

# OPS C++ User's Manual

Mike Giles, Istvan Reguly, Gihan Mudalige

December 2015

# Contents

|          |   |          |
|----------|---|----------|
| <b>1</b> | <b>Introduction</b>                             | <b>4</b> |
| <b>2</b> | <b>Key concepts and structure</b>               | <b>4</b> |
| <b>3</b> | <b>OPS C++ API</b>                              | <b>6</b> |
| 3.1      | Initialisation and termination routines         | 6        |
|          | ops_init  | 6        |
|          | ops_decl_block                                  | 6        |
|          | ops_decl_block_hdf5                             | 6        |
|          | ops_decl_dat                                    | 6        |
|          | ops_decl_dat_hdf5                               | 7        |
|          | ops_decl_const                                  | 7        |
|          | ops_update_const                                | 8        |
|          | ops_decl_halo                                   | 8        |
|          | ops_decl_halo_hdf5                              | 8        |
|          | ops_decl_halo_group                             | 8        |
|          | ops_decl_reduction_handle                       | 9        |
|          | ops_partition                                   | 9        |
|          | ops_diagnostic_output                           | 9        |
|          | ops_printf                                      | 9        |
|          | ops_timers                                      | 9        |
|          | ops_fetch_block_hdf5_file                       | 9        |
|          | ops_fetch_stencil_hdf5_file                     | 10       |
|          | ops_fetch_dat_hdf5_file                         | 10       |
|          | ops_fetch_dat                                   | 10       |
|          | ops_print_dat_to_txtfile                        | 10       |
|          | ops_timing_output                               | 10       |
|          | ops_exit  | 10       |
| 3.2      | Halo exchange                                   | 11       |
|          | ops_halo_transfer                               | 11       |
| 3.3      | Parallel loop syntax                            | 12       |
|          | ops_par_loop                                    | 12       |
|          | ops_arg_gbl                                     | 12       |
|          | ops_arg_reduce                                  | 12       |
|          | ops_arg_dat                                     | 13       |
|          | ops_arg_idx                                     | 13       |
| 3.4      | Stencils  | 14       |
|          | ops_decl_stencil                                | 14       |
|          | ops_decl_strided_stencil                        | 14       |
|          | ops_decl_stencil_hdf5                           | 14       |
| 3.5      | Checkpointing                                   | 15       |
|          | ops_checkpointing_init                          | 15       |
|          | ops_checkpointing_manual_datlist                | 16       |
|          | ops_checkpointing_fastfw                        | 16       |
|          | ops_checkpointing_manual_datlist_fastfw         | 16       |
|          | ops_checkpointing_manual_datlist_fastfw_trigger | 16       |

|   |                           |    |
|---|---------------------------|----|
| 4 | Tiling for Cache-blocking | 17 |
| 5 | OPS User Kernels          | 18 |

# 1 Introduction

OPS is a high-level framework with associated libraries and preprocessors to generate parallel executables for applications on **multi-block structured grids**. Multi-block structured grids consists of an unstructured collection of structured meshes/grids. This document describes the OPS C++ API, which supports the development of single-block and multi-block structured meshes.

Many of the API and library follows the structure of the OP2 high-level library for unstructured mesh applications [1]. However the structured mesh domain is distinct from the unstructured mesh applications domain due to the implicit connectivity between neighbouring mesh elements (such as vertices, cells) in structured meshes/grids. The key idea is that operations involve looping over a “rectangular” multi-dimensional set of grid points using one or more “stencils” to access data. In multi-block grids, we have several structured blocks. The connectivity between the faces of different blocks can be quite complex, and in particular they may not be oriented in the same way, i.e. an  $i, j$  face of one block may correspond to the  $j, k$  face of another block. This is awkward and hard to handle simply.

To clarify some of the important issues in designing the API, we note here some needs connected with a 3D application:

- When looping over the interior with loop indices  $i, j, k$ , often there are 1D arrays which are referenced using just one of the indices.
- To implement boundary conditions, we often loop over a 2D face, accessing both the 3D dataset and data from a 2D dataset.
- To implement periodic boundary conditions using dummy “halo” points, we sometimes have to copy one plane of boundary data to another. e.g. if the first dimension has size  $I$  then we might copy the plane  $i = I - 2$  to plane  $i = 0$ , and plane  $i = 1$  to plane  $i = I - 1$ .
- In multigrid, we are working with two grids with one having twice as many points as the other in each direction. To handle this we require a stencil with a non-unit stride.
- In multi-block grids, we have several structured blocks. The connectivity between the faces of different blocks can be quite complex, and in particular they may not be oriented in the same way, i.e. an  $i, j$  face of one block may correspond to the  $j, k$  face of another block. This is awkward and hard to handle simply.

