# The **P**DI **D**ata **I**nterface

# IO, restart, in situ or coupling: PDI, a single interface to decouple data handling from computation in numerical simulation

October 4[th], 2023 – SISMA

Julien Bigot[2], François-Xavier Morvan[2], Pierre-Antoine Raclius[2], Yushan Wang[2]

➢ We want it easy to use

➢ We want it fast

➢ We want a portable library

➢ We want large language support

➢ We want parallelization independent file format

➢ We want a portable file format

➢ We want to leverage the underlying hardware

➢ We want…

PDI

- ➢ We want it easy to use
- ➢ We wan~~t~~
- ➢ We wan~~t~~

> Handling I/O is complex
> Optimizing I/O is a job on its own

- ➢ We want large language support
- ➢ We want parallelization independent file format
- ➢ We want a portable file format
- ➢ We want to leverage the underlying hardware
- ➢ We want…

➢ We want it easy to use

➢ We war...

➢ We war...

➢ We want large language support

➢ W...

➢ We want a portable file format

➢ We want to...

➢ We want…

**Handling I/O is complex**
**Optimizing I/O is a job on its own**

**Complex but common problem,**
**A community with dedicated expert**

**Let's use libraries**

FTI

MPI I/O

pNetCDF

SIONlib

FlowVR

Hercule

Adios

NetCDF

HDF5

XIOS

IME

POSIX

SCR

pHDF5

Damaris

PaDaWAn

Choosing the best library: a problem on its own

The best library depends on…

➢ The code specifics, the type of I/O

  ➢ Parallelism level, replicated / distributed data, I/O frequency, …

  ➢ Initialization data reading, result writing (small or large), checkpoint writing, coupling related I/O

➢ The specific execution

  ➢ Small case / large case, debug / production, …

➢ The specific hardware available

  ➢ I/O bandwidth, intermediate storage, …

## Choosing the best library: a problem on its own
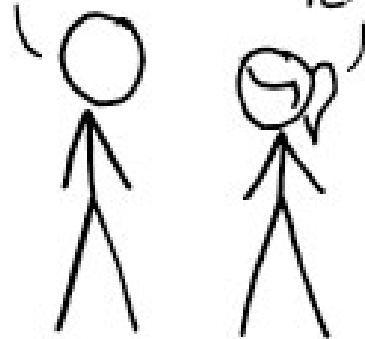
The best library depends on…

➢ The code specifics, the type of I/O

  ➢ Parallelism le[...] [...]equency, …

  ➢ Initialization data reading, result writing (small or large), checkpoint writing, coupling related I/O

➢ The specific execution

**Not one-size fits all library**

➢ The specific hardware available

**Many codes end-up with an IO abstraction layer**

  ➢ I/O bandwidth, intermediate storage, …

© XKCD
https://xkcd.com/927/

© XKCD
https://xkcd.com/927/

No!

PDI is an Interface...
just an interface!

PDI



Main loop

Initialization

Finalization

- Data assimilation
- Coupling inputs,
- …

Main loop

- Intermediate results,
- Checkpoint,
- Post-processing,
- Coupling outputs,

- parameters reading,
- data initialization,
- …

- Final results,
- Final checkpoint,
- …

Initialization

Finalization

Similar from the code point of view:
➢ Import or export data

But… different libraries needed

- parameters reading,
- data initialization,
- …

**Main loop**

- Intermediate results,
- Checkpoint,
- Post-processing,
- Coupling outputs,

- Data assimilation
- Coupling inputs,
- …

- Final results,
- Final checkpoint,
- …

Initialization

Finalization

Similar from the code point of view:

➢ Import or export data

But… different libraries needed

- Data assimilation
- Coupling inputs,
- …

Main loop

- parameters rea
- data initializati
- …

**The data-handling problem**

- Final results,
- Final checkpoint,
- …

Initialization

Finalization

**MPI**

| Code | Code | ... | Code |

**?**

API
API
API
API
API
API
API
API
API

**HPC Library**

➢ PDI annotations: a purely declarative API

➢ Plugins for access to existing libraries

- ➢ PDI annotations: a purely declarative API

- ➢ PDI YAML spec. tree: What to do with data

- ➢ Plugins for access to existing libraries

```
plugins:
  decl_hdf5:
    - file: meta${pcoord[0]}x${pcoord[1]}.h5
      write: [ dsize, psize ]
```

# What is PDI?

**MPI**

| Code | Code | . . . | Code |
|------|------|-------|------|
| PDI | PDI | . . . | PDI |

Data references

PDI Data Store

```
plugins:
  decl_hdf5:
    - file: meta${pcoord[0]}x${pcoord[1]}.h5
      write: [ dsize, psize ]
```

| plugin | plugin | . . . | plugin |
|--------|--------|-------|--------|

API
API
API

API    API    API

HPC Library

➤ PDI data store: a dict of buffer references

  ➢ Name ⇒ unique identifier

  ➢ Reference

    • Ownership & locking information

      • RW-lock: Single Writer / Multiple Readers

      • Memory ownership : Strong or Semi-weak

  • Type ⇒ memory layout and interpretation

  • Buffer address ⇒ pointer to user memory (CPU/GPU[WIP])

PDI Data Store

```
#pragma pdi metadata
int buffer_size;
#pragma pdi size:[$buffer_size+1]
double *main_buffer;
```

➤ Data type: memory layout & sematics

- ➢ Annotations (C/C++), fully automatic (Python), or YAML (Fortran)
- ➢ MPI / HDF5 inspired model: scalar / array / record

