



FAKULTA INFORMAČNÍCH TECHNOLOGIÍ VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ

Formální jazyky a překladače

Dokumentace projektu do IFJ a IAL

Tým 066, varianta II

Roman Ondráček	xondra58	28%	vedoucí
Pavel Raur	xraurp00	28%	
František Jeřábek	xjerab25	28%	
Radim Lipka	xlipka02	16%	

Seznam implementovaných rozšíření

BOOLOP, BASE, FUNEXP, IFTHEN TABUNARY

Obsah

1	Úvod	2
2	Rozbor částí překladače	2
2.1	Lexikální analýza	2
2.2	Syntaktická analýza	4
2.2.1	LL–Gramatika a LL–Tabulka	4
2.2.2	Precedenční syntaktická analýza	6
2.3	Sémantická analýza	7
2.4	Generátor cílového kódu	7
2.5	Tabulka symbolů a další datové struktury	7
3	Implementovaná rozšíření	8
3.1	BOOLOP	8
3.2	BASE	8
3.3	FUNEXP	8
3.4	IFTHEN	8
3.5	TABUNARY	8
4	Týmová spolupráce	9
4.1	Rozdělení práce	9
5	Závěr	9
6	Použitá literatura, nástroje, programy	9

1 Úvod

Tento dokument byl vytvořen jako dokumentace společného projektu do předmětů Formální jazyky a překladače a Algoritmy. Jsou zde popsány postupy implementace jednotlivých komponent překladače a problémy spojené s jejich implementací.

2 Rozbor částí překladače

Zadáním bylo vytvořit překladač jazyka **IFJ19**, který je podmnožinou jazyka Python 3, do cílového jazyka **IFJcode19**. Jelikož se jedná o variantu projektu II, bylo nutné implementovat tabulku symbolů implementovat jako tabulku s rozptýlenými položkami.

Implementovaný překladač se dělí na několik dílčích modulů:

- Lexikální analyzátor
- Syntaktický analyzátor
- Sémantický analyzátor
- Generátor cílového kódu

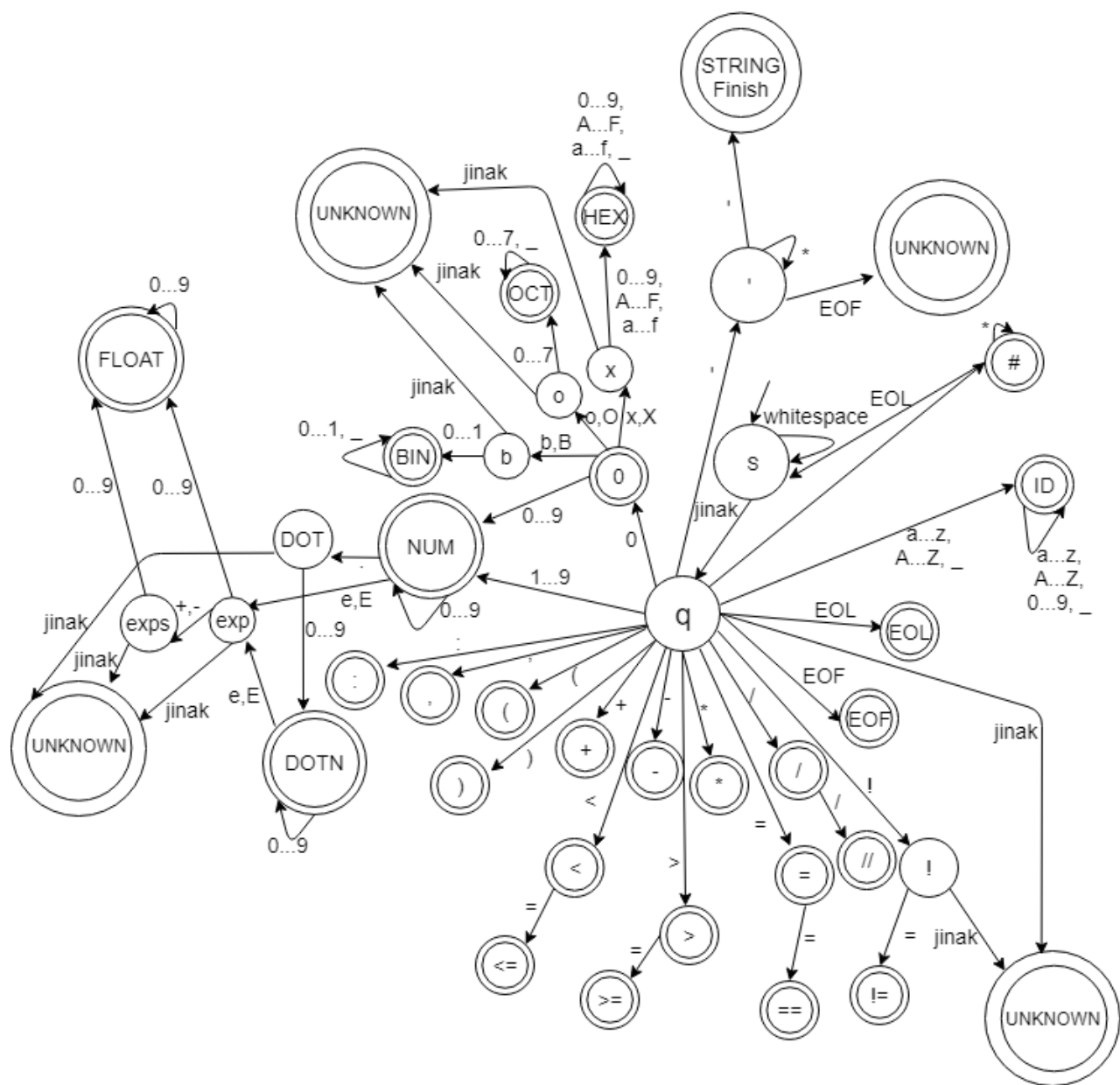
Při samotném překladu se tedy dostane ke slovu jako první lexikální analyzátor, poté přijde na řadu syntaktický analyzátor se sémantickou analýzou a následně generátor cílového kódu.

