

# Advanced MPI



SCIENTIFIC &  
DATA-INTENSIVE COMPUTING

Luca Tornatore - I.N.A.F. 

Advanced HPC 2024-2025 @ Università di Trieste

# Outline



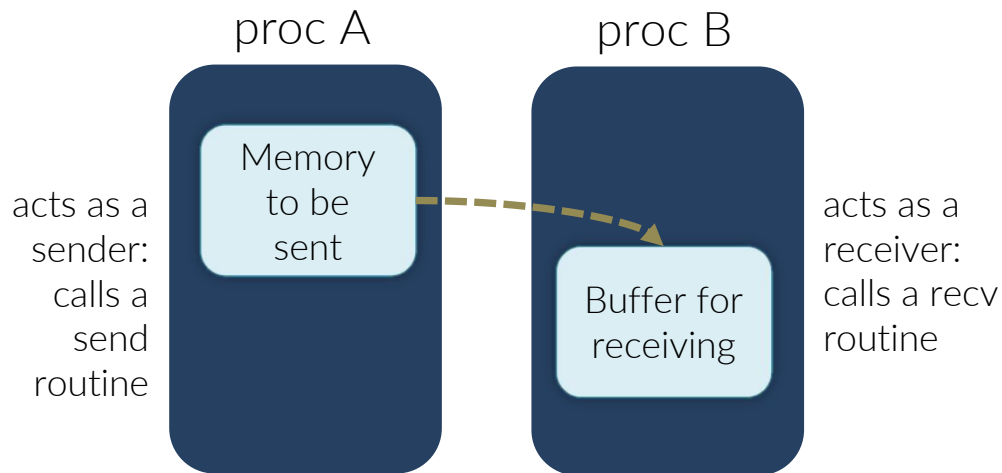
Advanced Usage of  
some **MPI** features  
on *topology* awareness,  
*shared* memory &  
*one-sided* communications

- Basics of **one-sided** communications
- Building a hierarchy of Communicators that reflects the **topology**
- Exchanging data in **shared-memory** windows among MPI processes

# Two-Sided communications

By its very nature, the message-passing paradigm is designed around the concept of cooperative exchange of informations among two or more processes whose address space are isolated and not directly inaccessible by other processes.

This model is very effective in protecting the memory access and in making clear what memory location will be modified and when that happens

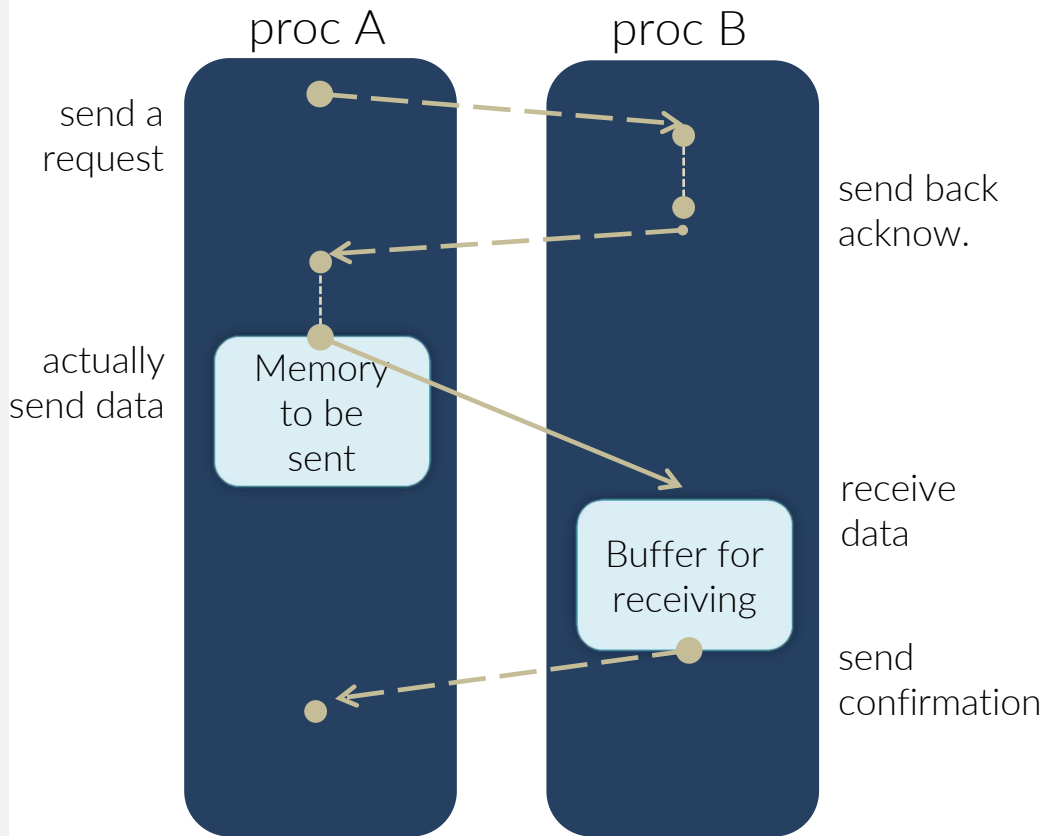


# Two-Sided communications

However, this model has also several cons.

The processes act as peers and must collaborate; as such, the “sender” actually send a request that must be accepted by the receiver.

Moreover, every send must be matched by a receive, and that make certain types of code more complex.



