

集合通信算法研究分析报告

summer 2025

WizHUA 
NUDT 计算机学院

最初写作于：2025 年 07 月 06 日

最后更新于：2025 年 07 月 09 日

目录

1. 摘要 3

2. 项目概述与研究背景 3

3. Open MPI 集合通信算法源码分析 3

 3.1. Open MPI 架构概述 3

 3.2. 集合通信框架分析 3

 3.2.1. 框架核心文件结构 4

 3.2.2. 核心文件分析 4

 3.3. 集合通信操作示例 5

 3.3.1. 用户代码调用 5

 3.3.2. 初始调用暂记 6

 3.3.3. 组件选择的核心机制 8

 3.3.4. 初始化阶段：组件选择机制 9

 3.3.5. MPI_Reduce 的调用过程 10

 3.3.6. 具体执行过程 12

 3.3.7. 调用链总结 12

 3.4. 主要集合通信算法实现 13

4. 集合通信参数配置分析 13

5. 数据集构建与应用场景设计 13

6. 机器学习建模与性能预测 13

7. 模型验证与检验 13

8. 结果分析与总结 13

9. 结论与展望 13

参考文献 14

章节 1. 摘要

章节 2. 项目概述与研究背景

章节 3. Open MPI 集合通信算法源码分析

基于 [The Open MPI Project](#)。 ([仓库链接](#))

3.1. Open MPI 架构概述

Open MPI 作为高性能计算领域的主流 MPI 实现，采用了分层的模块化组件架构(Modular Component Architecture, MCA)设计。整个架构分为三个核心层次：

- **OPAL**(Open Portability Access Layer) 提供操作系统和硬件抽象
- **OMPI** 实现 MPI 标准的核心功能
- **OSHMEM** 提供 OpenSHMEM 支持

这种分层设计确保了代码的模块化和可移植性，并且为集合通信算法的实现和优化提供了灵活的框架支持。

集合通信的核心实现位于 `ompi/mca/coll/` 目录下，包含了 `base`、`basic`、`tuned`、`han`、`xhc` 等多个专门化组件。这种组件化设计通过 `mca_coll_base_comm_select()` 机制，能够根据消息大小、进程数量、网络拓扑等运行时参数动态选择最优算法。

在此次任务中，我们专注于研究消息大小和进程数量这两个核心因素对集合通信算法选择和性能的影响，暂不涉及网络拓扑、硬件特性等复杂环境因素。因此，我们在源码分析部分将重点关注以下三个核心组件：

- `base` 组件 - 提供基础算法实现和算法选择框架
- `basic` 组件 - 包含简单可靠的参考算法实现
- `tuned` 组件 - 集成多种优化算法和智能选择机制

MCA 架构的另一个关键特性是其参数化配置系统，通过 MCA 参数可以在运行时动态调整算法选择策略、消息分片大小、通信拓扑等关键参数，同时 `MPI_T` 接口提供了运行时性能监控和参数调优的能力。这种设计不仅为我们的参数配置分析提供了完整的参数空间，也为机器学习模型的训练数据收集和在线预测部署提供了技术基础。通过深入分析这些组件的实现机制和参数影响，我们可以系统地理解集合通信性能的影响因素，为后续的数据集构建、特征工程和预测模型设计奠定坚实的理论基础。

3.2. 集合通信框架分析

集合通信框架的核心实现位于 [ompi/mca/coll/base/](#) 目录，采用动态组件选择机制为每个通信子配置最优的集合通信实现。

3.2.1. 框架核心文件结构

框架的关键文件包括：

- `coll_base_functions.h` - 定义基础算法接口

该部分定义所有集合通信操作的函数原型和参数宏（`typedef enum COLLTYPE`）；提供算法实现的标准化接口；并声明各种拓扑结构的缓存机制，提供通用的工具函数和数据结构，如二进制树(binary tree)、二项式树(binomial tree)、k 进制树(k-nomial tree)、链式拓扑(chained tree)、流水线拓扑(pipeline) 等。

