

Süntaksanalüüs

Süntaksanalüüs

- **Süntaksanalüüs** kontrollib programmi struktuuri vastavust keele grammatikale:
 - saab sisendina, skanneri poolt genereeritud, lekseemide jada;
 - väljastab programmi esitava (abstraktse) süntaksipuu;
 - süntaktiliste vigade korral, teeb kindlaks nende asukoha;
 - ... teavitab võimalikest vea põhjustest;
 - ... püüab veast toibuda ja jätkata analüüsi (et järgnevaid vigu avastada).
- Süntaksanalüüsi kutsutakse **parsimiseks** (**parsing**) ning vastavat analüsaatorit nimetatakse **parseriks** (**parser**).

Grammatikad

- Keelte süntaksi kirjeldatakse reeglina kontekstivaba grammatika abil.
- **Grammatika** on nelik $G = \langle N, T, P, S \rangle$, kus
 - N on lõplik **mitteterminaalide** tähestik;
 - T on lõplik **terminaalsümbolite** tähestik;
 - $N \cap T = \emptyset$ ja $V = N \cup T$;
 - $P \subset \{ \alpha \rightarrow \beta \mid \alpha \in V^+, \beta \in V^* \}$ on lõplik **produksioonireeglite** hulk;
 - $S \in N$ on **alg sümbol**.
- Grammatika on **kontekstivaba** (**context-free**), kui produktsioonireeglid on kujul $A \rightarrow \alpha$, kus $A \in N$ ja $\alpha \in V^*$.

Grammatikad

- Jada $w \in V^*$ nimetatakse **lausevormiks** (sentential form).
- Lausevorm $v \in V^*$ on **otsetuletatav** (directly derivable) lausevormist $u \in V^*$ (tähistus $u \Longrightarrow v$), kui leiduvad $w_1, w_2, \alpha, \beta \in V^*$ sellised, et $u = w_1\alpha w_2$, $v = w_1\beta w_2$ ja $\alpha \rightarrow \beta \in P$.
- Relatsiooni \Longrightarrow refleksiivset transitiivset sulundit (tähistus \Longrightarrow^*) nimetatakse **derivatsiooniks** (derivation) ehk **tuletuseks**.
- Grammatika $G = \langle N, T, P, S \rangle$ genereerib **keele**

$$L(G) = \{w \in T^* \mid S \Longrightarrow^* w\}$$

- Grammatikad G_1 ja G_2 on **ekvivalentsed**, kui $L(G_1) = L(G_2)$.

Grammatikad

Chomsky hierarhia:

	Produksioonid	Keelte tüüp	Automaat
L_0	$\alpha \rightarrow \beta$	Semi-Thue süsteemid	Turingi masin
L_1	$\alpha A \beta \rightarrow \alpha \gamma \beta$	Kontekstist sõltuvad keeled	Tõkestatud Turingi masin
L_2	$A \rightarrow \alpha$	Kontekstivabad keeled	Magasinmäluga automaat
L_3	$A \rightarrow w, A \rightarrow wB$	Regulaarsed keeled	Lõplik automaat
(L_4)	$A \rightarrow w$	Lõplikud keeled	Tsükliteta lõplik automaat

kus $A, B \in N$, $\alpha, \beta, \gamma \in V^*$ ja $w \in T^*$.

Lemma: Chomsky hierarhia on range; so.:

$$(L_4) \subset L_3 \subset L_2 \subset L_1 \subset L_0$$

Kontekstivabad grammatikad

- Edaspidi käsitleme ainult kontekstivabu grammatikaid.
- Kontekstivabade grammatikate produktsioonireegleid esitatakse tavaliselt **Backus-Naur'i kujul** (BNF).
- Näide: olgu $N = \{\text{Exp}\}$ ja $T = \{+, *, (,), id\}$, siis

$$\begin{array}{l} \text{Exp} \rightarrow \text{Exp} + \text{Exp} \\ | \text{Exp} * \text{Exp} \\ | (\text{Exp}) \\ | id \end{array}$$

esitab produktsioonireeglite hulka

$$P = \{ \text{Exp} \rightarrow \text{Exp} + \text{Exp}, \text{Exp} \rightarrow (\text{Exp}), \\ \text{Exp} \rightarrow \text{Exp} * \text{Exp}, \text{Exp} \rightarrow id \}.$$

Kontekstivabad grammatikad

- Mitteterminaal A on **produktiivne** (productive), kui leidub $w \in T^*$ selline, et $A \Longrightarrow^* w$.
- Mitteterminaal A on **saavutatav** (reachable), kui leiduvad lausevormid $u, v \in V^*$ sellised, et $S \Longrightarrow^* uAv$.
- KV-grammatika $G = \langle N, T, P, S \rangle$ on **taandatud** (reduced), kui tema iga mitteterminaal on produktiivne ja saavutatav.
- **Lemma:** Iga KV-grammatika saab teisendada temaga ekvivalentseks taandatud KV-grammatikaks.

Kontekstivabad grammatikad

- Reeglina saab ühte ja sama lauset tuletada paljudel eri viisidel.
- Kanoonilised derivatsioonid:
 - **vasakderivatsioon** – derivatsiooni igal sammul asendatakse vasakpoolseim mitteterminaal;
 - **parenderivatsioon** – derivatsiooni igal sammul asendatakse parempoolseim mitteterminaal.
- Näide:

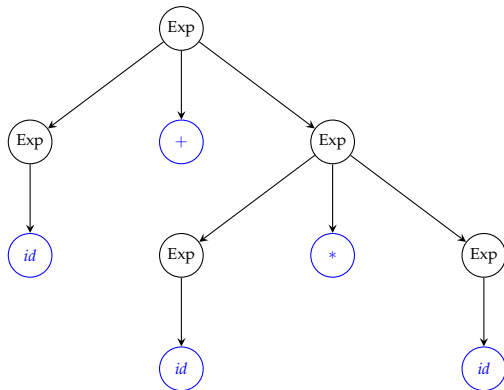
$$\begin{aligned} \text{Exp} &\Rightarrow_{lm} \text{Exp} + \text{Exp} \\ &\Rightarrow_{lm} id + \text{Exp} \\ &\Rightarrow_{lm} id + \text{Exp} * \text{Exp} \\ &\Rightarrow_{lm} id + id * \text{Exp} \\ &\Rightarrow_{lm} id + id * id \end{aligned}$$
$$\begin{aligned} \text{Exp} &\Rightarrow_{rm} \text{Exp} + \text{Exp} \\ &\Rightarrow_{rm} \text{Exp} + \text{Exp} * \text{Exp} \\ &\Rightarrow_{rm} \text{Exp} + \text{Exp} * id \\ &\Rightarrow_{rm} \text{Exp} + id * id \\ &\Rightarrow_{rm} id + id * id \end{aligned}$$