The latest proposal is to handle all of these different requirements through stencil definitions.

## 2 Key concepts and structure

An OPS applicaiton can generally be divided into two key parts: initialisation and parallel exeution. During the initialisation phase, one or more blocks (ops\_block) are defined: these only have a dimensionality (i.e. 1D, 2D, etc.), and serve to group datasets together. Datasets are defined on a block, and have a specific size (in each dimension of the block), which may be slightly different across different datasets (e.g. staggered grids), in some directions they may be degenerate (a size of 1), or they can represent data associated with different multigrid levels (where their size if a multiple or a fraction of other datasets). Datasets can be declared with empty (NULL) pointers, then OPS will allocate the appropriate amount of memory, may be passed non-NULL pointers (currently only supported in non-MPI environments), in which case OPS will assume the memory is large enough for the data and the block halo, and there are HDF5 dataset declaration routines

which allow the distributed reading of datasets from HDF5 files. The concept of blocks is necessary to group datasets together, as in a multi-block problem, in a distributed memory environment, OPS needs to be able to determine how to decompose the problem.

The initialisation phase usually also consists of defining the stencils to be used later on (though they can be defined later as well), which describe the data access patterns used in parallel loops. Stencils are always relative to the “current” point; e.g. if at iteration  $(i,j)$ , we wish to access  $(i-1,j)$  and  $(i,j)$ , then the stencil will have two points:  $(0, 0, -1, 0)$ . To support degenerate datasets (where in one of the dimensions the dataset’s size is 1), as well as for multigrid, there are special strided, restriction, and prolongation stencils: they differ from normal stencils in that as one steps through a grid in a parallel loop, the stepping is done with a non-unit stride for these datasets. For example, in a 2D problem, if we have a degenerate dataset called `xcoords`, size  $(N,1)$ , then we will need a stencil with stride  $(1,0)$  to access it in a regular 2D loop.

Finally, the initialisation phase may declare a number of global constants - these are variables in global scope that can be accessed from within user kernels, without having to pass them in explicitly. These may be scalars or small arrays, generally for values that do not change during execution, though they may be updated during execution with repeated calls to `ops_decl_const`.

The initialisation phase is terminated by a call to `ops_partition`.

The bulk of the application consists of parallel loops, implemented using calls to `ops_parallel_loop`. These constructs work with datasets, passed through the opaque `ops_dat` handles declared during the initialisation phase. The iterations of parallel loops are semantically independent, and it is the responsibility of the user to enforce this: the order in which iterations are executed cannot affect the result (within the limits of floating point precision). Parallel loops are defined on a block, with a prescribed iteration range that is always defined from the perspective of the dataset written/modified (the sizes of datasets, particularly in multigrid situations may be very different). Datasets are passed in using `ops_arg_dat`, and during execution, a pointer to the values at the current grid point will be passed to the user kernel. This pointer needs to be dereferenced using the `OPS_ACC` macros, where the user may specify the relative offset to access the grid point’s neighbours (which accesses have to match the declared stencil). Datasets written may only be accessed with a one-point, zero-offset stencil (otherwise the parallel semantics may be violated).

Other than datasets, one can pass in read-only scalars or small arrays that are iteration space invariant with `ops_arg_gbl` (typically weights,  $\delta t$ , etc. which may be different in different loops). The current iteration index can also be passed in with `ops_arg_idx`, which will pass a globally consistent index to the user kernel (i.e. also under MPI).

Reductions in loops are done using the `ops_arg_reduce` argument, which takes a reduction handle as an argument. The result of the reduction can then be acquired using a separate call to `ops_reduction_result`. The semantics are the following: a reduction handle after it was declared is in an “uninitialised” state. The first time it is used as an argument to a loop, its type is determined (increment/min/max), and is initialised appropriately  $(0, \text{inf}, -\text{inf})$ , and subsequent uses of the handle in parallel loops are combined together, up until the point, where the result is acquired using `ops_reduction_result`, which then sets it back to an uninitialised state. This also implies, that different parallel loops, which all use the same reduction handle, but are otherwise independent, are independent and their partial reduction results can be combined together associatively and commutatively.

OPS takes responsibility for all data, its movement and the execution of parallel loops. With different execution hardware and optimisations, this means OPS will re-organise data as well as execution (potentially across different loops), and therefore any data accesses or manipulation may only be done through the OPS API.

