

Projektová dokumentace
Implementace překladače jazyka IFJ20
Formální jazyky a překladače

Tým 002, varianta 2

Jan Polišenský	xpolis04	40 %
Nina Štefeková	xstefe11	30 %
Martin Benovič	xbenov00	20 %
Sabína Švecková	xsveck00	10 %

8. prosince 2020

Obsah

1	Úvod	2
2	Rozdělení práce v týmu	2
2.1	Odůvodnění odchylek od rovnoměrnosti bodů	2
3	Implementace	2
3.1	Lexikální analýza	3
3.2	Syntaktická a sémantická analýza	4
3.3	Zpracování výrazů	6
3.4	Tabulka symbolů	7
3.5	Generování kódu	7

1 Úvod

Cílem tohoto projektu bylo vytvořit překladač v jazyce C zpracovávající jazyk IFJ20, který vychází z programovacího jazyka GO. Tento překladač jsme implementovali jako soubor kooperujících nezávislých modulů lexikální analýzy, syntaktické analýzy, sémantické analýzy a generování kódu.

Ačkoli je tato míra nezávislost jednotlivých modulů netradiční, umožňuje pohodlnou změnu funkcionality jedné části bez zásahu do jiné. Toho se v průběhu implementace ukázalo jako značná výhoda oproti klasické jednopružodové či dvouprůžodové implementaci.

2 Rozdělení práce v týmu

Jan Polišenský - návrh architektury, parser, generátor kódu, zpracování výrazů, tvorba gramatiky

Martin Benovič - dokumentace, testování, LL gramatika

Nina Štefeková - scanner, kód vestavěných funkcí

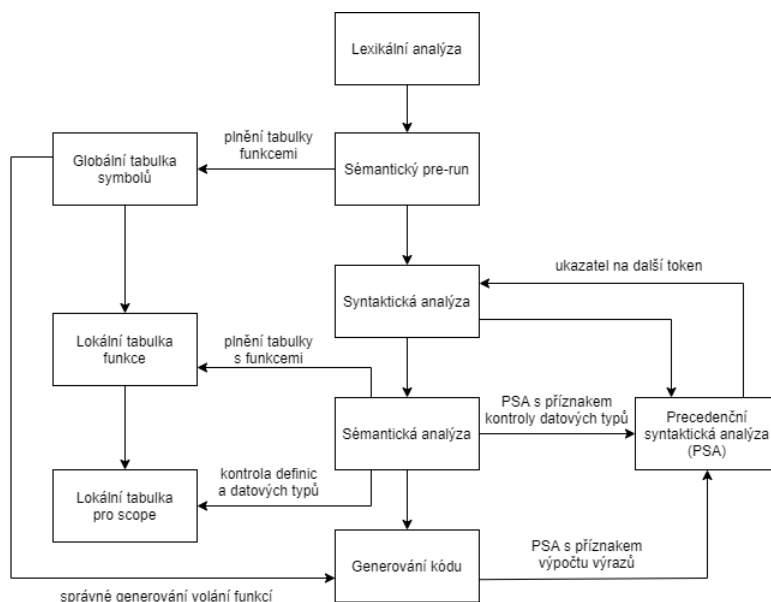
Sabína Švecková - scanner

2.1 Odůvodnění odchylek od rovnoměrnosti bodů

Složitá komunikace a absence prezenčních schůzek zapříčinily nerovnoměrné rozdělení v odpracovaných částech.

3 Implementace

Výsledná implementace vychází z níže uvedeného diagramu, vstupní kód je postupně sekvenčně zpracováván jednotlivými částmi našeho překladače. Hlavními bloky jsou pak lexikální analýza, sémantický prerun, syntaktická analýza a výsledné generování kódu.



