



IFJ, IAL 2020 - Dokumentace  
**Implementace překladače jazyka IFJ20**

Tým 48 - Varianta II

Implementovaná rozšíření: UNARY, FUNEXP

<b>Marek Filip</b>	( <b>xfilip46</b> )	– 27%
Vojtěch Bůbela	(xbubel08)	– 23%
Vojtěch Fiala	(xfiala61)	– 23%
Ondřej Míchal	(xmicha80)	– 27%

# Obsah

<b>1 Úvod</b>	<b>2</b>
<b>2 Práce v týmu</b>	<b>2</b>
2.1 Komunikace . . . . .	2
2.2 Verzovací systém . . . . .	2
2.3 Rozdělení práce . . . . .	2
<b>3 Implementace jednotlivých částí projektu</b>	<b>3</b>
3.1 Lexikální analýza . . . . .	3
<b>4 Syntaktická analýza</b>	<b>4</b>
<b>5 Sémantická analýza</b>	<b>5</b>
5.1 Generování cílového kódu . . . . .	6
<b>6 Datové struktury</b>	<b>7</b>
6.1 Tabulka symbolů . . . . .	7
6.2 Lineární seznam - 3AC . . . . .	8
<b>7 Závěr</b>	<b>8</b>
<b>8 Literatura</b>	<b>9</b>
<b>9 Příloha – Deterministický konečný automat pro lexikální analýzu</b>	<b>10</b>
<b>10 Příloha – LL gramatika</b>	<b>11</b>
<b>11 Příloha – LL tabulka</b>	<b>13</b>
<b>12 Příloha – Precedenční tabulka</b>	<b>14</b>

# 1 Úvod

Tato dokumentace popisuje implementaci překladače imperativního jazyka IFJ20, který je založen na jazyce Go. Implementace překladače byla zadána jako projekt do předmětů IFJ a IAL.

## 2 Práce v týmu

### 2.1 Komunikace

Týmová komunikace probíhala převážně osobně a nebo skrz aplikaci Discord. Náš tým měl každý týden pravidelné schůzky, na kterých jsme probírali pokrok, jakého jsme za uplynulý týden dosáhli a určovali si milníky, kterých chceme dosáhnout do další schůzky.

### 2.2 Verzovací systém

Pro správu verzí projektu jsme použili verzovací systém git a platformu GitHub.

### 2.3 Rozdělení práce

Základní rozdělení práce na jednotlivých částech projektu proběhlo na prvním meetingu. Následně jsme si pak s jednotlivými problémy ještě navzájem vypomáhali. Konkrétní rozdělení práce popisuje následující tabulka:

<b>Marek Filip (xfilip46)</b>	Organizační záležitosti, syntaktická analýza, sémantická analýza, zásobníková struktura, testování, dokumentace
Vojtěch Bůbela (xbubel08)	Lexikální analýza, generování cílového kódu, testování, dokumentace
Vojtěch Fiala (xfiala61)	Lexikální analýza, generování cílového kódu, testování, dokumentace, prezentace
Ondřej Míchal (xmicha80)	Lexikální analýza, sémantická analýza, tabulka symbolů, testování, dokumentace, git podpora

Na mírně nerovnoměrném rozdělení bodů jsme se společně dohodli na základě toho, že mezi výkonností jednotlivých členů byly rozdíly a shodli jsme se, že je tímto způsobem dorovnáme.

### 3 Implementace jednotlivých částí projektu

#### 3.1 Lexikální analýza

První krok při tvorbě překladače je lexikální analýza. Pro její implementaci jsme vytvořili deterministický konečný automat založený na diagramu 1.

Ze standardního vstupu jsou načítány jednotlivé znaky (v případě, že se jedná o bílý znak či komentář, je zahozen), které dále DKA analyzuje, přidává do struktury `string`<sup>1</sup> a skládá z nich tokeny. Token je reprezentován strukturou `token_t`, která se skládá z typu a atributu.

