

Reconocimiento de Escritura

Daniel Keysers

Image Understanding and Pattern Recognition
German Research Center for Artificial Intelligence (DFKI)
Kaiserslautern, Germany

Mar-2007

OCR - Introduction

OCR fonts

Tesseract

Sources of OCR Errors

OCR - Introduction

OCR fonts

Tesseract

Sources of OCR Errors

- ▶ OCR = Optical Character Recognition
- ▶ steady progress in OCR since the mid-fifties
- ▶ 1975: IBM Optical Page Reader cost over 3,000,000 US\$ and displaced several dozen keypunch operators
- ▶ applications
 - ▶ home or office use
 - ▶ forms processing
 - ▶ address reading
 - ▶ conversion of large archives of text to computer-readable form

Although researchers have worked on the problem of OCR for at least thirty years, there has been a renewed interest in OCR technology in the recent years. This is partly due to

- ▶ the increasing need for efficient information storage and retrieval,
- ▶ the increasing need for cross-language information access, and
- ▶ the dramatic drop in scanner prices
- ▶ [large scale digitization projects]

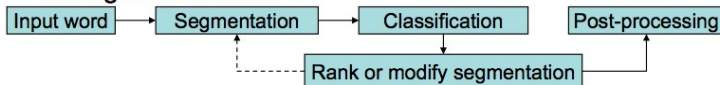
(Kanungo, 1998)

- OCR paradigms: [Casey 96]

- Dissection (Segmentation driven OCR):



- Recognition driven:

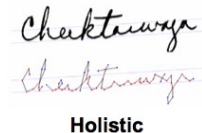


- Holistic:



Input	Window	Remainder	
mm	m m	1 null	Reject
mm	r r	m m	Accept

Recognition driven



building blocks of an OCR system:

- ▶ noise removal
- ▶ skew estimation
- ▶ page segmentation & line detection (= layout analysis)
- ▶ segmentation
- ▶ character classifier
- ▶ language modeling

What are the main differences to a handwriting recognition system?

even 99% accuracy
= 30 errors on a typical printed page of 3000 characters

“in almost every application, either the OCR results must be corrected by a human operator or a significant fraction of the documents are rejected in favor of operator entry” (Nagy⁺ 2000)

2 Browser and Design Testing

There are multiple implementations of HTML rendering engines; some common ones are Microsoft's Internet Explorer, Mozilla's Gecko, Apple's Safari, Opera's browser, and KDE's KHTML. Each of these render web pages differently due to bugs and incomplete specifications of web standards. Common defects are missing text, text that is unintentionally rendered overlapping, text that unintentionally overlaps graphical elements, bad font substitutions, bad spacing, and unreadable choices of foreground and background colors.

Our approach to this problem is to render the HTML into an image-based representation and then subject the image-based representation to OCR (including layout analysis). Common rendering problems can be detected by comparing the HTML input against the OCR output. For example, incorrect rendering due to missing text, overlapping text, bad font substitutions, and text in invisible colors can be flagged by detecting text that is present in the original HTML but missing in the OCR output. Incorrect layouts can be detected by comparing the paragraph structure and reading order of the original HTML against the layout analysis output.

We automate this process by using user interface scripting support. Our initial prototype is implemented on Apple Macintosh OS X, where we use a combination of AppleScript, the Firefox ScreenGrab extension, and the Safari Snagit extension to automatically capture web pages with different browsers, versions, and browser settings, and to send the captured page to be processed by the OCR system; analogous technologies exist for Windows and Linux. This way, a large collection of web pages can be rendered, analyzed, and verified without the need for operator intervention in different browsers and browser versions. The approach can detect HTML layout problems without prior assumptions about layout engines and for browsers returning a wide variety of layouts; the approach correctly distinguishes incorrect and correct layouts even in the presence of JavaScript and style sheets.

The same approach can be used for checking HTML layouts against design rules. Design rules for HTML are intended to ensure readability, accessibility, and easy interpretation by readers, as well as to ensure correct representations of organizational identity. Design rules specify such features as minimum font sizes, acceptable fonts and color choices, and minimum spacings between logical groupings of page elements.

