

Dieses Poster bietet eine Übersicht über die Ergebnisse und Vorgehensweisen der am Lehrstuhl für Angewandte Computerlinguistik (ACoLi) der Goethe Universität Frankfurt eingereichten Masterthesis "Verfahren zur Wortsegmentierung nicht-alphabetischer Schriften" von Timo Homburg.

Problemstellung: Wortsegmentierung bzw. Worterkennung ist in der Computerlinguistik eine elementare Aufgabe für die maschinelle Analyse eines Textes. Ausgehend von einer Zerlegung eines Textes in bekannte Einheiten, sogenannte Tokens, ist es möglich eine umfassende linguistische Analyse anzustreben. In der Masterthesis werden Algorithmen vorgestellt, die für die Worttrennung im Chinesischen entwickelt wurden und auf der akkadischen Keilschrift ausgetestet. Auf dem Poster werden die einzelnen Typen von Segmentierungsalgorithmen sowie der Evaluationsprozess, die Visualisierung und die Ergebnisse dargestellt.

Der Lehrstuhl für Angewandte Computerlinguistik (ACoLi) an der Goethe Universität Frankfurt am Main wurde im Januar 2013 gegründet um die Aktivitäten der Digital Humanities im Bereich Natural Language Processing (NLP) zu unterstützen. Aufbauend auf vorheriger Forschung von Juniorprofessor Dr. Christian Chiarcos entwickelt ACoLi technische Infrastrukturen zur Analyse, zum Aufbau und zur automatischen Verarbeitung von linguistischen Daten und unterhält Infrastrukturen zur Speicherung, Abfrage und Visualisierung von linguistischen Experimenten. Ein Interessensgebiet der Forschung ist die Analyse von außereuropäischen und historischen Sprachen, welche mit diesem Poster aufgegriffen werden.

Segmentierung:

Satz: 这是一张有趣的海报

Worte: 这 是 一 张 有 趣 的 海 报

(Zhèshì yízhāng yǒuqù de hǎibào)
(Das ist ein interessantes Poster)

Regelbasierte Verfahren und Baseline

Baseline Verfahren:

- Verfahren mit eingeschränkter Funktionalität
- Referenzpunkt zur Evaluation fortgeschrittener Methoden

Original(O): [Image of original Chinese characters]

AVG(L = 1): [Image of segmentation result with L=1]

AVG(L = 2): [Image of segmentation result with L=2]

Beispiel Mittlere Wortlänge: Mittlere Wortlänge des Korpus als Baselinemethode der Klassifikation

Regelbasierte Verfahren:

- Händisch entwickelte Regeln werden verwendet
- Scoringssysteme auf Zeichenbasis werden verwendet
- Scoringssysteme gewichten Zeichen(folgen) auf Segmentierbarkeit

Original(O): [Image of original Chinese characters]

Prefix/Suffix: [Image of segmentation result]

Zeichenliste: [List of characters]

Frequenzen: [Frequency table]

Prefix/Suffix Segmentierung: Klassifiziere mittels Frequenzen auf Zeichenebene Beginn/Mittel/End und Einzelzeichens

Wörterbuchbasierte Verfahren

Idee:

- Leite von bekannten Worten eine Segmentierung her
- Arbeite auf Grundlage von ggf. aus Korpora generierten linguistisch angereicherten Wörterbüchern
- Linguistische Anreicherungen: Wortfrequenzen und Zeichenfrequenzen pro Position, Vorhergehende/Nachfolgende Worte, Zeichenvarianz und Wortkategorien

Beispiele:

- **MaxMatch:** Versucht von Beginn zu Ende des Textes jeweils das längste Wort abzutrennen
- **MinWCMATCH:** Minimiere die Anzahl der Worte im gegebenen Text
- **MaxProbMatch:** Maximiere die Gesamtwahrscheinlichkeit aller Worte im jeweiligen Text

Original(O): [Image of original Chinese characters]

MinWCMATCH: [Image of segmentation result]

MaxMatch: [Image of segmentation result]

Wörterbuch: [Image of dictionary]

Zeichenliste: [List of characters]

Überwachtes maschinelles Lernen

Idee:

- Annotiere eine Klasse zu jedem Zeichen des Textes
- Klassifiziere auf Grundlage von FeatureSets aus dem Chinesischen
- Verwende im Chinesischen erfolgreiche Verfahren des maschinellen Lernens