- `coll_base_comm_select.c` - 实现组件选择机制

该部分为每个通信子动态选择最优的集合通信组件；处理组件优先级和兼容性检查；支持运行时组件偏好设置（通过 `comm->super.s_info` 等机制）。

- `coll_base_util.h` - 工具函数定义

该部分支持配置文件解析和参数处理；提供调试和监控支持。

3.2.2. 核心文件分析

具体而言，粗略分析这部分代码可以观察到：

1. 框架支持 MPI 标准定义的 22 种集合通信操作，通过 `COLLTYPE` 枚举类型进行分类管理：

```
1  typedef enum COLLTYPE {
2      ALLGATHER = 0,          ALLGATHERV,          ALLREDUCE,
3      ALLTOALL,              ALLTOALLV,          ALLTOALLW,
4      BARRIER,              BCAST,                EXSCAN,
5      GATHER,                 GATHERV,            REDUCE,
6      REDUCESCATTER,          REDUCESCATTERBLOCK, SCAN,
7      SCATTER,                 SCATTERV,            NEIGHBOR_ALLGATHER,
8      NEIGHBOR_ALLGATHERV, NEIGHBOR_ALLTOALL,  NEIGHBOR_ALLTOALLV,
9      NEIGHBOR_ALLTOALLW,  COLLCOUNT
10 } COLLTYPE_T;
```

代码 3.1 `coll_base_functions.h` 集合通信操作类型枚举定义

2. 每种集合通信操作都提供三个层次的接口：

- 阻塞接口：如 `BCAST_ARGS`，标准的同步集合通信
- 非阻塞接口：如 `IBCAST_ARGS`，支持异步执行和重叠计算
- 持久化接口：如 `BCAST_INIT_ARGS`，支持 MPI-4 的持久化集合通信

3. 为每个操作提供了丰富的算法变体实现，以 Broadcast 操作为例，框架提供了 10 种不同的算法实现：

```
1  /* Bcast */
2  int ompi_coll_base_bcast_intra_generic(BCAST_ARGS, uint32_t
3  count_by_segment, ompi_coll_tree_t* tree);
4  int ompi_coll_base_bcast_intra_basic_linear(BCAST_ARGS);
5  int ompi_coll_base_bcast_intra_chain(BCAST_ARGS, uint32_t segsize,
6  int32_t chains);
7  int ompi_coll_base_bcast_intra_pipeline(BCAST_ARGS, uint32_t segsize);
8  int ompi_coll_base_bcast_intra_binomial(BCAST_ARGS, uint32_t segsize);
9  int ompi_coll_base_bcast_intra_bintree(BCAST_ARGS, uint32_t segsize);
10 int ompi_coll_base_bcast_intra_split_bintree(BCAST_ARGS, uint32_t
11 segsize);
12 int ompi_coll_base_bcast_intra_knomial(BCAST_ARGS, uint32_t segsize,
13 int radix);
14 int ompi_coll_base_bcast_intra_scatter_allgather(BCAST_ARGS, uint32_t
15 segsize);
16 int ompi_coll_base_bcast_intra_scatter_allgather_ring(BCAST_ARGS,
17 uint32_t segsize);
```

代码 3.2 coll_base_functions.h 中的 bcast 的算法实现变体

同样地，Allreduce 操作提供了 7 种算法，Allgather 提供了 8 种算法，覆盖了从延迟优化到带宽优化的各种应用场景。

4. 多数算法函数支持参数化配置，如：

- segsize 参数：控制消息分段大小，影响内存使用和流水线效率
- radix 参数：控制树形算法的分支数，平衡通信轮数和并发度
- max_requests 参数：控制并发请求数量，影响内存和网络资源使用

5. 以及其它相关拓展和配置接口（如拓扑感知的算法优化等），此处略。

通过上面的分析，结合官方文档 [The Modular Component Architecture \(MCA\) – Open MPI 5.0.x documentation](#)，我们可以较好的理解 Open MPI 集合通信框架，并通过配置参数使用特定的算法实现来优化通信操作的性能。具体的参数配置分析见于[章节 4. 集合通信参数配置分析](#)。