3 Phishing and Search Engine Spam

Phishing (www.antiphishing.org) is a problem in which an adversary attempts to obtain personal and private information by creating E-mails and web pages that belong to a trusted organization (e.g., the recipient's bank), but actually transmit the information to the adversary. Closely related to phishing is search engine spam, where a web site will present HTML content to a search engine that gives indications to the search engine that the web site contains relevant information about a popular topic but actually renders as

2 Browser and Design Testing

There are multiple implementations of HTML rendering engines; some common ones are Microsoft's Internet Explorer, Mozilla's Gecko, Apple's Safari, Opera's browser, and KDE's KHTML. Each of these render web pages differently due to bugs and incomplete specifications of web standards. Common defects are missing text, text that is unintentionally rendered overlapping, text that unintentionally overlaps graphical elements, bad font substitutions, bad spacing, and unreadable choices of foreground and background colors.

Our approach to this problem is to render the HTML into an image-based representation and then subject the image-based representation to OCR (including layout analysis)...

44

LECTURE 4. EUROPEAN OPTIONS IN COMPLETE MARKETS

Indeed, it follows from (3.5') in view of (4.30) that

$$\begin{aligned} (4.40) \quad E_{t-1}^Q(U) X_{t-1}^{S_0} &= (1-\alpha)C + \sum_{i=1}^N E_{t-1}^Q(U) \gamma_i^S S_{t-1} (p_i - r_k) \\ &= (1-\alpha)C + \sum_{i=1}^N E_{t-1}^Q(U) S_{t-1} \gamma_i^S (p_i - r_k) \\ &= \sum_{i=1}^N E_{t-1}^Q(U) S_{t-1} \gamma_i^S C (p_i - r_k) \\ &= (1-\alpha)C - (1-\alpha)C + M_{t-1}^Q = C M_{t-1}^Q. \end{aligned}$$

From (4.40),

$$E_{t-1}^Q(U) X_{t-1}^{S_0} = E_{t-1}^Q(U) f - C \mathbf{1}_{\{S_{t-1} < K\}}$$

and hence

$$E_{t-1}^Q(U) X_{t-1}^{S_0} \geq E_{t-1}^Q(U) f - C.$$

The last inequality means that $s_{t-1} \in \text{SF}(f, N)$.

Further, it follows from (4.38) that

$$(4.41) \quad \mathbf{P}^*(X_{t-1}^{S_0} \geq f) = \mathbf{P}^*(f - C E_N(U) \mathbf{1}_{\{S_{t-1} < K\}} \geq f) \\ = \mathbf{P}^*(\mathbf{1}_{\{S_{t-1} < K\}} \leq 0) = \mathbf{P}^*(S_{t-1} \geq K) = (1-\alpha).$$

Finally, we get from (4.38) and (4.41) that

$$(4.42) \quad \begin{aligned} \mathbf{P}^*(X_{t-1}^{S_0} \geq f) &= \mathbf{E}^*[\mathbf{1}_{\{X_{t-1}^{S_0} \geq f\}} E_N] \\ &\geq \mathbf{E}^*[\mathbf{1}_{\{X_{t-1}^{S_0} \geq f\}} \mathbf{1}_{\{S_{t-1} \geq K\}} E_N] \\ &\geq \lambda(1-\alpha) \geq 1-\alpha. \end{aligned}$$

The relations (4.41) and (4.42) show that the condition (4.35) holds for the strategy s_{t-1} , and hence s_{t-1} is an α -($1-\alpha$)- C -(f, N)-hedge.

What has been obtained shows that it is possible to hedge a contingent claim with a specified probability ($1-\alpha$). Further, the initial funds can be reduced by the amount αC , though with a risk α the accepted contingent claim cannot be repaid.

