

## **Poglavje 3**

# **SINTAKSA**

### **3.1 Osnovni sintaktični objekti**

Uporabljali bomo dve vrsti objektov: nize in abstraktno sintaksno drevo.

Nizi so primerna podatkovna struktura za linearno predstavitev jezika, niso pa primerni za analizo.

Abstraktno sintaksno drevo prikažejo hierarhično strukturo sintakse zato je primerno za semantično analizo.

#### **3.1.1 Simboli**

Uporabljali bomo množico simbolov, ki bodo služili v različnih vlogah: znaki, spremenljivke, imena polj, ...

Simbolom bomo včasih rekli imena ali atomi ali identifikatorji v odvisnosti od navad v konkretnem delovnem okolju.

O simbolih bomo razmišljali kot o atomih, ki nimajo strukture razen identitete.

Napišemo  $x$  sym kar pomeni, da je  $x$  simbol in predpostavljamo, da imamo neskončno mnogo simbolov za to identiteto.

Sodba  $x \# y$ , kjer je  $x$  sym in  $y$  sym, pove, da sta  $x$  in  $y$  različna simbola.

Uporabljali bomo različne razrede simbolov. V splošnem bomo predpostavili, da sta katerakoli razreda simbolov med sabo različna, tako da ne bo zmešnjav med njimi.

### 3.1.2 Nizi nad abecedo

Abeceda je končna množica znakov  $c_1$  sym, ...,  $c_n$  sym. Naj  $\Sigma$  označuje končno množico takšnih sodb.

Forma sodbe za  $\Sigma \vdash s$  str, ki definira nize nad abecedo  $\Sigma$ , je induktivno definirana z naslednjimi pravili.

$$\overline{\Sigma \vdash \varepsilon \text{ str}} \quad (3.1)$$

$$\frac{\Sigma \vdash c \text{ sym} \quad \Sigma \vdash s \text{ str}}{\Sigma \vdash c \cdot s \text{ str}} \quad (3.2)$$

Niz je torej sekvenca znakov oz. null v primeru, da ne vsebuje znakov. Pogosto ne omenjamo  $\Sigma$ , če je abeceda jasna iz konteksta.

V primeru nizov uporabljamo indukcijo na sledeč način. Ko hočemo pokazati lastnost  $P$  nad nizom  $s$  je zadosti, da vidimo sledeče:

- $\varepsilon P$
- Ob predpostavki  $s P$  in  $c$  sym pokažemo  $c \cdot s P$

To včasih imenujemo princip *indukcije nad nizi*. Metoda je ekvivalentna indukciji na dolžino niza le da ni potrebno definirati dolžine niza.

Naslednje pravilo induktivno definira izjavo  $s_1 \hat{s}_2 = s \text{ str}$ , ki pravi, da je  $s$  konkatencija  $s_1$  in  $s_2$ .

$$\overline{\epsilon \hat{s} \text{ str}} \quad (3.3)$$

$$\frac{s_1 \hat{s}_2 = s \text{ str}}{(c \cdot s_1) \hat{s}_2 = c \cdot s \text{ str}} \quad (3.4)$$

Enostavno je pokazati z indukcijo, da pravila definirajo totalno funkcijo ( $\forall \exists !$ ) na dveh argumentih.

Nizi so običajno zapisani kot sekvenca simbolov npr. "abcd" je predstavitev niza  $a \cdot (b \cdot (c \cdot (d \cdot e)))$ .

### 3.1.3 Abstraktna sintaksna drevesa

Abstraktno sintaksno drevo (okr. ast) je urejeno drevo v katerem so vozlišča označena s simboli, ki jih imenujemo *operatorji*.

Signatura  $\Omega$  je končna množica sodb oblike  $ar(o) = n$ , kjer velja  $o \text{ sym}$  in  $n \text{ nat}$ . Te sodbe priredijo števnost  $n$  operatorjem  $o$ , tako da če  $\Omega \vdash ar(o) = n$  in  $\Omega \vdash ar(o) = n'$  potem  $n = n'$ .

Razred abstraktnih sintaksnih dreves nad signaturo  $\Omega$  je induktivno definiran na sledeč način.

$$\frac{\Omega \vdash ar(o) = n \quad a_1 \text{ ast} \dots a_n \text{ ast}}{o(a_1, \dots, a_n) \text{ ast}} \quad (3.5)$$

Osnovni primer induktivne definicije je operator s števnostjo nič. V tem primeru je premisa prazna. Takšna vozlišča drevesa so listi.

### Strukturna indukcija

Princip strukturne indukcije je specializacija principa indukcije pravil na pravila, ki definirajo abstraktna sintaksna drevesa nad signaturami.

Če želimo pokazati  $P(a \text{ ast})$  je zadosti, da pokažemo da je  $P$  zaprta glede na pravilo 3.5. Z drugimi besedami, če  $\Omega \vdash ar(o) = n$ , potem moramo pokazati, da

$$\text{if } P(a_1 \text{ ast}), \dots, P(a_n \text{ ast}) \text{ then } P(o(a_1, \dots, a_n) \text{ ast}). \quad (3.6)$$

V primeru, da je  $n = 0$  se dokaz reducira na to, da pokažemo  $P(o())$ .

**Primer 3.1.1** Poglejmo si sledečo induktivno definicijo višine ast.

$$\frac{hgt(a_1) = h_1 \dots hgt(a_n) = h_n \quad \max(h_1, \dots, h_n) = h}{hgt(o(a_1, \dots, a_n)) = succ(h)} \quad (3.7)$$

Z uporabo strukturne indukcije lahko pokažemo  $(\forall, \exists!)$ , kar pomeni, da ima vsak ast unikatno višino.

Predpostavimo lahko, da imajo vsa ast višine  $1 \leq i \leq n$  unikatno višino  $hgt(a_i) = h_i$ . Pokažemo lako, da je maksimum vseh unikatno definiran in je torej celotna višina unikatna.  $\square$

### Spremenljivke in substitucija

V praksi pogosto potrebujemo abstraktna sintaksna drevesa s spremenljivkami kot vrzeli, ki jih je mogoče napolniti z drugimi drevesi.

Spremenljivke so instancirane z uporabo substitucije spremenljivk ast z dugim drevesom.

