

**Dokumentace k projektu z předmětů  
IFJ a IAL Implementace překladače  
imperativního jazyka IFJ23  
Tým xkukht01, varianta – vv-BVS**

Kukhta Myron (xkukht01) 39%

6. prosince 2023

Malashchuk Vladyslav (xmalas04) 33%

Burylov Volodymyr (xburyl00) 28%

Pikulin Artemii (xpikul03) 0%

# 1 Úvod

Cílem projektu bylo vytvořit program na C, který načte zdrojový kód napsaný na IFJ23 a přeložil jste ho do jazzuka IFJ23code. Vracenou hodnotou programu je chybový kód (vrátí 0, pokud překlad proběhl bez chyb). Jazyk IFJ23 je zjednodušenou podmnožinou jazyka Swift2, což je moderní a stále populárnější jazyk používaný především při vývoji pro zařízení Apple.

## 2 Práce v týmu

### 2.1 Rozdělení práce

Rozdělili jsme práci mezi čtyři členy týmu. Každý dělal svou část, v případě potřeby mu pomáhali ostatní členové týmu. Celkové testování projektu jsme prováděli společně. Práce byla komplikována kvůli neplnění svých povinností jedním z členů týmu, který ji v posledních dnech opustil a nevykonával svou část práce(generátor kódu).

### 2.2 Způsob spolupráce

Při pracích na projektu jsme aktivně využívali systém pro správu verzí Git. Naše zdrojové kódy byly umístěny v repozitáři na online platformě GitHub, což zajistilo, že každý člen týmu měl snadný přístup k nejaktuálnější verzi projektu.

Pro specifické části projektu jsme vytvořili individuální větve, kde jsme prováděli testování a provedli nezbytné úpravy. Před začleněním do hlavní větve byl nezbytný tzv. "pull request", následovaný pečlivým posouzením kódu (code review) a schválením od jednoho či více členů týmu. Tímto postupem jsme zaručili kvalitu a integritu našeho projektu.

## 3 Návrh a implementace

### 3.1 Lexikální analýza

Lexikální analyzátor je implementován jako konečný automat, jak lze vidět ve zdrojovém souboru scanner.c a přidruženém hlavičkovém souboru scanner.h. Hlavní funkce, `get_token`, je zodpovědná za rozpoznávání tokenů a jejich zápis do struktury `Token`. Tato struktura obsahuje typ tokenu, atributy a informaci o tom, zda za ním následuje nový řádek. Atributy tokenu mohou být `integer`, `floating`, `string(Dynamic_string_t)`, `keyword(Klíčové slovo)` v závislosti na tokenu.