2.1 Lexikální analýza

Úkol lexikálního analyzátoru je načítat jednotlivé lexikální jednotky (lexémy) a převádět je na jednotlivé tokeny, které reprezentují daný lexém v dalších částech překladače.

Část lexikální analýzy překladače je naimplementována v modulu **scanner.c** pomocí *konečného stavového automatu*, který byl pro tuto část překladače navrhnut. Nerozpoznatelný vstupní token je označen `T_UNKNOWN` a token reprezentující interní chybu překladače, například chybu alokace paměti `T_ERROR`.

Důležitou součástí lexikální analýzy překladače jazyka IFJ19 je také zásobníkový automat, který je použit pro práci s odsazením jednotlivých řádků vstupního souboru.



Obrázek 1: Automat lexikální analýzy

2.2 Syntaktická analýza

Implementace syntaktické analýzy se nachází v modulu **parser.c**. Jedná se o implementaci metodou shora dolů, konkrétněji metodou *rekurzivního sestupu*, která je založena na LL–gramatice a LL–tabulce.

Vstupem syntaktické analýzy jsou jednotlivé tokeny z lexikální analýzy.

Jedním z úkolů syntaktické analýzy je naplnění tabulky symbolů.

2.2.1 LL–Gramatika a LL–Tabulka

1. $\langle \text{code} \rangle \rightarrow \langle \text{body} \rangle \langle \text{eols} \rangle \text{EOF}$
2. $\langle \text{body} \rangle \rightarrow \langle \text{definitions} \rangle \langle \text{statements} \rangle$
3. $\langle \text{definitions} \rangle \rightarrow \langle \text{eols} \rangle \langle \text{definition} \rangle \langle \text{definitions} \rangle$
4. $\langle \text{definition} \rangle \rightarrow \epsilon$
5. $\langle \text{definition} \rangle \rightarrow \text{DEF IDENTIFIER} (\langle \text{function_params} \rangle) : \langle \text{eols} \rangle \text{INDENT} \langle \text{statements} \rangle \text{DEDENT}$
6. $\langle \text{definition} \rangle \rightarrow \langle \text{function_call} \rangle$
7. $\langle \text{function_call} \rangle \rightarrow \text{IDENTIFIER} (\langle \text{function_params} \rangle)$
8. $\langle \text{function_params} \rangle \rightarrow \epsilon$
9. $\langle \text{function_params} \rangle \rightarrow \langle \text{function_param} \rangle \langle \text{function_nparam} \rangle$
10. $\langle \text{function_nparam} \rangle \rightarrow , \langle \text{function_param} \rangle \langle \text{function_nparam} \rangle$
11. $\langle \text{function_nparam} \rangle \rightarrow \epsilon$
12. $\langle \text{function_param} \rangle \rightarrow \text{IDENTIFIER}$
13. $\langle \text{statements} \rangle \rightarrow \epsilon$
14. $\langle \text{statements} \rangle \rightarrow \langle \text{statement} \rangle \text{EOL} \langle \text{eols} \rangle \langle \text{statements} \rangle$
15. $\langle \text{statement} \rangle \rightarrow \langle \text{returnRule} \rangle$
16. $\langle \text{statement} \rangle \rightarrow \langle \text{condition} \rangle$
17. $\langle \text{statement} \rangle \rightarrow \langle \text{assignment} \rangle$
18. $\langle \text{statement} \rangle \rightarrow \langle \text{whileRule} \rangle$
19. $\langle \text{statement} \rangle \rightarrow \text{PASS}$
