

# **PhastaIO Specification version 2011-11-10**

PetaApps Group  
Rensselaer Polytechnic Institute  
Troy, New York 12180

## **Overview**

PhataIO is a library in the Phasta CFD system for performance reading and writing of application data in “blocked” parallel approach. This document is meant to serve both as a user’s guide and engineering specification document for PhastaIO. We begin with a PhastaIO example followed by a Man-page like specification for each PhastaIO routine.

# 1 File Format

## 1.1 New file structure

The following parameters are given before writing:

1. # fields;
2. # parts in each file; (or total number of parts)
3. # files;

Note: Global part ID (GPID) associated with each part would change with each data block written and come with field tag. (But we won't have it before writing)

We will discuss later on how to pass these parameters, the options include:

1. Use the interface of current version of phastaIO.
2. Set up probably two new function to do this: one to create a string to field number mapping in writing, and the other to access this mapping information when reading.
3. The file consists of mainly two parts: master header and body. In master header, we store information like, # fields and their names, # parts in this file, # files. The most important one is the offset table which is used for looking up the field data in this file. The basic idea is illustrated by the following simple example and corresponding figure.

In the simple example below, we have 2 processors, each processor deals with 2 parts and each part contains 3 fields (e.g., solution, time derivative and error). These two processors write to one file in parallel. In multiple files case, each file is mostly the same as this example except that different processor will have different rank name. **Note: use global part ID GPID instead of rank when writing/reading parts..**

## 1.2 Pseudo-code

```
2 processors: proc 0 and proc 1;
2 parts:      proc 0 deals with part 0 and part 1; Proc 1 deals with part 2 and part 3;
3 fields:      field 0, field 1 and field 2.
```

```
Call open file function: (create one file)
Setup(#fields, #files, # parts/file)
Loop over fields: (0-2, solution, time derivative solution and error)
    Loop over parts: (proc 0: 0-1, proc 1: 2-3.)
        Compute GPID for this part data(it's changing);
        Call write header function;
        Call write datablock function;
    End loop
End loop
Close file
```

## 1.3 Master Header and Body Structure: 1 File Case

Below is the structure for the Master Header and Body in Figure 1 for the 2 processor, 2 parts per processor with 3 fields example from above assuming only 1 file for both processors.

The formula for finding the index into the offset table is computed as follows:

$$offset\_index = part_i + (field_j \times nppf) \quad (1)$$

## Master Header:

MPI\_IO\_Tag

nFields = 3

solution: 0

time derivative: 1

error: 2

nParts per file = 4

nFiles = 1

Address table:

(0, 0)

(1, 0)

(2, 0)

(3, 0)

(0, 1)

(1, 1)

.

.

.

(3, 2)

## File Body

Header:	part 0, field 0
Data:	part 0, field 0
Header:	part 1, field 0
Data:	part 1, field 0
Header:	part 2, field 0
Data:	part 2, field 0
Header:	part 3, field 0
Data:	part 3, field 0
Header:	part 0, field 1
Data:	part 0, field 1
Header:	part 1, field 1
Data:	part 1, field 1
Header:	part 2, field 1
Data:	part 2, field 1
Header:	part 3, field 1
Data:	part 3, field 1
Header:	part 0, field 2
Data:	part 0, field 2
Header:	part 1, field 2
Data:	part 1, field 2
Header:	part 2, field 2
Data:	part 2, field 2
Header:	part 3, field 2
Data:	part 3, field 2
Proc 0	
Proc 1	

Figure 1: Master Header and Body structure for 1 File case.

As an example, let's consider what entry (0, 0) means. The first index is the global part ID,  $part_i$ . The second is the field

index,  $field_j$ . Thus,  $(0, 0)$  returns the offset index into the file for  $part_0, field_0$ .  $(2, 1)$  returns the offset index into the file for  $part_2, field_1$ .

Last, we plan to use @ as string separator for the field tag. That is a field tag of “solution@2” means solution field on part 2 would mean  $part_i = 2, field_j = 0$  in the offset table of the example file, again assuming only 1 file for both processors.

**Note:** we decided to increase the master header size from previously defined 2MB to 4MB. This is because 2MB was too small to keep the offset data table of some cases when file number is small (i.e. file size gets bigger and offset table gets bigger too). This will not fit all cases with very small file number (e.g. 3B elements mesh with 1 file) but it will be sufficient for typical cases and we can further increase the size when we decide it is needed.

