

# Remote Memory Access

---

Getting started with RMA



# Reusing this material



This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License.

<http://creativecommons.org/licenses/by-nc-sa/4.0/>

This means you are free to copy and redistribute the material and adapt and build on the material under the following terms: You must give appropriate credit, provide a link to the license and indicate if changes were made. If you adapt or build on the material you must distribute your work under the same license as the original.

Note that this presentation contains images owned by others. Please seek their permission before reusing these images.

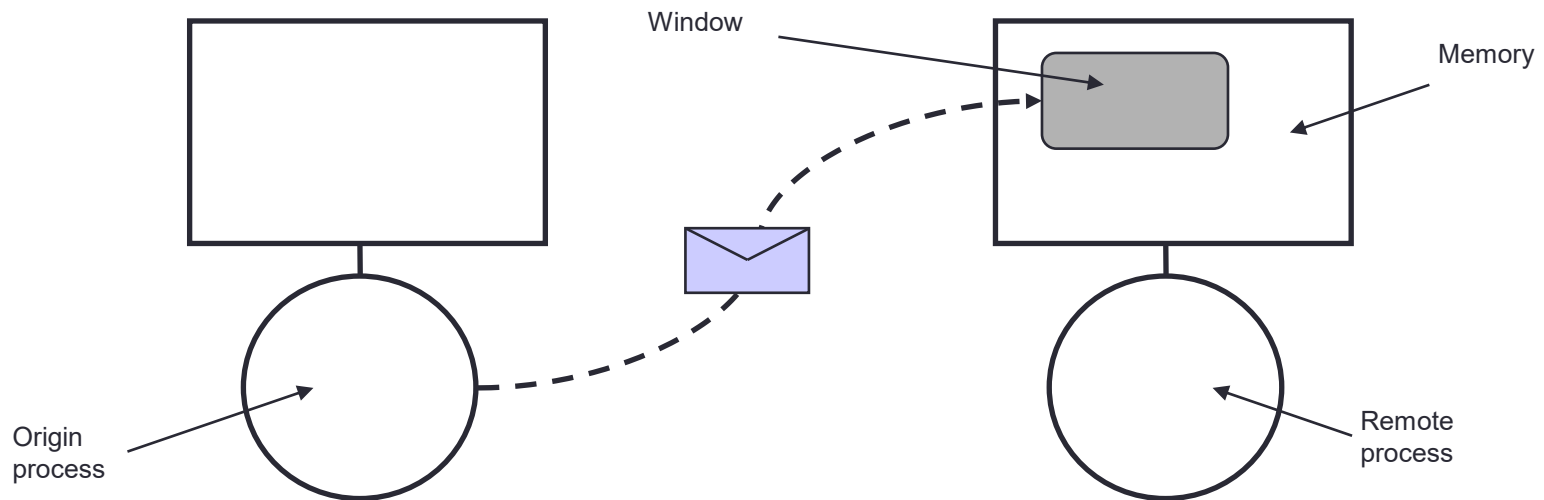
# Outline

- MPI RMA Basic Concepts
  - Why RMA?
  - Terminology
  - Program flow
- Getting started with RMA
  - Management of windows
  - Fence synchronization
  - Moving data around
- Practical
  - Modifying P2P code to use RMA

# MPI RMA Concepts

# Single-Sided Model

- Remote memory can be read or written directly using library calls



- Remote process does not actively participate
  - No matching receive (or send) needs to be performed
  - Synchronisation is now a major issue

# Motivation

- Why extend the basic message-passing model?
- Hardware
  - Many supercomputer networks support Remote Memory Access (RMA) in hardware
  - This is the fundamental model for SMP systems
  - Many users started to use RMA calls for efficiency
    - Lead to the development of non-portable parallel applications
- Software
  - Many algorithms naturally single-sided
    - e.g., sparse matrix-vector
  - Matching send/receive pairs requires extra programming
  - Even worse if communication structure changes
    - e.g., adaptive decomposition

# Why RMA

- One-sided comms functions are an interface to MPI RMA
  - I think “one sided” is a confusing term because, as we will see, whilst the communication calls themselves are one sided often the synchronisation is issued on both sides
- Is a natural fit for some codes
- Can provide a performance/scalability increase for codes
  - Programmability reasons
  - Hardware (interconnect) reasons
  - But is not a silver bullet!

