

Dokumentace k projektu do předmětů IFJ a IAL
Implementace překladače imperativního jazyka IFJ20

Tým 86, varianta II

Samuel Valaštín	xvalas10	25%
Jonáš Tichý	xtichy29	25%
Daniel Gavenda	xgaven08	25%
Miroslav Štěpánek	xstepa68	25%

9. prosince 2020

Obsah

1 Úvod	3
2 Použité datové struktury	3
2.1 Tabulka s rozptýlenými položkami	3
2.2 Lineárně vázaný seznam	3
2.3 Parse tree	3
2.4 AST	3
3 Lexikální analýza	3
4 Syntaktická analýza	3
4.1 Analýza za využití rekurzivního sestupu	4
4.2 Precedenční syntaktická analýza	4
4.3 Tvorba AST	4
4.4 Konverze výrazu do AST	4
5 Sémantická analýza	4
5.1 Analýza výrazů	4
6 Generování kódu	5
6.1 Struktura výsledného programu	5
6.2 Řešení různého zanoření proměnných	5
6.3 Generování výrazů	5
7 Práce v týmu	5
7.1 Způsob práce	5
7.2 Rozdělení práce	5
7.3 Komunikace	6
8 Závěr	6
9 Přílohy	8

1 Úvod

Cílem projektu je vytvoření překladače, který ze standardního vstupu načítá zdrojový program jazyka IFJ20 (podmnožina jazyka GO) a na standardní výstup generuje mezikód v jazyku IFJcode20.

2 Použité datové struktury

2.1 Tabulka s rozptýlenými položkami

Základem tabulky s rozptýlenými položkami je efektivní hashovací funkce, u které jsme se inspirovali `sdbm` algoritmem¹. Pro implementaci jsme zvolili metodu explicitně zřetězených prvků. V programu tato struktura slouží k ukládání údajů o deklaracích proměnných a funkcí. Pro umožnění práce s rámci a zastíňování proměnných jsme se rozhodli vytvořit pro každý rámec unikátní tabulku symbolů o velikosti 509 prvků. Toto číslo bylo zvoleno z důvodu adekvátní velikosti a faktu, že jde o prvočíslo.

2.2 Lineárně vázaný seznam

Tato struktura je využívána v syntaktické analýze pro ukládání jednotlivých tokenů. Při sémantické analýze a generování kódu slouží k zachování hierarchie vnořených rámců, reprezentovaných tabulkami symbolů.

2.3 Parse tree

Tento strom tvořený uzly `PTNode` explicitně reprezentuje zdrojový kód a je později použit pro tvorbu AST.

2.4 AST

AST slouží k co nejjednodušší reprezentaci zdrojového kódu, přičemž stále zachovává jeho význam. Tohoto zjednodušení je dosaženo vynecháním nepotřebných tokenů, jako jsou například odřádkování, závorky nebo klíčová slova. Uzly, které vytvářejí vlastní rámec, obsahují ve své struktuře také ukazatel na příslušnou tabulku symbolů. Jednotlivé uzly jsou reprezentovány strukturou `ASTNode`.

3 Lexikální analýza

První částí při tvorbě překladače byla implementace lexikální analýzy. Hlavní funkcí tohoto modulu je `getToken`. Prostřednictvím této funkce načítáme znaky ze standardního vstupu a transformujeme je na tokeny. Struktura tokenu se skládá z typu a atributu. Atribut je typu `union` a může nabývat hodnoty podle typu tokenu. Pokud je token typu `string` nebo identifikátor, tak nabývá hodnoty `char*`. Pro typ `integer` nabývá hodnoty `int64_t`, pro typ `float64` hodnoty `double` a v ostatních případech je atribut prázdný. Lexikální analyzátor je implementovaný nekonečným cyklem, ve kterém pomocí `switch` vybíráme stavy automatu na základě načteného znaku. Pokud načtený znak není v daném stavu povolený, pak funkce vrací příslušnou návratovou hodnotu volající funkci, která zajistí bezpečné ukončení programu. V opačném případě se načítají další znaky a přechází se mezi stavy dokud neskončí v koncovém stavu automatu. Pro zpracování escape sekvencí využíváme buffer, který načítá slovo v hexadecimální soustavě a následně ho převede do ASCII.