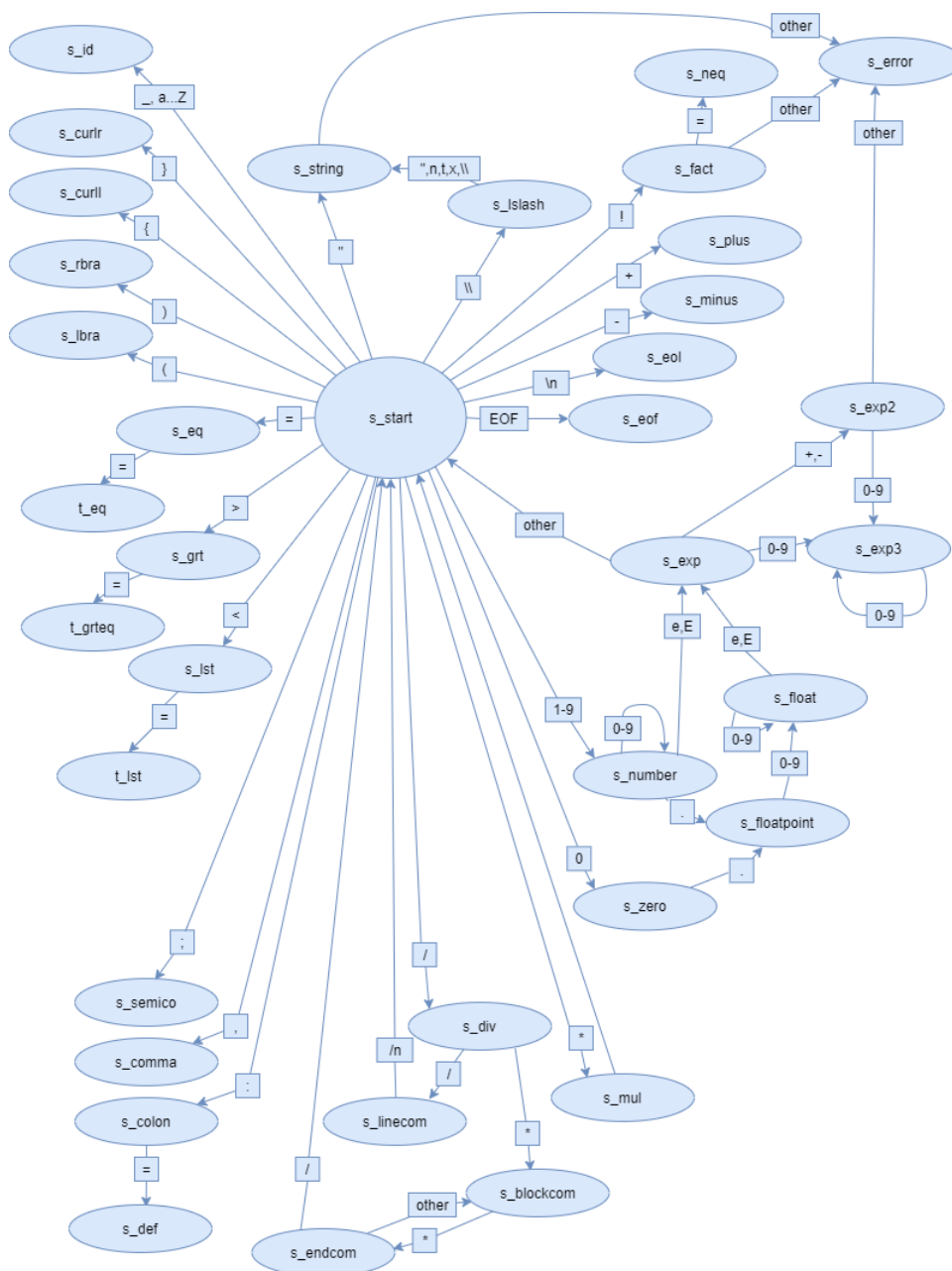
Obrázek 1: Schéma překladače

3.1 Lexikální analýza

Prvním krokem při tvorbě překladače byla implementace lexikální analýzy. Lexikální analyzátor jsme implementovali jako deteministický konečný automat. Před implementací jsme si automat navrhli, ale při samotné implementaci jsme ještě doplnili další stavy.

Automat jsme implementovali v jazyku C jako nekonečný **switch**, kde každý **case** odpovídá jednomu stavu automatu. Tento automat rozděluje zdrojový kód na jednotlivé tokeny. Tokeny jsou různé symboly, operátory, čísla, řetězce, identifikátory a klíčová slova. Pokud narazí na lexikální chybu, vrátí token **t_error**, který vyvolává chybu 1 při načítání tokenů v první části syntaktické analýzy.

Níže je pak uveden graf konečného automatu.



Obrázek 2: Schéma konečného automatu

3.2 Syntaktická a sémantická analýza

Před samotným spuštěním syntaktické analýzy je prováděn sémantický prerun. Jsou zde načítány tokeny ze scanneru a následně jsou zřetězovány do dvojsměrně vázaného lineárního seznamu, toto provázání tokenů nám umožňuje se v průběhu překladu libovolně vracet či posouvat. Jsou zde také prováděny jednoduché sémantické akce a plnění globální tabulky symbolů. Tato funkce se nachází v souboru `preprocesor.c`.

Syntaktická a sémantická analýza je pak jádro celého překladače a nachází se v souboru `parser.c`. Syntaktická analýza je implementována metodou rekurzivního sestupu který se řídí LL-gramatikou. Sémantické akce jsou pak prováděny až po dokončení syntaktické analýzy a jsou realizovány jako samostatný běh.

