

VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ
FAKULTA INFORMAČNÍCH TECHNOLOGIÍ

Formální jazyky a překladače 2016/2017

Dokumentace interpretu imperativního jazyka IFJ16

Tým 013, varianta b/3/I

Rozšíření: BOOLOP

| | | | |
|----------------|----------|-----|--------------|
| Šuba Adam | xsubaa00 | 25% | vedoucí týmu |
| Paliesek Jakub | xpalie00 | 25% | |
| Šuhaj Peter | xsuhaj02 | 25% | |
| Tóth Adrián | xtotha01 | 25% | |

11. 12. 2016

Obsah

| | | |
|----------|--|----------|
| 1 | Úvod | 2 |
| 2 | Návrh | 2 |
| 3 | Interpret | 3 |
| 3.1 | Lexikální analyzátor | 3 |
| 3.2 | Syntaktický analyzátor | 3 |
| 3.3 | Sémantická analýza | 4 |
| 3.4 | Tabulka symbolů | 4 |
| 3.5 | Interpret | 5 |
| 4 | Algoritmy | 5 |
| 4.1 | Vyhledávání podřetězce použitím Boyer-Mooreova algoritmu (b) | 5 |
| 4.2 | Řazení použitím algoritmu Shell sort (3) | 6 |
| 4.3 | Implementace tabulky symbolů použitím binárního vyhledávacího stromu (I) | 6 |
| 5 | Práce v týmu | 6 |
| 5.1 | Rozdělení týmu | 6 |
| 5.2 | Způsob práce | 6 |
| 6 | Závěr | 6 |
| 7 | Literatura | 7 |
| 8 | Příloha | 8 |
| 8.1 | LL-gramatika | 8 |
| 8.2 | Precedenční tabulka | 10 |
| 8.3 | Diagram konečného automatu lexikální analýzy | 11 |

1 Úvod

Dokumentace popisuje vývoj interpretu imperativního jazyka IFJ16. Jazyk IFJ16 je podmnožinou jazyka Java SE 8, což je staticky typovaný, objektově-orientovaný jazyk. Úlohou našeho interpretu je kontrola zdrojového kódu napsaného v jazyku IFJ16 a jeho následná interpretace.

Zvolená varianta zadání Tým 013 má parametry b/3/I:

- (b) vyhledávání podřetězce s použitím Boyer-Mooreova algoritmu
- (3) řazení s použitím algoritmu Shell sort
- (I) implementace tabulky symbolů použitím binárního vyhledávacího stromu

Dokumentace je rozdělená na kapitoly, které blíže popisují jednotlivé části interpretu a způsob jejich řešení. Taktéž obsahuje i popis použitých algoritmů, rozdělení úloh a způsob práce v týmu. V závěru dokumentace se nachází shrnutí naší práce.

Příloha obsahuje:

- diagram konečného automatu lexikální analýzy
- LL-gramatika pro syntaktickou analýzu shora dolů
- precedenční tabulka pro syntaktickou analýzu zdola nahoru

2 Návrh

Plánování procesu vývoje interpretu nám bylo značně zjednodušeno závazně zadanými metodami implementace syntaktické analýzy. V případě metody shora dolů jsme se rozhodli pro rekursivní sestup.

Samotný vývoj jsme naplánovali na 6 fází, s cílem mít při pokusném odevzdání plně funkční produkt.

1. Vývoj lexikálního analyzátoru
2. Vývoj syntaktického analyzátoru
3. Příprava tabulky symbolů
4. Vývoj sémantické analýzy a její integrace do syntaktického analyzátoru
5. Vývoj interpretu
6. Dokončení a vyladění celého produktu, tvorba dokumentace

Vývoj lexikálního analyzátoru zahrnuje přípravu konečného automatu rozpoznávajícího lexémy jazyka IFJ16 a jeho implementaci. První částí vývoje syntaktického analyzátoru je implementace rekursivního sestupu. To vyžaduje tvorbu LL-gramatiky a LL-tabulky. Druhou částí je implementace precedenční analýzy, která vyžaduje tvorbu precedenční tabulky. V mezichase je možné navrhnout a implementovat tabulku symbolů a připravit sémantickou analýzu pro integraci do syntaktického analyzátoru. Jakmile je hotová kontrola správnosti vstupního programu a jsou generovány instrukce, je možné přistoupit

k implementaci interpretu. V průběhu vývoje jednotlivých modulů jsou připravovány testy jak samotných modulů, tak jejich spolupráce s ostatními moduly. Závěrečná fáze zahrnuje odstranění co nejvíce chyb a tvorba dokumentace.

