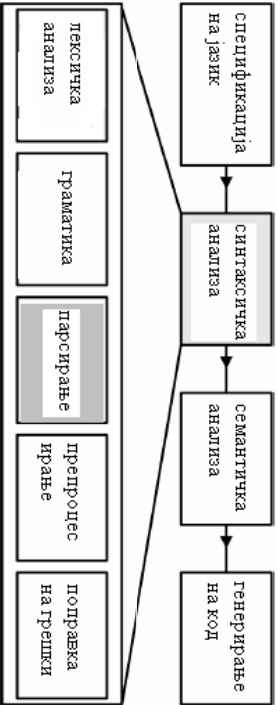


Парсирање

Вовед



Вовед

- Конструиравме граматика за формален јазик
 - За генерирање на валидна реченица треба да се напишат синтаксни дрва и синтаксни дијаграми
 - Компјутерите не работат врз основа на овие дијаграми туку врз основа на разбирање (парсирање) на речениците
- Програмата која ги парсира речениците се нарекува парсер.
 - Некој парсер градат синтаксни дрва за речениците и така ги разбираат
 - Тие се генерираат од конкретна реченица, а не од синтаксен дијаграм

Вовса

- Парсерот е програма која го чита внесениот текст и кажува дали тој ги следи правилата на граматиката или не
 - Ако не треба да каже зошто
- Друг начин на работа на парсерот е да каже дали реченицата може да биде генерирана од граматиката.
- Само посложените компајлери градат синтаксно дрво во нивната меморија
 - градбата на дрвата е илустративно за работата на компајлерот
- Поедноставување на синтаксно дрво со креирање на апстрактно синтаксно дрво (покомпактно)
 - Апстрактното синтаксно дрво е битно за фазата на семантичка анализа и генерирање на код

Префиксен код

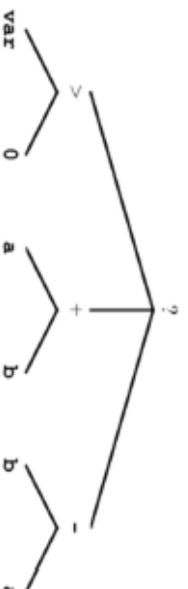
- Префиксната нотација е воведена од Полскиот научник Ј. Lukasiewicz.
- Нотација за која не требаат загради
 - Перфектна за излез од компајлерот
- Пример: $1 + 2 * 3$ станува $+ 1 * 2 3$
 $1 / 2 - 3$ станува $- / 1 2 3$
- При пишување на граматиката треба да се води сметка за асоцијативноста или за предимство на оператори
- Понатамошните фази не треба да водат сметка за тоа
 - Парсерот треба да го конвертира текстот во некој меѓу формат (меѓу јазик) во кој тоа ќе биде вметнато аунтоматично (the parser should convert the input text to an intermediate format)
 - Ова се решава со недвосмислено синтаксно дрво или со префиксна нотација

Префиксен код

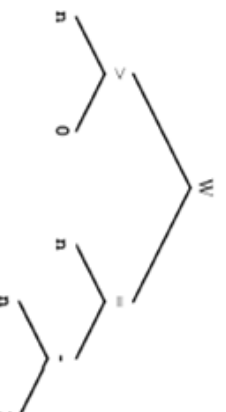
- Со префиксната нотација може да се претстават и комплексни конструкции како *if ... then* и *while ... do*
- Пример
 - if (var > 0) { a + b } else { b – a } станува ?>var'0'+a'b'-b'a'
 - while (n > 0) { n = n – 1 } станува W>n'0'=n'-n'1'
- Апострофите (') се користат за да се одвојат променливите
- Префиксната нотација всушност е слично на дрво и со користење на рекурзивни функции може да се претвори во дрво
- Може да се конструира многу едноставен компајлер кој би превел математички изрази во префиксен код
- Потоа може да се направи втора програма-евалuator-кој може да го интерпретира овој код и да го прорачуна резултатот.

