

## CS146 - Assignment #5: more shell scripting

In this assignment, you are going to learn how to write “shell” scripts that have complex embedded awk(1) scripts. That is, the file will technically be a Bourne (or Bash) shell script, but one (or more) of the command-lines inside the shell script will be using awk; in turn, the awk program can be arbitrarily long and complicated. In the end, half or more of the characters in the file may actually be awk code. For example, say we want to evaluate mathematical expressions on the command-line. While Bash can handle the 4 basic operations (add, subtract, multiply, divide):

```
$ echo $(((15+5)/5)
4
```

it knows nothing about basic mathematical functions:

```
$ echo $((cos(0))) # hoping to get cosine(0), which is 1
bash: cos(0): syntax error in expression
```

Luckily, awk knows about such simple functions:

```
$ (echo 0; echo 1; echo 3.1415) | awk '{print cos($1)}'
1
0.540302
-1
$
```

So if we want to evaluate mathematical functions on the Unix command line, we can write a little shell script around awk, as follows, which shows a script called “fexpr” (meaning “floating point” version of the existing expr(1) program):

```
$ cat fexpr
#!/bin/sh
exec awk "BEGIN{print $*; exit}"
$ fexpr 'cos(0)'
1
$ fexpr 'exp(1)'
2.71828
$ fexpr 'cos(1)^exp(1)' # raise cos(1) to the power of 2.71828
0.187599
$
```

A few things to note: while normally awk executes the program you give it on every input line, in this case we want it to evaluate one expression, and exit. The BEGIN tag is always executed by awk before any lines are read. We put “\$@” in double quotes, so the Bourne Shell will expand that to everything given on the command line of fexpr; and then we just exit, so awk doesn’t sit around waiting for any input lines. To save on the expense of one fork() system call, we tell the Bourne shell to **exec** awk rather than forking it as a child process.

That one-line Bourne shell script demonstrates the power of combining Bourne shell with awk; you might think that it would be easy to have the interpreter file use *only* awk, but the problem is that awk has no easy way to grab all its command line arguments at once, interpret them all together as one expression, and print the resulting value.<sup>1</sup> So instead, we let Bourne do the heavy lifting here and just use \$\* inside double-quotes, and BAM! we’re done. Note also that this is one of the very few cases where we want to use “\$\*” rather than “\$@”.

---

<sup>1</sup>**Bonus:** figure out how to do it using awk as the interpreter file rather than Bourne. That means the first line of the file must start with “#!/bin/awk ...”—there could be more on that first line, but I’ll let you figure it out. Note *this is hard*, which is why I’m assigning it as a bonus. I’m not ever sure it’s *possible*, so don’t spend too much time on it.

# 1 rename: programmatically rename files using regular expressions

Let's say we're in a directory and we want to rename every file ending in `.c`, to end in `.C` (i.e., capitalize the 'c'). Here's one way to do it:

```
$ for i in *.c; do b='basename "$i" .c'; mv "$b.c" "$b.C"; done
```

Things get more complicated if you want more complex renaming. Say for example you just created exactly one thousand simulation output files named `output.XXX`, where `XXX` runs from 000 through 999 inclusive. These files used a simulation parameter equal to 10, but you want to re-do the simulation changing that parameter to 20. So it would be useful now to rename all the files to look like `output.p10.XXX`, and then the new ones will be called `output.p20.XXX`. Then it gets harder. One way would be to do this:

```
for i in output.*; do #..... hmmmmm, now what??? Maybe not a loop....
ls output.* | sed 's/output\.//' | awk '{printf "mv output.%d output.p10.%d\n",$1,$1}' | sh
# Yes, the above works. But also note you don't have to use '/' to separate the source
# from the target in sed. Whatever follows the initial 's' is defined to be your
# separator, so the above could just as easily have been:
..... sed 's;output\.;;' ..... # use a semi-colon separator.
```

That last one provides a general hint: when algorithmically renaming files, regular expressions are far more powerful than shell globbing. So let's formalize that.

Create a program called `rename`. Its first argument is a `sed` substitution command, and all remaining arguments are existing filenames that will be renamed using the substitution command. For simplicity the user will leave out the initial 's', and the first character in the first argument shall be used as the separator. Thus, both examples above can be done in one line each:

```
rename ';\.c$;.C;' *.c
rename ';\.;.p10.;;' output.*
```

Note in the second case, only the first period will be matched. If the user wants to match *every* occurrence of the source regular expression, the user can append a 'g' to the end of the substitution string, so for example "`rename ';\.c$;.Cg;' *`" will change *every* lower-case 'a' to an upper case one in all filenames in the current directory.

