

# Projektna dokumentacija

Tema projekta bila je implementacija jednostavnog kompajlera. Implementiran je korišćenjem flex-a i bison-a. Za potrebe projekta kao izvorni jezik korišćen je miniC (uproštena verzija C programskog jezika) a kao ciljni jezik RISC-V assembler. Implementacija miniC kompajlera je podeljena na uobičajene faze kompajliranja: leksičku, sintaksnu, semantičku analizu i generisanje koda.

## Leksička analiza

Leksički analizator (skener) je program namenjen za leksičku analizu (skeniranje) teksta. Zadatak skenera je da pročita ulazni tekst i da ga podeli na simbole. Skener radi tako što čita tekst, karakter po karakter, i pokušava da prepozna neku od zadatih reči. Ukoliko naiđe na simbol čija pojava nije predviđena, treba da prijavi leksičku grešku. Beline (u koje spadaju prazno mesto, tabulator i karakter za novi red) se preskaču, odnosno prepoznaju i ignorišu, jer imaju ulogu separatora simbola. Za izradu ovog projekta za generisanje skenera koristi se program flex. Flex generiše skener na osnovu pravila zadatih u flex specifikaciji. Specifikaciju kreira korisnik u posebnoj datoteci koja po konvenciji ima ekstenziju **.l**. Ova datoteka sadrži pravila koja opisuju reči koje skener treba da prepozna. Pravila se zadaju u obliku regularnih izraza. Svakom regularnom izrazu moguće je pridružiti akciju (u obliku C koda) koja će se izvršiti kada skener prepozna dati regularni izraz. Uobičajena akcija je vraćanje oznake simbola, odnosno tokena za taj simbol. Token je numerička oznaka klase (vrste) simbola. Na primer: token **NUMBER** označava bilo koji broj (i broj 2 i broj 359), dok token **IF** označava jedino ključnu reč **if**. Prilikom izvršavanja, skener će tražiti pojavu stringa u ulaznom tekstu koji odgovara nekom regularnom izrazu. Ako ga pronađe, izvršiće akciju (kod) koju je korisnik pridružio tom regularnom izrazu. Skener za miniC jezik je implementiran u datoteci **micko.l**. Prvi deo specifikacije započinje uključivanjem dve opcije koje utiču na to kako će flex napraviti skener:

`% option noyywrap yylineno`

**norrywrap** opcija znači da će skener skenirati samo jednu datoteku a **yylineno** opcija znači da će skener koristiti globalnu promenljivu flex-a **yylineno** koja sadrži broj trenutno skenirane linije. Za neke simbole se, osim tokena, prosleđuje još neka vrednost koja bliže opisuje simbol (npr. za simbol

broj se uz token **NUMBER** šalje i vrednost konkretnog broja, recimo 123). U ovakvim situacijama, za smeštanje vrednosti, koja može biti različitog tipa - za različite simbole, se koristi unija. Za miniC skener, potrebna je unija koja sadrži jedan celobrojni tip (**int i**) i jedan pokazivač na karakter (**char \*s**). Promenljiva tipa unije, preko koje se, uz token, prosleđuje vrednost, se u flex-u zove **yylval**. Niz pravila koja su definisana u skeneru:

```
[ \t\n]+          { /* skip */ }

"int"             { yynval.i = INT; return _TYPE; }
"unsigned"        { yynval.i = UINT; return _TYPE; }
"if"              { return _IF; }
"else"            { return _ELSE; }
"return"          { return _RETURN; }

"for"             { return _FOR; }
"++"              { return _INC; }
"while"           { return _WHILE; }
"?"              { return _QMARK; }
":"              { return _COLON; }

"("               { return _LPAREN; }
")"               { return _RPAREN; }
"{"               { return _LBRACKET; }
"}"               { return _RBRACKET; }
";"              { return _SEMICOLON; }
"="              { return _ASSIGN; }

"+"              { yynval.i = ADD; return _AROP; }
"-"              { yynval.i = SUB; return _AROP; }

"<"              { yynval.i = LT; return _RELOP; }
">"              { yynval.i = GT; return _RELOP; }
"<="             { yynval.i = LE; return _RELOP; }
">="             { yynval.i = GE; return _RELOP; }
"=="             { yynval.i = EQ; return _RELOP; }
"!="             { yynval.i = NE; return _RELOP; }

[a-zA-Z][a-zA-Z0-9]* { yynval.s = strdup(yytext);
                        return _ID; }
[+-]?[0-9]{1,10}     { yynval.s = strdup(yytext);
                        return _INT_NUMBER; }
[0-9]{1,10}[uU]      { yynval.s = strdup(yytext);
                        yynval.s[yyleng-1] = 0;
                        return _UINT_NUMBER; }

\\.\\.\\.          { /* skip */ }
.                  { printf("line %d: LEXICAL ERROR on char %c\n", yylineno, *yytext); }
```

Prvo pravilo preskače beline. Pravila koja se nalaze u crvenom kvadratu prepoznavaju ključne reči **if**, **else**, **return**, **for**, **while**, delimitere **(, )**, **{, }**, separatore **;**, operator dodele **=** i simbole ternarnog operatora **? i :** i za sve njih

je dovoljno vezati samo tokene(nisu potrebne dodatne vrednosti koje bi ih detaljnije denisale). Dva pravila iznad kvadrata prepoznaju označene i neoznačene celobrojne tipove. Za oba simbola se vezuje isti token (**\_TYPE**), ali različite vrednosti uz token. Vrednosti su konstante (**INT** i **UINT**) koje opisuju prepoznate tipove. Vrednost se prosleđuje kroz promenljivu **yyval**. Kako je konstanta celobrojnog tipa, njena vrednost se smešta u polje i promenljive **yyval**. Zatim slede pravila za prepoznavanje aritmetičkih i relacionih operatora. Za ove operatore se, osim tokena, prosleđuje još i odgovarajuća konstanta kao vrednost simbola. Konstanta je potrebna da bi se operatori razlikovali međusobno, jer imaju isti token (na primer: konstanta **ADD** za operator **+**). Zatim slede pravila za prepoznavanje identifikatora i literala. U ovim slučajevima se, pored tokena, prosleđuje još i string simbola, koji se smešta u polje **s** (tipa **char\***) promenljive **yyval**. String prepoznatog simbola flex čuva u promenljivoj **yytext**. Regularni izraz za identifikator opisuje string koji počinje slovom (malim ili velikim: **[a-zA-Z]**, a iza njega se može naći slovo (malo ili veliko) ili cifra, 0 ili više puta: **[a-zA-Z0-9]\***. Minimalni identifikator se, znači, sastoji od jednog slova. Regularni izraz za označeni celobrojni literal opisuje string od minimalno jedne do maksimalno deset cifara: **[0-9]{1,10}**, ispred kojih se može, a ne mora, naći predznak plus ili minus: **[+-]?**. Regularni izraz za neoznačeni celobrojni literal ne prepoznaje predznak, već dozvoljava da se iza broja navede malo ili veliko slovo **u**: **[uU]**. Linijski komentari se prepoznaju pomoću regularnog izraza koji prepoznaje dva znaka slash na početku stringa: **\\** (može i **//**), a zatim bilo koji znak osim **\n** znaka, 0 ili više puta: **.\***. Poslednje pravilo koje sadrži operator **.** (bilo koji znak osim **\n** znaka) služi da prepozna leksičku grešku. Bilo koji znak koji ne pripada nijednom prethodnom pravilu predstavlja grešku, koju treba saopštiti korisniku.

## Sintaksna analiza

Sintaksa jezika opisuje pravila po kojima se kombinuju simboli jezika (npr. u deklaraciji promenljive, prvo se piše ime tipa, zatim ime promenljive i na kraju ;). Sintaksa se opisuje gramatikom, obično pomoću BNF notacije. Sintaksna analiza ima zadatak da proveriti da li je ulazni tekst sintaksno ispravan. Ona to čini tako što preuzima niz tokena od skenera i proverava da li su tokeni navedeni u ispravnom redosledu. Ako jesu, to znači da je ulazni tekst napisan u skladu sa pravilima gramatike korišćenog jezika, tj. da je sintaksno ispravan. U suprotnom, sintaksna analiza treba da prijavi sintaksnu grešku i da nastavi analizu. Opisani proces se vrši u delu kompajlera koji se zove parser i naziva se parsiranje. Za implementaciju parsera, često se koristi generator parsera

**bison**. Prvi korak u generisanju parsera je priprema njegove specifikacije. Specifikaciju kreira korisnik u posebnoj datoteci koja po konvenciji ima ekstenziju **.y**. Ova datoteka sadrži gramatiku koju parser treba da prepozna. Pravila se zadaju u BNF obliku. Svakom pravilu je moguće pridružiti akciju (u obliku C koda) koja će se izvršiti kada parser prepozna dato pravilo. Prilikom izvršavanja, parser komunicira sa skenerom tako što od njega traži sledeći token iz ulaznog teksta. Kada dobije token, parser proverava da li je njegova pojava na datom mestu dozvoljena (tj. da li je u skladu sa gramatikom). Ako jeste nastaviće parsiranje, a ukoliko nije, prijavice grešku, pokušaće da se oporavi od greške i da nastavi parsiranje. U momentu kada prepozna celo pravilo, parser će izvršiti akciju koja je pridružena tom pravilu. Parser za miniC je implementiran je u datoteci **micko.y**. Sve konstante koje su potrebne skeneru i parseru su smeštene u posebnoj datoteci **defs.h** radi lakšeg održavanja koda.

```
//tipovi podataka
enum types { NO_TYPE, INT, UINT };

//vrste simbola (moze ih biti maksimalno 32)
enum kinds { NO_KIND = 0x1, REG = 0x2, LIT = 0x4,
             ... FUN = 0x8, VAR = 0x10, PAR = 0x20 };

//konstante aritmetickih operatora
enum arops { ADD, SUB, MUL, DIV, AROP_NUMBER };
//stringovi za generisanje aritmetickih naredbi
static char *ar_instructions[] = { "ADDS", "SUBS", "MULS", "DIVS",
                                   ... "ADDU", "SUBU", "MULU", "DIVU" };

//konstante relacionih operatora
enum relops { LT, GT, LE, GE, EQ, NE, RELOP_NUMBER };
//stringovi za JMP naredbu
static char* jumps[]={"JLTS", "JGTS", "JLES", "JGES", "JEQ ", "JNE ",
                     ... "JLTU", "JGTU", "JLEU", "JGEU", "JEQ ", "JNE " };

static char* opp_jumps[]={"JGES", "JLES", "JGTS", "JLTS", "JNE ", "JEQ ",
                        ... "JGEU", "JLEU", "JGTU", "JLTU", "JNE ", "JEQ "};
```

