

# VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ

## FAKULTA INFORMAČNÍCH TECHNOLOGIÍ



Dokumentácia k projektu z predmetu IFJ a IAL  
Implementácia prekladača jazyka IFJ20  
Tým 085, varianta 1

**Vedúci tímu:**

Marián Zimmermann, xzimme03 – 33%

**Ďalší členovia:**

Vanessa Jóriová, xjorio00 – 34%

Magdaléna Bellayová, xbella01 – 33%

Matušík Adrián, xmatus35 – 0%

# Úvod

Táto dokumentácia popisuje návrh a implementáciu prekladača jazyka IFJ20, ktorý je podmnožinou jazyka Go. Prekladač načítava zdrojový kód zapísaný jazykom IFJ20 a prekladá ho do cieľového jazyka IFJcode20 – medzikódu. Pre implementáciu sme si zvolili variantu 1, využívajúcu štruktúru binárneho stromu.

Riešenie projektu sa dá rozdeliť na 4 hlavné (neskôr popísané) úseky:

- Lexikálny analyzátor (scanner)
- Syntaktická a sémantická analýza metódou rekurzívneho zostupu (parser)
- Precedenčná analýza výrazov
- Generovanie medzikódu

## Rozdelenie práce v tíme

Marián Zimmermann – lexikálna analýza, generovanie kódu

Magdaléna Bellayová – precedenčná analýza

Vanessa Jóriová – syntaktická a sémantická analýza, generovanie kódu

Adrián Matušík – člen sa do vypracovania projektu nezapojil

## Práca v tíme

Na projekte sme pracovali spoločne, každý člen tímu si vzal na starosti jednu podmnožinu projektu s tým, že mu ostatní v prípade potreby pomáhali. Situáciu nám skomplikovalo vypadnutie štvrtého člena, ktorý mal na starosti generovanie kódu, ktoré sa teda rozdelilo medzi ostatných.

Ako verzovací systém sme využili Git a ako vzdialený repozitár Github.

Keďže nám aktuálna situácia nedovolila osobné stretnutia, väčšina komunikácie prebiehala písomne (cez skupinovú konverzáciu) alebo hlasovo na Discorde.

## Lexikálny analyzátor

Lexikálny analyzátor (scanner) je implementovaný ako deterministický konečný automat. Pozostáva z cyklu while a switchu do ktorého sa postupne posielajú znaky (ignoruje biele znaky a komentáre) a na základe nich vyhodnocuje o aký token sa jedná po prípade ukladá jeho hodnotu do stringu strtmp.

## Syntaktický a sémantický analyzátor

Syntaktický a sémantický analyzátor (neskôr len *parser*) je ústrednou a hlavnou časťou nášho programu, ktorá sa spúšťa hneď na začiatku a svoju činnosť ukončuje až po spracovaní vstupného kódu v jazyku IFJ20. Počas svojej činnosti úzko spolupracuje s ostatnými modulmi:

- Žiada scanner o tokeny, pomocou ktorých vyhodnocuje, či je zdrojový kód napísaný v jazyku popísanom vytvorenou gramatikou
- Pokiaľ narazí na výraz, predá riadenie precedenčnej analýze, ktorá výraz vyhodnotí

- Priamo z parseru sa generuje cieľový medzikód

Syntaktická a sémantická analýza prebiehajú súbežne. Zdrojový kód sa spracuje jedným behom, syntaktická analýza prijíma tokeny od lexikálneho analyzátoru a rekurzívnym zostupom vyhodnocuje, či je zdrojový súbor gramaticky správne. V prípade, že syntaktická analýza narazí na výraz, predá riadenie precedenčnej analýze, ktorá tokeny až po koniec výrazu spracuje.

