

Bottom-Up Syntax Analysis

staršie - metódy

Ján Šturc

Zima 2010

Posunovo redukčná schéma

Shift-Reduce Parsing

Grammar:

1. $S \rightarrow a A B e$
2. $A \rightarrow A b c$
3. $A \rightarrow b$
4. $B \rightarrow d$

Reduced sentence:

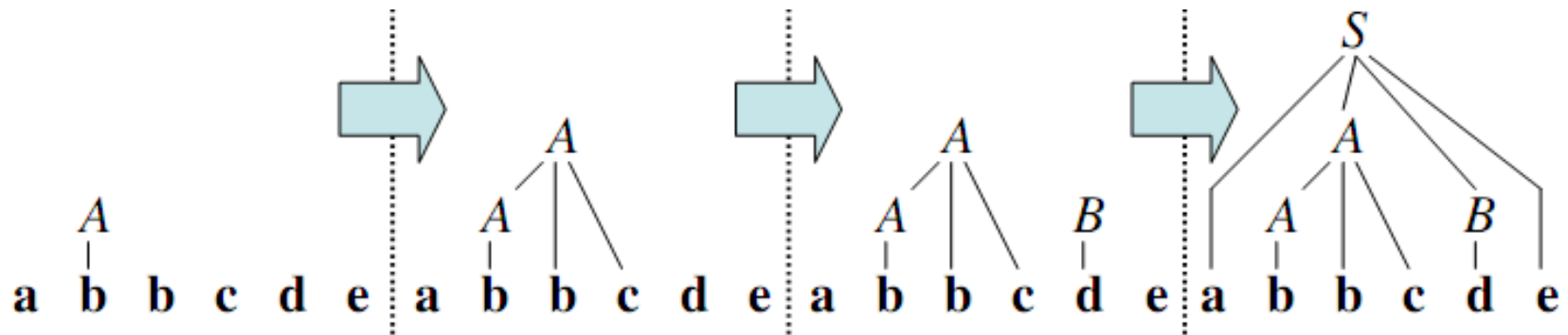
a b b c d e
a A b c d e
a A d e
a A B e
S

Shift-reduce corresponds to a rightmost derivation:

$S \Rightarrow_{rm} a A B e$
 $\Rightarrow_{rm} a A d e$
 $\Rightarrow_{rm} a A b c d e$
 $\Rightarrow_{rm} a b b c d e$

These production's match right-hand sides: 3 2 4 1

Shift-reduce parsing produces rightmost derivation in the reverse order.



Handles

A **handle** is a substring of grammar symbols in *a right-sentential form* that matches a *right-hand side of a production*

Grammar:

1. $S \rightarrow a A B e$

2. $A \rightarrow A b c$

3. $A \rightarrow b$

4. $B \rightarrow d$

Reducing a sentence:

$a \underline{b} b c d e$

$a \underline{A b c} d e$

$a A \underline{d} e$

$\underline{a A B e}$

S

Handles:

b

Abc

d

aABe

In the right-sentential form $a A \mathbf{b} c d e$ the string \mathbf{b} is not a handle because the result $a A \mathbf{A} d e$ is not a sentential form (cannot be derived from starting symbol S by the given grammar). Further reductions would fail.

Posunovo redukčný automat

- Je to zásobníkový automat, ktorý na základe vstupého symbolu a stavu, v ktorom sa nachádza vykoná jednu z akcií
 1. **Shift**: posunie symbol zo vstupu na zásobník a posunie vstup
 2. **Reduce**: redukuje symboly na vrchole zásobníku
 3. **Accept**: akceptuje
 4. Ak nemôže vykonať ani jednu z uvedených akcií, vykoná spracovanie chyby: **Error**.

Činnost posunovo redukčního automatu

Grammar:

- $E \rightarrow E + E$
- $E \rightarrow E * E$
- $E \rightarrow (E)$
- $E \rightarrow \text{id}$

Stack	Input	Action
\$	id+id*id\$	shift
\$ <u>id</u>	+id*id\$	reduce $E \rightarrow \text{id}$
\$ <u>E</u>	+id*id\$	shift
\$ <u>E+</u>	id*id\$	shift
\$ <u>E+id</u>	*id\$	reduce $E \rightarrow \text{id}$
\$ <u>E+E</u>	*id\$	shift (or reduce?)
\$ <u>E+E*</u>	id\$	shift
\$ <u>E+E*id</u>	\$	reduce $E \rightarrow \text{id}$
\$ <u>E+E*E</u>	\$	reduce $E \rightarrow E * E$
\$ <u>E+E</u>	\$	reduce $E \rightarrow E + E$
\$ <u>E</u>	\$	accept

Handles
to reduce

A conflict.
How to
resolve it?

