

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 *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 $u \in V^*$ on *otsetuletatav* (directly derivable) lausevormist $v \in V^*$ (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 $N = \{ Exp \}$ ja $T = \{ +, *, (,), id \}$, siis

$$\begin{array}{l}
 Exp ::= Exp + Exp \\
 \quad | Exp * Exp \\
 \quad | (Exp) \\
 \quad | id
 \end{array}$$

esitab produktsioonireeglite hulka

$$\begin{aligned}
 P = \{ & Exp \rightarrow Exp + Exp, \quad Exp \rightarrow (Exp), \\
 & Exp \rightarrow Exp * Exp, \quad Exp \rightarrow id \}.
 \end{aligned}$$

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;
 - *paremderivatsioon* — derivatsiooni igal sammul asendatakse parempoolseim mitteterminaal.
- Näide:

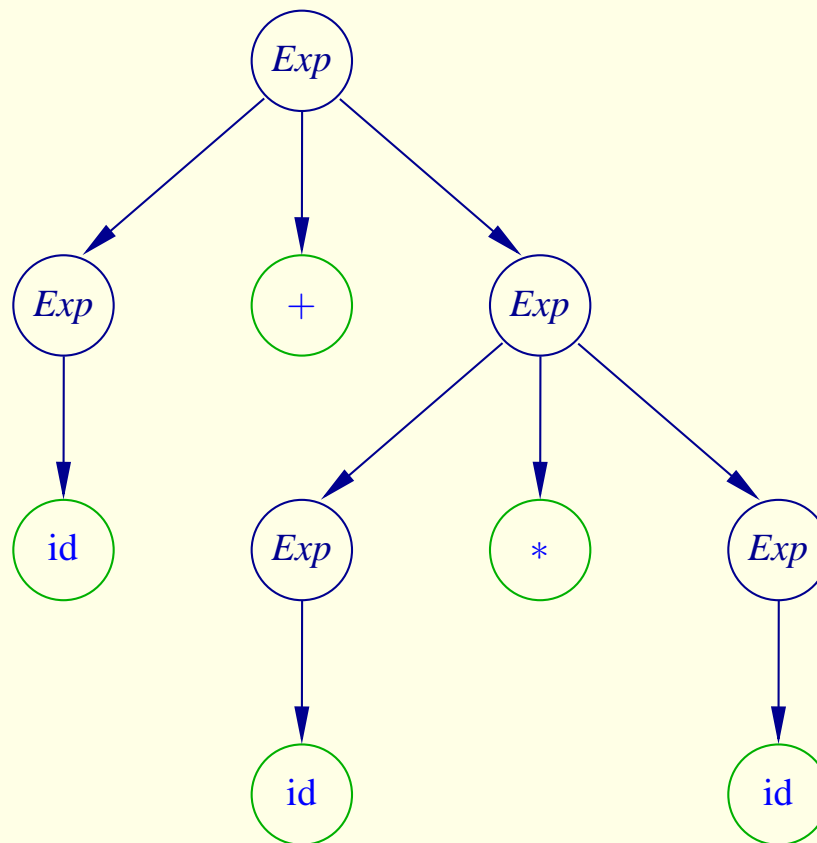
$Exp \implies_{lm} Exp + Exp$	$Exp \implies_{rm} Exp + Exp$
$\implies_{lm} id + Exp$	$\implies_{rm} Exp + Exp * Exp$
$\implies_{lm} id + Exp * Exp$	$\implies_{rm} Exp + Exp * id$
$\implies_{lm} id + id * Exp$	$\implies_{rm} Exp + id * id$
$\implies_{lm} id + id * id$	$\implies_{rm} id + id * id$