3.3. 集合通信操作示例

为了更好地理解 Open MPI 集合通信框架的工作原理，以一个具体的 MPI_Reduce 操作为例，粗略分析从用户调用到底层算法执行的几个关键过程。

3.3.1. 用户代码调用

在用户的代码中调用 MPI 函数，假设 8 进程对一个 int 数据执行 Reduce 操作，有如下示例代码：

```
1 // size = 8
2 int main(int argc, char** argv) {
```

```
3     MPI_Init(&argc, &argv);
4
5     int rank, size;
6     MPI_Comm_rank(MPI_COMM_WORLD, &rank);
7     MPI_Comm_size(MPI_COMM_WORLD, &size);
8
9     int local_data = rank + 1;
10    int result = 0;
11
12    // 执行Reduce操作, 求和到进程0
13    MPI_Reduce(&local_data, &result, 1, MPI_INT, MPI_SUM, 0, MPI_COMM_WORLD);
14
15    if (rank == 0) {
16        printf("Sum result: %d\n", result);
17    }
18
19    MPI_Finalize();
20    return 0;
21 }
```

代码 3.3 MPI_Reduce 示例代码

3.3.2. 初始调用暂记

该节存在较多困惑尚未解决, 仅作参考用。实际的操作示例描述从[初始化阶段: 组件选择机制](#)开始

当 MPI_Init() 调用时, 首先被 monitoring_prof.c 中的监控层拦截:

```
1 int MPI_Init(int* argc, char*** argv)
2 {
3     // 1. 调用实际的MPI实现
4     result = PMPI_Init(argc, argv);
5
6     // 2. 获取通信子基本信息
7     PMPI_Comm_size(MPI_COMM_WORLD, &comm_world_size);
8     PMPI_Comm_rank(MPI_COMM_WORLD, &comm_world_rank);
9
10    // 3. 初始化MPI_T监控接口
11    MPI_T_init_thread(MPI_THREAD_SINGLE, &provided);
12    MPI_T_pvar_session_create(&session);
13
14    // 4. 注册集合通信监控变量
15    init_monitoring_result("coll_monitoring_messages_count", &coll_counts);
16    init_monitoring_result("coll_monitoring_messages_size", &coll_sizes);
17    start_monitoring_result(&coll_counts);
18    start_monitoring_result(&coll_sizes);
19 }
```

代码 3.4 MPI_Init() 的过程

补充: 这里所指的“实际的 MPI 实现”, 区别于此处 monitoring_prof.c 中的 MPI_Init() :

1. 监控库定义: MPI_Init (拦截版本)
2. 真实库定义: MPI_Init (真实实现)

3. 通过 `#pragma weak: PMPI_Init -> MPI_Init` (实际实现), 实际实现在 `ompi/mpi/c/init.c.in` 中通过模板文件实现, 在编译时生成实际的 c 代码。¹

4. 通过 `LD_PRELOAD`: 用户调用先到监控版本

`init.c.in` 如下:

```

1  PROTOTYPE INT init(INT_OUT argc, ARGV argv)
2  {
3      int err;
4      int provided;
5      int required = MPI_THREAD_SINGLE;
6
7      // 1. 检查环境变量设置的线程级别
8      if (OMPI_SUCCESS > ompi_getenv_mpi_thread_level(&required)) {
9          required = MPI_THREAD_MULTIPLE;
10     }
11
12     // 2. 调用后端初始化函数
13     if (NULL != argc && NULL != argv) {
14         err = ompi_mpi_init(*argc, *argv, required, &provided, false);
15     } else {
16         err = ompi_mpi_init(0, NULL, required, &provided, false);
17     }
18
19     // 3. 错误处理
20     if (MPI_SUCCESS != err) {
21         // return ompi_errhandler_invoke(NULL, NULL,
22         OMPI_ERRHANDLER_TYPE_COMM,
23                                     err < 0 ?
24         ompi_errcode_get_mpi_code(err) : err,
25                                     FUNC_NAME);
26     }
27
28     // 4. 初始化性能计数器
29     SPC_INIT();
30     return MPI_SUCCESS;
31 }

```

代码 3.5 模板生成的 `MPI_Init` 实现

在 `ompi_mpi_init.c` 中的核心初始化过程如下:

```

1  int ompi_mpi_init(int argc, char **argv, int requested, int *provided, bool
2  reinit_ok)
3  {
4      // 1. 状态检查和线程级别设置
5      ompi_mpi_thread_level(requested, provided);
6
7      // 2. 初始化MPI实例
8      ret = ompi_mpi_instance_init(*provided, &ompi_mpi_info_null.info.super,

