

Lexikalische Analyse

- **aus Einführung:**

Scanner

- Programmteil, der lexikalische Analyse durchführt.
- Unterprogramm des Parsers (Syntaxanalyse)
- übergibt Parser Token (Bezeichner der Lexemklassen):
IDENT, NUMBER, ...

Bsp.: bez1 = 354*bez2;
 -> IDE"="NUM"*"IDE";"

- für Parser Token ausreichend
 - bei Codeerzeugung müssen aber Bezeichnernamen und Zahlwerte bekannt sein
- > da Scanner einziger Teil des Compilers, der Quellprogramm durchläuft, liefert er weitere Informationen als Attribute

Token & Attribute

- **Aufgaben des Scanners:**

- Zerlegung des Quellprogramms in Folge von Paaren aus
 - Token Bezeichnung der Lexemklassen
 - Attribute Lexemwerte
 - für Bezeichner: Name
 - für Zahlen: Werte
 - nur, wenn nötig

Bsp: bez1 = 354*bez2;

-> (IDENT,bez1) ("=",) (NUMBER,354) ("*",)
 (IDENT,bez2) (";",)

- **Verwendung der Scannerausgabe:**

- Token: Parser
- Attribute: Codeerzeugung

Analyseverfahren

- **bisher:**
Beschreibung der Lexemstrukturen durch reguläre Ausdrücke
- **RA**
 - Def.: ε, a **sind RAs für jedes Eingabezeichen** a
Seien a, b reguläre Ausdrücke. Dann auch
 - $a|b$: a oder b
 - ab : a konkateniert mit b
 - a^* : ε oder a kein mal, einmal oder beliebig oft mit sich selbst konkateniert (Kleene Stern)
 - Bsp.: Siehe Beschreibung von Bezeichnern (Erweiterungen)
 - deklarative Beschreibung
 - können in Prologprogrammierung verwendet werden
 - nicht in imperativen oder objektorientierten Sprachen

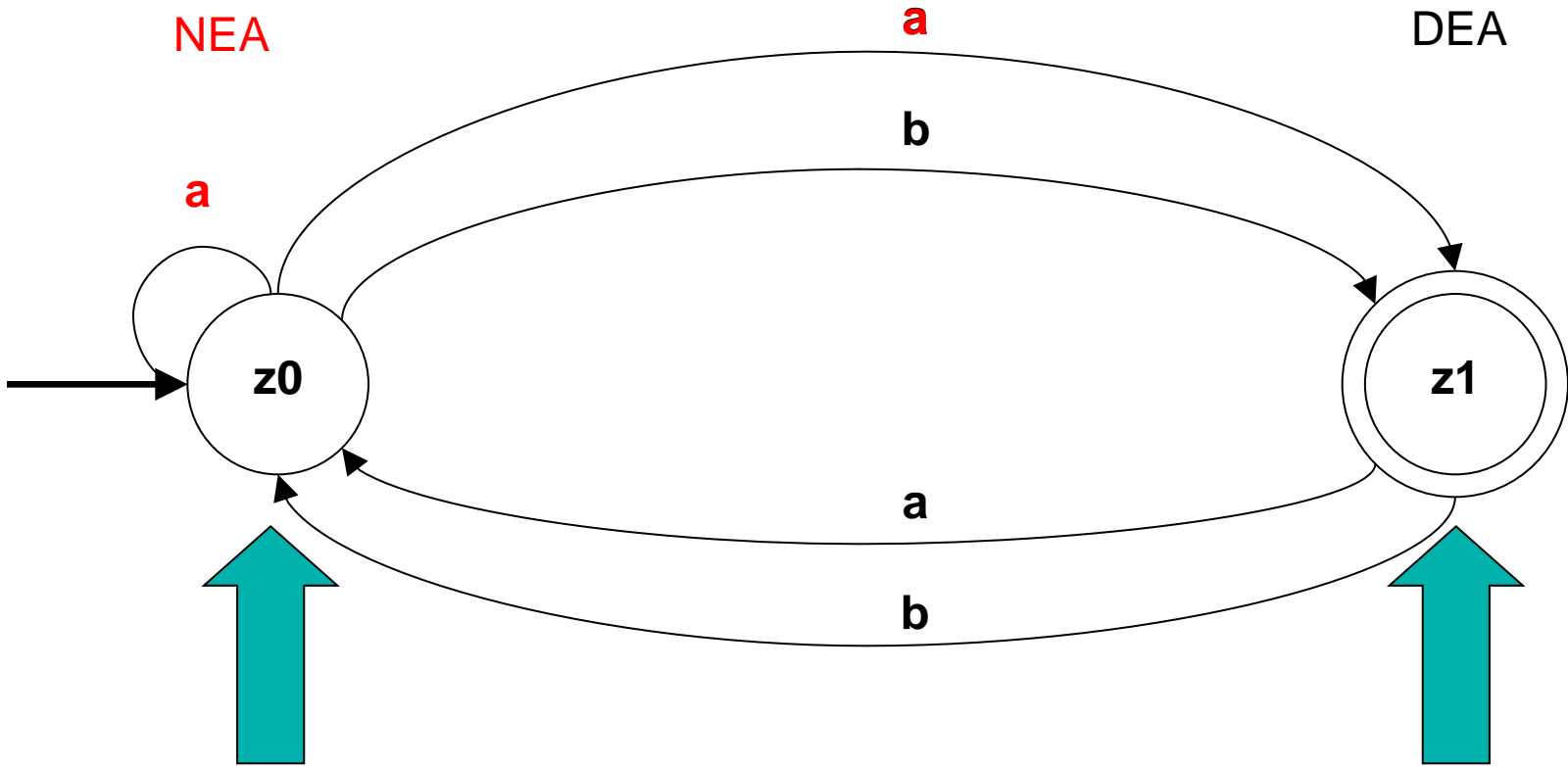
Scannerprogrammierung in C oder Java?

