

Enhancing performance and reliability of Network File System

*Note: Sub-titles are not captured in Xplore and should not be used

1st Aswin Babu Karuvally
dept. of computer applications
College of engineering, Trivandrum
Trivandrum, India
aswinbabuk@gmail.com

2nd Basith Hameem
dept. of computer applications
College of Engineering, Trivandrum
Trivandrum, India
basithhameem@cet.ac.in

3rd Ann Jerin Sundar
dept. of computer applications
College of Engineering Trivandrum
Trivandrum, India
annjerinajs@gmail.com

4th Prof. John Prakash Joseph
dept. of Computer Applications
College of Engineering, Trivandrum
Trivandrum, India
john@cet.ac.in

Abstract—Network File System is a commonly used distributed file system, allowing the user to access and manipulate storage on remote computers as if they were part of the local machine. Network File System is notoriously slow in its default configuration. This is accentuated if the NFS environment has more than a dozen clients. When configured to deliver faster speeds by turning on asynchronous mode, the system suffers from higher risk of data corruption and loss.

This paper proposes a number of modifications to the Network File System allowing the file system to provide very high performance, while minimizing the risk of data loss and corruption. Further, the proposed system behaves better in congested networks by consuming less bandwidth ensuring decent speeds even during periods of heavy network traffic.

Index Terms—UNIX, NFS, Performance, Data loss, Data corruption, Reliability, Optimize, Speed, Tweak, Link, Journal

I. INTRODUCTION

Network File System is a distributed File System protocol primarily used by the UNIX family of Operating Systems. It allows users to mount, access and manipulate disk partitions or directories on a remote computer, as if the said partition or directory was a part of the local machine. Network File System was developed as an open standard by SUN Microsystems in 1984.

NFS is widely used in Local Area Networks to conveniently share data and provides users the ability to access their files across the network. Sometimes, a directory access protocol such as LDAP is combined with NFS allowing the users to login to their user account from any computer on the network.

The main drawback of NFS is the slow read and write speeds it offers with the default setup. Though NFS offers a number of parameters in its configuration files to increase the performance, these either do not affect the performance much or increases the chance of data corruption and loss. Thus

the users are forced to run the system with the default, slow configuration.

In many environments, this leads to lost human productivity as interactive computing becomes impossible and Operating System needing access to data on NFS share often ends up freezing the whole computer. This also bottlenecks the CPU of the computer, thereby wasting precious computing resources.

This paper proposes a number of changes to the Network File System protocol which increases the performance of Network File System while reducing the risk of data corruption and loss. The proposed system also ensures decent speeds in congested network as it consumes less bandwidth than the original NFS implementation and protects the ability of the NFS server to provide access to files during period of peak network activity.

II. BACKGROUND WORK

A. Maintaining the Integrity of the Specifications

The IEEEtran class file is used to format your paper and style the text. All margins, column widths, line spaces, and text fonts are prescribed; please do not alter them. You may note peculiarities. For example, the head margin measures proportionately more than is customary. This measurement and others are deliberate, using specifications that anticipate your paper as one part of the entire proceedings, and not as an independent document. Please do not revise any of the current designations.

III. RESEARCH METHODOLOGY

The whole NFS system is tried to benchmark using virtual machine and nfs version 4.2 is used for simulation. One of the main reason to use virtual machine is that, they provide a number of networking hardware specifications which allow us to simulate target machine network hardware properties.

Identify applicable funding agency here. If none, delete this.

Virtual machine is free to use, fast and easy to implement. A 100 Mbps bridged adapter is setup in the virtual machine to directly to outside network. In other words NAT (Network Address Translation) is not used because in NAT the translation is done by the CPU. This make sure that CPU will not bottleneck the system simulation. Benchmarking the system was a tedious process because of the absence of proper benchmarking tool. The benchmarking was properly completed by combing some of the available tools. Bonnie ++, Phoronix, Dbench and dd-the command line utility was used. The main disadvantages of first three test suits was caching effect. This was effectively avoided by using the command line utility "dd". Dbench is a powerful benchmarking tool which itself can simulate upto 512 clients by its own. Benchmarking was done through proper combination of all these tools.

```
ddif = /dev/zero of = /home/user1/tmpfile bsize = 1K count = 2048000 conv = fdatasync
```

The "conv=fdatasync" option flushes each 1 KB of data from cache and this ensures that caching effect on virtual machine and nfs performance is zero.

