

# PR #24297 完整报告

sgl-project/sglang

Rename SGLANG\_USE\_JIT\_ALL\_REDUCE to SGLANG\_OPT\_USE\_CUSTOM\_ALL\_REDUCE\_V2

合并时间: 2026-05-03 14:59

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/24297>

## 执行摘要

- 一句话: 重命名 JIT all-reduce 环境变量并迁移至 envs 模块
- 推荐动作: 值得精读的点: 了解 SGLang 项目中环境变量的集中管理模式 (Envs 类) 以及弃用处理的惯用模式 (`_print_deprecated_env`)。该 PR 展示了小范围代码整洁重构的标准流程。

## 功能与动机

PR 描述指出: 'Rename the JIT all-reduce opt-in env var and migrate to the envs module for consistency with other recent allreduce flags.' 该变更目的是与其他 all-reduce 标志保持命名和访问方式的一致性。

## 实现拆解

1. 在 `Envs` 类中注册新环境变量: 在 `python/sglang/srt/environ.py` 的 `Env` 类中添加 `SGLANG_OPT_USE_CUSTOM_ALL_REDUCE_V2 = EnvBool(False)` 定义。
2. 更新读取逻辑: 在 `python/sglang/srt/distributed/device_communicators/custom_all_reduce.py` 的 `dispatch_custom_allreduce` 函数中, 将原来通过 `get_bool_env_var("SGLANG_USE_JIT_ALL_REDUCE", default="false")` 读取环境变量的方式改为通过 `envs.SGLANG_OPT_USE_CUSTOM_ALL_REDUCE_V2.get()` 访问。
3. 添加弃用兼容: 在 `_convert_SGL_to_SGLANG` 函数中调用 `_print_deprecated_env("SGLANG_USE_JIT_ALL_REDUCE", "SGLANG_OPT_USE_CUSTOM_ALL_REDUCE_V2")`, 当旧环境变量存在时自动设置新变量并输出弃用警告。
4. 更新文档字符串: 在 `custom_all_reduce.py` 中同步更新了函数注释中的环境变量名。

关键文件:

- `python/sglang/srt/distributed/device_communicators/custom_all_reduce.py` (模块 分布式通信; 类别 `source`; 类型 `core-logic`; 符号 `dispatch_custom_allreduce`): 核心变更文件: 环境变量读取方式从 `get_bool_env_var` 改为通过 `envs` 模块访问, 并更新了文档字符串。
- `python/sglang/srt/environ.py` (模块 环境配置; 类别 `source`; 类型 `core-logic`; 符号 `Env, _convert_SGL_to_SGLANG`): 新增环境变量定义和弃用兼容映射, 确保向后兼容。

关键符号: `dispatch_custom_allreduce`, `_convert_SGL_to_SGLANG`

## 关键源码片段

`python/sglang/srt/distributed/device_communicators/custom_all_reduce.py`

核心变更文件: 环境变量读取方式从 `get_bool_env_var` 改为通过 `envs` 模块访问, 并更新了文档字符串。

```
# python/sglang/srt/distributed/device_communicators/custom_all_reduce.py ( 关键函数 dispatch_
custom_allreduce)
def dispatch_custom_allreduce():
    """Return the CustomAllreduce class to use (aiter on ROCm if enabled).

    On AMD with 1-stage AR enabled, use sglang's CustomAllreduce.
    Otherwise use AiterCustomAllreduce if available.

    Set SGLANG_OPT_USE_CUSTOM_ALL_REDUCE_V2=1 to use the JIT-compiled v2
    implementation.
    """
    # 使用集中式 envs 模块读取新变量名, 替代原来的 get_bool_env_var + 硬编码字符串
    if _is_cuda and envs.SGLANG_OPT_USE_CUSTOM_ALL_REDUCE_V2.get():
        from .custom_all_reduce_v2 import CustomAllReduceV2

        logger.debug("[AR] Using CustomAllReduceV2 (JIT-compiled)")
        return CustomAllReduceV2

    # 以下逻辑保持不变 (NVIDIA CUDA/MUSA / AMD 分支)
    if _is_cuda or _is_musa:
        return CustomAllreduce

    assert _is_hip
    # ... AMD 分支代码略
```

`python/sglang/srt/environ.py`

新增环境变量定义和弃用兼容映射, 确保向后兼容。

```
# python/sglang/srt/environ.py ( 在 Envs 类中 )
class Envs:
    # ... 其他定义 ...
    # Deterministic inference
    SGLANG_ENABLE_DETERMINISTIC_INFERENCE = EnvBool(False)
    SGLANG_USE_1STAGE_ALLREDUCE = EnvBool(False)
    # 新增: JIT compiled v2 custom all-reduce 的开关变量, 默认关闭
    SGLANG_OPT_USE_CUSTOM_ALL_REDUCE_V2 = EnvBool(False)
    # ...

# 函数 _convert_SGL_to_SGLANG 中新增弃用映射
def _convert_SGL_to_SGLANG():
    # 其他映射 ...
```

```
# 添加旧变量到新变量的兼容映射，用户设置旧变量时会自动迁移并收到警告
_print_deprecated_env(
    "SGLANG_USE_JIT_ALL_REDUCE", "SGLANG_OPT_USE_CUSTOM_ALL_REDUCE_V2"
)
# ... 其他映射 ...
```

## 评论区精华

Review 评论 (gemini-code-assist[bot]) 指出: 'The renaming ... is a breaking change for users currently relying on the old environment variable. To maintain backward compatibility ... please add a deprecation mapping.' 该建议已被采纳, 第二个提交正是添加了 `_print_deprecated_env` 调用。

- 添加弃用映射以保持向后兼容 (correctness): 已采纳: 第二个提交中在 `_convert_SGL_to_SGLANG` 函数中添加了 `_print_deprecated_env("SGLANG_USE_JIT_ALL_REDUCE", "SGLANG_OPT_USE_CUSTOM_ALL_REDUCE_V2")`。

## 风险与影响

- 风险: 风险极低。变量名更改本身是向后兼容的 (通过弃用映射), 且改动仅涉及环境变量读取路径和名称, 不涉及 all-reduce 算法逻辑本身。经过审查后添加的弃用映射进一步降低了用户升级时的断裂风险。
- 影响: 影响范围小, 仅涉及两个文件。用户如果之前设置了 `SGLANG_USE_JIT_ALL_REDUCE=1`, 在升级后会收到弃用警告但功能仍正常工作。推荐用户迁移到新变量名。
- 风险标记: 暂无

## 关联脉络

- 暂无明显关联 PR