

# Tutorial on the MPL interface to MPI

Victor Eijkhout

PEARC 2023



While the C API to MPI is usable from C++, it feels very unidiomatic for that language. Message Passing Layer (MPL) is a modern C++17 interface to MPI. It is both idiomatic and elegant, simplifying many calling sequences.

<https://github.com/rabauke/mpl>



1 pingpongbuffer 19

2 scangather 35

3 scangather 36

4 setdiff 57



# Part I

## Basics



For doing the exercises:

```
module load mpl
```

which defines TACC\_MPL\_INC, TACC\_MPL\_DIR



To compile MPL programs, add a line

```
1  #include <mpl/mpl.hpp>
```

to your file. You need to add a path to your compile line:

```
1  mpicxx -o mpiprogram -I${MPL_LOCATION}/include mympiprogram.cpp
```

where MPL\_LOCATION is system-dependent.



There is no initialization or finalize call.

*Implementation note: Initialization is done at the first  
`mpl::environment` method call, such as `comm_world`.*



The naive way of declaring a communicator would be:

```
1 // commrank.cxx
2 mpl::communicator comm_world =
3   mpl::environment::comm_world();
```

calling the predefined environment method `comm_world`.

However, if the variable will always correspond to the world communicator, it is better to make it `const` and declare it to be a reference:

```
1 const mpl::communicator &comm_world =
2   mpl::environment::comm_world();
```





The `processor_name` call is an environment method returning a `std::string`:

```
1  std::string mpl::environment::processor_name ();
```



The rank of a process (by `mpl::communicator::rank`) and the size of a communicator (by `mpl::communicator::size`) are both methods of the `communicator` class:

```
1  const mpl::communicator &comm_world =  
2      mpl::environment::comm_world();  
3  int procid = comm_world.rank();  
4  int nprocs = comm_world.size();
```



The `environment` namespace has the equivalents of `MPI_COMM_WORLD` and `MPI_COMM_SELF`:

```
1  const communicator& mpl::environment::comm_world();  
2  const communicator& mpl::environment::comm_self();
```

Uses of `MPI_COMM_NULL` are handled differently.



Communicators can be duplicated but only during initialization. Copy assignment has been deleted. Thus:

```
1  // LEGAL:  
2  mpl::communicator init = comm;  
3  // WRONG:  
4  mpl::communicator init;  
5  init = comm;
```



Pass communicators by reference to avoid communicator duplication:

```
1 // commpass.cxx
2 // BAD! this does a MPI_Comm_dup.
3 void comm_val( const mpl::communicator comm );
4
5 // correct!
6 void comm_ref( const mpl::communicator &comm );
```



The communicator class has its copy operator deleted; however, copy initialization exists:

```
1 // commcompare.cxx
2 const mpl::communicator &comm =
3     mpl::environment::comm_world();
4 cout << "same: " << boolalpha << (comm==comm) << endl;
5
6 mpl::communicator copy =
7     mpl::environment::comm_world();
8 cout << "copy: " << boolalpha << (comm==copy) << endl;
9
10 mpl::communicator init = comm;
11 cout << "init: " << boolalpha << (init==comm) << endl;
```

(This outputs true/false/false respectively.)

*Implementation note: The copy initializer performs an  
MPI\_Comm\_dup.*



# Part II

## Collectives



Collectives have many polymorphic variants, for instance for 'in place', and buffer handling.

Operators are handled through functors.





*Buffer* type handling is done through polymorphism (templating and ADL):

no explicit indication of types.

Scalars are handled as such:

```
1 float x,y;  
2 comm.bcast( 0,x ); // note: root first  
3 comm.allreduce( mpl::plus<float>(), x,y ); // op first
```

where the reduction function needs to be compatible with the type of the buffer.



If your buffer is a `std::vector` you need to take the `.data()` component of it:

```
1  vector<float> xx(2),yy(2);  
2  comm.allreduce( mpl::plus<float>(),  
3      xx.data(), yy.data(), mpl::contiguous_layout<float>(2) );
```

The `contiguous_layout` is a 'derived type'; this will be discussed in more detail elsewhere (see note 65 and later). For now, interpret it as a way of indicating the count/type part of a buffer specification.



