

Skannerite elust

Vesal Vojdani
(TÜ Arvutiteaduse Instituut)

- **Regulaaravaldised** üle (lõpliku) tähestiku Σ

$$E ::= \varepsilon \mid a \mid (E E) \mid (E \mid E) \mid E^*$$

kus $a \in \Sigma$.

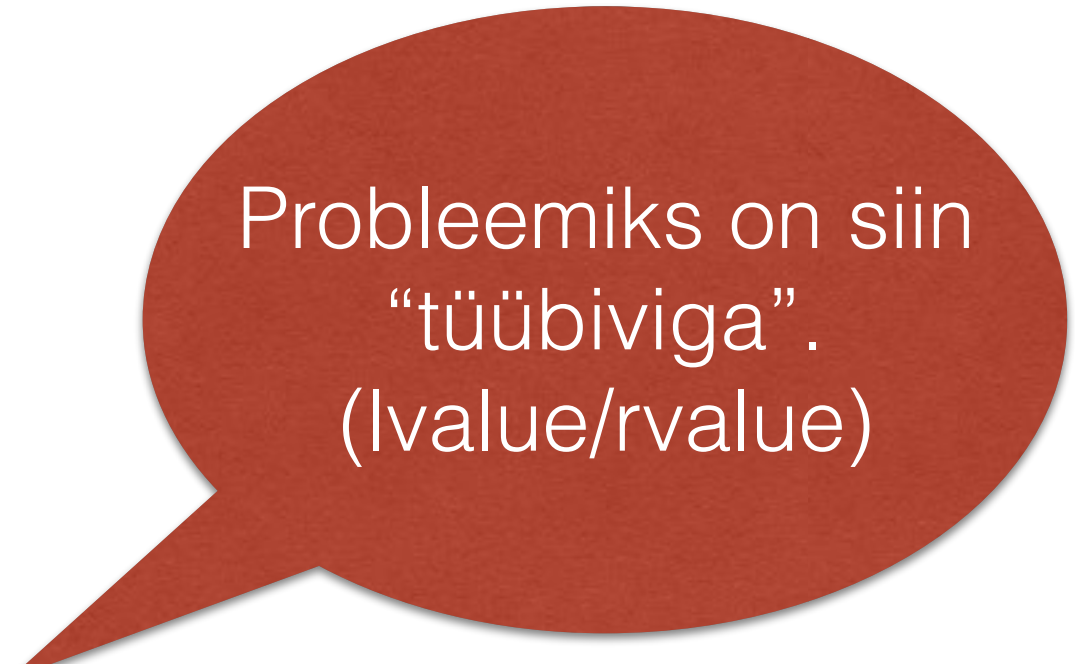
- Regulaaravaldis E defineerib **keele** $L(E) \subseteq \Sigma^*$

$$S \in L(E)$$

Regulaaravaldis defineerib keele, aga leksimine...

Leksimine!

