# Handwriting and Gesture Recognition

Interactive Systems Laboratories

KIT
Karlsruher Institut für Technologie
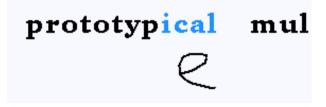
Interactive Systems Labs

# Speech, Handwriting, Text

# Gestures

## Delete Words and Characters:

prototype multimodal listening typewriter

prototype multimodal listening typewriter

prototypical mul

## Indicate Cursor Position:

prototype multimodal

prototype multimodal

## Select Characters:

prototypical mul

## Partial Word Correction:

prototypical mul

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Writing and Language

- About 6000 living languages exist
  - India:  > 1600
  - South America:  > 1000
  - Africa:  > 1000
  - Europe:  < 70

- 90% of world population speak one of the 100 widely used languages

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Writing and Language

- Only for 13% of the living languages, a written language exists


- Chinese, English, Spanish, Russian, Hindi, German cover ~50% of the world population


- Number of written languages ever used : ~660

KIT
Karlsruher Institut für Technologie
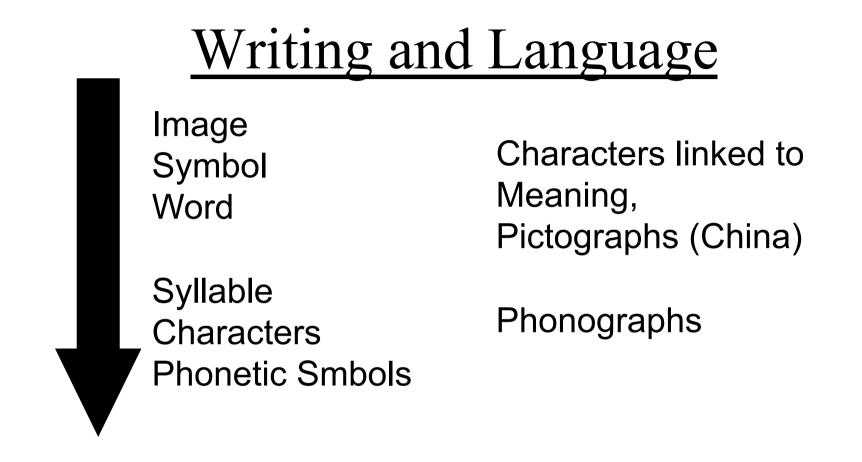
# Writing and Language

- First cave drawings:
  more than 30.000 years ago

- First writing systems:
  ~ 5000 b.C.    Sumerer in Mesopotamien
  ~ 7000 b.C.    South Eastern Europe
- Latin Alphabet:
  600 b.C.

KIT
Karlsruher Institut für Technologie

# Writing and Language

- Hundreds of written languages use the latin alphabet

- The latin alphabet is based on the sound of words (a significant departure)

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Writing and Language

Image
Symbol
Word

Characters linked to Meaning, Pictographs (China)

Syllable
Characters
Phonetic Smbols

Phonographs

Reduction of the Number of Symbols by way Phonetic Sound-based Systems (since ~ 1500 BC)

KIT
Karlsruher Institut für Technologie

# Language Input by Speech

- Fast Input of Long Texts (dictation, essays,…)

- Data Input under Devided Attention (Driving,
  Operating Machinery, ..)

- Hands, Eye Busy Situations
  (Surgery, Construction, Human Postal Sorting...)

KIT
Karlsruher Institut für Technologie

# Why Handwriting ??

Problems:

- It is Slower than Speech and Typing
- Recognition comes with Errors

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Handwriting

## Data Input without Keyboard and Mouse

- Cell Phones
- Personal Digital Assistants
- Palmtops
- WebPads
- Outdoor-Activities
- ...

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Handwriting

Input of brief messages:


- Keyboard substitute (?!)
- Mouse substitute


Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Handwriting

Communication in noisy environments:

- Factories
- Conventions
- Discos
- ....

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Handwriting

## Silent Communication:

- Meetings
- Presentations
- Military operations
- ...

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Handwriting

## Communication of Confidential Data:

- Personal Data
- Credit Card Numbers
- Codewords
- ...

Interactive Systems Labs

# Handwriting

Under Conditions, were verbal communication is not possible:

- Under water
- Handicapped people
- ...

KIT
Karlsruher Institut für Technologie

# Handwriting

Input of spatial data:

- Forms
- Mathematic Formulas
- Crossword puzzles
- ...

Interactive Systems Labs

# Handwriting

## Symbolic Data

- Graphs, Tables
- Symbolic gestures
- ...

KIT
Karlsruher Institut für Technologie

# Handwriting

Input of Biometric Data
Person Verification and Identification

- Signatures
- Writing Style

Interactive Systems Labs

# Handwriting

## Error Correction and Comments:

- Annotation and Modification of Documents
- Correcting Voice Recognition Errors
- ...

Interactive Systems Labs

# Handwriting

A Parallel, Alternate Input Modality !!

Redundance
Naturalness
Robustness
Flexibility

Interactive Systems Labs

# Handwriting Recognition

A Handwriting Recognizer transforms handwritten input in a computer readable format (e.g. ASCII)

KIT
Karlsruher Institut für Technologie

# Handwriting Recognition

# Applications

KIT
Karlsruher Institut für Technologie

# Postal Sorting

Mail:

World-wide:  ~ 1 billion mails a day

US:  ~ 43 Million mails  a day (Germany ?)

of which 30 million can be machine-processed (3 x Mount Everest)

Most of it „Junk Mail" which (luckily) has textual labels

Interactive Systems Labs

# Postal Sorting

Throughput:

Human:      3800 - 5000  Mails per hour

Machine:        30 000  Mails per hour

KIT
Karlsruher Institut für Technologie

# Handwriting Recognition

- Other Applications:
  - Processing Forms (UPS, FedEx, …)
  - PDA's, Palmtops
  - Graphic Tablets

Interactive Systems Labs

# On-Line vs. Off-line Recognition

Two types of Applications and Systems:

- Off-Line:
    - Computer input by scanning
    - Handwriting is stored as binary greyscale image
    - Writer doesn't need special hardware
      (paper and pencil are enough)
    - Data capturing can be done any time
- On-Line:
    - Need Tablet and Pen
    - Collect x,y Coordinates as a function of time
    - Use Temporal Information
    - Better Performance

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Off-line Handwriting Recognition

- Possible applications include
  - check reading
  - postal address reading
  - document analysis, ...

