

# Cell-MPI

Mastering the Cell Broadband Engine architecture through a Boost based parallel communication library

**digiteo**  
Research in science and technology of information



cea **list**

Sebastian Schaetz, Joel Falcou,  
Lionel Lacassagne

Digiteo Foundation, LRI - University  
Paris South XI, CEA LIST

May 17, 2011

## Why a talk about a library for the Cell in 2011

- Heterogeneous and composable architectures are not uncommon, powerful and worth studying.

## Why a talk about a library for the Cell in 2011

- Heterogeneous and composable architectures are not uncommon, powerful and worth studying.
- We present useful concepts that apply to all of them.

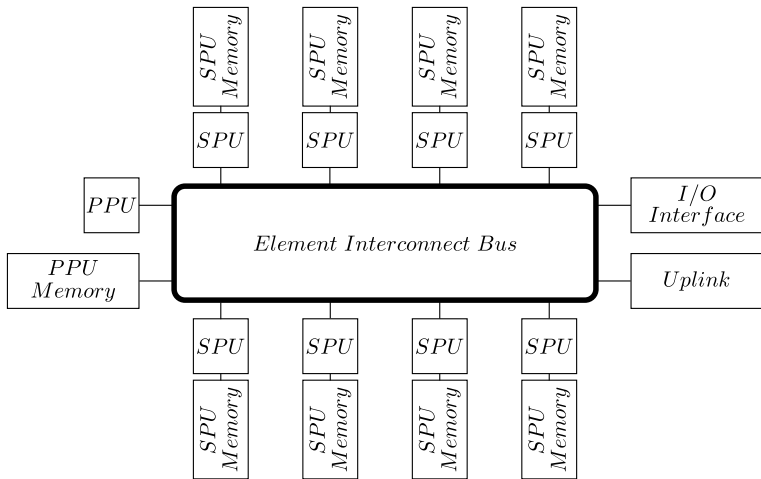
## Why a talk about a library for the Cell in 2011

- Heterogeneous and composable architectures are not uncommon, powerful and worth studying.
- We present useful concepts that apply to all of them.
- We illustrate the lessons we learned as we used Boost libraries on a constricted platform and

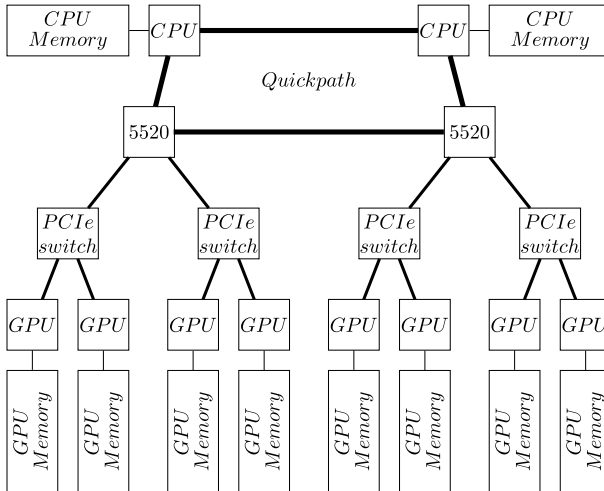
## Why a talk about a library for the Cell in 2011

- Heterogeneous and composable architectures are not uncommon, powerful and worth studying.
- We present useful concepts that apply to all of them.
- We illustrate the lessons we learned as we used Boost libraries on a constricted platform and
- elaborate what choices we had to make and why we made them as we created a Boost-like library for this platform.

# Cell Broadband Engine - Schematic



# A similar architecture - Multi-GPU Schematic



## Cell Broadband Engine - The good stuff

- Power architecture core paired with up to 8 streamlined vector co-processors: 204.8 GFlops/s (single) 102.4 GFlops/s (double)
- High data transfer bandwidth: theoretical 204.8 GB/s
- Good performance/watt (0.87 double precision GFlops/s per Watt for IBM BladeCenter QS22)

Due to these advantages, the CBE is a good fit for multimedia and vector processing applications as well as scientific computation.



## Cell Broadband Engine - The bad stuff

- Distributed system on one chip, explicit communication necessary
- SPE Memory limitations
  - 256kB for code and data per SPE
  - no overflow detection
- Communication intricacies
  - packet size
  - address alignment
  - explicit DMA
- Optimization for speed
  - SIMD (assembler-like)
  - convoluted pipeline mechanism

Due to these restrictions, the complexity of programming the CBE is comparable to writing code for embedded systems.

## Writing code for the CBE

- PPE and SPE entry points in separate main functions

## Writing code for the CBE

- PPE and SPE entry points in separate main functions
- Compilers: ppu-gcc, spu-gcc

## Writing code for the CBE

- PPE and SPE entry points in separate main functions
- Compilers: ppu-gcc, spu-gcc
- SPE object file passed to ppu-embedspu to generate library exports symbol that is accessible from PPE code

## Writing code for the CBE

- PPE and SPE entry points in separate main functions
- Compilers: ppu-gcc, spu-gcc
- SPE object file passed to ppu-embedspu to generate library exports symbol that is accessible from PPE code
- PPE creates thread for each SPE and loads the symbol

## Writing code for the CBE

- PPE and SPE entry points in separate main functions
- Compilers: ppu-gcc, spu-gcc
- SPE object file passed to ppu-embedspu to generate library exports symbol that is accessible from PPE code
- PPE creates thread for each SPE and loads the symbol
- Argument passed to thread is accessible in SPE through argument vector

## Writing code for the CBE

- PPE and SPE entry points in separate main functions
- Compilers: ppu-gcc, spu-gcc
- SPE object file passed to ppu-embedspu to generate library exports symbol that is accessible from PPE code
- PPE creates thread for each SPE and loads the symbol
- Argument passed to thread is accessible in SPE through argument vector
- Usual approach: argument is pointer to structure in main memory; structure is loaded to SPE through explicit DMA call:

## Writing code for the CBE

- PPE and SPE entry points in separate main functions
- Compilers: ppu-gcc, spu-gcc
- SPE object file passed to ppu-embedspu to generate library exports symbol that is accessible from PPE code
- PPE creates thread for each SPE and loads the symbol
- Argument passed to thread is accessible in SPE through argument vector
- Usual approach: argument is pointer to structure in main memory; structure is loaded to SPE through explicit DMA call:

```
1  /* DMA control block information from system memory. */
2  mfc_get((void*)&parms, parm_ptr, (sizeof(parms)+15)&~0xF, tag, td, rd);
3  mfc_write_tag_mask(1<<tag);
4  mfc_read_tag_status_all(); /* Wait for DMA to complete */
```



## Writing code for the CBE - continued

- "Getting started" can be tedious when developing for the CELL since compilation procedure and startup are not trivial

## Writing code for the CBE - continued

- "Getting started" can be tedious when developing for the CELL since compilation procedure and startup are not trivial
- CMake to the rescue:

## Writing code for the CBE - continued

- "Getting started" can be tedious when developing for the CELL since compilation procedure and startup are not trivial
- CMake to the rescue:
  - Great tool to simplify basically any build-related steps

## Writing code for the CBE - continued

- "Getting started" can be tedious when developing for the CELL since compilation procedure and startup are not trivial
- CMake to the rescue:
  - Great tool to simplify basically any build-related steps
  - Find all required libraries and binaries on the system

## Writing code for the CBE - continued

- "Getting started" can be tedious when developing for the CELL since compilation procedure and startup are not trivial
- CMake to the rescue:
  - Great tool to simplify basically any build-related steps
  - Find all required libraries and binaries on the system
  - Low-level macros: `ACTIVATE_PPE_COMPILER()`, `ACTIVATE_SPE_COMPILER()`

## Writing code for the CBE - continued

- "Getting started" can be tedious when developing for the CELL since compilation procedure and startup are not trivial
  