Kontekstivabad grammatikad

- Iga derivatsioon määrab üheselt **süntaksipuu** (*syntax-tree*, *parse-tree*), mis on järjestatud tippudega puu, kus:
 - puu juur on märgendatud algsümboliga S ;
 - vahetipud on märgendatud mitteterminaalidega;
 - lehed on märgendatud terminaalsümboliga või tühisümboliga ϵ ;
 - kui vahetipp on märgendatud mitteterminaaliga A ja tema alampuude (vasakult paremale) t_1, \dots, t_n juured on märgendatud vastavalt A_1, \dots, A_n , siis $A \rightarrow A_1 \dots A_n \in P$.
- Tuletatud lause saadakse lugedes puu lehtede märgendid vasakult paremale.
- Süntaksipuu määrab üheselt kasutatud produktsioonireeglid, kuid mitte nende rakendamise järjekorda.

Kontekstivabad grammatikad

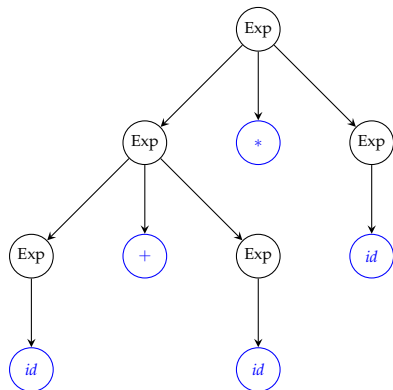
Näide: eelnevalt toodud vasak- ja parenderivatsioonile vastab mõlemal juhul sama süntaksipuu



Kontekstivabad grammatikad

NB! Ühel lausel võib olla mitu erinevat süntaksipuud!

$\text{Exp} \Rightarrow_{lm} \text{Exp} * \text{Exp}$
 $\Rightarrow_{lm} \text{Exp} + \text{Exp} * \text{Exp}$
 $\Rightarrow_{lm} id + \text{Exp} * \text{Exp}$
 $\Rightarrow_{lm} id + id * \text{Exp}$
 $\Rightarrow_{lm} id + id * id$



Kontekstivabad grammatikad

- KV-grammatika on **mitmene** (**ambiguous**), kui ühe lause jaoks leidub mitu süntaksipuud.
- Iga süntaksipuu jaoks leidub täpselt üks vasak- ja üks parenderivatsioon; seega:
 - ühesel lausel on täpselt üks vasak- ja üks parenderivatsioon;
 - mitmesel lausel on vähemalt kaks vasak- ja parenderivatsiooni.
- Lause erinevad süntaksipuud vastavad reeglina lause semantilistelt erinevatele tõlgendusvõimalustele.
- Mitmese grammatika saab teatud juhtudel (aga mitte alati) teisendada ekvivalentseks üheseks grammatikaks.

Kontekstivabad grammatikad

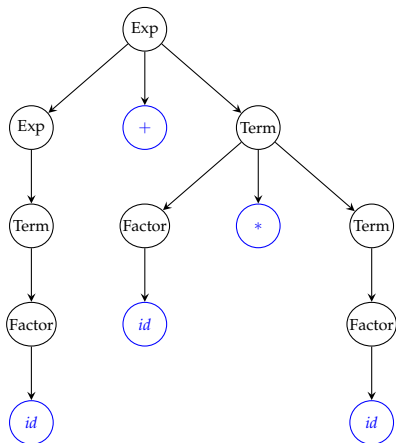
- Mitmesuse eemaldamine – binaarsed operaatorid:
 - iga prioriteeditaseme jaoks toome sisse uue mitteterminaali;
 - vasakassotsiatiivse operaatori korral kasutame vasakrekursiivset; paremassotsiatiivse korral paremrekursiivset produktsioonireeglit;
 - tugevama prioriteediga operaatoritele vastavad reeglid paigutame "sügavamale".
- Näide:

$$\text{Exp} \rightarrow \text{Exp} + \text{Term}$$
$$| \text{Term}$$
$$\text{Term} \rightarrow \text{Factor} * \text{Term}$$
$$| \text{Factor}$$
$$\text{Factor} \rightarrow (\text{Exp})$$
$$| \textit{id}$$

Kontekstivabad grammatikad

Näide:

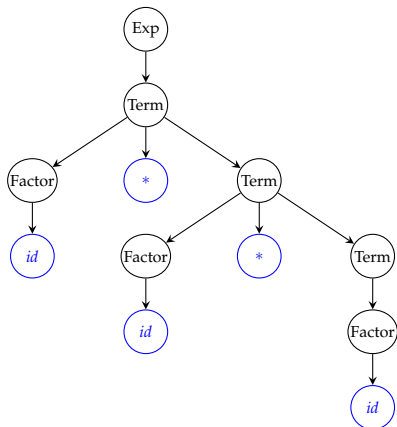
$\text{Exp} \Rightarrow_{lm} \text{Exp} + \text{Term}$
 $\Rightarrow_{lm} \text{Term} + \text{Term}$
 $\Rightarrow_{lm} \text{Factor} + \text{Term}$
 $\Rightarrow_{lm} id + \text{Term}$
 $\Rightarrow_{lm} id + \text{Factor} * \text{Term}$
 $\Rightarrow_{lm} id + id * \text{Term}$
 $\Rightarrow_{lm} id + id * \text{Factor}$
 $\Rightarrow_{lm} id + id * id$



Kontekstivabad grammatikad

Näide:

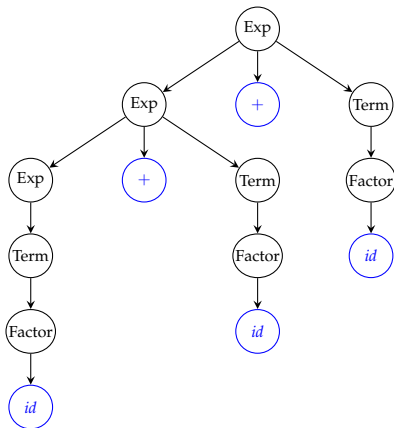
$\text{Exp} \Rightarrow_{lm} \text{Term}$
 $\Rightarrow_{lm} \text{Factor} * \text{Term}$
 $\Rightarrow_{lm} id * \text{Term}$
 $\Rightarrow_{lm} id * \text{Factor} * \text{Term}$
 $\Rightarrow_{lm} id * id * \text{Term}$
 $\Rightarrow_{lm} id * id * \text{Factor}$
 $\Rightarrow_{lm} id * id * id$



Kontekstivabad grammatikad

Näide:

$\text{Exp} \Rightarrow_{lm} \text{Exp} + \text{Term}$
 $\Rightarrow_{lm} \text{Exp} + \text{Term} + \text{Term}$
 $\Rightarrow_{lm} \text{Term} + \text{Term} + \text{Term}$
 $\Rightarrow_{lm} \text{Factor} + \text{Term} + \text{Term}$
 $\Rightarrow_{lm} id + \text{Term} + \text{Term}$
 $\Rightarrow_{lm} id + \text{Factor} + \text{Term}$
 $\Rightarrow_{lm} id + id + \text{Term}$
 $\Rightarrow_{lm} id + id + \text{Factor}$
 $\Rightarrow_{lm} id + id + id$



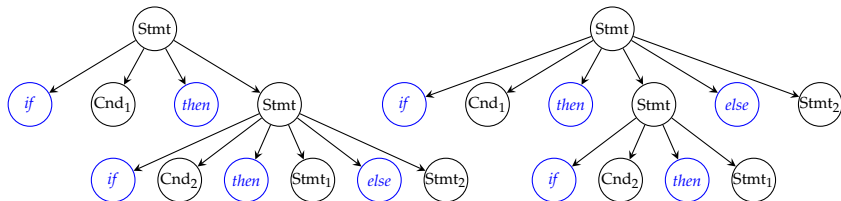
Kontekstivabad grammatikad

- Mitmesuse eemaldamine – tingimuslauseid:

$$\begin{array}{l} \text{Stmt} \rightarrow \textit{if} \text{ Cnd } \textit{then} \text{ Stmt} \\ \quad \quad | \textit{if} \text{ Cnd } \textit{then} \text{ Stmt } \textit{else} \text{ Stmt} \\ \quad \quad | \text{Other} \end{array}$$

- Järgnev lausevorm omab kahte süntaksipuud:

if Cnd₁ *then if* Cnd₂ *then* Stmt₁ *else* Stmt₂

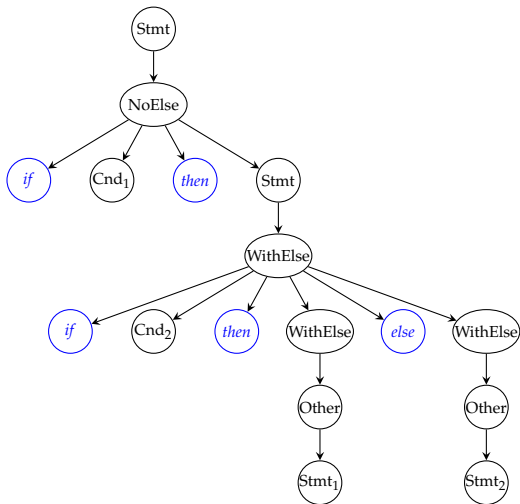


Kontekstivabad grammatikad

Tavaliselt loetakse korrektseks esimest süntaksipuud; so. *else* kuulub alati kõige sisemise võimaliku tingimuslause juurde:

Stmt	→	WithElse
		NoElse
WithElse	→	<i>if</i> Cnd <i>then</i> WithElse <i>else</i> WithElse
		Other
NoElse	→	<i>if</i> Cnd <i>then</i> Stmt
		<i>if</i> Cnd <i>then</i> WithElse <i>else</i> NoElse

Kontekstivabad grammatikad



Parsimistehnikad

Ülalt-alla parsimine (top-down parsing):

- alustab süntaksipuu ehitamist juurest ning kasvatab seda lehtede suunas;
- igal sammul valib produktsioonireegli ning üritab seda sobitada sisendsõnaga;
- mittesobiva reegli korral toimub tagasipöördumine (**backtracking**);
- reeglina annab vasakpoolse derivatsiooni.

Alt-üles parsimine (bottom-up parsing)

- alustab süntaksipuu ehitamist lehtedest ning kasvatab seda juure suunas;
- rakendab sobivaid reegleid paremalt vasakule kuni jõuab algsümbolini;
- reeglina annab parempoolse derivatsiooni.

Ülalt-alla parsimine

Ülalt-alla parsimise üldine algoritm:

- konstrueerime juurtipu, märgendame selle algsümboliga ja jätkame puu konstrueerimist lehteded suunas sügavuti vasakult paremale;
- kui vaatluse all olev tipp on mitteterминаал A , siis valime mingi produktsioonireegli kujul $A \rightarrow \alpha$, konstrueerime reegli paremale poolele vastavad tipud, ning jätkame vasakpoolseima alamtipuga;
- kui tipp on terminaalsümbol mis ei sobi sisensümboliga, siis teostame tagasipöördumise produktsioonireegli valikuni mis selle tipu konstrueeris ja jätkame sealt uuesti valides teise reegli;
- kui tipp on terminaalsümbol mis sobib sisensümboliga, siis jätkame järgmise vasakpoolseima läbivaatamata tipuga.

