

Wprowadzenie do informatyki



Jerzy Nawrocki
Wydział Informatyki
Politechnika Poznańska
jerzy.nawrocki@put.poznan.pl

Przetwarzanie tekstów i AWK

Wprowadzenie do informatyki

Problem konwersji plików

```
FName:Jurek SName:Busz Salary 585
FName:Alek SName:Gor Salary 700
```

↓

FName	SName	Salary
Jurek	Busz	585
Alek	Gor	700

Przetwarzanie tekstów i AWK (2)

Wprowadzenie do informatyki

Problem konwersji plików

```
#include <stdio.h>
#include <stdlib.h>
FILE *fin;
char token[200];
char gettoken(void)
{int i=0;
char c;
do {c =getc(fin);
if (c == EOF) return (EOF);
} while (c < ' ');
```

Rozwiązanie w C:
≈ 40 linii kodu

Przetwarzanie tekstów i AWK (3)

Wprowadzenie do informatyki

Problem konwersji plików


Rozwiązanie w AWK

```
BEGIN {FS=":| ";}
NR == 1 {print $1, "\t", $3, "\t", $5;}
{gsub(/,,".", $6); print $2, "\t", $4, "\t", $6;}
```

Przetwarzanie tekstów i AWK (4)

Wprowadzenie do informatyki

Powstanie języka



Bell Labs, Murray Hill (New Jersey), Foto: http://en.wikipedia.org/wiki/Bell_Labs

Bell Labs, New Jersey (USA), 1977

AWK: Aho, Weinberger, Kernighan
Platformy: Unix, MS DOS/Windows
Podobienstwo do C

Wprowadzenie do informatyki

Cel wykładu



Przedstawić:

- Inny paradygmat programowania (programowanie regułowe)
- Wyrażenia regularne
- Podstawy programowania w języku AWK

Przetwarzanie tekstów i AWK (6)

Wprowadzenie do informatyki

Plan wykładu



- **Idea języka AWK**
- **Najprostsze programy**
- **Wzorce wiersza**
- **Wyrażenia regularne**
- **Zmienne**

Przetwarzanie tekstów i AWK (7)

Wprowadzenie do informatyki

Idea języka AWK – Plik wejściowy

Jerzy Nawrocki	43089	I1
Jan Kowalski	43780	I2
Adam Malinowski	43990	I1

Pole

Wiersz

Pola: \$1, \$2, \$3, ...

Przetwarzanie tekstów i AWK (8)

Wprowadzenie do informatyki

Idea języka AWK – Schemat programu

```
wzorzec1 {instrukcje1}
wzorzec2 {instrukcje2}
... ..
```

Reguła przetwarzania

Przetwarzanie tekstów i AWK (9)

Wprowadzenie do informatyki

Zasada działania


```
Jerzy Nawrocki
Jan Kowalski
Adam Malinowski
```

```
wzorzec1 {instrukcje1}
wzorzec2 {instrukcje2}
... ..
```

Przetwarzanie tekstów i AWK (10)

Wprowadzenie do informatyki

Plan wykładu



- **Idea języka AWK**
- **Najprostsze programy**
- **Wzorce wiersza**
- **Wyrażenia regularne**
- **Zmienne**

Przetwarzanie tekstów i AWK (11)


Wprowadzenie do informatyki

Najprostsze programy

Jerzy Nawrocki	43089	I1
Jan Kowalski	43780	I2
Adam Malinowski	43990	I1

Ile pól na wyjściu?

```
$4=="I1" { print $2, $1; }
```



Przetwarzanie tekstów i AWK (12)

Wprowadzenie do informatyki

Najprostsze programy

Ile pól na wyjściu?

Jerzy Nawrocki	43089	I1
Jan Kowalski	43780	I2
Adam Malinowski	43990	I1

```
$4=="I1"
```

Przetwarzanie tekstów i AWK (13)

Wprowadzenie do informatyki

Najprostsze programy

Jakie pole najpierw?


Jerzy Nawrocki	43089	I1
Jan Kowalski	43780	I2
Adam Malinowski	43990	I1

```
{ print $2, $1; }
```

Przetwarzanie tekstów i AWK (14)

Wprowadzenie do informatyki

Plan wykładu



- Idea języka AWK
- Najprostsze programy
- Wzorce wiersza
- Wyrażenia regularne
- Zmienne

Przetwarzanie tekstów i AWK (15)

Wprowadzenie do informatyki

Wzorce wiersza

- Początek i koniec tekstu
- Relacje
- Wzorce złożone
- ~~Wzorce zakresu~~
- *Wyrażenia regularne*

Przetwarzanie tekstów i AWK (16)