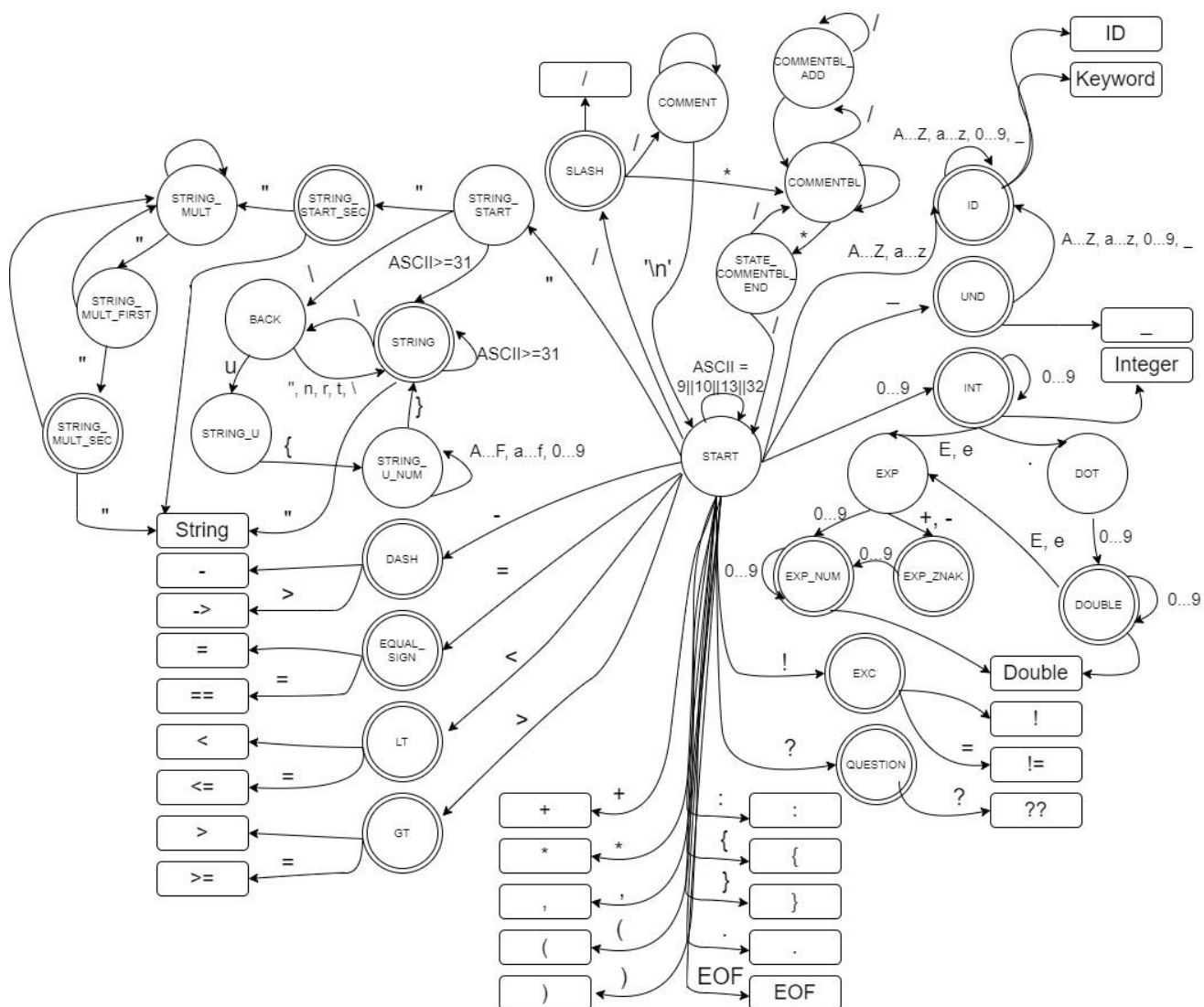
Konečný automat je implementován pomocí příkazu `switch-case`, kde každý stav odpovídá určitému znaku nebo skupině znaků. Funkce `get_token` čte znaky a přechází mezi stavy, dokud nedojde ke `SCAN_END`. V případě neshody s očekávanými znaky program signalizuje chybu (návratová hodnota 1).

Pokud funkce `get_token` identifikuje token typu `TOKEN_ID`, používá dvě další funkce: `keyword_or_id` a `keyword_or_not`. První kontroluje, zda je identifikátor klíčovým slovem, a případně přiřazuje odpovídající hodnotu atributu. Druhá se stará o případ, kdy na konci identifikátoru následuje otazník (?). Pokud je to klíčové slovo, přiřazuje hodnotu atributu `keyword` `KW_INT_NIL`, `KW_DOUBLE_NIL` nebo `KW_STRING_NIL`.

Kromě toho jsou v kódu vytvořeny pomocné proměnné, jako je `State_string_help_str` pro zpracování escape sekvencí (`\u{dd}`), `helpcomblock` pro počítání otevřených komentářů v bloku a `Dynamic_string_t` pro dočasné uložení dat tokenu pro další zpracování.

Během vývoje jsem se setkal s problémy při zápisu do znakového řetězce pomocí escape sekvence `'\u{dd}'`, kde je vyžadována práce s hexadecimálními číslicemi. Tento problém jsem vyřešil pomocí pomocných znaků a funkce `strtol`, která umožňuje převést číselný řetězec čísla zapsaného v hexadecimálním systému do desítkové soustavy.