20. $\langle \text{returnRule} \rangle \rightarrow \text{RETURN} \langle \text{return_expression} \rangle$
21. $\langle \text{return_expression} \rangle \rightarrow \epsilon$
22. $\langle \text{return_expression} \rangle \rightarrow (\langle \text{return_expression} \rangle)$

23. $\langle \text{return_expression} \rangle \rightarrow \langle \text{expression} \rangle$
24. $\langle \text{condition} \rangle \rightarrow \text{IF } \langle \text{condition_expression} \rangle : \text{EOL } \langle \text{eols} \rangle \text{ INDENT } \langle \text{statement} \rangle \text{ DEDENT } \langle \text{else_condition} \rangle$
25. $\langle \text{else_condition} \rangle \rightarrow \epsilon$
26. $\langle \text{else_condition} \rangle \rightarrow \text{ELSE} : \text{EOL } \langle \text{eols} \rangle \text{ INDENT } \langle \text{statement} \rangle \text{ DEDENT}$
27. $\langle \text{condition_expression} \rangle \rightarrow (\langle \text{condition_expression} \rangle)$
28. $\langle \text{condition_expression} \rangle \rightarrow \langle \text{expression} \rangle$
29. $\langle \text{assignment} \rangle \rightarrow \text{IDENTIFIER} = \langle \text{expression} \rangle$
30. $\langle \text{whileRule} \rangle \rightarrow \text{WHILE } \langle \text{condition_expression} \rangle : \text{EOL } \langle \text{eols} \rangle \text{ INDENT } \langle \text{statement} \rangle \text{ DEDENT}$
31. $\langle \text{eols} \rangle \rightarrow \epsilon$
32. $\langle \text{eols} \rangle \rightarrow \text{EOL } \langle \text{eols} \rangle$

	DEF	ID	()	INDENT	DEDENT	RETURN	IF	ELSE	..	WHILE	PASS	EOF	EOL	=	,
$\langle \text{code} \rangle$	1	1												1		
$\langle \text{body} \rangle$	2	2												2		
$\langle \text{defs} \rangle$	3	3												3		
$\langle \text{def} \rangle$	5	6														
$\langle \text{f_call} \rangle$		7														
$\langle \text{f_params} \rangle$		9														
$\langle \text{f_nparam} \rangle$		10														10
$\langle \text{f_param} \rangle$		12														
$\langle \text{stats} \rangle$																
$\langle \text{stat} \rangle$		17					15	16			18	19				
$\langle \text{retRule} \rangle$							20									
$\langle \text{ret_exp} \rangle$			22													
$\langle \text{condition} \rangle$								24								
$\langle \text{else_cond} \rangle$									26							
$\langle \text{cond_exp} \rangle$			27													
$\langle \text{assignment} \rangle$		29														
$\langle \text{whileRule} \rangle$											30					
$\langle \text{eols} \rangle$	31	31	31	31	31	31	31	31	31	31	31	31	31	32	31	31

Tabulka 1: LL–Tabulka

2.2.2 Precedenční syntaktická analýza

Precedenční syntaktická analýza je využita pro zpracování výrazů a je řízena precedenční tabulkou. Stejně jako syntaktická analýza je precedenční analýza implementována v modulu **parser.c**.

Jelikož bylo implementováno i rozšíření TABUNARY, bylo potřeba se vypořádat s rozdílem mezi unárními operátory plus a mínus a jejich binární podobou, protože do této části syntaktické analýzy přicházejí jako od sebe nerozpoznatelné tokeny.