➤ "Data" vs. "Metadata"

- ➢ PDI only handles the pointer for "data"
  - Minimal overhead
- ➢ PDI keeps a copy of "metadata"
  - Can be used in $-expressions

Kevin Barre

➤ PDI data store: a dict of buffer references

- ➤ Name ⇒ unique identifier
- ➤ Reference
  - Ownership & locking information
    - RW-lock: Single Writer / Multiple Readers
    - Memory ownership : Strong or Semi-weak
  - Type ⇒ memory layout and interpretation
  - Buffer address ⇒ pointer to user memory (CPU/GPU[WIP])
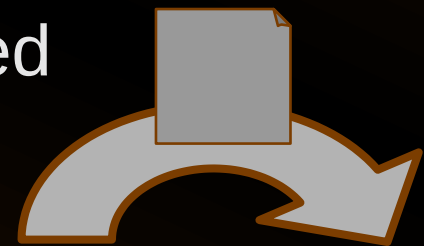
PDI Data Store

➢ **PDI data store: a dict of buffer references**

  ➢ Name ⇒ unique identifier

  ➢ Reference

   • Ownership & locking information
     ◦ RW-lock: Single Writer / Multiple Readers
     ◦ Memory ownership : Strong or Semi-weak

   • Type ⇒ memory layout and interpretation

   • Buffer address ⇒ pointer to user memory (CPU/GPU[WIP])

PDI Data Store

➢ **Notification system: plugins register to be called**

  ➢ On data share / access

  ➢ On arbitrary locations in code (named "events")

```
/** Initializes PDI */
PDI_status_t PDI_init(PC_tree_t yaml_conf);

/** Finalizes PDI */
PDI_status_t PDI_finalize();
```

a C / C++ API
Also available for:
➢ Fortran
➢ Python

➢ **Init** takes the specification tree as parameter

  ➢ The YAML is parsed using the *paraconf* library

➢ **Finalize** releases all PDI-related resources

```
typedef enum { PDI_IN, PDI_OUT, PDI_INOUT } PDI_inout_t;

// A data buffer is ready (filled)
PDI_status_t PDI_share(const char *name, void *data, PDI_inout_t access);

// A buffer will be reused
PDI_status_t PDI_reclaim(const char *name);
```

a C / C++ API
Also available for:
➢ Fortran
➢ Python

➢ **Share**

➢ A buffer is in a coherent consistent state

➢ Reference the buffer in PDI store

➢ **Reclaim**

➢ The buffer will be reused for a different use

➢ Un-reference the buffer in PDI store

```
double* data_buffer = malloc( buffer_size*sizeof(double) );

while ( !computation_finished )
{
    compute_the_value_of( data_buffer, /*...*/ );
    PDI_share("main_buffer", data_buffer, PDI_OUT);
    do_something_without_data_buffer();
    do_something_reading( data_buffer, /*...*/ );
    PDI_reclaim("main_buffer");
    update_the_value_of( data_buffer, /*...*/ );
}
```

buffer is shared
- between here
  ...
- and here

➢ Creates a "shared region" in code where
  ➢ Data referenced in PDI store
  ➢ Plugins can use it

➢ Code should refrain from
  ➢ modifying it (**PDI_IN|OUT**)
  ➢ accessing it (**PDI_IN**)

```c
typedef enum { PDI_IN, PDI_OUT, PDI_INOUT } PDI_inout_t;

// Combine share & expose for 1 piece of data
PDI_status_t PDI_expose(const char *name, void *data, PDI_inout_t access);

// When there is no data… (an interesting location has been reached)
PDI_status_t PDI_event(const char* event);

// And when there is multiple data
PDI_status_t PDI_multi_expose(const char *event_name,
    void *data, PDI_inout_t access,
    …,          …,
    NULL );
```

a C / C++ API
Also available for:
➢ Fortran
➢ Python

➢ Expose = share + reclaim

➢ Events: similar to exposing empty data

➢ Multi-expose:
  ➢ All share
  ➢ An event
  ➢ All reclaims

## In code

➢ Write & annotate your code

➢ Annotate buffers availability (share / reclaim)

➢ Compile and… DONE! (on the code side)

## In YAML

➢ Use pre-made plugins or write your own code to choose I/O libraries, describe behavior

   ➢ React to events

   ➢ Access data in the store

```c
PDI_expose("buffer_size", &buffer_size, PDI_OUT);
double* data_buffer = malloc( buffer_size*sizeof(double) );

while ( iteration_id < max_iteration_id )
{
    compute_the_value_of( data_buffer, /*...*/ );
    PDI_share("main_buffer", data_buffer, PDI_OUT);
    do_something_reading( data_buffer, /*...*/ );
    PDI_reclaim("main_buffer");
}
```

