

# PR #23486 完整报告

sgl-project/sglang

docs(cookbook): add Qwen3.6-27B dense variant

合并时间: 2026-04-23 01:22

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/23486>

## PR 分析报告: docs(cookbook): add Qwen3.6-27B dense variant

### 执行摘要

本次 PR 为 Qwen3.6 模型系列新增的 27B 密集变体更新了文档和部署工具，包括 cookbook 页面重写和交互式代码片段扩展，使用户能便捷地部署和使用新变体。变更主要涉及文档内容，技术风险低，但需注意引用的测试数据与官方结果的差异。

### 功能与动机

Qwen3.6 近期发布了 27B 密集变体 (Qwen3.6-27B / Qwen3.6-27B-FP8)，与现有的 35B-A3B MoE 变体并存。为保持文档与模型发布同步，需要更新相关材料以覆盖这一新变体。PR 作者在 body 中明确说明：“更新 cookbook 页面和部署代码片段以覆盖两者”，并“重写介绍 / 可用模型 / 硬件要求章节以覆盖两个变体”。这确保了用户能获得准确的部署指导，提升文档的实用性和时效性。

### 实现拆解

- 部署代码片段改造 (`docs_new/src/snippets/autoregressive/qwen36-deployment.jsx`)
  - 在配置对象 `options` 中添加 `modelSize` 字段，提供“35B-A3B (MoE)”和“27B (Dense)”两个单选选项。
  - 重构 `modelConfigs` 对象，从按硬件平铺改为按模型大小嵌套，每个变体包含 `baseName` (用于路径生成) 和硬件 / 量化配置。
  - 更新 `generateCommand` 函数，读取用户选择的 `modelSize`，动态生成对应的 `--model-path` 参数 (例如 `Qwen/Qwen3.6-27B`)。
  - 在遍历选项生成命令行时，跳过 `modelSize` 字段，避免其产生无效参数。

关键实现片段展示了配置结构和命令生成逻辑：

```
const modelConfigs = { '35b-a3b': {
  baseName: '35B-A3B', h100: { bf16: { tp: 1, mem: 0.8 }, fp8: { tp: 1, mem: 0.8 } }, // ... 其他硬件配置
}, '27b': { baseName: '27B', h100: { bf16: { tp: 1, mem: 0.8 }, fp8: { tp: 1, mem: 0.8 } }, // ... 其他硬件配置
}, };
const generateCommand = () => { const { hardware, modelSize, quantization, speculative } = values;
const sizeConfig = modelConfigs[modelSize];
const hwConfig = sizeConfig?.[hardware]?.[quantization];
const modelName = `Qwen/Qwen3.6-${sizeConfig.baseName}${quantSuffix}`; // ... 生成完整命令
};
```

- cookbook 文档更新 (`docs_new/cookbook/autoregressive/Qwen/Qwen3.6.mdx`) - 更新元描述，从单一变体扩展为系列描述。
  - 重写模型介绍，明确列出两个变体及其特点 (稀疏

MoE vs 密集)。 - 调整“关键特性”部分，反映架构的多样性。 - 在“可用模型”和“硬件要求”表格中添加 27B 变体的信息。

2. 清理过时内容 - 移除了 `sglang[all]` 安装提示，替换为通用的 `uv pip install sglang`，避免引入不相关的扩散 / 追踪依赖，与安装文档保持一致。

## [docs\\_new/src/snippets/autoregressive/qwen36-deployment.jsx](#)

这是交互式部署代码片段的核心文件，负责生成用户部署命令。本次变更添加了模型大小选择逻辑，直接影响用户部署体验。

### 关键源码片段

## [docs\\_new/src/snippets/autoregressive/qwen36-deployment.jsx](#)

这是交互式部署代码片段的核心文件，负责生成用户部署命令。本次变更添加了模型大小选择逻辑，直接影响用户部署体验。

```
export const Qwen36Deployment = () => {
  const options = {
    // ... 其他配置 (如 hardware) 保持不变
    modelSize: {
      name: 'modelSize',
      title: 'Model Size',
      items: [
        { id: '35b-a3b', label: '35B-A3B (MoE)', default: true }, // 默认选中 MoE 变体
        { id: '27b', label: '27B (Dense)', default: false }, // 新增 27B 密集变体选项
      ],
    },
    // ... 其他配置 (如 quantization、reasoning 等)
  };

  // 模型配置现在按 modelSize 嵌套，每个变体包含 baseName 和硬件 / 量化配置
  const modelConfigs = {
    '35b-a3b': {
      baseName: '35B-A3B', // 用于生成模型路径的基础名称
      h100: { bf16: { tp: 1, mem: 0.8 }, fp8: { tp: 1, mem: 0.8 } },
      h200: { bf16: { tp: 1, mem: 0.8 }, fp8: { tp: 1, mem: 0.8 } },
      b200: { bf16: { tp: 1, mem: 0.8 }, fp8: { tp: 1, mem: 0.8 } },
    },
    '27b': {
      baseName: '27B', // 27B 变体的基础名称
      h100: { bf16: { tp: 1, mem: 0.8 }, fp8: { tp: 1, mem: 0.8 } },
      h200: { bf16: { tp: 1, mem: 0.8 }, fp8: { tp: 1, mem: 0.8 } },
      b200: { bf16: { tp: 1, mem: 0.8 }, fp8: { tp: 1, mem: 0.8 } },
    },
  };

  const generateCommand = () => {
    const { hardware, modelSize, quantization, speculative } = values;
    const sizeConfig = modelConfigs[modelSize]; // 根据选择的 modelSize 获取配置
```

```
const hwConfig = sizeConfig?.[hardware]?.[quantization]; // 嵌套获取硬件和量化配置
if (!hwConfig) {
  return '# Please select a valid hardware and quantization combination';
}

const quantSuffix = quantization === 'fp8' ? '-FP8' : '';
const modelName = `Qwen/Qwen3.6-${sizeConfig.baseName}${quantSuffix}`; //
// 动态生成模型路径
// ... 后续命令生成逻辑
};
};
```

## 评论区精华

本次 PR 没有 review 评论，直接合并。但在关联的 Issue #23467 中，有评论者提出疑问：

“MMM U 结果与官方结果 (82.9) 差异很大。这个测试是在 sglang 主分支上进行的吗？”

这暗示文档中引用的测试数据可能需要进一步验证，但该讨论未在 PR 中展开，也未影响合并决策。PR 作者在测试计划中已说明差异在统计置信区间内，且 FP8 变体经 #23467 修复后达到 BF16 同等精度。

## 风险与影响

- 风险：主要风险在于文档内容的准确性。MMM U 测试结果 (55.1% 和 53.0%) 与官方数据 (82.9%) 存在显著差异，可能误导用户对模型性能的预期。不过，作者已解释该差异在 Wilson 95% CI 范围内，且 FP8 变体无精度回归。部署代码片段的新增逻辑简单，配置错误风险可控。
- 影响：正面影响显著。用户现在能获得完整的 Qwen3.6 变体信息，并使用交互式工具生成正确的部署命令，提升了文档的实用性和用户体验。对系统无运行时影响，不涉及核心代码变更。

## 关联脉络

- 与 PR #23467 紧密相关：该 PR 修复了 FP8 量化配置中的模块路径匹配错误，确保了 Qwen3.6-27B-FP8 变体的正确加载。本次文档更新基于该修复进行验证，两者共同支持新变体的端到端部署流程。
- 从近期历史 PR 看，文档更新类变更（如 #23459 更新 NPU 最佳实践）常见于模型或硬件支持扩展后，反映了项目对文档及时性的重视。