3 Interpret

Implementaci interpretu lze rozdělit do 5 logických celků. Hlavní komponentou je syntaktický analyzátor, neboli „parser“, který řídí činnost lexikálního analyzátoru, „scanneru“, a také provádí sémantické akce sémantického analyzátoru. Sémantická analýza také vkládá informace o funkcích a proměnných do tabulky symbolů. Poslední komponentou je samotný interpret, který provádí interpretaci instrukcí, vygenerovaných během kontroly vstupního programu.

Kontrola správnosti programu je prováděna ve dvou průchodech, v prvním se některé části zdrojového kódu úplně ignorují (zejména výrazy, které jsou závislé na sémantických akcích, více dále). Účelem prvního průchodu je naplnění tabulky symbolů.

3.1 Lexikální analyzátor

Jediným úkolem lexikálního analyzátoru je na žádost syntaktického analyzátoru přechíst další symbol ze zdrojového souboru, klasifikovat jej a předat ve formě tokenu zpět syntaktickému analyzátoru.

Implementace je založena na konečném automatu. Bílé znaky a komentáře jsou ze zdrojového kódu odignorovány. Pokud řetězec dojde do koncového stavu reprezentující nějaký lexém a následující znak nepatří mezi znaky, které může daný lexém obsahovat a zároveň znak patří mezi znaky které mohou bezprostředně za daným lexémem následovat, je typ vráceného tokenu nastaven na typ odpovídající tomuto lexému. Některým tokenům je ještě přidělen atribut, obsahující užitečná data získaná z lexému. Pokud je v některém stavu přečten znak, který není v daném stavu podporován, popřípadě pokud je v koncovém stavu přečten znak, který nemůže následovat bezprostředně za daným lexémem, dojde k chybě. Poslední načtený znak, který rozhodl o ukončení konečného automatu, je vrácen zpět do vstupního souboru.

Poznámka:

Při druhém průchodu syntaktické analýzy již není lexikální analyzátor činný. Při každém získání tokenu během prvního průchodu je daný token uložen do datového typu fronta a v druhém průchodu jsou tokeny brány z této fronty.

3.2 Syntaktický analyzátor

Mozkem celé kontroly správnosti zdrojového kódu je syntaktický analyzátor. Samotnou syntaktickou analýzu můžeme rozdělit do dvou částí: kontrola rekurzivním sestupem a kontrola správnosti výrazů precedenční analýzou.

Rekurzivní sestup je metoda syntaktické analýzy shora dolů, která implementuje pravidla LL-gramatiky (viz příloha) na základě LL tabulky. Podle aktuálního stavu analyzátoru a vstupního tokenu je rozhodnuto, které pravidlo se zvolí. Metoda využívá rekurzivního volání funkce implementující dané pravidlo.

Pro kontrolu správnosti výrazů se využívá precedenční analýza, která je metodou syntaktické analýzy

zdola nahoru. Neboť bylo nutné nějakým způsobem napojit precedenční analýzu na rekurzivní sestup, přistoupili jsme k zavedení speciálního pseudoterminálu `E` do naší LL-gramatiky. V místě kontroly, kde by se měl nacházet tento terminál, se analyzátor přepne do režimu, kdy pouze žádá další a další tokeny a ukládá je do datového typu fronta. Tento režim je ukončen v momentu, kdy narazí na token, který není ve výrazech podporován. Takto vytvořená fronta je poté předána precedenční analýze ke kontrole. Speciálním případem, kdy je nutné lehce pozměnit chování tohoto procesu, jsou výrazy podmínek ve větvení a v cyklu. Zde je nutné počítat levé a pravé závorky tak, aby byl analyzátor schopen poznat, kdy nastal konec podmínky a fronta je naplněná. Samotná precedenční analýza je implementována podle algoritmu vyučovaném v předmětu IFJ. Pro účely precedenční analýzy, byly k výčtovému typu `token_t` přidány položky `token_nonterminal` a `token_stop`.

