

A

AWK SUMMARY

This appendix contains a summary of the awk language. In syntactic rules, components enclosed in brackets [...] are optional.

Command-line

`awk [-Fs] 'program' optional list of filenames`

`awk [-Fs] -f progfile optional list of filenames`

The option `-Fs` sets the field separator variable `FS` to `s`. If there are no filenames, the standard input is read. A filename can be of the form `var=text`, in which case it is treated as an assignment of `text` to the variable `var`, performed at the time when that argument would be accessed as a file.

AWK programs

An awk program is a sequence of pattern-action statements and function definitions. A pattern-action statement has the form:

`pattern { action }`

An omitted pattern matches all input lines; an omitted action prints a matched line.

A function definition has the form:

`function name(parameter-list) { statement }`

Pattern-action statements and function definitions are separated by newlines or semicolons and can be intermixed.

Patterns

`BEGIN`

`END`

`expression`

`/regular expression/`

`pattern && pattern`

`pattern || pattern`

`!pattern`

`(pattern)`

`pattern, pattern`

The last pattern is a range pattern, which cannot be part of another pattern. Similarly, `BEGIN` and `END` do not combine with other patterns.

Actions

An action is a sequence of statements of the following kinds:

```
break
continue
delete array-element
do statement while (expression)
exit [expression]
expression
if (expression) statement [else statement]
input-output statement
for (expression; expression; expression) statement
for (variable in array) statement
next
return [expression]
while (expression) statement
{ statements }
```

A semicolon by itself denotes the empty statement. In an if-else statement, the first *statement* must be terminated by a semicolon or enclosed in braces if it appears on the same line as else. Similarly, in a do statement, *statement* must be terminated by a semicolon or enclosed in braces if it appears on the same line as while.

Program format

Statements are separated by newlines or semicolons or both. Blank lines may be inserted before or after any statement, pattern-action statement, or function definition. Spaces and tabs may be inserted around operators and operands. A long statement may be broken by a backslash. In addition, a statement may be broken without a backslash after a comma, left brace, &&, ||, do, else, and the right parenthesis in an if or for statement. A comment beginning with # can be put at the end of any line.

Input-output

close(<i>expr</i>)	close file or pipe denoted by <i>expr</i>
getline	set \$0 from next input record; set NF, NR, FNR
getline < <i>file</i>	set \$0 from next record of <i>file</i> ; set NF
getline <i>var</i>	set <i>var</i> from next input record; set NR, FNR
getline <i>var</i> < <i>file</i>	set <i>var</i> from next record of <i>file</i>
print	print current record
print <i>expr-list</i>	print expressions in <i>expr-list</i>
print <i>expr-list</i> > <i>file</i>	print expressions on <i>file</i>
printf <i>fmt</i> , <i>expr-list</i>	format and print
printf <i>fmt</i> , <i>expr-list</i> > <i>file</i>	format and print on <i>file</i>
system(<i>cmd-line</i>)	execute command <i>cmd-line</i> , return status

The *expr-list* following print and the *fmt*, *expr-list* following printf may be parenthesized. In print and printf, >>*file* appends to the *file*, and | *command* writes on a pipe. Similarly, *command* | getline pipes into getline. The function getline returns 0 on end of file, and -1 on error.

Printf format conversions

These conversions are recognized in `printf` and `sprintf` statements.

<code>%c</code>	ASCII character
<code>%d</code>	decimal number
<code>%e</code>	<code>[-]d.ddddddE[+-]dd</code>
<code>%f</code>	<code>[-]ddd.dddddd</code>
<code>%g</code>	e or f conversion, whichever is shorter, with nonsignificant zeros suppressed
<code>%o</code>	unsigned octal number
<code>%s</code>	string
<code>%x</code>	unsigned hexadecimal number
<code>%%</code>	print a %; no argument is converted

Additional parameters may lie between the % and the control letter:

<code>-</code>	left-justify expression in its field
<code>width</code>	pad field to this width as needed; leading 0 pads with zeros
<code>.prec</code>	maximum string width or digits to right of decimal point

Built-in variables

The following built-in variables can be used in any expression:

<code>ARGC</code>	number of command-line arguments
<code>ARGV</code>	array of command-line arguments (<code>ARGV[0..ARGC-1]</code>)
<code>FILENAME</code>	name of current input file
<code>FNR</code>	input record number in current file
<code>FS</code>	input field separator (default blank)
<code>NF</code>	number of fields in current input record
<code>NR</code>	input record number since beginning
<code>OFMT</code>	output format for numbers (default <code>"%.6g"</code>)
<code>OFS</code>	output field separator (default blank)
<code>ORS</code>	output record separator (default newline)
<code>RLENGTH</code>	length of string matched by regular expression in match
<code>RS</code>	input record separator (default newline)
<code>RSTART</code>	beginning position of string matched by match
<code>SUBSEP</code>	separator for array subscripts of form <code>[i,j,...]</code> (default <code>"\034"</code>)

`ARGC` and `ARGV` include the name of the invoking program (usually `awk`) but not the program arguments or options. `RSTART` is also the value returned by `match`.

The current input record is named `$0`. The fields in the current input record are named `$1`, `$2`, ..., `$NF`.

Built-in string functions

In the following string functions, *s* and *t* represent strings, *r* a regular expression, and *i* and *n* integers.

An `&` in the replacement string *s* in `sub` and `gsub` is replaced by the matched string; `\&` yields a literal ampersand.

<code>gsub(<i>r,s,t</i>)</code>	globally substitute <i>s</i> for each substring of <i>t</i> matched by <i>r</i> , return number of substitutions; if <i>t</i> is omitted, \$0 is used
<code>index(<i>s,t</i>)</code>	return the index of <i>t</i> in <i>s</i> , or 0 if <i>s</i> does not contain <i>t</i>
<code>length(<i>s</i>)</code>	return the length of <i>s</i>
<code>match(<i>s,r</i>)</code>	return index of where <i>s</i> matches <i>r</i> or 0 if there is no match; set RSTART and RLENGTH
<code>split(<i>s,a,fs</i>)</code>	split <i>s</i> into array <i>a</i> on <i>fs</i> , return number of fields; if <i>fs</i> is omitted, FS is used in its place
<code>sprintf(<i>fmt,expr-list</i>)</code>	return <i>expr-list</i> formatted according to <i>fmt</i>
<code>sub(<i>r,s,t</i>)</code>	like <code>gsub</code> except only the first matched substring is replaced
<code>substr(<i>s,i,n</i>)</code>	return the <i>n</i> -character substring of <i>s</i> starting at <i>i</i> ; if <i>n</i> is omitted, return the suffix of <i>s</i> starting at <i>i</i>

Built-in arithmetic functions

<code>atan2(<i>y,x</i>)</code>	arctangent of <i>y/x</i> in radians in the range $-\pi$ to π
<code>cos(<i>x</i>)</code>	cosine (angle in radians)
<code>exp(<i>x</i>)</code>	exponential e^x
<code>int(<i>x</i>)</code>	truncate to integer
<code>log(<i>x</i>)</code>	natural logarithm
<code>rand()</code>	pseudo-random number <i>r</i> , $0 \leq r < 1$
<code>sin(<i>x</i>)</code>	sine (angle in radians)
<code>sqrt(<i>x</i>)</code>	square root
<code>srand(<i>x</i>)</code>	set new seed for random number generator; uses time of day if no <i>x</i> given

Expression operators (increasing in precedence)

Expressions may be combined with the following operators:

<code>= += -= *= /= %= ^=</code>	assignment
<code>?:</code>	conditional expression
<code> </code>	logical OR
<code>&&</code>	logical AND
<code>in</code>	array membership
<code>~ !~</code>	regular expression match, negated match
<code>< <= > >= != ==</code>	relationals
	string concatenation (no explicit operator)
<code>+ -</code>	add, subtract
<code>* / %</code>	multiply, divide, mod
<code>+ - !</code>	unary plus, unary minus, logical NOT
<code>^</code>	exponentiation
<code>++ --</code>	increment, decrement (prefix and postfix)
<code>\$</code>	field

All operators are left associative, except assignment, `?:`, and `^`, which are right associative. Any expression may be parenthesized.

