

# Fakulta Informačních Technologií Vysoké Učení Technické v Brně



## Dokumentace k projektu IFJ a IAL **Interpret jazyka IFJ16**

Tým 004, Varianta a/2/II

11. prosince 2016

Vedoucí: **Javorka Martin (xjavor18) 20%**

Marušin Marek (xmarus08) 20%

Utkin Kirill (xutkin00) 20%

Mrva Marián (xmrvam01) 20%

Shapochkin Victor (xshapo00) 20%

# Obsah

<b>1 Úvod</b>	<b>1</b>
<b>2 Řešení projektu</b>	<b>1</b>
2.1 Lexikální analyzátor	1
2.2 Syntaktický a sémantický analyzátor	2
2.3 Precedenční analýza	2
2.4 Interpret	3
<b>3 Algoritmy a datové struktury</b>	<b>3</b>
3.1 Heap sort	3
3.2 Knuth-Morris-Pratt	3
3.3 Tabulka symbolů	3
3.4 Garbage collector	4
3.5 Rozšíření	4
3.5.1 Boolean	4
<b>4 Práce v týmu</b>	<b>4</b>
4.1 Rozdělení práce	4
<b>5 Závěr</b>	<b>5</b>
<b>6 Metriky kódu</b>	<b>5</b>
6.1 Literatura	5
<b>A Konečný automat pro lexikální analýzu</b>	<b>6</b>
<b>B LL - gramatika</b>	<b>7</b>
<b>C Precedenční tabulka</b>	<b>9</b>

# 1 Úvod

Tato dokumentace popisuje návrh a implementaci interpretu jazyka IFJ16, který je zjednodušenou podmnožinou jazyka Java SE 8. Zvolili jsme si variantu zadání a/2/II, ve které máme naimplementovat funkci find pomocí algoritmu Knuth-Morris-Prattův, funkci sort pomocí algoritmu Heap Sort a tabulku symbolů pomocí tabulky s rozptýlenými položkami.

## 2 Řešení projektu

Projekt jsme se rozhodli rozdělit na několik částí:

- Lexikální analyzátor
- Syntaktický analyzátor
- Semantický analyzátor
- Věstavené funkce
- Interpret
- Testování

Syntaktický analyzátor je “srdcem” celého interpretu, řídí ostatní části a zároveň provádí i semantickou kontrolu. Lexikální analyzátor mu zasílá tokeny, které potom analyzuje. Když je zpracován celý program a je správný, tak se potom spustí interpret, který na základě instrukční pásky provede výslednou interpretaci programu.

### 2.1 Lexikální analyzátor

Lexikální analyzátor (přip. scanner) je implementován konečným automatem (obrázek č. 1), který na vstupu dostane zdrojový soubor s posloupností znaků, odstraní všechny nedůležité části jako komentáře a bílé znaky a rozdělí je na tzv. **tokeny**.

Token (tToken) je struktura, která obsahuje informaci o lexéme do které se postupně přerazují jí znaky za běhu scanneru. Pokud se dostane lexéma do konečného stavu automatu bez zajištění chyby (vyhovuje IFJ16), tuto informaci o lexéme dostane parser, který si token požádal pomocí funkce getNextToken. Mezi informací o lexéme patří “stav”, “data” potřebné pro parser a několik pomocných atributů pro správný běh v scanneru. Pokud scanner vrátí token s errorovým stavem, parser ho vyhodnotí, vypíše požadovaný error-ový kód a zastaví program.

Velmi důležitou částí lexikálního analyzátoru je rozpoznávání a ošetřování řetězců, oktálních čísel, speciálních znaků a čísel, které mají největší počet stavů.

Do scanneru jsme měli implementovat buffer, který uchovává u sebe historii tokenů, používá ho parser když potřebuje dostat určitý token, který už v minulosti dostal.