4 Syntaktická analýza

V této části projektu používáme pro analýzu výrazů precedenční syntaktickou analýzu a pro vše ostatní metodu rekurzivního sestupu za použití námi vytvořené LL – gramatiky. Komunikace s lexikálním analyzátozem probíhá pomocí již dříve zmiňované funkce `getToken`. Jednotlivé tokeny se ukládají do jednostranně vázaného seznamu, pro možnost rychlého a efektivního přístupu k dříve načteným tokenům.

¹<http://www.cse.yorku.ca/~oz/hash.html>

V rámci rekurzivního sestupu se kromě syntaktické analýzy provádí také tvorba parse tree, který slouží k přesné reprezentaci zdrojového kódu. Tento strom je později použit k tvorbě AST.

4.1 Analýza za využití rekurzivního sestupu

Metoda rekurzivního sestupu za použití LL–gramatiky je v naší implementaci řešena souborem konstantních globálních proměnných v souboru `grammar.h`. Pomocí ukazatelů můžeme takto efektivně reprezentovat naši LL–gramatiku bez nutnosti implementace nadbytečného počtu funkcí. Celá syntaktická analýza je tedy reprezentována konstantami v `grammar.h` a jedinou funkcí `recursive` v souboru `parser.c`, která provádí rekurzivní sestup nad gramatikou a kontroluje validitu načtených tokenů. Zbytek funkcí v tomto souboru slouží buď k precedenční syntaktické analýze, tvorbě AST nebo uvolňování alokované paměti.

4.2 Precedenční syntaktická analýza

Zpracování výrazů probíhá způsobem bottom-up ve funkci `precExp`, která je volána pokaždé, když rekurzivní funkce najde možné místo výskytu výrazu. Zde je kontrolována syntaktická správnost výrazu a dále je zde volána funkce `precTable`, která na základě precedenční tabulky rozhoduje, zda bude proveden `shift`, `reduce` nebo se jedná o neplatný stav. Funkce `precShift` do seznamu vloží zarážku, jejíž pozice se nachází před posledním prvkem, který není uzal. Pokud žádný prvek není nalezen, vloží se zarážka na začátek celého seznamu. Následně je vložen i samotný token. Funkce `precReduce` pak redukuje tokeny do struktury `PTNode`, dokud nenarazí na již zmíněnou zarážku. Vše probíhá až do chvíle, kdy nastane některý z chybových stavů nebo se podaří výraz zredukovat do jediné `PTNode`, ve které jsou pak dílčí operace seřazeny podle priority.

4.3 Tvorba AST

Proběhne-li kontrola syntaxe bez nalezení chyby, pak přichází na řadu tvorba AST. Tento strom slouží k co nestručnější reprezentaci zdrojového kódu, která je významově identická. K tvorbě využíváme rekurzivního sestupu nad dříve vytvořeným `parse tree`, který zjednodušujeme vynecháním nepotřebných tokenů. Díky tomuto zjednodušení dosáhneme vyšší rychlosti při následných průchodech stromu při sémantické analýze a generování kódu.

4.4 Konverze výrazu do AST

Výrazy jsou za pomoci rekurzivní funkce `PTtoAST` překonvertovány ze struktury `PTNode` do struktury `ASTNode`. Dále jsou odstraněny všechny závorky, které byly u precedence nutné, ale dále již nejsou žádoucí.

