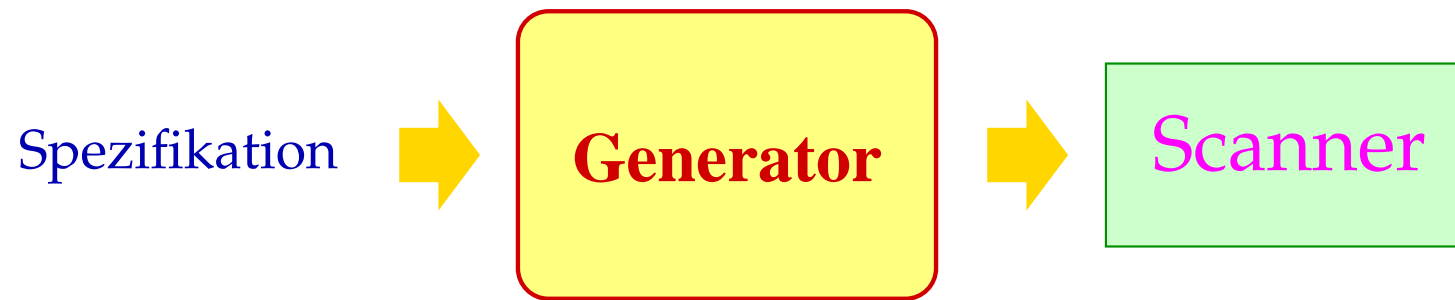
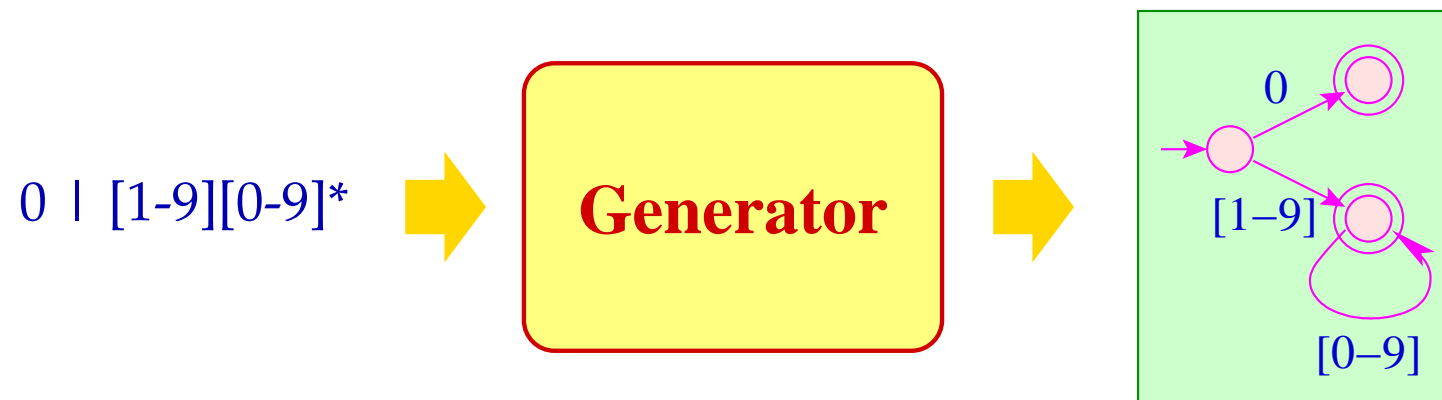


... in unserem Fall:



... in unserem Fall:



Spezifikation von Token-Klassen: Reguläre Ausdrücke;

Generierte Implementierung: Endliche Automaten + X :-)

1.1 Grundlagen: Reguläre Ausdrücke

- Programmtext benutzt ein endliches **Alphabet** Σ von Eingabe-Zeichen, z.B. ASCII :-)
- Die Menge der Textabschnitte einer Token-Klasse ist i.a. **regulär**.
- Reguläre Sprachen kann man mithilfe **regulärer Ausdrücke** spezifizieren.

1.1 Grundlagen: Reguläre Ausdrücke

- Programmtext benutzt ein endliches **Alphabet** Σ von Eingabe-Zeichen, z.B. ASCII :-)
- Die Menge der Textabschnitte einer Token-Klasse ist i.a. **regulär**.
- Reguläre Sprachen kann man mithilfe **regulärer Ausdrücke** spezifizieren.

Die Menge \mathcal{E}_Σ der (nicht-leeren) **regulären Ausdrücke** ist die kleinste Menge \mathcal{E} mit:

- $\epsilon \in \mathcal{E}$ (ϵ neues Symbol nicht aus Σ);
- $a \in \mathcal{E}$ für alle $a \in \Sigma$;
- $(e_1 \mid e_2), (e_1 \cdot e_2), e_1^* \in \mathcal{E}$ sofern $e_1, e_2 \in \mathcal{E}$.



Stephen Kleene, Madison Wisconsin, 1909-1994

Beispiele:

$$((a \cdot b^*) \cdot a)$$

$$(a \mid b)$$

$$((a \cdot b) \cdot (a \cdot b))$$

Beispiele:

$$((a \cdot b^*) \cdot a)$$

$$(a \mid b)$$

$$((a \cdot b) \cdot (a \cdot b))$$

Achtung:

- Wir unterscheiden zwischen Zeichen $a, 0, |, \dots$ und **Meta-Zeichen** $(, |,), \dots$
- Um (hässliche) Klammern zu sparen, benutzen wir **Operator-Präzedenzen**:

$$* > \cdot > |$$

und lassen “.” weg :-)

Beispiele:

$$((a \cdot b^*) \cdot a)$$

$$(a \mid b)$$

$$((a \cdot b) \cdot (a \cdot b))$$

Achtung:

- Wir unterscheiden zwischen Zeichen $a, 0, |, \dots$ und **Meta-Zeichen** $(, |,), \dots$
- Um (hässliche) Klammern zu sparen, benutzen wir **Operator-Präzedenzen**:

$$* > \cdot > |$$

und lassen “.” weg :-)

- Reale Spezifikations-Sprachen bieten zusätzliche Konstrukte wie:

$$e? \equiv (\epsilon \mid e)$$

$$e^+ \equiv (e \cdot e^*)$$

und verzichten auf “ ϵ ” :-)

Spezifikationen benötigen eine Semantik :-)

Im Beispiel:

Spezifikation	Semantik
ab^*a	$\{ab^n a \mid n \geq 0\}$
$a \mid b$	$\{a, b\}$
$abab$	$\{abab\}$

Für $e \in \mathcal{E}_\Sigma$ definieren wir die spezifizierte Sprache $\llbracket e \rrbracket \subseteq \Sigma^*$ induktiv durch:

$$\llbracket \epsilon \rrbracket = \{\epsilon\}$$

$$\llbracket a \rrbracket = \{a\}$$

$$\llbracket e^* \rrbracket = (\llbracket e \rrbracket)^*$$

$$\llbracket e_1 \mid e_2 \rrbracket = \llbracket e_1 \rrbracket \cup \llbracket e_2 \rrbracket$$

$$\llbracket e_1 \cdot e_2 \rrbracket = \llbracket e_1 \rrbracket \cdot \llbracket e_2 \rrbracket$$

Beachte:

- Die Operatoren $(_)*, \cup, \cdot$ sind die entsprechenden Operationen auf Wort-Mengen:

$$(L)^* = \{w_1 \dots w_k \mid k \geq 0, w_i \in L\}$$

$$L_1 \cdot L_2 = \{w_1 w_2 \mid w_1 \in L_1, w_2 \in L_2\}$$

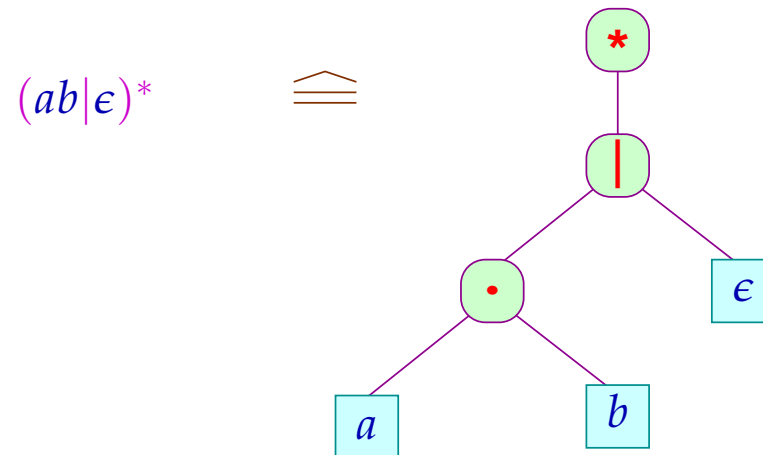
Beachte:

- Die Operatoren $(_)*, \cup, \cdot$ sind die entsprechenden Operationen auf Wort-Mengen:

$$(L)^* = \{w_1 \dots w_k \mid k \geq 0, w_i \in L\}$$

$$L_1 \cdot L_2 = \{w_1 w_2 \mid w_1 \in L_1, w_2 \in L_2\}$$

- Reguläre Ausdrücke stellen wir intern als **markierte geordnete Bäume** dar:



Innere Knoten: Operator-Anwendungen;

Blätter: einzelne Zeichen oder ϵ .

Finger-Übung:

Zu jedem regulären Ausdruck e können wir einen Ausdruck e' (evt. mit “?”) konstruieren so dass:

- $\llbracket e \rrbracket = \llbracket e' \rrbracket$;
- Falls $\llbracket e \rrbracket = \{\epsilon\}$, dann ist $e' \equiv \epsilon$;
- Falls $\llbracket e \rrbracket \neq \{\epsilon\}$, dann enthält e' kein “ ϵ ”.

Finger-Übung:

Zu jedem regulären Ausdruck e können wir einen Ausdruck e' (evt. mit “?”) konstruieren so dass:

- $\llbracket e \rrbracket = \llbracket e' \rrbracket$;
- Falls $\llbracket e \rrbracket = \{\epsilon\}$, dann ist $e' \equiv \epsilon$;
- Falls $\llbracket e \rrbracket \neq \{\epsilon\}$, dann enthält e' kein “ ϵ ”.

Konstruktion:

Wir definieren eine Transformation \mathcal{T} von regulären Ausdrücken durch:

$$\begin{aligned}
\mathcal{T}[\epsilon] &= \epsilon \\
\mathcal{T}[a] &= a \\
\mathcal{T}[e_1|e_2] &= \text{case } (\mathcal{T}[e_1], \mathcal{T}[e_2]) \text{ of } \begin{array}{l} (\epsilon, \epsilon) : \epsilon \\ | (e'_1, \epsilon) : e'_1? \\ | (\epsilon, e'_2) : e'_2? \\ | (e'_1, e'_2) : (e'_1 | e'_2) \end{array} \\
\mathcal{T}[e_1 \cdot e_2] &= \text{case } (\mathcal{T}[e_1], \mathcal{T}[e_2]) \text{ of } \begin{array}{l} (\epsilon, \epsilon) : \epsilon \\ | (e'_1, \epsilon) : e'_1 \\ | (\epsilon, e'_2) : e'_2 \\ | (e'_1, e'_2) : (e'_1 \cdot e'_2) \end{array} \\
\mathcal{T}[e^*] &= \text{case } \mathcal{T}[e] \text{ of } \begin{array}{l} \epsilon : \epsilon \\ | e_1 : e_1^* \end{array} \\
\mathcal{T}[e?] &= \text{case } \mathcal{T}[e] \text{ of } \begin{array}{l} \epsilon : \epsilon \\ | e_1 : e_1? \end{array}
\end{aligned}$$

Unsere Anwendung:

Identifizier in Java:

le = [a-zA-Z_\\$]

di = [0-9]

Id = {le} ({le} | {di})*

Unsere Anwendung:

Identifizier in Java:

le = [a-zA-Z_\\$]

di = [0-9]

Id = {le} ({le} | {di})*

Bemerkungen:

- “le” und “di” sind **Zeichenklassen**.
- **Definierte Namen** werden in “{”, “}” eingeschlossen.
- Zeichen werden von **Meta-Zeichen** durch “\” unterschieden.