Neboť neimplementujeme rozšíření FUNEXP je pro nás funkce `ifj16.print()` speciální syntaktickou konstrukcí, neboť jí lze jako parametr dát jednoduchou konkatenaci, jejíž výsledkem je vždy datový typ string. Z tohoto důvodu jsme se rozhodli nezahrnovat tohle chování do LL-gramatiky, ale vyřešit to jako speciální případ až při samotné implementaci analyzátoru. V druhém průchodu zdrojového souboru již díky sémantickým akcím analyzátor ví, že volaná funkce je `ifj16.print()` a na základě této informace se přepne do speciálního režimu, podobného přípravě na precedenční analýzu. Rozdíl je v tom, že povolené tokeny jsou pouze proměnné, literály a operátor konkatenace „+“. Při ukončení tohoto režimu je zkontrolováno, zda-li alespoň jeden token byl typu string. Poté je konkatenace vyhodnocena precedenční analýzou. V naší implementaci je tedy tato jednoduchá konkatenace vyhodnocována zleva doprava podle stejných pravidel jako jiné výrazy. Výsledkem `3 + 5 + "\n"` je tedy `"8\n"` a nikoliv `"35\n"`.

3.3 Sémantická analýza

Sémantická analýza je soubor funkcí a kontrol volaných během syntaktické analýzy v obou průchodech. V 1. průchodu do tabulky symbolů vkládá informace o třídách, statických funkcích, statických proměnných, argumentech funkcí a kontroluje, zda nedošlo k redefinici. Těla funkcí nejsou v 1. průchodu sémanticky zpracována. V druhé fázi jsou pak generovány instrukce inicializačních výrazů statických proměnných a do tabulky symbolů jsou vkládány lokální proměnné. Rovněž jsou v této fázi zpracována těla funkcí, což zahrnuje kontrolu typů, kontrolu, zda byl identifikátor při jeho použití (u volání funkce, přiřazení) již definován a následné generování instrukcí, pokud nastaly chyby zmíněné dříve.

V případě sémantické správnosti kódu je během precedenční analýzy vždy při redukci části výrazu na neterminál generována aritmetická instrukce, případně i dodatečné konverzní instrukce, pokud je daná implicitní konverze podporována. Výstupem sémantické akce generující aritmetickou instrukci je pomocná proměnná, která je během precedenční analýzy přiřazena jako atribut cílovému neterminálu a může být dále použita při jiné redukci.

3.4 Tabulka symbolů

Tabulka symbolů (dále jen TS) je podle zadání implementovaná pomocí binárního vyhledávacího stromu. Vzhledem k tomu, že v jazyce IFJ16 existují různé rozsahy platnosti identifikátorů a třídy umožňují vytvářet oddělené jmenné prostory, navrhli jsme TS jako tříúrovňovou kaskádu vyhledávacích struktur. První úroveň je jediná tabulka tříd. Každá třída má potom vlastní tabulku členů (společnou pro statické proměnné a funkce). Funkce mají dále vlastní tabulky lokálních proměnných.

Hodnoty statických proměnných jsou uloženy přímo v odpovídajících položkách TS. Hodnoty lokálních

proměnných nejsou, kvůli možnému rekurzivnímu volání funkcí, uložené přímo v TS, ale jejich hodnoty si ukládá interpret na svém zásobníku během vykonávání instrukcí.

Různé typy proměnných (statická, lokální, pomocná) a literály jsou v naší implementaci reprezentovány různými strukturami, ale při sémantických akcích často není potřeba mezi nimi rozlišovat. Kvůli tomuto jsme zavedli jednu extra strukturu, která je umístěna jako první člen ostatních struktur. Obsahuje informaci o tom, na kterou strukturu lze ukazatel na tuto univerzální strukturu dále přetypovat v případě potřeby. Je zde využita vlastnost jazyka C, že ukazatel na strukturu ukazuje na svého prvního člena a naopak.¹ Kromě toho ještě obsahuje informaci o datovém typu, který může sémantická analýza zkontrolovat bez nutnosti přetypování.

