



European  
Commission

Horizon 2020  
European Union funding  
for Research & Innovation

## WP4: I/O & Data Flow

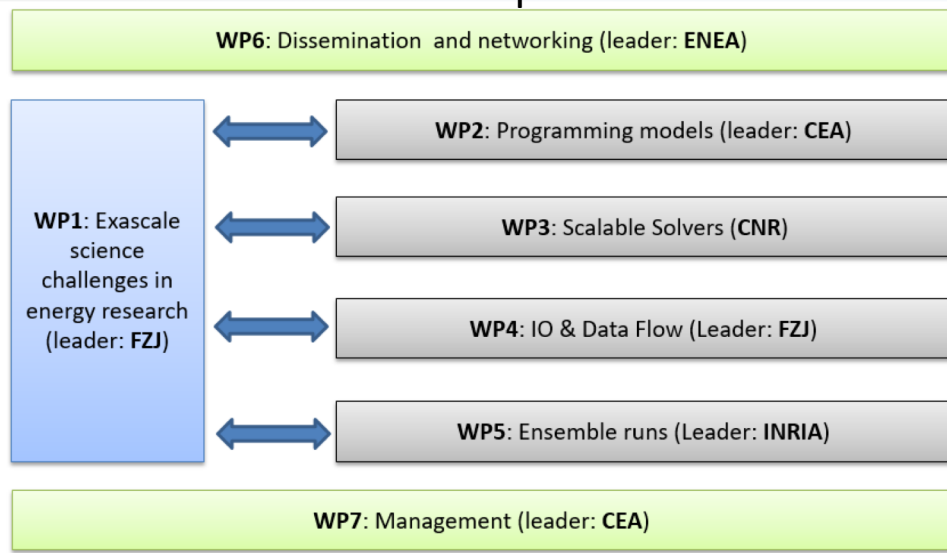
F.Iannone, F.Ambrosino, A.Funel, ENEA

Report Project – 12/02/2020

# EoCoE<sup>2</sup>: Energy Oriented Center of Excellence : toward exascale for energy

- ✓ HPC apps: Wind, Meteo, Materials, Water, Fusion
- ✓ Start date: 2019 – End Date: 2021

## Work plan



## PM efforts

	Partner number	WP1	WP2	WP3	WP4	WP5	WP6	WP7	Total Person-Months per Participant
CEA	1	56	70	2	24	30	5	30	217
FZJ	2	75	36	6	39	17	11	19	203
ENEA	3	30	0	0	7	0	25	1	63
BSC	4	23	42	12	27	0	1	1	106
CNRS	5	16	0	27	24	0	0	1	68
INRIA	6	14	14	18	2	40	1	1	90
CERFACS	7	12	0	33	0	0	1	1	47
MPG	8	39	15	2	0	0	1	1	58
FRAUNHOFER	9	24	0	0	0	0	1	1	26
FAU	10	0	30	0	0	0	1	1	32
CNR	11	0	0	48	0	0	0	0	48
UNITN	12	28	0	0	0	0	1	1	30
PSNC	13	0	0	0	48	18	24	0	90
ULB	14	0	0	33	0	0	1	1	35
UBAH	15	30	0	3	0	0	1	1	35
CIEMAT	16	7	0	0	0	0	2	1	10
IFPEN	17	3	18	0	0	0	1	1	23
DDN	18	0	0	0	24	0	1	1	26
<b>TOTAL Person Months</b>		<b>357</b>	<b>225</b>	<b>184</b>	<b>195</b>	<b>105</b>	<b>78</b>	<b>63</b>	<b>1 207</b>

- ✓ Project funds: 1207 PMs – 9.2 MEuro
- ✓ ENEA efforts : 63 PMs – 560 kEuro
- ✓ ENEA in T4.2: I/O Refactoring & optimization – T4.2.2: Gysela I/O optimization – 7 PMs

