

# PR #23669 完整报告

sgl-project/sglang

Clean up noisy startup warnings from third-party deps

合并时间: 2026-04-27 18:10

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/23669>

## 执行摘要

- 一句话: 清理启动时第三方库噪音警告
- 推荐动作: 值得阅读, 尤其是 `hf_transformers_patches.py` 中的日志压制技巧和 `common.py` 中的 API 迁移方法, 可作为处理第三方库警告和 transformers 升级的参考模式。

## 功能与动机

Server startup logs are cluttered with noisy warnings from third-party dependencies (torchao, transformers v5) that are irrelevant to the user. These changes silence the noise without changing any behavior.

## 实现拆解

1. 压制 torchao 警告: 在 `hf_transformers_patches.py::_patch_removed_symbols` 中, 导入 `modeling_llama` 前临时将 torchao 日志级别设为 ERROR, 避免其版本不匹配警告。
2. 替换弃用配置键: 在 `common.py::get_hf_text_config` 中将 `config.torch_dtype` 替换为 `config.dtype`, 消除 transformers v5 弃用警告。
3. 替换图片处理器基类: 在多模态处理器 `base_processor.py` 和 `ernie45_vl.py` 中将 `BaseImageProcessorFast` 替换为 `BaseImageProcessor`, 适配 transformers v5 移除 Fast 后缀。
4. 降低平台检测日志级别: 在 `platforms/__init__.py` 中将 "No platform detected" 从 warning 降为 debug, 减少非必要输出。
5. 新增诊断文档: 创建 `.claude/skills/clean-startup-log/SKILL.md`, 记录日志清理方法与排查流程。测试方面: PR body 要求验证无警告且多模态模型正常, 但未包含自动化测试用例。

关键文件:

- `python/sglang/srt/utils/hf_transformers_patches.py` (模块 补丁层; 类别 source; 类型 dependency-wiring; 符号 `_patch_removed_symbols`): 核心变更: 压制 torchao 启动警告
- `.claude/skills/clean-startup-log/SKILL.md` (模块 文档; 类别 docs; 类型 documentation; 符号 `_tracing_import, TraceHandler, emit`): 新增技能文档, 记录日志清理方法

- python/sglang/srt/utils/hf\_transformers/common.py (模块 配置层; 类别 source; 类型 core-logic; 符号 get\_hf\_text\_config) : 替换 torch\_dtype 为 dtype 消除弃用警告
- python/sglang/srt/multimodal/processors/base\_processor.py (模块 多模态处理; 类别 source; 类型 dependency-wiring; 符号 process\_mm\_data) : 替换 BaseImageProcessorFast 引用
- python/sglang/srt/multimodal/processors/ernie45\_vl.py (模块 多模态处理; 类别 source ; 类型 dependency-wiring; 符号 process\_mm\_data) : 替换 BaseImageProcessorFast 引用
- python/sglang/srt/multimodal/processors/kimi\_k25.py (模块 多模态处理; 类别 source ; 类型 core-logic) : 更新注释说明基类变更
- python/sglang/srt/platforms/\_\_init\_\_.py (模块 平台检测; 类别 source; 类型 core-logic) : 降低平台检测日志级别

