

PR #24511 完整报告

sgl-project/sglang

Support per-call extras and dataclass transform input in dumper grafter

合并时间: 2026-05-06 16:57

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/24511>

执行摘要

- 一句话: dumper grafter 支持 per-call extras 和 dataclass 输入
- 推荐动作: 建议开发者精读此 PR, 特别是 GraftTransformInput 的设计模式和 per-call extras 的实现思路, 可作为内部调试工具 API 设计的参考。合并前应修复 review 中提出的进程组 rank 问题, 否则可能导致静默错误。

功能与动机

原用户变换函数签名(received_list,target)扩展困难, 每次新增字段都需要修改所有用户代码。改用 dataclass 模式可向后兼容地添加 fields。同时 per-call extras 机制让用户能在单次 dump 中传递上下文信息到变换函数, 提高灵活性。

实现拆解

1. 在 dumper.py 中定义 GraftTransformInput dataclass, 包含 tags、received_list、received_extras_list 和 target 字段。
2. 修改 maybe_intercept 方法, 新增 extras 参数, 将 all-gather 的贡献改为 (value, extras) 元组, 接收方构建 received_extras_list 并传入 GraftTransformInput。
3. 修改 dump 和 _dump_inner 方法, 透传 grafter_extras 参数。
4. 修改用户变换函数的调用逻辑, 将原来展开的参数替换为 GraftTransformInput 实例。
5. 更新测试用例: 增加 test_extras_flow_to_recv_transform 和 test_extras_default_none_flow, 同时调整现有测试以适应新的函数签名。

关键文件:

- python/sglang/srt/debug_utils/dumper.py (模块 调试工具; 类别 source; 类型 core-logic; 符号 GraftTransformInput, transform, maybe_intercept, _default_transform): 核心逻辑实现, 包含 GraftTransformInput dataclass 定义、maybe_intercept 方法修改以及 per-call extras 和 dataclass 输入的支持。
- test/registered/debug_utils/test_dumper.py (模块 测试; 类别 test; 类型 test-coverage; 符号 test_extras_flow_to_recv_transform, _test_extras_func, test_extras_default_none_flow, _test_extras_none_func): 新增测试覆盖 extras 传递和默认 None 行为, 同时更新现有测试以适应新的变换函数签名。

关键符号: maybe_intercept, _default_transform, dump, _dump_inner, GraftTransformInput

关键源码片段

python/sglang/srt/debug_utils/dumper.py

核心逻辑实现, 包含 GraftTransformInput dataclass 定义、maybe_intercept 方法修改以及 per-call extras 和 dataclass 输入的支持。

```
# python/sglang/srt/debug_utils/dumper.py
```

```
@dataclass
```

```
class GraftTransformInput:
```

```
    """用户变换函数的单一输入参数。
```

```
    User transforms have signature:
```

```
    def transform(graft_input: GraftTransformInput) -> torch.Tensor: ...
```

```
    使用 dataclass 封装后, 未来可安全新增字段而不破坏已有用户代码。
```

```
    """
```

```
    # Full dumper.dump tags dict (name + recompute_status + extra_kwargs + ctx).
```

```
    tags: "dict[str, Any]"
```

```
    # One tensor per sender rank, in sender-rank order.
```

```
    received_list: "list[torch.Tensor]"
```

```
    # Parallel list of per-sender `grafter_extras` (the dict passed to
```

```
    # dumper.dump on each sender; None if the sender omitted it).
```

```
    received_extras_list: "list[Optional[dict]]"
```

```
    # Recv side's local tensor that will be copy_'d into.
```

```
    target: "torch.Tensor"
```

```
class _Grafter:
```

```
    # ...
```

```
    def maybe_intercept(
```

```
        self,
```

```
        *,
```

```
        value,
```

```
        tags: dict,
```

```
        extras: Optional[dict] = None, # 新增: per-call extension dict
```

```
    ) -> None:
```

```
        # ... 原有逻辑
```

```
        # all-gather: sender 贡献 (value, extras) 元组, recv 为 None
```

```
        my_contribution = (value, extras) if is_send else None
```

```
        gathered: list = [None] * total_world
```

```
        dist.all_gather_object(gathered, my_contribution, group=self._pg)
```

```

if not is_send:
    # recv 端: 拆分 gathered 为 value_list 和 extras_list
    value_list = [item[0] for item in gathered if item is not None]
    extras_list = [item[1] for item in gathered if item is not None]
    # 构造 GraftTransformInput 并调用 transform
    graft_input = GraftTransformInput(
        tags=tags,
        received_list=value_list,
        received_extras_list=extras_list,
        target=value,
    )
    override = transform(graft_input)
    # ... 执行 copy_