IV. WORKING AND TYPICAL PERFORMANCE OF NFS

Network File System is a distributed file system originally developed by SUN Microsystems in 1984. A distributed file system differ from normal file systems in the sense that they operate over the network and allow sharing of storage resources. NFS was developed as an open standard. It was initially implemented for UNIX but is now compatible with a wide array of Operating Systems. NFS has gone through four major revisions, with the first publicly available version being v2. All the experiments for this research have been conducted using NFS v4.2

NFS uses the Client-Server architecture. The NFS Server makes available a disk partition or directory on the network which can then be mounted just like a local storage device by the clients. For applications running on the client machine, NFS is just another file system and requires zero modifications inorder to work. This is made possible by an abstraction layer called Virtual File System or VFS. VFS defines what operations can be done on the filesystem regardless of the file system type. When an application deals with the UNIX file system, it is actually dealing with VFS. VFS receives data regarding the target file from the application, then hands it over to the actual file system, which in this case is NFS.

Once NFS receives the data, it is transferred to the server with the help of Remote Procedure Call (RPC) in External Data Representation (XDR) format. RPC allows the NFS client to execute instructions on the server. XDR is a data representation standard that provides a uniform data format which can be understood by a variety of computers. This is one of the factors that provide NFS with cross platform compatibilty. The working of NFS protocol has been shown in figure.

VFS-NFS handover Figure

NFS by default runs in what is called server side synchronous mode. In this mode, when the NFS client receives a

write operation, it connects to the server, requests a write and transfers the data. Once the transfer is complete, the server syncs the data to its disks. After completing the sync operation, server returns an acknowledgment message back to the client. The issue here is that, the client has to wait till it receives the acknowledgment. It cannot perform any additional write to the server till the acknowledgment arrives. This considerably slows down the system. Further, clients are often configured to work in client side synchronous mode. In this mode, the client is forced to write the data to the server as soon as it receives the request. This often causes the system to crawl.

The performance of an NFS system with client and server side synchronous mode turned on was benchmarked with client side options in /etc/fstab set as [rw, sync, hard, intr 0 0] and the server side options in /etc/exports set as [rw, no-root-squash, subtree-check]. The benchmarks were conducted using dbench utility part of Phoronix Test Suite. This resulted in performance of 0.94 MB/s. The server and client systems in this case were equipped with 100mbps network adapters. The speeds obtained are thus clearly sub-optimal. The performance graph obtained from the benchmark has been shown in figure.

V. PARAMETERS AFFECTING PERFORMANCE

Testing environment Testing environment is set up in a HP-Z640 Desktop Workstation running RHEL (RedHat Enterprise Linux). A NFS server and six clients were setup on a virtual machine. Dbench is specifically meant for SMB/NFS benchmarking. Using the above tools for benchmarking, various tests were run on the previously setup environment. There are various factors that can be manipulated in the configuration of NFS server side and client side. NFS client side configuration is done in "/etc/fstab" whereas NFS server side configuration is done in "/etc/exports". Various combinations of these

TABLE I
NFS OPTIONS-CLIENT SIDE

SlNo.	Option	Description
1	rw	Read/Write
2	syn	Sync file system with the server
3	hard	NFS requests are retried indefinitely
4	intr	Provided for backward compatibility
5	nfsvers	Specifies the nfs versions
6	rsize	Maximum number of bytes when reading data
7	wsize	Maximum number of bytes when writing data
8	udp	Specifies the connection to UDP
9	async	Asynchronous write

TABLE II
NFS OPTIONS-SERVER SIDE

SlNo.	Option	Description
1	rw	Read/Write
2	no-root-squash	Turn off root squashing
3	subtree-check	Specified directory/its subrectory for access
4	async	Synchronous write
5	sync	Asynchronous write

options were experimentally simulated in the above mentioned

environment. Test span was from 40 minutes to 12 hours. Test span depends upon the combination of tools we use and the options we enforce in `"/etc/fstab"` and `"/etc/exports"`. It is found that there is considerable decrease in performance when UDP used in client side. Write speed decrement 0.94 MB/s to 0.77 MB/s. More importantly there is sharp increase in the performance by turning on `"async"` mode on both server side and client side.

VI. PERFORMANCE WITH SERVER SIDE ASYNC

In server side synchronous mode, the server waits till the data has been written to its disk before returning the acknowledgment message to the client. Server side asynchronous mode changes the behaviour of NFS server such that it returns the acknowledgment message as soon as the client completes the transfer of data. This has tremendous impact on the performance of the system.