关键符号: \_patch\_removed\_symbols, get\_hf\_text\_config

## 关键源码片段

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

核心变更: 压制 torchao 启动警告

```
def _patch_removed_symbols():
    """Re-export symbols removed in transformers v5.4.0."""
    # LlamaFlashAttention2 compat
    try:
        import logging

        # 导入 modeling_llama 会触发深层导入链:
        # modeling_llama -> quantizers -> torchao,
        # torchao 会发出不兼容 torch 版本的警告,
        # 在这里临时抑制。
        _torchao_logger = logging.getLogger("torchao")
        _prev_level = _torchao_logger.level
        _torchao_logger.setLevel(logging.ERROR)
    try:
        from transformers.models.llama import modeling_llama
    finally:
        _torchao_logger.setLevel(_prev_level)

    if not hasattr(modeling_llama, "LlamaFlashAttention2"):
        if hasattr(modeling_llama, "LlamaAttention"):
            modeling_llama.LlamaFlashAttention2 = modeling_llama.LlamaAttention
except ImportError:
    logger.warning(
        "Could not import transformers.models.llama.modeling_llama; "
        "LlamaFlashAttention2 compat patch not applied."
    )

# is_flash_attn_greater_or_equal_2_10 compat
```

```

try:
    import transformers.utils as _u

    if not hasattr(_u, "is_flash_attn_greater_or_equal_2_10"):
        if hasattr(_u, "is_flash_attn_greater_or_equal"):
            _u.is_flash_attn_greater_or_equal_2_10 = (
                lambda: _u.is_flash_attn_greater_or_equal("2.10.0")
            )
        else:
            _u.is_flash_attn_greater_or_equal_2_10 = lambda: False
except ImportError:
    logger.warning(
        "Could not import transformers.utils; "
        "is_flash_attn_greater_or_equal_2_10 compat patch not applied."
    )

```

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

替换 torch\_dtype 为 dtype 消除弃用警告

```

def get_hf_text_config(config: PretrainedConfig):
    """Get the sub config relevant to LLM for multi modal models."""
    # Skip for custom Llava models
    if config.architectures is not None:
        class_name = config.architectures[0]
        if class_name.startswith("Llava") and class_name.endswith("ForCausalLM"):
            setattr(config, "dtype", torch.float16)
            return config

    text_config = None

    # Convert dict sub-configs to PretrainedConfig
    parent_dtype = getattr(config, "dtype", None) # 使用 dtype 替代已弃用的 torch_dtype
    for _attr in ("text_config", "llm_config", "language_config", "thinker_config"):
        _sub = getattr(config, _attr, None)
        if isinstance(_sub, dict):
            _converted = PretrainedConfig(**_sub)
            if getattr(_converted, "dtype", None) is None and parent_dtype is not None:
                _converted.dtype = parent_dtype # 保证子配置继承父配置的 dtype
                setattr(config, _attr, _converted)

    # Priority: thinker_config > llm_config > language_config > text_config
    if hasattr(config, "thinker_config"):
        thinker_config = config.thinker_config
        if hasattr(thinker_config, "text_config"):
            setattr(
                thinker_config,
                "dtype", # 原为 torch_dtype
                getattr(thinker_config, "dtype", None), # 原为 torch_dtype
            )

```

```
        text_config = thinker_config.text_config
    else:
        text_config = thinker_config
    elif hasattr(config, "llm_config"):
        assert hasattr(config.llm_config, "num_attention_heads")
        text_config = config.llm_config
    elif hasattr(config, "language_config"):
        text_config = config.language_config
    elif hasattr(config, "text_config"):
        assert hasattr(config.text_config, "num_attention_heads")
        text_config = config.text_config

    normalize_rope_scaling_compat(config)

    if text_config is not None:
        return _patch_text_config(config, text_config)
    return config
```

## 评论区精华

本 PR 无 review 讨论，作者直接合并。仅有 bot 提示 quota 已满和 author 的 /tag-and-rerun-ci 命令。

- 暂无高价值评论线程

## 风险与影响

- 风险：风险极低。
- torch\_dtype → dtype: transformers v5 中 torch\_dtype 仍可用但触发弃用警告，替换为 dtype 完全向后兼容。
- BaseImageProcessorFast → BaseImageProcessor: 前者是后者的子类，isinstance 检查父类同样匹配，多模态处理行为不变。
- 日志级别调整不影响任何功能逻辑。
- torchao 日志压制仅作用于特定导入时刻，不影响后续 torchao 功能。
- 影响：对用户：启动日志显著干净，无功能性影响。对系统：无性能回归。对团队：降低 log 噪音，便于问题诊断；新增的 SKILL.md 文档可作为后续日志清理的模板。
- 风险标记：低风险，向后兼容，无行为变更

## 关联脉络

- PR #23525 Upgrade transformers from 5.5.4 to 5.6.0: 本 PR 的 API 替换（torch\_dtype → dtype、BaseImageProcessorFast → BaseImageProcessor）是 transformers 升级的配套清理。