

# PR #7278 完整报告

PaddlePaddle/FastDeploy

[Docs]add dsk-3.2 doc

合并时间: 2026-04-09 17:28

原文链接: <http://prhub.com.cn/PaddlePaddle/FastDeploy/pull/7278>

## 执行摘要

本次 PR 为 FastDeploy 添加了 DeepSeek-V3.2 模型的部署文档，覆盖中英文版本，包括环境准备和多个部署示例。通过修正拼写错误和变量名不一致，提升了文档的准确性和可执行性。该变更对用户部署新模型有直接帮助，风险较低，属于常规文档维护。

## 功能与动机

动机是响应用户对 DeepSeek-V3.2 模型部署指南的需求。尽管 PR body 未填写具体描述，但标题“add dsk-3.2 doc”和文件变更表明，目标是为这一新模型提供详细的部署说明，确保用户能快速上手使用 FastDeploy 进行部署。

## 实现拆解

实现涉及两个关键文件：

- docs/best\_practices/DeepSeek-V3.md：新增 DeepSeek-V3.2 章节，内容包括：
  - 硬件支持：H800 80GB GPU 在 block\_wise\_fp8 量化下需 16 卡。
  - 安装指南：指向现有 GPU 安装文档。
  - 部署示例：三个示例展示不同配置，如使用 DSA\_ATTN 后端、多服务器设置等。

```
shell export FD_ATTENTION_BACKEND="DSA_ATTN" python -m fastdeploy.entrypoints.openai.multi_api_server \ --quantization block_wise_fp8 \ --data-parallel-size 16
```
- docs/zh/best\_practices/DeepSeek-V3.md：同步中文内容，结构相同，确保中文用户获得一致信息。

## 评论区精华

review 中，fastdeploy-bot 自动化工具指出了文档中的关键问题：

- 拼写错误：如将 'block\_wise\_fp8' 误写为 'black\_wise\_fp8'，fastdeploy-bot 评论道“拼写错误会误导用户”。
- 变量名不一致：示例中变量名大小写混合（如 \$MODEL\_PATH 与 \$model\_path），fastdeploy-bot 建议“修改为一致的大写变量名，避免命令执行失败”。这些反馈被及时采纳并在提交中修正，无进一步争议，体现了自动化工具在文档质量保障中的作用。

## 风险与影响

风险：主要风险是文档准确性，如拼写错误可能导致用户误解量化类型；变量名不一致可能使示例命令失败。但这些已通过 review 修正，风险缓解。无代码变更，故无技术回归、性能或安全风险。影响：对用户影响积极，提供了新模型的部署指南，降低学习成本。对系统无直接影响。对团队，完善了文档覆盖，提升了整体用户体验。

## 关联脉络

与近期 PR #7267 (“[Docs] Update docs for release/2.5”) 关联，两者都属于文档维护工作，反映了团队对文档更新的持续投入。从历史 PR 看，FastDeploy 近期侧重于模型支持优化（如 MoE、Attention 模块）和文档完善，本 PR 是这一趋势的延续，为新增模型提供配套文档。