

Zero-blocks optimisation in MFS

Uladzislau Sobal
the University of Warsaw
ws374078@students.mimuw.edu.pl

ABSTRACT

This paper introduces and explains a simple optimisation for the Minix File System (MFS) that improves efficiency of storage for blocks filled with zeroes.

1. INTRODUCTION

Some files contain a lot of zeroes, so a natural question to ask would be if we can be more efficient around these zeroes. To answer this question, let's get our hands dirty with the workings of MFS (and ext, since they are similar).

2. MINIX FILE SYSTEM

Both MFS and ext use i-nodes. I-node is a structure used to store data. It's basically an array of pointers. The first 12 pointers are direct pointers to data. The next 256 pointers are pointers to pointers to data (indirect blocks). After that there are doubly indirect blocks (pointers to pointers to pointers) and trebly indirect blocks (pointers to pointers to pointers to pointers). What's interesting is if we are trying to read one of these blocks and it's not initialised (the pointer is equal to 0), the system just acts as if it was filled with zeroes. You can test this behaviour by creating a file, and dd'ing into it some random numbers with huge seek. Something like:

```
dd if=/dev/random bs=512 of=out seek=100000000
```

You will get a huge file that isn't even supposed to fit on the hard drive, and first 10^8 blocks will be filled with zeroes.

However, the system doesn't do anything of the kind when writing zeroes, it simply initialises a block and writes zeros. This is what I changed.

3. CHANGING MFS SOURCE CODE

To test the optimisation I altered standard MFS implementation. The file that I changed was `/fs/mfs/read.c`. There is a huge multi-purpose function called `rw_chunk`. I altered so that it first checks the data it is about to write, and, if the data is all zeroes, doesn't create the block. Implementation

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SOMIMUW '17 Warsaw, Poland

© 2017 ACM. ISBN 123-4567-24-567/08/06.

DOI: 10.475/123_4

is very simple, and the diff is only a couple of dozens of lines long.

4. OPTIMISATION RESULTS

To check how much we gain from this optimisation, I did 4 tests on different types of data - code, books, films and minix image. I copied to a clean mfs and the patched one a folder containing code in different languages, a folder with books and documents, a folder containing several episodes of a TV-series and .img file of a cloned system. Here are the results:

Table 1: Time

Data type	Time without patch	Time with patch
Code	16.226s	16.142s
TV series	2m12.663s	2m32.435s
Documents	10.542s	11.973s
minix.img	9.593s	13.733s

Table 2: Size

Data type	Size without patch	Size with patch
Code	219792	219328
TV series	4029608	4029608
Documents	249072	248920
minix.img	373344	371408

All tests were run on a machine with CPU Intel i7-4510U, SSD. Virtual machines were run with kvm enabled.

5. CONCLUSION

As you can see from the results, the optimisation is not always useful. It does improve space usage a little bit, but it also takes more time, and the tradeoff is not good enough to include the patch into the standard implementation. However, if you have a special case where data has a lot of zeroes, and time is not critical for you, you could use a patch that writes zeroes more effectively. Since the data written in this way can be read by the unpatched version of the file system, the most sensible way to use the optimisation is to estimate how much you gain and run the patched write only when needed.