



Implementace překladače imperativního jazyka IFJ23

Tým xjakub41, Varianta: TRP-izp

Milan Jakubec (xjakub41) 25%

Jir Tesar (xtesar44) 25%

Norman Babiak (xbabia01) 25%

Jakub Racek (xracek12) 25%

Brno, 6. prosince 2023

Obsah

1	Rozdělení práce	2
2	Implementace	3
2.1	Lexikální analýza	3
2.2	Syntaktická analýza	3
2.3	Semantická analýza	4
2.4	Analýza výrazu	4
2.5	Generování kódu	5
3	Datové struktury	5
3.1	Tabulka symbolů	5
3.2	Zásobník	5
3.2.1	Zásobník pro analýzu výrazu	5
3.2.2	Zásobník pro tabulku symbolů	6
3.3	Dynamické pole	6
4	Závěr	7
5	FSM pro lexikální analýzu	8
5.1	Diagram	8
6	LL	9
6.1	LL gramatika	9
6.2	LL tabulka	12
7	Precedenční analýza	13
7.1	Gramatika pro výrazy	13
7.2	Precedenční tabulka	14

1 Rozdělení práce

Milan Jakubec

- Navrh gramatiky
- Implementace rekurzivní syntaktické analýzy
- Implementace generatoru (obecně)
- Dokumentace

Jiří Tesař

- Lexikální analýza
- Dokumentace

Norman Babiak

- Implementace analýzy výrazu
- Implementace generatoru (obecně, výrazy)
- Implementace zásobníku
- Dokumentace

Jakub Ráček

- Implementace semantické analýzy
- Implementace tabulky symbolů
- Implementace zásobníku
- Dokumentace

2 Implementace

Jednotlivé moduly prekladače jsou rozděleny do samostatných souborů. Vyjma semanticke analýzy rekurzivního syntaktického parseru, která se nachází společně se samotným parserem v jednom souboru. Deklarace struktur a funkcí (vyjma pomocných funkcí) se nachází ve stejnojmenných hlavičkových souborech, v nichž jsou uvedeny i Doxygen komentáře popisující jejich funkci.

2.1 Lexikální analýza

Lexikální analyzátor prekladače je implementován za pomoci konečného automatu (FSM). Tento automat lze ve zdrojových souborech nalézt pod názvem **scanner.c**, jeho diagram se pak nachází dále v této dokumentaci. Vstup se tokenizuje do datové struktury *token_t*, tato struktura a její další pomocné struktury se nacházejí v souboru **token.c**. Každý token v sobě má svůj typ (typy tokenů jsou dle nově pomocného enumu *token_type*), atributy a boolean *eol*, který pomáhá zjistit, zda se token nachází na novém řádku, či nikoliv.

Pro samotné načtení tokenů je klíčová funkce *get_next_token()*. Implementována je i možnost práce s dynamickými stringy, k nalezení v souboru **str.c**. Pokud jsou tedy součástí vstupního zdrojového kódu řetězcové hodnoty jako například řetězcový literal nebo identifikátor, jsou funkce ze **str.c** použity k akumulaci jednotlivých znaků do kompletního řetězce, který je pak přiřazen jako atribut příslušného tokenu. V případě výskytu chyby se vypíše chybová hláška a provádění se ukončí s odpovídající návratovou hodnotou. K tomu se využívá externí modul **error.c**.

2.2 Syntaktická analýza

Syntaktická analýza parseru je implementována na základě doporučení dle rekurzivního sestupu za použití poměrně rozsáhlé LL gramatiky. Implementaci této gramatiky resp. jejích pravidel lze nalézt v souboru **parser.c**. Samotná LL gramatika včetně LL tabulky je k dispozici v pozdějších částech tohoto dokumentu. Jelikož je parser středobodem celého překladu, je to právě on, kdo využívá lexikální analyzátor a inicializuje v sobě veškeré potřebné struktury, včetně instance generatoru, který během překladu rovnou generuje výsledný mezikód. Jednotlivá pravidla parseru jsou implementována jako funkce, která má jako parametr pokladce instanci parseru, a potom dle

potřeby různé pomocné struktury (pro parametry, argumenty funkce). Samotný parser potom pomocí předpřipravené kontroly svůj stav, jako například zda se nachází ve funkci, v cyklu, a podobně. Tokeny parser konzumuje pomocí vlastní funkce `load_token()`, která využívá rozhraní lexikálního analyzátoru. V každé funkci se také po konzumaci tokenu kontroluje, zda byl na vstupu přijat syntakticky správný a očekávaný token, pakliže nebyl, vypíše se syntaktická chyba a program se okamžitě ukončí (viz. sekce "Lexikální analýza").

