# Revisiting HDF5 SWMR and file versioning features

Elena Pourmal Scot Breitenfeld elena.pourmal@lifeboat.llc brtnfld@hdfgroup.org





## HDF5 Prototypes

- Single writer / multiple reader access to HDF5 files (VFD SWMR)
  - Available in the HDF5 repo branches "feature/vfd\_swmr\_beta\_2" or "feature/vfd\_swmr"
- HDF5 file provenance and revision/version control
  - Released in 1.14.0

#### We need your feedback on both features to finish productization!

May 26 - 28, 2025



## Acknowledgement

- VFD SWMR work is supported by the U.S. Department of Energy, Office of Science under award number DE-SC0025856 for Phase I SBIR project "Toward Robust HDF5"
- The HDF Group



### VFD SWMR

Full access to HDF5 file under construction

### Current Single Writer/Multiple Reader (SWMR)

- Introduced in HDF5 1.10.0
- Implementation has many limitations
  - Based on POSIX semantics
    - No support for NSF
  - Doesn't allow modification of the file structure
    - Only dataset's raw data can be modified
    - Groups, datasets, attributes cannot be added
  - New data may not be visible for some time
  - No support for parallel applications
  - Order of writer and reader is critical



#### <u>Lifeboat</u>

May 26 - 28, 2025





### Full Single Writer/Multiple Reader (VFD SWMR)

#### Goals

- SWMR feature
  - Allow modification of a file structure
  - Guarantee max time for new data visibility
  - No POSIX semantics requirements
  - Support NFS and Object Store
  - Lift requirement for writer/reader ordering
  - Support for VL types
  - Support for parallel applications
  - Create a tool to recover corrupted HDF5 file
- HDF5 Architecture improvement
  - Remove complex flush dependencies



#### Lifeboat

May 26 - 28, 2025

The HDF Gro

#### <u>\_\_\_\_</u>

## Overview of VFD SWMR

- See "RFC: VFD SWMR" for design details.
- Major implementation concepts and requirements:
  - Tick describes periods of operations for both writer and reader: if "t" is desired maximum latency (user defined parameter) from write to visibility on the reader side, tick is defined to be "t/3"; tick length is measured in tenths of a second.
  - VFD SWMR configuration parameters and HDF5 metadata file :
    - At the end of each tick HDF5 metadata cache is flushed to page buffer that keeps track of modified metadata pages
    - Modified metadata pages are written to storage to a metadata file
    - Written metadata pages are released in metadata file after max\_lag ticks
    - Metadata file is shared between writer and readers
    - Readers satisfy metadata I/O requests from the metadata file
  - See requirements for file creation and access properties for writer and readers on the next slides.

#### <u>Lifeboat</u>



### **VFD SWMR Architecture**



Lifeboat

May 26 - 28, 2025

HUG 25

#### 9

## VFD SWMR Programming Model - Writer

- Set creation property to use paged allocation: H5Pset\_file\_space\_strategy(fcpl, H5F\_FSPACE\_STRATEGY\_PAGE, false, 1)
- Set access properties to:
  - use latest (1.14.\*) file format with H5Pset\_libver\_bounds
  - enable page buffering with H5Pset\_buffer\_size
  - set VFD SWMR configuration properties in the H5F\_vfd\_swmr\_config\_t structure; see H5Fpublic.h for description

```
config.tick_len = 4;
config.max_lag = 7;
config.writer = true;
config.md_file_path = "./";
config.md_file_name = "my_md_file";
H5Pset_vfd_swmr_config (fapl, &config);
```

Create file with specified creation and access properties

<u>Lifeboat</u>





# VFD SWMR Programming Model - Reader

- Set access properties with the same parameters as for writer:
  - Use latest (1.14.\*) file format with H5Pset\_libver\_bounds
  - Enable page buffering with H5Pset\_buffer\_size
  - Set VFD SWMR configuration properties

```
config.tick_len = 4;
config.max_lag = 7;
config.writer = false;
config.md_file_path = "./";
config.md_file_name = "my_md_file";
H5Pset_vfd_swmr_config (fapl, &config);
```