Unsere Anwendung:

Identifizier in Java:

le = [a-zA-Z_\\$]

di = [0-9]

Id = {le} ({le}|{di})*

Gleitkommazahlen:

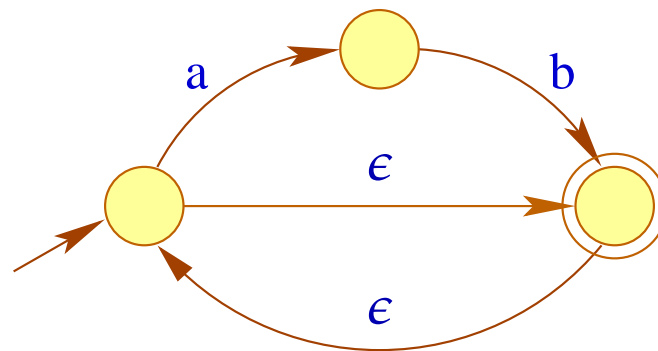
Float = {di}* (\.{di}|{di}\.) {di}* ((e|E)(\+|\-)?{di}+)?

Bemerkungen:

- “le” und “di” sind **Zeichenklassen**.
- **Definierte Namen** werden in “{”, “}” eingeschlossen.
- Zeichen werden von **Meta-Zeichen** durch “\” unterschieden.

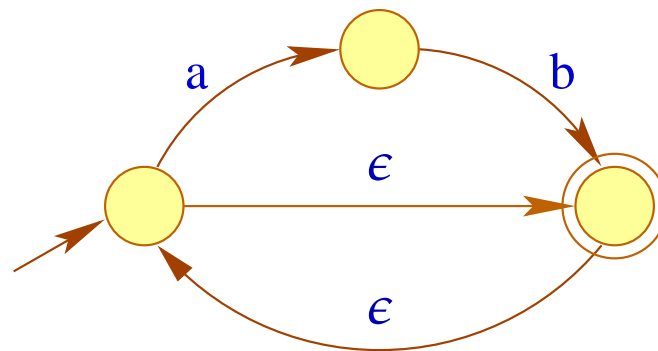
1.2 Grundlagen: Endliche Automaten

Beispiel:



1.2 Grundlagen: Endliche Automaten

Beispiel:



Knoten: Zustände;

Kanten: Übergänge;

Beschriftungen: konsumierter Input :-)



Michael O. Rabin, Stanford University



Dana S. Scott, Carnegie Mellon
University, Pittsburgh

Formal ist ein nicht-deterministischer endlicher Automat mit ϵ -Übergängen (ϵ -NFA) ein Tupel $A = (Q, \Sigma, \delta, I, F)$ wobei:

- Q eine endliche Menge von Zuständen;
- Σ ein endliches Eingabe-Alphabet;
- $I \subseteq Q$ die Menge der Anfangszustände;
- $F \subseteq Q$ die Menge der Endzustände und
- δ die Menge der Übergänge (die Übergangs-Relation) ist.

Formal ist ein nicht-deterministischer endlicher Automat mit ϵ -Übergängen (ϵ -NFA) ein Tupel $A = (Q, \Sigma, \delta, I, F)$ wobei:

- Q eine endliche Menge von Zuständen;
- Σ ein endliches Eingabe-Alphabet;
- $I \subseteq Q$ die Menge der Anfangszustände;
- $F \subseteq Q$ die Menge der Endzustände und
- δ die Menge der Übergänge (die Übergangs-Relation) ist.

Für ϵ -NFAs ist:

$$\delta \subseteq Q \times (\Sigma \cup \{\epsilon\}) \times Q$$

Formal ist ein **nicht-deterministischer** endlicher Automat mit ϵ -Übergängen (**ϵ -NFA**) ein Tupel $A = (Q, \Sigma, \delta, I, F)$ wobei:

- Q eine endliche Menge von Zuständen;
- Σ ein endliches Eingabe-Alphabet;
- $I \subseteq Q$ die Menge der Anfangszustände;
- $F \subseteq Q$ die Menge der Endzustände und
- δ die Menge der Übergänge (die Übergangs-Relation) ist.

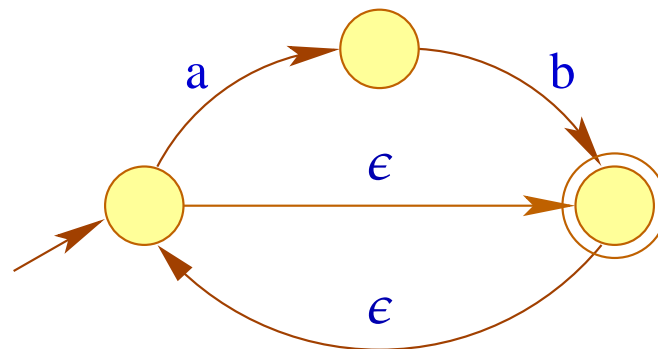
Für **ϵ -NFAs** ist:

$$\delta \subseteq Q \times (\Sigma \cup \{\epsilon\}) \times Q$$

- Gibt es keine ϵ -Übergänge (p, ϵ, q) , ist A ein **NFA**.
- Ist $\delta : Q \times \Sigma \rightarrow Q$ eine Funktion und $\#I = 1$, heißt A **deterministisch (DFA)**.

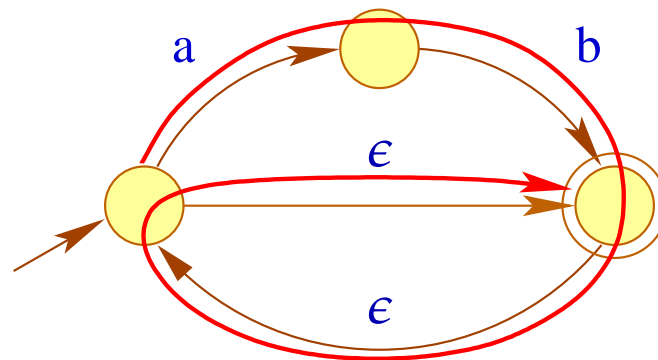
Akzeptierung

- Berechnungen sind Pfade im Graphen.
- akzeptierende Berechnungen führen von I nach F .
- Ein akzeptiertes Wort ist die Beschriftung eines akzeptierenden Pfades ...



Akzeptierung

- Berechnungen sind Pfade im Graphen.
- akzeptierende Berechnungen führen von I nach F .
- Ein akzeptiertes Wort ist die Beschriftung eines akzeptierenden Pfades ...



- Dazu definieren wir den **transitiven Abschluss** δ^* von δ als kleinste Menge δ' mit:

$$(p, \epsilon, p) \in \delta' \text{ und}$$

$$(p, xw, q) \in \delta' \text{ sofern } (p, x, p_1) \in \delta \text{ und } (p_1, w, q) \in \delta'.$$

δ^* beschreibt für je zwei Zustände, mit welchen Wörtern man vom einen zum andern kommt :-)

- Die Menge aller akzeptierten Worte, d.h. die von A akzeptierte Sprache können wir kurz beschreiben als:

$$\mathcal{L}(A) = \{w \in \Sigma^* \mid \exists i \in I, f \in F : (i, w, f) \in \delta^*\}$$

Satz:

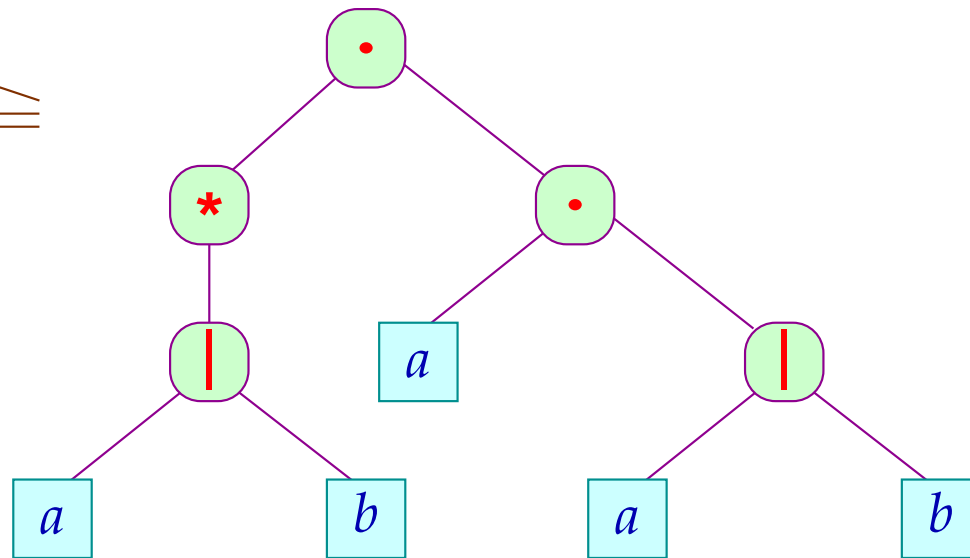
Für jeden regulären Ausdruck e kann (in linearer Zeit :-) ein ϵ -NFA konstruiert werden, der die Sprache $\llbracket e \rrbracket$ akzeptiert.

Idee:

Der Automat verfolgt (konzeptionell mithilfe einer Marke “ \bullet ”), wohin man in e mit der Eingabe w gelangen kann.

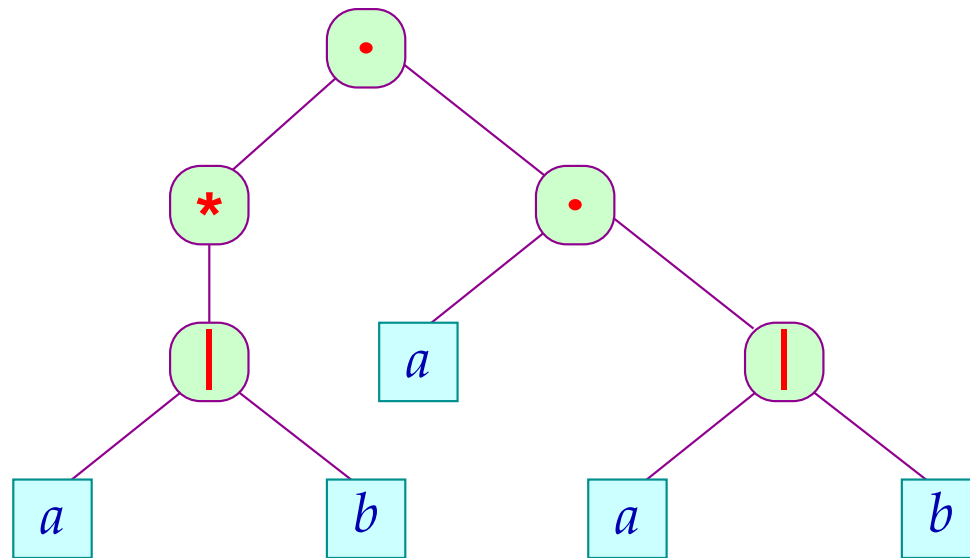
Beispiel:

$$(a|b)^*a(a|b)$$



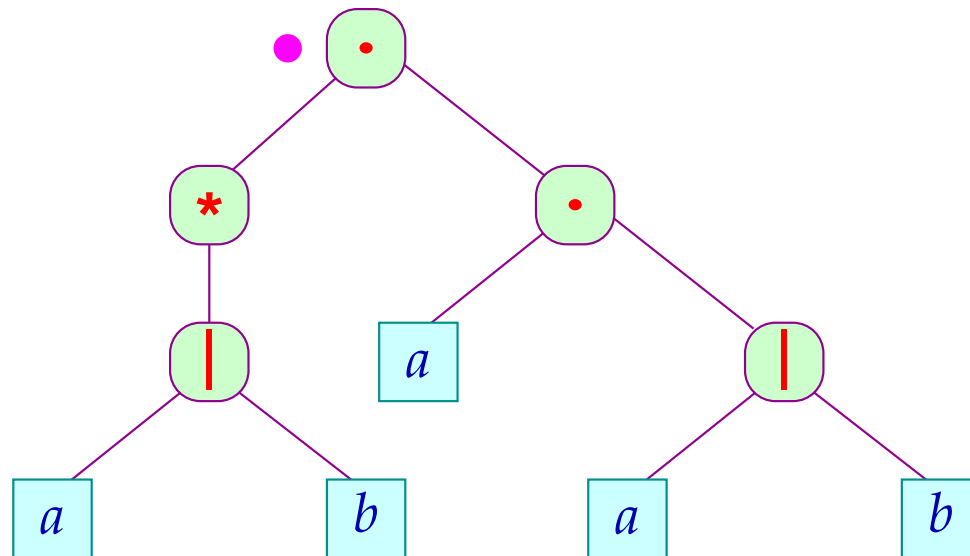
Beispiel:

$w = bbaa$:



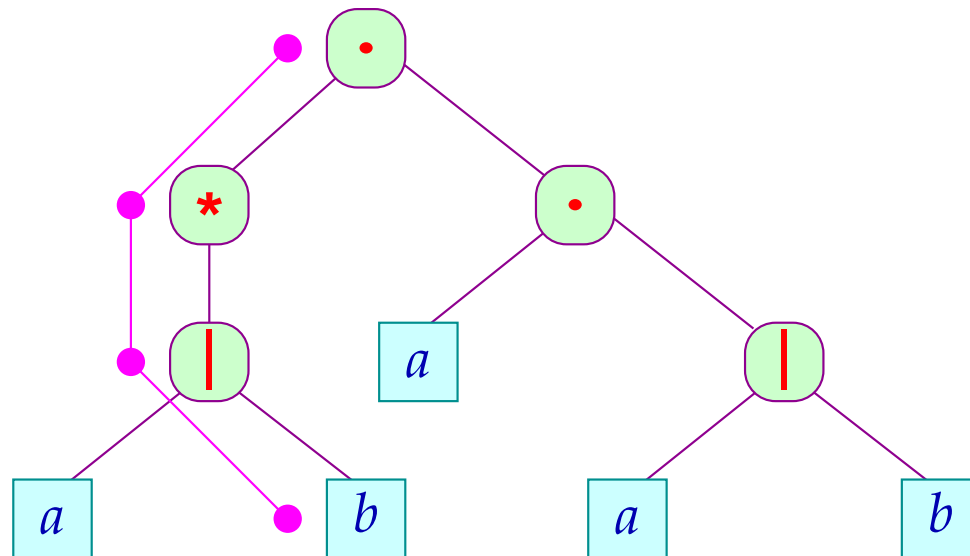
Beispiel:

$w = bbaa$:



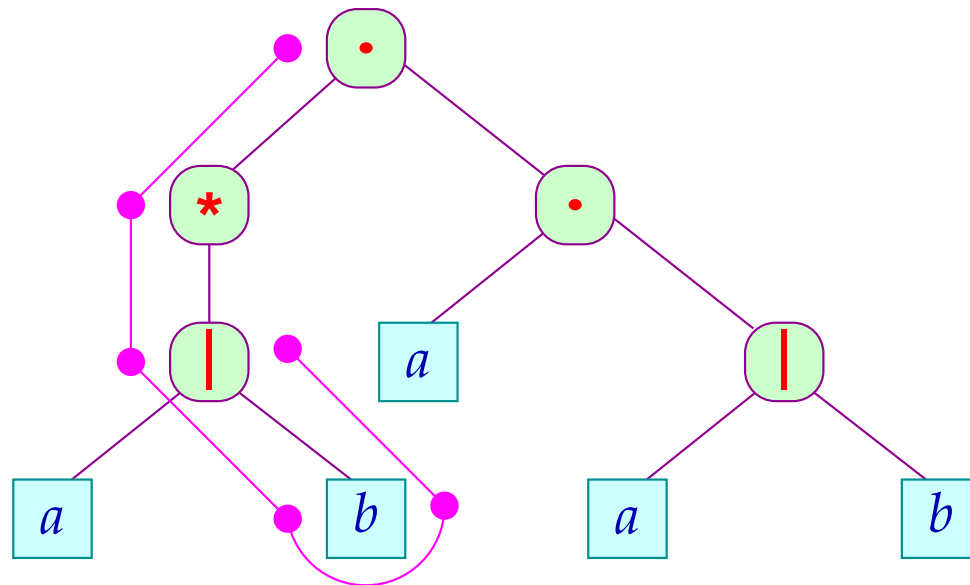
Beispiel:

$w = bbaa$:



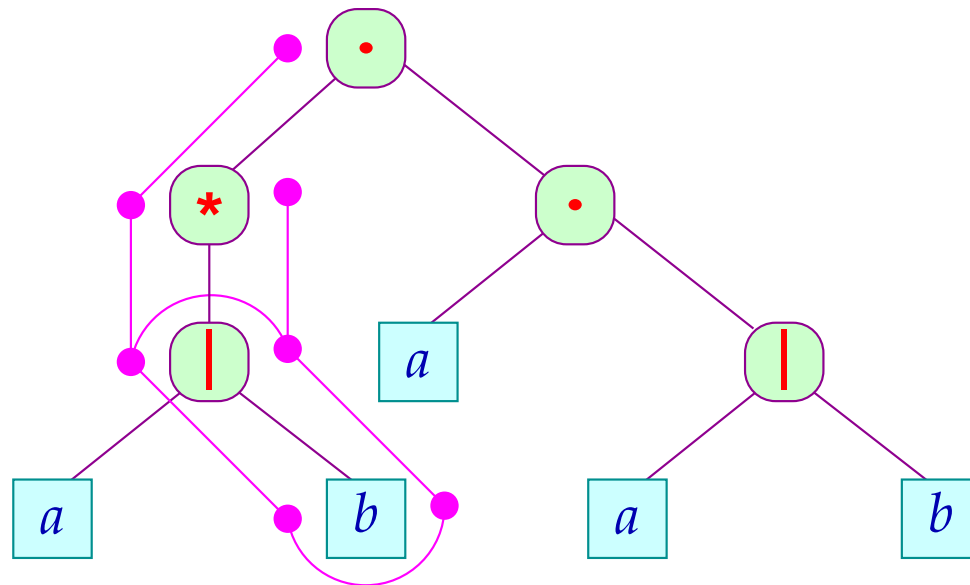
Beispiel:

$w = bbaa$:



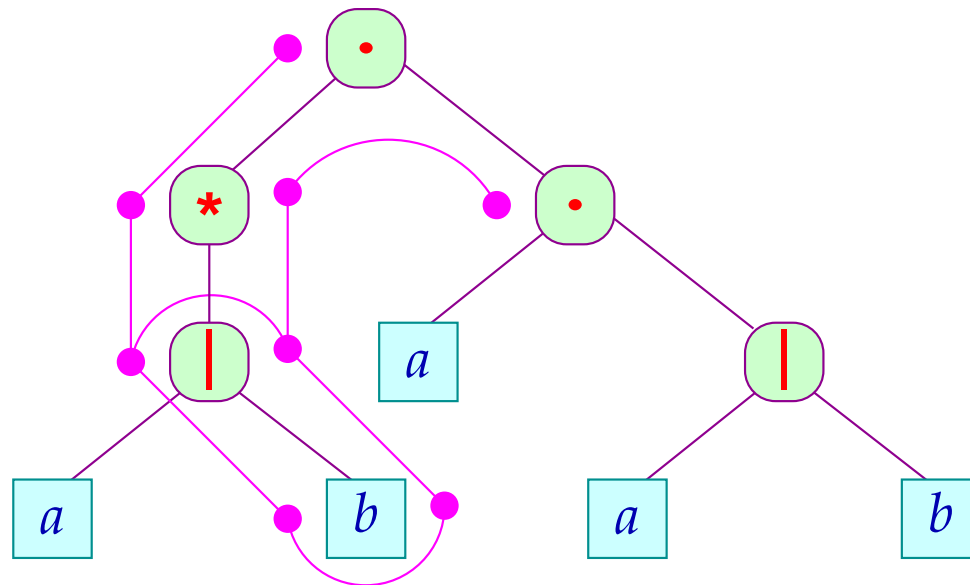
Beispiel:

$w = bbaa$:



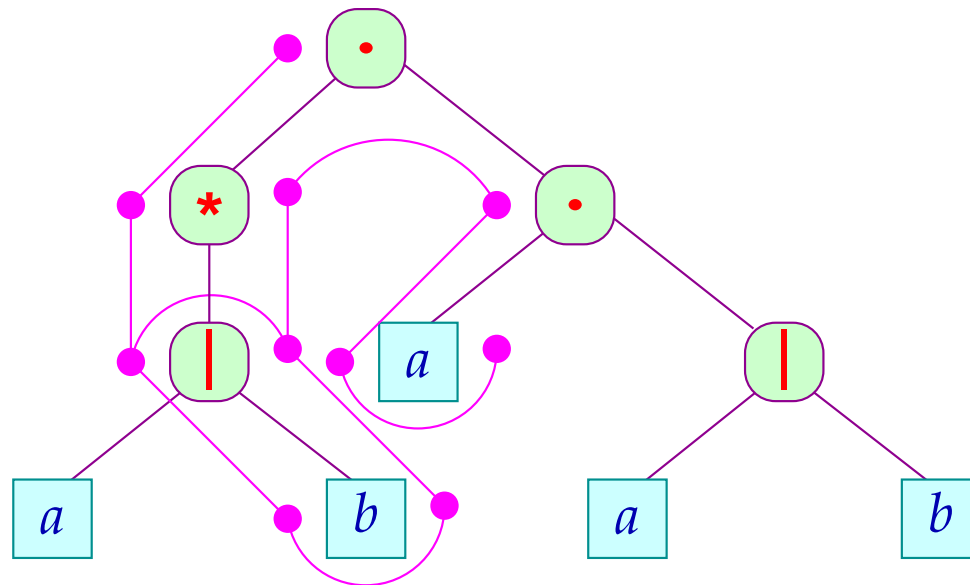
Beispiel:

$w = bbaa$:



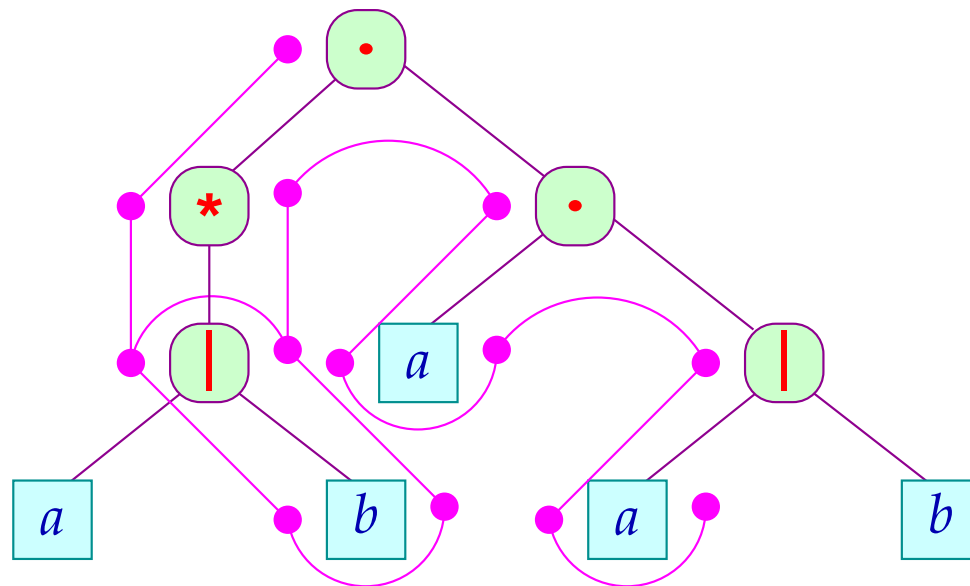
Beispiel:

$w = bbaa$:



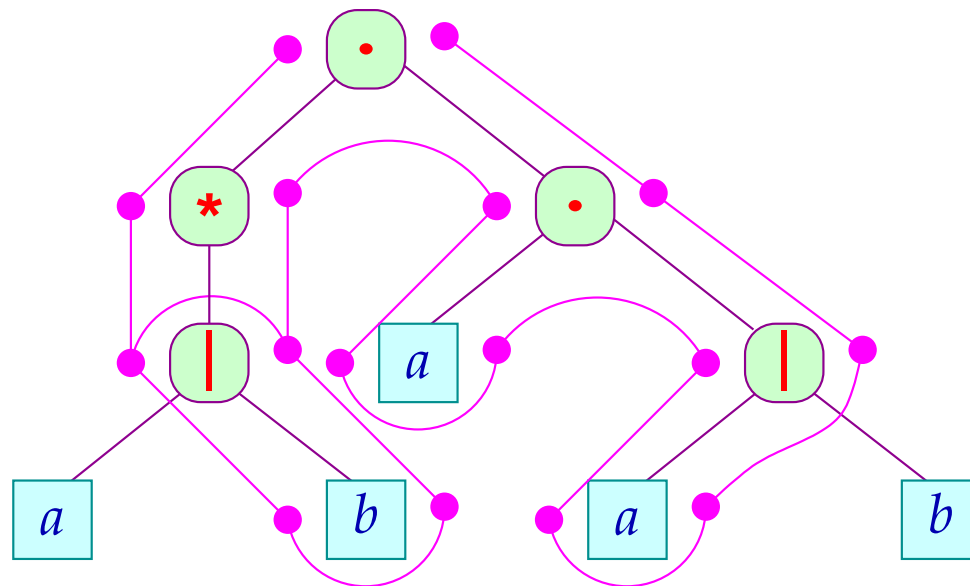
Beispiel:

$w = bbaa$:



Beispiel:

$w = bbaa$:

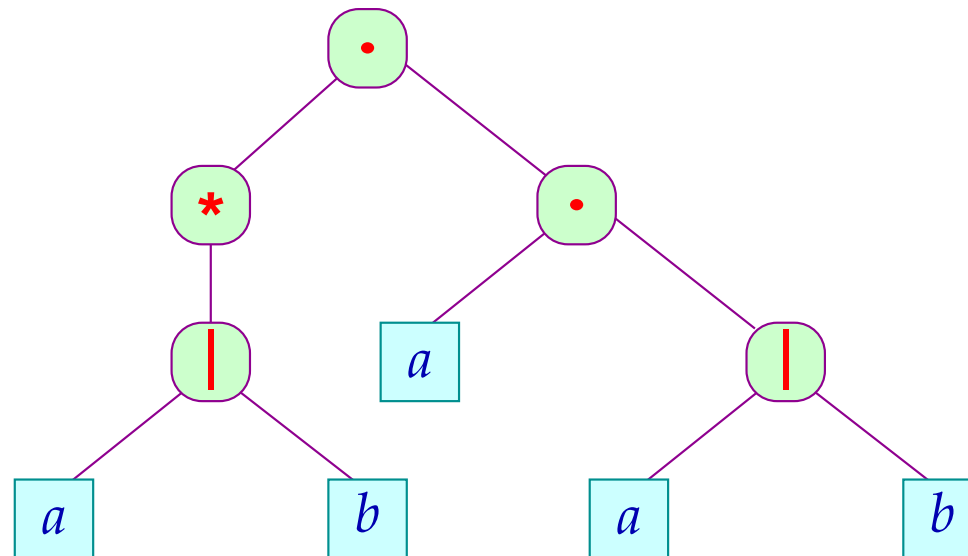


Beachte:

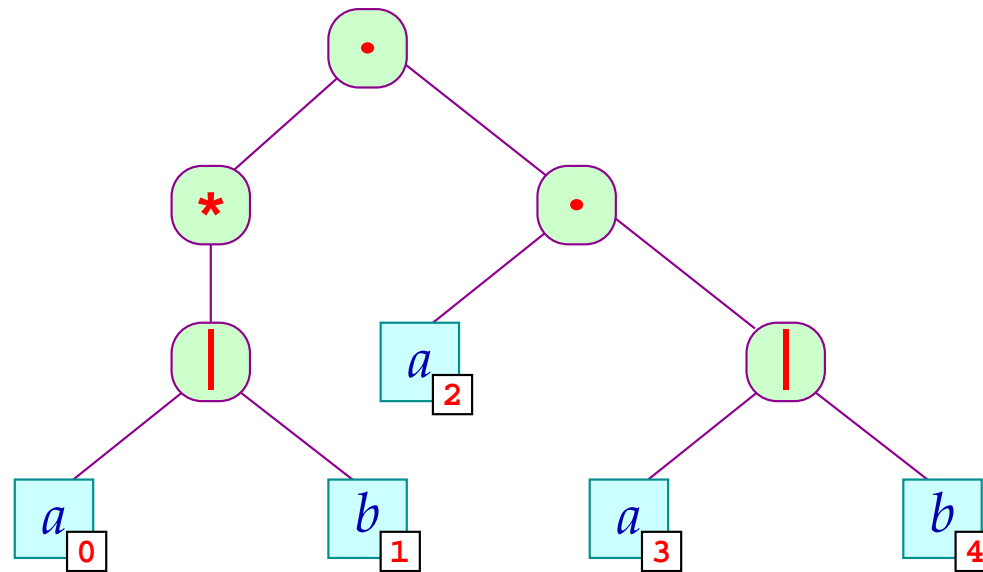
- Gelesen wird nur an den Blättern.
- Die Navigation im Baum erfolgt ohne Lesen, d.h. mit ϵ -Übergängen.
- Für eine formale Konstruktion müssen wir die Knoten im Baum **bezeichnen**.
- Dazu benutzen wir (hier) einfach den dargestellten **Teilausdruck** :-)
- Leider gibt es eventuell mehrere gleiche Teilausdrücke :-)

\implies Wir numerieren die Blätter durch ...

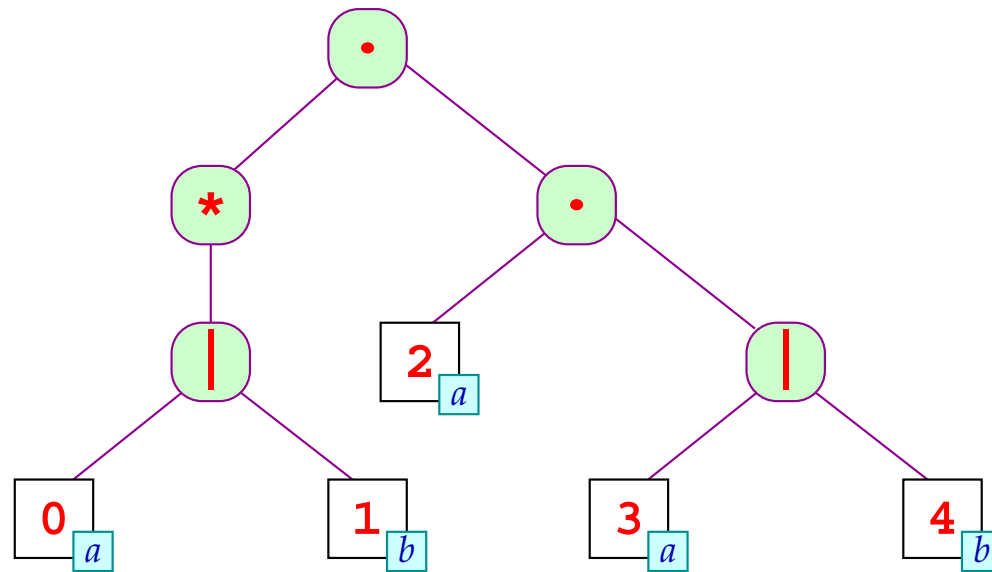
... im Beispiel:



... im Beispiel:



... im Beispiel:



Die Konstruktion:

Zustände: $\bullet r, r \bullet$ r Knoten von e ;

Anfangszustand: $\bullet e$;

Endzustand: $e \bullet$;

Übergangsrelation:

Für Blätter $r \equiv \boxed{i \mid x}$ benötigen wir: $(\bullet r, x, r \bullet)$.

Die übrigen Übergänge sind:

r	Übergänge
$r_1 \mid r_2$	$(\bullet r, \epsilon, \bullet r_1)$ $(\bullet r, \epsilon, \bullet r_2)$ $(r_1 \bullet, \epsilon, r \bullet)$ $(r_2 \bullet, \epsilon, r \bullet)$
$r_1 \cdot r_2$	$(\bullet r, \epsilon, \bullet r_1)$ $(r_1 \bullet, \epsilon, \bullet r_2)$ $(r_2 \bullet, \epsilon, r \bullet)$

r	Übergänge
r_1^*	$(\bullet r, \epsilon, r \bullet)$ $(\bullet r, \epsilon, \bullet r_1)$ $(r_1 \bullet, \epsilon, \bullet r_1)$ $(r_1 \bullet, \epsilon, r \bullet)$
$r_1?$	$(\bullet r, \epsilon, r \bullet)$ $(\bullet r, \epsilon, \bullet r_1)$ $(r_1 \bullet, \epsilon, r \bullet)$

Diskussion:

- Die meisten Übergänge dienen dazu, im Ausdruck zu navigieren :-)
- Der Automat ist i.a. nichtdeterministisch :-)

Diskussion:

- Die meisten Übergänge dienen dazu, im Ausdruck zu navigieren :-)
- Der Automat ist i.a. nichtdeterministisch :-)

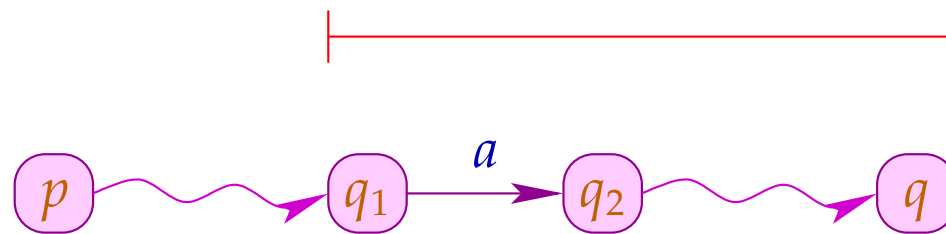


Strategie:

- (1) Beseitigung der ϵ -Übergänge;
- (2) Beseitigung des Nichtdeterminismus :-)

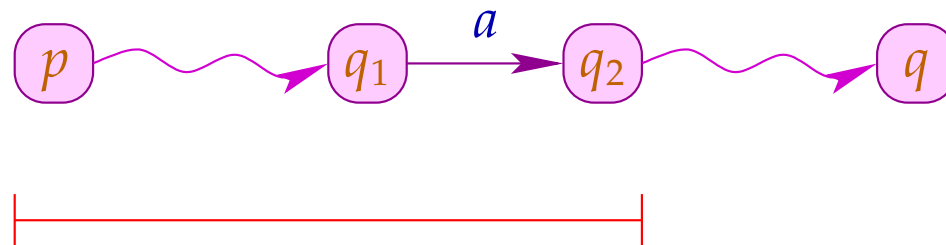
Beseitigung von ϵ -Übergängen:

Zwei einfache Ansätze:



Beseitigung von ϵ -Übergängen:

Zwei einfache Ansätze:



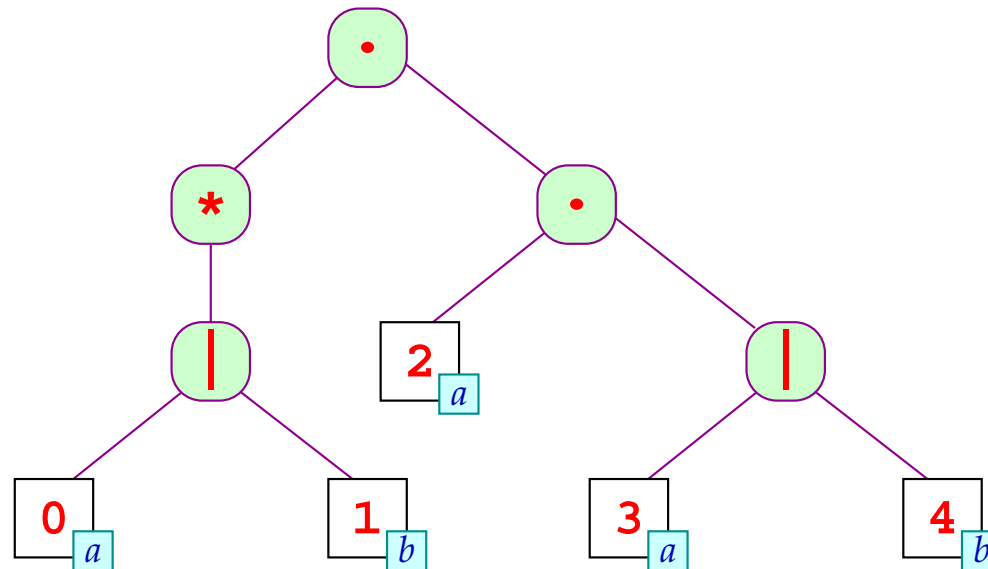
Wir benutzen hier den zweiten Ansatz.

