

Fsck - The UNIX File System Check Program

Marshall Kirk McKusick

Computer Systems Research Group
Computer Science Division
Department of Electrical Engineering and Computer Science
University of California, Berkeley
Berkeley, CA 94720

T. J. Kowalski

Bell Laboratories
Murray Hill, New Jersey 07974

ABSTRACT

This document reflects the use of fsck with the 4.2BSD and 4.3BSD file system organization. This is a revision of the original paper written by T. J. Kowalski.

File System Check Program (fsck) is an interactive file system check and repair program. Fsck uses the redundant structural information in the UNIX file system to perform several consistency checks. If an inconsistency is detected, it is reported to the operator, who may elect to fix or ignore each inconsistency. These inconsistencies result from the permanent interruption of the file system updates, which are performed every time a file is modified. Unless there has been a hardware failure, fsck is able to repair corrupted file systems using procedures based upon the order in which UNIX honors these file system update requests.

The purpose of this document is to describe the normal updating of the file system, to discuss the possible causes of file system corruption, and

UNIX is a trademark of Bell Laboratories.
This work was done under grants from the National Science Foundation under grant MCS80-05144, and the Defense Advance Research Projects Agency (DoD) under Arpa Order No. 4031 monitored by Naval Electronic System Command under Contract No. N00039-82-C-0235.

to present the corrective actions implemented by fsck. Both the program and the interaction between the program and the operator are described.

Revised July 16, 1985

TABLE OF CONTENTS

- 1. Introduction
- 2. Overview of the file system
 - 2.1. Superblock
 - 2.2. Summary Information
 - 2.3. Cylinder groups
 - 2.4. Fragments
 - 2.5. Updates to the file system
- 3. Fixing corrupted file systems
 - 3.1. Detecting and correcting corruption
 - 3.2. Super block checking
 - 3.3. Free block checking
 - 3.4. Checking the inode state
 - 3.5. Inode links
 - 3.6. Inode data size
 - 3.7. Checking the data associated with an inode
 - 3.8. File system connectivity
- Acknowledgements
- References
- 4. Appendix A
 - 4.1. Conventions
 - 4.2. Initialization
 - 4.3. Phase 1 - Check Blocks and Sizes
 - 4.4. Phase 1b - Rescan for more Dups
 - 4.5. Phase 2 - Check Pathnames
 - 4.6. Phase 3 - Check Connectivity
 - 4.7. Phase 4 - Check Reference Counts
 - 4.8. Phase 5 - Check Cyl groups
 - 4.9. Cleanup

1. Introduction

This document reflects the use of fsck with the 4.2BSD and 4.3BSD file system organization. This is a revision of the original paper written by T. J. Kowalski.

When a UNIX operating system is brought up, a consistency check of the file systems should always be performed. This precautionary measure helps to insure a reliable environment for file storage on disk. If an inconsistency is discovered, corrective action must be taken. Fsck runs in two modes. Normally it is run non-interactively by the system after a normal boot. When running in this mode, it will only make changes to the file system that are known to always be correct. If an unexpected inconsistency is found fsck will exit with a non-zero exit status, leaving the system running single-user. Typically the operator then runs fsck interactively. When running in this mode, each problem is listed followed by a suggested corrective action. The operator must decide whether or not the suggested correction should be made.

The purpose of this memo is to dispel the mystique surrounding file system inconsistencies. It first describes the updating of the file system (the calm before the storm) and then describes file system corruption (the storm). Finally, the set of deterministic corrective actions used by fsck (the Coast Guard to the rescue) is presented.

2. Overview of the file system

The file system is discussed in detail in [Mckusick84]; this section gives a brief overview.

2.1. Superblock

A file system is described by its super-block. The super-block is built when the file system is created (newfs(8)) and never changes. The super-block contains the basic parameters of the file system, such as the number of data blocks it contains and a count of the maximum number of files. Because the super-block contains critical data, newfs replicates it to protect against catastrophic loss. The default super block always resides at a fixed offset from the beginning of the file system's disk partition. The redundant super blocks are not referenced unless a head crash or other hard disk error causes the default super-block to be unusable. The redundant blocks are sprinkled throughout the disk partition.

Within the file system are files. Certain files are distinguished as directories and contain collections of pointers to files that may themselves be directories. Every file has a descriptor associated with it called an inode.

The inode contains information describing ownership of the file, time stamps indicating modification and access times for the file, and an array of indices pointing to the data blocks for the file. In this section, we assume that the first 12 blocks of the file are directly referenced by values stored in the inode structure itself. The inode structure may also contain references to indirect blocks containing further data block indices. In a file system with a 4096 byte block size, a singly indirect block contains 1024 further block addresses, a doubly indirect block contains 1024 addresses of further single indirect blocks, and a triply indirect block contains 1024 addresses of further doubly indirect blocks (the triple indirect block is never needed in practice).

In order to create files with up to 2^{32} bytes, using only two levels of indirection, the minimum size of a file system block is 4096 bytes. The size of file system blocks can be any power of two greater than or equal to 4096. The block size of the file system is maintained in the super-block, so it is possible for file systems of different block sizes to be accessible simultaneously on the same system. The block size must be decided when newfs creates the file system; the block size cannot be subsequently changed without rebuilding the file system.

2.2. Summary information

Associated with the super block is non replicated summary information. The summary information changes as the file system is modified. The summary information contains the number of blocks, fragments, inodes and directories in the file system.

2.3. Cylinder groups

The file system partitions the disk into one or more areas called cylinder groups. A cylinder group is comprised of one or more consecutive cylinders on a disk. Each cylinder group includes inode slots for files, a block map describing available blocks in the cylinder group, and summary information describing the usage of data blocks within the cylinder group. A fixed number of inodes is allocated for each cylinder group when the file system is created. The current policy is to allocate one inode for each 2048 bytes of disk space; this is expected to be far more inodes than will ever be needed.

All the cylinder group bookkeeping information could be placed at the beginning of each cylinder group. However if

The actual number may vary from system to system, but is usually in the range 5-13.

this approach were used, all the redundant information would be on the top platter. A single hardware failure that destroyed the top platter could cause the loss of all copies of the redundant super-blocks. Thus the cylinder group book-keeping information begins at a floating offset from the beginning of the cylinder group. The offset for the $i+1$ st cylinder group is about one track further from the beginning of the cylinder group than it was for the i th cylinder group. In this way, the redundant information spirals down into the pack; any single track, cylinder, or platter can be lost without losing all copies of the super-blocks. Except for the first cylinder group, the space between the beginning of the cylinder group and the beginning of the cylinder group information stores data.