Konflikty – Conflicts

- Shift-reduce and reduce-reduce conflicts are caused by
 - The limitations of the parsing method (even when the grammar is unambiguous)
 - Ambiguity of the grammar

Operátorovo precedenčné gramatiky

- Neobsahujú ϵ -ové pravidlá
- Žiadna pravá strana neobsahuje dva po sebe idúce neterminály.
- Ideálne striedajú sa terminálne a neterminálne symboly
- Neterminálne symboly nie sú dôležité. Jeden neterminál stačí.

Príklad: aritmetický výraz

$$E \rightarrow E + E \mid E - E \mid E * E \mid E / E \mid E \uparrow E \mid (E) \mid id$$

Precedenčné relácie

- Definujeme tri rôzne precedenčné relácie:
 - $a \succ b$ a „má väčšiu precedenciu ako“ b
 - $a \doteq b$ a „má rovnakú precedenciu ako“ b
 - $a \preceq b$ a „má menšiu precedenciu ako“ b
- Relácia medzi dvojicou terminálov sa určuje:
 - na základe matematickej intuície o prioritě a asociatívnosti operátorov
 - Zostrojením jednoznačnej gramatiky a použitím mechanickej metódy na odvodenie precedenčných relácií
- Precedenčné relácie slúžia aj na ohraničenie handle
 - \preceq začiatok handle
 - \succ koniec handle
 - Ak $a \doteq b$, potom sa symboly a a b vyskytujú v tej istej handle

Precedenčná posunovo redukčná schéma

- Konfigurácia posunovo redukčného automatu pre precedenčnú gramatiku je tvaru:
 $\$ \beta_0 a_1 \beta_1 \dots a_m \beta_m, a_{m+1} a_{m+2} \dots a_n \$$, kde β_i je ε alebo neterminál a a_i je terminál.
- Ak $a_m \succ a_{m+1}$ (posledný operátor na zásobníku má väčšiu prioritu ako operátor na vstupe) redukujeme.
- Pri redukcií vyberáme symboly zo zásobníku (pop) pokiaľ neplatí $a_k \prec a_{k+1}$ (operátor zostávajúci na zásobníku má menšiu prioritu ako posledný vybraný operátor).
- Neterminály neovplyvňujú parsovanie. Pre popis gramatiky stačí jeden a pri parsovaní ich nepotrebuje vôbec.

Algoritmus precedenčnej analýzy

set ip na začiatok vstupu w\$;

loop

if top=\$ **and** ip↑ = \$ **then** EXIT

else { a := top; b := ip ↑;

if a < b **or** a = b **then** push; advance(ip)

elseif a > b **then** /*reduce*/

{ **repeat** pop

until top < symbol last popped;

output(rule)

}

else error

}

pool

Intuitívna konštrukcia precedenčných relácií

1. O_1 má väčšiu prioritu ako O_2
 - $O_1 \succ O_2$
 - $O_2 \prec O_1$
 - Príklad: $+ \prec *, * \succ +$
2. O_1 a O_2 majú rovnakú prioritu (alebo O_1 a O_2 je ten istý operátor)
 - $O_1 \succ O_2, O_2 \succ O_1$ ak operátory sú ľavo asociatívne
 - $O_1 \prec O_2, O_2 \prec O_1$ ak operátory sú pravo asociatívne
 - Príklad: $+ \succ +, + \succ -, - \succ +, - \succ -, \uparrow \prec \uparrow$

Operátory a ostatné tokeny

- Nech O je ľubovoľný operátor. Potom
 - $O \prec \text{id}, \quad \text{id} \succ O, \quad \$ \prec O \quad O \succ \$,$
 - $O \prec (, \quad (\prec O$
 - $O \succ), \quad) \succ O$
- Relácie neobsahujúce operátory
 - $(\doteq) \quad \$ \prec (\quad \$ \prec \text{id}$
 - $(\prec (\quad \text{id} \succ \$ \quad) \succ \$$
 - $(\prec \text{id} \quad \text{id} \succ) \quad) \succ)$

Precedenčná tabuľka – príklad

Gramatika:

$E \rightarrow E + E \mid E - E \mid E * E \mid E / E \mid E \uparrow E \mid (E) \mid id$

	+	-	*	/	↑	id	()	\$
+	·>	·>	<·	<·	<·	<·	<·	·>	·>
-	·>	·>	<·	<·	<·	<·	<·	·>	·>
*	·>	·>	·>	·>	<·	<·	<·	·>	·>
/	·>	·>	·>	·>	<·	<·	<·	·>	·>
↑	·>	·>	·>	·>	<·	<·	<·	·>	·>
id	·>	·>	·>	·>	·>		·>	·>	·>
(<·	<·	<·	<·	<·	<·	<·	·	
)	·>	·>	·>	·>	·>		·>		·>
\$	<·	<·	<·	<·	<·	<·	<·		

Kompresia precedenčnej tabuľky

dvojica funkcií f, g .

