

PR #24008 完整报告

sgl-project/sglang

[diffusion] fix: align encoder of flux klein with official

合并时间: 2026-04-29 22:42

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/24008>

执行摘要

- 一句话: Flux Klein 文本编码器对齐官方, 引入 masked causal attention
- 推荐动作: 该 PR 值得精读, 特别是 `_masked_causal_attention` 的实现展示了如何在不依赖 FlashAttention 内置 masking 时手动处理变长批处理注意力, 对扩散模型文本编码器设计有参考价值。

功能与动机

使 Flux Klein 模型的文本编码器与官方实现一致, 修复先前实现中注意力掩码不匹配的问题, 从而提升生成图像质量。

实现拆解

1. 在 `Qwen3Attention.forward` 中添加可选参数 `attention_lengths`, 并将原有 `self.attn` 调用替换为新增的 `_masked_causal_attention` 方法。
2. 实现 `_masked_causal_attention`: 当传入的 `attention_lengths` 为 `None` 或所有长度等于序列长度时退化为标准 causal attention; 否则逐 batch 对有效部分使用 FlashAttention 的 causal 模式, 对填充部分通过 PyTorch SDPA (非因果) 计算, 并处理 GQA 键值重复。
3. 在 `Qwen3DecoderLayer.forward` 中添加 `attention_lengths` 参数并传递给 `self_attn`。
4. 在 `Qwen3ForCausalLM.forward` 中根据 `attention_mask` 计算每个 batch 的有效长度并向下传递。
5. 更新测试配置: 收紧 `flux_2_klein_image_t2i` 的一致性阈值 (`clip/ssim/psnr/mean_abs_diff` 均显著提高), 删除不再支持的 `flux_2_klein_ti2i_2_gpus` 测试用例及其相关的性能基线、精度配置等。

关键文件:

- `python/sglang/multimodal_gen/runtime/models/encoders/qwen3.py` (模块 文本编码器; 类别 source; 类型 data-contract; 符号 `_masked_causal_attention`): 核心编码器改动, 添加 masked causal attention 支持变长序列, 与官方对齐。
- `python/sglang/multimodal_gen/test/server/perf_baselines.json` (模块 性能基线; 类别 test; 类型 test-coverage): 删除 `flux_2_klein_ti2i_2_gpus` 的性能基线, 反映测试用例移除。

- python/sglang/multimodal_gen/test/server/consistency_threshold.json (模块 一致性阈值; 类别 test; 类型 test-coverage) : 收紧 flux_2_klein_image_t2i 的一致性阈值, 反映编码器对齐后的质量提升。
- python/sglang/multimodal_gen/test/server/gpu_cases.py (模块 GPU 用例; 类别 test; 类型 test-coverage) : 移除 flux_2_klein_ti2i_2_gpus 测试用例, 该场景不再支持。
- python/sglang/multimodal_gen/test/server/accuracy_config.py (模块 精度配置; 类别 test; 类型 test-coverage) : 移除 flux_2_klein_ti2i_2_gpus 的精度配置项。
- python/sglang/multimodal_gen/test/server/accuracy_testcase_configs.py (模块 测试配置; 类别 test; 类型 test-coverage) : 从精度测试列表中移除 flux_2_klein_ti2i_2_gpus。

关键符号: Qwen3Attention.forward, Qwen3Attention._masked_causal_attention, Qwen3DecoderLayer.forward, Qwen3ForCausalLM.forward

评论区精华

该 PR 无公开 review 讨论。

- 暂无高价值评论线程

风险与影响

- 风险: 主要风险在于新增的 masked causal attention 逻辑是否正确处理变长批处理场景, 特别是当 attention_lengths 与实际序列长度不一致时的边界情况。此外, 移除 2-gpu 测试用例可能降低对多 GPU 场景的覆盖, 需确认是否通过其他测试保障。测试阈值调整虽反映质量提升, 但也可能降低 CI 对微小退化的敏感度。
- 影响: 影响范围: 仅影响使用 Flux Klein 模型 (black-forest-labs/FLUX.2-klein-4B) 的用户, 文本编码质量应提升, 生成图像更符合预期。系统层面, 新增参数向后兼容 (默认 None), 不会影响已有功能。测试配置更新确保 CI 通过, 但移除了 2-gpu 用例, 后续若添加多 GPU 需重新评估。
- 风险标记: 注意力掩码逻辑变更, 测试用例删除

关联脉络

- 暂无明显关联 PR