3.5 Interpret

Interpret obsahuje jednu velkou funkci, jejímž vstupem je instrukční páska. Funkce se pomocí cyklu posouvá po instrukční pásce až do dosažení konce, tj. na instrukční pásce je NULL.

Podle operačního kódu instrukce se zvolí příslušná větev interpretu a pomocná dekodovací funkce získá z instrukce operandy (jejich hodnotu, informaci o tom, jestli jsou inicializované apod.). Pro rámce funkcí má interpret svůj vlastní zásobník. Při volání funkce se vytvoří nový rámec, vyplní se potřebnými hodnotami (argumenty funkce, místo pro návratovou hodnotu, ukazatel na instrukci pro návrat) a vloží se na zásobník. Při návratu z funkce je návratová hodnota uložena na místa v rámci pro ni vymezeného a interpretace pokračuje od poznačené instrukce.

4 Algoritmy

Všechny popsané algoritmy jsou implementovány v souboru `ial.c`

4.1 Vyhledávání podřetězce použitím Boyer-Mooreova algoritmu (b)

Boyer-Moorův algoritmus je neefektivnější algoritmus na vyhledávání podřetězce v řetězci. Vychází z klasického porovnávání hledaného vzorku patřičně přiloženému k prohledávanému řetězci, ale s tím rozdílem, že změna zarovnání vzorku s řetězcem se neděje po jednom, ale po více znacích. Porovnávání znaků přiloženého vzorku se děje od konce k počátku.

K určení, o kolik znaků lze vzorek posunout (kolik zarovnání lze přeskočit), se používají dvě heuristiky. První z nich je heuristika „bad character“, která se při neshodě snaží vzorek zarovnat tak, aby ke shodě došlo. Pokud to není možné (vzorek již daný znak neobsahuje), vzorek se zarovná až za daný znak v řetězci. Druhou heuristikou je „good suffix“, která se v případě neshody po několika shodných znacích snaží vzorek zarovnat tak, aby nedošlo k porušení již nalezené shody. V naší implementaci jsme se rozhodli využít pouze první heuristiku.

¹<http://www.open-std.org/jtc1/sc22/wg14/www/docs/n1124.pdf>, s. 103

4.2 Řazení použitím algoritmu Shell sort (3)

Shell sort je řadící algoritmus, který pracuje na principu bublinového vkládání. Neřadí prvky nacházející se výhradně vedle sebe, ale i poslovnosti oddělené určitým počtem prvků (krokem). V naší implementaci je na počátku tento krok stanoven na polovinu počtu prvků. Při každé další iteraci je hodnota kroku snížena na polovinu. V různých implementacích existují jiné způsoby určení těchto kroků, od čeho se odvíjí jejich časová složitost. Algoritmus skončí jakmile krok dosáhne hodnoty 0. Shell sort pracuje „in situ“ a je nestabilní, což nám nevadí. Ve srovnání s jinými algoritmy je pomalejší než Quicksort a rychlejší než Heapsort. Složitost naší implementace je $O(n^2)$.

4.3 Implementace tabulky symbolů použitím binárního vyhledávacího stromu (I)

Binární vyhledávací strom je binární strom, platí tedy, že je složen z kořene a levého a pravého podstromu. Všechny uzly nalevo od kořene mají klíče menší než kořen a všechny uzly napravo mají klíče větší než kořen. Vyhledávacím klíčem je v naší implementaci identifikátor symbolu. Funkce pro práci s binárními vyhledávacími stromy (v našem případě vyhledání uzlu, vložení uzlu a zrušení celého stromu průchodem PostOrder) byly implementovány rekurzivně. Implementace samotné tabulky symbolů byla popsána detailně v kapitole 3.4.

