

# PR #26026 完整报告

sgl-project/sglang

[bug fix] Fix 3 issues when using Gemma4 MTP

合并时间: 2026-05-23 18:16

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/26026>

## 执行摘要

- 一句话: 修复 Gemma4 MTP 三个初始化崩溃问题
- 推荐动作: 此 PR 值得关注, 它展示了处理模型初始化兼容性问题的典型模式:
  - 当子类跳过了父类的 `__init__` 时, 需要显式设置父类依赖的成员变量。
  - 对于 MoE 模型的 Dense 变体, 需要安全处理 `num_experts` 缺失的情况。
  - 硬件特定的自动后端选择应结合量化类型共同判断, 避免对不支持的后端进行硬编码。建议其他模型后端自动选择逻辑也参考此模式。

## 功能与动机

Gemma4 MTP assistant 模型在初始化时崩溃, 无法正常启动服务。PR body 中引用了两个堆栈跟踪:

- `AttributeError: 'Gemma4AssistantForCausalLM' object has no attribute 'pp_group'` —— 因为 MTP 类跳过了 `Gemma4CausalLM.__init__`, 未设置 `pp_group`。
- `TypeError: 'NoneType' object cannot be interpreted as an integer` —— 因为 `num_experts` 为 `None`, 但仍调用 `FusedMoE.make_expert_params_mapping`。
- 第三个问题: BF16 checkpoint 不应使用 `flashinfer` 作为 MoE runner 后端, 但当前代码对所有 `moe_runner_backend == "auto"` 的情况都选择了 `flashinfer_trtllm`。

## 实现拆解

1. 修复 MTP 类缺失 `pp_group` 属性(`python/sglang/srt/models/gemma4_mtp.py`):
  - 新增导入 `from sglang.srt.distributed import get_pp_group`。
  - 在 `Gemma4AssistantForCausalLM.__init__` 中显式调用 `self.pp_group = get_pp_group()`, 确保 `pp_group` 在模型初始化时被正确设置, 避免后续 `tie_weights` 等操作中因父类 `__init__` 被跳过而引发的 `AttributeError`。
2. 修复 Dense 子类中 `num_experts` 为 `None` 的问题(`python/sglang/srt/models/gemma4_causal.py`):
  - 将 `num_experts = self.config.num_experts` 改为 `num_experts = getattr(self.config, "num_experts", None) or 0`, 安全处理 MoE 属性缺失的情况。
  - 将 `per_expert_params_mapping` 的构造包裹在条件判断 `if num_experts: ... else []` 中, 当 `num_experts` 为 0 时返回空列表, 避免调用 `FusedMoE.make_expert_params_m`

apping 时因 None 传参而崩溃。

### 3. 修复 BF16 下 MoE runner 后端的自动选择(`python/sglang/srt/server_args.py`):

- 在 SM100 的 MoE 后端自动选择逻辑中，增加 `if self.get_model_config().quantization == "modelopt_fp4"`: 条件判断，只有当量化配置为 `modelopt_fp4` 时，才设置 `self.moe_runner_backend = "flashinfer_trtllm"` 并更新 `self.quantization = "modelopt_fp4"`。对于非 FP4 的其他量化类型（如 BF16），保持默认 MoE runner 后端，避免使用不支持的 flashinfer 后端。

### 4. 测试与验证：通过运行 #24552 的测试验证修复效果，PR 提交前服务器在初始化阶段崩溃，修复后服务可正常启动。本次 PR 未新增测试文件，但后续 CI 覆盖了 Gemma4 MTP 相关的回归场景。

关键文件：

- `python/sglang/srt/models/gemma4_causal.py`（模块 模型层；类别 source；类型 core-logic；符号 `load_weights`）：核心修复：使用 `getattr` 安全获取 `num_experts`，并条件构造 `per_expert_params_mapping`，避免 Dense 子类因 `num_experts=None` 崩溃。
- `python/sglang/srt/server_args.py`（模块 服务配置；类别 source；类型 core-logic）：修复 Gemma4 在 SM100 上使用 BF16 时 MoE 后端自动选择错误的问题，增加了量化类型判断。
- `python/sglang/srt/models/gemma4_mtp.py`（模块 模型层；类别 source；类型 entrypoint；符号 `Gemma4AssistantForCausalLM.init`）：修复 MTP 类因跳过父类 `__init__` 而缺失 `pp_group` 属性的问题。

