



---

# AFM Migration: The Road To Perdition

Spectrum Scale Users Group – UK Meeting  
9<sup>th</sup>-10<sup>th</sup> May 2017

Mark Roberts (AWE)  
Laurence Horrocks-Barlow (OCF)



# GPFS Systems

- Legacy System
  - GPFS 3.5 (x86 m35), 40GbE attached
  - Home (DDN SFA12K-40) 10 servers per site
    - 4 MB block size, ~1.8PiB, 100+ mil resident files, ~1.5 PiB migrated (HSM) files
  - Group (IBM DCS9900) 2 servers per site
    - 1 MB block size, ~650TiB in ~20 mil resident files
- New System
  - ESS Power8 (GL6), dual 40GbE
  - 6TB SAS disks (data), Metadata (GL4)
  - All disks are FIPs140-2 encrypted
  - Home (4 PB), 16 MB block size
  - Group (2 PB), 8 MB block
- Both systems are NSD replicated.....Not AFM-DR



# Migration Options

- Requirements for migration
  - Minimum downtime (24\*7 production class)
  - Gain benefits of new 4.2.1 (ESS) file system features
- Possible methods
  - NSD migration - does not meet requirements due to restrictions
  - Active File Management (AFM)
    - Utilise native GPFS transfer
    - IBM support
    - Complicated compared to rsync
  - rsync
    - Need patches to support GPFS (officially supported??)
    - Slow
  - TSM dump & restore.....enough said !