2.4. Fragments

To avoid waste in storing small files, the file system space allocator divides a single file system block into one or more fragments. The fragmentation of the file system is specified when the file system is created; each file system block can be optionally broken into 2, 4, or 8 addressable fragments. The lower bound on the size of these fragments is constrained by the disk sector size; typically 512 bytes is the lower bound on fragment size. The block map associated with each cylinder group records the space availability at the fragment level. Aligned fragments are examined to determine block availability.

On a file system with a block size of 4096 bytes and a fragment size of 1024 bytes, a file is represented by zero or more 4096 byte blocks of data, and possibly a single fragmented block. If a file system block must be fragmented to obtain space for a small amount of data, the remainder of the block is made available for allocation to other files. For example, consider an 11000 byte file stored on a 4096/1024 byte file system. This file uses two full size blocks and a 3072 byte fragment. If no fragments with at least 3072 bytes are available when the file is created, a full size block is split yielding the necessary 3072 byte fragment and an unused 1024 byte fragment. This remaining fragment can be allocated to another file, as needed.

2.5. Updates to the file system

Every working day hundreds of files are created, modified, and removed. Every time a file is modified, the operating system performs a series of file system updates. These updates, when written on disk, yield a consistent file system. The file system stages all modifications of critical information; modification can either be completed or cleanly backed out after a crash. Knowing the information that is first written to the file system, deterministic procedures can be developed to repair a corrupted file system.

To understand this process, the order that the update requests were being honored must first be understood.

When a user program does an operation to change the file system, such as a write, the data to be written is copied into an internal in-core buffer in the kernel. Normally, the disk update is handled asynchronously; the user process is allowed to proceed even though the data has not yet been written to the disk. The data, along with the inode information reflecting the change, is eventually written out to disk. The real disk write may not happen until long after the write system call has returned. Thus at any given time, the file system, as it resides on the disk, lags the state of the file system represented by the in-core information.

The disk information is updated to reflect the in-core information when the buffer is required for another use, when a sync(2) is done (at 30 second intervals) by update(8), or by manual operator intervention with the sync(8) command. If the system is halted without writing out the in-core information, the file system on the disk will be in an inconsistent state.

If all updates are done asynchronously, several serious inconsistencies can arise. One inconsistency is that a block may be claimed by two inodes. Such an inconsistency can occur when the system is halted before the pointer to the block in the old inode has been cleared in the copy of the old inode on the disk, and after the pointer to the block in the new inode has been written out to the copy of the new inode on the disk. Here, there is no deterministic method for deciding which inode should really claim the block. A similar problem can arise with a multiply claimed inode.

The problem with asynchronous inode updates can be avoided by doing all inode deallocations synchronously. Consequently, inodes and indirect blocks are written to the disk synchronously (i.e. the process blocks until the information is really written to disk) when they are being deallocated. Similarly inodes are kept consistent by synchronously deleting, adding, or changing directory entries.

3. Fixing corrupted file systems

A file system can become corrupted in several ways. The most common of these ways are improper shutdown procedures and hardware failures.

File systems may become corrupted during an unclean halt. This happens when proper shutdown procedures are not observed, physically write-protecting a mounted file system, or a mounted file system is taken off-line. The most common

operator procedural failure is forgetting to sync the system before halting the CPU.

File systems may become further corrupted if proper startup procedures are not observed, e.g., not checking a file system for inconsistencies, and not repairing inconsistencies. Allowing a corrupted file system to be used (and, thus, to be modified further) can be disastrous.

Any piece of hardware can fail at any time. Failures can be as subtle as a bad block on a disk pack, or as blatant as a non-functional disk-controller.

3.1. Detecting and correcting corruption

Normally fsck is run non-interactively. In this mode it will only fix corruptions that are expected to occur from an unclean halt. These actions are a proper subset of the actions that fsck will take when it is running interactively. Throughout this paper we assume that fsck is being run interactively, and all possible errors can be encountered. When an inconsistency is discovered in this mode, fsck reports the inconsistency for the operator to choose a corrective action.

A quiescent^{|=} file system may be checked for structural integrity by performing consistency checks on the redundant data intrinsic to a file system. The redundant data is either read from the file system, or computed from other known values. The file system must be in a quiescent state when fsck is run, since fsck is a multi-pass program.

In the following sections, we discuss methods to discover inconsistencies and possible corrective actions for the cylinder group blocks, the inodes, the indirect blocks, and the data blocks containing directory entries.

3.2. Super-block checking

The most commonly corrupted item in a file system is the summary information associated with the super-block. The summary information is prone to corruption because it is modified with every change to the file system's blocks or inodes, and is usually corrupted after an unclean halt.

The super-block is checked for inconsistencies involving file-system size, number of inodes, free-block count, and the free-inode count. The file-system size must be larger than the number of blocks used by the super-block and the number of blocks used by the list of inodes. The file-system size and layout information are the most critical

^{|=} I.e., unmounted and not being written on.

pieces of information for fsck. While there is no way to actually check these sizes, since they are statically determined by newfs, fsck can check that these sizes are within reasonable bounds. All other file system checks require that these sizes be correct. If fsck detects corruption in the static parameters of the default super-block, fsck requests the operator to specify the location of an alternate super-block.

3.3. Free block checking

Fsck checks that all the blocks marked as free in the cylinder group block maps are not claimed by any files. When all the blocks have been initially accounted for, fsck checks that the number of free blocks plus the number of blocks claimed by the inodes equals the total number of blocks in the file system.

If anything is wrong with the block allocation maps, fsck will rebuild them, based on the list it has computed of allocated blocks.

The summary information associated with the super-block counts the total number of free blocks within the file system. Fsck compares this count to the number of free blocks it found within the file system. If the two counts do not agree, then fsck replaces the incorrect count in the summary information by the actual free-block count.

The summary information counts the total number of free inodes within the file system. Fsck compares this count to the number of free inodes it found within the file system. If the two counts do not agree, then fsck replaces the incorrect count in the summary information by the actual free-inode count.

