HDF5 in the ECP-EQSIM Project

Houjun Tang Lawrence Berkeley National Lab





The <a>EarthQuake <a>SIMulation (EQSIM)

awrence

aboratory

lational

BERKELEY LAB



HDF5 integration in EQSIM workflow

BERKELEY LA



EXASCALE COMPUTIN PROJECT

Effective data management is critical

- 1 simulation run can generate 200+ TB data (uncompressed).
- Effective (lossy) compression is essential for data management.



Significant amount of ground motion data





FOM run on 3600 GPU-accelerated Summit nodes

- 391 billion total grid points, 381899 simulation steps (90 seconds).
- Top surface: 68577x45729, ~3 billion points, 1.75m horizontal grid size.
- Down-sampled every 100 steps, 3819 steps written to GPFS.
- ZFP accuracy mode, 1e-2





Lessons learned & best practice

- Custom-made formats are not ideal when data needs to be shared to many people with different backgrounds.
- HDF5 format is useful and effective for managing large datasets.
 - Cross-platform, multi-language support (C, Python, MATLAB).
 - Self-describing, stores metadata together with data.
- HDF5 library can be efficiently used for parallel I/O
 - Requires knowledge of tuning parameters selection on different systems.
 - Compression filters can be enabled with a few lines of code.







Thanks!

Questions?