**AFM chosen for speed despite complexity**

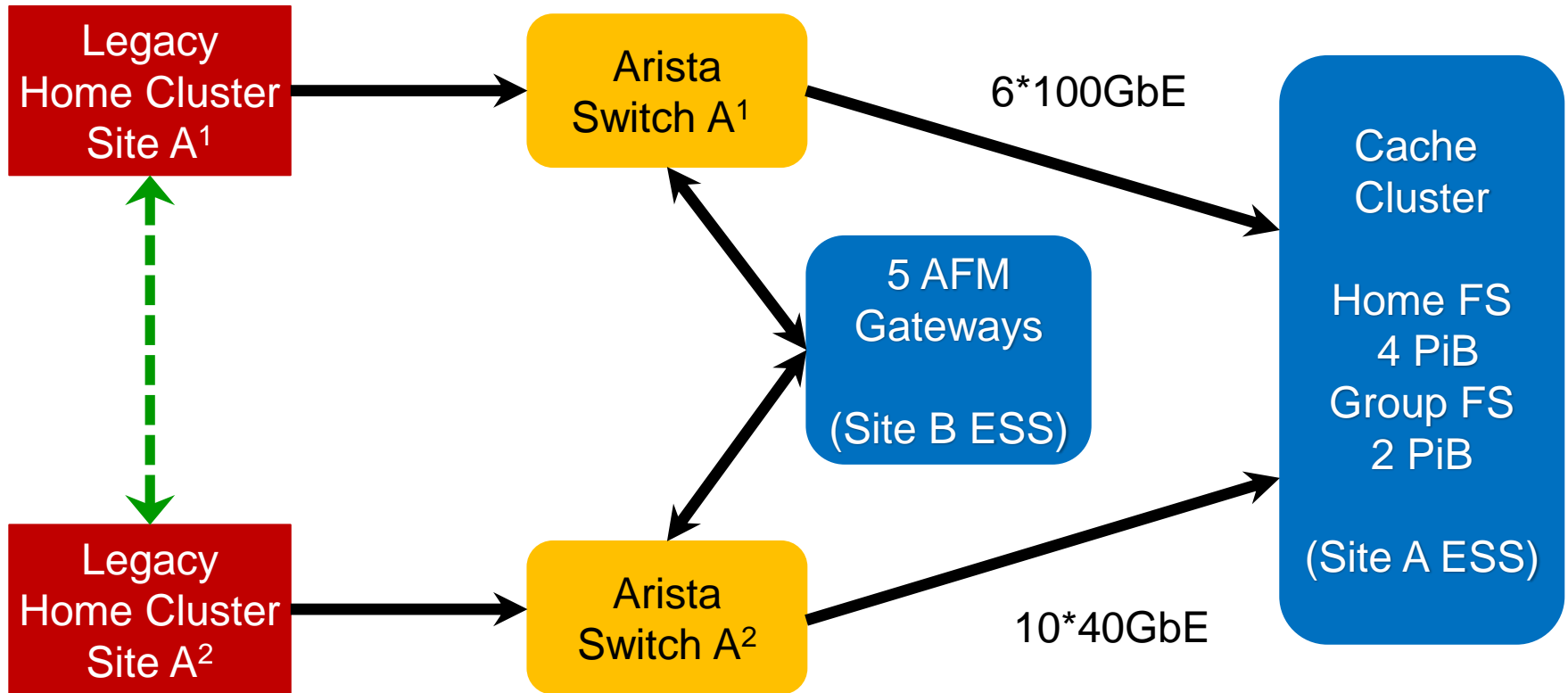


# Overview

- Legacy system upgrade
  - 3.5.0.32 -> 4.1.1.0 -> 4.1.1.9
    - FS version 3.5 to 4.1.1.4
  - 4.1.1.9 -> 4.2.1.0 -> 4.2.1.2
    - FS version 4.1.14 to 4.2.0
  - Small data files now being stored in metadata inode
  - Worked perfectly.....no issues
- Reality check.....no spare hardware available
  - Had to use ESS site B nodes as AFM gateways
    - Site B ESS disks shutdown
    - Vanilla GPFS 4.2.1.2 installed (probably not supported !)

**Legacy system out of support – fingers crossed .....**

# ESS - AFM Setup





# Migration procedure

- Home cluster
  - Disable TSM backups (run selective if needed), user dir “locked”
  - ILM policy to separate resident and migrated files into filelists
  - Migrated files need to be deleted first (restored later)
- Cache cluster
  - Independent fileset (local update mode)
  - Home cluster ILM filelists used to perform metadata & data prefetch
  - After transfer check AFM fileset for uncached/partial/dirty files (ILM)
  - When AFM transfer clean :
    - Fix up directory times (rsync)
    - Disassociate with home cluster (convert to std fileset)
    - Backup new data (mmbackup)
    - Restore migrated files from home cluster TSM client then backup again

**Group FS finished in 7-8 days (not 24\*7)**



## Site A failure

- Network switch configs corrupted
  - Zero communications until resolved
- DS3500 TSM DB/disk cache .....
- Lights on but no-one's home !
- Home FS corruption
  - mmfsck (logrecovery assert)
    - Tracked issue to directory structure (--skip-directory-check)
  - Unable to run ILM
    - No ILM filelists for AFM, no in-house backup (pre-mmbackup clone)
    - Fallback to in-house GNU find (-type m) & post processing for false positives
    - Verify via HSM dsmls command (not pretty)
  - mmfs daemon crashes when searching affected dirs (understandable)

**Only as strong as your weakest link !**



# Network Issues

- Single user transfer (default settings)
  - ~225 TiB, >500K files, average file size 536 MiB, largest file size 1.11 TiB
  - Listfile metadata prefetch took 12m30s
  - Listfile data prefetch took 8m57s
- Lots of uncached directories/files
  - Used a dedicated pipe, single transfer and overnight !
  - 1st migration attempt (~12.5hrs)
    - Significant errors on inter-switch links (1\*40GbE and 1\*100GbE)
  - 2nd migration attempt (~12hrs)
    - Additional bad 40GbE link detected and removed
  - 3rd migration attempt (~11.5hrs)
    - No significant errors, no uncached files encountered

**Appreciate your network team !**

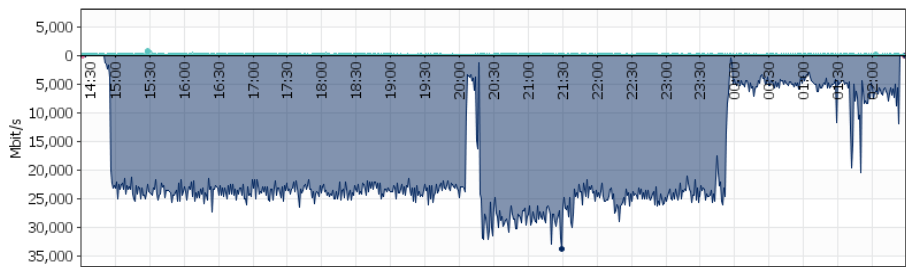




# AFM monitoring

- No nice method to work how much data was migrated :(
  - Had to go look at inter-switch link counters
  - Observed drop @20:05 happens at same point regardless of time of day

