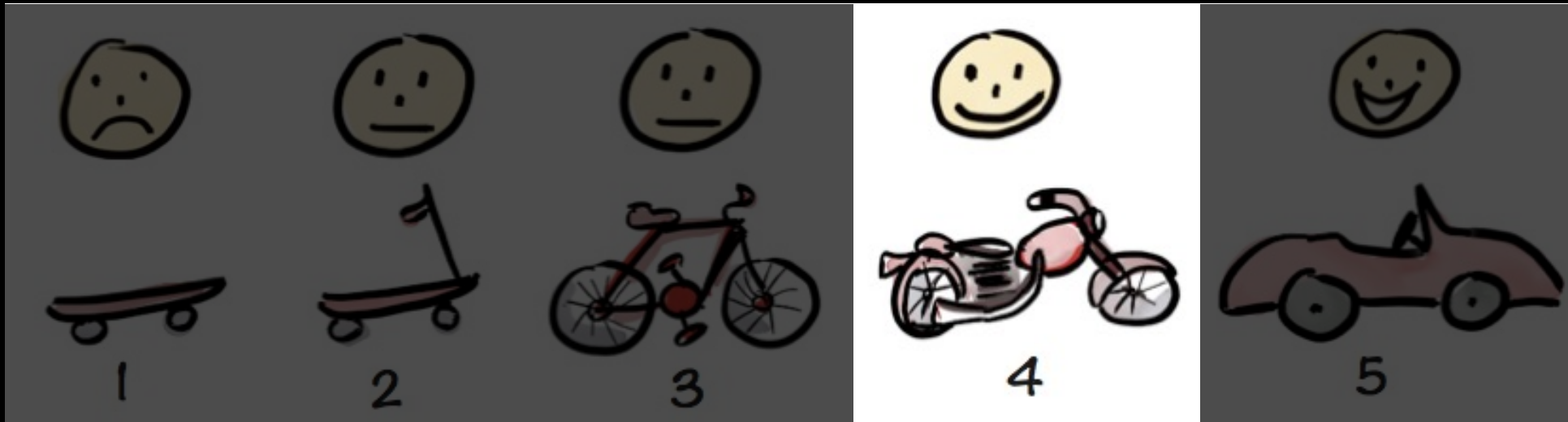


Errungenschaften der letzten Vorlesung

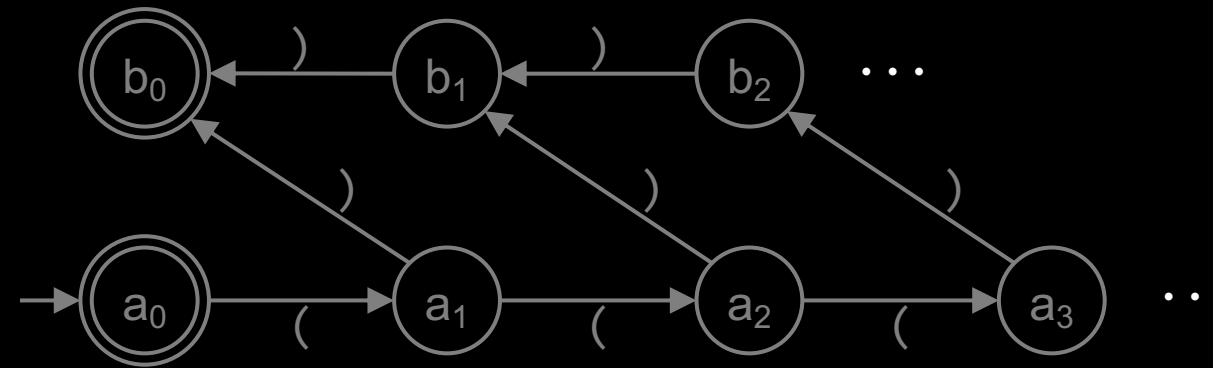


Die Syntax einer formalen Sprache wird mit formalen Grammatiken spezifiziert.