Zur Konstruktion von Parsern werden wir später den ersten benutzen :-)

1. Schritt:

$$\text{empty}[r] = t \quad \text{gdw.} \quad \epsilon \in [[r]]$$

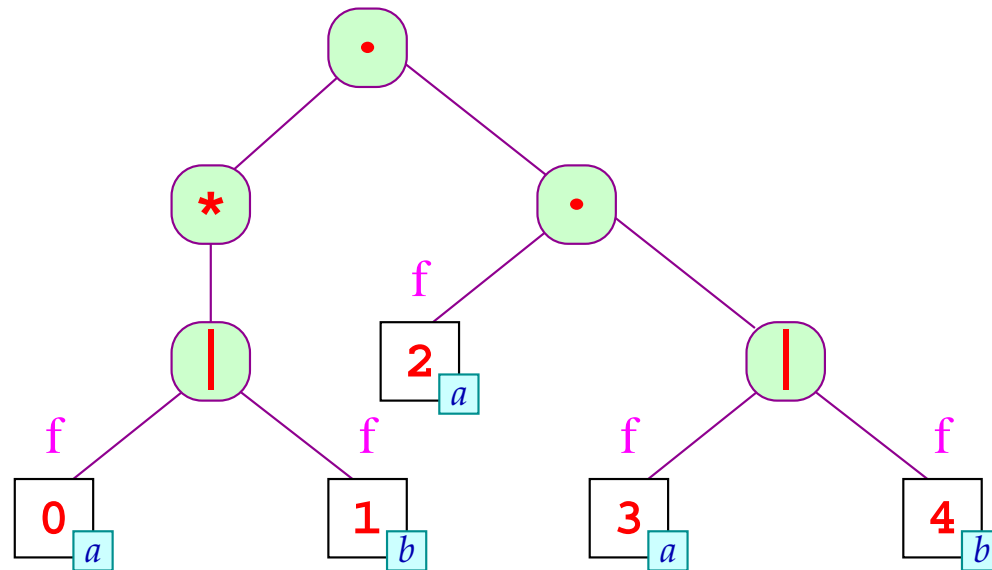
... im Beispiel:



1. Schritt:

$$\text{empty}[r] = t \quad \text{gdw.} \quad \epsilon \in [[r]]$$

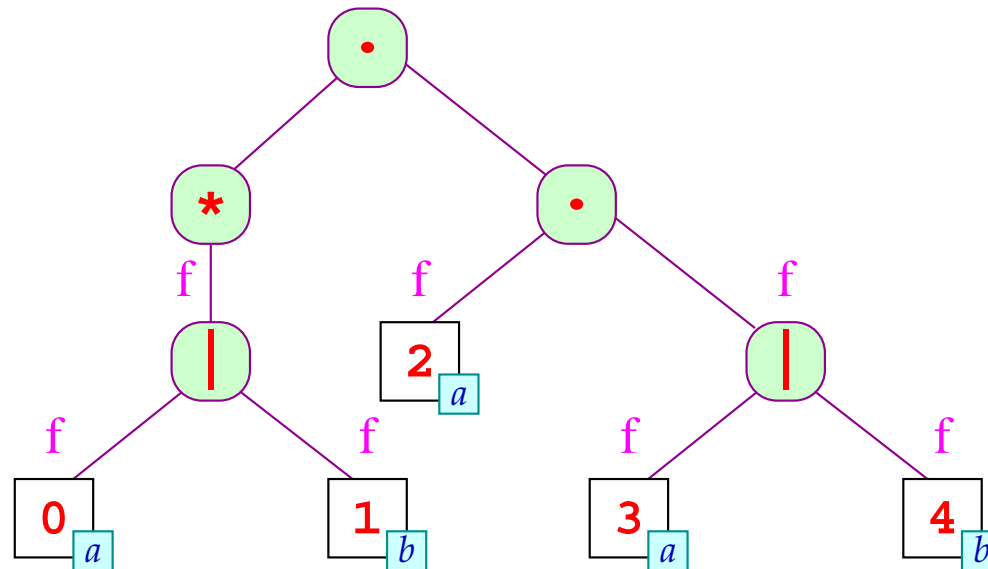
... im Beispiel:



1. Schritt:

$$\text{empty}[r] = t \quad \text{gdw.} \quad \epsilon \in [[r]]$$

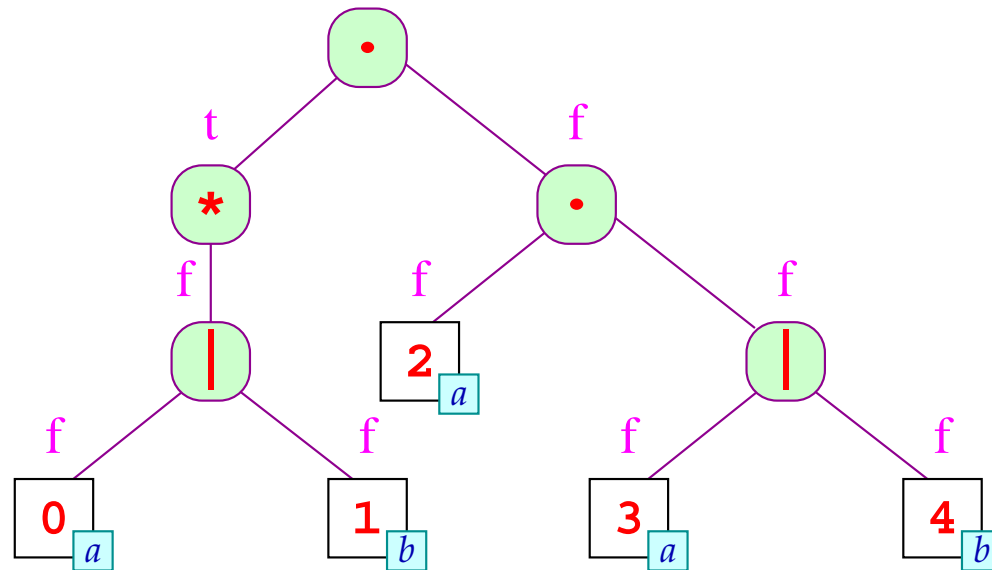
... im Beispiel:



1. Schritt:

$$\text{empty}[r] = t \quad \text{gdw.} \quad \epsilon \in [[r]]$$

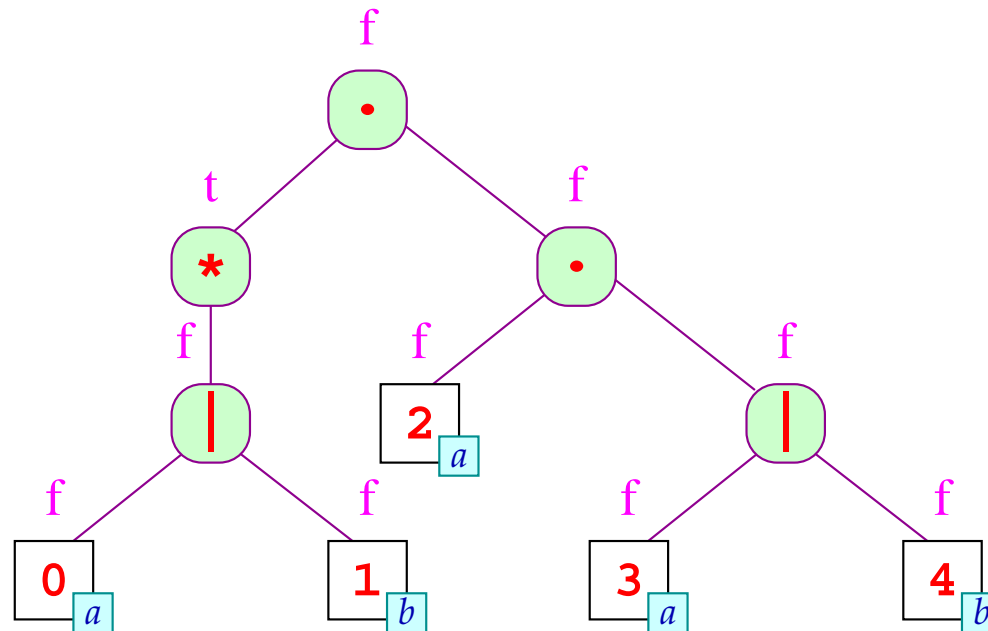
... im Beispiel:



1. Schritt:

$$\text{empty}[r] = t \quad \text{gdw.} \quad \epsilon \in [[r]]$$

... im Beispiel:



Implementierung:

DFS post-order Traversierung

Für Blätter $r \equiv \boxed{i \mid x}$ ist $\text{empty}[r] = (x \equiv \epsilon)$.

Andernfalls:

$$\text{empty}[r_1 \mid r_2] = \text{empty}[r_1] \vee \text{empty}[r_2]$$

$$\text{empty}[r_1 \cdot r_2] = \text{empty}[r_1] \wedge \text{empty}[r_2]$$

$$\text{empty}[r_1^*] = t$$

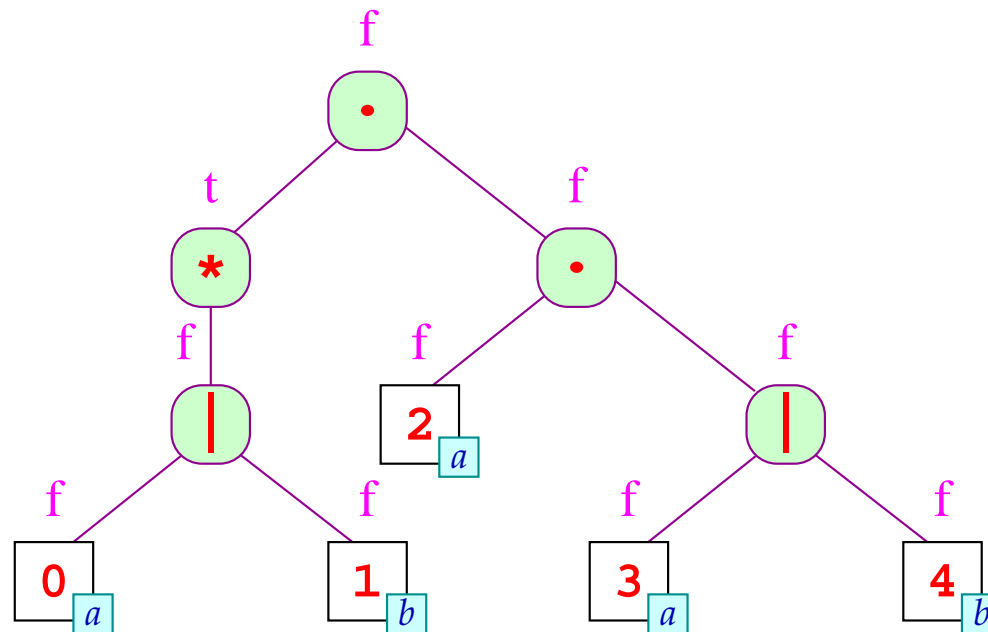
$$\text{empty}[r_1?] = t$$

2. Schritt:

Die Menge erster Blätter:

$$\text{first}[r] = \{i \text{ in } r \mid (\bullet r, \epsilon, \bullet \boxed{i \mid x}) \in \delta^*, x \neq \epsilon\}$$