- **Implementierung durch DEA (det. endl. Automaten)**

Endliche Automaten (Wdh.)

nicht deterministisch
NEA

deterministisch
DEA



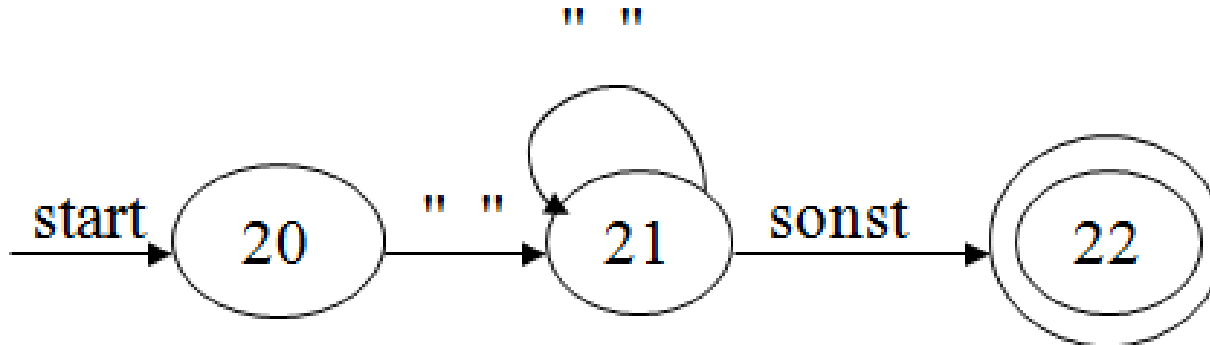
Anfangszustand

Endzustand

RAs für DEA & NEA?