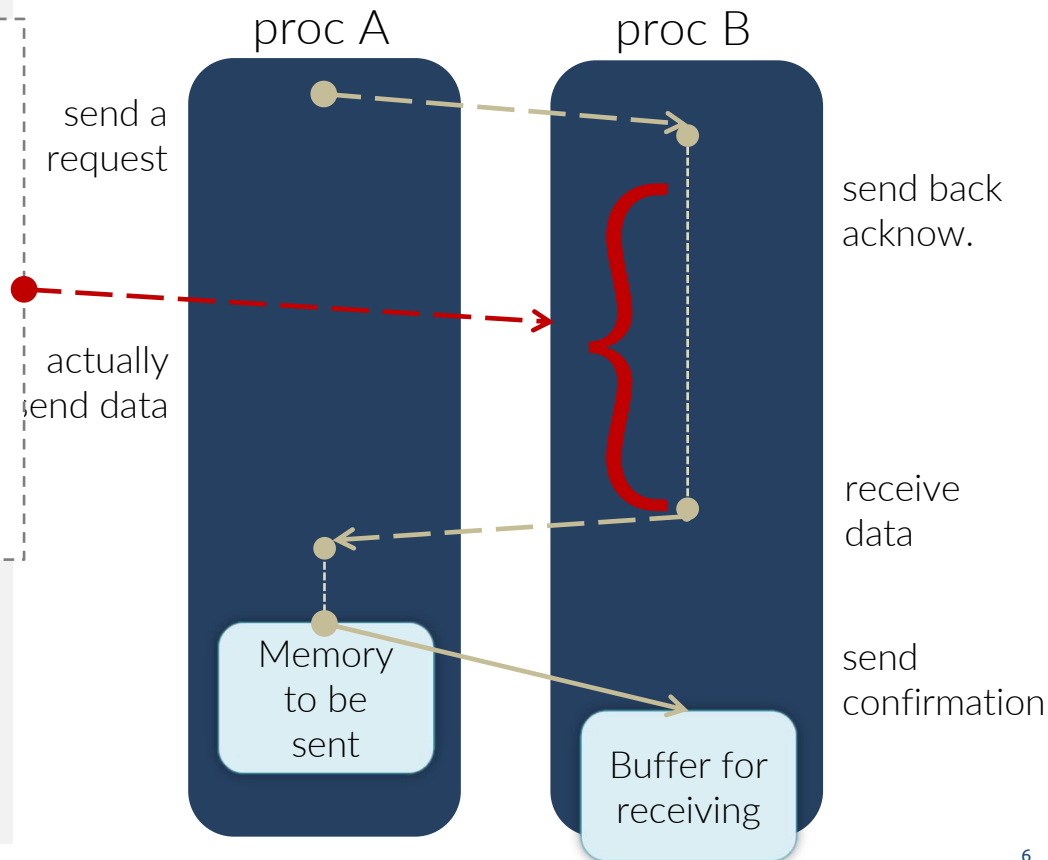
# Two-Sided communications

However, this is not a cons.

The processes collaborate to actually send and receive data.

Moreover, every send must be matched by a receive, and that makes certain types of code more complex.

This delay may be large, depending on the status of process B. Even the Send operation will be delayed as well, because it requires the cooperation of both processes.

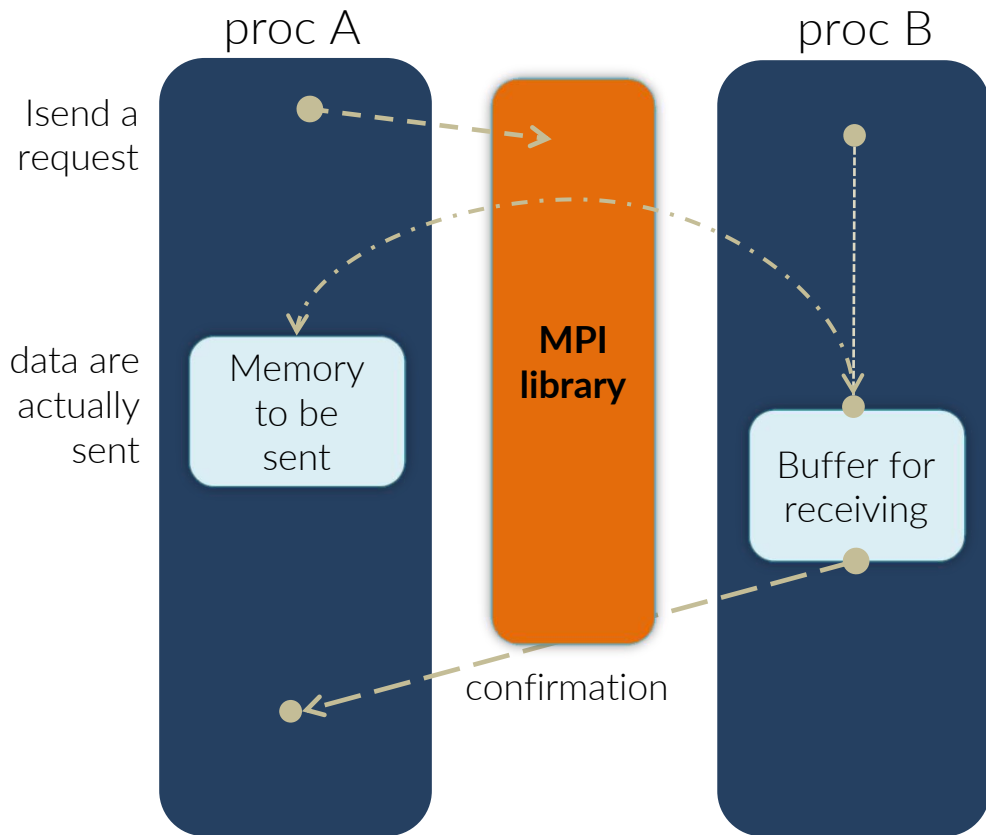


# Two-Sided communications

The *non-blocking* routines mitigate the difficulties linked to the synchronization: the MPI library manages the operations after the sending process posts his request for sending data.

The sending process may even overlook the confirmation about the data reception has concluded, if not needed by its semantics.

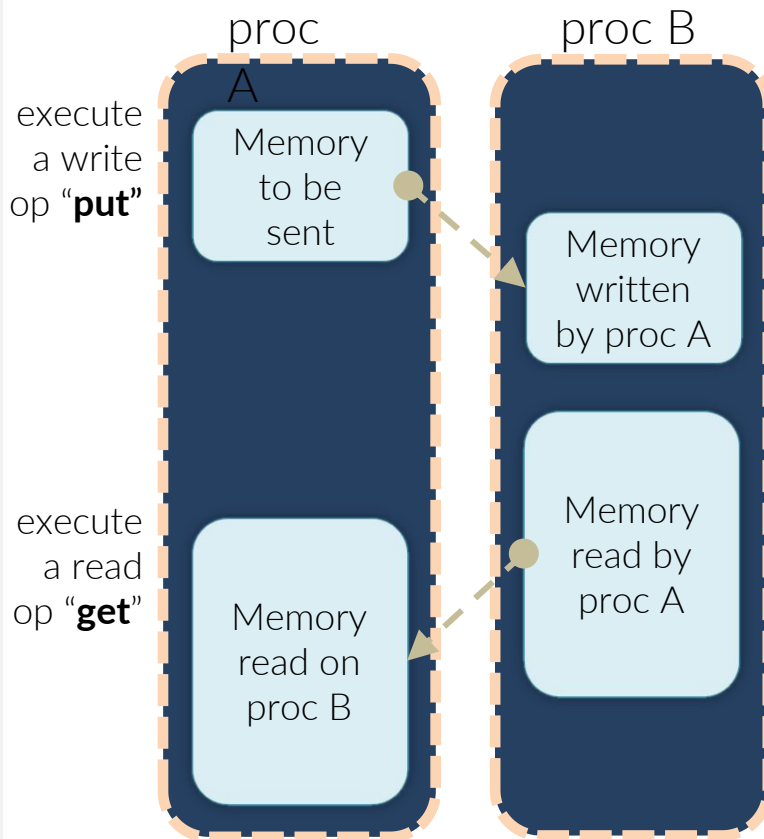
While this makes easier to implement several algorithms, it does not change the fact that the two process need to cooperate.



# One-Sided communications

On the other hand, if **Remote Memory Access (RMA)** was possible, the protocol may be much more relaxed, at the cost of a much larger burden - on the shoulder of the developer:

to ensure a correct synchronization of the operations and the absence of data races or situations with an undefined behaviour.



Proc B does not need to collaborate. In fact, it may even not "know" that proc A is writing or reading.

