

Customer-Oriented Storage Performance Management

Dany Felzenszwalbe Intel

Legal Notices

This presentation is for informational purposes only. INTEL MAKES NO WARRANTIES, EXPRESS OR IMPLIED, IN THIS SUMMARY.

Intel and the Intel logo are trademarks of Intel Corporation in the U.S. and/or other countries.

* Other names and brands may be claimed as the property of others.

Copyright © 2014, Intel Corporation. All rights reserved.



About the Presenter

Dany Felzenszwalbe

- Joined Intel in 2001
- □ Working in PDIT / EC
 - Product Development IT
 - **Engineering** Computing
- Primary focus is NAS/NFS data and storage solutions



Felsenschwalbe

Cliff swallow



Agenda

- Intel Design NAS environment overview
- Performance monitoring challenges
- Performance management overview
- Reactive performance management
- Proactive performance management
- Lessons learnt
- Call for Action



Intel's Design NAS Environment



43 Petabytes of NFS capacity

- On > 1000 fileservers
- Multiple vendors / platforms
- In 44 locations



91% in the largest 10 data centers

- 58% in the largest 4 data centers
- From 1 to 190 fileservers in a site

More than 1000 projects

- About 30,000 Users (aka Customers)
- 20-30 million batch jobs running weekly



... and

"everything"

has to run

"FAST"… !







NFS Performance

- 4 performance (and cost) tiers are defined for NFS
- The design teams/activities are expected to use the appropriate tiers
- When performance is "bad" the tier "doesn't matter"
- Understanding when it happens is challenging
- Severe NFS performance issues could impact a project's TTM



Cost

 Until 2007 we relied on customers as our Monitoring

SD @

2014 Storage Developer Conference. © Intel. All Rights Reserved.

Performance

Customer-Oriented Storage Performance Management





Resources Utilization Monitoring



- Monitoring the fileservers' CPU utilization, diskdrives utilization, network collisions and such
- Used Cricket (Open Source from SourceForge) for tracking and as an alerting mechanism
- No correlation between high utilization and users reporting performance issues





Client-Side User-Experience Monitoring



- Simulate user activity to identify problems
 - Slow mounts and reads
- Using home grown monitoring framework
 - 2 clients
 - Different subnets
 - Duration matters
 - Time of day
 - Server type/tier

HARMLESS – event is 90 seconds long CRICITICAL – event is 300 seconds long FATAL – event is 600 seconds long



Custom Use-Case Performance Monitoring



- In spite of the successful monitoring capabilities some undetected events remained painful
- Still needed to identify problems before customers
- Cloning with GIT and a similar tool were chosen and monitors were created



Latencies Monitoring



Latency is a key metric that should provide visibility into customer impact



- Fileserver reported latencies should correlate with the user-experience monitoring and should also match other performance issues ("misses")
- Different platforms offer different capabilities
- We have been experimenting with "disk" latencies and NFS operations latencies





SD[®]

Analysis and Resolution



We know that a NFS server is slow – now what ?



We need to promptly clear the impact We need the same solutions for any platform



I 5 ℃ ↑ ↓	🔏 👻 🗢 NFS slowness - FILESERVER10.intel.com (FATAL) - reaso	on copy - sent by MONITORING_CLIENT2 - Message (HTML)	? 🖻 – 🗆 🗙					
FILE MESSAGE IN	SERT OPTIONS FORMAT TEXT REVIEW							
	🚫 🖂 📴 Meeting 🎽 temp	💼 Rules 🐐 🏹 😪 Mark Unread 🔐 🏦 Find	Q = A					
	L← L→ G‡IM - C→ To Manager -	PoneNote Categorize - Related -						
Junk - Delete Reply	All More - Team Email - More	Actions * Policy * Follow Up * Select *	Zoom Start Inking					
Delete	Respond Quick Steps 🕞	Move Tags 🖬 Editing	Zoom Ink 🔺					
Sun 8/10/201	4 10:43 PM							
root Fileserver name, severity (duration), copy/mount								
NFS slow	ness - FILESERVER10.intel.com (FATAL) - reason copy	/ - sent by MONITORING_CLIENT2						
To nfsadmin	anvf							
Retention Policy Mail Cloud -	Inbox (60 days)	Expires 10/9/2014						
			_					
Model VENDOR A MC	DEL VV (tier 2)	tion and groups	<u> </u>					
Business groups wit	th data on fileserver: Jeiver model,	tier and groups						
ABC	usir	ig it						
BCD		0						
SUMMARY								
NES	operations							
read	6589 39.92%							
write	57 0.35%	VVORKIOAD						
lookup	916 5.55%	(OPS min)						
getattr	3303 20.01%	(OPS mix)						
access	5640 34.17%							
username1 OPS 10	0735 weight 3350.771914							
(read=6546 write=0	lookup=734 getattr=1774 ac ess=1670 create=0) remove=0 readdir=0 readdirplus=0 setattr=0 rmd	ir=0 symlink=0)					
pool_1 1								
pool_2 271								
pool_3 405	Iop-hitter user with OPS mix							
disk6	and Patch jobs amounts and	/nfs/site/disks/some_project_disk477	8650					
disk0	and batch jobs amounts and	/his/site/disks/some_project_disk059	1025					
disk0	^a paths being used	/nfs/site/disks/some_project_disk058	155					
disk0	a paths being used /nfs/site/disks/some_project_disk030							
disk0	s	/nfs/site/disks/some_project_disk001	82					
disku	some_project_diskUS1	/nis/site/disks/some_project_diskUS1	12					
username2 OPS 217	73 weight 277.846433							
(read=0 write=0 loc	kup=36 getattr=515 access=1618 create=0 rem	ove=0 readdir=0 readdirplus=0 setattr=0 rmdir=0	symlink=0)					
pool_1 262								
pool_2 334								
boor ² 183								
disk0	some project disk058	/nfs/site/disks/some project disk058	1181					
disk0	some project disk117	/nfs/site/disks/some project disk117	231					
disk0	some project_disk051 /nfs/site/disks/some project_disk051							
disk0	some project disk030	/nfs/site/disks/some project disk030	208					
disk0	some_project_disk059	/nfs/site/disks/some_project_disk059	192					
disk0	some_project_disk001 /nfs/site/disks/some_project_disk001 103							
disk2	some_project_disk071	/nfs/site/disks/some_project_disk071	36					
•			•					