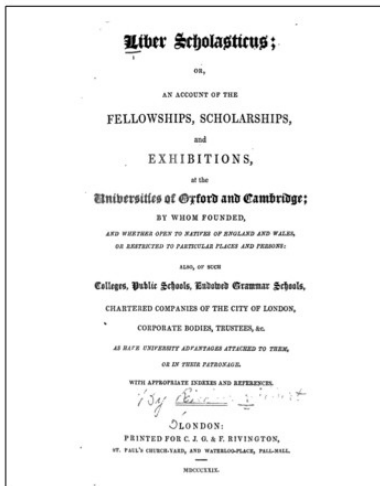
PROBLEMS

4.1. Prove that on a no-arbitrage (B, S)-market we have for a standard European option to buy (sell) that $C(N_1) \geq C(N_2)$ (respectively, $P(N_2) \geq P(N_1)$) when $N_2 \geq N_1$.

4.2. Prove that the fair price $C = C(N, S_0, K)$ of a standard European option to buy, where N is the exercise time, S_0 is the initial price of a share, and K is the exercise price, has the following properties:

- $C(S_0, K)$ is monotone in S_0 and K ;
- $C(S_0, K)$ is convex in S_0 and K ;
- $C(\lambda S_0, \lambda K) = \lambda C(S_0, K)$ for $\lambda > 0$.

Indeed, it follows from (3.5') in view of (4.39) that (4.40) $S_{t-1} = (1-\alpha)C - (1-\alpha)C +$
From (4.40), $\diamond^A(u)x^A = \diamond$ and hence the last inequality means that $\text{TTQ} \in \text{SF}(f, N)$. Further, it follows from (4.38) that (4.41) $\mathbf{P}^*(X_{t-1}^{S_0} \geq f) = \mathbf{P}^*(f - C \diamond N(U) \mathbf{1}_{\{Z_N\}} \geq f) = \mathbf{P}^*(\mathbf{1}_{\{Z_N < 0\}} = \mathbf{P}^*(Z_N > A) = (1-\alpha)$. Finally, we get from (4.38) and (4.41) that (4.42) $\mathbf{P}^*(X_{t-1}^{S_0} \geq f) = \mathbf{E}^*[W \diamond > A | (1-\alpha) > 1-\alpha]$. The relations (4.41) and (4.42) show that the condition (4.35) holds for the strategy n_a , and hence TTQ is an α -($1-\alpha$)- C -(f, N)-hedge. What has been obtained shows that it is possible to hedge a contingent claim with a specified probability ($1-\alpha$). Further, the initial funds can be reduced by the amount αC , though with a risk α the accepted contingent claim cannot be repaid. PROBLEMS 4.1. Prove that on a no-arbitrage (B, S)-market we have for a standard European option to buy (sell) that $C(N_1) \geq C(N_2)$ (respectively, $P(N_2) \geq P(N_1)$) when $N_2 \geq N_1$. 4.2. Prove that the fair price $C = C(N, S_0, K)$ of a standard European option to buy, where N is the exercise time, S_0 is the initial price of a share, and K is the exercise price, has the following properties: a) $C(S_0, K)$ is monotone in S_0 and K ; b) $C(S_0, K)$ is convex in S_0 and K ; c) $C(\lambda S_0, \lambda K) = \lambda C(S_0, K)$ for $\lambda > 0$.



OR,
AN ACCOUNT OF THE
FELLOWSHIPS, SCHOLARSHIPS,
and
EXHIBITIONS,
at the
atttttonvitto of <C2><A9>Tforfc anft
<E2><82><AC>amftitUO
BY WHOM FOUNDED,
J.VJ> UHKriKK OPEff TO IfATIFES OF
SNOLAND AND WALES,
Ott RKITRICTEU TO PARTICULAR
PLACES AND PERSONS;
ALSO, OF SUCH
CoKrgs, IJutir \$rf)ool6, Kniutotti (Grammar
5rf)ool
CHABTERED COMPANIES OF THE CITY
OF LONDON,
CORPORATE BODIES, TRUSTEES, &c.
At BarS OXIrESSJTY ADrANTAOES
ATTACHED TO TBEX,
OS IN THEM PATRONAGE.
WITH APPROPRIATE INDEXES AND
REFERENCES.
^LONDON:
PRINTED FOR C. J. O. & F. RIV1NGTON,
. PAI L's Clif RCII.YAlin, AMI
WATERLOO.PLACE, PALUMALL.
MDCCCXXIX.

