

PR #40654 完整报告

vllm-project/vllm

[Core] Avoid seq_lens_cpu GPU->CPU sync

合并时间: 2026-04-24 08:35

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/40654>

执行摘要

- 一句话: 避免 GPU→CPU 同步, 引入 `seq_lens_cpu_upper_bound`
- 推荐动作: 此 PR 值得精读, 尤其是从事 speculative decoding 或 attention 后端开发的工程师。设计决策: 用 CPU 计算的上界替代 GPU 张量访问, 是一种典型的异步优化模式。建议关注 `eagle.py` 中减法操作的实现, 确认其无同步。

功能与动机

主要动机是消除 prefill 阶段的 CPU 同步 (特别是 DS3.2 场景)。作者在 issue 评论中说明: 'eliminating the current prefill cpu sync for DS3.2', 并希望将 V1 特有的 `optimistic_seq_lens_cpu` 泛化到 `CommonAttentionMetadata` 中, 同时废弃 `seq_lens_cpu` 属性以避免隐式同步。

实现拆解

实现主要分为以下步骤:

1. 在 `InputBatch` 中新增 `seq_lens_cpu_upper_bound` 字段 (`vllm/v1/worker/gpu/model_runner.py`) - 在 `prepare_inputs` 方法中, 通过 `np.add(num_computed_tokens_np, num_scheduled_tokens)` 计算上界, 并转换为 Tensor 存入 `InputBatch`。- 该上界是纯 CPU 计算的, 不依赖 GPU 数据, 避免了同步。
2. 在 `CommonAttentionMetadata` 中新增 `seq_lens_cpu_upper_bound` 字段 (涉及 `vllm/v1/attention/backends/utils.py`、`vllm/v1/worker/gpu/attn_utils.py` 等) - 在构建元数据时传入该字段, 替代原先使用 `seq_lens_cpu` 计算 `max_seq_len` 等操作。
3. 在各 attention 后端中改用上界 (`cross_attention.py`, `mha_attention.py`, `flex_attention.py`, `mha_indexer.py`) - 在交叉注意力中, 利用 `seq_lens_cpu_upper_bound - query_lens_cpu` 计算 `num_computed_tokens_cpu`, 避免从 GPU 读取。- 在 MLA 注意力中, 用上界计算 `context_lens_cpu`。
4. 在 `ubatch_utils.py` 中调整切片逻辑 - 读取 `_seq_lens_cpu` 时先检查是否为 None, 避免触发属性同步。- 改用 `seq_lens_cpu_upper_bound` 计算 `max_seq_len`。
5. 在 `postprocess` 和 `postprocess_pool` 中乐观推进 CPU 镜像 (`model_runner.py`) - 每次后处理后, 将 `num_computed_tokens_np` 加上 `num_scheduled_tokens`, 确保上界在下一轮仍然有效。

6. 在 speculative decoding 组件中适配 (`eagle.py`, `dflash.py`, `llm_base_proposer.py`) - 在 `eagle.py` 中, 对 `seq_lens_cpu_upper_bound` 进行减去 rejected tokens 的修正, 保持上界正确。

关键文件:

- `vllm/v1/worker/ubatch_utils.py` (模块 `UBatch`; 类别 `source`; 类型 `core-logic`; 符号 `_make_metadata_with_slice`): 核心切片逻辑: 改用 `_seq_lens_cpu` 和 `seq_lens_cpu_upper_bound`, 避免同步并计算 `max_seq_len`。
- `vllm/v1/worker/gpu/model_runner.py` (模块 `ModelRunner`; 类别 `source`; 类型 `data-contract`; 符号 `prepare_inputs`, `postprocess`, `postprocess_pool`): 主入口: 计算并传递 `seq_lens_cpu_upper_bound`, 并在 `postprocess` 中推进 CPU 镜像。
- `vllm/v1/spec_decode/eagle.py` (模块 `SpecDecode`; 类别 `source`; 类型 `core-logic`; 符号 `prepare_inputs`): `spec decode` 关键路径: 修正上界以反映 rejected tokens, 减少过估计。
- `vllm/model_executor/layers/attention/cross_attention.py` (模块 `交叉注意力`; 类别 `source`; 类型 `data-contract`; 符号 `CrossAttentionBuilder.build`): 交叉注意力: 用上界计算 `num_computed_tokens`, 避免读取 GPU 属性。
- `vllm/model_executor/layers/attention/mla_attention.py` (模块 `MLA 注意力`; 类别 `source`; 类型 `data-contract`; 符号 `MLABackend.build`): `MLA 注意力`: 类似地使用上界计算 `context lens`。

关键符号: `prepare_inputs`, `postprocess`, `postprocess_pool`, `prepare_attn`, `_make_metadata_with_slice`, `build (cross attention)`, `build (MLA attention)`

