

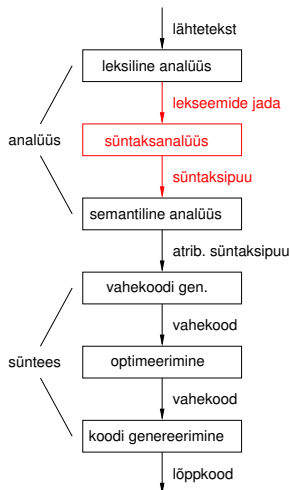
Süntaksanalüüs

Sissejuhatus

Kontekstivabad grammatikad

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 **algsümbol**.
- 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 $u \Rightarrow v$), 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 \Rightarrow refleksiivset transitiivset sulundit (tähistus \Rightarrow^*) nimetatakse **derivatsiooniks** (derivation) ehk **tuletuseks**.
- Grammatika $G = \langle N, T, P, S \rangle$ genereerib **keele**

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

- Grammatikad G_1 ja G_2 on **ekvivalentseted**, 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	Tõkestatud Turingi masin
L_2	$A \rightarrow \alpha$	Kontekstivabad	Magasinmäluga automaat
L_3	$A \rightarrow w, A \rightarrow wB$	Regulaarsed	Lõplik automaat
(L_4)	$A \rightarrow w$	Lõplikud	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

- 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:

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

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

Kontekstivabad grammatikad

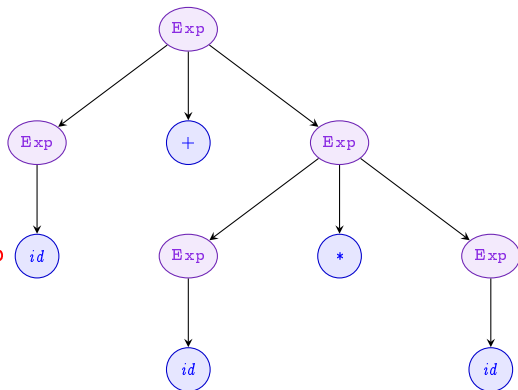
- Iga derivatsioon määrab üheselt **süntaksipuu** (**syntax-tree**, **parse-tree**), mis on järjestatud tippudega puu, kus:
 - puu juur on märgendatud algsümboliga S ;
 - vahetipud on märgendatud mitteterminaalidega;
 - lehed on märgendatud terminaalsümboliga või tühisümboliga ϵ ;
 - kui vahetipp on 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 produktsiooni-reeglid, kuid mitte nende rakendamise järjekorda.

Kontekstivabad grammatikad

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

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

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



Kontekstivabad grammatikad

NB! Ühel lausel võib olla mitu erinevat süntaksipuud!

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

