



VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ

BRNO UNIVERSITY OF TECHNOLOGY

FAKULTA INFORMAČNÍCH TECHNOLOGIÍ

FACULTY OF INFORMATION TECHNOLOGY

ÚSTAV INFORMAČNÍCH SYSTÉMŮ

DEPARTMENT OF INFORMATION SYSTEMS

**GRAMATICKÉ SYSTÉMY A SYNTAXÍ ŘÍZENÝ
PŘEKLAD ZALOŽENÝ NA NICH**

GRAMMAR SYSTEMS AND SYNTAX-DIRECTED TRANSLATION BASED ON THEM

DIPLOMOVÁ PRÁCE

MASTER'S THESIS

AUTOR PRÁCE

AUTHOR

Bc. JAROSLAV HANDLÍŘ

VEDOUCÍ PRÁCE

SUPERVISOR

Prof. RNDr. ALEXANDER MEDUNA, CSc.

BRNO 2017

Abstrakt

Práce zkoumá teorii formálních jazyků v oblasti bezkontextových gramatik. Zaměřuje se především na možnosti a modely spolupráce více gramatik při řešení společného problému. V těchto souvislostech představuje gramatické systémy, které byly navrženy jako formální prostředek pro popis distribuovaného a paralelního zpracování. Po uvedení do dané problematiky se práce zaměřuje na praktické uplatnění těchto mechanismů při syntaxí řízeném překladu, a proto je druhá část práce věnována implementaci dynamického syntaktického analyzátoru, který během analýzy aplikuje více gramatik. S ohledem na co největší uživatelskou přívětivost a možné didaktické použití je aplikace implementována pomocí moderních webových technologií HTML5, JavaScript, AngularJS, CSS3, LESS a další.

Abstract

The thesis examines the theory of formal languages in the field of context-free grammars. It focuses mainly on the possibilities and models of collaborating grammars to solve a common problem. In this context, it presents grammatical systems that have been designed as a formal means for describing distributed and parallel processing. After introducing to the problematics, the thesis focuses on the practical use of these mechanisms in the translation controlled syntax, and therefore the second part of the thesis deals with the implementation of a dynamic syntactic analyzer that uses more grammars during the analysis. With respect to the greatest user friendliness and the possible didactic use, the application is implemented using modern web technologies HTML5, JavaScript, AngularJS, CSS3, LESS and more.

Klíčová slova

bezkontextové gramatiky, gramatické systémy, CD, PC, syntaktická analýza, precedenční analýza, syntaxí řízený překlad, dynamická gramatika, HTML5, JavaScript, AngularJS, CSS3, LESS

Keywords

context-free grammars, grammar systems, CD, PC, syntax analysis, precedence analysis, syntax-directed translation, dynamic grammar, HTML5, JavaScript, AngularJS, CSS3, LESS

Citace

HANDLÍŘ, Jaroslav. *Gramatické systémy a syntaxí řízený překlad založený na nich*. Brno, 2017. Diplomová práce. Vysoké učení technické v Brně, Fakulta informačních technologií. Vedoucí práce prof. RNDr. Alexander Meduna, CSc.

Gramatické systémy a syntaxí řízený překlad založený na nich

Prohlášení

Prohlašuji, že jsem tuto diplomovou práci vypracoval samostatně pod vedením prof. RNDr. Alexandera Meduny, CSc. Uvedl jsem všechny literární prameny a publikace, ze kterých jsem čerpal.

.....
Jaroslav Handlř
15. května 2017

Poděkování

Tímto bych rád poděkoval svému vedoucímu diplomové práce prof. RNDr. Alexanderu Medunovi, CSc., za jeho čas a zájem, odborný dohled a vedení v průběhu celé práce.

Obsah

1	Úvod.....	2
2	Základní formální prostředky.....	3
3	Regulární gramatiky.....	4
3.1	Konečné automaty	5
4	Bezkontextové gramatiky	6
4.1	Derivační stromy a nejednoznačnost	7
4.2	Zásobníkové automaty.....	9
5	Gramatické systémy.....	10
5.1	CD gramatické systémy	11
5.1.1	Formální definice.....	11
5.1.2	Derivační módy	12
5.1.3	Generovaný jazyk	13
5.1.4	Příklady.....	13
5.1.5	Třídy jazyků CD gramatických systémů a jejich generativní síla	15
5.1.6	Hybridní CD gramatické systémy.....	16
5.1.7	Třídy jazyků hybridních CD gramatických systémů a jejich generativní síla	17
5.2	PC gramatické systémy.....	18
5.2.1	Formální definice.....	18
5.2.2	Derivační kroky	18
5.2.3	Generovaný jazyk	20
5.2.4	Varianty PC gramatických systémů.....	20
5.2.5	Příklad.....	21
5.2.6	Třídy jazyků PC gramatických systémů a jejich generativní síla	22
6	Syntaxí řízený překlad	23
6.1	LL gramatiky bez epsilon pravidel	26
6.1.1	Množina First.....	27
6.2	LL gramatiky s epsilon pravidly	29
6.2.1	Množina Empty	30
6.2.2	Množina First.....	32
6.2.3	Algoritmus $\text{First}(X_1X_2 \dots X_n)$	34
6.2.4	Algoritmus $\text{Empty}(X_1X_2 \dots X_n)$	34
6.2.5	Množina Follow.....	35
6.2.6	Množina Predict.....	36
6.3	Rekurzivní sestup.....	38
6.4	Prediktivní syntaktická analýza	39
6.5	LR syntaktický analyzátor	40
6.6	Precedenční syntaktický analyzátor	41
7	Aplikace	43
7.1	Uživatelské rozhraní	47
7.2	Implementace.....	48
7.3	Nasazení a vývoj.....	49
8	Závěr	50

1 Úvod

Práce se věnuje teorii formálních jazyků v oblasti bezkontextových gramatik, které jsou velmi významné především při analýze a překladu programovacích jazyků. Umožňují tak nejen formálně popsat a specifikovat daný programovací jazyk, ale zároveň provádět jeho syntaktickou analýzu a především pak transformaci většinou velmi abstraktního jazyka, který je však pro člověka dobře čitelný, do jazyka symbolických adres nebo přímo strojového kódu, aby mohl být daný program vykonávaný počítačem. Přestože lze popsat programovací jazyk jednou komplexní gramatikou, může být v mnoha ohledech užitečné zapojit do činnosti více gramatik, z nichž každá dokáže zpracovat např. pouze část jazyka nebo dokonce může být vstupní program tvořený více různými jazyky. Jedná se o spolupráci více gramatik při řešení společného problému. V běžném světě se tato situace vyskytuje denně. Tým lidí se podílí na jednom projektu, který se snaží společně a co nejefektivněji dokončit. V praxi se tyto přístupy často uplatňují i v oblasti překladu jazyků, ale zpravidla se klade důraz na funkčnost a dále se již nezkoumá teoretické pozadí a nehledají se formalismy. Avšak formální prostředky, které se snaží tuto problematiku zachytit a rozvinout, již dlouhou dobu existují. Řeč je o tzv. gramatických systémech, kterým je věnováno jádro této práce.

Gramatické systémy vznikly ve druhé polovině 20. století a definují různé modely spolupracujících gramatik, které v době jejich vzniku nemohli autoři z dnešního pohledu plně docenit. Tyto systémy zavádějí dva hlavní typy, které umožňují popis distribuovaného a paralelního zpracování. První typ se nazývá kooperující distribuované gramatické systémy, zkráceně CD gramatické systémy (z anglického *Cooperating Distributed*), a druhý typ se jmenuje paralelně komunikující gramatické systémy, zkráceně PC gramatické systémy (z anglického *Parallel Communicating*).

CD gramatické systémy jsou distribuovanou variantou a jednotlivé gramatiky systému (nazývány komponenty systému) pracují sekvenčně. Postupně se střídají při zpracování větné formy, kterou společně sdílí, dokud nenaleznou její řešení. Výsledek je tak dán dílčím přispěním všech gramatik, které mezi sebou přepínají podle různých pravidel. Obecnou analogii k tomuto systému je možné v literatuře nalézt např. u tzv. „problému tabule“ (anglicky *The Black Board Model*), kdy se více lidí střídá u tabule. Společně tak řeší daný problém a každý z účastníků přispěje svými nápady. Za tímto systémem se zároveň skrývá i velký potenciál distribuce výkonu, který v době vzniku těchto systémů nebyl tolik využitelný. Naopak dnes, kdy je internet samozřejmostí, většina služeb se přesunuje na cloud a data i výkon lze snadno distribuovat díky moderním sítím po celém světě, je myšlenka CD gramatického systému velmi aktuální.

PC gramatické systémy pracují podobným způsobem, ale na rozdíl od předchozí sekvenční distribuované varianty jsou plně paralelní. Jednotlivé komponenty systému (gramatiky) pracují paralelně a každá komponenta zpracovává svoji vlastní větnou formu. Mezi sebou mohou komunikovat a předávat si informace pomocí tzv. komunikačních symbolů. Tyto systémy tak dávají prostor nejen pro řízení výpočtu (může existovat hlavní komponenta, která si žádá práci ostatních, jedná se potom o centralizovaný režim), ale především tak vznikl formální prostředek pro zavedení paralelismu. Paralelní zpracování je z dnešního pohledu druhým velmi významným trendem, protože dnes neexistuje obyčejný chytrý telefon, který by nedisponoval vícejádrovým procesorem. Otevírá se tím další velká oblast využití těchto systémů.

První polovina práce pojednává o obecných základech formálních jazyků a bezkontextových gramatik, které jsou aplikovány v gramatických systémech. Poté již následuje podrobné

vysvětlení obou typů gramatických systémů doplněné jejich obecnými definicemi, variantami a možnostmi spolu s praktickými příklady ukazujícími způsob jejich použití.

Druhá část práce se věnuje využití gramatických systémů při syntaktické analýze. Jednotlivé kapitoly popisují různé možnosti přístupu k syntaktické analýze a syntaxi řízenému překladu, přičemž podrobněji jsou popsány ty metody, které byly použity v související praktické práci.

Praktická část práce kombinuje přístupy syntaktické analýzy shora dolů s analýzou zdola nahoru pomocí centralizovaného PC gramatického systému do moderní webové aplikace, která uživateli umožňuje nejen velmi jednoduchým způsobem provést syntaktickou analýzu svého vlastního vstupního kódu v prostředí internetového prohlížeče bez nutnosti cokoli instalovat a konfigurovat, ale zároveň dává možnost si přímo v prohlížeči definovat velmi snadno vlastní gramatiku, která bude následně použita při analýze. Výstup syntaktické analýzy informuje uživatele o tom, jestli analýza proběhla v pořádku, příp. kde se pravděpodobně ve vstupním textu nachází chyba. Dále se uživatel dozví posloupnost aplikovaných pravidel gramatiky spolu s odpovídající LL tabulkou. Celá aplikace má tedy především didaktický charakter, který by měl uživateli pomoci k lepšímu pochopení bezkontextových gramatik, syntaktické analýzy a gramatických systémů. Zároveň aplikace představuje použití moderních vývojářských přístupů a webových technologií v praxi. Pro implementaci byly použity HTML5, kaskádové styly CSS3 a LESS, klientský JavaScript s frameworkem AngularJS, pro responzivitu a další styly framework Bootstrap. Tímto však výčet použitých technologií a nástrojů nekončí a jejich podrobný popis lze nalézt v kapitole 7. Za zmínku také stojí sekce příloh, která obsahuje ukázky reálné aplikace a jejího použití.

2 Základní formální prostředky

Tato kapitola seznamuje čtenáře se základními prostředky formálních jazyků, které jsou použity v celé následující práci, jež si klade za cíl nejen prohloubení znalostí, ale i srozumitelné vysvětlení základních principů pro čtenáře, kteří zatím spíše tápou. Zkušenější čtenáři mohou tedy této kapitole věnovat menší pozornost.

Abeceda je konečná, neprázdná množina elementů, které se nazývají symboly. Značí se Σ ¹. Abecedu, která obsahuje tři symboly x, y, z , lze zapsat následovně:

$$\Sigma = \{a, b, c\}.$$

Řetězec je potom posloupnost symbolů dané abecedy, která může mít nulovou délku. Takový řetězec se značí ϵ ² a platí, že ϵ je řetězec nad abecedou Σ . ϵ je tedy řetězec nulové délky, tzn. prázdný řetězec. Jinak platí, že pokud je x řetězec nad Σ a $a \in \Sigma$, potom i xa je řetězec nad abecedou Σ .

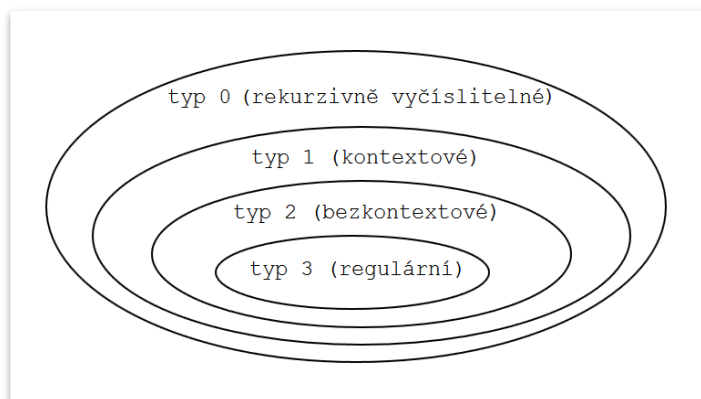
Jazykem se nazývá každá podmnožina řetězců nad danou abecedou, tedy necht' Σ^* značí množinu všech řetězců nad abecedou Σ , potom každá podmnožina $L \subseteq \Sigma^*$ je jazyk nad abecedou Σ .

Chomského hierarchie definuje kaskádu jednotlivých rodin jazyků, kde nadřazený typ má větší vyjadřovací sílu, a rozšiřuje tak typ pod ním. Bezkontextové gramatiky jsou rodinou jazyků typu 2 a rozšiřují tak nejjednodušší typ jazyků, kterým jsou regulární jazyky. Regulární jazyky jsou samotným základem a ve většině publikací jimi začíná výklad. Popisují se na nich nejzákladnější principy

¹ řecké velké písmeno sigma

² řecké malé písmeno epsilon

operací nad formálními jazyky. Tato práce je sice věnována bezkontextovým gramatikám, ale přesto následující kapitola obsáhne stručný úvod do regulárních jazyků a gramatik. Regulární jazyky a především konečné automaty jsou jedním z jejich formálních prostředků a budou zapotřebí v praktické části. Podrobné informace lze nalézt např. v publikaci [2].



Obrázek 1: Chomského hierarchie jazyků

3 Regulární gramatiky

Úvodem je třeba ještě vysvětlit, co je to *terminál* a *neterminál*. Za *terminál* či terminální symbol se označuje znak nebo řetězec z dané abecedy, zatímco *neterminál* představuje ostatní znaky, řetězce a symboly použité při definování pravidel gramatiky, čemuž je blíže věnována následující kapitola.

Regulární jazyky jsou nejjednodušší třídou jazyků v Chomského hierarchii a označují se typem 3. Obsahují všechny konečné jazyky, tj. jazyky, jejichž řetězce mají konečnou délku, a jednodušší případy nekonečných jazyků. Jednoduchým příkladem může být jazyk nad abecedou $\Sigma = \{a, b\}$, který obsahuje nejdříve sekvenci symbolů *a* a poté sekvenci symbolů *b*. Formální zápis takového jazyka by vypadal následovně:

$$L = \{a^n b^m : n \geq 0, m \geq 0\}$$

Jedním z formalismů regulárních jazyků jsou regulární výrazy. Jedná se o sadu pravidel a operací, které umožňují popsat regulární jazyky. Přestože regulární výrazy v rámci této práce přímo použity nejsou, jejich principy se prolínají i do bezkontextových gramatik, a je proto vhodné si je alespoň ve stručnosti definovat.

Rozlišují se tři základní operace nad řetězci, potažmo jazyky, a to

- *konkatenace*, značí se operátorem \cdot a provádí zřetězení dvou řetězců,
- *sjednocení*, značí se $+$ a umožňuje výběr,
- *iterace*, která se značí $*$,

a jejich vzájemná priorita je $* > \cdot > +$.

Regulární výrazy nad abecedou Σ a jazyky, které značí, lze definovat následovně:

- \emptyset je regulární výraz značící prázdnou množinu (prázdný jazyk)

- ϵ je regulární výraz značící jazyk $\{\epsilon\}$
- a , kde $a \in \Sigma$, je regulární výraz značící jazyk $\{a\}$
- Necht' r a s jsou regulární výrazy značící jazyky L_r a L_s , potom:
 - $(r \cdot s)$ je regulární výraz značící jazyk $L = L_r L_s$
 - $(r + s)$ je regulární výraz značící jazyk $L = L_r \cup L_s$
 - (r^*) je regulární výraz značící jazyk $L = L_r^*$

Regulární jazyk, jehož příklad je uveden na začátku této kapitoly, lze definovat následovně:

Necht' L je jazyk, potom je L regulární jazyk, pokud existuje regulární výraz r , který tento jazyk značí.

Neformálně řečeno, jazyk L je jazykem regulárním, pokud je možno pro tento jazyk najít regulární výraz, který jej popisuje. Dalšími možnostmi, jak specifikovat regulární jazyky, jsou regulární gramatiky nebo konečné automaty, jejichž definice zakončí tuto doplňující kapitolu.

Regulární gramatika G je čtveřice $G = (N, T, P, S)$, kde:

- N je abeceda neterminálů
- T je abeceda terminálů a platí, že $N \cap T = \emptyset$
- P je konečná množina pravidel tvaru:
 - $A \rightarrow xB$, kde $A, B \in N, x \in T$ nebo
 - $A \rightarrow x$, kde $A \in N, x \in T$ případně
 - $S \rightarrow \epsilon$, pokud se S neobjevuje na pravé straně žádného pravidla
- $S \in N$ je počáteční (startující) neterminál

3.1 Konečné automaty

Konečný automat je velmi užitečný nástroj, který umožňuje nejen přímou implementaci lexikální analýzy, o které bude zmínka později, ale využívá se v celé řadě aplikací praktických problémů, kde se systém nachází v určitém stavu a za daných podmínek přechází do stavu následujícího. Mimo to se velmi často používá grafická reprezentace pro názorné definování daného regulárního jazyka nebo obecně modelovaného problému.

Konečný automat je pětice $M = (Q, \Sigma, R, s, F)$, kde:

- Q je konečná neprázdná množina stavů
- Σ je vstupní abeceda
- R je konečná množina pravidel tvaru $pa \rightarrow q$, kde $p, q \in Q, a \in \Sigma \cup \{\epsilon\}$
- $s \in Q$ je počáteční stav
- $F \subseteq Q$ je množina koncových stavů

Konečný automat pracuje tak, že má k dispozici sadu stavů, ve kterých se může nacházet, a čtecí hlavu, pomocí níž čte symboly vstupní pásky po jednom zleva doprava. Na začátku je inicializován na počáteční stav s a čtecí hlava se nachází na první pozici pásky, tedy prvním symbolu. Poté postupně čte symboly z pásky a pro každý přečtený symbol a aktuální stav vyhledá v množině

pravidel, do jakého následujícího stavu může přejít přečtením tohoto symbolu. Pokud přečetl všechny symboly a dostal se do jednoho z koncových stavů f , pak může daný jazyk přijmout a tento jazyk je regulární.

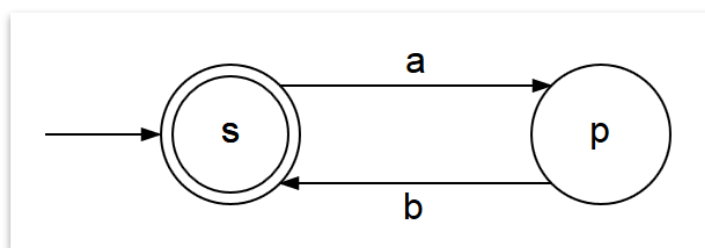
Formální definice jazyka přijímaného konečným automatem je následující. Necht' $M = (Q, \Sigma, R, s, F)$ je konečný automat. Jazyk přijímaný konečným automatem M , $L(M)$, je definován:

$$L(M) = \{w: w \in \Sigma^*, sw \vdash^* f, f \in F\}$$

Příklad 1.: Mějme konečný automat $M = (\{s, q\}, \{a, b\}, R, s, \{s\})$, kde R obsahuje pravidla:

- $sa \rightarrow q$
- $qb \rightarrow s$

Jazyk přijímaný tímto automatem je $L(M) = \{(ab)^n : n \geq 0\}$ a grafická reprezentace vypadá následovně:



Obrázek 2: Konečný automat reprezentující jazyk $(ab)^*$

4 Bezkontextové gramatiky

Následující kapitoly vychází z poznatků předchozích sekcí a postupně staví na tomto základu, takže nebude již vše vysvětlováno tak podrobně, jako tomu bylo v kapitole 2 a 3.

Bezkontextové jazyky, jak již název napovídá, nejsou omezovány vnitřními vazbami nebo vztahy na základě obsahu, tj. kontextu, a jediné, co určuje syntaktickou strukturu bezkontextového jazyka, jsou pravidla gramatiky.

Bezkontextová gramatika G je čtveřice $G = (N, T, P, S)$, kde:

- N je abeceda neterminálů
- T je abeceda terminálů a platí, že $N \cap T = \emptyset$
- P je konečná množina pravidel tvaru $A \rightarrow x$, kde $A \in N, x \in (N \cup T)^*$
- $S \in N$ je počáteční (startující) neterminál

Derivačním krokem se označuje akce, kdy se v dané větné formě provede změna řetězce za použití některého z pravidel v množině P . Obecná definice potom zní následovně. Necht' $G = (N, T, P, S)$ je bezkontextovou gramatikou, $u, v \in (N \cup T)^*$, tedy u, v jsou řetězce terminálních a neterminálních symbolů, které mohou být i prázdné, a $p = A \rightarrow x \in P$. Potom uAv přímo derivuje uxv za použití pravidla p v gramatice G , což lze formálně zapsat $uAv \Rightarrow uxv [p]$ či jednoduše $uAv \Rightarrow uxv$. Tato

definice neříká nic jiného, než že je možné přepsat neterminál A ve větě uAv řetězcem terminálů a neterminálů, popř. prázdným řetězcem (došlo by ke smazání neterminálu A), pokud v gramatice existuje příslušné pravidlo. Aplikováním více derivačních kroků vzniká sekvence derivačních kroků, ale její definice zde již nebude zmíněna, bližší podrobnosti lze nalézt např. v publikaci [2].

Jazyk generovaný bezkontextovou gramatikou je takový, jehož všechny řetězce w lze vygenerovat pomocí sekvence derivačních kroků z počátečního neterminálu S . Formální zápis vypadá následovně. Necht' $G = (N, T, P, S)$ je bezkontextová gramatika, pak jazyk generovaný bezkontextovou gramatikou G , $L(G)$, je definován:

$$L(G) = \{w: w \in T^*, S \Rightarrow w\}$$

Postupně tedy dochází k expanzi a přepisování neterminálů za pomoci pravidel gramatiky počínaje startujícím neterminálem S , až vznikne výsledný řetězec terminálů.

4.1 Derivační stromy a nejednoznačnost

Derivační strom je grafická reprezentace jedné či více derivací aplikovaných na danou větnou formu v podobě orientovaného grafu. Jeho nespornou výhodou je jasné a přehledné zobrazení provedených derivačních kroků. Mimo to umožňuje lépe pochopit některé problémy, které vyplývají z obecnosti bezkontextových gramatik a které by nemusely být hned zřejmé, obzvláště bez grafického znázornění.

Řeč je o nejednoznačnosti bezkontextových gramatik. Obecně bezkontextové gramatiky umožňují provést v každém derivačním kroku přepsání kteréhokoliv neterminálu za použití libovolného pravidla z množiny pravidel, které lze v daném kroku aplikovat na zpracovávanou větnou formu. V teoretickém světě by to nevadilo, naopak díky tomu vznikne v mnoha případech přehlednější gramatika, ale z hlediska praxe je to problém, protože různé posloupnosti aplikovaných pravidel na jedné a té samé větné formě mohou vygenerovat různé derivační stromy, což vede na již zmíněnou nejednoznačnost, a tedy i nedeterminističnost. V praxi je však potřeba mít jasně dané kroky, protože počítač nedokáže pracovat nedeterministicky, tedy hádat, co by bylo nejvhodnější. Určitě existují speciální případy a experimentální systémy, které se o to snaží, ale pro běžnou aplikaci syntaktických analyzátorů a překladačů je vyžadován deterministický přístup. Následující dva příklady se pokusí demonstrovat problém nejednoznačnosti na jednoduché gramatice pro základní matematické operace.

Příklad 2.: Necht' G je gramatika $G = (\{E\}, \{+, *, i\}, P, E)$, kde P obsahuje následující pravidla:

- $E \rightarrow E * E$
- $E \rightarrow E + E$
- $E \rightarrow i$

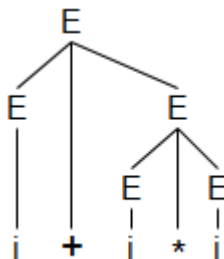
Vygenerujte řetězec $i + i * i$.