Regular expressions

The regular expression metacharacters are

`\ ^ $. [] ! () * + ?`

The following table summarizes regular expressions and the strings they match:

<code>c</code>	matches the nonmetacharacter <i>c</i>
<code>\c</code>	matches the escape sequence or literal character <i>c</i>
<code>^</code>	matches the beginning of a string
<code>\$</code>	matches the end of a string
<code>.</code>	matches any single character
<code>[abc...]</code>	character class: matches any of <i>abc...</i>
<code>[^abc...]</code>	negated class: matches any single character but <i>abc...</i>
<code>r₁ r₂</code>	alternation: matches any string matched by <i>r₁</i> or <i>r₂</i>
<code>(r₁)(r₂)</code>	concatenation: matches <i>xy</i> where <i>r₁</i> matches <i>x</i> and <i>r₂</i> matches <i>y</i>
<code>(r)*</code>	matches zero or more consecutive strings matched by <i>r</i>
<code>(r)+</code>	matches one or more consecutive strings matched by <i>r</i>
<code>(r)?</code>	matches the null string or one string matched by <i>r</i>
<code>(r)</code>	grouping: matches the same strings as <i>r</i>

The operators are listed in increasing precedence. Redundant parentheses in regular expressions may be omitted as long as the precedence of operators is respected.

Escape sequences

These sequences have special meanings in strings and regular expressions.

<code>\b</code>	backspace
<code>\f</code>	formfeed
<code>\n</code>	newline
<code>\r</code>	carriage return
<code>\t</code>	tab
<code>\ddd</code>	octal value <i>ddd</i> , where <i>ddd</i> is 1 to 3 digits between 0 and 7
<code>\c</code>	any other character <i>c</i> literally, e.g., <code>\"</code> for <code>"</code> and <code>\\</code> for <code>\</code>

Limits

Any particular implementation of awk enforces some limits. Here are typical values:

- 100 fields
- 3000 characters per input record
- 3000 characters per output record
- 1024 characters per field
- 3000 characters per `printf` string
- 400 characters maximum literal string
- 400 characters in character class
- 15 open files
- 1 pipe
- double-precision floating point

Numbers are limited to what can be represented on the local machine, e.g., 10^{-38} .. 10^{38} ; numbers outside this range will have string values only.

Initialization, comparison, and type coercion

Each variable and field can potentially be a string or a number or both at any time. When a variable is set by an assignment

```
var = expr
```

its type is set to that of the expression. ("Assignment" includes +=, -=, etc.) An arithmetic expression is of type number, a concatenation is of type string, and so on. If the assignment is a simple copy, as in `v1 = v2`, then the type of `v1` is set to that of `v2`.

In comparisons, if both operands are numeric, the comparison is made numerically. Otherwise, operands are coerced to string if necessary, and the comparison is made on strings. The type of any expression can be coerced to numeric by subterfuges such as

```
expr + 0
```

and to string by

```
expr ""
```

(i.e., concatenation with a null string). The numeric value of an arbitrary string is the numeric value of its numeric prefix.

Uninitialized variables have the numeric value 0 and the string value "". Accordingly, if `x` is uninitialized,

```
if (x) ...
```

is false, and

```
if (!x) ...
```

```
if (x == 0) ...
```

```
if (x == "") ...
```

are all true. But note that

```
if (x == "0") ...
```

is false.

The type of a field is determined by context when possible; for example,

```
$1++
```

implies that `$1` must be coerced to numeric if necessary, and

```
$1 = $1 ", " $2
```

implies that `$1` and `$2` will be coerced to strings if necessary.

In contexts where types cannot be reliably determined, e.g.,

```
if ($1 == $2) ...
```

the type of each field is determined on input. All fields are strings; in addition, each field that contains only a number is also considered numeric.

Fields that are explicitly null have the string value ""; they are not numeric. Non-existent fields (i.e., fields past `NF`) and `$0` for blank lines are treated this way too.

As it is for fields, so it is for array elements created by `split`.

Mentioning a variable in an expression causes it to exist, with the values 0 and "" as described above. Thus, if `arr[i]` does not currently exist,

```
if (arr[i] == "") ...
```

causes it to exist with the value "" and thus the `if` is satisfied. The test

```
if (i in arr) ...
```

determines if `arr[i]` exists without the side effect of creating it.