An NFS system with server side asynchronous mode was benchmarked with the client side options in `/etc/fstab` set as `[rw, sync, hard, intr 0 0]` and server side options in `/etc/exports` set as `[rw,no-root-squash,no-subtree-check, async]`. The test was conducted using `dbench` tool - part of `phoronix` test suite, and resulted in 28.72 MB/s, a huge increase in performance boost, compared with the earlier server side synchronous mode, which returned 0.94 MB/s. The performance graph obtained has been shown in figure. Further benchmarks were also conducted with `dd` utility, with a block size of 1KB, file size of 2GB and `fdatasync` option turned on. The block size of 1K ensures the computer writes only 1KB of data at a time and the `fdatasync` option flushes the cache after each block is written. The aim was to understand the performance of the system without the contribution of client side caches. The benchmark resulted in 7.1 MB/s which was still significantly higher than results recorded with server side asynchronous mode. Performance Figure. The higher performance comes from the fact that the clients do not have to wait till the server syncs the data with its disks. Clients can transfer data to the server, then get on with other tasks such as writing additional files. This also means that more number of clients can access the server in unit time. Still, with the current protocol, it is not advisable to leave server side `async` turned on due to the possibility of data loss and corruption.

A. Reliability concerns with Server side ASYNC

One side effect of enabling server side asynchronous mode is that, more clients can access the server in unit time. This in turn can cause a write queue to form on the client side. That is, writing a file to the server gives no guarantee that it has been written to permanent storage. In case the server crashes immediately after some data has been transferred to it, be it software crash or hardware failure, it can cause data corruption or loss. The more worrying fact is that, it is not just a single computer which loses data. Data loss can occur to most of the clients which has written to the server shortly before the crash.

If a server crash occurs, a client has no means of protecting itself from the data loss. A client has no NFS cache that is permanent in nature. Even if the client has the lost file in its primary memory, there is a high chance of losing it. This is because, if prior to crash the server was serving a critical configuration file, the application dependant on the file can crash. If the application is part of the Operating System, it can bring down the whole system. The latter is often the case with environments where home directory is served by an NFS share.

B. Fixing server side ASYNC behaviour

Once VFS handovers the write request to NFS client, it transfers the received data to the server rather than writing it to the local storage. In case of data loss, the file cannot be recovered, as the only copy of the file was in server's memory. The solution is to create a buffer in the client's local storage such that, a copy of all the data written to server will be kept with the client.

The buffer is a predefined storage area in the client's secondary memory. It acts like a ring buffer with a flexible memory size. The oldest files are deleted once the buffer reaches a predefined size. NFS client will maintain a plaintext file in the buffer containing names of each file in the buffer, its path and hash generated from each file. During an NFS write operation, the client stores a copy of the file in the buffer area and updates the metadata file with the information regarding the new file. The hash is calculated whenever a file is written to the buffer. To minimize the CPU overhead for hash calculation, one should implement a lightweight hashing algorithm such as QUARK or PHOTON. Figure shows the working and structure of the buffer.

Figure of the working and structure of buffer + metadata file

Optional (Sample implementation of buffer deletion strategy)

In case of a server crash, the server creates a list of corrupted or lost files. During the first boot after the crash, the server requests the metadata file from each client that connects to it. Once the server receives the metadata, it calculates the hash for the local copy of each file that is listed in the metadata. If a file is missing or if the hashes do not match, they are added to `retransmit-list`, a list of files to be retransmitted from the client. Once the metadata file from a client is fully scanned, the `retransmit list` is sent to the client. The client in turn transmits a new copy of each file in the `retransmit list`. Figure shows the error recovery process.

C. Equations

Number equations consecutively. To make your equations more compact, you may use the solidus (`/`), the `exp` function, or appropriate exponents. Italicize Roman symbols for quantities and variables, but not Greek symbols. Use a long dash rather than a hyphen for a minus sign. Punctuate equations with commas or periods when they are part of a sentence, as in:

$$a + b = \gamma \quad (1)$$

Be sure that the symbols in your equation have been defined before or immediately following the equation. Use “(1)”, not “Eq. (1)” or “equation (1)”, except at the beginning of a sentence: “Equation (1) is . . .”

D. *LaTeX-Specific Advice*

Please use “soft” (e.g., `\eqref{Eq}`) cross references instead of “hard” references (e.g., (1)). That will make it possible to combine sections, add equations, or change the order of figures or citations without having to go through the file line by line.