- Input consists of scanned <span style="color:red">bitmaps</span> without any temporal information

- Eventually location of handwriting needs to be found (document analysis)

- Stroke order doesn´t influence recognition

- But: problems through overlapping or touching characters and noisy input

KIT
Karlsruher Institut für Technologie

# Document Analysis



Figure 9. Original.

Figure 10. Result.

Figure 11. Original image with 15° skew.

Figure 12. Segmentation result.

# Character Recognition



Interactive Systems Labs

# Evaluating Handwriting Recognition Systems

- Comparing different handwriting recognizers is difficult

- Performance depends on

  - recognition task (e.g. isolated characters, words, unconstrained text)

  - writing style (e.g. printed, mixed, cursive)

  - size of dictionary

  - intended end user(s)

    - single writer (allows writer dependent system)

    - multi-writer

    - Writer-independent (requires writer independent system)

Interactive Systems Labs

# Handwriting Recognition Tasks

 boxed upper case letters

 boxed letters

 printed words

 mixed words

 cursive words

 unconstrained text

# Writing Styles

| Printed | Mixed | Cursive |
|---|---|---|
| carolyn | seller | proofs |
| cobler | hungarian | hampers |
| cluff | resignations | sember |

# LCD Graphics-Tablets



Interactive Systems Labs

# LCD Graphics-Tablets



Interactive Systems Labs

# LCD Graphics-Tablets



Interactive Systems Labs

# On-line Data Capturing

Analog Graphic tablet:
(UltraPad, ArtPad, PenPartner of Wacom

Digital-analog Graphic Tablet:
( Wacom Intuos Serie )

Tablet PC

Interactive Systems Labs

# On-line Handwriting Recognition

- Special Hardware needed (Graphic tablets)

- Interaction with Computer

- Simultaneous Writing und Capturing of Handwriting

- Handwriting stored as point sequence over time (x,y,t).

# On-line Information

1. (x,y,t)

2. Pressure

3. Tilt

4. Pen_down- and Pen_Up

5. Velocity (of 1.)

# Handwriting on a Smartboard



Interactive Systems Labs

# Handwriting - [erkennung]

## Example: Forensic,
##       Identification

KIT
Karlsruher Institut für Technologie

# Forensic

- Finding signatures in a database

- Comparison of signatures

KIT
Karlsruher Institut für Technologie

# Identification (Biometrie)

Signature Verification:

• measuring ballistic movements

• mostly based on non-visible features of the signature (e.g. pressure, acceleration, ...)

• special pens exists

Interactive Systems Labs

# Other Applications

- Mimio

- CrossPad

KIT
Karlsruher Institut für Technologie

# Processing Handwritten Material

|  | Off-line | On-line |
|---|---|---|
| Handwriting-recognition | Mail automization, Form Reading (Optical Character Recognition, OCR) | PDAs<br>Pen Computer<br>Graphic-Tablets<br>SmartBoards<br>CrossPad<br>Mimio |
| Forensic + Identification | Verification and Comparison of signatures (Document analysis) | Signature Verification (Biometrics) |

Interactive Systems Labs

# On-line & Off-line

Static Information is independent of the stroke
sequence:

E

KIT
Karlsruher Institut für Technologie

# On-line & Off-line

Dynamic Information Simplifies Segmentation:

# On-line & Off-line

Dynamic
Information
**→ Trivial →**
Static
Information

Static
Information
**→ Difficult →**
Dynamic
Information

Interactive Systems Labs

# Ergonomics

Mouse ⟷ Pen

# Ergonomics

## Ergonomic parameters:

- Wrist rotation

- Angle of Hand at wrist

- Angle between Fingers and Palm

- Distance between Fingers

- Underarm Rotation

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Ergonomics

## Normal position:

- No rotation of hand  at wrist

- No angle of hand at wrist

- No angle between fingers and palm

- no distance between fingers

- No rotation of underarm

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Ergonomics

Mouse: very big horizontal rotation at wrist

    (specifically rotation in the direction of the small finger)

# Ergonomics

Mouse:  Vertical rotation at wrist

Moving mouse forward and backward

# Ergonomic Advantages of Pens

- Less deviation from the normal position (especially for wrist rotation)

- Makes use of the fingers fine-motoric

- Complete Arm Movements

- Action at Tip of Pen (LCD Tabletts)

# On-Line Handwriting Recognition

KIT
Karlsruher Institut für Technologie

Interactive Systems Labs

# Design Parameters

- Handwriting Style

- Size of Dictionary

- Writer dependent / independent

- National Particuliarities (e.g., r's, ...)

- Left-handed / Right-Handed

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# History

- Online-Recognition started end of the 50s. (Off-line recognition already earlier)

- Handwriting recognition for mail sorting:
  mid 90s
  (blockletters already earlier)

- PDAs with Handwriting Recognition:
  Beginning of 90s (Newton of Apple)

- Palm Graffiti – Speed, Accuracy... Naturalness ?

