

Methoden



der automatischen Spracherkennung

WIE SAGE ICH ES MEINEM COMPUTER?

Foto: Aussenkyfer

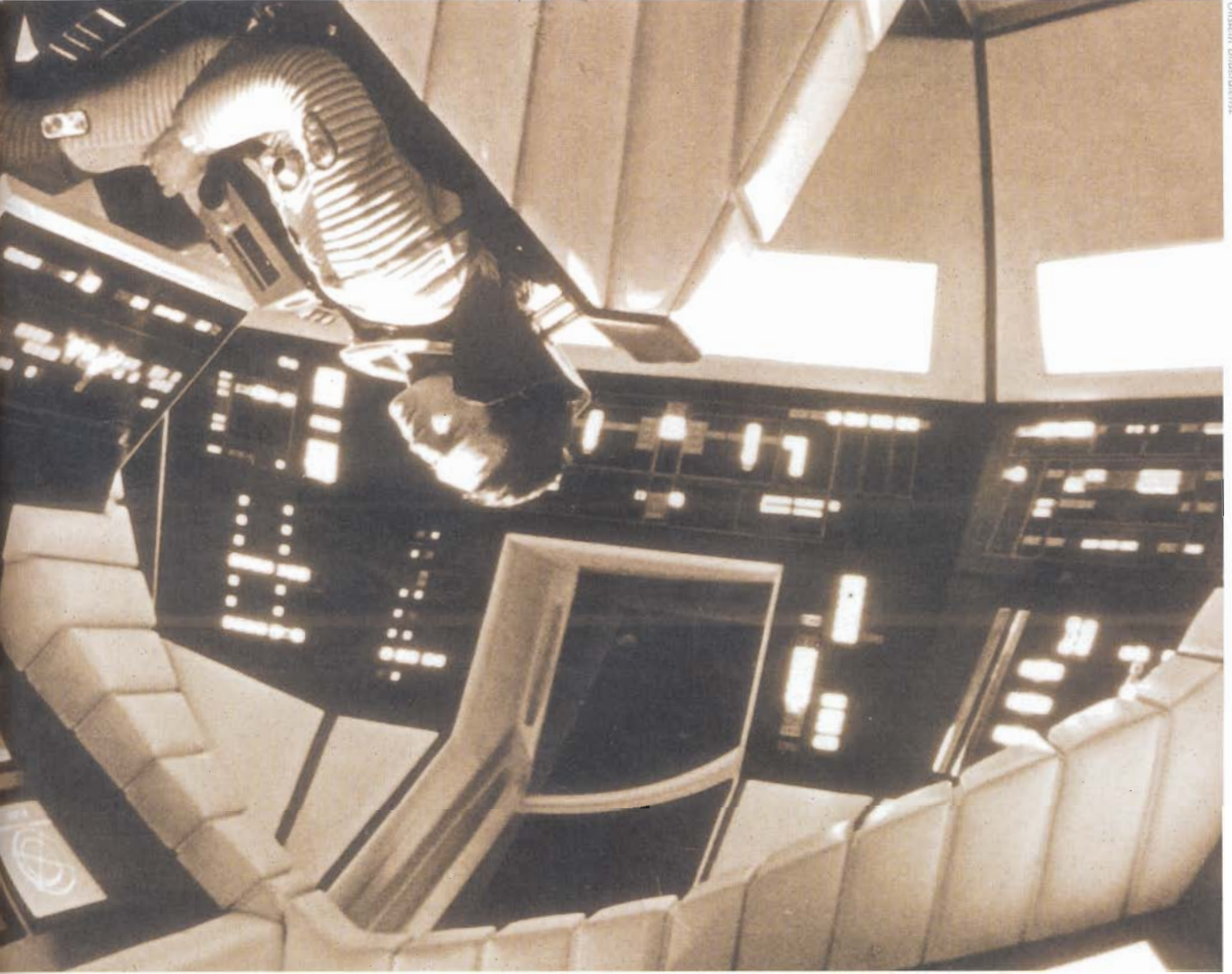


In „2001 Odyssee im Weltall“, einem der besten Science-Fiction Filme aller Zeiten, spielt ein Computer eine der Hauptrollen. Der Rechner „HAL“ kann sich mit dem Raumfahrer unterhalten, Schach spielen, von seinen Lippen ablesen und reagiert sogar emotional. Die Crew spricht mit „HAL“ wie mit einem Kollegen. Wenn aber der mörderische Computer abgeschaltet wird, geht sein Sprachvermögen langsam verloren, bis er in einer dramatischen Szene schließlich verstummt. Einen Computer wie „HAL“ zu bauen, einen Computer, der Sprache versteht, ist ein alter Traum der Wissenschaftler. Dabei geht es nicht darum, einen Computer zu entwickeln, der alle Nuancen der Sprache verstehen und deuten kann – es geht um etwas viel Einfacheres: Um die Diktiermaschine der Zukunft, einen Computer, der unsere Briefe aufnehmen und transkribieren kann. Systeme für die automatische Spracherkennung bringen