### 3.1.1 Diagram konečného automatu



## 3.2 Syntaktická analýza

V projektu kompilátoru pro Swift hraje syntaktická analýza klíčovou roli při přeměně vstupního zdrojového kódu do strukturovaného formátu v souladu se syntaxí jazyka. Tento proces zahrnuje interakci mezi parserem a skenerem, kde parser žádá o tokeny od skeneru pomocí funkce `get_token`.

Lexikální analýza prováděná skenerem převádí symboly zdrojového kódu na tokeny. Při úspěšné lexikální analýze skener vrátí token (lexém ze vstupního kódu). Parser na základě LL gramatiky, začíná rekurzivně zpracovávat pravidla odpovídající aktuálnímu stavu. Pokud aktuální stav parseru nesouhlasí s očekávaným tokenem, vrátíme syntaktickou chybu. Tento proces se opakuje až do úplného rozboru vstupního kódu.

### 3.2.1 LL-gramatika(tabulka)

	let	var	func	ID	LITERAL	if	else	while	{	}	(	)	:	,	Int	Int?	Double	Double?	String	String?	nil	return	=	EOF	e	expression	_	->	\$
<SWIFT_PROG>	3	3	2	3			3	3																1					
<STATEMENT_LIST>	6	7		8			4	5																		9			
<IF_EXPRESSION>	11																										10		
<VAR_DEFINITION>				12						11													10						
<ID_TYPE_DEFITION>													13												14				
<ID_ACTION>										16													15						
<ASSIGNMENT>				18																							17		
<ARG_LIST>				19																					20				
<ARG>				21																					22				
<ARG_NAME>				23																					23				
<NEXT_ARG>															24										25				
<FUNCTION_DEFINITION>			25																										
<PARAM_LIST>				26																					27		26		
<PARAM>				28																							28		
<FUNCTION_STATEMENT_LIST>	29	29		29			29	29														30							
<RETURN_EXPRESSION>																									32	31			
<NAME_PARAM>				33																								34	
<ID_PARAM>				35																								36	
<NEXT_PARAM>															37										38				
<FUNC_TYPE_DEFINITION>																								40				39	
<TERM>				41	42																	42							
TYPE															43	44	45	46	47	48									

### 3.2.2 LL-gramatika(pravidla)

- <SWIFT\_PROG> -> EOF
- <SWIFT\_PROG> -> <FUNCTION\_DEFINITION> <SWIFT\_PROG>
- <SWIFT\_PROG> -> <STATEMENT\_LIST> <SWIFT\_PROG>
- <STATEMENT\_LIST> -> if <IF\_EXPRESSION> { <STATEMENT\_LIST> } else { <STATEMENT\_LIST> } <STATEMENT\_LIST>
- <STATEMENT\_LIST> -> while EXPRESSION { <STATEMENT\_LIST> } <STATEMENT\_LIST>
- <STATEMENT\_LIST> -> let <VAR\_DEFINITION>
- <STATEMENT\_LIST> -> var <VAR\_DEFINITION>
- <STATEMENT\_LIST> -> ID <ID\_ACTION>
- <STATEMENT\_LIST> -> e
- <IF\_EXPRESSION> -> EXPRESSION
- <IF\_EXPRESSION> -> LET ID
- <VAR\_DEFINITION> -> ID <ID\_TYPE\_DEFINITION> <ASSIGNMENT> <local\_scope>
- <ID\_TYPE\_DEFINITION> -> : <TYPE>

