# SMECY C (or SME-C? ☺)

—

# C99 with pragma and API for parallel execution, processor mapping and communication generation

Remi.Barrere@thalesgroup.com
Ronan.Keryell@hpc-project.com

May 28, 2011

## 1 Introduction

In the SMECY project we want to use C source code as a portable intermediate representation (IR) between tools from high-level tools down to lower-level tools because of its good trade-off between expressiveness and readability, without compromising portability.

The targets envisioned in SMECY are heterogeneous multicore systems with shared memory or not, with various hardware accelerator, such as ASIC, ASIP, FPGA, GPU, specialized vector processors, partially reconfigurable accelerators...

Unfortunately, since it is undecidable to get high-level properties from such programs, we use decorations to help tools to understand program behaviour and generate codes for some hardware targets. We try to keep clear decorations, easy to understand, so that SMECY C can also used as a programming language.

A SMECY program contains various functions that may be executed on various processors, accelerators, GPU... that may consume and produce data from different physical memory spaces.

Since we want to express also performance on given platforms, we keep the opportunity to have platform-specific pragma and API or specialized intrinsic types and functions, for example to express use of special hardware accelerator functions or operations.

The hardware specific pragma and intrinsics are to be defined between software and hardware partners involved in various use cases. But it may not possible to address all the programming models and platforms envisioned in the project.

### 1.1 SMECY-C programming model

The programming model is based on C processes, with a virtual shared memory and threads *à la* OpenMP. Since we may have quite asynchronous processes in a real application description or at the back-end level in the execution model, we can cope indeed with different C processes communicating with an API.

We add mapping information stating on which hardware part a function is to be placed and run.

The programming model exposed in the following is based on an OpenMP SMP model because of its (rather) simple readability, elegance, old background and wide acceptance. We make the hypothesis that a SMECY program is a correct OpenMP program that can be executed in sequential with a C OpenMP-free compiler (just by ignoring OpenMP `#pragma`) and in parallel on a SMP target (such an $x86$ machine) by using an OpenMP compiler *with the same semantics*. Since we use C (by opposition to a DSL) as an internal representation, we choose this behaviour to stick to standard behaviour as much as possible. This is known as the sequential equivalence. Since we can cope different results for performance reasons from executing in parallel non associative floating point operations, we deal with only *weak* sequential equivalence instead of *strong* sequential equivalence.

Of course, this model is incompatible with a real hardware target envisioned in SMECY, so we need to add hints in the code explaining memory dependencies at the function call levels. Since it is quite difficult to describe general dependencies, we approximate memory dependencies with rectangles (and more generally hyperparallelepipede in any dimension) that can be read, written or both. We think these abstractions are good trade-offs between expressiveness (and what a programmer can endure...) and hardware capabilities. Even if there is some similitude with HPF (High Performance Fortran) we do not deal with strides[1].

With this information the tools can guess the communication to generate between the different functions and memory spaces to emulate the global OpenMP memory semantics.

The neat side effect is that we have the same global program executed on all the platforms (sequential, real OpenMP and SMECY) with the same semantics and we can see the sequential version as a functional simulator of the SMECY application and the real OpenMP version as a parallelized quicker version of this simulator.

It is also easy to debug the application, but also all the SMECY tools used or developed in the project.

To be able to address real hardware from the C level with special needs:

- to specify hardware register names;

- define input/output routines specifically but in portable way;

- define fixed point computations;

- specific data size;

- accumulator register (DSP...);

- different memory spaces that can be chosen specifically to optimize storage and speed (DSP, hardware accelerators with scratch pad memory...);

- saturated arithmetic.

---

[1]Indeed, by using some higher dimensional arrays than the arrays used in the application, you can express them... So may be we can express them in the syntax?

For all these we rely on the TR 18037 Embedded C standard, supported for example by ACE tools.

To describe different processes communicating together in an asynchronous way, we do not have anymore a sequential equivalence and then do not use pragma to express this. So we use a simple communication, synchronization and threading API. Since we target embedded systems with a light efficient implementation, we can rely on such standard API as the ones of the MultiCore Association: MCAPI (communication), MRAPI (synchronization) and MTAPI (threading).

As modern programs need a clear documentation, such as with Doxygen marking style providing meta information on different elements of the program, we can use this (hopefully correct) information to help the compiling process itself.