2.3 Sémantická analýza

Sémantická analýza je prováděna staticky, využívá přitom toho, že jazyk IFJ23 je silně staticky typovaný. Logika analýzy je tedy užze spojena s analýzou syntaktickou a generátorem. Obzvláště kriticky je zde zásobník, kterým jsme se rozhodli řešit rozsah platnosti. Zásobník obsahuje tabulky symbolů s jejich proměnných a funkcí, ze kterých na vrcholu je vždy aktuální rozsah platnosti, a na dně je globální rozsah platnosti. V analýze se sleduje i současný stav parseru, reprezentovaný několika proměnnými. Jedná se převážně o typové a kontextové kontroly. Typové kontroly jsou především při přiřazení do proměnných, volání funkcí, a zpracování výrazů. Pro uchování a kontrolu parametru funkce mimo jiné slouží struktura dynamického pole. Kontroly kontextové zamezují práci s nedefinovanými či neinicializovanými proměnnými a redefinovanými.

2.4 Analýza výrazů

Analýza výrazů probíhá v souboru `expression.c`, který využívá strukturu syntaktické analýzy a precedenční tabulku. Na začátku je inicializován zásobník a struktura `analysis_t` pro pozdější sémantickou analýzu výrazu. Na zásobník je uložen počáteční znak "\$", poté je načten token z vrcholu zásobníku a následující token ze vstupu pomocí parserové funkce `load_token()`. Načtené tokeny potom vyhodnotí pomocná funkce `precedence()`. Funkce vrátí výsledek volající funkci, `prec_analysis()`, která se rozhodne, jestli načtený token ze vstupu uložit na vrchol zásobníku, nebo výraz v zásobníku zredukuje na jediný netriviální. Redukce jednotlivých terminálů se vyhodnotí po sémantické analýze podle gramatiky pro výrazy (viz. sekce 7.1). Funkce `reduce` rovnou generuje výsledný mezikód. Pro ukončení analýzy výrazu slouží funkce `handle_upcoming()`, která rozhodne, zda je další token součástí výrazu, nebo už má předat kontrolu zpátky do struktury správným datovým typem výsledku výrazu, s kterým pak struktura dále pracuje.

2.5 Generování kódu

Generování kódu je řešeno v souborech **gen.c** a **buildin.c**, kde v prvním souboru je k nalezení základní struktura generovaného kódu, tzn. generované deklarace promenných, generované smyčky, podmíněných větví atd., a v druhém souboru se nachází definice vestavených funkcí. Všechny parametry funkcí a vlastně i vše ostatní se předává přes zásobník.

3 Datové struktury

3.1 Tabulka symbolů

Pro implementaci tabulky symbolů jsme zvolili variantu TRP-izp. Parser pracuje se zásobníkem tabulek symbolů, kde každá tabulka reprezentuje jiný kontext. Tabulka je postavena na znalostech z předmětu IJC a byla přizpůsobena variantě s otevřenou adresací (implicitního zřetězení položek) na základě prezentace předmětu IAL. V základu se jedná o dynamické kruhové pole, kde každá položka reprezentuje funkci anebo proměnnou existující v aktuálním kontextu. Pro hashovací funkci jsme zvolili řešení ze zadané letosní IJC projektu a dle publikovaného článku[1] jsme využili variantu hashovací funkce "sdbm".

3.2 Zásobník

Soubor **stack.c** obsahuje implementaci zásobníku, který zahrnuje základní funkce, jako jsou **stack_pop()**, **stack_init()**, **stack_push()**, a **stack_free()**. Dale obsahuje další pomocné funkce, které jsou využívány při analýze výrazu a tabulky symbolů.