Ülalt-alla parsimine

Näide:

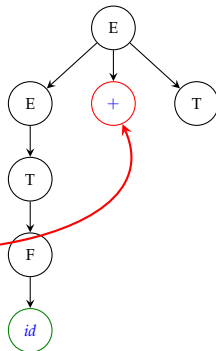
E	→	E + T	• <i>id</i> – <i>num</i> * <i>id</i>
		E – T	E.1 • <i>id</i> – <i>num</i> * <i>id</i>
		T	E.3 • <i>id</i> – <i>num</i> * <i>id</i>
T	→	T * F	T.3 • <i>id</i> – <i>num</i> * <i>id</i>
		T / F	F.2 • <i>id</i> – <i>num</i> * <i>id</i>
		F	<i>id</i> • – <i>num</i> * <i>id</i>
F	→	(E)	<i>id</i> – • <i>num</i> * <i>id</i>
		<i>id</i>	T.1 <i>id</i> – • <i>num</i> * <i>id</i>
		<i>num</i>	T.3 <i>id</i> – • <i>num</i> * <i>id</i>
			F.3 <i>id</i> – • <i>num</i> * <i>id</i>
			<i>id</i> – <i>num</i> • * <i>id</i>
			<i>id</i> – <i>num</i> * • <i>id</i>
			F.2 <i>id</i> – <i>num</i> * • <i>id</i>
			<i>id</i> – <i>num</i> * <i>id</i> •

Ülalt-alla parsimine

Näide:

$E \rightarrow E + T$
| $E - T$
| T
 $T \rightarrow T * F$
| T / F
| F
 $F \rightarrow (E)$
| id
| num

• $id - num * id$
E.1 • $id - num * id$
E.3 • $id - num * id$
T.3 • $id - num * id$
F.2 • $id - num * id$
 $id \bullet - num * id$
 $id - \bullet num * id$
T.1 $id - \bullet num * id$
T.3 $id - \bullet num * id$
F.3 $id - \bullet num * id$
 $id - num \bullet * id$
 $id - num * \bullet id$
F.2 $id - num * \bullet id$
 $id - num * id \bullet$

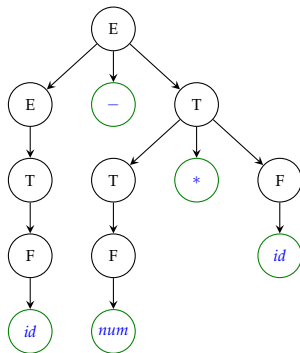


Ülalt-alla parsimine

Näide:

$E \rightarrow E + T$
| $E - T$
| T
 $T \rightarrow T * F$
| T / F
| F
 $F \rightarrow (E)$
| id
| num

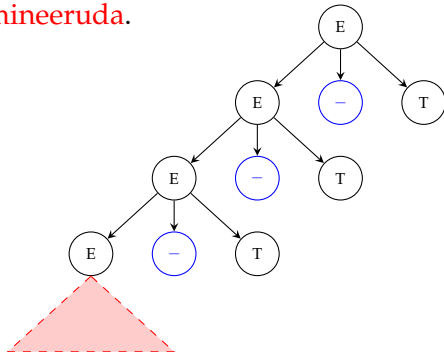
• $id - num * id$
E.2 • $id - num * id$
E.3 • $id - num * id$
T.3 • $id - num * id$
F.2 • $id - num * id$
 $id \bullet - num * id$
 $id - \bullet num * id$
T.1 $id - \bullet num * id$
T.3 $id - \bullet num * id$
F.3 $id - \bullet num * id$
 $id - num \bullet * id$
 $id - num * \bullet id$
F.2 $id - num * \bullet id$
 $id - num * id \bullet$



Ülalt-alla parsimine

- Parsimise efektiivsus sõltub tugevasti õige produktsioonireegli valikust.
- Vale reegli valimine põhjustab hilisema tagasipöördumise vajaduse.
- Vasakrekursiivsete reeglite olemasolul võib ülalt-alla parsimine **mittetermineeruda**.

- $id - num * id$
- E.2 ● $id - num * id$
- E.2 ● $id - num * id$
- E.2 ● $id - num * id$
- ...



Vasakrekursioon

- Grammatika on **vasakrekursiivne**, kui leidub mitteterminal $A \in N$ selline, et

$$A \Longrightarrow^+ A\alpha,$$

kus $\alpha \in V^*$.

- Vasakrekursioon on **vahetu**, kui leidub reegel kujul $A \rightarrow A\alpha$.
- Vastasel korral on vasakrekursioon **kaudne**.

Vasakrekursiooni elimineerimine

Vahetu vasakrekursiooni eemaldamine:

- Toome sisse uue mitteterminali ning asendame vasakrekursiooni paremrekursiooniga

$$A \rightarrow A \alpha \mid \beta \quad \longrightarrow \quad \begin{array}{l} A \rightarrow \beta A' \\ A' \rightarrow \alpha A' \mid \varepsilon \end{array}$$

- Üldjuhul

$$A \rightarrow A \alpha_1 \mid A \alpha_2 \mid \dots \mid \beta_1 \mid \beta_2 \mid \dots$$
$$\longrightarrow \quad \begin{array}{l} A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \\ A' \rightarrow \alpha_2 A' \mid \alpha_2 A' \mid \dots \mid \varepsilon \end{array}$$

Vasakrekursiooni elimineerimine

Näide:

$$\begin{array}{l} E \rightarrow E + T \\ \quad | \quad E - T \\ \quad | \quad T \end{array}$$



$$\begin{array}{l} E \rightarrow T E' \\ E' \rightarrow + T E' \\ \quad | \quad - T E' \\ \quad | \quad \varepsilon \end{array}$$

$$\begin{array}{l} T \rightarrow T * F \\ \quad | \quad T / F \\ \quad | \quad F \end{array}$$



$$\begin{array}{l} T \rightarrow F T' \\ T' \rightarrow * F T' \\ \quad | \quad / F T' \\ \quad | \quad \varepsilon \end{array}$$

- Uus grammatika genereerib sama keele, kuid on natuke vähem intuiitivne.
- Mõlemates grammatikates on operaatorid vasakassotsiatiivsed.

