



Dokumentácia k projektu z predmetov IFJ a
IAL

Implementace překladače imperativního jazyka IFJ22

Tým xbucka00, varianta TRP

Timotej Bučka	xbucka00	25%
Daniel Mačura	xmacur09	25%
Tobiáš Štec	xstect00	25%
Patrik Frůstök	xfrust00	25%

1.december.2022

Obsah

Spôsoby riešenia jednotlivých častí projektu	2
1. Lexikálna analýza	3
1.1. Konečný stavový automat lexikálnej analýzy	4
2. Syntaktická a sémantická analýza	6
2.1. Implementačné detaily	6
2.2. Gramatické pravidla	8
2.3. LL(1) tabuľka	9
2.4. Spracovanie výrazov	10
2.5. Precedenčná tabuľka pre výrazy	11
2.6. Zjednodušená precedenčná tabuľka pre výrazy	12
3. Tabuľka symbolov	13
4. Generovanie kódu	13
5. Vlastné telo prekladača	14
Práce v tímu	15
6. Rozdelenie práce medzi členmi tímu	15

Spôsoby riešenia jednotlivých častí projektu

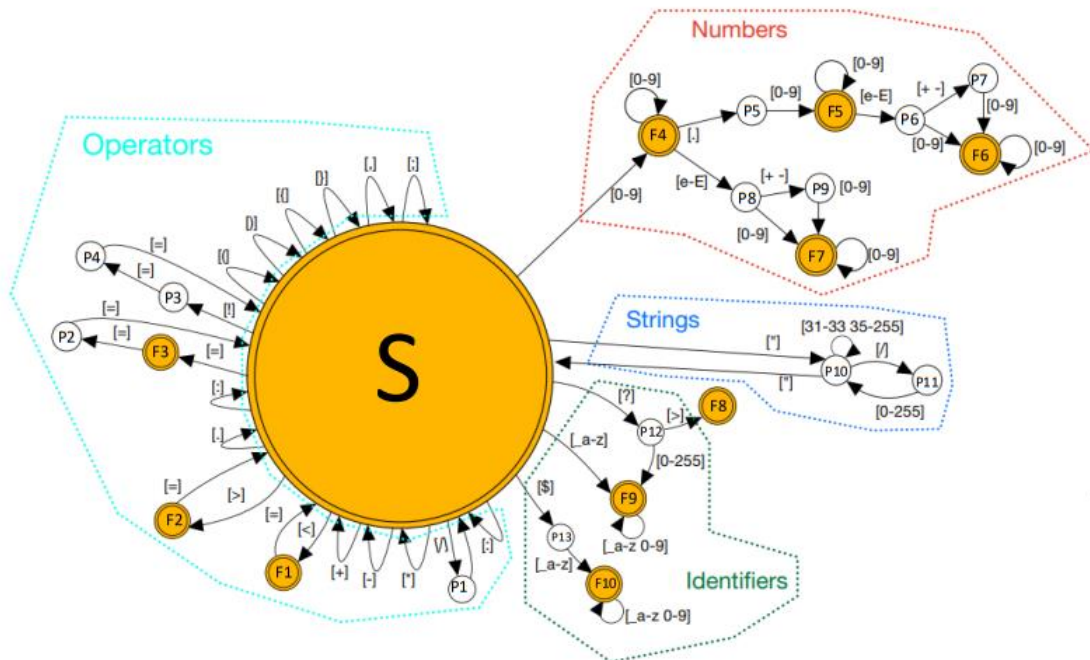
1. Lexikálna analýza

Lexikálnu analýzu sme implementovali s pomocou podporných štruktúr a funkcií zo súborov ***lexer.c*** a ***lexer.h***. V súbore ***token.h*** máme definovanú dátovú štruktúru `token`, ktorá sa skladá z typu (`ID`) a atributu hodnoty (`VAL`), ktorá uchováva reťazcovú reprezentáciu daného tokenu, resp. identifikátor kľúčového slova.

Hlavnou funkciou našej lexikálnej analýzy je `lexer_next_token`, ktorá je volaná z makra `next_tok`. V tomto makre sa získaný token pridá do nášho dvojito viazaného zoznamu `dll`. Funkcia `lexer_next_token`, načíta zo súboru `token` po tokene, a postupne každému pridá jeho typ (`ID`), na základe deterministického konečného automatu (**viz. 1.1**) a jeho stavov, implementovaného priamo vo funkcii ako `switch`, ktorý sa opakuje pokiaľ zo vstupu nie je načítaný znak `EOF`. Jednotlivé prípady `case` predstavujú už samotné stavy automatu. Ak jeden z načítaných znakov nie je validný, program vráti hodnotu `LEXICAL_ERR` (`==1`) a ukončí proces. Inak sa prechádza do ďalších stavov a načítajú sa ďalšie znaky, až do chvíle kedy už máme hotový ďalší token, ktorý vraciame a ukončujeme túto funkciu.