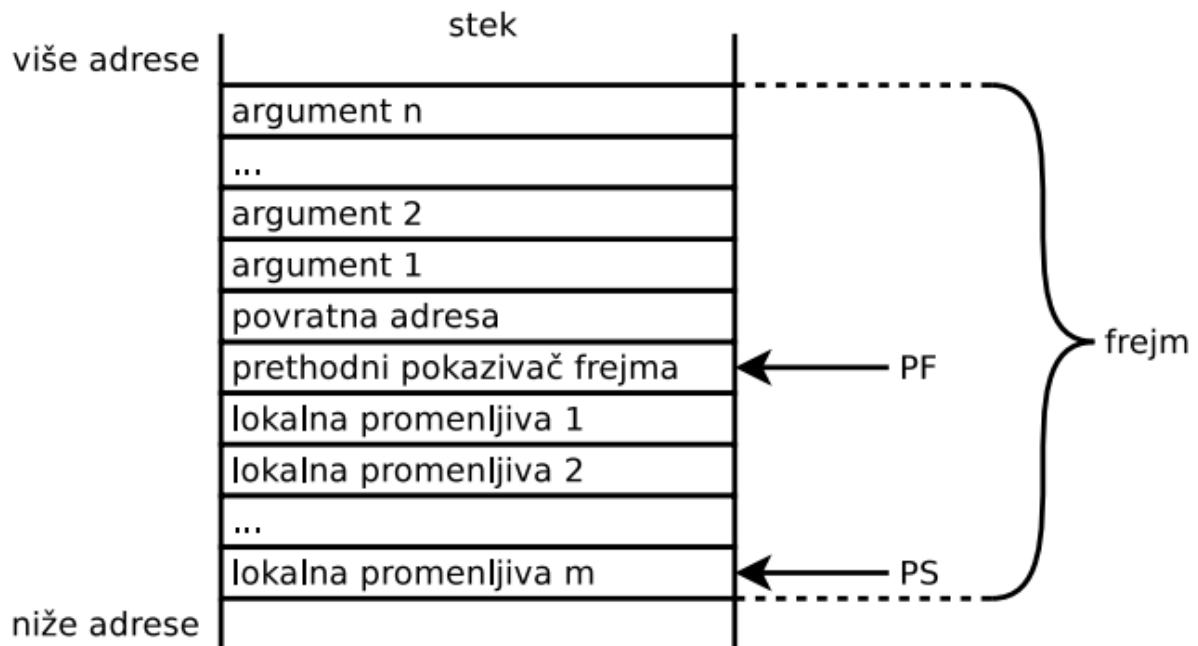
Enumeracije uokvirene crvenom bojom su bitne za implementaciju parsera. Ostale enumeracije su potrebne za generisanje asemblerskog koda osim enumeracije **kinds** koja se koristi u semnatičke svrhe. Enumeracija **types** sadrži konstante koje opisuju tipove. miniC jezik podržava samo **int** i **unsigned** tipove. Enumeracija **arops** sadrži konstante koje opisuju aritmetičke operatore sabiranja, oduzimanja, množenja i deljenja, a enumeracija **relops** sadrži konstante koje opisuju relacione operatore (<, >, <=, >=, == i !=). Tokeni koje šalje skener, a koje prima parser, se moraju denisati pomoću ključne reči **%token**. Uz svaki token se može denisati tip

vrednosti koja se prosleđuje sa njim, između operatora **<i>**. Tokeni koji imaju vrednosti celobrojnog tipa su **\_TYPE**, **\_AROP** i **\_RELOP** (njihova vrednost je konstanta koja opisuje vrstu tipa, aritmetičkog, odnosno relacionog operatora), a tokeni koji imaju vrednost tipa string su **\_ID**, **\_INT\_NUMBER** i **\_UINT\_NUMBER** (njihova vrednost je string imena ili broja). Funkcija `main()` jedino poziva parser (tj. funkciju **`yyparse()`**). Čim naiđe na prvu sintaksnu grešku, parser će završiti parsiranje.

## Semantička analiza

Semantička analiza je faza kompajliranja u kojoj se proverava značenje (semantika) programskog teksta (koda). Kod koji je sintaksno ispravan, ne mora biti i semantički ispravan. Na primer, ako se promenljiva tipa `int` poredi sa promenljivom tipa `boolean`, kompajler treba da prijavi semantičku grešku. Semantička analiza se, najčešće, implementira u parseru, tako što se pravilima dodaju semantičke provere. Za ovo je potrebno poznavati organizaciju memorije i arhitekturu računara za koju se generiše kod. Globalni identifikatori su statični postoje za sve vreme izvršavanja programa i za njih se mogu rezervisati memorijske lokacije u toku kompajliranja. Lokalni identifikatori su dinamični postoje samo za vreme izvršavanja funkcija i za njih se zauzimaju memorijske lokacije na početku izvršavanja funkcija. Ove lokacije se oslobađaju na kraju izvršavanja funkcija, pa se zato lokalnim identifikatorima dodeljuju memorijske lokacije sa steka. Deo steka koji se zauzima za izvršavanje neke funkcije se zove (stek) frejm. Stek frejm, redom, sadrži:

- argumente poziva funkcije (od poslednjeg ka prvom),
- povratnu adresu (od koje će se nastaviti izvršavanje koda nakon poziva funkcije),
- pokazivač prethodnog stek frejma (frejma funkcije iz koje se poziva ova funkcija) i
- lokalne promenljive (od prve ka poslednjoj).



Pokazivač frejma **PF** pokazuje na početak frejma, a pokazivač **PS** pokazuje na vrh steka. **PF** ostaje isti, dok se **PS** menja tokom izvršavanja funkcije. Pošto će generator koda generisati kod na RISC-V asemblerskom jeziku, potrebno je upoznati se sa arhitekturom naredbi i registara. Ukupno ima 32 registara (za potrebe projekta korisceno je 30). Oznaka registra se sastoji od oznake x i rednog broja registra: x0, x1, ..., x31. Registri od x5 do x9 i od x18 do x31 imaju opštu namenu i služe kao radni registri. Registri x10 i x11 su rezervisani za smeštanje argumenata koji će funkcija da primi i za povratnu vrednost funkcije, od x12 do x17 su registri za smeštanje argumenta koji ce funkcija da primi, registar x2 služi kao pokazivač steka (**PS**), a registar x0 sadrži u sebi vrednost 0. Tabela simbola (**symbol table**) je struktura podataka koja čuva identikatore iz kompajliranog programa i sve informacije o tim identikatorima. Za semantičku analizu, ona je neophodna, jer sadrži sve informacije na osnovu kojih se mogu sprovesti semantičke provere. Na primer, u sledećem programu, postoji samo main() funkcija, sa jednim parametrom i jednom lokalnom promenljivom:

```
int main ( int p )
{
    int x ;
    x = 100;
    return p + x ;
}
```

Tokom kompajliranja ove funkcije, u tabelu simbola treba zapisati da postoji:

- identikator **main** koji predstavlja funkciju, čiji povratni tip je **int** i koja ima jedan parametar tipa **int**
- identikator **p** tipa **int**, koji predstavlja parametar, prvi po redu
- identikator **x** tipa **int**, koji predstavlja lokalnu promenljivu, prvu po redu
- literal **100** tipa **int**.

Za ovaj program, sadržaj tabele simbola (na kraju parsiranja funkcije `main()`) bi bio:

| STRING<br>SIMBOLA | VRSTA<br>SIMBOLA | TIP<br>SIMBOLA | ATRIBUT<br>SIMBOLA | ATRIBUT<br>SIMBOLA |
|-------------------|------------------|----------------|--------------------|--------------------|
| main              | FUN              | INT_TYPE       | 1                  | INT_TYPE           |
| p                 | PAR              | INT_TYPE       | 1                  | -                  |
| x                 | VAR              | INT_TYPE       | 1                  | -                  |
| 100               | LIT              | INT_TYPE       | -                  | -                  |