Vasakrekursiooni elimineerimine

Kaudse vasakrekursiooni eemaldamine:

- Näide:

$$\begin{aligned}A_1 &\rightarrow A_2 \alpha \mid \beta \\A_2 &\rightarrow A_1 \gamma \mid A_2 \delta\end{aligned}$$

- Teisendame kaudse vasakrekursiooni vahetuks.
- Asendame mitteterminali A_2 reeglite paremates pooltes kõik A_1 esinemised tema definitsiooniga:

$$\begin{aligned}A_1 &\rightarrow A_2 \alpha \mid \beta \\A_2 &\rightarrow A_2 \alpha \gamma \mid \beta \gamma \mid A_2 \delta\end{aligned}$$

- Eemaldame vahetu vasakrekursiooni:

$$\begin{aligned}A_1 &\rightarrow A_2 \alpha \mid \beta \\A_2 &\rightarrow \beta \gamma A'_2 \\A'_2 &\rightarrow \alpha \gamma A'_2 \mid \delta A'_2 \mid \varepsilon\end{aligned}$$

Vasakrekursiooni elimineerimine

Üldine algoritm kaudse vasakrekursiooni eemaldamiseks:

Järjestame mitteterminalid mingisse järjekorda A_1, \dots, A_n

for $i \leftarrow 1$ **to** n

for $j \leftarrow 1$ **to** $i - 1$

Asendame iga produktsiooni kujul $A_i \rightarrow A_j \alpha$

produktsioonidega $A_i \rightarrow \beta_1 \alpha \mid \beta_2 \alpha \mid \dots \mid \beta_k \alpha$,

kus $A_j \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_k$ on kõik A_j produktsioonid

Eemaldame vahetu vasakrekursiooni mitteterminali A_i
produktsioonidest

NB! Eeldab, et algses grammatikas pole ei ε -produktsioone, ega tsükleid (so. $A_i \Longrightarrow^+ A_i$).

Vasakrekursiooni elimineerimine

Näide:

$$\begin{array}{l} A \rightarrow C \alpha \mid A \alpha \\ B \rightarrow A \beta \mid \gamma \\ C \rightarrow B \delta \mid \varphi \end{array}$$



$$\begin{array}{l} A \rightarrow C \alpha A' \\ B \rightarrow A \beta \mid \gamma \\ C \rightarrow B \delta \mid \varphi \\ A' \rightarrow \alpha A' \mid \varepsilon \end{array}$$



$$\begin{array}{l} A \rightarrow C \alpha A' \\ B \rightarrow C \alpha A' \beta \mid \gamma \\ C \rightarrow B \delta \mid \varphi \\ A' \rightarrow \alpha A' \mid \varepsilon \end{array}$$



$$\begin{array}{l} A \rightarrow C \alpha A' \\ B \rightarrow C \alpha A' \beta \mid \gamma \\ C \rightarrow C \alpha A' \beta \delta \mid \gamma \delta \mid \varphi \\ A' \rightarrow \alpha A' \mid \varepsilon \end{array}$$



$$\begin{array}{l} A \rightarrow C \alpha A' \\ B \rightarrow C \alpha A' \beta \mid \gamma \\ C \rightarrow \alpha A' \beta \delta C' \\ A' \rightarrow \alpha A' \mid \varepsilon \\ C' \rightarrow \gamma \delta C' \mid \varphi C' \mid \varepsilon \end{array}$$



$$\begin{array}{l} A \rightarrow C \alpha A' \\ C \rightarrow \alpha A' \beta \delta C' \\ A' \rightarrow \alpha A' \mid \varepsilon \\ C' \rightarrow \gamma \delta C' \mid \varphi C' \mid \varepsilon \end{array}$$

Ennustav parsimine

- Vale produktsioonireegli valimine põhjustab tagasipöördumise.
- Tihti on võimalik valitava reegli õigsuse üle otsustada mingi arvu sisendsümbolite ettevaatamise teel.
- Üldjuhul on tarvis ette vaadata piiramata arv sümboleid.
 - Näit.: Cocke-Younger-Kasami või Earley algoritmid.
- **Ennustav parsimine** (**predictive parsing**) on ülalt-alla parsimine, kus alati on võimalik ige reegel valida nii, et pole vaja tagasipöörduda.
 - Grammatika peab olema selline, et järgmine sisendsümbol (või mingi fikseeritud arv sümboleid) määrab unikaalselt ära valitava reegli.

Ennustav parsimine

- Iga lausevormi $\alpha \in (N \cup T)^*$ jaoks defineerime hulga:

$$\begin{aligned} \text{first}(\alpha) &= \{a \in T \mid \alpha \Longrightarrow^* a \beta\} \\ &\cup \{\varepsilon \mid \alpha \Longrightarrow^* \varepsilon\} \end{aligned}$$

kus $\beta \in (N \cup T)^*$.

- Iga mitteterminali $A \in N$ jaoks defineerime hulga:

$$\begin{aligned} \text{follow}(A) &= \{a \in T \mid S \Longrightarrow^* \alpha A a \beta\} \\ &\cup \{\$ \mid S \Longrightarrow^* \alpha A\} \end{aligned}$$

kus $\alpha, \beta \in (N \cup T)^*$ ja $\$$ on spetsiaalne sisendilõpu marker.

Ennustav parsimine

Näide:

$$\begin{aligned} S &\rightarrow ABC \\ A &\rightarrow aA \mid \varepsilon \end{aligned}$$

$$\begin{aligned} B &\rightarrow b \mid \varepsilon \\ C &\rightarrow c \mid d \end{aligned}$$

$$\begin{aligned} \mathit{first}(C) &= \{c, d\} \\ \mathit{first}(B) &= \{b, \varepsilon\} \\ \mathit{first}(A) &= \{a, \varepsilon\} \\ \mathit{first}(S) &= \mathit{first}(ABC) \\ &= (\mathit{first}(A) \setminus \{\varepsilon\}) \\ &\quad \cup (\mathit{first}(B) \setminus \{\varepsilon\}) \\ &\quad \cup \mathit{first}(C) \\ &= \{a, b, c, d\} \end{aligned}$$

$$\begin{aligned} \mathit{follow}(C) &= \{\$ \} \\ \mathit{follow}(B) &= \mathit{first}(C) \\ &= \{c, d\} \\ \mathit{follow}(A) &= (\mathit{first}(B) \setminus \{\varepsilon\}) \\ &\quad \cup \mathit{first}(C) \\ &= \{b, c, d\} \end{aligned}$$