## 2.2 Syntaktický a sémantický analyzátor

Syntaktický analyzátor (případně parser) je jeden z najrozšířenějších a nejdůležitějších modulů interpretu. Má za úlohu hned několik hlavních operací nad uživatelským vstupem. Parser kontroluje konstrukce vstupního souboru, jinými slovy má kontrolovat syntaktické správnosti a sémantickou stranu programu. Vyvolá se hned po spuštění interpretu a inicializaci globálních proměn. Syntaktický a sémantický analyzátor celý čas komunikuje s modulem lexikálního analyzátoru od kterého si vyžaduje stále nový token na jeho zpracování. Jazyk IFJ16 je podmnožina jazyka JAVA SE 8, proto je naše syntaktická a sémantická analýza inspirovaná dvojfázovým přechodem - to znamená, že modul parser se volá dvakrát, přitom si pamatuje ve které fázi přechod se nachází. V prvním přechodu se kontroluje syntaktická správnost vstupního souboru a zapisuje proměny, třídy a funkce do tabulky symbolů. Druhá fáze přechodu zabezpečí sémantickou kontrolu, kde víme zjistit pomocí tabulky symbolů redefinici identifikátorů, případně jich nedeklaraci, datové typy a podobně. Obě fáze dvojfázového přechodu ŠÚ jsou založené na rekurzivním sestupu, který je popsán níže. Tabulka symbolů je důležitá část interpretu, do které parser mnohokrát zapisuje a předčítá. Z toho parser musí rozpoznat globální tabulku symbolů, a tabulky symbolů nižších úrovní, které jsou uloženy právě v uvedené dřív globálně tabulce. Na základě těchto operací dokáže modul parser vyhodnotit i typovou kontrolu. Tím zjistíme jestli nám sedí jednotlivé typy proměn, které musí být kompatibilní, počty parametrů ve funkcích a podobně.

Kromě jiného parser v druhém přechodu generuje instrukce vnitřního kódu, které v nejbližší fázi modul interpretu zpracovává a interpretuje. Z hlediska modulu parser se jedna zejména o generování instrukcí skoku, volání funkcí, instrukcí přiřazení a podobně. Všechny vygenerované instrukce parser ukládá na instrukční pásku, zatímco jednu instrukci reprezentuje tříadresný kód.

Rekurzivní sestup uvedený dřív je založený na LL-gramatice 1 (viz B.), podle které jsme implementovali celý modul parser. Modul je implementován jako sada funkcí, které představují jednotlivé neterminály. Naopak jednotlivé tokeny jsou vnímány jako terminální symboly LL gramatiky. V gramatice se nachází taky speciální symbol <expression>, který vyjadřuje předání řízení precedenční analýze, která zpracovává výrazy pro podstatné zjednodušení syntaktické a sémantické analýzy. V případě že se nedá podle získaného tokenu určit pravidlo gramatiky, parser detekuje chybu a na výstup posílá příslušný návratový kód na výstup programu, pokud se v žádné z fází dvojfázového přechodu parseru nenajde syntaktická nebo sémantická chyba.

Úloha syntaktické a sémantické analýzy tímto končí a řazení se předává modulu interpret.

## 2.3 Precedenční analýza

Na zpracování výrazů je použita syntaktická analýza “zdola nahoru”. Pokud parser narazí na výraz, předává řízení precedenční analýze.

Precedenční analýza je implementovaná pomocí precedenční tabulky a zásobníku. Na zásobník se ukládají terminaly a neterminály. Terminaly mohou být identifikatory, konstanty a operátory. Neterminál je výraz, speciální neterminál \$ (dollar) určující konec výrazu, a stav '<' z precedenční tabulky.

Precedenční analýza na základě pravidla vytvoří novou položku v tabulce symbolů jako výsledek výrazu a přidá tříadresný kód na instrukční pásku. Precedenční analýza odhalí chybný výraz (syntaktickou chybu) nenajdením vhodného pravidla.

## 2.4 Interpret

Interpret byl poslední částí projektu, má za úkol zpracovávat výstup programu a vytvářet cílový kód v případě úspěšného ukončení sémantické a syntaktické analýzy. Interpret postupně vykonává instrukci (třídresní kód) z instrukční pásy. Je volán z mainu po skončení syntaktické analýzy ve dvou krocích. V prvním kroku vyhodnotí všechny instrukce pro statické proměny, které jsou na samostatné instrukční pásce. Ve druhém kroku interpret vyhodnotí instrukce, které se nachází ve funkcích. Nejdřív si přemapujeme identifikátory (proměny a konstanty) z tabulky symbolů na lokální rámec. Pro interpretování jednotlivých se nejdřív ověří, jestli se identifikátory nachází v rámci. V případě že se identifikátor nenajde, hledá se v tabulce symbol pro třídu, ve které se funkce nachází.