```

¹报告中对此类实现相关的问题只作了简单的追踪, 未进行进一步的探究


```

8             MPI_ERRORS_ARE_FATAL,
9             &ompi_mpi_instance_default,
10            argc, argv);
11
12    // 3. 初始化通信子子系统
13    if (OMPI_SUCCESS != (ret = ompi_comm_init_mpi3())) {
14        error = "ompi_mpi_init: ompi_comm_init_mpi3 failed";
15        goto error;
16    }
17
18    // 4. 为预定义通信子选择集合通信组件
19    if (OMPI_SUCCESS != (ret = mca_coll_base_comm_select(MPI_COMM_WORLD))) {
20        error = "mca_coll_base_comm_select(MPI_COMM_WORLD) failed";
21        goto error;
22    }
23
24    if (OMPI_SUCCESS != (ret = mca_coll_base_comm_select(MPI_COMM_SELF))) {
25        error = "mca_coll_base_comm_select(MPI_COMM_SELF) failed";
26        goto error;
27    }
28
29    // 5. 标记初始化完成
30    opal_atomic_swap_32(&ompi_mpi_state, OMPI_MPI_STATE_INIT_COMPLETED);
31 }

```

代码 3.6 ompi_mpi_init.c 中的核心初始化过程

3.3.3. 组件选择的核心机制

在 `mca_coll_base_comm_select()` 中执行具体的组件选择:

```

1  int mca_coll_base_comm_select(ompi_communicator_t *comm)
2  {
3      // 1. 初始化通信子的集合通信结构
4      comm->c_coll = (mca_coll_base_comm_coll_t*)calloc(1,
5      sizeof(mca_coll_base_comm_coll_t));
6
7      // 2. 查询所有可用的集合通信组件 (basic, tuned, han, xhc等)
8      selectable =
9      check_components(&ompi_coll_base_framework.framework_components, comm);
10
11     // 3. 按优先级排序并启用最高优先级组件
12     for (item = opal_list_remove_first(selectable);
13         NULL != item;
14         item = opal_list_remove_first(selectable)) {
15         mca_coll_base_avail_coll_t *avail = (mca_coll_base_avail_coll_t *)
16         item;
17         ret = avail->ac_module->coll_module_enable(avail->ac_module, comm);
18
19         if (OMPI_SUCCESS == ret) {
20             // 4. 设置函数指针到具体实现
21             if (NULL == comm->c_coll->coll_reduce) {

```



```

20         comm->c_coll->coll_reduce = avail->ac_module->coll_reduce;
21         comm->c_coll->coll_reduce_module = avail->ac_module;
22     }
23     if (NULL == comm->c_coll->coll_allgather) {
24         comm->c_coll->coll_allgather = avail->ac_module-
>coll_allgather;
25         comm->c_coll->coll_allgather_module = avail->ac_module;
26     }
27     // ... 为所有集合通信操作设置函数指针
28
29     opal_list_append(comm->c_coll->module_list, &avail->super);
30 }
31 }
32
33 // 5. 验证完整性 - 确保所有必需的集合通信操作都有实现
34 #define CHECK_NULL(what, comm, func) \
35     ((what) = # func, NULL == (comm)->c_coll->coll_ ## func)
36
37     if (CHECK_NULL(which_func, comm, allgather) ||
38         CHECK_NULL(which_func, comm, allreduce) ||
39         CHECK_NULL(which_func, comm, reduce) ||
40         // ... 检查其他操作
41     ) {
42         opal_show_help("help-mca-coll-base.txt",
43             "comm-select:no-function-available", true,
44             which_func);
45         return OMPI_ERR_NOT_FOUND;
46     }
47     return OMPI_SUCCESS;
48 }