V precedenční analýze je hojně využit zásobník. Každý příchozí token je porovnáván s aktuálním obsahem vrcholu zásobníku a podle precedenční tabulky je prováděna náležitá operace. Ty jsou typicky dvě, buď je zpracováván token pouze uložen na zásobník nebo je uvolněn aktuální vrchol zásobníku. Ten je redukován podle jednoho z předem definovaných pravidel. Výraz se takto postupně zpracovává a postupně redukuje, dokud není dosažen konec zpracovávaného výrazu, který je symbolizován znakem \$. Stejný znak je na začátku precedenční analýzy vložen na vrchol zásobníku, takže posledním porovnáním u každého výrazu, který je bez chyby, je porovnání \$ s \$. Po dokončení všech redukcí a po nalezení konce výrazu tedy vzniká výsledný strom.

$E \rightarrow id$	$E \rightarrow E + E$	$E \rightarrow E < E$
$E \rightarrow (E)$	$E \rightarrow E - E$	$E \rightarrow E > E$
$E \rightarrow id()$	$E \rightarrow E * E$	$E \rightarrow E \leq E$
$E \rightarrow id(E)$	$E \rightarrow E / E$	$E \rightarrow E \geq E$
$E \rightarrow id(E, \dots)$	$E \rightarrow E // E$	$E \rightarrow E \text{ and } E$
$E \rightarrow E == E$	$E \rightarrow E != E$	$E \rightarrow E \text{ or } E$
$E \rightarrow \text{not } E$		

Tabulka 2: Redukční pravidla

	*	/	//	+	-	<	>	<=	>=	==	and	or	not	!=	()	id	\$
*	>	>	>	>	>	>	>	>	>	>	>	>	<	>	<	>	<	>
/	>	>	>	>	>	>	>	>	>	>	>	>	<	>	<	>	<	>
//	>	>	>	>	>	>	>	>	>	>	>	>	<	>	<	>	<	>
+	<	<	<	>	>	>	>	>	>	>	>	>	<	>	<	>	<	>
-	<	<	<	>	>	>	>	>	>	>	>	>	<	>	<	>	<	>
<	<	<	<	<	<	-	-	-	-	-	>	>	<	-	<	>	<	>
>	<	<	<	<	<	-	-	-	-	-	>	>	<	-	<	>	<	>
<=	<	<	<	<	<	-	-	-	-	-	>	>	<	-	<	>	<	>
>=	<	<	<	<	<	-	-	-	-	-	>	>	<	-	<	>	<	>
==	<	<	<	<	<	-	-	-	-	-	>	>	<	-	<	>	<	>
and	<	<	<	<	<	<	<	<	<	<	>	>	<	<	<	>	<	>
or	<	<	<	<	<	<	<	<	<	<	<	>	<	<	<	>	<	>
not	>	>	>	>	>	>	>	>	>	>	>	>	<	>	<	>	<	>
!=	<	<	<	<	<	-	-	-	-	-	>	>	<	-	<	>	<	>
(<	<	<	<	<	<	<	<	<	<	<	<	<	<	<	=	<	-
)	>	>	>	>	>	>	>	>	>	>	>	>	>	>	-	>	-	>
id	>	>	>	>	>	>	>	>	>	>	>	>	>	>	=	>	-	>
\$	<	<	<	<	<	<	<	<	<	<	<	<	<	<	<	-	<	-

Tabulka 3: Tabulka precedenční syntaktické analýzy

2.3 Sémantická analýza

Sémantické kontroly jsou přidruženy k rekurzivnímu sestupu syntaktické analýzy. Nainplementovány jsou v modulu **semantic_analysis.c**. Největší část sémantických kontrol spočívá v kontrole datových typů výrazů, které vstupují do různých operací. V případě nejednotnosti těchto datových typů poté probíhá implicitní přetypování. Jelikož byla implementována podpora rozšíření BOOLOP, kromě implicitního přetypování z typu *int* na *float*, je nutné podporovat i přetypování z typu *bool* na *int* a *float*.

2.4 Generátor cílového kódu

Generátor cílového kódu byl implementován v modulu **inter_code_generator.c**. Úkolem tohoto modulu je vytvářet cílový kód, tedy v našem případě kód **IFJcode19**.