➢ Write data in the HDF5 format

➢ Heavily relies on
  ➢ $-expressions
  ➢ default configuration values

➢ Makes
  ➢ *Simple* things **easy**
  ➢ *Complex* things **possible**
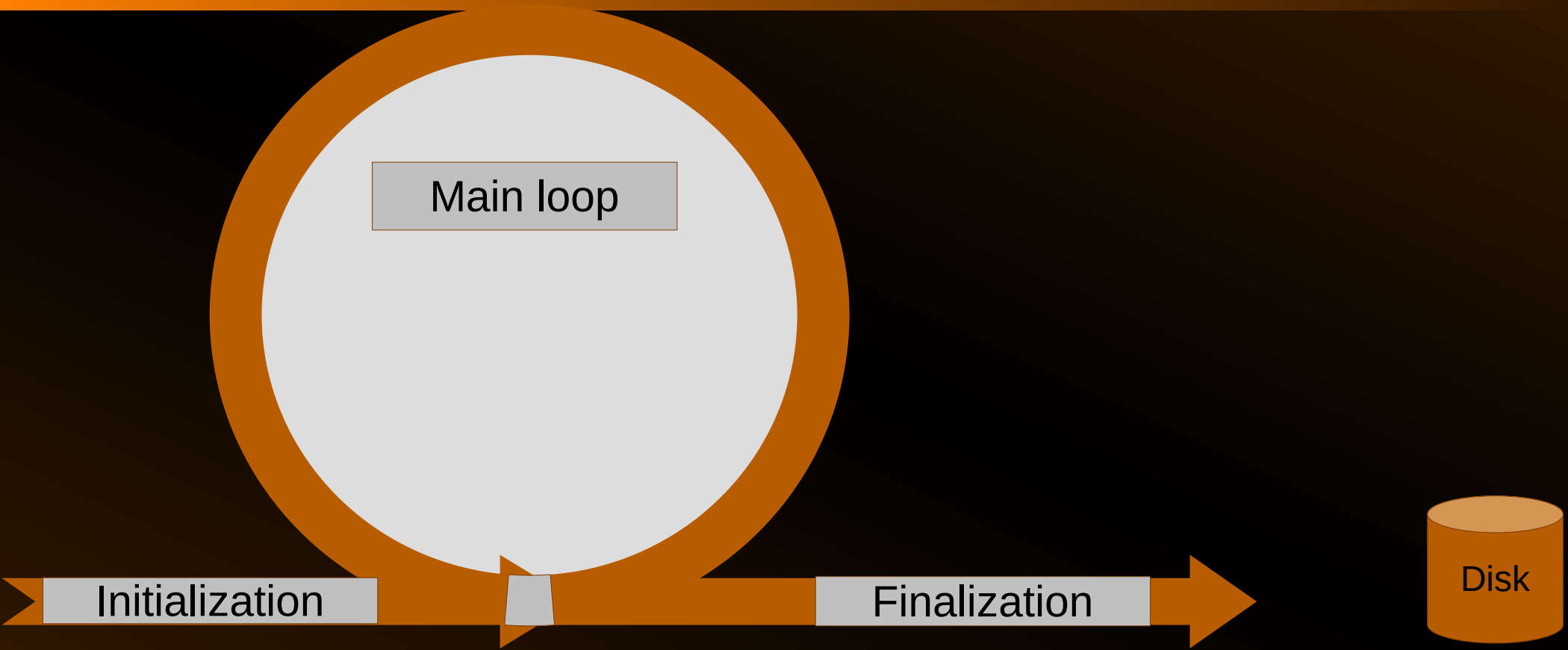
```yaml
plugins:
  decl_hdf5:
    file: 'my_file_${iteration_id}x${rank}.h5'
    write: main_buffer
```

➢ Simple to just dump data as HDF5
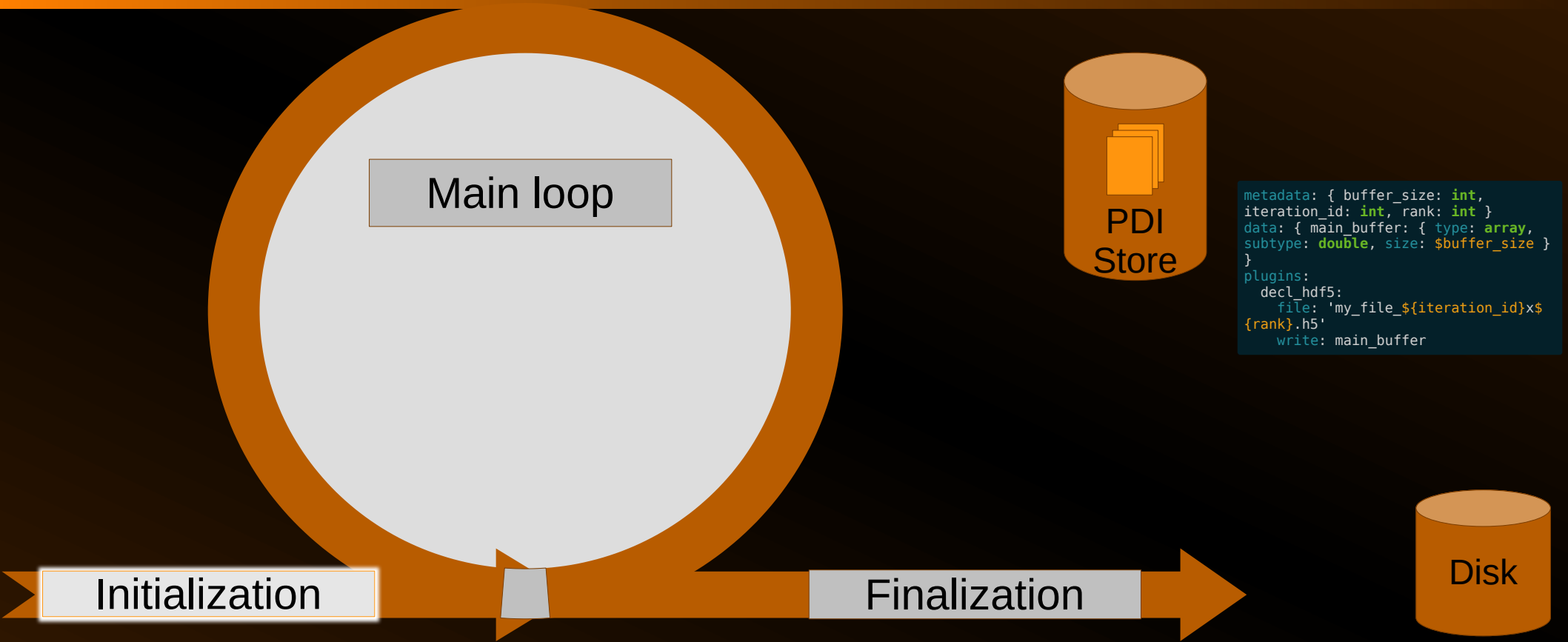
PDI

```yaml
plugins:
  decl_hdf5:
    file: 'my_file.h5'
    when: '$iteration_id % 100 = 0 & $iteration_id < 10000'
    datasets:
      main_dset:
        type: array
        subtype: double
        Size: [ '($buffer_size - 2) * $np', 100 ]
    write:
      main_buffer:
        memory_selection: { start: 1, size: '$buffer_size - 2' }
        dataset: main_dset
        dataset_selection:
          start: [ '($buffer_size - 2) * $iteration_id', '$iteration_id/100' ]
          size: [ '$buffer_size - 2', 1 ]
    communicator: $MPI_COMM_WORLD
  mpi:
```

