

# PR #14385 完整报告

sgl-project/sglang

[CPU] Implement MXFP4 Gemm kernels for intel AMX to support GPT OSS series.

合并时间: 2026-03-30 14:44

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/14385>

## 执行摘要

本 PR 为 sglang 仓库的 CPU 后端实现了 MXFP4 量化 GEMM 内核，针对 Intel AMX 和 AVX512 优化，旨在支持 GPT OSS 系列大型语言模型推理。通过扩展类型 dispatch、添加权重打包 / 解包逻辑和提供 tinygemm 接口，显著增强了 quantization 能力，但需关注 AVX512 代码的正确性和性能测试覆盖。

## 功能与动机

动机源自支持 GPT OSS 20B 和 120B 等模型的 MXFP4 量化推理需求。PR body 明确说明 'provide mxfp4 support for intel amx backend'，Issue 评论中进一步指出要 'enable gpt oss 20B and 120B'。这反映了项目在扩展低精度推理支持以优化资源使用和性能。

## 实现拆解

实现主要涉及以下模块：

- 类型扩展：修改 sgl-kernel/csrc/cpu/common.h，在 CPU\_DISPATCH\_PACKED\_TYPES 宏中添加 uint8\_t 支持，以处理 mxfp4/int4 类型。
- 权重处理：在 sgl-kernel/csrc/cpu/gemm.cpp 中，新增 pack\_vnni<uint8\_t> 函数实现 32-way VNNI 格式打包，并更新 convert\_weight\_packed 函数以适配 kByte 类型。
- 内核优化：在 sgl-kernel/csrc/cpu/gemm\_fp8.cpp 中，添加 unpack\_B 函数用于 mxfp4 解包，使用 AVX512 intrinsics 将 4-bit 值转换为 bfloat16，并扩展 tinygemm\_kernel 模板。
- 数值转换：在 sgl-kernel/csrc/cpu/vec.h 中引入 CVT\_MXFP4\_TO\_BF16 宏，基于查找表实现高效 MXFP4 到 bfloat16 转换。
- API 集成：在 sgl-kernel/csrc/cpu/torch\_extension\_cpu.cpp 中暴露 mxfp4\_scaled\_mm\_cpu 和 convert\_scale\_packed 函数，提供上层调用接口。

关键代码片段示例（来自 gemm\_fp8.cpp）：  
`inline void unpack_B(at::BFloat16* __restrict__ Btmp, const uint8_t* __restrict__ packed_B, int64_t N, int64_t K, int64_t ldb, int64_t ldb_tmp, const uint8_t* __restrict__ scale){ // AVX512 实现，使用预取和 intrinsics 优化`

## 评论区精华

没有正式的 review 讨论，但 Issue 评论中作者 mingfeima 提到：

"provide a mxfp4 moe kernel based on tinygemm interface and also other necessary changes to enable gpt oss 20B and 120B." 这表明实现已直接集成 MoE 支持，无公开争议，侧重于功能交付。

## 风险与影响

风险：

- AVX512 代码复杂性可能引入数值错误，尤其是在边界条件下。
- uint8\_t 类型同时表示 mxfp4 和 int4，使用中需清晰文档以避免混淆。
- 缺少性能基准测试，仅提及 accuracy tests，可能掩盖潜在性能回归。

影响：

- 用户可受益于 MXFP4 量化，降低 GPT OSS 模型的内存占用，提升 CPU 推理速度。
- 系统 quantization 栈得到扩展，为未来低精度优化（如 int4）铺平道路。
- 团队需维护新内核，可能增加代码库复杂性和测试需求。

## 关联脉络

与历史 PR 的关联揭示 quantization 优化脉络：

- PR #21625 涉及 MXFP8 Gemm 测试改进，共享 quantization 测试模式。
- PR #19835 修复 MXFP8 Triton 路径，反映在低精度量化方面的持续投入。这些 PR 共同显示仓库正积极扩展 quantization 支持，以适配多样硬件和模型需求，本 PR 是这一趋势在 CPU 和 MXFP4 领域的具体体现。