I globalni i lokalni identikatori mogu biti smešteni u istu tabelu simbola. Globalni su prisutni u tabeli simbola sve vreme kompajliranja, dok su lokalni prisutni samo u toku kompajliranja funkcije kojoj pripadaju. Izlaskom iz opsega vidljivosti, iz tabele simbola se brišu identikatori vidljivi (samo) u tom opsegu. Na primer, na kraju parsiranja funkcije, brišu se svi lokalni identikatori funkcije. Pošto u tabeli simbola istovremeno mogu postojati identični globalni i lokalni identikatori, radi njihovog razlikovanja, u tabeli mora biti naznačena vrsta identikatora (da li je funkcija, promenljiva, parametar, argument, ...). Za svaki identikator mora postojati i oznaka njegovog tipa (da li je u pitanju tip `int`, `unsigned`, `char`, ...). Za parametre i lokalne promenljive mora postojati njihova relativna pozicija na steku u odnosu na početak frejma, tj. njihov redni broj (ova informacija je potrebna za generisanje koda). Takođe, radi uniformnosti, u tabeli simbola mogu biti smeštene i oznake radnih registara. One se smeštaju u tabelu u vreme njene inicijalizacije. Registri od `x5` do `x7` i od `x28` do `x31` se smeštaju u element koji počinje sa indeksom 0 a završavaju se sa indeksom 6, a registri `x8`, `x9` i od `x18` do `x25` u element koji počinje sa indeksom 7 a završava se sa indeksom 16. Registri su prisutni u tabeli simbola sve vreme kompajliranja. Registri se koriste u fazi generisanja koda,

kada se u njih smeštaju međurezultati izraza. Koristi se vrlo jednostavna implementacija tabele simbola: tabela je organizovana kao niz struktura koje opisuju osobine identifikatora. U datoteci **defs.h** nalazi se enumeracija **kinds** koja sadrži konstante koje opisuju vrste simbola. Vrednosti konstanti su stepeni broja 2 (1, 2, 4, 8, ... tj. sadrže samo po 1 setovan bit) da bi se nad njima mogla primeniti bitwise logika. U ovu datoteku smeštena su i dva makroa koja olakšavaju prijavu greške ili upozorenja: **err(args ...)** i **warn(args ...)**, kao i nekoliko konstanti konstante za simulaciju bool tipa (**TRUE**, **FALSE**), dužina tabele simbola, tj. broj elemenata tabele: **SYMBOL\_TABLE\_LENGTH**, dužina bafera za smeštanje poruka o greškama: **CHAR\_BUFFER\_LENGTH**, oznaka da simbol nema atribut **NO\_ATR**, redni broj poslednjeg radnog registra **LAST\_WORKING\_REG**, redni broj registra koji čuva povratnu vrednost funkcije **FUN\_REG**, pomoćni registri **TMP\_REG** i **HELP\_REG**.

```
# define bool int
# define TRUE 1
# define FALSE 0
# define SYMBOL_TABLE_LENGTH 64
# define NO_ATR 0
# define LAST_WORKING_REG 16
# define FUN_REG 21
# define TMP_REG 17
# define HELP_REG 18
# define CHAR_BUFFER_LENGTH 128
```

Struktura jednog elementa tabele simbola se sastoji od polja(datoteka **syntab.h**):

- **name** - string identifikatora
- **kind** - vrsta simbola - sadrži vrednost iz enumeracije **kinds**
- **type** - tip simbola - sadrži vrednost iz enumeracije **types**
- **atr1** - atribut - sadrži različite vrednosti za različite vrste simbola
  - za lokalnu promenljivu, u ovom se smešta redni broj promenljive
  - za parametar, u ovom polju se smešta redni broj parametra
  - za funkciju, u ovom polju se čuva broj parametara
  - za ostale identifikatore, ovo polje nije popunjeno.
- **atr2** - atribut - sadrži različite vrednosti za različite vrste simbola
  - za funkciju, u ovom polju se čuva tip parametra (vrednost iz enumeracije **types**)



```
// Element tabele simbola
typedef struct sym_entry {
    char *   name;           // ime simbola
    unsigned kind;          // vrsta simbola
    unsigned type;          // tip vrednosti simbola
    unsigned atr1;          // dodatni atribut simbola
    unsigned atr2;          // dodatni atribut simbola
} SYMBOL_ENTRY;
```

Operacije za rad sa tabelom simbola su implementirane u datotekama **symtab.h** i **symtab.c**. Funkcija koja služi za ubacivanje jednog simbola u tabelu simbola se zove **insert\_symbol()**. Parametri funkcije su string, vrsta, tip i atributi simbola. Funkcija vraća indeks ubačenog elementa u tabeli simbola.

```
int insert_symbol ( char * name , unsigned kind , unsigned type ,
unsigned atr1 , unsigned atr2 ) ;
```

Funkcija koja služi za ubacivanje novog literala u tabelu simbola se zove **insert\_literal()**. Parametri funkcije su string i tip literala. Funkcija će ubaciti novi literal ako on već ne postoji u tabeli simbola. U suprotnom neće prijavljivati grešku, jer je dovoljna jedna instanca literala u tabeli (svi literali se tretiraju kao lokalni za funkciju).

```
int insert_literal ( char * str , unsigned type ) ;
```

Funkcija koja pretražuje tabelu simbola se zove **lookup\_symbol()**. Parametri funkcije su ime i vrsta simbola koji se traži. Funkcija vraća indeks elementa tabele simbola na kom je pronašla simbol, ili -1 ukoliko ga nije našla. Pretraga tabele simbola je potrebna u semantičkim proverama, kada se npr. proverava da li je promenljiva u nekom izrazu prethodno deklarirana ili ne.

```
int lookup_symbol ( char * name , unsigned kind ) ;
```

Set i get metode kojima se menjaju i čitaju vrednosti određenih polja elementa tabele simbola su:

```
void set_name ( int index , char * name ) ;
char * get_name ( int index ) ;
void set_kind ( int index , unsigned kind ) ;
unsigned get_kind ( int index ) ;
void set_type ( int index , unsigned type ) ;
unsigned get_type ( int index ) ;
```

```
void set_atr1 ( int index , unsigned atr1 ) ;  
unsigned get_atr1 ( int index ) ;  
void set_atr2 ( int index , unsigned atr2 ) ;  
unsigned get_atr2 ( int index ) ;
```

Funkcija koja briše deo tabele simbola i koja se koristi za brisanje lokalnih simbola funkcije se zove **clear\_symbols()**. Parametar ove funkcije je indeks prvog elementa koji treba da se obriše. Funkcija će obrisati sve simbole koji se nalaze između ovog indeksa i poslednjeg popunjenog elementa tabele simbola.

```
void clear_symbols ( unsigned begin_index );
```

Funkcije za rad sa tabelom simbola u celini su funkcija za prikaz svih popunjenih elemenata tabele: **print\_symtab()**, funkcija za inicijalizaciju tabele simbola: **init\_symtab()** i funkcija za brisanje svih elemenata tabele simbola: **clear\_symtab()**.

```
void print_symtab ( void ) ;  
void init_symtab ( void ) ;  
void clear_symtab ( void ) ;
```

miniC parser sa semantičkim proverama koristi prethodno opisanu tabelu simbola a implemantiran je u datoteci **micko.y**. Prvi deo specifikacije parsera je proširen denicijama promenljivih koje koristi parser u toku provere semantike. To su broj grašaka **error\_count**, broj upozorenja **warning\_count**, broj lokalnih promenljivih **var\_num**, indeks u tabeli simbola na kom se nalazi simbol trenutno parsirane funkcije **fun\_idx** i indeks u tabeli simbola na kom se nalazi ime funkcije čiji poziv se trenutno parsira **fcall\_idx**:

```
int error_count = 0;  
int warning_count = 0;  
int var_num = 0;  
int fun_idx = -1;  
int fcall_idx = -1;
```

U nastavku su denisani i tipovi pojmova: iza ključne reči **%type** se navodi tip iz unije ( **<i>** ili **<s>**), a zatim imena svih pojmova koji će biti tog tipa. Ako je tip pojma **<i>** to znači da će njegova vrednost biti celobrojnog tipa (**int**), a ako je tipa **<s>** onda će njegova vrednost biti tipa string (**char \***).

**%type <i> num\_exp exp literal**

**%type <i> function\_call argument rel\_exp if\_part cond\_exp**

Dva tokena: **ONLY\_IF** i **\_ELSE** su denisana upotrebom deklaracije

**%nonassoc** koja kaže da ovi tokeni nemaju asocijativnost. Ova dva tokena se koriste za potrebe razrešavanja dvosmislenosti if-else iskaza.

**%nonassoc ONLY\_IF**

**%nonassoc \_ELSE**

Pravila u drugom delu specikacije parsera su proširena akcijama koje sadrže semantičke provere. Semantička provera, koju treba obaviti nakon prepoznavanja celog programa, je da li postoji **main()** funkcija. Provera se vrši pretragom tabele simbola. Ako se u tabeli ne pronađe simbol main koji predstavlja funkciju (**FUN**), prijavi se semantička greška.

```
program
: function_list
{
    if(lookup_symbol("main", FUN) == NO_INDEX)
        err("undefined reference to 'main'");
}
;
```

Tokom parsiranja denicije funkcije vrši se nekoliko semantičkih akcija. Čim pristigne ime funkcije, ono se smešta u tabelu simbola zajedno sa informacijama o tipu povratne vrednosti funkcije (vrednost tokena **\_TYPE** koja nosi konstantu **INT** ili **UINT**) i vrsti identikatora (konstanta **FUN**). Promenljiva **fun\_idx** dobija indeks elementa u tabeli simbola u koji je smešteno ime funkcije. Nakon parsiranja celog tela funkcije (nakon prepoznavanja pojma body) iz tabele simbola se brišu svi simboli koji su lokalni za funkciju. Ovo sme da se radi (i potrebno je), jer je opseg vidljivosti lokalnih simbola samo do kraja funkcije u kojoj su deklarisan. Ovim je omogućeno da svaka funkcija ima svoje lokalne promenljive, koje mogu biti (potpuno) iste kao u nekoj drugoj funkciji. Funkcija **clear\_symbols()**, koja briše lokalne simbole, se poziva sa argumentom **fun\_idx+1**, jer brisanje tabele simbola treba početi nakon imena funkcije. Ime funkcije je globalni identikator i zato treba da bude prisutan u tabeli simbola tokom parsiranja celog programa. Promenljiva **var\_num**, koja broji lokalne promenljive jedne funkcije, se resetuje.

```

function
: _TYPE_ID
{
    fun_idx = lookup_symbol($2, FUN);
    if(fun_idx == NO_INDEX)
        fun_idx = insert_symbol($2, FUN, $1, NO_ATR, NO_ATR);
    else
        err("redefinition of function '%s'", $2);
}
_LPAREN parameter _RPAREN body
{
    clear_symbols(fun_idx + 1);
    var_num = 0;
}
;