14. <ID\_TYPE\_DEFINITION> -> e
15. <ID\_ACTION> -> = <ASSIGNMENT>
16. <ID\_ACTION> -> (<ARG\_LIST>)
17. <ASSIGNMENT> -> EXPRESSION
18. <ASSIGNMENT> -> ID(<ARG\_LIST>)
19. <ARG\_LIST> -> <ARG><NEXT\_ARG>
20. <ARG\_LIST> -> e
21. <ARG> -> <ARG\_NAME> <TERM>
22. <ARG\_NAME> -> ID :
23. <ARG\_NAME> -> e
24. <NEXT\_ARG> -> ,<ARG><NEXT\_ARG>
24. <NEXT\_ARG> -> e
25. <FUNCTION\_DEFINITION> -> FUNC ID (<PARAM\_LIST>) <FUNC\_TYPE\_DEFINITION> {  
<FUNCTION\_STATEMENT\_LIST> }
26. <PARAM\_LIST> -> <PARAM> <NEXT\_PARAM>
27. <PARAM\_LIST> -> e
28. <PARAM> -> <NAME\_PARAM> <ID\_PARAM> : <TYPE>
29. <FUNCTION\_STATEMENT\_LIST> -> <STATEMENT\_LIST><FUNCTION\_STATEMENT\_LIST>
30. <FUNCTION\_STATEMENT\_LIST> -> return <RETURN\_EXPRESSION>  
<FUNCTION\_STATEMENT\_LIST>
31. <RETURN\_EXPRESSION> -> EXPRESSION
32. <RETURN\_EXPRESSION> -> e
33. <NAME\_PARAM> > ID
34. <NAME\_PARAM> -> \_
35. <ID\_PARAM> -> ID
36. <ID\_PARAM> -> \_
37. <NEXT\_PARAM> -> ,<PARAM><NEXT\_PARAM>
38. <NEXT\_PARAM> -> e
39. <FUNC\_TYPE\_DEFINITION> -> -> <TYPE>
40. <FUNC\_TYPE\_DEFINITION> -> e
41. <TERM> -> ID
42. <TERM> -> LITERAL
43. <TYPE> -> INT
44. <TYPE> -> INT?
45. <TYPE> -> DOUBLE
46. <TYPE> -> DOUBLE?
47. <TYPE> -> STRING
48. <TYPE> -> STRING?

### 3.3 Sémantická analýza

Během procesu sémantické analýzy proouváme tabulku symbolů, implementovanou ve formě vyváženého binárního stromu podle výšky. Tabulku kterou uchováváme na zásobníku. První tabulka symbolů v dolní části zásobníku obsahuje informace o globálních proměnných a funkcích. Nad ní jsou umístěny další tabulky, zodpovědné za lokální oblasti viditelnosti, kde jsou uchovávány lokální proměnné, které mohou překrývat globální proměnné. Každá proměnná a funkce jsou uloženy v uzle stromu a obsahují důležité informace, jako je název, datový typ, definice, inicializace atd. Tento přístup nám umožňuje efektivně sledovat, zda všechny proměnné a funkce zapojené v kódu byly správně inicializovány a odpovídají očekávanému typu.

Hlavní složitost naší analýzy spočívá v použití jednopřechodového přístupu kódu, což přináší určité obtíže při práci funkcemi, které používáme do jejich definice. Jsem musel vytvořit správný algoritmus, abych předpokládal správný typ funkcí a jejich parametry.

### 3.4 Zpracování výrazu

Zpracování výrazu je implementováno v souboru expression.c a provádí se funkcí rule\_expression, kterou volá syntaktický analyzátor. Výrazy se zpracovávají pomocí precedenční tabulky a zásobníku symbolů. Zpracování jednotlivých výrazů začíná inicializací zásobníku symbolů a uložením symbolu dolaru na jeho vrchol. Pak token po tokenu začíná funkce analyzovat vstupní výraz. Z každého tokenu funkce bere potřebné informace (symbol, datový typ - pokud je typ tokenu TOKEN\_IDE (identifikátor), pak ho najde v tabulkách symbolů a předá typ jeho hodnoty). Pak na základě nejbližšího k vrcholu zásobníku terminálu a vstupního symbolu, který je získán z následujícího tokenu, funkce provádí jednu z 4 operací - shift, reduce, equal, error (ne znamená 100% chybu, ale provede několik kontrol), všechny možné kombinace vstupních symbolů a prvních terminálů na vrcholu jsou popsány precedenční tabulkou.