*well, for consistency reasons, it is advisable that it "knows" that somebody may be writing*

# One-Sided communications

The central idea of one-sided communication models is to disentangle data exchange from the need of processes synchronization

- data movement possible without requiring that the “collaboration” (e.g. sync) of remote process
- Each process exposes a range of its memory (a “window”)
- Other processes can directly read from or write to the exposed memory window



# One-Sided communications

RMA and Shared-Memory are two different concepts.

In **RMA** the players offer a “window” to the other players to access precise memory locations and that access happens in well-defined moments that are circumscribed by *fences*, *epochs* or *locks*.

The **Shared-Memory** is more general: the players share the memory and the access, either in reading or writing, is possible on the entire (shared)memory and the correctness of the operations is entire responsibility of the programmer.

# One-Sided communications

At abstract level, RMA workflow is something like:

```
create_the_memory_window();
```

```
advertise_an_epoch_of_remote_write_access();
```

```
..
```

```
perform_remote_write_access()
```

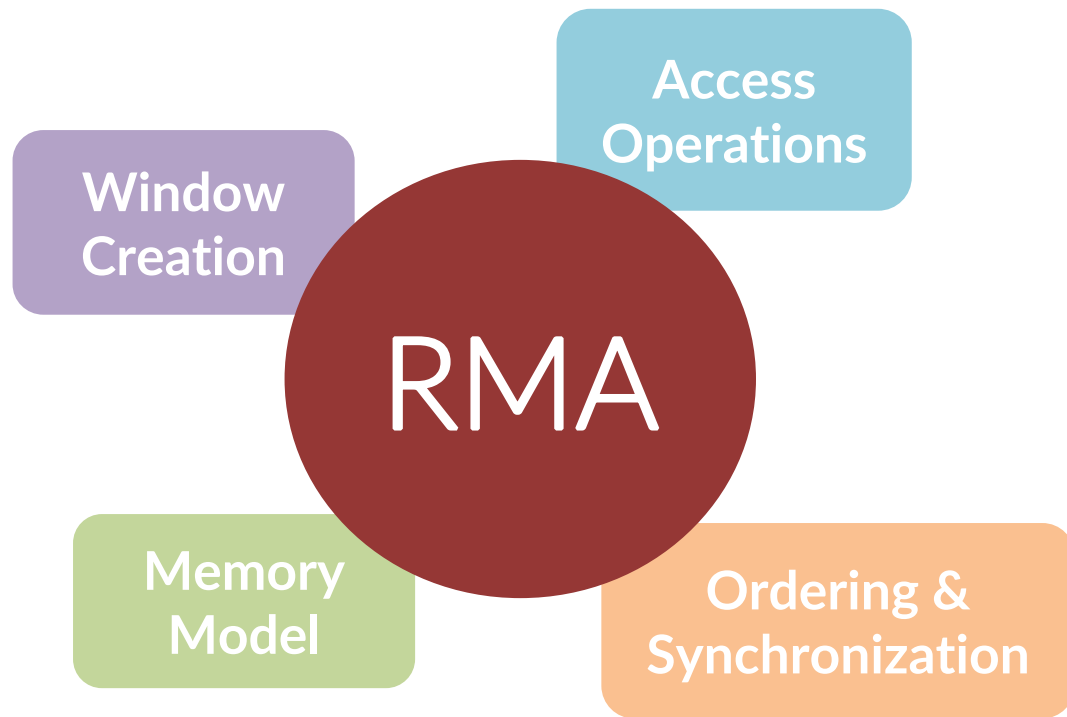
```
..
```

```
close_the_epoch_of_remote_write_access();
```

```
..
```

```
close_the_memory_window();
```