```

代码 3.7 集合通信组件选择过程

3.3.4. 初始化阶段：组件选择机制

在 `MPI_Init()` 调用时，系统为 `MPI_COMM_WORLD` 选择合适的集合通信组件，这将影响后续 `MPI_Reduce` 的实现方式。

从 `init.c.in` 模板开始的调用链：

```

用户调用: MPI_Init(&argc, &argv)
    ↓
模板生成: init.c.in → MPI_Init()
    ↓
核心初始化: ompi_mpi_init(*argc, *argv, required, &provided, false)
    ↓
通信子初始化: ompi_comm_init_mpi3() [创建MPI_COMM_WORLD]
    ↓
组件选择: mca_coll_base_comm_select(MPI_COMM_WORLD)

```

在组件选择过程中，系统为 `MPI_COMM_WORLD` 设置 `reduce` 函数指针：

```

1 // 在mca_coll_base_comm_select()中
2 if (OMPI_SUCCESS == ret) {
3     // 关键: 为reduce操作设置函数指针
4     if (NULL == comm->c_coll->coll_reduce) {
5         comm->c_coll->coll_reduce = selected_module->coll_reduce;
6         comm->c_coll->coll_reduce_module = selected_module;
7     }
8     // ... 其他集合通信操作的设置
9 }

```

代码 3.8 为 MPI_Reduce 设置函数指针

假设系统选择了 basic 组件, 则根据通信子的大小决定通信算法, 譬如:

```

1 // 在basic组件模块启用时
2 if (ompi_comm_size(comm) <= mca_coll_basic_crossover) {
3     // 小规模通信子: 使用线性算法
4     BASIC_INSTALL_COLL_API(comm, basic_module, reduce,
5                             ompi_coll_base_reduce_intra_basic_linear);
6 } else {
7     // 大规模通信子: 使用对数算法
8     BASIC_INSTALL_COLL_API(comm, basic_module, reduce,
9                             mca_coll_basic_reduce_log_intra);
10 }

```

代码 3.9 basic 组件的通信算法选择机制

对于 8 进程的情况, `mca_coll_basic_crossover (config = 4) < 8`, 则会选择对数算法。

3.3.5. MPI_Reduce 的调用过程

当用户调用

```
MPI_Reduce(&local_data, &result, 1, MPI_INT, MPI_SUM, 0, MPI_COMM_WORLD)
```

第 1 步: MPI 接口层

```

1 // 在ompi/mpi/c/reduce.c中
2 int MPI_Reduce(const void *sendbuf, void *recvbuf, int count,
3               MPI_Datatype datatype, MPI_Op op, int root, MPI_Comm comm)
4 {
5     // 参数验证: 检查buffer、datatype、root等参数
6     if (MPI_PARAM_CHECK) { /* ... */ }
7
8     // 关键调用: 通过函数指针调用已选定的reduce实现
9     return comm->c_coll->coll_reduce(sendbuf, recvbuf, count, datatype,
10                                     op, root, comm,
11                                     comm->c_coll->coll_reduce_module);
12 }

```

代码 3.10 MPI 接口层的分发

第 2 步: 组件实现层

假设在组件选择过程中选择了 basic 组件²，调用转入 `mca_coll_basic_reduce_log_intra`：

```

1 // 在ompi/mca/coll/basic/coll_basic_reduce.c中
2 int mca_coll_basic_reduce_log_intra(const void *sbuf, void *rbuf,
3                                     size_t count,
4                                     struct ompi_datatype_t *dtype,
5                                     struct ompi_op_t *op,
6                                     int root,
7                                     struct ompi_communicator_t *comm,
8                                     mca_coll_base_module_t *module)
9 {
10     int size = ompi_comm_size(comm); // 8
11     int rank = ompi_comm_rank(comm);
12     int dim = comm->c_cube_dim;      // log2(8) = 3
13
14     // 检查操作是否可交换
15     if (!ompi_op_is_commute(op)) {
16         // MPI_SUM是可交换的，所以不会走这个分支
17         return ompi_coll_base_reduce_intra_basic_linear(sbuf, rbuf, count,
18 dtype, op, root, comm, module);
19     }
20
21     // 使用超立方体算法，执行log(N)轮通信
22     int vrank = (rank - root + size) % size; // 虚拟rank，以root为0重新编号
23
24     for (int i = 0, mask = 1; i < dim; ++i, mask <=< 1) {
25         if (vrank & mask) {
26             // 高位进程向低位进程发送并停止
27             int peer = ((vrank & ~mask) + root) % size;
28             MCA_PML_CALL(send(snd_buffer, count, dtype, peer,
29 MCA_COLL_BASE_TAG_REDUCE,
30 MCA_PML_BASE_SEND_STANDARD, comm));
31             break;
32         } else {
33             // 低位进程接收并归约
34             int peer = vrank | mask;
35             if (peer < size) {
36                 peer = (peer + root) % size;
37                 MCA_PML_CALL(recv(rcv_buffer, count, dtype, peer,
38 MCA_COLL_BASE_TAG_REDUCE, comm,
39 MPI_STATUS_IGNORE));
40                 // 执行归约操作
41                 ompi_op_reduce(op, rcv_buffer, rbuf, count, dtype);
42             }
43         }
44     }
45 }