5 Sémantická analýza

Tato část je tvořena dvojím průchodem AST. V prvním průchodu dochází k naplnění tabulky symbolů globálního rámce definicemi uživatelských funkcí. Tento krok je důležitý pro zajištění možnosti volání funkcí před jejich samotnou definicí ve zdrojovém kódu. Následující průchod provádí kompletní sémantickou analýzu a naplňování příslušných tabulek symbolů. Dochází zde ke kontrole správného volání funkcí, definic, přiřazení, podmínek, cyklů a příkazů `return` pro zajištění bezchybnosti kódu před předáním stromu generátoru mezikódu. Zdrojové kódy sémantické analýzy jsou umístěny v samostatném souboru `semantics.c`, pro větší přehlednost ve struktuře programu.

5.1 Analýza výrazů

Sémantika výrazu je kontrolována funkcí `checkExpression`, která postupuje rekurzivně po uzlech AST, dokud neurčí datové typy všech operandů výrazu. Tyto typy mezi sebou porovná a zkontroluje, zda jsou

identické. Dále je také kontrolováno, zda jsou operátory přítomné ve výrazu kompatibilní s datovým typem `string`, zda není jmenovatel při dělení nulová konstanta nebo zda se ve výrazu nachází pouze jeden relační operátor, který musí být kořenem celého AST stromu.

6 Generování kódu

Vstupem při generování kódu je ukazatel na kořen abstraktního syntaktického stromu. Výstupem je mezikód `IFJcode20`, který je generovaný na standartní výstup.

6.1 Struktura výsledného programu

Mezikód začíná prologem a skokem na hlavní funkci `main`. Následně se generují jednotlivé funkce. Vzhledem k zadání pracujeme s běžným lokálním rámcem. V případě volání funkce vytváříme dočasný rámec, do kterého ukládáme parametry funkce. Vestavěné funkce generujeme jen v případě, že jsou volané ve zdrojovém kódu vstupního programu. Jejich samotné generování se vykoná až po vygenerování vstupního kódu. Speciálním případem vestavěné funkce je `print`, který generujeme, vzhledem k proměnlivému počtu parametrů a snížení ceny interpretace kódu, pro každé volání zvlášť.

6.2 Řešení různého zanoření proměnných

Vzhledem k možnému problému interpretace, při výskytu proměnných se stejnými identifikátory, v rozdílných zanořeních, na začátku každé funkce předem definujeme všechny proměnné. Název proměnné se skládá ze jména proměnné, za kterým následuje „*“ a hloubka zanoření. Výjimku tvoří definice pomocných proměnných při vícenásobném přiřazení, které ukládáme jako `tmpVar` a číslo proměnné.

6.3 Generování výrazů

Vyhodnocování výrazů probíhá rekurzivně, přičemž používáme instrukce, které pracují se zásobníkem. V případě výskytu operace „+“ a „/“ se před vyhodnocením zanoříme a určíme typ literálu nebo identifikátoru pomocí tabulky symbolů.

7 Práce v týmu

7.1 Způsob práce

Jako vývojové prostředí jsme si určili `Visual Studio Code`, které nabízí jednoduchou implementaci verzovacího systému `Git`. Jako společný vzdálený repozitář jsme používali webovou službu `GitHub`. Vzhledem k náročnosti projektu jsme v průběhu vývoje využívali párové programování, přičemž jsme se rozdělili na slovenskou a českou část.

7.2 Rozdělení práce

- Daniel Gavenda – lexikální analýza, generátor kódu
- Samuel Valaštín – lexikální analýza, generátor kódu
- Jonáš Tichý – syntaktická analýza, sémantická analýza, AST, tabulka symbolů, dokumentace
- Miroslav Štěpánek – syntaktická analýza, syntaktická a sémantická analýza výrazů, dokumentace

7.3 Komunikace

Na projektu jsme začali pracovat týden po zveřejnění zadání. Stihli jsme uskutečnit jednu schůzi, na které jsme si přiblížili detaily zadání a rozdělili jsme si práci na projektu. Vzhledem k situaci jsme přesunuli komunikaci na vlastní server Discord. Každý týden jsme uspořádali konferenční hovor, při kterém jsme si sdělili nový postup ve vývoji projektu.

