

MarFS Metadata Scaling

PDSW WIP Report 2016

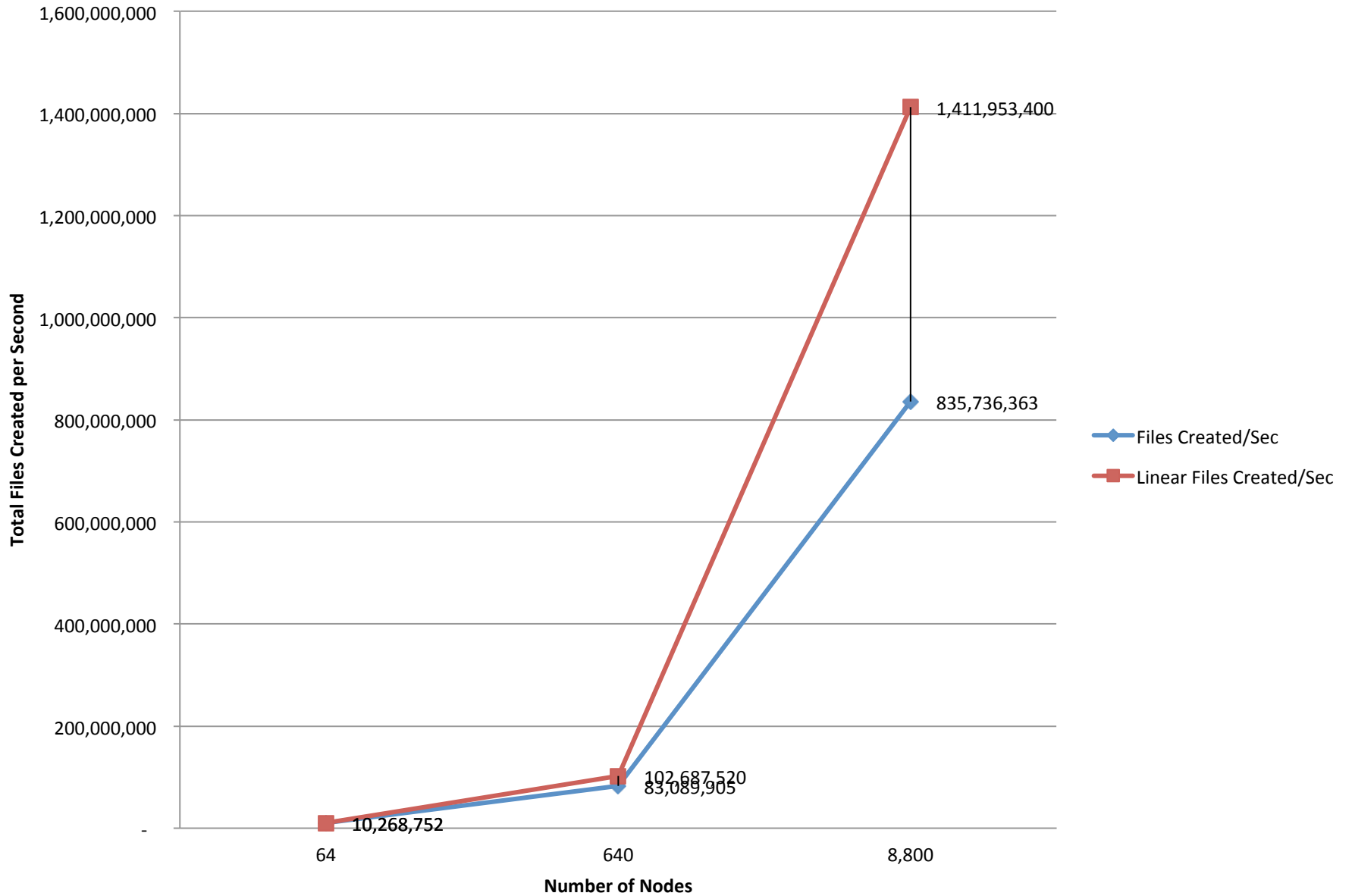
David Bonnie, Hsing-Bung Chen, Gary Grider, Jeffrey Inman, Brett
Kettering, William Vining