Poglejmo si najprej nekaj dogovorov v povezavi z notacijo. Naj bo  $\mathcal{X} = x_1 \text{ ast}, \dots, x_n \text{ ast}$  množica parametrov in hkrati seznam hipotez  $\{x_1, \dots, x_n\} \mid x_1 \text{ ast}, \dots, x_n \text{ ast}$ , kjer velja tudi  $x_1 \text{ sym}, \dots, x_n \text{ sym}$ . Naprej, naj izraz  $x \# \mathcal{X}$  pomeni  $x \notin \{x_1, \dots, x_n\}$ .

Z uporabo te notacije lahko napišemo sodbo  $\mathcal{X} \vdash a \text{ ast}$ , ki je induktivno definirana z naslednjimi pravili:

$$\overline{\mathcal{X}, x \text{ ast} \vdash x \text{ ast}} \quad (3.8)$$

$$\frac{\Omega \vdash ar(o) = n \quad \mathcal{X} \vdash a_1 \text{ ast} \dots \mathcal{X} \vdash a_n \text{ ast}}{\mathcal{X} \vdash o(a_1, \dots, a_n) \text{ ast}} \quad (3.9)$$

Po principu indukcije pravil je za dokaz  $P(\mathcal{X} \vdash a \text{ ast})$  zadosti, da pokažemo:

1.  $P(\mathcal{X}, x \text{ ast} \vdash x \text{ ast})$ .
2. Če  $\Omega \vdash ar(o) = n$  in če  $P(\mathcal{X} \vdash a_1 \text{ ast}), \dots, P(\mathcal{X} \vdash a_n \text{ ast})$ , potem  $P(\mathcal{X} \vdash o(a_1, \dots, a_n) \text{ ast})$ .

Parametri  $\mathcal{X}$  se obravnavajo kot atomični objekti, vendar ima vsak svoje abstraktno sintaksno drevo.

Definiramo sodbo  $\mathcal{X} \vdash [a/x]b = c$ , kar pomeni da je  $c$  rezultat substitucije  $a$  za  $x$  v  $b$  z uporabo naslednjih pravil.

$$\overline{\mathcal{X}, x \text{ ast} \vdash [a/x]x = a} \quad (3.10)$$

$$\frac{x \# y}{\mathcal{X}, x \text{ ast}, y \text{ ast} \vdash [a/x]y = y} \quad (3.11)$$

$$\frac{\mathcal{X} \vdash [a/x]b_1 = c_1 \dots \mathcal{X} \vdash [a/x]b_n = c_n}{\mathcal{X} \vdash [a/x]o(b_1, \dots, b_n) = o(c_1, \dots, c_n)} \quad (3.12)$$

Rezultat substitucije je enolično definiran z argumenti. Kot posledico lahko napišemo  $[a/x]b$  za unikatno  $c$ , tako da  $[a/x]b = c$ .

**Izrek 3.1.1** Če velja  $\mathcal{X} \vdash a$  ast in  $\mathcal{X}, x \text{ ast} \vdash b$  ast, kjer  $x \notin \mathcal{X}$ , potem obstaja enoličen  $c$ , tako da  $\mathcal{X} \vdash [a/x]b = c$  in  $\mathcal{X} \vdash c$  ast.

*Dokaz.* Trditev dokažemo s strukturno indukcijo na  $b$  relativno na kontekst  $\mathcal{X}, x \text{ ast}$ .

Imamo tri primere za obravnavo:

1. Ker velja  $\mathcal{X}, x \text{ ast} \vdash x$  ast, moramo pokazati da obstaja unikatno  $c$  tako da  $\mathcal{X} \vdash [a/x]x = c$ . Z uporabo pravila 3.10a vidimo da je zamenjava  $c$  z  $a$  potrebna in zadostna.
2. Če  $\mathcal{X}, x \text{ ast}, y \text{ ast} \vdash y$  ast za nek  $y \notin \mathcal{X}$ , potem po pravilu 3.10b zamenjava  $c$  z  $y$  je potrebna in zadostna.
3. Končno, če  $b = o(b_1, \dots, b_n)$ , potem velja po indukciji, da obstaja enoličen  $c_1, \dots, c_n$  tako da  $\mathcal{X} \vdash [a/x]b_1 = c_1, \dots, \mathcal{X} \vdash [a/x]b_n = c_n$ . Po pravilu 3.10c je edina možna izbira za  $c$  drevo  $o(c_1, \dots, c_n)$ , kar zadostuje.  $\square$

## 3.2 Konkretna sintaksa

Konkretna sintaksa jezika se uporablja za predstavitev gramatike z navadnim tekstom. Velikokrat uporabimo konkretno sintakso za izboljšanje berljivosti in izločanje dvoumnosti jezika.

V tej sekciji predstavimo glavne metode za specifikacijo konkretne sintakse. Uporabljali bomo jezik izrazov  $\mathcal{L}(\text{nat str})$ , ki predstavi elementarno aritmetiko na naravnih številih, enostavne operacije na nizih in enostavno delo s spremenljivkami.

Tip $t$	::=	num	num	
		str	str	
Izraz $e$	::=	$x$	$x$	
		num[ $n$ ]	$n$	
		str[ $s$ ]	$s$	
		plus( $e_1; e_2$ )	$e_1 + e_2$	(3.13)
		times( $e_1; e_2$ )	$e_1 * e_2$	
		cat( $e_1; e_2$ )	$e_1 \hat{e}_2$	
		len( $e$ )	$ e $	
		let( $e_1; x.e_2$ )	let $x$ be $e_1$ in $e_2$	

Jezik  $\mathcal{L}(\text{nat str})$  vsebuje osnovne operacije za delo z naravnimi števili in nizi. Edini sestavljen gradnik jezika je let stavek, ki izvede substitucijo spremenljivke v stavku.

### 3.2.1 Leksikalna struktura

Prva faza sintaktičnega procesiranja je pretvorba iz predstavitve osnovane na znakih v predstavitev osnovani na simbolih. Ta proces običajno imenujemo *leksikalna analiza*. Osnovna ideja leksikalne analize je združiti znake v simbole, ki so osnova za naslednje faze analize.

Na primer, niz znakov 467 se spremeni v celo število 467. Podobno predstavimo identifikator sestavljen iz niza znakov "temp" v spremenljivko z imenom *temp*.

Med leksikalno analizo počistimo tudi vse "bele znake" (presledke, tabulatorje, itd.) in komentarje, ki se ne uporabljajo v nadaljni analizi.