Typ tokenu může být buď identifikátor, literál, datový typ, operátor a nebo speciální typ – konec řádku, konec souboru a nebo neplatný typ. Jestliže se jedná o identifikátor, je jeho název dále porovnán s tabulkou klíčových slov a na základě výsledku tohoto srovnání se rozhodne, zda se jedná o klíčové slovo či název proměnné a nebo funkce.

Atribut tokenu je tvořen datovou strukturou `union` a v případě, že je typ tokenu literál, obsahuje jeho hodnotu.

Výsledný token je poté předán k dalšímu zpracování.

---

<sup>1</sup>Struktura `string` pochází z knihovny `str.h`, která je dostupná na stránkách projektu předmětu IFJ

## 4 Syntaktická analýza

Při návrhu jsme čerpali z populární mezinárodní knihy známé pod pseudonymem „Dragon Book“ [1]. Zádání nám určilo využití LL(1) gramatiky pro analýzu shora dolů a metodu precedenční tabulky pro analýzu zespodu nahoru. Jako první se sepsali pravidla LL(1) gramatiky. Pro správný postup byly použity instrukce získané z Dragon Book. Revize této gramatiky však probíhala několikrát, z důvodu neúplného pochopení zadání či špatného prvotního návrhu.

Sestavená gramatika (včetně gramatiky pro výrazy pro první nulté odevzdání) pak byla využita pro sestavení LL tabulky. Pro sestavení tabulky jsme postupovali dle instrukcí získaných ze záznamu democvičení IFJ ze dne 18. 10. 2011. Správnost LL tabulky (pro účely implementace) jsme pak zkontrolovali pomocí webové aplikace autorů Radima Kocmana a Dušana Koláře<sup>2</sup>.

Analýza zhora dolů je pak implementována metodou rekurzivního sestupu, kde se gramatická pravidla rozvíjí dle volání jiných funkcí nebo rekurzivního volání sebe sama. Tento postup se zdál jednodušší z důvodu banální implementace a absence zásobníku. Díky rozšíření FUNEXP jsme nemuseli řešit předávání kontroly nějakým ”hackem” a výsledný kód by měl proto působit „čistěji“.

Precedenční analýzu jsme započali po nultém pokusném odevzdání. V týmu jsme se dohodli na zvolení rozšíření FUNEXP a UNARY, takže se již od začátku počítalo se složitějším návrhem a následnou implementací. Rozšíření FUNEXP znamenalo hodně implementačních detailů a problému při různém stupni zanoření volání funkce a jejího správného vyhodnocení vzhledem na ostatní operátory včetně čárky, z pohledu tabulky to znamenalo přidat operátor `,`, který má nejmenší precedenci s tím, že na některých místech tabulky přibýlo více `-` operátorů.

Pro precedenční analýzu je použita struktura zásobníku (dále stack). Stack je implementoval obecně, díky chytrému využití maker pro abstrakci datového typu obsaženého ve stacku a díky shell skriptu generující nové soubory jsme mohli tento stack s lehkostí využít i pro potřeby sémantické analýzy a generování kódu. Na základu syntaktické analýzy je vybudována analýza sémantická.

---

<sup>2</sup><http://www.fit.vutbr.cz/~ikocman/llkptg/>

## 5 Sémantická analýza

Implementace sémantické analýzy blízce následuje syntaktickou analýzu. V úvodní části se značným problémem prokázala být kombinace analýzy zhora dolů a zdola nahoru. Důvodem byla tvrdá hranice mezi informacemi, které každá část analýzy má.

V rámci analýzy se ihned generuje tříadresný kód. Tento přístup byl zvolen s vidinou ulehčené práce. V průběhu implementace se ale objevila potřeba pro alespoň omezené použití abstraktního syntaktického stromu. Během analýzy se bohatě využívá několik datových struktur:

- tabulka symbolů uživatelských funkcí
- tabulka symbolů interních funkcí
- zásobník volaných funkcí
- zásobník dočasných identifikátorů
- zásobník proměnných na levé straně přiřazení
- zásobník typů, které by volaná funkce měla vrátit ve výrazu

V rámci vyhodnocování výrazů se každý výsledek redukce uložil do dočasné proměnné. Tento přístup umožnil už poměrně lehké propojení obou částí analýzy (zhora dolů a zdola nahoru). Zároveň ale tento přístup přináší značnou režii.

