

PR #24757 完整报告

sgl-project/sglang

Optimize ngram decode id computation

合并时间: 2026-06-02 17:37

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/24757>

执行摘要

- 一句话: 新增 ngram decode 专用 kernel, decode 计算加速 54%~86%
- 推荐动作: 值得精读。该 PR 展示了如何通过分析 decode 语义简化 kernel 实现以获得显著性能提升, 设计模式 (thread-per-output、边界检查) 对其他 token-wise 操作有借鉴意义。建议关注 ComputeNGramIdsDecodeKernel 的循环展开策略和 grid 大小选择。

功能与动机

在 decode 阶段, 每个请求只产生一个 token, 原始 kernel 需要传入 tokens、exclusive_req_len_sums 等参数并执行 cumsum, 这些开销可以消除。新 kernel 直接从 ne_token_table 按 row_index 和 column_start 读取最后一个 token, 省去额外计算和参数传递, 显著降低延迟。

实现拆解

1. 新增 CUDA kernel: 在 python/sglang/jit_kernel/csrc/ngram_embedding.cuh 中添加 ComputeNGramIdsDecodeKernel, 采用 thread-per-output 模式, 每个线程计算一个 ngram ID; 简化参数列表, 不再需要 tokens 和 exclusive_req_len_sums, 直接通过 row_indices 和 column_starts 从 ne_token_table 读取对应 token。
2. 注册 JIT 接口: 在 python/sglang/jit_kernel/ngram_embedding.py 中定义 compute_n_gram_ids_decode 函数, 使用 @debug_kernel_api 装饰, 并在 _jit_ngram_embedding_module 中注册对应的 CUDA wrapper。
3. 修改 forward 路由: 在 python/sglang/srt/layers/n_gram_embedding.py 的 NgramEmbedding.forward 中根据 forward_batch.forward_mode.is_decode() 分支: decode 模式调用新 kernel, extend 模式仍走原来的 compute_n_gram_ids 路径 (含 cumsum)。
4. 添加正确性测试: 新增 python/sglang/jit_kernel/tests/test_ngram_embedding.py, 参数化 batch size 为 [1,2,17,128,1024], 对比新 kernel 与原 kernel 输出是否严格一致 (torch.testing.assert_close)。
5. 添加性能基准: 新增 python/sglang/jit_kernel/benchmark/bench_ngram_compute_decode.py, 使用 Triton perf_report 对比 general 与 decode 路径在不同 batch size 下的延迟, 结果已计入 PR body。

关键文件:

- python/sglang/jit_kernel/csrc/ngram_embedding.cuh (模块 CUDA 内核; 类别 other; 类型 core-logic; 符号 ComputeNGramIdsDecodeKernel, kDecodeBlockSize, kMaxComputeNGramIdsDecodeBlocks, NgramEmbeddingKernel::compute_n_gram_ids_decode) : 核心 CUDA kernel 实现, 新增 ComputeNGramIdsDecodeKernel 和对应的静态方法 compute_n_gram_ids_decode, 是性能优化的核心。
- python/sglang/jit_kernel/ngram_embedding.py (模块 Python 封装; 类别 source; 类型 core-logic; 符号 compute_n_gram_ids_decode, _jit_ngram_embedding_module) : Python 端定义 compute_n_gram_ids_decode 函数并注册到 JIT 模块, 是连接 kernel 和上层调用的桥梁。
- python/sglang/srt/layers/n_gram_embedding.py (模块 模型层; 类别 source; 类型 dependency-wiring) : 修改 NgramEmbedding.forward 方法: decode 模式路由到新 kernel, extend 保持原样, 是性能提升生效的关键入口。
- python/sglang/jit_kernel/tests/test_ngram_embedding.py (模块 测试; 类别 test; 类型 test-coverage; 符号 _make_ngram_params, test_compute_n_gram_ids_decode_matches_general) : 参数化测试新增 kernel 的正确性, 对比 general 与 decode 路径的输出, 确保五个 batch size 下完全一致。
- python/sglang/jit_kernel/benchmark/bench_ngram_compute_decode.py (模块 基准测试; 类别 source; 类型 core-logic; 符号 _make_ngram_params, benchmark, fn) : 性能基准脚本, 使用 Triton perf_report 对比新旧路径延迟, 量化加速效果, 结果用于 PR 验证。

关键符号: compute_n_gram_ids_decode, NgramEmbeddingKernel::compute_n_gram_ids_decode, ComputeNGramIdsDecodeKernel, NgramEmbedding.forward

