

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 dünn 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.

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**

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**

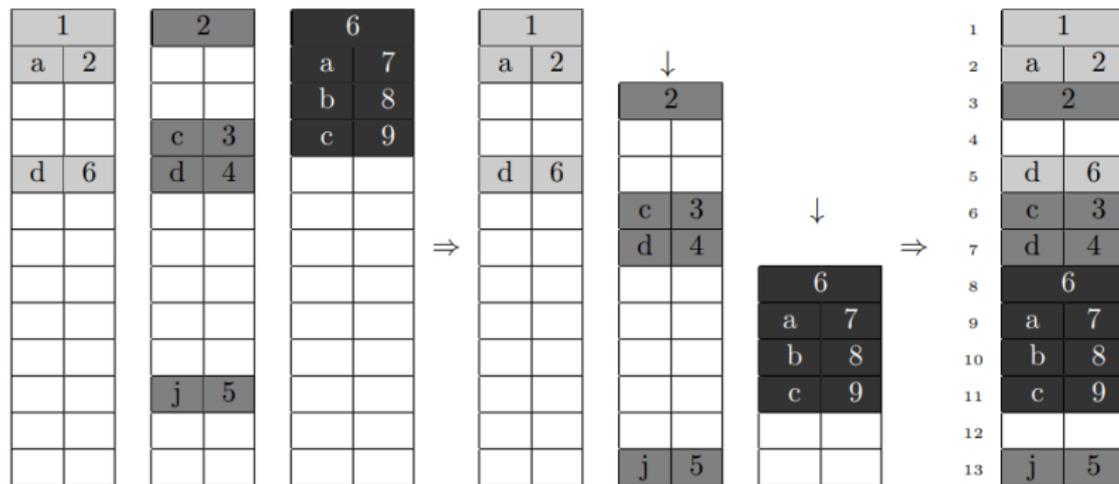


Figure: Quelle Perezautomatenscript

Tarjan-Tabelle - Beispiel

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**

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.

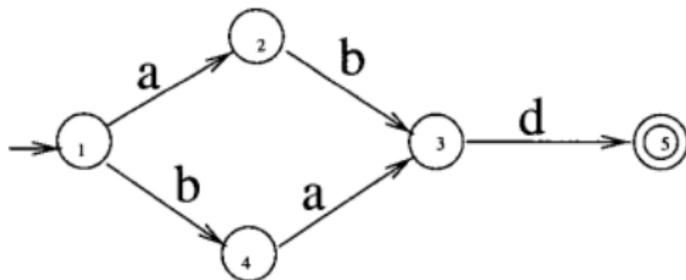


Figure: Quelle Daciuk Paper

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

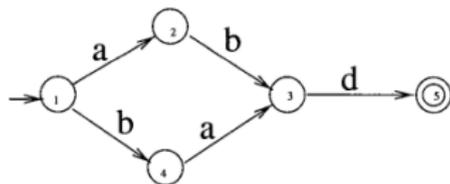


Figure: Quelle Daciuk Paper

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

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.