Subject: Re: the last line of a large file Posted by Conor on Fri, 10 Aug 2007 15:51:17 GMT

View Forum Message <> Reply to Message

On Aug 10, 11:00 am, Carsten Lechte <c...@toppoint.de> wrote:

- > Conor wrote:
- >> lol! Really! What in the world is the point of putting the number of
- >> lines at the end of the file?

- > One legitimate reason would be that sometimes you only know how much
- > data you have until after you have processed it all, especially if the
- > data sets are so large that you only ever have a small subset in RAM.

>

- A legitimate example are zip archives, where the table of contents is
- > written to the end of the file, because the the compressed sizes of
- > the archive members cannot be known in advance, and it would double
- > the running time to determine the compressed size beforehand, it would
- > furthermore use twice the disk space to re-write the file with the
- > contents in front, it would be impossible to keep the whole archive
- > in RAM before writing it, and finally, one may not be able leave space
- > for the contents table at the beginning of the file, to be filled in
- > later, because one would have to know how long the table will be
- > beforehand...

>

- Of course, this does not mean that the original poster's data has a
- > legitimate reason for being organised like this.

>

- For the original poster's problem, one idea is to get the file size
- > in bytes, skip to position file size-1000, read that small chunk and
- > parse it for the desired metadata. This might even be faster than
- > actually counting the lines with FILE LINES, but it is probably only
- > worth it if the metadata contains more useful information that just
- > the number of lines in the file.

> chl

As an actual suggestion, if the file\_lines dosesn't take too long you can always just count the number of lines and break down the file into manageable chunks. Imagine for a moment that the following file has 1,000,000 lines and your computer can only make arrays with 10,000 rows at a time (which you would know in advanced). You might do something like this:

```
max_size = 10000
num_rows = file_lines(file) ; 1,000,000
num_parts = num_rows/max_size; 10 parts
num cols = 10
```

```
data = fltarr(max_size,numcols)
```

for i=0,num\_parts-1 do begin readf,lun,data ; do something with data, then read the next chunk endfor

There's a couple other things you can do. For starters, if you don't already know it you can calucluate the number of columns in the file by reading in the first line, using strsplit, and then rewinding the file to the beginning. Also, I haven't included it in the above code, but you'll have to keep track of the last line still. In this case what you would probably do is calculate how many lines you want to read in the last chunk of data, and worry about it then. For instance, imagine the same example but now the line has 75,000 lines and you don't want to read the last one:

```
max_size = 10000
num_rows = file_lines(file) ; 75,000
num_parts = ceil(num_rows/max_size) ; 8 parts
num_cols = 10
last_read = max_size - (num_parts*max_size - num_rows) - 1 ; 4999
data = fltarr(max_size,numcols)
for i=0,num_parts-1 do begin
if i eq num_parts-1 then begin
    readf,lun,data
endif else begin
    data = fltarr(last_read,numcols)
    readf,lun,data
endelse
; do something with data, then read the next chunk
endfor
```

Not exactly elegant, but it should work for your problem.