Jak již bylo zmíněno výše, bezkontextové gramatiky neříkají, které neterminály se mají rozgenerovat v jakém pořadí, pokud jich větná forma obsahuje více. Generování začne tedy ze startujícího neterminálu a poté se vybírají neterminály v různém pořadí. Toto je první polovina problému nejednoznačnosti, která je ale snadno řešitelná. Lze zavést omezení, že se v každém derivačním kroku

zpracuje vždy určitý neterminál. V praxi se tak zavádí tzv. nejlevější, resp. nejpravější derivace, která v každém derivačním kroku rozgeneruje vždy nejlevější, resp. nejpravější neterminál. Protože se jedná o speciální případ obecné derivace, lze bez újmy na obecnosti prohlásit, že aplikování vždy nejlevější nebo nejpravější derivace nemá vliv na výsledek, ale zavede to jistý řád a jednotnost při zpracování. Podrobnější vysvětlení a důkazy lze nalézt v publikaci [2].

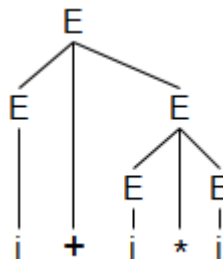
Náhodný výběr neterminálu

$$\begin{aligned} \underline{E} &\Rightarrow E + \underline{E} \\ &\Rightarrow E + \underline{E} * E \\ &\Rightarrow \underline{E} + i * E \\ &\Rightarrow i + i * \underline{E} \\ &\Rightarrow i + i * i \end{aligned}$$



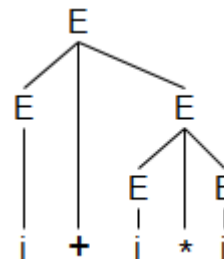
Nejlevější neterminál

$$\begin{aligned} \underline{E} &\Rightarrow \underline{E} + E \\ &\Rightarrow \underline{E} + E * E \\ &\Rightarrow i + \underline{E} * E \\ &\Rightarrow i + i * \underline{E} \\ &\Rightarrow i + i * i \end{aligned}$$



Nejpravější neterminál

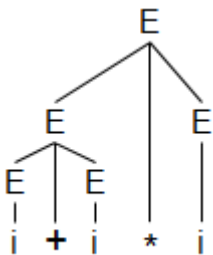
$$\begin{aligned} \underline{E} &\Rightarrow E + \underline{E} \\ &\Rightarrow E + E * \underline{E} \\ &\Rightarrow E + \underline{E} * i \\ &\Rightarrow \underline{E} + i * i \\ &\Rightarrow i + i * i \end{aligned}$$



I když se posloupnost provedených derivací liší, tak výsledný derivační strom je pro všechna pořadí (omezení) stejný. Tento výsledek je uspokojivý, protože se dá použít nejlevější nebo nejpravější derivace beze změny výsledku. Bohužel ale, jak lze vidět vzápětí, problém nejednoznačnosti spočívá v tom, že gramatika dovoluje použít libovolné pravidlo z množiny aktuálně dostupných pravidel. Pokud se zahájí derivace pravidlem s operací krát místo plus (gramatika tento krok dovolí, jedná se o naprosto legitimní derivaci), pak je možné vygenerovat stejný řetězec terminálů jako v předchozí variantě.

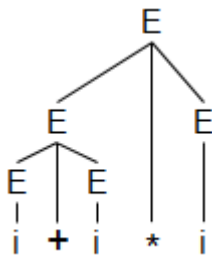
Náhodný výběr neterminálu

$$\begin{aligned} \underline{E} &\Rightarrow E * \underline{E} \\ &\Rightarrow \underline{E} * i \\ &\Rightarrow E + \underline{E} * i \\ &\Rightarrow \underline{E} + i * i \\ &\Rightarrow i + i * i \end{aligned}$$



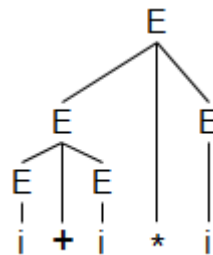
Nejlevější neterminál

$$\begin{aligned} \underline{E} &\Rightarrow \underline{E} * E \\ &\Rightarrow \underline{E} + E * E \\ &\Rightarrow i + \underline{E} * E \\ &\Rightarrow i + i * \underline{E} \\ &\Rightarrow i + i * i \end{aligned}$$



Nejpravější neterminál

$$\begin{aligned} \underline{E} &\Rightarrow E * \underline{E} \\ &\Rightarrow \underline{E} * i \\ &\Rightarrow E + \underline{E} * i \\ &\Rightarrow \underline{E} + i * i \\ &\Rightarrow i + i * i \end{aligned}$$



Opět se neliší výsledný derivační strom pro různé pořadí zpracovaných neterminálů, ale to jen díky tomu, že se zachovala návaznost aplikovaných pravidel. Díky grafické reprezentaci, jak bylo zmíněno výše, je na první pohled patrné, že pro stejný řetězec terminálů existují dva různé derivační stromy, a tedy gramatika G je nejednoznačná.

Pro předejití nejednoznačnosti gramatiky, potažmo jazyka, je nutné zavést determinismus. K dalšímu výkladu jsou potřeba zásobníkové automaty, resp. deterministické zásobníkové automaty, kterým bude věnována následující podkapitola.

4.2 Zásobníkové automaty

Stejně jako u regulárních jazyků byl jedním z formálních aparátů konečný automat, tak u bezkontextových jazyků se jedná o zásobníkový automat. Principiálně se od konečného automatu příliš neliší. Opět má množinu stavů, ve kterých se může nacházet, a ze vstupní pásky čte symboly. Oproti konečnému automatu ale pracuje navíc se zásobníkem. Následující stav je tedy dán aktuálním stavem, symbolem čteným z pásky a symbolem na zásobníku. Podle těchto tří hodnot rozhoduje, které pravidlo z množiny pravidel aplikovat v aktuálním kroku. Formální definice zásobníkového automatu je následující:

Zásobníkový automat je sedmice $M = (Q, \Sigma, \Gamma, R, s, S, F)$, kde:

- Q je konečná množina stavů
- Σ je vstupní abeceda
- Γ je zásobníková abeceda
- R je konečná množina pravidel tvaru $Apa \rightarrow wq$, kde $A \in \Gamma$, $p, q \in Q$, $a \in \Sigma \cup \{\varepsilon\}$, $w \in \Gamma^*$
- $s \in Q$ je počáteční stav
- $S \in \Gamma$ je počáteční symbol na zásobníku
- $F \subseteq Q$ je množina koncových stavů

Na rozdíl od konečného automatu zásobníkový automat může přijímat jazyk třemi způsoby, a to *přechodem do koncového stavu* (1), tento případ je totožný s konečným automatem, *vyprázdněním zásobníku* (2) nebo kombinací, tj. *přechodem do koncového stavu a zároveň vyprázdněním zásobníku* (3), přičemž všechny tyto varianty mají ekvivalentní vyjadřovací sílu. Formální definice jazyka přijímaného zásobníkovým automatem je následující. Necht' $M = (Q, \Sigma, \Gamma, R, s, S, F)$ je zásobníkový automat. Jazyk přijímaný zásobníkovým automatem M , $L(M)$, je definován:

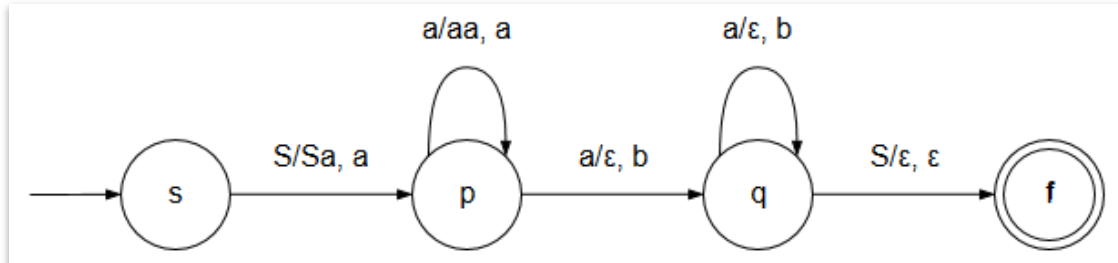
1. $L(M)_f = \{w: w \in \Sigma^*, Ssw \vdash^* zf, z \in \Gamma^*, f \in F\}$
2. $L(M)_\varepsilon = \{w: w \in \Sigma^*, Ssw \vdash^* zf, z = \varepsilon, f \in Q\}$
3. $L(M)_{f\varepsilon} = \{w: w \in \Sigma^*, Ssw \vdash^* zf, z = \varepsilon, f \in F\}$

Příklad 1.: Mějme zásobníkový automat $M = (\{s, p, q, f\}, \{a, b\}, \{a, S\}, R, s, S, \{f\})$, kde R obsahuje pravidla:

- $Ssa \rightarrow Sap$
- $apa \rightarrow aap$
- $apb \rightarrow q$

- $aqb \rightarrow q$
- $Sq \rightarrow f$

Jazyk přijímaný tímto automatem je $L(M) = \{a^n b^n : n \geq 1\}$. Tento zásobníkový automat je navíc schopen přijmout všemi třemi způsoby, které jsou popsány výše, protože v okamžiku, kdy přejde do koncového stavu, má zároveň i prázdný zásobník. Grafická reprezentace vypadá následovně:



Obrázek 3: Zásobníkový automat reprezentující jazyk $a^n b^n$, $n \geq 1$

Deterministický zásobníkový automat je pro uplatnění bezkontextových gramatik z praktického hlediska velmi podstatný nástroj. Drobnou úpravou definice obecného zásobníkového automatu tak vznikne aparát, u kterého se zamezilo náhodnému výběru či rozhodování, tedy nedeterminismu, a opakovaný výpočet proběhne vždy stejným způsobem. Je tedy nutné zavést omezení, aby automat nemohl provést v každém kroku výpočtu více než jeden přechod. Tedy formálně musí platit, že pro každé pravidlo tvaru $Apa \rightarrow wq \in R$ množina $R - \{Apa \rightarrow wq\}$ neobsahuje žádné pravidlo s levou stranou Apa nebo Ap .

Bohužel ne všechny bezkontextové gramatiky lze převést na ekvivalentní deterministickou formu, tedy třída *deterministických bezkontextových jazyků* \subset třída jazyků přijímaných *obecnými (nedeterministickými) zásobníkovými automaty*.

5 Gramatické systémy

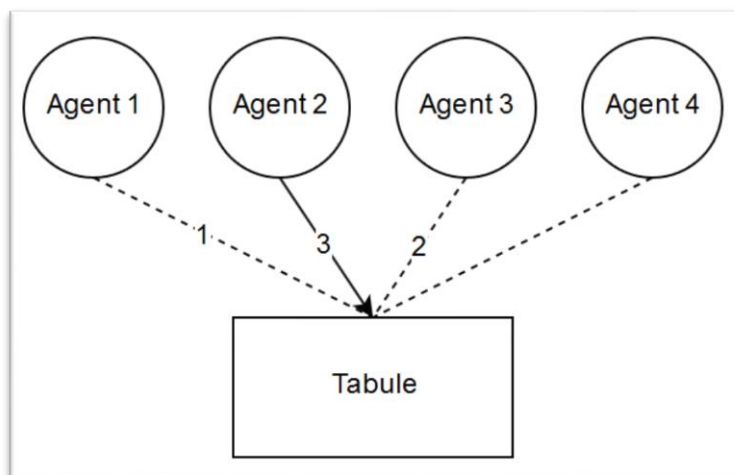
Gramatické systémy vznikly v průběhu druhé poloviny minulého století jako nová oblast formálních jazyků. Hlavní snahou bylo nalezení nástrojů a metod, které již v praxi existovaly dříve a běžně se používaly, ale nebyly pro ně vytvořeny formalismy a teoretický základ, o který by se dalo opřít při dalším zkoumání. V principu se gramatické systémy snaží pokrýt dvě oblasti informačních technologií, které nabraly velkého rozmachu právě ke konci minulého století a dnes je jejich vývoj stále velmi aktuální. Jedná se o distribuované a paralelní zpracování. V mnoha ohledech je výhodné i při běžné syntaktické analýze použít více gramatik, které se podílí na výpočtu. Pokud se k tomu zapojí možnosti distribuce výkonu nebo paralelního provádění, získá se nejen silný nástroje v rámci formálních jazyků, ale velké možnosti využití výpočetních zdrojů.

Nejlépe představitelnou variantou gramatického systému je syntaktická analýza, kdy se kombinuje přístup analýzy shora dolů s analýzou zdola nahoru. Dochází tak k přepínání dvou gramatik, které spolu komunikují a spolupracují při analýze, což umožňuje použít např. dvě jednodušší, přehlednější gramatiky, než se snažit celou analýzu provést pouze jednou gramatikou.

Následující dvě podkapitoly představí popořadě kooperující distribuované gramatické systémy, tzv. CD gramatické systémy, a paralelní gramatické systémy, tzv. PC gramatické systémy.

5.1 CD gramatické systémy

Gramatické systémy s označením CD definují model distribuované spolupráce gramatik, tedy kooperující distribuované gramatické systémy (z angličtiny *Cooperating Distributed Grammar Systems*). Tento systém si lze představit jako stolní hru. Hráči mají jasně dané pořadí a postupně se střídají. Když se hráč dostane na řadu, záleží na pravidlech hry, kdy musí svoji aktivitu předat dalšímu na řadě. Pravidla mohou být nastavena tak, že každý hráč provede přesně daný počet tahů (jeden nebo více) nebo provádí tahy tak dlouho, dokud může pokračovat. Poté předá hru dalšímu hráči, který je na řadě. Odborné články přirovnávají tento systém např. týmové spolupráci agentů, kteří řeší společný problém. Velmi často uváděný příklad je tzv. „problém tabule“, neboli *The Blackboard Model* či *The Blackboard System*. Skupina lidí řeší společný problém u tabule, kde se střídají, a každý z účastníků přispívá svými nápady, dokud dohromady problém nevyřeší.



Obrázek 4: Ukázka spolupráce agentů u tabule

CD gramatické systémy pracují sekvenčně na jedné sdílené větné formě. Každá gramatika, která je na řadě, provede modifikaci (příp. i více modifikací) větné formy a předá řízení další gramatice. Jednotlivé gramatiky jsou označovány jako tzv. komponenty systému. Následující podkapitoly uvedou jednotlivé definice a formalismy CD gramatických systémů.

5.1.1 Formální definice

V této práci jsou systémy definovány nad bezkontextovými gramatikami, které byly uvedeny v předchozí kapitole. CD gramatický systém Γ stupně n , $n \geq 1$, je $(n + 3)$ -tice, kde n značí počet množin pravidel systému, což lze chápat jako počet gramatik systému, které se budou podílet na výpočtu. Jednotlivé množiny pravidel (gramatiky) sdílí množiny terminálů, neterminálů a startující neterminál.

$\Gamma = (N, T, S, P_1, \dots, P_n)$, kde:

- N – abeceda neterminálů
- T – abeceda terminálů, $N \cap T = \emptyset$
- S – startující neterminál, $S \in N$
- $P_1 \dots P_n$ – konečná množina bezkontextových pravidel tvaru $a \rightarrow x$, $a \in N$, $x \in (N \cup T)^*$ (tzv. komponenty systému)

Formálně lze zapsat tento systém i takto:

$$\Gamma = G_1, \dots, G_n$$
$$G_i = (N, T, P_i, S) \text{ pro } 1 \leq i \leq n,$$

což vyplývá z popisu uvedeného výše. Zpravidla požadujeme, aby abecedy terminálů a neterminálů plus startující neterminál byly pro všechny participující gramatiky stejné, a ty se tak lišily pouze množinou pravidel.

5.1.2 Derivační módy

Derivační mód CD gramatického systému udává, jakým způsobem dochází k předávání řízení mezi jednotlivými gramatikami. Zavádí tedy pravidla, která určují, jak dlouho (kolik derivačních kroků) aktuálně aktivní gramatika může nebo musí provést, aby mohla předat řízení další gramatice, tj. komponentě systému.

Rozlišují se dva základní módy. První se nazývá ukončovací mód (anglicky *Terminating Mode*) a funguje velmi prostě. Každá komponenta pracuje (provádí přímé derivační kroky) tak dlouho, dokud může na větnou formu aplikovat některé ze svých pravidel. Tedy jednoduše řečeno „pracuje, dokud může“ a poté předá řízení další komponentě v pořadí.

Naopak druhý mód je striktnější. Nazývá se k -krokový mód (anglicky *K-Step Mode*) a určuje nebo omezuje počet derivačních kroků, které musí každá komponenta dodržet. Zavádí tedy tři podvarianty, z nichž každá má svoje označení:

- **právě k kroků** ($= k$) – každá komponenta musí provést přesný počet derivačních kroků, aby mohla předat řízení další komponentě. Pokud toto omezení nemůže splnit, pak daný jazyk není jazykem tohoto systému.
- **maximálně k kroků** ($\leq k$) – tato varianta je nejmírnější, protože komponenta může provést libovolný počet derivačních kroků, který nepřesáhne stanovené omezení k .
- **alespoň k kroků** ($\geq k$) – není tak striktní jako první varianta, ale vynucuje minimální počet derivačních kroků, které musí každá komponenta provést.

Okrajově je možné se setkat ještě se třetím módem, který je označován hvězdičkou (*). Tento mód je nejbenevolentnější a říká, že každá komponenta může pracovat tak dlouho, dokud sama potřebuje, což je velmi nedeterministické.

Postupně byly vysvětleny základní principy jednotlivých derivačních módů CD gramatických systémů a nyní následuje ukázka formálních definic a způsob jejich zápisu.

Terminating Mode (ukončující mód)

$x_i \Rightarrow^t y$, x derivuje y v i -té komponentě v t -módu, právě když:

1. x derivuje y v G_i ($x \Rightarrow^* y$ v $G_i = (N, T, P_i, S)$) a
2. y nederivuje/neexistuje přímá derivace y pro jakékoliv z nad abecedou $N \cup T$,
 $y \not\Rightarrow z$ pro $\forall z \in (N \cup T)^*$

K-Step Mode

1. k -step derivace $x_i \Rightarrow^{=k} y$ pokud $x \Rightarrow^k y$ v G_i
2. nejvýše k -step derivace $x_i \Rightarrow^{\leq k} y$ pokud $x \Rightarrow^j y$ v G_i pro $j \leq k$
3. alespoň k -step derivace $x_i \Rightarrow^{\geq k} y$ pokud $x \Rightarrow^j y$ v G_i pro $j \geq k$

Obecně se provádí přímé derivační kroky, dokud nevznikne řetězec terminálů a nelze pokračovat dále, v tomto okamžiku výpočet končí.

5.1.3 Generovaný jazyk

Spojením všech derivačních módů dohromady vzniká množina derivačních módů D .

$$D = \{*, t\} \cup \{\leq k, = k, \geq k : k = 1, 2, \dots\}$$

Pokud se k množině derivačních módů přidá gramatika, vznikne množina všech možných derivací F .

$$F(G_j, u, f) = \{v : u_j \Rightarrow^f v\}, \text{ kde } j \in \{1, \dots, n\}, f \in D, u \in V^*$$

Pro lepší pochopení, V je množina všech řetězců takových, že se z řetězce u módem f v j -té komponentě provede derivace z u do v .

Nechť $\Gamma = (N, T, S, P_1, \dots, P_n)$ je CD gramatický systém v derivačním módu $f \in D$, potom jazyk generovaný tímto módem je následující:

$$L_f(\Gamma) = \{w \in T^* : \text{kde jsou } v_0, v_1, \dots, v_n \text{ tak, že}$$

$$v_i \in F(G_{j_i}, v_{i-1}, f), i = 1, \dots, m, j_i \in \{1, \dots, n\}, v_0 = S, v_m = w \text{ pro } m \geq 1\}.$$

5.1.4 Příklady

Nyní následuje několik základních příkladů pro praktickou ukázkou, jak pracují CD gramatické systémy a jaké jazyky mohou generovat.

Příklad 1.: $L_t(\Gamma) = ?$ pro CD gramatický systém $\Gamma = (\{S, A\}, \{a\}, S, P_1, P_2, P_3)$ s komponentami:
 $P_1 = \{S \rightarrow AA\}$

$$P_2 = \{A \rightarrow S\}$$

$$P_3 = \{A \rightarrow a\}$$

Ze zadání je na první pohled patrné, že se jedná o ukončující mód. Tedy každá komponenta bude pracovat tak dlouho, dokud může aplikovat některé ze svých pravidel. Výpočet musí zahájit první komponenta P_1 , protože žádná jiná nemá pravidla pro rozgenerování startujícího neterminálu. Provede tedy derivační krok z S do AA a musí skončit, nemá již žádné další pravidlo, které by mohla aplikovat.

$$S \Rightarrow_{P_1} AA$$

Komponenta P_1 bude mít možnost ovlivnit větnou formu pouze po činnosti komponenty P_2 , takže na bezprostředně dalším výpočtu se budou podílet komponenty P_2 nebo P_3 . Komponenta P_3 má sílu ukončit výpočet v okamžiku, kdy se celá větná forma skládá pouze z neterminálů A , což je právě tento případ. Derivuje neterminály A na výsledné terminály, a protože v ukončujícím módu komponenta pracuje, dokud může, bude provádět derivace tak dlouho, dokud nezůstane řetězec terminálů, a výpočet tím skončí. Pro lepší ukázkou bude tedy pracovat komponenta P_2 , po které bude muset převzít řízení opět komponenta P_1 .

$$AA \Rightarrow_{P_2} SA \Rightarrow_{P_2} SS \Rightarrow_{P_1} AAAA$$

Pokud by nyní pokračovala komponenta P_2 , opakovala by se předchozí sekvence, proto bude ukázáno ukončení výpočtu komponentou P_3 .

$$AAAA \Rightarrow_{P_3} aAAA \Rightarrow_{P_3} aaAA \Rightarrow_{P_3} aaaa \Rightarrow_{P_3} aaaa$$

Vznikl tedy řetězec složený ze čtyř terminálů a , čímž výpočet končí. Je zřejmé, že podle toho, kolikrát bude pracovat komponenta P_2 , takovou délku bude mít výsledný řetězec terminálů a , přičemž samotná délka je omezena na mocniny dvou. V nejjednodušším případě skončí výpočet po třech derivačních krocích, pokud po komponentě P_1 bude pracovat rovnou komponenta P_3 a výsledkem bude řetězec aa . Výsledný jazyk tohoto systému je tedy:

$$\underline{\underline{L_t(\Gamma) = \{a^{2^n} : n \geq 1\}}}$$

Příklad 2.: $L_{=2}(\Gamma) = ?$ CD gramatický systém $\Gamma = (\{S, A, A', B, B'\}, \{a, b, c\}, S, P_1, P_2)$, kde:

$$P_1 = \{S \rightarrow S, S \rightarrow AB, A' \rightarrow A, B' \rightarrow B\}$$

$$P_2 = \{A \rightarrow aA'b, B \rightarrow cB', A \rightarrow ab, B \rightarrow c\}$$

Tento systém je o něco složitější než v prvním příkladu. Obsahuje dvě komponenty, z nichž každá obsahuje čtyři pravidla. Zároveň se už nejedná o ukončující mód, ale jak je podle požadovaného jazyka vidět, systém se nachází v k-step módu. Konkrétně tedy v první variantě tohoto módu, a to 2-step módu (= 2). Každá komponenta musí provést přesně dva derivační kroky, aby mohla předat řízení druhé komponentě. Pokud by nemohla tuto podmínku splnit, pak končí neúspěchem.

Samotné derivace probíhají stejně, jak to bylo detailně ukázáno v prvním příkladu, proto zde nebude vysvětlen celý výpočet krok po kroku, ale bude uveden pouze základní princip a možný postup. Generování musí začít komponenta P_1 , která má opět jako jediná pravidla pro zpracování startujícího neterminálu. První pravidlo komponenty P_1 , $S \rightarrow S$, je z praktického pohledu zbytečné, ale nesplnili bychom bez něho podmínku systému (= 2). Aplikují se tedy první a druhé pravidlo komponenty P_1 a předá se řízení komponentě P_2 . Ta opět provede dva derivační kroky a vrátí řízení zpět komponentě P_1 . Vždy se musí dodržet omezující podmínky systému, a proto se aplikují buďto dvě terminální nebo dvě neterminální pravidla. Výsledkem bude řetězec složený popořadě ze stejného počtu písmen a, b, c . Délka řetězce se odvíjí od toho, kolikrát se komponenty prostřídají.

$$\begin{aligned} S &\Rightarrow_{P_1} S \Rightarrow_{P_1} AB \Rightarrow_{P_2} aA'bB \Rightarrow_{P_2} aA'bcB' \Rightarrow_{P_1} aAbcB' \Rightarrow_{P_1} aAbcB \Rightarrow_{P_2} \\ &\Rightarrow_{P_2} aaA'bbcB \Rightarrow_{P_2} aaA'bbccB' \Rightarrow_{P_1} aaAbbccB' \Rightarrow_{P_1} aaAbbccB \Rightarrow_{P_2} \dots \\ &\Rightarrow_{P_1} a^n Ab^n c^n B \Rightarrow_{P_2} a^{n+1} b^{n+1} c^n B \Rightarrow_{P_2} a^{n+1} b^{n+1} c^{n+1} \end{aligned}$$

$$\underline{\underline{L_{=2}(\Gamma) = \{a^n b^n c^n : n \geq 1\}}}$$

Příklad 3.: $L_{=3}(\Gamma) = ?$ pro CD gramatický systém z příkladu 2.