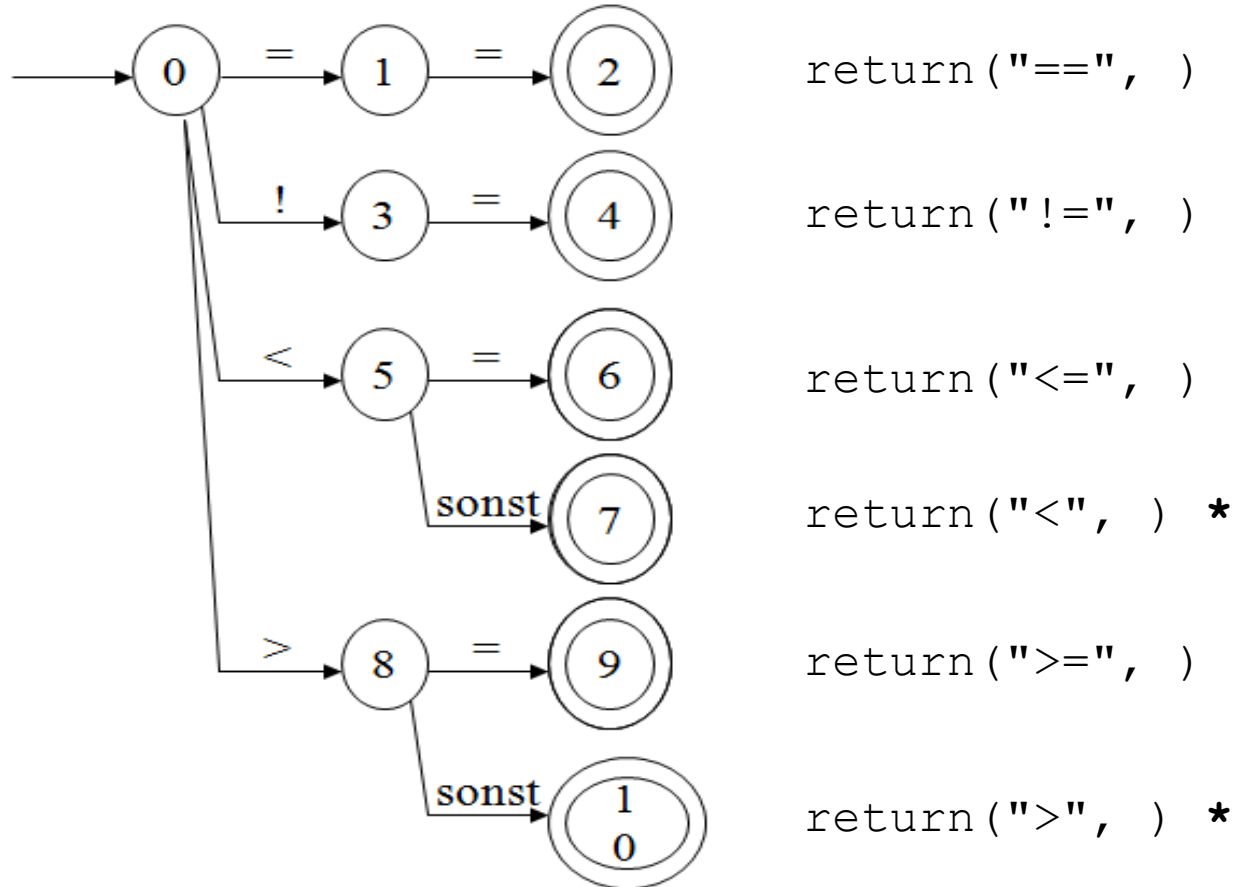
Händische Scannerkonstruktion (EAs)

- pro Lexemklasse ein oder mehrere erkennende EAs
- Transformation in Programmteile
- Sukzessives Durchlaufen mit **look-ahead** auf der Eingabe (Lexemtrennung!), da sonst kein Zurücksetzen möglich
- **Bsp.:** **EA für Blank**



- Hier: keine Rückgabe eines Tokens an den Parser (Überlesen)
- Scanner muss bis zum nächsten Token weiterarbeiten
- Implement. des Scanners durch Umsetzung der EAs in PSn

EA für CompOPs



`return ("==",)`

`return ("!=",)`

`return ("<=",)`

`return ("<",) *`

`return (">=",)`

`return (">",) *`

- **sonst:** Übergang mit jedem anderen Zeichen
- *****: Rücksetzen des look-ahead Zeigers um 1 Zeichen

automatische Scannerkonstruktion

- Grundlage: **RAs -> DEAs**
- **in javacc:**

Konfigurationsdatei

Test.jj $\xrightarrow{\text{javacc}}$

... ,
TestTokenManager.java

Implementierung des EAs
wird automatisch übersetzt
bei Übersetzung von
Test.java

- Spezifikation der RAs in Konfigurationsdatei -> Einführung
- **unter Unix:**
Scanner- und Parser-Erzeuger getrennt in **lex, yacc**
-> Konfigurationsdatei von lex enthält nur RAs

