

Scalable Distributed Directory Implementation on Orange File System

Shuangyang Yang
Walter B. Ligon III
Elaine C. Quarles
Clemson University

Outline

- Orange File System
- Motivation for Distributed Directory
- Design and Implementation
- Micro-benchmark Performance
- Summary and Conclusion
- Future Work

What is OrangeFS?

- High performance, parallel file system
- Continuation of PVFS
- Concentration on production quality features
- Open source, GPL licensed
- More details at <http://orangeefs.org/>



Huge Directories

- Millions of files under the same directory
- Scalability issues
- Needed for:
 - Gene sequencing
 - Image processing
 - Data mining
 - Real-time application monitoring

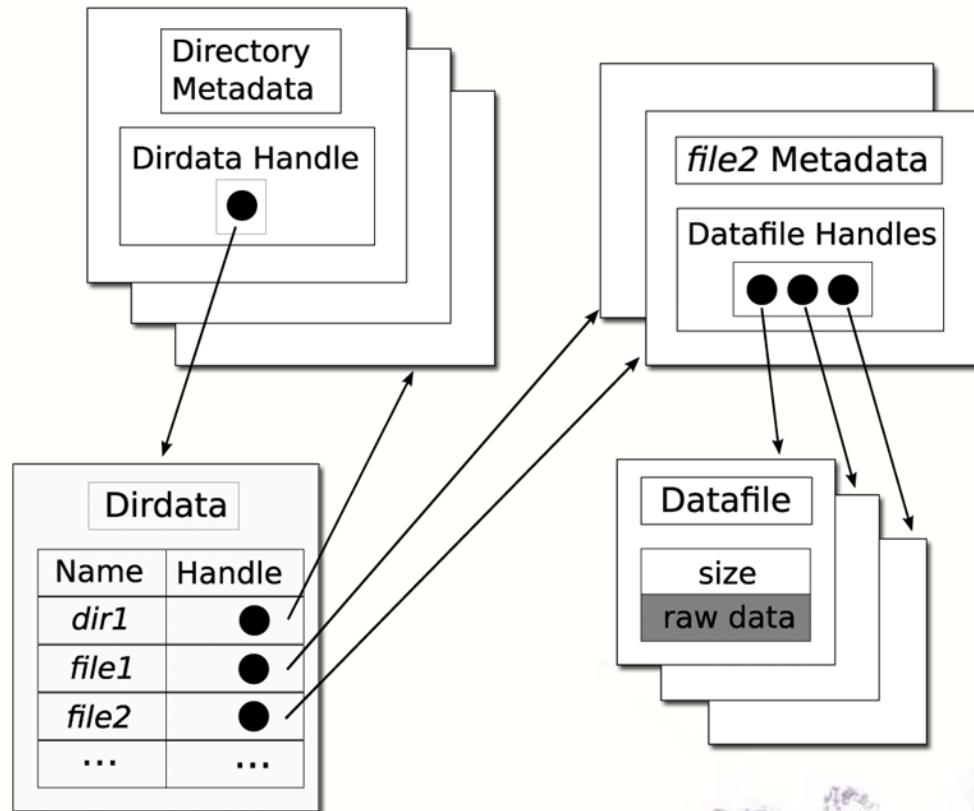
Motivation for Distributed Directories

- OrangeFS stripes contents of data files among servers, BUT
- All directory entries for any given directory are held on a single server, THUS
- Creating the potential for hot spots

OrangeFS Structure

Terminology:

