

PR #38405 完整报告

vllm-project/vllm

[Frontend] Add multimodal support to /inference/v1/generate endpoint

合并时间: 2026-04-18 11:31

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/38405>

执行摘要

- 一句话: 为解耦推理端点添加多模态支持, 实现渲染到生成的零客户端转换。
- 推荐动作: 建议开发者和架构师精读 `vllm/entrypoints/serve/disagg/mm_serde.py` 的序列化实现, 关注整数张量处理和 `Msgpack` 配置, 以及 `tests/entrypoints/serve/disagg/test_serving_multimodal_tokens.py` 的端到端测试设计, 以理解多模态数据流的集成方式。

功能与动机

PR body 说明: 'Adds multi-modal support to the disaggregated `/inference/v1/generate` endpoint, enabling coordinators to pass pre-processed multi-modal features (e.g. `pixel_values`, `image_grid_thw`) directly to vLLM workers over HTTP.' 目的是支持多模态推理的解耦流程, 减少客户端预处理负担。

实现拆解

1. 添加序列化工具: 在 `vllm/entrypoints/serve/disagg/mm_serde.py` 中新增 `encode_mm_kwargs_item` 和 `decode_mm_kwargs_item` 函数, 使用 `MsgpackEncoder` 和 `MsgpackDecoder` 将 `MultiModalKwargsItem` 序列化为 `base64` 字符串, 支持整数和浮点张量。 - 原因: 为了在 HTTP 请求中高效传输张量数据。 - 影响: 为生成端点提供了数据反序列化能力, 并确保整数张量 (如 `image_grid_thw`) 的兼容性。
2. 修改服务逻辑: 在 `vllm/entrypoints/serve/disagg/serving.py` 的 `serve_tokens` 函数中, 当请求包含 `features` 时, 解析 `MultiModalFeatures`, 将 `PlaceholderRangeInfo` 转换为 `PlaceholderRange`, 并反序列化 `kwargs_data` 以构建 `mm_input`。 - 原因: 集成多模态特征到现有的推理流程中。 - 影响: 使得生成端点能够处理从渲染端点传来的预处理数据, 实现零客户端转换。
3. 添加测试配套: 新增测试文件如 `tests/entrypoints/serve/disagg/test_serving_multimodal_tokens.py` (端到端测试) 和 `tests/entrypoints/openai/test_mm_serde.py` (序列化测试), 并更新 `tests/utils/_test_serial_utils.py` 以支持整数数据类型测试。 - 原因: 确保功能正确性和健壮性, 覆盖边界情况。 - 影响: 提高了代码覆盖率和可靠性, 验证了 Qwen3-VL 模型的多模态能力。
4. 提供示例和文档: 新增 `examples/online_serving/disaggregated_serving/example_mm_serve.py` 示例脚本, 展示从渲染到生成的完整流程, 并更新相关协议文件。 - 原因: 帮助用户理解和使用新功能。 - 影响: 改善了开发者体验, 促进了功能采纳。

关键文件:

- `examples/online_serving/disaggregated_serving/example_mm_serve.py` (模块 示例脚本; 类别 `source`; 类型 `core-logic`; 符号 `make_data_url`, `main`): 提供了端到端多模态服务示例, 展示从渲染到生成的完整流程, 是功能使用的最佳实践。
- `vllm/entrypoints/serve/disagg/mm_serde.py` (模块 序列化工具; 类别 `source`; 类型 `core-logic`; 符号 `encode_mm_kwargs_item`, `decode_mm_kwargs_item`): 核心序列化工具, 负责多模态特征的编码和解码, 是实现 HTTP 传输张量数据的关键模块。
- `vllm/entrypoints/serve/disagg/serving.py` (模块 服务入口; 类别 `source`; 类型 `dependency-wiring`): 修改了生成端点的服务逻辑, 集成多模态特征处理, 是功能落地的核心入口。
- `tests/entrypoints/serve/disagg/test_serving_multimodal_tokens.py` (模块 端到端测试; 类别 `test`; 类型 `test-coverage`; 符号 `test_image`, `server`, `client`, `test_render_to_generate_roundtrip`): 端到端测试文件, 验证多模态渲染到生成的完整流程, 确保功能正确性和稳定性。

关键符号: `make_data_url`, `test_render_to_generate_roundtrip`,
`encode_mm_kwargs_item`, `decode_mm_kwargs_item`, `serve_tokens`

关键源码片段

[examples/online_serving/disaggregated_serving/example_mm_serve.py](#)

