

On Exploiting gzip's Content-Dependent Compression

Matteo Franzil, Luis Augusto Dias Knob

1. Introduction

Despite the development of more recent compression techniques like **bzip2** and **xz**, **gzip** is still a well-liked UNIX compression tool. Gzip is employed as the default compression technique in the majority of Linux and UNIX distributions. Although it has a somewhat poor compression ratio when compared to other utilities, this can be justified by its simplicity and quick compression and decompression times [1, 4]. This is also true for some tools, like **containerd**, which employs gzip as its standard compression technique for its image layers [2, 3].

In this report, we explore the possibility of exploiting gzip's algorithm for artificially increasing the decompression time of a file. By creating files filled with semi-random data generated with various methods, we show that compression and decompression times may vary significantly depending on the content of the file and the compression level used. Indeed, we show that the decompression time of a file can be increased by a factor of 3 in the worst case, when compared to the decompression time of a file containing English text.

2. Results

2.1. System setup

We run our experiments on a machine with a 4-core Intel Xeon Silver 4112 CPU @ 2.60GHz, 64GB of RAM, running Ubuntu Server 20.04.2 LTS. We used **gzip** version 1.12 on both machines. Tests were run in isolation and with CPU pinning, in order to reduce the impact of other processes on the results, and were run 5 times to reduce the impact of noise.

Our tests comprised the following steps:

- 1) Generate a file of 100MB in size, with a specific content
- 2) Compress the file with **gzip**
- 3) Decompress the file with **gzip**

We measured the time taken by each step, the size of the compressed and uncompressed files, the CPU usage, and the compression ratios. We repeated this test for each of **gzip**'s compression levels (1-9).

2.2. Popular tools

We first started by analyzing the compression and decompression times of files generated with some popular random data generators. We generated files with the following tools [5]:

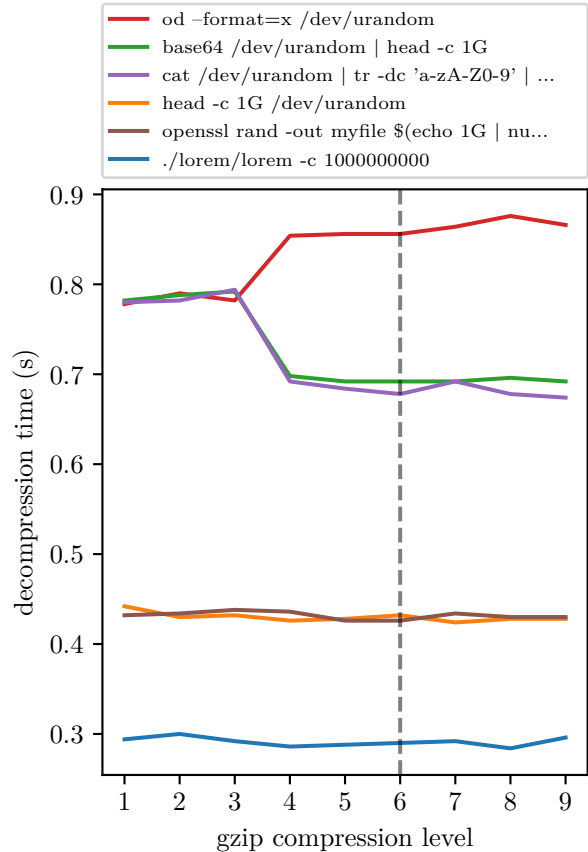


Figure 1: Decompression times for files generated with popular tools.

- `od --format=x /dev/urandom | head -c 1G`
- `base64 /dev/urandom | head -c 1G`
- `cat /dev/urandom | tr -dc 'a-zA-Z0-9' | head -c 1G`
- `openssl rand -out myfile "$(echo 1G | numfmt --from=iec)"`
- `lorem -c 1000000000`

We then compressed and decompressed these files with **gzip**, using the methods described above. The decompression times are shown in Figure 1.

We can see that the decompression times vary significantly depending on the tool used to generate the file. For example, the file generated by **lorem**, containing English text, is decompressed in 3 seconds. **openssl** and direct **/dev/urandom** output are decompressed in 4.5 seconds. Finally, the files generated by **base64** and **tr** require between 7 and 8 seconds

to be decompressed, depending on the compression level used. Finally, `od`'s output is decompressed in 9 seconds.

2.3. `od`'s output

Wishing to understand why `od`'s output was decompressed slower than the others, we decided to fully leverage `od`'s various output formats. We thus generated files with the following output formats:

- `x` (hexadecimal)
- `a`, `c` (ASCII both named and unnamed)
- `d1`, `d2`, `d4`, `d8` (decimal)
- `f` (floating point)
- `o` (octal)
- `u1`, `u2`, `u4`, `u8` (unsigned decimal)

We decided to use the decimal and unsigned decimal format variants with 1, 2, 4, and 8 bytes in order to see if the size of the numbers (and the minus sign in decimal formats) had any impact on the decompression time. It must be remembered we voluntarily un-padded the results of `od`, and thus, the files are comprised of a single line of content with no spaces or newlines.

We repeated the same tests as before, and the results are shown in Figure 2, although with 100MB files instead of 1GB files, due to the long time required to generate the files.

We can see that the decompression times vary significantly depending on the output format chosen. For our nefarious purposes, the best output formats (at level 6, the default one) are `a` and `c` (ASCII named and unnamed), although `x` almost matches them.