Závěrečný příklad ukazuje situaci, kdy není možné splnit zadané podmínky systému. Použije se stejný systém z předchozího příkladu a pouze se upraví podmínka z *přesně dvou kroků* na *tři kroky*. Vznikne tedy 3-step mód (= 3). Opět začne komponenta P_1 , ale po aplikaci prvních dvou pravidel není již v komponentě další pravidlo, které by bylo možné aplikovat, a nelze tak splnit podmínku přesně tří derivačních kroků pro každou komponentu. Stejný výsledek by byl i pro variantu $k \geq 3$.

$$S \Rightarrow_{P_1} S \Rightarrow_{P_1} AB \Rightarrow_{P_1} \text{neexistuje}$$

$$\underline{\underline{L_{=3}(\Gamma) = L_{\geq 3}(\Gamma) = \emptyset}}$$

5.1.5 Třídy jazyků CD gramatických systémů a jejich generativní síla

Pro univerzální definování jednotlivých tříd jazyků CD gramatického systému bylo zavedeno následující značení:

$$CD_x^y(f), \text{ kde}$$

- f – derivační mód, $f \in D$
- y – pokuch chybí, potom nepovoluje ε -pravidla
 - ε , pak jsou ε -pravidla povolena
- $x = n$, nejvyšší stupeň CD gramatického systému, $n \geq 1$
 - ∞ , počet komponent systému není omezen

$CD_{\infty}(=)$ sjednocení všech tříd jazyků $CD_{\infty}(=k)$ pro $k = 1, 2, \dots$

$CD_{\infty}(\geq)$ sjednocení všech tříd jazyků $CD_{\infty}(\geq k)$ pro $k = 1, 2, \dots$

Generativní síla CD gramatických systémů se liší podle zavedených omezení či počtu komponent. Následujících pět teorémů zachycuje vztahy jednotlivých variant CD gramatických systémů a jejich generativní síly vůči ostatním typům.

1. $CD_{\infty}^y(f) = \mathcal{L}(CF)$, pro všechny $f \in \{= 1, \geq 1, *\} \cup \{\leq k: k \geq 1\}$
2. $\mathcal{L}(CF) = CD_1^y(f) \subset CD_2^y(f) \subseteq CD_r^y(f) \subseteq CD_{\infty}^y(f) \subseteq \mathcal{L}(M)$,
pro všechny $f \in \{= k, \geq k: k \geq 2\}$, $r \geq 3$
3. $CD_r^y(\geq k) \subseteq CD_r^y(\geq k+1)$
4. $CD_{\infty}^y(\geq) \subseteq CD_{\infty}^y(=)$
5. $\mathcal{L}(CF) = CD_1^y(t) = CD_2^y(t) \subset CD_3^y(t) = CD_{\infty}^y(t) = \mathcal{L}(ETOL)$

CD gramatický systém, který není omezen počtem komponent, připuštěním epsilon pravidel a který může použít libovolný k -step derivační mód s jedním či více kroky, má generativní sílu shodnou s třídou bezkontextových jazyků (teorém 1).

Zavede-li se omezení na k -step derivační mód se dvěma kroky a více (epsilon pravidla nejsou opět omezena), pak s počtem komponent roste generativní síla počínaje silou třídy bezkontextových jazyků a konče třídou maticových jazyků (teorém 2).

Pokud se vezmou dvě třídy jazyků CD gramatických systémů se stejným počtem komponent v módu „alespoň k kroků“, pak je třída s omezením na alespoň k kroků podmnožinou třídy s omezením alespoň $(k+1)$ kroků (teorém 3).

Podobně je na tom třída jazyků CD gramatických systémů v módu „alespoň k kroků“, která je podmnožinou třídě v módu „přesně k kroků“ (teorém 4).

Pokud se jedná o ukončující mód, potom generativní síla tříd jazyků CD gramatických systémů v ukončujícím módu roste s počtem komponent. Při omezení na jednu či dvě komponenty dosahuje síly shodné s třídou bezkontextových jazyků. Od tří komponent a více je síla stejná s generativní silou třídy ETOL jazyků (teorém 5). Podrobnější informace lze nalézt v publikaci [7].

5.1.6 Hybridní CD gramatické systémy

Zatímco klasické CD gramatické systémy definovaly jeden globální derivační mód pro všechny komponenty (gramatiky) systému, hybridní CD gramatické systémy umožňují definovat vlastní derivační mód pro každou komponentu zvlášť. Díky tomu lze každé komponentě na základě jejich pravidel a vhodnosti použití vybrat vlastní mód, ve kterém bude pracovat. Tedy jedna komponenta může pracovat v ukončujícím módu, zatímco druhá bude přepínat přesně po dvou krocích atd. Formální definice hybridního CD gramatického systému se liší pouze v rozšíření jednotlivých komponent o daný derivační mód.

$$\Gamma = (N, T, S, (P_1, f_1), \dots, (P_n, f_n)), \text{ kde:}$$

- N – abeceda neterminálů
- T – abeceda terminálů, platí $N \cap T = \emptyset$

- S – startující symbol, $S \in N$
- $P_1 \dots P_n$ – konečná množina bezkontextových pravidel tvaru $a \rightarrow x$, $a \in N$, $x \in (N \cup T)^*$
- $f_1 \dots f_n$ – mód i -té komponenty, $f_i \in D$ pro všechna $i \in \{1, \dots, n\}$

Formální definice jazyka generovaného hybridním CD gramatickým systémem:

$$L_f(\Gamma) = \{w \in T^* : \text{kde jsou } v_0, v_1, \dots, v_n \text{ tak, že}$$

$$v_i \in F(G_{j_i}, v_{i-1}, f_j), i = 1, \dots, m, j_i \in \{1, \dots, n\}, v_0 = S, v_m = w \text{ pro nějaké } m \geq 1\}.$$

5.1.7 Třídy jazyků hybridních CD gramatických systémů a jejich generativní síla

Obecný zápis pro jednotlivé třídy jazyků hybridních CD gramatických je rozšířením předchozí varianty pro klasické CD gramatické systémy:

$$XCD_{x,v}^y(f), \text{ kde}$$

- f – derivační mód, $f \in D$
- y – pokud chybí, potom nepovoluje ε -pravidla
 - ε , pak jsou ε -pravidla povolena
- $x = n$, nejvyšší stupeň CD gramatického systému, $n \geq 1$
 - ∞ , počet komponent systému není omezen
- $v = m$, každá komponenta P_i obsahuje maximálně m pravidel, $m > 1$
 - ∞ nebo nic, počet pravidel není omezen
- $X = \text{nic}$, nedeterministický systém
 - D , deterministický (pro každé pravidlo z P_i platí, že pro každý neterminál A existuje nejvýše jedna pravá strana; $A \rightarrow u, A \rightarrow w \in P_i$, pak $u = w$)
 - H , hybridní systém (nepíše se (f) , protože neexistuje globální mód)

I pro hybridní CD gramatické systémy lze uvést sadu teorémů, které demonstrují jejich generativní sílu vůči ostatním třídám jazyků v rámci Chomského hierarchie jazyků. Rozbor jednotlivých teorémů je analogický s předchozím vysvětlením klasických CD gramatických systémů, pouze s rozdílem, že je nutné sledovat jednotlivé atributy z definice tříd hybridních CD gramatických systémů.

1. $CD_{\infty, \infty}(f) = CD_{\infty, 1}(f) = \mathcal{L}(CF)$, pro všechny $f \in \{= 1, \geq 1, *\} \cup \{\leq k : k \geq 1\}$
2. $\mathcal{L}(CF) \subset CD_{\infty, 1}^{\varepsilon}(t) \subset CD_{\infty, 2}^{\varepsilon}(t) \subseteq CD_{\infty, 3}^{\varepsilon}(t) \subseteq CD_{\infty, 4}^{\varepsilon}(t) \subseteq CD_{\infty, 5}^{\varepsilon}(t) = CD_{\infty, \infty}^{\varepsilon}(t) = \mathcal{L}(ETOL)$
3. $CD_{n, m}(f) \subset CD_{n+1, m}(f)$, $f \in \{*, t\}$
4. $CD_{n, m}(f) \subset CD_{n, m+1}(f)$, $f \in \{*, t\}$
5. $\mathcal{L}(CF) = HCD_1 \subset HCD_2 \subseteq HCD_3 \subseteq HCD_4 = HCD_{\infty} = \mathcal{L}(M)$
6. $\mathcal{L}(ETOL) \subset HCD_4$
7. $CD_{\infty}(=) \subset HCD_3$

5.2 PC gramatické systémy

Paralelně komunikující gramatické systémy, zkráceně PC gramatické systémy (z angličtiny *Parallel Communicating Grammar Systems*), jak již z názvu vyplývá, pracují paralelně. To je zásadní rozdíl oproti CD gramatickým systémům, které pracují sekvenčně a postupně se střídají ve zpracování jedné společné větné formy. U PC gramatických systémů pracují všechny komponenty současně na své vlastní větné formě. Důležité je i druhé slovo názvu, „komunikující“, protože nejenže komponenty pracují paralelně, ale mohou mezi sebou i komunikovat pomocí tzv. komunikačních symbolů. Tímto způsobem si mezi sebou předávají mezivýsledky nebo si může jedna komponenta vyžádat práci jiné komponenty. Tohoto principu se využívá např. při syntaktické analýze. V mnoha případech může být užitečné a velmi praktické mít pro analýzu či překlad více různých gramatik, ať už pro různé části vstupního kódu nebo obecně může být analyzovaný kód tvořen více různými jazyky. Pak je vhodné aplikovat v průběhu analýzy odpovídající gramatiky. Zároveň však při komunikaci mezi komponentami hrozí riziko uvážnutí, tzv. deadlock, kterému je nutné se vyvarovat.

5.2.1 Formální definice

PC gramatický systém Γ stupně n , $n \geq 1$, je $(n + 3)$ -tice stejně jako CD gramatický systém, avšak svým zápisem je spíše podobný hybridnímu CD gramatickému systému. n značí opět počet množin pravidel systému, přičemž každá množina P má svůj vlastní startovací neterminál.

$$\Gamma = (N, K, T, (S_1, P_1), \dots, (S_n, P_n)), \text{ kde:}$$

- N – abeceda neterminálů
- K – konečná množina komunikačních symbolů, $K \in \{Q_1, \dots, Q_n\}$
- T – abeceda terminálů
- S_i – startující symbol i -té komponenty, $S_i \in N$ pro $i = 1, \dots, n$
- P_i – konečná množina pravidel tvaru $A \rightarrow x$, kde $A \in N$ a $x \in (N \cup T \cup K)^*$ pro $i = 1, \dots, n$
- Dále platí, že $K \cap N = \emptyset$, $K \cap T = \emptyset$

5.2.2 Derivační kroky

U PC gramatických systémů se nerozlišují derivační módy, jako tomu bylo u CD gramatických systémů, ale při výpočtu se střídají dva druhy derivačních kroků. Oba druhy derivačních kroků jsou velmi specifické a mají svůj speciální účel. Primárním derivačním krokem, dalo by se říci i režimem, je tzv. **generující krok** (g-Step), během kterého systém provádí klasické generativní kroky z (x_1, \dots, x_n) do (y_1, \dots, y_n) , tedy (x_1, \dots, x_n) přímo derivuje (y_1, \dots, y_n) . Tento režim je totožný s derivačními kroky, které provádí CD gramatické systémy. n -tice (x_1, \dots, x_n) a (y_1, \dots, y_n) , kde $x_i, y_i \in (N \cup T \cup K)^*$ pro všechna $i: 1 \leq i \leq n$, se nazývají konfigurace systému.

Zatím se PC gramatické systémy, až na paralelní provádění, tolik neliší od CD gramatických systémů. Významný rozdíl nastává při přepnutí systému do druhého režimu, tzv. **komunikační krok** (c-Step). K tomu dochází v okamžiku, kdy některá z komponent vygeneruje komunikační symbol. Zastaví se provádění generujících kroků ve všech komponentách a spustí se obluka komunikačního

kroku, ve které se provede přepsání vygenerovaného komunikačního symbolu uvnitř větné formy (konfigurace) komponenty, která tento symbol vygenerovala, řetězcem (konfigurací) komponenty, na kterou komunikační symbol odkazuje. Odkazovaná komponenta však nesmí mít ve své větné formě žádné komunikační symboly. Určitě se to zdá trochu zmatené, ale dále následuje definice jednotlivých režimů a příklady, ze kterých bude jejich fungování pochopitelnější. Důležité je, že komunikační krok (režim) má vyšší prioritu než generativní krok, takže jsou vždy nejdříve zpracovány komunikační symboly a až poté systém pokračuje v generativním kroku. Lze si tedy komunikační krok představit jako obsluhu přerušení, pokud se uvažuje analogii z operačních systémů.

g-Step

Pokud $x_i \Rightarrow y_i \vee G_i = (N \cup K, T, P_i, S_i)$ nebo

$x_i = y_i \in T^*$ pro všechny $1 \leq i \leq n$, pak

$$(x_1, \dots, x_n)_g \Rightarrow (y_1, \dots, y_n)$$

c-Step

Nastav $z_i = x_i$ pro všechny $i = 1, \dots, n$

Pro každé $i = 1, \dots, n$ pokud

$$\text{abeceda}(x_i) \cap K \neq \emptyset$$

a zároveň pro každé $Q_j \vee x_i$

$$\text{abeceda}(x_j) \cap K = \emptyset$$

potom pro každé $Q_j \vee x_i$

1. Nastav $z_j = S_j$
2. Nahraď Q_j za $x_j \vee x_i$
3. Nastav do z_i výsledek z kroku 2

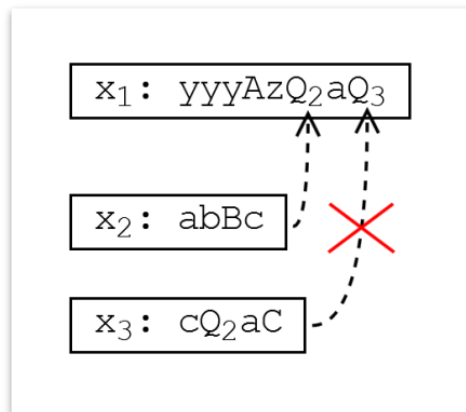
Proveď

$$(x_1, \dots, x_n)_c \Rightarrow (y_1, \dots, y_n)$$

pomocí $y_i = z_i$ pro všechny $i = 1, \dots, n$

Při generativním kroku, jak již bylo zmíněno výše, se provádí klasické derivační kroky z x_i do y_i pro všechny komponenty. V případě, že x_i již tvoří řetězce terminálů, provede se pouze simulace derivačního kroku, ale větná forma dané komponenty zůstane nezměněna. Ostatní komponenty mohou provést různé změny svých větných forem, zatímco tato „čeká“.

Naproti tomu komunikační krok provádí naprosto odlišné úkony a jeho algoritmus je daleko složitější, proto zde bude důkladně vysvětlen. Při úvodní inicializaci se do pole pomocných proměnných (z_i) uloží větné formy (x_i , pro všechna $i = 1, \dots, n$) jednotlivých komponent. Následně se v hlavním těle algoritmu procházejí větné formy. Pokud některá větná forma (x_i) obsahuje jeden či více komunikačních symbolů (Q_j) a zároveň komponenta, na kterou tento symbol odkazuje (x_j), neobsahuje žádný komunikační symbol, pak se pro tuto komponentu může vykonat komunikační krok. První podmínka je jasná, nelze provádět komunikační krok v komponentě, která neobsahuje žádný komunikační symbol. Druhá podmínka předchází zanořování při zpracování a především možnému zacyklení, protože pokud by komponenta (x_j), na kterou odkazuje komunikační symbol (Q_j) obsahovala také komunikační symboly, musely by se pro ni také v nejbližších iteracích vykonat komunikační kroky a v případě, že by komponenta odkazovala sama na sebe, došlo by k zacyklení.



Obrázek 5: Příklad porušení podmínky výskytu K symbolů. Požadavek Q_3 nebude v tomto kroku uspokojen.

Pokud je tedy podmínka splněna, potom se pro každou komponentu (x_j) nejdříve upraví pomocné pole tak, že do odpovídajících proměnných nastaví startující symbol příslušné komponenty ($z_j = S_j$). Tento krok se provádí u PC gramatických systémů s návraty (jeho další varianty budou vysvětleny později). Dále je to již přímočaré, přepíší se všechny komunikační symboly (Q_j) v původní větě formě (x_i) odpovídajícími větnými formami (x_j) a výsledek se uloží do (z_i). Až se zpracují všechny komponenty a komunikační symboly v tomto komunikačním kroku, přepíše se celé pomocné pole ($\forall z_i$) na výstup (do proměnných y_i).

5.2.3 Generovaný jazyk

PC gramatický systém, jak již bylo uvedeno výše, provádí buďto generativní kroky, které odpovídají klasickým derivačním krokům, nebo speciální komunikační kroky, které se snaží vypořádat s komunikačními symboly. Formálně lze generativní krok zapsat jako přechod z konfigurace (x_1, \dots, x_n) do konfigurace (y_1, \dots, y_n), tedy $(x_1, \dots, x_n)_g \Rightarrow (y_1, \dots, y_n)$, u komunikačního kroku je zápis obdobný, $(x_1, \dots, x_n)_c \Rightarrow (y_1, \dots, y_n)$. Obecně tedy konfigurace (x_1, \dots, x_n) provádí přímý derivační krok do konfigurace (y_1, \dots, y_n), $(x_1, \dots, x_n) \Rightarrow (y_1, \dots, y_n)$. Vychází-li se z počáteční konfigurace (S_1, \dots, S_n), která derivuje konfiguraci ($x, \alpha_2, \dots, \alpha_n$), pak je výsledným jazykem generovaným tímto PC gramatickým systémem množina všech řetězců terminálů, které se vygenerují v první komponentě x , ostatní komponenty ($\alpha_2, \dots, \alpha_n$) mohou obsahovat libovolné řetězce.

$$L(\Gamma) = \{x \in T^* : (S_1, \dots, S_n) \Rightarrow^* (x, \alpha_2, \dots, \alpha_n), \alpha_i \in (N \cup T \cup K)^*, \text{ pro všechny } i = 2, \dots, n\}$$

5.2.4 Varianty PC gramatických systémů

Centralizovaný PC gramatický systém je jedním ze speciálních případů klasického PC gramatického systému. Zavádí omezení, že pouze první komponenta P_1 může generovat komunikační symboly. Toto omezení má za přínos prevenci uváznutí, protože pouze do první komponenty mohou být vkládány větné formy ostatních komponent místo příslušných komunikačních symbolů a nemůže se stát, že by požadavek nemohl být uspokojen, protože by jiná komponenta obsahovala komunikační symboly. Celý systém je tedy řízen první komponentou, což se formálně zapíše následovně:

Nechť $\Gamma = (N, K, T, (S_1, P_1), \dots, (S_n, P_n))$ je PC gramatický systém. Γ je centralizovaný, pokud pro všechny $A \rightarrow x \in P_i$, kde $i = 2, \dots, n$, $\text{abeceda}(x) \cap K = \emptyset$.

Ve výkladu algoritmu komunikačního kroku byl zmíněn systém s návraty. Dle definice algoritmu se jedná o výchozí režim a vzniká v bodě 1 ($z_j = S_j$), kdy dojde k přepsání kopírované komponenty jejím startujícím neterminálem. U systému s návraty tedy nedochází pouze k nahrazení komunikačního symbolu v aktuálně zpracovávané komponentě, ale je zároveň i „resetována“ komponenta, která je kopírována. Může se pak v praxi stát, že první komponenta již obsahuje řetězec terminálů, který patří do jazyka tohoto systému, ale jiná komponenta si vyžádá jeho vkopírování (vygeneruje odpovídající komunikační symbol – toto je možné pouze u necentralizovaného systému), tím je první komponenta nahrazena startujícím neterminálem a její rozgenerování začíná od začátku. Jak je patrné z tohoto příkladu, lze kombinovat ne/centralizovaný systém s režimem s návraty a bez návratů. Jazyk generovaný systémem s režimem s návraty se značí $L_r(\Gamma)$.

Naopak PC gramatický systém bez návratů „neresetuje“ kopírovanou komponentu startujícím symbolem, ale pokračuje ve zpracování stávající větné formy. Je však nutné upravit definici algoritmu pro komunikační krok a odstranit z ní bod 1 ($z_j = S_j$). Jazyk generovaný tímto systémem se značí $L_{nr}(\Gamma)$.

5.2.5 Příklad

Příklad 1.: $L_r(\Gamma) = ?$ pro PC gramatický systém

$\Gamma = (\{S_1, S'_1, S_2, S_3\}, K, \{a, b, c\}, (S_1, P_1), (S_2, P_2), (S_3, P_3))$ s komponentami:

$$\begin{aligned} P_1 &= \{S_1 \rightarrow abc, S_1 \rightarrow a^2b^2c^2, S_1 \rightarrow aS'_1, S_1 \rightarrow a^3Q_2, S'_1 \rightarrow aS'_1, S'_1 \rightarrow a^3Q_2, \\ &\quad S_2 \rightarrow b^2Q_3, S_3 \rightarrow c\} \\ P_2 &= \{S_2 \rightarrow b S_2\} \\ P_3 &= \{S_3 \rightarrow c S_3\} \end{aligned}$$

Ze zadání lze na první pohled zjistit některé vlastnosti systému. Index r určuje, že se jedná o systém s návraty. Po bližším prozkoumání jednotlivých pravidel je patrné, že komunikační symboly generuje pouze první komponenta, takže se tedy jedná o centralizovaný systém. Samotný výpočet začíná obdobně jako u CD gramatických systémů, pouze s tím rozdílem, že je nutné provést derivaci paralelně pro každou komponentu. U komponenty P_1 dochází při startu k výběru ze tří pravidel, ale aby nebyl výpočet ukončen hned při prvním kroku, použije se pro ukázkou třetí pravidlo. Naopak u komponent P_2 a P_3 je generování přímočaré, obě komponenty mají každá pouze po jednom pravidle, které neustále produkují nové terminály. Start bude tedy následující:

$$(S_1, S_2, S_3) \Rightarrow (aS'_1, b S_2, c S_3)$$

Nevygeneroval se žádný komunikační symbol, takže systém zůstává v generativním kroku. Každá z komponent má „svůj“ terminál, který může aplikací vhodných pravidel generovat do nekonečna. Provede se tedy simulace nekonečného generování, aby byly pokryty všechny možnosti výsledného jazyka. Pro komponentu P_1 to znamená opakovat páté pravidlo, pro komponenty P_2 a P_3 stačí opakovat jejich jediná pravidla. Výsledkem jsou tak n -té mocniny jednotlivých terminálů:

$$(aS'_1, bS_2, cS_3) \Rightarrow^* (a^n S'_1, b^n S_2, c^n S_3)$$

Nyní je nutné pro komponentu P_1 vybrat jiné pravidlo, aby se výpočet posunul dále. Jediná možnost je šesté pravidlo v pořadí. Zbylé komponenty pokračují stejným způsobem.

$$(a^n S'_1, b^n S_2, c^n S_3) \Rightarrow (a^{n+3} Q_2, b^{n+1} S_2, c^{n+1} S_3)$$

Je vidět, že se v komponentě P_1 vygeneroval první komunikační symbol. Systém musí tedy přejít do komunikačního kroku a pokusit se tento komunikační symbol zpracovat. Jedná se o symbol Q_2 , což znamená, že odkazuje na komponentu P_2 , a vzhledem k tomu, že komponenta P_2 neobsahuje žádný komunikační symbol, nic nebrání v provedení komunikačního kroku. Důležité je ještě připomenout, že se jedná o systém s návraty, takže symbol Q_2 bude nahrazen větňou formou komponenty P_2 , a ta bude následně „resetována“ startujícím neterminálem.

$$(a^{n+3} Q_2, b^{n+1} S_2, c^{n+1} S_3) \Rightarrow (a^{n+3} b^{n+1} S_2, S_2, c^{n+1} S_3)$$

Tímto byly vyřešeny všechny komunikační symboly, tedy jeden, a následuje generativní krok. V komponentě P_1 je možné aplikovat pouze sedmé pravidlo, které opět vygeneruje komunikační symbol. V tomto případě se však jedná o požadavek na komponentu P_3 , protože symbol je Q_3 . Opět dojde k přechodu do komunikačního kroku, symbol Q_3 bude nahrazen větňou formou komponenty P_3 , a ta následně „resetována“ svým startujícím neterminálem.

$$(a^{n+3} b^{n+1} S_2, S_2, c^{n+1} S_3) \Rightarrow (a^{n+3} b^{n+3} Q_3, bS_2, c^{n+2} S_3) \Rightarrow (a^{n+3} b^{n+3} c^{n+2} S_3, bS_2, S_3)$$

Nyní lze v komponentě P_1 aplikovat pouze poslední pravidlo, čímž generování končí, protože vznikl řetězec terminálů. Zároveň se systém nachází v centralizovaném režimu, takže žádná jiná komponenta nemůže „resetovat“ první komponentu.

$$(a^{n+3} b^{n+3} c^{n+2} S_3, bS_2, S_3) \Rightarrow (a^{n+3} b^{n+3} c^{n+3}, b^2 S_2, c S_3)$$

Na závěr musí být jazyk doplněn o výsledky pravidel, kterými mohl systém ukončit generování hned v prvních krocích, a výsledný jazyk bude vypadat následovně:

$$\underline{\underline{L_r(\Gamma) = \{a^n b^n c^n : n \geq 1\}}}$$

5.2.6 Třídy jazyků PC gramatických systémů a jejich generativní síla

Pro PC gramatické systémy bylo zavedeno následující univerzální značení jednotlivých typů systémů a jejich režimů:

$$XPC_n Y, \text{ kde}$$

- $X - N$, systém bez návratů (non-returning mode)
– C , centralizovaný PC gramatický systém
- n – počet komponent systému, $n \geq 1$
- Y – typ pravidel (REG, LIN, CF)

Vzorové příklady: CPC_2REG , NPC_2LIN , $NCPC_\infty$

Následující teorémy zachycují několik vztahů mezi třídami jazyků PC gramatických systému a ostatních typů tříd jazyků:

1. $PC_nREG \subset PC_{n+1}REG$ pro $n > 1$
2. $CPC_nREG \subset CPC_nLIN \subset CPC_nCF$ pro $n > 1$
3. $NPC_\infty CF \subseteq PC_\infty CF$
4. $\mathcal{L}(M) \subset PC_\infty CF$
5. $\mathcal{L}(LIN) \subset PC_\infty REG$

První teorém popisuje vzrůstající sílu PC gramatických systémů s regulárními pravidly při zvětšujícím se počtu komponent. Čím více komponent je použito, tím větší se získá generativní síla. Druhý teorém není nijak překvapivý, popisuje vztah mezi centralizovanými PC gramatickými systémy s n komponentami, přičemž čím silnější třída pravidel se použije, tím silnější třída jazyků je získána. Třetí teorém uvažuje neomezený počet komponent a říká, že při použití bezkontextových pravidel je třída jazyků u systému bez návratů podmnožinou třídy jazyků systému s návraty. Zbylé dva teorémy porovnávají maticové a lineární jazyky s PC gramatickými systémy s neomezeným počtem komponent a bezkontextovými, resp. regulárními pravidly. Podrobnější informace s dalšími teorémy a důkazy lze nalézt v publikaci [7].