## 3 OPS C++ API

### 3.1 Initialisation and termination routines

**void ops\_init(int argc, char \*\*argv, int diags\_level)**

This routine must be called before all other OPS routines.

|                    |   |
|--------------------|---|
| <b>argc, argv</b>  | the usual command line arguments  |
| <b>diags_level</b> | an integer which defines the level of debugging diagnostics and reporting to be performed |

Currently, higher **diags\_level**s does the following checks

**diags\_level** = 1 : no diagnostics, default to achieve best runtime performance.

**diags\_level** > 1 : print block decomposition and **ops\_par\_loop** timing breakdown.

**diags\_level** > 4 : print intra-block halo buffer allocation feedback (for OPS internal development only)

**diags\_level** > 5 : check if intra-block halo MPI sends depth match MPI receives depth (for OPS internal development only)

**ops\_block ops\_decl\_block(int dims, char \*name)**

This routine defines a structured grid block.

|             |                                    |
|-------------|------------------------------------|
| <b>dims</b> | dimension of the block             |
| <b>name</b> | a name used for output diagnostics |

**ops\_block ops\_decl\_block\_hdf5(int dims, char \*name, char \*file)**

This routine reads the details of a structured grid block from a named HDF5 file

|             |   |
|-------------|---|
| <b>dims</b> | dimension of the block                                  |
| <b>name</b> | a name used for output diagnostics                      |
| <b>file</b> | hdf5 file to read and obtain the block information from |

Although this routine does not read in any extra information about the block from the named HDF5 file than what is already specified in the arguments, it is included here for error checking (e.g. check if blocks defined in an HDF5 file is matching with the declared arguments in an application) and completeness.

**ops\_dat ops\_decl\_dat(ops\_block block, int dim, int\* size, int \*base, int \*d\_m, int \*d\_p, T \*data, char \*type, char \*name)**

This routine defines a dataset.

|              |  |
|--------------|--|
| <b>block</b> | structured block   |
| <b>dim</b>   | dimension of dataset (number of items per grid element)                                  |
| <b>size</b>  | size in each dimension of the block  |
| <b>base</b>  | base indices in each dimension of the block  |
| <b>d_m</b>   | padding from the face in the negative direction for each dimension (used for block halo) |
| <b>d_p</b>   | padding from the face in the positive direction for each dimension (used for block halo) |
| <b>data</b>  | input data of type T   |
| <b>type</b>  | the name of type used for output diagnostics (e.g. “double”, “float”)                    |
| <b>name</b>  | a name used for output diagnostics   |

The **size** allows to declare different sized data arrays on a given **block**. **d\_m** and **d\_p** are depth of the “block halos” that are used to indicate the offset from the edge of a block (in both the negative and positive directions of each dimension).

**ops\_dat ops\_decl\_dat\_hdf5(ops\_block block, int dim, char \*type, char \*name, char \*file)**

This routine defines a dataset to be read in from a named hdf5 file

|              |   |
|--------------|---|
| <b>block</b> | structured block  |
| <b>dim</b>   | dimension of dataset (number of items per grid element)               |
| <b>type</b>  | the name of type used for output diagnostics (e.g. “double”, “float”) |
| <b>name</b>  | name of the dat used for output diagnostics                           |
| <b>file</b>  | hdf5 file to read and obtain the data from                            |

**void ops\_decl\_const(char const \* name, int dim, char const \* type, T \* data )**

This routine defines a global constant: a variable in global scope. Global constants need to be declared upfront so that they can be correctly handled for different parallelizations. For e.g CUDA on GPUs. Once defined they remain unchanged throughout the program, unless changed by a call to ops.update\_const(..)

|             |   |
|-------------|---|
| <b>name</b> | a name used to identify the constant                                  |
| <b>dim</b>  | dimension of dataset (number of items per element)                    |
| <b>type</b> | the name of type used for output diagnostics (e.g. “double”, “float”) |
| <b>data</b> | pointer to input data of type T                                       |

**void ops\_update\_const(char const \* name, int dim, char const \* type, T \* data)**

This routine updates/changes the value of a constant

|             |   |
|-------------|---|
| <b>name</b> | a name used to identify the constant                                  |
| <b>dim</b>  | dimension of dataset (number of items per element)                    |
| <b>type</b> | the name of type used for output diagnostics (e.g. "double", "float") |
| <b>data</b> | pointer to new values for constant of type T                          |

