

**awk** is a standard Unix utility which has its own script language to provide data extraction, transformations, and reports. Its name comes from the last names of its creators (Aho, Weinberger, Kernighan). The GNU version of awk is named **gawk**.

**Example 1:** We want the product ID (first field) and unit price (fourth field) from the inventory.txt file. Sample data line:

## Using a Program File

As shown above, the **-f** switch is used to specify a program file which allows for more complex capabilities. `awk` provides C-like **if**, counting **for**, **while**, and **printf** action commands. It also supports a **for in** to iterate over the contents of an array.

Some special conditions:

- BEGIN** executes the action before the records are read. In the actions, we typically initialize variables and print column headings.
- END** executes the action after the records are read. It is very common to print totals in the action corresponding to an **END** condition.

**Example 4:** Show the product ID, inventory quantity, and unit price for any products having an inventory quantity (second field) greater than 200. Also show a column heading.

```
$ cat >example4
```

```
BEGIN {printf("%-6s %4s %-10s\n", "ID", "QTY", "UNIT PRICE");}
$2 > 200 {printf("%6s %4d %8.2f\n", $1, $2, $4);}
```

CTRL-D

```
$ gawk -f example4 inventory.txt
```

```
ID      QTY  UNIT PRICE
SBB001  300    14.95
SBG002  400    14.95
NHC001  300     9.95
```

**Example 5:** Show the product ID and the product description for every product. Since a product description can be many words, we will use a counting for loop beginning at \$6 and ending at NF. Note: you can reference \$i.

```
$ cat >example5
```

```
{
    # output the product ID
    ??
    # output each word in the product description
    ??
    printf("\n");
}
```

CTRL-D

```
$ gawk -f example5 inventory.txt
```

```
PPF001 Popeil Pocket Fisherman
SBB001 Snuggie Brown
SBG002 Snuggie Green
BOM001 Bass-O-Matic
MCW001 Miracle Car Wax
TTP001 Topsy Turvy Planter
NHC001 Electric Nose Hair Clipper
SSX001 Secret Seal
```

The `awk` arithmetic operators are from the C programming language.

The type of comparison (numeric or string) is based on the operands. If both are numeric, a numeric comparison is done. Otherwise, a string comparison is used.

**Example 6:** For products having a unit price greater than \$10 and more than 100 items in inventory, print the ID, inventory quantity, unit price, and gross value (product of inventory quantity and unit price). Also print the total gross value for all products meeting the criteria.

```
$ cat >example6
```

```
BEGIN {
    printf("%-6s %4s %-10s %-12s\n", "ID", "QTY", "UNIT PRICE",
        "GROSS PRICE");
    total = 0;
}
{
    if ( $2 > 100 && $4 > 10.0 )
    {
        gross = $2 * $4;
        printf("%6s %4d %8.2f %10.2f\n", $1, $2, $4, gross);
        total = total + gross;
    }
}
END { printf("%-6s %4s %-10s %10.2f\n", " ", " ", " ", total);}
```

CTRL-D

	<pre>\$ gawk -f example6 inventory.txt ID      QTY  UNIT PRICE  GROSS PRICE SBB001  300    14.95      4485.00 SBG002  400    14.95      5980.00 SSX001  150    29.95      4492.50                                 14957.50</pre> <p>How can we rewrite this to not include the if-statement? ??</p>
<p><b>Exercise#1: List the login names for anyone in the faculty group (group 1000).</b> Use the following to get all users:</p> <pre>\$ getent passwd</pre> <p>It returns records that look like this:</p> <pre>krobbins:x:512:1000:Kay A. Robbins:/home/krobbins:/bin/csh maynard:x:511:1000:Hugh B. Maynard:/home/maynard:/bin/tcsh clark:x:1000:1000:Larry Clark:/home/clark:/bin/tcsh abc123:x:5035:1001:Bob:/home/abc123:/bin/tcsh ...</pre> <p>The faculty group is group 1000 which is in the 4th field.</p> <p>How will your awk program get its input? ??</p> <p>How do we specify a different field separator? ??</p>	<pre>\$ ?? ?? ??</pre>
<p><b>Range Patterns</b></p> <p>awk supports ranges of lines as a pattern just like sed:</p> <pre>/pat1/,/pat2/</pre>	<p><b>Example 7:</b> Analyze the code in cs1713p0.c, counting number of comment lines, number of blank lines, and number of code lines.</p> <pre>\$ cat &gt; example7 BEGIN { blankCount = 0; commentCount=0} /^\\\/\\\/,\\/\\\/ { commentCount++ } /^[ \t]*\\\/\\\/ { commentCount++} /^[ \t]*\$/ {blankCount++} END {print "Total Lines:", NR;       print "Comment lines:", commentCount;       print "Blank lines:", blankCount;       print "Code: ", NR - commentCount - blankCount; }</pre> <p>CTRL-D</p> <pre>\$ gawk -f example7 cs1713p0.c Total Lines: 77 Comment lines: 26 Blank lines: 9 Code: 42</pre>
<p><b>Associative Arrays</b></p> <p>Awk supports associative arrays (i.e., hash tables). The key for an associative array can be a character string. To assign a value:</p> <pre>array[key] = value;</pre> <p>To check whether an entry exists:</p>	<p><b>Example 8:</b> print the total of purchase items for each item requested. Examine invCommand.txt. We are only interested in the ORDER ITEM records.</p> <pre>\$ cat &gt;example8 \$1 == "ORDER" &amp;&amp; \$2 == "ITEM"{     if (\$3 in invM)         invM[\$3] += \$4;     else         invM[\$3] = \$4;</pre>