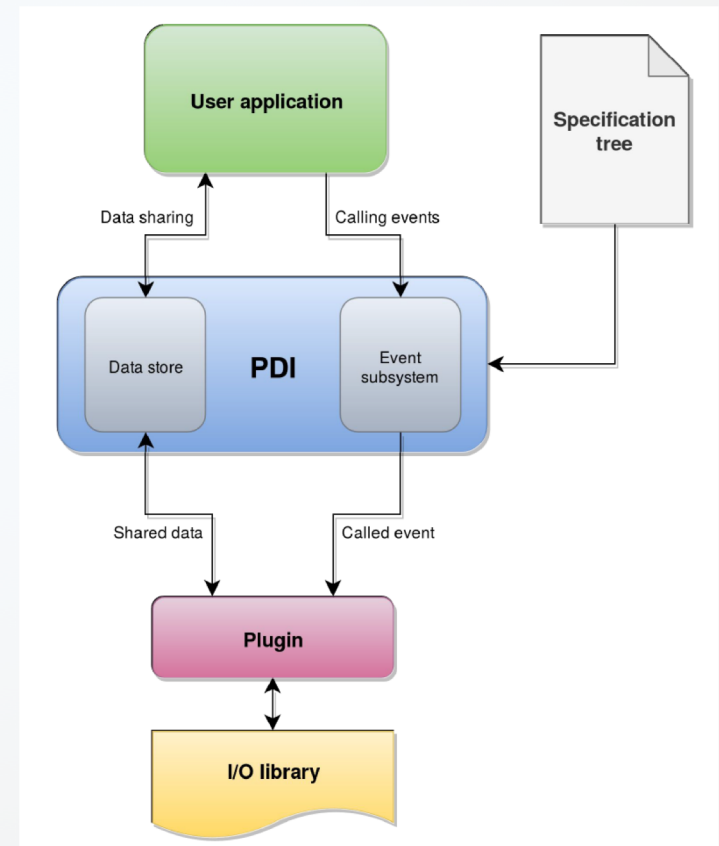
# EoCoE<sup>2</sup>: *Energy Oriented Center of Excellence : toward exascale for energy*

## ✓ WP 4 : I/O & Data Flow

- Task 4.2: I/O refactoring and optimization
  - Task 4.2.2: Gysela I/O optimization (Fusion) (Gysela-X not ready yet)

## ✓ PDI: Parallel Data Interface

- Middleware to interface I/O Library (Posix,HDF5,SionLib)
- C++/MPI comm
- Parser YAML allows to define data structure
- API:
  - PDI\_init
  - PDI\_share
  - PDI\_reclaim
  - PDI\_release
  - PDI\_expose
  - PDI\_access
  - PDI\_event
  - PDI\_multiexpose
  - PDI\_finalize

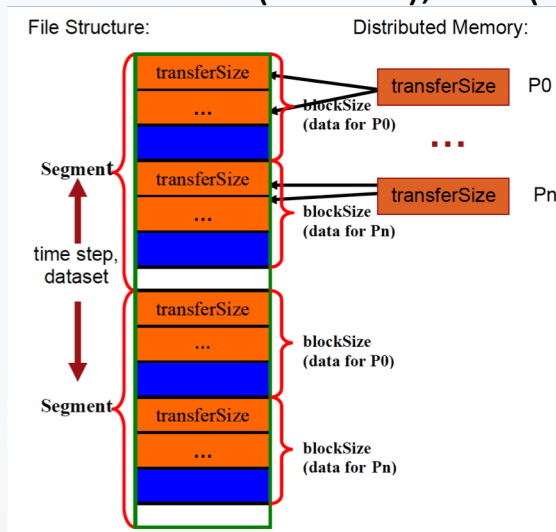


# EoCoE<sup>2</sup>: *Energy Oriented Center of Excellence : toward exascale for energy*

## ✓ ENEA activities:

- integrate PDI into the standard I/O benchmarking tool IOR waiting for GyselaX deployment
- test performances of PDI integrated in IOR with POSIX interface
- make an extensive benchmarks session using JUBE to automatically test, many parameters (e.g. file system block size, shared/independent access, collective I/O, cache effects etc.)
- find and report any bugs. Comparison of runs with and without PDI to give an estimation of the overhead especially for large I/O amount of data and/or large number of tasks

✓ Currently: installed standalone PDI (v. 0.5.1), IOR (v. 3.3.0) JUBE (v. 2.1.4).