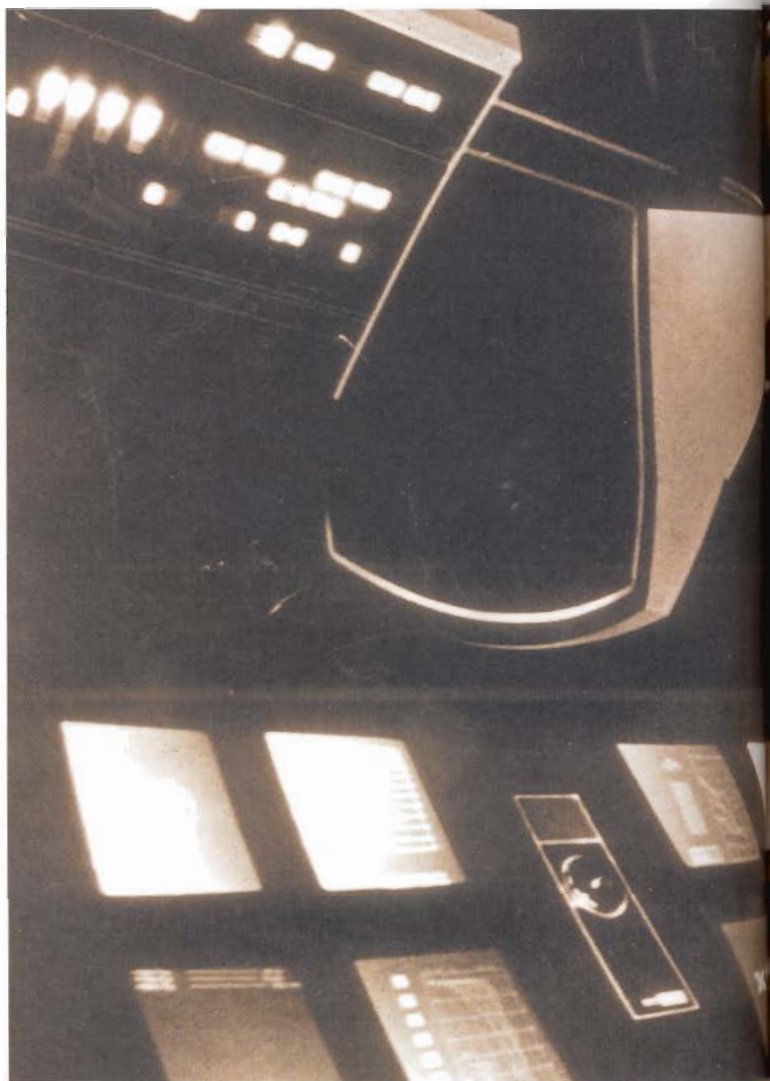
den Computer näher zu den Menschen, weil der Benutzer nicht mehr über eine Tastatur mit dem Computer kommuniziert, sondern in seinem eigenen Code, das heißt mit Worten. Spracherkennung sind aber noch weit davon entfernt, wie „HAL“ zu arbeiten. Fast perfekte automatische Spracherkennung bleibt noch ein offenes Problem und ist ein sehr aktives Forschungsfeld.

Dabei unterscheidet man zwei Arten der Spracherkennung: Sprecherabhängig und Sprecherunabhängig. Im ersten Fall wird das System auf die Stimme und Diktion einer Person trainiert, und andere Personen können das System erst nach einer Anpassung verwenden. Das Ziel der heutigen Forschung ist, Systeme zu bauen, die nicht maßgeschneidert für eine Person sind, das heißt vom Sprecher unabhängig. Dabei soll das System auch mit fließend gesprochenen Sätzen umgehen können. Der Benutzer soll nicht jedes Wort mit ... einem ... künstlichen ... Pause ... voneinander ... trennen, sondern soll natürlich reden können.

Der Rechner „Hal“



Ulrichen Bildarchiv



In diesem Beitrag möchte ich erläutern, warum automatische Spracherkennung immer noch ein schwieriges Problem ist und welche Methoden in der Informatik verwendet werden, um Computern das Verstehen der Sprache beizubringen. Automatische Spracherkennung ist ein wahrlich interdisziplinäres Feld, wo Methoden aus der Linguistik, der Signalverarbeitung und der Mustererkennung zum Einsatz kommen.