Rechts- und linkslineare Grammatiken lassen sich mit endlichen Automaten verarbeiten. Für andere Grammatiken muss der Syntaxbaum noch manuell ermittelt werden. Alle Token, für die ein endlicher Akzeptor existiert, können automatisch zu einem Tokenizer / Lexer kombiniert werden. Zu jedem NEA kann mithilfe der Potenzmengenkonstruktion ein äquivalenter DEA berechnet werden. Zu jedem DEA kann ein minimaler DEA mit geringstmöglicher Zustandszahl berechnet werden.

Welche Sprachen sind nicht linear?

$$L = \{a^k b^k \mid k \in \mathbb{N}_0\} = \{\varepsilon, ab, aabb, aaabbb, \dots\}$$



*formaler Beweis in Meyer's Foliensatz
Folie 46*

Ein endlicher Automat kann nicht zählen. Er benötigt deshalb für das nächstlängere Wort **jeweils zwei weitere Zustände**. Weil der Parameter k hier nicht **beschränkt** ist, benötigt der Automat zur Erkennung der Sprache L **unendlich viele Zustände**. Wie der Name „endlicher Automat“ bereits suggeriert, hat ein endlicher Automat aber nur endlich viele Zustände. Er kann die Sprache L deshalb nicht erkennen. L ist demnach auch **nicht linear**.

Welche Sprachen sind nicht linear?

$$L = \{a^k b^k \mid k \in \mathbb{N}_0\} = \{\varepsilon, ab, aabb, aaabbb, \dots\}$$

Vorsicht! Wir wissen bereits, dass die Sprache L nicht linear ist. Tatsächlich lässt sich die Sprache aber durch die nachfolgende Grammatik erzeugen:

$a : ' (' b ;$	rechtslineare (Produktions-)Regel
$a : \varepsilon ;$	rechts- und linkslinear
$b : a ') ' ;$	linkslinear

Eine lineare Grammatik ist entweder linkslinear oder rechtslinear. Eine rechtslineare Grammatik enthält ausschließlich rechtslineare Produktionen. Eine linkslineare enthält analog nur linkslineare Produktionen. Im Beispiel oben wurden Produktionen beider Typen verwendet. Die Grammatik ist deshalb kontextsensitiv und nicht linear.

context-free grammar (CFG)