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

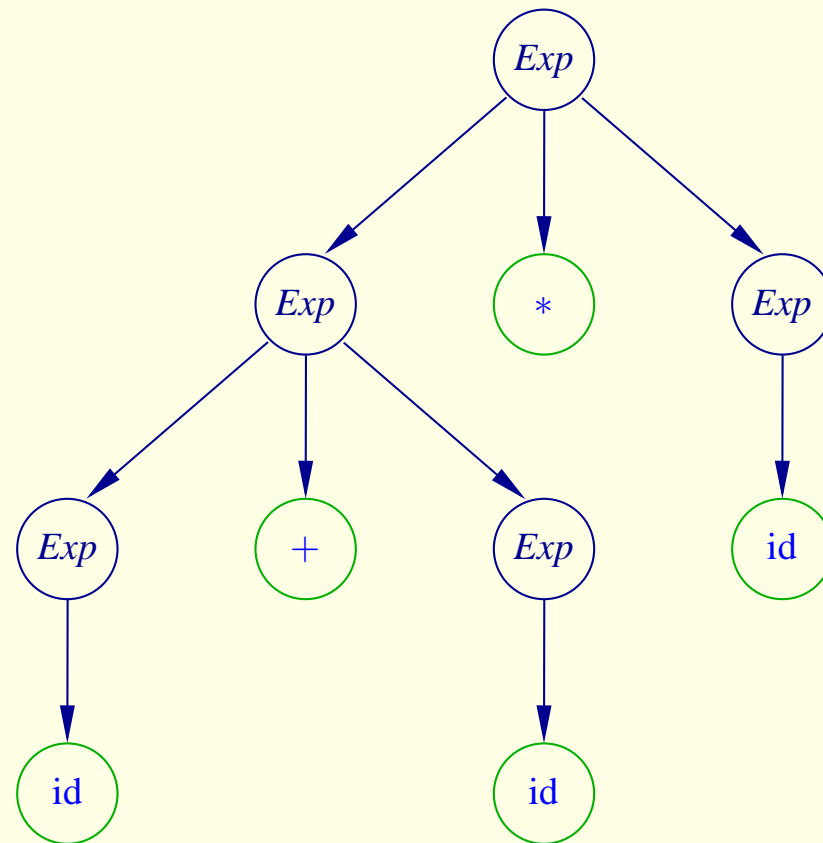
- Näide: eelnevalt toodud vasak- ja parenderivatsioonile vastab mõlemal juhul sama süntaksipuu



Kontekstivabad grammatikad

- **NB!** Ühel lausel $w \in L(G)$ võib olla mitu erinevat süntaksipuud!
- Näide:

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



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 semantiliselt erinevatele tõlgendusvõimalustele.
- Mitmese grammatika saab teatud juhtudel (aga mitte alati) teisendada ekvivalentseks üheseks grammatikaks.