- CMake to the rescue:
  - Great tool to simplify basically any build-related steps
  - Find all required libraries and binaries on the system
  - Low-level macros: `ACTIVATE_PPE_COMPILER()`, `ACTIVATE_SPE_COMPILER()`
  - `ADD_SPE_MODULE(target symbol file0 file1 ... fileN)`

## Writing code for the CBE - continued

- "Getting started" can be tedious when developing for the CELL since compilation procedure and startup are not trivial
  
- CMake to the rescue:
  - Great tool to simplify basically any build-related steps
  - Find all required libraries and binaries on the system
  - Low-level macros: `ACTIVATE_PPE_COMPILER()`, `ACTIVATE_SPE_COMPILER()`
  - `ADD_SPE_MODULE(target symbol file0 file1 ... fileN)`
  
- C++ and Boost to the rescue:

## Writing code for the CBE - continued

- "Getting started" can be tedious when developing for the CELL since compilation procedure and startup are not trivial
  
- CMake to the rescue:
  - Great tool to simplify basically any build-related steps
  - Find all required libraries and binaries on the system
  - Low-level macros: `ACTIVATE_PPE_COMPILER()`, `ACTIVATE_SPE_COMPILER()`
  - `ADD_SPE_MODULE(target symbol file0 file1 ... fileN)`
  
- C++ and Boost to the rescue:
  - Wrap recurring boilerplate code in clearly laid out functions and classes



## Writing code for the CBE - continued

- "Getting started" can be tedious when developing for the CELL since compilation procedure and startup are not trivial
- CMake to the rescue:
  - Great tool to simplify basically any build-related steps
  - Find all required libraries and binaries on the system
  - Low-level macros: `ACTIVATE_PPE_COMPILER()`, `ACTIVATE_SPE_COMPILER()`
  - `ADD_SPE_MODULE(target symbol file0 file1 ... fileN)`
- C++ and Boost to the rescue:
  - Wrap recurring boilerplate code in clearly laid out functions and classes
  - A kernel function should be declared and behave like a free function

## Cell-MPI Bootstrapping

- Launching a kernel function passing a data structure

```
struct mydatastruct { int x; int y; int z; };
```

## Cell-MPI Bootstrapping

- Launching a kernel function passing a data structure

```
struct mydatastruct { int x; int y; int z; };
```

- Kernel is defined with:

```
1 | BEGIN_CELL_KERNEL()  
2 | {  
3 |     mydatastruct * ptr;  
4 |     SPE_Custom(ptr);  
5 |     RETURN((ptr->x + ptr->y) * ptr->z);  
6 | }  
7 | END_CELL_KERNEL()
```

## Cell-MPI Bootstrapping

- Launching a kernel function passing a data structure

```
struct mydatastruct { int x; int y; int z; };
```

- Kernel is defined with:

```
1 | BEGIN_CELL_KERNEL()  
2 | {  
3 |     mydatastruct * ptr;  
4 |     SPE_Custom(ptr);  
5 |     RETURN((ptr->x + ptr->y) * ptr->z);  
6 | }  
7 | END_CELL_KERNEL()
```

- In PPE code the kernel is registered with  
PPE\_REGISTER\_KERNEL(kernel);

## Cell-MPI Bootstrapping

- Launching a kernel function passing a data structure

```
struct mydatastruct { int x; int y; int z; };
```

- Kernel is defined with:

```
1 | BEGIN_CELL_KERNEL()  
2 | {  
3 |     mydatastruct * ptr;  
4 |     SPE_Custom(ptr);  
5 |     RETURN((ptr->x + ptr->y) * ptr->z);  
6 | }  
7 | END_CELL_KERNEL()
```

- In PPE code the kernel is registered with  
PPE\_REGISTER\_KERNEL(kernel);
- The runtime is initialized with PPE\_Init();

## Cell-MPI Bootstrapping - continued

- The kernel is then called asynchronously:

```
1 | mydatastruct mydata(1, 5, 7);  
2 | PPE_Run(kernel, mydata);
```

## Cell-MPI Bootstrapping - continued

- The kernel is then called asynchronously:

```
1 | mydatastruct mydata(1, 5, 7);  
2 | PPE_Run(kernel, mydata);
```

- The PPE can wait for kernel completion: PPE\_Sync();

## Cell-MPI Bootstrapping - continued

- The kernel is then called asynchronously:

```
1 | mydatastruct mydata(1, 5, 7);  
2 | PPE_Run(kernel, mydata);
```

- The PPE can wait for kernel completion: PPE\_Sync();
- and access the kernels return value:

```
1 | int returnvalues[CBE_MPI_NUM_SPE];  
2 | PPE_Return(&returnvalues[0]);
```



## Cell-MPI Bootstrapping - continued

- The kernel is then called asynchronously:

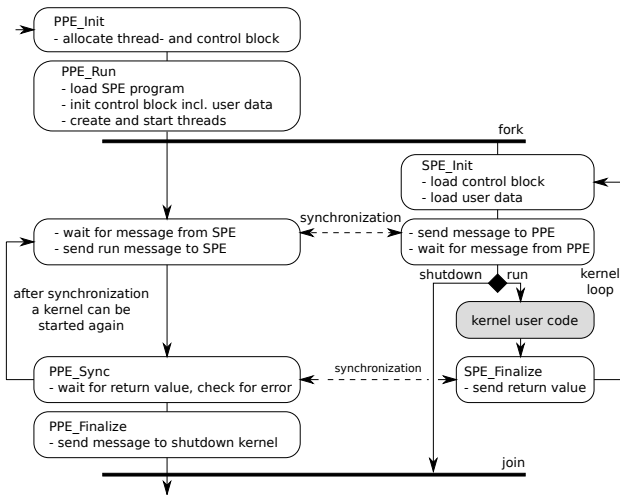
```
1 | mydatastruct mydata(1, 5, 7);  
2 | PPE_Run(kernel, mydata);
```

- The PPE can wait for kernel completion: `PPE_Sync()`;
- and access the kernels return value:

```
1 | int returnvalues[CBE_MPI_NUM_SPE];  
2 | PPE_Return(&returnvalues[0]);
```

- The runtime is finalized with `PPE_Finalize()`;

# Cell-MPI Bootstrapping Mechanism



## Cell-MPI Bootstrapping - Boostified

- A kernel can be declared in both PPE and SPE code with:

```
1 | SPE_FUNCTION(kernel_, kernel, (int x) (int y) (int z) );
```

## Cell-MPI Bootstrapping - Boostified

- A kernel can be declared in both PPE and SPE code with:

```
1 | SPE_FUNCTION(kernel_, kernel, (int x) (int y) (int z) );
```

- and implemented as a free function in SPE code

```
1 | int kernel(int x, int y, int z)  
2 | {  
3 |     return (x+y) * z;  
4 | }
```

## Cell-MPI Bootstrapping - Boostified

- A kernel can be declared in both PPE and SPE code with:

```
1 | SPE_FUNCTION(kernel_, kernel, (int x) (int y) (int z) );
```

- and implemented as a free function in SPE code

```
1 | int kernel(int x, int y, int z)  
2 | {  
3 |     return (x+y) * z;  
4 | }
```

- It can then be called as a free function from PPE code:

```
1 | int * returnvalues = kernel(2, 5, 7);
```

## Cell-MPI Bootstrapping - Boostified

- A kernel can be declared in both PPE and SPE code with:

```
1 | SPE_FUNCTION(kernel_, kernel, (int x) (int y) (int z) );
```

- and implemented as a free function in SPE code

```
1 | int kernel(int x, int y, int z)  
2 | {  
3 |     return (x+y) * z;  
4 | }
```

- It can then be called as a free function from PPE code:

```
1 | int * returnvalues = kernel(2, 5, 7);
```

- or asynchronously:

```
1 | kernel_async(2, 5, 7);  
2 | PPE_Sync();
```

So we do C++ but...