# Key concepts in RMA

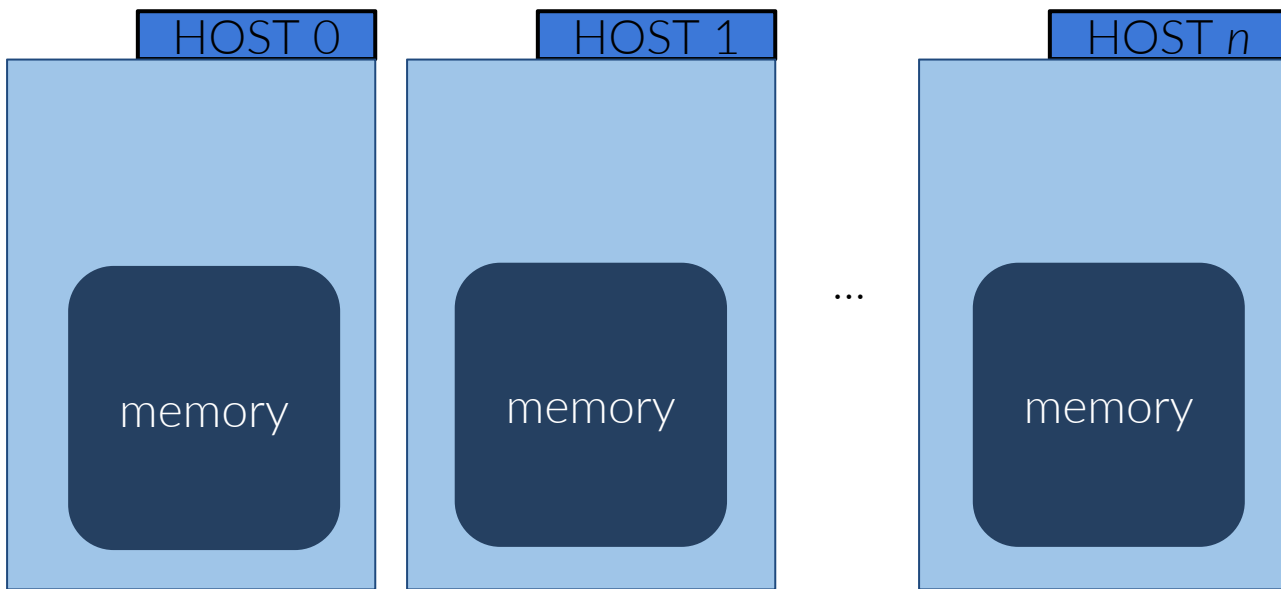


# Key concepts in RMA

## Creating Memory Windows

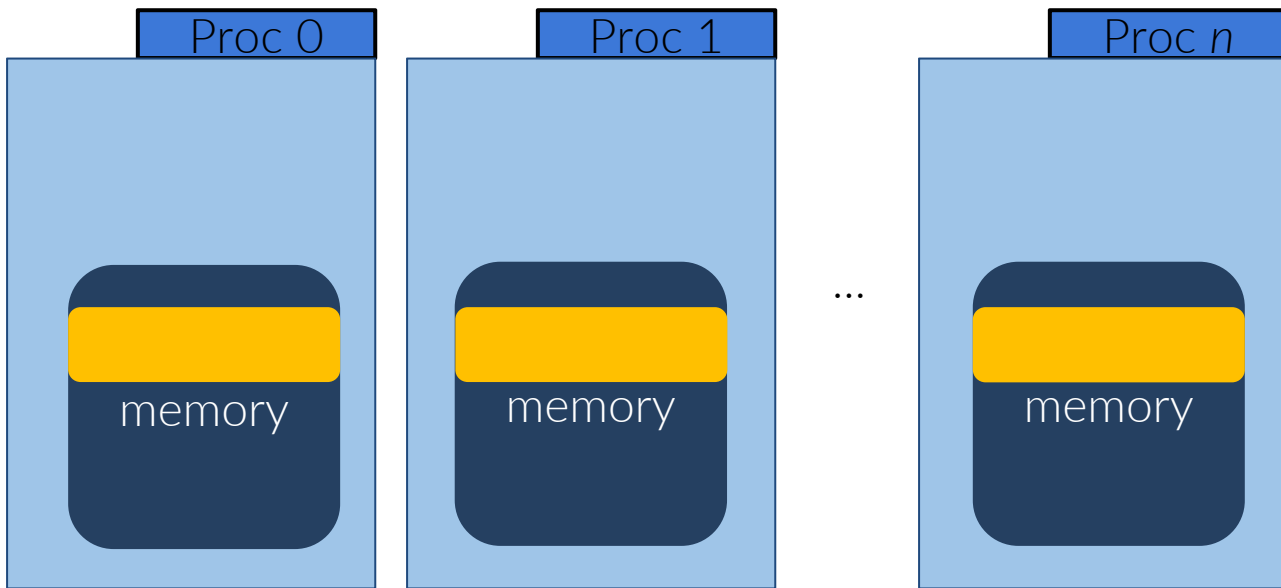
# Creating memory windows

- The memory of each process is only locally accessible to each MPI process



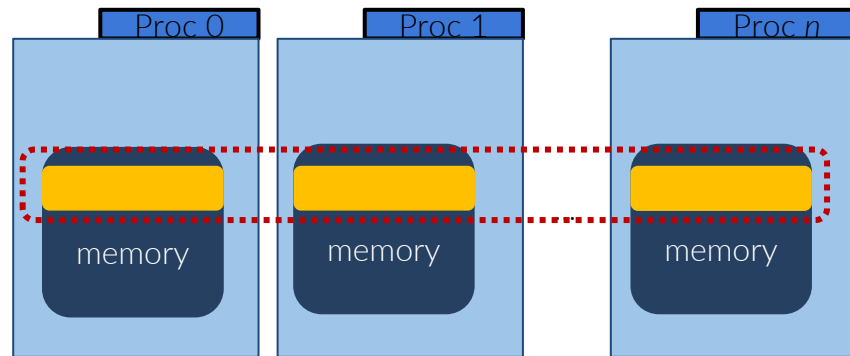
# Creating memory windows

- The programmer has to make an **explicit call to MPI** routines and declare that a region of memory is accessible from remote by the processes **within a given communicator**



# Creating memory windows

- The memory of each process is only locally accessible to each MPI process
- The programmer has to make an explicit call to MPI routines and declare that a region of memory is accessible from remote by the processes *within a given communicator*
  - the region made accessible by RMA is called “a window”
  - the **window creation is a collective operation**
- Once a window is created, the processes are able to write on/read from it remotely **without any collaboration** with the process that owns the memory



# Creating memory windows

There are **four possible ways** to create a window

- **MPI\_Win\_allocate**  
allocates a memory region *and* makes it available - the memory will be released when you *free* the window
- **MPI\_Win\_create**  
a memory region already exists, it will be made a window remotely accessible
- **MPI\_Win\_create\_dynamic**  
creates a window disjointly from a precise memory region; memory buffers will be dynamically added to / removed from the window at any time
- **MPI\_Win\_allocate\_shared**  
allocates memory in a shared memory nodes *and* creates a window linked to it



# Creating memory windows

## Note on memory buffers to be linked to a window and memory allocation

Some MPI implementations may be more efficient if the memory address is *aligned* to the memory pages (and hence to cache lines).

A further important detail may be that the size of the memory buffer is a multiple of the page size.

On `posix` system you may use `posix_memalign` or the C11 `memalign`. An alternative is to use `MPI_Alloc_mem` or `MPI_Win_allocate`.

# Creating memory windows

```
MPI_Win_allocate ( MPI_Aint size, int disp_unit,  
                  MPI_Info info, MPI_Comm comm, void *baseptr,  
                  MPI_Win *win )
```

<b>size</b>	the size of the memory buffer, in bytes (integer)
<b>disp_unit</b>	local units for access offset, in bytes (pos. integer)
<b>info</b>	a handle to an info argument
<b>comm</b>	an handle to a comm object
<b>baseptr</b>	the pointer to the allocated mem buffer (returned by the call)
<b>win</b>	a pointer to a window object (returned by the call)

# Creating memory windows

```
MPI_Win_allocate ( MPI_Aint size, int disp_unit,
                  MPI_Info info, MPI_Comm comm, void *baseptr,
                  MPI_Win *win )
```

Note the type

```
int      N;
double *data;
MPI_win mywin;
```

```
...
MPI_Win_allocate ( N*sizeof(double),
                  sizeof(double), MPI_INFO_NULL, MPI_COMM_WORLD,
                  &data, &mywin );
```

```
... // use data
MPI_Win_free ( &mywin );
```

local data size of double type

always valid not to specify anything here; this is an object to pass hints that may be useful to optimize memory management

It is collective within this communicator

# Creating memory windows

```
MPI_Win_allocate ( MPI_Aint size, int disp_unit,  
                  MPI_Info info, MPI_Comm comm, void *baseptr,  
                  MPI_Win *win )
```

```
int      N;  
double *data;  
MPI_win mywin;
```

```
...  
MPI_Win_allocate ( N*sizeof(double),  
                  sizeof(double), MPI_INFO_NULL, MPI_COMM_WORLD,  
                  &data, &mywin );
```

```
... // use data  
MPI_Win_free ( &mywin );
```

The call itself is a collective, but the window's size can be different at every MPI process



# Creating memory windows

```
MPI_Win_allocate ( MPI_Aint size, int disp_unit,  
                  MPI_Info info, MPI_Comm comm, void *baseptr,  
                  MPI_Win *win )
```

```
typedef struct { double d; int j;  
                char *buffer; } data_t;
```

```
data_t *dat;  
MPI_Win mywin;
```

...

```
MPI_Win_allocate ( N*sizeof(data_t), 1, MPI_INFO_NULL,  
MPI_COMM_WORLD, &data, &mywin );
```

```
... // use data  
MPI_Win_free ( &mywin );
```

```
typedef struct { double d; int j;  
                char *buffer; } data_t;
```

```
data_t *dat;  
MPI_Win mywin;
```

...

```
MPI_Win_allocate ( N*sizeof(data_t), sizeof(data_t),  
MPI_INFO_NULL, MPI_COMM_WORLD, &data, &mywin );
```

```
... // use data  
MPI_Win_free ( &mywin );
```