3.4. Checking the inode state

An individual inode is not as likely to be corrupted as the allocation information. However, because of the great number of active inodes, a few of the inodes are usually corrupted.

The list of inodes in the file system is checked sequentially starting with inode 2 (inode 0 marks unused inodes; inode 1 is saved for future generations) and progressing through the last inode in the file system. The state of each inode is checked for inconsistencies involving format and type, link count, duplicate blocks, bad blocks, and inode size.

Each inode contains a mode word. This mode word describes the type and state of the inode. Inodes must be one of six types: regular inode, directory inode, symbolic

link inode, special block inode, special character inode, or socket inode. Inodes may be found in one of three allocation states: unallocated, allocated, and neither unallocated nor allocated. This last state suggests an incorrectly formatted inode. An inode can get in this state if bad data is written into the inode list. The only possible corrective action is for fsck is to clear the inode.

3.5. Inode links

Each inode counts the total number of directory entries linked to the inode. Fsck verifies the link count of each inode by starting at the root of the file system, and descending through the directory structure. The actual link count for each inode is calculated during the descent.

If the stored link count is non-zero and the actual link count is zero, then no directory entry appears for the inode. If this happens, fsck will place the disconnected file in the lost+found directory. If the stored and actual link counts are non-zero and unequal, a directory entry may have been added or removed without the inode being updated. If this happens, fsck replaces the incorrect stored link count by the actual link count.

Each inode contains a list, or pointers to lists (indirect blocks), of all the blocks claimed by the inode. Since indirect blocks are owned by an inode, inconsistencies in indirect blocks directly affect the inode that owns it.

Fsck compares each block number claimed by an inode against a list of already allocated blocks. If another inode already claims a block number, then the block number is added to a list of duplicate blocks. Otherwise, the list of allocated blocks is updated to include the block number.

If there are any duplicate blocks, fsck will perform a partial second pass over the inode list to find the inode of the duplicated block. The second pass is needed, since without examining the files associated with these inodes for correct content, not enough information is available to determine which inode is corrupted and should be cleared. If this condition does arise (only hardware failure will cause it), then the inode with the earliest modify time is usually incorrect, and should be cleared. If this happens, fsck prompts the operator to clear both inodes. The operator must decide which one should be kept and which one should be cleared.

Fsck checks the range of each block number claimed by an inode. If the block number is lower than the first data block in the file system, or greater than the last data block, then the block number is a bad block number. Many bad blocks in an inode are usually caused by an indirect

block that was not written to the file system, a condition which can only occur if there has been a hardware failure. If an inode contains bad block numbers, fsck prompts the operator to clear it.

3.6. Inode data size

Each inode contains a count of the number of data blocks that it contains. The number of actual data blocks is the sum of the allocated data blocks and the indirect blocks. Fsck computes the actual number of data blocks and compares that block count against the actual number of blocks the inode claims. If an inode contains an incorrect count fsck prompts the operator to fix it.

Each inode contains a thirty-two bit size field. The size is the number of data bytes in the file associated with the inode. The consistency of the byte size field is roughly checked by computing from the size field the maximum number of blocks that should be associated with the inode, and comparing that expected block count against the actual number of blocks the inode claims.

3.7. Checking the data associated with an inode

An inode can directly or indirectly reference three kinds of data blocks. All referenced blocks must be the same kind. The three types of data blocks are: plain data blocks, symbolic link data blocks, and directory data blocks. Plain data blocks contain the information stored in a file; symbolic link data blocks contain the path name stored in a link. Directory data blocks contain directory entries. Fsck can only check the validity of directory data blocks.

Each directory data block is checked for several types of inconsistencies. These inconsistencies include directory inode numbers pointing to unallocated inodes, directory inode numbers that are greater than the number of inodes in the file system, incorrect directory inode numbers for ``.`' and ``..'', and directories that are not attached to the file system. If the inode number in a directory data block references an unallocated inode, then fsck will remove that directory entry. Again, this condition can only arise when there has been a hardware failure.

If a directory entry inode number references outside the inode list, then fsck will remove that directory entry. This condition occurs if bad data is written into a directory data block.

The directory inode number entry for ``.`' must be the first entry in the directory data block. The inode number for ``.`' must reference itself; e.g., it must equal the

inode number for the directory data block. The directory inode number entry for ``..' must be the second entry in the directory data block. Its value must equal the inode number for the parent of the directory entry (or the inode number of the directory data block if the directory is the root directory). If the directory inode numbers are incorrect, fsck will replace them with the correct values. If there are multiple hard links to a directory, the first one encountered is considered the real parent to which ``..' should point; fsckP recommends deletion for the subsequently discovered names.

3.8. File system connectivity

Fsck checks the general connectivity of the file system. If directories are not linked into the file system, then fsck links the directory back into the file system in the lost+found directory. This condition only occurs when there has been a hardware failure.

Acknowledgements

I thank Bill Joy, Sam Leffler, Robert Elz and Dennis Ritchie for their suggestions and help in implementing the new file system. Thanks also to Robert Henry for his editorial input to get this document together. Finally we thank our sponsors, the National Science Foundation under grant MCS80-05144, and the Defense Advance Research Projects Agency (DoD) under Arpa Order No. 4031 monitored by Naval Electronic System Command under Contract No. N00039-82-C-0235. (Kirk McKusick, July 1983)

I would like to thank Larry A. Wehr for advice that lead to the first version of fsck and Rick B. Brandt for adapting fsck to UNIX/TS. (T. Kowalski, July 1979)

References

- [Dolotta78] Dolotta, T. A., and Olsson, S. B. eds., UNIX User's Manual, Edition 1.1, January 1978.
- [Joy83] Joy, W., Cooper, E., Fabry, R., Leffler, S., McKusick, M., and Mosher, D. 4.2BSD System Manual, University of California at Berkeley, Computer Systems Research Group Technical Report #4, 1982.
- [McKusick84] McKusick, M., Joy, W., Leffler, S., and Fabry, R. A Fast File System for UNIX, ACM Transactions on Computer Systems 2,

3. pp. 181-197, August 1984.

[Ritchie78] Ritchie, D. M., and Thompson, K., The UNIX Time-Sharing System, The Bell System Technical Journal 57, 6 (July-August 1978, Part 2), pp. 1905-29.