➢ Possible to do complex rearranging of data in parallel

Main loop

Initialization

Finalization

Disk

Main loop

PDI
Store

```
metadata: { buffer_size: int,
iteration_id: int, rank: int }
data: { main_buffer: { type: array,
subtype: double, size: $buffer_size }
}
plugins:
  decl_hdf5:
    file: 'my_file_${iteration_id}x$
{rank}.h5'
    write: main_buffer
```

Initialization

Finalization

Disk

Main loop

PDI
Store

```
metadata: { buffer_size: int,
iteration_id: int, rank: int }
data: { main_buffer: { type: array,
subtype: double, size: $buffer_size }
}
plugins:
  decl_hdf5:
    file: 'my_file_${iteration_id}x$
{rank}.h5'
    write: main_buffer
```

Decl'HDF5
plugin

Initialization

Finalization

Disk

PDI

Data references

Main loop

PDI Store

```
metadata: { buffer_size: int,
iteration_id: int, rank: int }
data: { main_buffer: { type: array,
subtype: double, size: $buffer_size }
}
plugins:
  decl_hdf5:
    file: 'my_file_${iteration_id}x$
{rank}.h5'
    write: main_buffer
```

Decl'HDF5 plugin

Initialization

Finalization

Disk

```
metadata: { buffer_size: int,
iteration_id: int, rank: int }
data: { main_buffer: { type: array,
subtype: double, size: $buffer_size }
}
plugins:
  decl_hdf5:
    file: 'my_file_${iteration_id}x$
{rank}.h5'
    write: main_buffer
```

Main loop

Data references

PDI Store

Events

Decl'HDF5 plugin

Initialization

Finalization

Disk

Main loop

Data references

PDI Store

```
metadata: { buffer_size: int,
iteration_id: int, rank: int }
data: { main_buffer: { type: array,
subtype: double, size: $buffer_size }
}
plugins:
  decl_hdf5:
    file: 'my_file_${iteration_id}x$
{rank}.h5'
    write: main_buffer
```

Events

Decl'HDF5 plugin

Initialization

Finalization

Disk

Main loop

Data references

PDI Store

```
metadata: { buffer_size: int,
iteration_id: int, rank: int }
data: { main_buffer: { type: array,
subtype: double, size: $buffer_size }
}
plugins:
  decl_hdf5:
    file: 'my_file_${iteration_id}x$
{rank}.h5'
    write: main_buffer
```

Events

Decl'HDF5 plugin

Initialization

Finalization

Disk

Main loop

Data references

PDI Store

```
metadata: { buffer_size: int,
iteration_id: int, rank: int }
data: { main_buffer: { type: array,
subtype: double, size: $buffer_size }
}
plugins:
  decl_hdf5:
    file: 'my_file_${iteration_id}x$
{rank}.h5'
    write: main_buffer
```

Events

Decl'HDF5 plugin

Initialization

Finalization

Disk

- IO libraries

  - HDF5 / parallel HDF5, NetCDF4 / pNetCDF4, SIONlib

- Special purpose IO

  - FTI, ADIOS / SENSEI

- Workflow integration

  - Dask w. Deisa, FlowVR, Melissa

- Your own code

  - $-expressions based language, Python, C, C++, Fortran

PDI

```
plugins:
  pycall:
    on_event:
      trigger_event_name: # event that triggers the call
        with: { iter: $iteration_id, original_data: $main_field }
        exec: |
          if iter<1000:
              new_data = original_data*4 # uses numpy
              pdi.expose('new_data', new_data, pdi.OUT);
```

➢ Let you call your own Python code

  ➢ Data is exposed as numpy arrays

  ➢ Numpy arrays can be re-exposed

    ⇒ In-process post-processing and data transformation

PDI

```yaml
plugins:
  user_code:
    on_event:
      trigger_event_name: # event that triggers the call
        function_name { in1: $iteration_id, in2: $main_field }
```

```c
void function_name(void)
{
    int* iter = NULL; PDI_access("in1", &iter, PDI_IN);
    double* main_field = NULL; PDI_access("in2", &iter, PDI_IN);
    // …
    PDI_release("in2");
    PDI_release("in1");
}
```