3.2.1 Zásobník pro analýzu výrazů

Pro analýzu výrazu jsou implementovány funkce jako **stack_push_token()**, která uloží na vrchol zásobníku přicházející token v analýze výrazu. Kromě znaku se tam také uloží jeho datový typ pro sémantickou analýzu. Další pomocnou funkcí je **stack_push_after()**, která vloží token po vrchní terminal zásobníku. Tato funkce slouží pro vložení tokenů určujících konec pro redukci výrazu. Potom **stack_top_token()**, která vrátí znak na vrcholu zásobníku, **stack_top_terminal()**, vrátí terminal na vrcholu zásobníku. Pomocná funkce **stack_count_after()** slouží pro spočítání znaku nad znakem určujícím konec pro redukci, a pro vložení těchto znaků do struktury **analysis_t**, pro

semantickou analýzu. Pro testovací účely je také implementována funkce `print_stack`.

3.2.2 Zásobník pro tabulku symbolů

Pro semantickou analýzu je implementováno několik pomocných funkcí pro práci se zásobníkem tabulek symbolů. Kromě funkcí jako je `stack_push_table()` nebo `stack_pop_table()` existují také funkce pro orientaci v kontextech, `stack_top_table()` reprezentující aktuální kontext a `stack_bottom_table()` reprezentující globální kontext. Pro prohledání stacku a nalezení promenných a funkcí slouží `stack_lookup_var()` a `stack_lookup_func()`. Mimo jiné existuje i funkce `stack_lookup_var_in_global()`, která pomáhá generatoru rozhodnout, jaký rámeček má použít.

3.3 Dynamické pole

V naší implementaci máme vlastní strukturu jménem `vector_t` pro práci s dynamickým polem. Vector je datová struktura dynamického pole, která v sobě uchovává parametry funkce. Každý zázpis funkce v tabulce symbolů má vlastní vector, který je mu přidělen při prvním zavolání nebo definici funkce a slouží pro semantickou kontrolu parametru v následujících voláních nebo definicích.

4 Závěr

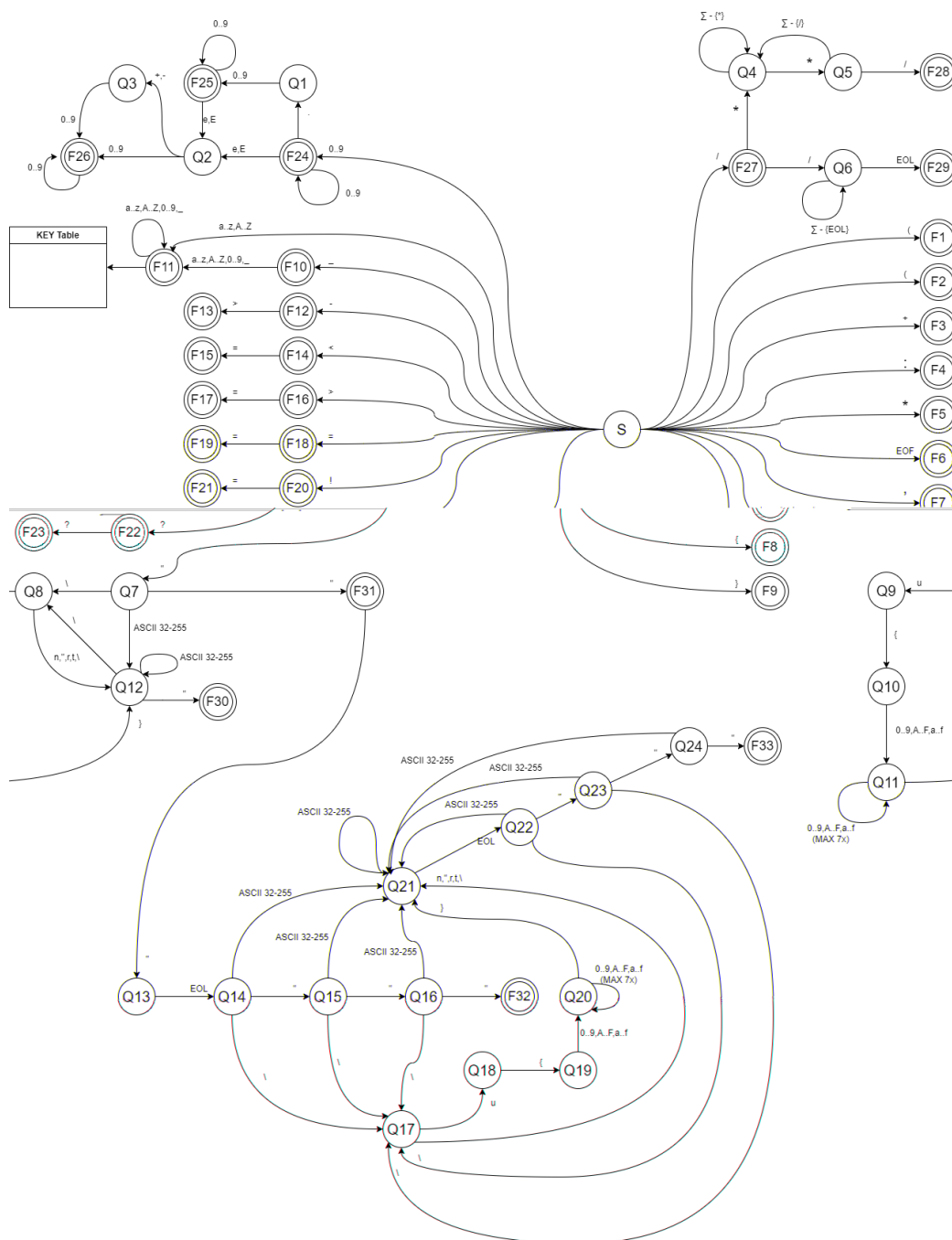
Práci na projektu jsme si navzdory jeho složitosti užili. Myslím si, že každý z nás se naučil spoustu nových a zajímavých věcí, nehlédneme na zvýšené schopnosti programovat a řešit problémy, přičemž pro každého člena týmu to samozřejmě platí převážně pro modul, který daný člen implementoval. Ale i tyto moduly spolu musely nějak fungovat, pracovat a komunikovat, a domnívám se, že to se nám povedlo zrealizovat dobře. Při řešení projektu u nás vystavaly během vývoje různé problémy, které jsme společnými silami byli schopni vyřešit. Velkou zásluhu na tom mají automatizované testy poskytnuté našimi skvělými kolegy z ročníku a samotný pan doktor Krivka, který při tvorbě prekládace byl všem velmi nápomocný. Ironicky jsme projekt do stavu, jaký jsme si představovali, dostali až den po druhém testovacím odevzdání, kdy se nám naraz povedlo opravit semantickou analýzu, naplno zprovoznit základní generování a vestavěné funkce. Závěrem bychom dodali, že velkou výzvou pro nás tím byl souboj s časem, pro každého z nás z různých důvodů, ale nakonec jsme se s tím, dle našeho názoru, takřka poprali statečně.