# Creating memory windows

```
MPI_Win_create ( void *base, MPI_Aint size, int disp_unit,  
                MPI_Info info, MPI_Comm comm, MPI_Win *win )
```

```
int      N;  
double *data;  
MPI_Win mywin;
```

```
MPI_Alloc_mem ( N*sizeof(double), MPI_INFO_NULL, &data );  
data[j] = ...;
```

```
...
```

```
MPI_Win_create ( data, N*sizeof(data_t), sizeof(double),  
                MPI_INFO_NULL, MPI_COMM_WORLD, &mywin );
```

```
... // use data
```

```
MPI_Win_free ( &mywin );
```

```
MPI_Free ( data );
```

# Creating memory windows

```
MPI_Win_create_dynamic ( MPI_Info info, MPI_Comm comm, MPI_Win *win )
```

```
int      N;  
double *data;  
MPI_win mywin;
```

```
MPI_Win_create_dynamic ( MPI_INFO_NULL, MPI_COMM_WORLD, &mywin );
```

```
MPI_Alloc_mem ( N*sizeof(double), MPI_INFO_NULL, &data );  
data[j] = ...;
```

```
...  
MPI_Win_attach ( mywin, data, N*sizeof(double) );
```

```
... // use data  
MPI_Win_detach ( mywin, data );  
MPI_Win_free ( &mywin );
```

# Setting info for memory windows

The `MPI_info` argument provides optimization hints to the runtime about the expected usage pattern of the window. The following info keys are predefined:

`"no_locks"` (boolean, default: false)

`"accumulate_ordering"` (string, default: `rar,raw,war,waw`)

`"accumulate_ops"` (string, default: `same_op_no_op`)

`"mpi_accumulate_granularity"` (integer, default 0)

You can set an info argument by

```
MPI_Info info;  
MPI_Info_create( &info );  
MPI_Info_set( info, "no_locks", "true" );  
MPI_Win_create( ..., info, ... )  
MPI_Info_free( &info );
```



# Setting info for memory windows

`"no_locks" (boolean, default: false)`

if set to **true**, then the implementation may assume that passive target synchronization (i.e., **MPI\_WIN\_LOCK**, **MPI\_WIN\_LOCK\_ALL**) will not be used on the given window. This implies that this window is not used for 3-party communication, and RMA can be implemented with no (less) asynchronous agent activity at this process.

`"accumulate_ordering" (string, default rar,raw,var,waw)`

controls the ordering of accumulate operations at the target.

`"accumulate_ops" (string, default: same_op_no_op)`

if set to **"same\_op"**, the implementation will assume that all concurrent accumulate calls to the same target address will use the same operator.

If set to **"same\_op\_no\_op"**, then the implementation will assume that all concurrent accumulate calls to the same target address will use the same operator or **MPI\_NO\_OP**.

This can eliminate the need to protect access for certain operators where the hardware can guarantee atomicity.

# Setting info for memory windows

```
"mpi_accumulate_granularity" (integer, default 0)"
```

provides a hint to implementations about the desired *synchronization granularity for accumulate operations*, i.e., the size of memory ranges in bytes for which the implementation should acquire a synchronization primitive to ensure atomicity of updates.

If the specified granularity is not divisible by the size of the type used in an accumulate operation, it should be treated as if it was the next multiple of the element size.

For example, a granularity of 1 byte should be treated as 8 in an accumulate operation using `MPI_UINT64_T`. By default, this info key is set to 0, which leaves the choice of synchronization granularity to the implementation. If specified, all MPI processes in the group of a window must supply the same value.

# Key concepts in RMA

## Data movement

# Move memory

You can **read**, **write**, **modify** data atomically

- **MPI\_PUT, MPI\_GET**

- **MPI\_Accumulate**

- **MPI\_Get\_accumulate**

- **MPI\_Compare\_and\_swap**

- **MPI\_Fetch\_and\_op**

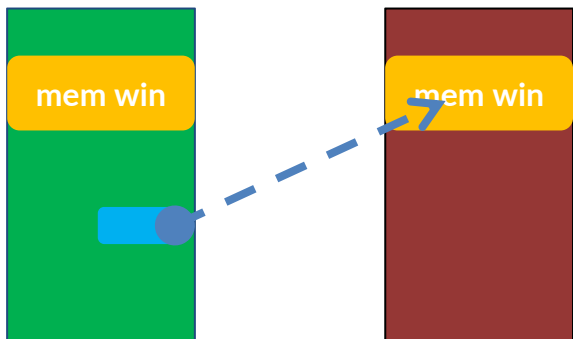
**atomic operations**

# Move memory: put and get

move data **from origin to target**

```
MPI_Put ( void *origin_addr, int origin_count,
MPI_Datatype origin_dtype,
int target_rank,
MPI_Aint target_disp, int target_count,
MPI_Datatype target_dtype,
MPI_Win win )
```

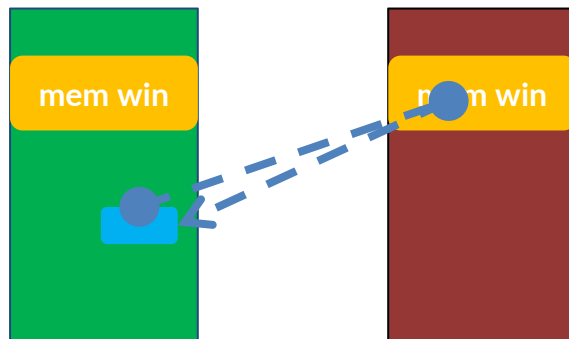
**origin** (calling proc)   **target** (owner of accessed mem)



move data **to origin from target**

```
MPI_Get ( void *origin_addr, int origin_count,
MPI_Datatype origin_dtype,
int target_rank,
MPI_Aint target_disp, int target_count,
MPI_Datatype target_dtype,
MPI_Win win )
```

**origin** (calling proc)   **target** (owner of accessed mem)

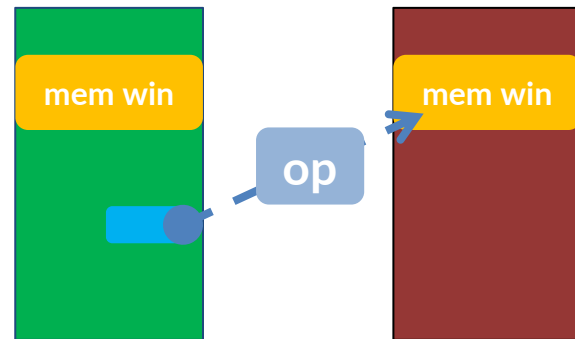


# Move memory: accumulate

```
MPI_Accumulate ( void *origin_addr, int origin_count, MPI_Datatype origin_dtype,
                 int target_rank,
                 MPI_Aint target_disp, int target_count, MPI_Datatype target_dtype,
                 MPI_Op op,
                 MPI_Win win )
```

This implement an **atomic** update operation

- perform the **op** reduction operation between the origin data and the target data
  - OP are the reduction operations: MPI\_SUM, MPI\_PROD, MPI\_OR, MPI\_NO\_OP, ...
  - user-defined operations are not allowed
- if **op=MPI\_REPLACE** you have an **atomic put**



# Why do we need atomics ?

```
if (Rank == 0) N = 1;
else N = 0;

MPI_Win_allocate ( N*sizeof(int), sizeof(int), MPI_INFO_NULL,
                  MPI_COMM_WORLD, &global_counter, &mywin );

...

int counter;
if (Rank > 0 )
{
    MPI_Get( &counter, 1, MPI_INT,      // origin
            0,                          // target rank
            0, 1, MPI_INT, mywin );    // target + win

    while ( counter < MAX ) {
        do_something();
        counter++;
        MPI_Put( &origin, 1, MPI_INT, 0, 0, 1, MPI_INT, mywin); }
}
```

Is this gonna work?