Interactive Systems Labs

# Graffiti

# Graffiti



**Satzzeichenmodus = einmal tippen**

Interactive Systems Labs

# Graffiti

## Zeichen mit Akzent

Zur Eingabe der folgenden Buchstaben
mit Akzent:

à á â ã ä å è é ê ë ì í î ï
ò ó ô õ ö ù ú û ü ÿ y ñ,

Schreiben Sie zuerst den jeweiligen
Buchstaben und dann den Akzent wie folgt:

Die folgenden Zeichen können mit Graffit ohne
Umschalten geschrieben werden:

ç æ

## Besondere Schriftzüge

Einzelheiten siehe
PalmPilot-Gebrauchsanweisung.

| | |
|---|---|
| | ShortCuts |
| | Befehlszeichen |
| | Cursor nach links |
| | Cursor nach rechts |
| | Nächstes Feld (Adreßbearbeitungs-bildschirm) |
| | Vorheriges Feld (Adreßbearbeitungs-bildschirm) |
| | Eintrag öffnen (Adreßbearbeitungs-bildschirm) |

Interactive Systems Labs

# Graffiti

# Handwriting types

Classification of
   Handwriting Types
   According to
   Tappert (IBM):

Newark
Spaced discrete characters

Lawndale
Run-on discretely written characters

Columbus
Pure cursive script writing

Libertyville
Mixed cursive and discrete

KIT
Karlsruher Institut für Technologie

# Recognition Rates

Recognition of single symbols:

Npen++:  0-9   96.5%

            A-Z  92.7%

            a-z  91.1%

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Recognition Rates

Recognition of words:

On-line recogniton rate is higher than off-line recognition rates:

Off-line:  95%  on      500 words
On-line:  95%  on   5.000 words
          90%  on  50.000 words

# Recognition Rates

Sentence Recognition:

Robust recognition of word sequences is not yet solved
 (Npen++: 86,6% with 20.000 Words).

Interactive Systems Labs

# Training and Classification

Interactive Systems Labs

# Training



? → **Training of Recognizer** → ?

↑

Labels (Meaning)

# Training

## Manuelle Segmentation:



Interactive Systems Labs

# Training

- Strokes

- Letters

- Words

- Sentences

- Texts

KIT
Karlsruher Institut für Technologie

# Training

- Strokes

- Letters

- Words

- Sentences

- Texts

Increasing
Dictionary,
Greater
Specificity

Increasing
cost of
Labeling,
More Training
Data

KIT
Karlsruher Institut für Technologie

# Strokes Level

- Analysis-by-Synthesis

- Rule-based Recognizer

- Syntactic Recognition

- Symbolic learning

Interactive Systems Labs

# Stroke Level

Decomposition / Identification of
atomic units, building blocks

# Strokes Level

## Finding Rules

# Stroke Level

Disadvantage: complex rule bases

A oder H ?

N oder V ?

KIT
Karlsruher Institut für Technologie

# Stroke Level

- Training of rules is difficult.

- Syntactic Approaches did not prove to work well

Interactive Systems Labs

# Letters Level

Training:



Sequence of (local) Feature vectors

Class „c"

Interactive Systems Labs

# Letters Level

**Features:** Direction r,
Slope s,
Curvature k,
...

➡ (r,s,k, ...)

No methods exist to find optimal features.

Mathematical methods exist to rank features according to their relevance.

KIT
Karlsruher Institut für Technologie

# Word Level

## Word Level Training (holistic approaches):

- Recognition on the whole word
- No segmentation on letter level

Advantage: No Segmentation, Detailed Modeling

Disadvantage:
- Recognition restricted to special dictionary
- Extension of dictionary requires new training
- Less training data per class

Interactive Systems Labs

**KIT**
Karlsruher Institut für Technologie

# Training

Use of unsegmented data on word and sentence :

Idea: Use pre-trained recognizer to segment higher levels

# Training

## Manual Segmentation:

*Buchstabenebene*

## Automatic Segmentation:

*Wortebene*

# Classification

# Classification

Pattern ➡️ **Classificator (Similarity Measure )** ➡️ Class

(Similarity)

# Classification

## Class:

 Group of feature vectors, which are typical for a part of a letter
e.g. all feature vectors, which are typical for the beginning of „d"

KIT
Karlsruher Institut für Technologie

# Letter Level

Training:

Sequence of (local) feature vectors

Class „d"

Interactive Systems Labs

# Letter Models

### Letter models:
 Representation of typical sequences of feature vectors for letters

KIT
Karlsruher Institut für Technologie

# Segmentation

Interactive Systems Labs

# Segmentation

Classification of Words and Texts requires
Segmentation:

- explicit Segmentation

- implicit Segmentation

Interactive Systems Labs

# Explicit Segmentation

Segmentation performed prior to recognition:

- spatial features  (e.g. stroke distance)

- temporal features
  (e.g. temporal distance between strokes)

KIT
Karlsruher Institut für Technologie

# Explicit Segmentation

Disadvantages:

• Segmentation is difficult to find
• Segmentation errors lead to classification errors

KIT
Karlsruher Institut für Technologie

# Explicit Segmentation

# Implicit Segmentation



Interactive Systems Labs

# Context



Interactive Systems Labs

# Context

## Context influences segmentation:

B

d

H2

Interactive Systems Labs

# Recognition of Single Symbols

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Recognition of single symbols

Problem:
- Deal with Deviations
- Find best Match

Dynamic Programming:
- Elastic Matching
- Editing distance
- Levenstein-Distance

KIT
Karlsruher Institut für Technologie

# Elastic Matching

INDUSTRY ⬅ distance ➡ INTEREST

Operations:

• Deleting
• Inserting
• Replacing

Interactive Systems Labs

# Elastic Matching

INDUSTRY

    Deletion of D

INUSTRY

    Deletion of U

INSTRY

    Replace Y with S

INSTRS

    Insertion  of E

INSTERS

    Insertion  of E

INSTERES

    Deletion of S

INTERES

    Insertion  of T

INTEREST

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Elastic Matching

$$d(a,b) = d(a^m, b^n)$$

$$d(a^i, b^{j)}) = \min \quad \begin{matrix} d(a^{i-1}, b^j) + w(a_i, \varnothing) \\ d(a^{i-1}, b^{j-1}) + w(a_i, b_j) \\ d(a^i, b^{j-1}) + w(\varnothing, b_j) \end{matrix}$$