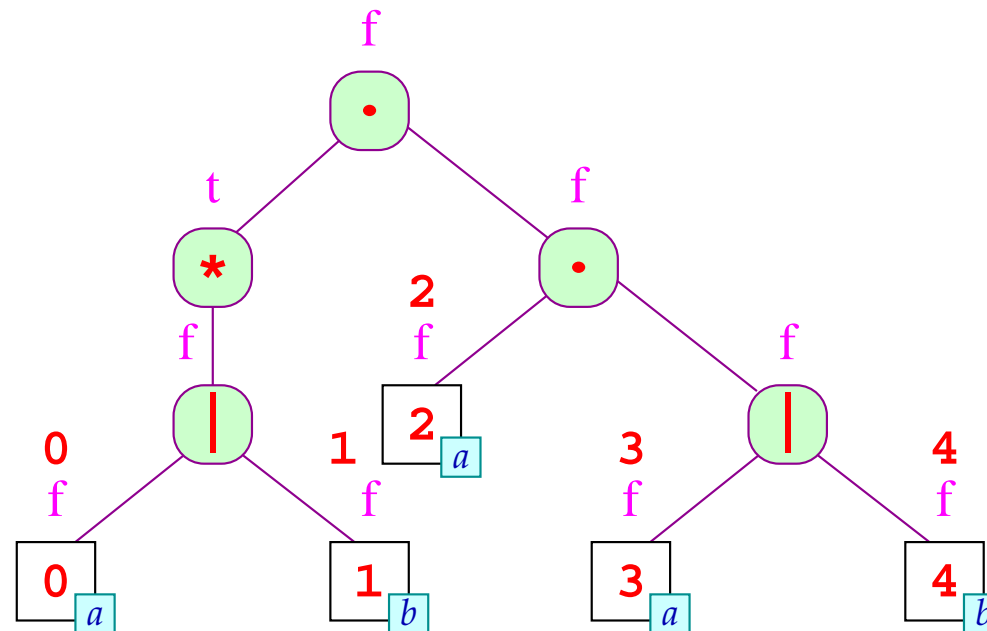
... im Beispiel:



2. Schritt:

Die Menge erster Blätter: $\text{first}[r] = \{i \text{ in } r \mid (\bullet r, \epsilon, \bullet \boxed{i \mid x}) \in \delta^*, x \neq \epsilon\}$

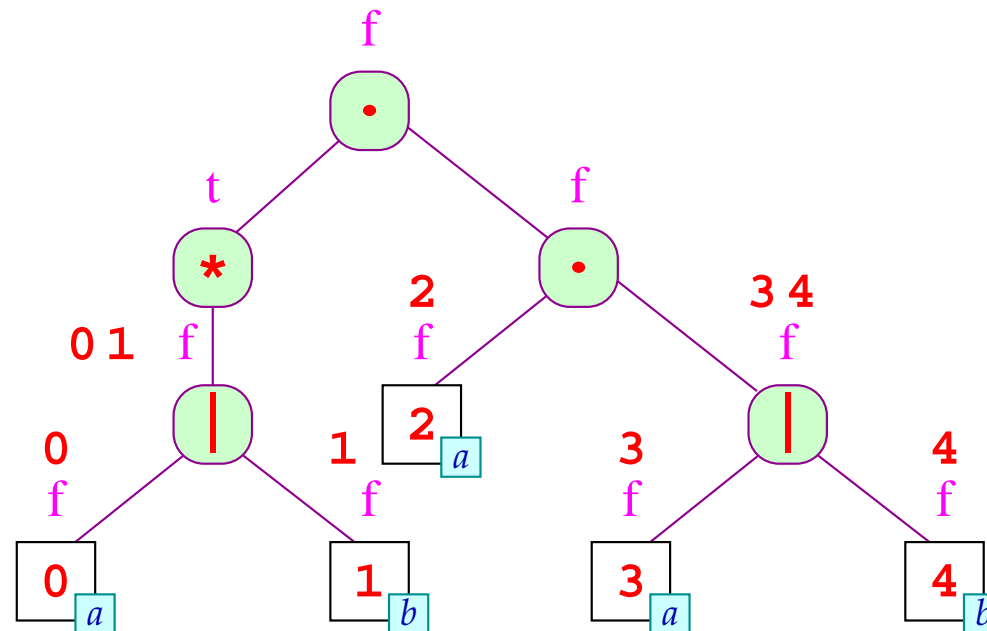
... im Beispiel:



2. Schritt:

Die Menge erster Blätter: $\text{first}[r] = \{i \text{ in } r \mid (\bullet r, \epsilon, \bullet \boxed{i} \mid x) \in \delta^*, x \neq \epsilon\}$

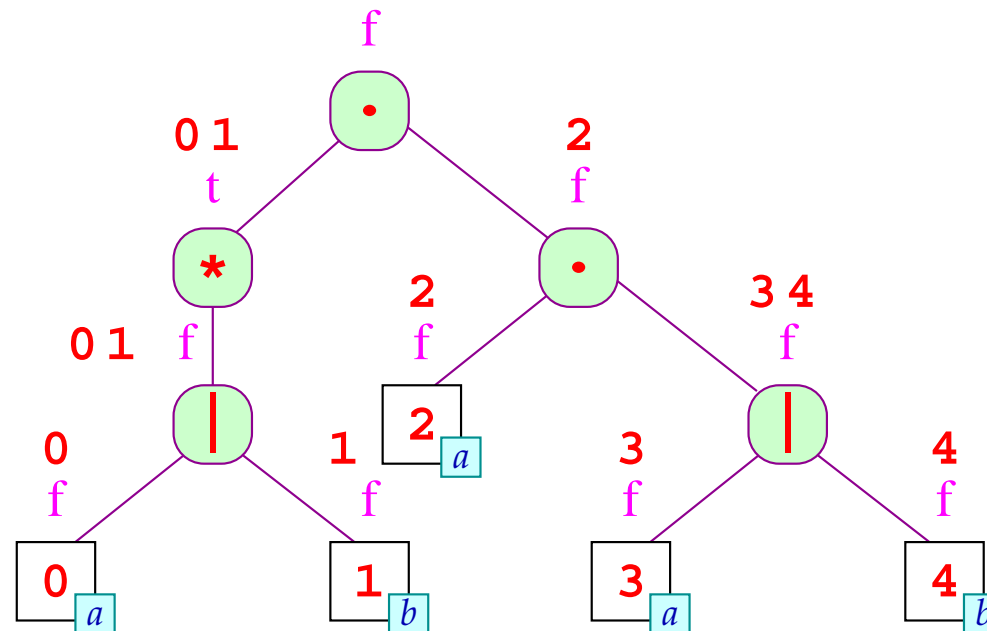
... im Beispiel:



2. Schritt:

Die Menge erster Blätter: $\text{first}[r] = \{i \text{ in } r \mid (\bullet r, \epsilon, \bullet \boxed{i \mid x}) \in \delta^*, x \neq \epsilon\}$

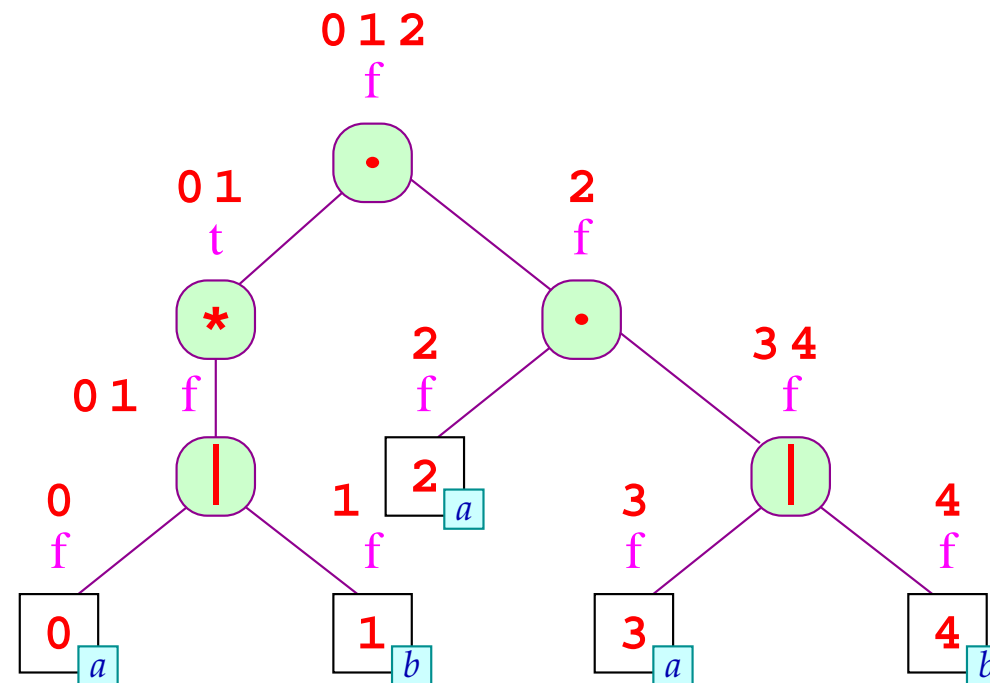
... im Beispiel:



2. Schritt:

Die Menge erster Blätter: $\text{first}[r] = \{i \text{ in } r \mid (\bullet r, \epsilon, \bullet \boxed{i \mid x}) \in \delta^*, x \neq \epsilon\}$

... im Beispiel:



Implementierung:

DFS post-order Traversierung

Für Blätter $r \equiv \boxed{i \mid x}$ ist $\text{first}[r] = \{i \mid x \neq \epsilon\}$.

Andernfalls:

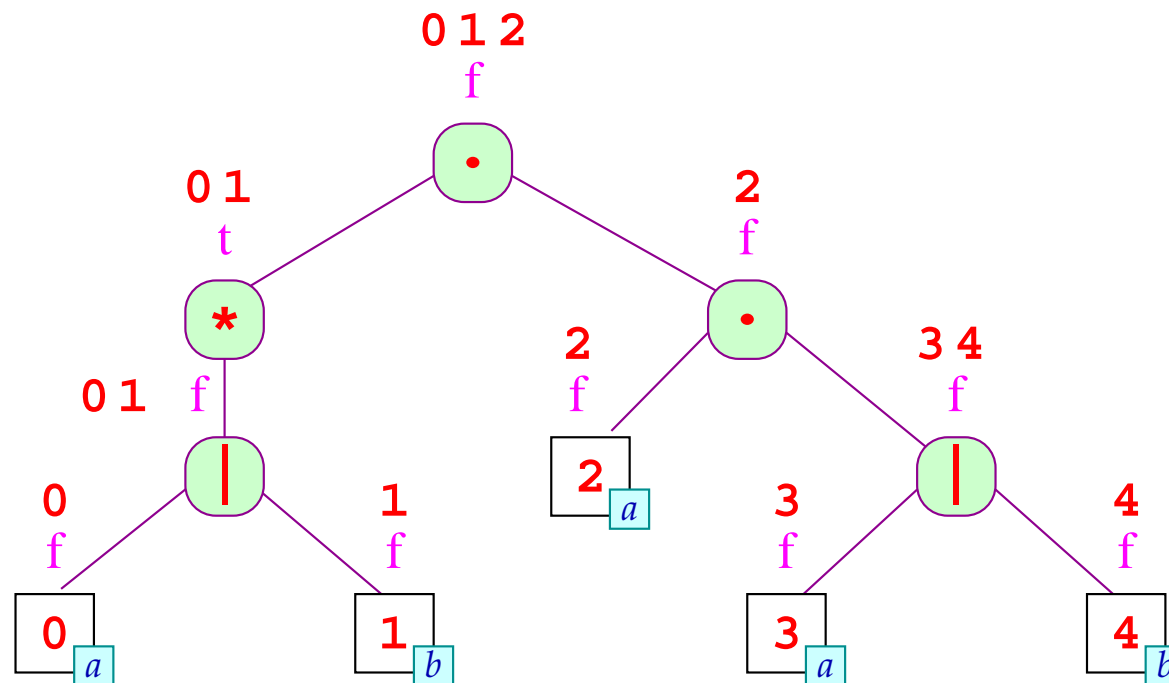
$$\begin{aligned}\text{first}[r_1 \mid r_2] &= \text{first}[r_1] \cup \text{first}[r_2] \\ \text{first}[r_1 \cdot r_2] &= \begin{cases} \text{first}[r_1] \cup \text{first}[r_2] & \text{falls } \text{empty}[r_1] = t \\ \text{first}[r_1] & \text{falls } \text{empty}[r_1] = f \end{cases} \\ \text{first}[r_1^*] &= \text{first}[r_1] \\ \text{first}[r_1?] &= \text{first}[r_1]\end{aligned}$$