Implement the layout for a pingpong with an array.



You can pass a C-style array as buffer, requiring a layout:

```
1 // collectarray.cxx
2 float rank2p2p1[2] = { 2*xrank, 2*xrank+1 };
3 mpl::contiguous_layout<float> p2layout(2);
4 comm_world.allreduce(mpl::plus<float>(), rank2p2p1, p2layout);
```



MPL point-to-point routines have a way of specifying the buffer(s) through a *begin* and *end* iterator.

```
1 // sendrange.cxx
2 vector<double> v(15);
3 comm_world.send(v.begin(), v.end(), 1); // send to rank 1
4 comm_world.recv(v.begin(), v.end(), 0); // receive from rank 0
```

Not available for collectives.



The usual reduction operators are given as templated operators:

```
1 float
2     xrank = static_cast<float>( comm_world.rank() ),
3     xreduce;
4 // separate recv buffer
5 comm_world.allreduce(mpl::plus<float>(), xrank, xreduce);
6 // in place
7 comm_world.allreduce(mpl::plus<float>(), xrank);
```

Note the parentheses after the operator. Also note that the operator comes first, not last.

Available: *max*, *min*, *plus*, *multiplies*, *logical\_and*, *logical\_or*, *logical\_xor*, *bit\_and*, *bit\_or*, *bit\_xor*.

*Implementation note: The reduction operator has to be compatible with  $T(T,T)>$*



```
template<typename T , typename F >
void mpl::communicator::allreduce
    ( F,const T &, T & ) const;
    ( F,const T *, T *,
      const contiguous_layout< T > & ) const;
    ( F,T & ) const;
    ( F,T *, const contiguous_layout< T > & ) const;
F : reduction function
T : type
```



```
void mpl::communicator::reduce
    // root, in place
    ( F f,int  root_rank,T &  sendrecvdata ) const
    ( F f,int  root_rank,T *  sendrecvdata,const contiguous_layout< T
    // non-root
    ( F f,int  root_rank,const T &  senddata ) const
    ( F f,int  root_rank,
        const T * senddata,const contiguous_layout< T > & l ) const
    // general
    ( F f,int  root_rank,const T &  senddata,T &  recvdata ) const
    ( F f,int  root_rank,
        const T *  senddata,T *  recvdata,const contiguous_layout< T
```





The broadcast call comes in two variants, with scalar argument and general layout:

```
1  template<typename T >
2  void mpl::communicator::bcast
3      ( int root_rank, T &data ) const;
4  void mpl::communicator::bcast
5      ( int root_rank, T *data, const layout< T > &l ) const;
```

Note that the root argument comes first.



```
template<typename T >
void mpl::communicator::bcast
    ( int  root, T &  data ) const
    ( int  root, T *  data, const layout< T > &  l ) const
```



Gathering (by `communicator::gather`)  
or scattering (by `communicator::scatter`)  
a single scalar takes a scalar argument and a raw array:

```
1 vector<float> v;  
2 float x;  
3 comm_world.scatter(0, v.data(), x);
```

If more than a single scalar is gathered, or scattered into, it becomes necessary to specify a layout:

```
1 vector<float> vrecv(2), vsend(2*nprocs);  
2 mpl::contiguous_layout<float> twonums(2);  
3 comm_world.scatter  
4   (0, vsend.data(), twonums, vrecv.data(), twonums );
```



There is a separate variant for non-root usage of rooted collectives:

```
1 // scangather.cxx
2 if (procno==0) {
3     comm_world.reduce
4         ( mpl::plus<int>(),0,
5           my_number_of_elements,total_number_of_elements );
6 } else {
7     comm_world.reduce
8         ( mpl::plus<int>(),0,my_number_of_elements );
9 }
```