KIT
Karlsruher Institut für Technologie

# Elastic Matching

KIT
Karlsruher Institut für Technologie

# Elastic Matching

# Word Recognition

- Problem

  - Word is a Sequence of Subword Units:
    A Sequence of Letters

  - Do not want to Train at the Word Level

    - Number of Units

    - Training Data Requirement

    - Problem of Adding New Words

KIT
Karlsruher Institut für Technologie

# Word Recognition

Data Capturing

Pre-Processing

Feature Extraction

Segmentation

Classification

Post-Processing

Interactive Systems Labs

**KIT**
Karlsruher Institut für Technologie

# Implicit Segmentation



Interactive Systems Labs

# Word Recognition

Data Capturing

Pre-Processing

Feature Extraction

Segmentation ⟷ Classification

Post-Processing

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Speech Recognition

- Recognizer Components:

# Word and Letter Models

Techniques that use letter models:


- Hidden Markov Models  (HMMs)
  - Model Training
  - Finding the best path (Segmentation)
  - Finding the probability  of a word


- Multi-State Time Delay Neural Networks (MS - TDNNs)

Interactive Systems Labs

# Words and Letter Models

Techniques to help guide the search:

- Language Models (e.g. Bi- or Trigrams)

- Dictionaries

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Word Model

*Suche*

| | |
|---|---|
| **Model A** | ------------------------------ |
| **Model S** | ++++-------------------------- |
| **Model U** | -----++++++------------------- |
| **Model C** | ----------++++--------------- |
| **Model H** | -------------+++++-------- |
| **Model E** | ----------+++----++++++ |
| **Model I** | ----------+++------------- |
| **Model L** | ---------------+++----------- |

Interactive Systems Labs

# Word Model

*Suche*

| | |
|---|---|
| **Model A** | ------------------------------- |
| **Model S** | **++++**----------------------- |
| **Model U** | -----+**+++**++--------------- |
| **Model C** | -----------+**++**+----------- |
| **Model H** | --------------+**+++**+-------- |
| **Model E** | -----------+++----**+++++**    <span style="color:red">Best Path</span> |
| **Model I** | -----------+++--------------- |
| **Model L** | ---------------+++----------- |

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Dictionaries

Model S  Model U  Model C  Model H  Model E

Computation of word probability through
concatenation of letter models

Disadvantage: classificator only recognizes words
from the dictionary

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Speech Recognition

- Recognizer Components:



Interactive Systems Labs

# Pre-Processing

## Size Normalization

Base-line computation:

# Pre-Processing

$$b_0, \ldots, b_n \quad b_i \in IR^d$$

$$b(t) := \sum_{i=0}^{n} b_i B_i^{\,n}(t), \quad t \in [0,1]$$

$$B_i^{\,n}(x) = \binom{n}{i} x^i (1-x)^{n-i}, \quad i = 0, \ldots, n.$$

# Pre-Processing

## Re-sampling:

temporally equidistant points

$\downarrow$

spatially equidistant points

# Pre-Processing

Base-line finding

Rotation-normalizing

Non-linear Interpolation

Smothing

Deletion of delayed Strokes

Re-sampling

Tilt-correction

Interactive Systems Labs

# MS-TDNN Architecture

# MS-TDNN Architecture

1                                                    T′

Sequenz der Merkmalsvektoren

normalisierte Eingabe

1

T′

globale Merkmale

Kontext-Bitmap

Ober-/Unterlänge

vertikale Position

Schreibrichtung

Krümmung

Stiftzustand

"Hut"

Neigung

Erscheinungsbild

Gewelltheit

Linearität

lokale Merkmale

Interactive Systems Labs

KIT

Karlsruher Institut für Technologie

# MS-TDNN Architecture



Modeling:

# MS-TDNN Architecture



Use of Language Models and Dictionaries

# MS-TDNN Architecture

Each word is represented by concatenating its letter models.

# On-line & Off-line

## Connection of two Pixels with a straight line

Idea: Use the horizontally right or diagonally situated pixel, which is closest to the exact

# On-line & Off-line

$$p_1 = (x_1, y_1), p_2 = (x_2, y_2)$$

$$0 \leq y_2 - y_1 \leq x_2 - x_1 \quad x_i, y_i \in IN$$

$$p_1 \qquad p_2$$

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

$$\Delta_x := x_2 - x_1; \quad \Delta_y := y_2 - y_1;$$

$$x := x_1; \quad y := y_1;$$

$$c_y := 2 * \Delta_y; \quad \Delta := c_y - \Delta_x; \quad c_x := \Delta - \Delta_x;$$

$$(x, y)$$

$$x := x + 1;$$

$$\Delta < 0 \qquad \Delta := \Delta + c_y$$

$$y := y + 1;$$

$$\Delta := \Delta + c_x;$$

$$x > x_2$$

# Human-Computer Interaction

Commands
Data

Mouse

Keyboard

Face Tracking

Speech Recognition

Handwriting Recognition

Gesture Recognition

Eye Tracking

Lipreading

Human

Computer

Multi-modal Interface

**Figure 9.** Original.

**Figure 10.** Result.

**Figure 11.** Original image with 15° skew.

**Figure 12.** Segmentation result.

# Handwriting Recognition

*Part II*

## Algorithms and Systems

Alex Waibel

Carnegie Mellon University

University of Karlsruhe

KIT
Karlsruher Institut für Technologie

# Off-line and On-line Handwriting Recognition

## Off-line recognition:



## On-line recognition:



Interactive Systems Labs

# On-line Handwriting Recognition

- Provides <span style="color:red">pen-based input</span> for human-computer interaction

- Pen-based input can be used
  - alone or
  - as part of a multi-modal interface

- Possible applications include
  - form filling
  - editing existing text
  - short notes
  - calendars

- Input consists of <span style="color:red">dynamic writing information</span> (e.g. temporal sequence of data points)