[Thompson78] Thompson, K., UNIX Implementation, The Bell System Technical Journal 57, 6 (July-August 1978, Part 2), pp. 1931-46.

4. Appendix A - Fsck Error Conditions

4.1. Conventions

Fsck is a multi-pass file system check program. Each file system pass invokes a different Phase of the fsck program. After the initial setup, fsck performs successive Phases over each file system, checking blocks and sizes, path-names, connectivity, reference counts, and the map of free blocks, (possibly rebuilding it), and performs some cleanup.

Normally fsck is run non-interactively to preen the file systems after an unclean halt. While preen'ing a file system, it will only fix corruptions that are expected to occur from an unclean halt. These actions are a proper subset of the actions that fsck will take when it is running interactively. Throughout this appendix many errors have several options that the operator can take. When an inconsistency is detected, fsck reports the error condition to the operator. If a response is required, fsck prints a prompt message and waits for a response. When preen'ing most errors are fatal. For those that are expected, the response taken is noted. This appendix explains the meaning of each error condition, the possible responses, and the related error conditions.

The error conditions are organized by the Phase of the fsck program in which they can occur. The error conditions that may occur in more than one Phase will be discussed in initialization.

4.2. Initialization

Before a file system check can be performed, certain tables have to be set up and certain files opened. This section concerns itself with the opening of files and the initialization of tables. This section lists error conditions resulting from command line options, memory requests, opening of files, status of files, file system size checks, and creation of the scratch file. All the initialization errors are fatal when the file system is being preen'ed.

C option?

C is not a legal option to fsck; legal options are -b, -y, -n, and -p. Fsck terminates on this error condition. See the fsck(8) manual entry for further detail.

cannot alloc NNN bytes for blockmap
cannot alloc NNN bytes for freemap
cannot alloc NNN bytes for statemap
cannot alloc NNN bytes for lncntp

Fsck's request for memory for its virtual memory tables failed. This should never happen. Fsck terminates on this error condition. See a guru.

Can't open checklist file: F

The file system checklist file F (usually /etc/fstab) can not be opened for reading. Fsck terminates on this error condition. Check access modes of F.

Can't stat root

Fsck's request for statistics about the root directory ``/' failed. This should never happen. Fsck terminates on this error condition. See a guru.

Can't stat F

Can't make sense out of name F

Fsck's request for statistics about the file system F failed. When running manually, it ignores this file system and continues checking the next file system given. Check access modes of F.

Can't open F

Fsck's request attempt to open the file system F failed. When running manually, it ignores this file system and continues checking the next file system given. Check access modes of F.

F: (NO WRITE)

Either the -n flag was specified or fsck's attempt to open the file system F for writing failed. When running manually, all the diagnostics are printed out, but no modifications are attempted to fix them.

file is not a block or character device; OK

You have given fsck a regular file name by mistake. Check the type of the file specified.

Possible responses to the OK prompt are:

YES ignore this error condition.

NO ignore this file system and continues checking the next file system given.

UNDEFINED OPTIMIZATION IN SUPERBLOCK (SET TO DEFAULT)

The superbblock optimization parameter is neither OPT_TIME nor OPT_SPACE.

Possible responses to the SET TO DEFAULT prompt are:

YES The superblock is set to request optimization to minimize running time of the system. (If optimization to minimize disk space utilization is desired, it can be set using tunefs(8).)

NO ignore this error condition.

IMPOSSIBLE MINFREE=D IN SUPERBLOCK (SET TO DEFAULT)
The superblock minimum space percentage is greater than 99% or less than 0%.

Possible responses to the SET TO DEFAULT prompt are:

YES The minfree parameter is set to 10%. (If some other percentage is desired, it can be set using tunefs(8).)

NO ignore this error condition.

One of the following messages will appear:

MAGIC NUMBER WRONG

NCG OUT OF RANGE

CPG OUT OF RANGE

NCYL DOES NOT JIVE WITH NCG*CPG

SIZE PREPOSTEROUSLY LARGE

TRASHED VALUES IN SUPER BLOCK

and will be followed by the message:

F: BAD SUPER BLOCK: B

USE -b OPTION TO FSCK TO SPECIFY LOCATION OF AN ALTERNATE
SUPER-BLOCK TO SUPPLY NEEDED INFORMATION; SEE fsck(8).

The super block has been corrupted. An alternative super block must be selected from among those listed by newfs (8) when the file system was created. For file systems with a blocksize less than 32K, specifying -b 32 is a good first choice.

INTERNAL INCONSISTENCY: M

Fsck's has had an internal panic, whose message is specified as M. This should never happen. See a guru.

CAN NOT SEEK: BLK B (CONTINUE)

Fsck's request for moving to a specified block number B in the file system failed. This should never happen. See a guru.

Possible responses to the CONTINUE prompt are:

YES attempt to continue to run the file system check.
Often, however the problem will persist. This error

condition will not allow a complete check of the file system. A second run of fsck should be made to re-check this file system. If the block was part of the virtual memory buffer cache, fsck will terminate with the message ``Fatal I/O error''.

NO terminate the program.

CAN NOT READ: BLK B (CONTINUE)

Fsck's request for reading a specified block number B in the file system failed. This should never happen. See a guru.

Possible responses to the CONTINUE prompt are:

YES attempt to continue to run the file system check. It will retry the read and print out the message:
THE FOLLOWING SECTORS COULD NOT BE READ: N
where N indicates the sectors that could not be read. If fsck ever tries to write back one of the blocks on which the read failed it will print the message:
WRITING ZERO'ED BLOCK N TO DISK
where N indicates the sector that was written with zero's. If the disk is experiencing hardware problems, the problem will persist. This error condition will not allow a complete check of the file system. A second run of fsck should be made to re-check this file system. If the block was part of the virtual memory buffer cache, fsck will terminate with the message ``Fatal I/O error''.

NO terminate the program.

CAN NOT WRITE: BLK B (CONTINUE)

Fsck's request for writing a specified block number B in the file system failed. The disk is write-protected; check the write protect lock on the drive. If that is not the problem, see a guru.

