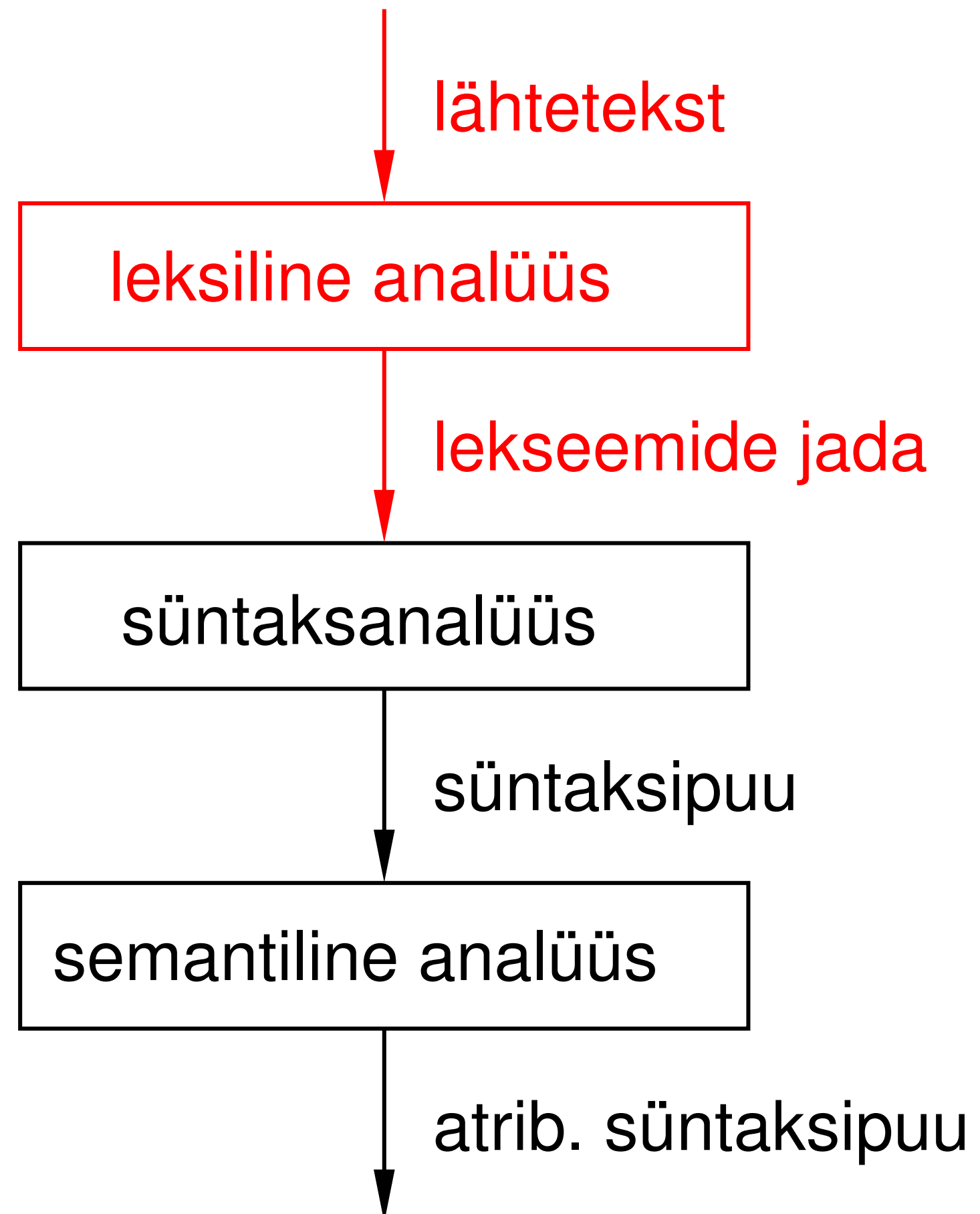


Automaadid, keeled, translaatorid

Leksiline analüüs
Regulaaravaldised

Leksiline analüüs



`x := y + 60`

↓

`muutuja(x),
omistamine,
muutuja(y),
liitmine,
arv(60)`

Lekser ja skanner

- Leksiline analüsaator (**lekser**) tükeldab sisendit ja loob lekseeme.