关键源码片段

python/sglang/jit_kernel/csrc/ngram_embedding.cuh

核心 CUDA kernel 实现, 新增 `ComputeNGramIdsDecodeKernel` 和对应的静态方法 `compute_n_gram_ids_decode`, 是性能优化的核心。

```
// ngram_embedding.cuh (partial) — 新增 decode 专用 kernel 及其 wrapper
```

```
constexpr int kDecodeBlockSize = 256;
constexpr int kMaxComputeNGramIdsDecodeBlocks = 65535;

// 线程 per output 的 ngram ID 计算 kernel
__global__ void ComputeNGramIdsDecodeKernel(
    int batch_size,
    int ne_n,
    int ne_k,
    const int* __restrict__ ne_weights,
    const int* __restrict__ ne_mods,
    const int* __restrict__ exclusive_ne_embedder_size_sums,
    const int* __restrict__ ne_token_table,
    int max_context_len,
    const int64_t* __restrict__ row_indices,
    const int* __restrict__ column_starts,
```

```

int* __restrict__ n_gram_ids) {
const int num_configs = (ne_n - 1) * ne_k;
const int total_outputs = batch_size * num_configs;
for (int output_idx = blockIdx.x * blockDim.x + threadIdx.x;
     output_idx < total_outputs;
     output_idx += blockDim.x * gridDim.x) {
const int req_id = output_idx / num_configs;
const int config_idx = output_idx - req_id * num_configs;
const int k_idx = config_idx % ne_k;
const int n_idx = config_idx / ne_k;
const int weight_offset = n_idx * ne_k * ne_n + k_idx * ne_n;
const int ne_mod = ne_mods[n_idx * ne_k + k_idx];
uint64_t n_gram_id = 0;
const int64_t req_token_table_offset = row_indices[req_id] * static_cast<int64_t>(max_
context_len);
const int64_t current_token_table_offset = req_token_table_offset + column_starts[req_id];
for (int j = 0; j < n_idx + 2; j++) {
if (current_token_table_offset - j < req_token_table_offset) { break; }
const int token = ne_token_table[current_token_table_offset - j];
if (token < 0) { break; } // 忽略被标记为 ignore 的 token
const uint64_t term = static_cast<uint64_t>(token) * static_cast<uint64_t>(ne_
weights[weight_offset + j]);
n_gram_id += term % ne_mod;
}
n_gram_id %= ne_mod;
n_gram_id += exclusive_ne_embedder_size_sums[n_idx * ne_k + k_idx];
n_gram_ids[output_idx] = static_cast<int>(n_gram_id);
}
}

```

// 在 NgramEmbeddingKernel 中的静态 wrapper, 供 TVM JIT 调用

```

static void compute_n_gram_ids_decode(
const int64_t ne_n,
const int64_t ne_k,
const tvm::ffi::TensorView ne_weights,
const tvm::ffi::TensorView ne_mods,
const tvm::ffi::TensorView exclusive_ne_embedder_size_sums,
const tvm::ffi::TensorView ne_token_table,
const tvm::ffi::TensorView row_indices,
const tvm::ffi::TensorView column_starts,
const tvm::ffi::TensorView n_gram_ids) {
// ... TensorMatcher 验证维度、类型、设备, 然后 launch kernel
auto stream = DeviceRuntime::Get(0)->GetStream();
int grid_size = min(kMaxComputeNGramIdsDecodeBlocks,
static_cast<int>((total_outputs + kDecodeBlockSize - 1) / kDecodeBlockSize));
ComputeNGramIdsDecodeKernel<<<grid_size, kDecodeBlockSize, 0, stream>>>(
batch_size, ne_n, ne_k, ...);
}