# Terminology

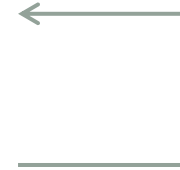
- Origin is process initiating the request (performs the call)
  - Irrespective of whether data is being retrieved or written
- Target is the process whose memory is accessed
  - By the origin, either remotely reading or writing to this
- All remote access performed on windows of memory
- All access calls non-blocking and issued inside an *epoch*
  - The epoch is what forces synchronisation of these calls



# RMA program flow

- Collectively initialise a window

- a) Start an RMA epoch (synchronisation)
- b) Issue communication calls
- c) Stop an RMA epoch (synchronisation)



Repeat as many  
times as you want

- Collectively free the window

# Getting started with RMA

Window management, fences and data movement

# Window creation

- A collective call, issued by all processes in the communicator

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

- Each process may specify completely different locations, sizes, displacement units and info arguments.
- You can specify no memory with a zero size and NULL base
- The same region of memory may appear in multiple windows that have been defined for a process. But concurrent communications to overlapping windows are disallowed.
- Performance may be improved by ensuring that the windows align with boundaries such as word or cache-line boundaries.

# Other window management

- Retrieving window attributes

```
int MPI_Win_get_attr(MPI_Win win, int win_keyval,  
    void *attribute_val, int *flag)
```

- win\_keyval is one of MPI\_WIN\_BASE, MPI\_WIN\_SIZE, MPI\_WIN\_DISP\_UNIT, MPI\_WIN\_CREATE\_FLAVOR, MPI\_WIN\_MODEL
- Attribute\_val if the attribute is available and in this case (flag is true), otherwise flag will be false

- Freeing a window

```
int MPI_Win_free(MPI_Win *win)
```

- All RMA calls must have been completed (i.e. the epoch stopped)

# Fences

- Synchronisation calls required to start and stop an epoch
  - Fences are the simplest way of doing this where global synchronisation phases alternate with global communication
- Most closely follows a barrier synchronisation
  - A (collective) fence is called at the start and stop of an epoch

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

`MPI_Win_fence(0, window);` ← *RMA can not be started until this first fence*  
*Communication calls go here*  
`MPI_Win_fence(0, window);` ← *All issued communication calls block here*  
← *Default value – no assertions*

# Fence attributes

- Attributes allow you to tell the MPI library more information for performance (but MPI implementations are allowed to ignore it!)
  - **MPI\_MODE\_NOSTORE** local window was not updated by local writes of any form since last synchronisation. *Can be different on processes*
  - **MPI\_MODE\_NOPUT** local window will not be updated by put/accumulate RMA operations until AFTER the next synchronisation call. *Can be different on processes*
  - **MPI\_MODE\_NOPRECEDE** fence does not complete any sequence of locally issues RMA calls. *Attribute must be given by all processes*
  - **MPI\_MODE\_NOSUCCEED** fence does not start any sequence of locally issued RMA calls. *Attribute must be given by all processes*
- Attributes can be or'd together, i.e.
  - `MPI_Win_fence((MPI_MODE_NOPUT | MPI_MODE_NOPRECEDE), window)` or `ior(MPI_MODE_NOPUT, MPI_MODE_NOPRECEDE)`

# RMA Communication calls

- Three general calls, all non-blocking:

- Get data from target's memory

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

- Put data into target's memory

```
int MPI_Put(const void *origin_addr, int origin_count,  
            MPI_Datatype origin_datatype, int target_rank,  
            MPI_Aint target_disp, int target_count,  
            MPI_Datatype target_datatype, MPI_Win win)
```

- Accumulate data in target's memory with some other data

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

# RMA communication comments

- Similarly to non-blocking P2P one must wait for synchronisation (i.e. end of the epoch) until accessing retrieved data (*get*) or overwriting written data (*put/accumulate*)
- `target_disp` is multiplied by window displacement unit, `origin_count` and `target_count` are in units of data type
- Undefined operations:
  - Local stores/reads with a remote PUT in an epoch
  - Several origin processes performing concurrent PUT to the same target location
  - Single origin process performing multiple PUTs to same target location in a single epoch
- Accumulate supports the `MPI_Reduce` operations, but NOT user defined operations. Also supports `MPI_REPLACE` - Effectively the same as a put.