Ennustav parsimine

- Kui grammatikas on reeglid $A \rightarrow \alpha$ ja $A \rightarrow \beta$ sellised, et $first(\alpha) \cap first(\beta) = \emptyset$, siis on esimese sisendsümboli põhjal võimalik otsustada kumba reeglit valida.
- **NB!** Kehtib ainult juhul, kui $\varepsilon \notin \{first(\alpha) \cup first(\beta)\}$.
- Vastasel korral tuleb vaadelda ka hulka $follow(A)$.
- Defineerime iga reegli $A \rightarrow \alpha$ jaoks hulga:

$$first^+(\alpha) = \begin{cases} (first(\alpha) \setminus \{\varepsilon\}) \cup follow(A), & \text{kui } \varepsilon \in first(\alpha) \\ first(\alpha), & \text{kui } \varepsilon \notin first(\alpha) \end{cases}$$

- Grammatika on **LL(1)**, kui iga (paarikaupa erineva) produktsioonireegli $A \rightarrow \alpha$ ja $A \rightarrow \beta$ korral

$$first^+(\alpha) \cap first^+(\beta) = \emptyset$$

Ennustav parsimine

- **NB!** LL(1) grammatika ei saa olla ei vasakrekursiivne ega mitmene!
- Näide:

$$S \rightarrow S a \mid \beta$$

- Kui $\beta \neq \varepsilon$
 - Siis $first(\beta) \subseteq first(S) = first(S a)$
 - Seega $first^+(S a) \cap first^+(\beta) \neq \emptyset$
- Kui $\beta = \varepsilon$
 - Siis $a \in first(S a)$ ja $a \in follow(S) = \{a, \$\}$
 - Seega $first^+(S a) \cap first^+(\beta) \neq \emptyset$

Ennustav parsimine

- Paljusid grammatikaid on võimalik teisendada LL(1) kujule kasutades:
 - vasakrekursiooni elimineerimist;
 - vasakfaktoriseerimist;
 - halvimal juhul üldistame natuke grammatikat (ning kontrollime eemaldatud kitsendusi pärast parsimist).
- Vasakfaktoriseerimine (left factoring) asendab ühise prefiksiga reeglid uute reeglitega, kus ühine prefiks on ainult ühes parempooles.
- Näide:

$$\begin{array}{l} A \rightarrow BaCD \\ \quad | BaCE \end{array} \quad \longrightarrow \quad \begin{array}{l} A \rightarrow BaCZ \\ Z \rightarrow D \\ \quad | E \end{array}$$

Ennustav parsimine

Vasakfaktoriseerimise algoritm:

- 1 Iga mitteterminali $A \in N$ korral leiame pikima prefixi α , mis esineb kahes või rohkemas A produktsioonireegli parempooles.
- 2 Kui $\alpha \neq \varepsilon$, siis asendame kõik A produktsioonid

$$A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \dots \mid \alpha \beta_n \mid \gamma$$

produktsioonidega

$$\begin{aligned} A &\rightarrow \alpha Z \mid \gamma \\ Z &\rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n \end{aligned}$$

kus $Z \in N$ on uus mitteterminaal.

- 3 Kordame protsessi, kuni ükski parempool ei oma ühist prefiksit.

Ennustav parsimine

- **Rekursiivselt laskuv** (**recursive descent**) parsimine on ennustava parsimise meetod, kus:
 - iga mitteterminali kohta realiseeritakse üks protseduur, mis tunneb ära sellele mitteterminalile vastavad lausevormid;
 - iga protseduur valib sõltuvalt sisendist reegli ning kutsub (rekursiivselt) välja reegli paremas pooles asuvatele mitteterminalidele vastavad protseduurid.
- Rekursiivselt laskuv parsimine on levinuim meetod (lihtsate keelte) parserite käsitsi kirjutamiseks.

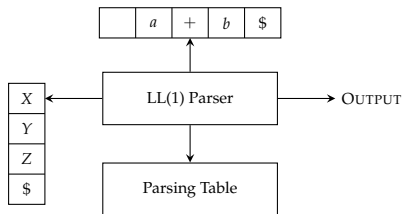
Ennustav parsimine

Näide: olgu A produktsioonireeglid $A \rightarrow aB \mid bCA \mid DE$

```
parseA() {  
  if token = a then {  
    token := nextWord(); parseB();  
  } else if token = b then {  
    token := nextWord(); parseC(); parseA();  
  } else {  
    parseD(); parseE();  
  }  
}
```


Ennustav parsimine

- Automaatselt genereeritavate LL(1) parserite korral kasutatakse tavaliselt **tabeljuhitavat LL(1) parsimist**:
 - moodustatakse maatriks M , kus read ja veerud on vastavalt indekseeritud mitteterminaalide ja terminalidega;
 - tabeli pesades on produktsioonireeglid, mida antud mitteterminaali ja sisendsümboli korral valida.
- Tabeljuhitava LL(1) parseri struktuur:



Ennustav parsimine

LL(1) parsimise algoritm:

```
push($); push(S);  
token := nextWord();  
while stack  $\neq$  empty do {  
  A := pop();  
  if  $A \in N$  then {  
    if  $M[A, token] = B_1 \dots B_n$  then {  
      push( $B_n$ ); ... ; push( $B_1$ );  
    } else error();  
  } else if  $A = token$  then {  
    token := nextWord(); pop();  
  } else error();  
}
```

Shift-reduce parsimine

Shift-reduce parsimine on üldine meetod alt-üles süntaksanalüüsiks:

- puu konstrueerimine toimub lehtedest juure suunas eesmärgiga "taandada" sisendstring alsümboliks;
- parsimise ajal on korraga terve mets puid, mis vastavad erinevatele, juba äratuntud, alamstringidele;
- kaks baasaktsiooni:
 - **shift** loeb uue sisendsümboli;
 - **reduce** taandab mingi produktsioonireegli parempoollele vastava (juba loetud/taandatud sisendsümbolite ja mitteterminalide) jada reegli vasakpoolle olevaks mitteterminaliks;
- konstruktsioon vastab parenderivatsioonile.