# Why do we need atomics ?

```
if (Rank == 0) N = 1;
else N = 0;

MPI_Win_allocate ( N*sizeof(int), sizeof(int), MPI_INFO_NULL,
                  MPI_COMM_WORLD, &global_counter, &mywin );

...

int counter;
if (Rank == 1 )
{
    MPI_Get( &counter, 1, MPI_INT, 0, 0, 1, MPI_INT, mywin );

    while ( there_is_something_todo ) {
        do_something();
        counter++;
        MPI_Put ( &counter, ... ); }
}
```

Is this gonna work?

note: the MPI standard does **not enforce the order** of execution of put and get operations

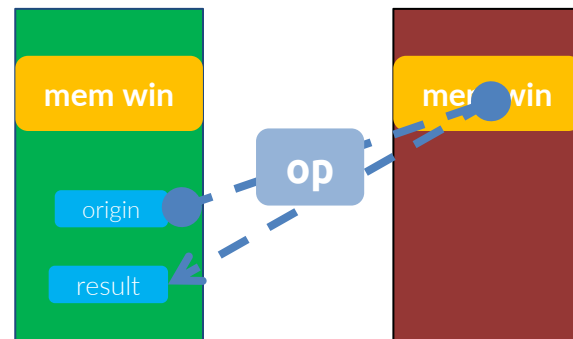


# Move memory: get\_accumulate

```
MPI_Get_accumulate ( void *origin_addr, int origin_count, MPI_Datatype origin_dtype,
                    void *result_addr, int result_count, MPI_Datatype result_dtype,
                    int target_rank,
                    MPI_Aint target_disp, int target_count, MPI_Datatype target_dtype,
                    MPI_Op op, MPI_Win win )
```

This implement an **atomic** read-modify-write operation

- perform the **op** reduction operation between the origin data and the target data
  - OP are the reduction operations: MPI\_SUM, MPI\_PROD, MPI\_OR, MPI\_NO\_OP, ...
  - user-defined operations are not allowed
- if **op**=MPI\_REPLACE you have an **atomic swap**
- if **op**=MPI\_NO\_OP you have an **atomic get**
- the result of the op is stored in the target buffer
- the value at target before op is stored in result buffer
- basic datatype must match



# Move memory: CAS and fop

```
MPI_Fetch_and_op ( void *origin_addr, void *target_addr,  
                  MPI_Datatype origin_dtype, int target_rank,  
                  MPI_Aint target_disp, MPI_Op op, MPI_Win win )
```

```
MPI_Compare_and_swap ( void *origin_addr, void *compare_addr, void *target_addr,  
                      MPI_Datatype origin_dtype, int target_rank,  
                      MPI_Aint target_disp, MPI_Win win )
```

FOP: it is an MPI\_Get\_accumulate but for 1 basic data at a time → more optimized

CAS: it is an atomic swap between origin and target if the target value is equal to compare value;

# Move memory: request-based

MPI\_Rput  
MPI\_Rget  
MPI\_Raccumulate  
MPI\_Rget\_accumulate

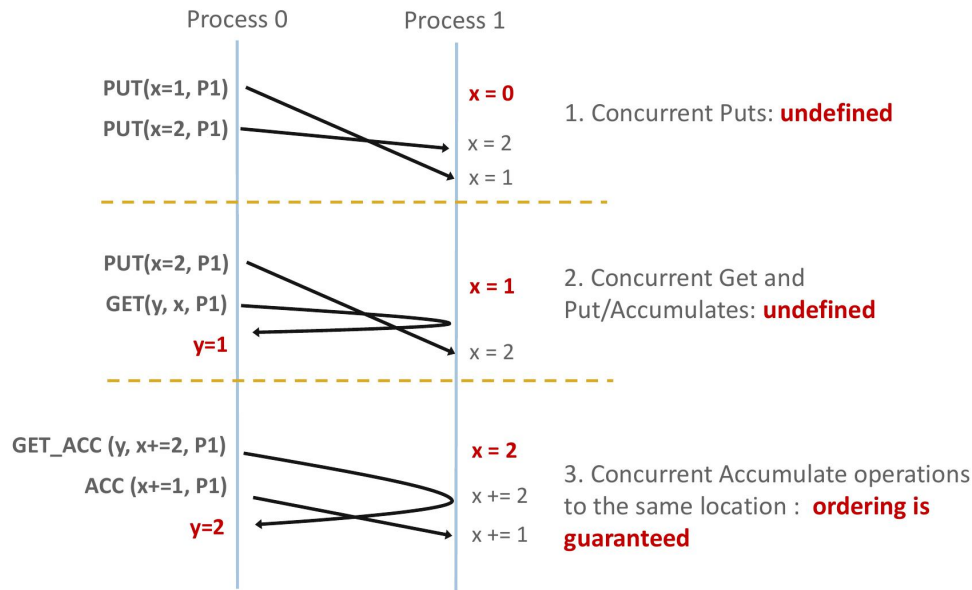
MPI offers also Request-based versions of put, get, accumulate, get\_accumulate.

I.e. routines, to be used only within *passive synchronization* (i.e. lock/unlock; see later), that return a request handle that can be managed using MPI\_Wait.

# Key concepts in RMA

## Ordering and Sync

# Ordering and Sync RMA ops



from "Advanced MPI programming", SC24

- MPI does not ensure any ordering for put and get
  - Concurrent puts and gets have an undefined result
  - A get concurrent with put/accumulate have an undefined result
- concurrent accumulate have a result defined accordingly to the order of operations.
 

Remind:

  - Accumulate with `op=MPI_REPLACE`
  - Get\_accumulate with `op=MPI_NO_OP`
- Accumulate ops from the same process are ordered by default
  - you can tell the MPI not to order, as optimization hint
  - you can ask RAW, WAR, RAR or WAW

*note: here above "concurrent" mean "concurrent on the same memory location"*

# Ordering and Sync RMA ops

All RMA routines are “non-blocking” routines: as such, to ensure the correctness of the operations, they need to be appropriately surrounded by synchronization calls

- to ensure that operations are completed
- to ensure that cache sync have been done

Access model:

- When a process is allowed to read/write remote memory ?
- When data written by process A are available for process B ?

