

Non-Collective Communicator Creation

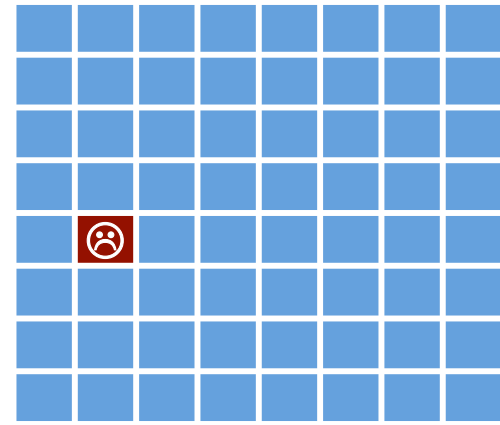
Ticket #286

```
int MPI_Comm_create_group(MPI_Comm comm,  
MPI_Group group, int tag, MPI_Comm *newcomm)
```

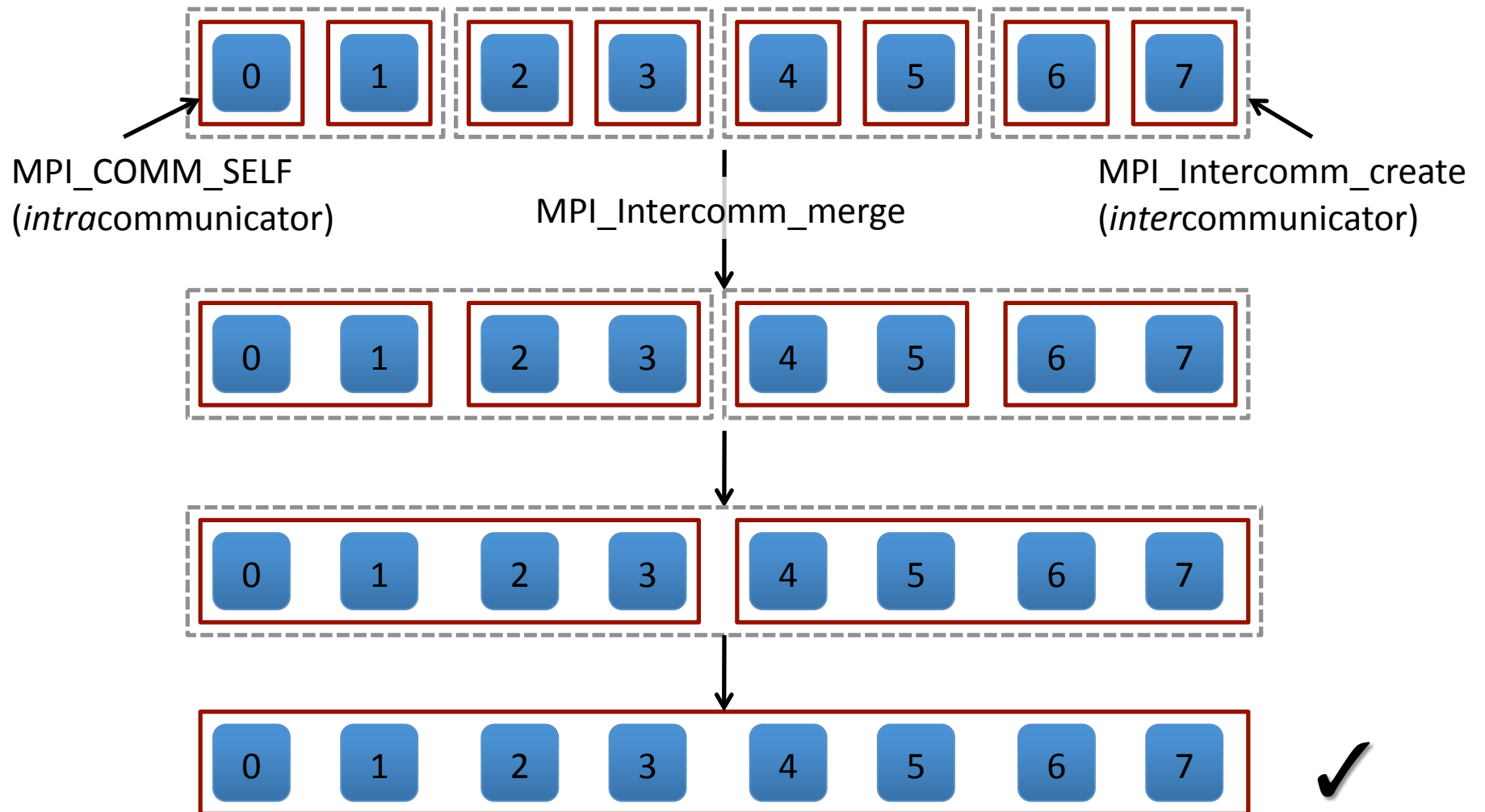
Non-Collective Communicator Creation in MPI. Dinan, et al., Euro MPI '11.