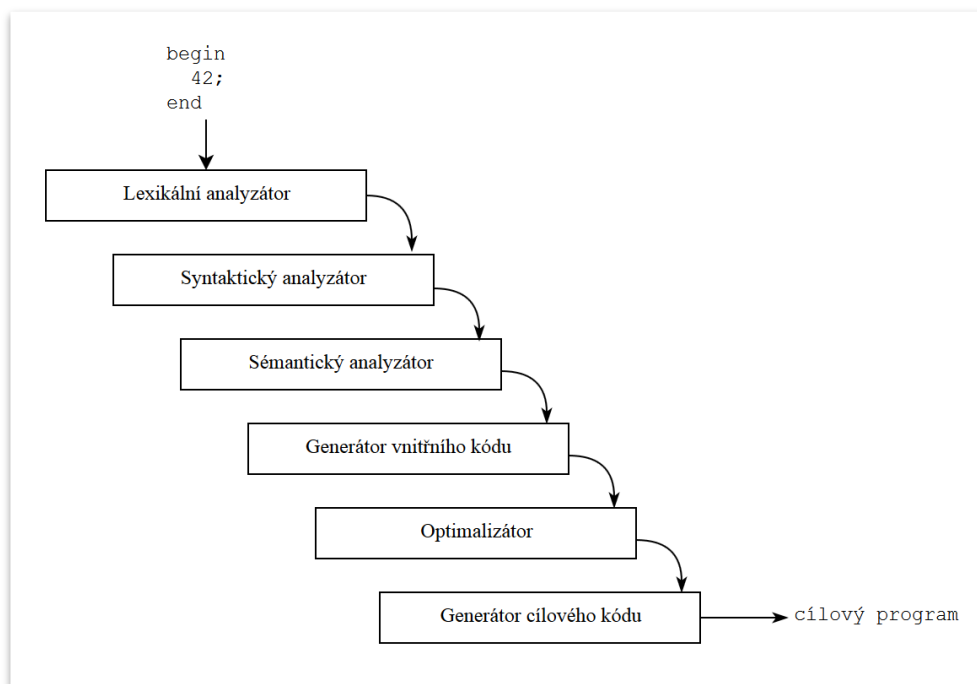
6 Syntaxí řízený překlad

Předchozí kapitoly pokrývaly především teorii. Práce je určená pro odborné čtenáře, ale přesto se úvod věnuje shrnutí obecnějšího základu, který by sice měl čtenář dobře znát, ale pro rychlé pochopení složitější látky je vhodné sjednotit i ty nejzákladnější stavební kameny. Nejdříve byly zopakovány obecné formální prostředky. Po nich následovaly regulární gramatiky, kterých se tato práce dotkne pouze v rámci implementace. Tvoří jádro lexikální analýzy, která je nezbytná jako přípravná fáze syntaktické analýzy. Poté následovala problematika bezkontextových gramatik a zásobníkových automatů. Bezkontextové gramatiky jsou velmi významné nejen z teoretického pohledu, ale umožňují definici a popis programovacích jazyků, takže jejich uplatnění v oblasti překladačů je velké. Následně již mohla navázat stěžejní část práce, a to samotné gramatické systémy, které formálně popsují součinnost gramatik a zavádějí různé modely jejich spolupráce.

Nyní je možné bezpečně přejít od teorie k praxi a zaměřit se na syntaxí řízený překlad a především pak jeho podčást, syntaktickou analýzu. Samotná syntaktická analýza slouží jako nástroj kontroly správné (syntaktické) struktury analyzovaného textu. Provádí se za pomoci požadované gramatiky definující jazyk, do kterého by měl kontrolovaný text patřit. Uvažujme tedy např. abecedu $\Sigma = \{\mathit{begin}, \mathit{,}, \mathit{int}, \mathit{id}, \mathit{end}, \dots\}$, gramatiku G , které odpovídá výsledný programovací jazyk L_G , a

zkoumáme, zda-li řetězec $x = \text{„begin 42; end“}$ patří do daného jazyka, $x \in L_G$, a tedy se jedná o syntakticky správně napsaný zdrojový kód, či nikoli, $x \notin L_G$. Syntaktická analýza tvoří základ překladu jazyků, a pokud je v celém procesu zapojena jako řídicí člen, jedná se o syntaxí řízený překlad. Překladač je obecně nástroj, který převádí zdrojový text (program napsaný ve zdrojovém jazyce) do podoby cílového textu (programu napsaného v cílovém jazyce), aniž by se změnil sémantický význam textu. U programů to znamená, že jsou oba programy, jak zdrojový, tak cílový, funkčně ekvivalentní. Zpravidla je v praxi žádoucí převést program z programovacího jazyka s vysokou abstrakcí, který je pro člověka dobře čitelný, do jazyka symbolických adres nebo přímo strojového kódu, aby ho mohl počítač snadno vykonávat.

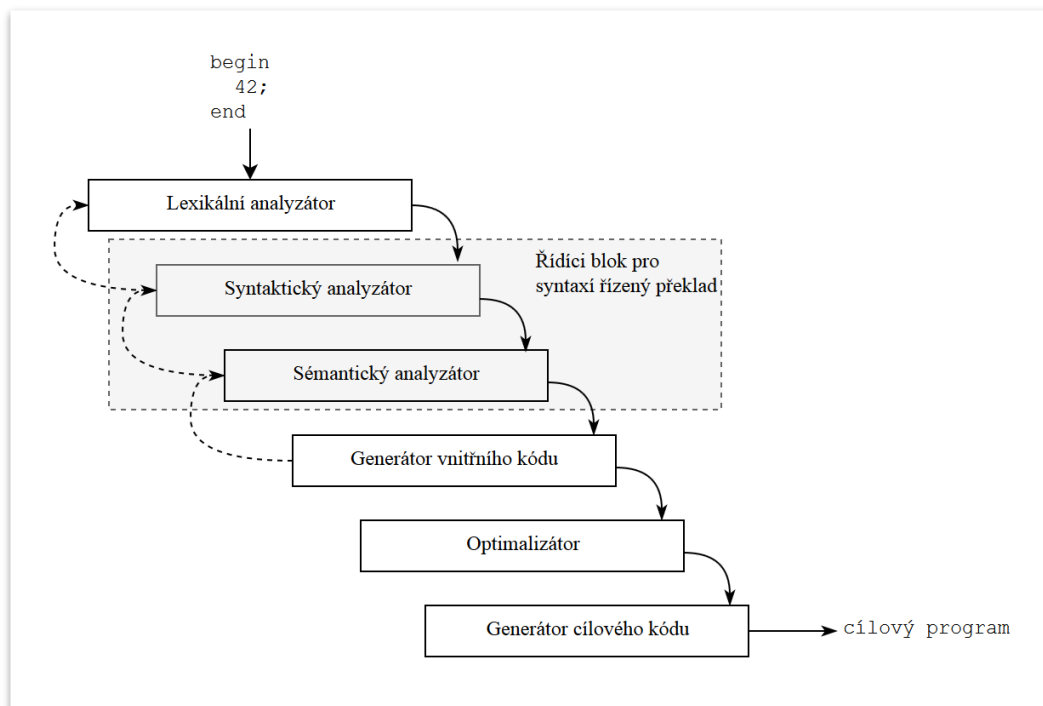
Obecný překladač si lze představit jako posloupnosti šesti na sebe navazujících bloků, které postupně zpracovávají vstupní text. Výstup prvního bloku je vstupem druhého bloku atd., až konečně poslední blok vygeneruje požadovaný cílový program.



Obrázek 6: Schéma obecného překladače

S ohledem na navazující praktickou část se tato práce zaměřuje pouze na syntaktickou analýzu, která je podmnožinou celého překladače a sestává se z prvních dvou bloků, tj. lexikálního analyzátoru a syntaktického analyzátoru. Další části překladače provádějí podrobnější kontroly kódu a připravují vše potřebné pro vygenerování cílového programu. Pokud je syntaxe vstupu v pořádku, potom sémantický analyzátor kontroluje datové typy, deklarace a inicializace proměnných a další náležitosti, na které již samotná syntaktická analýza nestačí. Následně zpravidla probíhá překlad do vnitřního (např. tříadresného) kódu, který je většinou univerzálně přenositelný a snadněji se nad ním provádějí volitelné optimalizace a generování výsledného kódu.

Pokud se jedná o syntaxí řízený překlad, bývají jednotlivé části překladače úzce propojené a komunikují spolu obousměrně během celého překladu, přičemž řídicím členem je syntaktický analyzátor.



Obrázek 7: Schéma překladače se syntaxí řízeným překladem

Při syntaxí řízeném překladu provádí syntaktický analyzátor kontrolu syntaktické struktury kódu a přitom volá a řídí ostatní bloky tak, aby byla aktuálně zpracovávaná část kódu rovnou zkontrolována i sémantickým analyzátozem a případně byl hned vygenerován odpovídající vnitřní kód. Tento přístup řízení je v praxi velmi častý. Není nutné, aby na sebe jednotlivé části překladače čekaly, ale analýza probíhá na více úrovních a umožňuje nejen úsporu času, ale i dřívější odhalení sémantických chyb, které by se jinak zjistily až po provedení celé syntaktické analýzy.

Lexikální analýza je nezbytná pro předpřípravu vstupního textu. Jejím vstupem je zdrojový program a výstupem řetězec tzv. tokenů. Jedná se o lexémy, tedy logické rozdělení vstupního kódu na lexikální jednotky jako jsou identifikátory, čísla, klíčová slova, operátory atd. Lexikální analýza vychází z povolené struktury vstupního textu (např. z jakých znaků a symbolů se může skládat identifikátor) a datových typů. Výsledné tokeny potom obsahují kromě samotného lexému i dodatečné informace, tzv. atributy, které usnadňují zbylý překlad. Lexikální analýza předchozího příkladu by mohla vypadat následovně:

Zdrojový program	Identifikace lexémů	Výsledné tokeny
Begin	begin	begin [keyword]
42;	42	42 [int]
End	;	;
	end	end [keyword]

Samotná implementace se provádí konečnými automaty, které byly právě z tohoto důvodu uvedeny na začátku práce.

Syntaktický analyzátor pracuje nad tokeny, které mu předá lexikální analyzátor, a snaží se k nim dle pravidel gramatiky sestavit odpovídající derivační strom. Pokud je takový strom nalezen, potom je vstupní text syntakticky správně. V opačném případě se jedná o syntaktickou chybu a

překlad končí. Pro sestavení derivačního stromu se uplatňují dva přístupy, a to shora dolů a zdola nahoru, které budou podrobně vysvětleny v následujících kapitolách.

Gramatické systémy lze zapojit do syntaktické analýzy mnoha způsoby dle jejich jednotlivých typů. Tato práce přistupuje k dané problematice z obecného hlediska a pokouší se aplikovat znalosti gramatických systémů na obecně známou kombinaci přístupů syntaktické analýzy shora dolů s analýzou zdola nahoru. Každý z těchto přístupů má svá specifika a je vhodnější pro jiný způsob použití. Pro účel této práce byla aplikována analýza shora dolů pro syntaktickou analýzu struktury kódu, zatímco analýza zdola nahoru pro zpracování matematických výrazů. Oba přístupy vyžadují vstupní text, neboli zdrojový kód, rozdělený na jednotlivé tokeny a hlavní rozdíl spočívá v tom, že se při syntaktické analýze shora dolů používá levý rozbor, tj. posloupnost pravidel nejlevější derivace vstupního řetězce, a simuluje se tak tvorbu derivačního stromu od kořene k listům, zatímco v případě analýzy zdola nahoru se provádí sestavení derivačního stromu od listů ke kořeni. V praxi to znamená, že je použít pravý rozbor, tj. reverzovaná posloupnost pravidel používající nejpravější derivaci vstupního řetězce. Pro oba přístupy existuje dále více variant, z nichž byla užitá tzv. prediktivní syntaktická analýza pro analýzu shora dolů. Tato metoda vyžaduje několik mezikroků vč. vytvoření tzv. LL tabulky, ale následná analýza je univerzální a nezávislá na konkrétní tabulce, což může být velmi užitečné. Pro analýzu zdola nahoru byla naopak použita jednodušší ze dvou variant, a to precedenční syntaktický analyzátor. Tato metoda je opět založena na sestavené speciální tabulce definující precedenci operátorů, závorek, identifikátorů a dalších symbolů a je velmi vhodná pro analýzu matematických výrazů. Tímto přístupem by se složitě popisovala struktura zdrojového kódu, ale díky spojení těchto dvou metod lze provést syntaktickou analýzu celého vstupního textu pomocí dvou jednodušších gramatik, což je jedna z hlavních myšlenek gramatických systémů a i jejich analogií, které zobrazují situace, kdy se více účastníků podílí na řešení společného problému a každý přispěje svými znalosti. Samostatně by celý problém nevyřešili, ale společně je to možné.

V následujících kapitolách jsou popsány oba přístupy, které byly použity ve výsledné aplikaci a porovnání s jejich alternativami. Samotný závěr práce je věnován implementaci a ukázce výsledné aplikace.

6.1 LL gramatiky bez epsilon pravidel

Jednou z hlavních překážek využití bezkontextových gramatik pro syntaktickou analýzu shora dolů je existence nejednoznačnosti, která umožňuje vytvoření více derivačních stromů pro stejnou větnou formu. Tato skutečnost vnáší do celého přístupu nedeterminističnost, což je při praktickém použití problém. Vznikl tak speciální typ gramatik, které nesou označení LL gramatiky. Název vznikl ze způsobu analýzy, která probíhá zleva doprava (*Left*) a konstruuje se nejlevější derivace věty (*Leftmost*). Nejdříve bude vysvětlena základní definice LL gramatiky bez epsilon pravidel a poté postupně popsáno vytvoření pomocných množin s následnou konstrukcí LL tabulky, která je nezbytná pro vytvoření prediktivního syntaktického analyzátoru. Definice LL gramatiky bez epsilon pravidel je následující:

Nechť $G = (N, T, P, S)$ je bezkontextová gramatika bez ε -pravidel. G je LL gramatika, pokud pro každé $a \in T$ a $A \in N$ existuje maximálně jedno pravidlo tvaru $A \rightarrow X_1 X_2 \dots X_n \in P$ a zároveň platí, že $a \in \text{First}(X_1 X_2 \dots X_n)$.

6.1.1 Množina First

LL gramatika bez epsilon pravidel je poměrně jednoduchá na použití, a proto bude vhodné na ní vysvětlit hlavní principy a poté pokračovat komplexnější verzí s epsilon pravidly. Hlavní myšlenkou množiny **First** a zároveň i algoritmu **First(x)** je výpočet všech terminálních symbolů, kterými může začínat větná forma derivovaná z **x**. Zjednodušeně řečeno, pokud je provedena derivace **x**, tak všechny terminální symboly, které se mohou objevit na prvním místě místo **x**, musí být vloženy do množiny **First(x)**. Pro terminální symboly bude jejich množina **First** obsahovat sebe sama, zatímco pro neterminály je nutné určit všechny možné terminály, které se z nich mohou derivovat. Tímto se daný algoritmus stává složitějším i přesto, že dosud nebyla použita epsilon pravidla, se kterými je výpočet ještě komplikovanější. Formální definice množiny **First** je následující:

Nechť $G = (N, T, P, S)$ je bezkontextová gramatika. Pro každé $x \in (N \cup T)^*$ je definována množina **First(x)** jako: $\text{First}(x) = \{a: a \in T, x \Rightarrow^* ay, y \in (N \cup T)^*\}$.

Algoritmus pro výpočet **First(x)** očekává na vstupu bezkontextovou gramatiku bez epsilon pravidel a jeho výstupem je množina **First(x)** pro každý terminál i neterminál.

1. Každému terminálu $a \in T$ vytvoř množinu $\text{First}(a) = \{a\}$
2. Aplikuj následující pravidlo, dokud bude možné měnit nějakou množinu **First**:
if $A \rightarrow X_1X_2 \dots X_n \in P$ **then** přidej $\text{First}(X_1)$ do $\text{First}(A)$

Příklad 1.: Mějme LL gramatiku pro jednoduchý program s následujícími pravidly:

- 1: $\langle \text{prog} \rangle \rightarrow \text{begin } \langle \text{state} \rangle$
- 2: $\langle \text{state} \rangle \rightarrow \langle \text{cmd} \rangle \langle \text{item} \rangle ; \langle \text{state} \rangle$
- 3: $\langle \text{state} \rangle \rightarrow \text{end}$
- 4: $\langle \text{cmd} \rangle \rightarrow \text{print}$
- 5: $\langle \text{cmd} \rangle \rightarrow \text{read}$
- 6: $\langle \text{item} \rangle \rightarrow \text{int}$
- 7: $\langle \text{item} \rangle \rightarrow \text{id}$

Jak lze vidět z jednotlivých pravidel, tento jazyk vyžaduje, aby program začal klíčovým slovem **begin** a skončil slovem **end**, mezi kterými může být libovolné množství příkazů **print** a **read**, které mají oba jako parametr buďto číslo (**int**) nebo proměnnou označenou jako **id**. Nejdříve se sestaví množina **First** pro všechny terminály a neterminály, a poté je již možné zkonstruovat LL tabulku.

Množiny **First** pro všechny terminály dle kroku 1.:

- $\text{First}(\text{begin}) = \{\text{begin}\}$
- $\text{First}(;) = \{;\}$
- $\text{First}(\text{print}) = \{\text{print}\}$
- $\text{First}(\text{read}) = \{\text{read}\}$
- $\text{First}(\text{end}) = \{\text{end}\}$
- $\text{First}(\text{int}) = \{\text{int}\}$
- $\text{First}(\text{id}) = \{\text{id}\}$

Výpočet množin **First** pro všechny neterminály dle kroku 2.:

$\langle \text{prog} \rangle \rightarrow \text{begin} \in P$: přidej $\text{First}(\text{begin})$ do $\text{First}(\langle \text{prog} \rangle)$
 $\text{First}(\langle \text{prog} \rangle) = \{\text{begin}\}$

$\langle \text{state} \rangle \rightarrow \langle \text{cmd} \rangle \in P$: přidej $\text{First}(\langle \text{cmd} \rangle)$ do $\text{First}(\langle \text{state} \rangle)$
 $\langle \text{state} \rangle \rightarrow \text{end} \in P$: přidej $\text{First}(\text{end})$ do $\text{First}(\langle \text{state} \rangle)$
 $\text{First}(\langle \text{state} \rangle) = \{\text{end}\}$

$\text{First}(\langle \text{cmd} \rangle)$ zatím nebylo zpracováno, takže se z ní do $\text{First}(\langle \text{state} \rangle)$ nic nepřidalo a bude nutné se v další iteraci k této množině pravděpodobně vrátit a doplnit ji. Tuto situaci ošetřuje podmínka druhého pravidla, která požaduje opakované iterování, dokud je možné něco měnit.

$\langle \text{cmd} \rangle \rightarrow \text{print} \in P$: přidej $\text{First}(\text{print})$ do $\text{First}(\langle \text{cmd} \rangle)$

$\langle \text{cmd} \rangle \rightarrow \text{read} \in P$: přidej $\text{First}(\text{read})$ do $\text{First}(\langle \text{cmd} \rangle)$

$\text{First}(\langle \text{cmd} \rangle) = \{\text{print}, \text{read}\}$

$\langle \text{item} \rangle \rightarrow \text{int} \in P$: přidej $\text{First}(\text{int})$ do $\text{First}(\langle \text{item} \rangle)$

$\langle \text{item} \rangle \rightarrow \text{id} \in P$: přidej $\text{First}(\text{id})$ do $\text{First}(\langle \text{item} \rangle)$

$\text{First}(\langle \text{item} \rangle) = \{\text{int}, \text{id}\}$

Nyní se provede druhá iteraci, ve které lze změnit už jenom výše zmíněnou množinu $\text{First}(\langle \text{state} \rangle)$.

$\langle \text{state} \rangle \rightarrow \langle \text{cmd} \rangle \in P$: přidej $\text{First}(\langle \text{cmd} \rangle)$ do $\text{First}(\langle \text{state} \rangle)$

$\text{First}(\langle \text{state} \rangle) = \{\text{end}, \text{print}, \text{read}\}$

Tímto jsou zpracovány všechny množiny **First**, žádná z nich již nejde změnit, takže algoritmus končí a následuje konstrukce LL tabulky. LL tabulka je velmi užitečný nástroj, který jednoznačně určuje, jaké pravidlo aplikovat na základě aktuálního neterminálu a terminálu (tokenu) na vstupu. Pokud se daná kombinace v tabulce nevyskytuje, nastává syntaktická chyba.

α	...	a	...
...			
A		$\alpha(A, a)$	
...			

Tabulka 1: Ukázka rozložení LL tabulky

Jak je možné vidět v tabulce, řádky indexují neterminály a sloupce terminály. Každá buňka pak obsahuje číslo pravidla, které se má aplikovat pro danou kombinaci. Formálně je definována následovně:

$\alpha(A, a) = A \rightarrow X_1 X_2 \dots X_n \in P$ pokud $a \in \text{First}(X_1)$,
jinak je $\alpha(A, a)$ prázdné \Rightarrow syntaktická chyba.

	begin	;	end	print	read	int	id
<prog>	1						
<state>			3	2	2		
<cmd>				4	5		
<item>						6	7

Tabulka 2: LL tabulka pro Příklad 1

6.2 LL gramatiky s epsilon pravidly

Pro obecnou představu je základní myšlenka konstrukce LL tabulky dostačující, ale pro využití systému v praxi by byly LL gramatiky bez epsilon pravidel velmi omezující, protože z reálného prostředí by popsaly jen malé množství systémů. Je proto nutné zavést LL gramatiky s epsilon pravidly. Jediné, v čem se liší od LL gramatik bez epsilon pravidel, je složitost sestavení LL tabulky. Ta u gramatik s epsilon pravidly vyžaduje několik pomocných množin, které navzájem využívají svoje data, a LL tabulka se sestaví pomocí poslední vytvořené množiny, takzvané množiny **Predict**. Nyní bude uvedena obecná definice LL gramatik s epsilon pravidly a poté představeny jednotlivé pomocné množiny. K ukázce poslouží trochu upravený příklad z předchozí části. Postupně budou sestrojeny dané množiny, až opět vznikne LL tabulka. Následně bude uveden příklad jejího použití při syntaktické analýze, která se pro sestavenou LL tabulku neliší, zda se jednalo o gramatiku s epsilon pravidly nebo bez.

Nechť $G = (N, T, P, S)$ je bezkontextová gramatika. G je LL gramatika s epsilon pravidly, pokud pro každé $a \in T$ a $A \in N$ existuje maximálně jedno pravidlo tvaru $A \rightarrow X_1X_2 \dots X_n \in P$ a zároveň platí, že $a \in \text{Predict}(X_1X_2 \dots X_n)$.