$f(a) < g(b)$, ak $a \prec b$

$f(a) = g(b)$, ak $a \doteq b$

$f(a) > g(b)$, ak $a \succ b$

Príklad:

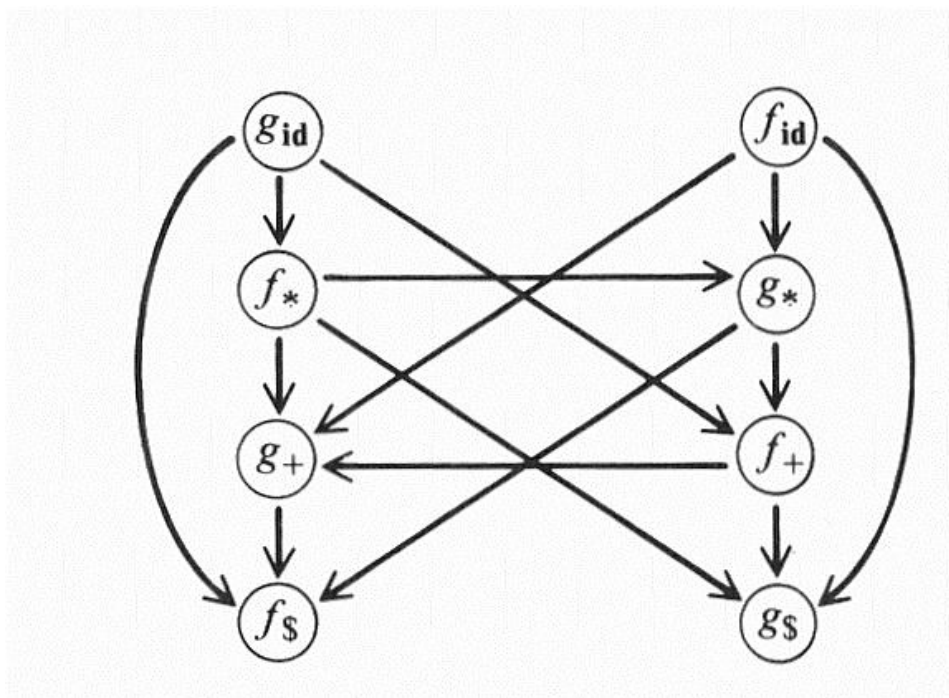
	+	-	*	/	↑	id	()	\$
f	2	2	4	4	4	0	6	6	0
g	1	1	3	3	5	5	0	5	0

Konštrukcia funkcií f a g

1. Vytvor symboly fa a ga pre všetky $a \in T \cup \{ \$ \}$
2. Zlučuj skupiny: Na počiatku je každý symbol jedna skupina.
 - Ak $a \doteq b$ zlúč skupinu fa a gb .
3. Zo skupín urob uzly orientovaného grafu.
4. Pridaj hrany
 - Ak $a \prec b \Rightarrow$ hrana $(gb) \rightarrow (fa)$
 - Ak $a \succ b \Rightarrow$ hrana $(fa) \rightarrow (gb)$.
5. Ak graf je cyklický precedenčné funkcie neexistujú.
6. Ak graf je acyklický
 - $f(a)$ } = dĺžka najdlhšej cesty vychádzajúcej z fa resp. ga
 - $g(a)$ }

Príklad

	+	*	id	\$
+	▷	◁	◁	▷
*	▷	▷	◁	▷
id	▷	▷	▷	▷
\$	◁	◁	◁	
f	2	4	4	0
g	1	3	5	0



Syntaktická analýza zdola nahor a vyhľadávanie reťazcov

- Bottom-Up parsing hľadá výskyt pravej strany pravidla vo vetnej forme (viable prefix).
- Idea je použiť štandardný algoritmus na paralelné vyhľadávanie vzoriek (Aho-Corasick, Dömölki)
- V príslušných automatoch stačí doplniť prácu so zásobníkom
- Keďže sa nechceme obmedziť na bezkontextové jazyky zavedieme zásobníky dva
 - vstupný (na vracanie už čiastočne spracovaného textu na vstup)
 - pracovný (normálny zásobník posunovo redukčného automatu)

