

Textextraktion aus nichttextuellen Datenquellen

Motivation

Im WWW sind viele nichttextuelle Datenquellen verfügbar

- PostScript (PS und PDF)
- Wordprozessor (z.B. DOC, RTF)
- Bilder (z.B. GIF, JPG)
- ...
- Potentiell hoher Informationsgehalt
- Werden nicht indiziert
- Sind nicht auffindbar

Gliederung

- Textextraktion aus PostScript-Dateien
- Textextraktion aus WWW-Images
- Bezug auf die PG

Das Problem mit PostScript

PostScript ist eine Programmiersprache.

Es gibt zwei Problemebenen:

1. Es muss nicht der gesamte Textinhalt in der ps-Datei stehen

a

```
/fib {  
  dup  
  1 eq { }  
    { dup 2 eq { pop 1 }  
      { dup 1 sub fib exch 2 sub fib add }  
    }  
  ifelse }  
ifelse  
} def  
  
(Fib[6] = ) show 6 fib ( ) cvs show
```

b

Fib[6] = 8

Das Problem mit PostScript (2)

2. Der Text ist fragmentiert und gemischt mit Formatierungen

```

Internal data in parentheses Word fragments No spaces
...
getinterval dup(Display)eq exch 0 4 getinterval(Next)eq or)(pop false)
...
(/usr/users/eksl/oates/papers/96/mlc/final/paper.dvi)
...
(abstract)98 973 y fr (Pinding)e(structure)i(in)em)o(ultiple)h(streams)
f(of)f(data)98 1018 y(is)30 b(an)e(imp)q(ortan)o(t)h(problem.)64
b(Consider)30 b(the)98 1064 y(streams)19 b(of)e(data)h(\01fo)o(wing)
g(from)f(a)g(rob)q(ot's)h(sen-)98 1110 y(sors.)f(the)g(monitors)h(in)
g(an)f(in)o(tensiv)o(e)i(care)g(unit.)98 1155 y(or)f(p)q(ario)q(dic)
j(measuremen)o(ts)f(of)e(v)n(arious)h(indica-)98 1201 y(tors)k(of)
g(the)g(health)h(of)e (the)i(econom)o(y)m(.).41 b(There)98 1247 y(is)17
b(clearly)h(utilit)o (y)g(in)f(determining)h(ho)o(w)d(curren)o(t)98
1292 y(and)g(past)h(v)n (alues)g(in)g(those)g(streams)h(are)e(related)
98 1338 y(to)22 h(re)h(v)n(alues.)45 b(W)m(e)22 orm)o(ulate)
g(the)h (prob y(lem)17 b(of)f(\014nding)i (e)g(in)g(m)o
(ultiple)g(str 1429 y(of)f(categorical)i(data (searc)o(h)g
(o)o(v)o(er)g(space)98 1475 y(of)24 b(dep)q(en q(s.)30
  
```

Abstract
 Finding structure in multiple streams of data is an important problem. Consider the streams of data flowing from a robot's sensors, the monitors in an intensive care unit, or periodic measurements of various indicators of the health of the economy. There is clearly utility in determining how current and past values in those streams are related to future values. We formulate the problem of finding structure in multiple streams of categorical data as search over the space of dependencies, unexpectedly frequent or

Abstract
 Finding structure in multiple streams of data is an important problem. Consider the streams of data flowing from a robot's sensors, the monitors in an intensive care unit, or periodic measurements of various indicators of the health of the economy. There is clearly utility in determining how current and past values in those streams are related to future values. We formulate the problem of finding structure in multiple streams of categorical data as search over the space of dependencies, unexpectedly frequent or

Einfacher Textextraktor

Der Lösungsansatz aus [1] redefiniert den PostScript **show**-Operator:
Die extrahierten ASCII-Zeichen werden in eine Datei umgeleitet