Pokud znak z tabulky je <(S), vloží menšítko před prvním terminálem, vloží načtený symbol na vrchol zásobníku (funkce shift) a zavolá funkci get\_token.

Pokud znak z tabulky je >(R), ověří se, že mezi < a > je validní výraz (funkce test\_rule pro nalezení vhodného pravidla a funkce test\_semantics pro ověření typové kompatibility, vhodnosti typů operandů pro konkrétní operace a možná i, pokud je to nutné, přetypování některých operandů) a zredukuje získaný výraz do neterminálu E s výsledným datovým typem (funkce reduce).

Pokud znak z tabulky je =, provede push a zavolá funkci get\_token.

Pokud znak z tabulky je "N", začne redukovat všechno, co se zjistalo v zásobníku až do okamžiku, kdy prvním terminálem v zásobníku bude dolar. Pak zkontroluje symbol následujícího tokenu a v případě korektnosti výrazu nastaví logickou proměnnou success na true, což způsobí východ z loopu.

Na konci zkontroluje typ výrazu a na základě toho provede přetypování.

### 3.4.1 Precedenční tabulka

	<b>+-</b>	<b>*/</b>	<b>!</b>	<b>??</b>	<b>r</b>	<b>(</b>	<b>)</b>	<b>i</b>	<b>\$</b>
<b>+-</b>	>	<	<	>	>	<	>	<	>
<b>*/</b>	>	>	<	>	>	<	>	<	>
<b>!</b>	>	>	#	>	>	<	>	#	>
<b>??</b>	<	<	<	<	<	<	>	<	>
<b>r</b>	<	<	<	>	#	<	>	<	>
<b>(</b>	<	<	<	<	<	<	=	<	#
<b>)</b>	>	>	>	>	>	#	>	#	>
<b>i</b>	>	>	>	>	>	#	>	#	>
<b>\$</b>	<	<	<	<	<	<	#	<	#

### 3.5 Generování cílového kódu

Chybí kvůli nedodržení práce jednoho z členů týmu.



## 4. Algoritmy

### 4.1 Zasobník

V rámci našeho kompilátoru pro Swift jsme implementovali efektivní systém pro správu tabulek symbolů pomocí struktury zásobníku. Tato organizace umožňuje pružně a hierarchicky uchovávat informace o symbolech během celého procesu sémantické analýzy. V úvodu inicializujeme dno zásobníku. Dno zásobníku slouží jako úložiště pro tabulku symbolů globální oblasti viditelnosti, obsahující informace o globálních proměnných a funkcích. Během průchodu kódem postupně přidáváme nové tabulky symbolů na vrchol zásobníku. Každá nová tabulka symbolů reprezentuje lokální oblast viditelnosti. Tabulky symbolů se na vrcholu zásobníku ukládají a odstraňují dle aktuálního kontextu analýzy.

### 4.2 Binární Výškově Vyvážený Strom

Tabulka symbolů je implementována jako uzly v binárním vyváženém stromu. Uzly obsahují informace o symbolech, jako je název, datový typ, definice, inicializace. Při přidání nového symbolu provádíme vyvažování stromu pomocí speciálně připravené funkce. Tato operace nám umožňuje optimalizovat časovou složitost při hledání symbolu v tabulce.

### 4.3 Dynamický řetězec

Dynamický řetězec je pomocná struktura, která se používá k zápisu znaků do řádku, když předem neznáme počet těchto znaků. Dynamický řetězec se skládá ze tří prvků: samotného řádku, délky řádku a velikosti vyhrazené pro řetězec paměti. Má funkce `ds_init` pro inicializaci řádku, `ds_add_char` pro přidání jednoho znaménka do řádku, `ds_add_chars` pro přidání více znaků do řádku, `ds_copy` pro kopírování jednoho dynamického řádku do druhého, `ds_to_string` pro kopírování dynamického řádku do `char` a `ds_free` pro uvolnění paměti zvýrazněné pod dynamickým řádkem.