Rozšíření FUNEXP přineslo komplikaci ve formě nejasnosti vracených hodnot z volání funkcí. Tyto nejasnosti byly částečně vyřešeny sledováním "hloubky" zanoření ve výrazu a na základě kontextu se pro funkci určí možná vrácená hodnota. I bez rozšíření FUNEXP je ale volání funkcí v ifj20 náročnější, protože nevyžaduje, aby funkce před voláním byla deklarována. Proto, kromě již zmíněného mechanismu, se každé funkční volání přidává na zásobník. Po zdánlivě úspěšném zpracování analýzy se každé funkční volání porovná s tabulkami symbolů, kde jsou dostupné signatury interních i uživatelských funkcí.

Interně vyhodnocování výrazů analýza rozumí relačním operacím, ale neumožňuje práci s nimi mimo výrazy v hlavičkách prvků řídících běh programu (if, for).

Inspirace pro sémantickou analýzu: [2]

## 5.1 Generování cílového kódu

Generování cílového kódu IFJ20code je posledním krokem v běhu našeho překladače. Na začátku generování kódu dochází k vygenerování všech podporovaných vestavěných funkcí na začátek výstupu. Mezikód je generován na základě instrukcí v podobě tříadresného kódu, které jsou generátoru předány jako položky v lineárně vázaném seznamu. Při tvorbě struktury tříadresného kódu jsme čerpali z [3].

Argumenty tříadresného kódu jsou typově rozděleny - prefixy `i, f, s, d, b` značí, o jaký typ argumentu se jedná (`int, float, string, identifier, bool`). Na základě těchto prefixů jsou složky tříadresného kódu zpracovány a převedeny do formátu vhodného k vygenerování instrukce.

Následně se na základě jednotlivých instrukcí získaných z tříadresného kódu generuje výsledný mezikód.

## 6 Datové struktury

### 6.1 Tabulka symbolů

Dle zadání je implementace založena na tabulce s rozptýlenými položkami. Při implementaci se vycházelo z kódu, který byl vytvořen v rámci předmětu IJC.

Implementace tabulky symbolů probíhala v čtyřech iteracích, kdy se kromě funkcionality a implementačních detailů řešily problémy s komplexností API. Komplexnost se ve výsledku vyřešila tak, že veřejné API je v souboru `syntable.h` a zbytek je v interní `syntable-private.h`.

Tabulka symbolů je komplexní struktura, která funguje na několika úrovních:

- Funkční bloky oddělují symboly v jednotlivých funkcích. Slouží také pro reprezentaci globálního rámce. Jsou implementovány jako jednosměrný lineární seznam. Funkční bloky obsahují identifikátor funkce, parametry, návratové hodnoty a zásobník rámců.
- Zásobník rámců je lineárním seznam, kdy nejvrchnější prvek představuje nejhlouběji zanořený rámec. Struktura, reprezentující rámec, obsahuje odkaz na vnější rámec a samotnou tabulku symbolů.
- Tabulka symbolů je tabulka s rozptýlenými položkami, která obsahuje klíče, které jsou asociovány se symboly.
- Symbol je struktura, představující identifikátor, který obsahuje token (viz Lexikální analýza 3.1).

Lineární seznam byl pro většinu struktur zvolen pro lehkost implementace a nízkou paměťovou režii. Samotná hashovací tabulka je schopná se realokovat, pokud její kapacita začne být naplňována. K realokaci dojde, pokud lineární seznam v polích tabulky dosáhne zhruba 70% maximální „povolené“ délky.

Po vytvoření tabulka symbolů obsahuje vždy jeden funkční blok (globální), v něm zásobník rámců a v něm jeden rámec.