```

python/sglang/jit_kernel/ngram_embedding.py

Python 端定义 `compute_n_gram_ids_decode` 函数并注册到 JIT 模块，是连接 kernel 和上层调用的桥梁。

```
# ngram_embedding.py (partial) — 新增 decode 专用 Python wrapper

@cache_once
def _jit_ngram_embedding_module() -> Module:
    return load_jit(
        "ngram_embedding",
        cuda_files=["ngram_embedding.cuh"],
        cuda_wrappers=[
            ("compute_n_gram_ids", "&NgramEmbeddingKernel::compute_n_gram_ids"),
            # 新增 decode kernel 注册
            (
                "compute_n_gram_ids_decode",
                "&NgramEmbeddingKernel::compute_n_gram_ids_decode",
            ),
            ("update_token_table", "&NgramEmbeddingKernel::update_token_table"),
        ],
    )

@debug_kernel_api
def compute_n_gram_ids_decode(
    ne_n: int,
    ne_k: int,
    ne_weights: torch.Tensor,
    ne_mods: torch.Tensor,
    exclusive_ne_embedder_size_sums: torch.Tensor,
    ne_token_table: torch.Tensor,
    row_indices: torch.Tensor,
    column_starts: torch.Tensor,
    n_gram_ids: torch.Tensor,
) -> None:
    """
    计算 ngram ID 的 decode 版本（每个请求一个 token）
    参数比 general 版本少：不需要 tokens 和 exclusive_req_len_sums
    """
    module = _jit_ngram_embedding_module()
    module.compute_n_gram_ids_decode(
        ne_n,
        ne_k,
        ne_weights,
        ne_mods,
        exclusive_ne_embedder_size_sums,
        ne_token_table,
        row_indices,
        column_starts,
```

```
        n_gram_ids,
    )
```

python/sclang/srt/layers/n_gram_embedding.py

修改 `NgramEmbedding.forward` 方法: decode 模式路由到新 kernel, extend 保持原样, 是性能提升生效的关键入口。

n_gram_embedding.py (partial) — forward 方法中的 decode 路由逻辑

```
def forward(self, input_ids: torch.Tensor, forward_batch: ForwardBatch):
    if (
        forward_batch.forward_mode.is_extend()
        or forward_batch.forward_mode.is_decode()
    ):
        ngram_embedding_info = forward_batch.ngram_embedding_info
        if forward_batch.forward_mode.is_decode():
            # decode 阶段: 使用简化 kernel, 无需 cumsum
            compute_n_gram_ids_decode(
                ne_n=self.over_embedding_n,
                ne_k=self.over_embedding_k,
                ne_weights=self.oe_weights,
                ne_mods=self.oe_mods,
                exclusive_ne_embedder_size_sums=self.exclusive_oe_embedder_size_sums,
                ne_token_table=ngram_embedding_info.token_table,
                row_indices=forward_batch.req_pool_indices,
                column_starts=ngram_embedding_info.column_starts,
                n_gram_ids=self.oe_n_gram_ids[: len(input_ids)],
            )
        else:
            # extend/prefill 阶段: 仍走原路径, 包含 cumsum
            torch.cumsum(
                ngram_embedding_info.req_lens,
                dim=0,
                dtype=torch.int32,
                out=self.exclusive_req_len_sums[1 : 1 + forward_batch.batch_size],
            )
            compute_n_gram_ids(
                ne_n=self.over_embedding_n,
                ne_k=self.over_embedding_k,
                ne_weights=self.oe_weights,
                ne_mods=self.oe_mods,
                tokens=input_ids.to(torch.int32),
                exclusive_ne_embedder_size_sums=self.exclusive_oe_embedder_size_sums,
                exclusive_req_len_sums=self.exclusive_req_len_sums[: forward_batch.batch_size + 1],
                ,
                ne_token_table=ngram_embedding_info.token_table,
                row_indices=forward_batch.req_pool_indices,
                column_starts=ngram_embedding_info.column_starts,
                n_gram_ids=self.oe_n_gram_ids[: len(input_ids)],
```

)
... 后续 ngram 嵌入计算

评论区精华

本 PR 的 review 评论均为机器人 /rerun-failed-ci 指令，未产生实质性技术讨论。两位 reviewer (kpham-sgl, yuan-luo) 均 approve，表明社区对变更质量和测试覆盖表示认可。

- 暂无高价值评论线程

风险与影响

- 风险：正确性风险低：新增 kernel 有参数化测试覆盖 5 种 batch size，并与原 kernel 逐元素比较。但未覆盖极端情形（如 column_starts=0、token_table 中有大量负 token），kernel 内置 break 条件处理负 token，理论上安全。性能风险：decode 路径替换后若存在未知 bug 可能导致静默错误，但测试验证通过。兼容性风险：无，decode 路径独立，extend 不受影响。
- 影响：对使用 NgramEmbedding 的模型推理任务（推测为某些 LLM 的 ngram 特征）的 decode 阶段延迟降低 50% 以上。对 extend/prefill 无影响。系统整体吞吐量可能因 decode 加速而略有提升。
- 风险标记：核心路径变更，新增 CUDA kernel，测试覆盖高

关联脉络

- 暂无明显关联 PR