Interpret prochází jednotlivé instrukci za sebou v cyklu. Výjimkou jsou instrukce skoku, které přímo odkazují na následující instrukci. Volání funkce je vyřešené rekurzivním voláním interpretu, kdy se mu předá instrukční páska pro volanou funkci a tabulka symbolů pro volanou funkci.

## 3 Algoritmy a datové struktury

V této kapitole budou popsány algoritmy, které byli implementovaný v našem projektu.

### 3.1 Heap sort

Tento algoritmus patří do skupiny třídících algoritmů a je jedním z nejefektivnějších algoritmu. Základní myšlenkou je využití datové struktury označované jako „halada“. Halada je binární strom o  $N$  úrovních, který má všechny uzly na všech  $(n-1)$  úrovních a na nejvzdálenější  $n$ -té úrovni má všechny uzly zleva, lze snadno implementovat polem. Pokud otec má index  $i$ , tak jsou jeho potomci na indexech  $2i$  (levý syn) a  $2i+1$  (pravý syn). Potom máme funkci zvanou Sift Down, která znovu staví hromadou porušenou v kořeni – proseje prvek v kořeni, který jako jediný porušuje pravidlo hromady. Algoritmus je velmi rychlý, ale není stabilní. Jeho náročná složitost je  $O(n \log n)$ .

### 3.2 Knuth-Morris-Pratt

Algoritmus se používá při vyhledávání podřetězce v jiném řetězci. Základní vyhledávací algoritmus se porovná podřetězec s řetězcem postupně po znaku. Má lineární časovou složitost  $O(m + n)$  a je význačný tím, že zachovává minimální možný počet celkových porovnání, tzn. neporovnává žádnou odpovídající dvojici (vzoru a řetězce) dvakrát. Před vyhledáváním vytvoříme si t.z vyhledávací tabulku, která bude mít pro každou pozici ve vzorku napsané číslo, které nám bude určovat, jaký prvek vzoru musíme porovnat s aktuálním znakem v řetězci. V případě že prvky neshodují, nemusí procházet už porovnané prvky a číslo z tabulky mu dá informace o tom, s kterým znakem vzoru má tento aktuální prvek řetězce porovnat před dalším posunutím.

### 3.3 Tabulka symbolů

Tuto tabulku symbolů jsme měli implementovat pomocí tabulky s rozptýlenými položkami. Pro hash tabulku používáme funkci djb2 (autor Bernstein). Tato tabulka je datovou strukturou pro

potřeby rychlého vyhledávání. Základem je obyčejné pole. Položkami pole jsou ukazatele na lineární seznamy, do nichž je možné uložit položky, které mají unikátní vyhledávací klíč. Důležitá část tabulky je tzv. hashovací funkce. Na vstupu dostává vyhledávací klíč (textový řetězec) a výstup je hash, který použijeme jako index do pole.

### 3.4 Garbage collector

Pro lepší správu paměti jsme se rozhodli využít garbage collector. Funguje tak, že vyhledává a uvolňuje úseky paměti, které již program (přip. proces) nepoužívá. Uvolnění paměti provádí funkce `freeGC()` na konci programu.

### 3.5 Rozšíření

Rozhodli jsme se v zadání implementovat jedno rozšíření.

#### 3.5.1 Boolean

Scanner pro rozšíření boolean v konečném automatu obsahuje navíc 2 klíčových slova (`true` a `false`) a jeden stav (`state`).

Precedenční analýza má přidáné nové pravidla pro operací nad typem boolean. Pro každé pravidlo je aj nový typ instrukci v tříadresním kódu. Zároveň byla rozšířená precedenční tabulka o tři nové operátory (`not`, `or` a `and`). Pro tyto operátory a pravidla je rozšířená typová kontrola, kde na dvou stranách výrazu očekává hodnotu typu boolean. Jsou upravené stavy určující přednost pro porovnávací operátory `==` (`equal`) a `!=` (`not equal`). Tyto operátory se teď můžou vyskytnout ve výrazu spolu s ostatními porovnávacími operátory (`<`, `>`, `<=`, `>=`). Parser rozšířil kontrolu tokenů o keyword `boolean` a rozšířil typovou kontrolu. Interpret přidal vykonávání instrukcí `not`, `and` a `or`.

