

HDF5: Toward a Community-Driven Open Source Project

October 14, 2020



Proprietary and Confidential.
© The HDF Group.

Elena Pourmal
Director of Engineering
The HDF Group

Outline



- State of HDF5
- The HDF Group efforts to revamp HDF5 as a *community-driven* Open Source project
- HDF5 Roadmap for 2021
- What is on *your* wish list?
 - Please use [Google doc](#) (see Lori's message in the chat window) to add your comments

State of HDF5

Year 2020

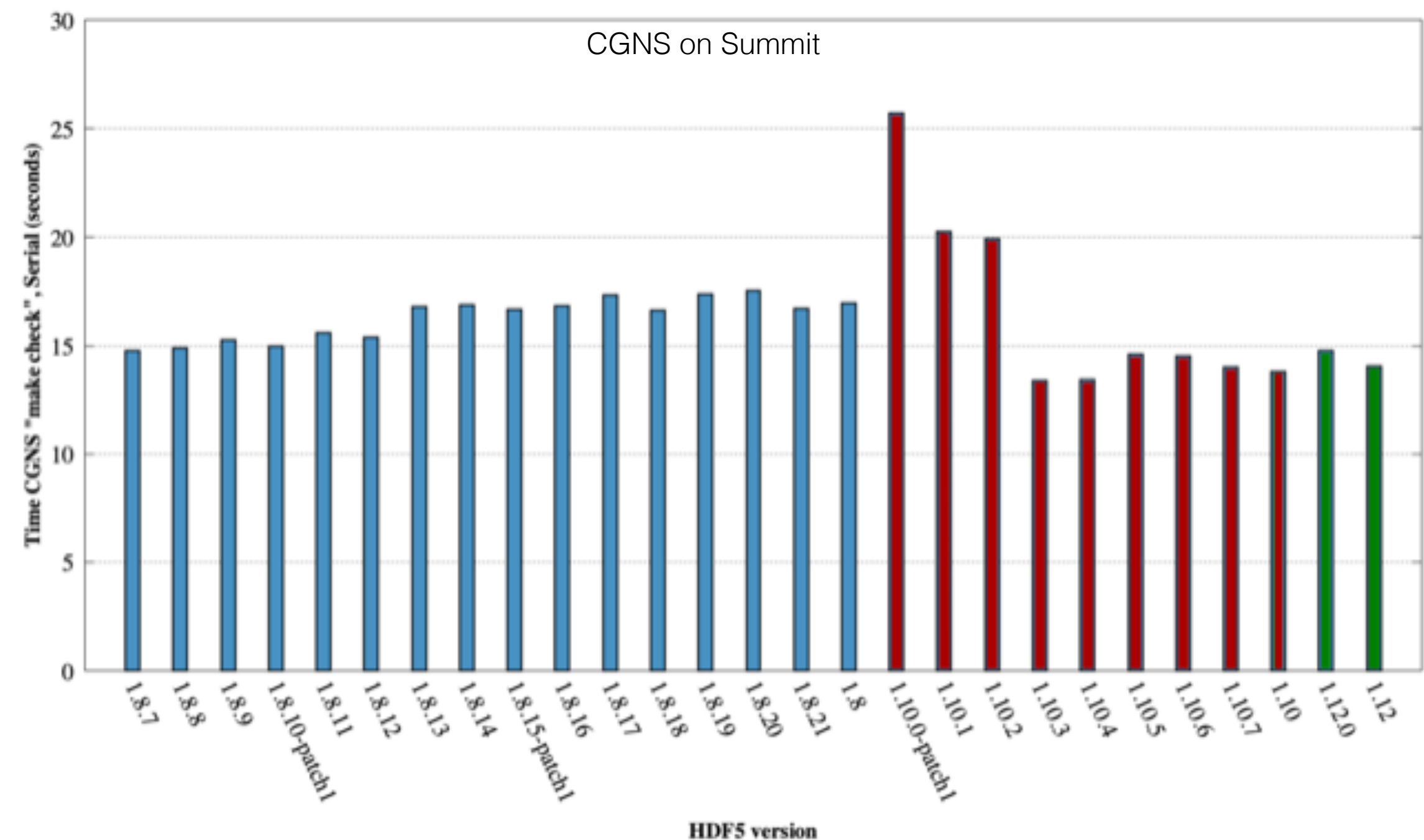
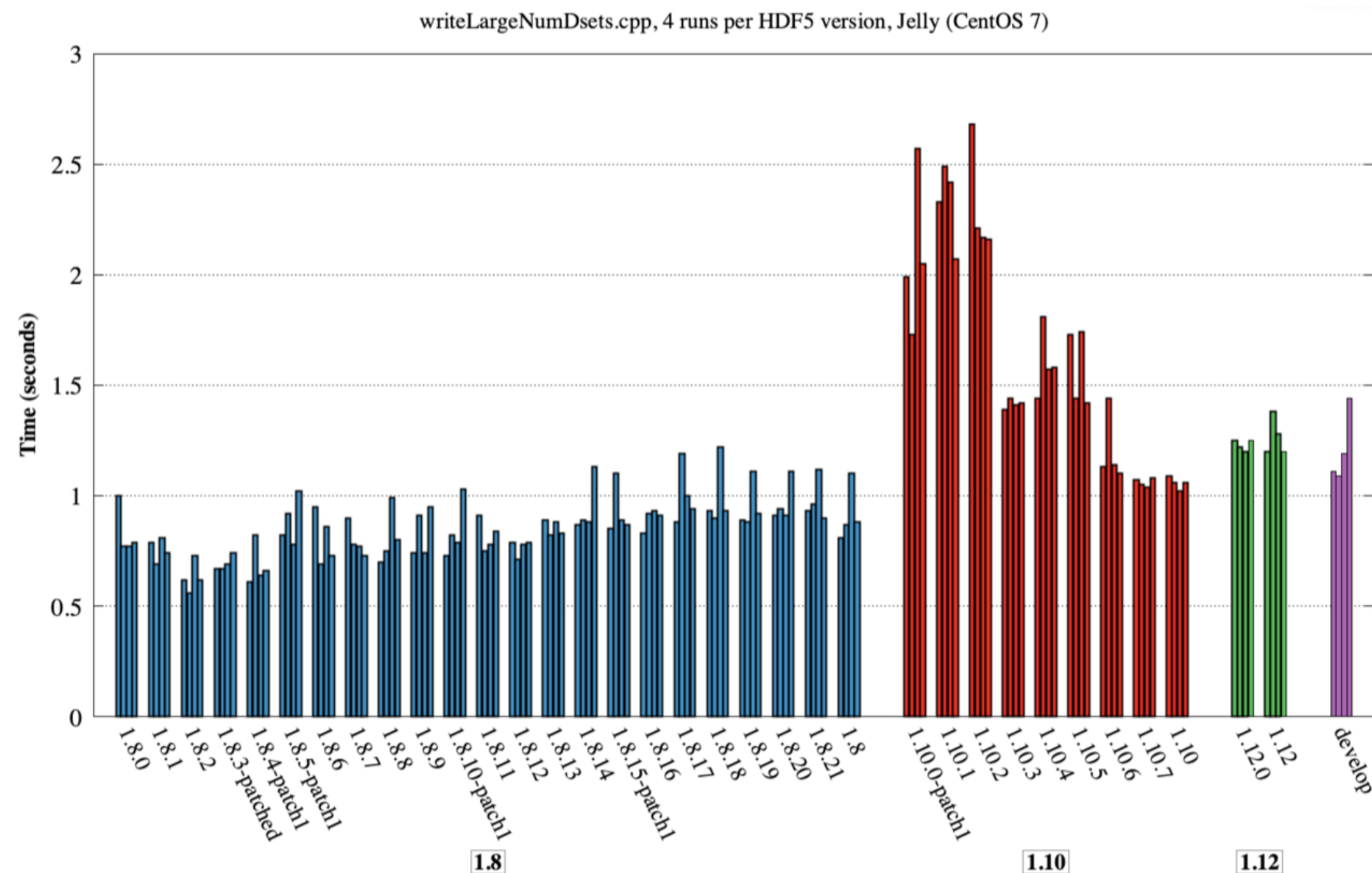
2020 Focus



