

# AWK Tham Khảo Nhanh

Pattern, field, array, hàm, xử lý văn bản

## Cơ Bản

### Chạy AWK

```
awk '{ print }' file.txt # in mọi dòng
awk '{ print $1 }' file.txt # in field đầu tiên
awk -F: '{ print $1 }' /etc/passwd # dấu phân tách tùy chỉnh
awk -f script.awk file.txt # chạy từ file
cmd | awk '{ print $2 }' # đầu vào qua pipe
```

### Cấu Trúc Chương Trình

<b>awk 'pattern { action }'</b>	Dạng cơ bản — action chạy khi pattern khớp
<b>BEGIN { ... }</b>	Chạy một lần trước khi xử lý đầu vào
<b>END { ... }</b>	Chạy một lần sau khi xử lý xong tất cả đầu vào
<b>Không có pattern</b>	Action chạy cho mỗi dòng
<b>Không có action</b>	Action mặc định là <b>{ print }</b>

## Pattern & Action

### Các Loại Pattern

```
awk '/error/' file.txt # khớp regex
awk '$3 > 100' file.txt # so sánh
awk 'NR >= 5 && NR <= 10' file.txt # phạm vi dòng
awk '/start/,/end/' file.txt # pattern phạm vi
```

### Tham Khảo Pattern

<b>/regex/</b>	Khớp dòng với regex
<b>\$1 ~ /pat/</b>	Field khớp regex
<b>\$1 !~ /pat/</b>	Field không khớp regex
<b>expr1, expr2</b>	Phạm vi: từ lần khớp đầu đến lần khớp thứ hai
<b>expr1 &amp;&amp; expr2</b>	AND logic
<b>expr1    expr2</b>	OR logic
<b>!expr</b>	NOT logic

## Biến

### Biến Built-in

<b>NR</b>	Số record (dòng) hiện tại
<b>NF</b>	Số field trong record hiện tại
<b>FS</b>	Dấu phân tách field đầu vào (mặc định: khoảng trắng)
<b>OF5</b>	Dấu phân tách field đầu ra (mặc định: dấu cách)
<b>RS</b>	Dấu phân tách record đầu vào (mặc định: newline)
<b>ORS</b>	Dấu phân tách record đầu ra (mặc định: newline)
<b>FILENAME</b>	Tên file đầu vào hiện tại
<b>FNR</b>	Số record trong file hiện tại

### Biến Người Dùng

```
awk '{ total += $1 } END { print total }' file.txt
awk -v threshold=50 '$1 > threshold' file.txt
awk 'BEGIN { count = 0 } /pat/ { count++ }
END { print count }' file.txt
```

## Fields

### Truy Cập Field

<b>\$0</b>	Toàn bộ dòng hiện tại
<b>\$1, \$2, ...</b>	Field thứ nhất, thứ hai, ...
<b>\$NF</b>	Field cuối cùng
<b>\$(NF-1)</b>	Field kế trước cuối

### Dấu Phân Tách Field

```
awk -F, '{ print $2 }' data.csv # dấu phẩy
awk -F'\t' '{ print $1 }' data.tsv # tab
awk 'BEGIN { FS = "[,:]" } { print $1 }' f # nhiều ký tự
awk 'BEGIN { OFS = "," } { print $1, $3 }' f # dấu phân tách đầu ra
```

## Luồng Điều Khiển

### Điều Kiện & Vòng Lặp

```
awk '{ if ($1 > 50) print "high"; else print "low" }' f
awk '{ for (i = 1; i <= NF; i++) print $i }' f
awk '{ i = 1; while (i <= NF) { print $i; i++ } }' f
awk '/skip/ { next } { print }' f # bỏ qua dòng khớp
```

### Câu Lệnh Điều Khiển

<b>if (cond) { ... } else { ... }</b>	Điều kiện
<b>for (i = 0; i &lt; n; i++) { ... }</b>	Vòng lặp for kiểu C
<b>for (key in array) { ... }</b>	Duyệt key của array
<b>while (cond) { ... }</b>	Vòng lặp while
<b>do { ... } while (cond)</b>	Vòng lặp do-while
<b>next</b>	Bỏ qua, chuyển sang record tiếp theo
<b>exit</b>	Dừng xử lý, chạy khối END