PROLOGO.

Voy á leerle unos manuscritos, que mas desvelos costó á mi padre el sustraerlos á tu curiosidad, que el escribirlos. Sé que cometo una imprudencia satisfaciendo un femenino deseo que te acarreará muchos dolores; pero contigo mas quiero pecar de tolerante que de severo. Profanaré con el secreto la memoria de mi buen padre, mas añadiré quilates á tu cariño: entre los respetos debidos á la memoria de un padre muerto, y el amor

â¬*MInv-

Toy aleertennof n^m^ritn. qaeva* desTdos eosto
4 mi padre d snstnerlos a tu oniosidad, qae d escri-
birlos. Se" que cometa ana impradtncia ilirfirirÂ«dn on
femenil deseo que te aearreara modiM dokns; pcro ew-
tigo mas quiero pecar de tolerant* que de wrcro. Pra-
fanart COD el secrete la memoria de mi boen padre.
mas anadirt qoilates a tu carioo: eatre 1Â« respeto* de-
bidos a. la memoria de on padre nmerlo, j d amor

commercial:

- ▶ Abbyy (FineReader)
- ▶ Nuance/Scansoft (OmniPage)
- ▶ Océ (RecoStar)
- ▶ Iris (ReadIris)
- ▶ many smaller vendors

open source:

- ▶ ocrad
- ▶ gocr
- ▶ Tesseract
- ▶ OCRopus

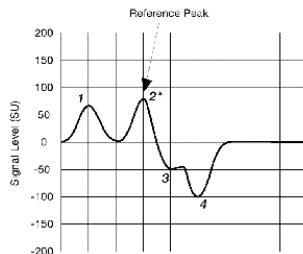
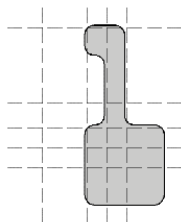
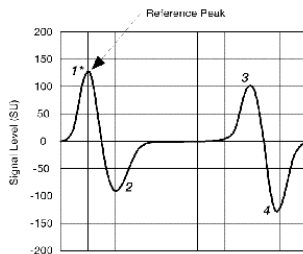
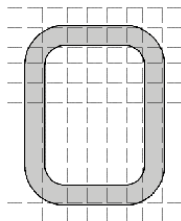
OCR - Introduction

OCR fonts

Tesseract

Sources of OCR Errors

If you can, try to make the problem easier. This is usually easier than improving the classifier.



two groups of OCR fonts

- ▶ Magnetic Ink Character Recognition (MICR)
- ▶ Optical Character Recognition (OCR)
- ▶ artificial distinction

history:

- ▶ financial world
- ▶ bar-code not easily readable by humans
- ▶ magnetic ink preferred

(source: D. Winter)



- ▶ first font used in automated banking
- ▶ digits and four special symbols used in banking"
"dash", "amount", "on*us" and "transit"
- ▶ The font was specially designed so that magnetic pulses would be read unambiguously. That is the reason for some of the heavy black features of some symbols. For this purpose a grid of 7 by 11 squares was used that either had to be white or black. The filling was done so that a magnetic scan would give a pulse signal that was very distinctive.



- ▶ digits, the letters, and five special symbols
- ▶ can also be seen as a barcode:
each symbol encoded by seven vertical bars separated by small or large spaces