## 1.4 A Multi-file example

Now for a multi-file example. Consider 3 fields, 2 files, 5 parts per processor (ppp) and 4 processors (meaning a total of 20 parts). This configuration could be read by the same configuration or by 20 processors (1 ppp, 10 processors per file), 10 processors (2 ppp, 5 processors per file), 4 processors (5 ppp, 2 processors per file). Please note, that a 5 processor configuration is not possible since while it would yield 4 parts-per-processors, we only have 2 files and so we would have an unequal distribution of parts per file which violates one of our core assumptions. For short hand we use the following:

- number of parts per file  $\rightarrow$  nppf
- number of parts per processor  $\rightarrow$  nppp

Note, that even though we have shown the address table as a two dimensional array we will actually keep it as a 1 dimensional array to avoid any language-based confusion (e.g., C will layout a such a table in row major ordering where as Fortran is column major by default). The one dimensional index is easily computed from the two keys (field number and part number).

file 1 restart.stepnum.0	file 2 restart.<stepnum>.1
MPI_IO_TAG for File 0	MPI_IO_TAG for File 1
nfields : 3	nfields : 3
solution:	solution:
time derivative:	time derivative:
error:	error:
nParts per file (nppf): 10	nppf : 10
nFiles: 2	nFiles : 2
(0,0)	(10,0)
(1,0)	(11,0)
(2,0)	(12,0)
(3,0)	(13,0)
(4,0)	(14,0)
(5,0)	(15,0)
(6,0)	(16,0)
(7,0)	(17,0)
(8,0)	(18,0)
(9,0)	(19,0)
(0,1)	(10,1)
.	.
.	.
.	.
(9,1)	(19,1)
(0,2)	(10,2)
.	.
.	.
.	.
(9,2)	(19,2)
HEADER(0,0)	HEADER(10,0)
DATA(0,0)	DATA(10,0)

.  
. ORDER NOT IMPORTANT, NOR IS HAVING 0 SIZE DATA BLOCKS  
.

We know the number of fields being written we can pre-allocate the table and access it in any  $(part_i, field_j)$  order.

## 1.5 Application Usage, Implementation and Caveat Details

ALL THE COMMENTS BELOW ARE ON A PER FILE BASIS. THERE IS A STRONG LINK BETWEEN INIT and OPEN. ONE TO ONE AND ONTO. ACTUALLY THE ONLY REASON THEY ARE TWO SEPARATE FUNCTIONS IS TO KEEP BACKWARD COMPATIBILITY.

We add a `FinalizePHMPIO` where local comms are released and file header structure is committed. That is, a file or group of files could be closed and then re-opened and effectively append to the file group. However, once the `FinalizePHMPIO` is called, those files are now committed and we loose our ability to append to those files.

`InitPHMPIO` has a global handle as one of its arguments, to indicate which file it is working on. This handle points to the struct that is holding all the data associated with the file that we are working with. `FinalizePHMPIO` takes this handle as an argument and frees memory of the struct associated with this file. This will allow multiple files to be dealt with at the same time. (The `MAX_PHASTA_FILES` macro is used to indicate max number of files we can work on at the same time)

Other implications are as follows:

1. Application must compute total number of fields to be written before `InitPHMPIO` is called.
2. `InitPHMPIO`: The arguments are: `nfields`, `nppf`, `nfiles`, and returns the handle for the struct. Since we know these parameters when Phasta is started, they can be directly passed into the PhastaI/O library via this function.
3. `openfile_`: We can keep current interface where the filename can still be computed before the file open function is called. See the discussion in the next subsection.
4. Application must concatenate the GPID to the field phrase before each field write with special separator “@”.
5. Write header and data as previously described.

Another key question here is: **Upon file opening, how does a given processor know which file to open especially for the multi-file case?** The answer is the file name will contain the proper information that will enable a processor to compute which file it should open. In particular the file name format will be: `< fileTypeString > . < stepnum > . < fnum >` where `fnum = int(mpi_rank/nfiles)`. **That is, we will assume/require that parts are ordered and sequentially/equally distributed in ranges across ranks.**

Some notes/caveats on using as well as implementing the new PhastaIO lib that address this question are as follows:

1. On file open/create, library will create local mpicomm for parts writing to a given file. This will be set for all fields and for each “set” of parts in the multiple ppp case since we have agreed to not vary ppp across processors. Thus NO CHANGE except sequential GPID’s in a given file rather than the MOD striping currently used in MPIO library.
2. Two functions to parse phrase into string, field# and GPID must be created. One for writing that “builds” the fields table and one for reading that uses the table as a map.
3. A C structure must be added to:
  - stored field string (only if new: e.g., in multiple parts-per-processor case same string will be encountered multiple times).
  - store all starting file addresses in a 64 bit integer (i.e., `long long`) in a way that the field-part locator (e.g., offset index) can be written after all parts/fields are written (when closing).
  - this field-part locator requires the order to appear sequential in the GPID index.
4. Space reserved in the header for the above field-part long-long addresses is exactly computable.
5. Currently we are planning header-then-data (sequentially one after the other). This can easily be relaxed later by storing all of the headers in a struct and, sorting them to match the table, and then writing them as a single block.