Please don’t use the `{eqnarray}` equation environment. Use `{align}` or `{IEEEeqnarray}` instead. The `{eqnarray}` environment leaves unsightly spaces around relation symbols.

Please note that the `{subequations}` environment in *LaTeX* will increment the main equation counter even when there are no equation numbers displayed. If you forget that, you might write an article in which the equation numbers skip from (17) to (20), causing the copy editors to wonder if you’ve discovered a new method of counting.

BIBTeX does not work by magic. It doesn’t get the bibliographic data from thin air but from .bib files. If you use *BIBTeX* to produce a bibliography you must send the .bib files.

LaTeX can’t read your mind. If you assign the same label to a subsection and a table, you might find that Table I has been cross referenced as Table IV-B3.

LaTeX does not have precognitive abilities. If you put a `\label` command before the command that updates the counter it’s supposed to be using, the label will pick up the last counter to be cross referenced instead. In particular, a `\label` command should not go before the caption of a figure or a table.

Do not use `\nonumber` inside the `{array}` environment. It will not stop equation numbers inside `{array}` (there won’t be any anyway) and it might stop a wanted equation number in the surrounding equation.

E. *Some Common Mistakes*

- The word “data” is plural, not singular.
- The subscript for the permeability of vacuum μ_0 , and other common scientific constants, is zero with subscript formatting, not a lowercase letter “o”.
- In American English, commas, semicolons, periods, question and exclamation marks are located within quotation marks only when a complete thought or name is cited, such as a title or full quotation. When quotation marks are used, instead of a bold or italic typeface, to highlight a word or phrase, punctuation should appear outside of the quotation marks. A parenthetical phrase or statement at the end of a sentence is punctuated outside of the closing parenthesis (like this). (A parenthetical sentence is punctuated within the parentheses.)
- A graph within a graph is an “inset”, not an “insert”. The word alternatively is preferred to the word “alternately” (unless you really mean something that alternates).

- Do not use the word “essentially” to mean “approximately” or “effectively”.
- In your paper title, if the words “that uses” can accurately replace the word “using”, capitalize the “u”; if not, keep using lower-cased.
- Be aware of the different meanings of the homophones “affect” and “effect”, “complement” and “compliment”, “discreet” and “discrete”, “principal” and “principle”.
- Do not confuse “imply” and “infer”.
- The prefix “non” is not a word; it should be joined to the word it modifies, usually without a hyphen.
- There is no period after the “et” in the Latin abbreviation “et al.”.
- The abbreviation “i.e.” means “that is”, and the abbreviation “e.g.” means “for example”.

An excellent style manual for science writers is [7].

F. *Authors and Affiliations*

The class file is designed for, but not limited to, six authors. A minimum of one author is required for all conference articles. Author names should be listed starting from left to right and then moving down to the next line. This is the author sequence that will be used in future citations and by indexing services. Names should not be listed in columns nor group by affiliation. Please keep your affiliations as succinct as possible (for example, do not differentiate among departments of the same organization).

G. *Identify the Headings*

Headings, or heads, are organizational devices that guide the reader through your paper. There are two types: component heads and text heads.

Component heads identify the different components of your paper and are not topically subordinate to each other. Examples include Acknowledgments and References and, for these, the correct style to use is “Heading 5”. Use “figure caption” for your Figure captions, and “table head” for your table title. Run-in heads, such as “Abstract”, will require you to apply a style (in this case, italic) in addition to the style provided by the drop down menu to differentiate the head from the text.

Text heads organize the topics on a relational, hierarchical basis. For example, the paper title is the primary text head because all subsequent material relates and elaborates on this one topic. If there are two or more sub-topics, the next level head (uppercase Roman numerals) should be used and, conversely, if there are not at least two sub-topics, then no subheads should be introduced.

H. *Figures and Tables*

a) *Positioning Figures and Tables:* Place figures and tables at the top and bottom of columns. Avoid placing them in the middle of columns. Large figures and tables may span across both columns. Figure captions should be below the figures; table heads should appear above the tables. Insert

TABLE III
NFS OPTIONS-CLIENT SIDE

SI.No.	Option	Description
1	rw	Read/Write
2	syn	Sync file system with the server
3	hard	NFS requests are retried indefinitely
4	intr	Provided for backward compatibility
5	nfsvers	Specifies the nfs versions
6	rsize	Maximum number of bytes when reading data
7	wsize	Maximum number of bytes when writing data
8	udp	Specifies the connection to UDP
9	async	Asynchronous write

^aSample of a Table footnote.