V definici se objevuje množina **Predict**, která je právě tou poslední množinou při vytváření LL tabulky. Kromě již známé množiny **First** je potřeba se seznámit ještě s dalšími dvěma množinami, a to **Empty** a **Follow**. Základním problémem LL gramatiky s epsilon pravidly je jistá neurčitost terminálů z množiny **First**. Zatímco v předchozím případě bylo možné se spolehnout, že daný terminál v množině **First** bude vždy, pokud byl pro tuto množinu vypočítán, tak nyní vzniká problém, že se tento terminál v konkrétní množině **First** může vyskytovat. Pokud totiž dojde k expanzi neterminálu pomocí pravidla, které obsahuje na pravé straně další neterminály a některé z nich (nebo všechny) derivují epsilon, pak může dojít k jejich vymazání a počáteční terminál může být symbol derivovaný z následujícího neterminálu v pravidle nebo může vypadnout celý, na počátku přepisovaný neterminál, a prvním terminálem v tomto místě bude až první terminál za tímto neterminálem.

Mějme pravidlo $r: A \rightarrow X_1X_2 \dots X_n \in P$ a aktuální větnou formu xAy . V případě, že by se jednalo o gramatiku bez epsilon pravidel, tak po derivaci A pomocí pravidla r by se jako první terminál ve větné formě objevilo místo A nějaké $a \in \text{First}(A)$, které by zároveň bylo ve $\text{First}(X_1)$ atd. Jenomže pokud $X_1X_2 \dots X_n$ může derivovat epsilon, pak mohou některá vypadnout a může dojít k situaci, kdy buď $a \in \text{First}(X_i)$ pro $2 \leq i \leq n$, nebo mohou být vymazána všechna, pak $a \in \text{First}(y)$.

Pro epsilon pravidla je výpočet mnohem složitější a je nutné prozkoumat gramatiku opravdu do hloubky, odhalit všechny možnosti a kombinace, které symboly mohou být kterými nahrazeny a

zastoupeny. K tomu pomůže postupné a systematické sestavení všech čtyř množin **Empty**, **First**, **Follow** a **Predict**.

Pro výpočet množin se zavádí pomocný symbol \$, který slouží k zakončení každého vstupního řetězce. Jedná se o nutnost způsobenou epsilon pravidly. Pokud by se celý vstupní řetězec derivoval na epsilon, je jisté, že symbol \$ zůstane, a definuje tak konec vstupu.

6.2.1 Množina Empty

Množina **Empty** je nejjednodušší a pomáhá určit, jestli daný symbol derivuje epsilon nebo ne. Je to tedy startující bod pro analýzu epsilon derivací, které nemusí být vidět v gramatice na první pohled, ale díky návaznosti pravidel může být daný neterminál nahrazen za epsilon i přes několik pravidel (derivací). Obecně je tedy **Empty(x)** množina, která obsahuje pouze jediný prvek, a to ϵ , pokud x derivuje ϵ , jinak je prázdná. Lze tak říci, že množina **Empty** není přímo množinou tak, jak je běžně chápána, ale pouze dvoustavovým atributem určujícím, jestli se může na místě daného symbolu objevit epsilon nebo ne. Ostatní množiny již obsahují více symbolů, a proto i **Empty(x)** se považuje za množinu pro jednotný přístup a manipulaci se všemi čtyřmi množinami. Formální definice množiny **Empty** je následující:

Nechť $G = (N, T, P, S)$ je bezkontextová gramatika. Pro každé $x \in (N \cup T)^*$ je definována množina **Empty(x)** jako: $\text{Empty}(x) = \{\epsilon\}$ pokud $x \Rightarrow^* \epsilon$ jinak $\text{Empty}(x) = \emptyset$.

Algoritmus pro výpočet **Empty(x)** očekává na vstupu bezkontextovou gramatiku a jeho výstupem je množina **Empty(x)** pro každý terminál i neterminál.

1. Každému terminálu $a \in T$ vytvoř množinu $\text{Empty}(a) = \emptyset$
2. Pro každý neterminál $A \in N$ rozhodni:
pokud $A \rightarrow \epsilon \in P$ pak $\text{Empty}(A) = \{\epsilon\}$, jinak $\text{Empty}(A) = \emptyset$.
3. Aplikuj následující pravidlo, dokud bude možné měni nějakou množinu **Empty**:
pokud $A \rightarrow X_1 X_2 \dots X_n \in P$ a zároveň $\text{Empty}(X_i) = \{\epsilon\}$ pro všechna $i = 1, \dots, n$,
pak $\text{Empty}(A) = \{\epsilon\}$

Příklad 2.: Mějme LL gramatiku s epsilon pravidly pro program s následujícími pravidly:

- 1: $\langle \text{prog} \rangle \rightarrow \text{begin } \langle \text{state} \rangle$
- 2: $\langle \text{state} \rangle \rightarrow \langle \text{cmd} \rangle \langle \text{item} \rangle ; \langle \text{state} \rangle$
- 3: $\langle \text{state} \rangle \rightarrow \text{end}$
- 4: $\langle \text{cmd} \rangle \rightarrow \text{print}$
- 5: $\langle \text{cmd} \rangle \rightarrow \text{read}$
- 6: $\langle \text{cmd} \rangle \rightarrow \epsilon$
- 7: $\langle \text{item} \rangle \rightarrow \text{int}$
- 8: $\langle \text{item} \rangle \rightarrow \text{id}$

Pro praktickou ukázkou je mírně upraven příklad z předchozí části pouze přidáním jednoho pravidla, pravidla č. 6, které umožňuje, aby neterminál $\langle \text{cmd} \rangle$ derivoval epsilon. Mohlo by se zdát, že se jedná o triviální změnu, která nemůže mít zásadní vliv na konstrukci LL tabulky, ale opak je pravdou. Může

zde dojít ke vzorové situaci, která byla popsána na začátku této kapitoly o LL gramatikách s epsilon pravidly. Pokud se totiž v pravidle č. 2 derivuje první neterminál $\langle \text{cmd} \rangle$ na epsilon, pak pro **First**($\langle \text{state} \rangle$) nemusí být první terminál ve **First**($\langle \text{cmd} \rangle$), ale až ve **First**($\langle \text{item} \rangle$). Pokud by byla gramatika složitější a existovalo pravidlo $\langle \text{item} \rangle \rightarrow \varepsilon$, potom by mohl být prvním terminálem až symbol ‘;’ atd. Je tedy nutné při výpočtu LL tabulky zohlednit tuto neurčitost a zahrnout do množin **First** všechny potenciální možnosti, když se někde objeví epsilon.

Množiny **Empty** pro všechny terminály dle kroku 1.:

$\text{Empty}(\text{begin}) = \emptyset$
 $\text{Empty}(\text{;}) = \emptyset$
 $\text{Empty}(\text{print}) = \emptyset$
 $\text{Empty}(\text{read}) = \emptyset$
 $\text{Empty}(\text{end}) = \emptyset$
 $\text{Empty}(\text{int}) = \emptyset$
 $\text{Empty}(\text{id}) = \emptyset$

Výpočet množin **First** pro všechny neterminály dle kroku 2.:

$\langle \text{prog} \rangle \rightarrow \text{begin} \in P:$ $\langle \text{prog} \rangle \rightarrow \varepsilon \notin P$
Empty($\langle \text{prog} \rangle$) = \emptyset

$\langle \text{state} \rangle \rightarrow \langle \text{cmd} \rangle \in P:$
 $\langle \text{state} \rangle \rightarrow \text{end} \in P:$ $\langle \text{state} \rangle \rightarrow \varepsilon \notin P$
Empty($\langle \text{state} \rangle$) = \emptyset

$\langle \text{cmd} \rangle \rightarrow \text{print} \in P:$
 $\langle \text{cmd} \rangle \rightarrow \text{read} \in P:$
 $\langle \text{cmd} \rangle \rightarrow \varepsilon \in P:$ $\langle \text{cmd} \rangle \rightarrow \varepsilon \in P$
Empty($\langle \text{cmd} \rangle$) = $\{\varepsilon\}$

$\langle \text{item} \rangle \rightarrow \text{int} \in P:$
 $\langle \text{item} \rangle \rightarrow \text{id} \in P:$ $\langle \text{item} \rangle \rightarrow \varepsilon \notin P$
Empty($\langle \text{item} \rangle$) = \emptyset

Po prvním průchodu přes neterminály gramatiky je patrné, že pouze pro neterminál $\langle \text{cmd} \rangle$ je množina **Empty**($\langle \text{cmd} \rangle$) = $\{\varepsilon\}$, ostatní neterminály epsilon přímo nederivují. Při provedení dalšího průchodu pro splnění kroku č. 3 algoritmu **Empty**, je zřejmé, že se již žádná množina **Empty** nezmění. Jediným kandidátem na změnu by byla množina **Empty**($\langle \text{state} \rangle$), protože pravidlo č. 2 obsahuje neterminál $\langle \text{cmd} \rangle$, jehož množina **Empty** ve druhém průchodu již obsahuje ε , ale pravá strana tohoto pravidla obsahuje další neprázdné neterminály a zároveň i jeden terminál, takže neplatí podmínka pro pravou stranu pravidla **Empty**(X_i) = $\{\varepsilon\}$ pro všechna $i = 1, \dots, n$, a žádná množina **Empty** již opravdu nejde změnit. Tímto byly vyřešeny všechny množiny **Empty** pro danou gramatiku a následuje modifikovaná verze množiny **First** pro LL gramatiky s epsilon pravidly.

6.2.2 Množina *First*

Význam a základní myšlenka množiny *First* byla uvedena již v předchozí sekci o LL gramatikách bez epsilon pravidel, takže není nutné zacházet do takové hloubky a lze přejít rovnou ke změněnému algoritmu. Největším problémem epsilon pravidel je, jak již bylo zmíněno, nutnost ošetření případu, kdy dojde ke smazání prvních neterminálů (jednoho či více) pravé strany pravidla (přepsáním na epsilon), a následující terminály a neterminály se tak posunou na jejich místo. Za normálních okolností by se vyskytovaly až za nimi. Musí se tedy do množin *First* zahrnout nejen terminály, které se mohou derivovat na místě prvního neterminálu na pravé straně jednotlivých pravidel, ale postihnout i případy, kdy je tento první neterminál smazán. Potom se musí zahrnout do dané množiny *First* terminály následujícího neterminálu. Může nastat i situace, že se tímto způsobem dojde až na konec pravé strany pravidla, pak je potřeba sjednotit všechny množiny *First* jednotlivých neterminálů vyskytujících se na pravé straně tohoto pravidla, protože nikdy není dopředu zřejmé, kolik z nich může být zleva popořadě přepsáno na epsilon. Situaci, kdy by byly přepsány na epsilon všechny (množina *Empty* neterminálu na levé straně pravidla by obsahovala epsilon), řeší později množina *Follow*.

Algoritmus pro výpočet *First(x)* očekává na vstupu bezkontextovou gramatiku a jeho výstupem je množina *First(x)* pro každý terminál i neterminál.

1. Každému terminálu $a \in T$ vytvoř množinu $First(a) = \{a\}$
2. Každému neterminálu $A \in N$ vytvoř množinu $First(A) = \emptyset$
3. Aplikuj následující pravidlo, dokud bude možné měnit nějakou množinu *First*:
pokud $A \rightarrow X_1X_2 \dots X_{k-1}X_k \dots X_n \in P$
 - a. přidej všechny symboly z $First(X_1)$ do $First(A)$
 - b. pokud $Empty(X_i) = \{\epsilon\}$ pro $i = 1, \dots, k-1$, $k \leq n$, přidej všechny symboly z $First(X_k)$ do $First(A)$

Algoritmus je velmi podobný původní variantě, pouze s tím rozdílem, že předtím se automaticky přiřadily $First(X_1)$ do $First(A)$, což se nyní stane také, ale zároveň se kontrolují všechna X_i zleva doprava a pokud tvoří skevenci, kde jednotlivá $Empty(X_i) = \{\epsilon\}$, pak jsou postupně přidávána $First(X_i)$ do $First(A)$, dokud podmínka nepřestane platit nebo se nedojde na konec pravé strany pravidla. Vždy je přidána i první množina $First(X_i)$ pro X_i , u kterého již neplatí $Empty(X_i) = \{\epsilon\}$, popř. X_n , pokud $X_n = X_k$.

Množiny *First* pro všechny terminály (1.):

```
First(begin) = {begin}
First(;) = {;}
First(print) = {print}
First(read) = {read}
First(end) = {end}
First(int) = {int}
First(id) = {id}
```

Množiny *First* pro všechny neterminály (2.):

```
First(<prog>) = ∅
First(<state>) = ∅
First(<cmd>) = ∅
First(<item>) = ∅
```

Po inicializaci všech množin **First** je možné pokračovat jejich výpočtem dle kroku 3, ke kterému jsou potřeba množiny **Empty** z předchozí kapitoly:

```
<prog> → begin <state> ∈ P:
    Empty(begin) ≠ {ε}
        přidej First(begin) do First(<prog>)
First(<prog>) = {begin}
```

```
<state> → <cmd> <item> ; <state> ∈ P:
    Empty(<cmd>) = {ε}
        přidej First(<cmd>) do First(<state>)
    Empty(<item>) ≠ {ε}
        přidej First(<item>) do First(<state>)
```

```
<state> → end ∈ P:
    Empty(end) ≠ {ε}
        přidej First(end) do First(<state>)
First(<state>) = {end}
```

```
<cmd> → print ∈ P:
    Empty(print) ≠ {ε}
        přidej First(print) do First(<cmd>)
```

```
<cmd> → read ∈ P:
    Empty(read) ≠ {ε}
        přidej First(read) do First(<cmd>)
First(<cmd>) = {print, read}
```

```
<item> → int ∈ P:
    Empty(int) ≠ {ε}
        přidej First(int) do First(<item>)
```

```
<item> → id ∈ P:
    Empty(id) ≠ {ε}
        přidej First(id) do First(<item>)
First(<item>) = {int, id}
```

Nyní se provede druhá iterace, ve které lze změnit opět už jenom množinu $First(\langle state \rangle)$, nyní však na ni bude mít vliv jak množina $First(\langle cmd \rangle)$, tak i změněná množina $First(\langle item \rangle)$:

```
<state> → <cmd> <item> ; <state> ∈ P:
    Empty(<cmd>) = {ε}
        přidej First(<cmd>) do First(<state>)
    Empty(<item>) ≠ {ε}
        přidej First(<item>) do First(<state>)
First(<state>) = {end, print, read, int, id}
```

Lze zde vidět zásadní rozdíl oproti předchozí ukázce bez epsilon pravidel. Množina $\text{First}(\langle \text{state} \rangle)$ se rozšířila o dva terminály int a id díky tomu, že v pravidle č. 2 může být první neterminál $\langle \text{cmd} \rangle$ smazán pomocí epsilon pravidla a na jeho místo se tak dostane druhý neterminál $\langle \text{item} \rangle$.

6.2.3 Algoritmus $\text{First}(X_1X_2\dots X_n)$

Pro výpočet zbývajících množin se zavádí rozšíření stávajících algoritmů $\text{Empty}(x)$ a $\text{First}(x)$, aby dokázaly vypočítat příslušné množiny pro celé pravé strany pravidel, takže ne jenom pro jeden symbol x , ale pro celou sekvenci symbolů $X_1X_2 \dots X_n$.

Algoritmus pro výpočet $\text{First}(X_1X_2 \dots X_n)$ očekává na vstupu bezkontextovou gramatiku $G = (N, T, P, S)$ a předpočítané množiny $\text{First}(X)$ a $\text{Empty}(X)$ pro každé $X \in N \cup T$, $x = X_1X_2 \dots X_n$, kde $x \in (N \cup T)^+$. Výstupem je potom množina $\text{First}(X_1X_2 \dots X_n)$. Dále platí, že $\text{First}(\varepsilon) = \emptyset$.

1. $\text{First}(X_1X_2 \dots X_n) = \text{First}(X_1)$
2. Aplikuj následující pravidlo, dokud bude možné měnit množinu $\text{First}(X_1X_2 \dots X_{k-1}X_k \dots X_n)$:
pokud $\text{Empty}(X_i) = \{\varepsilon\}$ pro $i = 1, \dots, k-1$, $k \leq n$, přidej všechny symboly z $\text{First}(X_k)$ do $\text{First}(X_1X_2 \dots X_{k-1}X_k \dots X_n)$

Algoritmus se příliš neliší od předchozí verze. Opět se v prvním kroku vypořádá s případem bez epsilon pravidel, kdy do výsledné množiny $\text{First}(X_1X_2 \dots X_n)$ přiřadí $\text{First}(X_1)$, a poté se již opakuje situace z předchozího algoritmu, kde postupně testuje množiny Empty pro jednotlivá X_i , a přidává jejich množiny First do $\text{First}(X_1X_2 \dots X_n)$, dokud nenarazí na první, u které již neplatí, že se její $\text{Empty}(X_i) = \{\varepsilon\}$. Její First ještě přidá do $\text{First}(X_1X_2 \dots X_n)$ a pokud už nemůže nic měnit, výpočet končí.

6.2.4 Algoritmus $\text{Empty}(X_1X_2\dots X_n)$

Rozšíření $\text{Empty}(X_1X_2 \dots X_n)$ bylo zavedeno ze stejného důvodu jako předchozí algoritmus pro $\text{First}(X_1X_2 \dots X_n)$. Bude potřeba u následujících množin a jeho výpočet je ze všech nejjednodušší. Umožňuje určit, jestli může být celá pravá strana pravidla prázdná (Empty) nebo se z ní vždy nějaký terminál vygeneruje.

Algoritmus pro výpočet $\text{Empty}(X_1X_2 \dots X_n)$ očekává na vstupu bezkontextovou gramatiku $G = (N, T, P, S)$ a předpočítanou množinu $\text{Empty}(X)$ pro každé $X \in N \cup T$, $x = X_1X_2 \dots X_n$, kde $x \in (N \cup T)^+$. Výstupem je potom množina $\text{Empty}(X_1X_2 \dots X_n)$. Dále platí, že $\text{Empty}(\varepsilon) = \{\varepsilon\}$.

1. Pokud $\text{Empty}(X_i) = \{\varepsilon\}$ pro $i = 1, \dots, n$, potom $\text{Empty}(X_1X_2 \dots X_n) = \{\varepsilon\}$
2. jinak $\text{Empty}(X_1X_2 \dots X_n) = \emptyset$

6.2.5 Množina Follow

Množina **Follow** umožňuje pokrýt ty případy, kdy se v dané větné formě díky epsilon pravidlu smaže celý neterminál a je potřeba předem vědět, co za ním může následovat, resp. se dostat na jeho místo. **Follow(A)** je tedy množina všech terminálů, které se mohou vyskytovat vpravo bezprostředně od **A** ve větné formě.

Nechť $G = (N, T, P, S)$ je bezkontextová gramatika. Pro všechna $A \in N$ je definována množina **Follow(A)** jako:

$$\mathbf{Follow}(A) = \{a : a \in T, S \Rightarrow^* xAay, x, y \in (N \cup T)^*\} \cup \{\$: S \Rightarrow^* xA, x \in (N \cup T)^*\}$$

Musí se tak pro všechny neterminály vypočítat jejich množiny **Follow**, tedy určit všechny možné terminály, které za nimi mohou následovat pro případ, že by byl daný neterminál smazán pomocí epsilon pravidla.

Algoritmus pro výpočet **Follow(A)** očekává na vstupu bezkontextovou gramatiku $G = (N, T, P, S)$ a předpočítané množiny **Empty(X)** a **First(X)**. Výstupem je potom množina **Follow(A)** pro každé $A \in N$.

1. **Follow(S) = {\$}**
2. Aplikuj následující pravidlo, dokud bude možné měnit nějakou množinu **Follow**:
pokud $A \rightarrow xBy \in P$ potom
 - a. Pokud $y \neq \epsilon$ potom přidej všechny symboly z **First(y)** do **Follow(B)**
 - b. Pokud **Empty(y) = {ε}** potom přidej všechny symboly z **Follow(A)** do **Follow(B)**

Objevuje se zde s již dříve zmíněným symbolem \$, který pomáhá určit konec vstupního řetězce, a proto se přiřazuje hned na začátku algoritmu do **Follow(S)**. To je vlastně nejvýstižnější popis množiny **Follow**. Předpokládá se, že je vstupní řetězec vždy zakončen symbolem \$. Pokud by byl celý vstupní řetězec prázdný, tedy existovalo by např. pravidlo $S \rightarrow \epsilon$, pak bezprostředně za celým řetězcem, který se může derivovat z **S** (např. celý vstupní program), vždy následuje symbol \$. V tomto případě dokonce místo neterminálu **S**, protože byl epsilon pravidlem smazán. Bod a) druhého kroku postihuje variantu, kdy se za neterminálem **B** vyskytují další symboly, potom se přidá jejich množina **First** do **Follow(B)**. Druhá nezávislá podmínka v bodě b) ověří, jestli mohou být symboly za neterminálem **B Empty**. Buďto mohou být smazány epsilon pravidlem nebo se jedná o první neterminál zprava, v tom případě platí dodatek **Empty(ε) = {ε}**. V obou případech je přidána množina **Follow(A)** do **Follow(B)**, aby byla k dispozici informace, co může následovat bezprostředně za neterminálem **B**. Tímto způsobem se určí množiny **Follow** pro všechny neterminály.

Výpočet množin **Follow** pro krok 1.:

$$\mathbf{Follow}(\langle \text{prog} \rangle) = \{\$\}$$

Nyní se bude provádět krok 2. pro všechny neterminály, dokud bude možné měnit některou z množin **Follow**:

$\langle \text{prog} \rangle \rightarrow \text{begin } \langle \text{state} \rangle \underline{\varepsilon}$
 $\text{Empty}(\varepsilon) = \{\varepsilon\}$
 přidej $\text{Follow}(\langle \text{prog} \rangle)$ do $\text{Follow}(\langle \text{state} \rangle)$
 $\text{Follow}(\langle \text{state} \rangle) = \{\$\}$

$\langle \text{state} \rangle \rightarrow \langle \text{cmd} \rangle \langle \text{item} \rangle ; \langle \text{state} \rangle \underline{\varepsilon}$
 $\text{Empty}(\varepsilon) = \{\varepsilon\}$
 přidej $\text{Follow}(\langle \text{state} \rangle)$ do $\text{Follow}(\langle \text{state} \rangle)$
 $\text{Follow}(\langle \text{state} \rangle) = \{\$\}$

$\langle \text{state} \rangle \rightarrow \langle \text{cmd} \rangle \langle \text{item} \rangle \underline{; \langle \text{state} \rangle}$
 $; \langle \text{state} \rangle \neq \varepsilon$
 přidej $\text{First}(; \langle \text{state} \rangle)$ do $\text{Follow}(\langle \text{item} \rangle)$
 $\text{Empty}(; \langle \text{state} \rangle) \neq \{\varepsilon\}$
 $\text{Follow}(\langle \text{item} \rangle) = \{;\}$

$\langle \text{state} \rangle \rightarrow \langle \text{cmd} \rangle \underline{\langle \text{item} \rangle ; \langle \text{state} \rangle}$
 $\langle \text{item} \rangle ; \langle \text{state} \rangle \neq \varepsilon$
 přidej $\text{First}(\langle \text{item} \rangle ; \langle \text{state} \rangle)$ do $\text{Follow}(\langle \text{cmd} \rangle)$
 $\text{Empty}(\langle \text{item} \rangle ; \langle \text{state} \rangle) \neq \{\varepsilon\}$
 $\text{Follow}(\langle \text{cmd} \rangle) = \{\text{int}, \text{id}\}$

6.2.6 Množina Predict

Množina **Predict** je poslední množinou, která je nutná pro sestavení LL tabulky, resp. ostatní množiny sloužily pouze jako podpora k sestavení množiny **Predict**, ze které se už přímo konstruuje LL tabulka pro LL gramatiku s epsilon pravidly. Množina **Predict** se váže na jednotlivá pravidla a pro každé pravidlo gramatiky umožňuje **Predict**($A \rightarrow x$) určit množinu všech terminálů, které mohou být aktuálně nejlevěji vygenerovány, pokud se pro libovolnou větnou formu použije pravidlo $A \rightarrow x$.