DEA-Methode

- **Ziel:**

- automatische Konstruktion RA -> DEA
- Implementierung von DEA in C oder Java

- **Umweg:** RA -> NEA -> DEA automatisch konstruierbar

- **1. RA -> NEA**

Konstruktion von **Kleene** induktiv über Aufbau der RAs
siehe Theoretische Informatik oder Literatur

- **2. NEA -> DEA**

Potenzmengenkonstruktion PMK

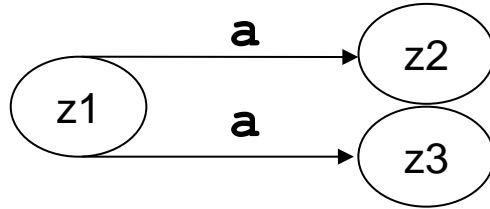
- **Problem**

DEA kann exponentiell viele Zustände haben
auch die hier vorgestellte verbesserte Version

Wdh.: PMK

• Idee:

– NEA



$$(z1, a, z2), (z1, a, z3) \in dRel$$

– DEA



$$dFun(\{z1\}, a) = \{z2, z3\}$$

• Formalisierung:

– Geg.: NEA $N = (States, \Sigma, dRel, z0, E)$

– Konstruiere DEA $D = (DStates, \Sigma, dFun, \{z0\}, ED)$ mit

$$DStates = P(States)$$

$$dFun(T, a) = \bigcup_{z \in T} \{z' \in States \mid (z, a, z') \in dRel\}$$

$$ED = \{T \in P(States) \mid \exists e \in T : e \in E\}$$

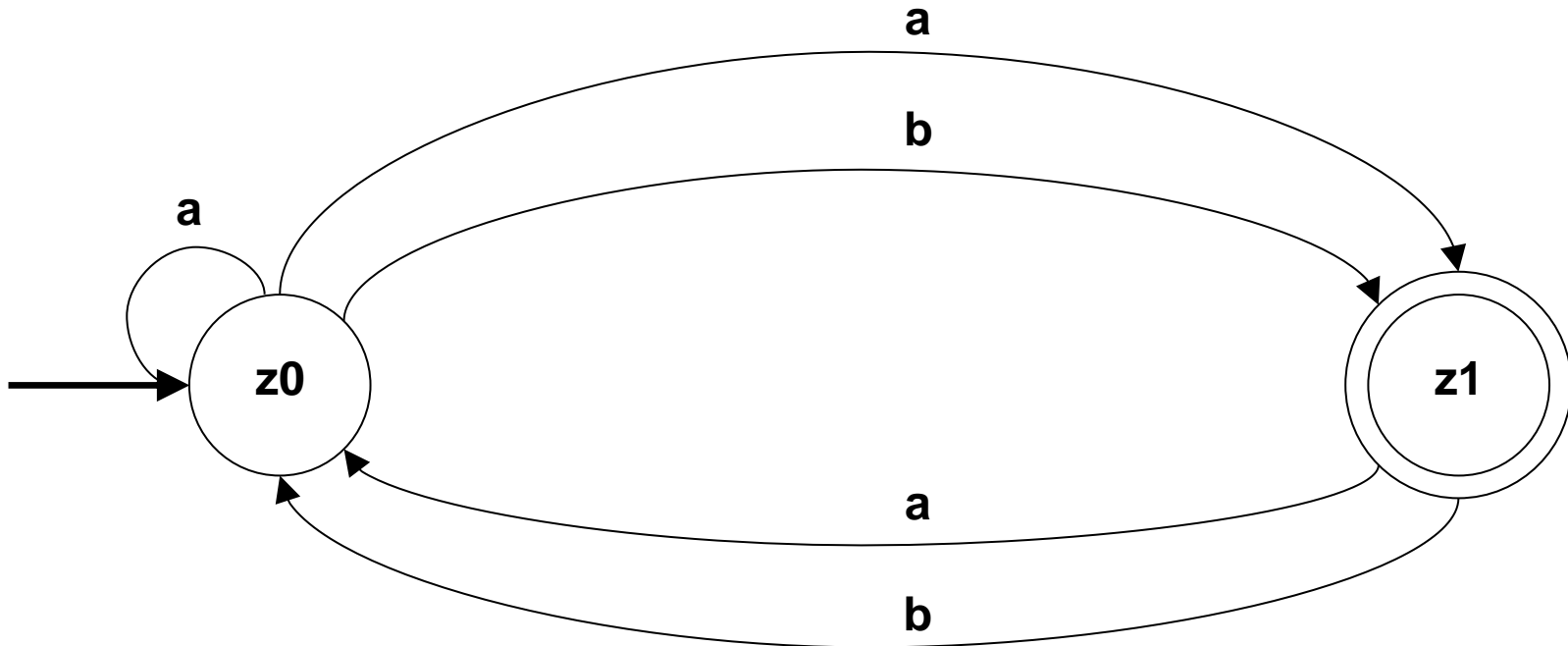
PMK-Algorithmus

- Füge nur erreichbare Zustände zu DStates ($\subseteq P(\text{States})$) so in TIWS gemacht:

```
DStates = {{z0}[ohne Marke]};
while T[ohne Marke] ∈ DStates
{