- Metadata
- Datafile
- Dirdata



Existing Techniques

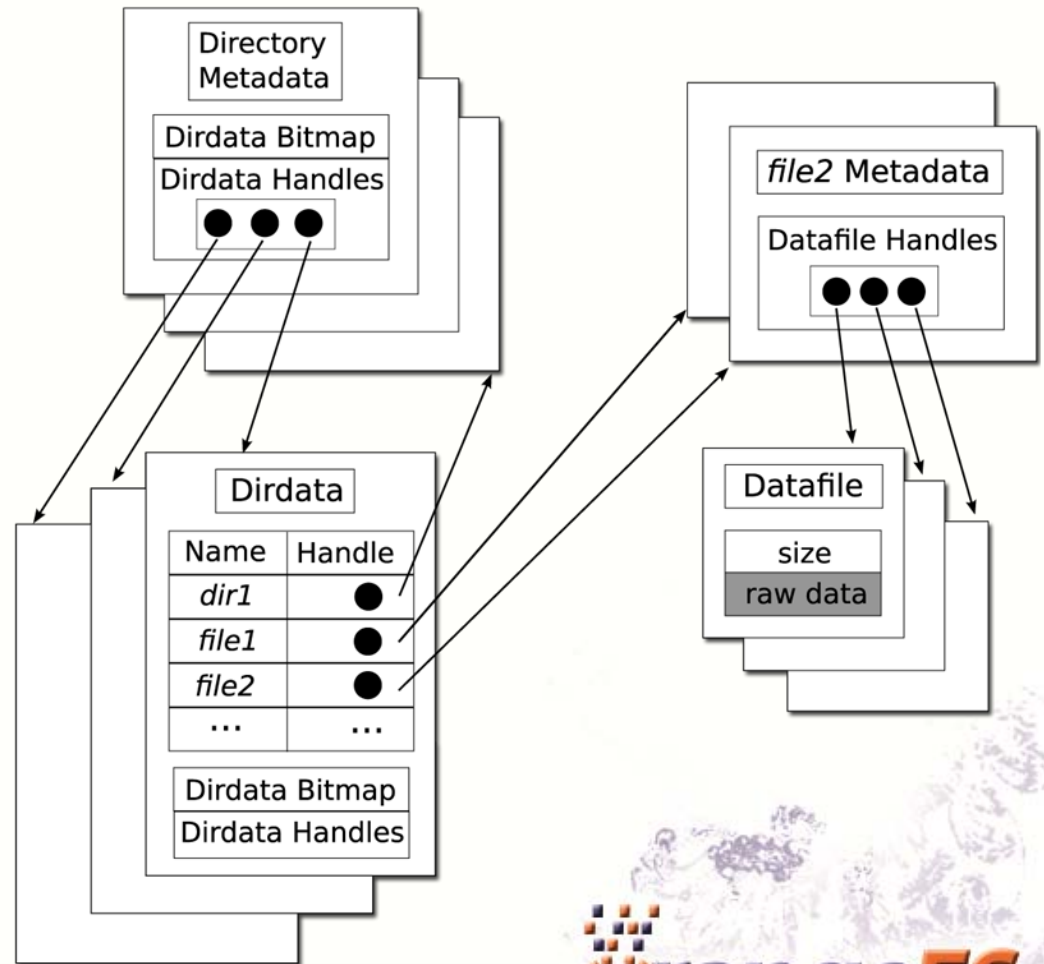
- Distribute directory entries in multiple blocks or partitions on multiple servers.
- GPFS: multiple disk blocks. Extensible hashing to lookup and grow.
- GIGA+: multiple fixed-size partitions. Partition bitmap to maintain status. Extensible hashing and incremental growth.

Goals of Current Work

- Distributed directory implementation based on GIGA+. Utilizing extensible hashing and partition bitmap representation.
- Different design decisions in order to be integrated into OrangeFS seamlessly.
- Evaluate performance of initial implementation