Префиксен код

- if (var > 0) { a + b } else { b – a }



- while (n > 0) { n = n – 1 }



Префиксен код

- Преискната нотација исто така помага во елиминирање на празни места и коментари без да се загуби смислата
- Програмата за евалуација работи рекурзивно
 - Почнува да го чита префиксниот стринг од лево кон десно и за секој оператор кој ќе го изброи се повикува самата себе за да ги добие операндите
 - Рекурзијата завршува кога ќе се најдат константите (иминја на променливи или литерални вредности).
- Евалуаторот работи на принцип што рекурзивната функција може да ја користи како стек во кој ги чува операторите кои сеуште не се евалуирани
- Значајно својство на префиксниот код и соодветното дрво е тоа што сега операторите веќе не се листови во дрвото туку внатрешни јазли
 - Со ова се штеди простор за зачувување на дрвото

Теорија на прасирање

- Како се конструира синтаксното (парсирачко) дрво од даден оригинален влез?
- Гледаме леви граматики
- Два вида на парсери: top-down и bottom-up.
 - Top-down парсерот почнува од стартниот симбол од граматиката и работи над конкретната реченица со замена на правилата на граматиката при што ги заменува нетерминалите до нетерминали
 - Овој метод е полесен но бара големи рестрикции на граматиката
- Bottom-up парсерот почнува со реченицата која ја земајува со нетерминалите од левата страна на правилата.
 - Поможен но потешок да се напише на рака, па за таа цел постојат алатки за конструкција на вакви bottom-up парсери од дадена граматика

Top-down парсирање

- Го зема најлевиот нетерминал и го заменува со тој од десната стран
 - ▢ Ако има повеќе избори ќе треба да се види со што да се замени, зависно од наредниот симбол кој следи
 - ▢ Ако има повеќе можни избори треба да се пробува по сите (постојат вакви парсери) и ако со некоја проба се успее сите нетерминали да се заменат со терминали, тигаш реченицата е добра. Ако во сите можности не се добие сите нетерминали да се заменат, треба компајлерот да јави грешка.
- За контекстно слободна граматика не е битно кој нетерминал ќе се замени, но обично се заменува најлевиот

Top-down парсирање

```
expression: factor restexpression .
restexpression : e.
restexpression : + factor restexpression .
restexpression : - factor restexpression .
factor : 0.
factor : 1.
factor : 2
factor : 3.
factor : 4.
factor : 5.
factor : 6.
factor : 7.
factor : 8.
factor : 9.

expression => factor restexpression
=> 1 restexpression
=> 1 +factor restexpression
=> 1 +2 restexpression
=> 1 +2 +factor restexpression
=> 1 +2 +factor restexpression
=> 1 +2 +3 restexpression
=> 1 +2 +3
```

Top-down парсирање

- Со оваа граматика не може да се испарсира изразот $1 + 2 * 3$
- Парсирањето ќе падне кога ќе се појави знакот *

$expression \Rightarrow_L 1 + 2 \text{ restexpression}$

- Парсерот ќе бара замена која почнува со *, и бидејќи не постои таква нема со што да направи замена и реченицата не може да се распознае.

Лево рекурзивни граматики

- Контекстно слободна граматика е $(V, \Sigma; S; P)$ лево рекурзивна ако има правило од облик $X\alpha \Rightarrow X\beta$
- X може бесконечно да се заменува

- Пример:

$expression: expression + term.$

$expression: expression - term.$

$expression: term.$

Лево рекурзивни граматики

- Левата рекурзија може да се отстрани
- Пример

expression: term + expression.

expression: term - expression.

expression: term.

- Операторите + и – сега се третираат како да се десно асоцијативни,
 - Кај одземање изразот погрешно да се испарсира
- Но сега не може да се избере дали со кој од овие правила да се направи замена, дали со +, дали со – или без знак. Може да се гледа неколку чекори наназад-но до каде?
- Може да искористиме лева факторизација
 - терм да се отстрани од почетокот и да се стави како посебно правило. Ова се нарекува “factoring out”. Се добива граматика со еден поглед на назад и е лево асоцијативна

expression:

term restexpression .

restexpression :

+ term restexpression.

restexpression :

- term restexpression.

restexpression : e.

Bottom-ур парсирање

- Bottom-ур парсерите работат обратно од top-down
- Почнуваат со целосниот внес на реченицата и работат со замена на терминалните со нетерминални користејќи ги правилата на граматиката во обратен правец.

Bottom-ур парсирање

expression: expression + expression.
expression: expression - expression.
expression: factor .
factor : 0.
factor : 1.
factor : 2.
factor : 3.
factor : 4.
factor : 5.
factor : 6.
factor : 7.
factor : 8.
factor : 9.