Bei der automatischen Spracherkennung geht es darum, ein akustisches Signal in den entsprechenden Text zu transformieren. Der erste Schritt besteht also darin, die Sprache mit einem **Signal** Mikrophon aufzunehmen. Wir wissen aus physiologischen Untersuchungen, dass Menschen akustische Signale in ihre elementaren Bestandteile zerlegen, das heißt die Frequenzmischung wird im Ohr und Gehirn analysiert. Eine ähnliche Art der Informationsverarbeitung kann man mit einem Spektrum durchführen: Dabei wird Sprache sichtbar

Das anschauliche und fundierte Fachbuch

Waldemar von Suchdoleiz (Hrsg.)
**Sprachentwicklungs-
 störung und Gehirn**
 Neurobiologische Grundlagen von Sprache
 und Sprachentwicklungsstörungen
 2001, 176 Seiten, kart.
 DM 57,90
 ISBN 3-17-016761-8

Kohlhammer
 www.kohlhammer-katalog.de

Der Herausgeber, Prof. W. v. Suchdoleiz, ist Leiter einer Forschungsabteilung und Spezialprechstunde für Entwicklungsstörungen an der Ludwig-Maximilians-Universität München, die weiteren Autoren sind ausgewiesene Fachexperten der Medizin und Psychologie.

Im Mittelpunkt des Buches steht die Frage, wie Störungen beim Spracherwerb zu erklären sind und welche Besonderheiten von Struktur und Funktion des Gehirns umschreibenden Sprachentwicklungsstörungen zugrunde liegen. Ausgehend von Kommunikationsstrategien im Tierreich und neuen Erkenntnissen der Psycholinguistik werden in einzelnen Kapiteln neurobiologische Aspekte von Sprache und Sprachentwicklungsstörungen anschaulich dargestellt.

Das Buch stellt die Frage, wie Störungen beim Spracherwerb zu erklären sind und welche Besonderheiten von Struktur und Funktion des Gehirns umschreibenden Sprachentwicklungsstörungen zugrunde liegen. Ausgehend von Kommunikationsstrategien im Tierreich und neuen Erkenntnissen der Psycholinguistik werden in einzelnen Kapiteln neurobiologische Aspekte von Sprache und Sprachentwicklungsstörungen anschaulich dargestellt.

FEINMECHANISCHE UND OPTISCHE SYSTEMTECHNIK

Optische Bänke
Präzisionsversteller
Positioniersysteme

OWIS GmbH
 Im Gaisgraben 7, D-79219 Staufen i. Br.
 Tel. ++49 (0) 7633/95 04-0, Fax ++49 (0) 7633/95 04-44
 http://www.owis-staufen.de, e-mail: info@owis-staufen.de

Zurück ins Leben

Wir helfen Ihnen aus Ängsten, Depressionen, Abhängigkeiten (inkl. Entgiftung) und Erschöpfungszuständen wie z.B. „Burn-Out“. Und führen Sie mit **individuellen**, sofortigen und kurzzeitigen Therapien zurück zu Ihren Stärken. Beihilfe und alle privaten Krankenkassen.

Aufnahme jederzeit in unseren Akutkliniken!

Klinik Berlin/Brandenburg ☎ 03 36 79/64 - 1 00
 Klinik Schwarzald ☎ 0 78 33 /79 20
 Klinik Weserbergland ☎ 0 57 54 /8 70



Die Oberbergkliniken
 PSYCHOTHERAPIE · PSYCHIATRIE · PSYCHOSOMATIK
 Dr. E. Gottschaldt, Dr. L. Schlüter-Dupont

