

# PR #43519 完整报告

vllm-project/vllm

Add model support for granite speech plus

合并时间: 2026-06-04 22:47

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/43519>

## 执行摘要

- 一句话: 支持 Granite Speech Plus 模型推理
- 推荐动作: 值得精读的设计模式: 通过提取工厂方法 `_build_encoder` 实现子类化复用, 避免了复制粘贴基类 `__init__`。同时, 该 PR 展示了在 vLLM 中添加新多模态模型的标准流程: 模型代码、注册、测试、文档联动。对于需要扩展语音模型变体的开发者有参考价值。

## 功能与动机

Adds support for the `GraniteSpeechPlus` architecture (`GraniteSpeechPlusForConditionalGeneration`) to vLLM, enabling inference for models such as `ibm-granite/granite-speech-4.1-2b-plus`.

## 实现拆解

1. 提取 `_build_encoder` 钩子 (`granite_speech.py`): 将 `__init__` 中直接实例化 `GraniteSpeechCTCEncoder` 的代码替换为调用 `self._build_encoder(...)`, 并在基类中提供默认实现返回 `GraniteSpeechCTCEncoder`, 为子类重写开放入口。
2. 新增 `granite_speech_plus.py`: 定义 `GraniteSpeechPlusCTCEncoder` (继承 `GraniteSpeechCTCEncoder`), 重写 `forward` 方法实现层隐藏状态拼接与中间层残差逻辑; 定义 `GraniteSpeechPlusForConditionalGeneration` (继承 `GraniteSpeechForConditionalGeneration`), 重写 `_build_encoder` 返回自定义编码器, 并通过装饰器注册多模态处理器。
3. 模型注册 (`registry.py`): 在 `GraniteSpeechForConditionalGeneration` 条目之后添加 `GraniteSpeechPlusForConditionalGeneration` 映射到模块 `granite_speech_plus`。
4. 测试集成: 将 `granite-speech-4.1-2b-plus` 加入现有 `test_granite_speech.py` 的 `models` 字典, 复用餐具类; 在 `tests/models/registry.py` 中添加 `HfExamplesInfo` 条目, 并设定 `min_transformers_version="5.8.0"` 以跳过低版本 CI 失败。
5. 文档更新 (`supported_models.md`): 在两处表格 (多模态模型与 ASR 模型) 中增加 `GraniteSpeechPlusForConditionalGeneration` 行。

关键文件:

- `vllm/model_executor/models/granite_speech_plus.py` (模块 语音模型; 类别 `source`; 类型 `data-contract`; 符号 `GraniteSpeechPlusCTCEncoder`, `forward`, `GraniteSpeechPlusForConditionalGeneration`, `_build_encoder`): 核心新增文件, 包含自

定义 CTC 编码器和生成模型类，实现层拼接与残差逻辑。

- `vllm/model_executor/models/granite_speech.py` (模块 语音模型; 类别 source; 类型 data-contract; 符号 `_build_encoder`) : 基类重构, 提取 `_build_encoder` 工厂方法, 使子类无须重复 `__init__` 即可更换编码器。
- `vllm/model_executor/models/registry.py` (模块 模型注册; 类别 source; 类型 data-contract) : 注册新架构到模型映射表, 使 vLLM 能根据 HF config 自动加载对应实现。
- `tests/models/multimodal/generation/test_granite_speech.py` (模块 测试套件; 类别 test; 类型 test-coverage) : 将新模型加入现有测试的 `models` 字典, 复用参数数量和 `logprobs` 对比逻辑。
- `tests/models/registry.py` (模块 测试注册; 类别 test; 类型 test-coverage) : 添加 `HfExamplesInfo` 条目, 使新模型出现在测试矩阵中, 并设定 `transformers` 版本门槛。
- `docs/models/supported_models.md` (模块 文档; 类别 docs; 类型 documentation) : 更新多模态和 ASR 模型表格, 方便用户查找可用模型。

关键符号: `_build_encoder`, `GraniteSpeechPlusCTCEncoder.forward`,  
`GraniteSpeechPlusForConditionalGeneration._build_encoder`

## 关键源码片段

### `vllm/model_executor/models/granite_speech_plus.py`

核心新增文件, 包含自定义 CTC 编码器和生成模型类, 实现层拼接与残差逻辑。

```
# vllm/model_executor/models/granite_speech_plus.py

class GraniteSpeechPlusCTCEncoder(GraniteSpeechCTCEncoder):
    def forward(self, hidden_states: torch.Tensor) -> torch.Tensor:
        hidden_states = self.input_linear(hidden_states)
        # cat_hidden_layers 为非负层索引, 这些中间层的 hidden states
        # 会在最终 hidden states 之前沿特征维度拼接, 其中 0 表示编码器输入。
        cat_layers = set(self.config.cat_hidden_layers or [])
        exported_hidden_states = []

        if 0 in cat_layers:
            exported_hidden_states.append(hidden_states)

        for idx, layer in enumerate(self.layers, start=1):
            hidden_states = layer(hidden_states, attention_dists=self.attention_dists)

            # 跳过最后一层, 因为其输出会在循环后自动附加, 避免重复。
            if idx in cat_layers and idx != self.num_layers:
                exported_hidden_states.append(hidden_states)

            # 中间层残差: 将第 N/2 层的输出经过 out -> softmax -> out_mid 后加回。
            if idx == self.num_layers // 2:
                hidden_states_mid = hidden_states.clone()
                hidden_states_mid, _ = self.out(hidden_states_mid)
                hidden_states_mid = self.softmax(hidden_states_mid)
```

```

        hidden_states_mid, _ = self.out_mid(hidden_states_mid)
        hidden_states += hidden_states_mid

    if exported_hidden_states:
        hidden_states = torch.cat([*exported_hidden_states, hidden_states], dim=-1)
    return hidden_states

@MULTIMODAL_REGISTRY.register_processor(
    GraniteSpeechMultiModalProcessor,
    info=GraniteSpeechMultiModalProcessingInfo,
    dummy_inputs=GraniteSpeechDummyInputsBuilder,
)
class GraniteSpeechPlusForConditionalGeneration(GraniteSpeechForConditionalGeneration):
    supported_languages = ISO639_1_SUPPORTED_LANGS

    def _build_encoder(self, config: PretrainedConfig,
                      quant_config: QuantizationConfig | None,
                      prefix: str) -> GraniteSpeechCTCEncoder:
        return GraniteSpeechPlusCTCEncoder(
            config=config, quant_config=quant_config, prefix=prefix)

