



# Dokumentace projektu

Implementace překladače imperativního jazyka  
**IFJ23**

Tým xsalta01, varianta vv-BVS

Autory:

Nadzeya Antsipenka xantsi00 - 23%

Lilit Movsesian xmovse00 - 23%

**Veranika Saltanova xsalta01** - 23%

Kirill Shchetiniuk xshche05 - 31%

6. prosince 2023

# Obsah

<b>1</b>	<b>Úvod</b>	<b>2</b>
<b>2</b>	<b>Implementace</b>	<b>2</b>
2.1	Lexikální analýza	2
2.2	Syntaktická analýza	2
2.2.1	Resení nedeterminismu v LL-tabulce	3
2.2.2	Rozhodování o ukončení výrazu	3
2.3	Sémantická analýza	3
2.4	Generování kódu	3
2.4.1	Začátek generování	3
2.4.2	Generování funkce	4
2.4.3	Generování výrazu	4
2.4.4	Generování volání funkce	4
2.5	Speciální algoritmy a datové struktury	4
2.5.1	Zásobník	4
2.5.2	Dynamické pole	4
2.5.3	Dynamický řetězec	5
2.5.4	Zdrojový soubor	5
2.5.5	Výškově vyvážený binární vyhledávací strom	5
2.6	Rozšíření na prémiové body	5
2.6.1	BOOLTHEN	5
2.6.2	CYCLES	5
2.6.3	FUNEXP	5
<b>3</b>	<b>Práce v týmu</b>	<b>6</b>
3.1	Rozdělení práce mezi členy týmu	6
<b>4</b>	<b>DKA</b>	<b>7</b>
<b>5</b>	<b>LL-gramatika</b>	<b>8</b>
<b>6</b>	<b>LL-tabulka</b>	<b>10</b>
<b>7</b>	<b>Precedenční tabulka</b>	<b>11</b>
<b>8</b>	<b>Použité zdroje</b>	<b>12</b>

# 1 Úvod

Cílem tohoto projektu bylo vytvořit program v jazyce C, který načte zdrojový kód zapsaný ve zdrojovém jazyce IFJ23 a přeloží jej do cílového jazyka IFJcode23.

## 2 Implementace

Překladač je sestaven z několika dílčích částí, které jsou popsány dále.

### 2.1 Lexikální analýza

Vypracování projektu jsme začali tvorbou diagramu deterministického konečného automatu (viz [DKA](#)), popisujícího lexikální analýzu. Na základě společně vypracovaného a odsouhlaseného diagramu jsme zahájili práce na implementaci projektu. Implementace lexikálního analyzátoru je v souboru `scanner.c`, pomocný modul pro práce s tokeny je implementován v souboru `token.c`. Lexikální analýza v projektu je provedena v funkci `source_code_to_tokens`, kde je realizován konečný automat, který převádí zdrojový kód na pole tokenů. Jednotlivé tokeny ukládáme jako inicializovanou strukturu `token_t`. Automat je realizován pomocí příkazu `switch`. Načtené symboly přidáváme do lexému a v každém konečném stavu lexém resetujeme. Identifikátory od klíčových slov oddělujeme porovnáním načteného lexému s klíčovými slovy. Skript zpracovává řádkové a blokové komentáře na základě načtených symbolů `'/'` a `'*'`, taktéž podporuje vnořené blokové komentáře. Pokud automat narazí na symbol odřádkování, nastaví u předešlého tokenu položku `has_newline_after` na hodnotu `true`. Pokud automat vrátí řetězcový literál, zavolá se funkce `verify_str`, která zpracovává řetězce a převádí je do formátu, který akceptuje interpret. Pokud řetězcový literál je víceřádkový, funkce `verify_str` zpracovává odřádkování a případně odstraňuje zbytečné bílé symboly na začátku řádku. Případné escape sekvence jsou převedeny na desítkový ASCII kód. Výsledkem provedené lexikální analýzy jsou tokeny, se kterými dál pracujeme v rámci projektu. Komunikace lexikálního analyzátoru s polem tokenů je provedena pomocí interface `TokenArray`.