Pri implementácii **syntaktickej analýzy** sme ako najzásadnejší problém vnímali implementáciu lexikálneho analyzátoru, ktorý načíta token po tokene bez možnosti token „vrátiť“. To čiastočne vyriešil mód *Peek*, ktorý nahliadol na ďalší token a uložil ho, aby bol vrátený pri ďalšom zavolaní. V niektorých prípadoch sa ale z tohto dôvodu program implementačne odchyľuje od LL-gramatiky, popr. sú využité iné spôsoby, ako rozhodnúť o rozložení neterminálu (sémantická analýza), podobné odchýlky sú ale vykonávané s dôrazom na to, aby nenarušili gramatiku popísanú LL gramatikou a prijímali rovnaký jazyk ako ten LL gramatikou popísaný.

Pri implementácii parseru sme sa pridržali tejto gramatiky (neterminály sú písané veľkým písmenom:

1. START  $\rightarrow$  PROLOG PROG
2. PROLOG  $\rightarrow$  package main
3. PROG  $\rightarrow$  FUNCTION PROG
4. PROG  $\rightarrow$  *eps*
5. FUNCTION  $\rightarrow$  func id ( ARGS) RETURN\_F {S\_BODY }
6. ARGS  $\rightarrow$  ARG ARG\_NEXT
7. ARG\_NEXT  $\rightarrow$  *eps*
8. ARG\_NEXT  $\rightarrow$  , ARG ARG\_NEXT
9. ARG  $\rightarrow$  id DATA\_TYPE
10. RETURN\_F  $\rightarrow$  ( RETURN\_TYPES)
11. RETURN\_F  $\rightarrow$  *eps*
12. RETURN\_TYPES  $\rightarrow$  RETURN RETURN\_NEXT
13. RETURN  $\rightarrow$  DATA\_TYPE
14. RETURN\_NEXT  $\rightarrow$  , RETURN RETURN\_NEXT
15. RETURN\_NEXT  $\rightarrow$  *eps*
16. S\_BODY  $\rightarrow$  STATEMENT S\_BODY
17. S\_BODY  $\rightarrow$  *eps*
18. STATEMENT  $\rightarrow$  DEFINITION
19. STATEMENT  $\rightarrow$  ASSIGN
20. STATEMENT  $\rightarrow$  IF
21. STATEMENT  $\rightarrow$  FOR
22. STATEMENT  $\rightarrow$  F\_CALL
23. STATEMENT  $\rightarrow$  RETURN\_S
24. DEFINITION  $\rightarrow$  id := <expr>
25. ASSIGN  $\rightarrow$  id ID\_NEXT = FUNC\_OR\_EXPR
26. FUNC\_OR\_EXPR  $\rightarrow$  <expr> EXPR\_NEXT
27. FUNC\_OR\_EXPR  $\rightarrow$  F\_CALL
28. ID\_NEXT  $\rightarrow$  , id ID\_NEXT
29. ID\_NEXT  $\rightarrow$  *eps*
30. EXPR\_NEXT  $\rightarrow$  , <expr> EXPR\_NEXT
31. EXPR\_NEXT  $\rightarrow$  *eps*
32. IF  $\rightarrow$  if <expr> { S\_BODY } else { S\_BODY }
33. FOR  $\rightarrow$  for FOR\_DEFINITION ; <expr> ; FOR\_ASSING { S\_BODY }
34. FOR\_DEFINITION  $\rightarrow$  DEFINITION
35. FOR\_DEFINITION  $\rightarrow$  *eps*

36.  $\text{FOR\_ASSIGN} \rightarrow \text{F\_ASSIGN}$
37.  $\text{F\_ASSIGN} \rightarrow \text{id ID\_NEXT} = \langle \text{expr} \rangle \text{EXPR\_NEXT}$
38.  $\text{FOR\_ASSIGN} \rightarrow \text{eps}$
39.  $\text{F\_CALL} \rightarrow \text{id}(\text{PARAMS})$
40.  $\text{PARAMS} \rightarrow \text{TERM TERM\_NEXT}$
41.  $\text{PARAMS} \rightarrow \text{eps}$
42.  $\text{TERM\_NEXT} \rightarrow , \text{TERM TERM\_NEXT}$
43.  $\text{TERM\_NEXT} \rightarrow \text{eps}$
44.  $\text{RETURN\_S} \rightarrow \text{return RET\_EXPR}$
45.  $\text{RET\_EXPR} \rightarrow \langle \text{expr} \rangle \text{EXPR\_NEXT}$
46.  $\text{RET\_EXPR} \rightarrow \text{eps}$
47.  $\text{TERM} \rightarrow \text{int\_value}$
48.  $\text{TERM} \rightarrow \text{float\_value}$
49.  $\text{TERM} \rightarrow \text{string\_value}$
50.  $\text{TERM} \rightarrow \text{id}$
51.  $\text{DATA\_TYPE} \rightarrow \text{int}$
52.  $\text{DATA\_TYPE} \rightarrow \text{float64}$
53.  $\text{DATA\_TYPE} \rightarrow \text{string}$

