

# PR #21569 完整报告

sgl-project/sglang

Upgrade transformers to 5.5.3 and refactor hf\_transformers\_utils into subpackage

合并时间: 2026-04-16 11:03

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/21569>

## 执行摘要

- 一句话: 将 transformers 升级至 5.5.3 并重构 hf\_transformers\_utils 为子包, 解决兼容性问题。
- 推荐动作: 建议技术管理者和工程师精读此 PR, 特别是 compat.py 中的补丁设计和 tokenizer.py 中的 TokenizersBackend 处理策略, 这些展示了在依赖升级中的兼容性保障技巧。

## 功能与动机

根据 PR body, 主要动机是升级 transformers 以获取最新功能, 并重构 hf\_transformers\_utils.py 以改善代码结构。具体表述为: 'upgrade pinned transformers from 5.3.0 to 5.5.3 with compatibility patches' 和 'refactor hf\_transformers\_utils.py into hf\_transformers/ subpackage'。

## 实现拆解

1. 依赖升级: 更新所有平台配置文件 (如 pyproject.toml、pyproject\_xpu.toml) 中的 transformers 版本从 5.3.0 到 5.5.3, 并添加必要的依赖如 easydict 和 addict 以支持远程代码导入验证。
2. 兼容性补丁: 新增 python/sglang/srt/utils/hf\_transformers/compat.py 模块, 包含猴子补丁解决 transformers v5 的回归问题, 如 rope 配置验证、移除符号 shim、图像处理器 kwargs 过滤等。
3. 代码重构: 将 hf\_transformers\_utils.py 拆分为多个子模块:
  - common.py: 共享助手函数, 如 download\_from\_hf 和 get\_rope\_config。
  - config.py: 配置加载逻辑, 包括 get\_config 函数。
  - tokenizer.py: tokenizer 加载工具, 新增 \_load\_tokenizer\_by\_declared\_class 处理 TokenizersBackend 问题。
  - processor.py: 处理器加载逻辑。
  - mistral\_utils.py: Mistral 模型特定工具, 从原位置迁移而来。
4. 向后兼容: 保持 hf\_transformers\_utils.py 作为薄层重新导出所有子模块符号, 确保现有导入路径不变。
5. 测试配套: 新增 test/registered/unit/utils/test\_hf\_transformers.py 测试文件, 验证子包行为、rope 配置助手和兼容性补丁。

关键文件：

- `python/sglang/srt/utils/hf_transformers/tokenizer.py` (模块 `Tokenizer` 加载; 类别 `source`; 类型 `core-logic`; 符号 `_load_tokenizer_by_declared_class`, `TokenizerWarningsFilter`, `filter`, `_resolve_tokenizer_name`) : 新增的 `tokenizer` 加载模块, 核心解决 `transformers v5` 中 `TokenizersBackend` 问题, 影响所有模型加载路径。
- `python/sglang/srt/utils/hf_transformers/compat.py` (模块 兼容性补丁; 类别 `source`; 类型 `dependency-wiring`; 符号 `apply_all`, `normalize_rope_scaling_compat`, `_patch`, `_ensure_gguf_version`) : 新增的兼容性补丁模块, 集中处理 `transformers v5` 的回归问题和远程代码兼容性, 确保升级后系统稳定。
- `python/sglang/srt/utils/hf_transformers/common.py` (模块 通用工具; 类别 `source`; 类型 `core-logic`; 符号 `download_from_hf`, `_resolve_local_or_cached_file`, `check_gguf_file`, `get_rope_config`) : 新增的共享助手模块, 包含下载、配置检测和 `rope` 配置等通用函数, 是子包的基础。
- `python/sglang/srt/utils/hf_transformers_utils.py` (模块 兼容性入口; 类别 `source`; 类型 `dependency-wiring`; 符号 `download_from_hf`, `get_rope_config`, `_patch_text_config`, `get_hf_text_config`) : 原始单文件模块被大幅修改为重新导出 `shim`, 保持向后兼容, 是变更的入口点。
- `test/registered/unit/utils/test_hf_transformers.py` (模块 单元测试; 类别 `test`; 类型 `test-coverage`; 符号 `TestNormalizeRopeScalingCompat`, `test_adds_type_from_rope_type`, `test_preserves_existing_type`, `test_no_op_when_no_rope_scaling`) : 新增的单元测试文件, 验证子包行为、`rope` 配置助手和兼容性补丁, 确保升级质量。

关键符号: `apply_all`, `normalize_rope_scaling_compat`,  
`_load_tokenizer_by_declared_class`, `get_config`, `download_from_hf`, `get_rope_config`