### 2.2 Syntaktická analýza

Dále jsme společně sestavili LL-gramatiku (viz [LL-gramatika](#)). Syntaktický analyzátor je implementován na základě metody rekurzivního sestupu. Pro implementaci jsme použili LL-tabulku (viz [LL-tabulka](#)). Implementaci syntaktického analyzátoru máme v souboru `parser.c`. Pro každý non-terminal jsme vytvořili funkci se `switch` příkazy, ve kterých se rekurzivně zpracovávají jednotlivá pravidla. K zpracování terminálu používáme funkci `match`, která porovnává současný token s požadovaným tokenem a nastaví následující token jako současný. Syntaktická analýza se provádí ve dvou průchodech. Během prvního průchodu se sbírají všechny definice a deklarace funkcí do tabulky symbolů. Ve druhém průchodu provádíme sémantické kontroly a generování kódu. Analýzu výrazu provádíme na základě precedenční analýzy. Implementace je založená na precedenční tabulce (viz [Precedenční tabulka](#)). Analýza výrazu

je implementována ve souboru `expr_parser.c`. Pro získání indexu v precedenční tabulce používáme funkci `map_token`. Poté podle dvou indexu vyhledáváme jednotlivou akci, kterou musíme provést. Akci 'shift' odpovídá hodnota '1', 'reduce' - '2', 'equal' - '3', 'error' - '4'. Pro kontrolu místa deklarace funkce máme proměnné `inside_func`, `inside_loop` a `inside_branch`. Deklarace funkcí je povolena pouze v hlavním těle programu mimo cykly a podmíněné příkazy. Pro kontrolu, že příkaz `return` je uvnitř funkce, používáme proměnnou `inside_func`. Pro kontrolu, že non-void funkce obsahuje alespoň jeden příkaz `return` používáme proměnnou `has_return`.

### 2.2.1 Resení nedeterminismu v LL-tabulce

Během návrhu gramatiky se objevil problém nedeterminismu při zpracování non-terminálu `<CALL_PARAM>`, tento problém jsme vyřešili tak, že nejprve zkusíme to zpracovat podle pravidla číslo 48. Pokud se to nepodaří, zkoušíme to zpracovat podle pravidla číslo 49.

### 2.2.2 Rozhodování o ukončení výrazu

Kvůli tomu, že v IFJ23 mohou být výrazy na více řádcích a zároveň používáme konec řádku jako oddělovač jednotlivých příkazů, potřebovali jsme nějak rozhodnout zda je konec řádku zároveň i koncem výrazu. Na to používáme první token na následujícím řádku, přidáváme jej do výrazu a sledujeme, zda analyzátor výrazu vrátí chybu. Způsobí-li přidání tokenu chybu, ukončujeme analýzu výrazu, jinak pokračujeme.

## 2.3 Sémantická analýza

Pro sémantickou analýzu používáme tabulku symbolů a zásobník tabulek symbolů. Tabulka symbolů je implementována v souboru `symtable.c`. Jednotlivé tabulky symbolů jsou implementovány pomocí výškově vyváženého binárního vyhledávacího stromu. V rámci sémantické analýzy ověřujeme správnost použití proměnných, konkrétně jestli jsou definovány a deklarovány. Pro uložení funkcí používáme vlastní tabulku symbolů, která je určena jenom pro funkce. Během inicializace tabulky symbolů přidáváme do ní vestavěné funkce. Pomocí tabulky symbolů pro funkce ověřujeme správnost volání funkce a správnost typu návratového příkazu. Správnost volání ověřujeme pomocí funkce `check_func_signature`, která porovnává požadovanou a získanou signatury volání.