Znakovna predstavitev je v največ primerih definirana z regularnimi izrazi. Leksikalna struktura  $\mathcal{L}(\text{nat str})$  je določena z naslednjimi pravili.

<i>Item</i>	<i>itm</i>	::=	<i>kwd</i>   <i>id</i>   <i>num</i>   <i>lit</i>   <i>spl</i>	
<i>Keyword</i>	<i>kwd</i>	::=	<i>l.e.t.ε</i>   <i>b.e.ε</i>   <i>i.n.ε</i>	
<i>Identifier</i>	<i>id</i>	::=	<i>ltr</i> ( <i>ltr</i>   <i>dig</i> )*	
<i>Numeral</i>	<i>num</i>	::=	<i>dig dig</i> *	
<i>Literal</i>	<i>lit</i>	::=	<i>qum</i> ( <i>ltr</i>   <i>dig</i> )* <i>qum</i>	(3.14)
<i>Special</i>	<i>spl</i>	::=	+   *   ^   (   )	
<i>Letter</i>	<i>ltr</i>	::=	<i>a</i>   <i>b</i>   ...	
<i>Digit</i>	<i>dig</i>	::=	0   1   ...	
<i>Quote</i>	<i>qum</i>	::=	"	

Leksikalni simbol je bodisi ključna beseda, identifikator, število, niz znakov ali posebni simbol. Simboli se med seboj razlikujejo po sestavi znakov, ki jih posamezna vrsta simbola vsebuje. Simboli, ki predstavljajo nize vsebujejo sekvence natisljivih znakov, ki so obdani z narekovaji.

Delo leksikalnega analizatorja je pretvorba nizov znakov v niz simbolov. Zgornja definicija leksikalne sintakse služi kot vodilo pri pretvorbi.

Vhoden niz znakov se pretvarja v niz simbolov z uporabo naslednjih pravil, ki definirajo simbolični jezik razčlenjevalnika oz. statično semantiko jezika.

$$\frac{s \text{ str}}{ID[s] \text{ tok}} \quad (3.15)$$

$$\frac{n \text{ nat}}{NUM[n] \text{ tok}} \quad (3.16)$$

$$\frac{s \text{ str}}{LIT[s] \text{ tok}} \quad (3.17)$$

$$\overline{LET \text{ tok}} \quad (3.18)$$

$$\overline{BE \text{ tok}} \quad (3.19)$$

$$\overline{IN \text{ tok}} \quad (3.20)$$

$$\overline{ADD \text{ tok}} \quad (3.21)$$

$$\overline{MUL \text{ tok}} \quad (3.22)$$



$$\overline{CAT tok} \quad (3.23)$$

$$\overline{LP tok} \quad (3.24)$$

$$\overline{RP tok} \quad (3.25)$$

$$\overline{VB tok} \quad (3.26)$$

Lesikalna analiza je definirana s pretvorbo nizov znakov v nize simbolov. Pretvorba je induktivno definirana s pomočjo naslednjih pravil:

$$\begin{array}{ll}
 s \text{ inp} \longleftrightarrow t \text{ tokstr} & \text{Scan input} \\
 s \text{ itm} \longleftrightarrow t \text{ tok} & \text{Scan an item} \\
 s \text{ kwd} \longleftrightarrow t \text{ tok} & \text{Scan a keyword} \\
 s \text{ id} \longleftrightarrow t \text{ tok} & \text{Scan an identifier} \\
 s \text{ num} \longleftrightarrow t \text{ tok} & \text{Scan a number} \\
 s \text{ spl} \longleftrightarrow t \text{ tok} & \text{Scan a symbol} \\
 s \text{ lit} \longleftrightarrow t \text{ tok} & \text{Scan a string literal}
 \end{array} \quad (3.27)$$

Definicija form sodb, ki sledijo, uporablja številne pomožne sodbe, ki ustrezajo klasifikaciji znakov in leksikalni strukturi jezika. Na primer,  $s \text{ whs}$  pravi da se niz  $s$  sestoji samo iz "belih znakov", izjava  $s \text{ lord}$  pravi, da se  $s$  sestoji iz alfanumeričnih znakov.

$$\overline{\epsilon \text{ inp} \longleftrightarrow \epsilon \text{ tokstr}} \quad (3.28)$$

$$\frac{s = s_1 \hat{\ } s_2 \hat{\ } s_3 \text{ str} \quad s_1 \text{ whs} \quad s_2 \text{ itm} \longleftrightarrow t \text{ tok} \quad s_3 \text{ inp} \longleftrightarrow ts \text{ tokstr}}{s \text{ inp} \longleftrightarrow t \cdot ts \text{ tokstr}} \quad (3.29)$$

$$\frac{s \text{ kwd} \longleftrightarrow t \text{ tok}}{s \text{ itm} \longleftrightarrow t \text{ tok}} \quad (3.30)$$

$$\frac{s \text{ id} \longleftrightarrow t \text{ tok}}{s \text{ itm} \longleftrightarrow t \text{ tok}} \quad (3.31)$$

$$\frac{s \text{ num} \longleftrightarrow t \text{ tok}}{s \text{ itm} \longleftrightarrow t \text{ tok}} \quad (3.32)$$

$$\frac{s \text{ lit} \longleftrightarrow t \text{ tok}}{s \text{ itm} \longleftrightarrow t \text{ tok}} \quad (3.33)$$

$$\frac{s \text{ spl} \longleftrightarrow t \text{ tok}}{s \text{ itm} \longleftrightarrow t \text{ tok}} \quad (3.34)$$

$$\frac{s = l \cdot e \cdot t \cdot \epsilon \text{ str}}{s \text{ kwd} \longleftrightarrow LET \text{ tok}} \quad (3.35)$$

$$\frac{s = b \cdot e \cdot \epsilon \text{ str}}{s \text{ kwd} \longleftrightarrow BE \text{ tok}} \quad (3.36)$$

$$\frac{s = i \cdot n \cdot \epsilon \text{ str}}{s \text{ kwd} \longleftrightarrow IN \text{ tok}} \quad (3.37)$$

$$\frac{s = s_1 \hat{s}_2 \text{ str} \quad s_1 \text{ ltr} \quad s_2 \text{ lord}}{s \text{ id} \longleftrightarrow ID[s] \text{ tok}} \quad (3.38)$$

