

# The MPI Message-passing Standard

## Practical use and implementation (VI)

SPD Course

24/03/2015

Massimo Coppola

Intracommunicators

# COLLECTIVE PRIMITIVES WITH COMM. AND COMPUTATION

- **int MPI\_Reduce(**  
    **const void\* sendbuf, void\* recvbuf, int count,**  
    **MPI\_Datatype datatype,**  
    **MPI\_Op op, int root, MPI\_Comm comm)**
- reduce operation across all processes of a Communicator
  - Reduces the elements in the same position of each process' buffer, leaving results in root's buffer
- count, datatype, op, root, comm arguments must match
  - If count == 1 we have a classical element-wise reduction
  - If count > 1 we have several reductions at the same time
- As with any collective, the communication pattern is implementation dependent (but is op commutative ?)
- MPI provides most basic operators
  - Operators are associative
  - Operators may be commutative → potential optimizations
  - Note that: floating point op.s may be seen as non-commutative
  - *Datatype* must be compatible with *op*

- Arithmetic operations
  - MPI\_MAX, MPI\_MIN, MPI\_SUM, MPI\_PROD
    - Generally allowed on MPI integral and floating point types (including complex)
- Logic (L) and bit wise (B) operations
  - MPI\_LAND, MPI\_LOR, MPI\_LXOR
    - Generally allowed on C integers and on logical types
  - MPI\_BAND, MPI\_BOR, MPI\_BXOR
    - Generally allowed on C/Fortran integers

# MPI\_MINLOC and MPI\_MAXLOC

- Operators defined on couples (value, index)
  - MPI\_MAXLOC, MPI\_MINLOC
  - Value is any integral or floating point type
  - Index is an integral type
  - chars used as integers require special attention
    - e.g. explicitly using MPI\_SIGNED\_CHAR / MPI\_UNSIGNED\_CHAR
- MINLOC : compute the global minimum of v and the index attached to it
- MAXLOC : compute the global maximum of v and the index attached to it
- Lexicographic order
  - when more values hit the minimum (maximum) the lower one is chosen
- Example application
  - pass (value, rank) to detect the rank of the process with the minimum/maximum value

# MPI couple types

- These couple types are supported by the MPI\_MINLOC and MPI\_MAXLOC operators
- MPI\_FLOAT\_INT - struct { float, int }
- MPI\_LONG\_INT - struct { long, int }
- MPI\_DOUBLE\_INT - struct { double, int }
- MPI\_SHORT\_INT - struct { short, int }
- MPI\_2INT - struct { int, int }
- MPI\_LONG\_DOUBLE\_INT - struct { long double, int }
  - this is an OPTIONAL type

- Operators are called within the reduction collective by the MPI library instances on the processes
- Each operators receives two local buffers and performs a reduction step on their contents
  - The buffers are possibly allocated by the library implementation as temporaries
  - Many operators are polymorphic, so they have to detect the type of data in the buffer
  - Datatype is a parameter passed from the collective down to the operator, but remember it is a handle
    - Easy case: MPI basic datatypes are globally known to MPI runtime and to the program
    - Besides, MPI standard operators are easy

- MPI operators (including user-defined ones) are used by all MPI collectives performing distributed computation
  - MPI\_REDUCE, MPI\_ALLREDUCE, MPI\_REDUCE\_SCATTER\_BLOCK, MPI\_REDUCE\_SCATTER, MPI\_SCAN, MPI\_EXSCAN
  - All the non blocking version of those collectives (since MPI-3)
  - MPI\_REDUCE\_LOCAL (special case actually designed for MPI implementers)



- MPI allows you to define your own operators
  - They can apply to basic and user-defined datatypes
- What do you need to do
  - (Possibly) provide relevant datatype definitions
  - Provide MPI with a definition of the operator
    - a compiled function with a specific signature
    - the operator definition this is local to each process
  - Detect and recognize the MPI\_Datatype within the operator code
    - To detect errors
    - If the operator needs to be polymorphic
  - Combine each couple of elements in the same position of the two input buffers
- Can call **no** MPI communication function
  - Can call MPI\_ABORT in case of error