Väčšina **sémantickej analýzy** prebieha tiež v parseri, kedy parser sémanticky okontroluje takmer všetko, k čomu má prístup, čím sa minimalizuje nutnosť behových kontrol.

IFJ20 je silne typovaný jazyk a to umožňuje statické priradenie typu novo definovanej premennej. Najprv sa vyhodnotí výraz za definíciou novej premennej, precedenčná analýza oboznámi parser o type novej premennej a premenná sa s informáciou o type a fakte, že sa jedná o premennú, nie o funkciu, uloží na vrchol zásobníka tabuliek symbolov.

Sémantické kontroly volania funkcie sú o niečo zložitejšie vďaka faktu, že funkcia môže byť zavolaná ešte pred definíciou. Implementácia lexikálneho analyzátora nepodporovala sémantickú kontrolu v dvoch behoch, ktorá by samotnú sémantickú analýzu zjednodušila, funkcie teda sémanticky kontrolujeme v jednom behu spôsobom:

- Pokiaľ sa narazí na definíciu funkcie ešte pred jej volaním, je názov funkcie uložený do tabuľky symbolov aj s informáciami o počte, type, názve parametrov a počte a type návratových hodnôt. Zaznačený je aj fakt, že bola definovaná. Definovanie funkcie s rovnakým menom vedie na chybu, rovnako ako aj následné volanie funkcie so zlými parametrami/návratovými hodnotami
- Pokiaľ sa narazí na volanie dosiaľ nedefinovanej funkcie, do tabuľky symbolov sú poznačené všetky informácie, ktoré z volania môžeme získať – typ parametru, návratových hodnôt a jej názov. Pokiaľ sú tieto informácie nekompletné (napr. \_ v návratových hodnotách) a narazí sa na ďalšie volanie, ktoré informácie rozšíri (volanie so známym typom návratovej hodnoty), informácie sú aktualizované. Volanie s konfliktnými informáciami vedie na sémantickú chybu. Pokiaľ sa neskôr v programe nájde definícia funkcie odpovedajúca zaznačeným informáciám, tabuľka symbolov je aktualizovaná o fakt, že bola funkcia definovaná – definícia nekompatibilná s volaniami vedie na chybu. Na konci programu sa vyhodnotí, či boli všetky funkcie definované – nedefinované funkcie spôsobia chybu.

Parser využíva funkcie na generovanie kódu, ktorým predáva informácie získané syntaktickou a sémantickou analýzou.

# Precedenčná analýza

Zavolá ju scanner, keď narazí na výraz. Tu sa berú tokeny, kým sa nenarazí na koniec výrazu, alebo nejaký chybu. Keď sa prečíta token tak sa podľa toho, čo je na vrchole zásobníka určí stav z precedenčnej tabuľky. Keď je v precedenčnej tabuľke stav LO tak sa vloží do zásobníka stav LO a potom až token. V prípade, že je stav HI tak sa ide skracovať výraz podľa pravidla. Do druhého zásobníka sa uložia znaky až po stav LO a súčasne sa vyberú z prvého zásobníka. Podľa počtu a typu operandov je následne vybrané pravidlo. Keď sa našlo pravidlo tak celá časť sa nahradí neterminálom, ktorý sa vloží do prvého zásobníka. Ak nastane stav EQ, tak sa len vloží token do zásobníka a ak nastane stav ER, tak sa vráti syntaktický chyba. Súčasný token sa porovnáva z najvrchnejším terminálom, kým ešte má, čo redukovať, ak nie a ešte súčasný token nebol vložený do zásobníka tak sa vloží. Ukončenie bez chyby nastane, keď príde koniec riadku a pomocou pravidiel sa bude vedieť na všetko na čo ešte neboli uplatnené pravidlá skrátiť tak, že zostane len jeden neterminál. Precedenčná analýza vracia TRUE ak prebehla v poriadku a nastavuje flag na typ výsledného neterminálu.