The architecture forces some restrictions especially on the SPE part of the library:

- Compilation without run-time type information

So we do C++ but...

The architecture forces some restrictions especially on the SPE part of the library:

- Compilation without run-time type information
- No dynamic memory allocation for predictable footprint



## So we do C++ but...

The architecture forces some restrictions especially on the SPE part of the library:

- Compilation without run-time type information
- No dynamic memory allocation for predictable footprint
- Custom, lightweight STL compatible allocators

## So we do C++ but...

The architecture forces some restrictions especially on the SPE part of the library:

- Compilation without run-time type information
- No dynamic memory allocation for predictable footprint
- Custom, lightweight STL compatible allocators
- Exception handling deactivated

## Exception emulation

Due to architecture limitations we emulate exceptions:

- An exception stops the kernel and notifies the PPE

## Exception emulation

Due to architecture limitations we emulate exceptions:

- An exception stops the kernel and notifies the PPE
- Only an error code is "thrown":

```
1 | #define THROW(errno) {spe_errno = errno ; SPE_Finalize(-1); exit(0);}
```

## Exception emulation

Due to architecture limitations we emulate exceptions:

- An exception stops the kernel and notifies the PPE
- Only an error code is "thrown":

```
1 | #define THROW(errno) {spe_errno = errno ; SPE_Finalize(-1); exit(0);}
```

- The PPE translates the error code into real exceptions:

```
1 | struct spe_error_bundle  
2 | { std::vector<spe_error_data> exception_info; };
```

## Exception emulation

Due to architecture limitations we emulate exceptions:

- An exception stops the kernel and notifies the PPE
- Only an error code is "thrown":

```
1 | #define THROW(errno) {spe_errno = errno ; SPE_Finalize(-1); exit(0);}
```

- The PPE translates the error code into real exceptions:

```
1 | struct spe_error_bundle  
2 | { std::vector<spe_error_data> exception_info; };
```

```
1 | typedef boost::error_info<struct tag_spe_error_info_bundle,  
2 |     spe_error_bundle> spe_error_info_bundle;
```

## Exception emulation

Due to architecture limitations we emulate exceptions:

- An exception stops the kernel and notifies the PPE
- Only an error code is "thrown":

```
1 | #define THROW(errno) {spe_errno = errno ; SPE_Finalize(-1); exit(0);}
```

- The PPE translates the error code into real exceptions:

```
1 | struct spe_error_bundle  
2 | { std::vector<spe_error_data> exception_info; };
```

```
1 | typedef boost::error_info<struct tag_spe_error_info_bundle,  
2 |     spe_error_bundle> spe_error_info_bundle;
```

```
1 | struct spe_runtime_exception : virtual boost::exception {};
```

## Exception emulation - continued

- If desired SPE exceptions can interrupt PPE execution



## Exception emulation - continued

- If desired SPE exceptions can interrupt PPE execution
- To define errors we use the same trick as in boosted bootstrapping:

## Exception emulation - continued

- If desired SPE exceptions can interrupt PPE execution
- To define errors we use the same trick as in boosted bootstrapping:
  - 1 | ERROR(MPI\_TAG\_MISMATCH, 7, "Send receive tag mismatch")
  - 2 | ERROR(BOOST\_FUNCTION\_BAD\_CALL, 12, "Bad boost function call")
  - 3 | ERROR(BAD\_ALLOC, 14, "bad alloc")

## Exception emulation - continued

- If desired SPE exceptions can interrupt PPE execution
- To define errors we use the same trick as in boosted bootstrapping:

```
1 | ERROR(MPI_TAG_MISMATCH, 7, "Send receive tag mismatch")  
2 | ERROR(BOOST_FUNCTION_BAD_CALL, 12, "Bad boost function call")  
3 | ERROR(BAD_ALLOC, 14, "bad alloc")
```

- Compiled with the SPE compiler (`#ifdef _SPE_`) generates:

```
1 | enum { MPI_TAG_MISMATCH = 7, BOOST_FUNCTION_BAD_CALL = 12,  
2 |         BAD_ALLOC = 14 };
```

## Exception emulation - continued

- If desired SPE exceptions can interrupt PPE execution
- To define errors we use the same trick as in boosted bootstrapping:

```

1 | ERROR(MPI_TAG_MISMATCH, 7, "Send receive tag mismatch")
2 | ERROR(BOOST_FUNCTION_BAD_CALL, 12, "Bad boost function call")
3 | ERROR(BAD_ALLOC, 14, "bad alloc")

```

- Compiled with the SPE compiler (`#ifdef _SPE_`) generates:

```

1 | enum { MPI_TAG_MISMATCH = 7, BOOST_FUNCTION_BAD_CALL = 12,
2 |        BAD_ALLOC = 14 };

```

- And with the PPE compiler generates a vector of objects:

```

1 | struct spe_error_struct
2 | { int id; const char * symbol; const char * message; };

```

# Unit Testing

## Unit Testing

- Boost.Test is great but builds don't fit SPEs:  
libboost\_unit\_test\_framework.so.1.45.0: 998kB

## Unit Testing

- Boost.Test is great but builds don't fit SPEs:  
libboost\_unit\_test\_framework.so.1.45.0: 998kB
- First idea: boost/detail/lightweight\_test.hpp

## Unit Testing

- Boost.Test is great but builds don't fit SPEs:  
libboost\_unit\_test\_framework.so.1.45.0: 998kB
- First idea: boost/detail/lightweight\_test.hpp  
misses a lot of the Boost.Test goodness



## Unit Testing

- Boost.Test is great but builds don't fit SPEs:  
libboost\_unit\_test\_framework.so.1.45.0: 998kB
- First idea: boost/detail/lightweight\_test.hpp  
misses a lot of the Boost.Test goodness

Enter SPE-Unit:

- Compromise between lightweight and feature-complete

## Unit Testing

- Boost.Test is great but builds don't fit SPEs:  
libboost\_unit\_test\_framework.so.1.45.0: 998kB
- First idea: boost/detail/lightweight\_test.hpp  
misses a lot of the Boost.Test goodness

Enter SPE-Unit:

- Compromise between lightweight and feature-complete
- Designed after Boost.Test

## Unit Testing - SPE Unit Features

- Only one test suite is available: `CBE_MPI_SPEUNIT_AUTO_TEST_SUITE()`;

## Unit Testing - SPE Unit Features

- Only one test suite is available: `CBE_MPI_SPEUNIT_AUTO_TEST_SUITE()`;
- Tests are started explicitly:

```
1 | uint32_t result = CBE_MPI_SPEUNIT_RUN_TEST_SUITE();  
2 | SET_RETURN_VALUE(result);
```

## Unit Testing - SPE Unit Features

- Only one test suite is available: `CBE_MPI_SPEUNIT_AUTO_TEST_SUITE()`;
- Tests are started explicitly:

```
1 | uint32_t result = CBE_MPI_SPEUNIT_RUN_TEST_SUITE();  
2 | SET_RETURN_VALUE(result);
```

- The powerful `AUTO_TEST_CASE_TEMPLATE(testname, T, typelist)` and a normal template `TEST_CASE_TEMPLATE(testname)` are included

## Unit Testing - SPE Unit Features

- Only one test suite is available: `CBE_MPI_SPEUNIT_AUTO_TEST_SUITE()`;
- Tests are started explicitly:

```
1 | uint32_t result = CBE_MPI_SPEUNIT_RUN_TEST_SUITE();  
2 | SET_RETURN_VALUE(result);
```

- The powerful `AUTO_TEST_CASE_TEMPLATE(testname, T, typelist)` and a normal template `TEST_CASE_TEMPLATE(testname)` are included
- Different test tool levels are supported: `WARN_*`, `CHECK_*`, `REQUIRE_*`