- Release major version of HDF5 1.12.0 that enables access to data in Object Store and the Cloud
- Address performance gap between 1.8.* and later maintenance releases
- Improve scalability of parallel applications
- Address Common Vulnerabilities and Exposures (CVE) issues reported to us
- Reduce number of compilation warnings
- Support for [Open Source SZIP Compression \(AEC\)](#) (from the [German Climate Computing Center](#))

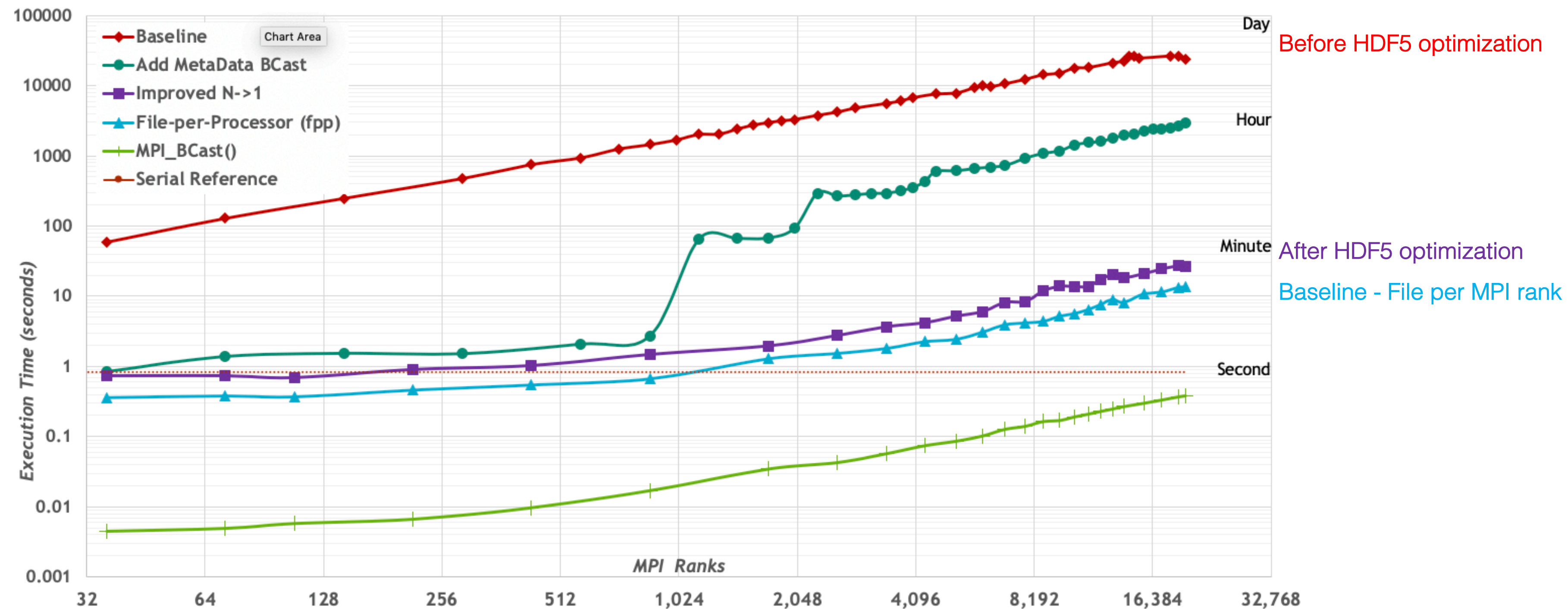
Closing Performance Gap

- Identified several bottlenecks caused by usage of
 - HDF5 property lists
 - Skip lists in metadata cache
 - Skip lists in ID look-ups
- Please share your benchmarks with us!