Príklad: Aho Corasickovej automat pre gramatiku aritmetického výrazu

$S \rightarrow E\$$

$E \rightarrow E \text{ addop } T \mid T$

$T \rightarrow T \text{ mulop } F \mid F$

$F \rightarrow (E) \mid -(E) \mid \text{id} \mid -\text{id}$

Follow množiny:

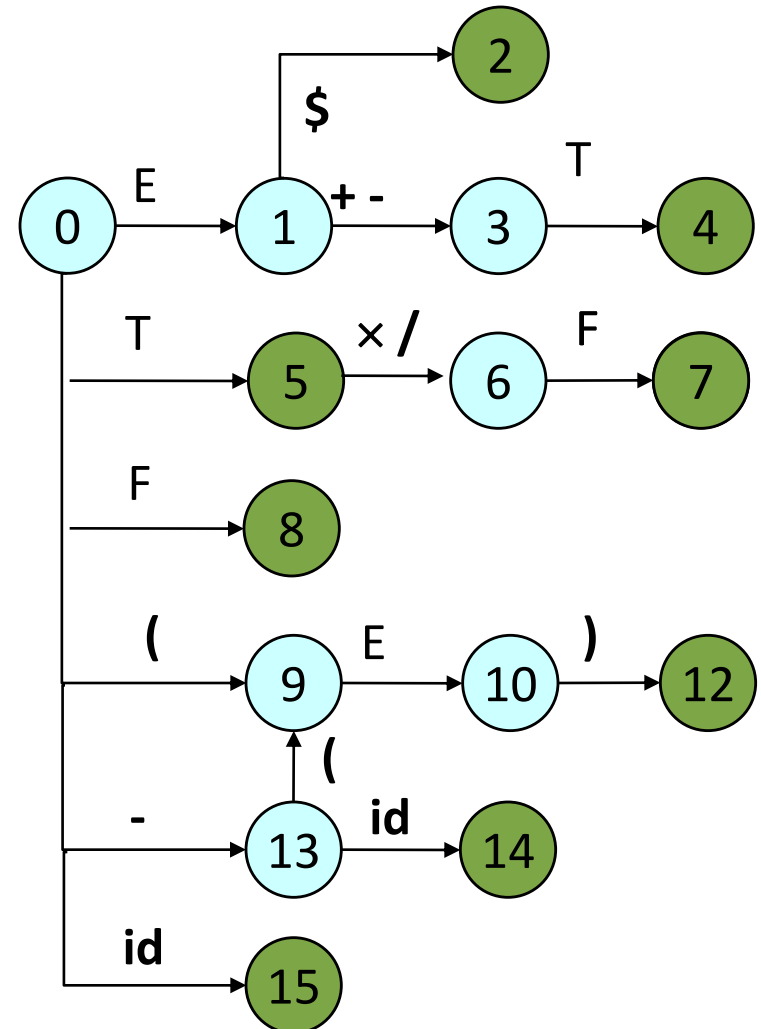
(E) $\{+, -,), \$\}$

(T) $\{+, -,), \$, \times, /\}$

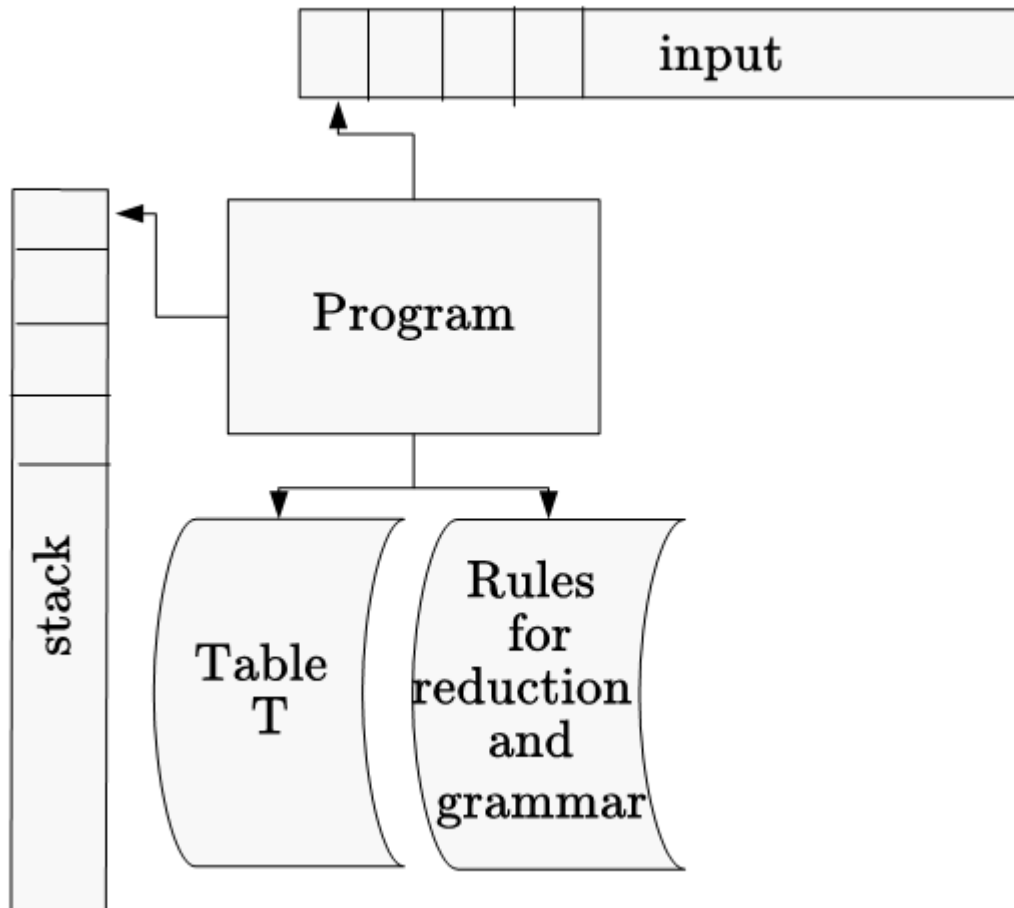
(F) $\{+, -,), \$, \times, /\}$

$\text{First}(X) = \{-, (, \text{id}\}$,

pre $X \in \{S, E, T, F\}$.



Dvozásobníkový automat



Program:

$q := 0;$

$\text{push}(\text{stack}, q);$

repeat

if $q \in F$ **then** Reduce;

else { $c := \text{pop}(\text{input});$

$\text{push}(\text{stack}, c);$

$q := T[q, c];$

$\text{push}(\text{stack}, q)$

}

until empty(input) ;

if empty(stack) **then** accept

else error;

Analýza kontext senzitivných jazykov

- Pravidlá sú tvaru: (1) $\alpha A \beta \rightarrow \alpha \omega \beta$.
- Algoritmus AC hľadá vzorku pravú stranu ($\alpha \omega \beta$ je to handle).
- Rovnako ako pri LR-like vkladáme do zásobníku symbol a stav.
- Shift:
 $s := \text{top}(\text{input}); \text{pop}(\text{input});$
 $q := T[q, s];$
 $\text{push}(\text{stack}, s); \text{push}(\text{stack}, q);$
- Reduce (podľa pravidla 1):
 for $i := |\beta| - 1$ **downto** 0 **do** { $\text{pop}(\text{stack}); s := \text{top}(\text{stack}); \text{pop}(\text{stack});$
 $\text{push}(\text{input}, s)$ } /* vráti na vstup β */
 for $i := |\omega| - 1$ **downto** 0 **do** { $\text{pop}(\text{stack}); \text{pop}(\text{stack})$ }
 /* vypopuje z pracovného zásobníku ω */
 $q := \text{top}(\text{stack}); \text{push}(\text{input } A);$ /* Dá na vstup A */

Konfigurácie

- Konfigurácie sú podobné konfiguráciám posunovo redukčného automatu uvidíme ich pre prípad redukcie:

- Pred redukciou:

$$s_1 q_1 \dots s_i q_i a_1 q_{i+1} \dots a_k q_{i+k} w_1 q_{i+k+1} \dots w_m q_{i+k+m} b_1 q_{i+k+m+1} \dots b_n q_{i+k+m+1}$$

$$\uparrow c_j \dots c_{j+z}$$

- Po redukcii

$$s_1 q_1 \dots s_i q_i a_1 q_{i+1} \dots a_k q_{i+k} \uparrow a b_1 \dots b_n c_j \dots c_{j+z}$$

- Pre lineárne ohraničené jazyky vstupný zásobník nikdy nebude väčší ako pôvodný vstup