LA-UR 16-28615

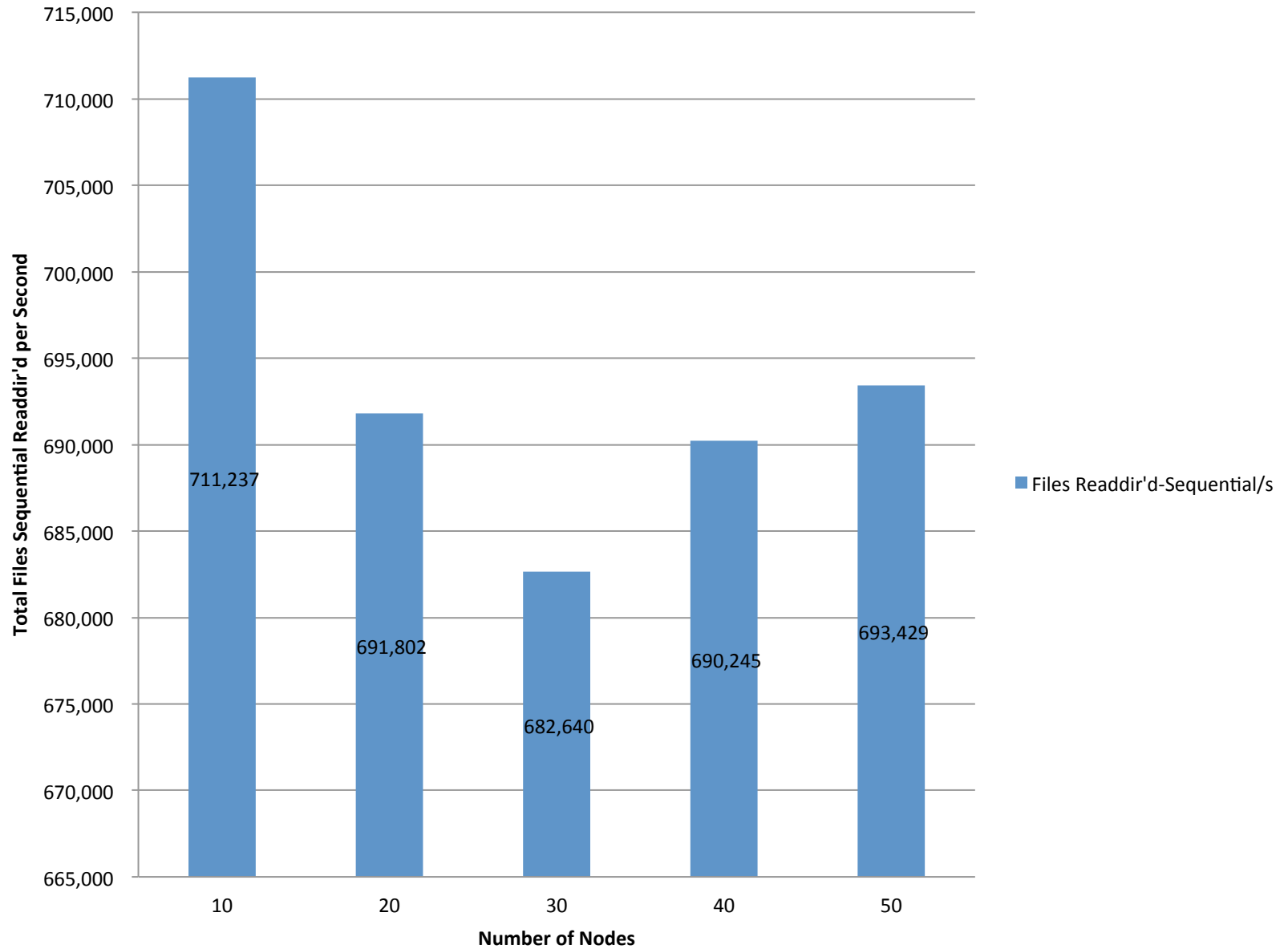
Metadata scaling components

- Deploy one drMDS per file system as rank 1 on first node
 - Make new directories & broadcast dir inode to fdMDSc's
- Deploy fsMDSc's on $\frac{1}{4}$ cores for each node in file system service
 - Handles its sharded part of distributed file metadata when broadcast commands are sent
- Deploy fsMDSp's on $\frac{1}{4}$ cores for each node in file system service
 - Handles its sharded part of distributed file metadata when command are sent to a specific fsMDSp.
- Deploy file system Clients on $\frac{1}{2}$ cores for each node in file system service
 - Execute file system operations, such as create