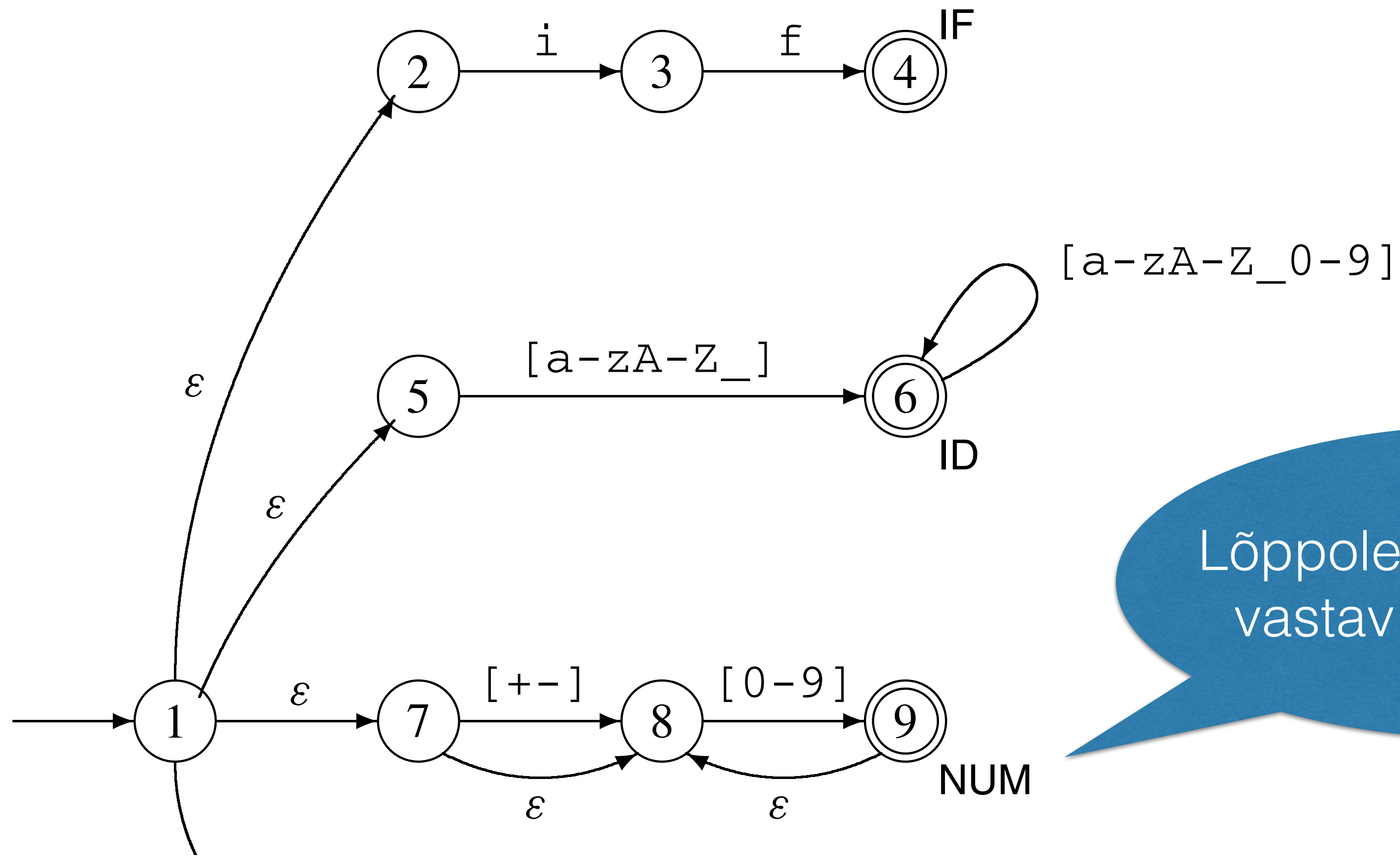
- **x+++++y**
<ID:x>, <Op:++>, <Op:++>, <Op:+>, <ID:y>
- **x+++ ++y**
<ID:x>, <Op:++>, <Op:+>, <Op:++>, <ID:y>



Probleemiks on siin
"tüübiviga".
(lvalue/rvalue)

Leksiline spetsifikatsioon

- Keyword: 'if' | ...
- Op: '++' | '+' | ...
- Identifier: [a-zA-Z_][a-zA-Z_0-9]*
- ...



Lõppolekutel on meeles vastav lekseemiklass

Kombineeritud keel

$$E = E_1 \mid E_2 \mid \dots \mid E_n$$

Kuidas kasutada

- Sisendiks on $x_1 \dots x_m$.
- Otsime sellist i , et alamsõne $w = x_1 \dots x_i$ kuuluks keelde $L(E)$.
- Kui leidub, siis peab olema alamkeel $L(E_j)$, kuhu sõne kuulub.
(Automaadi lõppolekus oli kirjas: E oli ju $E_1 \mid E_2 \mid \dots \mid E_n$.)
- Eemaldame sisendist sõne w ja jätkame kuni sisend on tühi.

Maximal Munch!

