

Das Kommandozeilenwerkzeug AWK

Alfred v. Aho, Peter J. Weinberger und Brian
W. Kernigham

Kathrin Kulmus

Technische Universität Chemnitz
Institut für Physik
THUS

10. Januar 2020

- **AWK & Shell** sind Skriptsprachen
- UNIX bietet viele Werkzeuge zur Datenverarbeitung:
sed, grep, paste, join, cut, diff, head, tail, sort, ...
- in **UNIX-Shell** können diese zu Programmen kombiniert werden
Glue-Language
- **AWK** ist vollständige Programmiersprache, die kaum auf dem UNIX Werkzeugkasten zurückgreift

Programm	Skript
wird <i>einmal</i> durch Compiler übersetzt	wird <i>jedesmal</i> durch Interpreter übersetzt

AWK-Skript wird *jedesmal* in internen Byte-Code übersetzt

- + Wiederholt ausgeführter Code wird nur *einmal* übersetzt und dann sehr schnell ausgeführt
- + keine Syntaxfehler zur Laufzeit mehr möglich
- kleine Programme: Übersetzungszeit länger als Laufzeit

Allgemeine Eigenschaften

- ▶ fasst Fähigkeiten von UNIX-Tools zusammen
- ▶ Automatische Speicherverwaltung
- ▶ Verwandt mit C
- ▶ wenige Datentypen: Text oder Dezimalzahl
- ▶ Keine Variablen-Deklaration nötig
- ▶ kann Reguläre Ausdrücke

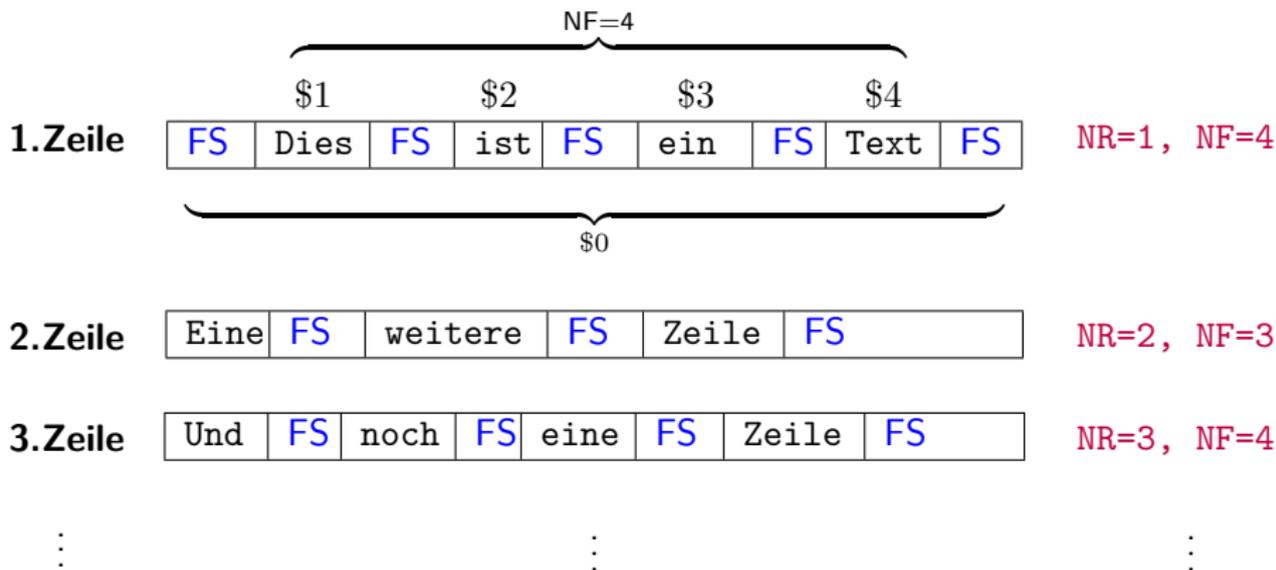
Typische Einsatzgebiete

- ▶ Datenauswertung und -aufbereitung
- ▶ Datentransformation, Texte umformatieren
- ▶ Datenvalidierung (z.B. Sytax, Semantik korrekt?)
- ▶ Prototyping und Realisierung von Programmen

```

awk  'BEGIN      { Aktion... }\  
      Bedingung1 { Aktion... }\  
      Bedingung2 { Aktion... }\  
      .....  
      BedingungN { Aktion... }\  
      END        { Aktion... }\  
file
  
```

- Interpreter ist `awk` „'“ schützen vor Interpretation durch Shell
- `BEGIN`, i.e. vor Verarbeitung der ersten Eingabezeile und `END`, i.e. nach Verarbeitung der letzten Eingabezeile sind optional
- Aktionen stehen in „{ }“
- mehrere Aktionen durch „;“ oder „\n“ trennen



Automatische Leseschleife:

- 0.) Datei Zeilenweise einlesen
- 1.) BEGIN Aktion ausführen
- 2.) Bedingungen der Reihe nach überprüfen
 - ▶ Bedingung für aktuelle Zeile (\$0) erfüllt → Aktion(en) ausführen
- 3.) alle Bedingungen geprüft
 - ▶ Nächste Zeile einlesen und Start bei 2.)
- 4.) END Aktion ausführen

 Keine Aktion: Zur Bedingung passende Zeilen werden ausgegeben

 Keine Bedingung: Aktion wird auf alle Zeilen angewandt

Operator	Bedeutung
(...)	Klammerung
\$	Feldzugriff(\$1,\$2, ... , \$NF)
++ --	Inkrement/Dekrement
^	Potenzierung x^y
+ - !	Unäres Plus/Minus, logisch Nicht
* / %	Multiplikation, Division, Modulo
+ -	Addition, Subtraktion
== != < (=) > (=)	Vergleich (Zahlen & Strings)
~ !~	Vergleich (Regex)
&&	logisches und / oder
= += -= *= /= %=	Operation + Zuweisung

- **Reguläre Ausdrücke** *regex* beschreiben eine Menge von Zeichenketten
- in **awk** durch „/regex/“ als Bedingung
- Folge von Nummern und/oder Buchstaben, Metazeichen, Zeichenlisten, Escape-Sequenzen
- Vorteil gegenüber Zeichenketten („foo“):
 - Zeichenketten schwieriger zu schreiben/lesen
 - *regex* in **awk** effizienter (wird intern ohne Umwandlung gespeichert)

regex	Bedeutung
<code>\n \t \v</code>	Zeilenumbruch, horizontaler/vertikaler Tabulator
<code>\</code>	Unterdrückt Sonderbedeutung: <code>/\„Hallo\“/</code>
<code>^, \$, .</code>	Zeilenanfang, -ende, beliebiges Zeichen <code>/^.\$/</code>
<code>[...], [^ ...]</code>	positive/negative Zeichenliste
<code>[:alnum:]</code>	„A-Za-z0-9“
<code>[:alpha:]</code>	„A-Za-z“
<code>[:lower:]</code>	„a-z“
<code>[:digit:]</code>	„0-9“
<code>\w, \W gawk</code>	„ <code>[:alnum:]-</code> “, „ <code>^[:alnum:]-</code> “

Standardvariable	Bedeutung
NF	Anzahl der Felder
NR	Anzahl der Zeilen insgesamt
FNR	Nummer der aktuell bearbeiteten Zeile
FS	Feldtrenner („ “)
RS	Zeilentrenner („\n“)
FILENAME	aktueller Dateiname

Standardfunktionen	Bedeutung
print <i>LISTE</i>	<i>LISTE</i> ausgeben
length(<i>TEXT</i>)	Länge von <i>TEXT</i>
substr(<i>TEXT</i> , <i>START</i> , <i>LEN</i>)	Teilstück von <i>TEXT</i>
sub(<i>REGEX</i> , <i>TEXT</i> , <i>VAR</i>)	<i>REGEX</i> in <i>VAR</i> durch <i>TEXT</i> 1x ersetzen
gsub(<i>REGEX</i> , <i>TEXT</i> , <i>VAR</i>)	Analog aber global

man awk, man gawk, info awk, info gawk

<https://www.ostc.de/awk.pdf>

https://www-user.tu-chemnitz.de/~hot/unix_linux_werkzeugkasten/awk.html

<https://linux-club.de/wiki/opensuse/Awk>

<http://mikiwiki.org/wiki/awk>



Beispiele

Bearbeite das Gedicht:

1. Welche zeilen sind leer? Entferne alle Leerzeilen, speichere den Text in jandl2.txt
2. Spiele mit dem Feldtrenner: ändere den Standartwert, gib die neue Anzahl der Felder und den Text ohne Feldtrenner aus.
 behandle das Gedicht als Einzeiler/Einspalter (Gib NF,NR als Test aus)
3. analysiere das Gedicht:
 Wie viele Wörter, Buchstaben, Zeichen, Zeilen?
4. Ersetze jedes o durch ein i.
5. Wie oft kommt das Wort mops vor?

1. Alle Spiele mit Auswärtssieg oder mehr als 30000 Zuschauern
2. Wo waren die meisten Zuschauer und wie viele?
3. Wo waren die Werbeeinnahmen je Zuschauer am geringsten, wie wenig?
4. Statistik:

Es gab NUM Spiele.

Zuschauer (gesamt) \pm Änderung (gesamt): ANZAHL \pm DIFF

Zuschauer pro Spiel (Durchschnitt) \pm Änderung (Durchschnitt): ANZAHL2
 \pm DIFF2

Tore (gesamt): TORE

Tore pro Spiel (Durchschnitt): TORE2