IOR file structure and Processors

- ✓ API development in IOR of a PDI interface (POSIX plugin) in IOR: aiori-PDI.c
  - PDI\_create, PDI\_Mknode, PDI\_Open, **PDI\_Xfer**, PDI\_Close, PDI\_Fsync, PDI\_GetFileSize

# EoCoE<sup>2</sup>: Energy Oriented Center of Excellence : toward exascale for energy

## ✓ Plugin: xfer\_pdi to interface PDI to POSIX functions read()/write()

```
int main(int argc,char *aa[])
{
  int bs,ts,sc;
  int fd;
  int n_block;
  char *buffer;
  if(argc!=5){
    printf("Usage : filename blocksize[MB] transfersize[MB] segmentcount \n");
    return -1;
  }
  fd=open(aa[1],O_CREAT|O_RDWR);
  if(fd==-1){
    printf("file create error.\n");
    return -1;
  }
  bs=atoi(aa[2]);
  ts=atoi(aa[3]);
  sc=atoi(aa[4]);
  printf("bs: %d\n",bs);
  n_block = bs/ts;
  buffer = malloc(sizeof(char)*ts*1024*1024);
  memset(buffer,1,ts*1024*1024);
  for (int i = 0; i < sc; i++) {
    for (int j = 0; j < n_block; j++) {
      printf("block: %d %d\n",i,j);
      write(fd,buffer,ts*1024*1024);
    }
  }
  close(fd);
}
```

### IOR POSIX benchmark

```
YAML file
metadata:
  ts: int64
  access: int
  file: int
  offset: int64
  return: int64
data:
  buffer: {type: array, subtype: char, size: $ts}
plugins:
  user_code:
    on_data:
      buffer:
        xfer_pdi:
          value: $buffer
```

```
133 void xfer_pdi()
134 {
135   // arguments of the function
136   int* access; PDI_access("access", (void*)&access, PDI_IN);
137   long* length; PDI_access("ts", (void*)&length, PDI_IN);
138   long* offset; PDI_access("offset", (void*)&offset, PDI_IN);
139   int* fd; PDI_access("file", (void*)&fd, PDI_IN);
140   long l=length;
141   char* pptr;
142   if (*access == READ) { //read
143     PDI_access("value", (void*)&pptr, PDI_OUT);
144   } else { // write
145     PDI_access("value", (void*)&pptr, PDI_IN);
146   }
147   long long rc;
148   long long remaining = (long long)*length;
149   // function body
150   if (*access == READ) {
151     if (verbose >= VERBOSE_4) {
152       fprintf(stdout,
153               "task %d reading from offset %lld\n",
154               rank,
155               *offset + length - remaining);
156     }
157     rc = read(*fd, pptr, *length);
158   } else {
159     if (verbose >= VERBOSE_4) {
160       fprintf(stdout,
161               "task %d writing to offset %lld\n",
162               rank,
163               *offset + length - remaining);
164     }
165     rc = write(*fd, pptr, *length);
166   }
167 }
168
169 // release all metadata
170 PDI_release("access");
171 PDI_release("file");
172 PDI_release("ts");
173 PDI_release("offset");
174 PDI_release("value");
175 // update return value
176 PDI_expose("return", &rc, PDI_OUT);
177
178 }
```

# EoCoE<sup>2</sup>: Energy Oriented Center of Excellence : toward exascale for energy

✓ Collaborative work: SLACK – PDI [https://app.slack.com/client/T9G3CB93N/DL5P6H8KY/user\\_profile/UL3537YHJ](https://app.slack.com/client/T9G3CB93N/DL5P6H8KY/user_profile/UL3537YHJ)

Slack needs your permission to enable desktop notifications.

