

# PR #22808 完整报告

sgl-project/sclang

[NPU] qwen3next low latency best practice docs.

合并时间: 2026-04-14 21:21

原文链接: <http://prhub.com.cn/sgl-project/sclang/pull/22808>

## 执行摘要

本次 PR 在 SGLang 的 Ascend NPU 最佳实践文档中添加了 Qwen3-Next 模型的低延迟配置方案，涵盖两个特定输入输出长度场景，旨在帮助用户优化 NPU 部署以降低延迟。变更纯属文档更新，风险较低，影响范围针对 NPU 平台用户，是 NPU 文档维护工作流程的一部分。

## 功能与动机

动机源于补充 Qwen3-Next 模型在 Ascend NPU 上的低延迟最佳实践文档，以支持用户在该平台上实现高效推理部署。PR body 中明确表述为“add qwen3next low latency best practice docs.”，反映出对 NPU 生态文档完善的持续需求。

## 实现拆解

实现仅修改一个文件: [docs/platforms/ascend/ascend\\_npu\\_best\\_practice.md](#)。关键改动包括:

- 在性能汇总表格中新增两行，分别对应 Qwen3-Next 模型在 1K+0.3K 和 6K+1.5K 输入输出长度下的 TPOT (时间每输出令牌) 数据。
- 添加两个配置块，详细说明部署命令和环境变量设置，例如: `shell export SGLANG_ENABLE_SPEC_V2=1 export DEEPEP_NORMAL_LONG_SEQ_ROUND=5`
- 提供性能测试方法和硬件信息，确保用户可复现结果。

## 评论区精华

Review 过程中无评论或讨论，仅由 sclang-npu-bot 自动批准，表明变更被视为常规文档更新，无需技术交锋。

## 风险与影响

- 风险: 主要风险是文档准确性，如配置参数错误可能导致用户部署失败；文档未经过代码级测试，依赖作者经验。
- 影响: 直接影响 Ascend NPU 用户，通过提供优化配置，潜在提升推理性能 (TPOT 14.21ms 至 15.62ms)；对系统无代码变更，属于文档增强。

## 关联脉络

与近期多个 NPU 文档 PR (如 #22804、#22799、#22795) 关联, 共同构成 NPU 平台文档维护序列。这些 PR 均聚焦于 Ascend NPU 功能描述和最佳实践更新, 显示团队对该平台文档的系统性完善, 以支持日益增长的 NPU 部署需求。