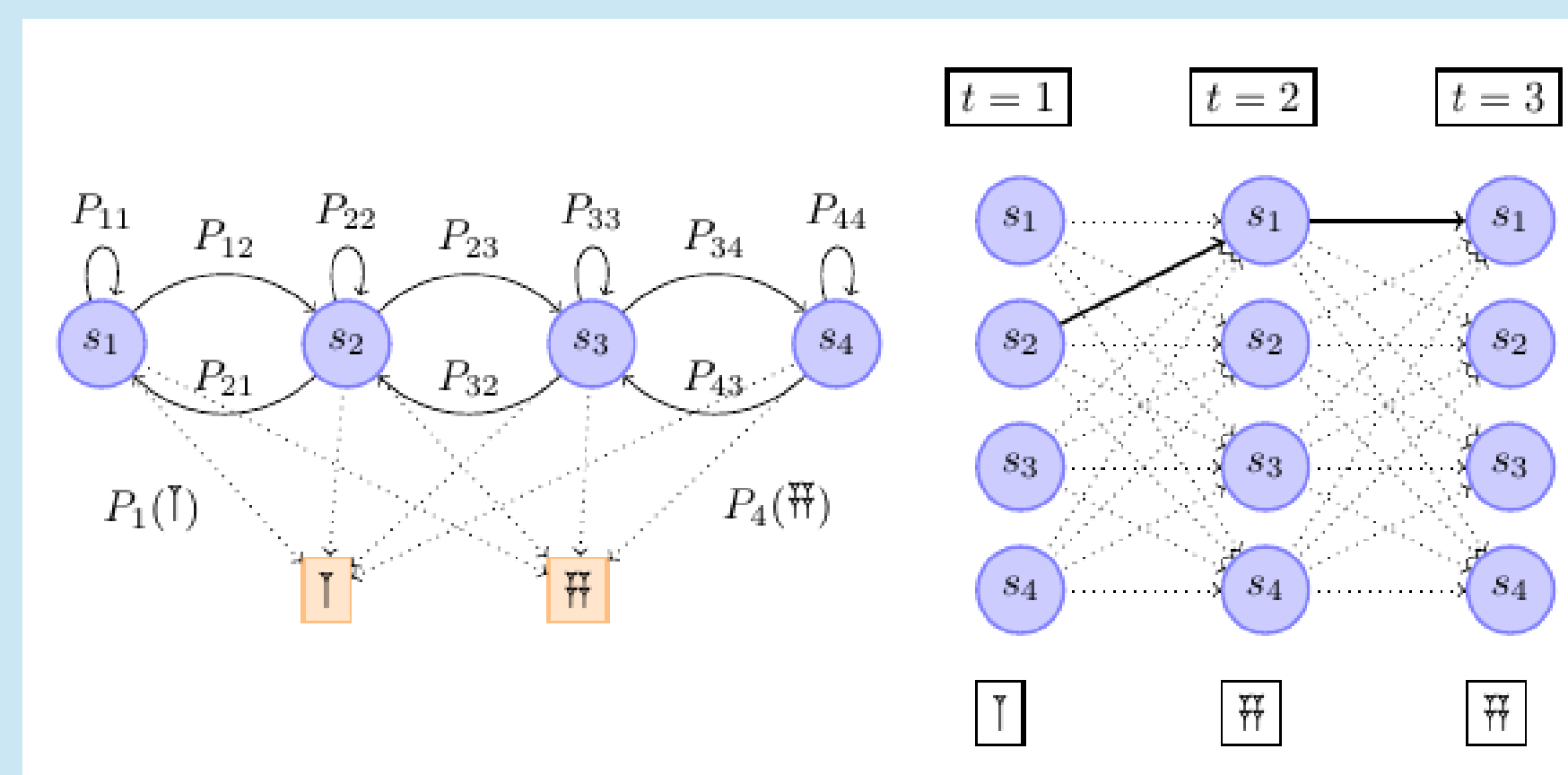
Verwendete Verfahren des maschinellen Lernens:

- Clusteringverfahren (kNN, kMeans)
- Entscheidungsbaumverfahren (C45)
- Statistische Verfahren (NaiveBayes, MaxEnt)
- Sequenzlabeling (HMM, CRF)
- Perzeptronlernen

我妈妈是中国人
S B E S B I E

Beispielannotation der Klassen:

Begin, Intermediate, End, Single Character



Beispielklassifikation mit HMM*:

Abbildung eines Hidden Markov Models für die Klassifikation mit

den Observations von zwei Keilschriftzeichen

*Grafik modifiziert von <https://gist.github.com/mblonde/472540>

Evaluation

Metriken:

- **Characterbasierte Metriken**
Analysiere die Richtigkeit der jeweiligen Entscheidung pro Zeichen
- **Wortbasierte Metriken**
Wie viele Worte wurden korrekt segmentiert?
- **Wortfenstermetriken**
Wie viele Wortfenster wurden richtig bzw. falsch segmentiert?
- **Editierdistanzmetriken**
Wie groß ist die Editierdistanz der generierten Wortgrenzen zu den richtigen Wortgrenzen?

Original: [Image of original Chinese characters]

Generiert: [Image of generated segmentation]

FalsePositive, FalseNegative, TrueNegative, TruePositive

$$Precision = \frac{TruePositive}{TruePositive + FalsePositive}$$

$$Recall = \frac{TruePositive}{TruePositive + FalseNegative}$$

Binäre Evaluation: Entscheide für jeden Character ob eine Wortgrenze danach richtig oder falsch gesetzt bzw. nicht gesetzt wurde.