Nechť $G = (N, T, P, S)$ je bezkontextová gramatika. Pro každé $A \rightarrow x \in P$ je definována množina **Predict**($A \rightarrow x$) následovně:

1. Pokud $\text{Empty}(x) = \{\varepsilon\}$ potom **Predict**($A \rightarrow x$) = $\text{First}(x) \cup \text{Follow}(A)$
2. jinak pokud $\text{Empty}(x) = \emptyset$ potom **Predict**($A \rightarrow x$) = $\text{First}(x)$

Nyní se vypočítá množiny **Predict** pro jednotlivá pravidla gramatiky našeho vzorového příkladu:

1: $\langle \text{prog} \rangle \rightarrow \text{begin } \langle \text{state} \rangle$
 $\text{Empty}(\text{begin } \langle \text{state} \rangle) = \emptyset$
 $\text{Predict}(1) = \text{First}(\text{begin } \langle \text{state} \rangle) = \text{First}(\text{begin}) = \{\text{begin}\}$

2: $\langle \text{state} \rangle \rightarrow \langle \text{cmd} \rangle \langle \text{item} \rangle ; \langle \text{state} \rangle$

$\text{Empty}(\langle \text{cmd} \rangle \langle \text{item} \rangle ; \langle \text{state} \rangle) = \emptyset$
 $\text{Predict}(2) = \text{First}(\langle \text{cmd} \rangle \langle \text{item} \rangle ; \langle \text{state} \rangle) =$
 $= \text{First}(\langle \text{cmd} \rangle \langle \text{item} \rangle) = \{\mathbf{print, read, int, id}\}$

3: $\langle \text{state} \rangle \rightarrow \text{end}$
 $\text{Empty}(\text{end}) = \emptyset, \text{Predict}(3) = \text{First}(\text{end}) = \{\mathbf{end}\}$

4: $\langle \text{cmd} \rangle \rightarrow \text{print}$
 $\text{Empty}(\text{print}) = \emptyset, \text{Predict}(4) = \text{First}(\text{print}) = \{\mathbf{print}\}$

5: $\langle \text{cmd} \rangle \rightarrow \text{read}$
 $\text{Empty}(\text{read}) = \emptyset, \text{Predict}(5) = \text{First}(\text{read}) = \{\mathbf{read}\}$

6: $\langle \text{cmd} \rangle \rightarrow \varepsilon$
 $\text{Empty}(\varepsilon) = \{\varepsilon\}$
 $\text{Predict}(6) := \text{First}(\varepsilon) \cup \text{Follow}(\langle \text{cmd} \rangle) = \{\mathbf{int, id}\}$

7: $\langle \text{item} \rangle \rightarrow \text{int}$
 $\text{Empty}(\text{int}) = \emptyset, \text{Predict}(7) = \text{First}(\text{int}) = \{\mathbf{int}\}$

8: $\langle \text{item} \rangle \rightarrow \text{id}$
 $\text{Empty}(\text{id}) = \emptyset, \text{Predict}(8) = \text{First}(\text{id}) = \{\mathbf{id}\}$

Na závěr výpočtu se provede sestavení výsledné LL tabulky pro tento vzorový příklad. K vyplnění tabulky postačí pouze množina **Predict**, ostatní množiny jsou pomocné. Celá tabulka je totožná s předchozí verzí, která již byla detailně popsána v předchozí kapitole. Liší se způsob výběru pravidla pro jednotlivé buňky a přibyl jeden sloupec pro zakončovací terminál \$:

$\alpha(A, a) = A \rightarrow X_1 X_2 \dots X_n \in P$ pokud $a \in \text{Predict}(A \rightarrow X_1 X_2 \dots X_n)$,
 jinak je $\alpha(A, a)$ prázdné \Rightarrow syntaktická chyba.

	begin	;	end	print	read	int	id	\$
$\langle \text{prog} \rangle$	1							
$\langle \text{state} \rangle$			3	2	2	2	2	
$\langle \text{cmd} \rangle$				4	5	6	6	
$\langle \text{item} \rangle$						7	8	

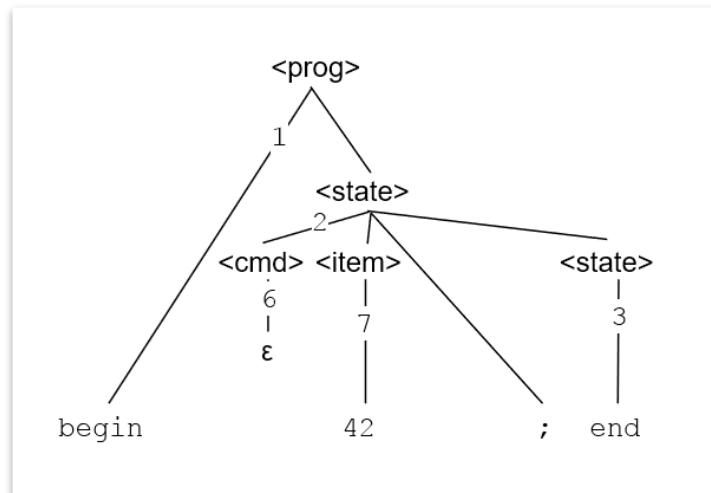
Tabulka 3: Výsledná LL tabulka pro Příklad 2

LL tabulka se oproti předchozí variantě bez epsilon pravidel rozšířila o případy, kdy dojde ke smazání neterminálu $\langle \text{cmd} \rangle$, který má jako jediný epsilon pravidlo, a syntaktický analyzátor tak díky tabulce ví, jak se s touto situací vypořádat, a které symboly očekávat.

Příklad 3.: Proveďte syntaktickou analýzu následujícího vstupního kódu pomocí LL tabulky sestavené výše:

```
begin
  42;
end
```

Následuje ukázka jednoduchého způsobu použití LL tabulky při syntaktické analýze. Zatím pouze intuitivně a velmi triviálně, později bude vysvětlen algoritmus prediktivní syntaktické analýzy, pomocí které lze provádět analýzu automaticky a nad libovolně velkou LL tabulkou.



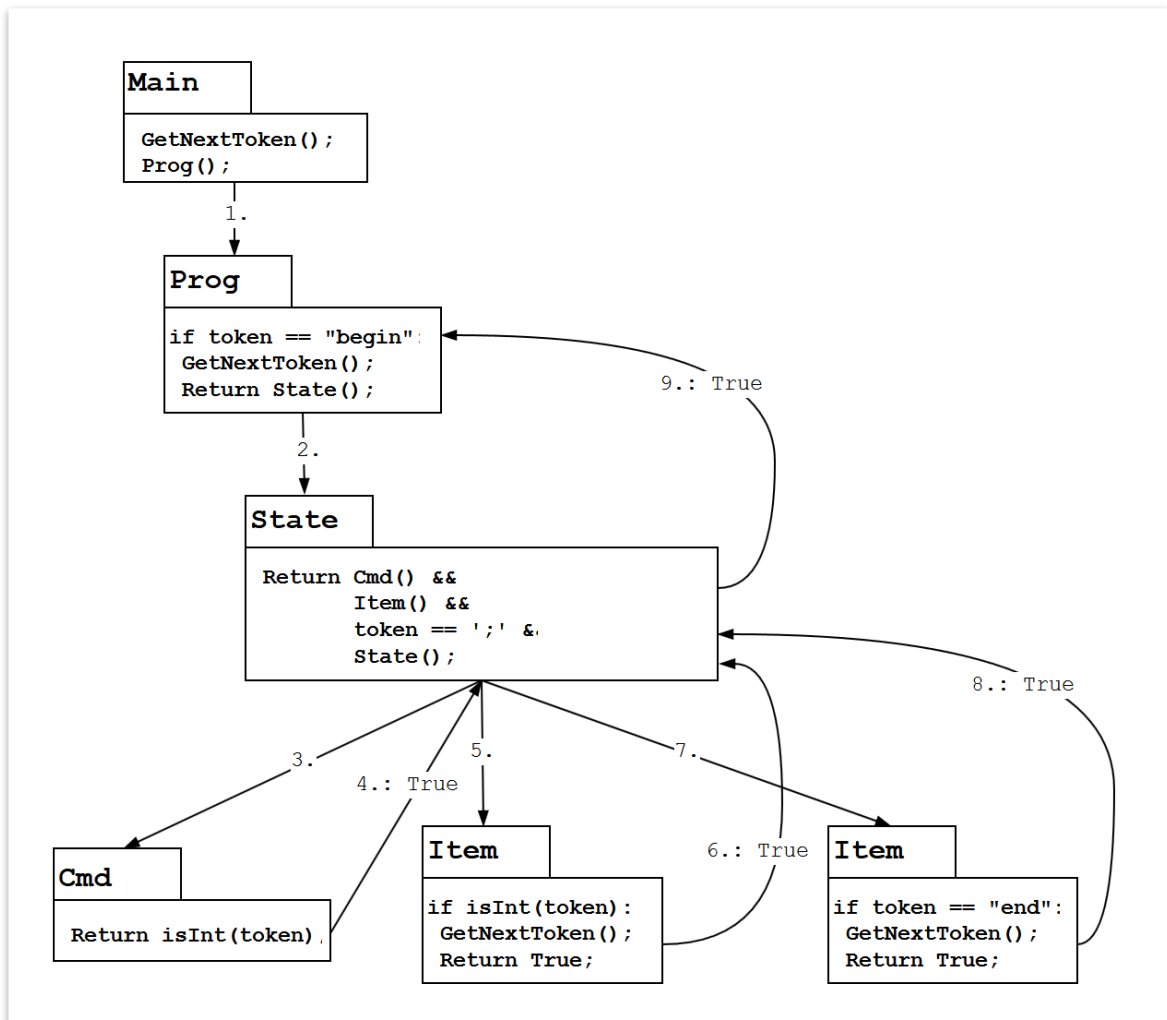
Obrázek 8: Derivační strom příkladu 3 s posloupností aplikovaných pravidel: 12673

6.3 Rekurzivní sestup

Rekurzivní sestup je jednoduchá metoda, která transformuje neterminály gramatiky do odpovídajících procedur, které obstarávají obsluhu jednotlivých akcí, a provádějí tak simulaci pravidel gramatiky. Z programátorského hlediska se jedná o velmi intuitivní přístup. Implementace rekurzivního sestupu zahrnuje velké množství podmínek a větvení pro kontrolu výskytu jednotlivých terminálů a podle toho se vzájemně rekurzivně volají jednotlivé procedury, které buď odhalí syntaktickou chybu a vrátí logickou hodnotu **false** nebo zavolají další procedury a postupným zanořováním kontrolují celý vstupní kód, dokud poslední metoda nevrátí **true**, resp. **false**, podle toho, zda byla syntaxe až do konce v pořádku. Výsledná hodnota se poté propaguje jednotlivými úrovněmi zanoření zpět do první metody, která odstartovala syntaktickou analýzu, a ta oznámí uživateli, zda je celý vstup v pořádku nebo se někde vyskytuje syntaktická chyba.

Jak již bylo uvedeno, jedná se o velmi jednoduchou a dobře pochopitelnou metodu, protože programátor implementuje jednotlivá pravidla gramatiky. Pro lepší kontrolu a přehlednost je vhodné opět vycházet z LL tabulky, protože při jejím sestavení se odhalí všechny možnosti a skryté kombinace, které do analýzy vnášejí epsilon pravidla, ale lze sestavit rekurzivní analyzátor i ze samotné gramatiky. Velkou nevýhodou tohoto přístupu je jednoúčelovost použití, protože implementace daného analyzátoru odpovídá konkrétní gramatice, a pokud by bylo třeba později upravit nebo úplně vyměnit gramatiku, musí se celý analyzátor přeprogramovat. Oproti tomu

prediktivní syntaktická analýza využívá při analýze LL tabulku a při změně gramatiky stačí pouze upravit LL tabulku a není potřeba provádět zásah do implementace. Z tohoto důvodu byla v práci upřednostněna implementace prediktivní syntaktické analýzy před rekurzivním sestupem, aby bylo možné využít jejich invariantních vlastností vůči gramatice, i když vždy vyžaduje složitou konstrukci LL tabulky.



Obrázek 9: Ilustrace volání metod rekurzivního sestupu pro vstupní řetězec „begin 42; end“

6.4 Prediktivní syntaktická analýza

Prediktivní syntaktický analyzátor je jeden z přístupů syntaktické analýzy shora dolů, který není oproti předchozí variantě rekurzivního sestupu tolik programátorsky intuitivní, ale na druhou stranu celou analýzu řídí jeden obecný algoritmus, který není implementačně závislý na dané gramatice, což přináší velkou výhodu, protože při vhodné implementaci je výsledný analyzátor znovupoužitelný pro různé gramatiky na rozdíl od rekurzivního sestupu.

Prediktivní syntaktický analyzátor se skládá ze svého jádra, které řídí syntaktickou analýzu a čte jednotlivé tokeny zpracované lexikální analýzou. Dále vyžaduje implementaci zásobníku a LL tabulky, která opovídá dané gramatice. Tyto dvě komponenty přidávají prediktivnímu analyzátoru na

složitosti a mohlo by se zdát, že je rekurzivní sestup jednodušší cestou, ale jak již bylo zmíněno, náročnější počáteční implementace se později vrátí ve flexibilitě a univerzálnosti celého syntaktického analyzátoru.

Syntaktická analýza probíhá shora dolů, takže výstupem analýzy je posloupnost pravidel pro vstupní řetězec, která je použita v nejlevější derivaci. Simuluje se tak sestavení derivačního stromu od kořene k listům.

Algoritmus prediktivní syntaktické analýzy očekává na vstupu LL tabulku pro zvolenou gramatiku $G = (N, T, P, S)$ a vstupní řetězec tokenů $x \in T^*$. Výstupem je potom levý rozbor pro x , pokud je x řetězcem jazyka $L(G)$, $x \in L(G)$, jinak syntaktický chyba.

1. Ulož symbol $\$$ na zásobník
2. Ulož symbol S na zásobník
3. Necht' X je vrchol zásobníku a a aktuální token
4. Pokud $X = \$$ a zároveň $a = \$$, potom **úspěch**; ukonči algoritmus
5. Pokud $X \in T$ a zároveň $X = a$, potom odeber X z vrcholu zásobníku a načti další a ze vstupního řetězce; vrať se na bod 3
6. Pokud $X \in N$ a zároveň $r: X \rightarrow x \in LL - \text{tabulka}[X, a]$, potom zaměň vrchol zásobníku X za řetězec symbolů x v opačném pořadí a zapiš r na výstup; vrať se na bod 3
7. **Syntaktická chyba**; ukonči algoritmus

6.5 LR syntaktický analyzátor

LR syntaktický analyzátor je jeden z nástrojů syntaktické analýzy, který umožňuje provádět syntaktickou analýzu zdola nahoru a simuluje tak konstrukci derivačního stromu od listů ke kořeni. Zkratka LR vznikla podle principu jeho fungování, protože čte řádky vstupního textu zleva doprava, shora dolů bez návratů (anglicky *Left to right*) a jeho výstupem je nejpravější derivace v převráceném pořadí (anglicky *Rightmost derivation in reverse*). Odtud tedy zkratka LR syntaktický analyzátor.

Při syntaktické analýze je nutné se typicky vypořádat se dvěma základními problémy, které již byly zmíněny, a to pravidla se stejnou pravou stranou, u kterých je problém deterministicky určit, které pravidlo se má aplikovat, a nejednoznačné gramatiky, které umožňují zkonstruovat více než jeden derivační strom pro stejný řetězec. Tyto problémy byly řešeny u analýzy shora dolů typicky převodem gramatiky na LL gramatiku, a poté již bylo možné použít rekurzivní sestup nebo prediktivní syntaktickou analýzu. LR syntaktický analyzátor používá ke své činnosti rozšířený zásobníkový automat a tabulku, která se skládá ze dvou částí, akční části a přechodová částí. Na výstupu pak vytváří pravý rozbor vstupního řetězce, což je reverzovaná posloupnost pravidel, která je použita v nejpravější derivaci. Samotnou analýzu obsluhuje podobně jako u prediktivní syntaktické analýzy jeden algoritmus, ale konstrukce rozšířeného zásobníku a především pak celé tabulky je velmi složitá, a proto je LR syntaktický analyzátor náročný na implementaci. Na druhou stranu má větší výpočetní sílu, než analyzátor shora dolů založený na LL gramatikách nebo jeho mnohem slabší alternativa pro analýzu zdola nahoru, precedenční syntaktický analyzátor.

6.6 Precedenční syntaktický analyzátor

Precedenční syntaktický analyzátor je nástroj pro syntaktickou analýzu zdola nahoru, který je řízen velmi jednoduchým algoritmem a rozhoduje se na základě tabulky, která definuje především vzájemnou precedenci a asociativitu operátorů používaných v gramatice. Odtud vznikl název precedenční tabulka. Výpočetní síla tohoto analyzátoru je nesrovnatelně malá s LR syntaktickým analyzátozem, ale na druhou stranu je velmi oblíbený, protože je snadný na implementaci a na rozdíl od analyzátorů shora dolů se dokáže velmi dobře vypořádat se zpracováním matematických výrazů. Gramatiku pro základní matematické operace lze popsat snadno několika pravidly, ale při analýze shora dolů dochází k problému nejednoznačnosti a více derivačním stromům, z nichž správným by byl pouze ten, který dodržuje precedenci operátorů, a právě s tímto se precedenční analyzátor snadno vypořádá. To je jeden z důvodů, proč práce v praktické části spojuje univerzální prediktivní syntaktický analyzátor shora dolů, který je složitý na konstrukci, ale ne tolik jako LR analyzátor, pro analýzu syntaktické struktury kódu spolu s precedenčním analyzátozem zdola nahoru pro zpracování výrazů. Tento přístup naplňuje jednu z myšlenek gramatických systémů, kdy více přispěvatelů s částečným řešením dokáže společně vyřešit celý problém, aniž by se muselo konstruovat jedno souhrnné složité řešení, což by byl v tomto případě LR analyzátor.

Jedním z důvodů, proč je síla precedenčního analyzátoru tak malá, jsou zavedená omezení, která požadují, aby v gramatice neexistovalo více pravidel se stejnou pravou stranou a zároveň gramatika neobsahovala epsilon pravidla.

Algoritmus precedenčního analyzátoru požaduje na vstupu precedenční tabulku pro gramatiku $G = (N, T, P, S)$ a vstupní řetězec tokenů $x \in T^*$. Výstupem je potom pravý rozbor pro x , pokud je x řetězcem jazyka $L(G)$, $x \in L(G)$, jinak syntaktický chyba.

1. Ulož symbol \$ na zásobník
2. Necht' a je terminál nejbližší vrcholu zásobníku a b je aktuální token na vstupu
3. Pokud $a = \$$ a zároveň $b = \$$, potom **úspěch**; ukonči algoritmus
4. Pokud $Tabulka[a, b] = '='$, potom ulož b na zásobník; načti další symbol b ze vstupu; vrať se na bod 2
5. Pokud $Tabulka[a, b] = '<'$, potom zaměň a za ' $a <$ ' na vrcholu zásobníku; ulož b na zásobník; načti další symbol b ze vstupu; vrať se na bod 2
6. Pokud $Tabulka[a, b] = '>'$ a zároveň ' $< y$ ' se nachází na vrcholu zásobníku a zároveň existuje pravidlo $r: A \rightarrow y \in P$, potom zaměň ' $< y$ ' za A na zásobníku; vypiš pravidlo r na výstup; vrať se na bod 2
7. **Syntaktická chyba**; ukonči algoritmus

Samotná precedenční tabulka je tvořena maticí průniku všech terminálů. Řádky tabulky reprezentují terminály na vrcholu zásobníku ($a_i..a_n$ pro $i = 1..n$) a záhlaví sloupců představuje terminály na vstupu ($a_j..a_n$ pro $j = 1..n$). Jejich vzájemný průnik (jednotlivé buňky) může nabývat hodnot $\{<, =, >, nic\}$. Místo teoretických ukázek bude popsána skutečná precedenční tabulka na gramatice, která byla využita při implementaci praktické části práce. Tato gramatika/tabulka umožňuje zpracovat matematické výrazy obsahující základní operace (+, -, *, /) a všechny porovnávací operátory (<, >, <=, >=, ==, !=). Dále pak dovoluje zkontrolovat libovolné zanoření kulatých závorek a symbol i reprezentuje identifikátor, kterým může být číselná hodnota nebo proměnná.

Nechť $G = (N, T, P, E)$ je gramatika pro precedenční syntaktický analyzátor použitý v praktické části práce, kde:

$$N = \{E\},$$

$$T = \{+, -, *, /, <, >, <=, >=, !=, (,), i, \$\},$$

$$P = \{ \begin{array}{ll} 1: E \rightarrow E + E, & 7: E \rightarrow E <= E, \\ 2: E \rightarrow E - E, & 8: E \rightarrow E >= E, \\ 3: E \rightarrow E * E, & 9: E \rightarrow E == E, \\ 4: E \rightarrow E / E, & 10: E \rightarrow E != E, \\ 5: E \rightarrow E < E, & 11: E \rightarrow (E), \\ 6: E \rightarrow E > E, & 12: E \rightarrow i \} \end{array}$$

Symbole na vstupu

	+	-	*	/	<	>	<=	=>	==	!=	()	i	\$
+	>	>	<	<	>	>	>	>	>	>	<	>	<	>
-	>	>	<	<	>	>	>	>	>	>	<	>	<	>
*	>	>	>	>	>	>	>	>	>	>	<	>	<	>
/	>	>	>	>	>	>	>	>	>	>	<	>	<	>
<	<	<	<	<	>	>	>	>	>	>	<	>	<	>
>	<	<	<	<	>	>	>	>	>	>	<	>	<	>
<=	<	<	<	<	>	>	>	>	>	>	<	>	<	>
=>	<	<	<	<	>	>	>	>	>	>	<	>	<	>
==	<	<	<	<	<	<	<	<	>	>	<	>	<	>
!=	<	<	<	<	<	<	<	<	>	>	<	>	<	>
(<	<	<	<	<	<	<	<	<	<	<	=	<	
)	>	>	>	>	>	>	>	>	>	>		>		>
i	>	>	>	>	>	>	>	>	>	>		>		>
\$	<	<	<	<	<	<	<	<	<	<	<		<	

Tabulka 4: Tabulka symbolů precedenční analýzy gramatiky G

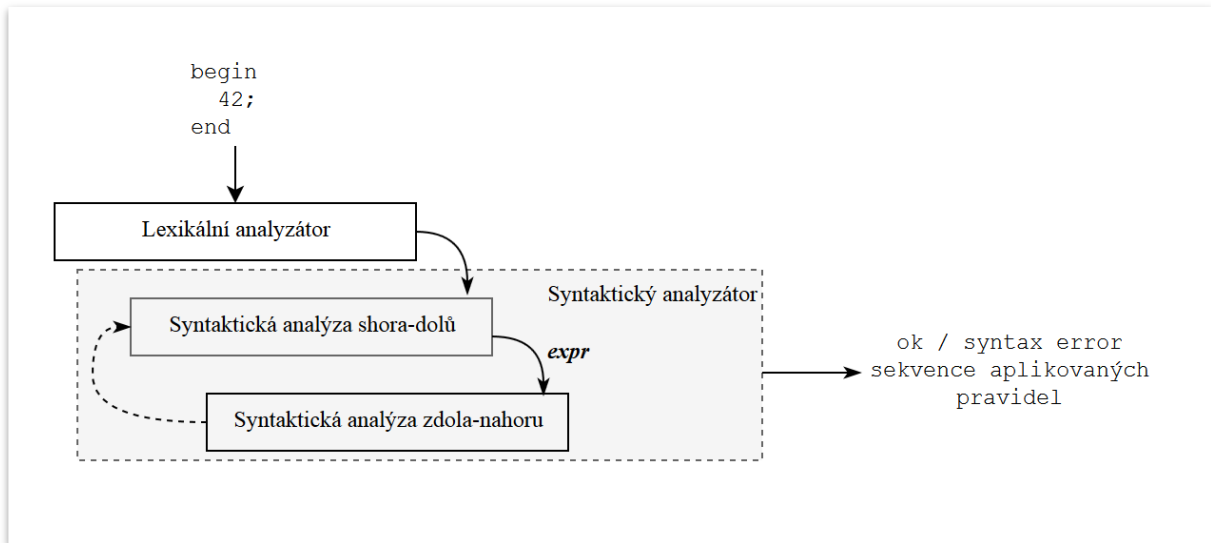
Precedenční tabulka zachycuje vztahy jednotlivých operátorů (a dalších symbolů) mezi sebou, aby je mohl analyzátor správně vyhodnotit a zpracovat v pořadí určeném jejich prioritou. Průnik operátorů se stejnou prioritou se řídí dle jejich asociativity. Levě asociativní operátory (v základní matematice to jsou všechny) mají hodnotu ' $>$ ', naopak pravě asociativní by měly opačnou šipku, ' $<$ '. Ostatní šipky porovnávají prioritu symbolů mezi sebou, přičemž se vždy porovnává priorita symbolu na

zásobníku vůči prioritě symbolu na vstupu, takže např. $' + < * '$ a naopak $' / > - '$. Prázdná buňka určuje kombinaci, kdy se symboly nemohou potkat vedle sebe, jako například dvě kulaté závorky $') ('$ nebo identifikátor a závorka $' i ('$. Jediný případ, kde se vyskytuje symbol rovnítka $' = '$, je při průniku dvou párových závorek $' () '$, které jsou následně při analýze zredukovány a hlídá se tak jejich správný počet.