➢ Let you call your own (C/Fortran) functions

  ➢ When performance matters

  ➢ To call library APIs not covered by plugins

```python
from sklearn.decomposition import IncrementalPCA
import yaml, json
import h5py
# load the simulation configuration
simu = yaml.load(open('simulation.yml'))
# Load data from HDF5
gtemp = h5py.File('data.hdf5',mode='r')['gtemp']
# process each time-step independently
for step in range(0, simu['timesteps']):
  pca = IncrementalPCA(n_components=2, copy=False,
                       svd_solver='randomized')
  pca.fit(gtemp[step,:,:])
  print(pca.explained_variance_)
```



Asahi, Y. & Fujii, K. & Heim, D. & Maeyama, S. & Garbet, X. & Grandgirard, V. & Sarazin, Y. & Dif-Pradalier, G. & Idomura, Y. & Yagi, M. (2021). "**Compressing the time series of five dimensional distribution function data from gyrokinetic simulation using principal component analysis**". *Physics of Plasmas*. 28. 012304. 10.1063/5.0023166.
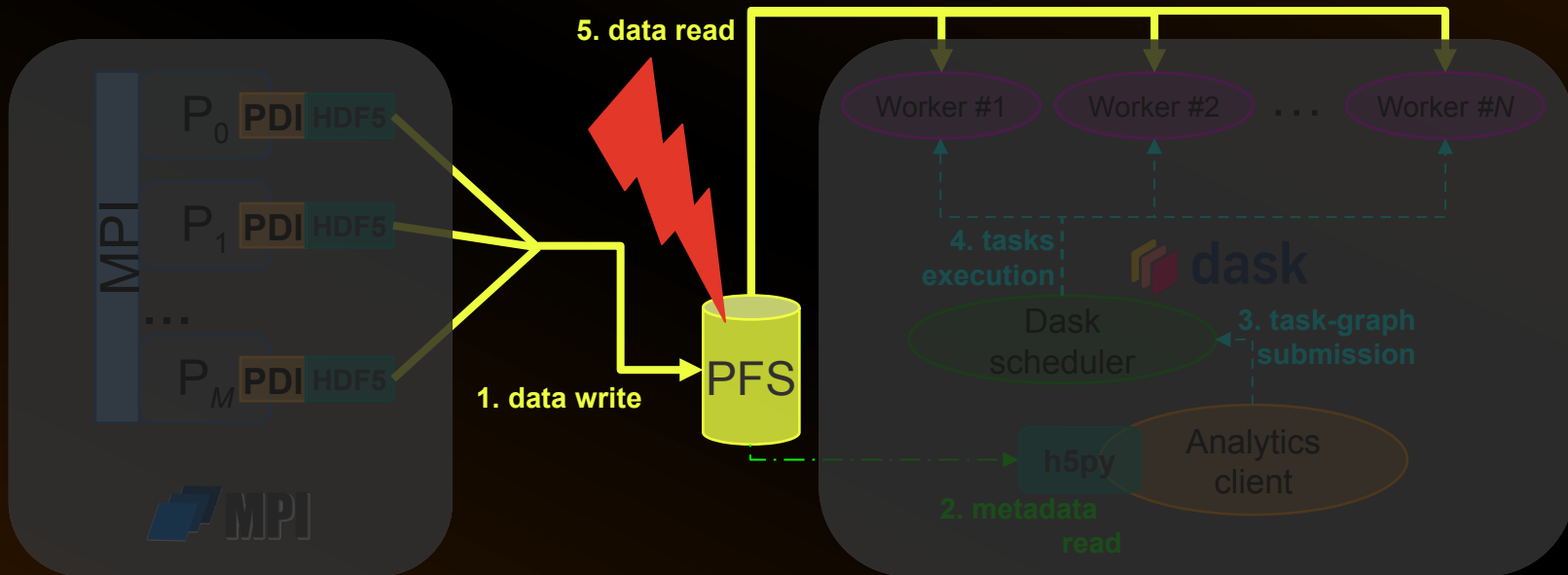
# Dask distributed?

➢ A scheduler/workers (+client) model to run work (each on its own process/node)

➢ A task-based model to describe work

➢ Many tools ported to dask for ease of use

  ➢ Numpy / SciPy

  ➢ Scikit-learn

  ➢ Pandas

  ➢ ...

```python
1  from sklearn.decomposition import IncrementalPCA
2  import yaml, json
3  import h5py
4  # load the simulation configuration
5  simu = yaml.load(open('simulation.yml'))
6  # Load data from HDF5
7  gtemp = h5py.File('data.hdf5',mode='r')['gtemp']
8  # process each time-step independently
9  for step in range(0, simu['timesteps']):
10   pca = IncrementalPCA(n_components=2, copy=False,
11                        svd_solver='randomized')
12   pca.fit(gtemp[step,:,:])
13   print(pca.explained_variance_)
```