Improving Parallel Performance

- New property to optimize read/write an entire dataset by all ranks in communicator



Courtesy Greg Sjaardema, Sandia National Labs

HDF5 1.12.0 Release



- New features highlights
 - The Virtual Object Layer (VOL)
 - An abstraction layer within the HDF5 library that enables different methods for accessing data and objects that conform to the HDF5 data model.
 - REST VOL to access data via HSDS in AWS and Azure Clouds (also, on prem. solutions like OpenIO)
 - DAOS VOL
 - VOLs under development (AIO, caching, log, ADIOS, PnetCDF)
 - HDF5 external references
 - Support attributes, and object and dataset selections that reside in another HDF5 file.
- Performance improvements
 - The hyperslab selection code was improved by an order of magnitude
 - Optimizations for parallel HDF5 applications released in 1.10.0 - 1.10.6 releases
 - Collective metadata storm reads/writes
 - Optimizations for open/close/flush
 - Reading of the same dataset by all the process collectively

HDF5 1.10.7 Release

- Performance improvements for parallel applications
- Split and Mirror VFDs were added
- Internal memory sanity checking was disabled in Autotools debug builds by default.
 - The sanity checking replaced C API malloc(3) and free(3) calls with HDF5 internal calls, which caused problems with the external filter plugins.
- Updated h5repack to follow external links and merge data
- Enhanced file locking control
 - Environment variable (highest)
 - APIs H5Pset(get)_file_locking()
 - Configure/CMake option (sets property default)
 - Library defaults (currently best-effort – Use file locking and ignore “disabled” errors)

HDF5 1.8.22 Release



- Coming before 12/15/2020
- Maintenance release with
 - CMake improvements
 - Fixes for CVEs and bugs
 - S3 and HDFS VFDs
 - For more info, see https://github.com/HDFGroup/hdf5/blob/hdf5_1_8/release_docs/RELEASE.txt
- **We are dropping support for 1.8.* in 2021**
 - Exact day will be announced on HDF FORUM

Miscellaneous information



- CMake minimum supported version is 3.12, if possible, move to 3.15
- We dropped support for Windows 7 and Visual Studio 2010, 2012, and 2013 after HDF5 1.10.7
- Pre-compiled binaries are now available from The HDF Group Website and from our ftp server <https://support.hdfgroup.org/ftp/HDF5/releases/>
- For each release we provide HDF5 compression plugins (encoding and decoding) along with the pre-compiled HDF5 binaries for Windows, macOS and Linux
 - BLOSC, Bit-shuffle, BZIP2, LZ4, LZF
- Plugin source will be in GitHub soon (stay tuned for the announcement)

Outreach



- This is our second HDF User Group meeting!
- Webinars and Blogs
 - <https://www.hdfgroup.org/category/webinar/>
 - <https://www.hdfgroup.org/blog/>
 - Please let us know if you are interested to present your work
 - Contact Lori Cooper lori.cooper@hdfgroup.org
- Tutorials
 - Tell us which topics you would like us to present in 2021
 - Contact help@hdfgroup.org or post on [FORUM](#) or contact Lori Cooper lori.cooper@hdfgroup.org

Other developments

- HDF5 is on GitHub <https://github.com/HDFGroup/hdf5> now
 - [HDF5 repo in Bitbucket](#) is not updated anymore!
 - Check doc directory in the source for the materials on the work flow, coding standards, etc.
- With the move to GitHub The HDF Group:
 - Adopted C coding style based on clang format
 - Enforced using Actions on GitHub
 - Provided formatter tool in the bin directory of the source code distribution
- HDF5 Documentation
 - We worked on addressing known problems
 - Search engines don't index HDF5 Docs served via Atlassian Confluence
 - Access to HDF [support portal](#) and especially to HDF5 RM is slow
 - Hard to contribute
 - Decision was made to revamp HDF5 RM in [Doxygen](#) (will be released in HDF5 1.12.1)
 - We worked on [prototyping](#) future HDF documentation solution

One Stop HDF5 Documentation



Browser tabs: v10.3.20 Time, HDF5 Users Group, Mail - Elena Pourn, Index of ftp://gam, hdf5 - master - AT, Tutorial/Parallel-h, HUG-Parallel-han, Parallel HDF5 Test, HDF5 Wiki

Address bar: https://hdf5.wiki/index.php/Main_Page

Search: binary file dump

Most Visited: Batter up: It's Natio..., New Tab Search, Calendar - Elena Po..., HDF5 SWMR Docu..., HDF Log In - Confluence, FDA-approved diet ..., A Report of Progres..., Getting Started, The New York Time..., HDF Home Page, ECU Login

Create account Log in

Main page Discussion Read View source View history Search HDF5 Wiki

Welcome to One Stop HDF5 Documentation!