SCHOOL OF DISTANCE EDUCATION

494 2407356, 2400288

Calicut University.P.O, 673635

466/SDE-A-ASST-2/2017/Admn

01.11.2017

പത്രക്കുറിപ്പ്

(എല്ലാ എഡിഷനുകളിലും പ്രസിദ്ധീകരിക്കുന്നതിനുവേണ്ടി)

കാലിക്കറ്റ് സർവ്വകലാശാല വിദൂരവിദ്യാഭ്യാസ വിഭാഗം നടത്തുന്ന സിദ്ധാന്തം
ജ്ഞാനന്തര ബിരുദ പ്രോഗ്രാമുകളിലേക്ക് 2017-2018 അദ്ധ്യയന വർഷത്തെ
വശനത്തിന് പിഴയടയ്ക്കൽ 15-11-2017 വരെയും 100/- രൂപ പിഴയടയ്ക്കൽ 20-11-2017
യ്ക്കും ഓൺലൈനായി അപേക്ഷിക്കാവുന്നതാണ്.

Sd/-
Deputy Registrar

Fig. 1. Example of a figure caption.

figures and tables after they are cited in the text. Use the abbreviation “Fig. 1”, even at the beginning of a sentence.

Figure Labels: Use 8 point Times New Roman for Figure labels. Use words rather than symbols or abbreviations when writing Figure axis labels to avoid confusing the reader. As an example, write the quantity “Magnetization”, or “Magnetization, M”, not just “M”. If including units in the label, present them within parentheses. Do not label axes only with units. In the example, write “Magnetization (A/m)” or “Magnetization {A[m(1)]}”, not just “A/m”. Do not label axes with a ratio of quantities and units. For example, write “Temperature (K)”, not “Temperature/K”.

ACKNOWLEDGMENT

The preferred spelling of the word “acknowledgment” in America is without an “e” after the “g”. Avoid the stilted expression “one of us (R. B. G.) thanks ...”. Instead, try “R. B. G. thanks...”. Put sponsor acknowledgments in the unnumbered footnote on the first page.

REFERENCES

Please number citations consecutively within brackets [1]. The sentence punctuation follows the bracket [2]. Refer simply to the reference number, as in [3]—do not use “Ref. [3]” or “reference [3]” except at the beginning of a sentence: “Reference [3] was the first ...”

Number footnotes separately in superscripts. Place the actual footnote at the bottom of the column in which it was cited. Do not put footnotes in the abstract or reference list. Use letters for table footnotes.

Unless there are six authors or more give all authors’ names; do not use “et al.”. Papers that have not been published, even if they have been submitted for publication, should be cited as “unpublished” [4]. Papers that have been accepted for publication should be cited as “in press” [5]. Capitalize only the first word in a paper title, except for proper nouns and element symbols.

For papers published in translation journals, please give the English citation first, followed by the original foreign-language citation [6].

REFERENCES

- [1] G. Eason, B. Noble, and I. N. Sneddon, “On certain integrals of Lipschitz-Hankel type involving products of Bessel functions,” *Phil. Trans. Roy. Soc. London*, vol. A247, pp. 529–551, April 1955.
- [2] J. Clerk Maxwell, *A Treatise on Electricity and Magnetism*, 3rd ed., vol. 2. Oxford: Clarendon, 1892, pp.68–73.
- [3] I. S. Jacobs and C. P. Bean, “Fine particles, thin films and exchange anisotropy,” in *Magnetism*, vol. III, G. T. Rado and H. Suhl, Eds. New York: Academic, 1963, pp. 271–350.
- [4] K. Elissa, “Title of paper if known,” unpublished.
- [5] R. Nicole, “Title of paper with only first word capitalized,” *J. Name Stand. Abbrev.*, in press.
- [6] Y. Yorozu, M. Hirano, K. Oka, and Y. Tagawa, “Electron spectroscopy studies on magneto-optical media and plastic substrate interface,” *IEEE Transl. J. Magn. Japan*, vol. 2, pp. 740–741, August 1987 [Digests 9th Annual Conf. Magnetism Japan, p. 301, 1982].
- [7] M. Young, *The Technical Writer’s Handbook*. Mill Valley, CA: University Science, 1989.