

Zusammenfassung zur Vorlesung Basismodul Computerlinguistik

Tarjan-Tabelle & Einfügen und Löschen im Lexikonautomaten

09.12.2021

- 1 Tarjan-Tabelle
- 2 Einfügen und Löschen eines Wortes im Lexikonautomaten

Effizientes Speichern von dürrtig besetzten Tabellen (sparse tables)

Grundidee

- Wir speichern die Zustände und Übergänge in einer Liste
- Wir speichern nur tatsächlich existierende Übergänge
- Die Liste hat möglichst wenige Lücken.

$$\Sigma = \{a, b, c, d\} \quad Q = \{1, 2, 3, 4, 5, 6\}$$

	1	2	3	4	5	6
a	2					3
b						4
c		3				5
d	6	4				

Tarjan-Tabelle

Bausteine

- Ein 'Kamm' in Länge des Alphabets als Eintrag in die Liste
- Tabelle mit 2 Spalten, um Label und Zielzustand zu speichern
- besteht aus **Zellen**, also **Zustandszellen** und **Übergangszellen**

$$\Sigma = \{a, b, c, d\}$$

$$\{a:1, b:2, c:3, d:4\}$$

$$1: a \rightarrow 2$$

$$d \rightarrow b$$

$$2: c \rightarrow 3$$

$$d \rightarrow 4$$

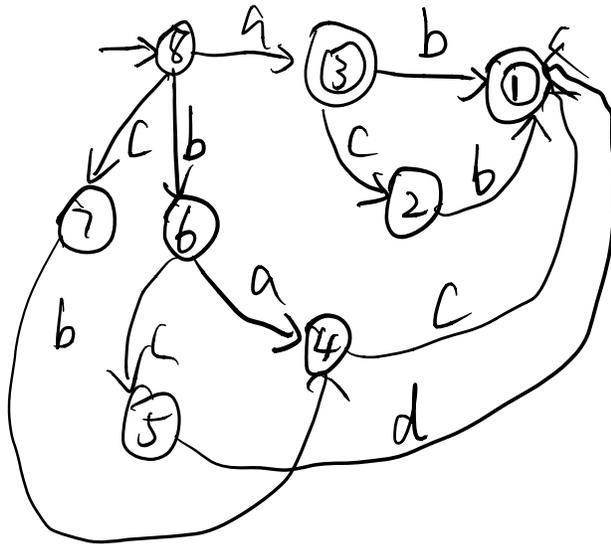
		1	
1	a	2	
2			
3			
4	d	b	

			2
1	-		
2			
3	c	3	
4	d	4	

Tarjan-Tabelle - Beispiel

a
ab
acb
bac
bcd
cbc

$\Sigma = \{a, b, c, d\}$
 $\{a:1, b:2, c:3, d:4\}$



1	F
2	NF
3	F
b	1
b	1
c	2
4	NF
5	NF
7	NF
c	1
b	7
d	1
6	NF
a	7
	8
c	8
8	NF
a	5
b	13
c	9

3
b 1
c 2

Einfügen und Löschen im Lexikonautomaten

Um aus einer **unsortierten** Liste einen Lexikonautomaten zu bauen.
⇒ Das entspricht dem wiederholten **Einfügen eines Wortes** in einen Lexikonautomaten.

- Der Automat ist vorher **minimiert**
- Das Wort wird eingefügt oder gelöscht
- Der Automat ist wieder **minimiert**

Einfügen eines Wortes

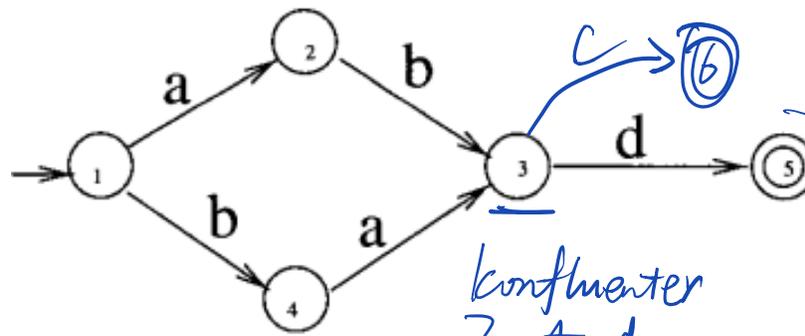
Wie fügen wir ein neues Wort im sortierten Fall hinzu?