2014 Storage Developer Conference. © Intel. All Rights Reserved.

SD @



Automatic Resolution



- Different approach "autosuspend"
 - No human intervention required
 - Starts at the first sign of slowness
 - Hooks on "yana" and suspends jobs by clients' traffic
- Tradeoff between NFS Slowness and Batch jobs suspending



Analytics



Can't see the forest for the trees

- Many sites, servers, admins, customers, jobs
- Needed trending and visualization
- Needed to focus on the "big" top-hitters

State 2004 45 MM Contained d Cont	Description Filescription Read Base 1 Water 2 Unit 1 Water 1 Water 1 Water Base 1 Water 1 Water 1 Water 1 Water 1 Water 1 Water Base 1 Water 1 Water	Top Marking between segmented bet belog bad en MES mess - Mensage (HTM) F = - (1 X A construction of the second sec
Hi Jane Dos. About 1861 of your Netbatch jobs caused high In order to prevent alsowness for other users and We support that the disk with high usage is offst Please asset that the first with high usage is offst Please constat your DARA and try to find a sy Please constat to 24.67 of you have any question Have a size day	BI Need	ed !	rec maters relations
Regards. For Computing Monitoring & Control Centrol For Computing Monitoring & Control Centrol Image: Computing Monitoring & Control Centrol For Control Centrol For Control Centrol Image: Computing Monitoring & Control Centrol For Control Centrol For Control Centrol Image: Computing Monitoring & Control Centrol For Control Centrol For Control Centrol Image: Control Centrol Image: Control Centrol For Control Centrol Image: Control Centrol Image: Control Centrol Image: Control Centrol Image: Control Centrol Image: Control Centrol Image: Control Centrol Image: Control Centrol Image: Control Centrol Image: Control Centrol Image: Control Centrol Image: Control Centrol Image: Control Centrol Image: Control Centrol Image: Control Centrol Image: Control Centrol Image: Control Centrol Image: Control Centrol Image: Control Centrol Image: Control Centrol Image: Control Centrol Image: Control Centrol Image: Control Centrol Image: Control Centrol Image: Control Centrol Image: Control Centrol Image: Centrol Centrol Image: Centrol Centrol Image: Control Image: Centrol Centrol Image: Cent	We suspect that the disk with high usage is infolded belower project_area01063 Please note that the Policy was approved by Engineering Computing and customer groups. Please constity over DARA and nyt for fand avy to over disk seck and 6 susses in the future. We also suggest your consider submitting your jobs in small amounts (through a Netbatch feeder) to prevent causing environmental slowness. Please constitute of any Please constitute of any Here a nice day Regards, Computing Monotring & Control Cented	the second	a readility lased accessed matters equilibred hard/accessed accessed matters equilibred hard/accessed accessed accessed hard/accessed accessed hard/accessed

Trending

 Reactive Performance Management - Real-time
 Proactive Performance Management

 Server
 Cient

 Besources
 User Experience

 Flargstein
 Custom Use-cases

 Phontering
 Custom Use-cases

 Postering
 Resolution

 Postering
 Resolution



SD @

Top-Hitters

Currently applied filters

en estera.