$$\frac{s = s_1 \hat{s}_2 \text{ str} \quad s_1 \text{ dig} \quad s_2 \text{ dgs} \quad s \text{ num} \longleftrightarrow n \text{ nat}}{s \text{ num} \longleftrightarrow NUM[n] \text{ tok}} \quad (3.39)$$

$$\frac{s = s_1 \hat{s}_2 \hat{s}_3 \text{ str} \quad s_1 \text{ qum} \quad s_2 \text{ lord} \quad s_3 \text{ qum}}{s \text{ lit} \longleftrightarrow LIT[s_2] \text{ tok}} \quad (3.40)$$

$$\frac{s = + \cdot \epsilon \text{ str}}{s \text{ spl} \longleftrightarrow ADD \text{ tok}} \quad (3.41)$$

$$\frac{s = * \cdot \epsilon \text{ str}}{s \text{ spl} \longleftrightarrow MUL \text{ tok}} \quad (3.42)$$

$$\frac{s = \wedge \cdot \epsilon \text{ str}}{s \text{ spl} \longleftrightarrow CAT \text{ tok}} \quad (3.43)$$

$$\frac{s = ( \cdot \epsilon \text{ str}}{s \text{ spl} \longleftrightarrow LP \text{ tok}} \quad (3.44)$$

$$\frac{s = ) \cdot \epsilon \text{ str}}{s \text{ spl} \longleftrightarrow RP \text{ tok}} \quad (3.45)$$

$$\frac{s = | \cdot \epsilon \text{ str}}{s \text{ spl} \longleftrightarrow VB \text{ tok}} \quad (3.46)$$

Pravilo 3.38 velja samo, če ne velja nobeno izmed pravil od 3.35 do 3.37. Gledano tehnično ima pravilo 3.38 premise, ki izključijo ključne besede kot možene identifikatorje.

### 3.2.2 Kontekstno neodvisne gramatike

Standardna metoda za definicijo konkretne sintakse jezika je uporaba kontekstno neodvisne gramatike (okr. KNG). Gramatika se sestoji iz treh komponent:

1. žetonov ali *terminalnih simbolov* nad katerimi je gramatika definirana,
2. sintaktičnih razredov ali *neterminalnih simbolov*, ki se razlikujejo od terminalnih simbolov in
3. pravil ali produkcij oblike  $A ::= \alpha$ , kjer je  $A$  neterminalen simbol in je  $\alpha$  niz terminalnih in neterminalnih simbolov.

Vsak sintaktični razred je kolekcija nizov žetonov. Pravila povejo kateri niz spada k katerem sintaktičnem razredu. Pri definiciji gramatike pogosto okrajšamo del produkcij tako, da združimo produkcije z istim sintaktičnim razredom.

$$\begin{array}{l}
 A ::= \alpha_1 \\
 \cdot \\
 \cdot \\
 \cdot \\
 A ::= \alpha_n
 \end{array}
 \tag{3.47}$$

Vse produkcije imajo enak levi del. Če sestavimo zgornje produkcije v eno samo dobimo:

$$A ::= \alpha_1 \mid \dots \mid \alpha_n, \tag{3.48}$$

Takšno pravilo definira množico alternativ za razvoj sintaktičnega razreda  $A$ , ki so ločene s pokončno črto  $\mid$ .

Kontekstno neodvisna gramatika določa hkratno induktivno definicijo množice sintaktičnih razredov. Vsak neterminalni simbol ali sintaktični razred vidimo kot sodbo  $s$   $A$  nad nizi terminalnih simbolov.

Neterminalni simboli imajo *dvojno vlogo*: a) vlogo spremenljivk v produkciji in b) vlogo lastnosti oz. razreda objektov. Poglejmo si zdaj kako prevedemo pravila izražena v KNG v produkcije.

Za prikaz prevajanja pravil KNG v pravila izpeljave uporabimo splošno obliko pravila KNG. Objekti  $s_1, \dots, s_{n+1}$  predstavljajo terminalne simbole, medtem ko predstavljajo  $A_1, \dots, A_n$  sintaktične razrede, ki se naprej razvijajo v terminalne simbole.

$$A ::= s_1 A_1 s_2 \dots s_n A_n s_{n+1} \quad (3.49)$$

Pravilo prevedemo tako, da uporabljamo sintaktične razrede  $A_i$  kot sodbe s katerimi priredimo tip izpeljanim terminalnim simbolom  $s'_i$ . Vsak sintaktični razred (neterminalni simbol) je definiran z induktivno definicijo. Množica takšnih pravil sestavlja induktivno definicijo *statične semantike* jezika.

$$\frac{s'_1 A_1 \dots s'_n A_n}{s_1 s'_1 s_2 \dots s_n s'_n s_{n+1} A}. \quad (3.50)$$

Spet se spomnimo, da je postavljanje nizov skupaj krajši zapis za konkatencijo nizov. Pravilo lahko napišemo takole.

$$\frac{s'_1 A_1 \dots s'_n A_n \quad s = s_1 \hat{\ } s'_1 \hat{\ } s_2 \hat{\ } \dots \hat{\ } s_n \hat{\ } s'_n \hat{\ } s_{n+1}}{s A} \quad (3.51)$$

Formulacija pove, da za  $s$  velja  $A$ , če lahko  $s$  razdelimo kot je opisano zgoraj:  $s'_i A_i$  za  $1 \leq i \leq n$ . Ker konkatencija nizov ni enolično inverzibilna, dekompozicija niza ni unikatna in imamo lahko več različnih načinov po katerih ovrednotimo pravilo.

### 3.2.3 Gramatična struktura

V okviru teorije programskih jezikov se bomo ukvarjali s statično oz. gramatično strukturo jezika—le-to bomo imenovali *statična semantika* jezika. Za specifikacijo statične strukture bomo uporabljali *abstraktno sintakso* jezika in se s tem izognili vsem problemom konkretne sintakse jezika.