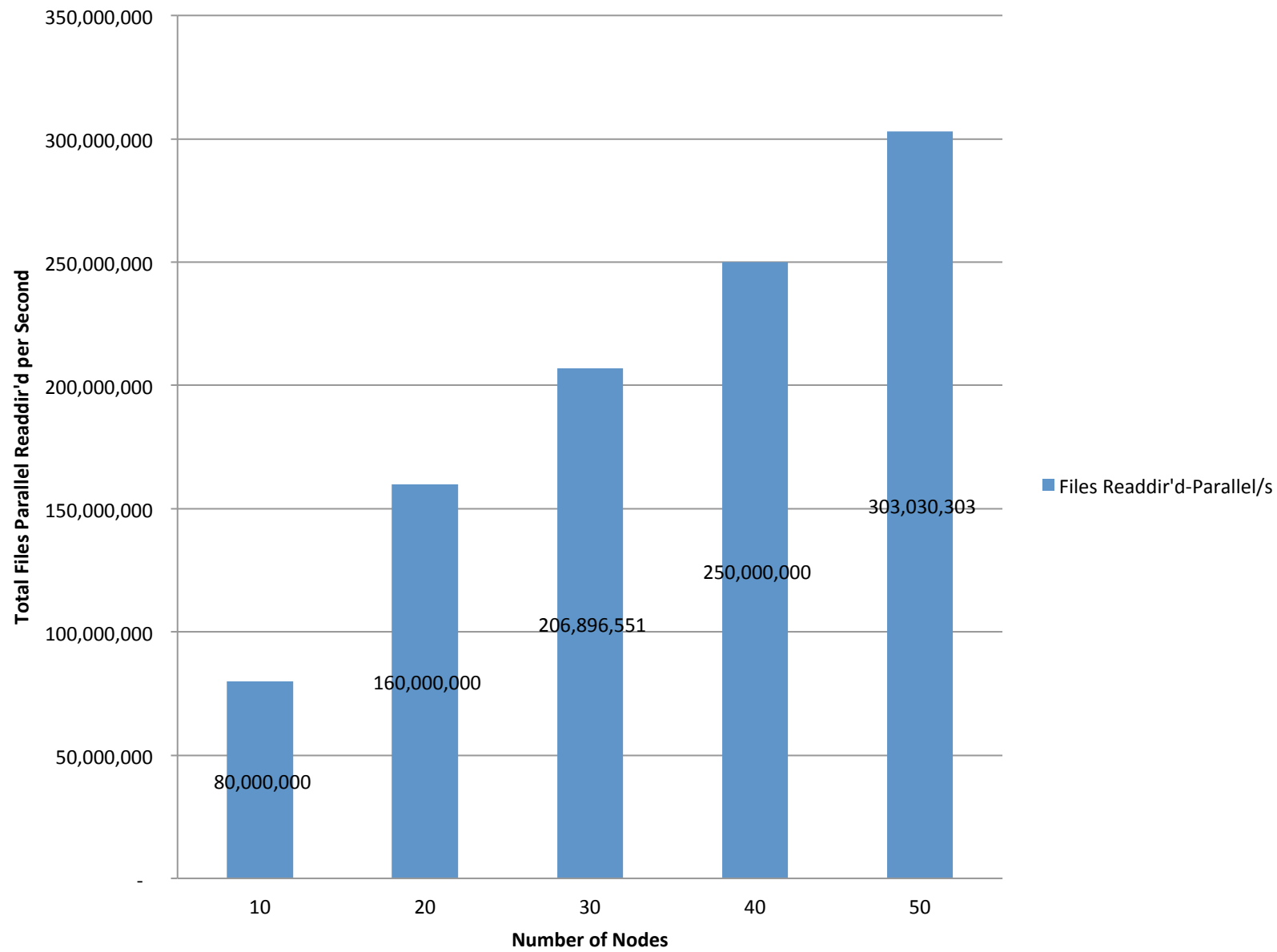
File Creation Rate by Node



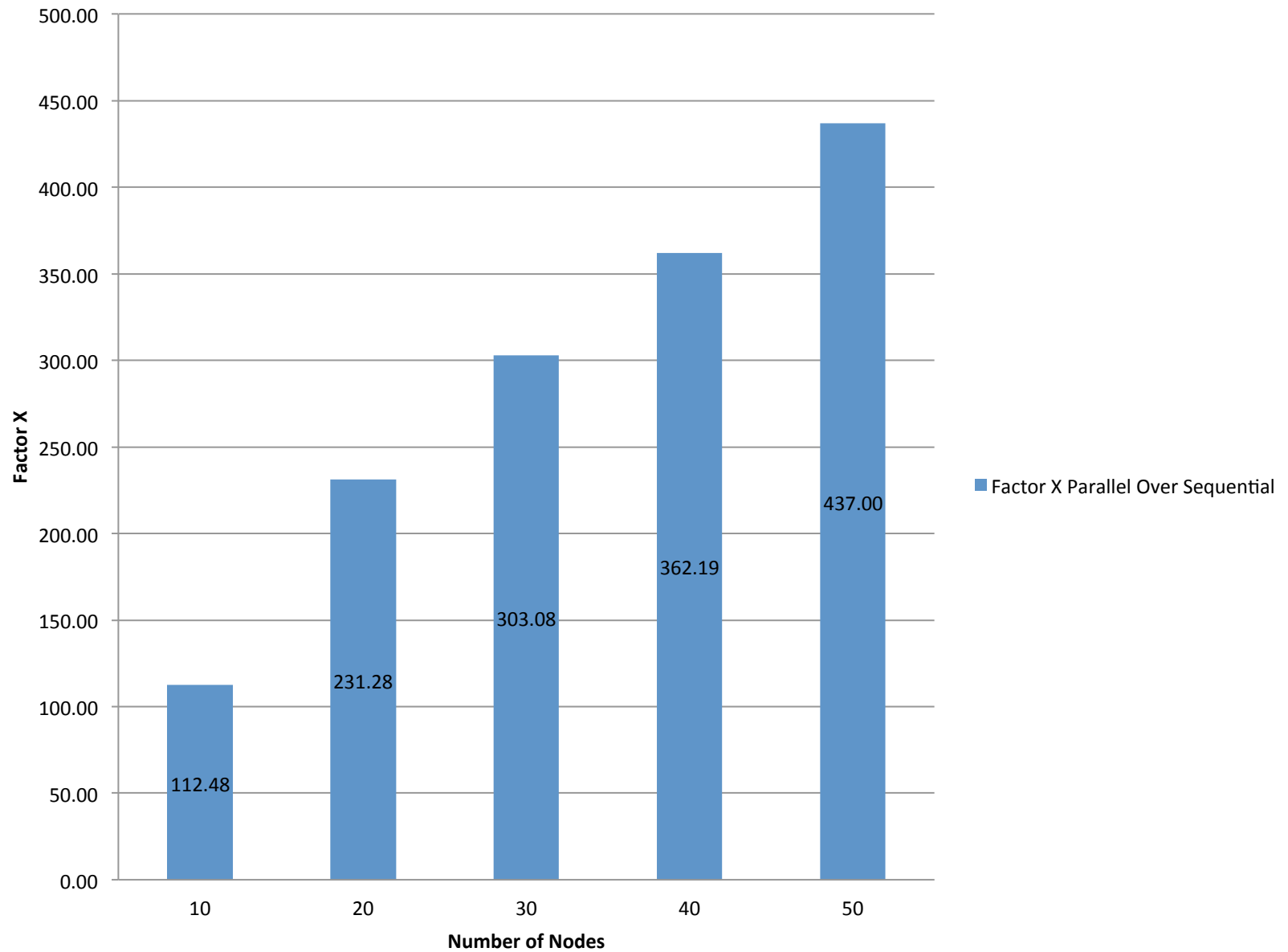
File Sequential Readdir Rate by Node



File Parallel Readdir Rate by Node



Factor of X that Parallel Readdir Rate is Greater than Sequential