**Karol**  
☆ | away | Karol Sierociński

November 12th, 2019

Tip: Try **X** **F** to search this channel.

now I have to change my goals.  
In order to integrate PDI in IOR I have to consider as benchmark the HDF5 interface not POSIX

**Karol** 11:53 AM  
I don't understand

**Francesco** 11:55 AM  
The initial goal was to compare IOR performances with POSIX interface versus IOR performances with PDI/posix plugin  
Now the new goal is to compare IOR performances with HDF5 interface versus IOR performances with PDI/DECL\_HDF5 plugin

**Karol** 11:57 AM  
ok, so what is the benchmark?

**Francesco** 11:59 AM  
IOR benchmark is the I/O throughput to a parallel filesystem  
What IOR do is to write/read buffer of bytes in the following format:  
12.01 each parallel task (p#) write a block size in one file  
image.png

**File Structure:**  
transferSize  
...  
transferSize  
...  
transferSize  
...  
transferSize  
...  
transferSize

**Distributed Memory:**  
blockSize (data for P0)  
...  
blockSize (data for Pn)

Segment  
time step, dataset  
Segment

Fig. 1. The design of the IOR benchmark for shared file type. Blocks are stored in separate files for the 1-file-per-processor mode of operation.

**Karol** 12:02 PM  
You can't specify this in the decl\_hdf5 plugin

**Francesco** 12:02 PM  
Message Karol

**Workspace Directory**

**francesco iannone** •  
ENEA

Status  
Set a status

Display name  
Francesco

Local time  
3:10 PM

Phone number  
(39)0694005124

# EoCoE<sup>2</sup>: Energy Oriented Center of Excellence : toward exascale for energy

## ✓ TEST ENVIRONMENT: Jube

### IORPDI.xml

```
<parameter name="api">POSIX,PDI</parameter>
<parameter name="blockSize">${transferSize}</parameter>
<parameter name="transferSize" type="int" mode="python">" , ".join(str(i) for i in [128*1024,1024**2,4*(1024**2),256*(1024**2)])</parameter>
<parameter name="segmentCount" type="int" mode="python">(2*(256*1024**2))/${transferSize}</parameter>
<parameter name="repetitions" type="int">3</parameter>
<parameter name="filePerProc" type="int">1</parameter>
<parameter name="keepFile" type="int">1</parameter>
<parameter name="testFile">${workdir}/testFile</parameter>
<parameter name="workdir">/gpqr3/eocoe/work/${jube_wp_id}</parameter> <!--originariamente era indicata come testdir-->
<parameter name="nodes" type="int">1,2,4,6,8</parameter>
<parameter name="taskspernode" type="int">1,2,4,8,16,32</parameter>
<parameter name="tasks" type="int" mode="python">${nodes}*${taskspernode}</parameter>
```

```
(start) 2019-12-06 16:29:35 (1575646175)
IOR-3.3-0-der: MPI Coordinated Test of Parallel I/O
Began      : Fri Dec 6 16:29:35 2019
Command line : /gpqr3/eocoe/1/ior-test/eocoe2/ior/build_openmpi_acc730/bin/lor -iqr_input_dir=/gpqr3/eocoe/work/24-k
Machine     : Linux cpe006k293.opfcdi.enea.it:3.10-514.26.2.6f7.x86_64 #1 SMP Tue Jul 4 15:04:05 UTC 2017 x86_64
Using synchronized MPI timer
Start time skew across all tasks: 0.00 sec
TestID     : 0
StartTime  : Fri Dec 6 16:29:35 2019
Path       : /gpqr3/eocoe/work/24
FS         : 726.0 TiB Used FS: 33.9% Inodes: 1918.8 Mi Used Inodes: 0.3%
Participating tasks : 1
Using reorder/Tasks -C (expecting block, not cyclic, task assignment)

Options:
api      : PDI
apiVersion :
test filename : /gpqr3/eocoe/work/24/testFile
access   : file-per-process
segments : independent
ordering in a file : sequential
task offset : 4096
clients per node : 1
repetitions : 10
xferSize  : 131072 bytes
blockSize : 131072 bytes
aggregate filesize : 512 MiB

Results:
Using Time Stamp 1575646175 (Ox5dea73df) for Data Signature
access bw(MiB/s) block(KiB) xfer(KiB) opens() wrfd(s) close(s) total(s) iter
[PD] [User-code] [16:29:35] *** info: Plugin loaded successfully
[PD] [16:29:35] *** info: Initialization successful
Commencing write performance test: Fri Dec 6 16:29:35 2019
[PD] [16:29:36] *** info: Finalization
[PD] [User-code] [16:29:36] *** info: Closing plugin
write 1478.95 128.00 128.00 0.008377 0.337223 0.000546 0.346191 0

Operation Max(MiB) Min(MiB) Mean(MiB) StdDev Max(OPs) Min(OPs) Mean(OPs) StdDev Mean(s) Test# #Tasks IPN Reps IPP (reord reordoff reordrand seed segcnt blksize xsize aggs(MB) API RetNum
write 1963.58 1478.95 1849.75 128.16 15708.64 11831.60 14788.00 1025.24 0.27839 0 1 1 10 1 1 1 0 0 4096 131072 131072 512.0 PDI 9999
read 2711.39 2670.33 2689.65 14.15 21691.11 21362.63 21517.21 113.23 0.19036 0 1 1 10 1 1 1 0 0 4096 131072 131072 512.0 PDI 9999

Summary of all tests:
Operation Max(MiB) Min(MiB) Mean(MiB) StdDev Max(OPs) Min(OPs) Mean(OPs) StdDev Mean(s) Test# #Tasks IPN Reps IPP (reord reordoff reordrand seed segcnt blksize xsize aggs(MB) API RetNum
write 1963.58 1478.95 1849.75 128.16 15708.64 11831.60 14788.00 1025.24 0.27839 0 1 1 10 1 1 1 0 0 4096 131072 131072 512.0 PDI 9999
read 2711.39 2670.33 2689.65 14.15 21691.11 21362.63 21517.21 113.23 0.19036 0 1 1 10 1 1 1 0 0 4096 131072 131072 512.0 PDI 9999
Finished : Fri Dec 6 16:29:40 2019
```