Since tools are to be oriented for a specific target, taking into account various hardware and compilation parameters that are to be kept orthogonally to application sources, those parameters are kept aside in some description files that flow between tools. These files may be represented with an XML syntax (using a SMECY naming space) or by even simpler format (JASON, YAML...).

## 1.2 Reference documents

Besides the SMECY documentation, the reader should be knowledgeable of some work of the ISO/IEC JTC1/SC22/WG14 committee on the C language standard:

- ISO/IEC 9899 - Programming languages - C (Technical Corrigendum 3 Cor. 3:2007(E)) `http://www.open-std.org/jtc1/sc22/wg14/www/docs/n1256.pdf`

- TR 18037: Embedded C `http://www.open-std.org/jtc1/sc22/wg14/www/docs/n1169.pdf`

- Future C1X standard `http://www.open-std.org/jtc1/sc22/wg14/www/docs/n1548.pdf`

and other standards such as

- MCA (MultiCore Association) API: MCAPI, MRAPI & MTAPI, `www.multicore-association.org`

- HPF (High Performance Fortran) `http://hpff.rice.edu/versions/hpf2`

- ISO/IEC TR 24717:2009, Information technology – Programming languages, their environments and system software interfaces – Collection classes for programming language COBOL $20xx$ `http://www.cobolstandard.info/j4/files/std.zip`

## 2 Intermediate representation use cases

### 2.1 Direct programming

A programmer can program her application directly in SMECY C with OpenMP and SMECY-specific pragma and API to target a SMECY platform.

It can be at rather high-level, by using only high-level pragma, or rather at a lower level, by using different communicating processes with the API and even specialized API and pragma for specific hardware.

A process written in C code with pragma and API can express a global host controlling process of an application or a local program in a specialized processor. And we may have many of such processes to express different producer and consumer Kahn processes interacting through a NoC in an asynchronous way.

## 2.2 System high-level synthesis

A compiler can take a sequential plain C or Matlab or another language code, analyze and parallelize the code by adding automatically parallel and mapping pragma. This can be seen as a high-level synthesis at system level.

A tool such as Par4All can do these kinds of transformation.

## 2.3 Hardware high-level synthesis

A compiler can take a program with SMECY pragma and compile any call with a mapping of a given kind into some hardware configuration or program to be executed instead of the function and an API call to use this hardware part from the host program.

Since the pragma are designed to be concretely compilable, such a tool should be easy to do with a simple compilation framework, such as ROSE Compiler.

The SMECY API and intrinsics are chosen to be mapped quite straightfor-wardly to real hardware functions by the back-end.

# 3 Exemples

## 3.1 Program with contiguous memory transfers

During C memory transfer, if we work on arrays with the last dimension taken as a whole, the memory is contiguous and the programs often work even there are some aliasing such as using a 2D array zone as a linearized 1D vector.

The following program exposes this kind of code where some work sharing is done by contiguous memory blocks.

```
1   /* To compile this program on Linux, try:

        make CFLAGS='-std=c99 -Wall' pragma_example

        To run:
6       ./pragma_example; echo $?
        It should print 0 if OK.

        You can even compile it to run on multicore SMP for free with

11      make CFLAGS='-std=c99 -fopenmp -Wall' pragma_example

        To verify there are really some clone() system calls that create the threads:
        strace -f ./pragma_example ; echo $?
```

26
```
#include <stdbool.h>

/* function Gen
```

```
31      Example of old C89 array use−case where the size is unknown. Note that
        this implies some nasty access linearization with array with more than
        1 dimension.
     */
    void Gen(int *out, int size) {
36    // Can be executed in parallel
    #pragma omp parallel for
      for (int i = 0; i < size; i++)
        out[i] = 0;
    }
41


    /* function Add

        Nice C99 array with dynamic size definition. Note this implies having
46      array size given first
    */
    void Add(int size, int in[size], int out[size]) {
      // Can be executed in parallel
    #pragma omp parallel for
51    for (int i = 0; i < size; i++)
        out[i] = in[i] + 1;
    }


56  /* function Test */
    bool Test(int size, int in[size]) {
      bool ok = true;
      /* Can be executed in parallel, ok is initialized from global value and
         at loop exit ok is the && operation between all the local ok
61       instances: */
    #pragma omp parallel for reduction(&&:ok)
      for (int i = 0; i < size; i++)
        /* We cannot have this simple code here:
           if (in[i] != 2)
66           exit(−1);
           because a loop or a fonction with exit() cannot be executed in parallel.

           Proof: there is a parallel execution interleaving that may execute
```