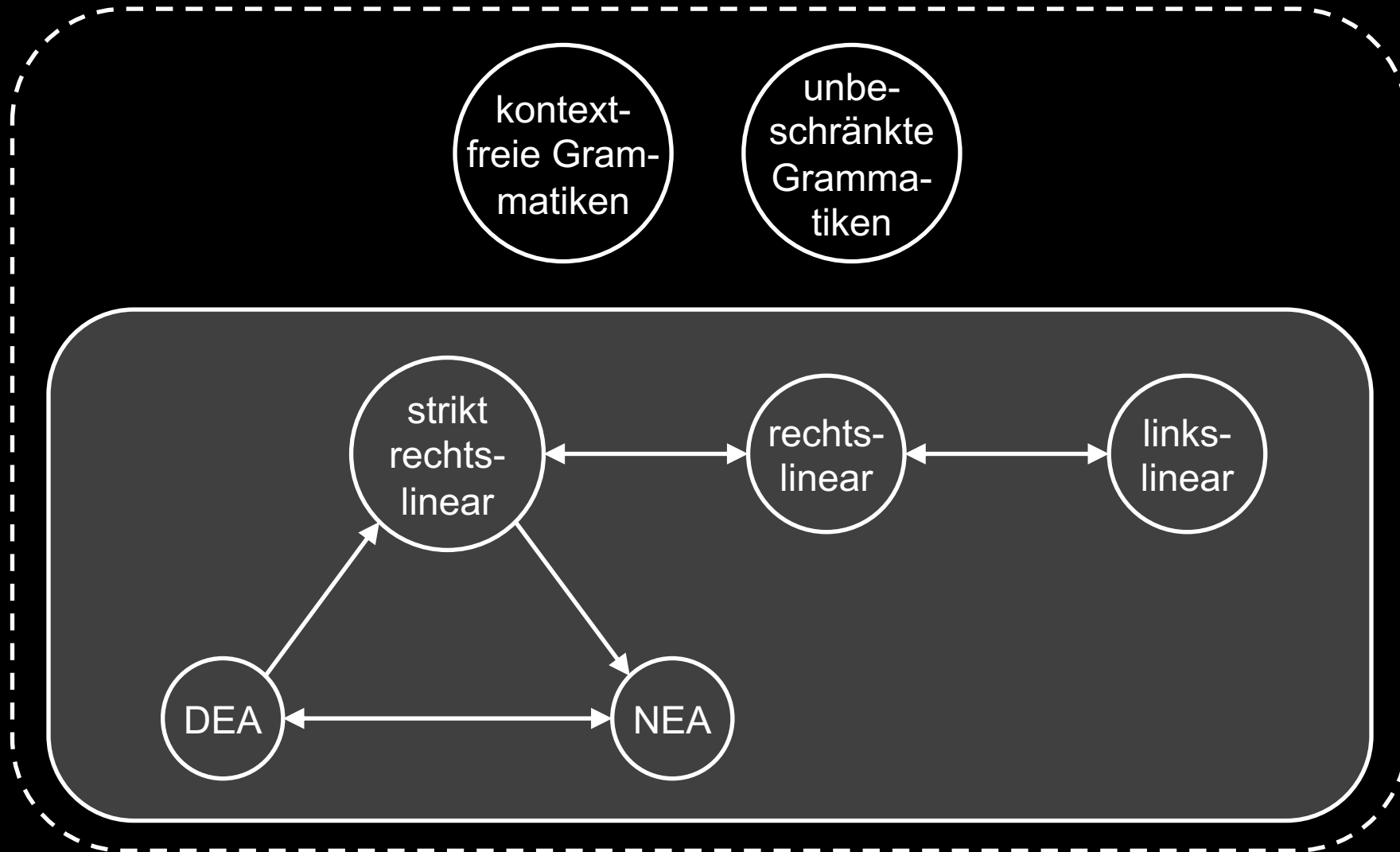
Eine kontextfreie Grammatik oder Typ-2-Grammatik ist eine Grammatik $G = (V_N, V_T, S, P)$, bei der alle Produktionen die folgende Form haben:

$$X \rightarrow w \quad \text{mit } X \in N \text{ und } w \in (N \cup T)^*$$

Eine Sprache heißt kontextfrei, wenn Sie von einer kontextfreien Grammatik erzeugt werden kann.

Alle Regeln, die wir bislang allgemein für Grammatiken getroffen haben, gelten selbstverständlich auch für diese Untergruppe.