Shift-reduce parsimone

Näide:

$S \rightarrow aABe$

$A \rightarrow bcA \mid c$

$B \rightarrow d$

shift

shift

shift

shift

reduce $A \rightarrow c$

reduce $A \rightarrow bcA$

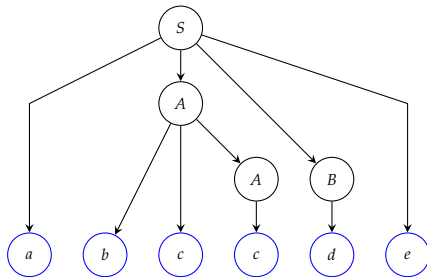
shift

reduce $B \rightarrow d$

shift

reduce $S \rightarrow aABe$

Input String:



$S \xRightarrow{rm} aABe \xRightarrow{rm} aAde$
 $\xRightarrow{rm} abcAde \xRightarrow{rm} abccde$

Shift-reduce parsimine

- Lausevormi nimetatakse **paremlausevormiks** (**right-sentential form**), kui ta on tuletatav parenderivatsiooni abil.
- Paremlausevormi γ **pide** (**handle**) on alamstring β , selline et $S \xRightarrow{*}_{rm} \delta A w \xRightarrow{rm} \delta \beta w = \gamma$, kus $\beta, \gamma, \delta \in V^*$ ja $w \in T^*$.
- Näide: olgu antud grammatika

$$\begin{aligned} S &\rightarrow a A B e \\ A &\rightarrow b c A \mid c \\ B &\rightarrow d \end{aligned}$$

ja parenderivatsioon

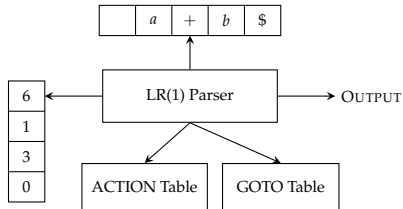
$$S \xRightarrow{rm} a A B e \xRightarrow{rm} a A d e \xRightarrow{rm} a b c A d e$$

Paremlausevormi $abcAde$ pide on bcA .

- Ühese grammatika korral on parenderivatsioon, ja seega ka pidemed, unikaalsed.

Shift-reduce parsimine

LR(1) parseri struktuur:



Shift-reduce parsimime

LR(1) parseri skelett:

```
push(Invalid); push( $s_0$ ); found := false;
token := nextWord();
while found  $\neq$  true do {
     $s$  := top();
    if ACTION[ $s$ , token] = reduce( $A \rightarrow \beta$ ) then
        pop( $2 * |\beta|$ );
     $s$  := top(); push( $A$ ); push(GOTO[ $s$ ,  $A$ ]);
    else if ACTION[ $s$ , token] = shift( $s_i$ ) then
        push(token); push( $s_i$ );
        token := nextWord();
    else if ACTION[ $s$ , token] = accept & token = $ then
        found := true;
    else report error;
}
report success;
```

Shift-reduce parsimine

- **LR(k)-element** (**LR(k)-item**) on paar $[A \rightarrow \alpha \cdot \beta, w]$, kus $A \rightarrow \alpha\beta$ on kontekstivaba grammatika produktsioonireegel ning $w \in T^*$ on sõna pikkusega $|w| \leq k$.
- Element $[A \rightarrow \cdot \beta\gamma, w]$ tähendab, et seni nähtud sisend on kooskõlas reegli $A \rightarrow \beta\gamma$ rakendamisega koheselt peale magasinini tipus olevat sümbolit.
- Element $[A \rightarrow \beta \cdot \gamma, w]$ tähendab, et seni nähtud sisend on kooskõlas reegli $A \rightarrow \beta\gamma$ rakendamisega ja parser on juba ära tundnud β .
- Element $[A \rightarrow \beta\gamma \cdot, w]$ tähendab, et parser on ära tundnud $\beta\gamma$ ja ettevaatamisstring w on kooskõlas selle redutseerimisega mitteterminaliks A .

Shift-reduce parsimine

- Näide: olgu antud grammatika

$$\begin{aligned} S &\rightarrow aAc \\ A &\rightarrow Ab \mid \varepsilon \end{aligned}$$

Tema LR(k)-elemendid ettevaatamisstringi w jaoks on:

$$\begin{array}{ll} [S \rightarrow \cdot aAc, w] & [A \rightarrow \cdot Ab, w] \\ [S \rightarrow a \cdot Ac, w] & [A \rightarrow A \cdot b, w] \\ [S \rightarrow aA \cdot c, w] & [A \rightarrow Ab \cdot, w] \\ [S \rightarrow aAc \cdot, w] & [A \rightarrow \cdot, w] \end{array}$$

- **NB!** Iga kontekstivaba grammatika LR(k)-elementide hulk on lõplik.

Shift-reduce parsimine

LR(1)-elementide sulundi leidmine

- Funktsioon $Closure(\mathcal{S})$ lisab LR(1)-elementide hulka \mathcal{S} need elemendid mis järelduvad seal juba olevatest.
 - Elemendist $[A \rightarrow \beta \cdot B\delta, a]$ järeldub iga reegli $B \rightarrow \tau$ ja iga sümboli $x \in first(\delta)$ jaoks element $[B \rightarrow \cdot\tau, x]$.

- Algoritm:

```
Closure( $\mathcal{S}$ ) {  
  while ( $\mathcal{S}$  is still changing) do {  
     $\forall [A \rightarrow \beta \cdot B\delta, a] \in \mathcal{S}$   
     $\forall B \rightarrow \tau \in P$   
     $\forall b \in first(\delta a)$   
    if  $[B \rightarrow \cdot\tau, b] \notin \mathcal{S}$  then  
       $\mathcal{S} := \mathcal{S} \cup \{[B \rightarrow \cdot\tau, b]\}$   
  }  
}
```

- Algoritmi termineerumine on garanteeritud, kuna elementide hulk on lõplik.