```

Nakon parsiranja pojma parameter, ime parametra (vrednost **\$2**) se ubacuje u tabelu simbola kao **PAR** sa tipom koji nosi token **\_TYPE** i sa rednim brojem 1 (polje **atr1** parametra čuva redni broj parametra). Tip parametra se zapisuje u polju **atr2** simbola funkcije u tabeli simbola (jer polje **atr2** za funkciju čuva tip njenog parametra). Ako nije bilo parametara u polje **atr1** smešta se podatak o broju parametara funkcije koji je u ovom slučaju 0, a ukoliko je detektovan jedan parametar u polje **atr1** smešta se 1.

```

parameter
: /* empty */
{ set_atr1(fun_idx, 0); }

| _TYPE_ID
{
    insert_symbol($2, PAR, $1, 1, NO_ATR);
    set_atr1(fun_idx, 1);
    set_atr2(fun_idx, $1);
}
;

```

Tokom parsiranja deklaracije promenljive (pojam **variable**), ime promenljive (vrednost **\$2**) se ubacuje u tabelu simbola kao **VAR** sa tipom koji nosi token **\_TYPE**. Vrednost brojača lokalnih promenljivih se inkrementira i zapisuje se kao atribut **atr1** ovog simbola u tabeli simbola, zato što lokalna promenljiva na mestu prvog atributa u tabeli simbola čuva svoj redni broj. Ukoliko je u tabeli simbola pronađen parametar sa istim imenom, prijavljuje se semantička greška zato što dva lokalna identikatora unutar funkcije ne mogu imati isto ime.

```

variable
: _TYPE _ID _SEMICOLON
{
    if(lookup_symbol($2, VAR|PAR) == NO_INDEX)
        insert_symbol($2, VAR, $1, ++var_num, NO_ATR);
    else
        err("redefinition of '%s'", $2);
}
;

```

Nakon parsiranja iskaza dodele(pojam **assignment\_statement**), proverava se da li je identikator koji se nalazi sa leve strane jednakosti promenljiva ili parametar. Ako nije, prijavljuje se semantička greška, jer je na tom mestu dozvoljeno samo ime promenljive ili parametra (a ne, recimo, ime funkcije). Dodatno, vrši se i provera tipova leve i desne strane jednakosti, jer bi oni trebali biti isti. Za preuzimanje tipova koristi se funkcija **get\_type()**. Njoj treba predati indeks u tabeli simbola onog simbola čiji tip treba da preuzme (u iskazu dodele to su: indeks pronađenog identikatora sa leve strane jednakosti i indeks u tabeli simbola gde se nalazi rezultat numeričkog izraza sa desne strane jednakosti (vrednost **\$3**)).

```

assignment_statement
: _ID _ASSIGN num_exp _SEMICOLON
{
    int idx = lookup_symbol($1, VAR|PAR);
    if(idx == NO_INDEX)
        err("invalid lvalue '%s' in assignment", $1);
    else
        if(get_type(idx) != get_type($3))
            err("incompatible types in assignment");
}
;

```

Nakon parsiranja prvog dela for petlje(pojam **for\_statement**), proverava se da li je identikator promenljiva ili parametar. Ako nije, prijavljuje se semantička greška, jer je na tom mestu dozvoljeno samo ime promenljive ili parametra. Dodatno, vrši se i provera tipova leve i desne strane jednakosti, jer bi oni trebali biti isti. Preuzimanje tipova vrši se na isti način kao i kod **assignment\_statement**. U poslednjem delu for petlje, ponovo se proverava da li je identikator promenljiva ili parametar i ako nije prijavljuje se semantička greška.

```

for_statement
: _FOR_LPAREN_ID_ASSIGN literal
{
    int i = lookup_symbol($3, VAR|PAR);
    if(i == NO_INDEX)
        err("'s' undeclared", $3);

    if(get_type(i) != get_type($5))
        err("incompatible types");
}
_SEMICOLON rel_exp
_SEMICOLON_ID_INC_RPAREN statement
{
    int i = lookup_symbol($11, VAR|PAR);
    if(i == NO_INDEX)
        err("'s' undeclared", $11);
}
;

```

Nakon parsiranja aritmetičkog izraza (pojam **num\_exp**), vrši se provera tipova operanada, pa ako nisu isti, prijavljuje se semantička greška. Provera se vrši pozivom funkcije **get\_type()** sa argumentima: indeks prvog operanda u tabeli simbola (**\$1**) i indeks drugog operanda u tabeli simbola (**\$3**).

```

num_exp
: exp
| num_exp _AROP exp
{
    if(get_type($1) != get_type($3))
        err("invalid operands: arithmetic operation");
}
;

```

Nakon parsiranja identifikatora u izrazima (**exp**), vrši se provera postojanja dotičnog imena, tj. njegovog prisustva u tabeli simbola. Ovu proveru obavlja funkcija **lookup\_symbol()**. Ako ga nema u tabeli, znači da nikada nije deklarisan, pa treba prijaviti semantičku gršku, jer se ne može koristiti promenljiva koja nije prethodno deklarirana. Ako se izraz nalazi u zagradama (pretposlednja varijanta pravila **exp**), njegova vrednost treba da postane i vrednost celog izraza. U slučaju proslednjeg pravila **exp** tu se proveravaju tipovi operanada, pa ako nisu isti, prijavljuje se semantička greška. Provera se vrši isto kao i kod **num\_exp**.

```

exp
: literal
| _ID
{
    $$ = lookup_symbol($1, VAR|PAR);
    if($$ == NO_INDEX)
        err("'$$' undeclared", $1);
}
| function_call
| _LPAREN num_exp _RPAREN
{ $$ = $2; }
| _LPAREN rel_exp _RPAREN _QMARK cond_exp _COLON cond_exp
{
    if(get_type($5) != get_type($7))
        err("expressions don't have same types");
}
;

```

Nakon parsiranja identifikatora u izrazima (**cond\_exp**), vrši se provera postojanja dotičnog imena, tj. njegovog prisustva u tabeli simbola. Ova provera obavlja se isto kao i kod **exp**.

```

cond_exp
: _ID
{
    if(($$ = lookup_symbol($1, VAR|PAR)) == NO_INDEX)
        err("'$$' undeclared", $1);
}
| literal
;

```

Kada se isparsira literal, treba ga ubaciti u tabelu simbola kao lokalni simbol, a pojmu **literal** dodeliti vrednost: indeks u tabeli simbola na kom je smeštena konstanta.

```

literal
: _INT_NUMBER
{ $$ = insert_literal($1, INT); }
| _UINT_NUMBER
{ $$ = insert_literal($1, UINT); }
;

```

Tokom parsiranja poziva funkcije (pojam **function\_call**) proverava se da li je navedeni identifikator ime neke od postojećih funkcija, pa ako nije, prijavljuje se greška da se poziva nepostojeća funkcija. Kako je usvojeno, povratna vrednost funkcije će biti smeštena u registru **%13**, pa je potrebno postaviti vrednost pojma **function\_call** na indeks 13 u tabeli simbola (gde se nalazi

registar %13). Registru %13 se postavlja tip povratne vrednosti funkcije koja se kasnije koristi u izrazima.

```
function_call
: _ID
{
    fcall_idx = lookup_symbol($1, FUN);
    if(fcall_idx == NO_INDEX)
        err("'<strong>%s' is not a function", $1);
}
_LPAREN argument _RPAREN
{
    if(get_atrl(fcall_idx) != $4)
        err("wrong number of arguments");
    set_type(FUN_REG, get_type(fcall_idx));
    $$ = FUN_REG;
}
;
```

Tokom parsiranja argumenta, treba proveriti da li je njegov tip isti kao tip odgovarajućeg parametra. Ako nije, treba prijaviti grešku. Tip parametra se dobija pozivom funkcije **get\_atr2(fcall\_idx)**, dok se tip argumenta dobija pozivom funkcije **get\_type(\$1)** (\$1 je indeks elementa u tabeli simbola koji sadrži rezultat numeričkog izraza koji je prosleđen na mestu argumenta). Vrednost pojma **argument** dobija vrednost 1 koja označava da je to prvi po redu (i jedini) argument. Ukoliko argumenta nema, vrednost ovog pojma postaje 0 (oznaka da nema argumenata).

```
argument
: /* empty */
{ $$ = 0; }

| num_exp
{
    if(get_atr2(fcall_idx) != get_type($1))
        err("incompatible type for argument");
    $$ = 1;
}
;
```

Tokom parsiranja relacionog izraza odnosno pojma **rel\_exp** potrebno je proveriti tipove operandada. Semantika jezika kaže da oni moraju biti isti. To se postiže pozivima funkcije **get\_type()** sa argumentima poziva: indeks prvog operanda u tabeli simbola (\$1) i indeks drugog operanda u tabeli simbola (\$3).



```

rel_exp
: num_exp _RELOP num_exp
{
    if(get_type($1) != get_type($3))
        err("invalid operands: relational operator");
}
;