5 Práce v týmu

5.1 Rozdělení týmu

- Adam Šuba – lexikální analyzátor, syntaktická analýza shora dolů a zdola nahoru, dokumentace, knihovna pro nekonečné řetězce, vedoucí
- Jakub Paliesek – sémantický analyzátor, tabulka symbolů, algoritmy, generátor 3-adresného kódu, knihovna pro nekonečné řetězce, dokumentace
- Peter Šuhaj – interpret, algoritmy, dokumentace
- Adrián Tóth – testy, algoritmy, dokumentace

5.2 Způsob práce

Pro práci byl zvolen systém správy verzí Git, hostovaný na serveru GitHub. Kromě osobní komunikace byla hlavním komunikačním kanálem skupinová konverzace na sociální síti Facebook.

6 Závěr

Tým pracoval spolehlivě a časový plán se nám dařilo plnit. Jediné období, kdy nebylo možné na projektu aktivně pracovat bylo během půlsestrálního období, s čímž jsme ovšem během plánování počítali a problémy to tedy nezpůsobilo. Při prvním pokusném odevzdání jsme měli až na několik drobných chyb funkční interpret.

7 Literatura

1. HONZÍK, Jan M., Algoritmy: Studijní opora. Verze: 16-B, Brno: Vysoké učení Technické, 2014
2. Shell sort. In: Wikipedia: the free encyclopedia [online]. San Francisco (CA): Wikimedia Foundation, 2001- [cit. 2016-12-11]. Dostupné z: https://cs.wikipedia.org/wiki/Shell_sort
3. ISO/IEC. ISO International Standard ISO/IEC 9899:TC2 – Programming Language C. Dostupné na <http://www.open-std.org/jtc1/sc22/wg14/www/docs/n1124.pdf>