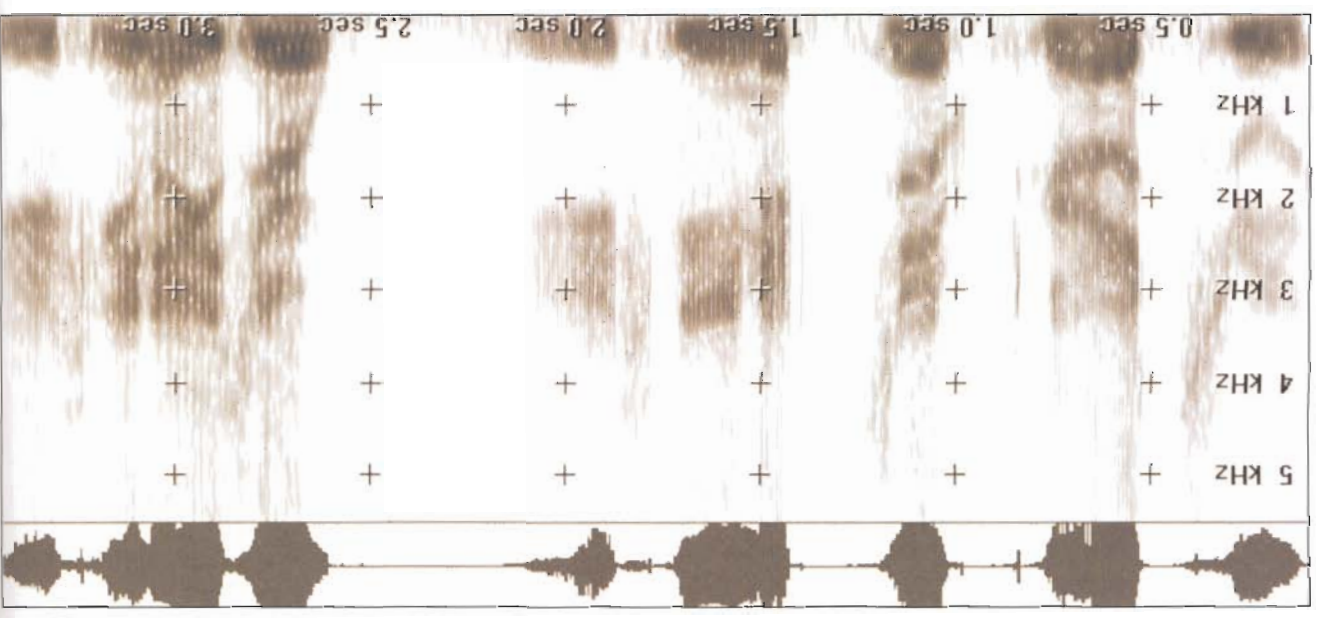
Info-Telefon
 0180/5257405
 0,24 DM/Min.
 www.oberbergkliniken.de



bar für das Auge gemacht (womit wir mehr über den Prozess der Spracherzeugung erfahren können) und ähnliche Daten werden dann dem Computer übergeben. Ein Spektrogramm ist eine Zerlegung des akustischen Signals in seine Frequenzkomponenten. Abbildung 1 zeigt ein Beispiel: In die waagerechte Richtung verläuft die Zeit. In die senkrechte Richtung wird mit Grauwerten die Intensität von jeder vorhandenen Frequenz angezeigt. Die gemusterten Frequenzen gehen von 0 bis 5 KHz. Dort, wo eine Frequenz besonders stark hervortritt, wird geschwärzt. Ist eine Frequenz nicht vorhanden, bleibt das entsprechende Feld hell. Ganz oben im Bild sieht man die Mikrophondaten: Schwingungen entsprechen den in drei Sekunden ausgesprochenen Worten. Darunter sieht man das Spektrogramm in der Zeit. Die senkrechten Streifen stammen aus dem

Musterungsprozess (die Frequenzen werden alle 50 Mikroskunden ermittelt). Man sieht sehr deutlich, dass sich gewisse waagerechte Muster bilden, dunkle Streifen, die sich nach oben (in die höheren Frequenzen) oder nach unten bewegen. Diese dunklen Streifen werden „Formanten“ genannt und ihre Position und Anzahl ist bereits ein wichtiges Indiz für die Art des Lautes, der ausgesprochen wurde. Gebübe Linguisten können sogar manchmal aus dem Spektrogramm auf die gesprochenen Worte schließen. Insbesondere Vokale zeigen eine deutliche Bänderstruktur in dem Spektrogramm. Abbildung 2 zeigt ein Beispiel des Spektrogramms, wenn „e-a-e-a“ ausgesprochen wird. Man sieht die Übergänge von einem Vokal in den anderen, und man sieht deutlich, dass die Formanten sich auf unterschiedlichen Höhen befinden.

Abbildung 1: Spektrogramm aus den Mikrofondaten



Das Phänomen des Übergangs der Formanten eines Lautes in den nächsten macht automatische Spracherkennung schwierig, da es nicht immer klar ist, wo ein Laut aufhört und wo der nächste beginnt. Andere Lauten, insbesondere solche wie „sch“ und „f“, erzeugen Spektrogramme, die eher wie Rauschen aussehen. Bei Konsonanten ist die Formantenstruktur nicht besonders ausgeprägt.