**ops\_halo ops\_decl\_halo(ops\_dat from, ops\_dat to, int \*iter\_size, int\* from\_base, int \*to\_base, int \*from\_dir, int \*to\_dir)**

This routine defines a halo relationship between two datasets defined on two different blocks.

|                  |   |
|------------------|---|
| <b>from</b>      | origin dataset  |
| <b>to</b>        | destination dataset   |
| <b>iter_size</b> | defines an iteration size (number of indices to iterate over in each direction) |
| <b>from_base</b> | indices of starting point in "from" dataset                                     |
| <b>to_base</b>   | indices of starting point in "to" dataset                                       |
| <b>from_dir</b>  | direction of incrementing for "from" for each dimension of <b>iter_size</b>     |
| <b>to_dir</b>    | direction of incrementing for "to" for each dimension of <b>iter_size</b>       |

A from\_dir [1,2] and a to\_dir [2,1] means that x in the first block goes to y in the second block, and y in first block goes to x in second block. A negative sign indicates that the axis is flipped. (Simple example: a transfer from (1:2,0:99,0:99) to (-1:0,0:99,0:99) would use iter\_size = [2,100,100], from\_base = [1,0,0], to\_base = [-1,0,0], from\_dir = [0,1,2], to\_dir = [0,1,2]. In more complex case this allows for transfers between blocks with different orientations.)

**ops\_halo ops\_decl\_halo\_hdf5(ops\_dat from, ops\_dat to, char\* file)**

This routine reads in a halo relationship between two datasets defined on two different blocks from a named HDF5 file

|             |  |
|-------------|--|
| <b>from</b> | origin dataset                             |
| <b>to</b>   | destination dataset                        |
| <b>file</b> | hdf5 file to read and obtain the data from |

**ops\_halo\_group ops\_decl\_halo\_group(int nhalos, ops\_halo \*halos)**

This routine defines a collection of halos. Semantically, when an exchange is triggered for all halos in a group, there is no order defined in which they are carried out.

|               |                                 |
|---------------|---------------------------------|
| <b>nhalos</b> | number of halos in <b>halos</b> |
| <b>halos</b>  | array of halos                  |



**ops\_reduction ops\_decl\_reduction\_handle(int size, char \*type, char \*name)**

This routine defines a reduction handle to be used in a parallel loop

|             |   |
|-------------|---|
| <b>size</b> | size of data in bytes   |
| <b>type</b> | the name of type used for output diagnostics (e.g. "double", "float") |
| <b>name</b> | name of the dat used for output diagnostics                           |

**void ops\_reduction\_result(ops\_reduction handle, T \*result)**

This routine returns the reduced value held by a reduction handle

|               |  |
|---------------|--|
| <b>handle</b> | the ops_reduction handle   |
| <b>result</b> | a pointer to write the results to, memory size has to match the declared |

**ops\_partition(char \*method)**

Triggers a multi-block partitioning across a distributed memory set of processes. (links to a dummy function for single node parallelizations). This routine should only be called after all the ops\_halo ops\_decl\_block and ops\_halo ops\_decl\_dat statements have been declared

|               |   |
|---------------|---|
| <b>method</b> | string describing the partitioning method. Currently this string is not used internally, but is simply a place-holder to indicate different partitioning methods in the future. |
|---------------|---|

**void ops\_diagnostic\_output()**

This routine prints out various useful bits of diagnostic info about sets, mappings and datasets. Usually used right after an ops\_partition() call to print out the details of the decomposition

**void ops\_printf(const char \* format, ...)**

This routine simply prints a variable number of arguments; it is created in place of the standard C printf function which would print the same on each MPI process

**void ops\_timers(double \*cpu, double \*et)**

gettimeofday() based timer to start/end timing blocks of code

|            |   |
|------------|---|
| <b>cpu</b> | variable to hold the CPU time at the time of invocation     |
| <b>et</b>  | variable to hold the elapsed time at the time of invocation |

**void ops\_fetch\_block\_hdf5\_file(ops\_block block, char \*file)**

Write the details of an ops\_block to a named HDF5 file. Can be used over MPI (puts the data in an ops\_dat into an HDF5 file using MPI I/O)

|              |                         |
|--------------|-------------------------|
| <b>block</b> | ops_block to be written |
| <b>file</b>  | hdf5 file to write to   |

**void ops\_fetch\_stencil\_hdf5\_file(ops\_stencil stencil, char \*file)**

Write the details of an ops\_block to a named HDF5 file. Can be used over MPI (puts the data in an ops\_dat into an HDF5 file using MPI I/O)

|                |                           |
|----------------|---------------------------|
| <b>stencil</b> | ops_stencil to be written |
| <b>file</b>    | hdf5 file to write to     |