Modellübersicht

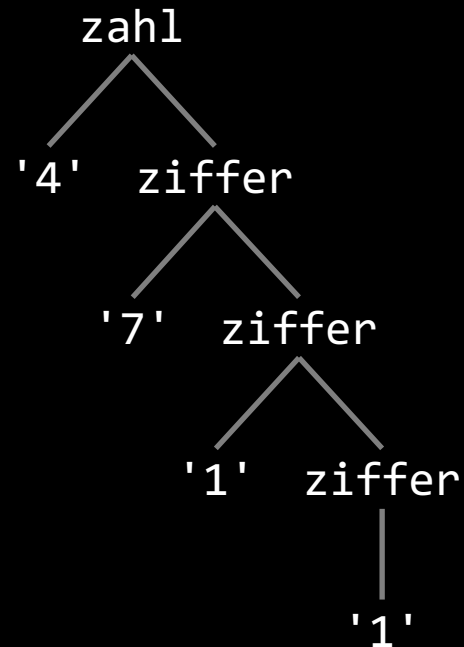


Unser Adventskalender

Typ	Name	Erlaubte Produktionen	Akzeptierender Automat	(Trennendes) Beispiel
3		$P \subseteq N \times (\{\varepsilon\} \cup T^* \cup T^* N)$ oder $P \subseteq N \times (\{\varepsilon\} \cup T^* \cup NT^*)$	endlicher Automat	$L = \{a^k b^l \mid k, l \in \mathbb{N}_0\}$
2	kontextfrei	$P \subseteq N \times (N \cup T)^*$		$L = \{a^k b^k \mid k \in \mathbb{N}_0\}$

Wiederholung

```
zahl : '1' ziffer | '2' ziffer | ... | '9' ziffer  
      | '0' | '1' | '2' | ... | '7' | '8' | '9' ;  
ziffer : '0' ziffer | '1' ziffer | ... | '9' ziffer  
         | '0' | '1' | '2' | ... | '7' | '8' | '9' ;
```



Syntaxbäume bei
Typ-3 Sprachen sind
immer listenartig

Wiederholung

$$e : e \text{ '+' } e \mid e \text{ '-' } e \mid e \text{ '*' } e \mid e \text{ '/' } e ;$$

$$e : \text{'0'} \mid \text{'1'} \mid \text{'2'} \mid \dots \mid \text{'7'} \mid \text{'8'} \mid \text{'9'} ;$$

$$\overset{1}{e} \Rightarrow e + \overset{7}{e} \Rightarrow 2 + \overset{3}{e} \Rightarrow 2 + \overset{8}{e} * e \Rightarrow 2 + 3 * \overset{9}{e} \Rightarrow 2 + 3 * 4$$

$$\overset{1}{e} \Rightarrow e + \overset{3}{e} \Rightarrow e + e * \overset{9}{e} \Rightarrow e + \overset{8}{e} * 4 \Rightarrow e + 3 * 4 \overset{7}{\Rightarrow} 2 + 3 * 4$$

Bei Typ-2 Sprachen gibt es immer

eine Linksableitung und eine Rechtsableitung.

Ableitung

Sei $G = (V_N, V_T, S, P)$ eine kontextfreie Grammatik mit $X \rightarrow w \in P$ und $u, u' \in (N \cup T)^*$.

Falls $u \in T^*$, heißt ein Ableitungsschritt $uXu' \Rightarrow uwu'$ mit Produktion $X \rightarrow w \in P$ Linksableitungsschritt, geschrieben $uXu' \xRightarrow{l} uwu'$.

$w_0 \Rightarrow w_1 \Rightarrow \dots \Rightarrow w_k$ heißt Linksableitung(sfolge), falls $w_0 \xRightarrow{l} w_1 \xRightarrow{l} \dots \xRightarrow{l} w_k$, geschrieben $w_0 \xRightarrow{*} w_k$.

Analog: Rechtsableitungsschritt \xRightarrow{r} und Rechtsableitung $\xRightarrow{*}$.

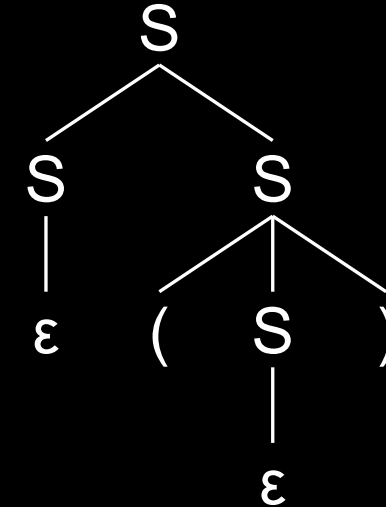
example: derivation

$$G = (\{S\}, \{ (,) \}, S, \{ S \rightarrow \varepsilon \mid (S) \mid SS \})$$

Linksableitung: $S \Rightarrow \underline{S}S \Rightarrow \underline{S} \Rightarrow (\underline{S}) \Rightarrow ()$

Rechtsableitung: $S \Rightarrow S\underline{S} \Rightarrow S(\underline{S}) \Rightarrow \underline{S}() \Rightarrow ()$

Weder noch: $S \Rightarrow S\underline{S} \Rightarrow \underline{S}(S) \Rightarrow (\underline{S}) \Rightarrow ()$



Alle hier gezeigten Ableitungen liefern den gleichen Ableitungsbaum.

