

# Prednáška 12 - Sémantická analýza

Ing. Viliam Hromada, PhD.

C-510  
Ústav informatiky a matematiky  
FEI STU

`viliam.hromada@stuba.sk`















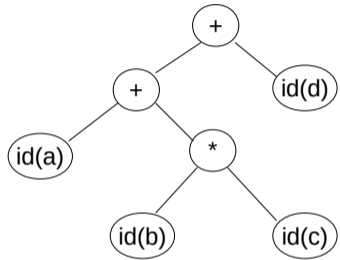
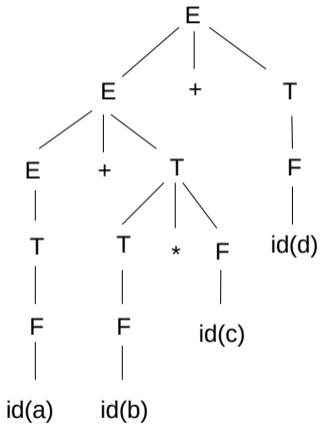








### Príklad AST č. 1







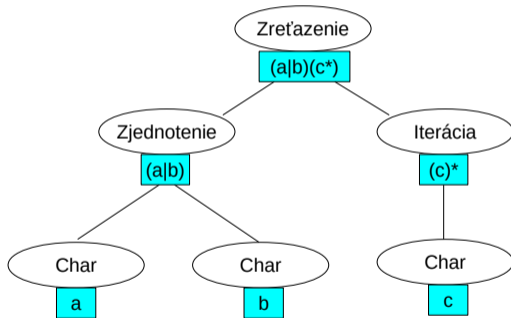








Iný abstraktný syntaktický strom by mohol vyzeráť nasledovne (sémantické atribúty sú teraz znázornené v modrých obdĺžnikoch a predstavujú regulárny výraz, ktorý je asociovaný s príslušným vrcholom AST).











# Medzikód

Jeho použitie má 2 hlavné výhody:

- Možnosť optimalizácie na úrovni medzikódu t.j. nezávislej od cieľového jazyka / platformy. Preto je táto optimalizácia ľahko prenositeľná medzi platformami a zväčša aj jednoduchšia ako optimalizácia na úrovni strojového jazyka.
- Jednoduchšia prenositeľnosť kompilátorov medzi jazykmi / platformami.





3 často uvažované formy medzikódu sú:

1. Abstraktný syntaktický strom
2. Postfixový zápis
3. Jazyk štvoric (niekedy nazývaný aj trojadresný kód, resp. three-address code)









# Jazyk štvoríc - typické štvorice

Štvorica	Popis
(operátor, <i>op1</i> , <i>op2</i> , výsledok)	$vysledok \leftarrow op1 \text{ operátor } op2$ , pričom $operátor \in \{+, -, *, /, AND, OR, =, \neq, <, <=, >, >=\}$ (binárne aritmetické, logické a relačné operátory)
(-, <i>op1</i> , výsledok)	$vysledok \leftarrow -op1$ (unárne mínus)
(NOT, <i>op1</i> , výsledok)	$vysledok \leftarrow \text{not}(op1)$ (negácia)
(:=, <i>op1</i> , <i>op2</i> , výsledok)	$vysledok \leftarrow op1, op2$ voliteľne obsahuje počet prenesených bajtov (presun/pri- radenie)
(FLOAT, <i>op1</i> , výsledok)	$vysledok \leftarrow \text{Float}(op1)$ (konverzia na typ Float)
(INTEGER, <i>op1</i> , výsledok)	$vysledok \leftarrow \text{Integer}(op1)$ (konverzia na typ Integer)
(LABEL, návestie)	návestie sa použije ako návestie nasle- dujúcej inštrukcie
(JUMP, návestie)	nepodmienený skok na inštrukciu s ná- vestím <i>návestie</i>
(JUMPT, <i>op1</i> , návestie)	podmienený skok na inštrukciu s náves- tím <i>návestie</i> , ak <i>op1</i> obsahuje logickú hodnotu TRUE
(JUMPF, <i>op1</i> , návestie)	podmienený skok na inštrukciu s náves- tím <i>návestie</i> , ak <i>op1</i> obsahuje logickú hodnotu FALSE
(RANGETEST, dol, hor, <i>op1</i> )	testuje, či $dol \leq op1 \leq hor$
(INDEX, pole, <i>i</i> , prvok)	$prvok \leftarrow \text{adresa prvku pola pole}[i]$