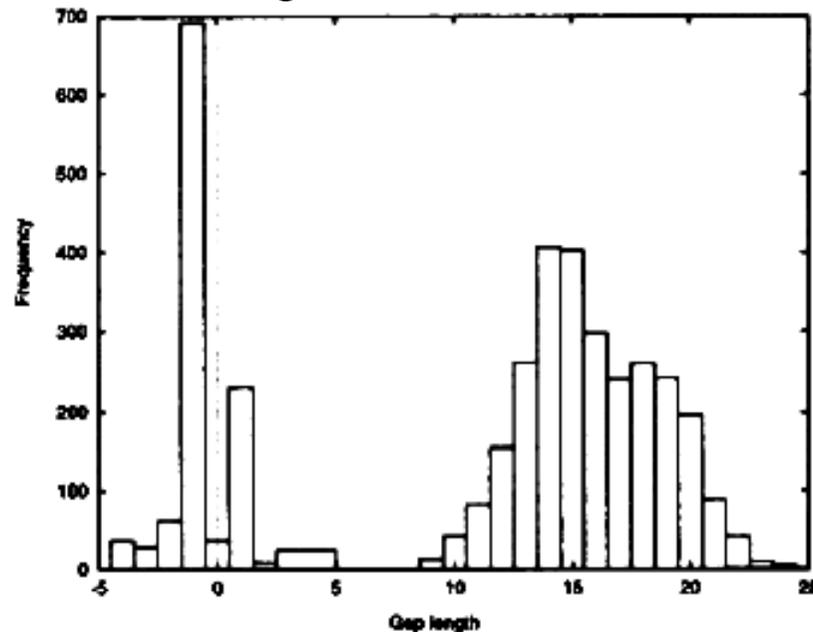
a	<code>/show (print) def</code>	Finding structure in multiple streams of data is an important problem. Consider the streams of data flowing from a robot's sensors, the monitors in an intensive care unit, or periodic measurements of various indicators of the health of the economy. There is clearly utility in determining how current and past values in those streams are related to future values
---	----------------------------------	--

Durch Verbesserung werden Wortfragmente getrennt ausgegeben

b	<code>/show (print () print) def</code>	Finding structure in multiple streams of data is an important problem. Consider the streams of data flowing from a robot's sensors, the monitors in an intensive care unit, or periodic measurements of various indicators of the health of the economy. There is clearly utility in determining how current and past values in those streams are related to future values
---	--	--

Einfacher Textextraktor (2)

Die Stellen zwischen Wortfragmenten werden heuristisch ermittelt.



```

c /X 0 def
/show {
  currentpoint pop
  X sub 5 gt ( ( ) print ) if
  dup print
  systemdict /show get exec
  currentpoint pop /X exch def
} def

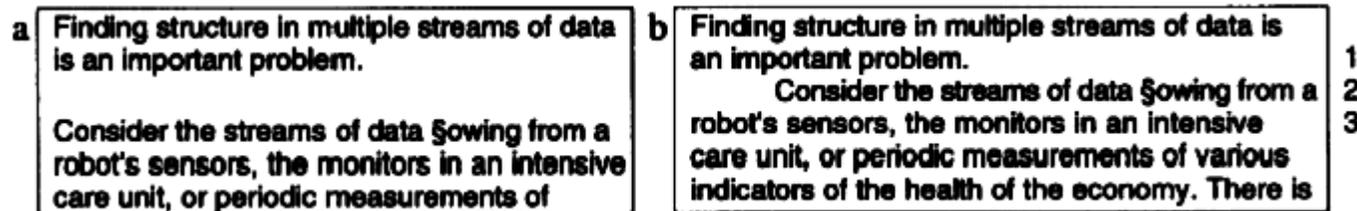
```

Finding structure in multiple streams of data is an important problem. Consider the streams of data flowing from a robot's sensors, the monitors in an intensive care unit, or periodic measurements of various indicators of the health of the economy. There is clearly utility in determining how current and past values in those streams are related to future values.

Verbesserter Textextraktor

Die Erweiterung (in Python) basiert auf dem ersten Ansatz und behandelt folgende Fälle:

- Abstände bei großen und kleinen Schriften
- Paragraph vs. Zeilenumbruch



- Nicht-ASCII Zeichen
- Worttrennung am Zeilenende
- Seitenumkehrung

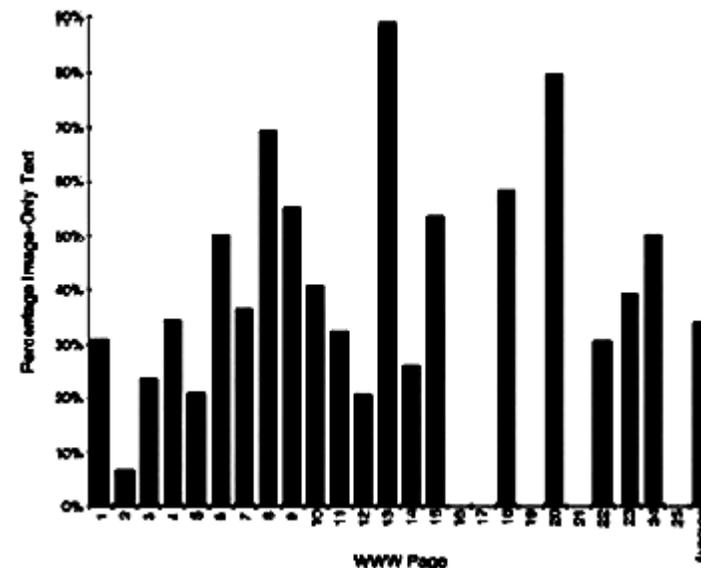
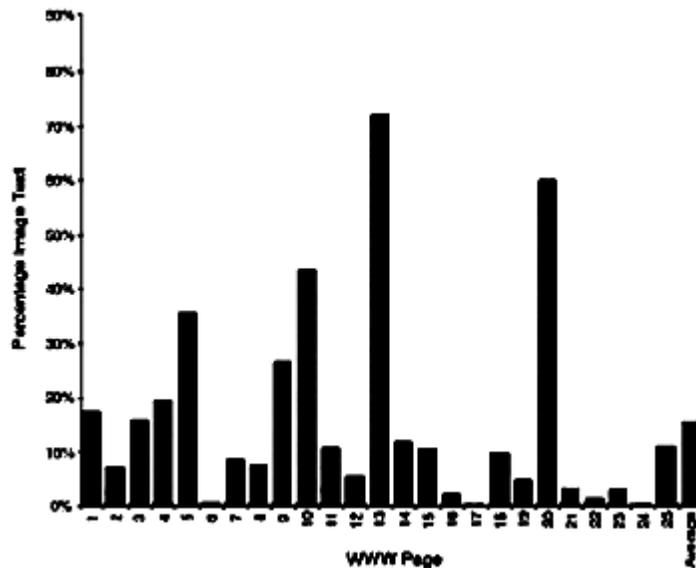
Andere Ansätze

- **ps2ascii.pl** Perl-Script, extrahiert geklammerten Text
- **ps2txt** C-Programm, extrahiert geklammerten Text, spezieller Code für durch dvips generierte PostScript-Dateien
- **ps2a.sh** komplexes PostScript-Programm, optimiert für Dateien, generierten aus T_EX
- **pstotext (DEC)** PostScript- und C-Programm von DEC, teuer, entspr. komplex und sehr langsam, aber konvertiert exzellent
- **ps2ascii** aus Ghostscript, nicht besonders zuverlässig
- **ps2html** Variante des ps2ascii, entwickelt an Johns Hopkins University für spezielle Dateien (PS aus QuarkXPress)
- **pstotext** funktioniert mit GhostScript, konvertiert auch PDF, aber mit weniger zuverlässigem Ergebnis
- **pdf2text**
pdf2html Bestandteil des Pakets xpdf

Textextraktion aus WWW-Images

Ergebnis einer kleinen Stichprobe (25 Websites) [2]:

- Im Durchschnitt 15 % der ganzen Textinformation einer Website steht in Images.
- Im Durchschnitt 34 % davon steht bei vielen Websites nur in Images, und nirgendwo sonst auf der Website.



Textextraktion aus WWW-Images (2)

GIF's und JPEG's sind im Internet am meisten verbreitet.

- GIF (Graphics Interchange Format): 256 Farben, verlustfrei

- JPEG: RGB-Farben, verlustbehaftet, besser für Fotos



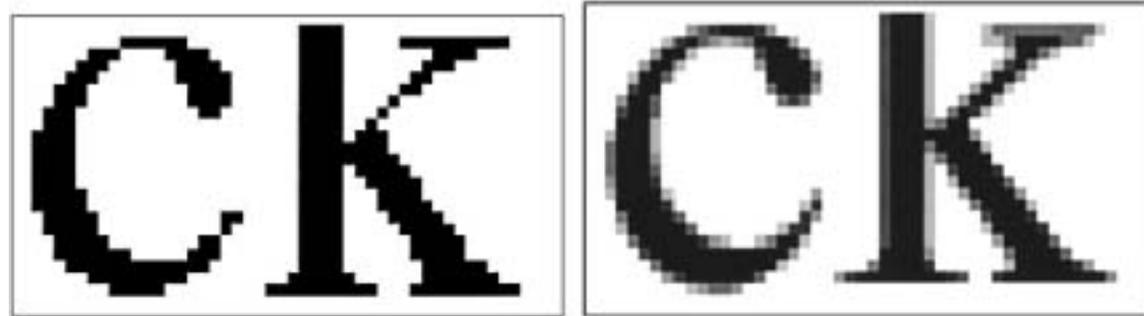
Abb.: JPEG's Kompressionsartefakte

Die Ausnutzung der besonderen Eigenschaften eines Graphikformats kann die Ergebnisse der Texterkennung verbessern.

