

# PR #35809 完整报告

vllm-project/vllm

[Models] Cohere Transcribe

合并时间: 2026-03-18 05:04

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/35809>

## 执行摘要

本 PR 添加了 Cohere ASR 语音转录模型支持到 vLLM，包括新增模型实现、自定义处理器和测试集成。它利用了近期变长编码器改进，并通过引入 `skip_decoder_start_token` 标志优化输入处理。影响范围涵盖多模态推理扩展，但存在硬编码路径和设备依赖风险，需后续更新。

## 功能与动机

此变更旨在支持 Cohere 的自动语音识别模型，解决语音转录场景的需求。根据 PR body，动机是“添加 Cohere Transcribe 到 vLLM”，并利用最近几个月对调度器做的变长编码器变更（如 PR 31058），以放松固定长度编码器输入限制。由于模型初始没有 Hugging Face 实现，因此添加了自定义处理器和配置，后续改为信任远程代码路径。

## 实现拆解

实现按模块拆解如下：

- 模型核心：新增 `vllm/model_executor/models/cohere_asr.py`，实现 `CohereASRForConditionalGeneration` 类，包含编码器（Conformer 架构）和解码器（Transformer），支持音频特征提取和变长输入处理。
- 处理器集成：新增 `vllm/transformers_utils/processors/cohere_asr.py`，定义 `CohereASRProcessor`，处理音频到梅尔频谱图的转换，并集成到多模态框架中。
- 输入处理增强：修改 `vllm/inputs/data.py` 中的 `build_enc_dec_inputs` 函数，引入 `skip_decoder_start_token` 参数，避免为 Cohere ASR 添加解码器起始令牌；相关更改传播到 `vllm/inputs/preprocess.py` 和 `vllm/renderers/base.py`。
- 配置与注册：更新 `vllm/model_executor/models/registry.py` 注册新模型；在 `vllm/transformers_utils/model_arch_config_convertor.py` 中添加 `CohereAsrModelArchConfigConvertor`，支持模型架构解析。
- 测试与示例：修改 `tests/entrypoints/openai/correctness/test_transcription_api_correctness.py` 以支持多模型 WER 测试；更新 `examples/offline_inference/audio_language.py` 添加离线推理示例，但路径暂硬编码。

## 评论区精华

Review 讨论中最有价值的交锋包括：

- 设备硬编码: gemini-code-assist[bot] 指出: “Hardcoding the device to 'cuda' will cause the model to fail on non-CUDA backends”。作者回应移除硬编码, 提升后端兼容性。
- 路径可移植性: gemini-code-assist[bot] 批评: “The model path is hardcoded to a local path ... makes the example code not portable”。作者标记为待更新, 但风险仍存。
- 设计优化: DarkLight1337 建议: “Let's make this based on a new flag in EncDecMultiModalProcessor”, 作者采纳并实现, 避免核心代码侵入。
- 代码共享: 作者解释为何不与 FireRedASR2 共享 Conformer: “short answer: not worth it ... fundamental structural difference”, 体现技术权衡。

## 风险与影响

技术风险:

- 硬编码路径 (如 /host/engines/vllm/audio/2b-release) 在示例和测试中, 导致其他环境运行失败, 需模型发布后修复。
- 修改 build\_enc\_dec\_inputs 等核心函数, 可能影响其他编码器 - 解码器模型, 需通过现有测试验证无回归。
- 设备处理初始硬编码, 虽已修复, 但类似模式需在代码库中审查。

影响范围:

- 用户可通过离线命令行和在线 API 使用新 ASR 模型, 扩展 vLLM 多模态能力。
- 系统层面集成变长编码器支持, 增强调度灵活性; 团队需维护新增代码, 但 review 促进代码质量提升。

## 关联脉络

与历史 PR 的关联揭示功能演进方向:

- 直接相关: PR 38120 “[Cohere] Enable Cohere-Transcribe” 可能为本 PR 的后续启用步骤, 涉及文档和测试完善。
- 基础改进: PR 31058、29268、29278 被引用, 提供了变长编码器调度器变更, 是本 PR 实现的前提, 显示 vLLM 在音频处理方向的持续投入。
- 多模态趋势: 结合近期 PR 如 38119 (添加 numpy 数组嵌入支持), 表明仓库正积极扩展多模态功能, Cohere ASR 加入是这一趋势的体现。