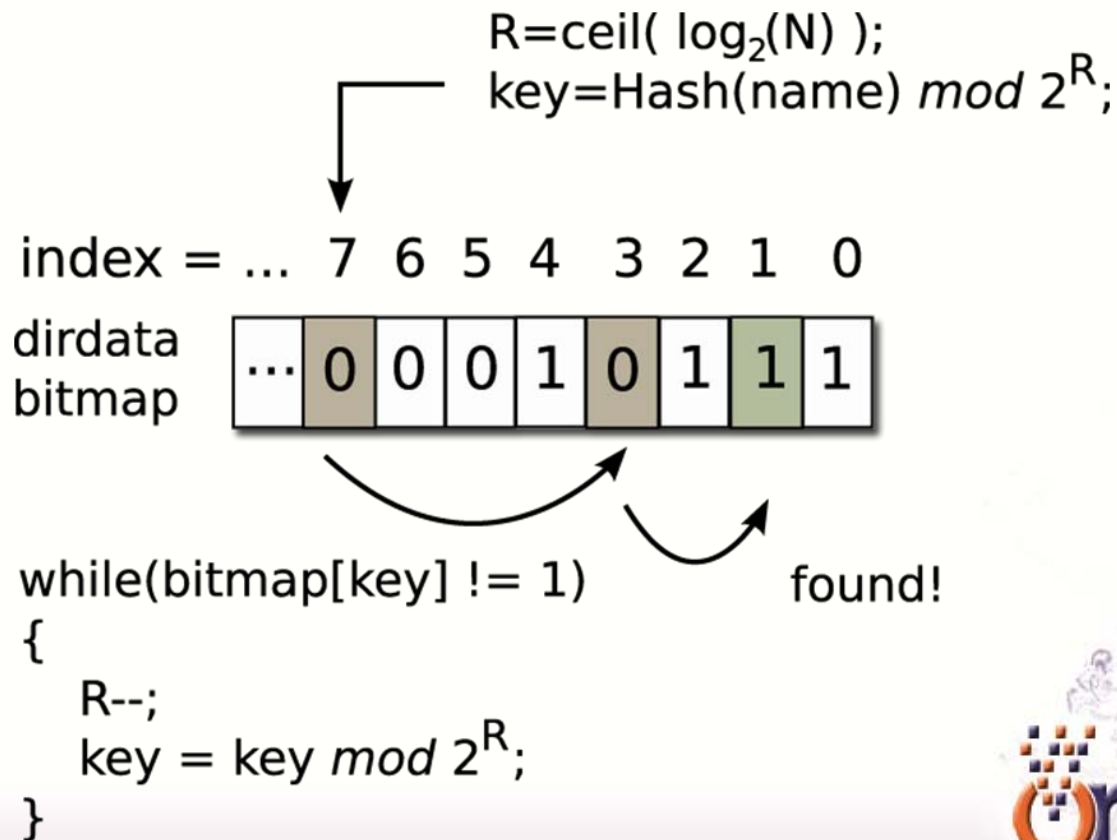
Modified Structure

- Multiple dirdata objects represented by handles.
- Dirdata bitmap
- Stored on both metadata object and dirdata object.



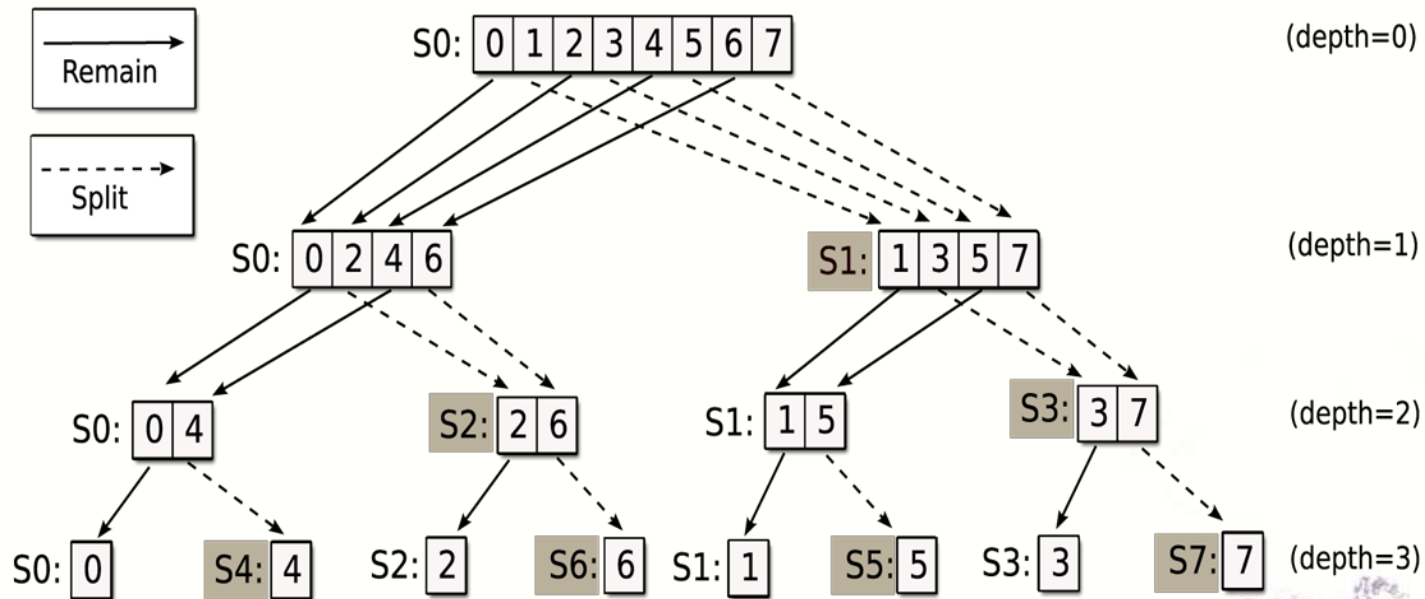
Lookup

- Map a directory entry to a dirdata object.