3. Schritt:

Die Menge nächster Blätter:

$$\text{next}[r] = \{i \mid (r \bullet, \epsilon, \bullet \boxed{i \mid x}) \in \delta^*, x \neq \epsilon\}$$

... im Beispiel:

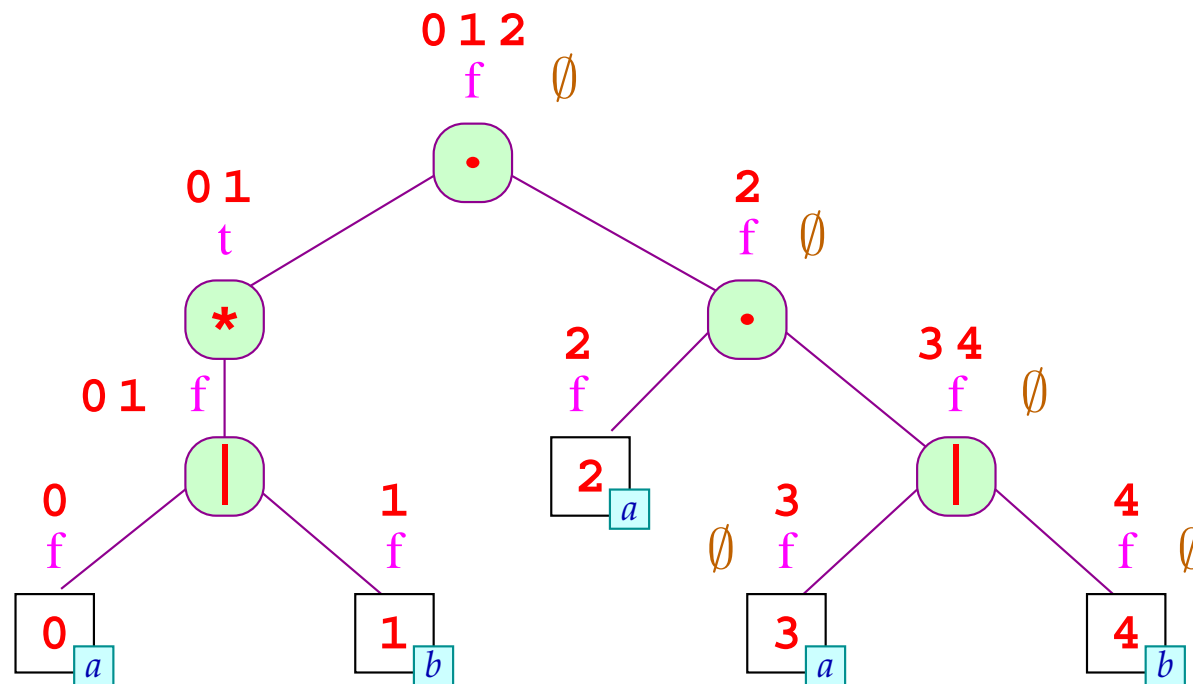


3. Schritt:

Die Menge nächster Blätter:

$$\text{next}[r] = \{i \mid (r \bullet, \epsilon, \bullet \boxed{i \mid x}) \in \delta^*, x \neq \epsilon\}$$

... im Beispiel:

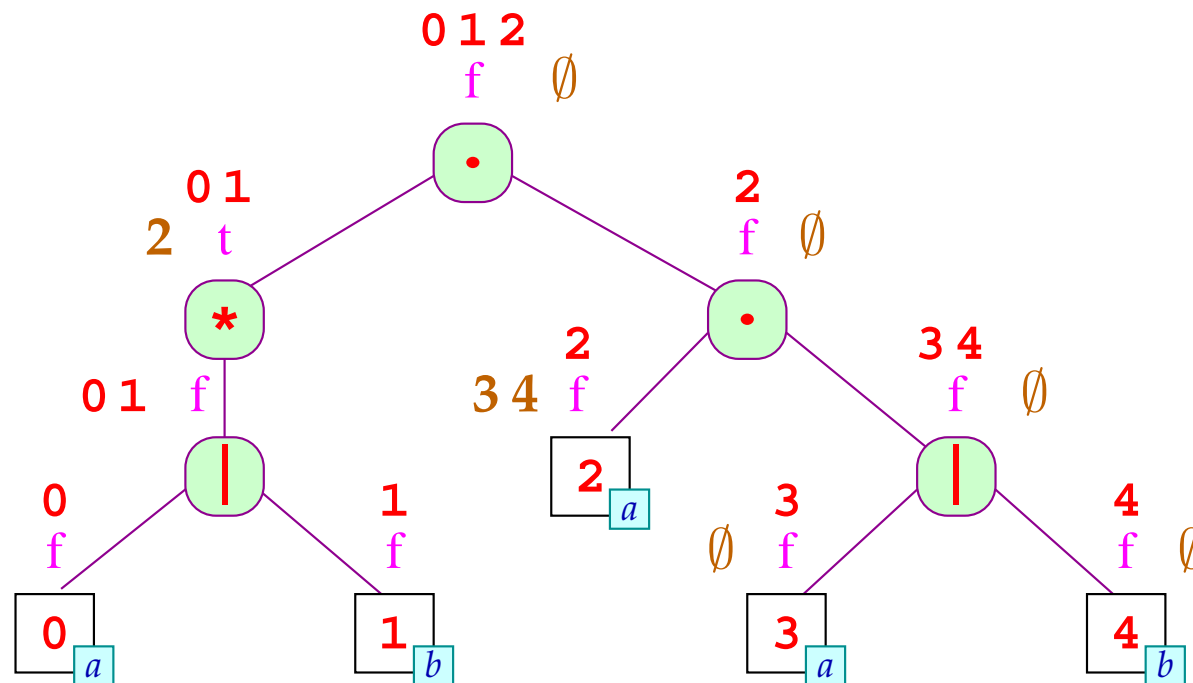


3. Schritt:

Die Menge nächster Blätter:

$$\text{next}[r] = \{i \mid (r \bullet, \epsilon, \bullet \boxed{i \mid x}) \in \delta^*, x \neq \epsilon\}$$

... im Beispiel:

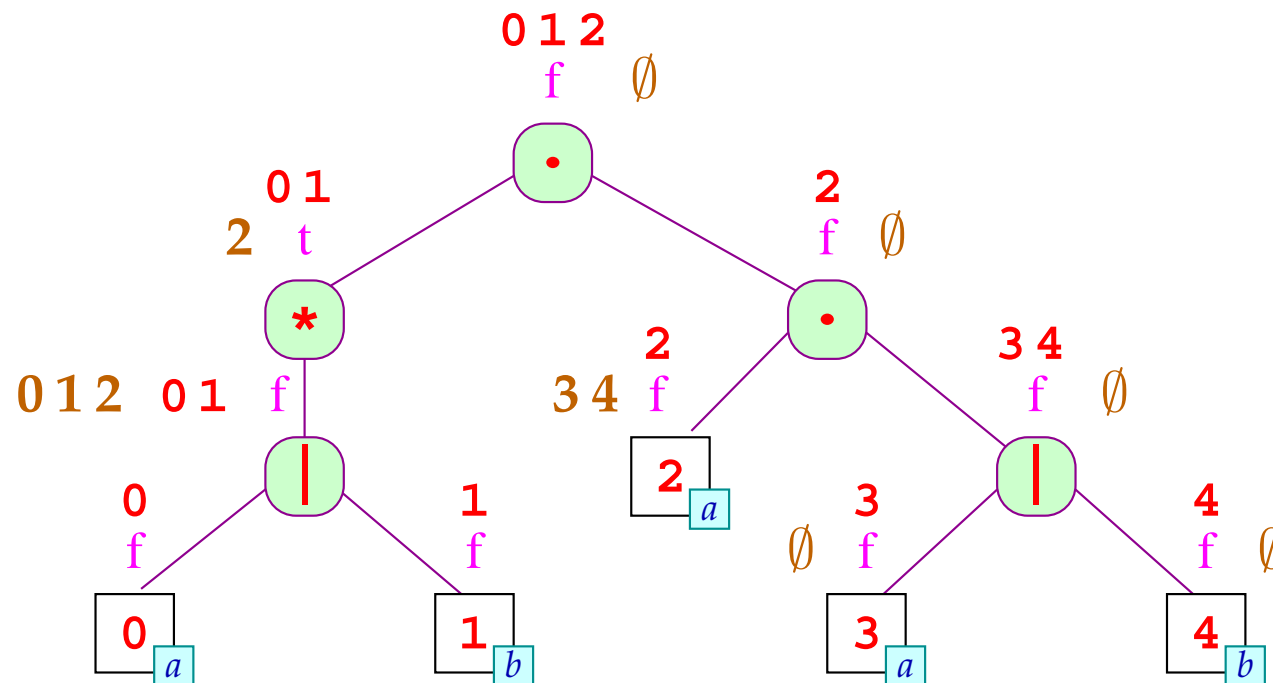


3. Schritt:

Die Menge nächster Blätter:

$$\text{next}[r] = \{i \mid (r \bullet, \epsilon, \bullet \boxed{i \mid x}) \in \delta^*, x \neq \epsilon\}$$

... im Beispiel:

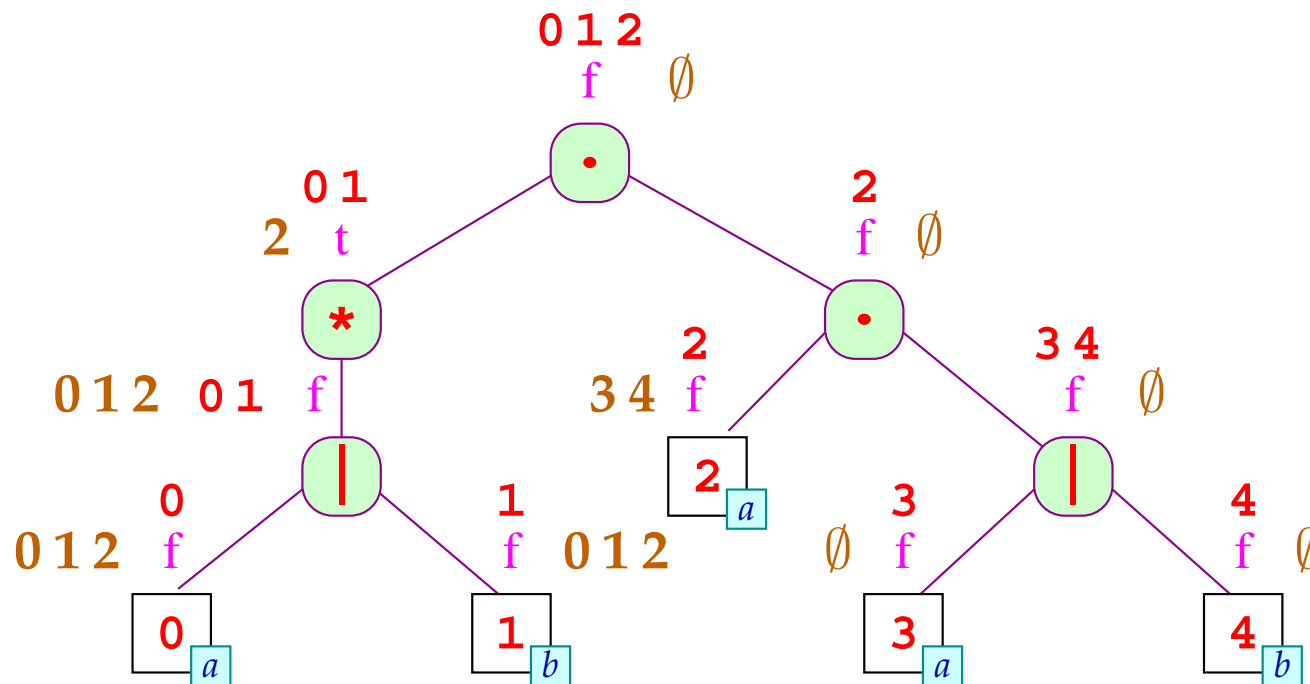


3. Schritt:

Die Menge nächster Blätter:

$$\text{next}[r] = \{i \mid (r \bullet, \epsilon, \bullet \boxed{i \ x}) \in \delta^*, x \neq \epsilon\}$$