6. Note that the first GPID for a given file is  $GPID_{start} = fnum * nppf$  (recall  $fnum = int(mpi\_rank / nfiles)$ ). This can be computed for each file and correct the GPID to starting address as follows  $offset_{index} = part_i + nppf * field_j$  but now  $part_i = GPID_i - GPID_{start}$ . Thus, the  $part_i$  is not purely a global part ID anymore but is a local position relative to the file on which this processor is operating on.

The Application steps for reading a file are as follows:

1. Application must know (or use system commands to find) the number of files.
2. It will know its number of processors so, as above, each processor can determine which files hold its part's data.
3. It will not typically know the number of fields but this will be set after the file is open where it will be read
4. When application wants to read a field it will have to open the appropriate file in read mode.
5. At open time, the available fields will be loaded into the field table struct so that later strings can be mapped to field#.
6. Again the only change will be that the GPID of the part desired will have to be concatenated (by application) onto the phrase (with ) for the field requested
7. A different function inside of header will then parse this phrase into the field phrase and gpids. The field phrase will be further mapped to a field number using the struct described previously.

## 2 queryphmpiio\_

### 2.1 Synopsis

```
#include<phastaIO.h>

void
queryphmpiio_( const char filename[],
               int *nfields,
               int *nppf );
```

### 2.2 Description

This function takes filename (Michel: should we add file descriptor number as well to distinguish different files opened at same time?), open it and read the header to validate endianness and format. If endianness is correct and format is MPI new format, it will read number of fields (nfields) and number of parts per file (nppf) and pass it back to PHASTA (using a pointer).

## 3 initphmpiio\_

### 3.1 Synopsis

```
#include<phastaIO.h>

int
initphmpiio_( int *nfields,
              int *nppf,
              int *nfiles,
              int *filehandle,
              char mode[] );
```

## 3.2 Description

This function takes the file handle and other numbers needed by header (nfields, nppf, nfiles), allocate and initialize the struct properly for this file and increases the global file descriptor index by 1. The reason we need mode in the argument is so that we need to differentiate “read” and “write”. For “read” case, we presume query function was called already and MasterHeaderSize is already known; for “write” case, we need to compute a proper MasterHeaderSize we need according to our dataset size.

## 4 finalizephmpiio\_

### 4.1 Synopsis

```
#include<phastaIO.h>
void
finalizephmpiio_( int *fileDescriptor );
```

### 4.2 Description

This function takes the file handle as argument and free the struct associated with the files pointed by file handle.

## 5 openfile\_

### 5.1 Synopsis

```
#include<phastaIO.h>

void
openfile_( const char filename[],
           const char mode[],
           int*   fileDescriptor ) {
```

### 5.2 Description

This is the file open function. The first argument `filename[]` is a string that contains path of the file we want to open; second argument is a constant string that specifies the `mode` we want to use, i.e. “read”, “write” and “append”; third argument is an integer value that contains the index of the file descriptor array(named “fileArray”).

This function would push a file handle into “fileArray” and give “fileDescriptor” a proper index value.

## 6 readheader\_

### 6.1 Synopsis

```
#include<phastaIO.h>

void
readheader_( int* fileDescriptor,
            const char keyphrase[],
            void* valueArray,
            int* nItems,
            const char datatype[],
            const char iotype[] )
```

## 6.2 Description

This is the header reading function. The first argument is an `index` of file descriptor array (named “fileArray”), which specifies the right file handle; the second argument is a constant string that specifies the `field type`, e.g. “solution”; the third argument is an array that contains `types of data`(only integers so far), e.g. “nshg” and “nv”; the fourth argument is number of types of data for “valueArray”; the fifth argument is `datatype` in datablock, i.e. integer or double; the sixth argument is `iotype`, i.e. binary or non-binary.

This function reads file header and gets information like field type, type of data in datablock, `datatype`(double, integer, etc), and so on.

## 7 readatablock\_

```
#include<phastaIO.h>

void
readatablock_( int*   fileDescriptor,
               const char keyphrase[],
               void* valueArray,
               int*   nItems,
               const char datatype[],
               const char iotype[] ) {
```

### 7.1 Description

This `readatablock` function has the same arguments as the `readheader` function. The first, second, fifth, sixth arguments specify exactly the same things as `readheader` function. The third argument is an array that stores the data read from datablock; the fourth argument specifies the size of the data array.