## 2.4 Generování kódu

Za cílový kód považujeme IFJcode23. Kód je generován na standardní výstup při druhém průchodu syntaktického analyzátoru. Jednotlivé funkce pro generování kódu jsou implementovány ve souboru `codegen.c`. Generování kódu provádíme hned na standardní výstup programu pomocí funkce `gen_line`.

### 2.4.1 Začátek generování

Na začátku generování kódu generujeme header a pomocné globální proměnné. Hned poté generujeme nepodmíněný skok na hlavní tělo programu. Mezi skokem a hlavním

tělem programu generujeme vestavěné funkce. Pro generace začátku programu voláme funkci `gen_header` a `gen_std_functions`.

### 2.4.2 Generování funkce

Před začátkem každé funkce generujeme nepodmíněný skok až za konec funkce pro obcházení funkce. Poté generujeme návěští, které odpovídá jménu funkci a vytváříme nový rámec. Dál generujeme nepodmíněný skok na definice lokálních proměnných, použitých během vykonání funkce. Hned za tímto skokem máme návěští pro návrat z definice lokálních proměnných a generujeme načtení parametrů v opačném pořadí z vrcholu zásobníku. Poté generujeme příkazy uvnitř funkce. Při naražení na příkaz `return`, v případě non-void funkce přiřazíme do proměnné `GF@$RET` výsledek návratového výrazu, generujeme zrušení rámce a návrat z funkce. V případě void funkce a naražení na poslední příkaz, generujeme hned za ním zrušení rámce a návrat z funkce. Za koncem funkce generujeme návěští pro definice lokálních proměnných. Generujeme definice proměnných a zpáteční skok do načtení parametru. Hned za tím generujeme návěští pro obcházení funkce.

### 2.4.3 Generování výrazu

Během syntaktické analýzy výrazu ukládáme všechny použité ve výrazech proměnné na zásobník. V případě funkce jako součásti výrazů generujeme volání funkce a dáváme návratovou hodnotu funkce (hodnotu proměnné `GF@$RET`) na zásobník.

### 2.4.4 Generování volání funkce

Rozhodli jsme, že budeme předávat všechny parametry funkcí přes zásobník v původním pořadí. Dále generujeme volání podle uvedeného ve zdrojovém kódu jména.

## 2.5 Speciální algoritmy a datové struktury

### 2.5.1 Zásobník

Jako pomocnou strukturu pro implementaci tabulky symbolů a syntaktického analyzátoru výrazu jsme použili zásobník. Zásobník je implementován ve souboru `stack.c` jako spojový seznam, kde vrchol zásobníku je první prvek seznamu. Zásobník umožňuje ukládání libovolného ukazatele na cokoliv. Operace nad zásobníkem umožňuje interface `Stack`.

### 2.5.2 Dynamické pole

Jako pomocnou strukturu pro různé případy jsme v souboru `list.c` implementovali dynamické pole, které umožňuje ukládání ukazatelů na jakékoliv datové typy. Dynamické pole je použito pro ukládání alokovaných ukazatelů pro další jednoduchou správu paměti. Operace nad polem umožňuje interface `DynamicArray`.

### 2.5.3 Dynamický řetězec

Jako alternativu poli znaku jsme implementovali dynamické řetězce pro jednoduchou práci s řetězci. Implementace dynamických řetězců je ve souboru `string_util.c`. Operace nad řetězcem umožňuje interface `String`.

### 2.5.4 Zdrojový soubor

Pro jednoduchou práci se zdrojovým kódem jsme implementovali datovou strukturu pro zdrojový kód. Implementace struktury je ve souboru `source_file.c`. Komunikace se zdrojovým kódem je provedena pomocí interface `SourceCode`.

### 2.5.5 Výškově vyvážený binární vyhledávací strom

Tabulky symbolů použité v rámci sémantické analýzy jsou implementovány pomocí výškově vyváženého binárního vyhledávacího stromu.