## 6.2 Lineární seznam - 3AC

Jednosměrný lineárně vázaný seznam je implementován ve 2 souborech – `three-address-code.c` a `three-address-code.h`. Je založen na domácí úloze č. 1 z předmětu IAL. Využíván je k předávání instrukcí tříadresného kódu. Skládá se z ukazatele na aktivní položku a na první položku v seznamu. Tělo jednotlivých položek obsahuje typ operandu, 3 argumenty (výsledek a 2 parametry operace) a ukazatel na další položku. Tělo tříadresného kódu je založeno na [3].

## 7 Závěr

Po zveřejnění zadání nás projekt zarazil svou rozsáhlostí a komplexitou. Projekt jsme vypracovávali od jeho zveřejnění, nicméně jeho kompletní dokončení se nám nepodařilo. Jednotlivé části jsme testovali za pomoci knihovny `gtest`<sup>3</sup> a větší celky pak za pomoci end-to-end testů a knihovny `bats`<sup>4</sup>. Při řešení projektu jsme využili znalosti získané na přednáškách předmětů IFJ a IAL, pro některé z nás to byla také první zkušenost s větším projektem. Jako tým jsme byli dohodnutí již předem, tedy k žádným problémům nedocházelo. Při řešení projektu se projevíly rozdíly ve schopnostech jednotlivých členů týmu a tyto rozdíly jsme se tedy rozhodli kompenzovat úpravou finálního rozdělení.

---

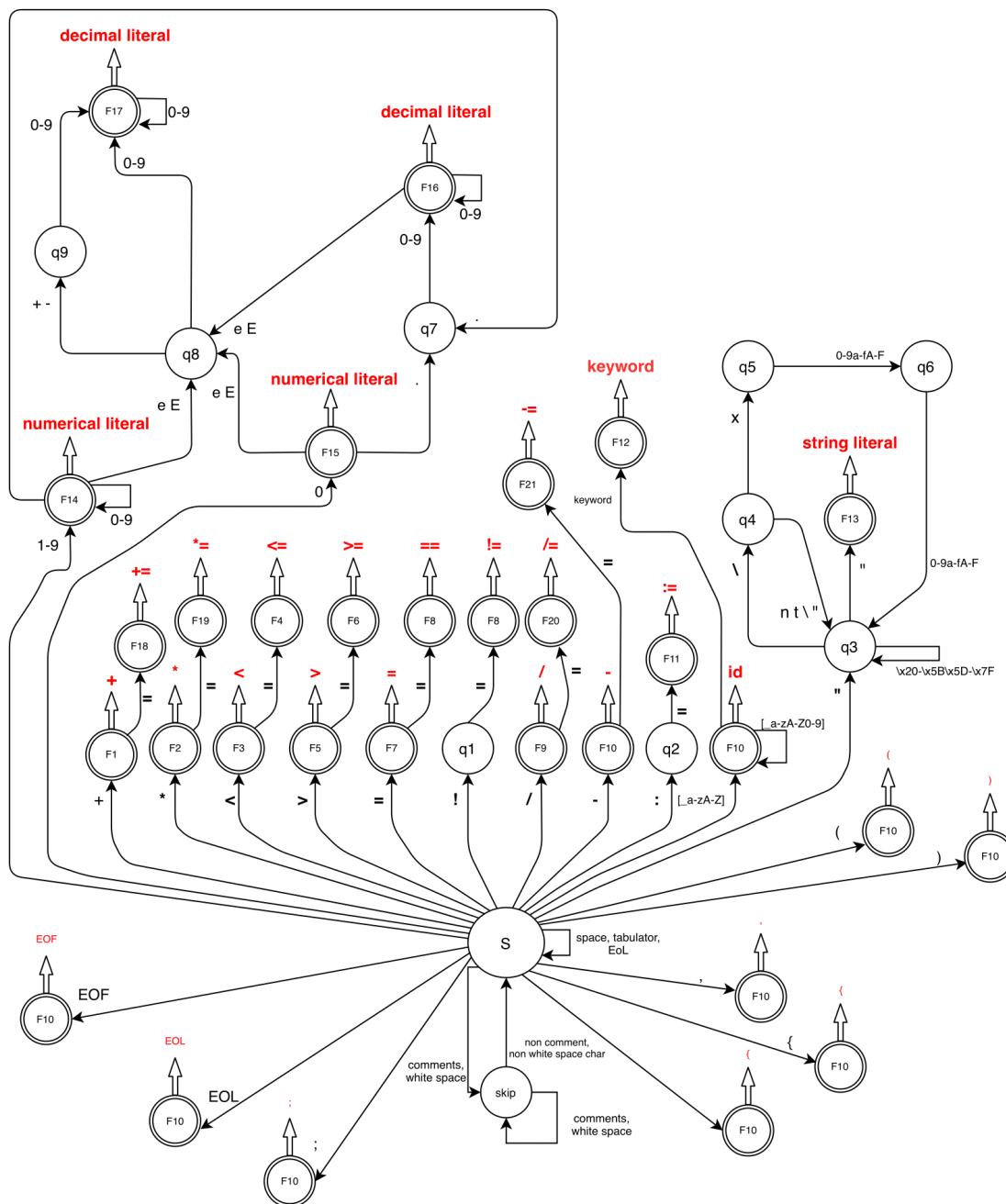
<sup>3</sup><https://github.com/google/googletest>