```

test/registered/debug_utils/test_dumper.py

新增测试覆盖 extras 传递和默认 None 行为，同时更新现有测试以适应新的变换函数签名。

```
# test/registered/debug_utils/test_dumper.py
```

```

def test_extras_flow_to_recv_transform(self, tmp_path: Path):
    """发送方携带 grafter_extras, 接收方变换函数读取并据此生成覆盖值。"""
    module_name = "_xform_uses_extras"
    (tmp_path / f"{module_name}.py").write_text(
        "import torch\n"
        "def transform(graft_input):\n"
        "    fill = graft_input.received_extras_list[0]['fill_value']\n"
        "    return torch.full_like(graft_input.target, fill)\n"
    )
    graft_port = find_available_port(29645)
    _run_graft_test(
        self._test_extras_func,
        graft_port=graft_port,
        group_name="grafter_extras",
        transform_dir=str(tmp_path),
        transform_path=f"{module_name}.transform",
    )

    @staticmethod
    def _test_extras_func(rank, graft_port, group_name, transform_dir, transform_path):
        sys.path.insert(0, transform_dir)
        grafter = _Grafter(
            config=_make_grafter_test_config(
                rank=rank,
                graft_port=graft_port,
                group_name=group_name,
                transform_path=transform_path,
            )
        )
    try:

```

```

if rank == 0:
    tensor = torch.tensor([1.0, 2.0, 3.0], device="cuda:0")
    grafter.maybe_intercept(
        value=tensor,
        tags={"name": "x"},
        extras={"fill_value": 42.0}, # 发送方传递 extras
    )
else:
    target = torch.zeros(3, device="cuda:1")
    grafter.maybe_intercept(value=target, tags={"name": "x"})
    # 接收方变换函数返回 [42, 42, 42]
    assert target.tolist() == [42.0, 42.0, 42.0], target.tolist()
finally:
    if grafter._pg is not None:
        dist.destroy_process_group(grafter._pg)

```

评论区精华

review 中 gemini-code-assist[bot] 指出：

- `_default_transform` 方法使用了默认进程组的 `dist.get_rank()` 而非 `graft` 进程组的 `self._pg`，可能导致 `rank` 和 `world size` 错误（高优先级）。
- 使用内部 PyTorch 分布式 API（如 `_new_process_group_helper`）存在维护风险，建议改用公共 API 或增加版本适配（中优先级）。两个问题均未在合并前修复。
- `_default_transform` 使用了默认进程组的 `rank`，而不是 `graft` 进程组的 `rank` (`correctness`): 未在合并前修复，需要作者后续处理。
- 使用内部 PyTorch 分布式 API 的维护风险 (`other`): 未在合并前修复，建议作者改用公共 API 或增加版本兼容层。

风险与影响

- 风险：
 1. 兼容性风险：用户变换函数签名从 `(received_list, target)` 改为 `(graft_input)`，所有现有用户代码必须更新，否则会立即报错。
 2. 正确性风险：`_default_transform` 仍使用默认进程组的 `rank`，当 `graft` 进程组与默认进程组不一致时，通信索引会错位，导致错误的的数据拷贝。
 3. 维护性风险：依赖内部 PyTorch API 可能在版本升级时失效，需持续跟踪。- 影响：影响范围仅限于使用 `dumper grafter` 功能的开发者（调试 / 对比场景）。变更会破坏现有用户变换脚本，但考虑到 `grafter` 仍是较新的实验性功能，用户基数小。新功能支持更灵活参数传递，为后续扩展提供了更干净的基础。- 风险标记：API 变更破坏兼容性，分布式通信索引可能错误，内部 API 依赖维护风险

关联脉络

- PR #24512 Enhance diff and tensor-info logging in dumper grafter: 同一功能线（`dumper grafter`），两个 PR 先后增强了 `dumper` 的日志和变换能力。

- PR #24513 Add e2e test with log snapshot in dumper grafter: 为 dumper grafter 添加端到端测试，与本 PR 的测试增强互补。