5

```
              some  computations  in  some  threads  with  a  greater  i  that  the  one
71            executing  the  exit()  done  on  another  thread.  So  the  causality  is
              not  respected.

              Anyway,  in  an  heterogenous  execution ,  just  think  about  how  to
              implement  the  exit()  operating  system  call  from  an
76            accelerator...  No  hope.  :−)

              So  use  a  reduction  instead  and  return  the  status  for  later
              inspection :
            */
81        ok &= (in[i] == 2);

        // Return  false  if  at  least  one  in[i]  is  not  equal  to  2:
        return ok;
      }
86


    /* main */
    int main(int argc, char* argv[]) {
      int tab[6][200];
91    // Gen  is  mapped  on  GPP  0,  it  produced  (out)  an  array  written  to  arg  1:
    #pragma smecy map(GPP, 0) arg(1, [6][200], out)
      /* Note  there  is  an  array  linearization  here,  since  we  give  a  2D  array
          to  Gen()  that  uses  it .  This  is  bad  programming  style ,  but  it  is  just
          to  show  it  can  be  handled  in  the  model  :−) */
96    Gen((int *) tab, 200*6);

      // Launch  different  things  in  parallel:
    #pragma omp parallel sections
      {
101       // Do  one  thing  in  parallel...
    #pragma omp section
        {
          /* Map  this  "Add"  call  to  PE  0,  arg  2  is  communicated  as  input  as  an
              array  of  "int  [3][200]",  and  after  execution  arg  3  is
106           communicated  out  as  an  array  of  "int  [3][200]"

              Note  the  aliasing  of  the  2  last  arguments.  Just  to  show  we  can
              handle  it .  :*/
    #pragma smecy map(PE, 0) arg(2, [3][200], in) arg(3, [3][200], out)
111       Add(200*3, (int *) tab, (int *) tab);
        }
        // ...with  another  thing
    #pragma omp section
        {
116       /* Map  this  "Add"  call  to  PE  1,  arg  2  is  communicated  as  input  as  an
              array  of  "int  [3][200]"  from  address  tab[3][0],  that  is  the
              second  half  of  tab,  and  after  execution  arg  3  is  communicated  out
              as  an  array  of  "int  [3][200]",  that  is  the  second  half  of  tab

121           Note  the  aliasing  of  the  2  last  arguments.  Just  to  show  we  can
              handle  it .  :*/
    #pragma smecy map(PE, 1) arg(2, [3][200], in) \
```

```
                              arg ( 3 , [ 3 ] [ 2 0 0 ] , out )
           Add(200∗3, &tab [3][0], &tab [3][0]);
126     }
      }

      // Launch different things in parallel:
#pragma omp parallel sections
131   {
#pragma omp section
       {
#pragma smecy map(PE, 2) arg ( 2 , [ 2 ] [ 2 0 0 ] , in ) arg ( 3 , [ 2 ] [ 2 0 0 ] , out )
          Add(200∗2, (int ∗) tab , (int ∗) tab );
136    }
#pragma omp section
       {
#pragma smecy map(PE, 3) arg ( 2 , [ 2 ] [ 2 0 0 ] , in ) arg ( 3 , [ 2 ] [ 2 0 0 ] , out )
          Add(200∗2, &tab [2][0], &tab [2][0]);
141    }
#pragma omp section
       {
#pragma smecy map(PE, 4) arg ( 2 , [ 2 ] [ 2 0 0 ] , in ) arg ( 3 , [ 2 ] [ 2 0 0 ] , out )
          Add(200∗2, &tab [4][0], &tab [4][0]);
146    }
    }
    // An example where arg 2 is just used as a whole implicitly:
#pragma smecy map(GPP, 0) arg ( 2 , in )
    bool result = Test(200∗6, (int ∗) tab );
151   // Return non 0 if the computation went wrong:
    return ! result ;
}
```