8 Závěr

Projekt se nám podařilo dokončit s dostatečnou časovou rezervou. Díky dobré týmové komunikaci jsme dokázali vyřešit mnoho problémů, které při vývoji nastaly a nasbírat cenné zkušenosti z oblasti vývoje a týmové spolupráce.

Odkazy

- [1] Jack Crenshaw. *Let's Build a Compiler*. [online]. [cit. 9.12.2020]. 1997-2008. URL: <https://compilers.iecc.com/crenshaw/>.
- [2] Ruslan Spivak. *Let's Build A Simple Interpreter*. [online]. [cit. 9.12.2020]. 2015. URL: <https://ruslanspivak.com/lsbasi-part1/>.
- [3] Hayo Thielecke. *Typed trees and tree walking in C with struct, union, enum, and switch*. [online]. [cit. 9.12.2020]. URL: <https://www.cs.bham.ac.uk/~hxt/2015/c-plus-plus/trees-in-c-slides.pdf>.

9 Přílohy

1. $\langle prog \rangle \rightarrow \langle multiEol \rangle \textbf{package} \langle multiEol \rangle \textbf{main EOL} \langle progBody \rangle$
2. $\langle progBody \rangle \rightarrow \langle multiEol \rangle \langle func \rangle \langle progBody \rangle$
3. $\langle progBody \rangle \rightarrow \langle multiEol \rangle$
4. $\langle func \rangle \rightarrow \textbf{func ID} \langle funcParams \rangle \langle funcRetTypes \rangle \{ \textbf{EOL} \langle body \rangle \}$
5. $\langle funcParams \rangle \rightarrow (\langle params \rangle)$
6. $\langle funcParams \rangle \rightarrow ()$
7. $\langle params \rangle \rightarrow \textbf{ID} \langle type \rangle , \langle params \rangle$
8. $\langle params \rangle \rightarrow \textbf{ID} \langle type \rangle$
9. $\langle funcRetTypes \rangle \rightarrow (\langle retTypes \rangle)$
10. $\langle funcRetTypes \rangle \rightarrow ()$
11. $\langle funcRetTypes \rangle \rightarrow \epsilon$
12. $\langle retTypes \rangle \rightarrow \langle type \rangle , \langle retTypes \rangle$
13. $\langle retTypes \rangle \rightarrow \langle type \rangle$
14. $\langle for \rangle \rightarrow \textbf{for} \langle forDec \rangle ; \langle expr \rangle ; \langle forInc \rangle \{ \textbf{EOL} \langle body \rangle \}$
15. $\langle forDef \rangle \rightarrow \langle define \rangle$
16. $\langle forDef \rangle \rightarrow \epsilon$
17. $\langle forInc \rangle \rightarrow \langle assign \rangle$
18. $\langle forInc \rangle \rightarrow \epsilon$
19. $\langle if \rangle \rightarrow \textbf{if} \langle expr \rangle \{ \textbf{EOL} \langle body \rangle \} \textbf{else} \{ \textbf{EOL} \langle body \rangle \}$
20. $\langle body \rangle \rightarrow \langle statement \rangle \textbf{EOL} \langle body \rangle$
21. $\langle body \rangle \rightarrow \epsilon$
22. $\langle statement \rangle \rightarrow \langle If \rangle$
23. $\langle statement \rangle \rightarrow \langle For \rangle$
24. $\langle statement \rangle \rightarrow \langle define \rangle$
25. $\langle statement \rangle \rightarrow \langle assign \rangle$
26. $\langle statement \rangle \rightarrow \langle funcCall \rangle$
27. $\langle statement \rangle \rightarrow \langle ret \rangle$
28. $\langle statement \rangle \rightarrow \epsilon$
29. $\langle funcCall \rangle \rightarrow \textbf{ID} (\langle funcCallArgs \rangle)$
30. $\langle funcCall \rangle \rightarrow \textbf{ID} ()$
31. $\langle funcCallArgs \rangle \rightarrow \langle operand \rangle , \langle funcCallArgs \rangle$