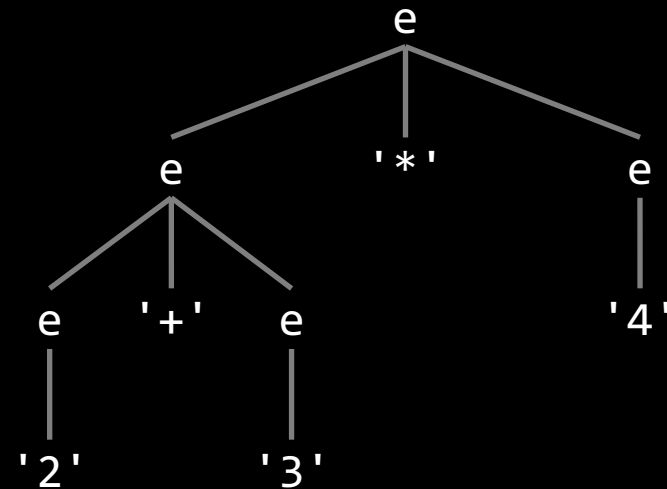
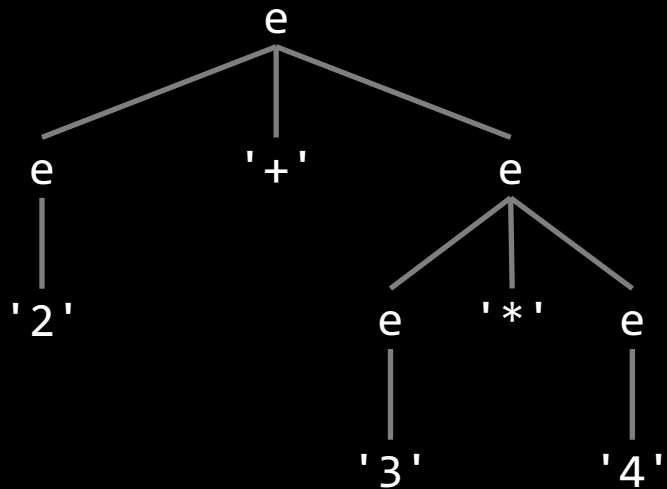
Ein Ableitungsbaum ist demnach „allgemeiner“ als eine Ableitung.

Aus jedem Ableitungsbaum ergibt sich genau eine Linksableitung und analog genau eine Rechtsableitung.

Also: Zu jeder Ableitung gibt es eine „äquivalente“ Links- bzw. Rechtsableitung.

Wiederholung

$$e : e \text{ '+' } e \mid e \text{ '-' } e \mid e \text{ '*' } e \mid e \text{ '/' } e ;$$

$$e : \text{'0'} \mid \text{'1'} \mid \text{'2'} \mid \dots \mid \text{'7'} \mid \text{'8'} \mid \text{'9'} ;$$


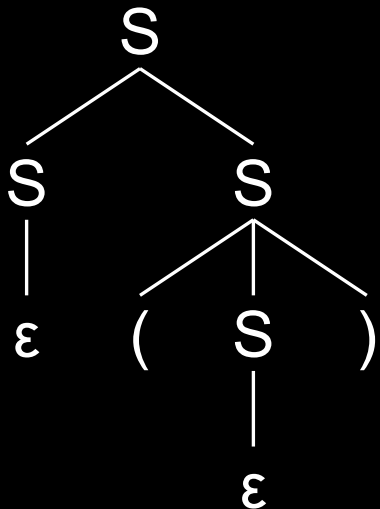
Eine Grammatik heißt **eindeutig**, wenn es für jedes Wort $w \in L(G)$ genau eine Linksableitung gibt. Nicht eindeutige Grammatiken nennt man auch **mehrdeutig**.