Špeciálny prípad $\beta = \varepsilon$

- Pravidlá tvaru (2) $\alpha A \rightarrow \alpha \omega$.
- Posun funguje rovnako ako v predošlom prípade nič sa nemení.
- Pri redukcii je prvý cyklus prázdny zostane:
for $i := |\omega| - 1$ **downto** **0** **do** { $\text{pop}(\text{stack}); \text{pop}(\text{stack})$ }
/ * vypopuje z pracovného zásobníku ω */
 $q := \text{top}(\text{stack}); \text{push}(\text{input } A);$ / * Dá na vstup A */
- Na vstup sa v každom prípade vracia jediný neterminál. Tento sa v nasledujúcom kroku presunie na zásobník.
- Ak si stav po presune predvypočítame (goto tabuľka) môže vstup byť jednosmerná páska ako pri xLR analýze.

Príklad – kontext senzitívna gramatika

$S \rightarrow E\$$
 $E \rightarrow E \text{ addop } E$
 $F \rightarrow F \text{ mulop } F$
 $E.2 \rightarrow F \{ \text{addop} \mid) \mid \$ \}$
 $F \rightarrow (E) \mid \text{id} \mid -\text{id}$

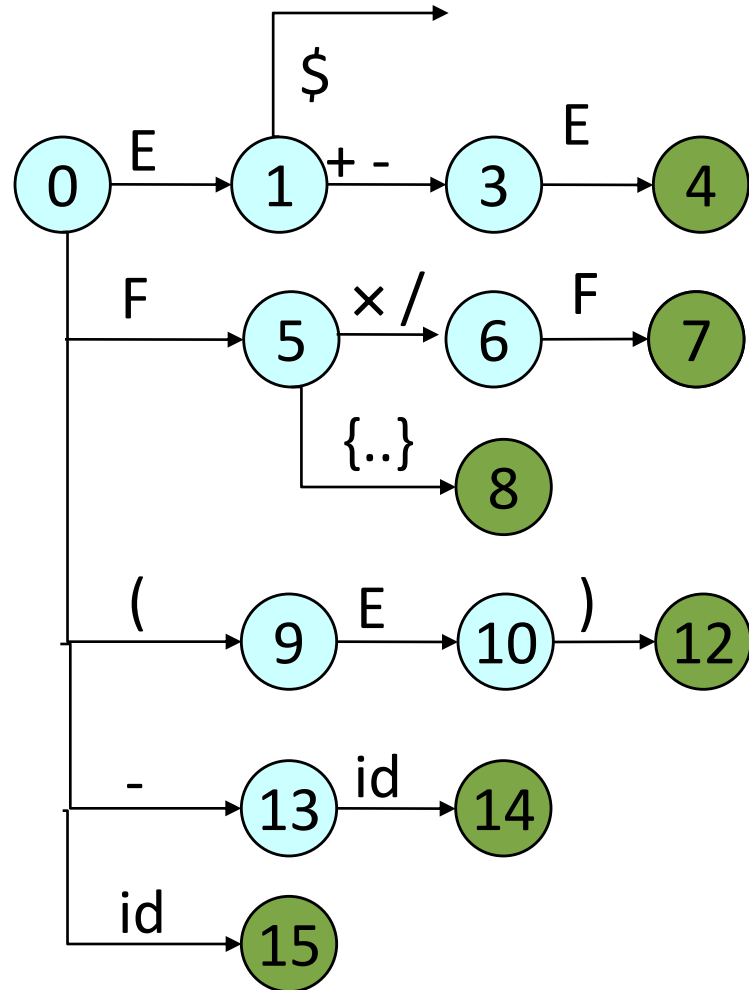
Skratky:

addop = {+|-}; mulop = {×|/}

Pravidlo $E.2 \rightarrow F \{ \text{addop} \mid) \mid \$ \}$ je
skratka za množinu pravidiel:

$\{ E + \rightarrow F +, E - \rightarrow F -, E) \rightarrow F), E \$ \rightarrow F \$ \}$

Skratky okrem toho, že skracujú zápis gramatiky sú aj efektívne implementované, znižujú aj rozpoznávací automat.