Main Page

TASKS 🧑	CONCEPTS 🧠	REFERENCE 📖
HDF5 Beginner New to HDF5? Start here! Cookbook (CB) Quick and easy recipes anyone can cook User Guides (UG) Tutorials on the most noteworthy features	Glossary Brief explanations of terms used throughout (Hyperspace?) Request for Comments (RFC) Methods, behaviors, research, or innovations Technical Notes (TN) Specific developments, techniques, procedures, or modifications	Reference Manual (RM) The nuts and bolts of HDF5. All of them! Release Information (REL) The poetry of the release cycle Specifications (SPEC) HDF5 is made exactly to it. File Format (FMT) No bit gets left behind.

Community Corner (CC)
Your Project Here

HDF5 Library Development (DEV)
Jump in on the deep end!

Other formats and versions of the HDF5 documentation can be found here (maintained) and here (unmaintained).

👉 Discuss RM requirements! 👈

News:

2020-07-28
A Bitbucket branch for the Doxygen-based RM was created. Periodic snapshots can be found here.

2020-07-24
A slightly refined Doxygen version is available

2020-07-16
Barbara is working on release information

2020-07-15
Werner is working on a Doxygen example

HDF5 Roadmap

2021

Releases



- Back to fixed releases schedule
 - HDF5 1.10.* releases in May and November
 - HDF5 1.12.* releases in June and December (exception, HDF5 1.12.1 will be released in January 2021)
- New features
 - VFD SWMR
 - Addresses deficiency of current SWMR implementation (e.g., allow to add or delete groups, dataset and attributes, support for VL types)
 - HDF5 Versioning VFD (aka "Onion" VFD)
 - Supports version control for file open/close session and provenance management
 - Parallel library enhancement – Sub-filing
 - Middle ground solution between single shared file and one file per process approach to improve I/O
 - Uses VFD approach
 - HDF5 file is striped across collections of sub-files
 - I/O requests are routed to I/O concentrators that access sub-files

Major focus

- Minimize HDF5 Performance gap between 1.8 and later releases
- Improve HDF5 documentation
 - Finish HDF5 RM conversion to Doxygen
 - Move on “One stop” documentation setup
 - Involve community in documentation development and maintenance
 - Publish documentation on HDF5 code standard and best practices to facilitate community contributions
- Address issues in HDF5 VOL architecture and productize HDF5 VOLs
 - Decrease VOL architecture overhead
 - Productize
 - THG supported VOLS (REST, DAOS, RADOS)
 - ECP VOLs as they mature
 - Implement REST VOL to directly access HDF5 Cloud-optimized files
- Continue outreach efforts

Community engagement



- Engage community to work toward
 - Multi-Threaded HDF5 library
 - Have a “proof of concept” limited implementation (e.g., reading contiguous and chunked datasets)
 - Develop a plan for full implementation
 - Attend the next talk and technical Webinar is planned for October 30 at 11:00 am Central
 - Full support for UTF-8 encoding and working with Unicode filenames
 - Support for 128 and 16-bit floating point numbers, Boolean and complex datatypes
- Open HDF5 compression filters repo maintained by THG
 - Repo contains HDF5 filters, compression libraries, and examples; source is built with CMake
 - We plan to move the project to GitHub, add to Spack builds and contribute back to the original filters repos.
 - Work with the h5py and other communities to coordinate distribution efforts?

The HDF Group wish list

- Address file corruption issues
 - Avoid corrupting file during catastrophic event (full disk, CTRL C, system crash)
- Long-standing improvement for parallel HDF5 including data aggregation
- Support for new storage
 - Column
 - Sparse
- Revamp chunk cache
- Performing VL types and data streaming into HDF5
- And the list continues.....
- What **you** would like to see?

HDF JIRA vs GitHub



- Register at The HDF Group Website to get access to JIRA
- Browse through the open issues and vote if the issue is important to you
- We will use GitHub for new issues and may bring some JIRA issues to GitHub.

THANK YOU!

Questions & Comments?

Questions? Suggestions?

Your turn now 😊