Asahi, Y. & Fujii, K. & Heim, D. & Maeyama, S. & Garbet, X. & Grandgirard, V. & Sarazin, Y. & Dif-Pradalier, G. & Idomura, Y. & Yagi, M. (2021). "**Compressing the time series of five dimensional distribution function data from gyrokinetic simulation using principal component analysis**". *Physics of Plasmas*. 28. 012304. 10.1063/5.0023166.

```python
1  import dask.array as da
2  from dask_ml.decomposition import IncrementalPCA
3  import yaml, json
4  import h5py
5  # Connect to Dask
6  sched = json.load(open('sched.json'))
7  client = dask.distributed.Client(sched["address"])
8  # load the simulation configuration
9  simu = yaml.load(open('simulation.yml'))
10 # Build a lazy array descriptor from HDF5
11 gtemp = h5py.File('data.hdf5',mode='r')['gtemp']
12 gtemp = da.from_array(gtemp, chunks=(1,4096,4096))
13 for step in range(0, simu['timesteps']):
14   pca = IncrementalPCA(n_components=2, copy=False,
15                        svd_solver='randomized')
16   pca.fit(gtemp[step,:,:])
17   print(pca.explained_variance_)
```

Asahi, Y. & Fujii, K. & Heim, D. & Maeyama, S. & Garbet, X. & Grandgirard, V. & Sarazin, Y. & Dif-Pradalier, G. & Idomura, Y. & Yagi, M. (2021). "**Compressing the time series of five dimensional distribution function data from gyrokinetic simulation using principal component analysis**". *Physics of Plasmas*. 28. 012304. 10.1063/5.0023166.

# Dask for post hoc analytics

```
plugins:
  decl_hdf5:
    file: 'my_file_${iteration_id}x${rank}.h5'
    write: main_buffer
```



5. data read

MPI

P$_0$ **PDI** HDF5

P$_1$ **PDI** HDF5

...

P$_M$ **PDI** HDF5

1. data write

PFS

Worker #1    Worker #2    ...    Worker #$N$

4. tasks execution

dask

Dask scheduler

3. task-graph submission

h5py    Analytics client

2. metadata read

- File-system requirements are huge
  - Let's run simulation & analysis at the same time
  - Erase files as soon as they are not required anymore

- File-system performance is still an issue

Yuuichi Asahi (JAEA)
Antoine Lavandier (MdlS)

# Dask in situ with Deisa



```
plugins:
  deisa:
    scheduler_file: "/home/user/xp/sched.json"
    transfer: { main_field: { when: "$iteration_id>0" } }
```

**5. data send**

**MPI**

$P_0$ **PDI** **DEISA Bridge**

$P_1$ **PDI** **DEISA Bridge**

...

$P_M$ **PDI** **DEISA Bridge**

**1. metadata send**

**2. metadata fetch**

**4. contract send**

Worker #1  Worker #2  ...  Worker #N

**5. tasks execution**

**dask**

Dask scheduler

**3. task-graph submission**

**DEISA Metadata adapter**

Analytics client

Amal Gueroudji, Julien Bigot, Bruno Raffin. **"DEISA: dask-enabled in situ analytics."** *HiPC 2021 - 28th International Conference on High Performance Computing, Data, and Analytics*, Dec 2021, virtual, India

Amal Gueroudji. "**Distributed Task-Based In Situ Data Analytics for High-Performance Simulations**". *PhD Thesis*, Université Grenoble Alpes [2020-..], 2023. English.

Amal Gueroudji (MdlS)

```
1   import dask.array as da
2   from dask_ml.decomposition import Increme
3   import yaml, json
4   import deisa
5   # Connect to Dask
6   sched = json.load(open('s
7   client = dask.distribu
8   # load the simulat
9   simu = yaml.lo
10  # Get data
11  gtemp
12  fo
                            , copy=False,
                        ='randomized')

                    ance_)
```

Used in production for grand-challenge on Adastra (CINES) #10 Top500
Multi-day full-scale run on the whole GPU partition

(a) Problem size/nodes (MB)

Corentin Roussel (MdlS)
Kai Keller (BSC)

➢ 4 versions of Gysela

  ➢ No checkpoint

  ➢ HDF5 checkpoints

  ➢ FTI fault-tolerance

  ➢ PDI (none / HDF5 / FTI / HDF5+FTI)

Execution time by MB of checkpointed data on 4 MareNostrum
Nodes with and without PDI

Corentin Roussel (MdIS)
Kai Keller (BSC)



Gysela Wallclock time in weak scaling on Curie (TGCC – France) with and without PDI

Checkpointed data ~2.1GB/node

PDI


(a) Instructions ($\times 10^{12}$)


(b) Instructions ($\times 10^{12}$)

Corentin Roussel (MdlS)
Kai Keller (BSC)


(c) Instructions ($\times 10^{12}$)

Memory usage during a Gysela execution with and without PDI on 4 nodes of MareNostrum (BSC – Spain)

# PDI In practice

- PDI is publicly available (BSD 3-clause license)
  - Regular releases since 2014
  - Packages available for Debian, Fedora, Ubuntu, Spack
  - Documentation & tutorials available @ https://pdi.dev/1.6/
  - Heavily tested & validated
    - more than 1500 tests…
    - …running on more than 12 platforms each
- Integration in production codes
  - Gysela, Parflow, ESIAS, Manta?, …
- Part of NumPEx software stack

# Conclusion

➢ An API for Data Coupling, Not an IO library

➢ A declarative annotation API

➢ Multiple plugins for and data processing

➢ Describe your IO from YAML

➢ Switch to in situ processing or more without even recompiling

➢ Your turn now!