## Unit Testing - SPE Unit Features

- Only one test suite is available: `CBE_MPI_SPEUNIT_AUTO_TEST_SUITE()`;
- Tests are started explicitly:

```
1 | uint32_t result = CBE_MPI_SPEUNIT_RUN_TEST_SUITE();  
2 | SET_RETURN_VALUE(result);
```

- The powerful `AUTO_TEST_CASE_TEMPLATE(testname, T, typelist)` and a normal template `TEST_CASE_TEMPLATE(testname)` are included
- Different test tool levels are supported: `WARN_*`, `CHECK_*`, `REQUIRE_*`
- Strings can be disabled to reduce overhead (silent mode)

## Unit Testing - SPE Unit Features

- Only one test suite is available: `CBE_MPI_SPEUNIT_AUTO_TEST_SUITE()`;
- Tests are started explicitly:

```
1 | uint32_t result = CBE_MPI_SPEUNIT_RUN_TEST_SUITE();  
2 | SET_RETURN_VALUE(result);
```

- The powerful `AUTO_TEST_CASE_TEMPLATE(testname, T, typelist)` and a normal template `TEST_CASE_TEMPLATE(testname)` are included
- Different test tool levels are supported: `WARN_*`, `CHECK_*`, `REQUIRE_*`
- Strings can be disabled to reduce overhead (silent mode)
- Emulated SPE exceptions can be validated with test tools like `CBE_MPI_REQUIRE_THROW`



## Unit Testing - Example

```
1 typedef boost::mpl::vector_c<int,1,2,4,8,16> aligned_alloc_alignments;
2
3 CBE_MPI_SPEUNIT_AUTO_TEST_SUITE();
4
5 CBE_MPI_SPEUNIT_AUTO_TEST_CASE_TEMPLATE( aligned_malloc_free_test, T,
6     aligned_alloc_alignments )
7 {
8     aligned_ptr<void,T::value> ptr = aligned_malloc<T::value>(T::value);
9     CBE_MPI_SPEUNIT_REQUIRE_EQUAL(is_aligned<T::value>(ptr.get()),true);
10    cbe_mpi::aligned_free(ptr);
11    CBE_MPI_SPEUNIT_REQUIRE_EQUAL(ptr.get(),((void*)(0)));
12 }
13
14 int kernel(void)
15 {
16     uint32_t result = CBE_MPI_SPEUNIT_RUN_TEST_SUITE();
17     SET_RETURN_VALUE(result);
18 }
```

## Data Transfer - Single Buffer

```
ii = in.get();  
oo = out.get();  
  
for(int i=0; i<iterations; i++) {  
  
    spe_ppe_get_c(in.get(), cd->inbuf1+(SPE_Rank()+i*SPE_Size()*slicesize*sizeof(float),  
        slicesize_padded*sizeof(float));  
  
    harris_simd(ii, oo, cd->slice_dimx, cd->slice_dimy, 0, PADY, buf1.get(), buf2.get(), buf3.get());  
  
    spe_ppe_put_c(cd->outbuf1+(SPE_Rank()+i*SPE_Size()*slicesize*sizeof(float) +  
        (cd->slice_dimx*PADY)*sizeof(float), oo, slicesize*sizeof(float));  
}
```

# Data Transfer - Single Buffer

```

ii = in.get();
oo = out.get();

for(int i=0; i<iterations; i++) {

    spe_ppe_get_c(in.get(), cd->inbuf1+(SPE_Rank()+i*SPE_Size())*slicesize*sizeof(float),
                 slicesize_padded*sizeof(float));

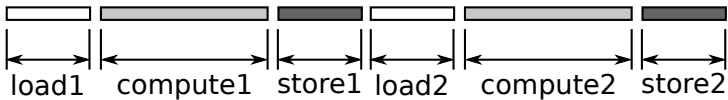
    harris_simd(ii, oo, cd->slice_dimx, cd->slice_dimy, 0, PADY, buf1.get(), buf2.get(), buf3.get());

    spe_ppe_put_c(cd->outbuf1+(SPE_Rank()+i*SPE_Size())*slicesize*sizeof(float) +
                 (cd->slice_dimx*PADY)*sizeof(float), oo, slicesize*sizeof(float));
}

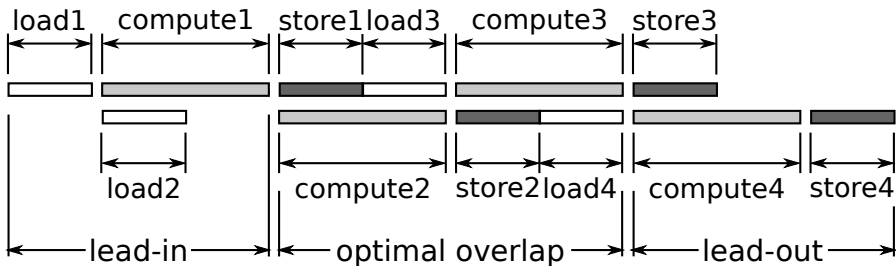
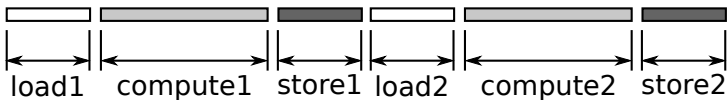
```

load  
 calc  
 store

# Data Transfer



# Data Transfer



# Data Transfer - Double Buffering

```
spe_ppe_get_async_c(in1.get(), cd->inbuf1+SPE_Rank()*slicesize*sizeof(float), slicesize_padded*sizeof(float), 9);
```

```
for(int i=0; i<iterations; i++) {
    if(i%2 == 0) {
        spe_ppe_get_async_c(in2.get(), cd->inbuf1+(SPE_Rank()+i+1)*SPE_Size()*slicesize*sizeof(float),
            slicesize_padded*sizeof(float), 10);

        dma_synchronize_c(9); dma_synchronize_c(11);
        ii = in1.get(); oo = out1.get();
    } else {
        spe_ppe_get_async_c(in1.get(), cd->inbuf1+(SPE_Rank()+i+1)*SPE_Size()*slicesize*sizeof(float),
            slicesize_padded*sizeof(float), 9);

        dma_synchronize_c(10); dma_synchronize_c(12);
        ii = in2.get(); oo = out2.get();
    }

    harris_simd(ii, oo, cd->slice_dimx, cd->slice_dimy, 0, PADY, buf1.get(), buf2.get(), buf3.get());

    if(i%2 == 0) {
        spe_ppe_put_async_c(cd->outbuf1+(SPE_Rank()+i*SPE_Size())* slicesize*sizeof(float) +
            (cd->slice_dimx*PADY)*sizeof(float), out1.get(), slicesize*sizeof(float), 11);
    } else {
        spe_ppe_put_async_c(cd->outbuf1+(SPE_Rank()+i*SPE_Size())* slicesize*sizeof(float) +
            (cd->slice_dimx*PADY)*sizeof(float), out2.get(), slicesize*sizeof(float), 12);
    }
}

spe_ppe_put_c(cd->outbuf1 + (SPE_Rank()+((iterations-1)*SPE_Size())) * slicesize*sizeof(float) +
    (cd->slice_dimx*PADY)*sizeof(float), oo, slicesize*sizeof(float));
```