Logically speaking, on every nonroot process, the gather call only has a send buffer. MPL supports this by having two variants that only specify the send data.

```
1  if (procno==0) {
2      vector<int> size_buffer(nprocs);
3      comm_world.gather
4          (
5              0,my_number_of_elements,size_buffer.data()
6          );
7  } else {
8      /*
9       * If you are not the root, do versions with only send buffers
10      */
11      comm_world.gather
12          ( 0,my_number_of_elements );
```



```
void mpl::communicator::gather
( int  root_rank, const T & senddata ) const
( int  root_rank, const T & senddata, T *  recvdata ) const
( int  root_rank, const T * senddata, const layout< T > &  sendl ) c
( int  root_rank, const T * senddata, const layout< T > &  sendl,
    T *  recvdata, const layout< T > &  recvl ) c
```



The in-place variant is activated by specifying only one instead of two buffer arguments.

```
1 float
2   xrank = static_cast<float>( comm_world.rank() ),
3   xreduce;
4 // separate recv buffer
5 comm_world.allreduce(mpl::plus<float>(), xrank, xreduce);
6 // in place
7 comm_world.allreduce(mpl::plus<float>(), xrank);
```

Reducing a buffer requires specification of a `contiguous_layout`:

```
1 // collectbuffer.cxx
2 float
3   xrank = static_cast<float>( comm_world.rank() );
4 vector<float> rank2p2p1{ 2*xrank, 2*xrank+1 }, reduce2p2p1{0,0};
5 mpl::contiguous_layout<float> two_floats(rank2p2p1.size());
6 comm_world.allreduce
7   (mpl::plus<float>(),
8    ↪rank2p2p1.data(), reduce2p2p1.data(), two_floats);
9 if ( iprint )
10   cout << "Got: " << reduce2p2p1.at(0) << ", "
11        << reduce2p2p1.at(1) << endl;
```



The size/displacement arrays for `MPI_Gatherv` / `MPI_Alltoallv` are handled through a `layouts` object, which is basically a vector of `layout` objects.

```
1  mpl::layouts<int> receive_layout;  
2  for ( int iproc=0,loc=0; iproc<nprocs; iproc++ ) {  
3      auto siz = size_buffer.at(iproc);  
4      receive_layout.push_back  
5          ( mpl::indexed_layout<int>( {{ siz,loc }} ) );  
6      loc += siz;  
7  }
```





As in the C/F interfaces, MPL interfaces to the scan routines have the same calling sequences as the 'Allreduce' routine.





- Let each process compute a random value  $n_{\text{local}}$ , and allocate an array of that length. Define

$$N = \sum n_{\text{local}}$$

- Fill the array with consecutive integers, so that all local arrays, laid end-to-end, contain the numbers  $0 \cdots N - 1$ . (See figure 34.)



Take the code from exercise 2 and extend it to gather all local buffers onto rank zero. Since the local arrays are of differing lengths, this requires `MPI_Gatherv`.

How do you construct the lengths and displacements arrays?



Arithmetic: *plus, multiplies, max, min.*

Logic: *logical\_and, logical\_or, logical\_xor.*

Bitwise: *bit\_and, bit\_or, bit\_xor.*



A user-defined operator can be a templated class with an `operator()`.

Example:

```
1 // reduceuser.cxx
2 template<typename T>
3 class lcm {
4 public:
5     T operator()(T a, T b) {
6         T zero=T();
7         T t((a/gcd(a, b))*b);
8         if (t<zero)
9             return -t;
10        return t;
11    }

1 comm_world.reduce(lcm<int>(), 0, v, result);
```

(The templated class can be a lambda expression)



You can also do the reduction by lambda:

```
1 comm_world.reduce
2   ( [] (int i,int j) -> int { return i+j; },
3       0,data );
```



Nonblocking collectives have the same argument list as the corresponding blocking variant, except that instead of a `void` result, they return an `irequest`. (See 54)

Wait calls are methods of the `irequest` object.

```
1 // ireducescalar.cxx
2 float x{1.}, sum;
3 auto reduce_request =
4     comm_world.ireduce(mpl::plus<float>(), 0, x, sum);
5 reduce_request.wait();
6 if (comm_world.rank()==0) {
7     std::cout << "sum = " << sum << '\n';
8 }
```





# Part III

## Point-to-point communication



- Scalar data type is handled through templating (and 'argument-dependent-lookup'): derived by the compiler.
- Count  $> 1$  is declared in the layout datatype.