<pre>if (key in array)     doSomething;</pre> <p>To iterate over the keys of the array:</p> <pre>for (key in array)     doSomething;</pre>	<pre>} END {     for (key in invM)         print key, invM[key]; }</pre> <p>CTRL-D</p> <pre>\$ gawk -f example8 invCommand.txt XXX001 20 SBG002 410 SVC001 3 APC001 1 NHC001 260 BOM001 2 PPF001 11 MCW001 62 SBB001 38</pre>
<p><b>Example 9: produce a shell script</b></p> <p>In this example, the output from awk will be a shell script.</p> <p>It is very common to have core dump files named "core" taking up lots of space throughout your directories. Can we use awk to help remove the files?</p> <p>As a system administrator, you would use "locate core" to find the files which would give us a huge result containing many files (not just mine).</p> <p>To simplify, we will simulate the use of locate by using find.</p> <p>First lets find core files using find:</p> <pre>\$ find ~ -name '*core*' /home/ssilvestro/cs1713/Pgm3/core /home/ssilvestro/cs4713/core /home/ssilvestro/cs2123/core /home/ssilvestro/.cache/compizconfig/core.pb /home/ssilvestro/core /home/ssilvestro/cs3423/score /home/ssilvestro/cs3423/core</pre> <p>Notice that some of the files are core, but others just have "core" somewhere in the name. We don't want to remove the other files. We will also confirm that each file starts with /home/ssilvestro/ since locate would have returned others.</p>	<p><b>Example 9:</b> produce a shell script to remove files</p> <p>What should our awk program do?</p> <pre>\$ cat &gt; example9 BEGIN { count = 0;} /^\home\/ssilvestro\/ &amp;&amp; /\core\$/ {      ??  } END ??</pre> <p>CTRL-D</p> <pre>\$ find ~ -name '*core*'   awk -f example9 - &gt; coreRm</pre> <pre>\$ bash coreRm removed 5 files</pre>

<p>Suppose we have awk output a shell script file that contains</p> <pre>rm /home/ssilvestro/cs1713/Pgm3/core rm /home/ssilvestro/cs4713/core rm /home/ssilvestro/cs2123/core rm /home/ssilvestro/core rm /home/ssilvestro/cs3423/core echo "removed 5 files"</pre> <p>We can then run that script file after verifying it.</p>	
<p><b>Passing arguments into awk code</b></p> <p>We can pass variable values into awk by specifying:  <code>awk options 'program' -v 'var=value' file1</code></p> <p>We can change example9 to allow us to pass in the user.</p> <p>The awk function match returns true if the functions the first argument matches the pattern specified in the second argument.</p>	<p><b>Example 10:</b> passing in a variable for the /home/user/  \$ cat &gt; example10  BEGIN { count = 0;}  /<code>/core</code>\$/ {      if ( match( \$0, arg1 ) )      {          print "rm", \$0;          count++;      }  }  END {print "echo removed", count, "files"}  CTRL-D  \$ find ~ -name '*core*'   awk -f example10 -v 'arg1=/home/ssilvestro/' -</p> <p>Why use match()?  ??</p>
<p><b>Special Variables</b></p> <ul style="list-style-type: none"> <li>FS    input field separator (defaulted to white space)</li> <li>OFS   output field separator (defaulted to blank)</li> <li>RS    input record separator (defaulted to \n)</li> <li>ORS   output record separator (defaulted to \n)</li> <li>NF    number of fields for the current line</li> <li>NR    record number of the current line</li> </ul>	
<p><b>Some built-in functions</b></p> <ul style="list-style-type: none"> <li>int(<i>val</i>)                returns the truncated integer value</li> <li>length(<i>val</i>)            returns the length of the value</li> <li>index(<i>str</i>,<i>match</i>)      returns the index of <i>match</i> in <i>str</i> or 0 if it isn't found</li> <li>substr(<i>str</i>,<i>pos</i>,<i>length</i>) returns the substring of <i>str</i> beginning at <i>pos</i> for <i>length</i> characters</li> <li>split(<i>string</i>, <i>array</i> [, <i>fieldsep</i>]) splits string into array (indexed at 1) according to fieldsep</li> <li>sub, gsub, gensub        see gawk manual for full details</li> </ul>	