## 2.6 Rozšíření na prémiové body

### 2.6.1 BOOLTHEN

Realizovali jsme podporu typu `Bool` a `Bool?`, booleovských literálů `true` a `false` a základních booleovských operátorů `!`, `&&` a `||`. Dále jsme podpořili rozšířený podmíněný příkaz `if - else if - else` a udělali `else if` a `else` nepovinnými.

### 2.6.2 CYCLES

Podpořili jsme cyklus `for`, příkazy `break` a `continue`. Rozsah cyklu `for` musí být celočíselnou hodnotou.

### 2.6.3 FUNEXP

Podpořili jsme volání funkce jako součást libovolných výrazu, výrazy mohou být v parametrech volání funkce.

### 3 Práce v týmu

Hned na začátku jsme rozhodli, že nad všemi součástmi projektu budeme pracovat společně, aby každý pochopil, co se děje v každé části.

Jako verzovací systém jsme použili Git, jako vzdálený repozitář GitHub. Pro spolupráci jsme použili technologie Code With Me od JetBrains. Jako vývojové prostředí jsme použili IDE CLion.

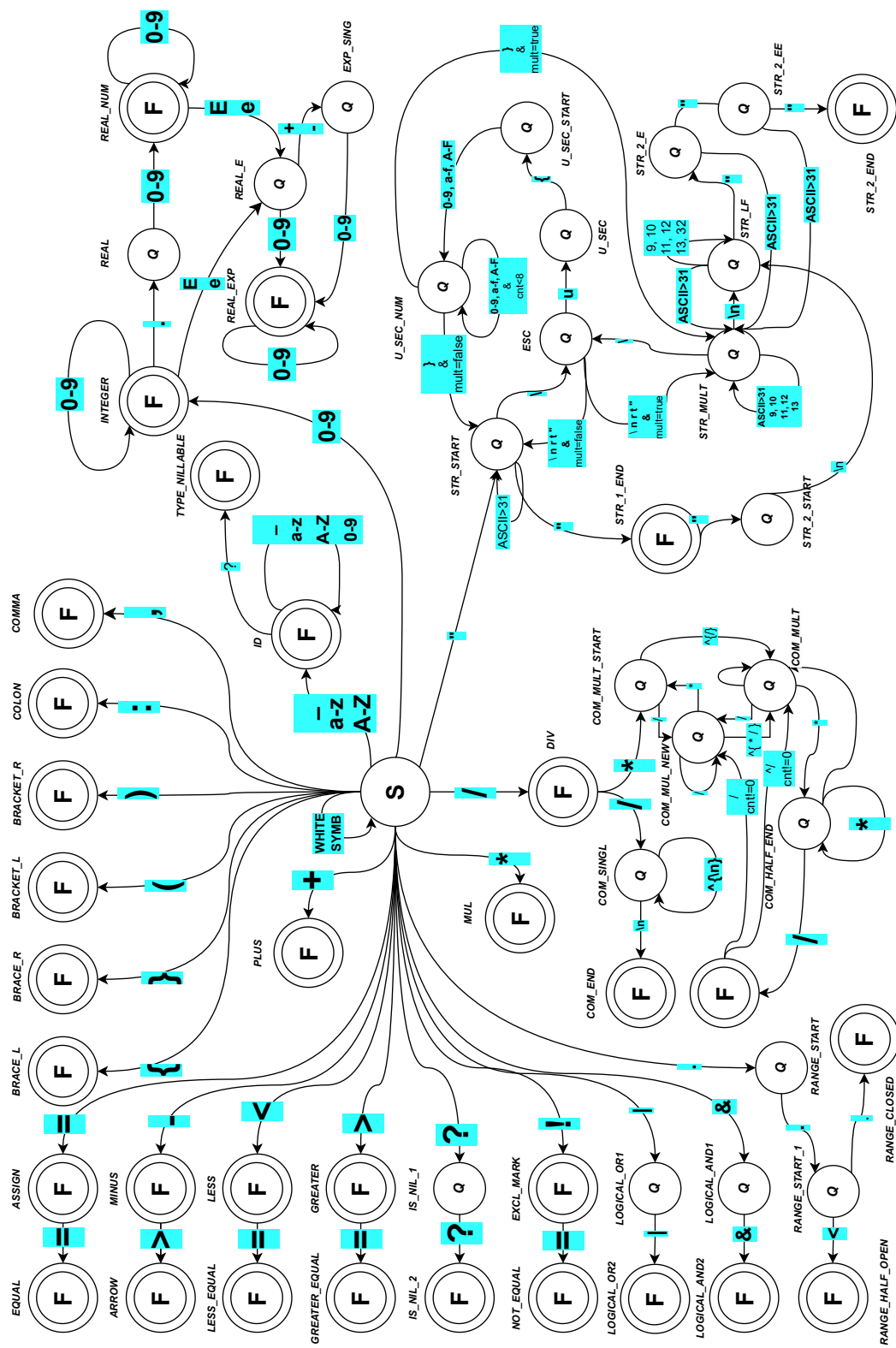
#### 3.1 Rozdělení práce mezi členy týmu