This function checks field type, data type, `iotype` first, then uses `fread`(for binary data) or `fscanf`(for integer and double) to get data into data array.

## 8 writeheader\_

```
#include<phastaIO.h>

void
writeheader_ ( int*   fileDescriptor,
               const char keyphrase[],
               void* valueArray,
               int*   nItems,
               int*   ndataItems,
               const char datatype[],
               const char iotype[] ) {
```

### 8.1 Description

This `writeheader` function has very similar arguments as `readheader` function, except that it has an extra argument “`ndataItems`”, which basically means number of data items in the datablock. This information is written in header right after field type.

This function gathers information like field type, types of data(in the header such as “nshg” and “nv”), total size of data in its data block, `datatype`, `iotype`, etc, and writes these information into file header in a certain sequence.

## 9 writedatablock\_

```
#include<phastaIO.h>
```



```

void
writedatablock_( int* fileDescriptor,
                 const char keyphrase[],
                 void* valueArray,
                 int* nItems,
                 const char datatype[],
                 const char iotype[] ) {

```

## 9.1 Description

This write datablock function has the same arguments as readdatablock function: the first argument is an index of file descriptor array (named “fileArray”), which specifies the right file handle; the second argument is a constant string that specifies the field type, e.g. “solution”; the third argument is an array that contains types of data (only integers so far), e.g. “nshg” and “nv”; the fourth argument is number of types of data for “valueArray”; the fifth argument is datatype in datablock, i.e. integer or double; the sixth argument is iotype, i.e. binary or non-binary.

This function checks field type, iotype, data type, and total size of data block, and then writes data into file with fwrite (for binary), or fprintf (for integer and double).

## 10 closefile\_

```

#include<phastaIO.h>

void
closefile_( int* fileDescriptor,
            const char mode[] ) {

```

### 10.1 Description

This is the file close function. The first argument is the index of file descriptor array as to which file is being closed; second argument, mode, is a constant string which specifies the mode, if mode is “write” or “append”, a file flush function would be called before file is finally closed.

## 11 Limitations with using PhastaIO

What are the types of things you cannot do with PhastaIO?

## 12 Motivating Example

Below is an example of how PhastaIO is used.

```

#include <stdio.h>
#include <iostream>
#include "phastaIO.h"
#include <string.h>
#include "rdtsc.h"
#include <stdlib.h>
#include "mpi.h"

using std::cout;
using std::endl;

int main(int argc, char *argv[]){

```

```

//These are used for timing the code:.....
unsigned long long timer_[6];
unsigned long long timer_span[5];
float time_span[5];

int number_of_groups = atoi(argv[3]);
MPI_Init(&argc,&argv);

int myrank, numprocs;
MPI_Comm_rank(MPI_COMM_WORLD, &myrank);
MPI_Comm_size(MPI_COMM_WORLD, &numprocs);

//split the mpi com_world
MPI_Comm local_comm;
//int number_of_groups=M;
int color=myrank % number_of_groups;
int new_size, newrank;
MPI_Comm_split(MPI_COMM_WORLD, color, myrank, &local_comm);
MPI_Comm_size(local_comm, &new_size);
MPI_Comm_rank(local_comm, &newrank);

int iterations = atoi(argv[1]);
int lstep = 0;
int numparts = atoi(argv[2]); // total number of parts

int i,irstin,irstout,*nshg,*nv,*sn,iarray[3],isize,nitems;
int ithree=3;
double **field_data;
int nchar = 255;
char inrfname[255], outrfname[255], fprefix[255], fieldtag[255], hdata[255];
char iotype[16];
strcpy(iotype,"binary");
int magic_number = 362436;
int* mptr = &magic_number;

// check range of parts that needs to be processed by this rank
int numparts_l = numparts/numprocs;

// note myrank starts from 0 (lowest one is 0)
// startpart will begin from 1 on lowest rank
// (also order of operations is imp. as numparts_l will be changed on last rank/proc. as below)
int startpart = myrank*numparts_l+1;

// last rank/proc. may get extra partial number of parts
if(myrank+1==numprocs)
    numparts_l = numparts-(numprocs-1)*numparts_l;

// startpart begins from 1 on rank 0
// (order of operations is imp. as numparts_l is corrected on last rank/proc.)
int endpart_plus1 = startpart+numparts_l;
if(numparts<numprocs || numparts_l==0) {
    if(myrank<numparts) {
        startpart = myrank+1;
        endpart_plus1 = startpart+1;
    }
    else { // some ranks will be idle
        startpart = 0;
        endpart_plus1 = 0;
    }
}