```

### vllm/model\_executor/models/granite\_speech.py

基类重构，提取 `_build_encoder` 工厂方法，使子类无须重复 `__init__` 即可更换编码器。

# vllm/model\_executor/models/granite\_speech.py ( 相关片段 )

```

class GraniteSpeechForConditionalGeneration(...):

    def __init__(self, *, vllm_config: VllmConfig, prefix: str = ""):
        super().__init__()
        config = vllm_config.model_config.hf_config
        quant_config = vllm_config.quant_config
        # ... 其他初始化 ...

        with self._mark_tower_model(vllm_config, "audio"):
            # 使用工厂方法创建编码器，子类可重写以返回不同的编码器实例。
            self.encoder = self._build_encoder(
                config=config.encoder_config,
                quant_config=quant_config,
                prefix=maybe_prefix(prefix, "encoder"),
            )
            # 投影器保持不变
            self.projector = GraniteSpeechEncoderProjector(
                config=config, quant_config=quant_config,
                cache_config=cache_config,
                prefix=maybe_prefix(prefix, "projector"),
            )

    def _build_encoder(self, config: PretrainedConfig,

```

```
quant_config: QuantizationConfig | None,  
prefix: str) -> "GraniteSpeechCTCEncoder":  
# 基类默认实现, 子类可覆盖。  
return GraniteSpeechCTCEncoder(  
    config=config, quant_config=quant_config, prefix=prefix)
```

## 评论区精华

1. 测试文件合并建议: alex-jw-brooks 建议将独立测试文件合并到现有 `test_granite_speech.py` 的 `models` 字典中, 避免重复测试代码。作者采纳并删除独立文件。
  2. 命名风格统一: alex-jw-brooks 指出 `_ISO639_1_SUPPORTED_LANGS` 前导下划线在基类中不使用, 建议移除以保持一致性。作者采纳。
  3. 注释与逻辑防护: alex-jw-brooks 要求对 `cat_hidden_layers` 添加注释说明其含义 (非负索引, 附加在最终 hidden states 之外), 并防止列表包含最后一层时重复追加。作者添加注释及 `idx != self.num_layers` 条件。
  4. Forward 逻辑正确性争议: gemini-code-assist[bot] 提出重复层处理、拼接顺序、in-place 更新顺序等潜在问题。作者逐一回应, 认为当前实现符合 Hugging Face 参考逻辑, 不需要重复层, 且顺序不影响数值结果。未引发进一步修改。
- 测试文件组织方式 (design): 采纳, 新模型复用现有测试逻辑。
  - 命名风格统一 (style): 移除前导下划线。
  - `cat_hidden_layers` 注释与防护逻辑 (correctness): 添加注释, 并在条件中加入 `idx != self.num_layers`。
  - Forward 逻辑正确性争议 (correctness): 作者解释后未修改, gemini 未进一步反驳。潜在风险较低。

## 风险与影响

- 风险:
  1. Transformers 版本依赖: 新模型要求 `transformers>=5.8.0`, CI 默认锁定的 `5.5.3` 无法识别 `model_type="granite_speech_plus"`。PR 通过测试注册中的 `min_transformers_version` 跳过低版本运行, 但用户使用低版本时可能遇到导入或配置错误。
  2. 新编码器 Forward 正确性: `GraniteSpeechPlusCTCEncoder.forward` 涉及层的状态拼接和中间残差, 逻辑与 Hugging Face 参考实现的对等性仅在单一配置 (`granite-speech-4.1-2b-plus`) 上验证。若未来出现含重复 `cat_hidden_layers` 索引的配置, 当前 `set` 去重可能改变语义。
  3. 测试覆盖局限: 仅添加了一个模型变体, 未覆盖边界情况 (如空 `cat_hidden_layers`、含重复索引、不同层数配置等)。- 影响: 用户: 可使用 `ibm-granite/granite-speech-4.1-2b-plus` 等 Granite Speech Plus 模型进行推理。系统: 新增模型不影响现有模型加载, 但需满足 `transformers` 版本约束。团队: 提供基于 `_build_encoder` 钩子的可扩展模式, 未来类似变体 (如其他编码器结构) 仅需重写该方法即可复用基类逻辑。- 风险标记: `transformers` 版本约束, 新编码器逻辑风险, 测试覆盖有限

## 关联脉络

- 暂无明显关联 PR