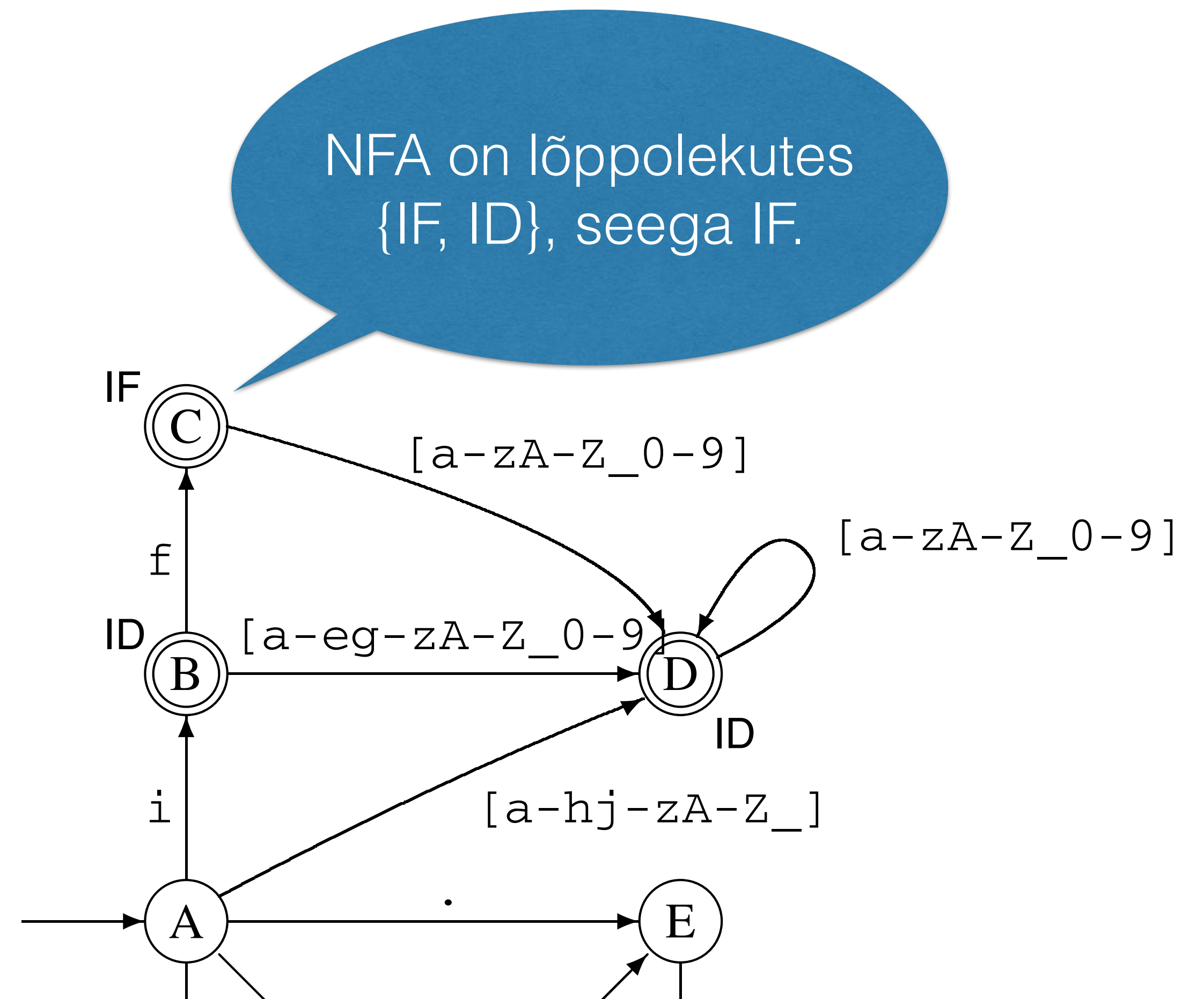
- Mis juhtub, kui sobivad kaks alamsõne??
- Võib ju juhtuda, et $x_1 \dots x_i$ ja $x_1 \dots x_j$ mõlemad kuuluvad keelde $L(E)$.
- Näiteks ‘++’ ja ‘+’ on mõlemad Java operaatorid.
- Loomulik konventsioon on valida **pikim alamsõne**, mis sobitub!

Milline lekseemiklass?

- Kas “**if**” on muutuja nimi või keyword?
- Prioriteedijärjekord: leksilise spetsifikatsiooni järjekord on oluline.
- NB! Kõigepealt valime pikim alamsõne ja alles siis valime lekseemiklass. (Seega, “**ifo**” on muutuja.)

Implementatsioonist

- Determiniseerime: jätame lõppolekutel meelde kõrgeima prioriteediga lekseem.
- Minimeerimine: lõppolekud on eristatavad, kui nad annavad erinevad lekseeme: {A,E}, {B,D}, {C}.



Pikim pigistus

- Automaat tuleb jooksutada nii kaugele kui saab, kasvõi sisendi lõpuni... (näiteks "+++" enam ei ole keeles)
- Hoida alles viimane lõppolek ja tagastada sellele vastav lekseem.
- Sisendis peab ka tagasi minema viimase sobitumise kohale
- Näite tagastame "++" ja jätkame sõnega "+".

Quiz!

Olgu järgmine leksiline spetsifikatsioon:

T1: a

T2: a^*b

Millise sõne puhul on leksimine kõige aeglasem?

Ja kui aeglane see on (keerukus, kui n on sisendi pikkus)?

1. a

2. $aaaaaaaaaa$

3. $aaaaaaaaaab$

Käsitis tehes (kodutöö)

```
StringBuilder sb = new StringBuilder();  
while (Character.isLetter(peek())) {  
    sb.append(peek());  
    consume();  
}  
String identOrIf = sb.toString();  
  
if (identOrIf.equals("if")) {  
    return new Token(IF);  
} else {  
    return new Token(IDENT, identOrIf);  
}
```

See muidugi ei ole enam puhas automaat

Tähtis on siin aru saada, mida lekser peab tegema (**maximal munch**), aga selleks me ei pea päris lekserit simuleerima.