- But: stroke order (within characters or words) influences recognition

Interactive Systems Labs

# Evaluating Handwriting Recognition Systems

- Comparing different handwriting recognizers is difficult

- Performance depends on
  - recognition task (e.g. isolated characters, words, unconstrained text)
  - writing style (e.g. printed, mixed, cursive)
  - size of dictionary
  - intended end user(s)
    - single writer (allows writer dependent system)
    - multi-writer
    - omni-writer (requires writer independent system)

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Handwriting Recognition Tasks

boxed upper case letters

boxed letters

printed words

mixed words

cursive words

unconstrained text

# Character Recognition



Interactive Systems Labs

# On-Line Recognition

- Simplest Case:
  - Single Characters (English, German)
  - Upper Case
  - Position Known
  - On at a time

Interactive Systems Labs

# Elastic Matching



Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Writing Styles

| Printed | Mixed | Cursive |
|---|---|---|
| carolyn | seller | proofs |
| cobler | hungarian | hampers |
| cluff | resignations | sember |

# The Unipen Project

- Initiated by the Technical Committee 11 of the IAPR to
  - constitute a sizable, quality database
  - evaluate the state of the art in on-line handwriting recognition by testing recognizers in the same conditions on several tasks of various difficulty
  - bring together researchers and developers in on-line handwriting recognition from Universities and Industry.

- About 40 participants

  AT&T, Apple Computer Inc., Bolt Beranek and Newman Inc., IBM (NY), Lexicus Corp (CA), NICI (Netherlands), Aachen Technical University (Germany), University of Karlsruhe (Germany), Hewlett-Packard Labs (UK), Philips Research Lab (Netherlands), ...

- About 4 million characters (approx. 500K words) of data donated by participants

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Recognition Strategies

- **Wholistic approach**

  – recognition is performed globally on the whole word

  – no attempt is made to identify characters individually


- **Analytical approach**

  – recognition is performed at an intermediate level

  – words are considered as sequences of smaller units (e.g. letters)

# Wholistic Strategies

Wholistic methods usually follow a two-step scheme:

❶ feature extraction

❶ global recognition by comparing the representation of the unknown word with those of references stored in the dictionary

Practical consequences:

– as letter segmentation is avoided and recognition performed in a global way it is tolerant to deformations that affect unconstrained cursive handwriting

– the recognition is constrained to a specific dictionary of words

– if training on word samples is required, the dictionary cannot be updated automatically from letter information and thus a training step is mandatory to expand or modify the dictionary

Suitable for applications

– with small dictionaries

– where the dictionary is statically defined and not likely to change

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Analytical Strategies

- Analytical strategies deal with several levels of representation
  - Feature level
  - One or more intermediate levels dealing with subparts of words
  - Integration at word level
- **Different kinds of subparts:** letters, graphemes, states, strokes ...
- Units of intermediate level usually are related to letters
- Letter-based recognition is independent from specific dictionary
- Dictionary can be replaced or modified without any retraining
- Analytical approaches fall into two main categories:
  - ❶ with explicit segmentation (input segmentation)
  - ❷ with implicit segmentation (output segmentation)

  which differ in the way (letter) segmentation is performed

# Stroke Level

Decomposition / Identification of
atomic units, building blocks

# Strokes Level

## Finding Rules

# Stroke Level

Disadvantage: complex rule bases

A oder H ?

N oder V ?