➢ Get the doc: https://pdi.dev/master/

➢ Join the fun on slack https://join.slack.pdi.dev/

WRITE BANDWIDTH WITH BLOCKSIZE OF 128 kiB

WRITE BANDWIDTH WITH BLOCKSIZE OF 256 MiB

IOR IO Benchmark PDI integration

Scaling with small (128k) & large (256M) data blocks
on CRESCO6

Francesco Iannone
(ENEA)

## Setup:

- Ruche cluster
  - 192 nodes (2 CPUs 20 cores each, 180 GB)
  - Omni-Path 100 Gbit/s
  - Spectrum Scale GPFS (IOs rate: 9 GB/s)
- Mini-app
  - 2D heat solver
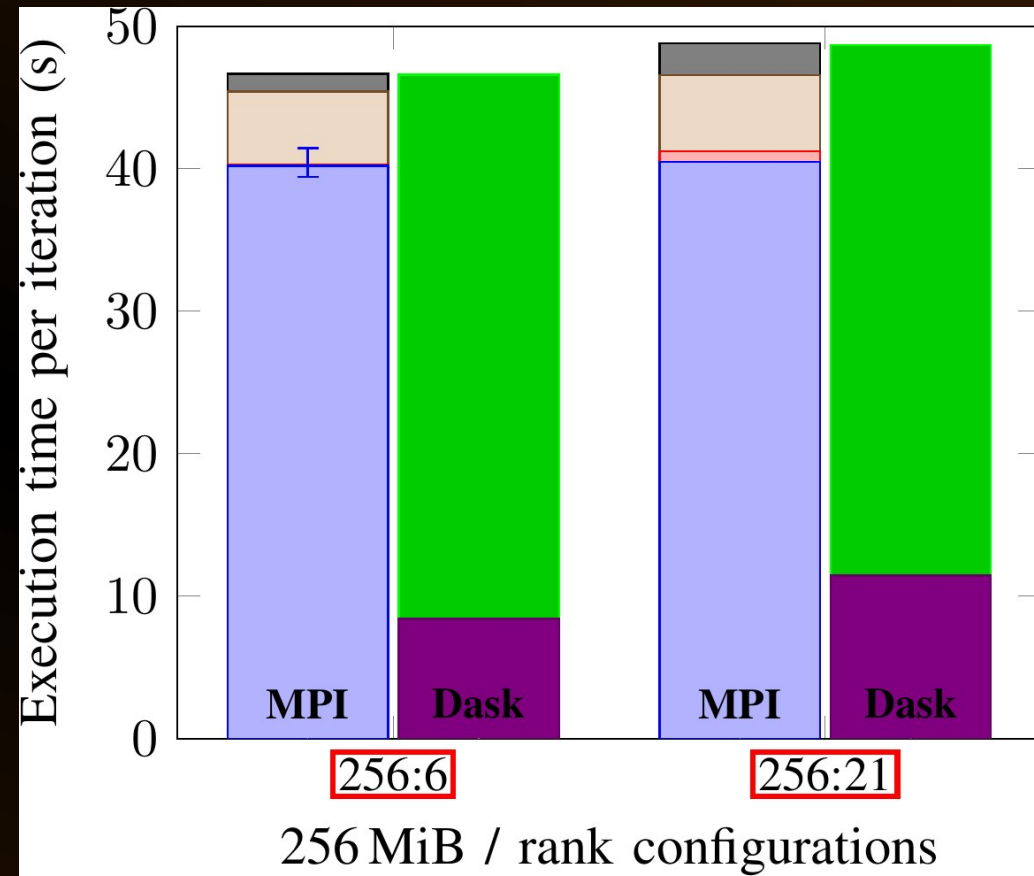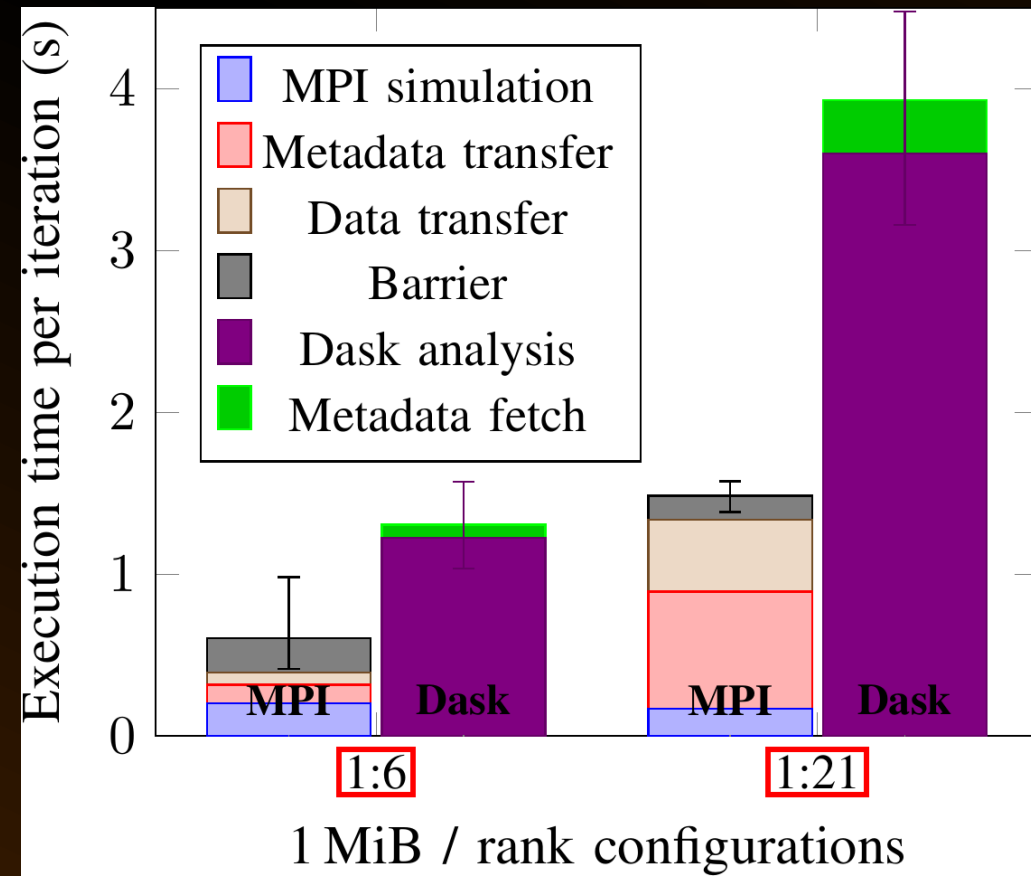  - Incremental Principal Component Analysis