## 3.2   Program with non-contiguous memory transfers

In the following example, we apply different computations on square pieces of the image, that do not have contiguous representation in memory. That is why we need to express restrictions on the use of the whole array.

```
/∗ To compile this program on Linux , try :
 2
     make CFLAGS='−std=c99 −Wall ' example_2D

     To run :
     ./ example_2D ; echo $?
 7   It should print 0 if OK.

     You can even compile it to run on multicore SMP for free with

     make CFLAGS='−std=c99 −fopenmp −Wall ' example_2D
12
     To verify there are really some clone () system calls that create the threads :
     strace −f ./ example_2D ; echo $?

     You can notice that the #pragma smecy are ignored ( the project is
17   on−going :−) ) but that the program produces already correct results in
```

```
          sequential execution and parallel OpenMP execution.

          Enjoy!

22        Ronan.Keryell@hpc-project.com
          for ARTEMIS SMECY European project.
     */

     #include <stdlib.h>
27   #include "example_helper.h"


     // Problem size
     enum { WIDTH = 500, HEIGHT = 200 };
32

     /* The main host program controlling and representing the whole
        application */
     int main(int argc, char* argv[]) {
37     int image[HEIGHT][WIDTH];
       unsigned char output[HEIGHT][WIDTH];

       // Initialize with some values
       init_image(WIDTH, HEIGHT, image);
42
     #pragma omp parallel sections
       {
         // On one processor
         // We rewrite a small part of image:
47   #pragma smecy map(PE, 0) arg(3, inout, [HEIGHT][WIDTH]           \
                             /[HEIGHT/3:HEIGHT/3 + HEIGHT/2 - 1]      \
                             [WIDTH/8:WIDTH/8 + HEIGHT/2 - 1])
         square_symmetry(WIDTH, HEIGHT, image, HEIGHT/2, WIDTH/8, HEIGHT/3);

52       // On another processor
     #pragma omp section
         // Here let the compiler to guess the array size
     #pragma smecy map(PE, 1) arg(3, inout, /[HEIGHT/4:HEIGHT/4 + HEIGHT/2 - 1] \
                                   [3*WIDTH/8:3*WIDTH/8 + HEIGHT/2 - 1])
57       square_symmetry(WIDTH, HEIGHT, image, HEIGHT/2, 3*WIDTH/4, HEIGHT/4);

         // On another processor
     #pragma omp section
         // Here let the compiler to guess the array size
62   #pragma smecy map(PE, 1) arg(3, inout, /[2*HEIGHT/5:2*HEIGHT/5 + HEIGHT/2 - 1]  \
                                   [WIDTH/2:WIDTH/2 + HEIGHT/2 - 1])
         square_symmetry(WIDTH, HEIGHT, image, HEIGHT/2, WIDTH/2, 2*HEIGHT/5);
       }
       // Here there is a synchronization because of the parallel part end
67
       // Since there
       normalize_to_char(WIDTH, HEIGHT, image, output);

       write_pgm_image("2D_example-output.pgm", WIDTH, HEIGHT, output);
```

```
72
       return EXIT_SUCCESS;
   }
```