```

Nakon prepoznavanja **return** iskaza, potrebno je proveriti tip izraza koji se vraća iz funkcije: semantika nalaže da ovaj tip mora biti isti kao povratni tip funkcije. Funkcija koja preuzima tipove **get\_type()** se poziva sa argumentima: indeks u tabeli simbola na kom se nalazi ime funkcije (**fun\_idx**, na kom stoji informacija o povratnom tipu funkcije) i indeks na kom se nalazi rezultat izraza koji se vraća iz funkcije.

```

return_statement
: _RETURN num_exp _SEMICOLON
{
    if(get_type(fun_idx) != get_type($2))
        err("incompatible types in return");
}
;

```

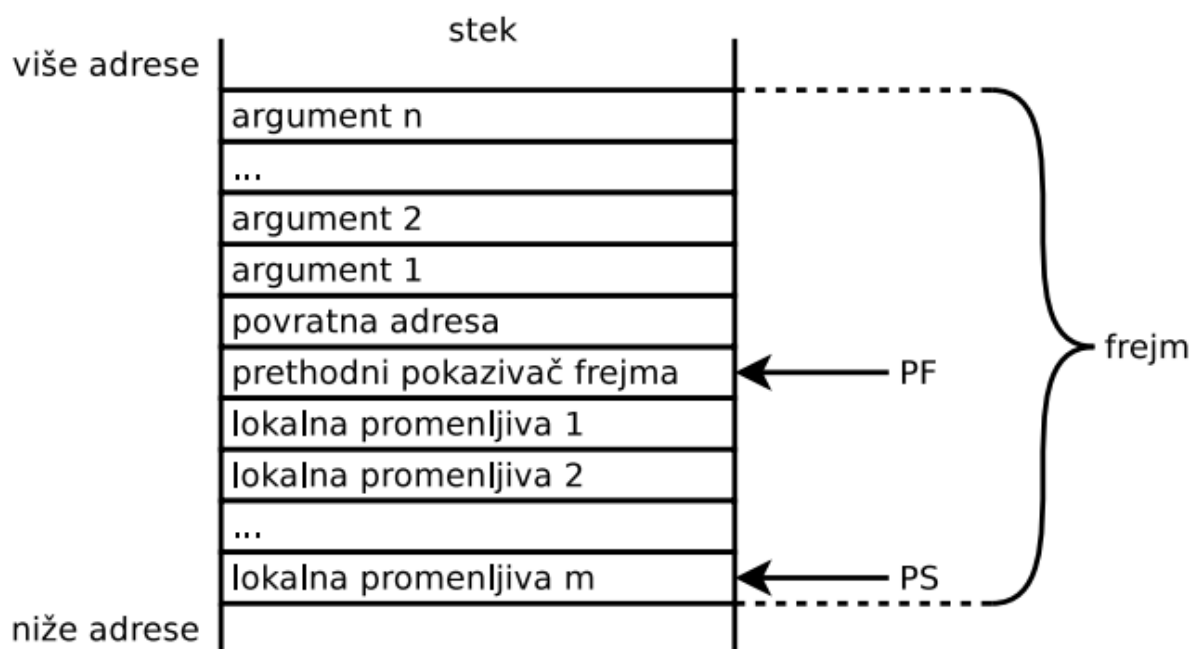
Funkcije **yyerror()** i **warning()** služe za prijavljivanje greške ili upozorenja korisniku.

Funkcija **main** vrši inicijalizaciju tabele simbola (**init\_symtab()**) pre početka parsiranja i brisanje tabele simbola (**clear\_symtab()**) nakon parsiranja. Na kraju, prijavljuje broj grešaka i upozorenja u toku parsiranja.

## Generisanje koda

Generisanje koda je faza kompajliranja u kojoj se proizvodi datoteka sa ekvivalentnim programom napisanim na ciljnom programskom jeziku. Ciljni jezik je jezik niskog nivoa, i može biti neki mašinski ili asemblerski jezik. Ciljni jezik na koji će u projektu biti preveden miniC programski jezik, je RISC-V asemblerski jezik. Podrazumeva se da registri i memorijske lokacije zauzimaju po 4 bajta. Ukupno ima 32 registara. Oznaka registra se sastoji od oznake x i rednog broja registra: x0, x1, ..., x31. Labele započinju malim slovom iza koga mogu da slede mala slova, cifre i podcrta `\_` (alfabet je 7 bitni ASCII). Iza labele se navodi dvotačka. Potrebno je imati šeme ekvivalencije. To su šabloni koji opisuju kako se određeni iskazi miniC jezika pišu na RISC-V asemblerskom jeziku. Ovi šabloni ne smeju izmeniti značenje polaznog miniC

iskaza. Primeri (mogućih) ekvivalencija između miniC koda i RISC-V asemblerskog koda, napisani su u dve kolone. U prvoj koloni se nalazi miniC kod koji se prevodi, a u drugoj se nalazi ekvivalentni asemblerski kod. Primeri sadrže i predloge načina generisanja ekvivalentnog koda. U svim navedenim primerima se podrazumeva da su **a**, **b** i **c** celobrojne označene lokalne promenljive, deklarisanе tim redom. Lokalne promenljive se smeštaju ispod pokazivača steka **sp**. Za pristup ovim lokacijama koristi se indeksno adresiranje pomoću registra **sp**. Tako je lokacija prve promenljive: **4(sp)** a druge **0(sp)** (memorijsku lokaciju adresira zbir vrednosti broja i sadržaja registra). Lokacije argumenata rade se identično kao i lokalne promenljive, pa se prvom argumentu pristupa na adresi **4(sp)**, drugom na adresi **0(sp)**, itd.



Primeru iskaza pridruživanja odgovara izgenerisani kod:

|                   |                       |
|-------------------|-----------------------|
| <b>a = b - a;</b> | <b>sub t2, t1, t0</b> |
|                   | <b>mv t0, t2</b>      |

Za smeštanje međurezultata izraza koriste se radni registri. Podrazumeva se da su svi radni registri na početku slobodni. U prvoj naredbi **sub** se zauzima (prvi slobodan) radni registar **t2** (jer su **t0** i **t1** registri zauzele promenljive **a** i **b**) i u njega se smešta razlika **b - a** (prvi operand naredbe **sub** je lokacija lokalne promenljive **b - t1**, drugi operand lokacija lokalne promenljive **a - t0**, a izlazni operand je registar **t2**). U naredbi **mv** se vrednost iz radnog registra **t2** kopira u lokaciju lokalne promenljive **a - t0**, i oslobađa se registar **t2**.

**a = ( a - b ) + ( a - c ) - c;**

```
sub    t3, t0, t1
sub    t4, t0, t2
add    t3, t3, t4
sub    t3, t3, t2
mv     t0, t3
```

U prvoj naredbi **sub** se zauzima radni registar **t3**, a u drugoj registar **t4**. U naredbi **add**(linija 3) se oslobađaju radni registri **t3** i **t4**, a ponovo se zauzima radni registar **t3**. U naredbi **sub** (linija 4) se oslobađa i ponovo zauzima radni registar **t3**, a u naredbi **mv** se oslobađa radni registar **t3**. Radni registar se zauzima za smeštanje rezultata svakog aritmetičkog izraza (**num\_exp**). Radni registar se oslobađa čim se preuzme njegova vrednost. Pošto se međurezultati izraza koriste u suprotnom redosledu od onog u kome su izračunati, radni registri, koji se koriste za smeštanje međurezultata, se zauzimaju i oslobađaju po principu steka. Kao "pokazivač steka registara" koristi se promenljiva **free\_reg\_num**, koja sadrži broj prvog slobodnog radnog registra. Zauzimanje radnog registra se sastoji od preuzimanja vrednosti promenljive **free\_reg\_num** i njenog inkrementiranja, a oslobađanje radnog registra se sastoji od dekrementiranja ove promenljive. Treba napomenuti da je broj registra istovremeno i indeks elementa tabele simbola. Broj zauzetog radnog registra služi kao vrednost sintaksnog pojma (**num\_exp**) koji odgovara izrazu čiji rezultat radni registar sadrži. Prekoračenje broja radnih registara (**free\_reg\_num > 16**) predstavlja fatalnu grešku u radu kompajlera. Za smeštanje povratne vrednosti funkcije se koristi registar **a0**. Funkcija, nakon poziva, lokalne promenljive i argumente čuva na stek frejmu. Primer ekvivalencije za deniciju funkcije:

```
int f ( int p ) {
    int a;
    return p + a;
}
```

```
f:
    addi    sp, sp, -4
    sw      t0, 0(sp)
f_body:
    add     t1, a0, t0
    mv      a0, t1
    j       f_exit
f_exit:
    lw      t0, 0(sp)
    addi    sp, sp, 4
    ret
```

Ekvivalentni asemblerski kod za funkciju započinje labelom koja ima isti identifikator kao funkcija, u ovom slučaju **f**:. Naredbom **addi**(linija 2) se zauzima prostor na steku za lokalnu promenljivu **a** (pokazivač steka se pomera jednu lokaciju, tj. 4 bajta niže), dok se pomoću naredbe **sw** postavlja vrednost u određenu lokaciju steka. Promenljiva **var\_num** služi kao brojač lokalnih promenljivih. Veličina prostora za lokalne promenljive se određuje kao **var\_num \* 4**. Prostor se zauzima samo ako funkcija ima lokalnih promenljivih (**var\_num > 0**). Ekvivalentni asemblerski kod tela funkcije započinje labelom koja sadrži ime funkcije sa postfiksom **\_body**, u ovom slučaju **f\_body**:. U naredbi **add** se računa povratna vrednost funkcije, a naredbom **mv** se ta vrednost smešta u registar **a0**. Ako funkcija ne sadrži return iskaz, kao povratna vrednost funkcije služi zatečeni sadržaj registra **a0**, koji je nepoznat u vreme defnisanja funkcije. Ekvivalentni asemblerski kod kraja funkcije započinje labelom koja sadrži ime funkcije sa postfiksom **\_exit**, u ovom slučaju **f\_exit**:. Naredbom **lw** se učitava vrednost iz memorije, a u naredbi **addi** se oslobađa prostor za lokalne promenljive (pokazivač steka se vraća u poziciju u kojoj je bio pre zauzimanja lokalnih promenljivih). Naredbom **ret** će doći do preusmeravanja toka izvršavanja na mesto odakle je funkcija pozvana. Primeru poziva funkcije (u okviru iskaza dodele) odgovara seldeći kod:

|                         |            |                   |
|-------------------------|------------|-------------------|
| <b>a = f ( a + b );</b> | <b>add</b> | <b>t2, t0, t1</b> |
|                         | <b>mv</b>  | <b>a0, t2</b>     |
|                         | <b>jal</b> | <b>f</b>          |
|                         | <b>mv</b>  | <b>t2, a0</b>     |
|                         | <b>mv</b>  | <b>t0, t2</b>     |