To avoid catastrophe, `rename` should refuse to over-write an existing file, but it will only print a warning when that happens, and continue to perform the rename operation on all files for which the destination doesn't over-write an existing file. Finally, if given the "-f" ("force") option before the substitution argument, then it will even over-write existing files with abandon. Thus, the man-page style SYNOPSIS is:

```
rename [-f] sed-substitution-command file1 file2 ...
```

Finally, to be maximally flexible, the files can actually be path names including slashes—which explains why we probably don't want to use '/' as the sed separator. So if the above examples were in some other directory, say "src" in the first case and "sim-output" in the second, this should work fine:

```
rename ';\.c$;.C;' src/*.c
rename ';\.;.p10.;;' sim-output/output.*
```

It is OK to assume that only the last element of the path (ie, none of the intervening directories) should match the regular expression. I think this would be hard to verify on-the-fly, so you're allowed to just assume it.

## 2 awkcel: use awk a bit like a spreadsheet

Write a Bourne shell script named `awkcel`. Yes, the name is a blatant rip-off of Excel(tm). It works on tab-separated files, which I recommend have the extension `.tsv`. (This need not be enforced; it will work on any file regardless of name, but the file *does*, in fact, need to be *tab* separated, and every tab is assumed to be a field separator.) Here is the first few lines of my solution:

```
$ cat awkcel
#!/bin/sh
# This is a front-end to awk(1), which allows NAMED COLUMNS to be accessed as variables.
# USAGE: awkcel {any standard awk program} FILENAME
# This script expects exactly two arguments.
# The first argument is an awk program that allows use of variables as described below.
# The second argument is an input file with a constant number of tab-separated columns.
# Any line that starts with a hash ("##") is assumed a comment, and discarded.
# The first (non-comment) line is the HEADER.
# The HEADER line contains tab-separated names of the columns.
# These column names must be valid awk(1) variable names.
# Each subsequent line must have the same number of tab-separated columns as the HEADER.
# COLUMNS SHOULD NOT BE EMPTY: behavior is undefined if any line has two adjacent tabs.
# Upon each line of input, values can be accessed using the name defined in the HEADER.
```

The above should be considered a *specification* of `awkcel` and must be abided by to get full marks. To enforce `awk(1)` using tabs to separate columns, use `awk's -F` option.

So for example, say the file `name-studnum.tsv` has the following:

```
$ cat name-studnum.tsv
name      studentNum  grade
hayes     00000      80.6
shawna    11111      85.4
dan       22222      83.7
$ awkcel 'studentNum=="00000"{print name, "is the prof"} \
        studentNum!="00000"{ \
            print name, "must be a TA or Reader or Student and got grade", \
            int(grade+0.5)}' name-studnum.tsv
hayes is the prof
shawna must be a TA or Reader or Student and got grade 85
dan must be a TA or Reader or Student 84
$
```

Some example tab-separated files are on openlab in `~wayne/pub/cs146/awkcel`. Based on those files, and your implementation of `awkcel`, answer the following questions. Your openlab submission should include all code for `awkcel` itself, and your write-up (in PDF) should show the *user-level* commands you used to answer the below questions.

1. The file `orthologs.tsv` contains lists of proteins (encoded by genes) that are identical across various species; the technical term is that two proteins that perform the same function, regardless of species, are called “orthologs”. On any given line, all proteins listed are orthologs of each other, and any species with an underscore “\_” as an entry has no such ortholog.
  - (a) Which pair of species has the greatest number of orthologs? Which pair has the least?
  - (b) Create a table that lists, for every pair of species in that file, the number of orthologs they have. (You can use a shell script to iterate over the pairs, but for any given pair you should use `awkcel` to compute the number of orthologs for that pair.)

2. The file `historical.2011.txt` contains the closing prices for a few hundred medium-to-large cap stock market symbols for most trading days of 2011, in date order. (Empty entries contain a single underscore.) Once you think your `awkcel` is working, use it to answer the following questions.
- (a) What were the closing prices of Amazon (AMZN) and Microsoft (MSFT) on July 1st, 2011? How about the 5th? Did they go up or down over the July 4th holiday?
  - (b) What was Amazon's worst day of 2011? (ie., worst 24-hour, close-to-close loss, in raw dollars, not percent)?
  - (c) Among all the companies in that file, which one suffered the worst single-day loss, as a percent of its price? Name the company, and the date. (Don't include any days surrounding days where the price is not listed, ie., is an underscore.)

**BONUS:** Try to write a program like `awkcel` in Python. How does its speed compare to `awkcel`? If you can make your Python version faster than *my* version of `awkcel` (which is just Bourne shell wrapped around `awk`), you automatically pass this course and get a 10% bonus on this assignment.