- Weak scaling
  - X + Y cores
  - X cores for MPI simu.
  - Y cores for Dask analytics
- No analytics
- vs. Post-hoc
- vs. DEISA

| Configuration | 128+16 | 256+32 | 512+64 |
|---|---|---|---|
| MPI processes | 128 | 256 | 512 |
| Dask workers | 16 | 32 | 64 |
| MPI nodes | 4 | 8 | 16 |
| Dask worker nodes | 1 | 2 | 4 |
| Global data size | 16 GiB | 32 GiB | 64 GiB |
| Dask generated tasks | 15210 | 29010 | 55150 |



Configurations, w. **(1)** no analytics, **(2)** DEISA, **(3)** post hoc

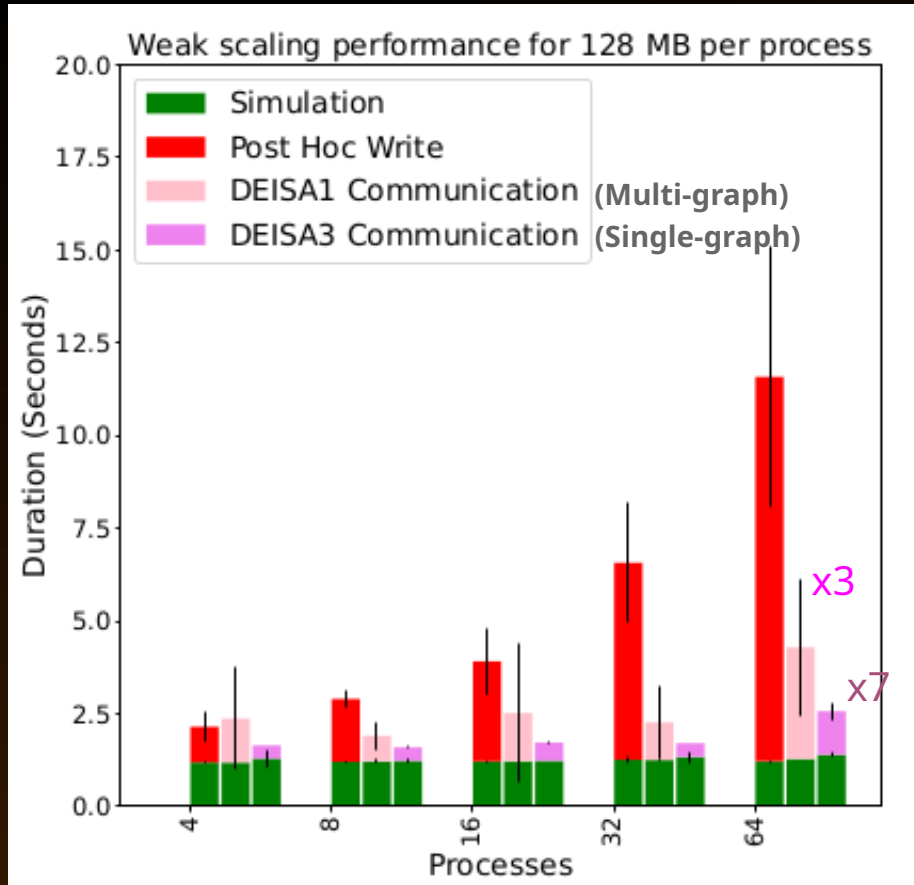mean: 430s
max: 445s
min: 421s

- IRENE supercomputer @ TGCC, France,
- Nodes:
    - 2x24-cores Intel Skylake@2.7GHz
    - 180GB RAM
- InfiniBand network (100Gb/s),
- Scratch disks: 300GB/s transfer rate
- Mini App 2D heat solver

| Parameter | Value |
|---|---|
| Number of runs | 3 |
| Number of iterations IPCA | 10 |
| Number of iteration Derivative | 12 |
| MPI nodes / Dask worker node | 2 |
| MPI process / MPI node | 2 |
| Dask worker / Dask worker node | 2 |
| Thread / Dask worker | 24 |
| MPI process / Dask worker | 2 |

| Configuration | XP1:128 MiB | XP1:256 MiB | XP1:512 MiB | XP1:1 GiB |
|---|---|---|---|---|
| MPI block size | 128 | 256 | 512 | 1 |
| Dask chunk size | 128 | 256 | 512 | 1 |
| MPI Nodes | [4, 8, 16, 32, 64, 128, 256] | | | |
| Dask Nodes | [2, 4, 8, 16, 32, 64, 128] | | | |

(c) Strong scaling results represented in hour-core for an 8 GiB problem size



(c) Strong scaling results represented in hour-core for a 8 GiB problem size

Multi-graph
-lot
metadata
-heartbit=5s

Single-graph
less
metadata
heartbit=∞