# MPI\_Op\_create

- `int MPI_Op_create(MPI_User_function* user_fn, int commute, MPI_Op* op)`
- Primitive for defining operators
- Takes a user function pointer as first argument
- Can specify non-commuting operators
- Returns the operator handle
- `MPI_Op_free` allows to free operators

# Operator signature

- **typedef void MPI\_User\_function(  
void\* invec,  
void\* inoutvec,  
int \*len,  
MPI\_Datatype \*datatype);**
- combines pointwise data from two buffers
  - results are placed in the second buffer
- The datatype handle comes from the reduction call and may not be known at compile time
  - For user-defined datatypes, polymorphic operators need to access a table of datatype handles that are defined by the program

# Example : rewriting MPI\_SUM

```
/* this follows the MPI_User_function typedef */
void my_sum_op(void * b_in, void * b_inout,
               int * count, MPI_Datatype * d) {
    if (d == MPI_INT) {
        for (i=0; i<count; i++) {
            ((int*)b_inout)[i]+=((int *)b_in)[i]; }
    } else if (d == MPI_FLOAT) {
        for (i=0; i<count; i++) {
            ((float*)b_inout)[i]+=((float *)b_in)[i]; }
    } else MPI_Abort (MPI_COMM_WORLD, -12345);
}
```

... ..

```
MPI_Op * op_sum;
MPI_Create_op (* my_sum_op, MPI_FALSE, op_sum)
```

- Very limited example: it only accepts INT and FLOAT types
- Can call specialized functions in each case (code reuse, hardware acceleration)

- Check the datatype
  - Compare the received datatype handle to a list of allowed handles, execute proper code
  - Simple if/else error if only one type is allowed
- Check and switch for polymorphic op.s
  - Operators that can handle several datatypes should use data structures in order to avoid the comparison overhead
    - E.g. an hash-map (perfect hash?) associating handles with code (function pointers) implementing each case of use of the operator
  - The overhead is usually negligible with respect to the communication overhead of a reduction or scan

# Scan

- ```
int MPI_Scan(const void* sendbuf, void* recvbuf,  
            int count,  
            MPI_Datatype datatype, MPI_Op op,  
            MPI_Comm comm)
```
- Applies a scan (parallel prefix) to the elements in corresponding position of the send buffers of the processes
- The scan works according to process rank.
  - Scan is the identity for process with rank 0
- If `count > 1` we have multiple scans within the same communication pattern

- `int MPI_Exscan(const void* sendbuf, void* recvbuf, int count, MPI_Datatype datatype, MPI_Op op, MPI_Comm comm)`
- Same as `MPI_Scan`, but results are accumulated on the *following* process
  - Process with rank  $i$  gets the parallel prefix result of data contributed from processes  $0.. i-1$
  - Process 0 receives no data and the receive buffer is not used by `MPI_Exscan`

# Other reduce collective operations

- `MPI_Allreduce`
  - Semantically equivalent to a reduce followed by a broadcast
  - May be implemented more efficiently, of course
- `MPI_Reduce_scatter_block`
  - Performs a reduction, then scatters the result buffer across the processes
  - Parameter `recvcount` is the number of elements received per process after the scatter
    - the overall reduction is computed on  $\text{recvcount} * N$  elements, where  $N$  is the communicator size.
- `MPI_Reduce_scatter`
  - Generalizes the `scatter_block` to a variable scatter (each process can receive a block of different size)
  - The `recvcount` is now an array of block sizes (the array is the same size as the communicator, see `MPI_Scatterv`)



- **With respect to MPI-3 standard**
  - Section 5.9 (Global reduction operators)
  - You can skip `reduce_local`