**void ops\_fetch\_dat\_hdf5\_file(ops\_dat dat, const char \*file)**

Write the details of an ops\_block to a named HDF5 file. Can be used over MPI (puts the data in an ops\_dat into an HDF5 file using MPI I/O)

|             |                       |
|-------------|-----------------------|
| <b>dat</b>  | ops_dat to be written |
| <b>file</b> | hdf5 file to write to |

**void ops\_fetch\_dat(ops\_dat dat, T \*data)**

Makes a copy of the data held internally by an ops\_dat and returns a pointer to that block of data. The memory block will be a flat 1D block of memory (similar to how any multi-dimensional data is held internally by OPS. Over MPI the pointer will point to a block of memory that holds the data held only within that MPI proc.

|             |   |
|-------------|---|
| <b>dat</b>  | ops_dat to fetch data from                              |
| <b>data</b> | pointer of type T pointing to memory block holding data |

**void ops\_print\_dat\_to\_txtfile(ops\_dat dat, char \*file)**

Write the details of an ops\_block to a named text file. When used under an MPI parallelization each MPI process will write its own data set separately to the text file. As such it does not use MPI I/O. The data can be viewed using a simple text editor

|             |                          |
|-------------|--------------------------|
| <b>dat</b>  | ops_dat to to be written |
| <b>file</b> | text file to write to    |

**void ops\_timing\_output(FILE \*os)**

Print OPS performance performance details to output stream

|           |  |
|-----------|--|
| <b>os</b> | output stream, use stdout to print to standard out |
|-----------|--|

**void ops\_exit()**

This routine must be called last to cleanly terminate the OPS computation.

## 3.2 Halo exchange

**void ops\_halo\_transfer(ops\_halo\_group group)**

This routine exchanges all halos in a halo group and will block execution of subsequent computations that depend on the exchanged data.

**group**                      the halo group

### 3.3 Parallel loop syntax

A parallel loop with N arguments has the following syntax:

```
void ops_par_loop( void (*kernel)(...),
                  char *name, ops_blk block, int dims, int *range,
                  ops_arg arg1, ops_arg arg2, ..., ops_arg argN )
```

|               |  |
|---------------|--|
| <b>kernel</b> | user's kernel function with N arguments              |
| <b>name</b>   | name of kernel function, used for output diagnostics |
| <b>block</b>  | the ops_block over which this loop executes          |
| <b>dims</b>   | dimension of loop iteration                          |
| <b>range</b>  | iteration range array                                |
| <b>args</b>   | arguments  |

The **ops\_arg** arguments in **ops\_par\_loop** are provided by one of the following routines, one for global constants and reductions, and the other for OPS datasets.

```
ops_arg ops_arg_gbl(T *data, int dim, char *type, ops_access acc)
```

Passes a scalar or small array that is invariant of the iteration space (not to be confused with ops\_decl\_const, which facilitates global scope variables).

|             |   |
|-------------|---|
| <b>data</b> | data array  |
| <b>dim</b>  | array dimension                                   |
| <b>type</b> | string representing the type of data held in data |
| <b>acc</b>  | access type                                       |

```
ops_arg ops_arg_reduce(ops_reduction handle, int dim, char *type, ops_access acc)
```

Passes a pointer to a variable that needs to be incremented (or swapped for min/max reduction) by the user kernel.

|               |   |
|---------------|---|
| <b>handle</b> | an ops_reduction handle                           |
| <b>dim</b>    | array dimension (according to type)               |
| <b>type</b>   | string representing the type of data held in data |
| <b>acc</b>    | access type                                       |

**ops\_arg ops\_arg\_dat(ops\_dat dat, ops\_stencil stencil, char \*type, ops\_access acc)**

Passes a pointer to the value(s) at the current grid point to the user kernel. the OPS\_ACC macro has to be used for dereferencing the pointer.

|                |  |
|----------------|--|
| <b>dat</b>     | dataset  |
| <b>stencil</b> | stencil for accessing data                           |
| <b>type</b>    | string representing the type of data held in dataset |
| <b>acc</b>     | access type  |

**ops\_arg ops\_arg\_idx()**

Give you an array of integers (in the user kernel) that have the index of the current grid point, i.e. idx[0] is the index in x, idx[1] is the index in y, etc. This is a globally consistent index, so even if the block is distributed across different MPI partitions, it gives you the same indexes. Generally used to generate initial geometry.