Non-Collective Communicator Creation

- Currently: Collective operation
 - Non-Collective: Create communicator collectively only on new members
1. Avoid bulk synchronization
 - Load balancing: Reassign processes from idle groups to active groups
 2. Overhead reduction
 - Multi-level parallelism, small communicators
 3. Recovery from failures
 - Not all ranks in parent can participate
 4. Compatibility with Global Arrays
 - Past: collectives using MPI Send/Recv



Non-Collective Communicator Creation Algorithm



Non-Collective Algorithm in Detail

INPUT: *group*, *comm*, *tag*

OUTPUT: *comm'*

REQUIRE: *group* is ordered by desired rank in *comm'* and is identical on all callers

LET: *grp_pids*[0..*|group|* - 1] = \mathbb{N} and *pids*[] be arrays of length *|group|*

MPI_Comm_rank(*comm*, &*rank*)

MPI_Group_rank(*group*, &*grp_rank*), MPI_Group_size(*group*, &*grp_size*)

MPI_Comm_dup(MPI_COMM_SELF, &*comm'*)

MPI_Comm_group(*comm*, &*parent_grp*)

MPI_Group_translate_ranks(*group*, *grp_size*, *grp_pids*, *parent_grp*, *pids*)

MPI_Group_free(&*parent_grp*)

} Translate group ranks to
ordered list of ranks on
parent communicator

for (*merge_sz* \leftarrow 1; *merge_sz* < *grp_size*; *merge_sz* \leftarrow *merge_sz* · 2) **do**

gid \leftarrow *grp_rank* / *merge_sz*, *comm_old* \leftarrow *comm'*

if *gid* mod 2 = 0 **then**

if ((*gid* + 1) · *merge_sz* < *grp_size* **then**

 MPI_Intercomm_create(*comm'*, 0, *comm*, *pids*[(*gid* + 1) · *merge_sz*], *tag*, &*ic*)

 MPI_Intercomm_merge(*ic*, 0 /* LOW */, &*comm'*)

end if

else

 MPI_Intercomm_create(*comm'*, 0, *comm*, *pids*[(*gid* - 1) · *merge_sz*], *tag*, &*ic*)

 MPI_Intercomm_merge(*ic*, 1 /* HIGH */, &*comm'*)

end if

if *comm'* \neq *comm_old* **then**

 MPI_Comm_free(&*ic*)

 MPI_Comm_free(&*comm_old*)