Bei der automatischen Spracherkennung sind viele unterschiedliche Ansätze getestet und implementiert worden. Eine populäre Methode besteht darin, die akustischen Signale für bestimmte Worte zu speichern und sie als Schablone zu benutzen. Dabei werden das Spektrogramm oder Varianten davon verwendet (wie zum Beispiel das cepstrum, mel cepstrum, usw.). In der Zeitachse werden zum Beispiel alle 100 ms zehn der Frequenzbänder gemustert und die Intensitätswerte gespeichert. So verwandelt sich ein Wort, das in einer halben Sekunde ausgesprochen wurde, in einen Datensatz mit 500 Einträgen. Wird jetzt ein neues Wort ausgesprochen, kann das Spektrogramm wieder berechnet werden und die Ähnlichkeit beider Worte kann getestet werden. Wir speichern also Schablonen für die Worte im Wörterbuch („ja“ oder „nein“ für ein einfaches Auskunfts-system) und vergleichen neue Worte mit diesen Schablonen. Es wird das Wort ausgewählt, das den kleinsten Abstand zu dem Signal hat. Diese Methode hat jedoch den Nachteil, dass nur eine begrenzte

Erste Erkennungsmethode: Schablonen

Prof. Dr. Raúl Rojas

Raúl Rojas, geboren in Mexiko Stadt, studierte Mathematik und Physik an der

Nationalen Technischen Universität Mexikos. Abschluss des Mathematikstudiums mit dem Master in Science.

Parallel dazu Master in Economics an der Autonomen Nationalen Universität

Mexikos. Promotion und Habilitation an der FU Berlin. Von 1994 bis 1997 C3 -

Professor für künstliche Intelligenz an der Universität Halle. Seit Ende 1997 C4 -

Professor an der Freien Universität Berlin. Rojas ist Herausgeber

mehrerer Bücher, u.a. „Neural Networks“, Springer-Verlag, 1996; „The First Computers“, MIT-Press, Cambridge 2000;

„Encyclopedia of Computers and Computer History“, Fitzroy-Deoborn, NY, 2001.



Foto: Aufschlager

Kontakt:

Fachbereich Mathematik und Informatik

Institut für Informatik

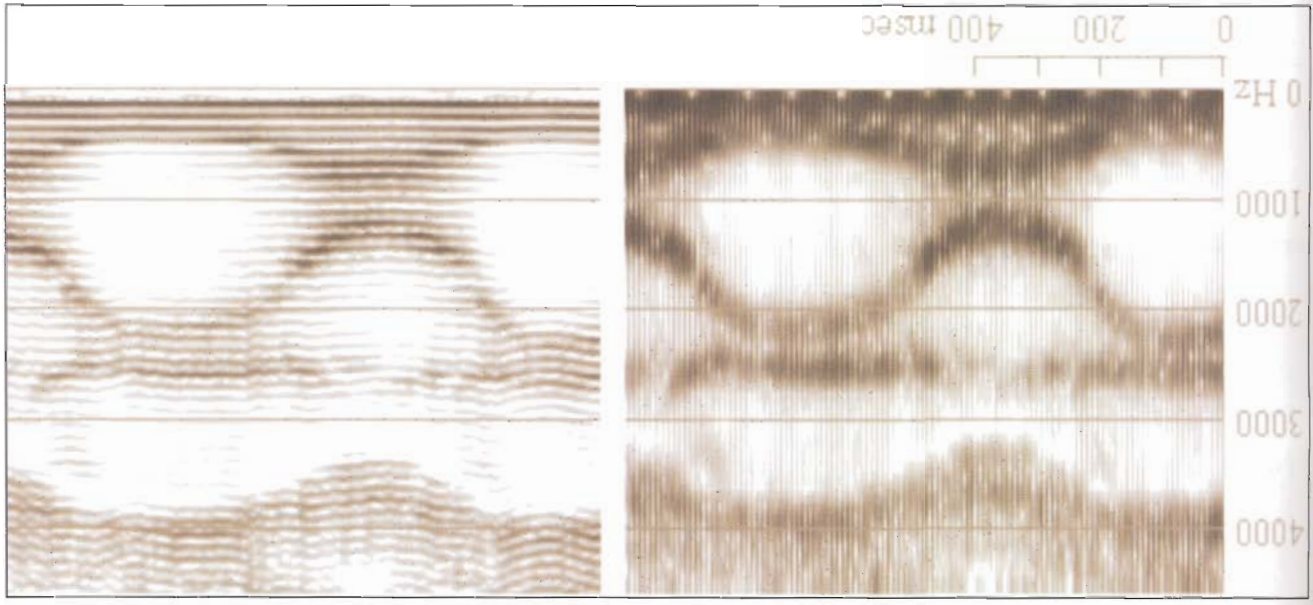
Takustr. 9

14195 Berlin

Tel. 030 / 83 87 51 30

E-Mail: rojas@inf.fu-berlin.de

Abbildung 2: Formanten der Vokale: e - a - e - a



Studentische Darlehenskasse e.V.