Kontekstivabad grammatikad

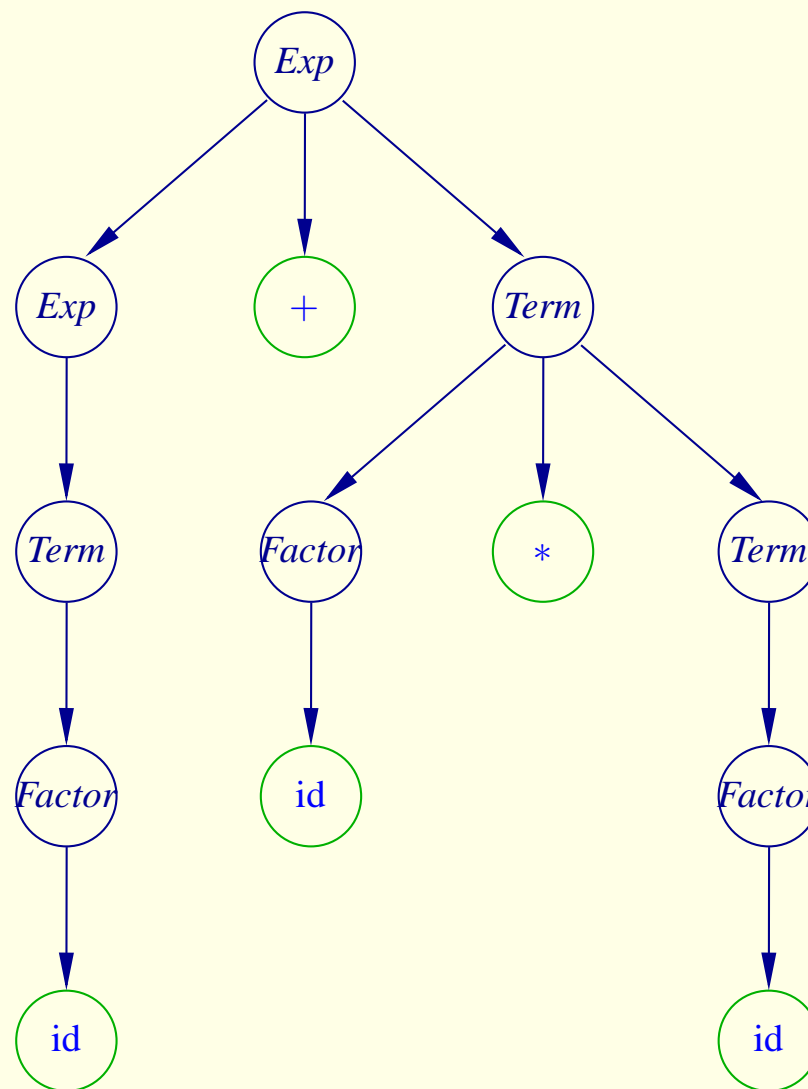
- Mitmesuse eemaldamine — binaarsed operaatorid:
 - iga prioriteeditaseme jaoks toome sisse uue mitteterminaali;
 - vasakassotsiatiivse operaatori korral kasutame vasakrekursiivset; paremassotsiatiivse korral paremrekursiivset produktsioonireeglit;
 - tugevama prioriteediga operaatoritele vastavad reeglid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

$$\begin{array}{lcl}
 \textit{Exp} & ::= & \textit{Exp} + \textit{Term} \\
 & & | \textit{Term} \\
 \textit{Term} & ::= & \textit{Factor} * \textit{Term} \\
 & & | \textit{Factor} \\
 \textit{Factor} & ::= & (\textit{Exp}) \\
 & & | \textit{id}
 \end{array}$$

Kontekstivabad grammatikad

Näide:

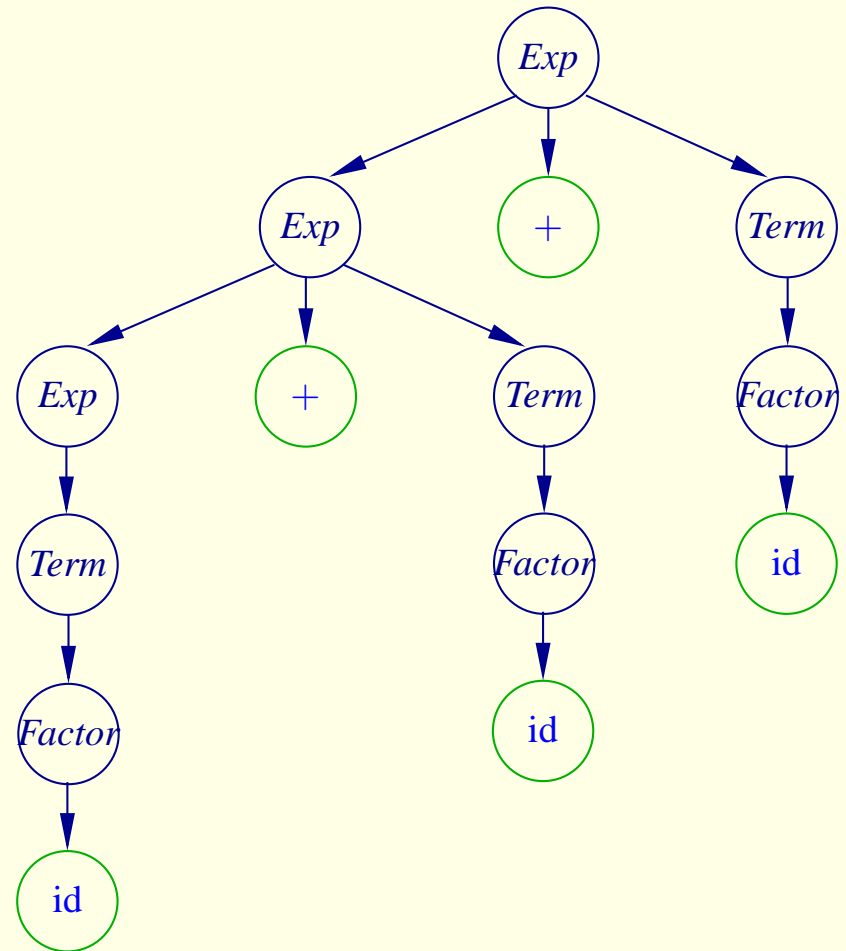
$Exp \Rightarrow_{lm} Exp + Term$
 $\Rightarrow_{lm} Term + Term$
 $\Rightarrow_{lm} Factor + Term$
 $\Rightarrow_{lm} id + Term$
 $\Rightarrow_{lm} id + Factor * Term$
 $\Rightarrow_{lm} id + id * Term$
 $\Rightarrow_{lm} id + id * Factor$
 $\Rightarrow_{lm} id + id * id$



Kontekstivabad grammatikad

Näide:

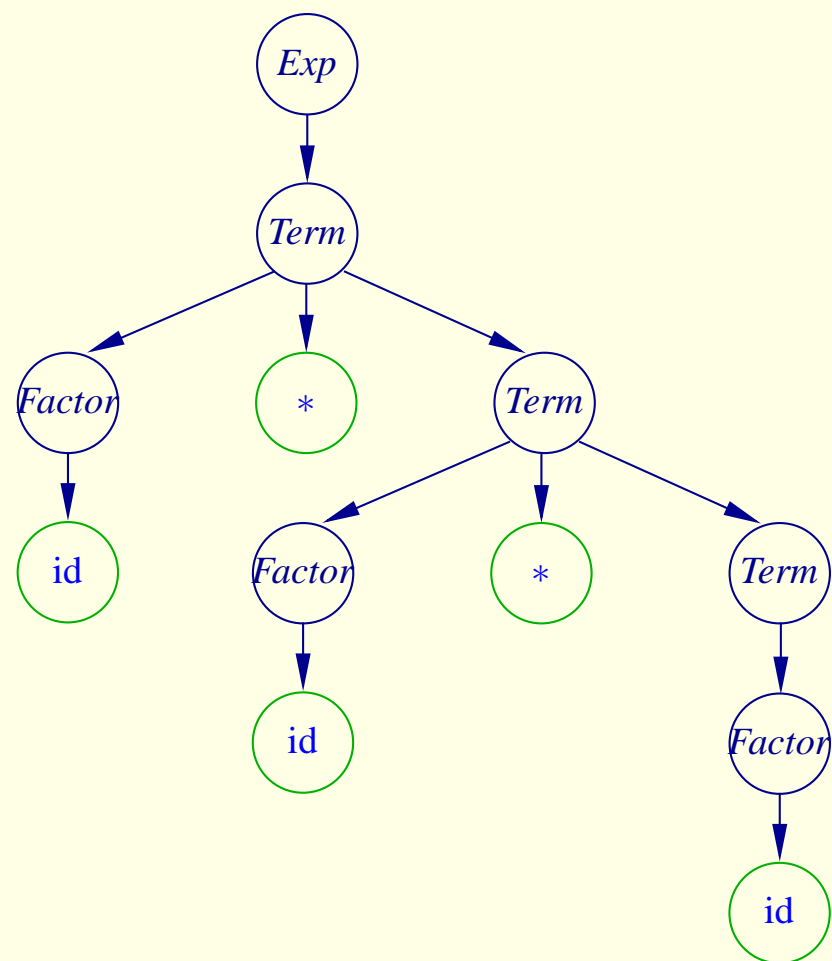
- $Exp \Rightarrow_{lm} Exp + Term$
- $\Rightarrow_{lm} Exp + Term + Term$
- $\Rightarrow_{lm} Term + Term + Term$
- $\Rightarrow_{lm} Factor + Term + Term$
- $\Rightarrow_{lm} id + Term + Term$
- $\Rightarrow_{lm} id + Factor + Term$
- $\Rightarrow_{lm} id + id + Term$
- $\Rightarrow_{lm} id + id +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



