

MD Performance Dashboard

MetaData Performance Project – Dashboard May 11, 2009

Release	Task	Bug	Status	Updates
1.6.x Branch – Resolve any regressions found and bring back to original pre-1.6 branch performance – COMPLETE				
	1.6.5 mdsrate performance regressions from 1.4.11 reported	18577	Completed	Patch (precreate optimizations) has been landed in 1.6. & 1.8
	Create rates improved significantly after reformatting the file system		Closed	Improved from 1500 to 5000 IOPS/s (close to average for a single MDS) after the reformat. Not enough conclusive data to file a bug.
1.8.x Branch – Double the current MD iop/s performance.				
	Unnecessary lookup RPC when doing non-open related creates	18534	In Progress	<ul style="list-style-type: none"> - Proof of concept patch created - Related kernel & client side patches being investigated - Oleg to run tests on 3/20 to test patch
	Asynchronous Close for Files	19017	In Progress	- This work will help improve on mdtest create tests, which do an open + close test
	Remove unnecessary fysnc's from mdtest test scripts		In Progress	<ul style="list-style-type: none"> - Recommendation is to run mdtest without the -y parameter. - Tests show 5x improvements, single client creates go up from 500/s to 2500/s
	Patches to reduce CPU boundedness for MDS- smp scalability improvements		In Progress	- Continue work on some patches Oleg had developed earlier.

MD Performance Dashboard

	LNET SMP scaling	15379	In Progress	<ul style="list-style-type: none"> - Prototype in Progress. test code expected Jun 09. - LNET self test runs on a 16 core SMP resulted in order of magnitude improvements for small messages(700,000 rpc's/s). - ptlrpc scaling work and test on actual mds (5/30)
2.0 Branch – Meet or Exceed GPFS Metadata Performance with Lustre CMD and Other Improvements				
	Size on MDS (SOM)		Awaiting Performance testing	<ul style="list-style-type: none"> - Preview version in 2.x - Productionized version with recovery and interoperability in 3.0
	Large readdir bulk RPCs for readdir efficiency		Not Started	
	Simple MD Write Back Cache		Not Started	- Can be used instead of the full blown WBC
	Batching of multiple separate metadata request into a single RPC		Not Started	- Dependant on MD-WBC
	Readdir+ system calls		Not Started	
	Create on Write (CROW)		Not Started	
	Changes to store the FID inside the directory		Not Started	
	Allow single clients to have multiple metadata modifying RPCs in flight		Not Started	
	pdirops (A locking protocol introduced in the VFS to allow for concurrent operations on a single directory inode)		Awaiting Performance testing	<ul style="list-style-type: none"> - Implemented in HEAD - Needed for parallel file creation rates in the same directory
3.x Timeframe				
	Clustered Meta Data		In Design	<ul style="list-style-type: none"> - Preview version of CMD to be available in 2.0 - All improvements for single server will benefit CMD - Linear scaling studies in progress.
	Meta Data Write Back Cache		In Design	- Design in progress.

MD Performance Dashboard

Performance Testing				
	- MD Baselining Test Design	18776	Completed	- Design created, reviewed by Senior engineers - Can use ramdisks for MDT's to simulate fast hw.
	MD Baseline Runs (1.6, 1.8 and 2.0)		Completed, Internal Reviews	- Baseline tests completed on a Sun Thor machine - Internal reviews in progress.
	MD Performance Runs(1.6, 1.8, 2.0)		In Progress	- Publish weekly metrics on MD improvements in metrics report.
	MD Baseline Run on non-Cray HW		Not Started	- Run baseline tests on Hyperion Cluster
Other Related				
	- Longhaul patch MD performance for WAN	18526	Not Started	- Increase number of in-flight RPC's and implement multi-slot transactional updates for state change operations. Read-only should be relatively easy, but the multi-slot tx updates is fairly complex and only slated for 2.0

Status Field - Key

B	Project Completed
G	Project on Track -- Quality/Requirements being met and within bounds
Y	Plan needs Close Monitoring -- Quality/Reqmnts at risk of being met or compromised
R	Recovery Plan Required -- Quality/Requirements will not be met or are compromised.
V	Product Enhancement