SD

14

Work Week: 2000 - 200 - 200

severity: severity - CRITICAL, FATAL, MINOR;

Duration by Fileserver Model

© *				٥	Ø *						
	ileserver	cell	tier	model	Business Groups	Duration		User	cell	Total hours of suspend	Jobs count
0			2		(in the second sec	6525	0	(ineselle	8	104636	32652
1	A 10	-	3	Bydan	1977 AN	4322	1		gd#	48643	6625
2	344	٠	4	404adi	(*******	1773	2	<i>e</i>	etas	48118	14470
3			1	400		1260	3	-		39635	41259
4		*	1			744	4			39108	8682
5		-	1	-		655	5	diagram.	Altimo	35319	12703
6	Nghania	-	3	andter		567	6	estilities.	4500	27834	11950
7			2	(jedža		540	7		4000	23950	6618
8	-		3	Survey.	*	470	8	and the set	#	18965	7157
9	acidette.		3			432	9	-	686	17729	11981
10	(>	1	**	*** ***	369	10	() design	etio	17207	8606
11	-		2		the second second	353	11	àng	51(sp.	10711	4639
12	(******	-	5	-	() And	239	12	William		9863	9891
13		****	1			238	13	diame	*	7782	10842
14	1		3			213	14	- Annalise	\$10	7592	237219
15			3	(alter a		213	15	(manifing	65	7293	7344
16		() (M)	1			205	16		8 00	5385	4945
17	(inches	-				199	17		40	4851	7889
18					· · · · · · · · · · · · · · · · · · ·						1042

The fileservers with the most slowness and the users with the longest job suspends are what we want to take action on !

2632 of 25 > > Page 1



Drill-down for action NFS Slowness - users & paths (by ops) Currently applied filters dsm_event_www: 201413 - 201420; severity: severity: califical, FATAL, MINOR;

businessgroup

9002.3

diffe.

10000

67.5

I STORES

and the state of t

1000

1000

1712

1000

Show Edit project w businessgroup ¥ cel X Apply All Filters otal of Slowness project qslot d_fileserver model Times as top-hitter 65.0 4000 CONTRACTOR OF A DESCRIPTION OF A DESCRIP 392 5 1000 1200 267 383 3

difference on

and the second

and the

a contra

101000-0

-

St. official

40000

A1453151

10000

22000

145

96

48

45

40

32

31

31

31

30

28

26

STATISTICS.

CONTRACTOR OFFICE

Confactoria in

a day of

1000

and the second

and the second second

a grint in the se

and a set

100

the second

and distant

-

Reactive Performance Management - Real-time

Client

User Experience

Custom Use-cases

Automatic Analysis

Resolution

Server

Resources

Utilization Filesystem Latency Monitoring

OPS Latency Monitoring

We can tell which problems are reoccurring by fileserver, specific export/path and activity

10000

-off-b

1000

1200

end 531 m

32.4

Contraction of the

-

17

14

.

3

1

3

.

5

1

1

2

1

2

2

1

2

Proactive Performance Management

Solutions Tool-Box And

Optimizatio

Lest Successful Refresh: 8/18/2014 2:56:45 PM (UTC + 3

Future won



cell: cell - celle celle

path

fileserver v

/efs/ /proj/

/ofs/ /proj/

/nfs/ proj/contine

Infish / disks/

/nfs/:/proj/

/nfs/~_proj/

/rfs/ Uproj/ 200

/nfs/ /disks

/rifs/ /proj/

/nfs/o /disks/ sep.p. o

/rfs/~__/disks/

/rfs/join/disks/

/nfs/ "disks/

/nfs/1/proj/

dsm event ww V.

cell

-

100

1010

25

witco:

100

Among A

1000

-

6212

mile-

100

020

104

Filters

eventr 2

(B. +

2

3

4

5

fileserver

seminates.

and the second

distantion in the

Cherry Mar

10000

Summer of

Contrario -

all the second

and the second

180

1000

Rest Statistics

ontento

root

Sec.

opposite the same

0.000

dillo-

Our Solutions Tool-Box



- Data migration to appropriate tier or platform
- Batch dispatch configurations slow ramp
- Local disk caching enablement
- Smarter Allocation consider slowness history
- Customer Flow changes
- Refresh of old HW
- Backups tuning







Future Plans

- Checking latencies usability
- Integration of new platforms
- Analyzing non-NFS traffic
- Monitoring SMB performance
- QoS implementation
- Automation for reoccurring top-hitters









Lessons Learnt



- User-experience monitoring
- Sometimes the simplest solution works best
- Automate everything, with caution
- NFS performance issues can be resolved
 Sometimes a downtime may be required



2014 Storage Developer Conference. © Intel. All Rights Reserved

Call For Action

- What we need from the NAS Vendors and/of the open source community:
 - Identify "Hot files/users/clients"
 - Alerts on counters with thresholds
 - **When** there is a performance problem ?
 - And Why would be ideal
 - CPU ? Cache ? Disks ? Network ?
- Standards and API ?





Any questions or feedback, please contact me at dany.felzenszwalbe@intel.com



2014 Storage Developer Conference. © Intel. All Rights Reserved.