## 3.3 Pipelined example

*TODO*

## 3.4 Remapping example

Some information can be in a given layout but needed in another layout to be
used by a specific hardware accelerator.

```
 1   #include <stdlib.h>
     #include "example_helper.h"

     // Problem size
     enum { WIDTH = 500, HEIGHT = 200, LINE_SIZE = 100 };
 6
     /* Apply some pixel value inversion in a 1D array
      */
     void
     invert_vector(int line_size,
11                   int input_line[line_size],
                     int output_line[line_size]) {
       for(int i = 0; i < line_size; i++)
         output_line[i] = 500 − input_line[i];
     }
16

     /* The main host program controlling and representing the whole
        application */
     int main(int argc, char* argv[]) {
21     int image[HEIGHT][WIDTH];
       unsigned char output[HEIGHT][WIDTH];

       // Initialize with some values
       init_image(WIDTH, HEIGHT, image);
26
       // Draw 70 horizontal lines and map operation on 8 PEs:
     #pragma omp parallel for num_threads(8)
       for(int proc = 0; proc < 70; proc++)
         // Each iteration is on a different PE in parallel:
31   #pragma smecy map(PE, proc & 7)                    \
                   arg(2, in, [1][LINE_SIZE])           \
                   arg(3, out, [1][LINE_SIZE])
         // Invert an horizontal line:
         invert_vector(LINE_SIZE,
36                     &image[HEIGHT − 20 − proc][WIDTH/2 + 2*proc],
                       &image[HEIGHT − 20 − proc][WIDTH/2 + 2*proc]);

       /* Here we guess we have 5 hardware accelerators and we launch
          operations on them: */
41   #pragma omp parallel for num_threads(5)
```

```
      for (int proc = 0; proc < 5; proc++) {
      /* This is need to express the fact that our accelerator only accept
         continuous data but we want apply them on non contiguous data in
         the array */
46    int input_line [LINE_SIZE];
      int output_line [LINE_SIZE];
      /* We need to remap data in the good shape. The compiler should use
         the remapping information to generate DMA transfer for example and
         remove input_line array */
51    SMECY_remap_int2D_to_int1D (HEIGHT, WIDTH, HEIGHT/3, 30 + 20*proc,
                                    LINE_SIZE, 1, image,
                                    LINE_SIZE, input_line);
      // Each iteration is on a different PE in parallel:
   #pragma smecy map(PE, proc) arg(2, in, [LINE_SIZE]) arg(3, out, [LINE_SIZE])
56    invert_vector (LINE_SIZE, input_line, output_line);
      SMECY_remap_int1D_to_int2D (LINE_SIZE, output_line,
                                    HEIGHT, WIDTH, HEIGHT/3, 30 + 20*proc,
                                    LINE_SIZE, 1, image);
   }
61
   // Convert int image to char image:
   normalize_to_char (WIDTH, HEIGHT, image, output);

   write_pgm_image ("remapping_example−output.pgm", WIDTH, HEIGHT, output);
66
   return EXIT_SUCCESS;
}
```

# 4   SMECY embedded C language

We take as input C99 (ISO/IEC 9899:2007) language with extension for embedded systems (TR 18037).

Refer to these document for more information.

# 5   Description of the SMECY directives

The generic format of SMECY code decorations are language dependent, because if here we describe a SMECY IR implementation based on C, it is indeed more general.

- In C/C++:

  #pragma smecy *clause[[,]clause]... newline*

  We can use \ at the end of line for continuation information.

- In Fortran:

  !$smecy *clause[[,]clause]... newline*

  Use & at the end of line for continuation information.

- In other languages: use `#pragma` equivalent, if not available, use comments *à la* Fortran. For example in Python:

  `#$smecy clause[[,]clause]... newline`

  Use also `&` at the end of line for continuation information

In implementations that support a preprocessor, the `_SMECY` macro name is defined to have the decimal value *yyyymm* where *yyyy* and *mm* are the year and month designations of the version of the SMECY API that the implementation supports. If this macro is the subject of a `#define` or a `#undef` preprocessing directive, the behavior is unspecified.

## 5.1 OpenMP support

SMECY is based on OpenMP. It is not clear yet what level of OpenMP is supported. Since a SMECY platform is made at least form a (SMP) control processor, any OpenMP compliant program can run on it anyway.

But the SMECY tools can use more or less information from available OpenMP decorations available in the code.

## 5.2 Mapping on hardware

The mapping of a function call on a specific piece of hardware can be specified with

```
#pragma smecy map(hardware[, unit])
some_function_call(...);
```

- `hardware` is a symbol representing a hardware component of a given target such as `CPU`, `GPP`, `GPU`, `PE`... They are target specific.

- `unit` is an optional instance number for a specific hardware part. This is typically an integer starting a 0. This hardware number can be an expression of the environment to be able to have a loop managing different accelerators.

We can add an `if(scalar-expression)` to predicate hardware launching according some run-time expression to choose between hardware or local software execution, as in OpenMP with the same syntax. The idea is to be able to do a software execution if the data to process is too small compared to the latency of an hardware accelerator.

Recursion is not supported on hardware-mapped functions. If there are functions called from hardware-mapped functions, they will be automatically inlined (so no recursion allowed in them either). If a function is mapped to a more programmable hardware (GPP), recursions in these called functions should be allowed.

## 5.3 Producer/consumer information

To generate hardware communications where there is only a function call, the compiler need to figure out what is the memory zone used to execute the function and then what memory zone in written by the function *and* that will be used later in the program[2]. From these information, copy-in and copy-out operations can be generates.