Possible responses to the CONTINUE prompt are:

YES attempt to continue to run the file system check. The write operation will be retried with the failed blocks indicated by the message:
THE FOLLOWING SECTORS COULD NOT BE WRITTEN: N
where N indicates the sectors that could not be written. If the disk is experiencing hardware problems, the problem will persist. This error condition will not allow a complete check of the file system. A second run of fsck should be made to re-check this file system. If the block was part of the virtual memory buffer cache, fsck will terminate with the message ``Fatal I/O error''.

NO terminate the program.

bad inode number DDD to ginode

An internal error has attempted to read non-existent inode DDD. This error causes fsck to exit. See a guru.

4.3. Phase 1 - Check Blocks and Sizes

This phase concerns itself with the inode list. This section lists error conditions resulting from checking inode types, setting up the zero-link-count table, examining inode block numbers for bad or duplicate blocks, checking inode size, and checking inode format. All errors in this phase except INCORRECT BLOCK COUNT and PARTIALLY TRUNCATED INODE are fatal if the file system is being preen'ed.

UNKNOWN FILE TYPE I=I (CLEAR)

The mode word of the inode I indicates that the inode is not a special block inode, special character inode, socket inode, regular inode, symbolic link, or directory inode.

Possible responses to the CLEAR prompt are:

YES de-allocate inode I by zeroing its contents. This will always invoke the UNALLOCATED error condition in Phase 2 for each directory entry pointing to this inode.

NO ignore this error condition.

PARTIALLY TRUNCATED INODE I=I (SALVAGE)

Fsck has found inode I whose size is shorter than the number of blocks allocated to it. This condition should only occur if the system crashes while in the midst of truncating a file. When preen'ing the file system, fsck completes the truncation to the specified size.

Possible responses to SALVAGE are:

YES complete the truncation to the size specified in the inode.

NO ignore this error condition.

LINK COUNT TABLE OVERFLOW (CONTINUE)

An internal table for fsck containing allocated inodes with a link count of zero cannot allocate more memory. Increase the virtual memory for fsck.

Possible responses to the CONTINUE prompt are:

YES continue with the program. This error condition will not allow a complete check of the file system. A second run of fsck should be made to re-check this file system. If another allocated inode with a zero link count is found, this error condition is repeated.

NO terminate the program.

B BAD I=I

Inode I contains block number B with a number lower than the number of the first data block in the file system or greater than the number of the last block in the file system. This error condition may invoke the EXCESSIVE BAD BLKS error condition in Phase 1 (see next paragraph) if inode I has too many block numbers outside the file system range. This error condition will always invoke the BAD/DUP error condition in Phase 2 and Phase 4.

EXCESSIVE BAD BLKS I=I (CONTINUE)

There is more than a tolerable number (usually 10) of blocks with a number lower than the number of the first data block in the file system or greater than the number of last block in the file system associated with inode I.

Possible responses to the CONTINUE prompt are:

YES ignore the rest of the blocks in this inode and continue checking with the next inode in the file system. This error condition will not allow a complete check of the file system. A second run of fsck should be made to re-check this file system.

NO terminate the program.

BAD STATE DDD TO BLKERR

An internal error has scrambled fsck's state map to have the impossible value DDD. Fsck exits immediately. See a guru.

B DUP I=I

Inode I contains block number B that is already claimed by another inode. This error condition may invoke the EXCESSIVE DUP BLKS error condition in Phase 1 if inode I has too many block numbers claimed by other inodes. This error condition will always invoke Phase 1b and the BAD/DUP error condition in Phase 2 and Phase 4.

EXCESSIVE DUP BLKS I=I (CONTINUE)

There is more than a tolerable number (usually 10) of blocks claimed by other inodes.

Possible responses to the CONTINUE prompt are:

YES ignore the rest of the blocks in this inode and continue checking with the next inode in the file system. This error condition will not allow a complete check of the file system. A second run of fsck should be made to re-check this file system.

NO terminate the program.

DUP TABLE OVERFLOW (CONTINUE)

An internal table in fsck containing duplicate block numbers cannot allocate any more space. Increase the amount of virtual memory available to fsck.

Possible responses to the CONTINUE prompt are:

YES continue with the program. This error condition will not allow a complete check of the file system. A second run of fsck should be made to re-check this file system. If another duplicate block is found, this error condition will repeat.

NO terminate the program.

PARTIALLY ALLOCATED INODE I=I (CLEAR)

Inode I is neither allocated nor unallocated.

Possible responses to the CLEAR prompt are:

YES de-allocate inode I by zeroing its contents.

NO ignore this error condition.

INCORRECT BLOCK COUNT I=I (X should be Y) (CORRECT)

The block count for inode I is X blocks, but should be Y blocks. When preen'ing the count is corrected.

Possible responses to the CORRECT prompt are:

YES replace the block count of inode I with Y.

NO ignore this error condition.

4.4. Phase 1B: Rescan for More Dups

When a duplicate block is found in the file system, the file system is rescanned to find the inode that previously claimed that block. This section lists the error condition when the duplicate block is found.

B DUP I=I

Inode I contains block number B that is already claimed by another inode. This error condition will always invoke the BAD/DUP error condition in Phase 2. You can determine which inodes have overlapping blocks by examining this error condition and the DUP error condition in Phase 1.

4.5. Phase 2 - Check Pathnames

This phase concerns itself with removing directory entries pointing to error conditioned inodes from Phase 1 and Phase 1b. This section lists error conditions resulting from root inode mode and status, directory inode pointers in range, and directory entries pointing to bad inodes, and directory integrity checks. All errors in this phase are fatal if the file system is being preen'ed, except for directories not being a multiple of the blocks size and extraneous hard links.

ROOT INODE UNALLOCATED (ALLOCATE)

The root inode (usually inode number 2) has no allocate mode bits. This should never happen.

Possible responses to the ALLOCATE prompt are:

YES allocate inode 2 as the root inode. The files and directories usually found in the root will be recovered in Phase 3 and put into lost+found. If the attempt to allocate the root fails, fsck will exit with the message:

CANNOT ALLOCATE ROOT INODE.

NO fsck will exit.

ROOT INODE NOT DIRECTORY (REALLOCATE)

The root inode (usually inode number 2) is not directory inode type.

Possible responses to the REALLOCATE prompt are:

YES clear the existing contents of the root inode and reallocate it. The files and directories usually found in the root will be recovered in Phase 3 and put into lost+found. If the attempt to allocate the root fails, fsck will exit with the message:

CANNOT ALLOCATE ROOT INODE.