KIT
Karlsruher Institut für Technologie

# Explicit Segmentation

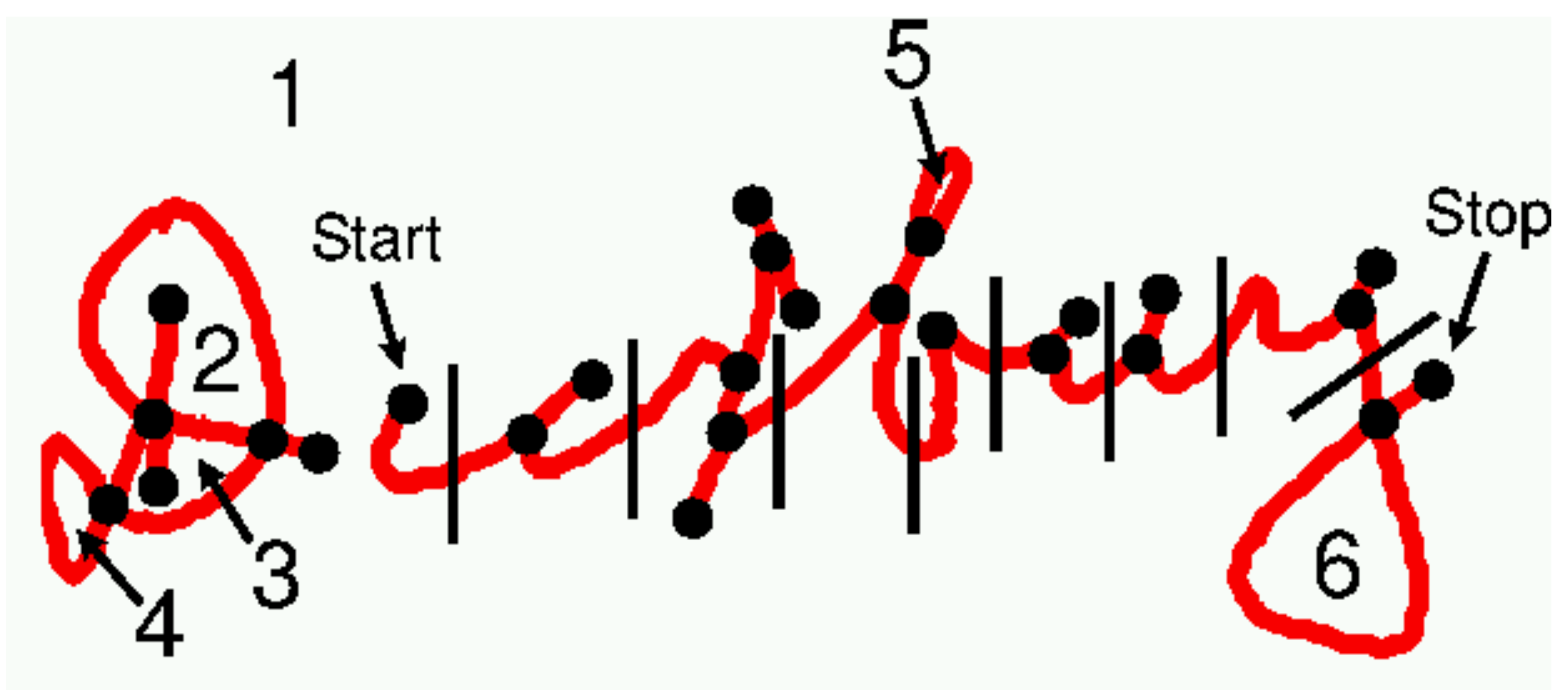

# Elastic Matching

# Input and Output Segmentation

# Input Segmentation

- Following three steps are performed:
  - ❶ external segmentation of the word into smaller units (e.g. letters)
  - ❷ individual recognition of these units
  - ❸ contextual post-processing using lexical, syntactic or semantic knowledge
- Major drawbacks of input segmentation:
  - segment boundaries are often difficult or impossible to find
  - erroneous segmentation may lead to incorrect recognition
- Contextual post-processing can be performed by
  - orthographic correction techniques using statistics of the dictionary (e.g. based on n-gram frequencies)
  - direct comparison with a dictionary (e.g. based on Edit Distance)
- Classical remark:

  *„it is necessary to segment to recognize, but it is also necessary to recognize to segment"*

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Output Segmentation

- A recognition-based segmentation is performed

- The decision about letter boundaries is delayed to the end of the recognition process
  - this avoids the problem of misrecognitions through early segmentation errors

- Contextual knowledge can be introduced
  - in a statistical way (e.g. letter n-gram)
  - by using dictionaries

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Training

## Manual Segmentation:

*Buchstabenebene*

## Automatic Segmentation:

*Wortebene*

# Speech Recognition

- Recognizer Components:

# NPen$^{++}$ - Cursive Handwriting Recognition

Preprocessing

coordinate sequence → Normalization → normalized coordinates → Feature Extraction → sequence of feature vectors

*Hello*

Recognition

recognized text ← Search ← character hypotheses ← TDNN

Hello

The system was designed ...

⇨ ... to be writer independent

⇨ ... to work with large vocabularies

⇨ ... to be fast enough for real-world applications

⇨ ... to make use of the dynamic writing information

# Normalization

- Removing undesired variability from the original pen trajectory
  - baseline normalization, deskewing
  - bezier normalization, smoothing
  - size normalization
  - resampling from temporal to spatial equidistance
  - removing delayed strokes

- The original dimension and temporal ordering of the input signal remain unchanged



Normalization

Sampling

Feature Extraction

Interactive Systems Labs

# Normalization Overview



Sampling

Baseline Normalization

Bezier Normalization

Smoothing

Deskewing

Resampling

Delayed Strokes

Feature Extraction

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Baseline Normalization

- According to the computed baselines the input pattern is rotated to a nearly horizontal orientation:

# Baseline Detection

- Using an EM (Expectation Maximization) algorithm the baseline, centerline, descenderline and ascenderline of the pattern are calculated simultaneously from the local minima and maxima:



- The second degree polynomials

$$f_i(x) = k(x - \widetilde{x})^2 + s(x - \widetilde{x}) + y_i \qquad (i = 0,...,3)$$

are used to approximate these lines, where the parameters $k$ (curvature), $s$ (slant) and $x$ (horizontal displacement) are shared among all four curves. The vertical displacements $y_i$ are given by $y_0 = b - d, y_1 = b, y_2 = b + c, y_3 = b + c + a$

Interactive Systems Labs

# Bezier Normalization and Smoothing

- A Bezier algorithm, which approximates missing data points, is used to compensate for sampling errors

- A moving average window is used for smoothing to remove sampling noise



Interactive Systems Labs

# Skew Normalization

- To ensure a nearly vertical orientation of all characters, the input pattern is normalized according to the skew angle

- The skew is computed from a histogram of all angles between a line segment and the x-axis multiplied with the length of this segment

Angle Histogram

0                    180

Skew Normalization

Interactive Systems Labs

# Size Normalization

- The pattern is rescaled with respect to its current <span style="color:red">coreheight</span>, which is the distance between the <span style="color:purple">baseline</span> and <span style="color:green">centerline</span>

- This ensures that all words have (nearly) the same character size



Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Resampling

- The spatial distance between two successive data points depends on
    - the general sampling rate of the used hardware
    - sampling errors
    - the current writing speed

- Therefore the sequence of data points is resampled from temporal to spatial equidistance



Resampling

Interactive Systems Labs

# Removing Delayed Strokes

- Delayed strokes like i-dots and t-strokes are removed from the sequence of data points, if they do not occur directly after the corresponding character in the temporal sequence



Interactive Systems Labs

# Normalization Overview



Sampling

Baseline Normalization

Bezier Normalization

Smoothing

Deskewing

Resampling

Delayed Strokes

Feature Extraction

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Feature Extraction

- Extraction of features along the normalized pen trajectory, yielding a temporal sequence of n-dimensional feature vectors

- Each feature vector consists of:
  - Local features: writing direction, curvature, position, pen up/down, lineness, aspect, curliness, slope, ...
  - Global features: context bitmaps

KIT
Karlsruher Institut für Technologie

# Feature Extraction

time

0                                                       T

final input representation

normalized
coordinate
sequence:

context bitmaps (a)

local features (b)

t-2   t-1   t   t+1 t+2

(a) context bitmaps

(b) local features

- writing direction

- curvature

- y position

- pen up/down

KIT

Karlsruher Institut für Technologie

# Computing State Hypotheses

Time Delay Neural Network (TDNN) to compute state
hypotheses over time given a feature vector sequence:

# Word Modeling

- Each word $w_i$ of the dictionary $W = \{w_1,...w_k\}$ is represented as its character sequence

$$w_i = c_{i1}c_{i2}...c_{ik}$$

- Each character $c_j$ itself is modelled by a three state hidden markov model

$$c_j = q_j^0\ q_j^1\ q_j^2$$

- The states $q_j^0$, $q_j^1$, $q_j^2$ model the initial, middle and final section of the character's coordinate sequence

- I.e. the final modelling of word $w_i$ is

$$w_i = q_{i0}\ q_{i1}\ ...\ q_{j3k}$$

(e.g. able = $a_0a_1a_2b_0b_1b_2l_0l_1l_2e_0e_1e_2$)

# Word Modelling

- Each word of the dictionary is represented as its character sequence, where each character itself is modeled by a three state hidden markov model
- The character's states model the initial, middle and final section of the character's coordinate sequence

E.g. the modeling for word "able" is

a            b            l            e

$a_0 \rightarrow a_1 \rightarrow a_2 \rightarrow$    $b_0 \rightarrow b_1 \rightarrow b_2 \rightarrow$    $l_0 \rightarrow l_1 \rightarrow l_2 \rightarrow$    $e_0 \rightarrow e_1 \rightarrow e_2 \rightarrow$

# MS-TDNN Architecture

# MS-TDNN Architecture



Use of Language Models and Dictionaries

# Flat Search Approach

Given the probabilities of the states the score for each word in the dictionary is defined to be a Viterbi approximation of the log likelihoods of the feature vector sequence:

# Decoding

The Viterbi Algorithm:

- Find the state sequence **Q** which maximizes **P(O, Q | λ )**
- Similar to Forward Algorithm except **MAX** instead of **SUM**

$$VP_t(i) = MAX_{q_0, \cdots q_{t-1}} \; P(O_1 O_2 \cdots O_t, q_t{=}i \mid \lambda \; )$$

Recursive Computation:

$$VP_t(j) = MAX_{i=0, \ldots, N} \; VP_{t-1}(i) \; a_{ij} b_j(O_t) \qquad t > 0$$

$$P(O, Q \mid \lambda \; ) = VP_T(S_N)$$

Save each maximum for backtrace at end

# Viterbi Trellis

$$\begin{bmatrix} A & 0.8 \\ B & 0.2 \end{bmatrix}$$

0.6

1.0

0.4

$$\begin{bmatrix} A & 0.3 \\ B & 0.7 \end{bmatrix}$$

**Initial**      **Final**

|  | **A** |  | **A** |  | **B** |  |
|---|---|---|---|---|---|---|
| | **t=0** | | **t=1** | | **t=2** | **t=3** |
| **state 1** | 1.0 | 0.6 * 0.8 | 0.48 | 0.6 * 0.8 | 0.23 | 0.6 * 0.2 | 0.03 |
| | | 0.4 * 0.3 | | 0.4 * 0.3 | | 0.4 * 0.7 | |
| **state 2** | 0.0 | 1.0 * 0.3 | 0.12 | 1.0 * 0.3 | 0.06 | 1.0 * 0.7 | 0.06 |

Interactive Systems Labs

# Flat Search

# Tree Architecture

# Number of Nodes

# Compressed Tree Architecture

- Linear lists of nodes can be merged into single nodes
- This results in fewer active nodes during search

# Tree Search Algorithm

Problem:

- The tree structure in itself does not yield enough of a benefit with respect to run-time efficiency

- There is still a linear scaling of run-time with the dictionary size

Solution:

- We give up evaluating all nodes (HMMs) of the search tree

- This leads to active and inactive nodes and a set of pruning rules which specify when to turn on an inactive node and when to turn off an active one

# Initialization and Basic Concept

- The search is initialized by setting all root nodes to be active and all other nodes (internal and word end nodes) to be inactive

- Then for each frame of the states layer of the neural network the algorithm goes through steps 1-4

**Step 0:** Initialization

**Step 1:** Evaluation

**Step 2:** Pruning

**Step 3:** Expansion

**Step 4:** Word Transition

# Tree Search Algorithm for One Frame

① Evaluation

- For each active node compute a viterbi step to find the accumulated scores $s_{ij}$ for the next frame
- Compute the best state score $s_i$ within each node and the best score
  - $s' = \max s_i$
- over all evaluated nodes

② Pruning

- Deactivate all currently active nodes in the search tree where the following criterion is fullfilled
  - $s_i < s' - beam$
- I.e. all nodes whose best accumulated score is below a certain threshold will become inactive in the next frame

③ Expansion

- For each currently active node test whether a transition from its last state to the first state of any child node leads to a higher accumulated score in the first state of that child  node
- If that holds and the new score is above the pruning threshold the child node is marked to be active in the next frame

④ Word Transition

- For each active word end node we test the transition from that node to any of the tree root nodes as we did in the expansion step above

KIT
Karlsruher Institut für Technologie

# Adjusting Beam Sizes

- The beam size influences both the recognition accuracy and recognition time:
  - smaller beam means increased speed and decreased accuracy
  - larger beam means decreased speed and increased accuracy

- Therefore the beam has to be adjusted according to the particular needs: e.g. maximum beam for evaluations, smaller beam for run-time version of the system

- For isolated word recognition (i.e. no word transitions allowed) only one single beam for all tree nodes is used

- For continuous recognition (sentence recognition) separate beams for different node types can be used. E.g.:

  a beam for all root and internal nodes

  a beam for all word end nodes

  a beam for word transitions

# Summary – On-Line Recognition

- A tree architecture reduces the number of nodes (HMMs) to be evaluated to approx. 1/3