Po načítaní tokenu sa prevedie kontrola, či sa nejedná o kľúčové slovo alebo identifikátor už predom známy v jazyku IFJ22.

1.1 Konečný stavový automat lexikálnej analýzy



Legend:

S - STATE_START

F1 - STATE STATE

F2 - STATE GT E

F3 - STATE_EQ_E

F4 - STATE_INTEGER_E

F5 - STATE_DOUBLE_E

F6 - STATE_DOUBLE_EXPONENT_E

F7 - STATE INTEGER EXPONENT E

F8 - STATE CLOSING TAG

F9 - STATE_IDENTIFIER_OR_KEYWORD_E

F10 - STATE_VARIABLE_E

P1 - STATE SLASH

P2 - STATE_EQEQ

P3 - STATE_NOT

P4 - STATE_NOT_EQ

P5 - STATE INTEGER AND SEPARATOR

P6 - STATE DOUBLE EXPONENT START

P7- STATE DOUBLE EXPONENT SIGN

P8 - STATE INTEGER EXPONENT STAR

P9 - STATE INTEGER EXPONENT SIGN

P10 - STATE QUOTATION CENTER E

P11 - STATE QUOTATION ESC

P12 - STATE QUESTIONMARK

P13 - STATE_VARIABLE_START

2. Syntaktická a sémantická analýza

Najdôležitejšiou časťou celého programu je syntaktická analýza.

2.1 Implementačné detaily

Hlavnú kostru syntaktickej a sémantickej analýzy tvoria funkcie definované v súbore ***analysis.c***.

Vo funkcii `run_analysis` dochádza k inicializácii globálnej a lokálnej tabuľky symbolov implementovanými ako TRP. Do globálnej tabuľky sú vo funkcii `preload_hash_table` postupne vkladané preddefinované funkcie.

Zvolili sme metódu rekurzívneho zostupu. Derivačné kroky simulujeme pomocou rekurzívneho volania funkcií na základe jednotlivých gramatických pravidiel. Za pomoci backtrackingu, ktorý je možné realizovať vďaka uloženiu tokenov v dvojito spojenom zozname `d11`, rozhodneme, ktoré gramatické pravidlo bude použité. Každá funkcia volá makro `next_tok`, ktoré sa rozšíri do volania funkcie `lexer_next_token` v situácii, kedy je aktívny element rovný poslednému elementu `d11`, inak je aktivita posunutá na ďalší element v `d11` a jeho hodnota je navrátená. Na základe získaného typu tokenu, daná funkcia rozhoduje o syntaktickej validite danej postupnosti tokenov. V prípade, že je očakávaný výraz, je zavolaná funkcia `parse_expression`, ktorá má na starosti syntaktickú kontrolu výrazov. Ak sa analýza nachádza v stave, kde je iný rámec (scope) ako rámec hlavnej funkcie, je inkrementovaná globálna premenná `BODYRECURSIONCOUNT`, ktorá zaistí, že sa bude pracovať s lokálnou tabuľkou symbolov miesto globálnej.

Pri deklarácii funkcie je do globálnej tabuľky symbolov pridaný prvok reprezentujúci danú funkciu s názvom, parametrami a

návratovým typom funkcie. To isté platí aj pre premenné, kde je daná premenná pridávaná vždy do správnej tabuľky symbolov (globálna alebo lokálna) určenej stavom premennej BODYRECURSIONCOUNT.

Za priebehu syntaktickej analýzy kontrolujeme jednotlivé sémantické pravidlá a predpoklady. Funkcia `parse_expression` prijíma ako jeden z parametrov ukazateľ na `data_type`, ktorého hodnotu mení na výsledný typ výrazu. Vo vnútri funkcie je kontrolovaná typová kompatibilita operátorov a pri výskyte premennej vo výraze dochádza k vyhľadaniu v tabuľke symbolov pomocou `hash_table_lookup`. Pri syntaktickej kontrole volania funkcie je zavolaná funkcia `hash_table_lookup` nad globálnou tabuľkou symbolov, a je skontrolovaný správny počet a typ argumentov.