8 Příloha

8.1 LL-gramatika

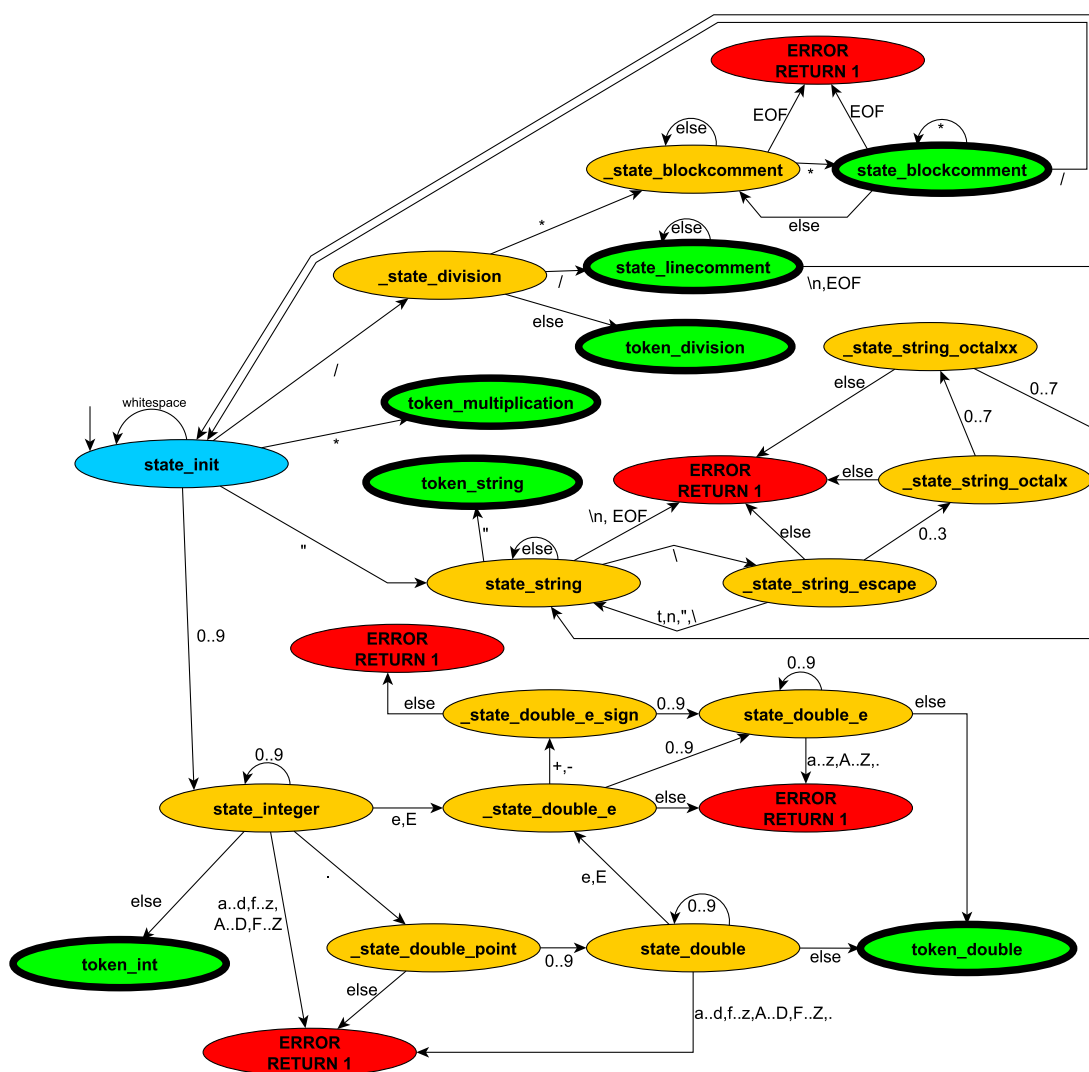
1. $\langle \text{c-list} \rangle \rightarrow \text{class id } \{ \langle \text{memb-list} \rangle \} \langle \text{c-list} \rangle$
2. $\langle \text{c-list} \rangle \rightarrow \$$
3. $\langle \text{memb-list} \rangle \rightarrow \langle \text{c-memb} \rangle \langle \text{memb-list} \rangle$
4. $\langle \text{memb-list} \rangle \rightarrow \epsilon$
5. $\langle \text{c-memb} \rangle \rightarrow \text{static } \langle \text{c-memb1} \rangle$
6. $\langle \text{c-memb1} \rangle \rightarrow \langle \text{type} \rangle \text{ id } \langle \text{c-memb2} \rangle$
7. $\langle \text{c-memb1} \rangle \rightarrow \text{void id } \langle \text{c-membfunc} \rangle$
8. $\langle \text{c-memb2} \rangle \rightarrow = \text{ E } ;$
9. $\langle \text{c-memb2} \rangle \rightarrow ;$
10. $\langle \text{c-memb2} \rangle \rightarrow \langle \text{c-membfunc} \rangle$
11. $\langle \text{c-membfunc} \rangle \rightarrow (\langle \text{fn-def-plist} \rangle) \{ \langle \text{fn-body} \rangle \}$
12. $\langle \text{type} \rangle \rightarrow \text{int}$
13. $\langle \text{type} \rangle \rightarrow \text{double}$
14. $\langle \text{type} \rangle \rightarrow \text{String}$
15. $\langle \text{type} \rangle \rightarrow \text{boolean}$
16. $\langle \text{fn-def-plist} \rangle \rightarrow \langle \text{par-def} \rangle \langle \text{fn-def-plist1} \rangle$
17. $\langle \text{fn-def-plist} \rangle \rightarrow \epsilon$
18. $\langle \text{fn-def-plist1} \rangle \rightarrow \langle \text{par-def} \rangle \langle \text{fn-def-plist1} \rangle$
19. $\langle \text{fn-def-plist1} \rangle \rightarrow \epsilon$
20. $\langle \text{par-def} \rangle \rightarrow \langle \text{type} \rangle \text{ id}$
21. $\langle \text{fn-body} \rangle \rightarrow \langle \text{stat} \rangle \langle \text{fn-body} \rangle$
22. $\langle \text{fn-body} \rangle \rightarrow \langle \text{type} \rangle \text{ id } \langle \text{opt-assign} \rangle ; \langle \text{fn-body} \rangle$
23. $\langle \text{fn-body} \rangle \rightarrow \epsilon$
24. $\langle \text{stat-com} \rangle \rightarrow \{ \langle \text{stat-list} \rangle \}$
25. $\langle \text{stat-list} \rangle \rightarrow \langle \text{stat} \rangle \langle \text{stat-list} \rangle$
26. $\langle \text{stat-list} \rangle \rightarrow \epsilon$
27. $\langle \text{stat} \rangle \rightarrow \langle \text{id} \rangle \langle \text{as-ca} \rangle ;$
28. $\langle \text{stat} \rangle \rightarrow \langle \text{stat-com} \rangle$
29. $\langle \text{stat} \rangle \rightarrow \text{if } (\text{ E }) \langle \text{stat-com} \rangle \text{ else } \langle \text{stat-com} \rangle$
30. $\langle \text{stat} \rangle \rightarrow \text{while } (\text{ E }) \langle \text{stat-com} \rangle$
31. $\langle \text{stat} \rangle \rightarrow \text{return } \langle \text{ret-val} \rangle ;$
32. $\langle \text{id} \rangle \rightarrow \text{id}$
33. $\langle \text{id} \rangle \rightarrow \text{fqid}$
34. $\langle \text{as-ca} \rangle \rightarrow (\langle \text{fn-plist} \rangle)$
35. $\langle \text{as-ca} \rangle \rightarrow = \langle \text{assign} \rangle$

