



Toward IOVerbs with the Memory-Centric Storage for Exascale (MCSE) Project



Table of contents

- 1 Introduction
- 2 Semantics
- 3 Conclusions

Goals

- Exploration of alternative storage APIs
- Development of IO-Verbs to express elementary parallel I/O
- ⇒ Uniform interface for memory and persistent storage
- Exploit different storage classes flexible in HPC-workflows (campaigns)
- Development of the Memory-Centric Storage System (MCS2)
- ⇒ Based upon IML and MCSE
- Migration path for existing applications

The Project

Key information

- 3-years funding by BMBF
- Funding Reference: 16ME0663K
- Start date: 10/2022

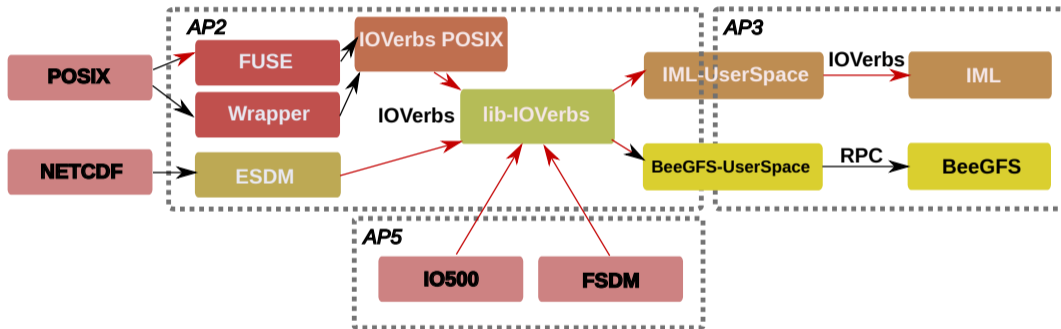
Partners

- University Göttingen (Coordinator)  GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN IN PUBLICA CONCORDIA
SINCE 1737
- Fraunhofer ITWM  Fraunhofer
ITWM
- Deutsches Zentrum für Luft- und Raumfahrt e.V.  DLR
- ThinkParQ  thinkparQ

What are IOVerbs?

- Analogon to IBVerbs but for I/O
- Describe I/O on the low-level and enabling parallel I/O
- Have a clear semantics
- Independent of storage system
- Allows high-level libraries to be layered on top
- Should be useful on local systems
- Maybe a prospective replacement of API in Linux VFS

IO Layers



Work Plan

AP7: Management

AP1: IOVERBS

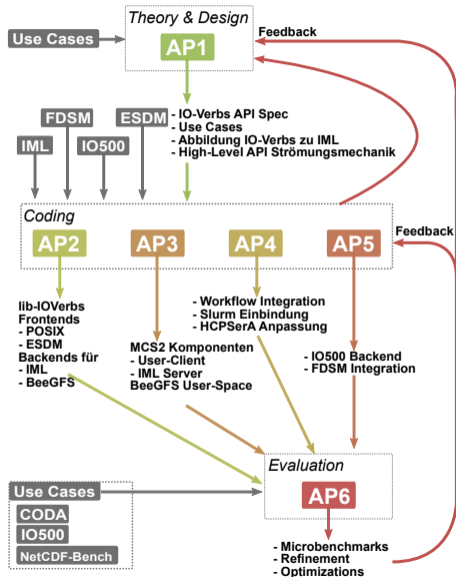
AP2: Implementation

AP3: MCS2

AP4: Campaign Storage

AP5: Anwendungsportierung

AP6: Evaluation



Outline

- 1 Introduction
- 2 Semantics**
- 3 Conclusions

Aspects for Semantics (Excerpt)

- Visibility: when will another proc see changes
 - Durability: when is data written back
 - Concurrency: what happens with concurrent access
 - Synchronous/Asynchronous API: App can proceed during I/O – background
 - In-transit: Allow computation during I/O for pre-/post-processing
 - Transparency: no hidden optimization
 - Optimistic vs. pessimistic approaches
- ⇒ There are many ways to describe semantics in the API

Motivating Example: Shallow Water Model

Facts about the model

- Stores data column-wise in memory
- Separates compute phase and IO phase¹

Existing NetCDF code for IO phase limits performance

```
1  size_t start[] = {0, 0};
2  size_t count[] = {nY, 1};
3  for(unsigned int col = 0; col < nX; col++) {
4      start[1] = col; //select col (dim "x")
5      nc_put_vara_float(dataFile, i_ncVariable, start, count,
6          &i_matrix[col+boundarySize[0]][boundarySize[2]]);
7  }
```