Tabulka 1 shrnuje rozdělení práce v týmu mezi jednotlivými členy.

Člen týmu	Přidělená práce
Nadzeya Antsipenka	lexikální analýza, syntaktická analýza, LL-gramatika, LL-tabulka, precedenční tabulka dokumentace, prezentace
Lilit Movsesian	lexikální analýza, syntaktická analýza, LL-gramatika, LL-tabulka, precedenční tabulka, tabulka symbolů (bvs), dokumentace, prezentace
<b>Veranika Saltanova</b>	vedení týmu, organizace práce, dohlížení na provádění práce, LL-gramatika, LL-tabulka, precedenční tabulka, syntaktická analýza, dokumentace
Kirill Shchetiniuk	generování cílového kódu, sémantická analýza, syntaktická analýza, syntaktická analýza výrazů, implementace pomocných datových struktur, LL-gramatika, LL-tabulka, precedenční tabulka, testování, rozšíření BOOLTHEN, CYCLES, FUNEXP

Tabulka 1: Rozdělení práce v týmu mezi jednotlivými členy

# 4 DKA



Obrázek 1: DKA



## 5 LL-gramatika

Jednotlivé non-terminály jsou označeny jako  $\langle \text{NON-TERMINAL} \rangle$ , terminály jsou označeny jako "terminal", konec souboru je označen jako \$, počáteční non-terminál je označen jako  $\langle S \rangle$ .