`kala := 42` \longrightarrow `kala | := | 42` \longrightarrow muutuja ("kala"),
omistamine,
arv(42)

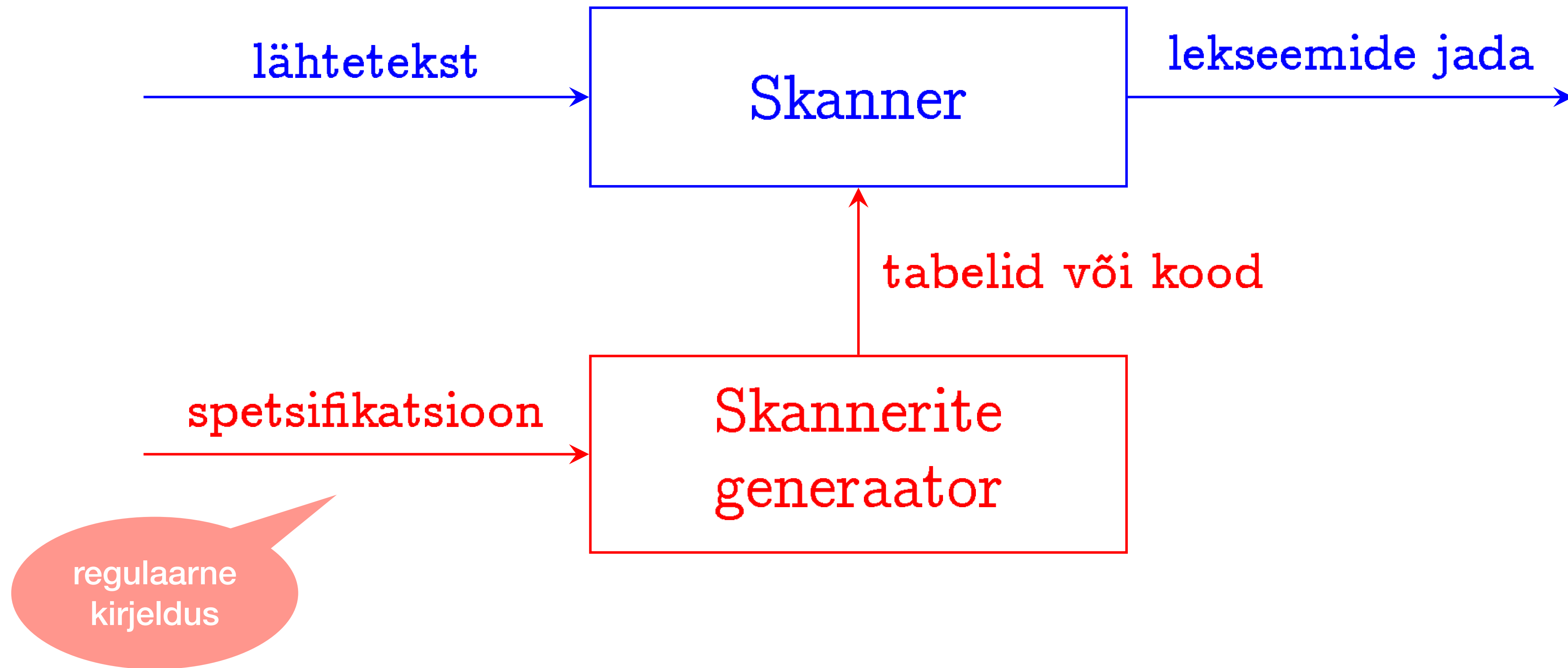
- Sõnade piiride leidmist nimetatakse ka **skannimiseks**.
- Lisaks on lekseemide loomisel viisakas konstante **väärtustada**.
Näiteks: Kui „4“ ja „2“ on kõrvuti, siis väärtuseks on $4*10+2 = 42$.
- Rangelt võttes võib neid samme eristada: **lekser** = **skanner** + **väärtustaja**.

Tihti aga kasutatakse sõnu
lekser ja skanner lihtsalt
sünonüümidenä.