- Open the file with specified access properties
- Reader needs to check if metadata file is available and then start read operations

Lifeboat

# VFD SWMR Demo

- Code is available from the HDF5 GitHub repo <u>https://github.com/HDFGroup/hdf5.git</u>, "feature/vfd\_swmr\_beta\_2" branch
  - See doc directory for VFD SWMR documentation
  - Check vfd-swmr-user-guide.md for building instructions and for the extensible datasets example; it compiles and runs as is.
- Demo
  - See README in the hdf5/examples directory
  - The files examples/credel.c and tools/lib/h5tools.c have been modified for the demo, but the updates have not been checked in yet.



## VFD SWMR Demo

credel.c VFD SWMR settings

H5Pget\_vfd\_swmr\_config(fapl, &config); config.version = H5F\_\_CURR\_VFD\_SWMR\_CONFIG\_VERSION; config.tick\_len = 4; config.max\_lag = 5; config.flush\_raw\_data = true; config.writer = true; config.md\_pages\_reserved = 128; config.pb\_expansion\_threshold = 50; config.maintain\_metadata\_file = true; strlcpy(config.md\_file\_path, "./", sizeof(config.md\_file\_path)); strlcpy(config.md\_file\_name, "my md file", sizeof(config.md file path));

<u>Lifeboat</u>



### Onion VFD

HDF5 file provenance and revision/version control

# Why HDF5 File Provenance?

- HDF5 data provenance is a major concern in cases when
  - Original data must be preserved
  - Changes are tracked and attributed
- HDF5 "undo" capability
  - The feature was requested for HDFView (reads and writes data)
- AI/ML applications with HDF5 used for:
  - Model checkpointing and serialization (e.g., Keras, TensorFlow, PyTorch)
  - Storing large datasets for training
  - Intermediate data storage for preprocessed data, model state and optimizer states to resume training from specific point
- RFC: Onion VFD, May 9, 2021
   <u>https://support.hdfgroup.org/releases/hdf5/documentation/rfc</u>

<u>Lifeboat</u>



## HDF5 File Provenance Storage Options (conceptual view)



# HDF5 Onion VFD Implemented Features

- Changes are tracked per file open/close
  - No "undo" capability
- Modification history is stored in a separate file (option 2) (*filename*.onion)
- Amended data is stored as fixed-size verbatim pages (size is power of 2) (vs. amended runs, i.e., only modified data)
- Metadata (header, index) can be aligned on page boundary
- Only the latest version can be modified, there is no branching.
- I/O on the onion file is performed with the same I/O driver as on the HDF5 file.
- There is no support for parallel applications.
- Each revision can be annotated.
- H5FDct1 function can be used to find onion info, revision history, to set or update revision annotation.

Lifeboat



## **Onion VFD APIs**

```
typedef struct H5FD onion fapl info t {
   uint8 t
                                       version;
   hid t
                                       backing fapl id;
   uint32 t
                                       page size;
   H5FD_onion_target_file_constant_t store_target;
   uint64 t
                                       revision num;
   uint8 t
                                       force_write_open;
                                       creation flags;
   uint8 t
    char
                                       comment[256];
} H5FD_onion_fapl_info_t;
```

#### Lifeboat





### Enhancements for productized version of Onion VFD

- Implement support for
  - Different storage options for history data
  - Different I/O drivers for history data when using options 2 and 3
  - Parallel HDF5 applications
- Other considerations:
  - Storage size
  - Onion file may become huge when working with HDF5 files
  - Branching
  - Looks like a useful feature for AI applications
  - Performance
  - Tools support
  - User and developer documentation
- Tell us about your needs and requirements!

#### <u>Lifeboat</u>



# Thank you!

Questions?

Lifeboat

May 26 - 28, 2025



20

HUG 25