# Data Transfer - Double Buffering

```
spe_ppe_get_async_c(in1.get(), cd->inbuf1+SPE_Rank()*slicesize*sizeof(float), slicesize_padded*sizeof(float), 9);
```

```
for(int i=0; i<iterations; i++) {
    if(i%2 == 0) {
        spe_ppe_get_async_c(in2.get(), cd->inbuf1+(SPE_Rank()+i+1)*SPE_Size()*slicesize*sizeof(float),
            slicesize_padded*sizeof(float), 10);

        dma_synchronize_c(9); dma_synchronize_c(11);
        ii = in1.get(); oo = out1.get();
    } else {
        spe_ppe_get_async_c(in1.get(), cd->inbuf1+(SPE_Rank()+i+1)*SPE_Size()*slicesize*sizeof(float),
            slicesize_padded*sizeof(float), 9);

        dma_synchronize_c(10); dma_synchronize_c(12);
        ii = in2.get(); oo = out2.get();
    }
}
```

lead-in

load,  
sync

```
harris_simd(ii, oo, cd->slice_dimx, cd->slice_dimy, 0, PADY, buf1.get(), buf2.get(), buf3.get());
```

calc

```
if(i%2 == 0) {
    spe_ppe_put_async_c(cd->outbuf1+(SPE_Rank()+i*SPE_Size())* slicesize*sizeof(float) +
        (cd->slice_dimx*PADY)*sizeof(float), out1.get(), slicesize*sizeof(float), 11);
} else {
    spe_ppe_put_async_c(cd->outbuf1+(SPE_Rank()+i*SPE_Size())* slicesize*sizeof(float) +
        (cd->slice_dimx*PADY)*sizeof(float), out2.get(), slicesize*sizeof(float), 12);
}
}
```

store

```
spe_ppe_put_c(cd->outbuf1 + (SPE_Rank()+((iterations-1)*SPE_Size())) * slicesize*sizeof(float) +
    (cd->slice_dimx*PADY)*sizeof(float), oo, slicesize*sizeof(float));
```

lead-out

# Double Buffering - Operations - Input Segment



## Double Buffering - Operations - Input Segment

- Start loading first segment (lead-in)

## Double Buffering - Operations - Input Segment

- Start loading first segment (lead-in)
  
- Start loading next segment

## Double Buffering - Operations - Input Segment

- Start loading first segment (lead-in)
- Start loading next segment
- Wait for segment to be ready for computation

## Double Buffering - Operations - Input Segment

- Start loading first segment (lead-in)
- Start loading next segment
- Wait for segment to be ready for computation
- Signal that computation on current segment is finished

## Double Buffering - Operations - Input Segment

- Start loading first segment (lead-in)  
`operator =()`
- Start loading next segment
- Wait for segment to be ready for computation
- Signal that computation on current segment is finished

## Double Buffering - Operations - Input Segment

- Start loading first segment (lead-in)  
`operator =()`
- Start loading next segment  
`operator ++(int)`
- Wait for segment to be ready for computation
- Signal that computation on current segment is finished

## Double Buffering - Operations - Input Segment

- Start loading first segment (lead-in)  
`operator =()`
- Start loading next segment  
`operator ++(int)`
- Wait for segment to be ready for computation  
`operator *()`
- Signal that computation on current segment is finished

## Double Buffering - Operations - Input Segment

- Start loading first segment (lead-in)  
`operator =()`
- Start loading next segment  
`operator ++(int)`
- Wait for segment to be ready for computation  
`operator *()`
- Signal that computation on current segment is finished  
`operator ++(int)`



## Double Buffering - Operations - Input Segment

- Start loading first segment (lead-in)  
`operator =()`
- Start loading next segment  
`operator ++(int)`
- Wait for segment to be ready for computation  
`operator *()`
- Signal that computation on current segment is finished  
`operator ++(int)`
- Check if end of data is reached

## Double Buffering - Operations - Input Segment

- Start loading first segment (lead-in)  
`operator =()`
- Start loading next segment  
`operator ++(int)`
- Wait for segment to be ready for computation  
`operator *()`
- Signal that computation on current segment is finished  
`operator ++(int)`
- Check if end of data is reached  
`operator ==()`

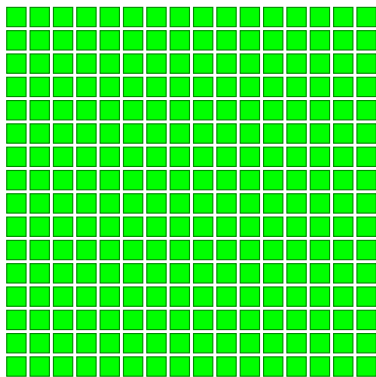
# Double Buffered Segmented Input Iterator

```
1  template<typename T> struct remote_segmented_input_iterator
2  {
3      // allocate required buffers
4      remote_segmented_input_iterator(...) {}
5
6      // start loading first buffer
7      void operator= (const addr64 & base_address_) { }
8
9      // wait for current segment to arrive and return pointer to it
10     T* operator *() {}
11
12     // start loading new data and increment current segment
13     inline void operator++(int) {}
14
15     // check if iterator has reached a position
16     bool operator ==(const addr64 & b) const {}
17 };
```

## Double Buffered Segmented Iterator Example

```
1 remote_segmented_input_iterator<float> it(depth,  
2   ssize, slicer(ssize));  
3 remote_segmented_output_iterator<float> ot(depth,  
4   ssize, slicer(ssize));  
5  
6 for(it = input, ot = output; /* lead-in */  
7   it!=input+overall_size; /* check end */  
8   it++, ot++) // load next, store current  
9 {  
10  float * in = *it; float * out = *ot; // synchronize  
11  harris_simd(in, out, cd->slice_dimx, cd->slice_dimy,  
12    0, PADY, buf1.get(), buf2.get(), buf3.get());  
13 }
```