32. $\langle funcCallArgs \rangle \rightarrow \langle operand \rangle$
33. $\langle multiExpr \rangle \rightarrow \langle expr \rangle , \langle multiExpr \rangle$
34. $\langle multiExpr \rangle \rightarrow \langle expr \rangle$
35. $\langle define \rangle \rightarrow \mathbf{ID} := \langle eolOpt \rangle \langle expr \rangle$
36. $\langle ret \rangle \rightarrow \mathbf{return} \langle multiExpr \rangle$
37. $\langle ret \rangle \rightarrow \mathbf{return}$
38. $\langle assign \rangle \rightarrow \langle multId \rangle = \langle eolOpt \rangle \langle funcCall \rangle$
39. $\langle assign \rangle \rightarrow \langle multId \rangle = \langle eolOpt \rangle \langle multiExpr \rangle$
40. $\langle lParenthsOpt \rangle \rightarrow (\langle multiEol \rangle \langle lParenthsOpt \rangle$
41. $\langle lParenthsOpt \rangle \rightarrow \epsilon$
42. $\langle rParenthsOpt \rangle \rightarrow) \langle rParenthsOpt \rangle$
43. $\langle rParenthsOpt \rangle \rightarrow \epsilon$
44. $\langle operand \rangle \rightarrow \mathbf{ID}$
45. $\langle operand \rangle \rightarrow \mathbf{intLit}$
46. $\langle operand \rangle \rightarrow \mathbf{fltLit}$
47. $\langle operand \rangle \rightarrow \mathbf{strLit}$
48. $\langle type \rangle \rightarrow \mathbf{int}$
49. $\langle type \rangle \rightarrow \mathbf{string}$
50. $\langle type \rangle \rightarrow \mathbf{float64}$
51. $\langle multId \rangle \rightarrow \mathbf{ID} , \langle multId \rangle$
52. $\langle multId \rangle \rightarrow \mathbf{ID}$
53. $\langle multiEol \rangle \rightarrow \mathbf{EOL} \langle multiEol \rangle$
54. $\langle multiEol \rangle \rightarrow \epsilon$
55. $\langle eolOpt \rangle \rightarrow \mathbf{EOL}$
56. $\langle eolOpt \rangle \rightarrow \epsilon$