2.5 Tabulka symbolů a další datové struktury

Tabulka symbolů je implementována jako tabulka s rozptýlenými položkami v modulu **symtable.c**. Tabulka symbolů uchovává informace o všech jednotkách vyskytujících se ve vstupním programu.

Pro práci s různě dlouhými vstupními řetězci byl v modulu **dynamic_string.c** implementován datový typ dynamický string. V lexikální analýze je také použit již dříve zmíněný zásobník v modulu **stack.c** pro práci s různými úrovněmi zanoření lexémů.

Precedenční analýza poté pracuje se zásobníkem z modulu **tree_element_stack.c**, na kterém probíhá porovnávání vrcholu tohoto zásobníku s nově přicházejícími tokeny.

3 Implementovaná rozšíření

3.1 BOOLOP

Pro implementaci rozšíření BOOLOP byla přidána další pravidla do precedenční analýzy, konkrétně pravidla pro **not**, **or**, **and** a **==** a zároveň další sémantická pravidla pro kontrolu datových typů vstupujících do těchto operací.

3.2 BASE

Implementace rozšíření BASE spočívá v přidání a vypořádání se s dalším způsobem zadávání číselných konstant. Do lexikální analýzy byly přidány stavy pro přijímání lexémů konstant zapsaných v binární, oktalové a hexadecimální soustavě, jejich převod na tokeny a předání dalším částem překladače.

3.3 FUNEXP

Rozšíření FUNEXP požadovalo podporu volání funkce jako součást výrazu a také výrazy jako parametry funkce při jejím volání. Do precedenční analýzy pro to byla přidány pravidla, která toto umožňují.

3.4 IFTHEN

V rozšíření IFTHEN se měl podporovat ternární operátor, **elif**, a konstrukce **if** bez části **else**. Právě poslední zmíněné rozšíření je naším překladačem podporováno.

3.5 TABUNARY

V rozšíření TABUNARY se mělo rozlišovat mezi unárním a binárním mínusem a plusem v precedenční analýze a podporou odsazování jak pomocí mezer, tak pomocí tabulátorů i prolínání těchto dvou způsobů odsazení. Část tohoto rozšíření, ve které se mělo rozlišovat právě mezi binárními a unárními operátory nebyla implementována. Pro část rozšíření, která se týká odsazování pomocí tabulátorů bylo třeba rozšířit zásobík u lexikální analýzy, aby bral v úvahu i tabulátory a ne jenom mezery.

4 Týmová spolupráce

Od začátku řešení projektu jsme se pravidelně scházeli nejdříve přibližně jedenkrát týdně a s blížícím se termínem odevzdání poté i vícekrát za týden. Pro rychlou komunikaci mezi sebou v reálném čase jsme si zvolili službu **Slack.com**, kde jsme ihned řešili všechny problémy, které nastaly.

Jako verzovací systém jsme použili **Git** hostovaný na stránce **GitHub**.

4.1 Rozdělení práce

- Roman Ondráček - implementace tabulky symbolů, implementace generátoru cílového kodu, pomocné abstraktní datové typy, testování
- Pavel Raur - implementace lexikální analýzy, implementace generátoru cílového kodu
- František Jeřábek - implementace syntaktické analýzy, implementace sémantické analýzy
- Radim Lipka - dokumentace, testování

5 Závěr

Cílem projektu bylo pochopit a prakticky si vyzkoušet metody a techniky vysvětlované v předmětech IFJ a IAL a ze všeho nejvíc pochopit, z čeho se skládá překladač a jak fungují jeho jednotlivé části.

Práce na projektu, ať už studium teorie, návrhy nebo samotná implementace nám přinesla řadu zkušeností, nových znalostí a rozšířila nám obzory na poli překladačů, jazyků a algoritmů.

6 Použitá literatura, nástroje, programy

- Alexandr Meduna a Roman Lukáš, Formální jazyky a překladače, prezentace k přednáškám
- Jan M. Honzík, Ivana Burgetová, Bohuslav Křena, Algoritmy, prezentace k přednáškám
- Zbyněk Křivka a Radim Kocman, Demonstrační cvičení IFJ, Implementace překladače IFJ19
- cppreference.com, <https://en.cppreference.com/>
- Manuálové stránky jazyka C
- GitHub, <https://github.com/>
- Slack, <https://slack.com/>