# Ordering and Sync RMA ops

All RMA routines are “non-blocking” routines: as such, to ensure the correctness of the operations, they need to be appropriately surrounded by synchronization calls

- to ensure that operations are completed
- to ensure that cache sync have been done

There two types of synchronization:

**active**      both origin and target have to call sync routines

**passive**      only the origin calls the sync routine

# Ordering and Sync RMA ops

## active

both origin and target have to call sync routines

**MPI\_Fence**

“active target”

A collective call that surround RMA routines to isolate different access types

**MPI\_Win\_post, MPI\_Win\_start,  
MPI\_Win\_complete,  
MPI\_Win\_wait**

“generalized active target”

Collectives that apply to a sub-group, to restrict the overhead of the needed communication

## passive

only the origin calls the sync routine

**MPI\_Win\_lock, MPI\_Win\_unlock**



# Ordering and Sync RMA ops

Data access happens between “epochs”

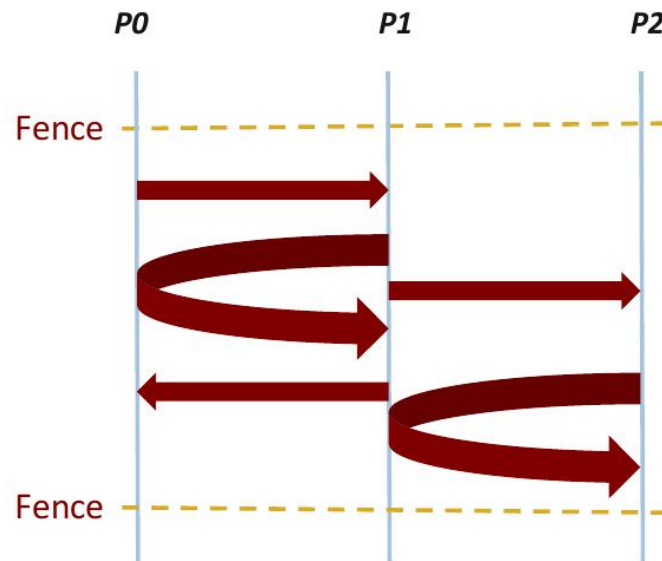
- **Access** *epochs*: a set of operations issued by origin processes
- **Exposure** *epochs*: remote processes can update a target’s memory window

Epochs define operation ordering and completion semantics

The synchronization model provides a way to establish epochs

# Fences

- A collective call, which opens and closes access and exposure epochs on *all* processes in the window
1. everyone in the windows post a **MPI\_Win\_fence** to open the epoch
  2. everyone is allowed to issue MPI\_Put and MPI\_Get
  3. everyone posts a **MPI\_Win\_fence** to close the epoch
  4. all operations completes at the exit of the second fence synchronization



from "Advanced MPI programming", SC24

# Fences: semantics

`MPI_Win_fence (int assert, MPI_Win win )`

The `assert` argument is used to provide optimization hints to the implementation:  
`assert == 0` is always valid.

Valid values may be combined with a bitwise OR operation (`assert1 | assert2`)

```
MPI_Win_fence ( 0, mywin );  
while ( there_is_something_todo ) {  
    ret = do_something( );  
    MPI_Accumulate( &ret, 1, MPI_INT, register[Rank], 0, 1, MPI_INT, MPI_SUM, mywin ); }  
MPI_Win_fence ( 0, mywin );
```

# Fences: assertions

`MPI_Win_fence (int assert, MPI_Win win )`

## **MPI\_MODE\_NOSTORE**

The local window **was not** updated by any local store since the last call to `MPI_Win_fence`. This assert refers to operations *before* the present fence call.

## **MPI\_MODE\_NOPUT**

The local window **will not** be remotely updated by put or accumulate between the present fence call and the next one. This assert involve *future* operations.

## **MPI\_MODE\_NOPRECEDE**

The called fence **will not** conclude any RMA calls made by the process calling the fence; then, no RMA calls should have been made between this call and the previous call (basically it says “no RMA to complete”).

## **MPI\_MODE\_NOSUCCEED**

the symmetric than before: no RMA calls **will be** made on this window before the next fence call (“no RMA to start”).

# Fences: assertions

```
MPI_Win_fence (int assert, MPI_Win win )
```

example

```
MPI_Win_fence(MPI_MODE_NOPRECEDE, win);
```

```
MPI_Put(..., target_a, ..., win);
```

```
MPI_Put(..., target_b, ..., win);
```

```
MPI_Win_fence(MPI_MODE_NOSTORE | MPI_MODE_NOPUT | MPI_MODE_NOSUCCEED,  
              win);
```

Similar to **Fence**, but not collective: origin and target declare themselves, in a sense

**Target:** declares **exposure** epoch

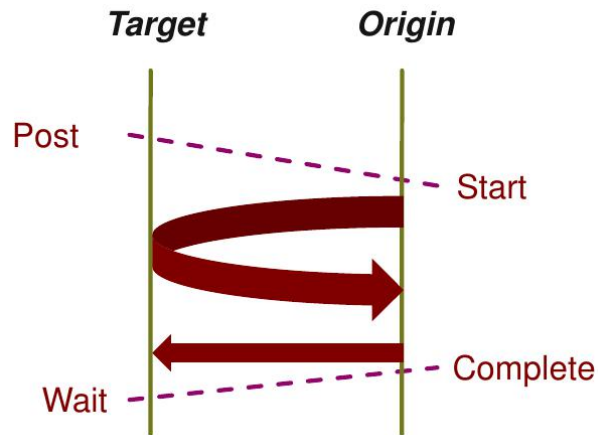
- **MPI\_Win\_post** initiates it
- **MPI\_Win\_complete** finalizes it

**Origin:** declares **access** epoch

- **MPI\_Win\_start** starts it
- **MPI\_Win\_wait** closes it

All synchronizations are allowed to block and defer to ensure the P-S/C-W ordering

Every process can be origin *and* target



from "Advanced MPI programming", SC24

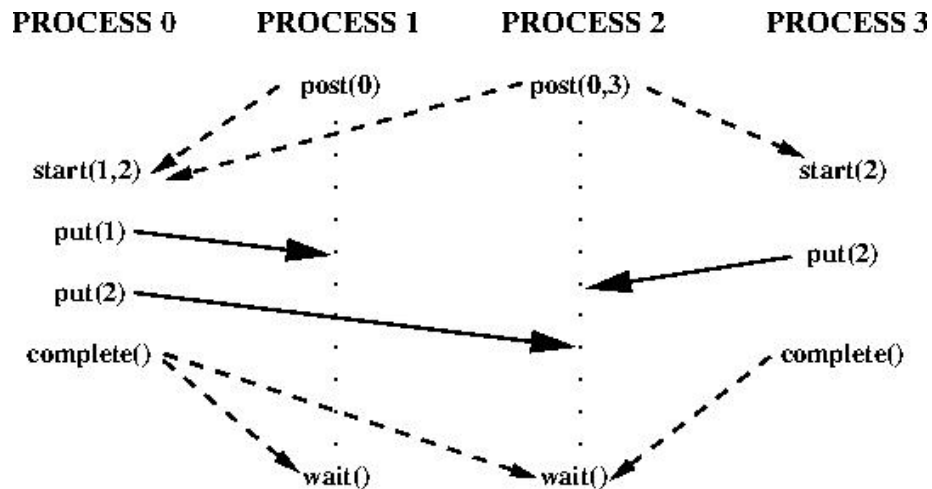
Similar to **Fence**, but not collective: origin and target declare themselves, in a sense

**Target:** declares **exposure** epoch

- `MPI_Win_post` initiates it
- `MPI_Win_complete` finalizes it

**Origin:** declares **access** epoch

- `MPI_Win_start` starts it
- `MPI_Win_wait` closes it



from MPI-forum: [link](#)

All synchronizations are allowed to block and defer to ensure the P-S/C-W ordering

Every process can be origin *and* target

```
MPI_Win_start    ( MPI_Group to_group, int assert, MPI_Win win )  
MPI_Win_complete ( MPI_Win win )
```

```
MPI_Win_post ( MPI_Group from_group, int assert, MPI_Win win )  
MPI_Win_wait ( MPI_Win win )
```

These routines are somehow equivalent to the fence call, but for the fact that they are not mandatorily executed by *all* the processes that are in the group of processes that created the window. They may be execute by a sub-group, even by 2 processes.