MPL uses a default value for the tag, and it can deduce the type of the buffer. Sending a scalar becomes:

```
1 // sendscalar.cxx
2 if (comm_world.rank()==0) {
3     double pi=3.14;
4     comm_world.send(pi, 1); // send to rank 1
5     cout << "sent: " << pi << '\n';
6 } else if (comm_world.rank()==1) {
7     double pi=0;
8     comm_world.recv(pi, 0); // receive from rank 0
9     cout << "got : " << pi << '\n';
10 }
```

(See also note 17.)



MPL can send *static arrays* without further layout specification:

```
1 // sendarray.cxx
2 double v[2][2][2];
3 comm_world.send(v, 1); // send to rank 1
4 comm_world.recv(v, 0); // receive from rank 0
```

Sending vectors uses a general mechanism:

```
1 // sendbuffer.cxx
2 std::vector<double> v(8);
3 mpl::contiguous_layout<double> v_layout(v.size());
4 comm_world.send(v.data(), v_layout, 1); // send to rank 1
5 comm_world.recv(v.data(), v_layout, 0); // receive from rank 0
```

(See also note 18.)



```
template<typename T >
void mpl::communicator::send
    ( const T scalar&,int dest,tag = tag(0) ) const
    ( const T *buffer,const layout< T > &,int dest,tag = tag(0) ) const
    ( iterT begin,iterT end,int dest,tag = tag(0) ) const
T : scalar type
begin : begin iterator
end : end iterator
```



```
template<typename T >
status mpl::communicator::recv
    ( T &,int,tag = tag(0) ) const inline
    ( T *,const layout< T > &,int,tag = tag(0) ) const
    ( iterT begin,iterT end,int source, tag t = tag(0) ) const
```



MPL differs from other Application Programmer Interfaces (APIs) in its treatment of tags: a tag is not directly an integer, but an object of class `tag_t`.

```
1 // sendrecv.cxx
2 mpl::tag_t t0(0);
3 comm_world.sendrecv
4   ( mydata,sendto,t0,
5     leftdata,recvfrom,t0 );
```

The `tag_t` class has a couple of methods such as `mpl::tag_t::any()` (for the `MPI_ANY_TAG` wildcard in receive calls) and `mpl::tag_t::up()` (maximal tag, found from the `MPI_TAG_UB` attribute).



Tag are *int* or an *enum* typ:

```
1     template<typename T >
2     tag_t (T t);
3     tag_t (int t);
4
```

Example:

```
1  // inttag.cxx
2  enum class pingpongtag : int { ping=1, pong=2 };
3  int pinger = 0, ponger = world.size()-1;
4  if (world.rank()==pinger) {
5      world.send(x, ponger, pingpongtag::ping);
6      world.recv(x, ponger, pingpongtag::pong);
7  } else if (world.rank()==ponger) {
8      world.recv(x, pinger, pingpongtag::ping);
9      world.send(x, pinger, pingpongtag::pong);
10 }
```





The constant `mpl::any_source` equals `MPI_ANY_SOURCE` (by *constexpr*).