Split Operation

- Overloaded dirdata object can send half of its hash space to another dirdata object.



Design Decisions

- One dirdata object on one metadata server, entries indexed by Berkeley DB
- Initial number of active dirdata objects is configurable
- Metadata object holds most up-to-date copy of dirdata bitmap

Implementation on OrangeFS

- Directory entries are distributed among multiple metadata servers.
- Dynamic splitting is close to completion.
- The scalable distributed directory feature is available as an experimental release at

<http://www.orangeFS.org/download/>

Micro-benchmark performances

- Modified version of UCAR metarates benchmark*.
- Measure throughput of file creation and removal under one directory concurrently.
- Conducted on the Palmetto Cluster housed by Clemson University.

* <http://www.cisl.ucar.edu/css/software/metarates/>

Palmetto Cluster

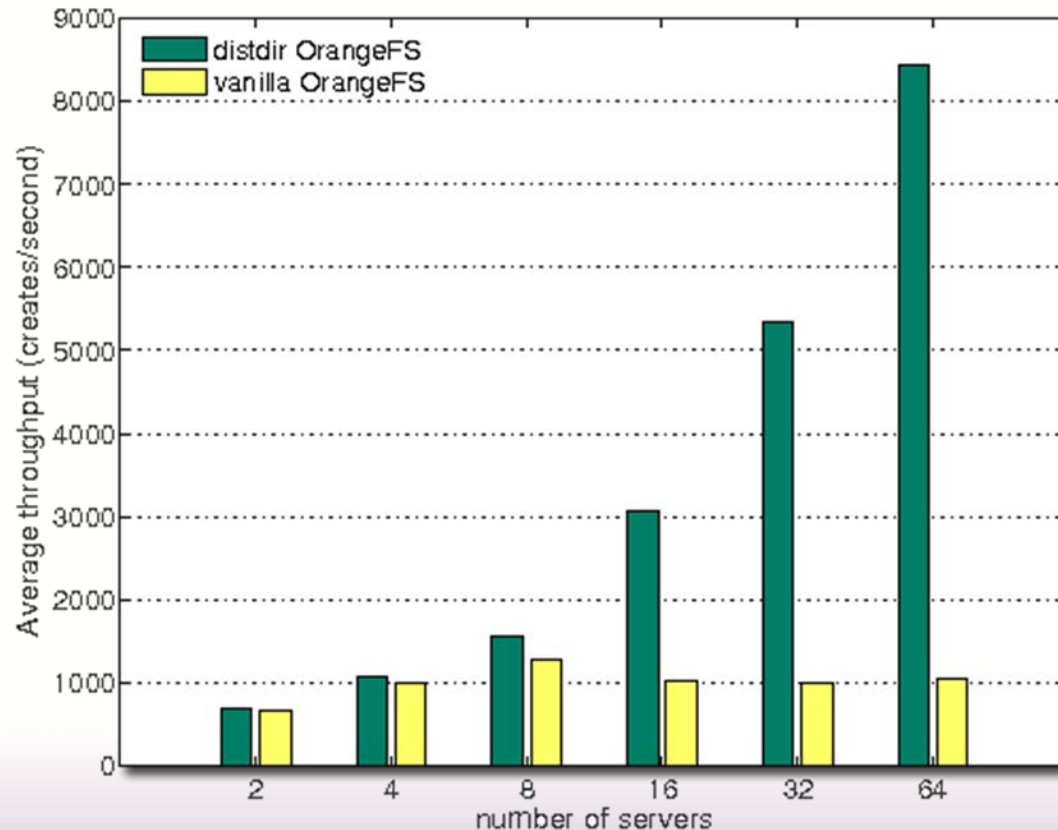
- 1541 nodes, 8 cores.
- 12/16 GB memory.
- Myrinet 10G interconnect.
- High throughput storage.
- More at