关键符号：`Gemma4AssistantForCausalLM.init`, `Gemma4ForCausalLM.load_weights`

## 关键源码片段

### `python/sglang/srt/models/gemma4_causal.py`

核心修复：使用 `getattr` 安全获取 `num_experts`，并条件构造 `per_expert_params_mapping`，避免 Dense 子类因 `num_experts=None` 崩溃。

```
# 关键变更：安全处理 num_experts
# Dense 子类（如 Gemma4 MTP assistant）复用此 load_weights
# 但 num_experts 可能为 None 或不存在，因此使用 getattr 默认 None
# 再通过 or 0 确保后续 make_expert_params_mapping 不会收到 None
num_experts = getattr(self.config, "num_experts", None) or 0

# 只有当 num_experts 非零（即模型确实是 MoE）时才生成 per-expert mapping
# 对于 Dense 子类，直接返回空列表，避免 TypeError
per_expert_params_mapping = (
    FusedMoE.make_expert_params_mapping(
        ckpt_gate_proj_name="gate_proj",
        ckpt_down_proj_name="down_proj",
        ckpt_up_proj_name="up_proj",
        num_experts=num_experts,
    )
    if num_experts
```

```
    else []
)
```

### python/sglang/srt/server\_args.py

修复 Gemma4 在 SM100 上使用 BF16 时 MoE 后端自动选择错误的问题，增加了量化类型判断。

```
# 关键变更：仅在 modelopt_fp4 量化时选择 flashinfer_trtllm
# 当前版本 flashinfer 不支持 bf16 checkpoint，因此自动选择不能硬编码
if is_sm100_supported() and self.moe_runner_backend == "auto":
    # 必须先检查量化类型，避免 BF16 模型错误使用 flashinfer
    if self.get_model_config().quantization == "modelopt_fp4":
        self.quantization = "modelopt_fp4"
        self.moe_runner_backend = "flashinfer_trtllm"
        logger.info(
            "Use flashinfer_trtllm as MoE runner backend on "
            "SM100 for Gemma-4 (modelopt_fp4)"
        )
```

### python/sglang/srt/models/gemma4\_mtp.py

修复 MTP 类因跳过父类 `__init__` 而缺失 `pp_group` 属性的问题。

```
class Gemma4AssistantForCausalLM(Gemma4ForCausalLM):
    """Gemma 4 MTP assistant: target embed + recurrent hidden through pre/post projection;
    own lm_head."""

    def __init__(
        self,
        config: PretrainedConfig,
        quant_config: Optional[QuantizationConfig] = None,
        prefix: str = "",
    ) -> None:
        text_config = copy.deepcopy(_get_text_config(config))
        text_config.num_kv_shared_layers = 0
        PreTrainedModel.__init__(self, config=text_config) # 跳过 Gemma4ForCausalLM.__init__
        self.assistant_config = config
        self.config = text_config
        self.quant_config = quant_config
        # 显式设置 pp_group，因为父类 __init__ 被跳过
        # 否则后续 tie_weights 中访问 self.pp_group.world_size 会抛出 AttributeError
        self.pp_group = get_pp_group()
        # ... 其余初始化不变
```

## 评论区精华

仅有一条审核评论：[kpham-sgl](#) 审阅者批准了该 PR，表示 "LGTM, thanks for the bug fix"。没有其他讨论或争议。

- 暂无高价值评论线程

## 风险与影响

- 风险:

1. 回归风险（低）： `gemma4_causal.py` 中的 `load_weights` 改为条件调用 `make_expert_params_mapping`，对原本使用 MoE 的 Gemma4 模型没有影响，因为 `num_experts` 仍然来自 `config.num_experts` 且非零。但若未来有新的 Dense 子类也继承此方法，需确保其 `num_experts` 属性正确。
2. 兼容性风险（低）： `server_args.py` 中的修改增加了 `quantization` 判断，仅影响 SM100 平台上的 `auto` 模式选择。对于非 SM100 或其他量化类型无影响。
3. 测试覆盖（中等）： PR 描述中提到通过运行 #24552 的测试验证，但未新增独立的单元测试。后续需要关注 CI 中 Gemma4 MTP 相关的回归测试是否已覆盖这些场景。

- 影响:

1. 用户影响（高）： 修复了 Gemma4 MTP 模型无法启动的严重 bug，所有尝试使用 Gemma4 MTP 的用户都将受益。
2. 系统影响（低）： 改动集中在模型初始化和配置选择逻辑，不涉及运行时推理路径，对已正常运行的 Gemma4 模型无影响。
3. 团队影响（低）： 改动较小，仅涉及 3 个文件，代码审查简单，已获得批准。 - 风险标记： 缺少测试覆盖

## 关联脉络

- PR #24552 test Gemma4 MTP: PR body 中提到此修复通过运行 #24552 的测试验证，是触发该 PR 的测试用例。