**10\*40GbE links**

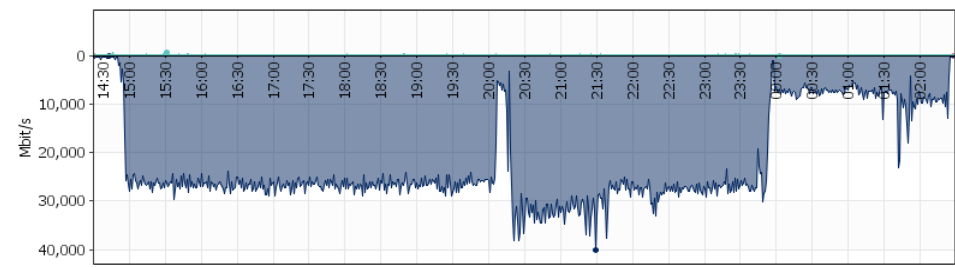


PRG Network Monitor 36.2.23.3270 17/02/2017

■ Downtime (%) ■ Traffic In (Mbit/s) ■ Traffic Out (Mbit/s) ■ Errors in (#/s) ■ Errors out (#/s) ■ Discards in (#/s) ■ Discards out (#/s)

Date Time	Traffic Total (volume)	Traffic Total (speed)	Traffic In (volume)	Traffic In (speed)	Traffic Out (volume)	Traffic Out (speed)	Errors in (v)
<b>Sums (of 720 values)</b>	100,271,349 MByte		1,413,077 MByte		98,858,271 MByte		0 #
<b>Averages (of 720 values)</b>	139,266 MByte	19,472 Mbit/s	1,963 MByte	274 Mbit/s	137,303 MByte	19,198 Mbit/s	0 #

**6\*100GbE links**



PRG Network Monitor 36.2.23.3270 17/02/2017

■ Downtime (%) ■ Traffic In (Mbit/s) ■ Traffic Out (Mbit/s) ■ Errors in (#/s) ■ Errors out (#/s) ■ Discards in (#/s) ■ Discards out (#/s)

Date Time	Traffic Total (volume)	Traffic Total (speed)	Traffic In (volume)	Traffic In (speed)	Traffic Out (volume)	Traffic Out (speed)	Errors in (v)
<b>Sums (of 720 values)</b>	114,295,913 MByte		963,227 MByte		113,332,686 MByte		0 #
<b>Averages (of 720 values)</b>	158,744 MByte	22,196 Mbit/s	1,338 MByte	187 Mbit/s	157,407 MByte	22,009 Mbit/s	0 #

**Same point on different switches – switch issue, legacy cluster, AFM ?**



# Issues

- mmfs daemon crash (4.2.1.x/4.2.2.2)
  - inodes exhausted during AFM transfer
- Inconsistent behaviour (4.2.1.x/4.2.2.2)
  - Empty directories not fetched
  - Fileset that first showed uncached entries for no reason
    - Fine once, it was unlinked and relinked
    - Other times it was resolved with another *find* run
  - Fileset that showing differing uncached entries between different runs
    - Fine once, gateways rebooted and transfer re-ran
  - Home cluster directories showed Archive attribute
    - Cache cluster showed Offline attribute



# Observations (1)

- Plan your AFM migrations
  - Use a reasonable number of concurrent AFM filesets
- Check the size of your file lists
  - Split if necessary across multiple gateways
- Watch you queue limits
  - `mmfasdm dump afm | grep QMem`
- Tune large filesets
  - `afmNumFlushThreads=16` or greater
- Independent fileset = double edged sword
  - Setup correct size prior to migration
  - Monitor and increase inodes as needed
    - Callback (`softQuotaExceeded`)
  - Watch your inode limits especially if you want to replicate the FS



## Observations (2)

- AFM gateways cannot run on AIX (restriction due to design)
- AFM documentation confusing.....getting better
- AFM migration forces use of independent fileset
- Some commands useful but hard to decipher
  - Watching single and parallel read threads via `mmfsadm dump afm`
- AFM perform any data verification ?
- AFM is mostly a black box - what is it doing ?
  - Provide IO rates to/from the cache filesets ?
  - How much data/files are behind, remaining data/files to be transferred ?
  - GUI monitoring - graph gateway usage (memory, queue length, I/O rates etc) ?
- AFM needs to resolve the use of `rsync` to fix directory times !

**Should we have just used `rsync/GNU parallel` instead ?**



## Our Gratitude to ...

Emily Barrett (IBM)

Stefano Gorino (CSCS)

Dean Hildebrand (IBM)

Venkateswara Puvvada (IBM)

Srikanth Srinivasan (IBM)