Syntaktická analýza končí neúspechom, ak je niektoré z pravidiel gramatiky porušené, a správna chybová hodnota je zapísaná do globálnej premennej `ERROR`. Hodnota `ERROR` je vrátená hlavnou funkciou `main` a v prípade že analýza prebehne bez vzniku syntaktických alebo sémantických chýb je vrátená hodnota `SUCCESS (0)`.

2.2 Gramatické pravidla

- (1) $BODY \rightarrow \varepsilon$
- (2) $BODY \rightarrow \text{"function" func_id "(" PARAMETERS ")" ":" TYPE "{" BODY "}" BODY}$
- (3) $BODY \rightarrow \text{"if" "(" EXPRESSION ")" "{" BODY "}" "else" "{" BODY "}" BODY}$
- (4) $BODY \rightarrow \text{"while" "(" EXPRESSION ")" "{" BODY "}" BODY}$
- (5) $BODY \rightarrow \text{"return" RETURN ";" BODY}$
- (6) $BODY \rightarrow \text{FUNCTION_CALL ";" BODY}$
- (7) $BODY \rightarrow \text{ASSIGNMENT ";" BODY}$
- (8) $BODY \rightarrow \text{EXPRESSION ";" BODY}$
- (9) $ASSIGNMENT \rightarrow \text{var_id "=" ASSIGNMENT_PRIME}$
- (10) $ASSIGNMENT_PRIME \rightarrow \text{EXPRESSION}$
- (11) $ASSIGNMENT_PRIME \rightarrow \text{FUNCTION_CALL}$
- (12) $FUNCTION_CALL \rightarrow \text{func_id "(" ARGUMENTS ")"}$
- (13) $RETURN \rightarrow \varepsilon$
- (14) $RETURN \rightarrow \text{EXPRESSION}$
- (15) $RETURN \rightarrow \text{FUNCTION_CALL}$
- (16) $TYPE \rightarrow \text{PREFIX C_TYPE}$
- (17) $PREFIX \rightarrow \varepsilon$
- (18) $PREFIX \rightarrow \text{"?"}$
- (19) $C_TYPE \rightarrow \text{"int"}$
- (20) $C_TYPE \rightarrow \text{"float"}$
- (21) $C_TYPE \rightarrow \text{"string"}$
- (22) $C_TYPE \rightarrow \text{"void"}$
- (23) $ARGUMENTS \rightarrow \varepsilon$
- (24) $ARGUMENTS \rightarrow \text{EXPRESSION ARGUMENTS_PRIME}$
- (25) $ARGUMENTS_PRIME \rightarrow \varepsilon$
- (26) $ARGUMENTS_PRIME \rightarrow \text{"," EXPRESSION ARGUMENTS_PRIME}$
- (27) $PARAMETERS \rightarrow \varepsilon$
- (28) $PARAMETERS \rightarrow \text{C_TYPE var_id PARAMETERS_PRIME}$
- (29) $PARAMETERS_PRIME \rightarrow \varepsilon$
- (30) $PARAMETERS_PRIME \rightarrow \text{"," C_TYPE var_id PARAMETERS_PRIME}$

2.3 LL(1) tabulka

[illegible]

2.4 Spracovanie výrazov

Spracovanie výrazov sa deje oddelene od spracovania syntaxe, a to precedenčnou analýzou. Funkcie starajúce sa o toto spracovanie nájdeme v súbore ***expressions.c***. Hlavnou funkciou je funkcia `parse_expression`, ktorá je volaná hlavným parserom v prípade, že parser narazí na výraz.

Dôležitou časťou prec. analýzy je zásobník reprezentovaný štruktúrou `expr_stack` pozostávajúcou z dĺžky zásobníka a zo štruktúry `expr_item`, ktorá reprezentuje jednotlivé prvky zásobníka. `expr_item` môže reprezentovať tri typy prvkov vyskytujúcich sa v precedenčnej analýze: znak dolár (znak neznámy v kontexte spracovania výrazu), terminál alebo neterminál.

Princíp fungovania precedenčnej analýzy je nasledovný: Na začiatku `parse_expression` je do zásobníku vložený znak dolár. Následne je cykle volané makro `next_tok` a na základe získaného tokenu je podľa precedenčnej tabuľky určené, ako sa má v analýze pokračovať. V tabuľke sa nachádzajú symboly:

- = - Vlož získaný token na zásobník a načítaj ďalší token
- < - Nastav najvrchnejšiemu termu na zásobníku príznak `breakpoint`, vlož získaný token na zásobník a načítaj ďalší token
- > - Vykonaj redukciu podľa známych pravidiel
- Prázdne políčko – chyba pri syntaktickej analýze

Analýza končí úspechom v prípade že načítaný token je interpretovaný ako dolár a na zásobník je v stave: dolár, neterminál (vrchol).