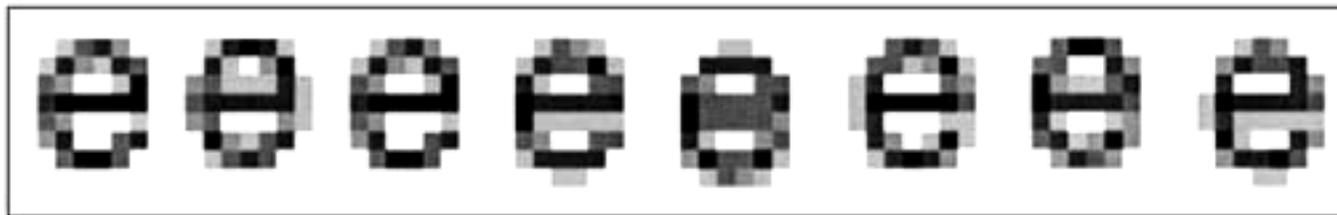
Allgemeine Probleme

- Kleine Auflösung (üblich 72 dpi)

- Anti-Aliasing



- Räumliche Mustereffekte (spatial sampling effects)



Allgemeine Probleme (2)

- Images mit schwer erkennbaren text



- Dynamische Images (GIF89a Standard)
- Images mit „Multizeichen“



Verfahren

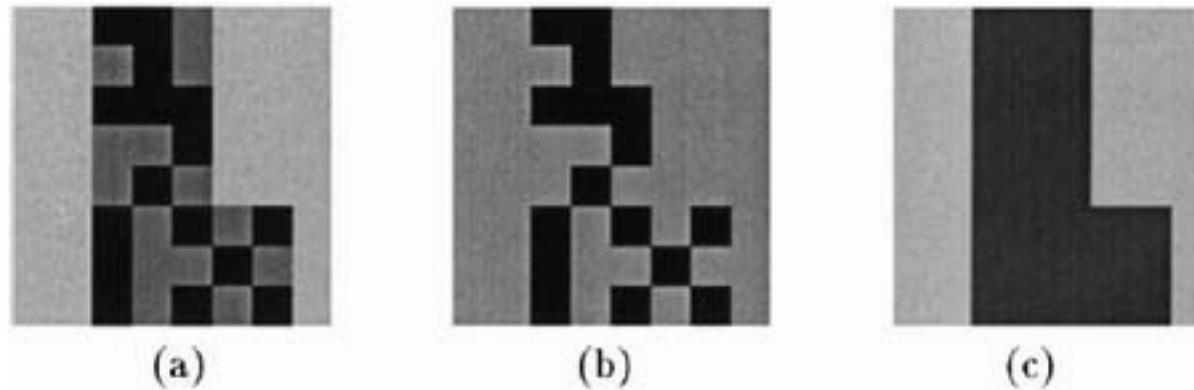
Aktuelle OCR-Technologie erkennt Text nur auf einem einheitlichen Hintergrund [3].

Alternative Verfahren benutzen bis zu drei Schritte:

- Farbclustering (color clustering)
- Zeichenentdeckung (character detection)
- Layoutanalyse (layout analysis)

Beispiel des Farbclusterings

- a) Buchstabe „L“
- b) Clusterung, basiert auf RGB-Distanz
- c) Clustering, basiert auf Kombination von RGB- /räumliche Distanz



Bezug auf die Projektgruppe

PostScript (PS und PDF):

- Public Domain Programme verfügbar
- Zur Verbesserung der Treffer der von uns eingesetzten Suchmaschine
- Umwandlung in HTML-Format (mit Qualitätsverlusten)

WWW-Images:

- Keine Public Domain Programme verfügbar
- Ggf. Ressourcen-intensiv

Literatur

- [1] G. G. Nevill-Manning, T. Reed, I.H. Witten. Extracting Text from PostScript. In *Software-Practice and Experince*, vol. 28(5), 481-491, 1998.
- [2] D. Lopresti, J. Zhou. Locating and Recognizing Text in WWW Images. In *Information Retrieval 2*, 177-206, 2000.
- [3] V. Wu, R. Manmatha, E. M. Riseman. In *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 21, No. 11, 1999
- [4] D. Byers. Full-text Indexing of Non-textual Resorces. In *Computer Networks and ISDN Systems 30* (1998), 141-148