7 Aplikace

Praktická část práce se zaměřuje na aplikování znalostí nabytých v teoretické části práce takovým způsobem, aby byly tyto vědomosti nejen přeneseny do praxe, ale zároveň napomohly uživateli k lepšímu pochopení celé problematiky gramatických systémů, překladačů a syntaktické analýzy. Vychází se z běžného problému syntaktické analýzy programovacích jazyků, které lze popsat složitou bezkontextovou gramatikou pro LR syntaktický analyzátor. Tento přístup je vhodný, pokud se celé gramatice věnuje velká pozornost a zároveň jsou k dispozici komerční nástroje, které dokáží výslednou komplexní gramatiku transformovat do funkčního analyzátoru, a udělají tak práci za nás. Pro profesionální použití je to velmi vhodná cesta, ale vyžaduje více zkušeností a znalostí. Druhý přístup nabízí použití více gramatik. Pomocí relativně jednoduché LL gramatiky lze popsat syntaktickou strukturu daného jazyka a následně s použitím LL tabulky vytvořit buďto jednoúčelový analyzátor pomocí rekurzivního sestupu nebo prediktivní syntaktický analyzátor, který pracuje nad LL tabulkou. Tato varianta je relativně snadná a dobře uchopitelná jak na pochopení, tak i na implementaci. Zůstává zde však problém s vyhodnocením matematických výrazů (často označovány jako *expressions*), protože syntaktická analýza shora dolů se špatně vypořádává s výrazy, ve kterých je nutné zohledňovat správné pořadí vyhodnocení operátorů na základě jejich priority (precedence) a dalšími záležitostmi, jako je např. libovolný počet kulatých závorek ve výrazech atd. Pro jejich zpracování je právě vhodnější přístup zdola nahoru, kterým pracuje LR syntaktický analyzátor, ale jak již bylo řečeno, vytvářet pouze kvůli výrazům jednu komplexní gramatiku, navíc s nutností implementovat složitý LR analyzátor, je pro nekomerční (příp. didaktické) účely zbytečné. Nabízí se tedy možnost nechat LL gramatiku pro popis syntaktické struktury programu a pro vyhodnocení výrazů navrhnout druhou gramatiku. Zde přichází opět v úvahu využít LR analyzátor, jehož gramatika by byla jednoduchá, protože by popisovala pouze výrazy, ale opět by byla nutná implementace celého analyzátoru, což by situaci vrátilo zpět na začátek k tomu, že je možné celý jazyk popsat gramatikou pro LR analyzátor. Druhá varianta je použít nejednoznačnou bezkontextovou gramatiku popisující pouze matematické výrazy, práci s operátory, čísly a příp. identifikátory, jejíž nejednoznačnost by za normálních okolností vadila, ale protože je použit precedenční syntaktický analyzátor, který se s tímto problémem dokáže vyrovnat, je vše v pořádku. Výsledná gramatika tak bude opravdu jednoduchá, velmi čitelná a přitom bez problému dokáže popsat požadovanou strukturu výrazů. Precedenční analyzátor je velmi jednoduchý nástroj pro syntaktickou analýzu zdola nahoru, který by někdo mohl označit za triviální a nepoužitelný. Pro složité gramatiky nebo pro celou syntaktickou analýzu by to byla pravda, ale jako jednoúčelový nástroj pro zpracování, v tomto případě matematických výrazů, je velmi vhodný. Zároveň je snadný na implementaci a z didaktických účelů donutí uživatele více přemýšlet nad vztahy mezi operátory, jejich prioritou, možnostech výskytu atd.

Výsledné dva protichůdné přístupy se vzájemně výborně doplňují. Samy o sobě jsou jednoduché a dohromady umožňují provést syntaktickou analýzu vstupního kódu požadovaného jazyka.



Obrázek 10: Schéma syntaktické analýzy aplikace

Jediným problémem tak zůstává, jak tyto dvě metody spojit, aby z nich vzniklo jedno fungující řešení. V praxi se zpravidla aplikuje přístup syntaxi řízeného překladu, kdy je syntaktický analyzátor hlavním jádrem překladače a řídí činnost ostatních komponent, jako je lexikální a sémantický analyzátor, generátor vnitřního kódu a další. Při zkoumání podobných přístupů v praxi by se pravděpodobně ukázalo, že hlavním řídicím prvkem je analýza shora dolů, která si nějakým způsobem volá metody precedenčního analyzátoru, a vše je implementované pouze za účelem funkčnosti, ale nikdo už nepřemýšlel nad tím, jestli lze toto propojení formálně specifikovat a blíže zkoumat. I zde je aplikován stejný model, protože dává smysl, aby analýzu řídil syntaktický analyzátor shora dolů. Hlídá a kontroluje strukturu celého kódu, takže ví, kde se mohou vyskytnout výrazy. Protože práce nezkoumá celý problém pouze po praktické stránce, ale především i po té teoretické, je zaveden v těchto místech do gramatiky pomocný pseudoterminál, nazvaný *expr* (od anglického *expression*). Lze si povšimnout, že spojení těchto dvou gramatik pomocným pseudoterminálem vede přímo ke gramatickým systémům. Možností je jistě více, ale nabízí se zvolení PC gramatického systému právě kvůli pomocnému pseudoterminálu *expr*, který nemá přímý význam ve struktuře jazyka, ale může posloužit právě jako komunikační symbol mezi dvěma komponentami PC gramatického systému. Tento pseudoterminál se může vyskytnout pouze v první komponentě, kterou je LL gramatika pro syntaktickou analýzu shora dolů. V gramatice precedenční analýzy se vyskytnout nemůže, protože by nedávalo příliš smysl, aby se při analýze výrazu znovu volala analýza výrazu. Tato analýza buďto skončí nebo se jedná o jeden složitější výraz, ale vždy se zpracuje najednou. Je tedy patrné, že celé řešení nejlépe pasuje na centralizovaný PC gramatický systém s návraty. Pouze první řídicí komponenta generuje komunikační symboly, pomocí kterých si vyžádá práci druhé komponenty, tedy precedenční analýzy, která se po vrácení výsledku resetuje do výchozího stavu (neukládá si rozpracované řešení) a opět čeká, až bude znovu zavolána.

Gramatika pro precedenční analyzátor byla podrobně představena již v dřívější kapitole 6.6 o precedenční analýze, takže zde již nebude znovu uvedena. V aplikaci byla použita přesně tato konkrétní gramatika bez dalších úprav. Naopak řídicí část syntaktické analýzy shora dolů je velmi

zajímavá, protože díky použití prediktivního syntaktického analyzátoru, který požaduje na vstupu již sestavenou LL tabulku, není tento analyzátor vázán na konkrétní gramatiku.

V rámci této praktické části byla vytvořena webová aplikace, která umožňuje zadat libovolný vstupní text, nad nímž provede syntaktickou analýzu a zobrazí uživateli posloupnost aplikovaných pravidel gramatiky. Pokud je zdrojový kód v pořádku, oznámí, že je vše *ok*. V opačném případě se výpis aplikovaných pravidel zastaví na výskytu první chyby a tato chyba je interpretována uživateli. Většinou se jedná o informaci, na kterém řádku a jakém tokenu došlo k chybě. Taková aplikace by sama o sobě byla velmi vhodnou didaktickou pomůckou, protože nevyžaduje žádnou instalaci nebo složité spouštění přes konzoli apod., jak tomu často bývá u aplikací v odborných pracích. Zároveň se uživatel dozví, jak byla pravidla postupně aplikována, a to pro obě gramatiky. Sekvence pravidel obsahuje čísla pravidel LL gramatiky při kontrole shora dolů, ale i výpis jednotlivých pravidel precedenční analýzy při přepnutí komponenty pro zpracování výrazů. Aplikace však zachází ještě dál a díky vhodné implementaci umožňuje vložit libovolnou LL gramatiku vč. epsilon pravidel. V jednotlivých kapitolách byla velmi podrobně vysvětlena konstrukce LL tabulky spolu s jednotlivými pomocnými množinami. Všechny tyto kroky jsou přeneseny do aplikace a umožňují dynamicky vygenerovat celou LL tabulku ze zadané gramatiky. Tato tabulka je potom předána syntaktickému analyzátoru při spuštění analýzy. Protože vytvoření LL tabulky a následná syntaktická analýza probíhá automaticky ze zadané LL gramatiky, je nutné klasickou definici bezkontextové gramatiky rozšířit o další atributy. Definice vstupní LL gramatiky G_{LL} je následující:

$$G = (N, T, P, S, Expr, \varepsilon, Int, Float, String, Id)$$

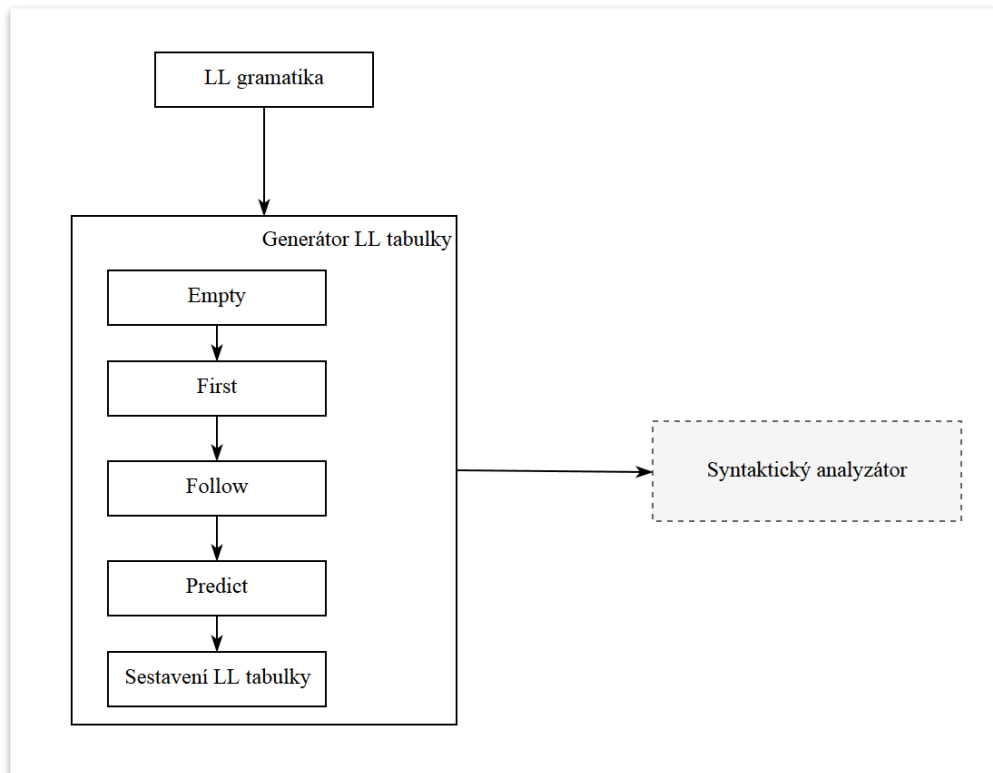
- N je konečná neprázdná množina neterminálních symbolů
- T je konečná neprázdná množina terminálních symbolů a platí, že $N \cap T = \emptyset$
- P je konečná množina pravidel tvaru $A \rightarrow x$, kde $A \in N, x \in (N \cup T \cup Expr)^+$ (epsilon pravidlo se vytvoří výběrem zástupného epsilon terminálu ε)
- $S \in N$ je počáteční (startující) neterminál

Následující atributy jsou volitelné a rozšiřují definici gramatiky G_{LL} :

- $Expr$ je pseudoterminál, který slouží pro přepnutí řízení syntaktické analýzy z analýzy shora dolů na analýzu zdola nahoru
- $\varepsilon \in T$ je vybraný terminál reprezentující epsilon pro vytvoření epsilon pravidla
- $Int \in T$ je vybraný terminál reprezentující celočíselný datový typ integer
- $Float \in T$ je vybraný terminál reprezentující číselný datový typ s plovoucí řádovou čárkou float
- $String \in T$ je vybraný terminál reprezentující datový typ řetězec
- $Id \in T$ je vybraný terminál reprezentující identifikátor

Rozšíření klasické bezkontextové gramatiky tímto způsobem bylo nezbytné, protože aplikace nabízí velkou flexibilitu, ale pro praktické použití je nutné specifikovat základní datové typy a další elementy jazyka, aby bylo možné provádět syntaktickou analýzu nad běžnými programovacími jazyky a dalšími podobnými konstrukcemi.

Následující schéma zachycuje jednotlivé bloky konstrukce LL tabulky tak, jak je LL gramatika postupně analyzována a výsledná LL tabulka je předána syntaktickému analyzátoru pro řízení analýzy shora dolů.



Obrázek 11: Schéma konstrukce LL tabulky

Obě gramatiky byly formálně zavedeny, přičemž LL gramatiku pro syntaktickou analýzu shora dolů lze označit G_{LL} a gramatiku pro precedenční analýzu zdola nahoru z předchozí kapitoly 6.6 G_{PR} . Pro již zmíněný centralizovaný PC gramatický systém s návraty bude jeho formální definice následující:

$$\Gamma = (N, K, T, (S_{LL}, P_{LL}), (E_{PR}, P_{PR})), \text{ kde:}$$

- $N = N_{LL} \cup N_{PR}$
- $K = \{Expr_{LL}\}$
- $T = T_{LL} \cup T_{PR}$
- pro všechny $A \rightarrow x \in P_{PR}$ platí, že $abeceda(x) \cap K = \emptyset$

Dynamická LL gramatika zvyšuje přidanou hodnotu aplikace, protože uživatel má možnost experimentovat. Může si vložit vlastní gramatiku a na základě výstupů ji libovolně modifikovat a upravovat v reálném čase. Vždy před další syntaktickou analýzou se znovu vygeneruje nová tabulka, aby zůstávala stále aktuální k zadané gramatice. Aplikace dále zobrazuje uživateli sestavenou LL tabulku, takže se díky aplikaci dozví nejen, jakým způsobem byl jeho vstupní kód na základě vlastní gramatiky analyzován, ale i jak vypadá příslušná LL tabulka. Podobný princip by bylo možné aplikovat i u precedenčního analyzátoru, ale pro lepší použití aplikace bylo cílem, aby měla alespoň nějaký pevný bod, a vyhodnocování matematických výrazů je vhodnou volbou, protože se pro běžné programovací jazyky neliší. Stačí si tedy načíst webovou stránku, kde je aplikace umístěná, vložit vlastní LL gramatiku spolu se vstupním kódem a již je možné sledovat vygenerovanou LL tabulku a výsledek analýzy.

7.1 Uživatelské rozhraní

Hlavním záměrem celé aplikace je jednoduchost, intuitivní použití a co největší přínos pro uživatele. Tyto aspekty se odráží již na první pohled v samotném designu, který vychází z čisté bílé barvy, na které dobře vyniknou důležité prvky. Po načtení se zobrazí jedna stránka (tento přístup je nazýván *single page* aplikace), která je vymezena záhlavím s názvem aplikace „Dynamic Parser“ a v pravém horním rohu se nachází informační tlačítko, které po stisknutí zobrazí dialogové okno s obecnými informacemi o aplikaci a precedenční tabulkou. Celá aplikace je v angličtině s výjimkou těchto obecných informací. Zápatí stránky je pak ohraničeno šedým pruhem obsahujícím kontakt a autora stránky.

Tělo stránky obsahuje vše ostatní. V pravém horním rohu se nachází tlačítka umožňující načíst vzorové příklady gramatiky a vstupního kódu. Pod mini je stránka rozdělena do dvou sloupců v poměru 2/3 a 1/3. Levá větší část obsahuje komponenty pro vložení gramatiky. Uživatel může definovat název gramatiky a následně vložit první pravidlo, které má samostatně vstup pro levou stranu a pravou stranu. Pod tímto řádkem se nachází tlačítko pro přidání dalších řádků. Jednotlivé řádky pravidel lze poté mazat tlačítkem **x**. Vstup pro vkládání pravidel není obyčejný textový HTML input, ale speciální komponenta `Tags-input`³, která umožňuje vkládat text formou vizuálně oddělených tagů (štítků), které v tomto případě přímo korespondují s jednotlivými terminály a neterminály. Pokud je vložen symbol do levé strany pravidla, lze vložit pouze jeden, automaticky se obarví modře a pokládá se za neterminál. Naopak v pravé straně pravidla je po vložení tagu uživatel vyzván k výběru, jestli se jedná o terminál, neterminál nebo expression (pseudoterminál pro komunikaci mezi gramatikami) reprezentující matematický výraz. Zároveň tato komponenta podle aktuálně napsaného textu zobrazuje nabídku již existujících tagů, takže v okamžiku, kdy je již napsaná část gramatiky, je její další doplňování o to rychlejší, protože opakující se tagy není nutné znovu psát, ale po zadání prvních písmen jsou automaticky nabídnuty ke vložení. Pravý sloupec obsahuje více předdefinovaných `Tags-input` komponent, které zobrazují výčty jednotlivých množin. Nachází se zde množina terminálů, neterminálů, expression pseudoterminál, aliasy pro epsilon, celá čísla, reálná čísla, identifikátory a další, tedy vše, co je nutné pro definování celé gramatiky, přičemž není potřeba vyplňovat všechna pole, ale záleží na komplexnosti pravidel. Jsou-li použity pouze např. proměnné, mohou zůstat vstupy definující terminály pro celá čísla, reálná čísla a řetězce prázdné. Tyto dva sloupce jsou úmyslně vedle sebe, protože na větších monitorech se tak zobrazí celá gramatika a vedle ní její konfiguraci. Pokud se obrazovka zmenší nebo se aplikace zobrazí např. na telefonu, přeskládají se sloupce díky responzivitě pod sebe a opět bude manipulace s aplikací pohodlnější.

Pod tímto blokem pro gramatiku se nachází sekce s LL tabulkou, kde je ve výchozím stavu umístěno pouze tlačítko pro generování samotné LL tabulky, ta je poté na tomto místě zobrazena. Dále navazuje poslední sekce, opět rozdělena na dva sloupce, ale tentokrát půl na půl. V levém sloupci je připraveno pole pro vstupní kód, který bude analyzován, spolu s tlačítkem, které analýzu vyvolá. Není nutné předem ručně generovat LL tabulku, při stisknutí tlačítka pro analýzu se vždy vygeneruje znovu. V pravém sloupci se po provedení analýzy zobrazí vstupní kód po jednotlivých očíslovaných řádcích a rozdělený na tokeny. Komentáře jsou obarveny zeleně a při analýze se přeskakují. Poté následuje výsledek analýzy (symbol ✓ pro úspěch nebo chybová zpráva) a posloupnost aplikovaných pravidel. Oba sloupce se opět při menším rozlišení obrazovky posunou pod sebe, aby nezabíraly tolik místa. Na velkém monitoru je pak vše přehledně vedle sebe, vstup i odpovídající výstup. Konkrétní ukázkou aplikace lze nalézt na konci práce v sekci příloh.

³ Komponenta pro automatické vytváření štítků. <http://mbenford.github.io/ngTagsInput/>

7.2 Implementace

Aplikace je vytvořena jako jednoduchá internetová stránka, tzv. single page aplikace (vše se nachází a funguje v rámci této jedné stránky), a implementuje celou řadu moderních webových technologií a přístupů. Samotná kostra je napsaná v HTML5⁴ a pro styly jsou použity jak CSS3 soubory, tak i modernější LESS⁵ soubory spolu s frameworkem Bootstrap⁶. Aby se aplikace dobře spravovala a vyvíjela, byla vytvořena jako projekt pro MS Visual Studio⁷ 2013 (a novější). V tomto vývojovém prostředí se nejen dobře orientuje v adresářové struktuře projektu, ale umožňuje konfiguraci dalších funkcionalit, jako je např. automatický publish profil pro snadnější nasazení aplikace.

Hlavní a jedinou stránkou aplikace je soubor index.html. Celá aplikační logika je implementovaná pouze pomocí JavaScriptu, takže lze říct, že se jedná o čistě klientskou aplikaci. Poté, co se stránka načte v prohlížeči, nevyžaduje internetové připojení a nikam se nepřipojuje, vše funguje lokálně v prostředí prohlížeče. Nejedná se však o čistý JavaScript, ale o moderní framework AngularJS⁸, který je složitější na pochopení, ale nabízí mnoho výhod, jako je například obousměrné vázání proměnných do webové stránky (tzv. two-way binding). Tyto výhody mají obecně vliv na výkon JavaScriptu, ale jeho degradace by byla citelná např. u informačního systému. U této aplikace, kde se neklade důraz na výkon a rychlost zpracování, to nijak omezuje, a lze tak využít všech výhod, které nabízí. Kolem AngularJS vznikla vývojářská komunita, takže existuje mnoho volně dostupných aplikací třetích stran, mezi které patří např. Tags-input použitý pro uživatelsky velmi přívětivé definování gramatiky pomocí štítků. Aby bylo možné jednotlivé moduly pro AngularJS lépe spravovat, využívá projekt balíčkovací systém NPM⁹ a pro automatický build JavaScriptu a stylů konzolový program gulp¹⁰. Při vývoji aplikace potom stačí spustit v konzoli příkaz gulp, který při jakémkoliv změně JavaScriptu nebo stylů znovu provede build aplikace, takže je celý klient během vývoje neustále aktuální. Před nasazením na produkci potom stačí spustit gulp s parametrem „-p“ a JavaScript se minifikuje do jednoho souboru app.js. Všechny soubory se styly (CSS i LESS) se zpracují podobným způsobem a uloží v souboru app.css. Styly využívají framework Bootstrap, jak již bylo zmíněno, a ikony jsou převzaty z Font Awesome¹¹.

Struktura aplikace obsahuje na nejvyšší úrovni soubor index.html a dva adresáře WebClient a Views. Ostatní soubory mají význam pouze pro Visual Studio a není nutné je zde rozebírat. Složka Views obsahuje html soubory, které představují dialogová okna zobrazující se v aplikaci. WebClient reprezentuje celého JavaScript klienta včetně všech stylů, modulů atd. Uvnitř této složky se nachází soubor package.json a dvě související složky node_modules a webpack, které obsahují moduly a konfigurace pro AngularJS a NPM. Dále jsou zde dva adresáře. První složka se jmenuje build a obsahuje výsledné minifikované soubory a fonty. A konečně složka src, ve které se nachází jednotlivé JavaScript soubory a styly. Základním souborem je main.js, který spojuje soubory dohromady. Samotná logika aplikace je potom rozdělena do kontrolérů umístěných ve složce controllers a pomocné funkce a konstanty jsou ve složce helpers. Veškeré styly použité v aplikaci se nachází ve složce style. Jsou zde uloženy soubory pro komponenty Bootstrapu, dialogová okna z modulu

⁴ Nový standard HTML. <https://www.w3.org/TR/html5/>

⁵ CSS soubory rozšířené o dynamické chování jako jsou proměnné, operace, funkce <http://lesscss.org/>

⁶ HTML, CSS, and JS framework pro responzivní design. <http://getbootstrap.com/>

⁷ Integrované vývojové prostředí firmy Microsoft. <https://www.visualstudio.com/cs/>

⁸ Moderní framework pro JavaScript. <https://angularjs.org/>

⁹ Balíčkovací systém pro JavaScript. <https://www.npmjs.com/>

¹⁰ Nástroj umožňující automatizaci různých úkonů pro vývoj (např. build)

<https://www.npmjs.com/package/gulp>

¹¹ Balík vektorových ikon. <http://fontawesome.io/>

ngDialog, Font Awesome a vlastní třídy pro samotnou aplikaci. Celý systém stylů a knihoven je udělán tímto způsobem, aby nebylo nutné nic stahovat z externích zdrojů z internetu, ale vše bylo při vývoji k dispozici offline a u výsledné aplikace vystavené na internet stačilo, aby si načetla své vlastní lokální zdroje.

Implementace výkonné logiky aplikace je uložena v jednotlivých kontrolérech. Na nejvyšší úrovni se nachází system-controller.js, jehož metody a proměnné jsou přístupné ostatním kontrolérům. Jsou zde umístěny metody pro základní práci s jednotlivými množinami, načítání terminálů, neterminálů atd. Zbývá funkcionalita je rozdělena na tři části. Celá lexikální analýza se nachází v souboru lexical-analysis.controllers.js. Sekvence metod pro vytvoření LL tabulky je v LLtable-controller.js a závěrečná syntaktická analýza je implementována v syntax-analysis-controller.js. Není nutné rozebírat implementaci a fungování jednotlivých metod, protože přesně kopírují funkcionalitu odpovídajících algoritmů, které byly podrobně probrány v jednotlivých kapitolách, a pouze ji upravují pro JavaScript. Za zmínku stojí, že lexikální analýza se neprovádí, jak tomu často bývá, průběžně, kdy syntaktická analýza řídí překlad a dotazuje se dle potřeby vždy na další token, ale provede se na začátku pro celý vstupní kód. Je to z toho důvodu, aby byl vstupní kód již ve formě tokenů celý k dispozici při dalších analýzách a bylo jej možné zobrazit v této podobě na výstupu. Po lexikální analýze se provede sestavení LL tabulky a poté již proběhne najednou celá syntaktická analýza.

7.3 Nasazení a vývoj

Jedná o čistě klientskou aplikaci s výchozím souborem index.html, takže pro nasazení stačí překopírovat obsah složky WebSite na server a vše bude fungovat automaticky samo. Toto platí pro případ, že je k dispozici publikovaná verze aplikace. V projektu pro MS Visual Studio je pro ulehčení práce vytvořený odpovídající publish profil.

