

VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ
FAKULTA INFORMAČNÍCH TECHNOLOGIÍ



Dokumentace projektu
Interpret jazyka IFJ15

Tým 100, varianta a/2/II

Vedoucí: Tomáš Strnka (xstrnk01) 22%

Členové: Petr Žufan (xzufan00) 22%

Martin Wojaczek (xwojac00) 22%

Martin Kvapil (xkvapi13) 22%

Petr Gibiš (xgibis00) 12%

v Brně 28. listopadu 2015

Obsah

1 Úvod.....	3
2 Řešení projektu	3
2.1 Lexikální analýza	3
2.2 Syntaktická a sémantická analýza	3
2.3 Instrukce a interpret.....	3
3 Algoritmy a datové struktury	4
3.1 Heap sort	4
3.2 Knuth-Morris-Prattův algoritmus.....	4
3.3 Tabulka s rozptýlenými položkami	5
4 Práce v týmu	5
4.1 Rozdělení práce	5
5 Závěr	5
6 Literatura.....	6
7 Metriky kódu.....	6
A Konečný automat - Lexikální analyzátor	7
B Precedenční tabulka	8
C LL Gramatika.....	9

1 Úvod

Úkolem projektu bylo vytvořit funkční interpret jazyka IFJ15, který byl inspirován podmnožinou jazyka C++ verzí 2011. Dále v dokumentaci popisujeme základní principy použité k řešení.

Zvolili jsme si variantu a/2/II, což znamená vyhledávání podřetězce pomocí Knuth-Morris-Pratt algoritmu, řazení za pomoci Heap sort algoritmu a užití tabulky s rozptýlenými položkami.

2 Řešení projektu

Kapitola popisující jednotlivé části interpretu.

2.1 Lexikální analýza

Činnost, kterou provádí tzv. lexikální analyzátor (scanner) je první částí interpretu. Lexikální analýza se stará o odstranění bílých znaků a komentářů a rozdělení vstupního souboru na jednotky neboli tokeny. Identifikátory, klíčová slova, operátory, čísla, řetězce. Lexikální jednotky jsou dále poskytnuty k dalšímu zpracování syntaktickému analyzátoru.

Reprezentace „scanneru“ pomocí konečného automatu (viz příloha A). Nejdůležitější funkcí lexikální analýzy je *get_next_token()*, která opakovaným voláním navrácí postupně všechny tokeny ze vstupního souboru. Když se ve vstupním souboru nachází lexikální chyba je navrácen token *LEX_ERROR* a načítání je ukončeno.

2.2 Syntaktická a sémantická analýza

Syntaktický analyzátor (parser) je nejdůležitějším členem celého kompilátoru. Je založený na rekurzivním sestupu. Postupně volá lexikální analyzátor pomocí funkce *GetNextToken()* a navrácené tokeny zpracovává. Aby syntaktická analýza fungovala správně, je potřeba vytvořit vhodně gramatická pravidla. Tyto pravidla reprezentuje LL gramatika (viz příloha C). Další částí parseru je precedenční syntaktická analýza, která kontroluje správnou syntax výrazů, podle precedenční syntaktické tabulky (viz příloha B). Pokud během kontroly parser zjistí, že program nedopovídá navržené LL gramatice, tak okamžitě skončí s chybou a dále již nepokračuje.

Během rekurzivního sestupu a precedenční syntaktické analýze se provádí kontrola sémantiky. Kontrola se zaměřuje na kontrolu deklarace proměnných. Další kontrola sémantiky probíhá až v interpretu. Zároveň s kontrolou syntaxe se generují instrukce a ukládají se symboly do globální a lokální tabulky symbolů.

2.3 Instrukce a interpret

Nezbytnou částí interpretu je generování a vykonávání instrukcí z instrukční pásky. Interpret pracuje se zásobníkem proměnných a s instrukcemi v tří-adresném kódu. Každý operand obsahuje jméno proměnné, po jejímž vyhledání získáme údaje jako například datový typ.

Interpret je závěrečnou fází zpracování zdrojového programu. Následuje po úspěšném skončení syntaktické analýzy, která již generovala posloupnost instrukcí tří-adresného kódu. Interpret vykonává

ještě některé sémantické kontroly inicializace a kompatibility datových typů. Instrukční páska je implementována jako lineárně zřetězený seznam. Interpret vyjme z globální tabulky symbolů instrukční pásku funkce *main* a zpracovává instrukce jednu po druhé. Zároveň na zásobník vloží kopii lokální tabulky symbolů, do které ukládá hodnoty proměnných. Když dojde na volání funkce, vytvoří se nová kopie proměnných na vrcholu zásobníku a instrukční páska dané funkce se vloží do právě zpracovávané instrukční pásky. Návrátová hodnota funkce je vždy uložena v proměnné *#RETURN*.