```

```

}

//printf("[%d] parts range : %d - %d\n",myrank,startpart,endpart_plus1-1);

// easier to avoid offset and allocate it as numparts
nshg = new int[numparts];
nv = new int[numparts];
sn = new int[numparts];
field_data = new double*[numparts];

sprintf(fprefix, "/gpfs/large/PHASTA/home/PHASzhou/Runs/AAAPatient2Exercise/M3/%d-procs_case", numparts);

strcpy(fieldtag, "solution");

// just to separate initialize and read/write phases
MPI_Barrier(MPI_COMM_WORLD);

if(myrank==0)
    cout << "Reading phase..." << endl;

// Read phase (field used currently is "solution")
for(i=startpart; i<endpart_plus1; i++){
    bzero((void*)inrfname, nchar);
    sprintf(inrfname, "%s/restart.%d.%d", fprefix, lstep, i);
    // cout << " [" << myrank << "] Opening " << inrfname << " (for read)..." << endl;
    openfile_(inrfname, "read", &irstin);

    /* read field data */
    iarray[0]=-1;
    readheader_(&irstin, fieldtag, (void*)iarray, &ithree, "double", iotype);
    nshg[i]=iarray[0];
    nv[i]=iarray[1];
    sn[i]=iarray[2];
    if(iarray[0]==-1){
        //cout << " [" << myrank << "] " << fieldtag << " not found... exiting" << endl;
        return 1;
    }
    isize=nshg[i]*nv[i];
    /* allocate the array to the right size */
    field_data[i]=new double[isize];

    // cout << " [" << myrank << "] Reading " << fieldtag << "..." << endl;
    /* now read the array */
    readdatablock_(&irstin, fieldtag, (void*)field_data[i], &isize, "double", iotype);
    // cout << " [" << myrank << "] Read " << fieldtag << " for " << nshg << " modes..." << endl;

    closefile_(&irstin, "read");
}

/* PROCESS READ DATA HERE */

// just to separate initialize and read/write phases
MPI_Barrier(MPI_COMM_WORLD);

double test_data = 28.922292;
field_data[startpart][5] = test_data;
double *new_field_data=new double[isize];

if(myrank==0)

```

```

    cout << "Writing phase..." << endl;

int descriptor;
char filename[255], path[255];
bzero((void*)filename,255);
sprintf(path, "/gpfs/large/PHASTA/home/PHASliuning/restart/%d", myrank%number_of_groups);
sprintf(filename,"%s/mytestdata.%d.%d",path,82,0);
openfile_new(filename, "read", &descriptor);

//New MPI IO stuff: .....
MPI_Barrier(MPI_COMM_WORLD);
timer_[0]=rdtsc();
writeheader_new( &descriptor, fieldtag, (void*)iarray, &ithree, &isize, "double", iotype);
MPI_Barrier(MPI_COMM_WORLD);
timer_[1]=rdtsc();
writedatablock_new( &descriptor, fieldtag, (void*)(field_data[startpart]), &isize, "double", iotype );
MPI_Barrier(MPI_COMM_WORLD);
timer_[2]=rdtsc();

timer_span[0] = timer_[1]-timer_[0];
timer_span[1] = timer_[2]-timer_[1];

/*
 * NOTE HARD CODE FOR 700 MHz BG/L -- BG/P is 850 MHz
 */
time_span[0]=(float)timer_span[0]/700.0;
time_span[1]=(float)timer_span[1]/700.0;

//Reading and checking part: .....
readdatablock_new(&descriptor,fieldtag, (void*)new_field_data,&isize,"double",iotype);
//readheader_new( &descriptor,fieldtag, (void*)iarray,&ithree,"double",iotype);
closefile_new(&descriptor, "read");

if (myrank==0){
    printf(" write header time is : %f \n", time_span[0]/1000.0);
    printf(" write datablock time is : %f \n", time_span[1]/1000.0);
}

if (myrank==6)
    printf("field_data is: %lf and new one is %lf\n", field_data[startpart][5], new_field_data[5]);

// just to separate finalize
MPI_Barrier(MPI_COMM_WORLD);

for(i=startpart;i<endpart_plus1;i++){
    delete [] field_data[i];
}
delete [] field_data;
delete [] sn;
delete [] nv;
delete [] nshg;

if(myrank==0)
    cout << "Finalize phase..." << endl;

MPI_Finalize();

return 0;
}

```

**Need a detail description of what is going on this example.**