## Jazyk štvoríc - príklad č. 2

Výraz  $D := b * b - 4 * a * c$ , kde  $D, b, c$  sú typu Float,  $a$  je typu Integer, 4 je konštanta typu Int, by zapísaný v jazyku štvoríc vyzeral:

(\*, b, b, t1)

(FLOAT, 4, t2)

(-, t2, t3)

(FLOAT, a, t4)

(\*, t3, t4, t5)

(\*, t5, c, t6)

(+, t1, t6, t7)

(:=, t7, sizeof(Float), D)





















## Sémantické záznamy

### Definícia

*Sémantický záznam zodpovedajúci ľavej strane pravidla označujeme symbolom  $$$$  a sémantické záznamy terminálnych a neterminálnych symbolov pravej strany pravidla označujeme zľava-doprava symbolmi  $\$1, \$2, \dots, \$n$ , pričom  $n$  je dĺžka pravej strany pravidla.*











## Sémantické podprogramy v Príklade - PRÍKAZY1

Sémantický podprogram PRÍKAZY1 sa vyvolá na konci 2. pravidla reťazca. V rámci tohoto pravidla sa generuje nový príkaz, preto sémantický podprogram zvýši počet príkazov generovaných z <príkazy> na pravej strane pravidla (\$3) o 1 a uloží ho v SZ asociovanom s neterminálom <príkazy> na ľavej strane, t.j. v \$\$:

**funkcia PRÍKAZY1**

$$$ \leftarrow \$3 + 1;$

**koniec funkcia**



## Sémantické podprogramy v Príklade - PRÍKAZY0

Sémantický podprogram PRÍKAZY0 sa vyvolá na konci 3. pravidla reťazca. V rámci tohoto pravidla sa negeneruje nový príkaz, preto sémantický podprogram uloží **nulu** do SZ asociovaného s neterminálom <príkazy> na ľavej strane pravidla, t.j. \$\$:

**funkcia PRÍKAZY0**

\$\$ ← 0;

**koniec funkcia**



# Znázornenie sémantického spracovania

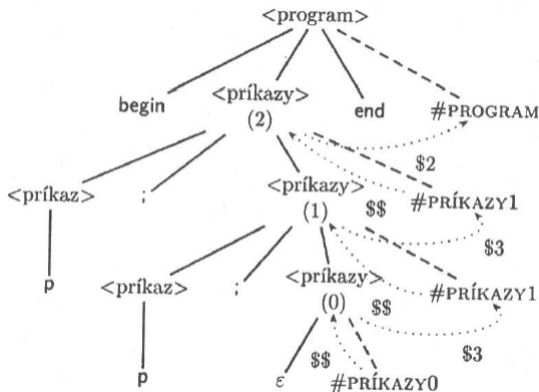


Figure: Prevzaté z Dедера, L.: Počítačové jazyky a ich spracovanie



## Výsledok príkladu

- Je dôležité si uvedomiť, že sémantiku sme si zvolili my, t.j. chceli sme spočítať, koľkokrát sa vo vstupnom reťazci "nachádza volanie programu", t.j. terminál  $p$ .
- Na to slúžili sémantické podprogramy PROGRAM, PRÍKAZY1, PRÍKAZY0.
- Všimnite si, že symboly gramatiky v sémantickom zásobníku medzi sebou komunikujú prostredníctvom svojich sémantických záznamov!
- Napr. pri redukcii podľa pravidla:  $\langle \text{príkazy} \rangle \rightarrow \langle \text{príkaz} \rangle; \langle \text{príkazy} \rangle$  sú vlastne príslušné SZ označené ako  $\$ \$ \rightarrow \$1 \$2 \$3$ .



# Výsledok príkladu

- Preto ak je v sémantickom podprograme PRÍKAZY1 napísané, že  $$$ = \$3 + 1$ , znamená to, že do SZ pre symbol \$\$ sa vloží hodnota SZ pre symbol \$3 zvýšená o 1.