3 Algoritmy a datové struktury

3.1 Heap sort

Neboli řazení hromadou. Podstatou řadící metody je implementace hromady polem. Je to implementace s implicitním zřetězením prvků binární stromové struktury hromady. Halda má rekurzivní vlastnost, že každý podstrom je taktéž haldou. Dále pro haldu platí, že kořenový prvek má vyšší hodnotu než jeho potomci. V reprezentaci je prvek označen indexem n a jeho potomci indexy $2n+1$ a $2n+2$. Funkce *siftdown* upravuje prvky haldy tak, že porovnává kořenový prvek haldy s jeho potomky a prohodí je v případě, že potomci mají vyšší hodnotu než kořenový prvek. Tato operace se opakuje i pro podhaldy které mají jako kořenový prvek potomky předchozí haldy. V případě výměny kořenového prvku je možné, že dojde k nerovnováze o úroveň méně, tehdy se *siftdown* se aplikuje i na danou úroveň. Takto program pokračuje, dokud nedojde ke kořenovému prvku bez potomků anebo konkrétní halda splňuje uspořádání. Program nakonec přečte prvky z uspořádané haldy.

Ukázka třídění haldy: 65843127

6	8
5 8	→ 7 6
4 3 1 2	5 3 1 2
7	4

V ukázce došlo k 3 krokům výměny. Prvek 8 a 6 byly prohozeny, a dále 7 a 4, kdy došlo k nerovnováze na vyšší úrovni, které vede k dalšímu prohození s prvkem 5.

3.2 Knuth-Morris-Prattův algoritmus

Pro vyhledávání podřetězce v řetězci při použití vestavěné funkce je použit Knuth-Morris-Prattův algoritmus. Ten vytváří masku hledaného podřetězce a určuje posun vyhledávání v případě neshody v porovnávání určitého písmene podřetězce s písmenem v prohledávaném textu.

Ukázka masky podřetězce:

1	2	3	4	5	6	7	8
A	B	A	C	A	B	D	B
0	0	1	0	1	2	0	0

Maska je pole indexů, které určuje na kterou pozici se má vyhledávací funkce vrátit v případě neshody. V našem případě (viz ukázka) je vidět význam tohoto algoritmu kdy se uvnitř podřetězce nachází stejná sekvence jako na začátku. Písmena s indexy 1, 2 a 5, 6 jsou stejné. Pokud nastane neshoda na indexu 7, pak víme, že předchozí znaky byly zkontrolovány a je tedy možné posunout rovnou na index 2. Následuje kontrola indexu 3, případně neúspěchu i tohoto indexu se podle masky vrátíme na index 0. Dochází ke zrychlení porovnávání ve smyslu, že díky přeskočení určitých částí textu, které jsou již zkontrolovány během předchozích kroků. Takto se podřetězec posunuje po délce celého textu, dokud nenarazí na konec řetězce. Časová složitost celého algoritmu je $O(m+n)$ díky vytvoření masky.

3.3 Tabulka s rozptýlenými položkami

Datová struktura tabulka s rozptýlenými položkami je použita pro uchování symbolů funkcí v překládaném kódu. Použita je tabulka s označením globální, pro symboly funkce, která obsahuje ukazatele na další tabulky, s označením lokální, pro symboly vně funkcí. Výhodou této datové struktury je rychlost vyhledávání symbolů v tabulce. Tabulka je reprezentována polem ukazatelů na jednosměrné seznamy. Každá položka obsahuje ukazatel na symbol a na další položku seznamu.

Pomocí hashovací funkce je pro každý symbol mapován klíčový řetězec na index, kterým je určeno rozložení v poli ukazatelů. Při vyhledávání se podle indexu určí seznam a ten sekvenčně se prohledá pro požadovaný symbol.

4 Práce v týmu

Před samotným zadáním projektu jsme si podle zadání z minulých let smluvili a rozdělili práci podle schopností a praxe každého z nás. Práce na projektu započaly koncem října, kdy jsme si založili Git repositář a byl zaznamenán první „commit“. Během řešení jsme se scházeli každý týden, abychom diskutovali o problémech jednotlivých částí a jak je následovně dát dohromady. Problémem byla velká časová náročnost projektu. Bojovali jsme do poslední chvíle.

4.1 Rozdělení práce

Tomáš Strnka	Vedoucí týmu, parser, dokumentace
Martin Wojaczek	Vyhodnocování výrazů, testování, parser
Petr Žufan	Interpret, struktury, testování
Petr Gibiš	Algoritmy, dokumentace
Martin Kvapil	Scanner, dokumentace

5 Závěr

Projekt se podařilo úspěšně dokončit a otestovat vzorovými zdrojovými kódy a několika vlastními. Nebylo to snadné a často jsme zažívali perné chvíle, obzvláště když ubíhaly poslední hodiny. Všichni to ale bereme jako velkou zkušenost pro naši budoucí kariéru.

