

# A Parallel I/O Abstraction for Partitioned Global Address Space Programming



### Felix Mößbauer

felix.moessbauer@campus.lmu.de Ludwig-Maximilians-Universität München





#### PGAS

- Data is distributed over units
- Global data structures know its layoutDASH

#### DASH

- C++ library similar to UPC++
- Implements PGAS approach
- Provides distributed data structures and algorithms



Data affinity

LUDWIG-MAXIMILIANS-UNIVERSITÄT MÜNCHEN

- data has well-defined owner but can be accessed by any unit
- data locality important for performance
- support for the *owner computes* execution model
- DASH:
  - unified access to local and remote data in global memory space



Array<int> arr(40)

Data affinity

- data has well-defined owner but can be accessed by any unit
- data locality important for performance
- support for the *owner computes* execution model
- DASH:
  - unified access to local and remote data in global memory space
  - and explicit views
     on local memory
     space

Array<int> arr(40)





### Why IO abstraction necessary

- IO is often boilerplate
- Parallel IO is tricky
- Re-implementation of IO part in each application

Abstraction idea

- Import / export distributed data structures
- various targets (HDF5, NetCDF, Image, ...)
- Conversion ?



- Implementation is portable (GPFS, NFS, ...)
- Minimal overhead
- Scalable

LUDWIG-MAXIMILIANS-UNIVERSITÄT MÜNCHEN

Simple interface

does not scale
<pre>// distributed data structure distdata data; if(rank == 0) {    load_data(&amp;data, file); }</pre>
<pre>MPI_Barrier();</pre>

**Design Goals** 





Similar to POSIX pipes

LUDWIG-

- Source  $\rightarrow$  IO Adapter  $\rightarrow$  Sink
- **Collective operation**

*n-dim container* OutputStream("file.hdf5") << dataset("array\_a") << [...] << array\_a;</pre> modifiers sink source







File

#### **IO Adapter Basics**

# Configuration using stream modifiers e.g.

- dataset

- store pattern
- restore pattern
- Good defaults
- IO Team is derived from container

```
namespace dio = dash::io::hdf5;
dash::NArray<int,3> data(10,50,30);
// [fill matrix]
dio::OutputStream os("filename.hdf5");
os << dio::dataset("path/to/dataset");
// kicker which passes configuration to io-adapter
os << data;</pre>
```

- Automatic type conversion (trivial types)
- User can pass custom type mapping function as lambda







- DASH: distributed datastructures
- HDF5

LUDWIG-

NÜNCHEN

- portable data format
- Independent of OS / Arch.
- Self-describing format
- Organized as a tree
- Dataset consists of
  - Tables
  - Metadata
  - Attributes



Use *h5dump < file >* to discover

#### HDF5 Sample File



**SPPEXA Workshop Japan 2017** 

LUDWIG-

#### **Implementation Aspects - Datastructure**

Container knows distribution

- Data is continuous memory
- Data domain is hyperrectangle
  - N-Dim rectangle
  - Sample: 2 dim
  - Propably underfilled blocks





#### **IO Adapter – Pattern Converter** MAXIMILIANS-UNIVERSITÄT

## **IO** Subtask

LUDWIG-

MÜNCHEN

- Has to be executed collectively
- Size of contribution might differ
- Order of tasks irrelevant for correctness
- No underfilled blocks





# Order tasks by local contribution **Drop all-empty tasks**



LUDWIG-

MÜNCHEN

### IO Adapter – Async IO (1)

## Observations

- IO tasks are long-running
- IO libs internally use MPI2-IO / ROMIO
- No support for non-blocking IO
- Solution
  - Use MPI multithreading
  - Create IO tasks as threads
  - Works for versions >= MPI2
  - Problem
    - Simultaneous IO tasks are horribly slow
    - Buggy in some MPI versions (e.g. mpich)

# IO Adapter – Async IO (2)

Simultaneous tasks are chained internally
Sync / async selection via launch-policies
Example:

```
os.flush();
```



#### Each pattern has tags

LUDWIG-MAXIMILIANS-UNIVERSITÄT

- Layout (blocked, row\_major, ...)
- Mapping (balanced, shifted, ...)
- Partitioning (balanced, regular, ...)
- Checked using type traits at compile time
  - Zero overhead at runtime
- At runtime using exceptions
  - Important for reading data
  - Not all pattern support underfilled blocks



- Minimal overhead (in pattern conversion)
- zero-copy of data

LUDWIG-

MÜNCHEN

- Performance depends on FS
- Large tiles are better
- SuperMUC Intelmpi, HDF5 on GPFS
  - Read: 1.2 GB/s per node
  - Write: 0.8 GB/s per node



#### **SPPEXA Workshop Japan 2017**

#### Tokyo, April,6 2017 | 19

All-Pairs latency benchmark

LUDWIG-MAXIMILIANS-UNIVERSITÄT MÜNCHEN

- Measures n-to-n latency multiple times
- Data is stored in Narray
- Post Processing
  - Import HDF5 data in R
  - Analyze and Plot

#### SuperMUC II, 20 Nodes



- Subsetting Data Structures
- Adios Bindings

- Instant support of various IO backends
- Converter DASH Pattern  $\rightarrow$  Adios Pattern
- DASH Coarray (Similar to CAF 2008)
  - Fulfills the Narray container concept  $\rightarrow$  IO works out of the box

#### Summary



LUDWIG-MAXIMILIANS-UNIVERSITÄT MÜNCHEN

# IO Abstraction

- Independent of data and grid layout
- Easy to use, but fast
- No special target FS

Use cases

- Import and export data
- Checkpointing

Reference Implementation part of DASH

#### **Acknowledgements**



LUDWIG-MAXIMILIANS-UNIVERSITÄT

> DASH on GitHub: https://github.com/dash-project/dash

IO stuff part of release 0.3.0

Funding





Bundesministerium für Bildung und Forschung

### The DASH Team

T. Fuchs (LMU), R. Kowalewski (LMU), D. Hünich (TUD), A. Knüpfer (TUD), J. Gracia (HLRS), C. Glass (HLRS), H. Zhou (HLRS), K. Idrees (HLRS), J. Schuchart (HLRS), F. Mößbauer (LMU), K. Fürlinger (LMU)