## Sémantické spracovanie - Kalkulačka

- Úlohou sémantickej analýzy bude v tomto príklade **vyhodnotiť výraz**.
- T.j. zadaný výraz sa rozparsuje a pre jednotlivé vrcholy derivačného stromu budú pridelené sémantické atribúty znamenajúce hodnotu podvýrazu.
- Hodnota sémantického atribútu pri koreni derivačného stromu bude potom hodnota celého výrazu.
- Výraz sa bude vyhodnocovať pomocou sémantických podprogramov.
- Akčné symboly sú vložené na **konci každého pravidla**, pričom im prislúchajú nasledovné sémantické akcie:



## Sémantické spracovanie - Kalkulačka

Č.	Pravidlo	Sémantická akcia
1	$S \rightarrow E$	Výsledok výrazu je \$1.
2	$E \rightarrow E + T$	$$$ \leftarrow \$1 + \$3$
3	$E \rightarrow E - T$	$$$ \leftarrow \$1 - \$3$
4	$E \rightarrow T$	$$$ \leftarrow \$1$
5	$T \rightarrow T * F$	$$$ \leftarrow \$1 * \$3$
6	$T \rightarrow T / F$	$$$ \leftarrow \$1 / \$3$
7	$T \rightarrow F$	$$$ \leftarrow \$1$
8	$F \rightarrow (E)$	$$$ \leftarrow \$2$
9	$F \rightarrow k$	$$$ \leftarrow$ konštanta vrátená lexikálnym analyzátorom



## Sémantické spracovanie - Kalkulačka

Simulujme výpočet pre vstupný reťazec  $5 + 3 * 6$  predstavujúci matematický výraz.

- Lexikálny analyzátor (LexA) rozdelí reťazec  $5 + 3 * 6$  na postupnosť lexém  $k + k * k$ .
- Syntaktický analyzátor (SynA) nájde derivačný strom postupnosti  $k + k * k$ .
- Sémantický analyzátor (SemA) vyhodnotí príslušné sémantické atribúty a nájde hodnotu vstupného výrazu (t.j. mal by vrátiť hodnotu 23).



## Obrázok sémantického spracovania podľa derivačného stromu:

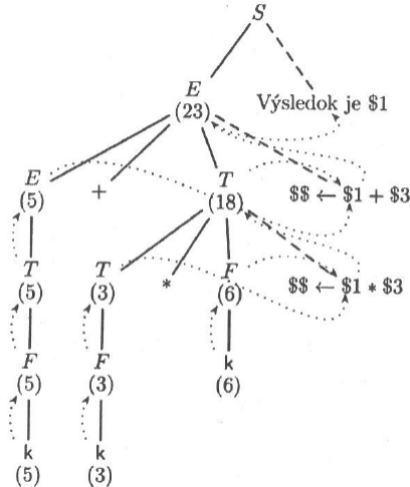


Figure: Prevzaté z Dедера, L.: Počítačové jazyky a ich spracovanie



## Tvorba jazykových procesorov pomocou nástrojov Flex a Bison

- Jazykový procesor, resp. kompilátor so syntaktickým analyzátorom ako ústrednou časťou, je možné vytvárať aj pomocou špeciálnych programov - generátorov jazykového procesora.
- Generátor jazykového procesora vytvorí jazykový procesor, ktorý obsahuje:
  - Lexikálny analyzátor na základe popisu lexém,
  - Syntaktický analyzátor na základe popisu bezkontextovej gramatiky,
  - Správu volaní sémantických podprogramov,
  - Správu sémantického zásobníka.
- Typický nástroj pre konštrukciu lexikálneho analyzátoru je **Flex** (vid' prednáška č. 6).
- Typický generátor jazykového procesora je nástroj **Bison**.



## Nástroj Bison

- Bison predstavuje generátor jazykového procesora vychádzajúci z nástroja Yacc.
- Generuje syntaktický analyzátor LALR(1), konflikty rieši pomocou priority asociatívnosti operátorov, resp. aj GLR prístupom.
- Predstavuje teda príklad sémantického spracovania pri analýze zdola-nahor.
- Spolupracuje s nástrojom Flex.
- Jeho vstupný súbor (zvyčajne s príponou `.y`) obsahuje popis gramatiky, ku ktorej sa jazykový procesor tvorí, spolu s časťami kódu v jazyku C, ktoré sa vykonávajú, keď sa rozpozná pravá strana príslušného pravidla (sémantické podprogramy).