6 Literatura

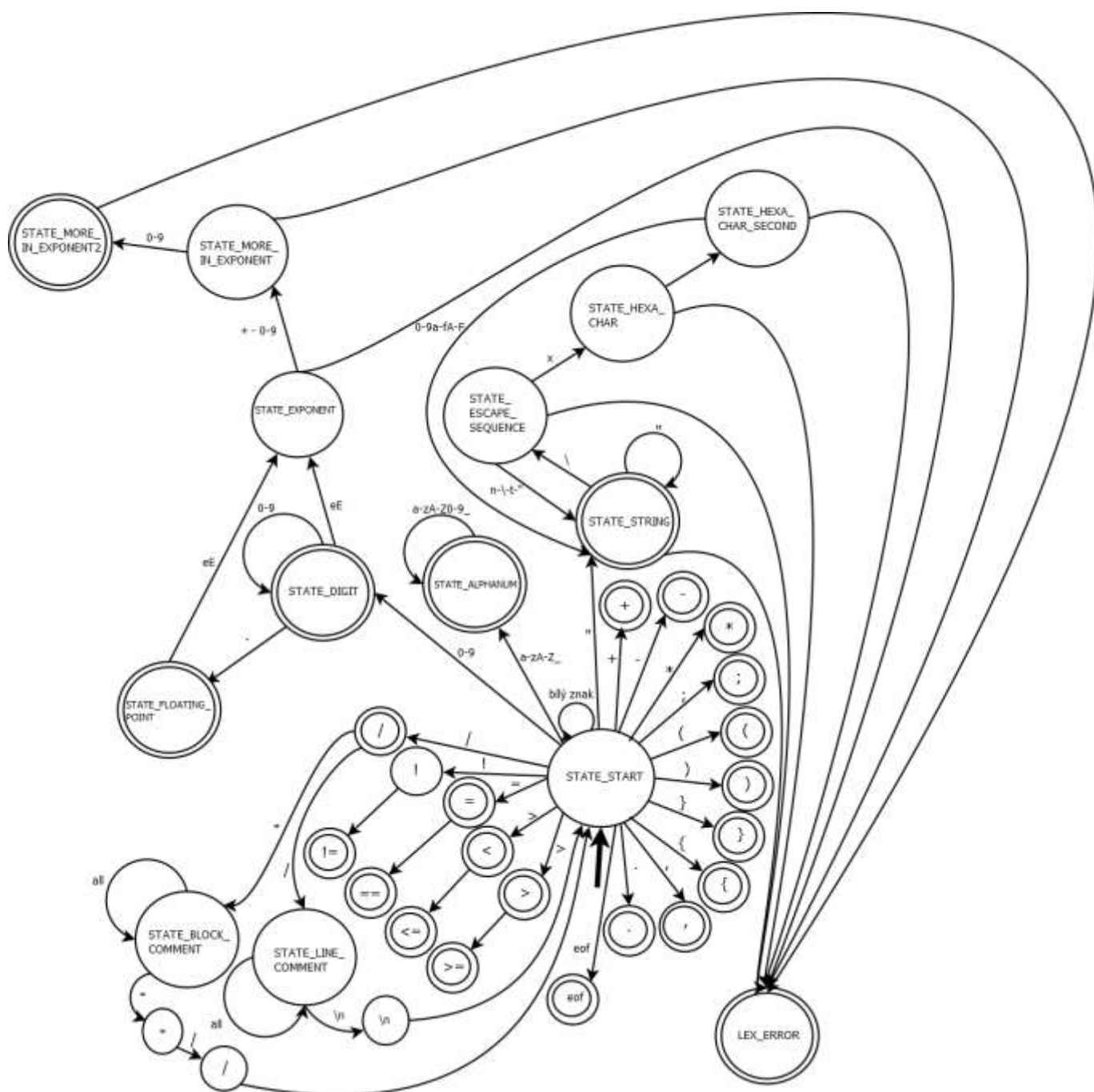
HONZÍK, Jan. Algoritmy – studijní opora. Brno: Vysoké učení technické, 2015.