The send-recv call in MPL has the same possibilities for specifying the send and receive buffer as the separate send and recv calls: scalar, layout, iterator. However, out of the nine conceivably possible routine signatures, only the versions are available where the send and receive buffer are specified the same way. Also, the send and receive tag need to be specified; they do not have default values.

```
1 // sendrecv.cxx
2 mpl::tag_t t0(0);
3 comm_world.sendrecv
4   ( mydata,sendto,t0,
5     leftdata,recvfrom,t0 );
```

```
1 // sendrecvarray.cxx
2 mpl::tag_t t0(0);
3 mpl::contiguous_layout<double>
4   ↪ twofloats(2);
5 comm_world.sendrecv
6   ( mydata,twofloats,sendto,t0,
7     ↪ leftdata,twofloats,recvfrom,t0
8     ↪ );
```



The `mpl::status_t` object is created by the receive (or wait) call:

```
1  mpl::contiguous_layout<double> target_layout(count);  
2  mpl::status_t recv_status =  
3      comm_world.recv(target.data(), target_layout, the_other);  
4  recv_count = recv_status.get_count<double>();
```



The `status` object can be queried:

```
1 int source = recv_status.source();
```



The `get_count` function is a method of the status object. The argument type is handled through templating:

```
1 // recvstatus.cxx
2 double pi=0;
3 auto s = comm_world.recv(pi, 0); // receive from rank 0
4 int c = s.get_count<double>();
5 std::cout << "got : " << c << " scalar(s): " << pi << '\n';
```



Nonblocking routines have an `irequest` as function result. Note: not a parameter passed by reference, as in the C interface. The various wait calls are methods of the `irequest` class.

```
1  double recv_data;  
2  mpl::irequest recv_request =  
3      comm_world.irecv( recv_data, sender );  
4  recv_request.wait();
```

You can not default-construct the request variable:

```
1  // DOES NOT COMPILE:  
2  mpl::irequest recv_request;  
3  recv_request = comm.irecv( ... );
```

This means that the normal sequence of first declaring, and then filling in, the request variable is not possible.

*Implementation note: The wait call always returns a `status_t` object; not assigning it means that the destructor is called on it.*



Instead of an array of requests, use an `irequest_pool` object, which acts like a vector of requests, meaning that you can *push* onto it.

```
1 // irecvsource.cxx
2 mpl::irequest_pool recv_requests;
3 for (int p=0; p<nprocs-1; p++) {
4     recv_requests.push( comm_world.irecv( recv_buffer[p], p ) );
5 }
```

You can not declare a pool of a fixed size and assign elements. (Why not? Can you find a way around it?)



The `irequest_pool` class has methods `waitany`, `waitall`, `testany`, `testall`, `waitsome`, `testsome`.

The 'any' methods return a `std::pair<mpl::test_result, size_t>`, where the `test_result` is an *enum class* with values:

- `completed` (for any/some/all completions),
- `no_completed` (for none),
- `no_active_requests` (if no more requests active).

```
1  auto [success, index] = recv_requests.waitany();
2  if ( success==mpl::test_result::completed ) {
3      auto recv_status = recv_requests.get_status(index);
```





Create two distributed arrays of positive integers. Take the set difference of the two: the first array needs to be transformed to remove from it those numbers that are in the second array.

How could you solve this with an `MPI_Allgather` call? Why is it not a good idea to do so? Solve this exercise instead with a circular bucket brigade algorithm.



Creating and attaching a buffer is done through `bsend_buffer` and a support routine `bsend_size` helps in calculating the buffer size:

```
1 // bufring.cxx
2 vector<float> sbuf(BUFLEN), rbuf(BUFLEN);
3 int size{
4     ↪comm_world.bsend_size<float>(mpl::contiguous_layout<float>(BUFLEN))
5     ↪};
6 mpl::bsend_buffer buff(size);
7 comm_world.bsend(sbuf.data(),mpl::contiguous_layout<float>(BUFLEN),
8     ↪next);
```

Constant: `mpl::bsend_overhead` is `constexpr`'d to the MPI constant `MPI_BSEND_OVERHEAD`.



There is a separate attach routine, but normally this is called by the constructor of the `bsend_buffer`. Likewise, the detach routine is called in the buffer destructor.

```
1 void mpl::environment::buffer_attach (void *buff, int size);  
2 std::pair< void *, int > mpl::environment::buffer_detach ();
```



MPL returns a `prequest` from persistent 'init' routines, rather than an `irequest` (MPL note 54):

```
1  template<typename T >  
2  prequest send_init (const T &data, int dest, tag t=tag(0)) const;
```

Likewise, there is a `prequest_pool` instead of an `irequest_pool` (note 55).