Käsitsi kodeeritud skanner

- Paljude keelte lekserid on käsitsi kirjutatud (sealhulgas OpenJDK).
- Nagu saame näha ([5. kodutöös](#)), siis on see üsna raske.
- Meie kasutame parema meelega lekseri generaatorit, mis loob lekseri automaatselt **regulaaravaldiste** põhjal.

Skannerite generaator



Näiteid regulaarsetest kirjeldustest

Identifikaatorid:

Letter $\rightarrow [a - z A - Z]$
Digit $\rightarrow [0 - 9]$
Identifier $\rightarrow \text{Letter} (\text{Letter} \mid \text{Digit})^*$

Arvkonstandid:

Sign $\rightarrow (+ \mid -)?$
Integer $\rightarrow 0 \mid \text{Sign} [1 - 9] \text{Digit}^*$
Decimal $\rightarrow \text{Integer} . \text{Digit}^+$
Real $\rightarrow (\text{Integer} \mid \text{Decimal}) \text{E Integer}$

Formaalne keel

- Vajame mugavat vahendit, millega defineerida sõnade hulkasid.
- Olgu meil lõplik hulk tähti Σ (tähestik).
Näiteks: eesti tähestik või hulk $\{0,1\}$.
- Kui võtta null või enam tähte hulgast Σ ja need kõrvuti paigutada, siis moodustub sõna! (Näiteks *kala* ja *1011001*).
- Σ^* tähistatakse kõikide sõnade hulka üle tähestiku Σ .
- Formaalne keel on suvaline Σ^* alamhulk!
Regulaaravaldiste abil saabki keelt defineerida (siis on ta regulaarne keel).

Regulaaravaldised

Avaldis	Nimi	Kirjeldus	Näide
a	Literaals	Üksiktäht meie tähestikust	“a”
ϵ	Tühisõna	Epsilon tähistab tühja sõna	“”
s t	Alternatsioon	Valik avaldiste sõnade vahel	“s” ja “t”
s t	Konkatenatsioon	Esimese avaldise sõna ja siis teise sõna	“st”
s*	Kleene'i tärn	Avaldise kordus	“”, “s”, “ss”, “sss”, ...

Ja neid saab muidugi kombineerida: avaldise (a|b)c keelde kuulub nii “ac” kui ka “bc”

Regulaaravaldise prioriteedid

- Nagu ka aritmeetikas, võib regulaaravaldises sulud ära jätta!
- Näiteks kirjutame avaldise $(a^* (b | \varepsilon)) | ((bc) a)$ asemel lihtsalt $a^* (b | \varepsilon) | bca$.
- Operaatorite prioriteedid on kahenavas järjekorras kordus, konkatenatsioon ja alternatsioon.

Mata	AKT
a^n	a^*
$a \times b$	ab
$a + b$	$a b$

Keele formaalne definitsioon

$$\begin{aligned}L(\varepsilon) &= \{\text{""}\} \\L(a) &= \{\text{"a"}\} \\L(E_1 E_2) &= \{vw \mid v \in L(E_1), w \in L(E_2)\} \\L(E_1 \mid E_2) &= L(E_1) \cup L(E_2) \\L(E^*) &= \{\text{""}\} \cup \{vw \mid v \in L(E), w \in L(E^*)\}\end{aligned}$$

Näited

Regulaaravaldis

$a \mid b$

abba

$(a \mid b)(c \mid d)$

ab^*a

$(ab)^*$

$(a \mid b)^*$

Defineeritav keel

$\{ "a", "b" \}$

$\{ "abba" \}$

$\{ "ac", "ad", "bc", "bd" \}$

$\{ "aa", "aba", "abba", "abbba", \dots \}$

$\{ "", "ab", "abab", "ababab", \dots \}$

$\{ "", "a", "b", "aa", "ab", "ba",$

$"aaa", "aab", "aba", "baa",$

$"abb", "bab", "bba", "bbb", \dots \}$

Kokkuvõte

- Leksiline analüüs
- Leksiline analüsaator ehk lekser ning selle genereerimine
- Formaalne keel
- Regulaaravaldised, nende koostamine ja nendele vastavad keeled