NO fsck will then prompt with FIX

Possible responses to the FIX prompt are:

YES replace the root inode's type to be a directory. If

the root inode's data blocks are not directory blocks, many error conditions will be produced.

NO terminate the program.

DUPS/BAD IN ROOT INODE (REALLOCATE)

Phase 1 or Phase 1b have found duplicate blocks or bad blocks in the root inode (usually inode number 2) for the file system.

Possible responses to the REALLOCATE prompt are:

YES clear the existing contents of the root inode and reallocate it. The files and directories usually found in the root will be recovered in Phase 3 and put into lost+found. If the attempt to allocate the root fails, fsck will exit with the message:
CANNOT ALLOCATE ROOT INODE.

NO fsck will then prompt with CONTINUE.

Possible responses to the CONTINUE prompt are:

YES ignore the DUPS/BAD error condition in the root inode and attempt to continue to run the file system check. If the root inode is not correct, then this may result in many other error conditions.

NO terminate the program.

NAME TOO LONG F

An excessively long path name has been found. This usually indicates loops in the file system name space. This can occur if the super user has made circular links to directories. The offending links must be removed (by a guru).

I OUT OF RANGE I=I NAME=F (REMOVE)

A directory entry F has an inode number I that is greater than the end of the inode list.

Possible responses to the REMOVE prompt are:

YES the directory entry F is removed.

NO ignore this error condition.

UNALLOCATED I=I OWNER=O MODE=M SIZE=S MTIME=T type=F (REMOVE)

A directory or file entry F points to an unallocated inode I. The owner O, mode M, size S, modify time T, and name F

are printed.

Possible responses to the REMOVE prompt are:

YES the directory entry F is removed.

NO ignore this error condition.

DUP/BAD I=I OWNER=O MODE=M SIZE=S MTIME=T type=F (REMOVE)
Phase 1 or Phase 1b have found duplicate blocks or bad blocks associated with directory or file entry F, inode I. The owner O, mode M, size S, modify time T, and directory name F are printed.

Possible responses to the REMOVE prompt are:

YES the directory entry F is removed.

NO ignore this error condition.

ZERO LENGTH DIRECTORY I=I OWNER=O MODE=M SIZE=S MTIME=T DIR=F (REMOVE)
A directory entry F has a size S that is zero. The owner O, mode M, size S, modify time T, and directory name F are printed.

Possible responses to the REMOVE prompt are:

YES the directory entry F is removed; this will always invoke the BAD/DUP error condition in Phase 4.

NO ignore this error condition.

DIRECTORY TOO SHORT I=I OWNER=O MODE=M SIZE=S MTIME=T DIR=F (FIX)
A directory F has been found whose size S is less than the minimum size directory. The owner O, mode M, size S, modify time T, and directory name F are printed.

Possible responses to the FIX prompt are:

YES increase the size of the directory to the minimum directory size.

NO ignore this directory.

DIRECTORY F LENGTH S NOT MULTIPLE OF B (ADJUST)
A directory F has been found with size S that is not a multiple of the directory blocksize B.

Possible responses to the ADJUST prompt are:

YES the length is rounded up to the appropriate block size. This error can occur on 4.2BSD file systems. Thus when preen'ing the file system only a warning is printed and the directory is adjusted.

NO ignore the error condition.

DIRECTORY CORRUPTED I=I OWNER=O MODE=M SIZE=S MTIME=T DIR=F (SALVAGE)

A directory with an inconsistent internal state has been found.

Possible responses to the FIX prompt are:

YES throw away all entries up to the next directory boundary (usually 512-byte) boundary. This drastic action can throw away up to 42 entries, and should be taken only after other recovery efforts have failed.

NO skip up to the next directory boundary and resume reading, but do not modify the directory.

BAD INODE NUMBER FOR `.' I=I OWNER=O MODE=M SIZE=S MTIME=T DIR=F (FIX)

A directory I has been found whose inode number for `.' does not equal I.

Possible responses to the FIX prompt are:

YES change the inode number for `.' to be equal to I.

NO leave the inode number for `.' unchanged.

MISSING `.' I=I OWNER=O MODE=M SIZE=S MTIME=T DIR=F (FIX)

A directory I has been found whose first entry is unallocated.

Possible responses to the FIX prompt are:

YES build an entry for `.' with inode number equal to I.

NO leave the directory unchanged.

MISSING `.' I=I OWNER=O MODE=M SIZE=S MTIME=T DIR=F
CANNOT FIX, FIRST ENTRY IN DIRECTORY CONTAINS F

A directory I has been found whose first entry is F. Fsck cannot resolve this problem. The file system should be mounted and the offending entry F moved elsewhere. The file

system should then be unmounted and fsck should be run again.

MISSING `.' I=I OWNER=O MODE=M SIZE=S MTIME=T DIR=F
CANNOT FIX, INSUFFICIENT SPACE TO ADD `.'
A directory I has been found whose first entry is not `.'. Fsck cannot resolve this problem as it should never happen. See a guru.

EXTRA `.' ENTRY I=I OWNER=O MODE=M SIZE=S MTIME=T DIR=F
(FIX)
A directory I has been found that has more than one entry for `.'.

Possible responses to the FIX prompt are:

YES remove the extra entry for `.'.

NO leave the directory unchanged.

BAD INODE NUMBER FOR `..' I=I OWNER=O MODE=M SIZE=S MTIME=T
DIR=F (FIX)
A directory I has been found whose inode number for `..' does does not equal the parent of I.

Possible responses to the FIX prompt are:

YES change the inode number for `..' to be equal to the parent of I (`..' in the root inode points to itself).

NO leave the inode number for `..' unchanged.

MISSING `..' I=I OWNER=O MODE=M SIZE=S MTIME=T DIR=F (FIX)
A directory I has been found whose second entry is unallocated.

Possible responses to the FIX prompt are:

YES build an entry for `..' with inode number equal to the parent of I (`..' in the root inode points to itself).

NO leave the directory unchanged.

MISSING `..' I=I OWNER=O MODE=M SIZE=S MTIME=T DIR=F
CANNOT FIX, SECOND ENTRY IN DIRECTORY CONTAINS F
A directory I has been found whose second entry is F. Fsck cannot resolve this problem. The file system should be

mounted and the offending entry F moved elsewhere. The file system should then be unmounted and fsck should be run again.