S statično semantiko bomo definirali osnovno sintakso jezika ter preverjali statičen izračun tipov izrazov jezika. Z *dinamično semantiko* jezika bomo predstavili različne



$$\frac{s1 \textit{id} \quad s2 \textit{exp} \quad s3 \textit{exp}}{LET \ s1 \ BE \ s2 \ IN \ s3 \ \textit{exp}} \quad (3.61)$$

$$\frac{n \ \textit{nat}}{NUM[n] \ \textit{num}} \quad (3.62)$$

$$\frac{s \ \textit{str}}{LIT[s] \ \textit{lit}} \quad (3.63)$$

$$\frac{s \ \textit{str}}{ID[s] \ \textit{id}} \quad (3.64)$$

Da bi poudarili vlogo konkatencije nizov bi lahko napisali pravilo 3.57 na naslednji način.

$$\frac{s = s1 \ \textit{MUL} \ s2 \ \textit{str} \quad s1 \ \textit{exp} \quad s2 \ \textit{exp}}{s \ \textit{exp}} \quad (3.65)$$

Izraz  $s \ \textit{exp}$  je izpeljiv, če je  $s$  konkatencija niza  $s_1$ , znaka za množenje in  $s_2$ , kjer velja  $s_1 \ \textit{exp}$  in  $s_2 \ \textit{exp}$ .

### 3.2.4 Kontekstno neodvisne gramatike in programski jeziki

Kontekstno neodvisne gramatike, ki so obravnavane v klasični šoli teorije avtomatov, jezikov in izračunljivosti [9, 8], predstavljajo *temeljni formalni jezik* za definicijo konkretne sintakse programskih jezikov.

V okviru teorije jezikov imamo množico specifičnih vrst gramatik, ki omogočajo enostavnejše prepoznavanje stavkov, ki so člani jezika dane gramatike. Na primer, *s-gramatike* so enostavne KNG, ki dovoljujejo za dan neterminalni simbol  $A$  eno samo produkcijo za dan začetek desne strani produkcije. Na ta način ima razčlenjevalni algoritem linearno kompleksnost.

Problemi, ki jih obravnava teorija jezikov so tudi dvoumnost jezikov in odpravljanje dvoumnosti v okviru KNG. Dvoumnosti se izognemo tako, da prepisemo gramatiko v drugo, nedvoumno obliko, ki dovoljuje eno samo drevo izpeljave za dan stavek.

Velikokrat želimo zapisati gramatiko v *normalni obliki*, preko katere lažje obravnavamo in vidimo lastnosti jezikov KNG. Najbolj znani normalni obliki gramatik sta Chomskijeva normalna oblika in Greibachova normalna oblika.

Normalne oblike ne vsebujejo marsikatere nadležne lastnosti gramatik kot so na primer produkcije enote in  $\lambda$ -produkcije, omogočajo pa tudi bolj enostavno teoretično in praktično analizo jezika<sup>1</sup>. Na primer, gramatika zapisana v Chomskijevi normalni obliki omogoča uporabo algoritma CYK za preverjanje članstva stavka v jeziku gramatike.

Kontekstno neodvisne gramatike torej služijo primarno kot orodje za specifikacijo in implementacijo programskih jezikov. Ena izmed zelo široko uporabljenih notacij za specifikacijo programskih jezikov je Backus-Naurjev zapis gramatike ali BNF je v osnovi enaka jezikom kontekstno neodvisnim jezikom.

Kontekstno neodvisni jeziki so bolj podrobno obravnavani v okviru področja *prevajalnikov in tolmačev*. To področje se ukvarja z množico zelo težkih problemov. Glavni problemi s katerimi se srečujemo pri implementaciji prevajalnikov in tolmačev so izločanje dvoumnosti gradnikov jezika, učinkovita implementacija razčlenjevanja in prevajanja jezika v bolj enostavne računalniške jezike.

Na primer, področje prevajalnikov in tolmačev se ukvarja s prevajanjem zelo kompleksnih gramatik, ki so še vedno v razredu KNG. V ta namen je bila razvita vrsta teoretičnih in praktičnih orodij za razčlenjevanje kot so npr. LL in LR gramatike, ki navkljub kompleksnim gradnikom jezika omogočajo razčlenjevanje v linearnem času [1]. Implementacija LL in LR gramatik je realizirana v okviru meta-programskih sistemih YACC in LEX.

Problemi, ki jih srečujemo pri razčlenjevanju in prevajanju jezikov presegajo okvir zapiskov o teoriji programskih jezikov, ki se primarno ukvarja s *semantiko programskih jezikov*. Jeziki so v okviru teorije programskih jezikov predstavljeni z abstraktno sintakso. Pričakuje se, da ni problemov s katerimi se ukvarja področje prevajalnikov programskih jezikov.

---

<sup>1</sup>Bolj podrobne rezultate si lahko ogledate v [9].

### 3.3 Abstraktna sintaksa

Konkretna sintaksa jezika definira linearno predstavitev jezikovnih konstruktov. Program je predstavljen kot niz simbolov.

Osnovni namen konkretne sintakse je izražanje konceptov programskih jezikov na katerih sloni dan jezik. Pod koncepti programskih jezikov so mišljene abstrakcije, ki služijo kot osnovni gradniki za definicijo programa. Primeri pogosto uporabljenih abstrakcij programskih jezikov so funkcije, iteracija, polimorfizem, itn.

Konkretna sintakso lahko torej vidimo kot realizacijo abstrakcij s katerimi so predstavljeni koncepti programskih jezikov. V okviru zasnove konkretne sintakse se ukvarjamo z berljivostjo jezika in primernostjo jezika za avtomatsko obdelavo. Področje, ki se ukvarja z razčlenjevanjem, prevajanjem in interpretacijo jezikov so *prevajalniki in tolmači* [1].

Pri matematični obravnavi jezika nas zanima pomen jezika ter formalne lastnosti jezika in ne toliko površinske predstavitve. Stavke jezika želimo obravnavati čim bolj neodvisno od konkretne sintakse jezika—zanima nas predvsem statična struktura stavkov izražena z abstraktno sintakso.

Abstraktna sintaksa jezika prikaže hierarhično in povezovalno strukturo jezika ter zane-mari linearno notacijo konkretne sintakse. Predstavitev programa v abstraktni sintaksi omogoča natančno predstavitev strukture stavka z uporabo dreves ali grafov.

Razčlenjevanje je proces prevajanja iz konkretne sintakse v abstraktno sintakso. Sestoji se iz analize linearne predstavitve jezika preko gramatike jezika ter prevajanja linearne predstavitve v abstraktno sintakso drevo ali podobno obliko abstraktnega zapisa.

V naslednji sekciji si bomo ogledali abstraktno sintakso jezika  $\mathcal{L}(\text{natstr})$ . Predstavljeno bo tudi pretvarjanje konkretne sintakse v abstraktno sintakso drevo.



