

MOVING TO AFS

SIMON WILKINSON
SCHOOL OF INFORMATICS
UNIVERSITY OF EDINBURGH
SIMON@SXW.ORG.UK

OVERVIEW

- ✻ Evaluation results
- ✻ Cell design
- ✻ Deployment process
- ✻ Deployment experiences

BACKGROUND ON INFORMATICS

- ✻ ~ 2000 active users, ~1500 hosts
- ✻ 20 Tb of centrally managed filestore
- ✻ Deployed Kerberos and LDAP infrastructure

OUR EXISTING FILESYSTEM

- ✻ NFS v3 based with Sun file servers and predominantly Linux clients
- ✻ AMD automounter providing identical filesystem on every machine
- ✻ Locally developed mechanisms to populate AMD filesystem maps, manage quotas, and do nightly mirroring
- ✻ Developed incrementally over many years.

WEAKNESSES

- ✱ Lack of security
 - ✱ Can't allow access from unmanaged machines
 - ✱ Can't allow access from beyond the firewall

WEAKNESSES

- ✱ Lack of portability
 - ✱ AMD infrastructure required significant modifications to off-the-shelf machines
- ✱ Lack of client availability for some systems

WEAKNESSES

- ✻ Lack of maintainability
 - ✻ Local glue required lots of effort just to keep running
 - ✻ Dealing with partition filling, and the resultant home directory moves
 - ✻ Fileserver failure leads to hung mounts, and lots of rebooting

CRITERIA

- ✻ Secure enough to permit access from foreign machines, and across firewalls
- ✻ Flexible ACL model
- ✻ Better performance
- ✻ Stability
- ✻ Linux and Solaris support required, Windows and Mac OS X desirable
- ✻ Easily scale to our client & data requirements
- ✻ No per-client licensing fees
- ✻ Preferably be a self-contained solution

CANDIDATES

☼ AFS

☼ CIFS

☼ Coda

☼ NFSv4

FEATURE COMPARISON

- ✻ On paper, most AFS features are present in NFSv4
- ✻ Critical absence is volume location independence
- ✻ Can't move filespace between servers without the user noticing
- ✻ No concept of a global namespace - still needs automounter glue!

EVALUATION

- ✻ AFS and NFSv4 feature sets very similar on paper, with NFSv4 leading the way
- ✻ However, NFSv4 “not quite ready yet” - few implementations of complete feature set
- ✻ Linux NFSv4 only did machine based authentication at mount time
- ✻ Bugs in NFSv4 implementation caused benchmarks to hang

BENCHMARKS

- ✱ Three benchmarks selected
 - ✱ iozone
 - ✱ blogbench
 - ✱ The Andrew Benchmark
- ✱ Only iozone and blogbench eventually used

BENCHMARKING RESULTS

- ✻ NFSv4 won the iozone one every time - by a small margin for files smaller than the AFS cache size
- ✻ Much more evenly matched with blogbench
- ✻ “Lies, damn lies, and statistics”

EVALUATION RESULTS

- ✻ NFSv4 just wasn't ready, and would still have required automounter madness.
- ✻ "Don't want our data to be their learning experience"
- ✻ OpenAFS met the majority of our criteria, with stability as an added bonus!

CELL DESIGN AUTHENTICATION

- ✱ AFS is tightly coupled with our authentication infrastructure
- ✱ Using RedHat's RH9 vintage pam_krb5 module (but planning on stopping)
- ✱ Using Doug Engert's pam_afs2 module (but looking at Russ's pam_openafs_session)

CELL DESIGN DIRECTORY

- ✿ Debated integrating pts with our existing LDAP directory
- ✿ Wrote some proof-of-concept code to backend pts with LDAP
- ✿ Decided that our LDAP service wasn't sufficiently reliable to do this in production
- ✿ Use 'standard' pts, with hooks into our account management system

CELL DESIGN - BACKUPS



CELL DESIGN ONLINE BACKUPS

- ✿ Our recent history makes us somewhat jumpy
- ✿ Off site disk mirrors was a requirement
- ✿ So, we use read-only user volumes
- ✿ All user volumes have an offsite RO copy which is released nightly.
- ✿ Backup volumes are still used to provide 'Yesterday' functionality, and tape backups ...

CELL DESIGN TAPE BACKUPS

- ✻ Finding a workable, scalable, tape backup system is a priority
- ✻ Currently embroiled in local politics
- ✻ At the moment, we just walk the AFS filespace and use our existing EBU licenses
- ✻ Not a very pretty bodge!

DEPLOYMENT EXPERIENCES

- ✻ Softly, softly ...
- ✻ Initially offered additional file space, rather than home directories, to the adventurous
- ✻ Gradually shifted computing staff home directories over
- ✻ Now creating all new users in AFS
- ✻ Starting to bulk move existing users

THINGS THAT MAKE OUR USERS SAD

- ✻ ACLs - especially the fact they are directory only
- ✻ Lack support for 'special' files such as devices or named pipes.
- ✻ Limits on maximum number of files per directory
- ✻ Linux's behaviour with sticky mode temporary directories

THINGS THAT CAUSED US PAIN

- ✻ .Xauthority files stored in home directories
- ✻ SSH public key files
- ✻ System daemons inheriting the PAG of the user starting them.
- ✻ Condor
- ✻ Beagle

SECURITY HURTS!

- ✻ Requirement to gain credentials before accessing files causes problems
 - ✻ Cron
 - ✻ Web servers
 - ✻ Condor and Grid Engine

SECURITY STILL HURTS

- ✻ Having to renew credentials is not popular
 - ✻ Long running jobs
 - ✻ Processes left running overnight
(Thunderbird, gnome-screensaver!)
- ✻ Unix applications aren't good at dealing with unexpected FS failure

REDUCE THE PAIN

- ✻ Get your filesystem credentials at login
- ✻ Renew them whenever you can (screensavers &c.)
- ✻ Don't have credentials expiring in the middle of the day
- ✻ Make sure all credentials renewal tools renew AFS tokens, too

LONG RUNNING JOBS

- ✻ Provide a mechanism for stashing credentials with a subset of permissions on the local disk
- ✻ Encourage people to use this to provide credentials for long running jobs
- ✻ k5start and k5renew are hugely useful tools
- ✻ Renewable tickets are great for medium-life jobs!

CONCLUSIONS

- ✻ Going well so far
- ✻ The crunch point is just around the corner!
- ✻ Softly, softly has perhaps been too soft
- ✻ Ensuring reliability before moving users, and responding rapidly to their concerns has been key

THANKS!

☼ There's a lot of good code and support out there!

QUESTIONS?