## Nástroj Bison

- Tieto fragmenty kódu (sémantické podprogramy) používajú symboliku \$\$, \$1, ... na prístup k sémantickým záznamom sémantického zásobníka.
- Výstupom z nástroja Bison je kód jazykového procesora v jazyku C v 2 súboroch s koncovkami: `.tab.c` a `.tab.h`. Tie sa dajú skompilovať kompilátorom jazyka C a prilinkovať k ďalším programom.





## Nástroj Bison - vstupný súbor

Vstup nástroja Bison predstavuje súbor s koncovkou `.y` obsahujúci tieto 3 časti:

1. Sekciu deklarácií
2. Sekciu pravidiel
3. Sekciu podprogramov

Tieto sekcie sú navzájom oddelené `%%`. Ich obsah si priblížime na príklade.



## Nástroj Bison - príklad

- Zostrojme jazykový procesor, ktorý vykoná spracovanie gramatiky zo slajdu č. 40 s rovnakou sémantickou analýzou, t.j. spočíta, koľko "príkazov"  $p$  obsahuje "program", t.j. vstupný reťazec. Gramatika:
  1.  $\langle \text{program} \rangle \rightarrow \text{begin} \langle \text{príkazy} \rangle \text{end}$
  2.  $\langle \text{príkazy} \rangle \rightarrow \langle \text{príkaz} \rangle ; \langle \text{príkazy} \rangle$
  3.  $\langle \text{príkazy} \rangle \rightarrow \epsilon$
  4.  $\langle \text{príkaz} \rangle \rightarrow p$



## Bison - sekcia deklarácií

Sekcia deklarácií vstupu nástroja Bison môže obsahovať:

- Kód (v jazyku C), ktorý sa vloží do výsledného súboru na jeho začiatok.
- Definíciu údajového typu sémantického záznamu - konštanta `YYSTYPE`. Ak sa nedefinuje, implicitne je typu `int`.
- Definície terminálnych symbolov gramatiky cez klauzulu `%token`
- Špecifikáciu počiatočného neterminálu gramatiky cez klauzulu `%start`
- Definíciu externej premennej `yytext` ktorá obsahuje reťazec získaný na vstupe zodpovedajúci príslušnej vrátenej lexéme.



## Bison - sekcia deklarácií

```
%{  
#include<stdio.h>  
#define YYSTYPE int  
extern char *yytext;  
void vypis(YYSTYPE argument);  
%}
```

```
%token _BEGIN  
%token _END  
%token _SEMICOLON  
%token _P  
%token NEZNAME
```

```
%start PROGRAM  
%%
```

# Bison - sekcia pravidiel

- Pravidlá sa zadávajú v tvare:

```

A: B1 . . . BN
   | C1 . . . CN
. . .
;

```

- pričom A je ľavá strana pravidiel, B1 . . . BN a C1 . . . CN sú pravé strany pravidla s rovnakou ľavou stranou. Všimnite si zvislú čiaru pri novom pravidle s rovnakou ľavou stranou a bodkočiarku, ktorá ukončuje pravidlá s rovnakou ľavou stranou.
- Všetko, čo nebolo deklarované ako %token v predchádzajúcej časti, je neterminál.
- Symboly **musia** byť oddelené medzerami.



# Bison - sekcia pravidiel

- Pravidlá sa zadávajú v tvare:

```

A: B1 ... BN
  | C1 ... CN
...
;
  
```

- pričom  $A$  je ľavá strana pravidiel,  $B1 \dots BN$  a  $C1 \dots CN$  sú pravé strany pravidla s rovnakou ľavou stranou. Všimnite si zvislú čiaru pri novom pravidle s rovnakou ľavou stranou a bodkočiarku, ktorá ukončuje pravidlá s rovnakou ľavou stranou.
- Všetko, čo nebolo deklarované ako `%token` v predchádzajúcej časti, je neterminál.
- Symboly **musia** byť oddelené medzerami.