<http://citi.clemson.edu/palmetto>



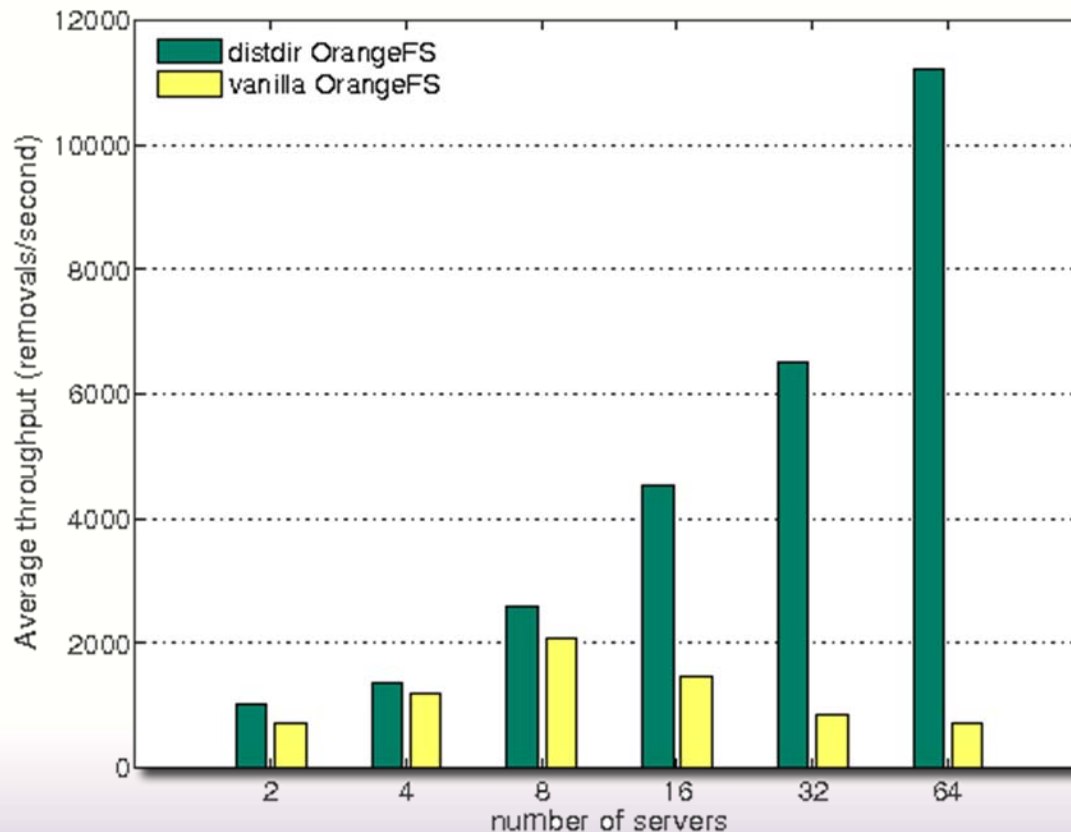
File Creation

- Average throughput can reach 8000+ creations/second with 64 servers, 128 clients.