2.4. Finding the optimal file

Finally, the graph in Figure 3 shows the decompression times for `od`'s best output formats along with the other tools.

2.5. Other results

To verify that our results were not specific to our machine, we also ran the same tests on a MacBook Pro (late 2020) with an M1 CPU, 16GB of RAM, and running macOS Ventura 13.3.1 (a). We used the same version of `gzip` as on Ubuntu (1.12). Results were similar, although the MacBook Pro was on average faster than the Ubuntu machine. We thus do not include the results in this report.

Furthermore, we tried various file sizes, between 100MB and 1GB, in order to verify there was also no correlation between the file size and the results. Again, we found that the results were comparable, and thus we only reported the results once for the 100MB files.

3. Conclusions

We have shown that the decompression time of a file can vary significantly depending on the content

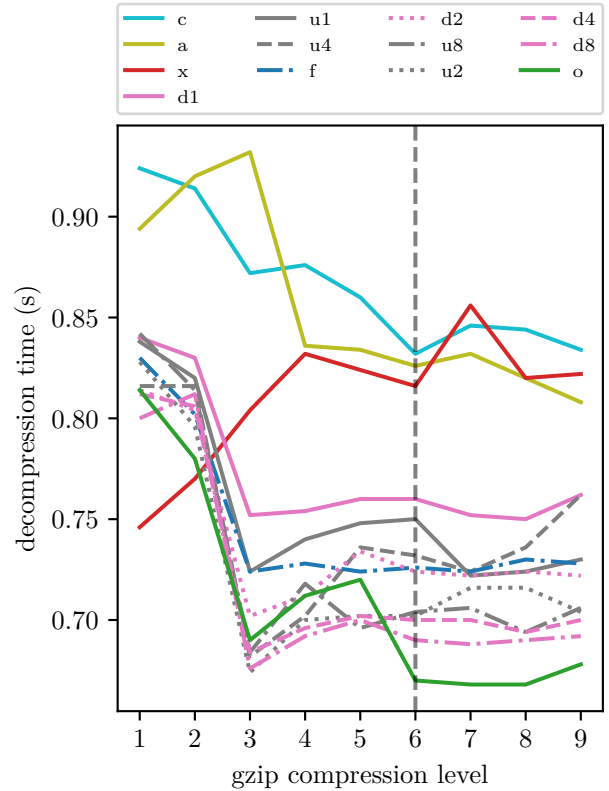


Figure 2: Decompression times for files generated with `od`'s various output formats.

of the file and the compression level used. Indeed, we have shown that the decompression time of a file can be increased by a factor of 3 when files are filled with data from `od`'s `a` and `c` output formats, when compared to the decompression time of a file containing English text.

We believe that this is a serious issue, as it can be used to artificially increase the decompression time of a file, which can be used to slow down systems that rely on `gzip` for decompression. For example, this could be used to slow down the unpacking of Docker images, which use `gzip` for compression of the various layers.

4. Bibliography

References

- [1] GNU Gzip.
- [2] Make image (layer) downloads faster by using pigz by sargun · Pull Request #35697 · moby/moby.
- [3] Support parallel decompression (pigz) by mxpv · Pull Request #2640 · containerd/containerd.
- [4] L. Peter Deutsch. DEFLATE Compressed Data Format Specification version 1.3. Request for Comments RFC 1951, Internet Engineering Task Force, May 1996.
- [5] Per Erik Strandberg. Lorem, December 2022.

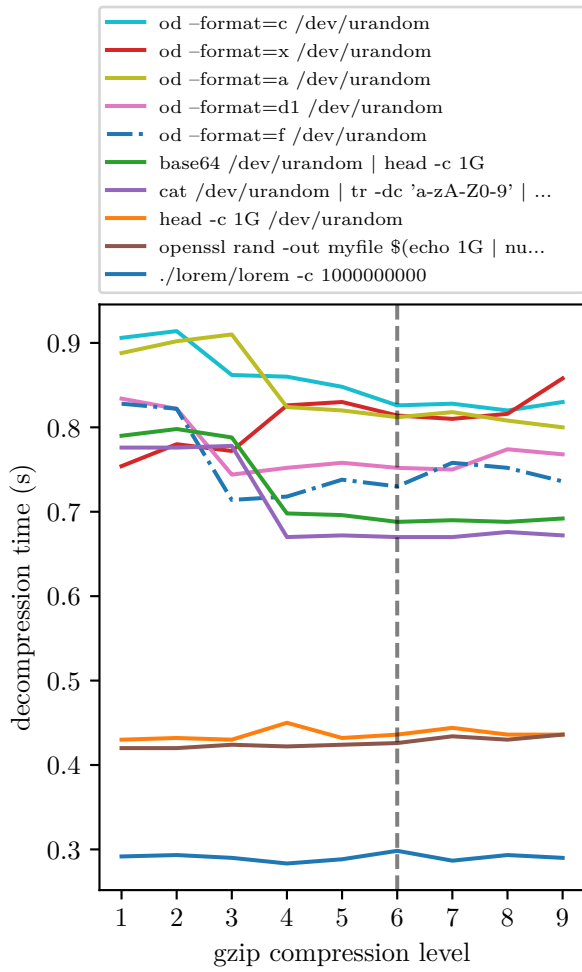


Figure 3: Decompression times for files generated with `od`'s best output formats and the other tools.

5. Appendix

The following are the full results of our tests for the last graph in the previous section.