## Bison - sekcia pravidiel

- Do pravidiel sa môžu vkladať kusy kódu v jazyku C uzatvorené do {...}, ktorý sa vykoná, keď bude rozpoznaná časť pravej strany pravidla pred blokom ohraničeným {...}.
- T.j. ak sa takýto blok kódu vloží na koniec pravidla, tak sa vykoná v prípade redukcie podľa príslušného pravidla.
- V podstate to nahrádza **akčné symboly**, resp. príslušné sémantické podprogramy.
- V sémantických podprogramoch je možné využívať:
  - Premenné `ytext` a `yyleng` z lexikálneho analyzátoru (rozpoznaný reťazec pre danú lexému, resp. jeho dĺžka)
  - Premenné `$$`, `$1`, ... pre symboly v pravidle, resp. príslušné sémantické záznamy v sémantickom zásobníku.
  - Volania funkcií definovaných v sekcii podprogramov.



## Bison - sekcia pravidiel

```
PROGRAM:  _BEGIN PRIKAZY _END  {vypis($2);}
;
PRIKAZY:  PRIKAZ _SEMICOLON PRIKAZY {$$ = $3 + 1;}
| /**/ {$$ = 0;}
;
PRIKAZ:  _P
;
%%
```





## Bison - sekcia podprogramov

- Tretia sekcia vstupného súboru nástroja Bison obsahuje definície funkcií, ktoré bude obsahovať výsledný jazykový procesor (výstup nástroja Bison).
- Sú tu teda aj definované funkcie, ktoré môžu byť volané v rámci sémantických podprogramov použitých v sekcii pravidiel.
- Samotný syntaktický analyzátor sa volá pomocou funkcie `yyparse()`, ktorá realizuje kompletnú syntaktickú a sémantickú analýzu. V prípade úspešnej syntaktickej analýzy vráti hodnotu 0 (nula), v prípade syntaktickej chyby vráti hodnotu 1.
- Taktiež je vhodné tu vložiť funkciu `main()`, v ktorej sa môže zavolať samotný analyzátor cez volanie `parse()`.



## Bison - sekcia podprogramov

```
void vypis(YSTYPE argument)
{
    printf("Program obsahuje %d prikazov\n", argument);
}
```

```
void main()
{
    if (yyparse() == 0)
        printf("Vyraz bol syntakticky spravny!\n");
    else
        printf("Vyraz nebol syntakticky spravny!\n");
    return ;
}
```



## Bison - vstup a výstup

- Vstupný súbor teda pozostáva z uvedených troch častí - ak ich spojíme dokopy a výsledok nazveme `program.y`, tak spracovanie vstupného súboru pomocou nástroja Bison vykonáme:

```
Bison -d -v program.y
```

- Prepínače `-d` a `-v` vygenerujú hlavičkový súbor `.tab.h`, resp. generujú viacero výpisov.
- Výsledkom sú súbory `program.tab.c` a `program.tab.h` obsahujúce kód výsledného analyzátora v jazyku C.
- Bison taktiež vygeneruje súbor `program.output` obsahujúci popis stavov automatu syntaktického analyzátora, ktorý rozpoznáva pravé strany pravidiel (viď. *LR(0)*-automat).



# Flex - vstupný súbor využívajúci výstup Bisonu

```
%{  
#include "program.tab.h"  
%}  
%%  
  
[\\t ]+  
";" return (_SEMICOLON);  
"p" return (_P);  
[b][e][g][i][n] return (_BEGIN);  
[e][n][d] return (_END);  
[\\n] yyterminate();  
. return (NEZNAME);  
%%
```





## Bison a Flex

- Spustíme Flex so vstupným Flex súborom (nech sa volá `program.l`)  
`Flex program.l`
- Flex vygeneruje kód lexikálneho analyzátoru `lex.yy.c`
- Súbor `lex.yy.c`, `program.tab.c`, `program.tab.h` spolu tvoria výsledný C-kód jazykového procesora.
- Podobných generátorov jazykových procesorov je viacero, pekný prehľad s použitými syntaktickými analyzátormi a s programovacím jazykom výstupu je na [https://en.wikipedia.org/wiki/Comparison\\_of\\_parser\\_generators](https://en.wikipedia.org/wiki/Comparison_of_parser_generators).