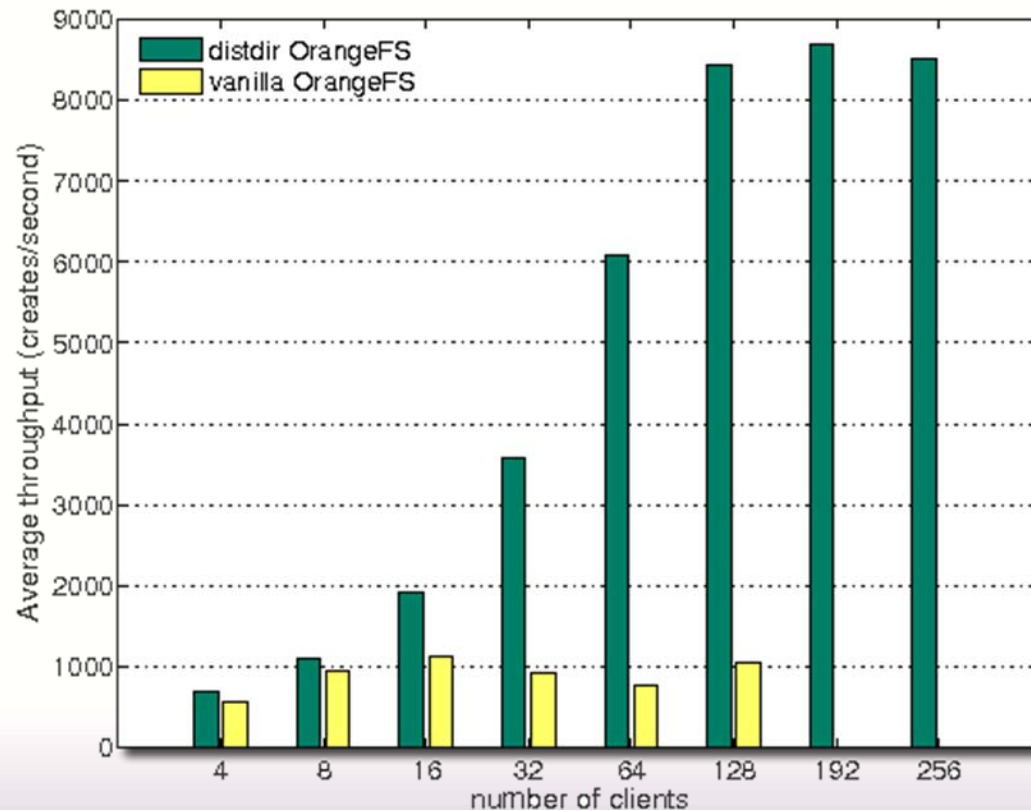
File Removal

- Average throughput can reach 11000+ removals/second with 64 servers, 128 clients.



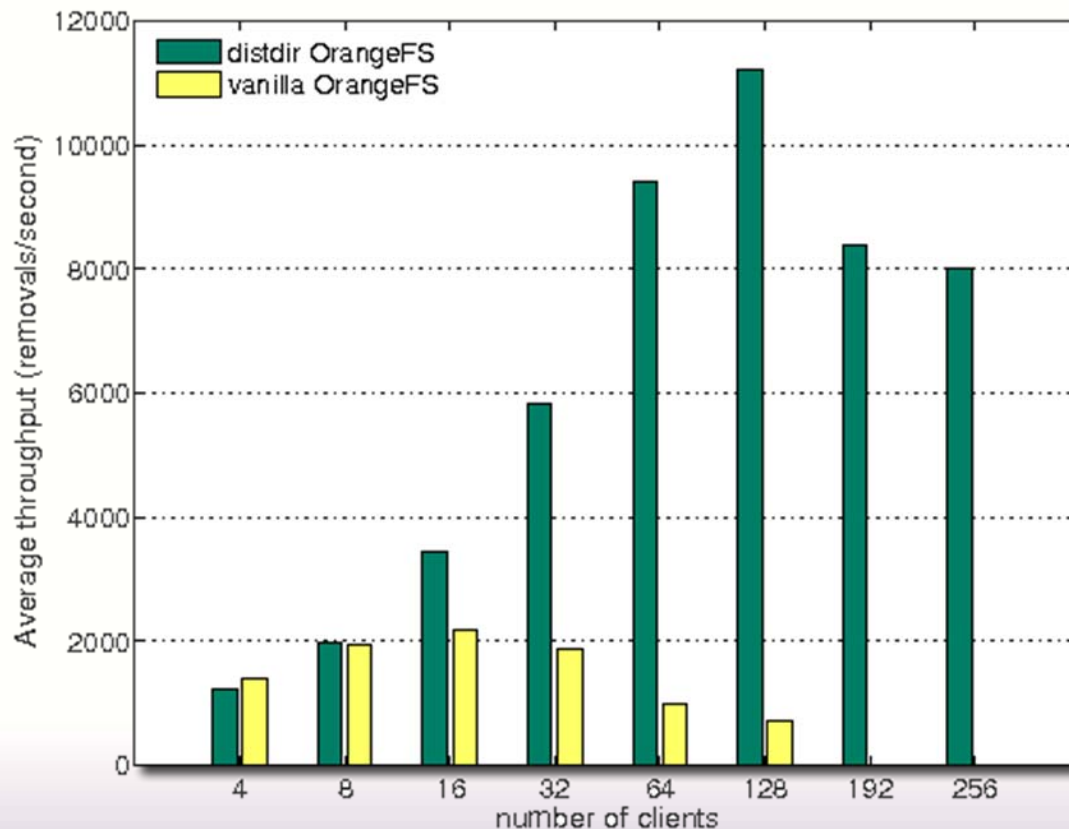
Creation with Variable Clients

- Maintain high throughput when saturated on 64 servers.



Removal with Variable Clients

- Maintain high throughput when saturated on 64 servers.



Summary and Conclusion

- The distributed directory implementation shows great scalability in creating and removing large numbers of files by multiple clients concurrently.
- Can reach 8000+ file creations/second and 11000+ file removals/second with 64 servers on Palmetto Cluster.

Future Work

- Finish implementation of splitting functionality and get ready for public release.
- Optimize performance, particularly using collective communication.
- Thoroughly evaluate performance of throughput, scalability, overheads and etc.

Acknowledgement