Využívajú sa tu aj funkcie generovania kódu, kde posielajú generátoru potrebné informácie na vygenerovanie kódu pre výrazy.

## Dátové štruktúry

### Tabuľka symbolov

Pri sémantickej analýze sa využíva zásobník tabuliek symbolov. Jednotlivé tabuľky symbolov sú implementované pomocou binárneho stromu. Uzly binárneho stromu si štruktúry obsahujúce kľúč (názov funkcie/premennej) a dáta – skutočnosť, či sa jedná o premennú alebo funkciu, či je funkcia definovaná, upravený názov rozšírený o návestie scopu využívaný pri generácii kódu, počet a typy parametrov a návratových hodnôt...

Jednotlivé tabuľky symbolov sú zviazané do zásobníka. Vkladať nové uzly sa dá len na aktuálny vrchol zásobníka (mimo volania dosiaľ nedefinovanej funkcie, kedy sa informácie značia do globálnej tabuľky symbolov – spodok zásobníka), vyhľadáva sa v zásobníku celom a vracia najvyššie položenú zhodu s kľúčom – to nám umožňuje prácu so scopmi v tele programu a dodržiavať viditeľnosť a prekryvanie premenných.

### String

Pre prácu s identifikátormi neobmedzenej dĺžky využívame knižnicu str.h, ktorú sme prevzali z IFJ stránok.

### Zásobník

Pri práci s precedenčnou analýzou sa využíva zásobník, kde si pre každý element ukladá symbol, typ, hodnotu a nasledovníka. Samotný zásobník má v sebe len pointer na vrchol zásobníka.