... im Beispiel:



Implementierung: DFS pre-order Traversierung ;-)

Für die Wurzel haben wir:

$$\text{next}[e] = \emptyset$$

Ansonsten machen wir eine Fallunterscheidung über den **Kontext**:

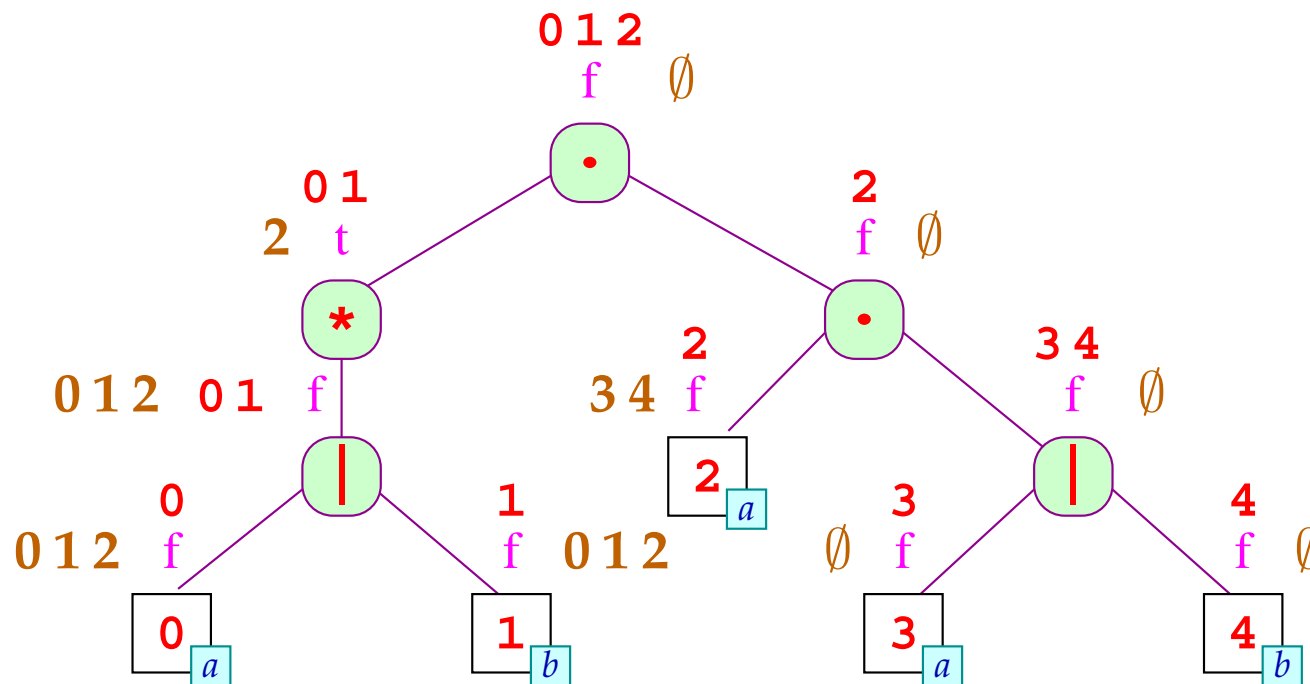
r	Regeln
$r_1 \mid r_2$	$\text{next}[r_1] = \text{next}[r]$ $\text{next}[r_2] = \text{next}[r]$
$r_1 \cdot r_2$	$\text{next}[r_1] = \begin{cases} \text{first}[r_2] \cup \text{next}[r] & \text{falls } \text{empty}[r_2] = t \\ \text{first}[r_2] & \text{falls } \text{empty}[r_2] = f \end{cases}$ $\text{next}[r_2] = \text{next}[r]$
r_1^*	$\text{next}[r_1] = \text{first}[r_1] \cup \text{next}[r]$
$r_1?$	$\text{next}[r_1] = \text{next}[r]$

4. Schritt:

Die Menge **letzter** Blätter:

$$\text{last}[r] = \{i \text{ in } r \mid (\boxed{i \ x} \bullet, \epsilon, r \bullet) \in \delta^*, x \neq \epsilon\}$$

... im Beispiel:

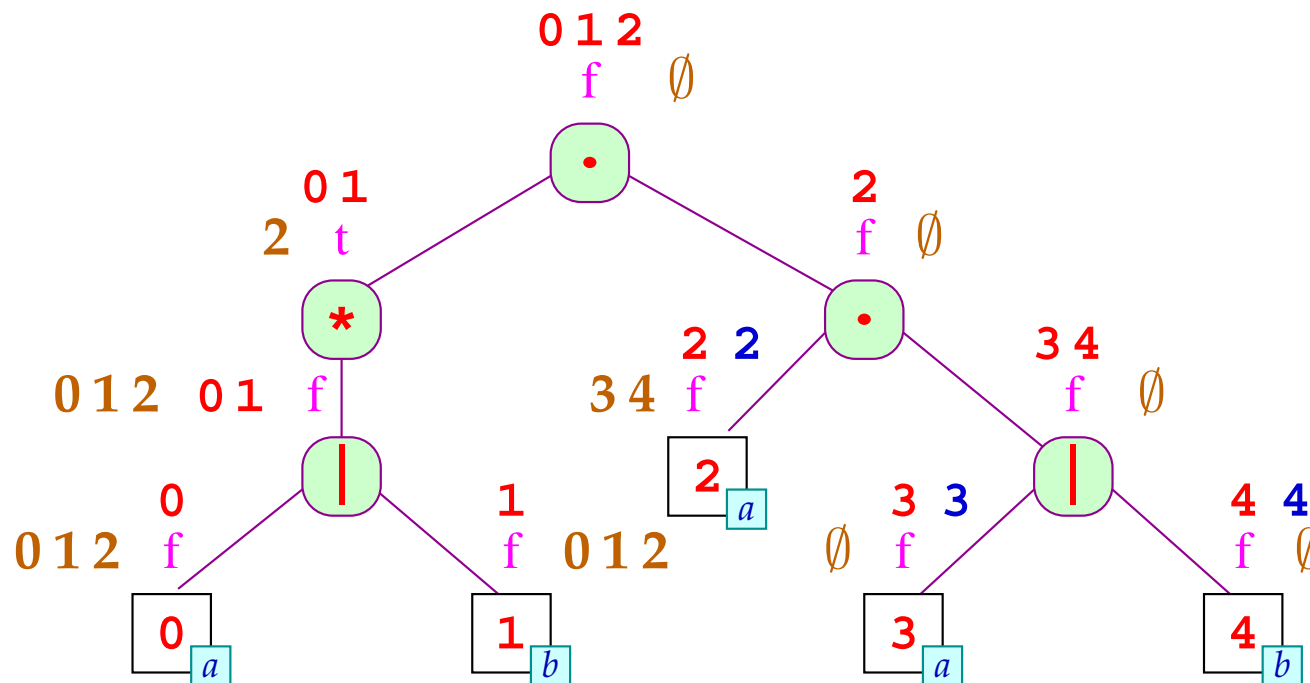


4. Schritt:

Die Menge **letzter** Blätter:

$$\text{last}[r] = \{i \text{ in } r \mid (\boxed{i \ x} \bullet, \epsilon, r \bullet) \in \delta^*, x \neq \epsilon\}$$

... im Beispiel:

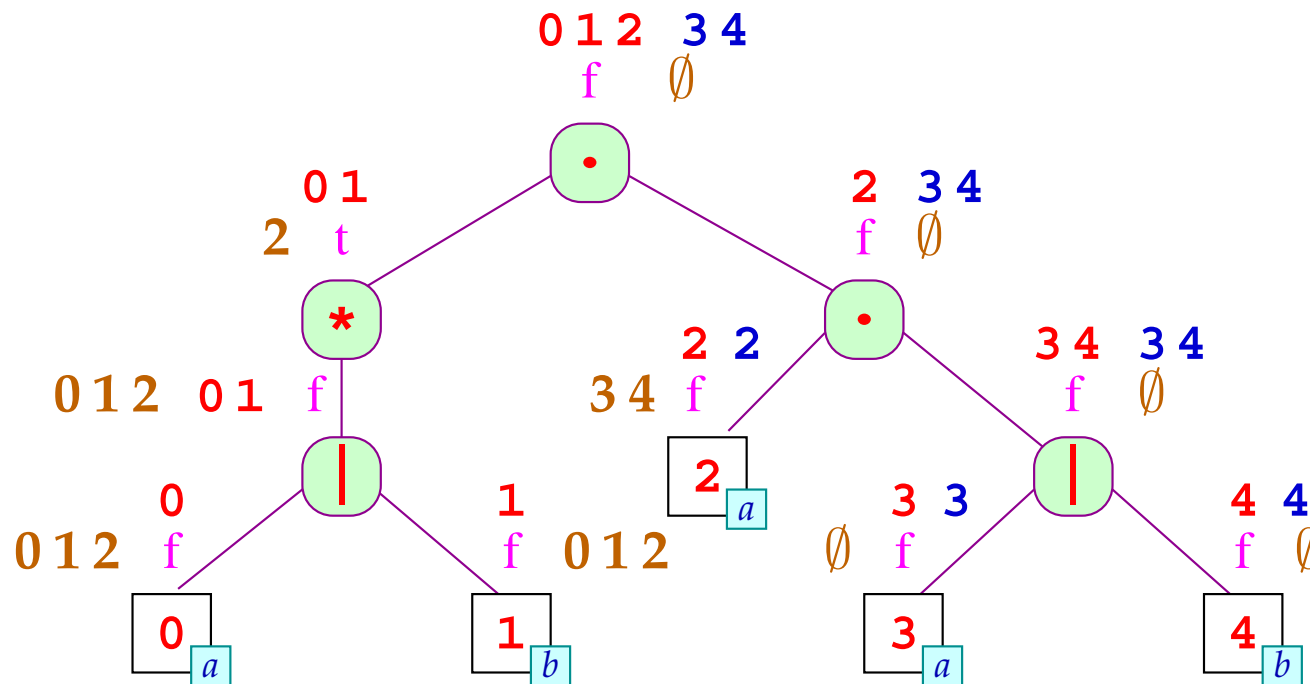


4. Schritt:

Die Menge **letzter** Blätter:

$$\text{last}[r] = \{i \text{ in } r \mid (\boxed{i \ x} \bullet, \epsilon, r \bullet) \in \delta^*, x \neq \epsilon\}$$

... im Beispiel:



Implementierung: DFS post-order Traversierung :-)

Für Blätter $r \equiv \boxed{i \mid x}$ ist $\text{last}[r] = \{i \mid x \neq \epsilon\}$.

Andernfalls:

$$\begin{aligned}\text{last}[r_1 \mid r_2] &= \text{last}[r_1] \cup \text{last}[r_2] \\ \text{last}[r_1 \cdot r_2] &= \begin{cases} \text{last}[r_1] \cup \text{last}[r_2] & \text{falls } \text{empty}[r_2] = t \\ \text{last}[r_2] & \text{falls } \text{empty}[r_2] = f \end{cases} \\ \text{last}[r_1^*] &= \text{last}[r_1] \\ \text{last}[r_1?] &= \text{last}[r_1]\end{aligned}$$

Integration:

Zustände: $\{\bullet e\} \cup \{i\bullet \mid i \text{ Blatt}\}$

Startzustand: $\bullet e$

Endzustände:

Falls $\text{empty}[e] = f$, dann $\text{last}[e]$. Andernfalls: $\{\bullet e\} \cup \text{last}[e]$.

Übergänge:

$(\bullet e, a, i\bullet)$ falls $i \in \text{first}[e]$ und i mit a beschriftet ist;

$(i\bullet, a, i'\bullet)$ falls $i' \in \text{next}[i]$ und i' mit a beschriftet ist.

Den resultierenden Automaten bezeichnen wir mit A_e .