1.  $\langle S \rangle \rightarrow \langle \text{CODE} \rangle \$$
2.  $\langle \text{CODE} \rangle \rightarrow \langle \text{VAR\_DECL} \rangle \text{"eol"} \langle \text{CODE} \rangle$
3.  $\langle \text{CODE} \rangle \rightarrow \langle \text{LET\_DECL} \rangle \text{"eol"} \langle \text{CODE} \rangle$
4.  $\langle \text{CODE} \rangle \rightarrow \langle \text{FUNC\_DECL} \rangle \text{"eol"} \langle \text{CODE} \rangle$
5.  $\langle \text{CODE} \rangle \rightarrow \langle \text{WHILE\_LOOP} \rangle \text{"eol"} \langle \text{CODE} \rangle$
6.  $\langle \text{CODE} \rangle \rightarrow \langle \text{FO\_LOOP} \rangle \text{"eol"} \langle \text{CODE} \rangle$
7.  $\langle \text{CODE} \rangle \rightarrow \langle \text{BRANCH} \rangle \text{"eol"} \langle \text{CODE} \rangle$
8.  $\langle \text{CODE} \rangle \rightarrow \langle \text{ID\_CALL\_OR\_ASSIGN} \rangle \text{"eol"} \langle \text{CODE} \rangle$
9.  $\langle \text{CODE} \rangle \rightarrow \langle \text{RETURN} \rangle$
10.  $\langle \text{CODE} \rangle \rightarrow \text{"break"}$
11.  $\langle \text{CODE} \rangle \rightarrow \text{"continue"}$
12.  $\langle \text{CODE} \rangle \rightarrow \varepsilon$
13.  $\langle \text{RETURN} \rangle \rightarrow \text{"return"} \langle \text{RET\_EXPR} \rangle$
14.  $\langle \text{RET\_EXPR} \rangle \rightarrow \langle \text{EXPR} \rangle$
15.  $\langle \text{RET\_EXPR} \rangle \rightarrow \varepsilon$
16.  $\langle \text{VAR\_DECL} \rangle \rightarrow \text{"var"} \text{"id"} \langle \text{VAR\_LET\_TYPE} \rangle \langle \text{VAR\_LET\_EXP} \rangle$
17.  $\langle \text{LET\_DECL} \rangle \rightarrow \text{"let"} \text{"id"} \langle \text{VAR\_LET\_TYPE} \rangle \langle \text{VAR\_LET\_EXP} \rangle$
18.  $\langle \text{VAR\_LET\_TYPE} \rangle \rightarrow \text{":"} \langle \text{TYPE} \rangle$
19.  $\langle \text{VAR\_LET\_TYPE} \rangle \rightarrow \varepsilon$
20.  $\langle \text{VAR\_LET\_EXP} \rangle \rightarrow \text{"="} \langle \text{EXPR} \rangle$
21.  $\langle \text{VAR\_LET\_EXP} \rangle \rightarrow \varepsilon$
22.  $\langle \text{FUNC\_DECL} \rangle \rightarrow \text{"func"} \text{"id"} \text{"("} \langle \text{PARAM\_LIST} \rangle \text{")"} \langle \text{RET\_TYPE} \rangle \text{"{"} \langle \text{CODE} \rangle \text{"}"}$
23.  $\langle \text{RET\_TYPE} \rangle \rightarrow \text{"->"} \langle \text{TYPE} \rangle$
24.  $\langle \text{RET\_TYPE} \rangle \rightarrow \varepsilon$
25.  $\langle \text{PARAM\_LIST} \rangle \rightarrow \langle \text{PARAM} \rangle \langle \text{NEXT\_PARAM} \rangle$
26.  $\langle \text{PARAM\_LIST} \rangle \rightarrow \varepsilon$
27.  $\langle \text{PARAM} \rangle \rightarrow \langle \text{PARAM\_NAME} \rangle \text{"id"} \text{":"} \langle \text{TYPE} \rangle$
28.  $\langle \text{PARAM\_NAME} \rangle \rightarrow \text{"id"}$
29.  $\langle \text{PARAM\_NAME} \rangle \rightarrow \text{"_"}$
30.  $\langle \text{NEXT\_PARAM} \rangle \rightarrow \text{","} \langle \text{PARAM} \rangle \langle \text{NEXT\_PARAM} \rangle$
31.  $\langle \text{NEXT\_PARAM} \rangle \rightarrow \varepsilon$
32.  $\langle \text{BRANCH} \rangle \rightarrow \text{"if"} \langle \text{BR\_EXPR} \rangle \text{"{"} \langle \text{CODE} \rangle \text{"}" } \langle \text{ELSE} \rangle$
33.  $\langle \text{BR\_EXPR} \rangle \rightarrow \langle \text{EXPR} \rangle$
34.  $\langle \text{BR\_EXPR} \rangle \rightarrow \text{"let"} \text{"id"}$
35.  $\langle \text{ELSE} \rangle \rightarrow \text{"else"} \langle \text{ELSE\_IF} \rangle$