Reference

- [1] YIGIT, Ozan - Hash Functions: <http://www.cse.yorku.ca/~oz/hash.html>.

5 FSM pro lexikální analýzu

5.1 Diagram



6 LL

6.1 LL gramatika

1. $\langle \text{program} \rangle \rightarrow \langle \text{statement} \rangle \langle \text{program} \rangle$
2. $\langle \text{program} \rangle \rightarrow \langle \text{function_definition} \rangle \langle \text{program} \rangle$
3. $\langle \text{program} \rangle \rightarrow \epsilon$

4. $\langle \text{statement_list} \rangle \rightarrow \langle \text{statement} \rangle \langle \text{statement_list} \rangle$
5. $\langle \text{statement_list} \rangle \rightarrow \epsilon$
6. $\langle \text{statement} \rangle \rightarrow \langle \text{variable_definition_let} \rangle$
7. $\langle \text{statement} \rangle \rightarrow \langle \text{variable_definition_var} \rangle$
8. $\langle \text{statement} \rangle \rightarrow \langle \text{assignment} \rangle$
9. $\langle \text{statement} \rangle \rightarrow \langle \text{conditional_statement} \rangle$
10. $\langle \text{statement} \rangle \rightarrow \langle \text{loop} \rangle$
11. $\langle \text{statement} \rangle \rightarrow \langle \text{function_call} \rangle$
12. $\langle \text{statement} \rangle \rightarrow \langle \text{return_statement} \rangle$