- The tree architecture itself does not improve recognition time compared to a flat search approach

- But combined with efficient pruning techniques the search space can be reduced dramatically

- The benefit in run-time is much higher than the small decrease in recognition accuracy

# Experiments (Isolated Words)

- Task:
    - isolated english words
    - writer independent
    - no restrictions on writing style
    - dictionary sizes ranging from 1,000 to 100,000 words

- Database:
    - collected at University of Karlsruhe and Carnegie Mellon University
    - mixture of german and english writers
    - 307 different writers (13,000 words)
        - 204 writers used for training (9,000 words)
        - 103 writers used for testing  (4,000 words)

KIT
Karlsruher Institut für Technologie

# Recognition Rates

Recognition of single symbols:

Npen++:   0-9    96.5%

              A-Z   92.7%

              a-z   91.1%

# Recognition Rates

Recognition of words:

On-line recogniton rate is higher than off-line recognition rates:

Off-line:  95%  on       500 words
On-line:  95%  on   5.000 words
                90%  on  50.000 words

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Results (Isolated Words)

Recognition time versus dictionary and beam sizes

Recognition accuracy versus dictionary and beam sizes

# Experiments (Sentences)

- Task:
  - sentences from Wall Street Journal (WSJ)
  - writer independent
  - no restrictions on writing style
  - 20,000 word WSJ dictionary

- Database:
  - collected at University of Karlsruhe and Carnegie Mellon University
  - mixture of german and english writers
  - training set:
    - ~ 20,000 isolated words
    - ~ 10,000 sentences (~ 70.000 words)
  - test set:
    - ~ 200 sentences from WSJ (from 40 different writers)

# Results (Sentences)

- 64% word recognition rate (without any language model)

- 85% word recognition rate (using a WSJ language model)

The word recognition rate is defined as
(#words - #insertions - #deletions -#substitutions) / #words

Remarks:

~ 90% of the errors are substitutions

~ 50% of these substitutions are caused by capitalization
errors and ending-"s"

# Handwriting Recognition

- On-line and Off-line recognition techniques merge:
  - extracting on-line information from static off-line data (?)
  - using off-line techniques in on-line systems
  - combining off-line and on-line recognizers (multiple experts)

- Issues:
  - Techniques: Stroke-based recognition, Context Dependence, Adaptations, ......
  - Benchmarks
  - Repair
  - The New Word Problem
  - Other symbol sets: Formulas, Gestures, Drawings, Icons, ..
  - Abbreviations
  - Combinations of symbol sets... *All* Pen Based Activity
  - Cursive off-line handwriting

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Systems

- Individual Characters
  - Many PDA's (particularly Japan, China)
  - Simplified Characters:  Graffiti

- Continuous Handwriting
  - TDNN Based:
    - Npen++
    - Similar Architectures: LeCun, Rumelhart
  - HMM Based:
    - BBN, …

# The Graffiti Recognition System



Interactive Systems Labs

# Off-Line Recognition

- Optical Character Recognition (OCR)
  - Mostly Printed Text

- Extract Text Sequence
  - Document/Scene Analysis
  - Extract/Scan Bitmap
  - Extraction and Computation of Features
  - Representation is Spatial, *not* Termporal

- Recognition Algorithm
  - Deal with Shift Invariance
  - Integrate Characters into Words (segmentation!)

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Off-line Handwriting Recognition

- Possible applications include
  - check reading
  - postal address reading
  - document analysis, ...

- Input consists of scanned <span style="color:red">bitmaps</span> without any temporal information

- Eventually location of handwriting needs to be found (document analysis)

- Stroke order doesn´t influence recognition

- But: problems through overlapping or touching characters and noisy input

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Document Analysis



**Figure 9.** Original.

**Figure 10.** Result.

**Figure 11.** Original image with 15° skew.

**Figure 12.** Segmentation result.

# Off-line Preprocessing

**Binary Connected Components**

**Baseline Normalization**

**Deskewing**

**Skeletonization**

**Approximation**

**Feature Extraction**

Interactive Systems Labs

# LeNet Architecture



Reference: Yann LeCun et al. „Handwritten Zipcode Recogntion with Multi-layer Networks", Proceedings of the ICPR-90, Atlantic City, 1990

Interactive Systems Labs

# Rummelhart Architecture



Reference: J. Keeler, D. E. Rummelhart. „A Self-Organizing Integrated Segmentation and Recognition Neural Net", Advances in Neural Information Processing Systems, Morgan Kaufman, 1991.

Interactive Systems Labs

# Comparison of Classifier Methods



- Reference: L. Bottou et al. „Comparison of Classifier Methods: A Case Study in Handwritten Digit Recognition", Proceedings of the ICPR-94, Jerusalem 1994)

Interactive Systems Labs

# Memory Requirements

# Training and Run Time



Interactive Systems Labs

# The Language Challenge
# Signs, Visual Text



Interactive Systems Labs

# Signs



Interactive Systems Labs

# Challenges of Sign Translation

- Detection
  - No motion information
  - No approximate shape and position information
  - No color reflectance information

- Recognition
  - Blured image
  - Low resolution
  - Character deformation
  - Reflections, Colors, Fonts

- Translation
  - Short sentence/Phrase
  - Abbreviation

Interactive Systems Labs

KIT
Karlsruher Institut für Technologie

# Sign Translator



Interactive Systems Labs

# Sign Translator - Architecture



Interactive Systems Labs

# Sign Detection

- A Hierarchical Approach
    - A multi-resolution edge detection algorithm
    - Adaptive searching in the neighborhood of initial candidates based on layout syntax
    - Layout analysis of the detected sign areas

# Character Recognition

- Intensity-based Approach
  - Feature from Gabor Transformation
  - LDA for the feature selection

- Result
  - Character set: 3755 Level 1 Chinese characters in Chinese national standard character set GB2312-80
  - Accuracy: 92.4%

# Some Results



Interactive Systems Labs

Interactive Systems Labs