### 3.3.1 Abstraktna sintaksna drevesa

Abstraktno sintaksno drevo jezika  $\mathcal{L}(\text{nat str})$  je definirano z naslednjo signaturo:

$$\begin{aligned}
 ar(\text{num}[n]) &= 0 \quad (n \text{ nat}) \\
 ar(\text{str}[s]) &= 0 \quad (s \text{ str}) \\
 ar(\text{id}[s]) &= 0 \quad (s \text{ str}) \\
 ar(\text{plus}) &= 2 \\
 ar(\text{times}) &= 2 \\
 ar(\text{cat}) &= 2 \\
 ar(\text{len}) &= 1 \\
 ar(\text{let}[s]) &= 2
 \end{aligned} \tag{3.66}$$

Vsak identifikator obravnavamo kot operator s števnostjo 0. Konstruktor 'let' je obravnavan kot družina operacij, ki imajo števnost 2 in so indeksirani z identifikatorjem, ki ga povezuje.

Pri specializaciji pravil za abstraktna sintaksna drevesa na predstavljeno signaturo dobimo naslednje induktivne definicije abstraktne sintakse  $\mathcal{L}(\text{nat str})$ :

$$\frac{n \text{ nat}}{\text{num}[n] \text{ ast}} \tag{3.67}$$

$$\frac{s \text{ str}}{\text{str}[s] \text{ ast}} \tag{3.68}$$

$$\frac{s \text{ str}}{\text{id}[s] \text{ ast}} \tag{3.69}$$

$$\frac{a_1 \text{ ast} \quad a_2 \text{ ast}}{\text{plus}(a_1; a_2) \text{ ast}} \tag{3.70}$$

$$\frac{a_1 \text{ ast} \quad a_2 \text{ ast}}{\text{times}(a_1; a_2) \text{ ast}} \tag{3.71}$$

$$\frac{a_1 \text{ ast} \quad a_2 \text{ ast}}{\text{cat}(a_1; a_2) \text{ ast}} \tag{3.72}$$

$$\frac{a \text{ ast}}{\text{len}(a) \text{ ast}} \tag{3.73}$$

$$\frac{s \text{ id} \quad a_1 \text{ ast} \quad a_2 \text{ ast}}{\text{let}[s](a_1; a_2) \text{ ast}} \quad (3.74)$$

Zadnje pravilo je specializacija pravila za 'let' stavek. Zahtevamo, da je prvi argument 'let' stavka identifikator.

### 3.3.2 Razčlenjevanje v abstraktno sintakšno drevo

Proces prevajanja iz konkretne sintakse v abstraktno sintakso se imenuje *razčlenjevanje* (angl. parsing).

Razčlenjevanje bomo obravnavali kot sodbo med konkretno in abstraktno sintakso jezika. Za to sodbo bo veljalo ( $\forall, \exists!$ ) za vse nize in ast.

Razčlenjevanje je torej parcialna funkcija med nizi in ast. Nedefinirana je za nize, ki nimajo pravilne gramatične strukture, preostale (pravilno strukturirane) nize pa enolično prevede v ast.

V nadaljevanju si bomo ogledali sodbe za razčlenjevanje jezika  $\mathcal{L}(\text{nat str})$ .

$$\begin{array}{ll} s \text{ prg} \longleftrightarrow a \text{ ast} & \text{Parse as a program} \\ s \text{ exp} \longleftrightarrow a \text{ ast} & \text{Parse as an expression} \\ s \text{ trm} \longleftrightarrow a \text{ ast} & \text{Parse as a term} \\ s \text{ fct} \longleftrightarrow a \text{ ast} & \text{Parse as a factor} \\ s \text{ num} \longleftrightarrow a \text{ ast} & \text{Parse as a number} \\ s \text{ lit} \longleftrightarrow a \text{ ast} & \text{Parse as a literal} \\ s \text{ id} \longleftrightarrow a \text{ ast} & \text{Parse as an identifier} \end{array} \quad (3.75)$$

Sodbe za razčlenjevanje so definirane induktivno na osnovi sledečih pravil:

$$\frac{n \text{ nat}}{\text{NUM}[n] \text{ num} \longleftrightarrow \text{num}[n] \text{ ast}} \quad (3.76)$$

$$\frac{s \text{ str}}{LIT[s] \text{ lit} \longleftrightarrow \text{str}[s] \text{ ast}} \quad (3.77)$$

$$\frac{s \text{ str}}{ID[s] \text{ id} \longleftrightarrow \text{id}[s] \text{ ast}} \quad (3.78)$$

$$\frac{s \text{ num} \longleftrightarrow a \text{ ast}}{s \text{ fct} \longleftrightarrow a \text{ ast}} \quad (3.79)$$

$$\frac{s \text{ lit} \longleftrightarrow a \text{ ast}}{s \text{ fct} \longleftrightarrow a \text{ ast}} \quad (3.80)$$

$$\frac{s \text{ id} \longleftrightarrow a \text{ ast}}{s \text{ fct} \longleftrightarrow a \text{ ast}} \quad (3.81)$$

$$\frac{s \text{ prg} \longleftrightarrow a \text{ ast}}{LP \text{ s } RP \text{ fct} \longleftrightarrow a \text{ ast}} \quad (3.82)$$

$$\frac{s \text{ fct} \longleftrightarrow a \text{ ast}}{s \text{ trm} \longleftrightarrow a \text{ ast}} \quad (3.83)$$

$$\frac{s_1 \text{ fct} \longleftrightarrow a_1 \text{ ast} \quad s_2 \text{ trm} \longleftrightarrow a_2 \text{ ast}}{s_1 \text{ MUL } s_2 \text{ trm} \longleftrightarrow \text{times}(a_1; a_2) \text{ ast}} \quad (3.84)$$

$$\frac{s \text{ fct} \longleftrightarrow a \text{ ast}}{VB \text{ s } VB \text{ trm} \longleftrightarrow \text{len}(a) \text{ ast}} \quad (3.85)$$

$$\frac{s \text{ trm} \longleftrightarrow a \text{ ast}}{s \text{ exp} \longleftrightarrow a \text{ ast}} \quad (3.86)$$