U prvoj liniji (u naredbi **add**) se računa vrednost argumenta poziva funkcije i smešta u radni registar (pretpostavka je da su svi radni registri slobodni). Naredbom **mv** se vrednost zbira smešta u registar za argument. Naredbom **jal** se poziva funkcija. Naredbom **mv** se isporučuje povratna vrednost funkcije. Ako se na mestu argumenta pojave literali ili promenljive, kao u slučaju poziva u sledećem primeru, onda nisu potrebni radni registri, jer se vrednosti literala i promenljivih direktno mogu smeštati na stek.

|                     |            |               |
|---------------------|------------|---------------|
| <b>a = f ( 1 );</b> | <b>li</b>  | <b>a0, 1</b>  |
|                     | <b>jal</b> | <b>f</b>      |
|                     | <b>mv</b>  | <b>t1, a0</b> |
|                     | <b>mv</b>  | <b>t0, t1</b> |

Primeru **if** iskaza odgovara sledeći kod:

|                        |                              |
|------------------------|------------------------------|
| <b>if ( a &lt; b )</b> | <b>if0:</b>                  |
| <b>a = 1;</b>          | <b>bge t0, t1, if_false0</b> |
| <b>else</b>            | <b>if_true0:</b>             |
| <b>a = 2;</b>          | <b>li t0, 1</b>              |
|                        | <b>j if_exit0</b>            |
|                        | <b>if_false0:</b>            |
|                        | <b>li t0, 2</b>              |
|                        | <b>if_exit0:</b>             |

U ekvivalentnom asemblerskom kodu se pojavljuju labele **if**, **if\_true**, **if\_false** i **if\_exit**. Neke od njih su potrebne, jer su ciljevi naredbi skokova (to su **if\_false** i **if\_exit**), dok su labele **if** i **if\_true** nepotrebne, ali zgodne za praćenje i proveru izgenerisanog koda. Labela **if** označava početak **if** iskaza, a labela **if\_exit** njegov kraj. Izvršavanje **if** iskaza započinje proverom uslova i naredbom skoka - **bge** (na **if\_true** ili **if\_false** labelu) u zavisnosti od ispunjenosti uslova. Neposredno iza labele **if\_true** se nalazi kod koji će se izvršiti ukoliko je uslov tačan i iza toga bezuslovan skok na kraj **if** iskaza (odnosno na **if\_exit** labelu), da izvršavanje ne bi “propalo” na **if\_false** labelu. Neposredno iza **if\_false** labele se nalazi kod koji će se izvršiti ukoliko uslov nije tačan. Labele u ekvivalentnom (i izgenerisanom) kodu moraju biti jedinstvene. Svaka labela se završava brojem koji sadrže promenljive **for\_num**, **if\_num**, **while\_num** i **ternary\_num** (u zavisnosti od toga da li se koristi **for**, **if**, **while** ili **ternarni operator**) (aktuelni broj labele). Jednoznačni brojevi se dobijaju inkrementiranjem datih promenljivih za svaki sledeći iskaz. Ukoliko bi se u prethodnom primeru izostavio **else** deo, u ekvivalentnom kodu bi nestao deo koda iza **if\_false** labele, kao na u sledećem kodu:

|                        |                              |
|------------------------|------------------------------|
| <b>if ( a &lt; b )</b> | <b>if0:</b>                  |
| <b>a = 1;</b>          | <b>bge t0, t1, if_false0</b> |
|                        | <b>if_true0:</b>             |
|                        | <b>li t0, 1</b>              |
|                        | <b>j if_exit0</b>            |
|                        | <b>if_false0:</b>            |
|                        | <b>if_exit0:</b>             |

Datoteke **codegen** (.c i .h) sarže funkcije za generisanje koda. Deo generisanja je smešten u ovim funkcijama, a deo u parseru (**micko.y**).

Funkcija **take\_reg()** zauzima prvi slobodni registar i vraća njegov indeks u tabeli simbola.

```
int take_reg(void);
```

Funkcija **free\_reg()** oslobađa poslednji zauzeti registar.

```
void free_reg(void);
```

Funkcija **free\_if\_reg()** oslobađa registar čiji je indeks prosleđen.

```
void free_if_reg(int reg_index);
```

Funkcija **gen\_sym\_name()** ispisuje simbol koji se nalazi na prosleđenom indeksu u tabeli simbola. Funkcija ispisuje simbol u odgovarajućem obliku, u zavisnosti od njegove vrste.

```
void gen_sym_name(int index);
```

Za generisanje naredbe poređenja - **CMP** naredba, predviđena je funkcija **gen\_cmp()** čiji su parametri indeksi operandi u tabeli simbola:

```
void gen_cmp ( int op1_index , int op2_index ) ;
```

Funkcija **gen\_mov\_risc()** generiše **mv** naredbu, a parametri su joj indeks ulaznog i indeks izlaznog operanda u tabeli simbola:

```
void gen_mov_risc ( int input_index , int output_index ) ;
```

Generisanje koda se realizuje u parseru, tako da se faza generisanja koda, praktično, odvija zajedno sa sintaksnom i semantičkom analizom. Sledi **micko.y** datoteka: parser sa semantičkim proverama i generisanjem hipotetskog asemblerskog koda. Parser pre početka parsiranja napravi **.asm** datoteku i zatim u nju generiše kod. Ukoliko parsiranje prođe bez greške, ova datoteka će biti validna. Ako se pojave greške u toku parsiranja, parser će obrisati ovu datoteku, jer sigurno nije validna. U prvom delu bison specifikacije dodate su promenljive: **for\_num**, **if\_num**, **while\_num** i **ternary\_num** - brojači labela i **output** - izlazni fajl u koji će biti izgenerisan hipotetski asemblerski program:

```
FILE *output;  
int for_num = -1;  
int if_num = -1;  
int while_num = -1;  
int ternary_num = -1;
```

Generisanje koda za deniciju funkcije podrazumeva generisanje:

1. labele koja se zove isto kao funkcija
2. tela funkcije (na čijem početku će se izgenerisati zauzimanje prostora za lokalne promenljive)
3. **exit** labele funkcije
4. oslobađanja prostora zauzetog za lokalne promenljive
5. naredbe **ret**, koja će sa steka skinuti povratnu adresu i dovesti do preusmeravanja toka izvršavanja na deo koda iz kog je funkcija pozvana - u slučaju da funkcija nije main. Ukoliko je main funkcija u pitanju umesto **ret** naredbe naći će 2 naredbe za izlazak iz programa

**li a7, 10**  
**ecall**

Prvi element generisanja koda treba da se izvrši na početku parsiranja funkcije, tj. odmah nakon prepoznavanja imena funkcije. Telo funkcije se generiše tokom parsiranja pojma body. Na kraju parsiranja funkcije se generišu elementi 4 i 5.

```
function
: _TYPE_ID
{
    fun_idx = lookup_symbol($2, FUN);
    if(fun_idx == NO_INDEX){
        int* param_types = (int*) malloc(sizeof(int)*128);

        fun_idx = insert_symbol($2, FUN, $1, NO_ATR, NO_ATR);

        parameter_map[fun_idx] = param_types;
    }
    else
        err("redefinition of function '%s'", $2);

    if(strcmp(get_name(fun_idx), "main") == 0){
        code("\n.data\n");
        code("strl:\t.string \"Result is \\\n\n");
        code("arr:\t.word 100\n\n");
        code(".text\n");
    }

    code("\n%s:", $2);
}
_LPAREN parameter_list _RPAREN body
{
    if(strcmp(get_name(fun_idx), "main") == 0)
        code("\n%s_exit:", $2);
    else
        code("\n%s_exit:\n", $2);

    if(var_num){
        int tmp = 0;

        for(int i = var_num;i > 0;i--){
            code("\n\t\tlw\t\t%s, %d(sp)", get_name(i - 1), tmp);
            tmp += 4;
            free_if_reg(i - 1);
        }

        code("\n\t\taddi\tsp, sp, %d\n", 4 * var_num);
    }

    if(strcmp(get_name(fun_idx), "main") != 0)
        code("\t\tret\n");

    if(strcmp(get_name(fun_idx), "main") == 0){
        code("\n\t\tli a7, 10");
        code("\n\t\tcall");
    }

    clear_symbols(fun_idx + 1);
    var_num = 0;
}
```