Kontekstivabad grammatikad

Näide:

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



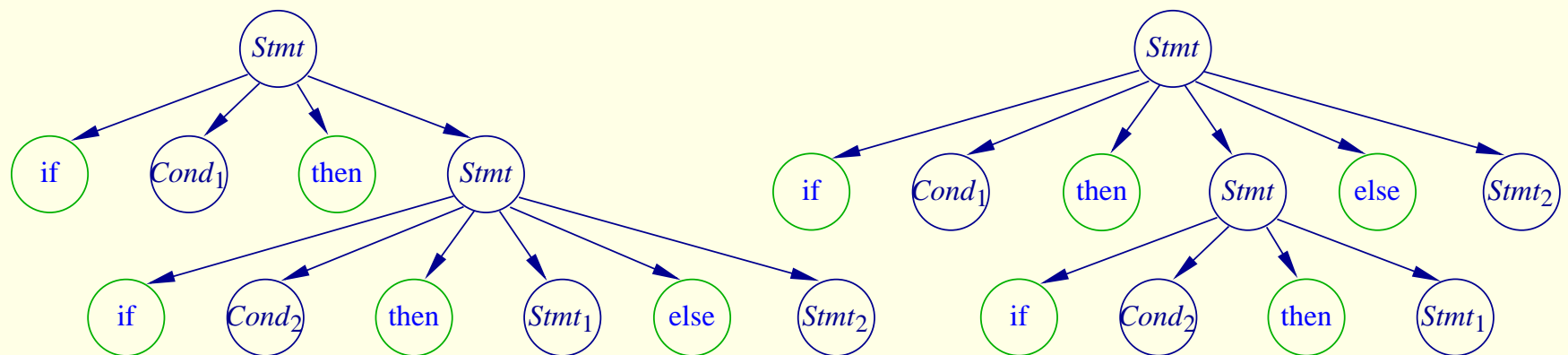
Kontekstivabad grammatikad

- Mitmesuse eemaldamine — tingimuslaused:

$$\begin{array}{l}
 Stmt ::= \text{if } Cond \text{ then } Stmt \\
 \quad | \text{if } Cond \text{ then } Stmt \text{ else } Stmt \\
 \quad | Other
 \end{array}$$

- Järgnev lausevorm omab kahte süntaksipuud:

if $Cond_1$ *then* *if* $Cond_2$ *then* $Stmt_1$ *else* $Stmt_2$



Kontekstivabad grammatikad

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

$$\begin{array}{lcl} \textit{Stmt} & ::= & \textit{WithElse} \\ & & | \textit{NoElse} \\ \textit{WithElse} & ::= & \textit{if Cond then WithElse else WithElse} \\ & & | \textit{Other} \\ \textit{NoElse} & ::= & \textit{if Cond then Stmt} \\ & & | \textit{if Cond then WithElse else NoElse} \end{array}$$

Kontekstivabad grammatikad