¹DSLs will help to separate those phases

ESDM Code for the Application (that worked)

```
1 int64_t offset[] = {(int64_t) timeStep, offsetY, offsetX};
2 int64_t size[] = {1, (int64_t) nY, (int64_t) nX};
3
4 esdm_wstream_float_t stream;
5 esdm_wstream_start(&stream, dset, 3, offset, size);
6 for(int y = 0; y < nY; y++) {
7     for(int x = 0; x < nX; x++) {
8         esdm_wstream_pack(stream, i_matrix
9             [x + boundarySize[0]][boundarySize[2] + y])
10         // this may trigger actual IO and postprocessing!
11     }
12 }
13 esdm_wstream_commit(stream);
```

- Ultimately, using DSLs an IO phase could switch between compute and "stream output" to minimize memory pressure (and trigger initial post processing)

IO500 Example - IOR Easy Write

- Each process operates on an (independent) file

```
1 fd = open(path, O_CREAT)
2 offset = 0
3 for I in SegmentCount:
4     write(fd, offset, Blocksize)
5     offset += Blocksize
6 close(fd)
```

An alternative formulation

```
1 fd = iov_open_async(path) // we can proceed *before* file is opened
2 offset = 0
3 iov_write_limit(fd, SegmentCount * Blocksize) // organizes space - note: Stonewall
4 for I in SegmentCount:
5     res = iov_write_excl_data_ready(fd, offset, Blocksize) // caching, no sync
6     offset += Blocksize // also user must check for error
7 iov_wait_completion(fd)
8 iov_destroy(fd)
```

IO500 Example - IOR Hard Write

- Processes operate on a shared file - write strided

An alternative formulation

```
1 fd = iov_open_async_coll(path) // collective
2 // includes amount of data that this proc and all procs write
3 iov_write_coll_limit(fd, SegmentCount*Blocksize, nproc*SegmentCount*Blocksize)
4 for I in SegmentCount:
5     res = iov_write_excl_data_ready(fd, offset, Blocksize) // exclusive: non-overlapping
6 iov_wait_completion(fd)
7 iov_destroy(fd)
```

IO500 Beispiel - MDTest Hard Write

■ Each process creates files ... (suboptimal anyway)

```
1 for f in file: // the recursive truth is too complex ;-)  
2     fd = open(path/f, O_CREAT)  
3     write(fd, 0, Blocksize)  
4     close(fd)
```

An alternative formulation

```
1 iov_limit_objects(no of files locally) //system could create container/return limits  
2 for f in file:  
3     fd = iov_open_async(path) // we can proceed *before* file is opened  
4     iov_write_limit(fd, Blocksize)  
5     iov_write_excl_data_ready(fd, 0, Blocksize)  
6     iov_destroy(fd)  
7     err = iov_check_state() // check async if there was an error so far, alternatively  
     ↪ use iov_register_err_handler() before  
8     if err != IOV_NO_ERR invoke error()  
9     iov_wait_completion_all()
```

Outline

- 1 Introduction
- 2 Semantics
- 3 Conclusions**

Conclusions

- MCSE will systematically explore parallel I/O semantics
- MCSE will develop IO-verbs
 - ▶ Addressing metadata and data
- MCSE develops an integrated solution for campaigns

PMDK Read (Simple case, no transactions)

- pmemobj API version 2.3
- Example Code/Github

```
1      struct my_root {
2          size_t len;
3          char buf[MAX_BUF_LEN];
4      };
5      void readIn(){
6          PMEMObjpool *pop;
7          pop = pmemobj_open(path, LAYOUT);
8          // Get the PMEMObj root
9          PMEMoid root = pmemobj_root(pop, sizeof(struct my_root));
10         // Pointer for structure at the root
11         struct my_root *rootp = pmemobj_direct(root);
12         printf("%s\n", rootp->buf);
13         pmemobj_close(pop); // Close PMEM object pool
14     }
```

PMDK Write

```
1 void writeOut(){
2     PMEMobjpool *pop;
3     pop = pmemobj_create(path, LAYOUT, PMEMOBJ_MIN_POOL, 0666);
4     // Get the PMEMObj root
5     PMEMoid root = pmemobj_root(pop, sizeof(struct my_root));
6     // Pointer for structure at the root
7     struct my_root *rootp = pmemobj_direct(root);
8     // code for write:
9     pmemobj_memcpy_persist(pop, rootp->buf, buf, strlen(buf));
10    // Assign the string length and persist it
11    rootp->len = strlen(buf);
12    pmemobj_persist(pop, &rootp->len, sizeof (rootp->len));
13    pmemobj_close(pop); // Close PMEM object pool
14 }
```