$$\frac{s_1 \text{ trm} \longleftrightarrow a_1 \text{ ast} \quad s_2 \text{ exp} \longleftrightarrow a_2 \text{ ast}}{s_1 \text{ ADD } s_2 \text{ exp} \longleftrightarrow \text{plus}(a_1; a_2) \text{ ast}} \quad (3.87)$$

$$\frac{s_1 \text{ trm} \longleftrightarrow a_1 \text{ ast} \quad s_2 \text{ exp} \longleftrightarrow a_2 \text{ ast}}{s_1 \text{ CAT } s_2 \text{ exp} \longleftrightarrow \text{cat}(a_1; a_2) \text{ ast}} \quad (3.88)$$

$$\frac{s \text{ exp} \longleftrightarrow a \text{ ast}}{s \text{ prg} \longleftrightarrow a \text{ ast}} \quad (3.89)$$

$$\frac{s_1 \text{ id} \longleftrightarrow \text{id}[s] \text{ ast} \quad s_2 \text{ exp} \longleftrightarrow a_2 \text{ ast} \quad s_3 \text{ prg} \longleftrightarrow a_3 \text{ ast}}{LET \text{ s}_1 \text{ BE } s_2 \text{ IN } s_3 \text{ prg} \longleftrightarrow \text{let}[s](a_2; a_3) \text{ ast}} \quad (3.90)$$

Uspešno razčlenjvanje pogojuje, da se niz žetonov razčleni v skladu s pravili gramatike, ki ni dvoumna, in da je rezultat razčlenjevanja dobro-definirano abstraktno sintaktično drevo.

**Izrek 3.3.1** Če  $s$  prg  $\longleftrightarrow a$  ast, potem  $s$  prg in  $a$  ast. Enako velja za vse ostale sodbe.

Velja celo več, če je niz generiran po pravilih gramatike, potem ga je mogoče razčleniti v ast.

**Izrek 3.3.2** Če  $s$  prg, potem obstaja enoličen  $a$  tako da  $s$  prg  $\longleftrightarrow a$  ast. Enako velja za ostale sodbe, ki opišejo razčlenjevanje: sodbe imajo lastnost  $(\forall\exists!)$  za vse dobro-definirane nize in ast.

Končno, vsak zapis abstraktne sintakse je mogoče preoblikovati v niz simbolov, ki se razčleni v skladu z danim ast.

**Izrek 3.3.3** Če velja  $a$  ast, potem obstaja (ni nujno enolično) niz  $s$  tako, da velja  $s$  prg in  $s$  prg  $\longleftrightarrow a$  ast. Z drugimi besedami, sodba za razčlenjevanje ima lastnost  $(\exists\forall)$ .

## 3.4 Primer enostavnega jezika

V tej sekciji bo predstavljen primer enostavnega jezika  $\mathcal{L}(\text{nat bool})$ , leksikalna struktura, leksikalna analiza ter gramatika  $\mathcal{L}(\text{nat bool})$  v obliki pravil.

### 3.4.1 Jezik $\mathcal{L}(\text{nat bool})$

Ogledali si bomo zelo enostaven jezik, ki vsebuje samo nekatere izmed sintaktičnih oblik:

- boolove konstante *true* in *false*,
- pogojni izraz *if – then – else*,

- numerično konstanto 0,
- aritmetične operacije *pred* in *succ* ter
- pogojni izraz *iszero*, ki preveri al ima izraz vrednost 0.

Gramatika  $\mathcal{L}(\text{nat bool})$  je definirana z naslednjimi izrazi, ki predstavljajo konkretno in abstraktno sintakso. Razlika med obema je minimalna v primeru jezika  $\mathcal{L}(\text{nat bool})$ .

$$\begin{array}{ll}
 \text{Tipi } t & ::= \text{ nat} \quad \text{nat} \\
 & \quad \text{bool} \quad \text{bool} \\
 \text{Izrazi } e & ::= \text{ bool}[t] \quad \text{true} \\
 & \quad \text{bool}[f] \quad \text{false} \\
 & \quad \text{if}(b, e_1, e_2) \quad \text{if } b \text{ then } e_1 \text{ else } e_2 \\
 & \quad \text{nat}[z] \quad 0 \\
 & \quad \text{succ}(e) \quad e + 1 \\
 & \quad \text{pred}(e) \quad e - 1 \\
 & \quad \text{iszero}(e) \quad \text{iszero}(e)
 \end{array} \tag{3.91}$$

Jezik  $\mathcal{L}(\text{nat bool})$  uporablja tipe *nat* in *bool*. Konkretni primerki teh dveh tipov so cela števila izražena v Churchevi obliki in boolove vrednosti *true* in *false*. Edini strukturiran jezikovni gradnik jezika  $\mathcal{L}(\text{nat bool})$  je *if* stavek, ki uporablja funkcijo *iszero*(*e*). Rezultat evaluacije izrazov so vedno enostavne vrednosti: boolove vrednosti ali cela števila. V naslednjem poglavju si bomo ogledali tudi pravila za evaluacijo izrazov jezika  $\mathcal{L}(\text{nat bool})$ .

**Primer 3.4.1** *Poglejmo si nekaj primerov izrazov v predstavljenem jeziku.*

$$\begin{array}{ll}
 & \text{succ(succ(succ(0)))} \\
 \blacktriangleright & 3 \\
 & \text{if false then 0 else 1} \\
 \blacktriangleright & 1 \\
 & \text{iszero(pred(succ(0)))} \\
 \blacktriangleright & \text{true}
 \end{array} \tag{3.92}$$

*Najprej, cela števila predstavimo z večkratnim vgnezdenjem operacije succ ter pred. Drugi primer prikaže uporabo if stavka. Tretji primer uporabi operacijo iszero nad celim številom.*  $\square$

### 3.4.2 Leksikalna analiza

#### Leksikalna struktura $\mathcal{L}(\text{nat bool})$

Znakovna predstavitev je v največjih primerih definirana z regularnimi izrazi. Leksikalna struktura  $\mathcal{L}(\text{nat bool})$  je določena z naslednjimi pravili.

