

PR #24662 完整报告

sgl-project/sglang

Breakable Cuda Graph Support for bs > 1

合并时间: 2026-05-11 13:28

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/24662>

执行摘要

- 一句话: 支持 $bs > 1$ 的可中断 CUDA 图执行
- 推荐动作: 值得精读。该 PR 体现了深刻的技术洞察: 通过重新划定 CUDA 图捕获边界, 使图与 batch size 解耦, 是使图化预填充支持多请求的关键设计。代码改动简洁 (仅 1 文件 +84/-57), 但思路值得借鉴。建议关注后续改进 layer_model 解析的多模型兼容性和测试覆盖。

功能与动机

原有的 BreakableCudaGraphRunner 捕获整个模型 forward (包括 logits_processor 和 pooler), 这些核的形状依赖于 batch size, 导致图只能在 $bs=1$ 时有效。为了支持多请求预填充时依然可以在 CUDA 图上执行, 需要将捕获边界缩小至不依赖 batch size 的内部 transformer 层。

实现拆解

实现分为四步:

1. 解析 layer_model (init): 在初始化时通过 `language_model = getattr(model_runner.model, 'language_model', model_runner.model)` 和后续的 `language_model.model` 解析出内部 transformer 栈模块 (与 PiecewiseCudaGraphRunner 的 patch_model 边界一致), 并赋值给 `self.layer_model`。此方法依赖模型属性名称, 存在鲁棒性问题 (见 review)。
2. 修改 _run_forward: 将 `model_runner.model.forward` 调用替换为 `self.layer_model.forward`, 只执行 transformer 栈。添加 `@torch.no_grad` 装饰器以匹配外部 ForCausalLM.forward 的无梯度模式。
3. 修改 _build_capture_forward_batch: 构建占位的 forward_batch 时使用 $bs=1$ (作为 attention/mamba 分段元数据形状的占位), 实际的 $bs > 1$ 元数据由 replay_prepare 在 replay 时注入。
4. 清理与简化: 移除原先为 bs 参数静态分配的 `static_seq_lens` 等一批形状为 (max_bs,) 的张量, 以及验证计数器 `replay/can_run_reject` 和相关日志。

无其他文件变更, 无测试文件配套 (但原机制已有覆盖)。

关键文件:

- python/sglang/srt/model_executor/breakable_cuda_graph_runner.py (模块 执行引擎; 类别 source; 类型 core-logic; 符号 init, _run_forward, _build_capture_forward_batch, _init_buffers) : 核心变更文件, 实现将 CUDA 图捕获边界缩小到 layer_model, 支持 bs>1。

关键符号: init, _run_forward, _build_capture_forward_batch, _init_buffers

关键源码片段

python/sglang/srt/model_executor/breakable_cuda_graph_runner.py

核心变更文件, 实现将 CUDA 图捕获边界缩小到 layer_model, 支持 bs>1。

```
# __init__ 中新增的 layer_model 解析 (替换原有的静态 bs 张量分配)
def __init__(self, model_runner: ModelRunner):
    # ... 前面代码不变

    # 解析内部 transformer 栈模块, 边界与 PCG patch_model 一致
    language_model = getattr(
        model_runner.model, "language_model", model_runner.model
    )
    self.layer_model = (
        language_model.model
        if hasattr(language_model, "model")
        and hasattr(language_model.model, "layers")
        else language_model
    )
    # 注意: 此方式依赖模型属性命名, 不够鲁棒, review 已指出

    # Memory pool (不变)
    if get_global_graph_memory_pool() is None:
        set_global_graph_memory_pool(self.device_module.graph_pool_handle())
    set_graph_pool_id(get_global_graph_memory_pool())

    # Warmup / capture (不变)
    self._warmup()
    self.device_module.synchronize()
    self.model_runner.tp_group.barrier()
    self._capture_all()

    self.raw_num_tokens = 0

# _run_forward 使用 layer_model.forward, 并添加 @torch.no_grad
@torch.no_grad() # 新增, 匹配外部的 torch.no_grad 装饰
def _run_forward(self, forward_batch, num_tokens):
    """只执行内部 transformer 栈前向, 避免 bs 相关形状固化。
    """
    forward_batch.dp_local_start_pos = forward_batch.dp_local_num_tokens = None
    set_dp_buffer_len(None, num_tokens, forward_batch.dp_padding_mode.is_max_len())
    set_is_extend_in_batch(False)
```

```
with set_forward_context(
    forward_batch,
    self.attention_layers,
    self.quant_config,
    self.moe_layers,
    self.moe_fusions,
):
    output = self.layer_model.forward( # 原为 model_runner.model.forward
        forward_batch.input_ids,
        forward_batch.positions,
        forward_batch,
    )
return output
```

评论区精华

Reviewer merrymercy 指出 `layer_model` 的解析方式依赖于字符串名称匹配（`model_runner.model.language_model.model`），非常脆弱，建议至少应在匹配失败时发出警告。PR 作者未公开回复此评论，未添加警告，但 PR 最终被合并。该设计决策属于技术债务，需后续改进。

- `layer_model` 解析依赖字符串名称匹配不鲁棒 (design): PR 作者未公开回应，未添加 warning，但 PR 最终被合并，该设计被接受为初期实现。

风险与影响

- 风险：
 1. 鲁棒性风险：解析 `layer_model` 的路径（`model_runner.model.language_model.model`）依赖模型内部属性名称，不同模型结构（如无 `language_model` 属性或层结构不同）可能导致解析失败，进而引发运行时错误。review 已指出此问题。
 2. 回归风险：改变了图捕获边界，影响所有使用 `BreakableCudaGraphRunner` 的预填充路径；虽然移除静态张量和计数器降低了复杂性，但可能遗漏某些模型（如多模态）的特殊处理。
 3. 性能权衡：`logits_processor/pooler` 在 `replay` 后 `eager` 执行，与原先全部在图中执行相比，会产生少许额外启动开销，但在多请求场景下被图重用收益覆盖。
 4. 测试覆盖：未发现直接对应的测试文件；虽然有 CI 覆盖，但新增的多请求预填充图执行场景可能未被充分测试。- 影响：对用户：使用 `breakable CUDA graph` 的预填充场景（如长文档预处理）现在可以同时处理多个请求（`bs>1`），显著提升吞吐，特别是对于具有多个并行预填充请求的高负载场景。对系统：需要维护 `layer_model` 的引用，初始化多一次属性查找，但性能收益大于成本。对团队：引入了依赖模型属性的脆弱设计，需在后续重构中通过更稳定的接口（如定义 `model.get_layer_stack()`）来解决。
- 风险标记：依赖模型结构，缺少测试覆盖，核心路径变更

关联脉络

- 暂无明显关联 PR