- Finde das **gemeinsame Präfix** und dessen **Split State**
- Füge das neue **Suffix** hinzu
- Minimiere den Pfad

Problem im unsortierten Fall: Es kann mehr als einen eingehenden **Pfad** zum Split State geben

- Dann würde das Suffix an **mehrere** Präfixe angehängt
- Dann würde man als Nebeneffekt ungewollt weitere Wörter einfügen.

abc



→ abc

ba c ungewollt

Figure: Quelle Daciuk Paper

Einfügen eines Wortes

Lösung: Es muss der Pfad des Wortes zuerst **isoliert** werden, um die Nebeneffekte zu verhindern.

→ **Konfluente** Zustände auf dem Präfixpfad müssen geklont / gesplittet werden

Konfluente Zustände

- Zustände mit mehreren eingehenden Übergängen

Klonen eines konfluenten Zustands

- geklonter Zustand hat den **entsprechend eingehenden** Übergang und **alle ausgehenden** Übergänge des ursprünglichen Zustands

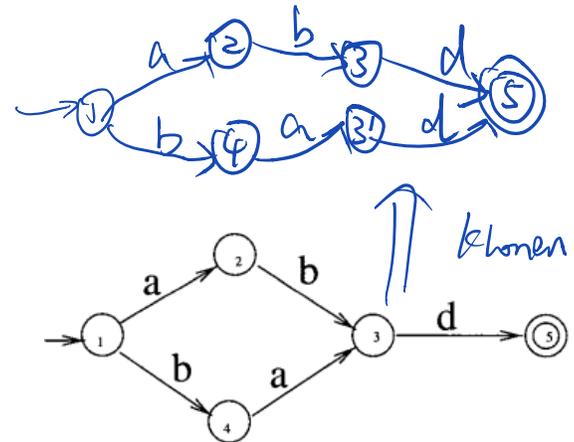


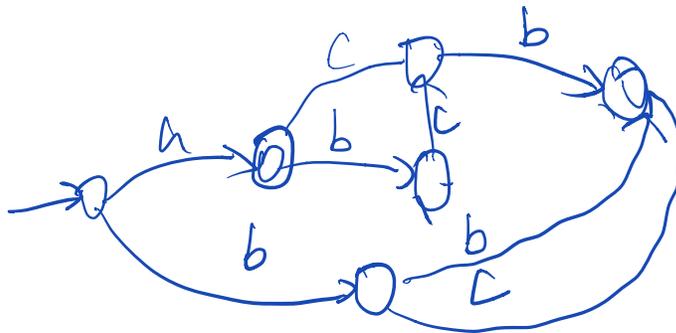
Figure: Quelle Daciuk Paper

Schritte des Einfügens

Einfügen

- Ermittle Split State
- Wenn auf dem Pfad **konfluente** Zustände erscheinen:
 - Klone ab dem ersten konfluenten Zustand **alle** Zustände bis einschließlich **Split State**
- Hänge Suffix an Split State an
- Minimiere

abc
bc
+a
+acb
+bb
+abcd



→ abc

Löschen

- Wenn es auf dem Pfad des zu löschenden Wortes **konfluente** Zustände gibt:
 - Klone ab diesem Zustand **alle** Zustände bis einschließlich **Finalzustand** des Wortes
- Wenn der Finalzustand des zu löschenden Wortes ausgehende Wörter hat:
 - Mache den Finalzustand nicht-final
- Sonst:
 - Lösche die Zustände in Rückwärtsrichtung, bis man auf einem Zustand tritt, der
 - **final** ist und / oder
 - mehr als **einen** ausgehenden Übergang hat
- Minimiere vom erreichten Zustand aus

Nächste Woche wird die Python-Implementierung des Daciuk-Algorithmus besprochen.