<sup>4</sup><https://github.com/bats-core/bats-core>

## 8 Literatura

- [1] AHO, A. V.; LAM, M. S.; SETHI, R.; aj.: *Compilers: Principles, Techniques, and Tools*, kapitola Syntax Analysis. Boston: Addison Wesley, 2006, ISBN 978-0321486813, str. 191.
- [2] AHO, A. V.; LAM, M. S.; SETHI, R.; aj.: *Compilers: Principles, Techniques, and Tools*, kapitola A simple Syntax-Directed translation. Boston: Addison Wesley, 2006, ISBN 978-0321486813, str. 64.
- [3] AHO, A. V.; LAM, M. S.; SETHI, R.; aj.: *Compilers: Principles, Techniques, and Tools*, kapitola Code Generation. Boston: Addison Wesley, 2006, ISBN 978-0321486813, str. 366.

## 9 Příloha – Deterministický konečný automat pro lexikální analýzu



Obrázek 1: Deterministický konečný automat

## 10 Příloha – LL gramatika

### Top-down grammar

1. **S** -> PROLOG FUNCS
2. **PROLOG** -> package id['main']
3. **FUNCS** -> func <NO\_EOL\_START> id PARAMS R\_PARAMS FUNCS  
| eps
4. **STMTS** -> STMT STMTS  
| eps
5. **STMT** -> id ID\_FIRST  
| if <NO\_EOL\_START> EXPR BLOCK else <NO\_EOL\_START> BLOCK <EOL>  
| for <NO\_EOL\_START> OPTDEF ; EXPR ; OPTASSIGN BLOCK <EOL>  
| RETURN <EOL>
6. ID\_FIRST -> DEF <EOL> | ASSIGN <EOL> | C\_PARAMS <EOL> | += EXPR | -= EXPR | \*=  
EXPR | /= EXPR

### Help rules

7. TYPE -> int | float64 | string
8. BLOCK -> { <NO\_EOL\_END> STMTS }

### Parameters definition

9. PARAMS -> ( PARAMS\_N
10. PARAMS\_N -> PARAM )  
| )
11. PARAM -> id TYPE PARAM\_N
12. PARAM\_N -> , id TYPE PARAM\_N  
| eps
13. R\_PARAMS -> ( R\_PARAMS\_N  
| eps
14. R\_PARAMS\_N -> R\_PARAM )  
| )
15. R\_PARAM -> TYPE R\_PARAM\_N
16. R\_PARAM\_N -> , TYPE R\_PARAM\_N  
| eps
17. C\_PARAMS -> ( C\_PARAMS\_N
18. C\_PARAMS\_N -> C\_PARAM )  
| )
19. C\_PARAM -> EXPR C\_PARAM\_N
20. C\_PARAM\_N -> , EXPR C\_PARAM\_N  
| eps

