensor(tensor.tensor_type)

# 第二遍：输出实际权重，MoE 权重拆分为每个专家
for tensor in reader.tensors:
    # 类似逻辑，根据 pattern 拆分为独立权重
    # ... (省略重复代码)

```

评论区精华

1. 代码重复问题: gemini-code-assist[bot] 指出 `materialize_gguf_weights` 和 `process_weights_after_loading` 中 w13/w2 的处理逻辑高度相似，建议提取辅助函数。该建议未在后续提交中显著重构，属于设计债务被记录。
2. 内联导入风格: ping1jing2 要求将 `materialize_gguf_weights` 方法内的 `from torch.nn.parameter import UninitializedParameter` 移到文件顶部。该问题已在最终代码中修复（文件顶部已有该导入）。
3. GPU 兼容性疑虑: TamirBaydasov 担心修改 `gguf_quant_weights_iterator` 会破坏 GPU 上的 GGUF 模型加载。TheKonka 回复称 GPU+GGUF+MoE 原有错误，修改后行为不变；GPU+GGUF+DENSE 正常。团队确认无回归。
4. 测试覆盖建议: ping1jing2 在 `qwen2_moe.py` 的 review 中建议添加单元测试。本 PR 仅包含集成测试，未添加单元测试。
 - 代码重复: `materialize_gguf_weights` 中 w13 和 w2 处理逻辑高度相似 (design): 建议未被采纳，在最终实现中仍存在重复代码。
 - 内联 import 风格问题 (style): 已在最终代码中修正：文件顶部添加了该导入语句。
 - GPU GGUF 加载兼容性风险 (correctness): 作者确认兼容性无变化，GPU+GGUF+MoE 的原有错误需另行修复。
 - 缺少针对 NPU GGUF 路径的单元测试 (testing): PR 仅包含端到端集成测试，未添加单元测试，建议未完全采纳。

风险与影响

- 风险:
 1. GPU GGUF MoE 兼容性: 作者指出 GPU+GGUF+MoE 组合在修改前后均报错，但未解释具体原因或提供修复方案。这可能导致用户在 GPU 上使用 GGUF MoE 模型时仍然失

败，尽管这可能是一个预先存在的问题。

2. 预去量化内存开销：预去量化策略在加载时将所有权重转为全精度，对于大模型会增加显存占用。用户需关注 `--mem-fraction-static` 参数的设置。
3. 模型兼容性有限：仅使用 Qwen3-14B 和 Qwen3-30B-A3B 进行了验证，其他 GGUF 模型（尤其是不同架构或量化类型）可能遇到未预期的错误。
4. 缺少单元测试：未针对 `ggml_dequantize_ascend` 等核心函数编写单元测试，仅依赖端到端集成测试，调试成本较高。
5. 代码维护债务：review 指出的代码重复问题未被修复，未来重构量化文件夹时可能增加技术债务。

• 影响：

- 用户影响：Ascend NPU 用户现在可以直接加载 GGUF 格式的量化模型，无需额外转换步骤。支持 Q4_K_M、Q8_0 等常见量化类型，覆盖密集和 MoE 模型架构。
- 系统影响：不影响 CUDA 和 MUSA 后端，仅当 `device='npu'` 时启用 Ascend 专用路径。预去量化过程发生在模型加载阶段，推理延迟不受影响。
- 团队影响：新增的 NPU 特异性代码增加了量化模块的维护分支。未来重构（如 Issue 中提及的量化文件夹重构）需考虑这些 Ascend 变体。
- 风险标记：新功能仅验证两种模型，GPU GGUF MoE 原有错误未修复，预去量化增加加载内存开销，缺少核心函数单元测试，代码重复债务被记录

关联脉络

- PR #23731 Fix Qwen3 MoE double-reduce when DP attention + EP + reduce_scatterv (#23729): 该 PR 同样修改了 `qwen3_moe.py`，修复了与本 PR 设置的 `norm_topk_prob=False` 可能相关的 MoE double-reduce 问题。两者对同一模型的 MoE 行为有叠加影响。