Wprowadzenie do informatyki

Początek i koniec tekstu

Jerzy Nawrocki	43089	I1
Jan Kowalski	43780	I2
Adam Malinowski	43990	I1

```
BEGIN { print "-----"; }
$4=="I2" { print $2, $1; }
END { print "*****"; }
```

Przetwarzanie tekstów i AWK (17)

Wprowadzenie do informatyki

Początek i koniec tekstu

Jerzy Nawrocki	43089	I1
Jan Kowalski	43780	I2
Adam Malinowski	43990	I1

```
END { print "*****"; }
$4=="I2" { print $2, $1; }
BEGIN { print "-----"; }
```

Przetwarzanie tekstów i AWK (18)

Wprowadzenie do informatyki

Relacje

12	11
2	11

`$1 > $2`

Przetwarzanie tekstów i AWK (19)

Wprowadzenie do informatyki

Wzorce złożone

|| lub (alternatywa)
`$1==1 || $2==1`

&& i (koniunkcja)
`$1==1 && $2==1`

! nie (zaprzeczenie)
`! $1==1`

Przetwarzanie tekstów i AWK (20)

Wprowadzenie do informatyki

Wzorce złożone


Jerzy	Adam	43089	I1
Adam	Kowalski	43780	I2
Adam	Malinowski	43990	I1

`$4=="I1" && $1=="Adam" { print $2, $1; }`

Przetwarzanie tekstów i AWK (21)

Wprowadzenie do informatyki

Plan wykładu



- Idea języka AWK
- Najprostsze programy
- Wzorce wiersza
- Wyrażenia regularne
- Zmienne

Przetwarzanie tekstów i AWK (22)

Wprowadzenie do informatyki

Stephen Kleene



1909-01-05, Connecticut, USA

1934: Dr, Princeton Univ., (Alonzo Church)

1935: Univ. of Wisconsin-Madison (USA)

1939-40: Inst. for Advanced Study, Princeton – teoria rekurencji

1990: National Medal of Sci.

1994-01-25, Madison

<http://www.math.wisc.edu/~gpslogic/>

Przetwarzanie tekstów i AWK (23)

Wprowadzenie do informatyki

Wyrażenia regularne

Wyrażenie arytmetyczne

Wartość: Tekst → Liczba

Wartość(`2-3 + 3`) = 9

Wyrażenie regularne

Wartość: Tekst → Zbiór Ciągów Znaków

Wartość(`/Ala | Ola/`) = {"Ala", "Ola"}

Przetwarzanie tekstów i AWK (24)

Wprowadzenie do informatyki

Wzorce z wyrażeniami regularnymi

Róża prawdziwa i sztuczna

Np. znak lub ciąg znaków

\$0, \$1, \$2, ..

Cały ciąg `Ciąg_zn ~ /^ wyr_reg $/`

Szydzi z prawdziwej sztuczna:
 - Krótkie twoje trwanie,
 Wdzięk pani wkrótce minie
 A mój pozostanie...
 - Tak - rzeczce wonna róża,
 Rumieniąc się skromnie –
 Ale patrząc na panią,
 Myśleć będą o mnie!

`$1 ~ /^A$/`

Wprowadzenie do informatyki

Wzorce z wyrażeniami regularnymi

Róża prawdziwa i sztuczna

Np. znak lub ciąg znaków

\$0, \$1, \$2, ..

Cały ciąg `Ciąg_zn ~ /^ wyr_reg $/`

Początek `Ciąg_zn ~ /^ wyr_reg /`

Szydzi z prawdziwej sztuczna:
 - Krótkie twoje trwanie,
 Wdzięk pani wkrótce minie
 A mój pozostanie...
 - Tak - rzeczce wonna róża,
 Rumieniąc się skromnie –
 Ale patrząc na panią,
 Myśleć będą o mnie!

`$1 ~ /^A/`

Wprowadzenie do informatyki

Wzorce z wyrażeniami regularnymi

Róża prawdziwa i sztuczna

Np. znak lub ciąg znaków

\$0, \$1, \$2, ..

Cały ciąg `Ciąg_zn ~ /^ wyr_reg $/`

Początek `Ciąg_zn ~ /^ wyr_reg /`

Koniec `Ciąg_zn ~ / wyr_reg $/`

Szydzi z prawdziwej sztuczna:
 - Krótkie twoje trwanie,
 Wdzięk pani wkrótce minie
 A mój pozostanie...
 - Tak - rzeczce wonna róża,
 Rumieniąc się skromnie –
 Ale patrząc na panią,
 Myśleć będą o mnie!

`$1 ~ /dzi$/`

Wprowadzenie do informatyki

Wzorce z wyrażeniami regularnymi

Róża prawdziwa i sztuczna

Np. znak lub ciąg znaków

\$0, \$1, \$2, ..