- 36. $\langle \text{assign} \rangle \rightarrow \text{E}$
- 37. $\langle \text{assign} \rangle \rightarrow \langle \text{id} \rangle (\langle \text{fn-plist} \rangle)$
- 38. $\langle \text{opt-assign} \rangle \rightarrow = \langle \text{assign} \rangle$
- 39. $\langle \text{opt-assign} \rangle \rightarrow \epsilon$
- 40. $\langle \text{fn-plist} \rangle \rightarrow \langle \text{val-id} \rangle \langle \text{fn-plist1} \rangle$
- 41. $\langle \text{fn-plist} \rangle \rightarrow \epsilon$
- 42. $\langle \text{fn-plist1} \rangle \rightarrow , \langle \text{val-id} \rangle \langle \text{fn-plist1} \rangle$
- 43. $\langle \text{fn-plist1} \rangle \rightarrow \epsilon$
- 44. $\langle \text{val-id} \rangle \rightarrow \langle \text{id} \rangle$
- 45. $\langle \text{val-id} \rangle \rightarrow \text{int-literal}$
- 46. $\langle \text{val-id} \rangle \rightarrow \text{double-literal}$
- 47. $\langle \text{val-id} \rangle \rightarrow \text{string-literal}$
- 48. $\langle \text{val-id} \rangle \rightarrow \text{boolean-literal}$
- 49. $\langle \text{ret-val} \rangle \rightarrow \text{E}$
- 50. $\langle \text{ret-val} \rangle \rightarrow \epsilon$