Vergleichsmaße:

- **Precision:** Wie viel Prozent der vom Algorithmus gesetzten Grenzen waren korrekt?
- **Recall:** Wie viel Prozent der vom Algorithmus gesetzten Grenzen entsprechen der erwarteten Segmentierung?
- **F-Score:** Gewichtetes harmonisches Mittel aus Precision und Recall

Original(O): [Image of original Chinese characters]

Generiert(G): [Image of generated segmentation]

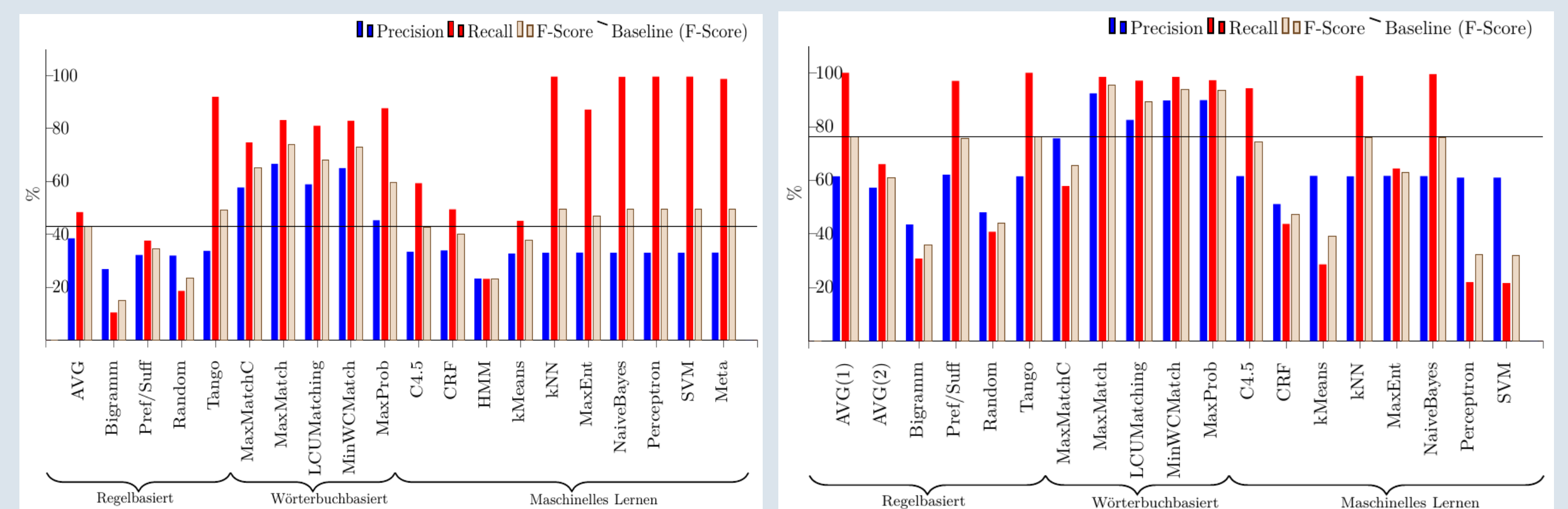
$WD(O, G) = \frac{1}{N} \sum_{i=1}^{N-k} \sum_{j=i+1}^N (|O_{ij} - G_{ij}| > 0)$

Wortfensterevaluation: Wende ein Scoringssystem für Wortfenster an um zu einer Aussage über die Güte der Segmentierung zu kommen.

Experimente und Ergebnisse

Aufbau der Experimente:

- 80% eines mittelbabylonischen Korpus als Trainingsmenge und 20% des Korpus als Testmenge
- 80% eines chinesischen Korpus als Trainingsmenge und 20% des Korpus als Testmenge



Mittelbabylonisch(Train) auf Mittelbabylonisch(Test)

Chinesisch(Train) auf Chinesisch(Test)

Ergebnisse:

- Baseline Verfahren erreichen im Chinesischen einen F-Score von >70% und im Akkadischen von ca. 45%
- Regelbasierte Verfahren können im Allgemeinen die Baseline nicht überbieten
- Wörterbuchbasierte Verfahren können in der Keilschrift mit einem F-Score von ca. 70% überzeugen
- Überwachtes Maschinelles Lernen erreicht meist mindestens die Baseline, scheitert jedoch an nichtoptimierten FeatureSets und zu geringen Trainingsmengen für die Keilschrift

Visualisierung

Original Text [Image of original Chinese text]

Generated Result [Image of generated segmentation]

Transliteration Text [Image of transliterated text]

Transliteration [Image of transliterated text]

Get Diffs

Legend:

- Correct segmentation (Green)
- Correct transliteration (Yellow)
- Wrong segmentation (Red)
- Missing segmentation (Blue)
- Wrong transliteration (Cyan)

Ziele der Visualisierung:

- Ein einfach zu bedienendes Werkzeug für Interessenten aller Fachbereiche
- Klare Strukturierung der Ergebnisse
→ Anzeige des Ergebnisses mit Signalfarben
- Rot: Falsche Segmentierung
- Grün: Richtige Segmentierung
- Blau: Fehlende Segmentierung