Ukoliko se u telu funkcije (pojam **body**) nalaze lokalne promenljive, potrebno je izgenerisati zauzimanje prostora za njih na steku. Ova akcija se može sprovesti tek kada se zna ukupan broj lokalnih promenljivih, a to je nakon parsiranja pojma **variable\_list**. Zauzimanje prostora se vrši pomeranjem stek pokazivača za onoliko lokacija koliko ima lokalnih promenljivih (svaka lokacija zauzima 4 bajta). Zato se u ovu svrhu generiše oduzimanje vrednosti **4\*var\_num** od registra **sp**. Zbog preglednosti izgenerisanog koda, na početku tela funkcije, generiše se **body** labela funkcije.

```
body
: _LBRACKET variable_list
{
    par_num = 0;
    if(var_num){
        code("\n\t\taddi\tsp, sp, -%d", 4 * var_num);

        int tmp = var_num;
        int reg = take_reg();

        for(int i = reg; i >= 0; i--){
            free_if_reg(i);
        }

        for(int i = 0; i < var_num; i++){
            reg = take_reg();

            code("\n\t\tsw\t\t");
            gen_sym_name(reg);
            code(", %d(sp)", (4 * tmp) - 4);
            tmp -= 1;
        }

        if(strcmp(get_name(fun_idx), "main") == 0)
            code("\n\t\tla\t\t%s, arr", get_name(HELP_REG));
        else{
            int reg = take_reg();

            for(int i = reg; i >= 0; i--){
                free_if_reg(i);
            }

            code("\n%s_body:", get_name(fun_idx));
        }
    }
    statement_list _RBRACKET
;
```

Kao ekvivalent iskaza dodele, generiše se **mv** naredba - ukoliko je **\$3** parametar, varijabla, registar ili **li** naredba - ukoliko je **\$3** literal. Funkciji

**gen\_mov\_risc()** se, kao prvi argument, prosleđuje **\$3** koji sadrži indeks elementa tabele simbola na kom se nalazi ulazni operand naredbe **mv**, tj. rezultat izraza sa desne strane znaka jednakosti. Drugi argument je indeks elementa tabele simbola na kom se nalazi identikator sa leve strane znaka jednakosti, koji predstavlja izlazni operand naredbe **mv**.

```
assignment_statement
: _ID _ASSIGN num_exp _SEMICOLON
{
    int idx = lookup_symbol($1, VAR|PAR);
    if(idx == NO_INDEX)
        err("invalid lvalue '%s' in assignment", $1);
    else
        if(get_type(idx) != get_type($3)){
            err("incompatible types in assignment");
        }

    if(get_kind($3) == LIT){
        code("\n\t\tli\t\t");
        gen_sym_name(idx);
        code(", ");
        gen_sym_name($3);
    }
    else if(get_kind($3) == PAR){
        code("\n\t\tmv\t\t");
        gen_sym_name(idx);
        code(", ");
        gen_sym_name($3);
    }
    else if(get_kind(idx) == PAR){
        code("\n\t\tmv\t\t");
        gen_sym_name(idx);
        code(", ");
        gen_sym_name($3);
    }
    else{
        gen_mov_risc(idx, $3);
    }
}
;
```

Za generisanje aritmetičkih operacija koriste se nizovi stringova **ar\_instructions[]** (u **defs.h**) koji sadrže varijante stringova (imena naredbi) za označene i neoznačene tipove podataka. Za indeksiranje ovog niza potrebna je oznaka aritmetičke operacije (konstanta **add**, **sub**, **mul** ili **div** (ove 4 konstante se koriste u slučaju da bilo koja od operacija odvijaju između 2 operanda), koja se nalazi u metapromenljivoj **\$2** a koja je denisana u enumeraciji **arops**) i oznaka tipa (bilo kog) operanda. Kada se ove dve informacije iskoriste, ona će dobiti string odgovarajuće aritmetičke naredbe.

Prvi i drugi (ulazni) operand naredbe se generišu pomoću funkcije **gen\_sym\_name()** jer imamo njihove indekse u tabeli simbola. Ukoliko je neki od ulaznih operanada bio smešten u registru, nakon preuzimanja vrednosti on se može osloboditi. U ovu svrhu se koristi funkcija **free\_if\_reg()**. Prilikom generisanja izlaznog operanda aritmetičke naredbe, zauzima se prvi slobodan registar pomoću funkcije **take\_reg()** i u njega smešta rezultat operacije. Vrednost **num\_exp** pojma biće indeks tog registra u tabeli simbola (kao mesto gde se nalazi rezultat parsiranog numeričkog izraza). Ukoliko se operacija odvija između operanda i literala onda se za to koriste drugi tip naredbi - za sabiranje **addi**, za oduzimanje isto **addi** s tim što ispred literala onda stoji znak -, množenje se odvija preko sabiranja a deljenje preko oduzimanja.

```

num_exp
: exp
| num_exp _AROP exp
{
    if(get_type($1) != get_type($3))
        err("invalid operands: arithmetic operation");

    int t1 = get_type($1);

    free_if_reg($3);
    free_if_reg($1);

    if(get_kind($3) == LIT && $2 == 0)
        code("\n\t\taddi\t");
    else if(get_kind($3) == LIT && $2 == 1)
        code("\n\t\taddi\t");
    else if(get_kind($3) == LIT && $2 == 2){
        multiplication_num++;

        code("\n");
        if(get_kind($1) == PAR){
            code("\n\t\tmv\t");
            gen_sym_name(TMP_REG);
            code(", ");
            gen_sym_name($$);
        }
        else
            gen_mov_risc(TMP_REG, $$);

        $$ = take_reg();
        set_type($$, t1);
        code("\n\t\tli\t");
        gen_sym_name($$);
        code(", 0\n");

        int reg1 = take_reg();
        code("\t\tli\t");
        gen_sym_name(reg1);
        set_type(reg1, t1);
        code(", 0\n");

        int reg2 = take_reg();
        code("\t\tli\t");
        gen_sym_name(reg2);
        set_type(reg2, t1);
        code(", %s\n", get_name($3));

        code("mul%d:\n", multiplication_num);

        code("\t\tadd\t");
        gen_sym_name($$);
        code(", ");
        gen_sym_name($$);
        code(", %s\n", get_name(TMP_REG));

        code("\t\taddi\t");
        gen_sym_name(reg1);
        code(", ");
        gen_sym_name(reg1);
        code(", 1\n");
    }
}

```

```

        code("\t\tblt\t\t");
        gen_sym_name(reg1);
        code(", ");
        gen_sym_name(reg2);
        code(", mul%d\n", multiplication_num);

        code("mul%d_end:", multiplication_num);

        free_if_reg(reg2);
        free_if_reg(reg1);
    }
    else if (get_kind($3) == LIT && $2 == 3){
        division_num++;

        code("\n");
        if (get_kind($1) == PAR){
            code("\n\t\tmv\t\t");
            gen_sym_name(TMP_REG);
            code(", ");
            gen_sym_name($$);
        }
        else
            gen_mov_risc(TMP_REG, $$);

        $$ = take_reg();
        set_type($$, t1);
        code("\n\t\tli\t\t");
        gen_sym_name($$);
        code(", 0\n");

        int reg1 = take_reg();
        code("\t\tmv\t\t");
        gen_sym_name(reg1);
        set_type(reg1, t1);
        code(", %s\n", get_name(TMP_REG));

        code("div%d:\n", division_num);

        code("\t\taddi\t\t");
        gen_sym_name(reg1);
        code(", ");
        gen_sym_name(reg1);
        code(", -%s\n", get_name($3));

        code("\t\tblt\t\t");
        gen_sym_name(reg1);
        code(", zero, div%d_exit\n", division_num);

        code("\t\taddi\t\t");
        gen_sym_name($$);
        code(", ");
        gen_sym_name($$);
        code(", 1\n");

        code("\t\tj\t\t\tdiv%d\n", division_num);
        code("div%d_exit:", division_num);

        free_if_reg(reg1);
    }
}

```

```

else
    code("\n\t\t%s\t\t", ar_instructions[$2 + (t1 - 1) * AROP_NUMBER]);

if((get_kind($3) != LIT && $2 != 2) ||
   (get_kind($3) != LIT && $2 == 2) ||
   (get_kind($3) == LIT && $2 == 0) ||
   (get_kind($3) == LIT && $2 == 1)){

    $$ = take_reg();
    gen_sym_name($$);
    set_type($$, t1);

    code(", ");
    gen_sym_name($1);
    code(", ");

    if(get_kind($3) == LIT && $2 == 1)
        code("-%s", get_name($3));
    else
        gen_sym_name($3);
}
}
;

```

U pravilu za izraze (**exp**), pravilo, koje opisuje poziv funkcije i pravilo za ternarni operator, sadrži generisanje koda. Kada se poziv funkcije nađe negde u izrazu, njenu povratnu vrednost treba kopirati u prvi slobodan registar. Ovo je neophodno, jer je moguće da se u izrazu ponovo zatekne poziv neke funkcije koji će njenu povratnu vrednost ponovo smestiti u registar **a0** (i time poništiti prethodnu vrednost). Zato povratne vrednosti funkcija u izrazima treba čuvati u registrima. Generisanje koda za ternarni izraz podrazumeva: zauzimanje prvog slobodnog registra, skok na **if\_false** labelu ukoliko relacioni izraz nije zadovoljen gde se pomera drugi **cond\_exp** u registar, generisanje **if\_true** labele i ukoliko je taj uslov ispunjen onda se prvi **cond\_exp** premešta u taj registar i skače se na **if\_exit** labelu. Vrednost pojma **exp** postaje indeks elementa u tabeli simbola, u kom se nalazi rezultat izvršavanja funkcije.

```

exp
: literal
| _ID
{
    $$ = lookup_symbol($1, VAR|PAR);
    if($$ == NO_INDEX)
        err("'s' undeclared", $1);
}
| function_call
{
    $$ = take_reg();
    gen_mov_risc($$, FUN_REG);
}
| _LPAREN num_exp _RPAREN
{ $$ = $2; }
| _LPAREN rel_exp _RPAREN _QMARK cond_exp _COLON cond_exp
{
    int out = take_reg();
    ++ternary_num;

    if(get_type($5) != get_type($7))
        err("expressions don't have same types in ternary exp");

    code("ternary_false%d", ternary_num);
    code("\nternary_true%d:", ternary_num);

    if(get_kind($5) == LIT){
        code("\n\t\tli\t\t");
        gen_sym_name(out);
        code(", ");
        gen_sym_name($5);
        set_type(out, get_type($5));
    }
    else
        gen_mov_risc(out, $5);

    code("\n\t\tj\t\tternary_exit%d", ternary_num);
    code("\nternary_false%d:", ternary_num);

    if(get_kind($7) == LIT){
        code("\n\t\tli\t\t");
        gen_sym_name(out);
        code(", ");
        gen_sym_name($7);
        set_type(out, get_type($7));
    }
    else
        gen_mov_risc(out, $7);

    code("\nternary_exit%d:", ternary_num);
    $$ = out;
}
;