Během sémantické analýzy jsou pro každou funkci a blok kódu vytvářeny lokální tabulky symbolů pomocí kterých jsou kontrolovány korektní datové typy. Funkce, jejich parametry a návratové hodnoty jsou pak uloženy v globální tabulce symbolů

Níže je pak uvedena LL-tabulka a LL-gramatika podle níž je implementovaná syntaktická analýza

	EOL	KEYWORD	EOF	COMMA	ID	LBRA	CHECH_FOR_DEF_FUNCTION	EQ	ASSIGN	EXPR	ε
<n_prog>	1	2									
<n_func>	3	5	4								
<n_param_n>				7	8						6
<n_retvals>						10					9
<n_retval>				12							11
<n_fretval>	14			13							15
<n_body>	16, 19	19, 21, 22			19, 20		17, 19				18, 19
<n_body_id>				25				23	24		
<n_body_comma>									26		
<n_body_id_var>	29			28			27			28	
<n_fretvals>	30	31									32
<n_expr>				34						33	
<n_func_call>							35				
<n_call_param>				37, 38						38	36
<n_if>		39									
<n_for>		40									
<n_def>					42						41
<n_assign>					44						43

Tabulka 1: LL tabulka

1. $\langle n_prog \rangle \rightarrow EOL \langle n_prog \rangle$
2. $\langle n_prog \rangle \rightarrow KEYWORD \ ID \ \langle n_func \rangle$
3. $\langle n_func \rangle \rightarrow EOL \langle n_func \rangle$
4. $\langle n_func \rangle \rightarrow EOF$
5. $\langle n_func \rangle \rightarrow KEYWORD \ ID \ LBRA \ \langle n_param_n \rangle \ RBRA \ \langle n_retvals \rangle \ CURLL \ \langle n_body \rangle$
 $\langle n_retvals \rangle \ CURLR \ \langle n_func \rangle$
6. $\langle n_param_n \rangle \rightarrow \varepsilon$
7. $\langle n_param_n \rangle \rightarrow COMMA \ ID \ KEYWORD \ \langle n_param_n \rangle$
8. $\langle n_param_n \rangle \rightarrow ID \ KEYWORD \ \langle n_param_n \rangle$
9. $\langle n_retvals \rangle \rightarrow \varepsilon$
10. $\langle n_retvals \rangle \rightarrow LBRA \ KEYWORD \ \langle n_retval \rangle \ RBRA$
11. $\langle n_retval \rangle \rightarrow \varepsilon$
12. $\langle n_retval \rangle \rightarrow COMMA \ KEYWORD \ \langle n_retval \rangle$
13. $\langle n_fretval \rangle \rightarrow COMMA \ \langle n_expr \rangle \ \langle n_fretval \rangle$
14. $\langle n_fretval \rangle \rightarrow EOL \ \langle n_fretval \rangle \ EOL$
15. $\langle n_fretval \rangle \rightarrow \varepsilon$
16. $\langle n_body \rangle \rightarrow EOL \ \langle n_body \rangle \ EOL$
17. $\langle n_body \rangle \rightarrow \langle n_func_call \rangle \ EOL \ \langle n_body \rangle$
18. $\langle n_body \rangle \rightarrow \varepsilon$
19. $\langle n_body \rangle \rightarrow \langle n_body \rangle \ EOL \ \langle n_expr \rangle$
20. $\langle n_body \rangle \rightarrow ID \ \langle n_body_id \rangle \ EOL \ \langle n_body \rangle$
21. $\langle n_body \rangle \rightarrow \langle n_if \rangle \ EOL \ \langle n_body \rangle$
22. $\langle n_body \rangle \rightarrow \langle n_for \rangle \ EOL \ \langle n_body \rangle$
23. $\langle n_body_id \rangle \rightarrow EQ \ \langle n_body_id_var \rangle$
24. $\langle n_body_id \rangle \rightarrow ASSIGN \ \langle n_body_id_var \rangle$
25. $\langle n_body_id \rangle \rightarrow COMMA \ ID \ \langle n_body_id \rangle$
26. $\langle n_body_comma \rangle \rightarrow ASSIGN \ \langle n_body_id_var \rangle$
27. $\langle n_body_id_var \rangle \rightarrow \langle n_func_call \rangle$
28. $\langle n_body_id_var \rangle \rightarrow \langle n_expr \rangle$
29. $\langle n_body_id_var \rangle \rightarrow EOL$
30. $\langle n_fretvals \rangle \rightarrow EOL$
31. $\langle n_fretvals \rangle \rightarrow KEYWORD \ \langle n_expr \rangle \ \langle n_fretval \rangle$
32. $\langle n_fretvals \rangle \rightarrow \varepsilon$
33. $\langle n_expr \rangle \rightarrow expr$
34. $\langle n_expr \rangle \rightarrow COMMA \ expr$
35. $\langle n_func_call \rangle \rightarrow CHECK_FOR_DEF_FUNCTION \ LBRA \ \langle n_call_param \rangle \ RBRA$
36. $\langle n_call_param \rangle \rightarrow \varepsilon$
37. $\langle n_call_param \rangle \rightarrow COMMA \ \langle n_expr \rangle \ \langle n_call_param \rangle$
38. $\langle n_call_param \rangle \rightarrow \langle n_expr \rangle \ \langle n_call_param \rangle$
39. $\langle n_if \rangle \rightarrow KEYWORD \ \langle n_expr \rangle \ CURLL \ \langle n_body \rangle \ CURLR \ KEYWORD \ CURLL \ \langle n_body \rangle \ CURLR$
40. $\langle n_for \rangle \rightarrow KEYWORD \ \langle n_def \rangle \ SEMICO \ \langle n_expr \rangle \ SEMICO \ \langle n_assign \rangle$
 $CURLL \ \langle n_body \rangle \ CURLR$
41. $\langle n_def \rangle \rightarrow \varepsilon$
42. $\langle n_def \rangle \rightarrow ID \ ASSIGN \ \langle n_expr \rangle$
43. $\langle n_assign \rangle \rightarrow \varepsilon$
44. $\langle n_assign \rangle \rightarrow ID \ EQ \ \langle n_expr \rangle$