$$\begin{array}{lll}
 \textit{Item} & \textit{itm} & ::= \textit{kwd} \mid \textit{boo} \mid \textit{zer} \mid \textit{spl} \\
 \textit{Keyword} & \textit{kwd} & ::= \textit{i.f.}\epsilon \mid \\
 & & \textit{t.h.e.n.}\epsilon \mid \\
 & & \textit{s.u.c.c.}\epsilon \mid \\
 & & \textit{p.r.e.d.}\epsilon \mid \\
 & & \textit{i.s.z.e.r.o.}\epsilon \mid \\
 \textit{Bool} & \textit{boo} & ::= \textit{true} \mid \textit{false} \\
 \textit{Special} & \textit{spl} & ::= ( \mid ) \\
 \textit{Zero} & \textit{zer} & ::= 0
 \end{array} \tag{3.93}$$

#### Jezik žetonov $\mathcal{L}(\text{nat bool})$

Vhoden niz znakov se pretvarja v niz simbolov (žetonov) z uporabo naslednjih pravil, ki definirajo simbolni jezik razčlenjevalnika oz. statično semantiko jezika.

$$\overline{\textit{ZERO tok}} \tag{3.94}$$

$$\frac{l \textit{ bool}}{\overline{\textit{BOOL}[l] tok}} \tag{3.95}$$

$$\overline{\textit{IF tok}} \tag{3.96}$$

$$\overline{\textit{THEN tok}} \tag{3.97}$$

$$\overline{\textit{ELSE tok}} \tag{3.98}$$

$$\overline{SUCC tok} \quad (3.99)$$

$$\overline{PRED tok} \quad (3.100)$$

$$\overline{ISZERO tok} \quad (3.101)$$

$$\overline{LP tok} \quad (3.102)$$

$$\overline{RP tok} \quad (3.103)$$

### Pretvorba $\mathcal{L}(\text{nat bool})$ v nize žetonov

Naslednja pravila predstavljata leksikalno razčlenjevanje stavkov jezika  $\mathcal{L}(\text{nat bool})$  v niz žetonov.

Pravila na začetku so definirana za usmerjanje pretvorbe izrazov  $\mathcal{L}(\text{nat bool})$  v žetone.

$$\overline{\epsilon inp \longleftrightarrow \epsilon tokstr} \quad (3.104)$$

$$\frac{s = s_1 \hat{\ } s_2 \hat{\ } s_3 str \quad s_1 whs \quad s_2 itm \longleftrightarrow t tok \quad s_3 inp \longleftrightarrow ts tokstr}{s inp \longleftrightarrow t . ts tokstr} \quad (3.105)$$

$$\frac{s kwd \longleftrightarrow t tok}{s itm \longleftrightarrow t tok} \quad (3.106)$$

$$\frac{s boo \longleftrightarrow t tok}{s itm \longleftrightarrow t tok} \quad (3.107)$$

$$\frac{s spl \longleftrightarrow t tok}{s itm \longleftrightarrow t tok} \quad (3.108)$$

$$\frac{s = i . f . \epsilon str}{s kwd \longleftrightarrow IF tok} \quad (3.109)$$

$$\frac{s = t . h . e . n . \epsilon str}{s kwd \longleftrightarrow THEN tok} \quad (3.110)$$

$$\frac{s = e.l.s.e \in str}{s \text{ kw}d \longleftrightarrow ELSE \text{ tok}} \quad (3.111)$$

$$\frac{s = t.r.u.e \in str}{s \text{ boo} \longleftrightarrow BOOL[t] \text{ tok}} \quad (3.112)$$

$$\frac{s = f.a.l.s.e \in str}{s \text{ boo} \longleftrightarrow BOOL[f] \text{ tok}} \quad (3.113)$$

$$\frac{s = 0 \in str}{s \text{ zer} \longleftrightarrow ZERO \text{ tok}} \quad (3.114)$$

$$\frac{s = (. \in str)}{s \text{ spl} \longleftrightarrow LP \text{ tok}} \quad (3.115)$$

$$\frac{s = ) \in str}{s \text{ spl} \longleftrightarrow RP \text{ tok}} \quad (3.116)$$

### 3.4.3 Razčlenjevanje $\mathcal{L}(\text{nat bool})$

V prejšnjem razdelku smo definirali kako se nizi znakov pretvarjajo v nize simbolov. Nize simbolov je naprej potrebno razčleniti v skladu z gramatično strukturo jezika  $\mathcal{L}(\text{nat bool})$ .

Najprej bomo definirali gramatično strukturo jezika v BNF obliki. Zapis bomo pretvorili v pravila s katerim je definirana statična struktura jezika  $\mathcal{L}(\text{nat bool})$ .

#### Gramatična struktura $\mathcal{L}(\text{nat bool})$

Konkretno sintakso  $\mathcal{L}(\text{nat bool})$  lahko definiramo s kontekstno neodvisno gramatiko nad simboli, ki so predstavljeni v 3.93. Spet ima gramatika ima en sam sintaktični razred, *exp*, ki je definiran z naslednjimi pravili:



$$\begin{array}{l}
\text{Expression } exp ::= zer \mid boo \mid \\
\quad \text{if } exp \text{ then } exp \text{ else } exp \mid \\
\quad \text{succ } exp \mid \text{pred } exp \mid \\
\quad \text{iszero } exp \mid LP \ exp \ RP \\
\text{Boolean } boo ::= BOOL[l] \quad (l \text{ bool}) \\
\text{Zero } zer ::= ZERO
\end{array} \tag{3.117}$$

Če prikažemo interpretacijo gramatike z induktivnimi definicijami, dobimo naslednja pravila.

$$\frac{s \ zer}{s \ exp} \tag{3.118}$$

$$\frac{s \ boo}{s \ exp} \tag{3.119}$$

$$\frac{s_1 \ exp \quad s_2 \ exp \quad s_3 \ exp}{\text{if } s_1 \ \text{then } s_2 \ \text{else } s_3} \tag{3.120}$$

$$\frac{s \ exp}{\text{succ } s \ exp} \tag{3.121}$$

$$\frac{s \ exp}{\text{pred } s \ exp} \tag{3.122}$$

$$\frac{s \ exp}{\text{pred } s \ exp} \tag{3.123}$$

$$\frac{s \ exp}{LP \ s \ RP \ exp} \tag{3.124}$$

$$\frac{l \ \text{bool}}{BOOL[l] \ boo} \tag{3.125}$$

$$\frac{}{ZERO \ zer} \tag{3.126}$$

### 3.5 Opombe

Poglavje vsebuje prevode izbranih sekcij iz učbenika *Practical Foundations for Programming Languages* [6] avtorja Roberta Harperja.