## 4 Práce v týmu

Během práce na projektu, jsme měli různé způsoby komunikaci mezi členy týmu. Jako hlavní komunikační prostředek byla zvolena aplikace **Slack**, ve které jsme vytvořili skupinu pro náš projekt, a kde řešili problémy vznikající během naší práce. Také jsme používali aplikaci **MeisterTask** pro lepší přehled úkolů, které máme v plánu a které už máme vypracované.

Pro uchovávání našeho projektu jsme využívali repository **BitBucket**.

Během celého semestru jsme měli několik schůzek. Na první schůzce jsme si rozdělili jednotlivé části projektu. Potom schůzky probíhali pouze několikrát pro posouzení jednotlivých částí zadání.

### 4.1 Rozdělení práce

**Javorka Martin** – parser, interpret, dokumentace

**Marušin Marek** – syntaktická a sémantická analýza (parser), generování vnitřních tříadresných instrukcí, návrh LL gramatiky, pomocné struktury, předávání zpětného tokenu do scanneru, dokumentace

**Utkin Kirill** – práce s řetězci, vřetavené funkce, garbage collector + testování, dokumentace

**Mrva Marián** – scanner, interpret, dokumentace

**Shapochkin Victor** – ial, tabulka symbolů, garbage collector + testování, dokumentace

## 5 Závěr

Práce nad projektem byla moc zajímavá a přinesla nám spoustu zkušeností do budoucna, především v jazyku C a vývoje aplikace ve velkém týmu. Projekt byl časově náročný, ale moc užitečný. Jsme velmi rádi, že byl vytvořen termín 2. pokusného odevzdání, díky kterému jsme stihli doladit funkčnost projektu před koncovým odevzdáním.

Výsledný program funguje dle zadání a má jedné rozšíření **BOOLEAN**.