```

代码 3.11 Basic 组件的 `mca_coll_basic_reduce_log_intra` 算法

²如果系统选择了不同的组件，MPI_Reduce 的执行方式会完全不同：

- **basic 组件**：使用简单的线性收集算法（如上所示）
 - **tuned 组件**：根据消息大小和进程数量选择二进制树、流水线等优化算法
 - **han 组件**：使用层次化算法，先在节点内归约，再在节点间归约
- 但对用户而言，调用接口完全相同，这正体现了 Open MPI 组件架构的优势。

3.3.6. 具体执行过程

对于上述对于 8 进程的 reduce 操作 (root=0), 超立方体算法的通信过程如下:

轮次 1 (mask=1): 处理 X 维度

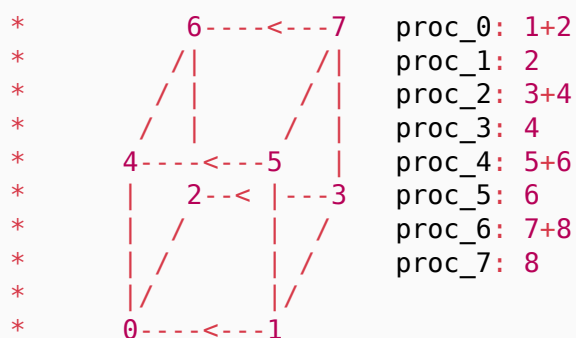
进程1 → 进程0: 发送数据2, 结果: 进程0有 (1+2)
 进程3 → 进程2: 发送数据4, 结果: 进程2有 (3+4)
 进程5 → 进程4: 发送数据6, 结果: 进程4有 (5+6)
 进程7 → 进程6: 发送数据8, 结果: 进程6有 (7+8)

轮次 2 (mask=2): 处理 Y 维度

进程2 → 进程0: 发送 (3+4), 结果: 进程0有 (1+2+3+4)
 进程6 → 进程4: 发送 (7+8), 结果: 进程4有 (5+6+7+8)

轮次 3 (mask=4): 处理 Z 维度

进程4 → 进程0: 发送 (5+6+7+8), 结果: 进程0有 (1+2+3+4+5+6+7+8)=36



代码 3.12 轮次 1 计算过程示意

3.3.7. 调用链总结

完整的 MPI_Reduce 调用链:

用户调用: `MPI_Reduce(&local_data, &result, 1, MPI_INT, MPI_SUM, 0, MPI_COMM_WORLD)`
 ↓
 MPI接口: `ompi/mpi/c/reduce.c::MPI_Reduce()`
 ↓
 函数指针分发: `comm->c_coll->coll_reduce()` [在MPI_Init时设置]
 ↓
 组件实现: `mca_coll_basic_reduce_intra()` [basic组件为例]
 ↓
 算法选择: 超立方体Reduce算法 [8进程, 基于crossover阈值]
 ↓
 底层通信: `MCA_PML_CALL(send/recv)` [点对点消息传递]
 ↓
 结果输出: 进程0获得最终结果36

3.4. 主要集合通信算法实现

该部分仅以

章节 4. 集合通信参数配置分析

章节 5. 数据集构建与应用场景设计

章节 6. 机器学习建模与性能预测

章节 7. 模型验证与检验

章节 8. 结果分析与总结

章节 9. 结论与展望

参考文献