Studentenhaus 1. Stock
 Hardenbergstraße 35, 10623 Berlin
 Tel.: 319 001-0
 E-Mail: mail@dakaberlin.de
 Internet: www.dakaberlin.de

Öffnungszeiten:
 Mo/Di/Da/Fr Mi
 10 - 12 Uhr 14 - 16 Uhr
 10 - 12 Uhr

In den Semesterferien:
 DVD
 10 - 12 Uhr

Wir finanzieren Studentendarlehen
 Wir finanzieren für Examenkandidaten und Doktoranden

Wer bekommt ein Darlehen?
 Studierende und Doktoranden der FU, TU, HU, FH, TFH. Die Rückzahlung beginnt im siebenten Monat nach Auszahlung der letzten Darlehensrate. Die Höhe der Rate ist abhängig von der Darlehenssumme.

Was hoch ist das Darlehen?
 Der Höchstbetrag für das Studienabschlussdarlehen liegt bei DM 15.600. Es wird in Monatsraten bis zu max. DM 1.300 ausbezahlt. Die Höhe des tatsächlichen Darlehens richtet sich nach den jeweiligen Studien- und Lebenshaltungskosten.

Wie läuft die Rückzahlung?
 Die Rückzahlung beginnt im siebenten Monat nach Auszahlung der letzten Darlehensrate. Die Höhe der Rate ist abhängig von der Darlehenssumme.

Was wir sonst noch brauchen...
 1. Zwei selbstschuldnerische Bürgschaften.
 2. Zwei Gutachten, die bestätigen, dass ein Studienabschluss in 12 Monaten realistisch ist.
 3. Ein Passfoto neueren Datums.
 4. Ein ausgefülltes Antragsformular.

Das Darlehen ist...

- ...elternumabhängig - nur Euer Einkommen zählt - auch für Langzeitstudenten
- Wir sind...
- ...studentisch
- ...flexibel
- ...schnell
- ...gemeinnützig
- ...für Euch da

-alle Mitarbeiter sind Studenten
 -lassen mit uns reden
 -zwei Bewilligungsstellen monatlich
 -seit 1950 als eingetragener Verein
 -Tel.: 319 001-0

offenunabhängige Studienabschlussdarlehen
 DAKA Finanzspritze
 studentisch
 elternumabhängig
 flexibel
 www.dakaberlin.de
 Tel 319 001-0

diese Daten die korrekte Transkription kennt, ist es relativ einfach, die Korrespondenz zwischen Text und segmentierten Daten herzustellen. Dies kann von einem mentierten Daten herzustellen. Wegen der hohen Anzahl der Parameter des Klassifikators und der großen Trainingsmenge nimmt der Lernvorgang viel Zeit in Anspruch. Durch schnelle parallele Hardware können aber die Lernzeiten auf ein vernünftiges Maß reduziert werden.

Im normalen Betrieb kann das Klassifizierungsnetz verwendet werden, um die Reihenfolge der Wahrscheinlichkeit der Phoneme zu ermitteln. Das heißt, bei einem Mikrophonsignal kann der Klassifikator für jeden Zeitpunkt bestimmen, wie groß die Wahrscheinlichkeit ist, dass beim Zeitpunkt t usw. ausgesprochen wurden. Dasselbe gilt für den Zeitpunkt $2, 3, \dots$. Der Sprachkennner muss dann die optimale Reihenfolge auswählen. Wir wissen zum Beispiel, dass ein "a" lang ausgesprochen werden kann. Das bedeutet: Mit großer Wahrscheinlichkeit kann auf ein "a" zum Zeitpunkt t auch ein "a" zum Zeitpunkt 2 folgen. Dagegen ist es schwierig, dass auf ein "t" zum Zeitpunkt 1 ein "k" zum Zeitpunkt 2 folgt. Der Sprachkennner baut also auf statistisches Wissen über die zu erhaltende Sprache auf. Linguisten liefern diese Daten und bilden Wortmodelle, womit die optimale Auswahl der Phoneme gemacht werden kann.

Abbildung 4 zeigt, wie die optimale Lautsequenz ausgewählt wird. Zu jedem Zeitpunkt gibt es eine Mischung

nach dem aktuellen Fenster. Dies ergibt insgesamt 13 mal 18, also 234 Eingabewerte für das Netz.