### 3.4 Stencils

The final ingredient is the stencil specification, for which we have two versions: simple and strided.

**ops\_stencil ops\_decl\_stencil(int dims, int points, int \*stencil, char \*name)**

|                |   |
|----------------|---|
| <b>dims</b>    | dimension of loop iteration                 |
| <b>points</b>  | number of points in the stencil             |
| <b>stencil</b> | stencil for accessing data                  |
| <b>name</b>    | string representing the name of the stencil |

**ops\_stencil ops\_decl\_strided\_stencil(int dims, int points,  
int \*stencil, int \*stride, char \*name)**

|                |   |
|----------------|---|
| <b>dims</b>    | dimension of loop iteration                 |
| <b>points</b>  | number of points in the stencil             |
| <b>stencil</b> | stencil for accessing data                  |
| <b>stride</b>  | stride for accessing data                   |
| <b>name</b>    | string representing the name of the stencil |

**ops\_stencil ops\_decl\_stencil\_hdf5(int dims, int points, char \*name, char\* file)**

|               |   |
|---------------|---|
| <b>dims</b>   | dimension of loop iteration                 |
| <b>points</b> | number of points in the stencil             |
| <b>name</b>   | string representing the name of the stencil |
| <b>file</b>   | hdf5 file to write to                       |

In the strided case, the semantics for the index of data to be accessed, for stencil point **p**, in dimension **m** are defined as:

**stride[m]\*loop\_index[m] + stencil[p\*dims+m],**

where **loop\_index[m]** is the iteration index (within the user-defined iteration space) in the different dimensions.

If, for one or more dimensions, both **stride[m]** and **stencil[p\*dims+m]** are zero, then one of the following must be true;

- the dataset being referenced has size 1 for these dimensions
- these dimensions are to be omitted and so the dataset has dimension equal to the number of remaining dimensions.

These two stencil definitions probably take care of all of the cases in the Introduction except for multiblock applications with interfaces with different orientations – this will need a third, even more general, stencil specification. The strided stencil will handle both multigrid (with a stride of 2 for example) and the boundary condition and reduced dimension applications (with a stride of 0 for the relevant dimensions).

### 3.5 Checkpointing

OPS supports the automatic checkpointing of applications. Using the API below, the user specifies the file name for the checkpoint and an average time interval between checkpoints, OPS will then automatically save all necessary information periodically that is required to fast-forward to the last checkpoint if a crash occurred. Currently, when re-launching after a crash, the same number of MPI processes have to be used. To enable checkpointing mode, the `OPS_CHECKPOINT` runtime argument has to be used.

**bool ops\_checkpointing\_init(const char \*filename, double interval, int options)**

Initialises the checkpointing system, has to be called after `ops_partition`. Returns true if the application launches in restore mode, false otherwise.

|                 |  |
|-----------------|--|
| <b>filename</b> | name of the file for checkpointing. In MPI, this will automatically be post-fixed with the rank ID.  |
| <b>interval</b> | average time (seconds) between checkpoints   |
| <b>options</b>  | <p>a combinations of flags, listed in <code>ops_checkpointing.h</code>:</p> <p><code>OPS_CHECKPOINT_INITPHASE</code> - indicates that there are a number of parallel loops at the very beginning of the simulations which should be excluded from any checkpoint; mainly because they initialise datasets that do not change during the main body of the execution. During restore mode these loops are executed as usual. An example would be the computation of the mesh geometry, which can be excluded from the checkpoint if it is re-computed when recovering and restoring a checkpoint. The API call <code>void ops_checkpointing_initphase_done()</code> indicates the end of this initial phase.</p> <p><code>OPS_CHECKPOINT_MANUAL_DATLIST</code> - Indicates that the user manually controls the location of the checkpoint, and explicitly specifies the list of <code>ops_dats</code> to be saved.</p> <p><code>OPS_CHECKPOINT_FASTFW</code> - Indicates that the user manually controls the location of the checkpoint, and it also enables fast-forwarding, by skipping the execution of the application (even though none of the parallel loops would actually execute, there may be significant work outside of those) up to the checkpoint.</p> |

OPS\_CHECKPOINT\_MANUAL - Indicates that when the corresponding API function is called, the checkpoint should be created. Assumes the presence of the above two options as well.