### 5.3.1 Function arguments

```
#pragma smecy arg(arg_id, arg_clause[, arg_clause]...)
some_function_call(...);
```

**Direction directive** defines how the data flows between the caller and the function:

- `in` the argument is used by the function;
- `out` the argument is used by the function;
- `inout` the argument is used and produced by the function;
- `unused` the argument is not used and produced by the function.

**Argument layout** specifies how the argument is used in the function.

- An optional *array_size_descriptor* such as `[n][m]` expressing that the data is used from the callee point of view as such an array starting at the address given in parameter. If not specified, all the caller argument is used;
- An optional /*array_range_descriptor* restriction such as `/[n:m][2][3:7]` expressing that the data is used from the callee point of view as an array with only this element ranges used. If not specified, all the array is used according to its size specified or not. If only some ranges are lacking, all the matching dimension is used. For example `/[4][]` matches the column 4 of an array.

The more precise this description is and the less data transfers occur.

### 5.3.2 Global variables

Right now we do not deal with sharing information through global variables, because it is more difficult to track. Only function parameters are used to exchange information.

But we can imagine to map global variables with this clause:

```
#pragma smecy global_var(var, arg_clause[, arg_clause]...)
```

## 5.4 Remapping specification

*TO FINISH*

Since we always want a sequential equivalence, that means that the sequential code representing the computation on an accelerator really consume The easy

HPF

---

[2]If a function produces something not used later, it is useless to get it back.

### 5.5   Hardware specific pragma

*To be defined in collaboration with the various hardware suppliers of the project (P2012, EdkDSP...).*

# 6   SMECY high-level API

## 6.1   OpenMP

Since we support OpenMP pragma, we also support OpenMP API that allows for example:

- getting/setting the number of threads;

- getting the number of available processors on the current domain;

- manipulating locks.

## 6.2   MultiCore Association API

We rely on the MCAPI, MTAPI and MRAPI for low memory footprint light-weight communications, threading and synchronization. Refer to their documentations for more information.

There is a reference implementation based on Linux *pthreads* that can be used as an example to port to the various available hardware.

This can be used to express communicating process, asynchronous communications, synchronization, etc.

## 6.3   NPL API

NPL is the API defined to program ST P2012 in a native way. Since it is rather at the same level of MCAPI/MRAPI/MTAPI, it should be easy to implement one above the other. Since the MCA APIs are standard, we thing that it is commercially interesting to provide a MCA APIs over NPL or other to widen P2012 usage.

*Define here NPL*

## 6.4   EdkDSP API

*Define here what is useful in the project.*

## 6.5   OpenCL

Since ST P2012 can be programmed in OpenCL which is also a programming API, a C process can use OpenCL orthogonally with other API. A kernel launching is done by defining the kernel source, the memory zone to use and to transfer and the different parameters of the kernel.

Refer to OpenCL documentation for more information.

# 7 Compilation

## 7.1 OpenMP support

From the programmer point of view it may be equivalent to have

- an OpenMP compiler generating SMECY API code such as the parallelism between SMECY target accelerators is run by OpenMP threads dealing each one with an hardware resource sequentially in a synchronous way;

- or a SMECY compiler understanding the OpenMP syntax and generating directly some parallel execution of SMECY accelerators in an asynchronous way.

## 7.2 Information available

### 7.2.1 Doxygen qualifier

In the Doxygen documentation mark-up language used in comments to detail various entities of a program, there are information such as in, out or inout qualifier on function parameters that can be used to generate the right communication with some hardware accelerators.

### 7.2.2 C qualifier

Qualifiers such as const attribute qualifier, address space names, register names, accumulator, saturation, etc. are used to generate the right target function or instruction.

### 7.2.3 Pragma

Of course the pragma information is heavily used in the compilation process

For example to generate correct communication, the mapping information is used, and from the used/defined information plus optional remapping information, a real communication with an API is used.

```
   void
2  invert_vector(int line_size,
                 int input_line[line_size],
                 int output_line[line_size]) {
     for(int i = 0; i < line_size; i++)
       output_line[i] = 500 - input_line[i];
7  }
     [...]
     int image[HEIGHT][WIDTH];
     [...]
   #pragma smecy map(PE, proc & 7)                    \
12              arg(2, in, [1][LINE_SIZE])           \
                arg(3, out, [1][LINE_SIZE])
     // Invert an horizontal line:
     invert_vector(LINE_SIZE,
                   &image[HEIGHT - 20 - proc][WIDTH/2 + 2*proc],
17                 &image[HEIGHT - 20 - proc][WIDTH/2 + 2*proc]);
```