13. $\langle \text{function_definition} \rangle \rightarrow \text{func } \langle \text{identifier} \rangle (\langle \text{parameter_list} \rangle)$
 $\quad \langle \text{function_return_type_and_body} \rangle$
14. $\langle \text{function_return_type_and_body} \rangle \rightarrow \{ \langle \text{void_function_body} \rangle \}$
15. $\langle \text{function_return_type_and_body} \rangle \rightarrow \text{"-"} \langle \text{type} \rangle$
 $\quad \{ \langle \text{nonvoid_function_body} \rangle \}$
16. $\langle \text{nonvoid_function_body} \rangle \rightarrow \langle \text{statement_list} \rangle$
17. $\langle \text{void_function_body} \rangle \rightarrow \langle \text{statement_list} \rangle$
18. $\langle \text{void_function_body} \rangle \rightarrow \epsilon$
19. $\langle \text{parameter_list} \rangle \rightarrow \langle \text{parameter} \rangle \langle \text{more_parameters} \rangle$
20. $\langle \text{parameter_list} \rangle \rightarrow \epsilon$

- 21. $\langle \text{more_parameters} \rangle \rightarrow , \langle \text{parameter} \rangle \langle \text{more_parameters} \rangle$
- 22. $\langle \text{more_parameters} \rangle \rightarrow \epsilon$
- 23. $\langle \text{parameter} \rangle \rightarrow \langle \text{no_name_parameter} \rangle$
- 24. $\langle \text{parameter} \rangle \rightarrow \langle \text{identifier_parameter} \rangle$
- 25. $\langle \text{no_name_parameter} \rangle \rightarrow _ \langle \text{identifier} \rangle : \langle \text{type} \rangle$
- 26. $\langle \text{identifier_parameter} \rangle \rightarrow \langle \text{identifier} \rangle \langle \text{rest_of_identifier_parameter} \rangle$
- 27. $\langle \text{rest_of_identifier_parameter} \rangle \rightarrow _ : \langle \text{type} \rangle$
- 28. $\langle \text{rest_of_identifier_parameter} \rangle \rightarrow \langle \text{identifier} \rangle : \langle \text{type} \rangle$

- 29. $\langle \text{type} \rangle \rightarrow \text{Double} \mid \text{Int} \mid \text{String} \mid \text{Double?} \mid \text{Int?} \mid \text{String?}$
- 30. $\langle \text{variable_definition_let} \rangle \rightarrow \text{let } \langle \text{identifier} \rangle \langle \text{definition_types} \rangle$
- 31. $\langle \text{variable_definition_var} \rangle \rightarrow \text{var } \langle \text{identifier} \rangle \langle \text{definition_types} \rangle$
- 32. $\langle \text{definition_types} \rangle \rightarrow \langle \text{type_def} \rangle$
- 33. $\langle \text{definition_types} \rangle \rightarrow \langle \text{initialization} \rangle$
- 34. $\langle \text{type_def} \rangle \rightarrow : \langle \text{type} \rangle \langle \text{type_def_follow} \rangle$
- 35. $\langle \text{type_def_follow} \rangle \rightarrow \langle \text{initialization} \rangle$
- 36. $\langle \text{type_def_follow} \rangle \rightarrow \epsilon$
- 37. $\langle \text{initialization} \rangle \rightarrow = \langle \text{init_type} \rangle$
- 38. $\langle \text{init_type} \rangle \rightarrow \langle \text{function_call} \rangle$
- 39. $\langle \text{init_type} \rangle \rightarrow \langle \text{expression} \rangle$

- 40. $\langle \text{assignment} \rangle \rightarrow \langle \text{identifier} \rangle = \langle \text{assignment_type} \rangle$
- 41. $\langle \text{assignment_type} \rangle \rightarrow \langle \text{function_call} \rangle$
- 42. $\langle \text{assignment_type} \rangle \rightarrow \langle \text{expression} \rangle$

- 43. $\langle \text{conditional_statement} \rangle \rightarrow \text{if } \langle \text{if_statement} \rangle$

44. `<if_statement>` -> `<classical_statement>`
45. `<if_statement>` -> `<variable_statement>`
46. `<classical_statement>` -> `<expression>` { `<statement_list>` } else { `<statement_list>` }
47. `<variable_statement>` -> let `<identifier>` { `<statement_list>` } else { `<statement_list>` }
48. `<loop>` -> while `<expression>` { `<statement_list>` }
49. `<function_call>` -> `<identifier>` (`<arguments>`)
50. `<arguments>` -> `<argument>` `<more_arguments>`
51. `<arguments>` -> ϵ
52. `<argument>` -> `<arg_name>` `<arg_value>`
53. `<arg_name>` -> `<identifier>` :
54. `<arg_name>` -> ϵ
55. `<arg_value>` -> `<identifier>` | `<int>` | `<double>` | `<string>` | `<nil>`
56. `<more_arguments>` -> , `<argument>` `<more_arguments>`
57. `<more_arguments>` -> ϵ
58. `<return_statement>` -> return `<returned_expression>`
59. `<returned_expression>` -> `<expression>`
60. `<returned_expression>` -> ϵ