36. <ELSE> ->  $\varepsilon$   
 37. <ELSE\_IF> -> "if" <BR\_EXPR> "{" <CODE> "}" <ELSE>  
 38. <ELSE\_IF> -> "{" <CODE> "}"  
 39. <ELSE\_IF> ->  $\varepsilon$   
  
 40. <WHILE\_LOOP> -> "while" <EXPR> "{" <CODE> "}"  
  
 41. <FOR\_LOOP> -> "for" <FOR\_ID> "in" <EXPR> <RANGE> "{" <CODE> "}"  
 42. <FOR\_ID> -> "id"  
 43. <FOR\_ID> -> "\_"  
 44. <RANGE> -> "... " <EXPR>  
 45. <RANGE> -> "..<" <EXPR>  
  
 46. <CALL\_PARAM\_LIST> -> <CALL\_PARAM> <NEXT\_CALL\_PARAM>  
 47. <CALL\_PARAM\_LIST> ->  $\varepsilon$   
 48. <CALL\_PARAM> -> "id" ":" <EXPR>  
 49. <CALL\_PARAM> -> <EXPR>  
 50. <NEXT\_CALL\_PARAM> -> "," <CALL\_PARAM> <NEXT\_CALL\_PARAM>  
 51. <NEXT\_CALL\_PARAM> ->  $\varepsilon$   
  
 52. <ID\_CALL\_OR\_ASSIGN> -> "id" <NEXT\_ID\_CALL\_OR\_ASSIGN>  
 53. <NEXT\_ID\_CALL\_OR\_ASSIGN> -> "(" <CALL\_PARAM\_LIST> ")"  
 54. <NEXT\_ID\_CALL\_OR\_ASSIGN> -> "=" <EXPR>  
  
 55. <TYPE> -> "Int"  
 56. <TYPE> -> "Double"  
 57. <TYPE> -> "String"  
 58. <TYPE> -> "Bool"  
  
 59. <EXPR> -> "id"  
 60. <EXPR> -> "int\_literal"  
 61. <EXPR> -> "double\_literal"  
 62. <EXPR> -> "string\_literal"  
 63. <EXPR> -> "true\_literal"  
 64. <EXPR> -> "false\_literal"  
 65. <EXPR> -> "nil\_literal"  
 66. <EXPR> -> "("  
 67. <EXPR> -> "!"

## 6 LL-tabulka

[illegible]

Obrázek 2: LL(1)-tabulka

## 7 Precedenční tabulka

	unwrap	not	mul	add	??	rel	&&		(	i	)	\$
unwrap		>	>	>	>	>	>	>			>	>
not	<	<	>	>	>	>	>	>	<	<	>	>
mul	<	<	>	>	>	>	>	>	<	<	>	>
add	<	<	<	>	>	>	>	>	<	<	>	>
??	<	<	<	<	<	>	>	>	<	<	>	>
rel	<	<	<	<	<		>	>	<	<	>	>
&&	<	<	<	<	<	<	>	>	<	<	>	>
	<	<	<	<	<	<	<	>	<	<	>	>
(	<	<	<	<	<	<	<	<	<	<	=	
i	>		>	>	>	>	>	>			>	>
)			>	>	>	>	>	>			>	>
\$	<	<	<	<	<	<	<	<	<	<		

Tabulka 2: Precedenční tabulka

Operátory v tabulce jsou v pořadí od nejvyšší priority do nejnižší (**unwrap** - nejvyšší, **||** - nejnižší). Operátory se stejnou prioritou a asociativitou jsme zahrnuli do následujících skupin:

```

unwrap = {!}
not = {!}
mul = {*, /}
add = {+, -}
rel = {==, !=, <, >, <=, >=}
```

## 8 Použité zdroje

1. Přednášky, materiály a záznamy z demonstračních cvičení z předmětu IFJ
2. <https://www.cs.princeton.edu/courses/archive/spring22/cos320/LL1/>
3. <https://www.fit.vutbr.cz/~ikocman/llkptg/>
4. <https://developer.apple.com/documentation/swift/>