    markiere T in DStates;
    for each a ∈ Σ
    {
        U = {}; //U ist Folgezustand von T mit a
        for each z ∈ T U ∪ = {z' | (z, a, z') ∈ dRel};
        if (U ∉ DStates) DStates ∪ = {U[ohne Marke]};
        dFun(T, a) = U;
    }
}
```

- In ED alle Mengen aus DStates, die Endzustand aus E enthalten

PMK-Beispiel

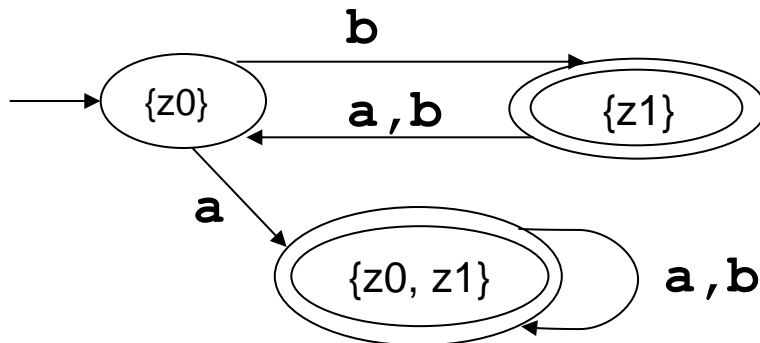


Formalisierung des NEA?

PMK-Beispiel

- geg.: NEA ($\{z_0, z_1\}, \{a, b\}, dRel, z_0, \{z_1\}$)
 mit $dRel = \{(z_0, a, z_0), (z_0, a, z_1), (z_0, b, z_1), (z_1, a, z_0), (z_1, b, z_0)\}$

- berechne: DEA ($DStates, \{a, b\}, dFun, \{z_0\}, ED$) Marke \sim
 - $DStates = \{\{z_0\}\} \rightarrow$
 - $dFun(\{z_0\}, a) = \{z_0, z_1\}$
 - $dFun(\{z_0\}, b) = \{z_1\}$
 - $DStates = \{\{z_0\}^\sim, \{z_0, z_1\}, \{z_1\}\} \rightarrow$
 - $dFun(\{z_1\}, a) = \{z_0\}$
 - $dFun(\{z_1\}, b) = \{z_0\}$
 - $DStates = \{\{z_0\}^\sim, \{z_0, z_1\}, \{z_1\}^\sim\} \rightarrow$
 - $dFun(\{z_0, z_1\}, a) = \{z_0, z_1\}$
 - $dFun(\{z_0, z_1\}, b) = \{z_0, z_1\}$
 - **$DStates = \{\{z_0\}^\sim, \{z_0, z_1\}^\sim, \{z_1\}^\sim\}$ $ED = \{\{z_1\}, \{z_0, z_1\}\}$**