Das gezeigte Netz wird trainiert, und zwar werden bekannte Sprachsignale segmentiert und das Netz präsentiert. Der "Trainer" weiß in jedem Moment, welches Phonem ausgesprochen wurde und verändert die Parameter des Klassifikators (im wesentlichen die Stärke der Verbindungskanten), um das Netz zu zwingen, eine "r" bei der Ausgabeleistung - die zum ausgesprochenen Phonem gehört - und "o" bei allen anderen Leitungen zu erzeugen. Dies wird mit Hunderten oder Tausenden von Phonemen wiederholt, bis das Netz gelernt hat, diese Daten zu klassifizieren. Dafür existieren die entsprechenden Algorithmen.

Interessant ist jetzt, dass sich das Netz bei der Präsentation eines unbekanntes Signals für die wahrscheinlichste Alternative entscheidet. Kommt zum Beispiel ein "a" über das Mikrophon, wird die Ausgabeleistung für diesen Laut am höchsten steigen, und es wird auf diese Weise angezeigt, dass wahrscheinlich ein "a" vorliegt. Der Klassifikator steigt aber selten bis zu "r" bei der Ausgabe, er liefert nur Werte zwischen o und r (z. B. $o,7$). Diese Werte können als die Wahrscheinlichkeit des Auftretens des Lautes "a" interpretiert werden. Andere Laute bekommen die Restwahrscheinlichkeit zugewiesen.

Die Art von Klassifikator ist ein Beispiel eines neuronalen Netzes. Für das Training des Netzes werden Stunden sprachlicher Daten verwendet (häufig werden diese Daten in CDs gepresst und verteilt). Da man für

Bestimmung des Pfades maximaler Wahrscheinlichkeit

Spracherkennungssysteme weit davon entfernt sind, den Sinn der ausgesprochenen Sätze zu verstehen. Es wird mit „brute force“ gerechnet, ähnlich wie Computer heute Schach spielen: Statt wie der Weltmeister auf einmal ein Muster zu erkennen, werden alle Möglichkeiten durchprobiert und die beste wird ausgewählt. Das reicht schon heute, um den Schach-Weltmeister zu schlagen – es reicht aber noch nicht, um Sprache fehlerfrei zu erkennen. An der Freien Universität arbeiten wir an diesem Problem und zwar mit einem Ansatz, der vielversprechend erscheint. Wir bauen ein geschlossenes System, in dem ein Computer gleichzeitig lernt, einen Text vorzulesen und das Vorgelesene zu erkennen und in Text zu verwandeln. Dabei versucht der Computer nah an einem menschlichen Vorleser zu bleiben. Damit kann der Computer sich selbst trainieren und die aufwändige linguistische Vorbereitung der Daten entfällt. Der Computer unterhält sich mit sich selbst.

Automatische Spracherkennung hat sich von den Unterverstärken in die Softwarehäuser bewegt. Es gibt bereits verschiedene Systeme, die für Diktate oder für Auskunfts-systeme verwendet werden können. Allerdings kann ein erfahrener Tipper schneller mit der Schreibmaschine als mit diesen Diktiermaschinen umgehen. Der Grund liegt in der noch großen Fehlerrate der Systeme: Korrekturen erfordern viel Zeit und stoppen den Fluss der Gedanken. Auch minimale Änderungen in der Umgebung (ein neues Mikrophon, Musik im Hintergrund, usw.) können die besten Spracherkennungsthema angreifen. Es sollte aber klar sein, dass heutige In dieser kurzen Darstellung haben wir lediglich das Thema angesprochen.

Garry Kasparow unterliegt IBM's Parallelrechner Deep Blue.