### Var definition and assign

21. DEF -> := EXPR
22. ASSIGN -> ID\_N = EXPRS

23.ID\_N -> , id ID\_N  
 | eps  
 24.EXPRS -> EXPR EXPR\_N  
 25.EXPR\_N -> , EXPR EXPR\_N  
 | eps

Helper 'for' rules

26.OPTDEF -> id DEF  
 | eps  
 27.OPTASSIGN -> id ASSIGN  
 | eps

Return statement

28.OPTRETURN -> RETURN  
 | eps  
 29.RETURN -> return OPTEXPRS  
 30.OPTEXPRS -> EXPRS  
 | eps

#### Bottom-up grammar

1. EXPR -> E
2.  $E \rightarrow id \mid lit$
3.  $E \rightarrow ( E )$
4.  $E \rightarrow E \text{ REL } E$
5.  $E \rightarrow E \text{ ADD } E$
6.  $E \rightarrow ( E )$
7.  $E \rightarrow \text{ADD } E$
8.  $E \rightarrow id ( )$
9.  $E \rightarrow id ( E )$
10.  $E \rightarrow id ( E, \dots )$

Get Index depending on the operation needed here

- 11.REL -> < | <= | > | >= | == | !=
- 12.ADD -> + | -
- 13.MUL -> \* | /

## 11 Příloha – LL tabulka

	package	id	func	if	expr	else	for	;	+=	-=	*=	/=	int	float64	string	{	}	)	,	(	:=	=	return	\$
\$	1																							
PROLOG	2																							
FUNCS		3																						4
PARAMS																22								
R_PARAMS																				28				
BLOCK																21								
STMTS		5		5			5										6							5
STMT		7		8			9																	10
ID_FIRST									14	15	16	17								12	13	11	12	
OPTDEF		47						48																
OPTASSIGN		49						50																
RETURN																								53
DEF																					40			
ASSIGN																			41			41		
C_PARAMS																				34				
TYPE													18	19	20									
PARAMS		23																						
PARAM		25																						
PARAMS_N																								
R_PARAM																								
R_PARAMS																				27	26			
R_PARAMS																				30	29			
R_PARAMS																				33	32			
C_PARAMS					35								31	31	31									
C_PARAMS_N					37																			
C_PARAM																								
ID_N																				39	38			
EXPRS					44															42			43	
EXPR_N		46		46			46	46									46			45				46
OPTRETURN																								51
OPTEXPRS					54																			

## 12 Příloha – Precedenční tabulka

	,	REL	ADD	MUL	(	)	id	lit	\$
,	=	<	<	<	<	=	<	<	
REL	>	>	<	<	<	>	<	<	>
ADD	>	>	> <	<	<	>	<	<	>
MUL	>	>	>	>	<	>	<	<	>
(	=	<	<	<	<	=	<	<	
)	>	>	>	>		>			>
id	>	>	>	>	=	>			>
lit	>	>	>	>		>			>
\$	<	<	<	<	<		<	<	

>|<      UNARY symbol → buď > nebo <

REL      >      <      >=      <=      '=='      !=

ADD      +      -

MUL      \*      /

Grammar rules

$E \rightarrow id \mid lit$       BASIC

$E \rightarrow ( E )$       BASIC

$E \rightarrow E \text{ REL } E$       BASIC

$E \rightarrow E \text{ ADD } E$       BASIC

$E \rightarrow ( E )$       BASIC

$E \rightarrow \text{ADD } E$       UNARY

$E \rightarrow id ( )$       FUNC

$E \rightarrow id ( E )$       FUNC

$E \rightarrow id ( E, \dots )$       FUNC