6.2 LL tabulka

	func	if	else	var	let	return	Int	Double	String	()	{	}	:	while	->	<VarID>	<FunID>	=	<Undes>	<Exp>	,	<Literal>
<program>	2	1		1	1										1		1	1					
<statement list>		4		4	4							5			4		4	4					
<statement>		9		7	6	12									10		8	11					
<function definition>	13																						
<function return type and body>												14				15							
<nonvoid function body>		16		16	16	16									16		16	16					
<void function body>		17		17	17	17						18			17		17	17					
<parameter list>											20						19			19			
<more parameters>											22											21	
<parameter>																	24			23			
<no name parameter>																				25			
<identifier parameter>																	26						
<rest of identifier parameter>																	28			27			
<type>							29	29	29														
<variable definition let>					30																		
<variable definition var>				31										32					33				
<definition types>														34									
<type def>																							
<type def follow>																							
<initialization>																							
<init type>																							
<assignment>																		38			39		
<assignment type>																	40						
<conditional statement>		43																41			42		
<if statement>					45																44		
<classical statement>																					46		
<variable statement>					47																		
<loop>																							
<function call>															48								
<arguments>											51							49					
<argument>																	50						50
<arg name>																	52						
<arg value>																	53						54
<more arguments>																							55
<return statement>											57											56	
<return expression>					58																59		

7 Precedenční analýza

7.1 Gramatika pro výrazy

1. $E \rightarrow i$
2. $E \rightarrow !E$
3. $E \rightarrow (E)$
4. $E \rightarrow E + E$
5. $E \rightarrow E - E$
6. $E \rightarrow E * E$
7. $E \rightarrow E / E$
8. $E \rightarrow E == E$
9. $E \rightarrow E != E$
10. $E \rightarrow E >= E$
11. $E \rightarrow E <= E$
12. $E \rightarrow E > E$
13. $E \rightarrow E < E$
14. $E \rightarrow E ?? E$

7.2 Precedenční tabulka

	+	-	*	/	==	!=	>=	>	<	??	!	()	ID	INT	FLO	STR	NIL	\$
+	>	>	>	<	>	>	>	>	>	>	>	<	>	<	<	<	<	<	>
-	>	>	>	<	>	>	>	>	>	>	>	<	>	<	<	<	<	<	>
*	>	>	>	>	>	>	>	>	>	>	<	<	>	<	<	<	<	<	>
/	>	>	>	>	>	>	>	>	>	>	<	<	>	<	<	<	<	<	>
==	<	<	<	<	>	>	>	>	>	>	<	<	>	<	<	<	<	<	>
!=	<	<	<	<	>	>	>	>	>	>	<	<	>	<	<	<	<	<	>
>=	<	<	<	<	>	>	>	>	>	>	<	<	>	<	<	<	<	<	>
<	<	<	<	<	>	>	>	>	>	>	<	<	>	<	<	<	<	<	>
??	<	<	<	<	>	>	>	>	>	>	<	<	>	<	<	<	<	<	>
!	>	>	>	>	>	>	>	>	>	>	<	<	>	<	<	<	<	<	>
(>	>	>	>	>	>	>	>	>	>	<	<	>	<	<	<	<	<	>
)	>	>	>	>	>	>	>	>	>	>	<	<	>	<	<	<	<	<	>
ID	>	>	>	>	>	>	>	>	>	>	<	<	>	<	<	<	<	<	>
INT	>	>	>	>	>	>	>	>	>	>	<	<	>	<	<	<	<	<	>
FLO	>	>	>	>	>	>	>	>	>	>	<	<	>	<	<	<	<	<	>
STR	>	>	>	>	>	>	>	>	>	>	<	<	>	<	<	<	<	<	>
NIL	>	>	>	>	>	>	>	>	>	>	<	<	>	<	<	<	<	<	>
\$	>	>	>	>	>	>	>	>	>	>	<	<	>	<	<	<	<	<	>

> Redukce, < Uložen na zásobník