## Double Buffered Segmented Iterator - Slicer



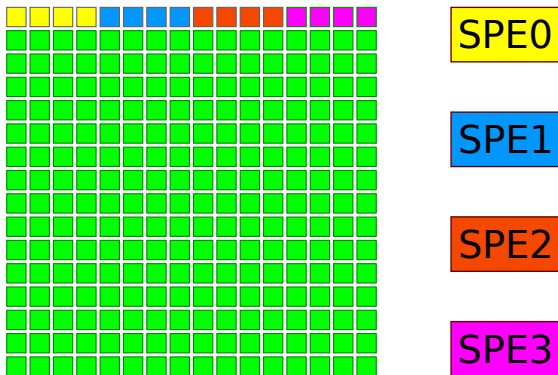
SPE0

SPE1

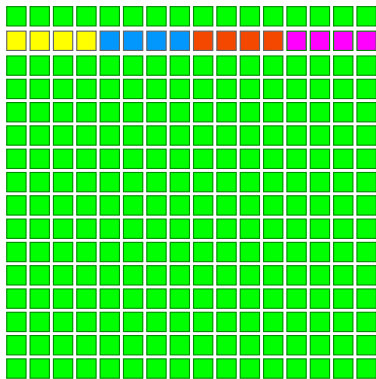
SPE2

SPE3

## Double Buffered Segmented Iterator - Slicer



## Double Buffered Segmented Iterator - Slicer



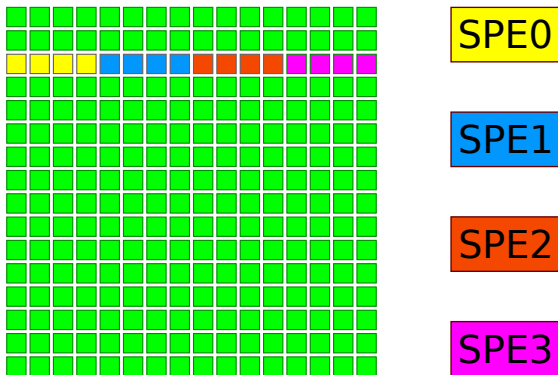
SPE0

SPE1

SPE2

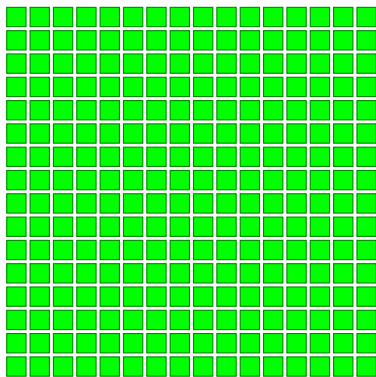
SPE3

## Double Buffered Segmented Iterator - Slicer





## Double Buffered Segmented Iterator - Slicer



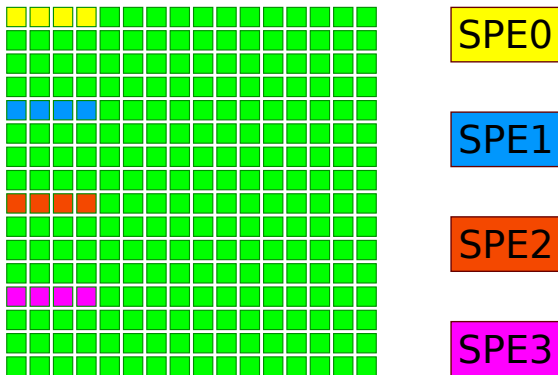
SPE0

SPE1

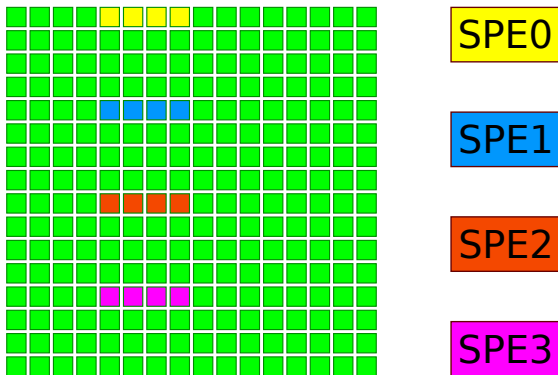
SPE2

SPE3

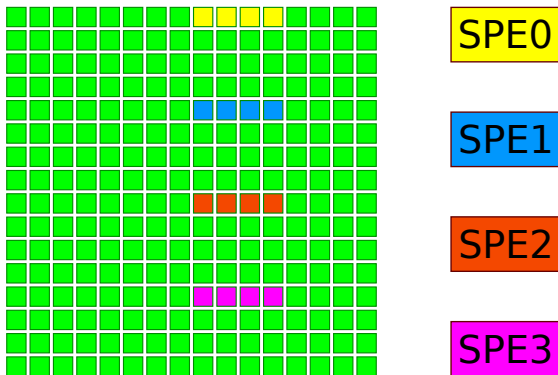
## Double Buffered Segmented Iterator - Slicer



## Double Buffered Segmented Iterator - Slicer



## Double Buffered Segmented Iterator - Slicer



## Multi-Buffered Segmented Iterator - Features

- `remote_vector<T>` for more expressive code:

```
1 // PPE:
2 std::vector<float> v(1024*1024); kernel(v);
3 // SPE:
4 kernel(remote_vector<float> v) {
5     remote_segmented_input_iterator<T> it(depth, ssize, slicer(ssize));
6     for(it = v.begin(); it!=v.end(); it++) {
7         float * in = *it;
8         /* computation */
9     }
10 }
```

## Multi-Buffered Segmented Iterator - Features

- `remote_vector<T>` for more expressive code:

```
1 // PPE:
2 std::vector<float> v(1024*1024); kernel(v);
3 // SPE:
4 kernel(remote_vector<float> v) {
5     remote_segmented_input_iterator<T> it(depth, ssize, slicer(ssize));
6     for(it = v.begin(); it!=v.end(); it++) {
7         float * in = *it;
8         /* computation */
9     }
10 }
```

- Read, write- and read-write Iterators with minimum buffer depth of 3

## Multi-Buffered Segmented Iterator - Features

- `remote_vector<T>` for more expressive code:

```
1 // PPE:  
2 std::vector<float> v(1024*1024); kernel(v);  
3 // SPE:  
4 kernel(remote_vector<float> v) {  
5     remote_segmented_input_iterator<T> it(depth, ssize, slicer(ssize));  
6     for(it = v.begin(); it!=v.end(); it++) {  
7         float * in = *it;  
8         /* computation */  
9     }  
10 }
```

- Read, write- and read-write Iterators with minimum buffer depth of 3
- Various slicers

## 2D Multi-Buffered Segmented Iterator

- Native 2D data transfer support through DMA lists



## 2D Multi-Buffered Segmented Iterator

- Native 2D data transfer support through DMA lists
- Difference to regular iterator:

## 2D Multi-Buffered Segmented Iterator

- Native 2D data transfer support through DMA lists
- Difference to regular iterator:
  - Slice size is 2D

## 2D Multi-Buffered Segmented Iterator

- Native 2D data transfer support through DMA lists
- Difference to regular iterator:
  - Slice size is 2D
  - Supports `remote_vector_2D`

## 2D Multi-Buffered Segmented Iterator

- Native 2D data transfer support through DMA lists
- Difference to regular iterator:
  - Slice size is 2D
  - Supports `remote_vector_2D`
  - Slicer takes 2D arguments:  
`slicer_2D(size_2d_t vector_dim, size_2d_t slice_dim);`

## 2D Multi-Buffered Segmented Iterator

- Native 2D data transfer support through DMA lists
- Difference to regular iterator:
  - Slice size is 2D
  - Supports `remote_vector_2D`
  - Slicer takes 2D arguments:  
`slicer_2D(size_2d_t vector_dim, size_2d_t slice_dim);`
- Ideal for image processing:

## 2D Multi-Buffered Segmented Iterator

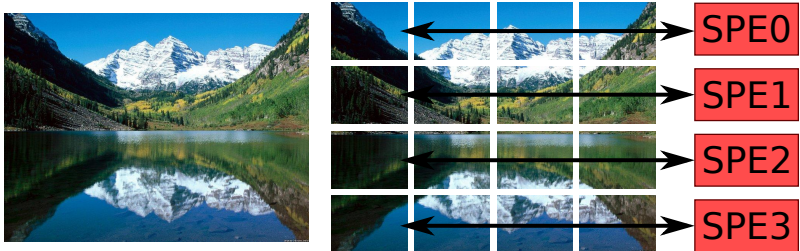
- Native 2D data transfer support through DMA lists
- Difference to regular iterator:
  - Slice size is 2D
  - Supports `remote_vector_2D`
  - Slicer takes 2D arguments:  
`slicer_2D(size_2d_t vector_dim, size_2d_t slice_dim);`
- Ideal for image processing:



## 2D Multi-Buffered Segmented Iterator

- Native 2D data transfer support through DMA lists
- Difference to regular iterator:
  - Slice size is 2D
  - Supports `remote_vector_2D`
  - Slicer takes 2D arguments:
 

```
slicer_2D(size_2d_t vector_dim, size_2d_t slice_dim);
```
- Ideal for image processing:



# High-Level Inter-SPE Communication: MPI



## High-Level Inter-SPE Communication: MPI

- Interprocess communication by message passing, SPEs send and receive message

## High-Level Inter-SPE Communication: MPI

- Interprocess communication by message passing, SPEs send and receive message
- API specification, used in high performance computing

## High-Level Inter-SPE Communication: MPI

- Interprocess communication by message passing, SPEs send and receive message
- API specification, used in high performance computing

Features:

- Virtual topology of processes

## High-Level Inter-SPE Communication: MPI

- Interprocess communication by message passing, SPEs send and receive message
- API specification, used in high performance computing

Features:

- Virtual topology of processes
- Synchronization

## High-Level Inter-SPE Communication: MPI

- Interprocess communication by message passing, SPEs send and receive message
- API specification, used in high performance computing