can be compiled into

```
      int image[HEIGHT][WIDTH];
      /* First prepare the PE #(proc & 7) hardware to execute invert_vector.
3         That may be used to load some program or microcode, reconfigure a
          FPGA, load/compile an OpenCL kernel... */
      SMECY_set(PE, proc & 7, invert_vector);
      /* Send the WIDTH integer as arg 1 on invert_vector hardware function on
          PE #(proc & 7): */
8     SMECY_send_arg(PE, proc & 7, invert_vector, 1, int, LINE_SIZE);
      /* Send a vector of int of size LINE_SIZE as arg 2 on invert_vector
          hardware function on PE #(proc & 7): */
      SMECY_send_arg_vector(PE, proc & 7, invert_vector, 2, int,
                            &image[HEIGHT − 20 − proc][WIDTH/2 + 2*proc], LINE_SIZE);
13    // Launch the hardware function or remote program:
      SMECY_launch(PE, 0, invert_vector);
      /* Get a vector of int of size LINE_SIZE as arg 3 on invert_vector
          hardware function on PE #(proc & 7): */
      SMECY_get_arg_vector(PE, proc & 7, invert_vector, 3, &image[HEIGHT − 20 − proc]
18                                                         [WIDTH/2 + 2*proc],
                            LINE_SIZE);
```

with low level macros described in § 8.

invert_vector() is either an already implemented function in a hardware library, or it is compiled by a target specific compiler with some callable interface.

For more complex calls needing remapping, such as:

```
      int input_line[LINE_SIZE];
      int output_line[LINE_SIZE];
      /* We need to remap data in the good shape. The compiler should use
          the remapping information to generate DMA transfer for example and
5         remove input_line array */
      SMECY_remap_int2D_to_int1D(HEIGHT, WIDTH, HEIGHT/3, 30 + 20*proc,
                                 LINE_SIZE, 1, image,
                                 LINE_SIZE, input_line);
      // Each iteration is on a different PE in parallel:
10    #pragma smecy map(PE, proc) arg(2, in, [LINE_SIZE]) arg(3, out, [LINE_SIZE])
      invert_vector(LINE_SIZE, input_line, output_line);
      SMECY_remap_int1D_to_int2D(LINE_SIZE, output_line,
                                 HEIGHT, WIDTH, HEIGHT/3, 30 + 20*proc,
                                 LINE_SIZE, 1, image);
```

is compiled by using other hardware interfaces involving more complex DMA:

```
1     // May not be useful if this function is already set:
      SMECY_set(PE, proc & 7, invert_vector);
      /* Send the WIDTH integer as arg 1 on invert_vector hardware function on
          PE #(proc & 7): */
      SMECY_send_arg(PE, proc & 7, invert_vector, 1, int, LINE_SIZE);
6     /* Send a vector of int of size LINE_SIZE as arg 2 on invert_vector
          hardware function on PE #(proc & 7) but read as a part of a 2D array: */
      smecy_send_arg_int_DMA_2D_PE_0_invert_vector(3, &image[HEIGHT − 20 − proc]
                                                      [WIDTH/2 + 2*proc],
                                                   HEIGHT, WIDTH, HEIGHT/3, 30 + 20*proc,
11                                                    LINE_SIZE, 1);
      SMECY_send_arg_DMA_2D_to_1D(PE, proc & 7, invert_vector, 2, int,
                                  &image[HEIGHT − 20 − proc][WIDTH/2 + 2*proc],
                                  HEIGHT, WIDTH, HEIGHT/3, 30 + 20*proc,
```

```
                                              LINE_SIZE, 1);
16   // Launch  the  hardware  function  or  remote  program:
     SMECY_launch(PE, 0, invert_vector);
     SMECY_get_arg_DMA_1D_to_2D(PE, proc & 7, invert_vector, 3, int,
                            &image[HEIGHT − 20 − proc][WIDTH/2 + 2*proc],
                            HEIGHT, WIDTH, HEIGHT/3, 30 + 20*proc,
21                                            LINE_SIZE, 1);
```

### 7.3 Geometric inference

In the OpenMP SMP model, there is a global memory (well, with a weak coherence model) that may not exist in the execution model of a given target. For example, even if 2 hardware accelerators exchange information through memory according to the high-level programming model, in the real world we may have 2 processors communicating through message passing with MCAPI on a NoC or 2 hardware accelerators connected through a pipeline.