MISSING `..' I=I OWNER=O MODE=M SIZE=S MTIME=T DIR=F
CANNOT FIX, INSUFFICIENT SPACE TO ADD `..'

A directory I has been found whose second entry is not `..'. Fsck cannot resolve this problem. The file system should be mounted and the second entry in the directory moved elsewhere. The file system should then be unmounted and fsck should be run again.

EXTRA `..' ENTRY I=I OWNER=O MODE=M SIZE=S MTIME=T DIR=F
(FIX)

A directory I has been found that has more than one entry for `..'.

Possible responses to the FIX prompt are:

YES remove the extra entry for `..'.

NO leave the directory unchanged.

N IS AN EXTRANEIOUS HARD LINK TO A DIRECTORY D (REMOVE)

Fsck has found a hard link, N, to a directory, D. When preen'ing the extraneous links are ignored.

Possible responses to the REMOVE prompt are:

YES delete the extraneous entry, N.

NO ignore the error condition.

BAD INODE S TO DESCEND

An internal error has caused an impossible state S to be passed to the routine that descends the file system directory structure. Fsck exits. See a guru.

BAD RETURN STATE S FROM DESCEND

An internal error has caused an impossible state S to be returned from the routine that descends the file system directory structure. Fsck exits. See a guru.

BAD STATE S FOR ROOT INODE

An internal error has caused an impossible state S to be assigned to the root inode. Fsck exits. See a guru.

4.6. Phase 3 - Check Connectivity

This phase concerns itself with the directory connectivity seen in Phase 2. This section lists error conditions resulting from unreferenced directories, and missing or full lost+found directories.

UNREF DIR I=I OWNER=O MODE=M SIZE=S MTIME=T (RECONNECT)

The directory inode I was not connected to a directory entry when the file system was traversed. The owner O, mode M, size S, and modify time T of directory inode I are printed. When preen'ing, the directory is reconnected if its size is non-zero, otherwise it is cleared.

Possible responses to the RECONNECT prompt are:

YES reconnect directory inode I to the file system in the directory for lost files (usually lost+found). This may invoke the lost+found error condition in Phase 3 if there are problems connecting directory inode I to lost+found. This may also invoke the CONNECTED error condition in Phase 3 if the link was successful.

NO ignore this error condition. This will always invoke the UNREF error condition in Phase 4.

NO lost+found DIRECTORY (CREATE)

There is no lost+found directory in the root directory of the file system; When preen'ing fsck tries to create a lost+found directory.

Possible responses to the CREATE prompt are:

YES create a lost+found directory in the root of the file system. This may raise the message:

NO SPACE LEFT IN / (EXPAND)

See below for the possible responses. Inability to create a lost+found directory generates the message:

SORRY. CANNOT CREATE lost+found DIRECTORY

and aborts the attempt to linkup the lost inode. This will always invoke the UNREF error condition in Phase 4.

NO abort the attempt to linkup the lost inode. This will always invoke the UNREF error condition in Phase 4.

lost+found IS NOT A DIRECTORY (REALLOCATE)

The entry for lost+found is not a directory.

Possible responses to the REALLOCATE prompt are:

YES allocate a directory inode, and change lost+found to reference it. The previous inode reference by the lost+found name is not cleared. Thus it will either be reclaimed as an UNREF'ed inode or have its link count ADJUST'ed later in this Phase. Inability to create a lost+found directory generates the message:
 SORRY. CANNOT CREATE lost+found DIRECTORY
 and aborts the attempt to linkup the lost inode. This will always invoke the UNREF error condition in Phase 4.

NO abort the attempt to linkup the lost inode. This will always invoke the UNREF error condition in Phase 4.

NO SPACE LEFT IN /lost+found (EXPAND)

There is no space to add another entry to the lost+found directory in the root directory of the file system. When preen'ing the lost+found directory is expanded.

Possible responses to the EXPAND prompt are:

YES the lost+found directory is expanded to make room for the new entry. If the attempted expansion fails fsck prints the message:
 SORRY. NO SPACE IN lost+found DIRECTORY
 and aborts the attempt to linkup the lost inode. This will always invoke the UNREF error condition in Phase 4. Clean out unnecessary entries in lost+found. This error is fatal if the file system is being preen'ed.

NO abort the attempt to linkup the lost inode. This will always invoke the UNREF error condition in Phase 4.

DIR I=I1 CONNECTED. PARENT WAS I=I2

This is an advisory message indicating a directory inode I1 was successfully connected to the lost+found directory. The parent inode I2 of the directory inode I1 is replaced by the inode number of the lost+found directory.

DIRECTORY F LENGTH S NOT MULTIPLE OF B (ADJUST)

A directory F has been found with size S that is not a multiple of the directory blocksize B (this can reoccur in Phase 3 if it is not adjusted in Phase 2).

Possible responses to the ADJUST prompt are:

YES the length is rounded up to the appropriate block size. This error can occur on 4.2BSD file systems. Thus when preen'ing the file system only a warning is printed and the directory is adjusted.

NO ignore the error condition.

BAD INODE S TO DESCEND

An internal error has caused an impossible state S to be passed to the routine that descends the file system directory structure. Fsck exits. See a guru.

4.7. Phase 4 - Check Reference Counts

This phase concerns itself with the link count information seen in Phase 2 and Phase 3. This section lists error conditions resulting from unreferenced files, missing or full lost+found directory, incorrect link counts for files, directories, symbolic links, or special files, unreferenced files, symbolic links, and directories, and bad or duplicate blocks in files, symbolic links, and directories. All errors in this phase are correctable if the file system is being preen'ed except running out of space in the lost+found directory.

UNREF FILE I=I OWNER=O MODE=M SIZE=S MTIME=T (RECONNECT)

Inode I was not connected to a directory entry when the file system was traversed. The owner O, mode M, size S, and modify time T of inode I are printed. When preen'ing the file is cleared if either its size or its link count is zero, otherwise it is reconnected.

Possible responses to the RECONNECT prompt are:

YES reconnect inode I to the file system in the directory for lost files (usually lost+found). This may invoke the lost+found error condition in Phase 4 if there are problems connecting inode I to lost+found.

NO ignore this error condition. This will always invoke the CLEAR error condition in Phase 4.