Shift-reduce parsimine

- Funktsioon $Goto(s, X)$, kus s on LR(1) parseri olek (elementide hulk) ja $X \in V$, leiab oleku kuhu parser, olles olekus s , satub pärast sümboli X äratundmist:

$$Goto(s, X) = Closure(\{[A \rightarrow \beta X \cdot \delta, a] \mid [A \rightarrow \beta \cdot X \delta, a] \in s\})$$

- Kanoonlise hulga konstrueerimine:

$s_0 := Closure(\{[S' \rightarrow \cdot S, \$]\});$

$\mathcal{S} := \{s_0\}; \Delta := \emptyset; k := 1;$

while (\mathcal{S} is still changing) **do** {

$\forall s_j \in \mathcal{S}, x \in V$

$s_k := Goto(s_j, x);$

$\Delta := \Delta \cup \{(s_j, x) \mapsto s_k\};$

if $s_k \notin \mathcal{S}$ **then**

$\mathcal{S} := \mathcal{S} \cup \{s_k\}; k := k + 1;$

}

Shift-reduce parsimine

Näide:

$$\begin{aligned} S &\rightarrow E \\ E &\rightarrow T - E \mid T \\ T &\rightarrow F * T \mid F \\ F &\rightarrow n \end{aligned}$$

Symbol	<i>first</i>
S	$\{n\}$
E	$\{n\}$
T	$\{n\}$
F	$\{n\}$
$-$	$\{-\}$
$*$	$\{*\}$
n	$\{n\}$

Algolek:

$$\begin{aligned} s_0 &= \text{Closure}(\{[S \rightarrow \cdot E, \$]\}) \\ &= \{ [S \rightarrow \cdot E, \$], [E \rightarrow \cdot T - E, \$], [E \rightarrow \cdot T, \$], \\ &\quad [T \rightarrow \cdot F * T, \$], [T \rightarrow \cdot F * T, -], [T \rightarrow \cdot F, \$], \\ &\quad [T \rightarrow \cdot F, -], [F \rightarrow \cdot n, \$], [F \rightarrow \cdot n, -], [F \rightarrow \cdot n, *] \} \end{aligned}$$

Shift-reduce parsimine

1. iteratsioon:

$$\begin{aligned} s_1 &= \text{Goto}(s_0, E) \\ &= \{ [S \rightarrow E\cdot, \$] \} \\ s_2 &= \text{Goto}(s_0, T) \\ &= \{ [E \rightarrow T\cdot -E, \$], [E \rightarrow T\cdot, \$] \} \\ s_3 &= \text{Goto}(s_0, F) \\ &= \{ [T \rightarrow F\cdot *T, \$], [T \rightarrow F\cdot *T, -], \\ &\quad [T \rightarrow F\cdot, \$], [T \rightarrow F\cdot, -] \} \\ s_4 &= \text{Goto}(s_0, n) \\ &= \{ [F \rightarrow n\cdot, \$], [F \rightarrow n\cdot, -], [F \rightarrow n\cdot, *] \} \end{aligned}$$

Shift-reduce parsimine

2. iteratsioon:

$$\begin{aligned} s_5 &= \text{Goto}(s_2, -) \\ &= \{ [E \rightarrow T - \cdot E, \$], [E \rightarrow \cdot T - E, \$], [E \rightarrow \cdot T, \$], \\ &\quad [T \rightarrow \cdot F * T, -], [T \rightarrow \cdot F, -], \\ &\quad [T \rightarrow \cdot F * T, \$], [T \rightarrow \cdot F, \$], \\ &\quad [F \rightarrow \cdot n, *], [F \rightarrow \cdot n, -], [F \rightarrow \cdot n, \$] \} \\ s_6 &= \text{Goto}(s_3, *) \\ &= \{ [T \rightarrow F * \cdot T, \$], [T \rightarrow F * \cdot T, -], \\ &\quad [T \rightarrow \cdot F * T, \$], [T \rightarrow \cdot F * T, -], \\ &\quad [T \rightarrow \cdot R, \$], [T \rightarrow \cdot F, -], \\ &\quad [F \rightarrow \cdot n, \$], [F \rightarrow \cdot n, -], [F \rightarrow \cdot n, *] \} \end{aligned}$$

3. iteratsioon:

$$\begin{aligned} s_7 &= \text{Goto}(s_5, E) \\ &= \{ [E \rightarrow T - E \cdot, \$] \} \\ s_8 &= \text{Goto}(s_6, T) \\ &= \{ [T \rightarrow F * T \cdot, \$], [T \rightarrow F * T \cdot, -] \} \end{aligned}$$

Shift-reduce parsimine

Üleminekurelatsioon Δ

State	<i>E</i>	<i>T</i>	<i>F</i>	-	*	<i>n</i>
0	1	2	3			4
1						
2				5		
3					6	
4						
5	7	2	3			4
6		8	3			4
7						
8						

Shift-reduce parsimime

LR(1)-tabelite genereerimine:

$\forall s_x \in \mathcal{S}$

$\forall i \in s_x$

if $i = [A \rightarrow \alpha \cdot a\beta, b], \Delta(s_x, a) = s_k, a \in T$ **then**

$ACTION[x, a] := \mathbf{shift}(k);$

else if $i = [S' \rightarrow S \cdot, \$]$ **then**

$ACTION[x, a] := \mathbf{accept};$

else if $i = [A \rightarrow \beta \cdot, a]$ **then**

$ACTION[x, a] := \mathbf{reduce}(A \rightarrow \beta);$

$\forall A \in N$

if $\Delta(s_x, A) = s_k$ **then**

$GOTO[x, A] := k;$

Shift-reduce parsimine

LR(1)-tabelid näitegrammatikale:

	ACTION				GOTO		
	<i>n</i>	–	*	\$	<i>E</i>	<i>T</i>	<i>F</i>
0	shift(4)				1	2	3
1				accept			
2		shift(5)		reduce(3)			
3		reduce(5)	shift(6)	reduce(5)			
4		reduce(6)	reduce(6)	reduce(6)			
5	shift(4)				7	2	3
6	shift(4)					8	3
7				reduce(2)			
8		reduce(4)		reduce(4)			