To solve this issue, we use the OpenMP global memory like a scoreboard memory that is used to symbolically relate all the data flows and generates various communication schemes.

Since the memory dependencies are expresses by hyperparallelepipedes, we can do some intersection analysis to compute if a communication is needed or not between 2 devices of the target.

## 8   SMECY low level hardware API

To call real hardware accelerators, few C macros are needed to interface a program running on some processor to a function running on another processor or in some hardware accelerator.

Since the implementation may depend also on the processor calling the macros (not the same IO bus will be used from an $x86$ or a DSP to call the same operator), a global preprocessing symbol must be defined by the compiler before using these macros, such as

**#define** SMECY_LOCAL_PROC x86

or

**#define** SMECY_LOCAL_PROC DSP

Few macros are necessary:

```
/* Prepare  a  processing  element  to  execute  a  function

    @param pe is  the  symbol  of  a  processing  element,  such  as  GPP,  DSP,  PE...
4   @param[in] instance  is  the  instance  number  of  the  processor  element  to  use
    @param func is  the  function  to  load  on  the  processor  element
*/
#define SMECY_set(pe, instance, func) ...

9  /* Send  a  scalar  argument  to  a  function  on  a  processing  element

    @param pe is  the  symbol  of  a  processing  element,  such  as  GPP,  DSP,  PE...
    @param[in] instance  is  the  instance  number  of  the  processor  element  to  use
```

```
        @param func is the function to load on the processor element
14      @param[in] arg is the argument instance to set
        @param type is the type of the scalar argument to send
        @param[in] val is the value of the argument to send
   */
   #define SMECY_send_arg(pe, instance, func, arg, type, val) ...
19
   /* Send a vector argument to a function on a processing element

        @param pe is the symbol of a processing element, such as GPP, DSP, PE...
        @param[in] instance is the instance number of the processor element to use
24      @param func is the function to load on the processor element
        @param[in] arg is the argument instance to set
        @param type is the type of the vector element to send
        @param[in] addr is the starting address of the vector to read from
                    caller memory
29      @param[in] size is the length of the vector
   */
   #define SMECY_send_arg_vector(pe, instance, func, arg, type, addr, size) ...

   /* Launch the hardware function or remote program using previously loaded
34      arguments

        @param pe is the symbol of a processing element, such as GPP, DSP, PE...
        @param[in] instance is the instance number of the processor element to use
        @param func is the function to load on the processor element
39
        A kernel can be launched several times without having to set/reset its function.
   */
   #define SMECY_launch(pe, instance, func) ...

44 /* Get the return value of a function on a processing element

        @param pe is the symbol of a processing element, such as GPP, DSP, PE...
        @param[in] instance is the instance number of the processor element to use
        @param func is the function to load on the processor element
49      @param type is the type of the scalar argument to send
        @return the value computed by the function
   */
   #define SMECY_get_return(pe, instance, func, type) ...

54 /* Get a vector value computed by a function on a processing element

        @param pe is the symbol of a processing element, such as GPP, DSP, PE...
        @param[in] instance is the instance number of the processor element to use
        @param func is the function to load on the processor element
59      @param[in] arg is the argument instance to retrieve
        @param type is the type of the vector element
        @param[out] addr is the starting address of the vector to write in
                    caller memory
        @param[in] size is the length of the vector
64 */
   #define SMECY_get_arg_vector(pe, instance, func, arg, type, addr, size) ...
```

```
   /* Reset a processing element to execute a function

69    @param pe is the symbol of a processing element, such as GPP, DSP, PE...
      @param[in] instance is the instance number of the processor element to use
      @param func is the function to unload from the processor element

      This is used for example to remove consuming resources to decrease
74    power. Giving here the function name may be useful for weird case to
      avoid having short−circuit between CLB in a FPGA during unconfiguring stuff
   */
   #define SMECY_reset(pe, instance, func) ...
```

There is also a macro for an asynchronous call and to wait for completion.