**void ops\_checkpointing\_manual\_datlist(int ndats, ops\_dat \*datlist)**

A use can call this routine at a point in the code to mark the location of a checkpoint. At this point, the list of datasets specified will be saved. The validity of what is saved is not checked by the checkpointing algorithm assuming that the user knows what data sets to be saved for full recovery. This routine should be called frequently (compared to check-pointing frequency) and it will trigger the creation of the checkpoint the first time it is called after the timeout occurs.

|                |  |
|----------------|--|
| <b>ndats</b>   | number of datasets to be saved               |
| <b>datlist</b> | arrays of <b>ops_dat</b> handles to be saved |

**bool ops\_checkpointing\_fastfw(int nbytes, char \*payload)**

A use can call this routine at a point in the code to mark the location of a checkpoint. At this point, the specified payload (e.g. iteration count, simulation time, etc.) along with the necessary datasets, as determined by the checkpointing algorithm will be saved. This routine should be called frequently (compared to checkpointing frequency), will trigger the creation of the checkpoint the first time it is called after the timeout occurs. In restore mode, will restore all datasets the first time it is called, and returns true indicating that the saved payload is returned in payload. Does not save reduction data.

|                |  |
|----------------|--|
| <b>nbytes</b>  | size of the payload in bytes                       |
| <b>payload</b> | pointer to memory into which the payload is packed |

**bool ops\_checkpointing\_manual\_datlist\_fastfw(int ndats, op\_dat \*datlist, int nbytes, char \*payload)**

Combines the manual datlist and fastfw calls.

|                |  |
|----------------|--|
| <b>ndats</b>   | number of datasets to be saved                     |
| <b>datlist</b> | arrays of <b>ops_dat</b> handles to be saved       |
| <b>nbytes</b>  | size of the payload in bytes                       |
| <b>payload</b> | pointer to memory into which the payload is packed |

**bool ops\_checkpointing\_manual\_datlist\_fastfw\_trigger(int ndats, opa\_dat \*datlist, int nbytes, char \*payload)**

With this routine it is possible to manually trigger checkpointing, instead of relying on the timeout process. as such it combines the manual datlist and fastfw calls, and triggers the creation of a checkpoint when called.



|                      |  |
|----------------------|--|
| <code>ndats</code>   | number of datasets to be saved                     |
| <code>datlist</code> | arrays of <code>ops_dat</code> handles to be saved |
| <code>nbytes</code>  | size of the payload in bytes                       |
| <code>payload</code> | pointer to memory into which the payload is packed |

The suggested use of these **manual** functions is of course when the optimal location for checkpointing is known - one of the ways to determine that is to use the built-in algorithm. More details of this will be reported in a tech-report on checkpointing, to be published later.

## 4 Tiling for Cache-blocking

OPS has a code generation (`ops_gen_mpi_lazy`) and build target for tiling. Currently this is only supported on a single node with OpenMP. Once compiled, to enable, use the `OPS_TILING` runtime parameter - this will look at the L3 cache size of your CPU and guess the correct tile size. If you want to alter the amount of cache to be used for the guess, use the `OPS_CACHE_SIZE=XX` runtime parameter, where the value is in Megabytes. To manually specify the tile sizes, use the `T1`, `T2`, and `T3` environment variables.

To test, please compile CloverLeaf under `apps/c/CloverLeaf`, modify `clover.in` to use a  $6144^2$  mesh, then run as follows: `OMP_NUM_THREADS=xx numactl -physnodebind=0 ./cloverleaf_tiled OPS_TILING`

## 5 OPS User Kernels

In OPS, the elemental operation carried out per mesh/grid point is specified as an outlined function called a *user kernel*. An example taken from the Cloverleaf application is given in Figure 1.