Cały ciąg `Ciąg_zn ~ /^ wyr_reg $/`

Początek `Ciąg_zn ~ /^ wyr_reg /`

Koniec `Ciąg_zn ~ / wyr_reg $/`

Podciąg `Ciąg_zn ~ / wyr_reg /`

Szydzi z prawdziwej sztuczna:
 - Krótkie twoje trwanie,
 Wdzięk pani wkrótce minie
 A mój pozostanie...
 - Tak - rzeczce wonna róża,
 Rumieniąc się skromnie –
 Ale patrząc na panią,
 Myśleć będą o mnie!

`$1 ~ /dzi/`

Wprowadzenie do informatyki

Wzorce z wyrażeniami regularnymi

Róża prawdziwa i sztuczna

Np. znak lub ciąg znaków

\$0, \$1, \$2, ..

Cały ciąg `Ciąg_zn ~ /^ wyr_reg $/`

Początek `Ciąg_zn ~ /^ wyr_reg /`

Koniec `Ciąg_zn ~ / wyr_reg $/`

Podciąg `Ciąg_zn ~ / wyr_reg /`

`$0 ~ / wyr_reg / = / wyr_reg /`

Szydzi z prawdziwej sztuczna:
 - Krótkie twoje trwanie,
 Wdzięk pani wkrótce minie
 A mój pozostanie...
 - Tak - rzeczce wonna róża,
 Rumieniąc się skromnie –
 Ale patrząc na panią,
 Myśleć będą o mnie!

`/dzi/`

Wprowadzenie do informatyki

Znaki specjalne

- `.` dowolny znak
- `[]` zbiór znaków
- `\n` nowa linia
- `\.` kropka
- `\"` znak cudzozyłowy
- `\\ddd` znak o kodzie oktalnym `ddd`

Przetwarzanie tekstów i AWK (30)

Wprowadzenie do informatyki

Znaki specjalne

Co robi ten program?

```
/^,$/
```

```
/[0123456789]/
```

```
/[0-9]/
```

Przetwarzanie tekstów i AWK (31)

Wprowadzenie do informatyki

Dopełnienie zbioru znaków

Co za różnica?

```
[^ ... ]
```

```
/[^0-9]/
```

```
/^0-9/
```

Przetwarzanie tekstów i AWK (32)

Wprowadzenie do informatyki

Alternatywa (lub)

Co robi ten program?

```
wyr_reg | wyr_reg
```

Ząbki listka
Nic tak nie potrafi gryźć
Jak cudzy laurowy liść.

```
/fi | ki/
```

Przetwarzanie tekstów i AWK (33)

Wprowadzenie do informatyki

Nawiasy i alternatywa

Nawiasy podnoszą priorytet operatorów:
 $3*(4 + 5) = 3*9 = 27$

Rozdzielność mnożenia względem dodawania:
 $3*(4 + 5) = 3*4 + 3*5 = 12 + 15 = 27$
 $(4 + 5)*3 = 4*3 + 5*3 = 12 + 15 = 27$

Rozdzielność konkatenacji względem alternatywy

```
/(t | k)/ = /ti | ki/
```

```
/Adam(e | i)k/ = /Adamek | Adami/
```

Przetwarzanie tekstów i AWK (34)

Wprowadzenie do informatyki

Operator powtarzania

Wypisz wszystkie wiersze, w których pierwsze pole jest liczbą.

```
1 2 3  
Rok 1984  
4 damy
```

```
$/1 ~ /[0-9]/
```

Przetwarzanie tekstów i AWK (35)

Wprowadzenie do informatyki

Operator powtarzania

Wypisz wszystkie wiersze, w których pierwsze pole jest liczbą.

```
1 2 3  
Rok 1984 2 raki  
4 damy
```

```
$/1 ~ /[0-9]/
```

```
1 2 3  
4 damy
```

Przetwarzanie tekstów i AWK (36)

Wprowadzenie do informatyki

Operator powtarzania *

Może wystąpić sekwencja w.

Musi przynajmniej raz wystąpić w.