## Hàm

### Hàm Tự Định Nghĩa

```
awk 'function max(a, b) {
    return (a > b) ? a : b
}
{ print max($1, $2) }' file.txt
```

### Hàm Số

<b>int(x)</b>	Cắt bỏ phần thập phân
<b>sqrt(x)</b>	Căn bậc hai
<b>sin(x), cos(x)</b>	Hàm lượng giác
<b>log(x), exp(x)</b>	Logarithm tự nhiên và hàm mũ
<b>rand()</b>	Số thực ngẫu nhiên trong khoảng 0 đến 1
<b>srand(seed)</b>	Khởi tạo seed cho bộ tạo số ngẫu nhiên

## Arrays

### Associative Arrays

```
awk '{ count[$1]++ }
END { for (k in count) print k, count[k] }' f
awk '{ arr[NR] = $0 }
END { for (i = NR; i >= 1; i--) print arr[i] }' f
```

### Thao Tác Array

<b>arr[key] = val</b>	Đặt phần tử
<b>arr[key]</b>	Lấy phần tử (tự tạo khi truy cập)
<b>key in arr</b>	Kiểm tra key có tồn tại không
<b>delete arr[key]</b>	Xóa một phần tử
<b>delete arr</b>	Xóa toàn bộ array
<b>for (k in arr)</b>	Duyệt qua các key (không theo thứ tự)
<b>length(arr)</b>	Số phần tử (gawk)

## Hàm Chuỗi

### Tham Khảo Chuỗi

<b>length(s)</b>	Độ dài chuỗi
<b>substr(s, start, len)</b>	Chuỗi con (chỉ số bắt đầu từ 1)
<b>index(s, target)</b>	Vị trí của target trong s (0 nếu không tìm thấy)
<b>split(s, arr, sep)</b>	Tách chuỗi thành array
<b>sub(pat, repl, s)</b>	Thay thế lần khớp đầu tiên
<b>gsub(pat, repl, s)</b>	Thay thế tất cả lần khớp
<b>match(s, pat)</b>	Vị trí khớp regex (đặt RSTART, RLENGTH)
<b>tolower(s) / toupper(s)</b>	Chuyển đổi chữ hoa/thường
<b>sprintf(fmt, ...)</b>	Định dạng chuỗi (như C printf)

## Ví Dụ Chuỗi

```
awk '{ gsub(/old/, "new"); print }' f # thay thế kiểu sed
awk '{ print toupper($0) }' f # viết hoa tất cả
awk '{ print substr($0, 1, 40) }' f # cắt bớt 40 ký tự
```

## I/O

### Đầu Ra

```
awk '{ print $1, $2 }' f # cách nhau bởi dấu cách
awk '{ printf "%s,%d\n", $1, $2 }' f # đầu ra có định dạng
awk '{ print $1 > "out.txt" }' f # chuyển hướng ra file
awk '{ print $1 >> "out.txt" }' f # thêm vào file
```

### Tham Khảo I/O

<b>print</b>	In với ORS (newline mặc định)
<b>printf fmt, ...</b>	In có định dạng (không có newline cuối)
<b>print &gt; file</b>	Chuyển hướng đầu ra ra file
<b>print &gt;&gt; file</b>	Thêm vào cuối file
<b>print   cmd</b>	Pipe đầu ra vào lệnh
<b>getline &lt; file</b>	Đọc một dòng từ file
<b>cmd   getline var</b>	Đọc đầu ra lệnh vào biến
<b>close(file)</b>	Đóng file hoặc pipe

## Các Pattern Thường Gặp

### One-Liners

```
awk '{ sum += $1 } END { print sum }' f # tổng cột
awk 'END { print NR }' f # đếm dòng
awk '!seen[$0]++' f # bỏ trùng lặp
awk 'NF' f # bỏ dòng trống
awk '{ print NF }' f # số field mỗi dòng
```

### Công Thức

<b>CSV to TSV</b>	<b>awk -F, 'BEGIN{OFS="\t"} {\$1=\$1; print}'</b>
<b>Tổng cột 2</b>	<b>awk '{ s += \$2 } END { print s }'</b>
<b>Top N dòng</b>	<b>awk 'NR &lt;= 10' (như head)</b>
<b>Đếm tần suất</b>	<b>awk '{ c[\$1]++ } END { for (k in c) print k, c[k] }'</b>
<b>Giữa các marker</b>	<b>awk '/BEGIN/,/END/'</b>
<b>In field thứ N</b>	<b>awk '{ print \$N }' (thay N)</b>