Zprovoznění aplikace pro vývoj na lokálním stroji není úplně jednoduché vzhledem k velkému množství použitých technologií, doplňků atd. Návod na instalaci celého vývojového prostředí je uveden v jednotlivých na sebe navazujících bodech, které by mělo stačit následovat. Vývojové prostředí se skládá ze dvou částí, první z nich je MS Visual Studio 2013/15, které ale není nezbytně nutné. Celý projekt je vytvořený pro toto prostředí, ale jednotlivé soubory lze editovat v libovolném textovém editoru. Důležitější a složitější je instalace webového klienta (návod je pro MS Windows).

1. Nainstalovat Node.js (<https://nodejs.org/dist/v0.12.7/>)
2. Přidat do systémových proměnných proměnnou GYP_MSVS_VERSION s hodnotou „2013“ (verze MS Visual Studia, některé použité balíčky tuto hodnotu potřebují)
3. Nainstalovat Python (verzi 2.7.10)
4. Přidat Python do Path (např. „C:\Python27\;C:\Python27\Scripts;“)
5. Pokud používáte Git, spustit git bash v adresáři WebClient, jinak klasicky cmd
6. Spustit příkaz: „npm install gulp –g“
7. Spustit příkaz: „npm install“

Příkazy pro vývoj:

„gulp“ Automatický build a tzv. livereload, který hlídá změny ve zdrojových souborech a automaticky provede build při každé změně.

„gulp –p“ Provede build JavaScriptu a stylů pro produkční prostředí, vytvoří tedy minifikované verze.

8 Závěr

Práce postupně obsáhla několik různých samostatných témat, které se ve výsledku spojily do jednoho funkčního celku a skvěle se vzájemně doplňují. Nejdříve se práce věnuje základu formálních jazyků, důležitým pojmům a formalismům. Tato kapitola, stejně jako ostatní, by jistě dokázala sama obsáhnout celou práci, ale bylo důležité zmínit alespoň to nejzákladnější pro lepší zasazení čtenáře do kontextu a sjednocení jeho znalostí s terminologií zbytku práce. Byly zde zmíněny tedy takové pojmy jako abeceda, řetězec, Chomského hierarchie a další. Po tomto základu mohla bez problému navázat kapitola o regulárních gramatikách. Význam této kapitoly je v kontextu celé práce trochu skrytý a někdo by mohl namítnout, že vložit tuto kapitolu do práce bylo zbytečné. Nicméně základy práce s řetězci a jazyky se velmi dobře vysvětlují nad regulárními jazyky, protože tvoří základní stavební kameny složitějších typů jazyků a úzce s nimi souvisejí konečné automaty, které jsou jedním z formálních prostředků regulárních jazyků. Zároveň byl konečný automat použit pro implementaci lexikální analýzy v praktické části práce, která není z pohledu práce významná, nicméně pro funkčnost aplikace je neopomenutelná. Následující kapitola navázala vyšší třídou jazyků Chomského hierarchie, a to bezkontextovými jazyky, resp. gramatikami, se kterými se poté pracuje ve zbytku celé práce. Bezkontextové gramatiky jsou velmi významné a dovolují popsat mimo jiné většinu programovacích jazyků, avšak vyskytuje se zde několik komplikací, z nichž jedna se týká nejednoznačnosti. K této situaci dochází, pokud existuje více pravidel gramatiky se stejnou levou stranou pravidla. I tento problém má své řešení, jako je např. zavedení determinismu nebo speciální LL gramatiky, které se využívají při syntaktické analýze shora dolů. S bezkontextovými gramatikami souvisí zásobníkové automaty, které jsou jedním z jejich formálních aparátů. Jedná se o konečné automaty rozšířené o zásobník a opět i tyto prostředky byly následně aplikovány v praktické části.

Po bezkontextových gramatikách přichází na řadu jádro práce, bezkontextové gramatické systémy. Tyto systémy vznikly v průběhu druhé poloviny 20. století a umožňují formálně popsat a definovat spolupráci více gramatik při řešení společného problému. V mnoha ohledech jsou podobné principy již dávno v praxi uplatňovány, ale nemají zavedené žádné formalismy, protože je hlavním cílem funkčnost a nedává se prostor teorii. Jenomže bez teoretického základu lze špatně provádět další výzkum a rozvoj takovýchto systémů, obzvláště pokud jde např. o validaci a verifikaci. Jedním z nejčastějších příkladů spolupráce více gramatik je syntaktická analýza. Lze provést analýzu vstupního kódu pomocí jedné komplexní gramatiky, ale to vyžaduje definici složité gramatiky a zároveň i komplexního nástroje, který ji dokáže zpracovat. V mnoha ohledech se jeví jako vhodnější přístup mít více gramatik, z nichž každá dokáže řešit část vstupního textu (jazyka) a společně zvládnou provést kompletní analýzu. Jako jeden konkrétní příklad za všechny je možné zmínit kombinaci syntaktické analýzy shora dolů pro kontrolu syntaktické stavby jazyka s analýzou zdola nahoru pro zpracování a vyhodnocení matematických výrazů, kde se klade velký důraz na správné pořadí zpracování matematických operací dle priority operátorů. Přesně pro tuto problematiku jsou gramatické systémy určeny, protože jejich základní myšlenka spočívá v tom, že úspěchu nemusí dosáhnout jeden složitý systémem, ale kombinace a spojení dílčích řešení. Vycházejí tak z velmi inspirativních modelů reálného světa, jako je například tzv. „problém tabule“ (anglicky *The Black Board Model*), ve kterém skupina lidí řeší společný problém u tabule a různě se střídají podle toho, kdo má nějaký nápad, kterým by mohl přispět, a posunout tak řešení blíže úspěchu. Vznikly dva základní typy gramatických systémů. První typ umožňuje distribuované zpracování a nazývá se kooperující distribuovaný gramatický systém, zkráceně CD gramatický systém (z anglického *Cooperating Distributed Grammar System*). Jednotlivé komponenty systému (takto se v terminologii gramatických systémů nazývají gramatiky daného systému) pracují sekvenčně na jedné sdílené větě

formě, kterou postupně upravují a dle různých pravidel se střídají při jejím zpracování. Naopak druhý typ gramatických systémů je představitelem paralelního zpracování a tento systém se nazývá paralelně komunikující gramatický systém, zkráceně PC gramatický systém (z anglického *Parallel Communicating Grammar System*). Komponenty tohoto systému pracují všechny současně, tedy paralelně, na svých vlastních větných formách. Každá komponenta tedy upravuje pouze svoji větnou formu, ale mohou mezi sebou komunikovat a předávat si mezivýsledky pomocí tzv. komunikačních symbolů. Velký potenciál gramatických systémů spočívá právě v těchto dvou přístupech, protože zavedly formální prostředky pro oblasti, které jsou v dnešní době stále více aktuální. Díky současným trendům v IT a komunikačních technologiích je paralelní a distribuované zpracování velkým tématem, a do budoucna lze očekávat jediné vzestupné tendence. Proto je výzkum v těchto oblastech, dle mého názoru, nejen důležitý, ale i nutný a gramatické systémy jsou jeden z prostředků, který by mohl přinést další možnosti, ať už v oblasti překladačů nebo i jiných odvětvích, jako je lingvistika, molekulární biologie a další.

Druhá polovina práce je věnována především praxi, a to v podobě syntaktické analýzy a syntaxí řízeném překladu. Jednotlivé kapitoly postupně popisují teorii nutnou pro implementaci syntaktického analyzátoru, porovnávají různé přístupy a možnosti. Mým cílem bylo vytvořit aplikaci, která kombinuje dva přístupy, tedy dvě různé gramatiky, syntaktické analýzy, a to syntaktickou analýzu shora dolů s analýzou zdola nahoru a uvažovat je v kontextu gramatických systémů. Pro syntaktickou analýzu shora dolů jsou velmi vhodné LL gramatiky, které sice mají menší vyjadřovací sílu, než obecné bezkontextové gramatiky, ale umožňují velmi snadno implementovat syntaktický analyzátor. Existují dva nejčastější přístupy, rekurzivní sestup a prediktivní syntaktická analýza. Oba přístupy vycházejí z tzv. LL tabulky, která na průniku všech terminálů s neterminály gramatiky definuje, jaké pravidlo z množiny pravidel gramatiky se má aplikovat pro danou kombinaci terminálu na vstupu s aktuálně zpracovávaným neterminálem na zásobníku. Ačkoli je rekurzivní sestup velmi intuitivní a jednoduchý (každý neterminál je reprezentován metodou a tyto metody se dle pravidel gramatiky mezi sebou rekurzivně volají, dokud nekontrolují celý vstupní řetězec), rozhodl jsem se pro prediktivní syntaktický analyzátor. Tento přístup sice vyžaduje implementaci zásobníku a LL tabulky, ale samotná analýza je řízena jedním algoritmem a výsledný analyzátor není nijak závislý na konkrétní LL tabulce, což se ukázalo velmi užitečné. Pro vytvoření LL tabulky je nutné postupně sestavit několik pomocných množin, jejichž podrobnému popisu jsou věnovány jednotlivé podkapitoly.

Syntaktická analýza zdola nahoru má také své dva zástupce, z nichž komplexním řešením je LR syntaktický analyzátor, který pokud má k dispozici komplexní gramatiku, dokáže provést analýzu celého vstupu. Je tak velmi vhodný pro profesionální a komerční řešení, ale je složitý na implementaci. Zároveň pokud by již byl použit, bylo by snazší ho aplikovat pro celou analýzu a ne pouze pro část. Proto je pro účely této práce vhodnější druhý přístup, tzv. precedenční syntaktický analyzátor. Jedná se o velmi jednoduchý aparát, který nemá příliš velkou vyjadřovací sílu, ale je velmi vhodný např. pro vyhodnocování matematických výrazů, což je přesně část kódu, se kterou má problém prediktivní syntaktická analýza shora dolů. Existují tak tedy dva přístupy, z nichž každý pracuje jiným způsobem a je vhodný pro analýzu jiné části vstupního kódu, což jde přesně ruku v ruce s myšlenkou gramatických systémů. Velmi se zde nabízí konkrétní řešení, a to centralizovaný PC gramatický systém. Tento systém je speciální v tom, že pouze hlavní řídicí komponenta může generovat komunikační symboly, a má tak možnost si dle potřeby vyžádat práci ostatních komponent. Řídicí komponenta je v tomto případě analýza shora dolů, která kontroluje syntaktickou strukturu celého jazyka a pro zpracování matematických výrazů si zavolá analýzu zdola nahoru, přičemž pro komunikaci mezi komponentami byl zaveden zvláštní pseudoterminál.

Výsledná aplikace tak umožňuje uživateli provádět syntaktickou analýzu jeho vlastního vstupního textu (předpokládá se nějaká forma programovacího jazyka) na základě spolupráce těchto dvou gramatik. Hlavním cílem celé aplikace je jednoduchost použití a co největší přidaná hodnota pro uživatele, proto se jedná o webovou aplikaci. Uživatel nemusí nic instalovat ani složitě spouštět, pouze se připojí na webovou stránku, kde je daná aplikace umístěná. Aplikace je navržena jako jednoduchá responzivní stránka (tzv. single page aplikace) s čistým bílým designem, aby byla co nejjednodušeji použitelná a přehledná. Uživatel má dále možnost snadno pomoci tagů (štítků) definovat svoji vlastní LL gramatiku pro analýzu shora dolů. To je možné díky tomu, že je analýza shora dolů implementovaná pomocí prediktivní syntaktické analýzy a není závislá na konkrétní LL tabulce. Uživatel tedy nadefinuje svoji LL gramatiku, vloží vstupní kód, nad kterým chce provést syntaktickou analýzu, a aplikace si již sama dynamicky sestaví LL tabulku a provede analýzu. Jako výsledek zobrazí uživateli úspěch nebo chybové hlášení a posloupnost aplikovaných pravidel jak LL gramatiky, tak i precedenčního analyzátoru. Uživatel má tak k dispozici nejen nástroj pro syntaktickou analýzu, ale uvidí způsob aplikace pravidel a vygenerovanou LL tabulku. Aplikace tedy slouží především pro didaktické použití a měla by uživateli napomoci k lepšímu pochopení problematiky překladačů, bezkontextových gramatik, syntaktické analýzy a v neposlední řadě gramatických systémů. Zároveň byla aplikace implementována pomocí celé řady moderních webových technologií, jako je HTML5, CSS3, LESS, Bootstrap, JavaScript, AngularJS a dalších, které jsou podrobněji popsány v kapitole o aplikaci. V sekci příloh se potom nachází reálná ukázka výsledné aplikace a jejího použití.

Práce zasahuje do velkého množství témat a otevírá mnoho dveří dalšímu výzkumu a rozvoji. V první řadě se nabízí teoretický i praktický výzkum paralelismu a distribuce výkonu gramatických systémů při syntaktické analýze. Zjišťovat, jaký dopad by mělo použití více gramatik (komponent), ze kterých by se sestavil a popsat komplexní gramatický systém. Tento systém v ideálním případě přenést do praxe, ať už pomocí paralelních výpočtů, pokud by se jednalo o PC gramatický systém, nebo distribuovat výkon mezi větší množství výpočetních stanic u CD gramatických systémů. Dalo by se tímto způsobem dosáhnout většího výkonu? Mělo by význam zavést tyto systémy do praxe vzhledem ke stávajícím komerčním přístupům? Okolo gramatických systémů je mnoho otázek, jejichž zodpovězení by mohlo přinést řadu nových a významných objevů nebo naopak zjištění, že jsou velmi zajímavé z teoretického pohledu („na papíře“), ale reálně nepoužitelné.

Další možností je rozvoj samotné aplikace, která díky implementaci pomocí moderních technologií tak rychle nezestárne a nabízí celou řadu způsobů, jak ji dále rozvinout. Jedna varianta by mohla být rozšíření výstupu analýzy o další podrobné informace, které by tak mohly poskytnout hlubší porozumění celé problematice. Nebo se zaměřit na dynamickou stránku a umožnit definici i druhé gramatiky, která pracuje zdola nahoru, což by mohlo být velmi zajímavé. Uživatel by mohl definovat kromě LL gramatiky i precedenční gramatiku, resp. precedenční tabulku. Případně by se dala nahradit precedenční analýza silnějším nástrojem, např. LR syntaktickým analyzátozem. Poté by šla zkoumat síla LL gramatiky vůči bezkontextové gramatice pro LR analyzátor. Pokud by se implementace provedla správně, bylo by velmi pravděpodobné, že by se dalo provést spojení obou analýz, jak je tomu teď, nebo by mohla být LR gramatika tak komplexně napsána, že by převzala sílu LL gramatiky a provedla pak analýzu sama. Stejně jako samotné gramatické systémy, tak i aplikace nabízí mnoho možností, jak v této práci pokračovat.

Literatura

- [1] MEDUNA, Alexander a Petr ZEMEK. *Regulated grammars and automata.*, 694 stran. ISBN 1493903683.
- [2] MEDUNA, Alexander. *Formal languages and computation: models and their applications.* 315 stran. ISBN 9781466513457.
- [3] MEDUNA, Alexander. *Elements of compiler design.* Boca Raton, FL: Auerbach Publications, 2008, 286 stran. ISBN 1420063235.
- [4] TECHET, Jiří, Tomáš MASOPUST a Alexander MEDUNA. Modern Formal Language Theory: Cooperating Distributed Grammar Systems [online]. Brno, 2007 [cit. 2017-02-01]. Dostupné z:
<http://www.fit.vutbr.cz/~meduna/work/lib/exe/fetch.php?media=lectures:phd:tid:frvs:09-cdgspres.pdf>
- [5] TECHET, Jiří, Tomáš MASOPUST a Alexander MEDUNA. Modern Formal Language Theory: Parallel Communicating Grammar Systems [online]. Brno, 2007 [cit. 2017-02-01]. Dostupné z:
<http://www.fit.vutbr.cz/~meduna/work/lib/exe/fetch.php?media=lectures:phd:tid:frvs:10-pcgspres.pdf>
- [6] ROZENBERG, Grzegorz. Cooperating grammar systems. *Mathematical Foundations of Computer Science 1978: Lecture Notes in Computer Science.* Springer Berlin Heidelberg, 1978, strany. 364-373. ISBN 9783540089216.
- [7] ROZENBERG, Grzegorz a Arto SALOMAA. *Handbook of formal languages: background and application.* Berlin: Springer-Verlag, 1997, 528 stran. ISBN 3540614869.
- [8] NIL, Penny. *The Blackboard Model of Problem Solving and the Evolution of Blackboard Architectures.* AI Magazine Volume 7 Number 2, 1986, strany 38–53
- [9] WEBBER, Adam. *Formal language: a practical introduction.* Wilsonville, Or.: Franklin, Beedle & Associates, 2008, 388 stran. ISBN 978-1590281970.
- [10] LINZ, Peter. *An introduction to formal languages and automata.* 5th ed. Sudbury, MA: Jones & Bartlett Learning, 2012, 437 stran. ISBN 978-1449615529.
- [11] BROOKSHEAR, J. Glenn. *Theory of computation: formal languages, automata, and complexity.* Redwood City, Calif.: Benjamin/Cummings Pub. Co., 1989, 320 stran. ISBN 978-0805301434.
- [12] SIPSER, Michael. *Introduction to the theory of computation.* 3rd Ed. Boston, MA: Course Technology Cengage Learning, 2012, 480 stran. ISBN 978-1133187790.
- [13] CSUHAI-VARJÚ, E. *Grammar systems: A grammatical approach to distribution and cooperation.* Psychology Press, 1994, 246 stran. ISBN 2881249574.
- [14] PAUN, Gheorghe a Arto SALOMAA. *New trends in formal languages: control, cooperation, and combinatorics.* New York: Springer, 1997, 474 stran. ISBN 3540628444.

Příloha A

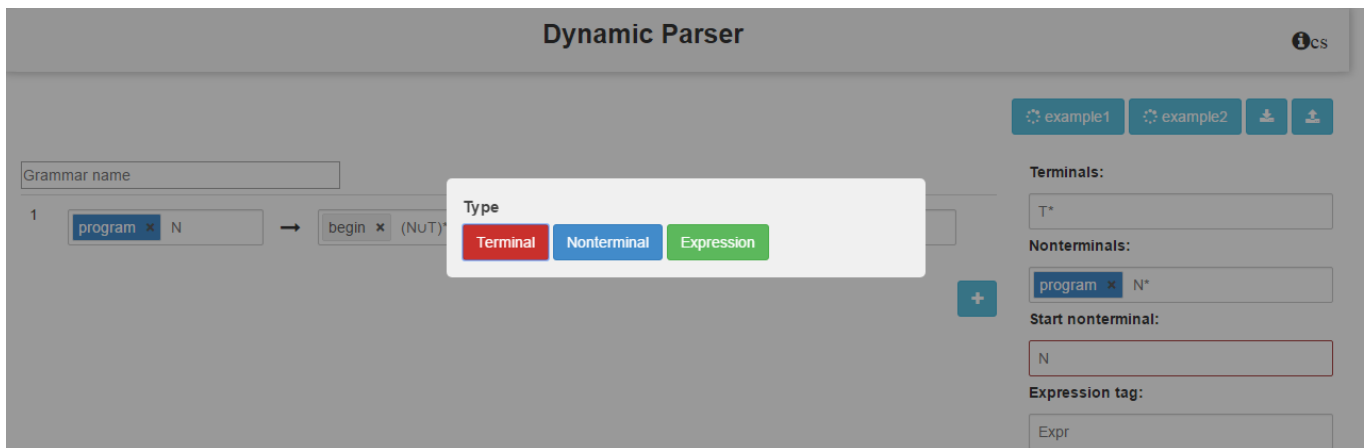
Obsah CD:

1. Text diplomové práce ve formátu PDF
2. Zdrojový text práce ve formátu ODT (docx)
3. Zdrojové soubory aplikace (manuál viz. kapitola 7.3)

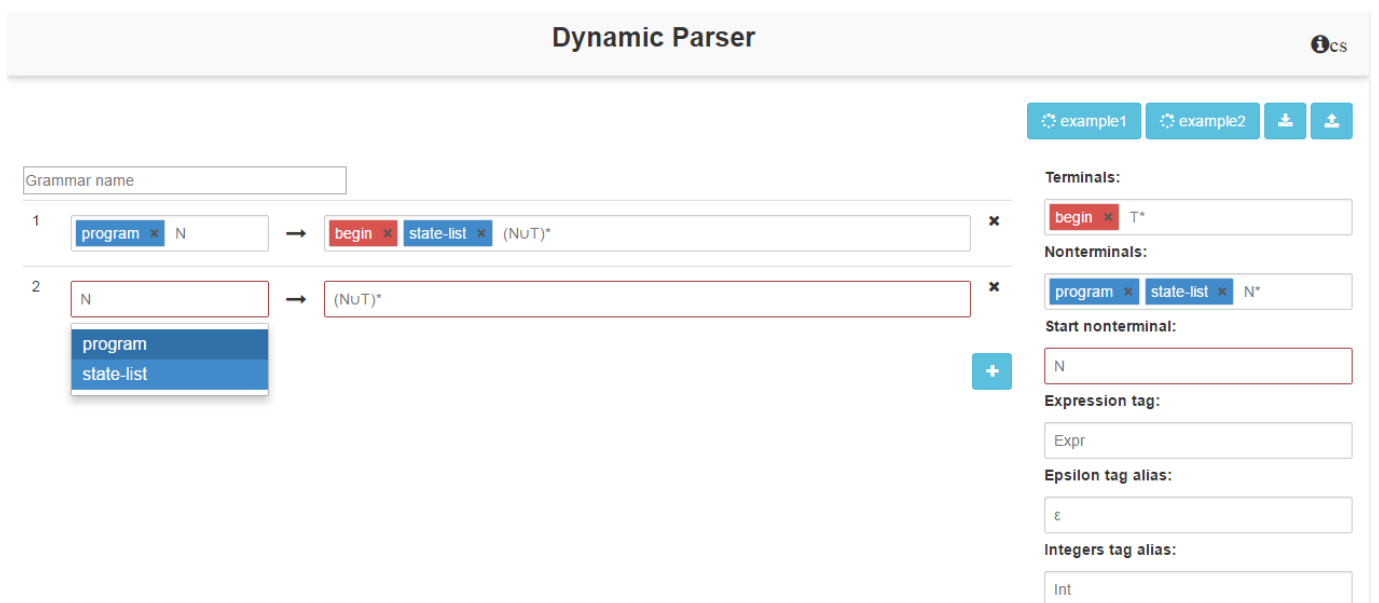
Příloha B

The screenshot shows the 'Dynamic Parser' web application. At the top, there is a header with the title 'Dynamic Parser' and a user icon 'i cs'. Below the header, there are two buttons: 'example1' and 'example2', along with download and upload icons. The main area is divided into two columns. The left column contains a 'Grammar name' input field and a rule editor. The rule editor shows a rule with index '1', a left-hand side 'N', and a right-hand side '(NUT)*'. A plus sign button is located below the rule editor. The right column contains several configuration fields: 'Terminals:' with a field containing 'T*'; 'Nonterminals:' with a field containing 'N*'; 'Start nonterminal:' with a field containing 'N'; 'Expression tag:' with a field containing 'Expr'; 'Epsilon tag alias:' with a field containing 'ε'; 'Integers tag alias:' with a field containing 'Int'; 'Floats tag alias:' with a field containing 'Float'; 'Strings tag alias: (e.g. "text")' with a field containing 'String'; and 'Identifier tag alias:' with a field containing 'Id'. Below these fields, there is a section for the 'LL Table' with a 'Generate LL Table' button. At the bottom, there is an 'Input code' section with a text area containing a comment: '/* ** Enter your own program code ** */'. Below the text area is an 'Analyze input' button. At the very bottom, there is a footer with the text: 'Dynamic Parser © Copyright 2017 Jarošlav Handlíř, FIT VUT Brno xhandl05@stud.fit.vutbr.cz'.

Obrázek 12: Výchozí vzhled aplikace po načtení stránky



Obrázek 13: Ukázka volby typu nového symbolu při psaní pravidel



Obrázek 14: Ukázka vytváření pravidel - nabídka již existujících tagů (neterminálů) pro pravou stranu pravidla

Dynamic Parser

ics

example1 example2 [upload] [download]

Example 1 (from Thesis, only top->down)

- 1 →
- 2 →
- 3 →
- 4 →
- 5 →
- 6 →
- 7 →
- 8 →

Terminals:

Nonterminals:

Start nonterminal:

Expression tag:

Epsilon tag alias:

Integers tag alias:

Floats tag alias:

Strings tag alias: (e.g. "text")

Identifier tag alias:

LL Table:

	begin	;	end	print	read	int	id	\$
prog	1							
state			3	2	2	2	2	
cmd				4	5	6	6	
item						7	8	

Generate LL Table

Input code:

```
begin
  42;
end
// end of code
```

Analyze input

Input in tokens:

- 1
- 2
- 3
- 4

Syntax analysis result: ✓

Sequence of applied rules:

- 1
- 2
- 6
- 7
- 3

Dynamic Parser © Copyright 2017 Jaroslav Handlíř, FIT VUT Brno
 xhandl05@stud.fit.vutbr.cz

Obrázek 15: Ukázka vzorového příkladu - kompletní gramatika, sestavená LL tabulka a provedená analýza