Chyby: diagnostika a recovery

- Pri hľadaní vzorky sa nemôže vyskytnúť chyba. Vzorka sa buď najde, alebo nenajde.
- Ak v trie automatu neexistuje, pokračovanie prejde sa do stavu, ktorý zodpovedá rozpoznaní najdlhšieho suffixu, t.j. v najhoršom prípade do počiatočného stavu.
- Pre syntaktickú analýzu môžeme „nie trie“ prechody doplniť citlivejšie, napr. podľa množín FIRST a Follow.
- Nie trie hrany zo stavu q_i do stavu q_j zodpovedajú korektnému odvodeniu: Ak zo stavu q_i vedie hrana označená neterminálnym symbolom X a cesta zo stavu 0 do stavu q_j je označená reťazcom r takým, že $X \Rightarrow r\alpha$.
- Ostatné nie trie hrany sú chybové.
- Chybovým hranám priradíme prechody do chybových stavov (stavov označených zápornými číslami).

Príklad: doplnenie tabuľky prechodov

$S \rightarrow E\$$
 $E \rightarrow E \text{ addop } E$
 $F \rightarrow F \text{ mulop } F$
 $E.2 \rightarrow F \{ \text{addop} \mid \} \mid \$ \}$
 $F \rightarrow (E) \mid \text{id} \mid -\text{id}$

$\text{First}(X) = \{-, (, \text{id}\},$
 pre $X \in \{S, E, F\}$.

Ostatné nechybové prechody:

(id	-	F	
3	9	15	13	5
6	9	15	13	
9	9	15	13	5
13	9	14	?	

Môžeme aj všetky poslať do stavu 0, bez prečítania vstupu.

Pridali sme pravidlo:

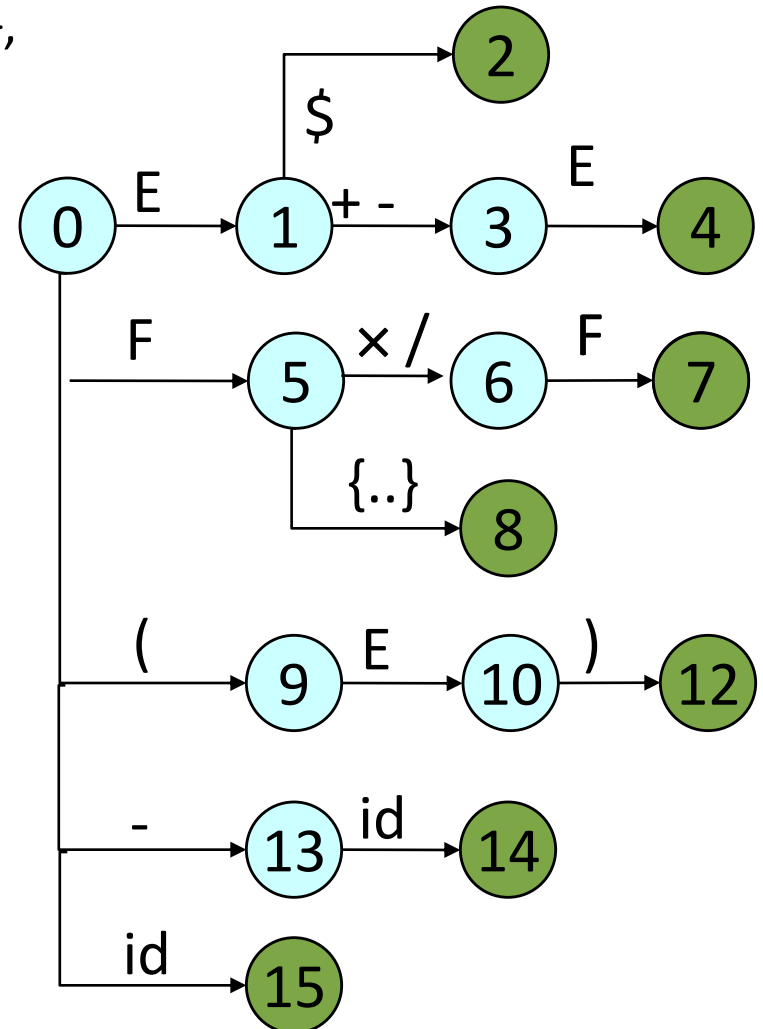
$F \rightarrow -(E)$

Chybové prechody:

1 nemali by vzniknúť

3 {+, mulop} operátor navyše

5 {id, (} chýba operátor ?



Ostatné chybové prechody

6	{+, mulop}	dva operátory za sebou
9	{+, mulop}	operátor za zátvorkou
10	{všetko okrem) }	chýba pravá zátvorka

- Poskytuje aspoň také možnosti spracovania chýb, ako klasické LR automaty.
- Chybové stavy
- Rozpoznanie vzorky sa automaticky zachytí na začiatkoch pravidiel (Panic mode, vstup sa presúva do zásobníku)
- Relikt (zombie) v zásobníku je vhodnejší pre globálnu opravu a hľadanie minimálnej chyby
- Daň za to je, že sa nemusíme dozvedieť o chybe včas, keď analyzovaný vstup prestane byť prefixom nejakého slova jazyka.