Mehrdeutigkeit

Eine Typ-2-Grammatik ist **mehrdeutig**, falls es für ein Wort zwei verschiedene Ableitungsbäume (Linksableitungen, Rechtsableitungen) gibt.

Sonst heißt die Grammatik **eindeutig**.

Eine kontextfreie Sprache ist (inhärent) **mehrdeutig**, falls jede sie erzeugende T2G mehrdeutig ist.

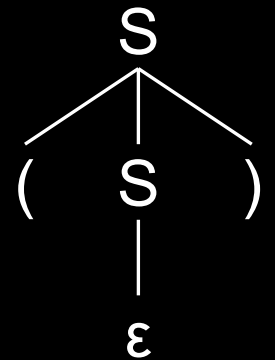


Die Grammatik $G = (\{S\}, \{(\,)\}, S, \{S \rightarrow \varepsilon \mid (S)SS\})$ ist mehrdeutig.

Die erzeugte formale Sprache ist aber nicht inhärent mehrdeutig.

Eine eindeutige Grammatik G' mit $L(G) = L(G')$ wäre

$$G' = (\{S, T\}, \{(\,)\}, S, \{S \rightarrow \varepsilon \mid ST, T \rightarrow (S)\})$$



Offene Fragen

- ~~1. Nicht jedes Token darf an jeder Stelle stehen. Wie lässt sich das regulieren?~~
- ~~2. Sind die von endlichen Automaten akzeptierten Sprachen identisch zu den von Grammatiken erzeugten Sprachen?~~
- ~~3. Lassen sich alle NEAs in DEAs umwandeln?~~
- ~~4. Sind rechtslineare und linkslineare Grammatiken verschieden?~~
- ~~5. Sind unbeschränkte Grammatiken mächtiger als die linearen?~~
- ~~6. Gibt es “kompaktere” Automaten?~~
7. Kann man die Erstellung eines endlichen Automaten automatisieren?

Wie sieht unsere gegenwärtige Implementierung aus?

Beispiel: StateMachineTableNumber.java

```
@Override
public void initStateTable() {
    compiler.State firstDigit = new compiler.State("firstDigit");
    firstDigit.addTransition('0', "number0");
    firstDigit.addTransition('1', "nextDigit");
    firstDigit.addTransition('2', "nextDigit");
    firstDigit.addTransition('3', "nextDigit");
    firstDigit.addTransition('4', "nextDigit");
    firstDigit.addTransition('5', "nextDigit");
    firstDigit.addTransition('6', "nextDigit");
    firstDigit.addTransition('7', "nextDigit");
    firstDigit.addTransition('8', "nextDigit");
    firstDigit.addTransition('9', "nextDigit");
    m_stateMap.put("firstDigit", firstDigit);

    compiler.State number0 = new compiler.State("number0");
    m_stateMap.put("number0", number0);

    compiler.State nextDigit = new compiler.State("nextDigit");
    nextDigit.addTransition('0', "nextDigit");
    nextDigit.addTransition('1', "nextDigit");
    nextDigit.addTransition('2', "nextDigit");
    nextDigit.addTransition('3', "nextDigit");
    nextDigit.addTransition('4', "nextDigit");
    nextDigit.addTransition('5', "nextDigit");
    nextDigit.addTransition('6', "nextDigit");
    nextDigit.addTransition('7', "nextDigit");
    nextDigit.addTransition('8', "nextDigit");
    nextDigit.addTransition('9', "nextDigit");
    m_stateMap.put("nextDigit", nextDigit);
}
```