### Features:

- Virtual topology of processes
- Synchronization
- Point to point communication

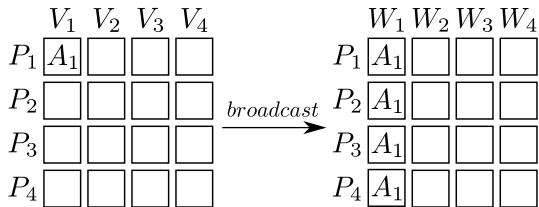
## High-Level Inter-SPE Communication: MPI

- Interprocess communication by message passing, SPEs send and receive message
- API specification, used in high performance computing

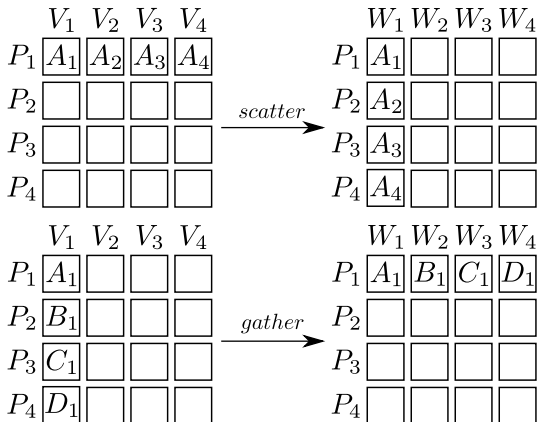
### Features:

- Virtual topology of processes
- Synchronization
- Point to point communication
- Collective communication

## MPI Collectives - Broadcast

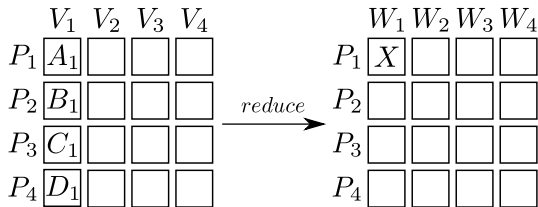


## MPI Collectives - Scatter and Gather

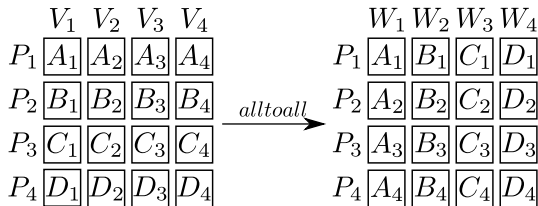




## MPI Collectives - Reduce and All to All



$$X = Op(A_1, Op(Op(B_1, Op(C_1, D_1))))$$



## MPI Interface - Example

```
1 communicator world;
2
3 if (world.rank() == 0)
4 {
5     char s1[] = "Hello";
6     world.send(1, 0, s1, sizeof(s1));
7     char s2[6];
8     world.recv(1, 1, s2, sizeof(s2));
9 }
10 else if (world.rank() == 1)
11 {
12     char s1[6];
13     world.recv(0, 0, s1, sizeof(s1));
14     char s2[] = "world";
15     world.send(0, 1, s2, sizeof(s2));
16 }
17 // Hello world from SPE 0, Hello world from SPE 1
```

# MPI Interface - Communicator

```
1 class communicator
2 {
3     void barrier();
4
5     template <typename T> void send(int dst, int tag, const T& value);
6     template <typename T> void send(int dst, int tag, const T* values, int n);
7     template <typename T> request isend(int dst, int tag, const T& value);
8     ...
9     template <typename T> status recv(int source, int tag, T& value);
10    template <typename T> status recv(int source, int tag, T* values, int n);
11    template <typename T> request irectv(int source, int tag, T& value);
12    ...
13    communicator include(uint16_t first, uint16_t last);
14    communicator exclude(uint16_t first, uint16_t last);
15    friend bool operator==(const communicator& c1, const communicator& c2);
16 };
```

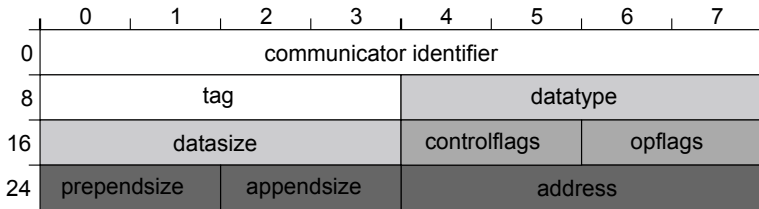
## MPI Interface - Request and Status

```
1 // represents current request
2 class request
3 {
4     request() {};
5     status wait();
6     boost::optional<status> test();
7 };
8
9 // represents status of a request
10 class status
11 {
12     int32_t source() const;
13     int32_t tag() const;
14     int32_t error() const;
15 };
```

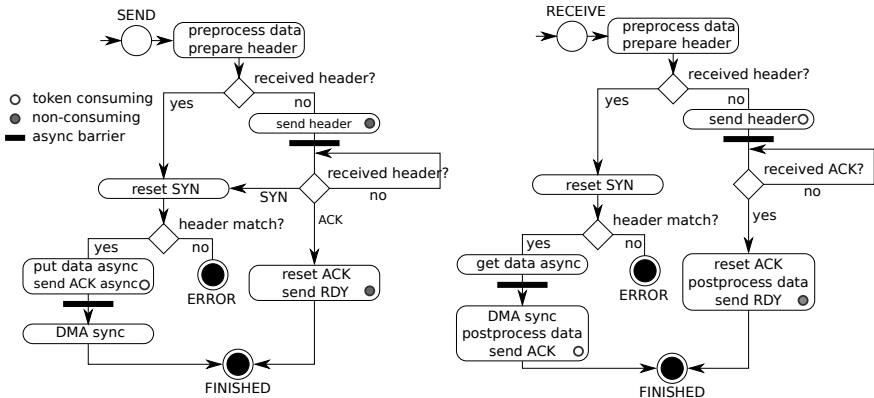
## MPI Interface - Collectives Interface

```
1  template<typename T, typename Op>
2  void reduce(const communicator & comm, const T & in,
3             T & out, Op op, int root);
4
5  template<typename T, typename Op>
6  void reduce(const communicator & comm, const T & in,
7             Op op, int root);
8
9  template<typename T, typename Op>
10 void reduce(const communicator & comm, const T * in,
11            int n, T * out, Op op, int root);
12
13 template<typename T, typename Op>
14 void reduce(const communicator & comm, const T * in,
15            int n, Op op, int root);
```

# MPI Header



# MPI Protocol



# MPI Types

We don't do Boost.Serialization but

- you may register your POD type:

```
1 | struct gps_position { /* POD */ };
2 | namespace cbe_mpi
3 | {
4 |     CBE_MPI_USER_POD_DATATYPE(gps_position);
5 | }
```

- or you may specialize send/receive methods:

```
1 | template <typename T>
2 | request isend(cbe_mpi::communicator & comm, int dst,
3 |     int tag, T data, int n);
4 |
5 | template <typename T>
6 | request irecv(cbe_mpi::communicator & comm, int src,
7 |     int tag, T data, int n);
```



## Registering POD Types

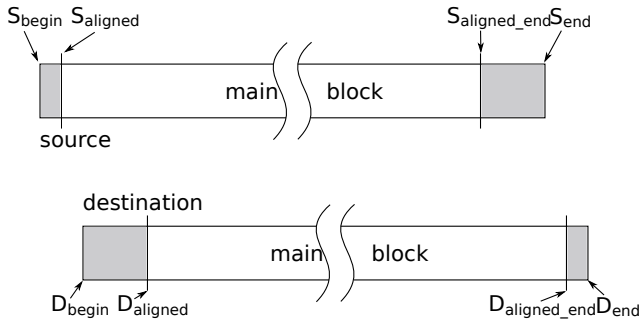
How we identify your type:

```
1  template<typename T>
2  struct cbe_mpi_user_pod_type_id { static void get() {} };
3
4  #define CBE_MPI_USER_POD_DATATYPE(CppType) \
5  template<> \
6  struct is_mpi_datatype< CppType > \
7  : boost::mpl::bool_<true> {}; \
8      \
9  inline int get_mpi_datatype(const CppType &) \
10 { \
11     return 0x80000000 | \
12         (int)&cbe_mpi_user_pod_type_id< CppType >::get(); \
13 }
```