Konkrétny príklad

- aritmetický výraz: id (id+id)\$
LR: E ↑(id+id)\$
SM: FE ↑\$
- Čo je lepšie vo všeobecnosti, mi nie je jasné;
 - Prístup LR (kompilátoru): Zistiť chybu najskôr ako sa dá.
 - Prístup SM (hladača vzoriek): Nájsť čo najviac, tak dlhých ako sa dá, odvoditeľných vzoriek.
- Nakoniec to môže dopadnúť rovnako.
 - LR metóda to urobí na dva krát:
 - Panic mode sa zachytí na (a rozpozná sa výraz.
 - Obom prístupom zostane rovnaký relikv v zásobníku a urobia phrase level recovery (missing operand)
 - Pri pozornom ošetrovaní chýb LR metóda nemusí redukovať (E)

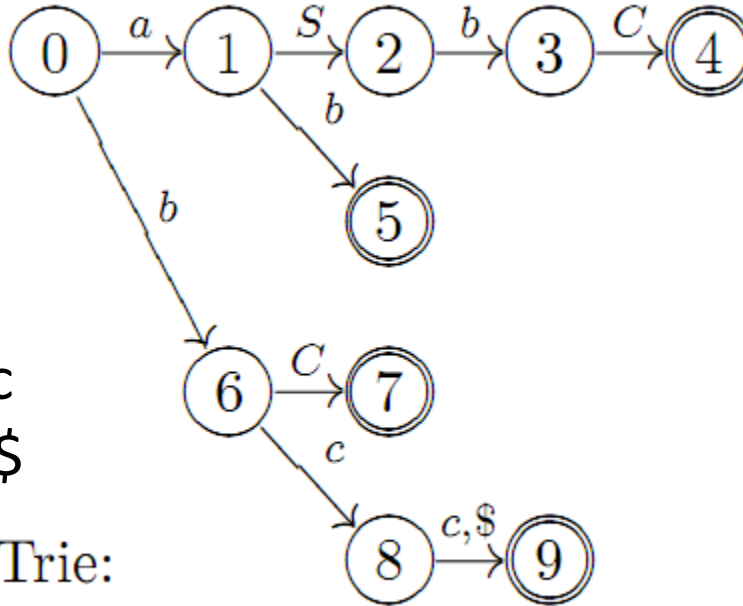
Rozšírenie na gramatiky bez obmedenia – prepisovacie systémy

- Tvar previdiel nie je dôležitý. Hľadáme pravé strany a nahradzujeme ľavými.
- Bez ohraničení na tvar pravidiel sa môže stať prepisy medzi vstupom a zásobníkom sa budú mnohokrát opakovať a analýza nebude lineárna.
- Je to prepisovací systém (Postov, Markovov, ... model algoritmu). Dá sa naprogramovať všetko. Už to nie je o kompilátoroch, ale o programovaní.
- Paradigma syntaxou (pravidlami) riadeného programovania
 - Módne a užívateľský príťažlivé, pokiaľ nevznikli grafické prostredia GUI.
 - Vychádzala z lineárneho modelu vstup a výstupu.

Príklad: $\{a^n b^n c^n : n \geq 1\}$

Grammar:

1. $S \rightarrow aSbC$
2. $aSb \rightarrow ab$
3. $Cb \rightarrow bC$
4. $bCc \rightarrow bcc$
5. $bC\$ \rightarrow bc\$$



	<i>a</i>	<i>b</i>	<i>c</i>	<i>\$</i>	<i>S</i>	<i>C</i>	<i>F</i>	<i>pattern</i>
0	1	6					0	
1	1	5			2		0	
2	1	3					0	
3	1	6				4	0	
4							1	<i>aSbC</i>
5							1	<i>ab</i>
6	1	6				7	0	
7							1	<i>bC</i>
8	1	6	9	9			0	
9							1	<i>bcc, bc\$</i>

Kritické je pravidlo 3.

Pre elimináciu jednej trojice abc musí prejsť cez všetky zvyšné b . Spôsobuje, že syntaktická analýza je kvadratickej zložitosti.

Skúste chybný vstup: $aaababcbccc$

Tu to zistí: $aaaSbS \uparrow Cbcc\$$;

Tu prestáva byť viable prefix: $aaaSb \uparrow abcbccc\$$