

PR #27236 完整报告

sgl-project/sglang

[diffusion] Speed up lossless realtime RGB transport

合并时间: 2026-06-04 16:59

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/27236>

执行摘要

- 一句话: gzip 压缩级别从 1 降为 0, 加速 RGB 传输
- 推荐动作: 该 PR 变更简洁高效, 性能收益显著且风险可控。值得精读以了解实时传输场景下的性能优化思路。

功能与动机

PR body 明确指出, 需要在实时 CPU 预算内压缩已经无损的帧负载 ('avoid spending realtime CPU budget compressing already lossless frame payloads'), 以加速实时 RGB 传输。

实现拆解

1. 在 `python/sglang/multimodal_gen/runtime/utils/realtime_video.py` 中新增模块级常量 `_RAW_RGB_DELTA_GZIP_LEVEL = 0` (第 30 行)。
2. 在 `build_delta_gzip_raw_rgb_payload` 函数中, 将 `zlib.compressobj` 的 `level` 参数由字面量 1 改为引用该常量 (第 51-53 行), 使压缩级别降为 0 (仅存储 / 不压缩)。
3. 添加注释说明: 保持 gzip 帧格式以实现无损传输, 同时不消耗实时压缩预算。
4. 纯源码改动, 无测试、配置或部署配套变更。

关键文件:

- `python/sglang/multimodal_gen/runtime/utils/realtime_video.py` (模块 多模态生成; 类别 source; 类型 core-logic): 核心变更文件: 新增常量 `_RAW_RGB_DELTA_GZIP_LEVEL = 0`, 并替换 `zlib.compressobj` 的压缩级别参数, 实现性能优化。

关键符号: `build_delta_gzip_raw_rgb_payload`

关键源码片段

`python/sglang/multimodal_gen/runtime/utils/realtime_video.py`

核心变更文件: 新增常量 `_RAW_RGB_DELTA_GZIP_LEVEL = 0`, 并替换 `zlib.compressobj` 的压缩级别参数, 实现性能优化。

```
# 定义 gzip 压缩级别为 0 (仅存储, 不压缩),  
# 避免在实时 CPU 预算内对已无损的帧负载进行压缩。  
_RAW_RGB_DELTA_GZIP_LEVEL = 0
```

```
def build_delta_gzip_raw_rgb_payload(
    frames: list[bytes],
    *,
    reference_frame: bytes | None = None,
) -> bytes:
    # ... 前置校验 ...

    previous = (
        np.frombuffer(reference_frame, dtype=np.uint8)
        if reference_frame is not None
        else None
    )
    # 使用 level=0 保持 gzip 帧格式，但不进行实际压缩，
    # 从而在保留无损传输能力的同时大幅降低 CPU 开销。
    compressor = zlib.compressobj(
        level=_RAW_RGB_DELTA_GZIP_LEVEL, method=zlib.DEFLATED, wbits=31
    )
    compressed_chunks = []
    for frame in frames:
        # ... delta 计算与压缩逻辑 ...
```

评论区精华

PR 无 review 评论，仅两个 bot 评论（配额限制和 CI 重跑指令），无技术讨论。

- 暂无高价值评论线程

风险与影响

- 风险：

1. 帧负载尺寸增加：PR body 指出 payload 从约 6.8-8.3 MiB 增至 13.7 MiB，但仍低于现有的 16 MiB 实时 payload 保护阈值，在可控范围内。
2. 兼容性风险：restore_delta_gzip_raw_rgb_payload 使用 zlib.decompress，gzip 格式（wbits=31）保持不变，仅压缩级别降为 0，因此接收端无需任何修改，可以正确解压。
3. 回归风险低：更改仅涉及一个常量的替换，逻辑不变。

- 影响：

1. 用户：实时 RGB 传输吞吐量提升约 22.4%，延迟显著降低（raw_payload_build 从约 226ms 降至约 20ms），无功能性损失。
2. 系统：网络带宽消耗增加（payload 大小翻倍约），但仍低于 16MiB 阈值，影响有限。
3. 团队：变更极小，易于理解和维护。 - 风险标记：负载尺寸增大，缺少测试配套

关联脉络

- 暂无明显关联 PR