$w = aaa$ in L
 $w = bb$ nicht in L

NEA-Methode

- **Problem bei DEA-Methode:**
 - u.U. exponentiell viele Zustände (Beispiel?)
- **Idee bei NEA-Methode:**
 - Simulation der PKM zur Laufzeit, d.h. eines Laufes durch Potenzmengenautomat bzgl. einer Eingabe
- **NEA-Methode:**
 - 2 Stacks zum Speichern einer Zustandsmenge und der jeweiligen Folgezustände
- **Eingaben:**
 - NEA $N = (\text{States}, \Sigma, \text{dRel}, z_0, E)$
 - Wort w
- **Ausgaben:**
 - Ja, falls w in zu erkennender Sprache
 - Nein, sonst

NEA-Methoden-Algorithmus

```
{
    T = {z0};
    a = nextchar(w);

    while (a != eof)
    {
        U = {}; //U ist Folgezustand von T mit a
        for each z ∈ T U ∪= {z' | (z,a,z') ∈ dRel};
        T = U;
        a = nextchar(w);
    }

    if (T ∩ E != {}) return "JA";
    else return "Nein";
}
```

NEA-Methoden-Beispiel

- geg.: NEA von Folie 4 ($\{z_0, z_1\}, \{a, b\}, dRel, z_0, \{z_1\}$)
mit $dRel = \{(z_0, a, z_0), (z_0, a, z_1), (z_0, b, z_1), (z_1, a, z_0), (z_1, b, z_0)\}$
- Lauf für $w = aaa$:
 - $T = \{z_0\}$ $nextchar(w) = a$
 - $T = \{z_0, z_1\}$ $nextchar(w) = a$ $T \cap E = \{z_1\}$
 - $T = \{z_0, z_1\}$ $nextchar(w) = a$ **-> "Ja"**
 - $T = \{z_0, z_1\}$ $nextchar(w) = eof$
- Lauf für $w = bb$:
 - $T = \{z_0\}$ $nextchar(w) = b$
 - $T = \{z_1\}$ $nextchar(w) = b$ $T \cap E = \emptyset$
 - $T = \{z_0\}$ $nextchar(w) = eof$ **-> "Nein"**

Komplexitätsvergleich

- **|States|** : Anzahl der Zustände des NEA
- **|w|** : Länge von w

	Platz	Zeit
DEA-Methode	$O(2^{ \text{States} })$	$O(w)$
NEA-Methode	$O(\text{States})$	$O(\text{States} * w)$

- **zur Zeit-Komplexität des NEA-Verfahrens**
 - while-Schleife wird |w| mal durchlaufen
 - for-Schleife wird |T| ($\leq |\text{States}|$) mal durchlaufen
 - $\rightarrow |w| * |T|$ Durchläufe $\leq O(|w| * |\text{States}|)$

Allgemeines Problem: Lexemende?

- `b`, `be`, `bez`, `bez1` alles Bezeichner
- Konvention: "**Principle of the Longest Match**"
d.h. Zeichenreihe bestimmt ein Lexem, wenn keine Verlängerung zu einer Lexemklasse gehört
- manchmal inkorrekt
- Bsp.: (Fortran)
`363`, `363.E1`, `.EQ.` alles Lexeme
Bei `363.EQ.363` Principle of the Longest Match
^ **lexikalischer Fehler**
- Lösung: Look-ahead von 3 Zeichen (mindert Effizienz)
- in javacc Standard-Look-ahead von 1

Zusammenfassung

- Scanner zerlegt Quellprogramm in Token und Attribute
- Struktur der Lexeme durch reg. Ausdrücke beschreibbar
- Angabe der Struktur in javacc durch reguläre Ausdrücke
- Analyse durch endliche Automaten, hierzu:
 - **Kleene:** Transformation der regulären Ausdrücke in NEA
 - **PMK:** Transformation NEA -> DEA
DEA-Methode -> **hoher Platzbedarf**
 - **Laufzeit-PMK:** Transformation NEA -> DEA zur Laufzeit
NEA-Methode -> **höherer Laufzeitbedarf**