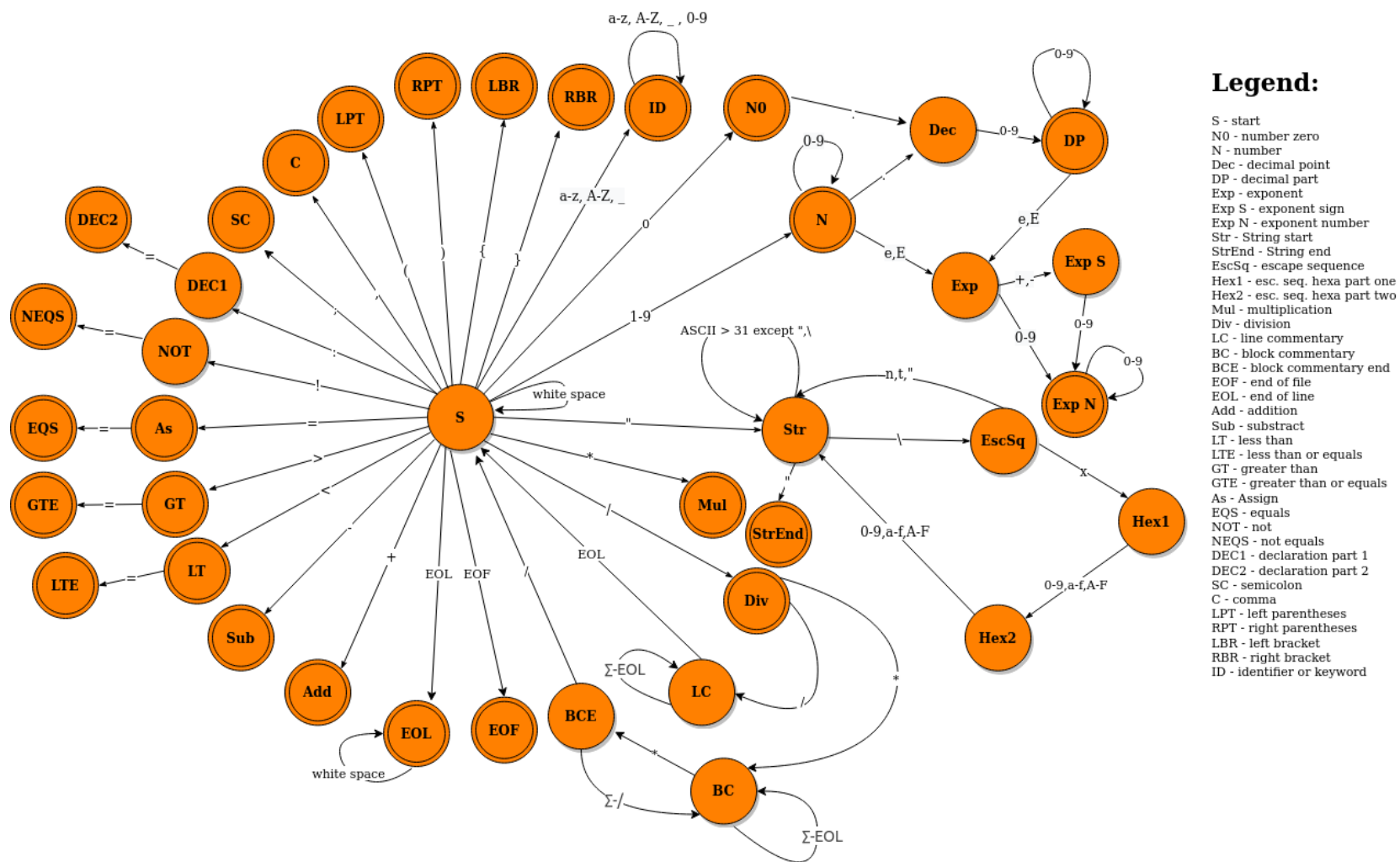
Tabulka 1: LL – gramatika

	+, -	*, /	Rel	()	Lit, ID	\$
+, -	>	<	>	<	>	<	>
*, /	>	>	>	<	>	<	>
Rel	<	<		<	>	<	>
(<	<	<	<	=	<	
)	>	>	>		>		>
Lit, ID	>	>	>		>		>
\$	<	<	<	<		<	

Tabulka 2: Precedenční tabulka

	EOL	ID	eps	func	for	if	return	int	string	float64	()	intLit	strLit	fltLit
< prog >	1		1												
< progBody >	2,3		2,3												
< func >				4											
< funcParams >											5,6				
< params >		7,8													
< funcRetTypes >			11								9,10				
< retTypes >								12,13	12,13	12,13					
< for >					14										
< forDef >	16	15													
< forInc >	18	17													
< if >						19									
< body >		20	21		20	20	20								
< statement >		24,25,26	28		23	22	27								
< funcCall >		29,30													
< funcCallArgs >		31,32											31,32	31,32	31,32
< multiExpr >			33,34												
< define >		35													
< ret >							36,37								
< assign >		38,39													
< lParenthsOpt >			41								40				
< rParenthsOpt >			43									42			
< operand >		44											45	47	46
< type >								48	49	50					
< multiId >		51,52													
< multiEol >	53		54												
< eolOpt >	55		56												

Tabulka 3: LL – tabulka



Obrázek 1: Diagram konečného automatu