```

Za generisanje poziva funkcije koristi se naredba **jal**. Vrednost pojma **function\_call** postaje registar **a0**, a njegov tip se postavlja na povratni tip funkcije. Dodatni kod iznad i ispod ispisa za **jal** služi da očuva postojeće vrednosti lokalnih varijabli pre nego što se uđe u funkciju na lokaciju **HELP\_REG** i nakon izvršavanja funkcije da vrati te vrednosti u lokalne varijable.

```
function_call
: _ID
{
    fcall_idx = lookup_symbol($1, FUN);
    if(fcall_idx == NO_INDEX)
        err("'s' is not a function", $1);

    arg_counter = 0;
}
_LPAREN argument_list _RPAREN
{
    if(get_atrl(fcall_idx) != arg_counter)
        err("wrong number of arguments to function 's'", get_name(fcall_idx));

    if(var_num){
        int tmp = var_num;
        code("\n");
        for(int i = 0; i < var_num; i++){
            code("\n\t\tsw\t\t");
            gen_sym_name(i);
            code(", %d(s)", (4 * tmp) - 4, get_name(HELP_REG));
            tmp -= 1;
        }
        code("\n");
    }

    code("\n\t\tjal\t\t", get_name(fcall_idx));

    if(var_num){
        int tmp = 0;
        code("\n");
        for(int i = var_num; i > 0; i--){
            code("\n\t\tlw\t\t", %d(s)", get_name(i - 1), tmp, get_name(HELP_REG));
            tmp += 4;
        }
        code("\n");
    }

    set_type(FUN_REG, get_type(fcall_idx));
    $$ = FUN_REG;

    arg_counter = 0;
}
;
```



U toku parsiranja pojma argument treba generisati naredbu **mv** ili **li** u zavisnosti od toga da li se radi o smeštanju registara ili varijabli u registar/re za argument ili literala u registar/re za argument.

```
argument
: num_exp
{
    if(parameter_map[fcall_idx][arg_counter] != get_type($1))
        err("incompatible type for argument in '%s'", get_name(fcall_idx));
    free_if_reg($1);

    if(get_kind($1) == VAR)
        gen_mov_risc(FUN_REG + arg_counter, $1);
    if(get_kind($1) == LIT){
        code("\n\t\tli\t\t");
        gen_sym_name(FUN_REG + arg_counter);
        code(", ");
        gen_sym_name($1);
    }
    if(get_kind($1) == REG){
        code("\n\t\tmv\t\t");
        gen_sym_name(FUN_REG + arg_counter);
        code(", ");
        gen_sym_name($1);
    }
    arg_counter += 1;
}
```

U okviru if iskaza treba izgenerisati više asemblerskih naredbi. Na početku if iskaza generiše se labela **ifX**, gde je **X** redni broj iskaza u programu. Brojač (**if\_num**) sa ovom vrednošću se na početku generisanja if iskaza inkrementira, da bi označio sledeći iskaz. Na mestu relacionog izraza će se izgenerisati skok sa suprotnim uslovom na odgovarajuću **if\_false** labelu. Ako je, na primer, u relaciji naveden operator **==**, izgenerisaće se skok **bne**. U ovu svrhu se koristi enumeracija **opp\_jumps[]** koja sadrži imena naredbi skokova ali u kontra redosledu od konstanti relacionih operatora(denisano u **defs.h**). U nastavku se generiše true labela sa brojem **if\_num**. Iza true labele, u toku parsiranja pojma statement, biće izgenerisani iskazi koji čine **then** telo if-a. Posle toga, treba izgenerisati bezuslovni skok **j** na kraj if iskaza, čime je obezbeđeno da se, odmah nakon then tela, neće izvršiti i **else** telo. Zatim sledi generisanje **if\_false** labele, iza koje će uslediti generisanje **else** tela (ako postoji).



Vrednost pojma **rel\_exp** postaje redni broj naredbe skoka u nizu **opp\_jumps[]**.

```

rel_exp
: num_exp _RELOP num_exp
{
    if(get_type($1) != get_type($3))
        err("invalid operands: relational operator");
    $$ = $2 + ((get_type($1) - 1) * RELOP_NUMBER);

    if(get_kind($3) == LIT && if_flag == 1){
        if_reg = take_reg();
        code("\n\t\tli\t\t%s, %s", get_name(if_reg), get_name($3));
    }
    else if(get_kind($3) == LIT && while_flag == 1){
        while_reg = take_reg();
        code("\n\t\tli\t\t%s, %s", get_name(while_reg), get_name($3));
    }
    else if(get_kind($3) == LIT && for_flag == 1){
        for_reg = take_reg();
        code("\n\t\tli\t\t%s, %s", get_name(for_reg), get_name($3));
    }

    code("\n\n\t\t%s\t\t", opp_jumps[$$]);
    gen_sym_name($1);
    code(", ");

    if(get_kind($3) == LIT && if_flag == 1){
        code("%s", get_name(if_reg));
        if_flag = 0;
    }
    else if(get_kind($3) == LIT && for_flag == 1){
        code("%s", get_name(for_reg));
        for_flag = 0;
    }
    else if(get_kind($3) == LIT && while_flag == 1){
        code("%s", get_name(while_reg));
        while_flag = 0;
    }
    else
        gen_sym_name($3);

    code(", ");
}
:

```

U okviru **return** iskaza treba izgenerisati **mv** ili **li** naredbu koja će rezultat izraza **num\_exp**(u zavisnosti da li je **num\_exp** literal, parametar, varijabla ili registar) prebaciti u registar **a0** (**FUN\_REG**). Iza toga se generiše безусловni skok (**j**) na **exit** labelu, koja označava kraj tela funkcije, čime se realizuje semantika **return** iskaza.



```

int main() {
    int synerr;
    init_symtab();
    output = fopen("output.asm", "w+");

    synerr = yyparse();

    clear_symtab();
    fclose(output);

    if(warning_count)
        printf("\n%d warning(s).\n", warning_count);

    if(error_count) {
        remove("output.asm");
        printf("\n%d error(s).\n", error_count);
    }

    if(synerr)
        return -1; //syntax error
    else if(error_count)
        return error_count & 127; //semantic errors
    else if(warning_count)
        return (warning_count & 127) + 127; //warnings
    else
        return 0; //OK
}

```

Za generisanje **for** iskaza radi se sledeće: literal se pomeri u promenljivu koja se koristi u **for** iskazu pomoću naredbe **li** pa se zatim izgeneriše labela **for**. Onda će se zbog prisustva **rel\_exp** izgenerisati naredba poređenja i skoka na **for\_exit** labelu u slučaju da **rel\_exp** nije ispunjen. Zatim se generiše kod pojam **statement** koji je već objašnjen. Nakon njega će se izgenerisati inkrement iskaz za koji se koristi **addi** i funkcija **gen\_sym\_name()**. Nakon njih će se izgenerisati skok na **for** labelu a zatim i sama **for\_exit** labela.

```

for_statement
: _FOR _LPAREN _ID _ASSIGN literal
{
    $<i>$ = ++for_num;
    for_flag = 1;

    int i = lookup_symbol($3, VAR|PAR);
    if(i == NO_INDEX)
        err("'%' undeclared", $3);

    if(get_type(i) != get_type($5))
        err("incompatible types");

    code("\n\n\t\tli\t\t");
    gen_sym_name(i);
    code(", ");
    gen_sym_name($5);

    code("\nfor%d:", for_num);
}
_SEMICOLON rel_exp
{
    code("for_exit%d", $<i>6);
}
_SEMICOLON _ID _INC _RPAREN statement
{
    int i = lookup_symbol($11, VAR|PAR);
    if(i == NO_INDEX)
        err("'%' undeclared", $11);

    code("\n\t\ttaddi\t");

    gen_sym_name(i);
    code(", ");
    gen_sym_name(i);
    code(", 1");

    code("\n\t\ttj\t\t\tfor%d", $<i>6);
    code("\nfor_exit%d:", $<i>6);

    free_if_reg(for_reg);
}
;

```

Generisanje koda za **while** petlju je slično kao i za **for** petlju pa čak i jednostavnije. Prvo će se izgenerisati **while** labela nakon koje će zbog prisustva **rel\_exp** biti izgenerisana naredba poređenja i skoka na **while\_exit** labelu ukoliko **rel\_exp** nije ispunjen. Nakon toga će se izgenerisati kod za pojam **statement**. I na kraju biće izgenerisan skok na **while** labelu nakon koje će biti izgenerisana sama **while\_exit** labela.

```

while_statement
: _WHILE
{
    .....
    $<i>$ = ++while_num;
    while_flag = 1;
    code("\nwhile%d:", while_num);
}
_LPAREN rel_exp
{
    .....
    code("while_exit%d", $<i>2);
}
_RPAREN statement
{
    .....
    code("\n\t\tj\t\twhile%d", $<i>2);
    code("\nwhile_exit%d:", $<i>2);
    free_if_reg(while_reg);
}
;

```