1 + 2 + 3 => factor + 2 + 3
=> expression + 2 + 3
=> expression + factor + 3
=> expression + expression + 3
=> expression + 3
=> expression + factor
=> expression + expression
=> expression

Множествата FIRST и FOLLOW

- Множеството FIRST на нетерминалот A е множество од сите терминали со кои може да почнува A.

$$\text{FIRST}(X) := \bigcup_{i=1}^n \text{PFIRST}(X_i)$$

- FIRST множеството на factor е {0; 1; 2; 3; 4; 5; 6; 7; 8; 9}.
- Множеството FOLLOW на нетерминалот A е множество од сите терминали кои може да следуваат директно после A.
 - Ова множество се бара само за нетерминалите кои се јавуваат од десната страна на правилата.

Множествата FIRST и FOLLOW

- Множеството FIRST на нетерминалот A е множество од сите терминали со кои може да почнува A.

$$\text{FIRST}(X) := \bigcup_{i=1}^n \text{PFIRST}(X_i)$$

- FIRST множеството на factor е {0; 1; 2; 3; 4; 5; 6; 7; 8; 9}.
- Множеството FOLLOW на нетерминалот A е множество од сите терминали кои може да следуваат директно после A.
 - Ова множество се бара само за нетерминалите кои се јавуваат од десната страна на правилата.

Множествата FIRST и FOLLOW

```
expression :      factor restexpression .
restexpression :  ε .
                | + factor restexpression
                | - factor restexpression .
factor :        0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

FOLLOW(expression) = {⊥}¹
FOLLOW(restexpression) = {⊥}
FOLLOW(factor) = {⊥, +, -}
```

Production	PFIRST
program:	statementlist RESULT =implication.
statementlist :	ϵ , $\{A \dots Z\}$ \emptyset
statement :	statement statementlist. $\{A \dots Z\}$
statement:	identifier =implication i. $\{A \dots Z\}$
implication :	conjunction restimplication. $\{\sim, (, 0, 1, A \dots Z\}$ \emptyset
restimplication :	ϵ , $\{-> \text{ conjunction restimplication.}\}$ $\{-> \text{ conjunction restimplication.}\}$
restimplication :	$\{-> \text{ conjunction restimplication.}\}$ $\{-> \text{ conjunction restimplication.}\}$
restimplication :	$\{-> \text{ conjunction restimplication.}\}$ $\{-> \text{ conjunction restimplication.}\}$
conjunction :	negation restconjunction. $\{\sim, (, 0, 1, A \dots Z\}$ \emptyset
restconjunction :	ϵ , $\{\& \text{ negation restconjunction.}\}$ $\{\& \text{ negation restconjunction.}\}$
restconjunction :	$\{\& \text{ negation restconjunction.}\}$ $\{\& \text{ negation restconjunction.}\}$
negation:	\sim negation. $\{\sim\}$
negation:	factor. $\{(, 0, 1, A \dots Z\}$ $\{\}\}$
factor :	(implication) . $\{\}\}$
factor :	identifier . $\{A \dots Z\}$ $\{1\}$
factor :	1. $\{1\}$
factor :	0. $\{0\}$
identifier :	A. $\{A\}$
identifier :	Z. $\{Z\}$

Table 6.1: PFIRST Sets for Logic Language

МНОЖЕСТВАТА FIRST и FOLLOW

Nonterminal	FIRST	FOLLOW
program	$\{A \dots Z\}$	$\{\perp\}$
statementlist	$\{A \dots Z\}$	$\{RESULT\}$
statement	$\{A \dots Z\}$	$\{A \dots Z, RESULT\}$
implication	$\{\sim, (, 0, 1, A \dots Z\}$	$\{;, \perp,)\}$
restimplication	$\{->, <-, <->\}$	$\{;, \perp,)\}$
conjunction	$\{\sim, (, 0, 1, A \dots Z\}$	$\{->, <-, <->, ;, \perp,)\}$
restconjunction	$\{\&, \}$	$\{->, <-, <->, ;, \perp,)\}$
negation	$\{\sim, (, 0, 1, A \dots Z\}$	$\{\&, , ->, <-, <->, ;, \perp,)\}$
factor:	$\{(, 0, 1, A \dots Z\}$	$\{\&, , ->, <-, <->, ;, \perp,)\}$
identifier:	$\{A \dots Z\}$	$\{=, \&, , ->, <-, <->, ;, \perp,)\}$

Table 6.2: FIRST and FOLLOW Sets for Logic Language