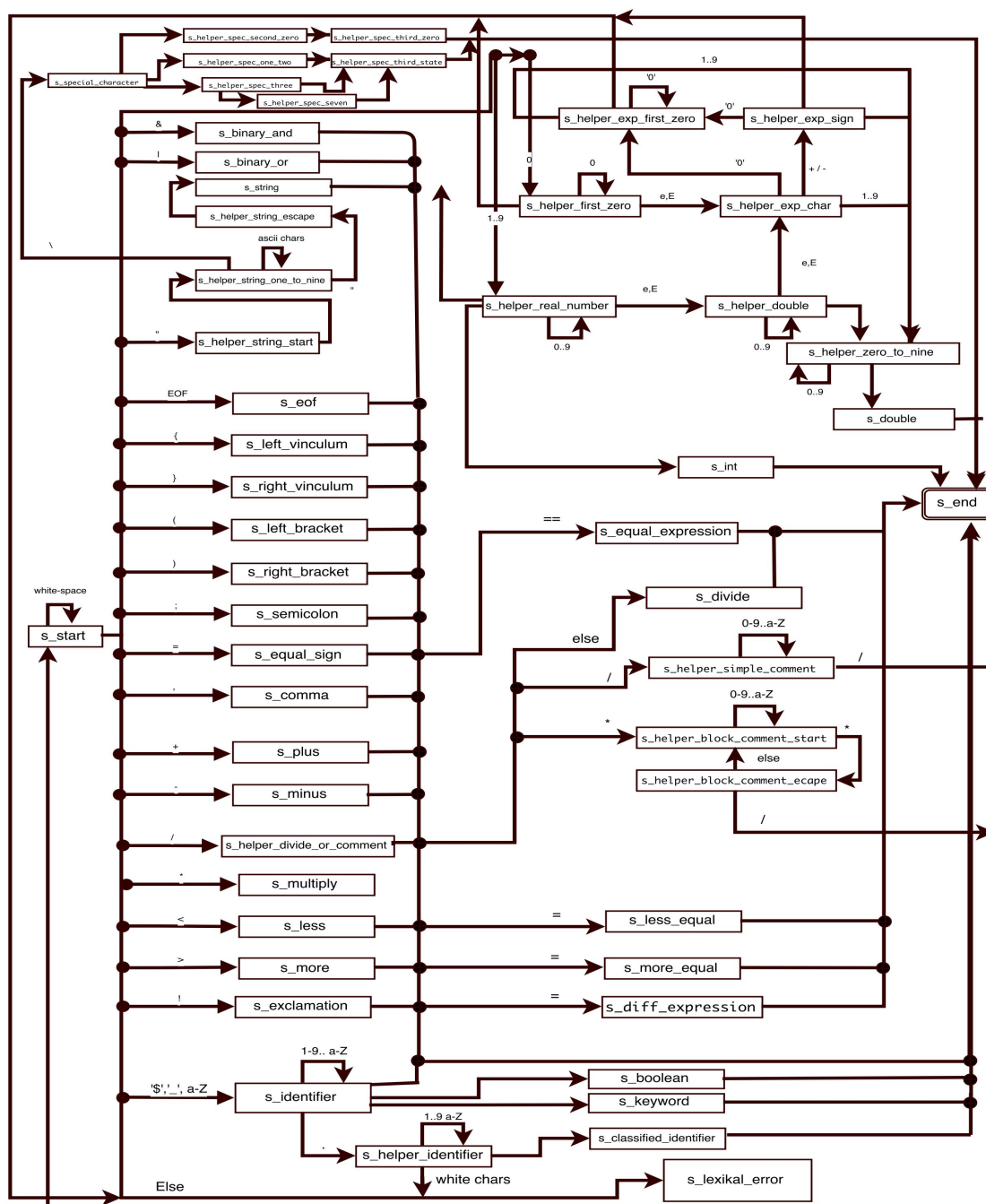
## 6 Metriky kódu

- Počet souborů: 34
- Řádků kódu: 11746

### 6.1 Literatura

- Přednášky, skripta a podklady z předmětu IFJ a IAL
- MEDUNA, Alexander. *Formal languages and computation: models and their applications*. New York: CRC Press, c2014. ISBN 978-1-4665-1345-7.

## A Konečný automat pro lexikální analýzu



Obrázek 1: Graf konečného automatu



## B LL – gramatika

1.  $\langle \text{prog} \rangle \rightarrow \langle \text{class\_list} \rangle \langle \text{eof} \rangle$
2.  $\langle \text{class\_list} \rangle \rightarrow \langle \text{class} \rangle \langle \text{class\_list} \rangle$
3.  $\langle \text{class\_list} \rangle \rightarrow \langle \text{e} \rangle$
4.  $\langle \text{class} \rangle \rightarrow \text{class id class\_body}$
5.  $\langle \text{class} \rangle \rightarrow \langle \text{e} \rangle$
6.  $\langle \text{class\_body} \rangle \rightarrow \langle \text{class\_body} \rangle \langle \text{class\_body} \rangle$
7.  $\langle \text{class\_body} \rangle \rightarrow \langle \text{e} \rangle$
8.  $\langle \text{class\_body} \rangle \rightarrow \text{static } \langle \text{type} \rangle \text{ id } \langle \text{class\_id\_list} \rangle$
9.  $\langle \text{class\_id\_list} \rangle \rightarrow ;$
10.  $\langle \text{class\_id\_list} \rangle \rightarrow = \langle \text{term} \rangle ;$
11.  $\langle \text{class\_id\_list} \rangle \rightarrow \langle \text{function\_definition} \rangle$
12.  $\langle \text{function\_definition} \rangle \rightarrow ( \langle \text{func\_arguments} \rangle ) \langle \text{body\_list} \rangle$
13.  $\langle \text{func\_arguments} \rangle \rightarrow \langle \text{type} \rangle \text{ id } \langle \text{func\_arguments2} \rangle$
14.  $\langle \text{func\_arguments} \rangle \rightarrow \langle \text{e} \rangle$
15.  $\langle \text{func\_arguments2} \rangle \rightarrow , \langle \text{type} \rangle \text{ id } \langle \text{func\_arguments2} \rangle$
16.  $\langle \text{func\_arguments2} \rangle \rightarrow \langle \text{e} \rangle$
17.  $\langle \text{body\_list} \rangle \rightarrow \langle \text{body} \rangle \langle \text{body\_list} \rangle$
18.  $\langle \text{body\_list} \rangle \rightarrow \langle \text{e} \rangle$
19.  $\langle \text{body} \rangle \rightarrow \langle \text{e} \rangle$
20.  $\langle \text{body} \rangle \rightarrow ;$
21.  $\langle \text{body} \rangle \rightarrow \text{return} \langle \text{term\_list} \rangle ;$
22.  $\langle \text{body} \rangle \rightarrow \text{if} ( \langle \text{term} \rangle ) \langle \text{body\_list} \rangle \text{else} \langle \text{body\_list} \rangle$
23.  $\langle \text{body} \rangle \rightarrow \text{while} ( \langle \text{term} \rangle ) \langle \text{body\_list} \rangle$
24.  $\langle \text{body} \rangle \rightarrow \langle \text{type} \rangle \text{ id } \langle \text{id\_list} \rangle ;$
25.  $\langle \text{body} \rangle \rightarrow \text{id} \langle \text{id\_statement} \rangle ;$
26.  $\langle \text{id\_statement} \rangle \rightarrow ( \langle \text{parlist} \rangle )$
27.  $\langle \text{id\_statement} \rangle \rightarrow = \text{id} ( \langle \text{id\_statement} \rangle )$