## 关键源码片段

### `python/sglang/srt/utils/hf_transformers/tokenizer.py`

新增的 `tokenizer` 加载模块, 核心解决 `transformers v5` 中 `TokenizersBackend` 问题, 影响所有模型加载路径。

```
# 关键函数: _load_tokenizer_by_declared_class
# 处理 transformers v5 中 AutoTokenizer 返回通用 TokenizersBackend 的问题
def _load_tokenizer_by_declared_class(tokenizer_name, *args, **kwargs):
    """Load tokenizer by the class declared in tokenizer_config.json.
    AutoTokenizer resolves to TokenizersBackend when the model's config
    model_type has no tokenizer class mapping (e.g. deepseek_vl_v2), even
    though tokenizer_config.json declares a standard class like
    LlamaTokenizerFast. Returns None if it cannot improve on AutoTokenizer.
    """
    import transformers
    try:
        revision = kwargs.get("revision") or kwargs.get("tokenizer_revision")
        config_file = _resolve_local_or_cached_file(
```

```

        tokenizer_name, "tokenizer_config.json", revision
    )
    with open(config_file) as f:
        tok_config = json.load(f)
        tok_class_name = tok_config.get("tokenizer_class")
except FileNotFoundError:
    return None
except (OSError, json.JSONDecodeError) as e:
    logger.debug(
        "Failed to read tokenizer_config.json for %s: %s", tokenizer_name, e
    )
    return None
if not tok_class_name:
    return None
# 跳过基础类, 如 PreTrainedTokenizer, 因为它们缺少必要方法
if tok_class_name in ("PreTrainedTokenizer", "PreTrainedTokenizerBase"):
    return None
tok_cls = getattr(transformers, tok_class_name, None)
if tok_cls is None and kwargs.get("trust_remote_code"):
    # 尝试通过 auto_map 动态加载自定义类
    try:
        auto_map = tok_config.get("auto_map", {})
        auto_tok_ref = auto_map.get("AutoTokenizer")
        if isinstance(auto_tok_ref, (list, tuple)):
            auto_tok_ref = auto_tok_ref[0]
        if auto_tok_ref:
            from transformers.dynamic_module_utils import get_class_from_dynamic_module
            tok_cls = get_class_from_dynamic_module(
                auto_tok_ref,
                tokenizer_name,
                code_revision=revision,
            )
        except (OSError, ImportError, ValueError, RuntimeError) as e:
            logger.debug("Dynamic module lookup for %s failed: %s", tok_class_name, e)
if tok_cls is None:
    return None
logger.info(
    "Loading tokenizer for %s directly as %s (bypassing AutoTokenizer)",
    tokenizer_name,
    tok_class_name,
)
try:
    return tok_cls.from_pretrained(tokenizer_name, *args, **kwargs)
except (OSError, ValueError, TypeError, ImportError) as e:
    logger.warning(
        "Direct load as %s failed for %s: %s. Falling back to AutoTokenizer result.",
        tok_class_name,
        tokenizer_name,
        e,
    )

```

```
)  
return None
```

## python/sglang/srt/utils/hf\_transformers/compat.py

新增的兼容性补丁模块，集中处理 transformers v5 的回归问题和远程代码兼容性，确保升级后系统稳定。

```
# 关键函数: normalize_rope_scaling_compat  
# 确保 rope_scaling 字典包含 "type" 键，以兼容远程代码模型  
def normalize_rope_scaling_compat(config) -> None:  
    """Ensure rope_scaling dicts have "type" alongside "rope_type".  
    Transformers v5 standardises rope_scaling to use "rope_type" and may  
    omit the legacy "type" key. Remote-code models (e.g. Kimi-VL) still  
    read rope_scaling["type"], causing a KeyError. This helper adds  
    "type" from "rope_type" whenever it is missing, recursively across  
    the config and all its sub-configs.  
    """  
  
    def _patch(cfg):  
        rs = getattr(cfg, "rope_scaling", None)  
        if isinstance(rs, dict) and "rope_type" in rs and "type" not in rs:  
            rs["type"] = rs["rope_type"]  
        # 递归处理子配置，如 text_config、vision_config 等  
        for attr in (  
            "text_config",  
            "llm_config",  
            "language_config",  
            "vision_config",  
            "thinker_config",  
        ):  
            sub = getattr(cfg, attr, None)  
            if sub is not None:  
                _patch(sub)  
    _patch(config)
```

## 评论区精华

Review 中重点讨论了两个问题：

1. `get_config` 错误处理：kpham-ssl 指出新代码中 `get_config` 函数丢失了 `ValueError` 处理器，可能导致 `deepseek_v32` 模型崩溃。JustinTong0323 回应并修复为捕获 (`ValueError`, `KeyError`) 并匹配字符串。
  2. `TokenizersBackend` 策略：kpham-ssl 建议在遇到通用 `TokenizersBackend` 时记录错误并崩溃，但 JustinTong0323 认为这过于激进，因为某些模型（如 `Nemotron-H`）能正常工作，并改为警告和回退机制。
- `get_config` 错误处理 (`correctness`): JustinTong0323 修复为捕获 (`ValueError`, `KeyError`) 并匹配字符串。
  - `TokenizersBackend` 处理策略 (`design`): 改为警告和回退机制，保留兼容性。

## 风险与影响

- 风险：技术风险包括：
- 回归风险：transformers 升级可能引入新 bug 或不兼容变更，影响模型加载和推理。
- 兼容性风险：兼容性补丁可能未覆盖所有边缘情况，导致特定模型或配置失败。
- 重构风险：代码拆分可能破坏内部依赖或导入路径，尽管有重新导出 shim，但复杂使用场景可能出错。
- 测试覆盖：新增测试可能未充分覆盖所有平台和模型组合。
- 影响：影响范围广泛：
- 用户：透明升级，但可能遇到新版本的兼容性问题，需要关注日志警告。
- 系统：所有使用 Hugging Face 模型加载的模块都受影响，包括多模态、NPU 支持等。
- 团队：代码结构改善，便于未来维护，但需要熟悉新子包结构。
- 风险标记：依赖版本跳跃，代码重构风险，兼容性补丁维护

## 关联脉络

- PR #22870 [AMD][MoRI] bump MoRI to v1.1.0: 类似依赖升级 PR，涉及更新关键依赖版本，可参考升级策略和兼容性处理。
- PR #22965 migrate CPU-only unit tests from openai\_server to unit/: 涉及代码重构和测试迁移，与本次 PR 中的代码拆分和测试配套有相似之处。