关键源码片段

`vllm/v1/worker/ubatch_utils.py`

核心切片逻辑: 改用 `_seq_lens_cpu` 和 `seq_lens_cpu_upper_bound`, 避免同步并计算 `max_seq_len`。

```
# vllm/v1/worker/ubatch_utils.py
# 改动核心: 直接读取内部字段 _seq_lens_cpu 和新增的 seq_lens_cpu_upper_bound,
# 避免触发 .seq_lens_cpu 属性引起的 GPU→CPU 同步。
seq_lens = attn_metadata.seq_lens[request_slice]
# 直接访问私有字段以绕过同步; 若为 None 则退化为 None
seq_lens_cpu = (
    attn_metadata._seq_lens_cpu[request_slice]
    if attn_metadata._seq_lens_cpu is not None
    else None
)
seq_lens_cpu_upper_bound = (
    attn_metadata.seq_lens_cpu_upper_bound[request_slice]
    if attn_metadata.seq_lens_cpu_upper_bound is not None
    else None
)
num_computed_tokens_cpu = (
```

```

    attn_metadata._num_computed_tokens_cpu[request_slice]
    if attn_metadata._num_computed_tokens_cpu is not None
    else None
)
# 对 split 的边界处理: 克隆并修改上界 (如果存在)
if splits_last_request:
    # ...
    if seq_lens_cpu is not None:
        seq_lens_cpu = seq_lens_cpu.clone()
        seq_lens_cpu[-1] -= tokens_skipped
    if seq_lens_cpu_upper_bound is not None:
        seq_lens_cpu_upper_bound = seq_lens_cpu_upper_bound.clone()
        seq_lens_cpu_upper_bound[-1] -= tokens_skipped

assert seq_lens_cpu_upper_bound is not None
max_seq_len = int(seq_lens_cpu_upper_bound.max()) # 使用上界而非 seq_lens_cpu

```

vllm/v1/worker/gpu/model_runner.py

主入口: 计算并传递 seq_lens_cpu_upper_bound, 并在 postprocess 中推进 CPU 镜像。

```

# vllm/v1/worker/gpu/model_runner.py
# prepare_inputs 中计算 CPU 上界
# 利用 numpy 纯 CPU 计算, 避免 GPU 张量访问
seq_lens_cpu_upper_bound_np = np.zeros(num_reqs_padded, dtype=np.int32)
np.add(
    self.req_states.num_computed_tokens_np[idx_mapping_np],
    num_scheduled_tokens,
    out=seq_lens_cpu_upper_bound_np[:num_reqs],
)
seq_lens_cpu_upper_bound = torch.from_numpy(seq_lens_cpu_upper_bound_np)
# 传入 InputBatch
return InputBatch(
    # ...
    seq_lens_cpu_upper_bound=seq_lens_cpu_upper_bound,
    # ...
)

# postprocess 中乐观推进 CPU 镜像, 保持上界正确
self.req_states.num_computed_tokens_np[idx_mapping_np] += (
    input_batch.num_scheduled_tokens
)

```

vllm/model_executor/layers/attention/cross_attention.py

交叉注意力: 用上界计算 num_computed_tokens, 避免读取 GPU 属性。

```

# vllm/model_executor/layers/attention/cross_attention.py
# 利用上界计算已计算的 token 数, 避免调用 .num_computed_tokens_cpu 属性 (触发同步)
query_lens_cpu = (
    common_attn_metadata.query_start_loc_cpu[1:]
)

```

```
- common_attn_metadata.query_start_loc_cpu[:-1]
)
assert common_attn_metadata.seq_lens_cpu_upper_bound is not None
num_computed_tokens_cpu = (
    common_attn_metadata.seq_lens_cpu_upper_bound - query_lens_cpu
)
num_cache_decodes = (num_computed_tokens_cpu > 0).sum().item()
```

评论区精华

主要讨论集中在 gemini-code-assist 对 eagle.py 中减法操作的潜在同步风险提出质疑，认为 `num_rejected_tokens` 可能是 GPU 张量，赋值给 `seq_lens_cpu_upper_bound` 会导致后续 `.item()/numpy()` 触发同步。作者直接回复 'Incorrect'，表明该问题已被充分考虑（`num_rejected_tokens` 已在 CPU 上或操作不会引入同步）。WoosukKwon 批准了该 PR。

- eagle.py 中上界减去 `num_rejected_tokens` 可能引入同步 (correctness): njhill 回复 'Incorrect'，表明该问题已考虑，`num_rejected_tokens` 在 CPU 上（来自 verification 逻辑），无同步风险。

风险与影响

- 风险：
 1. 核心路径变更（所有 attention 后端和 model runner）：引入 `seq_lens_cpu_upper_bound` 后，若在极端情况下上界过紧（如 speculative decoding 中 rejection 较多），可能导致 workspace 分配不足或 kernel dispatch 错误。但逻辑上上界始终 \geq 真实长度，不会低估。
 2. 兼容性风险：`seq_lens_cpu` 属性仍被部分后端使用（如 flash attention），但 PR 已确保主要路径（cross, MLA, flex）都已迁移。遗留路径若触发同步不会导致错误，只会性能回退。
 3. 测试覆盖：测试文件有相应调整，但未发现对边界条件（如空请求、dummy run）的专项测试。
 - 影响：用户影响：降低推理延迟，特别是在 long context 或 speculative decoding 场景下，消除 GPU→CPU 同步等待。系统影响：减少跨设备数据移动，提升吞吐量。团队影响：统一了上界计算模式，后续开发者应优先使用 `seq_lens_cpu_upper_bound` 而非 `seq_lens_cpu`。
 - 风险标记：核心路径变更，跨模块适配

关联脉络

- 暂无明显关联 PR