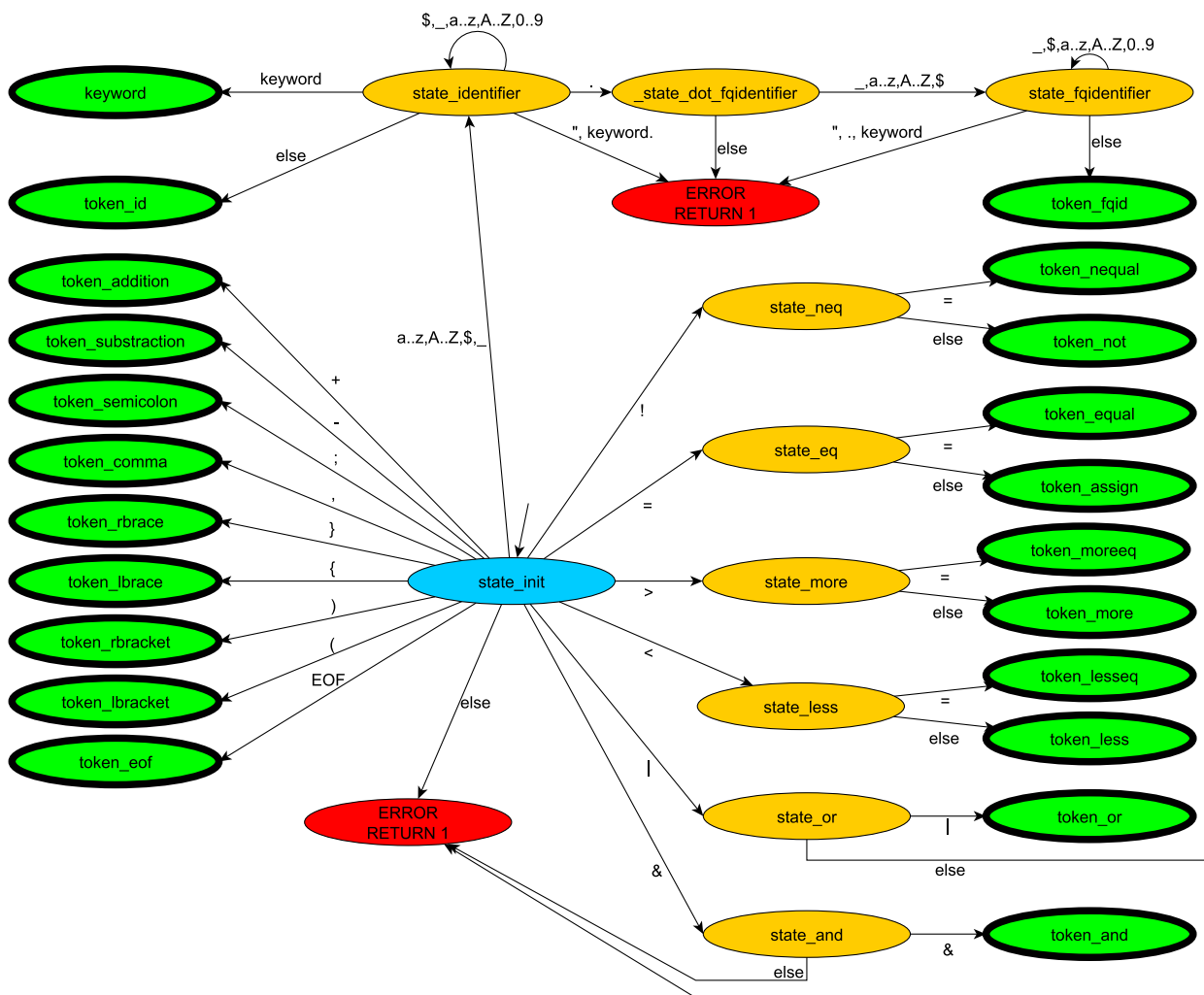
8.2 Precedenční tabulka

| | + | - | * | / | < | > | <= | >= | == | != | (|) | && | | ! | id | \$ |
|----|---|---|---|---|---|---|----|----|----|----|---|---|----|---|---|----|----|
| + | > | > | < | < | > | > | > | > | > | > | < | > | > | > | < | < | > |
| - | > | > | < | < | > | > | > | > | > | > | < | > | > | > | < | < | > |
| * | > | > | > | > | > | > | > | > | > | > | < | > | > | > | < | < | > |
| / | > | > | > | > | > | > | > | > | > | > | < | > | > | > | < | < | > |
| < | < | < | < | < | | | | | > | > | < | > | > | > | < | < | > |
| > | < | < | < | < | | | | | > | > | < | > | > | > | < | < | > |
| <= | < | < | < | < | | | | | > | > | < | > | > | > | < | < | > |
| >= | < | < | < | < | | | | | > | > | < | > | > | > | < | < | > |
| == | < | < | < | < | < | < | < | < | > | > | < | > | > | > | < | < | > |
| != | < | < | < | < | < | < | < | < | > | > | < | > | > | > | < | < | > |
| (| < | < | < | < | < | < | < | < | < | < | < | = | < | < | < | < | |
|) | > | > | > | > | > | > | > | > | > | > | | > | > | > | > | | > |
| && | < | < | < | < | < | < | < | < | < | < | < | > | > | > | < | < | > |
| | < | < | < | < | < | < | < | < | < | < | < | > | < | > | < | < | > |
| ! | > | > | > | > | > | > | > | > | > | > | < | > | > | > | < | < | > |
| id | > | > | > | > | > | > | > | > | > | > | | > | > | > | > | | > |
| \$ | < | < | < | < | < | < | < | < | < | < | < | | < | < | < | < | |

8.3 Diagram konečného automatu lexikální analýzy



Obrázek 1: Diagram konenečného automatu, část 1.



Obrázek 2: Diagram konenečného automatu, část 2.