7 Metriky kódu

Počet řádků zdrojového textu: 6065

Velikost spustitelného programu: 161 KiB

A Konečný automat - Lexikální analyzátor



B Precedenční tabulka

	+	-	*	/	<	>	<=	>=	==	!=	()	i	\$	s	d	f
+	>	>	<	<	>	>	>	>	>	>	<	>	<	>	<	<	<
-	>	>	<	<	>	>	>	>	>	>	<	>	<	>	0	<	<
*	>	>	>	>	>	>	>	>	>	>	<	>	<	>	0	<	<
/	>	>	>	>	>	>	>	>	>	>	<	>	<	>	0	<	<
<	<	<	<	<	>	>	>	>	>	>	<	>	<	>	<	<	<
>	<	<	<	<	>	>	>	>	>	>	<	>	<	>	<	<	<
<=	<	<	<	<	>	>	>	>	>	>	<	>	<	>	<	<	<
>=	<	<	<	<	>	>	>	>	>	>	<	>	<	>	<	<	<
==	<	<	<	<	>	>	>	>	>	>	<	>	<	>	<	<	<
!=	<	<	<	<	>	>	>	>	>	>	<	>	<	>	<	<	<
(<	<	<	<	<	<	<	<	<	<	<	=	<	0	<	<	<
)	>	>	>	>	>	>	>	>	>	>	0	>	0	>	0	0	0
i	>	>	>	>	>	>	>	>	>	>	0	>	0	>	0	0	0
\$	<	<	<	<	<	<	<	<	<	<	<	0	<	0	<	<	<
s	>	0	0	0	>	>	>	>	>	>	0	>	0	>	0	0	0
d	>	>	>	>	>	>	>	>	>	>	0	>	0	>	0	0	0
f	>	>	>	>	>	>	>	>	>	>	0	>	0	>	0	0	0

> – redukce, = – výjimka, < – handle, 0 – nedefinováno

C LL Gramatika

1. $\text{PROGRAM} \rightarrow \text{FUNC_DCLR}$
2. $\text{FUNC_DCLR} \rightarrow \text{TYPE id (PARAM) SELECT FUNC_DCLR}$
3. $\text{PARAM} \rightarrow \text{TYPE id PARAM_N}$
4. $\text{PARAM} \rightarrow \varepsilon$
5. $\text{PARAM_N} \rightarrow , \text{TYPE id PARAM_N}$
6. $\text{PARAM_N} \rightarrow \varepsilon$
7. $\text{SELECT} \rightarrow \text{BODY}$
8. $\text{SELECT} \rightarrow ;$
9. $\text{BODY} \rightarrow \{ \text{STMNT} \}$
10. $\text{STMNT} \rightarrow \text{BODY STMNT}$
11. $\text{STMNT} \rightarrow \text{IF STMNT}$
12. $\text{STMNT} \rightarrow \text{FOR STMNT}$
13. $\text{STMNT} \rightarrow \text{CIN STMNT}$
14. $\text{STMNT} \rightarrow \text{COUT STMNT}$
15. $\text{STMNT} \rightarrow \text{RETURN STMNT}$
16. $\text{STMNT} \rightarrow \text{PROM STMNT}$
17. $\text{FOR} \rightarrow \text{for (TYPE id I_PROM ; EXPR ; id = EXPR) STMNT}$
18. $\text{IF} \rightarrow \text{if (EXPR) STMNT else STMNT}$
19. $\text{CIN} \rightarrow \text{cin} \gg \text{id ID_N ;}$
20. $\text{COUT} \rightarrow \text{cout} \ll \text{TERM TERM_N ;}$
21. $\text{RETURN} \rightarrow \text{return EXPR ;}$
22. $\text{ID_N} \rightarrow \gg \text{id ID_N}$
23. $\text{ID_N} \rightarrow \varepsilon$
24. $\text{TERM_N} \rightarrow \ll \text{TERM TERM_N}$
25. $\text{TERM_N} \rightarrow \varepsilon$
26. $\text{I_PROM} \rightarrow = \text{EXPR}$
27. $\text{I_PROM} \rightarrow \varepsilon$
28. $\text{PROM} \rightarrow \text{id = CALL_DEC ;}$
29. $\text{PROM} \rightarrow \text{TYPE id I_PROM ;}$
30. $\text{CALL_DEC} \rightarrow \text{id (LIST_PAR)}$

31. $\text{CALL_DEC} \rightarrow \text{EXPR}$

32. $\text{LIST_PAR} \rightarrow \text{TERM LIST_PAR_N}$

33. $\text{LIST_PAR_N} \rightarrow , \text{TERM LIST_PAR_N}$

34. $\text{LIST_PAR_N} \rightarrow \varepsilon$

35. $\text{LIST_PAR} \rightarrow \varepsilon$

36. $\text{TYPE} \rightarrow \text{int}$

37. $\text{TYPE} \rightarrow \text{string}$

38. $\text{TYPE} \rightarrow \text{double}$

39. $\text{TYPE} \rightarrow \text{auto}$

40. $\text{TERM} \rightarrow \text{term}$