Foto: Ulfstein

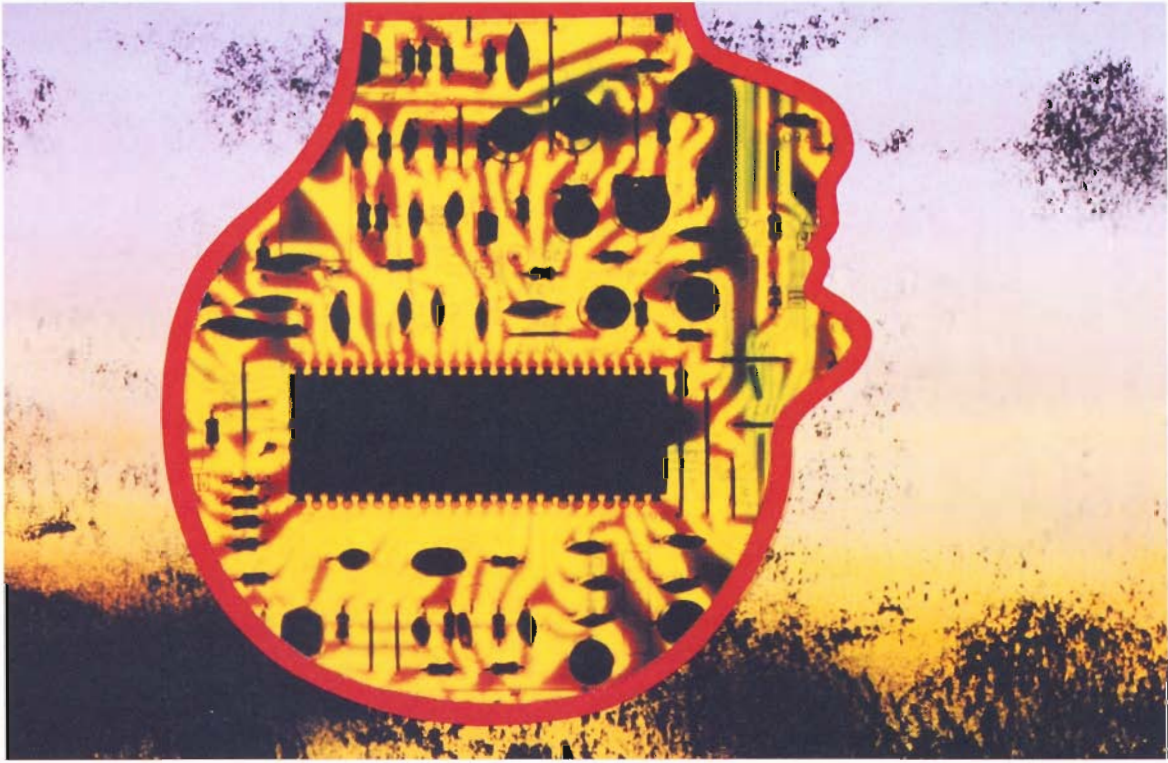


Illustration: Unicorn

systeme in die Knie zwingen. Aber auch so wird erwartet, dass bis 2005 der Markt für Spracherkennungssoftware stetig zunehmen wird, und es wird bis 2005 mit Umsätzen in Milliardenhöhe gerechnet.

Blinde und behinderte Benutzer können sicherlich bereits jetzt von Spracherkennungssoftware profitieren. Es ist zu erwarten, dass für diese Benutzer immer bessere Systeme erstellt werden.

Das größte Problem für Spracherkennungssysteme ist die soziale Akzeptanz. Wenn ich mit meinem Computer kommuniziere, will ich wirklich mit dem Computer sprechen? Ist es nicht einfacher, Befehle mit der Maus oder mit der Tastatur einzugeben? Stören sprachliche Befehle am Computer meine Mitarbeiter nicht? Wahrscheinlich ist es so, dass für Computer am Arbeitsplatz die Eingabe von Befehlen über ein

Die Zukunft der automatischen Spracherkennung

„künstliches“ Medium wie Tastatur und Maus ergonomischer ist. Anders verhält es sich mit den kleineren tragbaren Computern, den Personal Assistants. Sie sind so klein, dass keine Tastatur mehr angeschlossen werden kann. Die Eingabe über Stifte ist auch nicht besonders schnell. Da diese Geräte in der Zukunft auch als Handy

verwendet werden, bietet es sich an, sie gleich mit Spracherkennung auszustatten. Es wird deswegen erwartet, dass Spracherkennung vor allem für diese kleineren Geräte eingesetzt wird und dass PDAs die treibende Kraft für der Entwicklung neuer Software werden. Spracherkennungssoftware für diese Geräte muss jedoch um ein Vielfaches robuster sein als Software für Benutzung in einem Büro, da die Umgebung viel lauter ist.

Literaturverzeichnis

Bourlard, H., Morgan, N. (1993), Connectionist Speech Recognition, Kluwer.

Rabiner, L., Bing-Hwang, J. (1993), Fundamentals of Speech Recognition, Prentice-Hall International, London.

Rojas, R. (1996), Neural Networks, Springer-Verlag, Berlin, New York.

überschritten.

Jedoch werden wir wahrscheinlich in den nächsten 20 Jahren immer noch nicht so wie in den Science-Fiction-Filmen mit den Computern sprechen können. Auch wenn alle Sprachverarbeitungsschritte beherrscht werden und Worte mit 99,9% Genauigkeit erkannt werden sollten, wird es noch große Defizite im Sprachverständnis geben. Sprache ist sehr abhängig vom Kontext und bereits kleine Nuancen können den Sinn eines Satzes umdrehen. Ob etwas ironisch oder im Ernst gemeint ist, wird der Computer noch nicht erkennen können. Was allerdings nach diesen 20 Jahren möglich wird, ist schwer zu sagen. Man sollte sich hüten, in der Informationstechnik Prognosen zu liefern, die länger als 10 Jahre gelten. Die Grenze habe ich hier zweifach überschritten.