提供了端到端多模态服务示例, 展示从渲染到生成的完整流程, 是功能使用的最佳实践。

```
import io
import pybase64 as base64
import requests
from PIL import Image
from transformers import AutoTokenizer

BASE_URL = "http://localhost:8000"
MODEL_NAME = "Qwen/Qwen3-VL-2B-Instruct"

def make_data_url(image: Image.Image) -> str:
    """Encode a PIL image as a base64 data URL."""
    buf = io.BytesIO()
    image.save(buf, format="PNG")
    b64 = base64.b64encode(buf.getvalue()).decode()
    return f"data:image/png;base64,{b64}" # 生成标准的data URL格式, 用于HTTP传输

def main():
    # 创建测试图像 (红色方块) 并编码
    image = Image.new("RGB", (224, 224), color=(255, 0, 0))
    data_url = make_data_url(image)
    print("Created 224x224 red test image")

    # 渲染请求: 预处理多模态消息为token IDs和特征
```

```

render_payload = {
    "model": MODEL_NAME,
    "messages": [
        {
            "role": "user",
            "content": [
                {"type": "image_url", "image_url": {"url": data_url}},
                {"type": "text", "text": "What color is this image? Answer in one word."},
            ],
        }
    ],
}
render_resp = requests.post(f"{BASE_URL}/v1/chat/completions/render", json=render_payload)
render_resp.raise_for_status()
render_data = render_resp.json()

# 生成请求：直接使用渲染输出，仅添加采样参数
generate_payload = render_data
generate_payload["sampling_params"] = {"max_tokens": 20, "temperature": 0.0}
gen_resp = requests.post(f"{BASE_URL}/inference/v1/generate", json=generate_payload)
gen_resp.raise_for_status()
gen_data = gen_resp.json()

# 解码并输出结果
output_ids = gen_data["choices"][0]["token_ids"]
tokenizer = AutoTokenizer.from_pretrained(MODEL_NAME)
text = tokenizer.decode(output_ids, skip_special_tokens=True)
print(f"Generated text: {text!r}")

if "red" in text.lower():
    print("Model correctly identified the red image.") # 验证模型正确识别图像颜色
else:
    print(f"WARNING: Expected 'red' in output, got: {text!r}")

if __name__ == "__main__":
    main() # 执行示例主函数

```

vllm/entrypoints/serve/disagg/mm_serde.py

核心序列化工具，负责多模态特征的编码和解码，是实现 HTTP 传输张量数据的关键模块。

```

from __future__ import annotations
import pybase64
from vllm.multimodal.inputs import MultiModalKwargsItem
from vllm.v1.serial_utils import MsgpackDecoder, MsgpackEncoder

_encoder = MsgpackEncoder(size_threshold=2**62) #
强制所有张量内联，避免外部存储以简化传输
_decoder = MsgpackDecoder(t=MultiModalKwargsItem) #
指定目标类型为MultiModalKwargsItem，确保反序列化正确性

```

```
def encode_mm_kwargs_item(item: MultiModalKwargsItem) -> str:
    """Serialize a MultiModalKwargsItem to a base64 string."""
    bufs = _encoder.encode(item) # 使用Msgpack编码为二进制缓冲区
    assert len(bufs) == 1, "All tensors should be inline" # 验证张量全部内联, 防止数据分片
    return pybase64.b64encode(bufs[0]).decode("ascii") # 转换为base64字符串, 便于JSON/
    HTTP传输
```

```
def decode_mm_kwargs_item(data: str) -> MultiModalKwargsItem:
    """Deserialize a base64 string back to a MultiModalKwargsItem."""
    raw = pybase64.b64decode(data) # 解码base64字符串为原始二进制数据
    return _decoder.decode(raw) # 使用Msgpack解码回MultiModalKwargsItem对象
```

评论区精华

- 整数测试范围改进: [gemini-code-assist\[bot\]](#) 指出 `tests/utils/_test_serial_utils.py` 中的整数范围太小, 可能导致测试不充分, 作者随后扩大范围以覆盖 `int32` 和 `int64` 的边界值。
- 代码风格建议: [NickLucche](#) 建议在 `vllm/entrypoints/serve/disagg/serving.py` 中使用海象操作符简化代码, 作者采纳并调整。
- 缓存跳过以避免 bug: [DarkLight1337](#) 提到需要添加 `skip_mm_cache=True` 以避免 bug #38543, 作者在 `render` 和 `generate` 端点中添加此参数, 确保多模态缓存的一致性和幂等性。
 - 整数测试范围改进 (correctness): 作者采纳建议, 修改 `_build_integer_tensor` 函数, 使用更接近 `int32` 和 `int64` 极限的范围, 提高了测试健壮性。
 - 缓存跳过以避免 bug (bugfix): 作者在 `render` 和 `generate` 端点中添加了 `skip_mm_cache=True` 参数, 解决了潜在的多模态缓存冲突问题。

风险与影响

- 风险: - 序列化兼容性风险: 新增的 `encode_mm_kwargs_item` 和 `decode_mm_kwargs_item` 依赖于 `MsgpackEncoder`, 如果序列化协议变更可能导致数据不兼容, 需注意版本控制。
- 整数溢出风险: 虽然测试已改进, 但整数张量序列化时仍可能因值范围过大引发溢出, 尤其在 `image_grid_thw` 等数据类型中。
- 缓存一致性问题: 添加 `skip_mm_cache=True` 解决了已知 bug, 但可能增加缓存未命中, 影响性能, 需监控多模态场景下的缓存效率。
- 影响: - 用户影响: 多模态模型开发者现在可以利用解耦端点进行高效推理, 减少了客户端预处理负担, 提升了用户体验。
- 系统影响: 扩展了 `vLLM` 的前端服务能力, 支持更复杂的多模态工作流, 可能增加 HTTP 负载和序列化开销。
- 团队影响: 增加了代码库的复杂性, 需维护新的序列化工具和服务逻辑, 但通过测试和示例降低了维护难度。
- 风险标记: 序列化兼容性风险, 整数溢出风险, 缓存一致性问题

关联脉络

- PR #40089 [Misc][UX] Map mimo reasoning and tooling parsers: 同样涉及多模态功能扩展, 添加了 MiMo-V2-Flash 模型的推理和工具解析器映射。
- PR #39870 [BugFix] Support custom tool parsers when tool_choice is required and named function: 涉及前端工具调用解析, 与当前 PR 在多模态和前端模块有功能关联。