# Part IV

## Derived Datatypes



MPL mostly handles datatypes through subclasses of the `layout` class. Layouts are MPL routines are templated over the data type.

```
1 // sendlong.cxx
2 mpl::contiguous_layout<long long> v_layout(v.size());
3 comm.send(v.data(), v_layout, 1); // send to rank 1
```

Also works with complex of float and double.

The data types, where MPL can infer their internal representation, are enumeration types, C arrays of constant size and the template classes `std::array`, `std::pair` and `std::tuple` of the C++ Standard Template Library. The only limitation is, that the C array and the mentioned template classes hold data elements of types that can be sent or received by MPL.



Should you need the `MPI_Datatype` object contained in an MPL layout, there is an access function `native_handle`.



Large buffer communication is supported:

```
1 // bigint.cxx
2 size_t s = static_cast<size_t>(1)<<(s0-1);
3 vector<char> buffer(s);
4 mpl::contiguous_layout<char> buffersize(s);
5 comm_world.send( buffer.data(), buffersize, processB );
```





In MPL type creation routines are in the main namespace, templated over the datatypes.

```
1 // vector.cxx
2 vector<double>
3   source(stride*count);
4 if (procno==sender) {
5   mpl::strided_vector_layout<double>
6     newvectortype(count,1,stripe);
7   comm_world.send
8     (source.data(),newvectortype,the_other);
9 }
```

The commit call is part of the type creation, and freeing is done in the destructor.



The MPL interface makes extensive use of `contiguous_layout`, as it is the main way to declare a nonscalar buffer; see note 18.



Contiguous layouts can only use predefined types or other contiguous layouts as their ‘old’ type. To make a contiguous type for other layouts, use `vector_layout`:

```
1 // contiguous.cxx
2 mpl::contiguous_layout<int> type1(7);
3 mpl::vector_layout<int> type2(8,type1);
```

(Contrast this with `strided_vector_layout`; note 68.)



MPL has the `strided_vector_layout` class as equivalent of the vector type:

```
1 // vector.cxx
2 vector<double>
3   source(stride*count);
4 if (procno==sender) {
5   mpl::strided_vector_layout<double>
6     newvectortype(count,1,stripe);
7   comm_world.send
8     (source.data(),newvectortype,the_other);
9 }
```

(See note 67 for nonstrided vectors.)



MPL point-to-point routines have a way of specifying the buffer(s) through a *begin* and *end* iterator.

```
1 // sendrange.cxx
2 vector<double> v(15);
3 comm_world.send(v.begin(), v.end(), 1); // send to rank 1
4 comm_world.recv(v.begin(), v.end(), 0); // receive from rank 0
```

Not available for collectives.



Noncontiguous iterable objects can be send with a `iterator_layout`:

```
1  std::list<int> v(20, 0);  
2  mpl::iterator_layout<int> l(v.begin(), v.end());  
3  comm_world.recv(&(*v.begin()), l, 0);
```



The templated `subarray_layout` class is constructed from a vector of triplets of global size / subblock size / first coordinate.

```
1  mpl::subarray_layout<int>(  
2    { {ny, ny_1, ny_0}, {nx, nx_1, nx_0} }  
3  );
```



In MPL, the `indexed_layout` is based on a vector of 2-tuples denoting block length / block location.

```
1 // indexed.cxx
2 const int count = 5;
3 mpl::contiguous_layout<int>
4     fiveints(count);
5 mpl::indexed_layout<int>
6     indexed_where{ { {1,2}, {1,3}, {1,5}, {1,7}, {1,11} } };
7
8 if (procno==sender) {
9     comm_world.send( source_buffer.data(),indexed_where, receiver );
10 } else if (procno==receiver) {
11     auto recv_status =
12         comm_world.recv( target_buffer.data(),fiveints, sender );
13     int recv_count = recv_status.get_count<int>();
14     assert(recv_count==count);
15 }
```