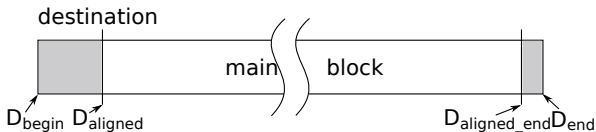
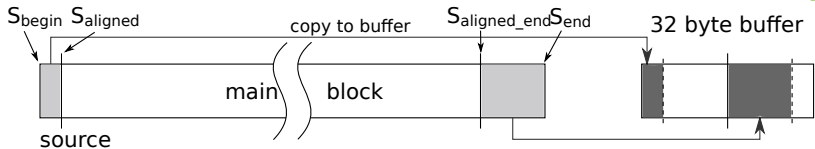
## Sending std::vector

```
1  template <typename T>
2  request isend(cbe_mpi::communicator com,
3              int dest, int tag, const std::vector<T> * values, int)
4  {
5      int vectorsize = values->size();
6      com.send(dest, tag, &vectorsize, 1);
7      return com.isend(dest, tag, &(*values)[0], vectorsize);
8  }
9
10 template <typename T>
11 request irecv(cbe_mpi::communicator com,
12             int source, int tag, std::vector<T> * values, int)
13 {
14     int vectorsize;
15     com.recv(source, tag, &vectorsize, 1);
16     values->resize(vectorsize);
17     return com.irecv(source, tag, &(*values)[0], vectorsize);
18 }
```

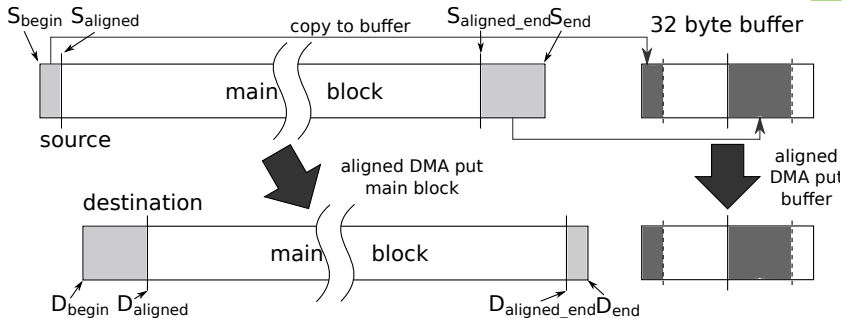
# MPI - Sending Unaligned Data



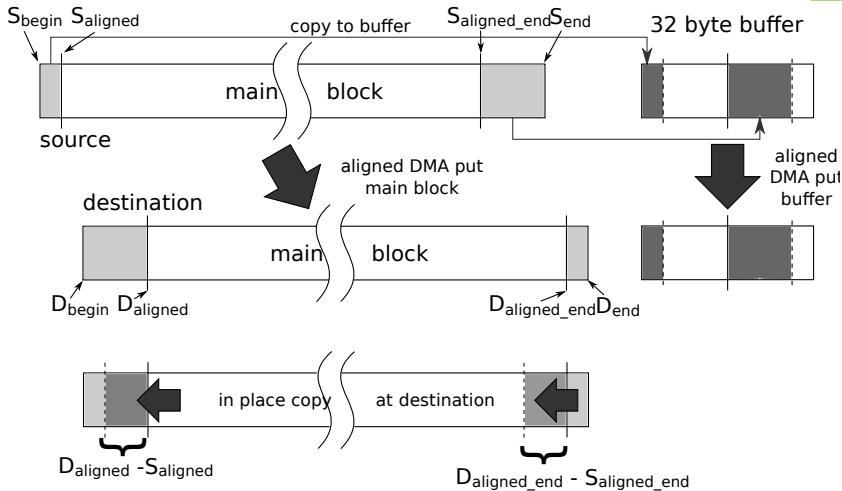
# MPI - Sending Unaligned Data



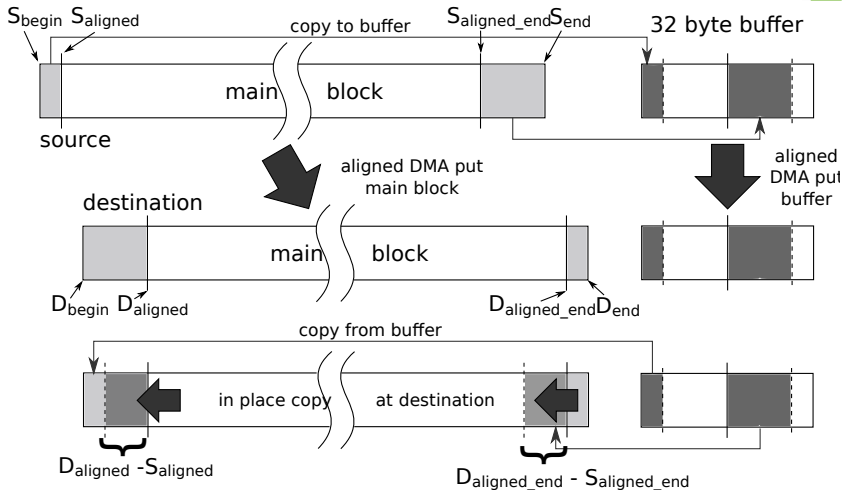
# MPI - Sending Unaligned Data



# MPI - Sending Unaligned Data



# MPI - Sending Unaligned Data



## Conclusion

- Build process can be simplified with CMake



## Conclusion

- Build process can be simplified with CMake
- Boilerplate code can be simplified with the help of Boost (e.g. PP)

## Conclusion

- Build process can be simplified with CMake
- Boilerplate code can be simplified with the help of Boost (e.g. PP)
- Ambiguity of functions or macros in different compilation units can be exploited

## Conclusion

- Build process can be simplified with CMake
- Boilerplate code can be simplified with the help of Boost (e.g. PP)
- Ambiguity of functions or macros in different compilation units can be exploited
  
- Optimal Boost solutions have to be adapted to fit embedded architecture

## Conclusion

- Build process can be simplified with CMake
- Boilerplate code can be simplified with the help of Boost (e.g. PP)
- Ambiguity of functions or macros in different compilation units can be exploited
  
- Optimal Boost solutions have to be adapted to fit embedded architecture
- Sweet spot between generic code and efficiency must be found

## Conclusion

- Build process can be simplified with CMake
- Boilerplate code can be simplified with the help of Boost (e.g. PP)
- Ambiguity of functions or macros in different compilation units can be exploited
  
- Optimal Boost solutions have to be adapted to fit embedded architecture
- Sweet spot between generic code and efficiency must be found
  
- Difficult low-level code can be wrapped nicely in C++ Concepts

## Conclusion

- Build process can be simplified with CMake
- Boilerplate code can be simplified with the help of Boost (e.g. PP)
- Ambiguity of functions or macros in different compilation units can be exploited
  
- Optimal Boost solutions have to be adapted to fit embedded architecture
- Sweet spot between generic code and efficiency must be found
  
- Difficult low-level code can be wrapped nicely in C++ Concepts
- C++ Concepts can be even more powerful on special purpose hardware

## Conclusion

- Build process can be simplified with CMake
- Boilerplate code can be simplified with the help of Boost (e.g. PP)
- Ambiguity of functions or macros in different compilation units can be exploited
  
- Optimal Boost solutions have to be adapted to fit embedded architecture
- Sweet spot between generic code and efficiency must be found
  
- Difficult low-level code can be wrapped nicely in C++ Concepts
- C++ Concepts can be even more powerful on special purpose hardware

Thank you for you kind attention.