

Home
Notes

Diffing very large files

Sometimes I want to find the difference between two files, line-by-line, using a program. There are several ways to accomplish this, and I prefer `vim`, using the `-d` command:

```
vim -d file1.dat file1_mod.dat
```

The advantage of this method is that `vim` will show the files side-by-side, using color highlights to indicate the differences between the two files.

A problem I have encountered is when the files are large. E.g. one of my files is 265MB in size, with 26 M lines. Running the plain `diff` command on these files takes <more than 20 mintes, after which I gave up>. Running `vimdiff` would take even longer!

Usually, I just want a quick check if there is a serious difference between two files. These files usually contain numerical data and will invariably contain small numerical differences, so I just want to glance over a broad enough spectrum of the files to get a feel for the differences.

To accomplish this, I can run `diff` on a subset of all lines in the two files. I can do that by combining the `head` command with `vim -d`:

```
vim -d <(head -n2000 file1.dat) <(head -n2000 file1_mod.dat)
```

This will show me the differences between the two files on only the first 2000 lines of head file. If I wanted some subset in the middle of the file, I could further pipe the data through the `tail` command:

```
vim -d <(head -n2000 file1.dat | tail -n1000) <(head -n2000 file1_mod.dat | tail -n1000)
```

This would show the difference in lines 1000-2000.