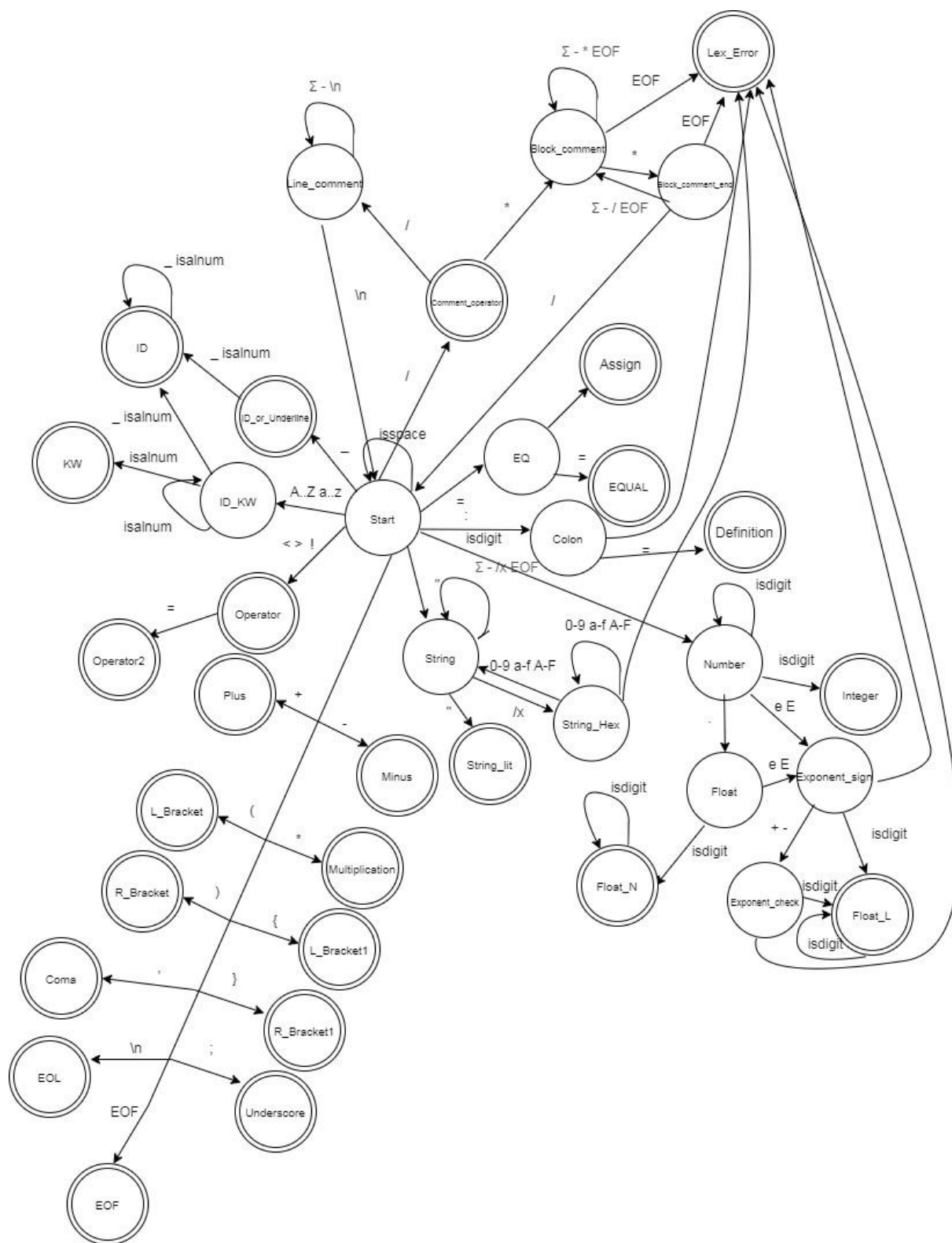
Vytvorili sme jednoduché funkcie na prácu so zásobníkom aby sme do neho vedeli vkladať aj vyberať hodnoty, zisťovať vrchol aj vrchný terminál, vyprázdňovať celý zásobník a zisťovať, či je prázdny.

## Záver

Aj keď sme boli oboznámení s rozsiahlosťou a náročnosťou projektu, jeho rozsah nás predsa prekvapil. Ako tím sme s programom podobnej náročnosti skúsenosti nemalo. Zadanie bolo náročné a opieralo sa o teoretické znalosti z predmetov IFJ a IAL – oba predmety nám však tieto znalosti v plnej miere poskytli.

Všetci sme sa zhodli na tom, že nás projekt aj napriek neľahkej implementácii a problémom, ktoré sme museli riešiť, veľa naučil a poskytol nám cenné skúsenosti do študijného aj profesionálneho života.

# Konečný automat



Operator = < ; >

Operator2 = >= ; <= ; !=

# LL-tabuľka

[illegible]