2.5 Precedenční tabulka výrazov

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

2.6 Zjednodušená precedenční tabulka výrazov

	\$	()	+ -	* /	.	rop	id
\$		<		<	<	<	<	<
(<	=	<	<	<	<	<
)	>		>	>	>	>	>	
+ -	>	<	>	>	<		>	<
* /	>	<	>	>	>		>	<
.	>	<	>			>	>	<
rop	>	<	>	<	<	<		<
id	>		>	>	>	>	>	

3. Tabuľka symbolov

V súboroch ***symtable.c*** a ***symtable.h*** sme implementovali tabuľku s rozptýlenými položkami (TRP) ako variantu zadania.

TRP sme sa rozhodli implementovať formou explicitného zreťazenia. Synonymá sme reťazili pomocou lineárne viazaných zoznamov.

Ako veľkosť mapovacieho poľa sme zvolili číslo prvočíslo **49 157**, kvôli jeho vhodným vlastnostiam¹ a použili ***polynomial rolling hash*** mapovaciu funkciu. Generovaný hash sme zkomprimovali za pomoci ***MAD compression***², aby sa s ním dalo indexovať v tabuľke.

Tabuľka obsahuje položky typu `table_item_data`, ktoré reprezentujú funkciu alebo premennú. Taktiež obsahujú názov, ktorý sa aj vyžíva pre gerenovanie hash funkcie

4. Generovanie kódu

Generovanie cieľového kódu **IFJcode22**, sme implementovali pomocou súboru ***generator.c***, ktorý je volaný v priebehu syntaktickej analýzy za chodu programu.

Ak sa token zhoduje s jednou z funkcií jazyka ***IFJ22***, zavolá sa príslušná funkcia definovaná ako makro, obsahujúca stringovú reprezentáciu danej funkcie v ***IFJcode22***. Okrem funkcií generujeme aj výstupný kód pre príkazy známe v jazyku ***IFJ22*** (napr. `if`, `while`, príkazy pre matematické operácie...).

Na záver, na základe vopred zadaného kódu, program vygeneruje a vypíše jeho trojadresný kód na štandardný výstup.

¹ <https://planetmath.org/goodhashtableprimes>

² <http://www.cs.emory.edu/~cheung/Courses/253/Syllabus/Map/hash-func.html>

5. Vlastné telo prekladača

V hlavnej funkcii **main** inicializujeme potrebné dátové štruktúry a ďalej voláme funkciu **run_analysis**, ktorá ďalej riadi činnosť prekladača. Na konci je vrátený správny kód chyby a uvoľnené predtým alokované miesto.

Práca v tíme

Na zložení tímu sme sa dohodli ešte počas prvého týždňa semestra, nakoľko sme sa už všetci poznali.

Na projekte sme začali pracovať ku koncu septembra. Z počiatku sme sa stretávali všetci spolu aby sme projektu a jednotlivým častiam pochopili všetci a mohli si tak rozdeliť našu prácu.

Popri súčasným prednáškam nám veľmi pomohli aj záznamy z minulého roku, vďaka ktorým sme mohli rýchlejšie postupovať.

Po spoločnom zhodnotení a pochopení častiam projektu sme si postupne začali prácu deliť a pracovať jednotlivo prípadne vo dvojiciach.

Prácu sme si rozdelili rovnomerne, každý člen teda dostal percentuálne hodnotenie 25%.

Po dokončení jednotlivých častí, sme si otestovali naše riešenie pokusným odovzdaním, po ktorom sme pokračovali v opravách chýb. Nakoniec sme sa všetci spojili a vypracovali ku projektu túto dokumentáciu.

6. Rozdelenie práce medzi členmi tímu

6.1 Timotej Bučka (xbucka00)

- Dokumentácia, syntaktická analýza, sémantická analýza, analýza výrazov

6.2 Patrik Früštök (xfrustoo)

- Dokumentácia, lexikálna analýza, tvorba testov, pridávanie komentárov

6.3 Daniel Mačura (xmacuro9)

- Dokumentácia, syntaktická analýza, sémantická analýza, generácia kódu

6.4 Tobiáš Štec (xstectoo)

- Dokumentácia, lexikálna analýza, syntaktická analýza, generácia kódu