## More questions for `awkcel`, this time about Spiral Galaxies

I have added more files to the directory `/home/wayne/pub/cs146/awkcell`. Note I've also added a shell script called `hawk`, which stands for "Hayes awk". It is simply a one-line shell script that runs normal `awk` but with a large number of new functions defined. Those functions are in the file `misc.awk` in the same directory. Feel free to peruse `misc.awk` for many examples of `awk` syntax. You are welcome to have your `awkcel` also read `misc.awk` (or better, just execute `hawk` rather than normal `awk`).

I have also added one very large (about 2.1GB) file called `SDSS+GZ1+SpArcFiRe+SFR.tsv`. It is an actual data file containing 341 columns and about 750,000 lines. Each line contains data for one galaxy. It is actual, real-life data on exactly 765,367 galaxies. The data are partly from the *Sloan Digital Sky Survey* (<http://sdss.org>), partly from the *Galaxy Zoo Project* (<http://zoo1.galaxyzoo.org>) and partly from *SpArcFiRe* (<http://sparcfire.ics.uci.edu>), a spiral galaxy analysis program written by one of my Ph.D. students a few years ago.

You don't need to know much about galaxies to do this assignment, but your `awkcel` needs to be pretty efficient to answer the following questions, since it's going to take quite awhile for `awk` to read through a 2.1GB file containing over 750,000 lines of 341 columns each. The only thing you need to know about galaxies is that some galaxies have spiral arms (for example the one that's on the background at Galaxy Zoo, or the one on our *SpArcFiRe* home page), and some are elliptical (ie., they have no spiral arm structure). Among spiral galaxies, some wind "clockwise", and others "anti-clockwise".

Among the columns in `SDSS+GZ1+SpArcFiRe+SFR.tsv`, two are named `P_CW` and `P_ACW`. They record the probability that the arms wind clockwise, and anti-clockwise, respectively. Using `awkcell`, define a variable `P_SP=P_CW+P_ACW`; this value tells us the probability that the galaxy shows observable spiral arm structure, regardless of the winding direction. Answer the following questions:

How many galaxies have `P_SP` exactly equal to 1? How about exactly equal to 0?

How many galaxies have `P_SP > 0.5`?

Print out a histogram of mean `totalArcLength` as a function of `P_SP` in increments of 0.1, giving 11 bins 0 through 10. Below is the exact `awkcel` command line to do this. (I've paragraphed it nicely but it could all go on one command line if you want.)

```
awkcel '{
    P_SP=P_CW+P_ACW;
    bin=int(10*P_SP);
    n[bin]++;
    numArcs[bin]+=totalArcLength;
}
END{
for(i=0;i<=10;i++)
    print i, n[i], numArcs[i]/n[i];
}' SDSS+GZ1+SpArcFiRe+SFR.tsv
```

### 3 Parallel in C

Write a C program named `parallel(1)`. Below is a short-form manual entry.

#### NAME

`parallel` - execute many independent Shell command lines in parallel, N at a time

#### SYNOPSIS

<sequence of independent Shell command lines> | `parallel [-s SHELL] [N]`

#### DESCRIPTION

Given an integer N and an arbitrarily large set of independent Shell command lines (here “independent” means that no line depends upon any other, so they can be executed in any order), maintain the property that exactly N of them are running simultaneously (except, of course, as the last N-1 command lines are executed, fewer than N will be running, until the last child exits). Thus, for example, if given 7 input lines and N=4, then start the first 4; each time any child exits (order does not matter), start the next in the sequence. Of course, once any 4 of these 7 have finished, only 3 will be left; then 2; then 1. Only once the last child exits, does `parallel` return. In other words, though `parallel` itself is starting processes asynchronously, `parallel` itself halts synchronously, so that it only exits once all of its children have exited. All processes run on the same host as `parallel` itself; that is, they are simply started under local shell processes.

If N is not given, then `parallel` determines the number of cores available; assuming that number is C, `parallel` acts as if N=C.

`parallel` runs each command separately, starting a new shell process for each. By default, `parallel` reads the SHELL environment variable to determine which shell to execute the commands under. The -s option allows the specification of a different shell. So for example if the SHELL environment variable is `/bin/csh`, then each command will run under a new `csh` process; if given “-s `/bin/bash`”, then `/bin/bash` will be used instead.

#### RETURN VALUE

`parallel`’s exit code is equal to the number of children that exited with a status != 0. That is, if all children exit successfully, `parallel` also exits with status 0. If one child failed for any reason, `parallel` exits with status 1, etc. Note that Unix only allows an 8 bit unsigned integer for the return status, so 255 is the maximum exit code allowed; if more than 255 children exited with errors, `parallel` returns 255.

#### SEE ALSO

`fork(2)`, `exec(2)`, `wait(2)`