The processes that *expose* their window initiate the *exposure epoch* with `MPI_Win_start` and ends it with `MPI_Win_wait`.

Instead, the processes that will *access* the windows initiate the *access epoch* by `MPI_Win_post` and close it by `MPI_Win_complete`.



# PSCW: assertions

`MPI_Win_start` (`MPI_Group` `to_group`, `int` `assert`, `MPI_Win` `win` )

`MPI_Win_post` (`MPI_Group` `from_group`, `int` `assert`, `MPI_Win` `win` )

`MPI_Win_post`'s asserts

( zero is always correct )

**`MPI_MODE_NOSTORE`**

The local window was not updated by any local store since the last call to `MPI_Win_complete`.

**`MPI_MODE_NOPUT`**

The local window *will not* be remotely updated by put or accumulate between the present fence call and the next matching `MPI_Win_complete`

**`MPI_MODE_NOCHECK`**

The matching `MPI_Win_start` have not been issued by a process that is an origin of an RMA with this process as a target (basically: no cross-RMAops). The matching `MPI_Win_start` *must* use the same assert.

`MPI_Win_start`'s asserts

( zero is always correct )

**`MPI_MODE_NOCHECK`**

Guarantees that the matching calls to `MPI_Win_post` have already been made

# PSCW: example

```
MPI_Group world_group;
MPI_Comm_group(MPI_COMM_WORLD, &world_group);
```

```
if ( Me == origin ) {
```

a group with multiple targets

```
    MPI_Group target_group;
    int ntargets = ...;
    int target_ranks[ntargets];
    MPI_Group_incl( world_group, ntargets, target_ranks, &target_group );
```

RMA ops

```
    MPI_Win_start( target_group, 0, win );
    // series of puts, gets, accumulate, etc.
    MPI_Put ( ... );
    MPI_Get ( ... );
    //
    MPI_Win_complete( win ); }
```

```
else {
```

a group with one origin

```
    MPI_Group origin_group;
    int norigins = 1;
    int origin_rank = origin;
    MPI_Group_incl( world_group, norigins, &origin_ranks, &origin_group );
```

Opens and  
closes exposure

```
    MPI_Win_post( origin_group, 0, win );
    MPI_Win_wait( win ); }
```

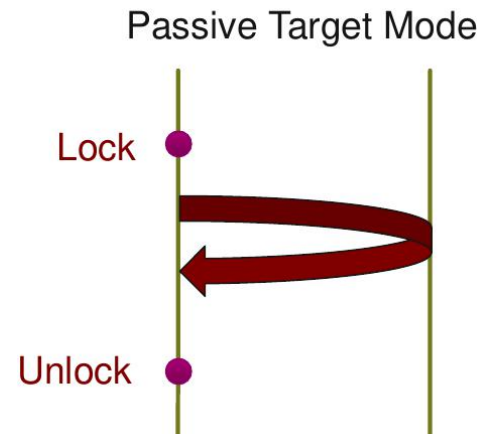
A snippet with one task being the origin of RMA ops towards multiple targets.

This can be easily generalized to a case in which many tasks are both origin and targets

## Passive mode

the target does **not participate** in operations

One-sided *asynchronous* communication



from "Advanced MPI programming", SC24

```
MPI_Win_lock (int lock_type, int rank, int assert, MPI_Win win )
```

```
MPI_Win_unlock (int rank, MPI_Win win )
```

```
MPI_Win_flush/flush_local (int rank, MPI_Win win )
```

Where `lock_type` can be:

`MPI_LOCK_SHARED`

concurrent access to the same target is allowed;

you are in charge of ensuring that no race conditions happen.


`MPI_LOCK_EXCLUSIVE`

concurrent access to the same target is **not** allowed

**NOTE:** “lock” has been an unfortunate choice for the name. That is not a mutual exclusion, it is more similar to just start/stop of RMA operations

**Flush:** complete operations on remote target process; data will be then available to the target task, or to other tasks

**Flush\_local:** locally complete operations to the target process



Advanced  
MPI


# Locks

A snippet to exemplify how to use the lock

```
MPI_Win win;

if (rank == 0) {
    // Rank 0 will perform the put, it does not need a window
    MPI_Win_create(NULL, 0, 1, MPI_INFO_NULL, MPI_COMM_WORLD, &win);
    // "locks" process 1
    MPI_Win_lock(MPI_LOCK_SHARED, 1, 0, win);
    // returns when succeeded
    MPI_Put(..., 1, ..., win);
    // returns when put has succeeded
    MPI_Win_unlock(1, win);
    MPI_Win_free(&win); }
else {
    // Rank 1 is the target, we need a window
    MPI_Win_create(..., ..., ..., MPI_INFO_NULL, MPI_COMM_WORLD, &win);
    //
    MPI_Win_free(&win);
}
```

Advanced HPC @ UniTs - 2024/2025



Luca Tomatore

54

```
MPI_Win_lock_all (int lock_type, int assert, MPI_Win win )
```

```
MPI_Win_unlock_all ( MPI_Win win )
```

```
MPI_Win_flush_all/flush_local_all (int rank, MPI_Win win )
```

**Lock\_all** starts an `MPI_LOCK_SHARED` on the group of tasks that participate in the window, whereas **unlock\_all** ends it.  
The routine is *not* collective

# Key concepts in RMA

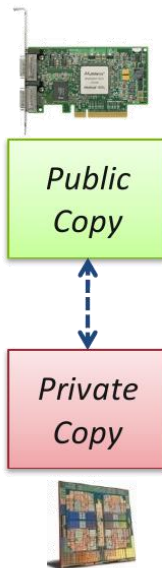
## Memory model

# Memory model in RMA

From MPI-3, two memory models are provided for RMA

- **Separate**, inherited from MPI-2  
designed to work on systems that do not provide cache coherence at hardware level.  
The “public” copy and the “private” memory are separated and MPI provides software coherence
- **Unified**  
there is only 1 copy of the window  
cache coherence must be at system level

Separate



Unified







# Memory model in RMA

```
int *model, flag;  
MPI_Win_get_attr(win, MPI_WIN_MODEL, &model,&flag);  
int is_unified = (*model == MPI_WIN_UNIFIED);
```

That's all folks, have fun

"So long  
and thanks  
for all the fish"



*That's all Folks!*