For the case where all block lengths are the same, use

`indexed_block_layout`:

```
1 // indexedblock.cxx
2 mpl::indexed_block_layout<int>
3   indexed_where( 1, {2,3,5,7,11} );
4 comm_world.send( source_buffer.data(), indexed_where, receiver );
```



One could describe the MPI struct type as a collection of displacements, to be applied to any set of items that conforms to the specifications. An MPL `heterogeneous_layout` on the other hand, incorporates the actual data. Thus you could write

```
1 // structscalar.cxx
2 char c; double x; int i;
3 if (procno==sender) {
4     c = 'x'; x = 2.4; i = 37; }
5 mpl::heterogeneous_layout object( c,x,i );
6 if (procno==sender)
7     comm_world.send( mpl::absolute,object,receiver );
8 else if (procno==receiver)
9     comm_world.recv( mpl::absolute,object,sender );
```

Here, the `absolute` indicates the lack of an implicit buffer: the layout is absolute rather than a relative description.



More complicated data than scalars takes more work:

```
1 // struct.cxx
2 char c; vector<double> x(2); int i;
3 if (procno==sender) {
4     c = 'x'; x[0] = 2.7; x[1] = 1.5; i = 37; }
5 mpl::heterogeneous_layout object
6     ( c,
7       mpl::make_absolute(x.data(),mpl::vector_layout<double>(2)),
8       i );
9 if (procno==sender) {
10     comm_world.send( mpl::absolute,object,receiver );
11 } else if (procno==receiver) {
12     comm_world.recv( mpl::absolute,object,sender );
13 }
```

Note the `make_absolute` in addition to `absolute` mentioned above.



Resizing a datatype does not give a new type, but does the resize 'in place':

```
1 void layout::resize(ssize_t lb, ssize_t extent);
```



# Part V

## Communicator manipulations



## Code:

```
1  const mpl::communicator &comm =
2      mpl::environment::comm_world();
3  MPI_Comm
4      world_extract = comm.native_handle(),
5      world_given = MPI_COMM_WORLD;
6  int result;
7  MPI_Comm_compare(world_extract, world_given, &result);
8  cout << "Compare raw comms: " << "\n"
9      << "identical: "
10     << (result==MPI_IDENT) << "\n"
11     << "congruent: "
12     << (result==MPI_CONGRUENT) << "\n"
13     << "unequal : "
14     << (result==MPI_UNEQUAL) << "\n";
```

## Output:

```
1  Compare raw comms:
2  identical: true
3  congruent: false
4  unequal : false
```



MPL does not have a routine for setting the error handler. Instead, use the *native\_handle* method to retrieve the embedded communiator.



In MPL, splitting a communicator is done as one of the overloads of the communicator constructor;

```
1 // commsplit.cxx
2 // create sub communicator modulo 2
3 int color2 = procno % 2;
4 mpl::communicator comm2
5     ( mpl::communicator::split, comm_world, color2 );
6 auto procno2 = comm2.rank();
7
8 // create sub communicator modulo 4 recursively
9 int color4 = procno2 % 2;
10 mpl::communicator
11     comm4( mpl::communicator::split, comm2, color4 );
12 auto procno4 = comm4.rank();
```

*Implementation note: The `communicator::split` identifier is an object of class `communicator::split_tag`, itself is an otherwise empty subclass of `communicator`:*

```
1 class split_tag {};
```

```
2 static constexpr split_tag split{};
```





Similar to ordinary communicator splitting (slide 80):

```
communicator::split_shared.
```

```
1 // commsplitttype.cxx
2 mpl::communicator shared_comm
3   ( mpl::communicator::split_shared_memory, world_comm );
4 int
5   onnode_procno = shared_comm.rank(),
6   onnode_nprocs = shared_comm.size();
```

But note: shared memory is currently not available, since windows are not (yet) implemented.



# Part VI

## Process topologies



There is a separate class `cartesian_communicator`.