---

```
1 void accelerate_kernel( const double *density0, const double *volume,
2                        double *stepbymass, const double *xvel0, double *xvel1,
3                        const double *xarea, const double *pressure,
4                        const double *yvel0, double *yvel1,
5                        const double *yarea, const double *viscosity) {
6
7     double nodal_mass;
8     //{0,0, -1,0, 0,-1, -1,-1};
9     nodal_mass = ( density0[OPS_ACC0(-1,-1)] * volume[OPS_ACC1(-1,-1)]
10 + density0[OPS_ACC0(0,-1)] * volume[OPS_ACC1(0,-1)]
11 + density0[OPS_ACC0(0,0)] * volume[OPS_ACC1(0,0)]
12 + density0[OPS_ACC0(-1,0)] * volume[OPS_ACC1(-1,0)] ) * 0.25;
13
14     stepbymass[OPS_ACC2(0,0)] = 0.5*dt/ nodal_mass;
15
16     //{0,0, -1,0, 0,-1, -1,-1};
17     //{0,0, 0,-1};
18     xvel1[OPS_ACC4(0,0)] = xvel0[OPS_ACC3(0,0)] - stepbymass[OPS_ACC2(0,0)] *
19 ( xarea[OPS_ACC5(0,0)] * ( pressure[OPS_ACC6(0,0)] - pressure[OPS_ACC6(-1,0)] ) +
20 xarea[OPS_ACC5(0,-1)] * ( pressure[OPS_ACC6(0,-1)] - pressure[OPS_ACC6(-1,-1)] ) );
21
22     //{0,0, -1,0, 0,-1, -1,-1};
23     //{0,0, -1,0};
24     yvel1[OPS_ACC8(0,0)] = yvel0[OPS_ACC7(0,0)] - stepbymass[OPS_ACC2(0,0)] *
25 ( yarea[OPS_ACC9(0,0)] * ( pressure[OPS_ACC6(0,0)] - pressure[OPS_ACC6(0,-1)] ) +
26 yarea[OPS_ACC9(-1,0)] * ( pressure[OPS_ACC6(-1,0)] - pressure[OPS_ACC6(-1,-1)] ) );
27
28     //{0,0, -1,0, 0,-1, -1,-1};
29     //{0,0, 0,-1};
30     xvel1[OPS_ACC4(0,0)] = xvel1[OPS_ACC4(0,0)] - stepbymass[OPS_ACC2(0,0)] *
31 ( xarea[OPS_ACC5(0,0)] * ( viscosity[OPS_ACC10(0,0)] - viscosity[OPS_ACC10(-1,0)] ) +
32 xarea[OPS_ACC5(0,-1)] * ( viscosity[OPS_ACC10(0,-1)] - viscosity[OPS_ACC10(-1,-1)] ) );
33
34     //{0,0, -1,0, 0,-1, -1,-1};
35     //{0,0, -1,0};
36     yvel1[OPS_ACC8(0,0)] = yvel1[OPS_ACC8(0,0)] - stepbymass[OPS_ACC2(0,0)] *
37 ( yarea[OPS_ACC9(0,0)] * ( viscosity[OPS_ACC10(0,0)] - viscosity[OPS_ACC10(0,-1)] ) +
38 yarea[OPS_ACC9(-1,0)] * ( viscosity[OPS_ACC10(-1,0)] - viscosity[OPS_ACC10(-1,-1)] ) );
39 }
```

---

Figure 1: example user kernel

This user kernel is then used in an `ops_par_loop` (Figure 2). The key aspect to note in the user kernel in Figure 1 is the use of the macros `OPS_ACC0`, `OPS_ACC1`, `OPS_ACC2` etc. These specifies the stencil in accessing the elements of the respective data arrays. For example in `OPS_ACC2`, 2 denotes the third function argument (argument number starts from 0) and (0, 0) denotes stencil. At compile time these macros will be expanded to give the correct array index

---

```

1  int rangexy_inner_plus1[] = {x_min,x_max+1,y_min,y_max+1};
2
3  ops_par_loop(accelerate_kernel, "accelerate_kernel", clover_grid, 2, rangexy_inner_plus1,
4  ops_arg_dat(density0, 1, S2D_00_M10_OM1_M1M1, "double", OPS_READ),
5  ops_arg_dat(volume, 1, S2D_00_M10_OM1_M1M1, "double", OPS_READ),
6  ops_arg_dat(work_array1, 1, S2D_00, "double", OPS_WRITE),
7  ops_arg_dat(xvel0, 1, S2D_00, "double", OPS_READ),
8  ops_arg_dat(xvel1, 1, S2D_00, "double", OPS_INC),
9  ops_arg_dat(xarea, 1, S2D_00_OM1, "double", OPS_READ),
10 ops_arg_dat(pressure, 1, S2D_00_M10_OM1_M1M1, "double", OPS_READ),
11 ops_arg_dat(yvel0, 1, S2D_00, "double", OPS_READ),
12 ops_arg_dat(yvel1, 1, S2D_00, "double", OPS_INC),
13 ops_arg_dat(yarea, 1, S2D_00_M10, "double", OPS_READ),
14 ops_arg_dat(viscosity, 1, S2D_00_M10_OM1_M1M1, "double", OPS_READ));

```

---

Figure 2: example `ops_par_loop`

## References

- [1] OP2 for Many-Core Platforms, 2013. <http://www.oerc.ox.ac.uk/projects/op2>