# Generic Simple Approach

- Declare local storage on each rank
- Create a window including all storage: `MPI_Win_create()`
  - replaces the communicator in subsequent RMA calls
- Access data in local storage using normal array operations
- Synchronise so everyone is ready: `MPI_Win_fence()`
  - Issue remote reads / writes to from / to data on other processes
    - `MPI_Get()` and `MPI_Put()`
- Synchronise so everyone is finished: `MPI_Win_fence()`
- Can now access data in local storage as normal

# Example

Based on an example at  
[cvw.cac.cornell.edu/MPIoneSided/fence](http://cvw.cac.cornell.edu/MPIoneSided/fence)

```
MPI_Win win;
int ctrlbuf[20], mybuf[20];
if (rank == 0) {
    MPI_Win_create(ctrlbuf, sizeof(int)*20, sizeof(int),
        MPI_INFO_NULL, comm, &win);
} else {
    MPI_Win_create(NULL, 0, 1, MPI_INFO_NULL, comm, &win);
}
if (rank == 0) initialise(ctrlbuf);

MPI_Win_fence(MPI_MODE_NOPRECEDE, win);
if (rank != 0) {
    MPI_Get(mybuf, 20, MPI_INT, 0, 0, 20, MPI_INT, win);
}
MPI_Win_fence(MPI_MODE_NOSUCCEED, win);

if (rank != 0) process(mybuf);
MPI_Win_free(&win);
```

*Rank 0 creates a window of 20 integers, displacement unit = 4 bytes (= 1 integer)*

*Other ranks create a window but attach no local memory*

*Fence, no preceding RMA calls*

*Non-zero ranks get the 20 integers from rank 0, disp 0*

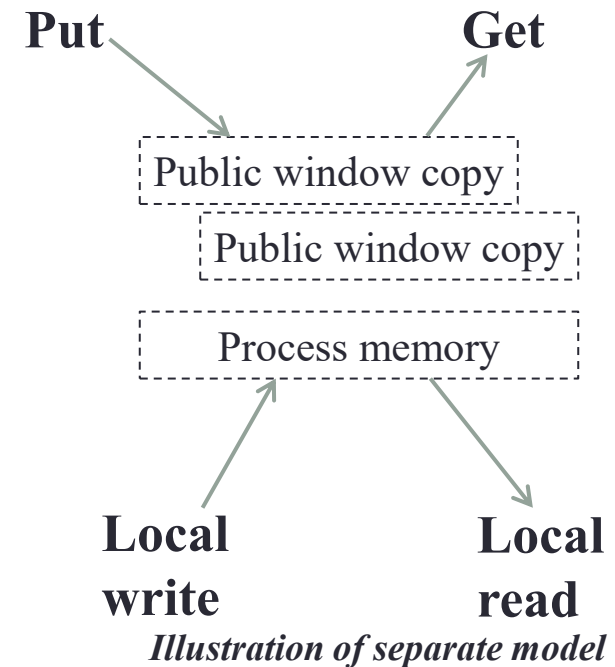
*Fence, complete all communications and no RMA calls in next epoch*

# RMA Memory model

- Public and private window copies
  - Public memory region is addressable by other processes (i.e. exposed main memory)
  - Private memory (i.e. transparent caches or communication buffers) which is only locally visible but elements from public memory might be stored.
- Coherent if updates to main memory are automatically reflected in private copy consistently
- Non-coherent if updates need to be explicitly synchronised

# RMA Memory model

- MPI therefore has two models
  - Unified if public and private copies are identical – used if possible, realistic on cache coherent machines. (*This was added in MPI v3*)
  - Separate if they are not, here there is only one copy of a variable in process memory but also a distinct public copy for each window that contains it. The old model
- In the separate model a suitable synchronisation call (i.e. end of an epoch) must be issued to make these consistent. In the unified model some synchronisation calls might be omitted for performance reasons
- The window attribute tells you which model it follows



# Additional functionality

- You should be aware that:
  - you *can* do point-to-point synchronisation in MPI RMA as well as global synchronisation with fences
    - it is called “Post-Start-Complete-Wait” - PSCW
    - you do not need to know any more details than this!
  - you *can* lock windows
    - this is called “passive target synchronisation” as the target process does not need to make any RMA calls
    - model is something like: lock / put / unlock
    - you do not need to know any more details than this!
- Covered in detail in additional lecture
  - provided for info but you **do not need to learn this extra material**

# Summary

- Model is quite simple
  - although syntax can be quite challenging
- Performance may not be very good
  - portability and flexibility requirements of MPI mean that latency may not be as small as you hoped
- However
  - windows are a key component of MPI shared-memory approach
  - see later ...