$w^+ = w | ww | www | \dots$

$w^* = \epsilon | w | ww | www | \dots$

Ciąg pusty

$w(\epsilon | w | ww | www | \dots) = w | ww | www | wwwww | \dots$

$x \epsilon = x$

$\epsilon x = x$

$w^+ = w w^* = w^* w$

$x w^* = x | x w^+$

$w^* x = x | w^+ x$

Przetwarzanie tekstów i AWK (43)

Wprowadzenie do informatyki

Zagadka

Co robi ten program?


```
Nadgodziny 2001/02
=====
Nawrocki 60
Complak 359
```

$\$2 \sim /^{[0-9]}+$/$

Przetwarzanie tekstów i AWK (44)

Wprowadzenie do informatyki

Plan wykładu



- Idea języka AWK
- Najprostsze programy
- Wzorce wiersza
- Wyrażenia regularne
- Zmienne

Przetwarzanie tekstów i AWK (45)

Wprowadzenie do informatyki

Zmienne

- Zmienne wprowadzone przez programistę
(*typ*: ciąg znaków;
wartość początkowa: ciąg pusty / zero)
- Zmienne wbudowane
(mają standardowe znaczenie)
- Zmienne polowe \$1, \$(i+j-1), ..

Przetwarzanie tekstów i AWK (46)

Wprowadzenie do informatyki

Przykładowe zmienne wbudowane

NF - liczba pól w wierszu
NR - numer wiersza
FILENAME - nazwa pliku

Przetwarzanie tekstów i AWK (47)

Wprowadzenie do informatyki

Zmienne

NR	NF	total
1	2	0
		2

Reguły sitwy

```
Zero do zera
A będzie kariera.
```

```
{total= total + NF;}
END {print "Pol: ", total;
print "Wierszy: ", NR;}
```

Przetwarzanie tekstów i AWK (48)

Wprowadzenie do informatyki



Podsumowanie

Przetwarzanie tekstów i AWK (49)

Wprowadzenie do informatyki

Cel wykładu

Przedstawić:

- Inny paradygmat programowania (programowanie regułowe)



Przetwarzanie tekstów i AWK (50)

Wprowadzenie do informatyki

Idea języka AWK – Plik wejściowy

Jerzy Nawrocki	43089	I1
Jan Kowalski	43780	I2
Adam Malinowski	43990	I1

Pole: \$1, \$2, \$3, ...

Przetwarzanie tekstów i AWK (51)

Wprowadzenie do informatyki

Idea języka AWK – Schemat programu

```
wzorzec1 {instrukcje1}
wzorzec2 {instrukcje2}
... ..
```

Reguła przetwarzania

Przetwarzanie tekstów i AWK (52)

Wprowadzenie do informatyki

Zasada działania

→ Jerzy Nawrocki
Jan Kowalski
Adam Malinowski

→ wzorzec1 {instrukcje1}
wzorzec2 {instrukcje2}
... ..


Przetwarzanie tekstów i AWK (53)

Wprowadzenie do informatyki

Cel wykładu

Przedstawić:

- Inny paradygmat programowania (programowanie regułowe)
- Wyrażenia regularne



Przetwarzanie tekstów i AWK (54)

Wprowadzenie do informatyki

Wyrażenia regularne

Wyrażenie arytmetyczne

Wartość: Tekst → Liczba

Wartość(2·3 + 3) = 9

Wyrażenie regularne

Wartość: Tekst → Zbiór Ciągów Znaków

Wartość(/Ala | Ola) = {"Ala", "Ola"}

Przetwarzanie tekstów i AWK (55)

Wprowadzenie do informatyki

Wyrażenia regularne

$\epsilon x = x \epsilon$

$L_1 \cdot L_2 = \{xy: x \in L_1, y \in L_2\}$

$L^0(r) = \{\epsilon\}$

$L^{n+1}(r) = L(r) \cdot L^n(r)$

$L^n(r) = \{\underbrace{xx\dots x}_n: x \in L(r)\}$

$L(r^*) = \cup L^n(r)$

$L(r^*) = \{\epsilon, x, xx, \dots, x^n: x \in L(r)\}$ $(01)^* = \{\epsilon, 01, 0101, 010101, \dots\}$

Przetwarzanie tekstów i AWK (56)

Wprowadzenie do informatyki

Cel wykładu




Przedstawić:

- Inny paradygmat programowania (programowanie regułowe)
- Wyrażenia regularne
- Podstawy programowania w języku AWK

Przetwarzanie tekstów i AWK (57)

Wprowadzenie do informatyki

Podsumowanie



Wreszcie!

`gawk -f prog.awk <in.txt >out.txt`


Inne mechanizmy AWK:

- Instrukcje złożone (if, while, ..)
- Tablice dynamiczne
- Funkcje wbudowane (gsub, ..)

Przetwarzanie tekstów i AWK (58)

Wprowadzenie do informatyki

Literatura



- A. Aho, B. Kernighan, P. Weinberger, *The AWK Programming Language*, Addison-Wesley, Reading, 1988.
- J. Nawrocki, W. Complak, Wprowadzenie do przetwarzania tekstów w języku AWK, *Pro Dialog 2* (1994), 23-46.

Przetwarzanie tekstów i AWK (59)