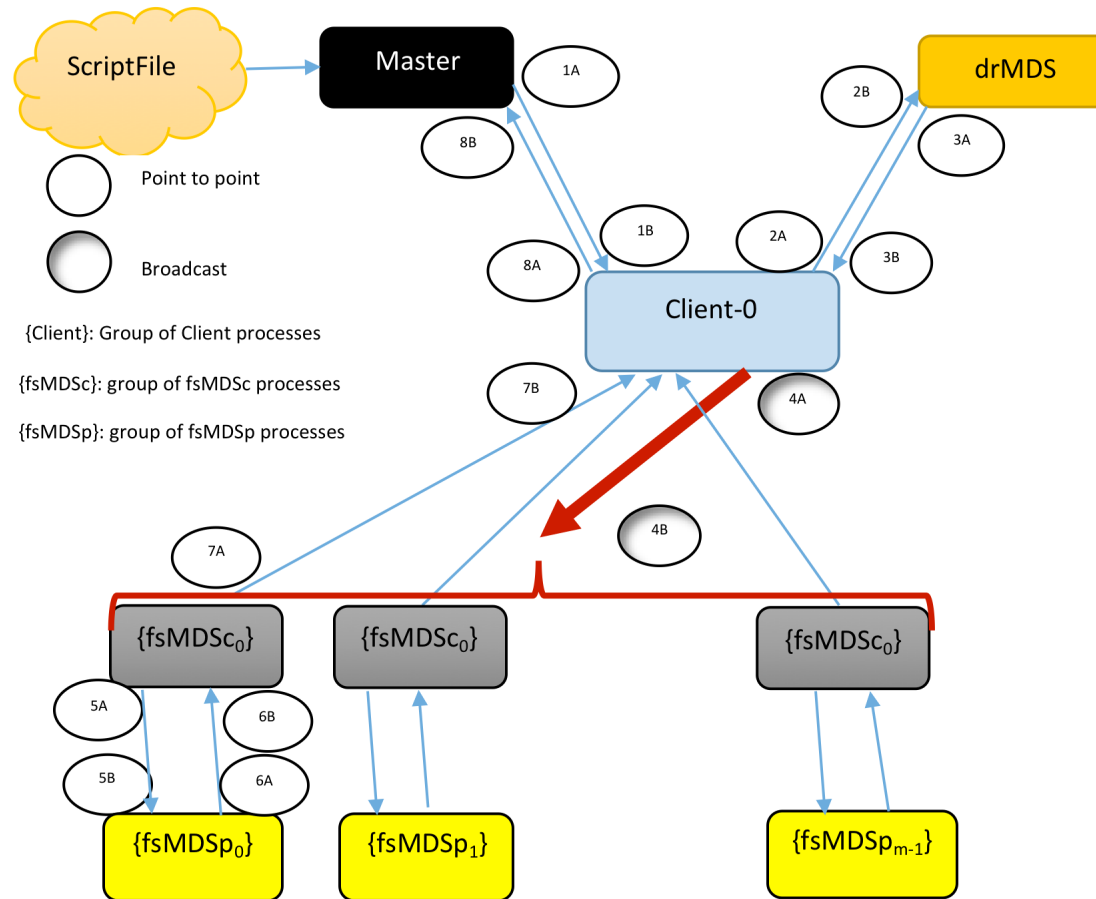


Background Information

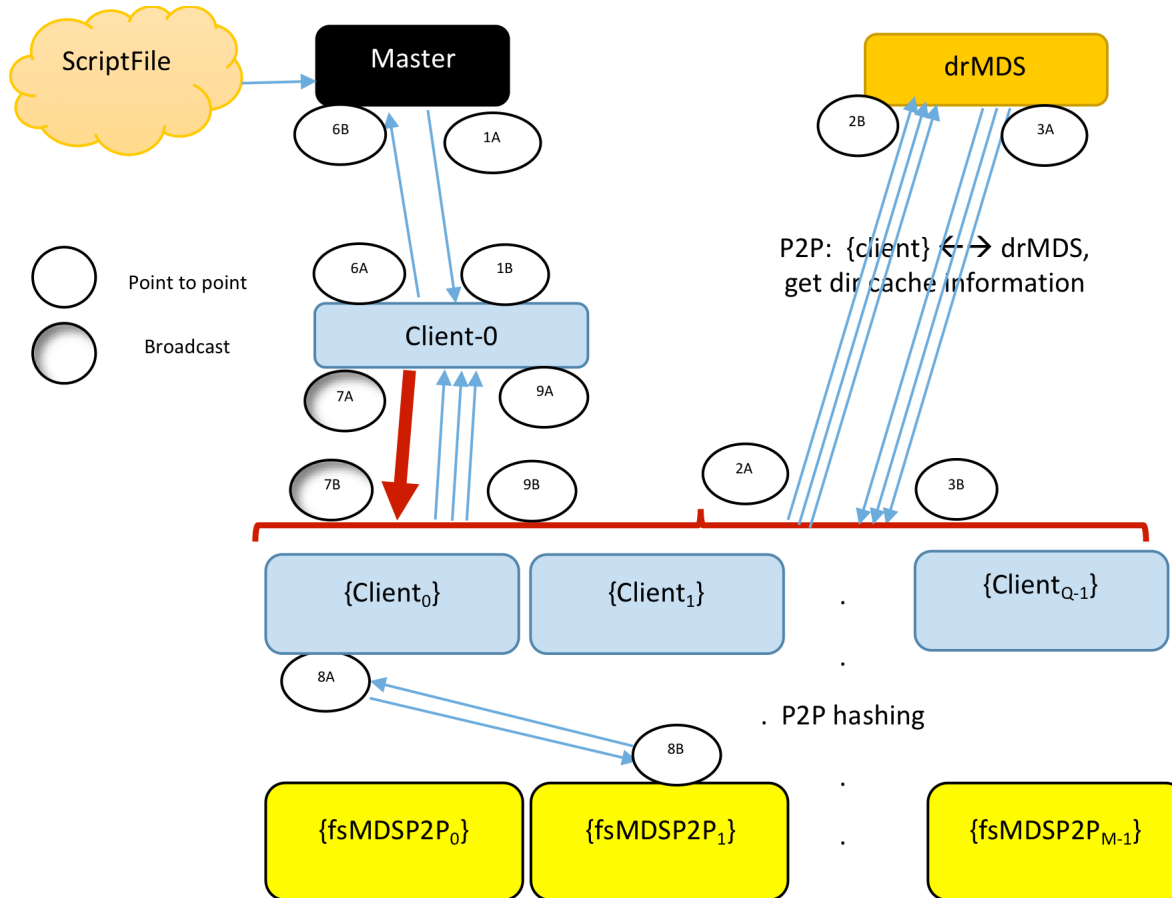
MARFS METADATA SCALING

MarFS Overview

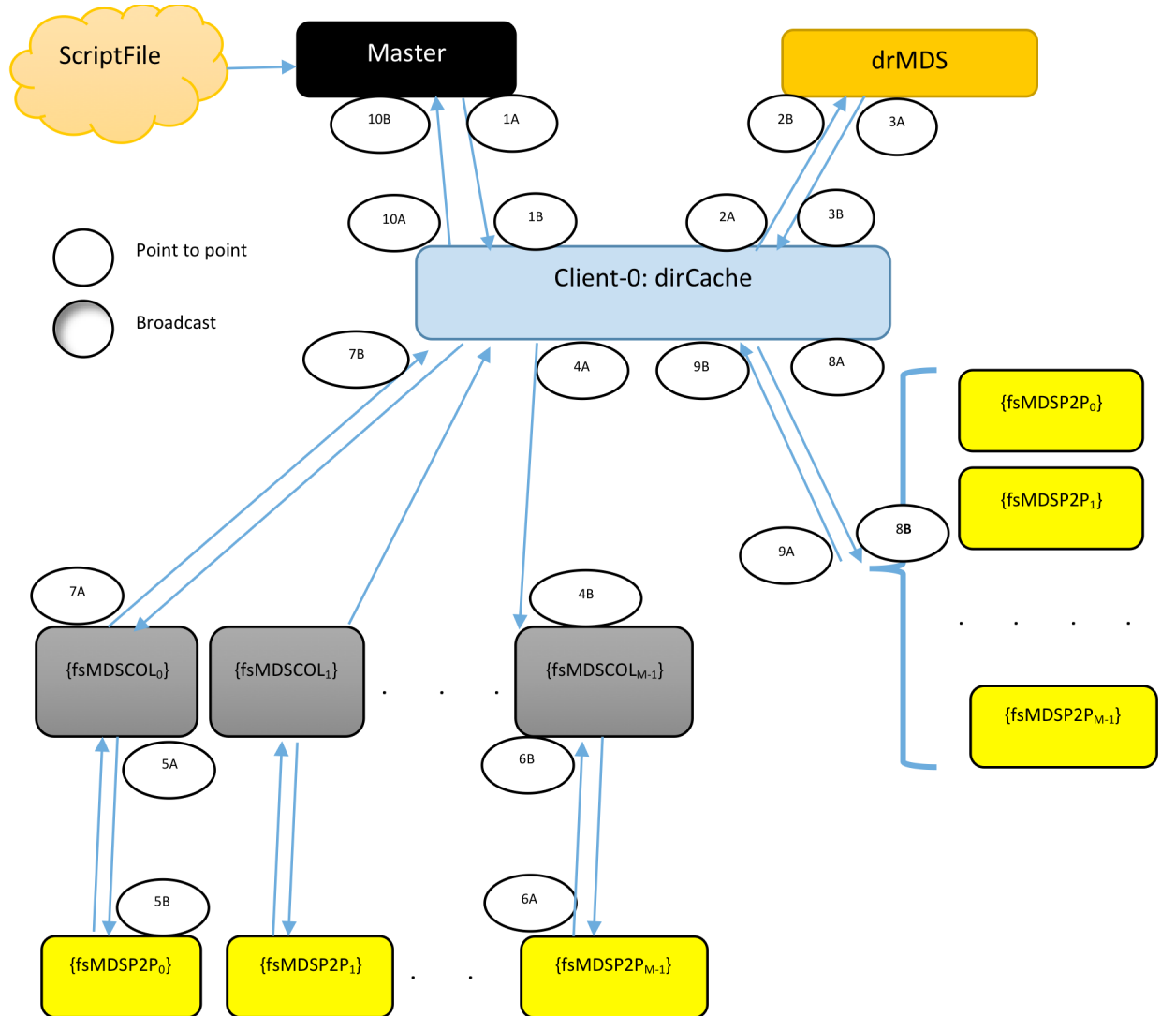
- Provides near-POSIX over cloud-style erasure and objects
 - Yields reliable storage on inexpensive disk
 - Supports legacy apps' files/folders/ownership/etc.
- Store large data sets for weeks to months on PFS, 1 TB/s
- Store data forever in archive, 10s GB/s
- Store large data sets for months to year'ish on MarFS, 100s GB/s
 - Data set $O(\text{PB})$, aggregate data $O(\text{EB})$
- Systems growing from $O(\text{M})$ cores/ $O(\text{PB})$ memory to $O(\text{B})$ cores/ $O(10\text{s PB})$ memory
 - Going to $O(\text{B})$ files per job in one directory and $O(10\text{s T})$ files per file system



Here's a picture of creating files



Here's a picture of sequential readdir



Here's a picture of parallel readdir