Additionally, there is a class `dimensions` that describes the shape of the Cartesian grid and its periodicity. The `size(int)` method retrieves the size in the given dimension.



The `dims_create` routine takes a `dimensions` object with only the dimensionality specified, and creates one with the sizes filled in.

```
1  mpl::cartesian_communicator::dimensions brick(3);  
2  brick = mpl::dims_create(nprocs,brick);
```

To have certain dimensions be periodic, the initial `dimensions` object needs to be create with periodicity values `periodic` or `non_periodic`.

```
1  mpl::cartesian_communicator::dimensions pbrick  
2    ( { mpl::cartesian_communicator::non_periodic,  
3        mpl::cartesian_communicator::periodic,  
4        mpl::cartesian_communicator::non_periodic } );  
5  pbrick = mpl::dims_create(nprocs,pbrick);
```



The actual Cartesian communicator has a constructor that takes a *dimensions* object as input.

```
1  mpl::cartesian_communicator cart_comm( comm_world,brick );
```



The *dimensions* object can be extracted from the communicator

```
1  mpl::cartesian_communicator::dimensions  
2      dimensions = cart_comm.get_dimensions();
```

after which dimensions and periodicities can be extracted:

```
1  // cartcoord.cxx  
2  int dsize = dimensions.size(idim);  
3  auto p = dimensions.periodicity(idim);
```



The `coordinates` method of the cartesian communicator returns a vector-like object describing the process coordinate:

```
1  for ( int ip=0; ip<nprocs; ip++ ) {  
2      mpl::cartesian_communicator::vector  
3          coord = cart_comm.coordinates(ip);  
4      print("[{:2}] coord: [",ip);  
5      for ( int id=0; id<dim; id++ )  
6          print("{},",coord[id]);  
7      print("]\n");
```



The routine `cartesian_communicator::shift` takes a dimension and a direction, and gives the source and destination as a `shifted_ranks` object, which is basically a tuple of two integers:

```
1  int pred,succ;
2  mpl::shift_ranks shifted = cart_comm.shift
3      ( /* dim: */ 1,/* up: */ 1 );
4  pred = shifted.source; succ = shifted.destination;
5  print("non-periodic: src={}, tgt={}\n",pred,succ);
```





The constructor `dist_graph_communicator`

```
1  dist_graph_communicator  
2      (const communicator &old_comm, const source_set &ss,  
3          const dest_set &ds, bool reorder = true);
```

is a wrapper around `MPI_Dist_graph_create_adjacent`.



Methods *indegree*, *outdegree* are wrappers around `MPI_Dist_graph_neighbors_count`. Sources and targets can be queried with *inneighbors* and *outneighbors*, which are wrappers around `MPI_Dist_graph_neighbors`.



# Part VII

## Other



The timing routines `wtime` and `wtick` and `wtime_is_global` are environment methods:

```
1  double  mpl::environment::wtime ();  
2  double  mpl::environment::wtick ();  
3  bool  mpl::environment::wtime_is_global ();
```



There is an *info* object in the *mpl* namespace:

```
1  mpl::info infoobject; // default constructor
```

Sample methods:

```
1  void set(const std::string &key, const std::string &value);  
2  [[nodiscard]] std::optional<std::string> value(const std::string  
    ↪ &key) const;
```



MPL always calls `MPI_Init_thread` requesting the highest level `MPI_THREAD_MULTIPLE`.

```
1  enum mpl::threading_modes {  
2      mpl::threading_modes::single = MPI_THREAD_SINGLE,  
3      mpl::threading_modes::funneled = MPI_THREAD_FUNNELED,  
4      mpl::threading_modes::serialized = MPI_THREAD_SERIALIZED,  
5      mpl::threading_modes::multiple = MPI_THREAD_MULTIPLE  
6  };  
7  threading_modes mpl::environment::threading_mode ();  
8  bool mpl::environment::is_thread_main ();
```



MPL is not a full MPI implementation

- One-sided communication
- Shared memory
- Process management