1.  $\text{START} \rightarrow \text{PROLOG PROG}$
2.  $\text{PROLOG} \rightarrow \text{package main}$
3.  $\text{PROG} \rightarrow \text{FUNCTION PROG}$
4.  $\text{PROG} \rightarrow \text{eps}$
5.  $\text{FUNCTION} \rightarrow \text{func id ( ARGS) RETURN\_F \{S\_BODY\}}$
6.  $\text{ARGS} \rightarrow \text{ARG ARG\_NEXT}$
7.  $\text{ARG\_NEXT} \rightarrow \text{eps}$
8.  $\text{ARG\_NEXT} \rightarrow , \text{ARG ARG\_NEXT}$
9.  $\text{ARG} \rightarrow \text{id DATA\_TYPE}$
10.  $\text{RETURN\_F} \rightarrow ( \text{RETURN\_TYPES})$
11.  $\text{RETURN\_F} \rightarrow \text{eps}$
12.  $\text{RETURN\_TYPES} \rightarrow \text{RETURN RETURN\_NEXT}$
13.  $\text{RETURN} \rightarrow \text{DATA\_TYPE}$
14.  $\text{RETURN\_NEXT} \rightarrow , \text{RETURN RETURN\_NEXT}$
15.  $\text{RETURN\_NEXT} \rightarrow \text{eps}$
16.  $\text{S\_BODY} \rightarrow \text{STATEMENT S\_BODY}$
17.  $\text{S\_BODY} \rightarrow \text{eps}$
18.  $\text{STATEMENT} \rightarrow \text{DEFINITION}$
19.  $\text{STATEMENT} \rightarrow \text{ASSIGN}$
20.  $\text{STATEMENT} \rightarrow \text{IF}$
21.  $\text{STATEMENT} \rightarrow \text{FOR}$
22.  $\text{STATEMENT} \rightarrow \text{F\_CALL}$
23.  $\text{STATEMENT} \rightarrow \text{RETURN\_S}$
24.  $\text{DEFINITION} \rightarrow \text{id} := \langle \text{expr} \rangle$
25.  $\text{ASSIGN} \rightarrow \text{id ID\_NEXT} = \text{FUNC\_OR\_EXPR}$
26.  $\text{FUNC\_OR\_EXPR} \rightarrow \langle \text{expr} \rangle \text{EXPR\_NEXT}$
27.  $\text{FUNC\_OR\_EXPR} \rightarrow \text{F\_CALL}$
28.  $\text{ID\_NEXT} \rightarrow , \text{id ID\_NEXT}$
29.  $\text{ID\_NEXT} \rightarrow \text{eps}$
30.  $\text{EXPR\_NEXT} \rightarrow , \langle \text{expr} \rangle \text{EXPR\_NEXT}$
31.  $\text{EXPR\_NEXT} \rightarrow \text{eps}$
32.  $\text{IF} \rightarrow \text{if} \langle \text{expr} \rangle \{ \text{S\_BODY} \} \text{ else } \{ \text{S\_BODY} \}$
33.  $\text{FOR} \rightarrow \text{for FOR\_DEFINITION} ; \langle \text{expr} \rangle ; \text{FOR\_ASSING} \{ \text{S\_BODY} \}$
34.  $\text{FOR\_DEFINITION} \rightarrow \text{DEFINITION}$
35.  $\text{FOR\_DEFINITION} \rightarrow \text{eps}$
36.  $\text{FOR\_ASSING} \rightarrow \text{F\_ASSIGN}$
37.  $\text{F\_ASSIGN} \rightarrow \text{id ID\_NEXT} = \langle \text{expr} \rangle \text{EXPR\_NEXT}$
38.  $\text{FOR\_ASSING} \rightarrow \text{eps}$
39.  $\text{F\_CALL} \rightarrow \text{id(PARAMS)}$
40.  $\text{PARAMS} \rightarrow \text{TERM TERM\_NEXT}$
41.  $\text{PARAMS} \rightarrow \text{eps}$
42.  $\text{TERM\_NEXT} \rightarrow , \text{TERM TERM\_NEXT}$
43.  $\text{TERM\_NEXT} \rightarrow \text{eps}$
44.  $\text{RETURN\_S} \rightarrow \text{return RET\_EXPR}$
45.  $\text{RET\_EXPR} \rightarrow \langle \text{expr} \rangle \text{EXPR\_NEXT}$
46.  $\text{RET\_EXPR} \rightarrow \text{eps}$
47.  $\text{TERM} \rightarrow \text{int\_value}$
48.  $\text{TERM} \rightarrow \text{float\_value}$
49.  $\text{TERM} \rightarrow \text{string\_value}$
50.  $\text{TERM} \rightarrow \text{id}$
51.  $\text{DATA\_TYPE} \rightarrow \text{int}$
52.  $\text{DATA\_TYPE} \rightarrow \text{float64}$
53.  $\text{DATA\_TYPE} \rightarrow \text{string}$

## Precedenčná tabuľka

[illegible]