(CLEAR)

The inode mentioned in the immediately previous error condition can not be reconnected. This cannot occur if the file system is being preen'ed, since lack of space to reconnect files is a fatal error.

Possible responses to the CLEAR prompt are:

YES de-allocate the inode mentioned in the immediately previous error condition by zeroing its contents.

NO ignore this error condition.

NO lost+found DIRECTORY (CREATE)

There is no lost+found directory in the root directory of the file system; When preen'ing fsck tries to create a lost+found directory.

Possible responses to the CREATE prompt are:

YES create a lost+found directory in the root of the file system. This may raise the message:

NO SPACE LEFT IN / (EXPAND)

See below for the possible responses. Inability to create a lost+found directory generates the message:

SORRY. CANNOT CREATE lost+found DIRECTORY

and aborts the attempt to linkup the lost inode. This will always invoke the UNREF error condition in Phase 4.

NO abort the attempt to linkup the lost inode. This will always invoke the UNREF error condition in Phase 4.

lost+found IS NOT A DIRECTORY (REALLOCATE)

The entry for lost+found is not a directory.

Possible responses to the REALLOCATE prompt are:

YES allocate a directory inode, and change lost+found to reference it. The previous inode reference by the lost+found name is not cleared. Thus it will either be reclaimed as an UNREF'ed inode or have its link count ADJUST'ed later in this Phase. Inability to create a lost+found directory generates the message:

SORRY. CANNOT CREATE lost+found DIRECTORY

and aborts the attempt to linkup the lost inode. This will always invoke the UNREF error condition in Phase 4.

NO abort the attempt to linkup the lost inode. This will always invoke the UNREF error condition in Phase 4.

NO SPACE LEFT IN /lost+found (EXPAND)

There is no space to add another entry to the lost+found directory in the root directory of the file system. When preen'ing the lost+found directory is expanded.

Possible responses to the EXPAND prompt are:

YES the lost+found directory is expanded to make room for the new entry. If the attempted expansion fails fsck prints the message:

SORRY. NO SPACE IN lost+found DIRECTORY

and aborts the attempt to linkup the lost inode. This will always invoke the UNREF error condition in Phase

4. Clean out unnecessary entries in lost+found. This error is fatal if the file system is being preen'ed.

NO abort the attempt to linkup the lost inode. This will always invoke the UNREF error condition in Phase 4.

LINK COUNT type I=I OWNER=O MODE=M SIZE=S MTIME=T COUNT=X
SHOULD BE Y (ADJUST)

The link count for inode I, is X but should be Y. The owner O, mode M, size S, and modify time T are printed. When preen'ing the link count is adjusted unless the number of references is increasing, a condition that should never occur unless precipitated by a hardware failure. When the number of references is increasing under preen mode, fsck exits with the message:

LINK COUNT INCREASING

Possible responses to the ADJUST prompt are:

YES replace the link count of file inode I with Y.

NO ignore this error condition.

UNREF type I=I OWNER=O MODE=M SIZE=S MTIME=T (CLEAR)

Inode I, was not connected to a directory entry when the file system was traversed. The owner O, mode M, size S, and modify time T of inode I are printed. When preen'ing, this is a file that was not connected because its size or link count was zero, hence it is cleared.

Possible responses to the CLEAR prompt are:

YES de-allocate inode I by zeroing its contents.

NO ignore this error condition.

BAD/DUP type I=I OWNER=O MODE=M SIZE=S MTIME=T (CLEAR)

Phase 1 or Phase 1b have found duplicate blocks or bad blocks associated with inode I. The owner O, mode M, size S, and modify time T of inode I are printed. This error cannot arise when the file system is being preen'ed, as it would have caused a fatal error earlier.

Possible responses to the CLEAR prompt are:

YES de-allocate inode I by zeroing its contents.

NO ignore this error condition.

4.8. Phase 5 - Check Cyl groups

This phase concerns itself with the free-block and used-inode maps. This section lists error conditions resulting from allocated blocks in the free-block maps, free blocks missing from free-block maps, and the total free-block count incorrect. It also lists error conditions resulting from free inodes in the used-inode maps, allocated inodes missing from used-inode maps, and the total used-inode count incorrect.

CG C: BAD MAGIC NUMBER

The magic number of cylinder group C is wrong. This usually indicates that the cylinder group maps have been destroyed. When running manually the cylinder group is marked as needing to be reconstructed. This error is fatal if the file system is being preen'ed.

BLK(S) MISSING IN BIT MAPS (SALVAGE)

A cylinder group block map is missing some free blocks. During preen'ing the maps are reconstructed.

Possible responses to the SALVAGE prompt are:

YES reconstruct the free block map.

NO ignore this error condition.

SUMMARY INFORMATION BAD (SALVAGE)

The summary information was found to be incorrect. When preen'ing, the summary information is recomputed.

Possible responses to the SALVAGE prompt are:

YES reconstruct the summary information.

NO ignore this error condition.

FREE BLK COUNT(S) WRONG IN SUPERBLOCK (SALVAGE)

The superblock free block information was found to be incorrect. When preen'ing, the superblock free block information is recomputed.

Possible responses to the SALVAGE prompt are:

YES reconstruct the superblock free block information.

NO ignore this error condition.

4.9. Cleanup

Once a file system has been checked, a few cleanup functions are performed. This section lists advisory messages about the file system and modify status of the file system.

V files, W used, X free (Y frags, Z blocks)

This is an advisory message indicating that the file system checked contained V files using W fragment sized blocks leaving X fragment sized blocks free in the file system. The numbers in parenthesis breaks the free count down into Y free fragments and Z free full sized blocks.

***** REBOOT UNIX *****

This is an advisory message indicating that the root file system has been modified by fsck. If UNIX is not rebooted immediately, the work done by fsck may be undone by the in-core copies of tables UNIX keeps. When preen'ing, fsck will exit with a code of 4. The standard auto-reboot script distributed with 4.3BSD interprets an exit code of 4 by issuing a reboot system call.

***** FILE SYSTEM WAS MODIFIED *****

This is an advisory message indicating that the current file system was modified by fsck. If this file system is mounted or is the current root file system, fsck should be halted and UNIX rebooted. If UNIX is not rebooted immediately, the work done by fsck may be undone by the in-core copies of tables UNIX keeps.