- 28.  $\langle \text{id\_statement} \rangle \rightarrow .\text{id}\langle \text{id\_statement} \rangle$
- 29.  $\langle \text{term\_list} \rangle \rightarrow \langle \text{term} \rangle$
- 30.  $\langle \text{term\_list} \rangle \rightarrow \langle \text{e} \rangle$
- 31.  $\langle \text{parlist} \rangle \rightarrow \langle \text{term} \rangle \langle \text{parlist2} \rangle$
- 32.  $\langle \text{parlist} \rangle \rightarrow \langle \text{e} \rangle$
- 33.  $\langle \text{parlist2} \rangle \rightarrow ,\langle \text{term} \rangle \langle \text{parlist2} \rangle$
- 34.  $\langle \text{parlist2} \rangle \rightarrow \langle \text{e} \rangle$
- 35.  $\langle \text{id\_list} \rangle \rightarrow \langle \text{e} \rangle$
- 36.  $\langle \text{id\_list} \rangle \rightarrow =\text{id}(\text{parlist})$
- 37.  $\langle \text{term} \rangle \rightarrow \text{id}$
- 38.  $\langle \text{term} \rangle \rightarrow \langle \text{string\_constant} \rangle$
- 39.  $\langle \text{term} \rangle \rightarrow \langle \text{int\_constant} \rangle$
- 40.  $\langle \text{term} \rangle \rightarrow \langle \text{float\_constant} \rangle$
- 41.  $\langle \text{term} \rangle \rightarrow \langle \text{expression} \rangle$
- 42.  $\langle \text{type} \rangle \rightarrow \text{null}$
- 43.  $\langle \text{type} \rangle \rightarrow \text{integer}$
- 44.  $\langle \text{type} \rangle \rightarrow \text{double}$
- 45.  $\langle \text{type} \rangle \rightarrow \text{string}$
- 46.  $\langle \text{type} \rangle \rightarrow \text{boolean}$
- 47.  $\langle \text{type} \rangle \rightarrow \text{void}$

## C Precedenční tabulka

	(	)	+	-	*	/	<	>	<=	>=	==	!=	i	\$	&&		!
(	<	=	<	<	<	<	<	<	<	<	<	<	<	err	<	<	<
)	err	>	>	>	>	>	>	>	>	>	>	>	err	>	>	>	>
+	<	>	>	>	<	<	>	>	>	>	>	>	<	>	<	>	<
-	<	>	>	>	<	<	>	>	>	>	>	>	<	>	<	>	<
*	<	>	>	>	>	>	>	>	>	>	>	>	<	>	>	>	<
/	<	>	>	>	>	>	>	>	>	>	>	>	<	>	>	>	<
<	<	>	<	<	<	<	err	err	err	err	>	>	<	>	>	>	<
>	<	>	<	<	<	<	err	err	err	err	>	>	<	>	>	>	<
<=	<	>	<	<	<	<	err	err	err	err	>	>	<	>	>	>	<
>=	<	>	<	<	<	<	err	err	err	err	>	>	<	>	>	>	<
==	<	>	<	<	<	<	<	<	<	<	err	err	<	>	<	<	<
!=	<	>	<	<	<	<	<	<	<	<	err	err	<	>	<	<	<
i	err	>	>	>	>	>	>	>	>	>	>	>	err	>	>	>	>
\$	<	>	<	<	<	<	<	<	<	<	<	<	<	err	<	<	<
&&	<	>	<	<	<	<	<	<	<	<	<	<	<	>	>	>	<
	<	>	<	<	<	<	<	<	<	<	<	<	<	>	<	>	<
!	<	>	>	>	>	>	>	>	>	>	>	>	<	>	>	>	<

Obrázek 2: Precedenční tabulka