end if

end for

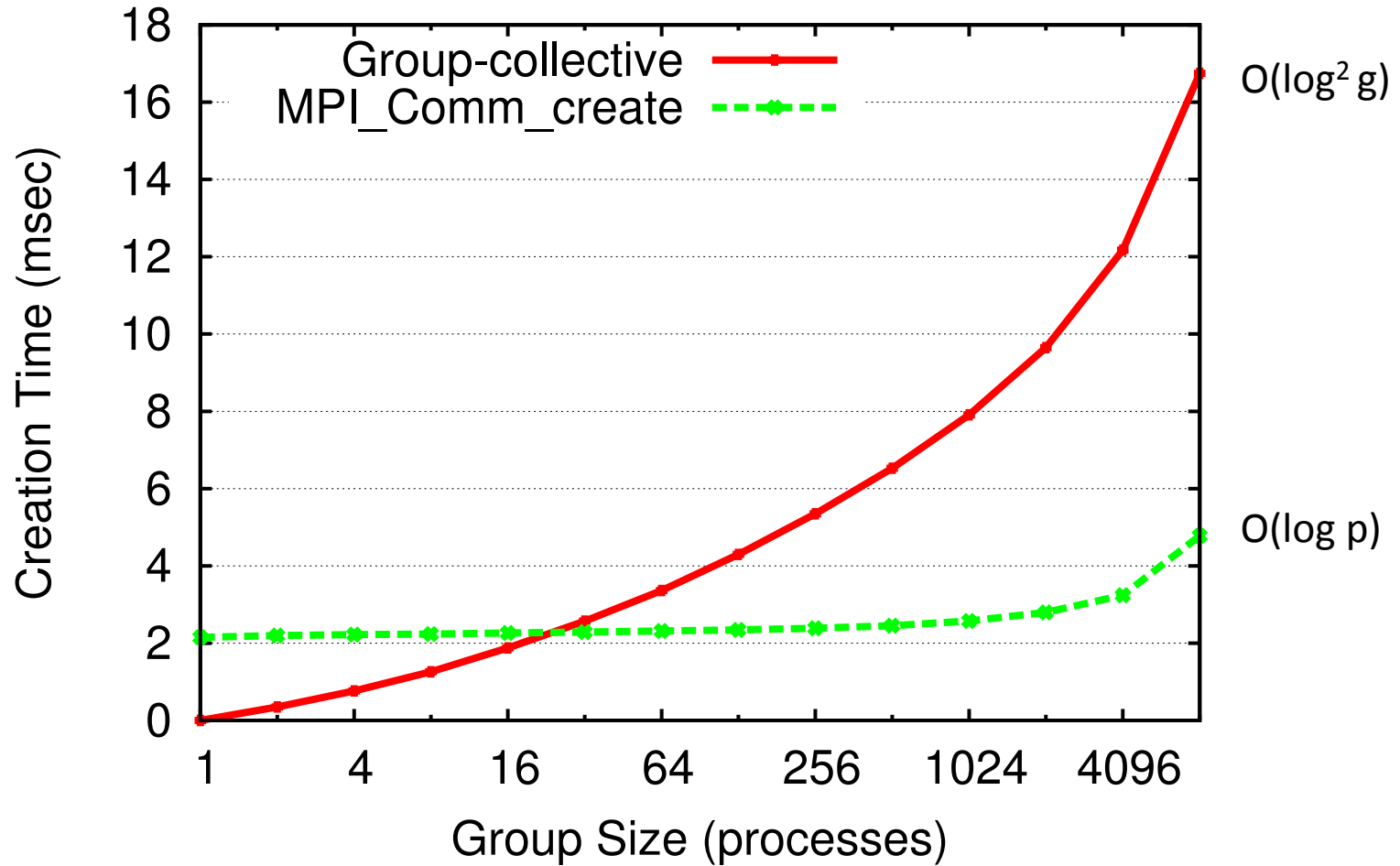
← Calculate my group ID

} Left group

} Right group

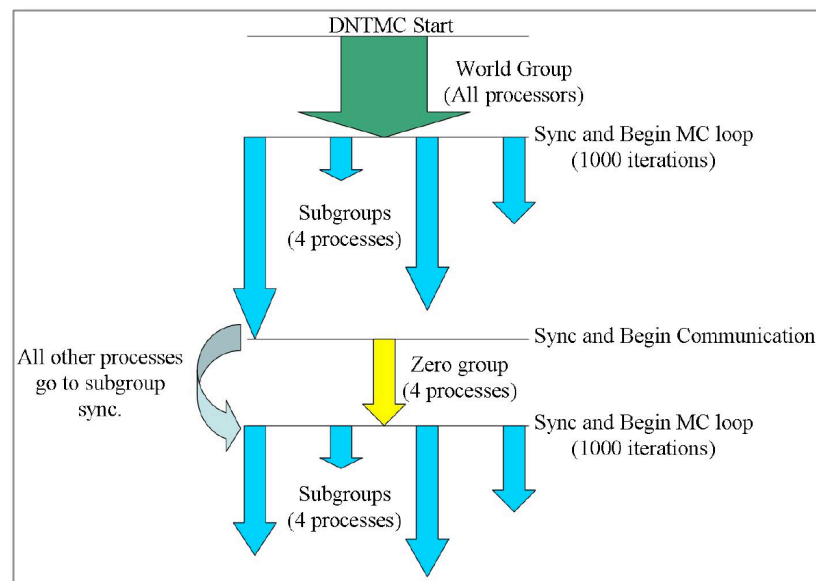


Evaluation: Microbenchmark



Case Study: Markov Chain Monte Carlo

- Dynamical nucleation theory Monte Carlo (DNTMC)
 - Part of NWChem
 - Markov chain Monte Carlo
- Multiple levels of parallelism
 - Multi-node “Walker” groups
 - Walker: N random samples
- Load imbalance
 - Sample computation (energy calculation) is irregular, can be rejected
- Regroup idle processes into active group
 - Preliminary results: 30% decrease in total application execution time



T L Windus et al 2008 *J. Phys.: Conf. Ser.* **125** 012017



- 
- Santorini straw vote to proceed with a formal reading:

Yes (16), No (0), Abstain (1)