nodes	taskspernode	tasks	API	Min(MiB/s)	Max(MiB/s)	Mean(MiB/s)	StdDev	Min(OPs)	Max(OPs)	Mean(OPs)	StdDev	Min(s)	Max(s)	Mean(s)	StdDev	Min(s)	Max(s)	Mean(s)	StdDev	Min(s)	Max(s)	Mean(s)	StdDev	
1	1	1	POSIX	128	4096	2560.00	2.88000	0.00000	0.00000	128	128	2092.0	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	128
1	1	1	PDI	1024	4096	9306.00	2.40000	0.00000	0.00000	1024	1024	1861.79	0.00020	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	1024
1	1	1	POSIX	4096	4096	536.00	2.76000	0.00000	0.00000	4096	4096	1019.99	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	4096
1	1	1	PDI	262144	3667	0.000010	3.40000	0.00000	0.00000	262144	262144	2066.67	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	262144
1	2	2	POSIX	128	10192	1336.00	3.76000	0.00000	0.00000	128	128	4131	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	128
1	2	2	PDI	1024	10089	2368.00	3.90000	0.00000	0.00000	1024	1024	3820	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	1024
1	4	4	POSIX	128	10098	2358.00	4.00000	0.00000	0.00000	128	128	3791	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	128
1	4	4	PDI	1024	10091	2358.00	4.00000	0.00000	0.00000	1024	1024	3754	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	1024
1	8	8	POSIX	128	10091	2358.00	4.00000	0.00000	0.00000	128	128	3754	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	128
1	8	8	PDI	1024	10098	2368.00	4.00000	0.00000	0.00000	1024	1024	3773	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	1024
1	16	16	POSIX	128	10098	2368.00	4.00000	0.00000	0.00000	128	128	3773	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	128
1	16	16	PDI	1024	10098	2368.00	4.00000	0.00000	0.00000	1024	1024	3773	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	1024
1	32	32	POSIX	128	10098	2368.00	4.00000	0.00000	0.00000	128	128	3773	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	128
1	32	32	PDI	1024	10098	2368.00	4.00000	0.00000	0.00000	1024	1024	3773	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	1024



# EoCoE<sup>2</sup>: *Energy Oriented Center of Excellence : toward exascale for energy*



## ✓ **Next actions:**

- New IOR POSIX/PDI benchmarks in a noiseless sessions
- Deliverable of WP4 with ENEA contribute on PDI integration
- Development of Gysela optimization (coordinate by JSC)
- Effort need: 1 developer for Gysela optimization