Tabulka 2: LL gramatika

3.3 Zpracování výrazů

Výrazy jsou zpracovávány pomocí precedenční syntaktické analýzy, jejíž implementace se nachází v souborech `psa_prerun.c` a `PSA.c`. PSA je v naší implementaci klíčovým komponentem, neboť zpracování výrazu je třeba jak na syntaktické, tak i na sémantické úrovni, PSA je také nutné na korektní generování kódu výrazů. Z těchto požadavků pak plyne poněkud obsáhlejší implementace PSA, kterou bylo nutné ještě doplnit preprocesorem výrazů umístěným v souboru `psa_prerun.c`. PSA je tedy možné spouštět ve 3 režimech:

1. kontrola syntaxe výrazů
2. kontrola sémantiky výrazů
3. generování kódu z výrazů

PSA pro správné zpracování výrazů využívá vlastní zásobník a lokální tabulky symbolů. Níže je pak uvedena precedenční tabulka symbolů využitá v PSA.

	+	-	*	/	R	()	id	int	float	str	\$
+	>	>	<	<	<	<	>	<	<	<	<	>
-	>	>	<	<	<	<	>	<	<	<	<	>
*	>	>	>	>	>	<	>	<	<	<	<	>
/	>	>	<	>	>	<	>	<	<	<	<	>
R	<	<	<	<		<	>	<	<	<	<	>
(<	<	<	<	<	<		<	<	<	<	
)	>	>	>	>	>		>					>
id	>	>	>	>	>		>					>
int	>	>	>	>	>		>					>
float	>	>	>	>	>		>					>
str	>	>	>	>	>		>					>
\$	<	<	<	<	<	<		<	<	<	<	

Tabulka 3: Precedenční tabulka

3.4 Tabulka symbolů

Tabulky symbolů jsou implementovány jako tabulky s rozptýlenými položkami s zřetěženými synonymy. Velikost mapovacího pole jsme dimenzovali aby nikdy nepřesáhla polovinu a je rovna prvočíslu 27457. Pro zřetězení synonym jsou využité lineární seznamy.

Funkce pro práci s tabulkou se nacházejí v souboru `symtable.c`. (kopírování tabulky, kontrola alokace tabulky, dealokace tabulky atd.).

3.5 Generování kódu

Hlavní funkce pro generování kódu se nachází v souboru `codegen.c`. Samotné generování pak probíhá v několika krocích. Nejprve je vygenerovaný kód vestavěných funkcí, tzv "base code", nachází se v souboru `code.c`.

Vstupní program je poté po jednotlivých řádcích procházen zdrojový kód a na základě prvního nebo i následujícího tokenu na každém řádu jsou volány jednotlivé funkce na generování konkrétních částí kódu (definice funkce, volání funkce). K vygenerování kódu z výrazů je pak použit precedenční syntaktický analyzátor s příznakem počítání výrazů. Ten k výpočtu výrazů využívá datový zásobník interpretu. Kvůli netradiční architektuře nastává v generátoru kódu problém s kontrolou platnosti jednotlivých proměnných v konkrétních blocích. Nejsou již totiž k dispozici lokální tabulky symbolů. Tento problém je řešen pomocí zásobníku na jehož vrchol se ukládá aktuálně platná proměnná daného jména.