2	3	4	5	6	7	8	9	:
;	<	=	>	?	@	A	B	C
D	E	F	G	H	I	J	K	L

- ▶ end of the sixties: full character recognition
- ▶ first font: OCR-A

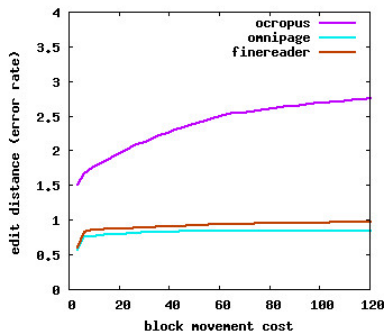
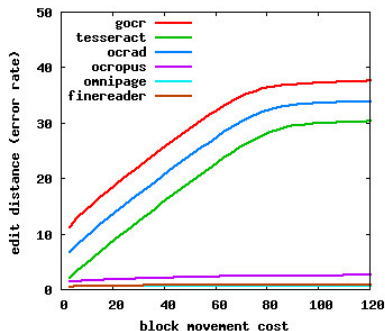
2	3	4	5	6	7	8	9	:
;	<	=	>	?	@	A	B	C
D	E	F	G	H	I	J	K	L

- ▶ accompanies EAN barcodes
- ▶ more symbols exist

- ▶ proposed open standard for representing OCR results
- ▶ motivation: existing formats have limitations in
 - ▶ multi lingual capabilities
 - ▶ typographic phenomena
 - ▶ separate formats for intermediate and final results
- ▶ goal: reuse as much existing technology as possible
- ▶ main idea: represent various aspects of OCR output
 - ▶ logical structuring
 - ▶ typesetting
 - ▶ character information
 - ▶ etc.
- ▶ HTML microformat

- ▶ usually: edit- or Levenshtein distance (cp. Speech Recognition)
- ▶ in the presence of reading order errors: allow block movements at certain cost (cp. Machine Translation)
- ▶ interesting problem: predict OCR accuracy from a set of document image measurements

- ▶ Layout Analysis → see previous lecture
- ▶ Character Recognition Engine:
 - ▶ based on Tesseract → discussed here
 - ▶ based on segmentation → see lecture on handwriting recognition
- ▶ flexible architecture

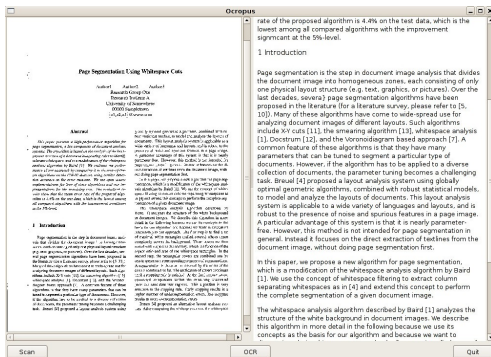


- ▶ 'OCRopus' is the DFKI open source OCR system
- ▶ framework for layout analysis and OCR integration
- ▶ DFKI layout analysis + DFKI or HP-labs 'Tesseract' OCR
- ▶ preliminary evaluation on 18 documents, 15K words
- ▶ goals: improving adaptivity to new fonts and layouts

OCRopus Open Source OCR System — Release



- ▶ first version (technology preview) to be released soon
- ▶ other application: screen OCR



antialiasing

binarization works poorly



colors

very small x-height

- ▶ motivation:
 - ▶ image-based HTML analysis
 - ▶ image-based cut and paste

- ▶ character recognizers:
 - ▶ HMMs
 - ▶ Tesseract

OCR - Introduction

OCR fonts

Tesseract

Sources of OCR Errors

R. Smith: 'An overview of the Tesseract OCR Engine', submitted to ICDAR 2007, personal communication.

- ▶ developed at HP between 1984 and 1994
- ▶ obtained good results at the 1995 UNLV Annual Test of OCR Accuracy
- ▶ then, development was stopped while other commercial OCR engines improved
- ▶ In late 2005, HP released Tesseract for open source.

- ▶ Layout Analysis was separate, therefore not included
- ▶ (OCROPUS closes this gap by including the possibility to use DFKI layout analysis with the Tesseract recognition engine)
- ▶ only supports US-ASCII
- ▶ connected component analysis
- ▶ stored as outlines ('blobs')
 - enable recognition of inverse text easily
- ▶ blobs → text-lines

distinguish fixed pitch and proportional text



fixed-pitch: chopping simple

**of 9.5% annually while the Fed-
erated junk fund returned 11.9%
*fear of financial collapse,***

proportional: Measure gaps in a limited vertical range between the baseline and mean line. Spaces close to the threshold are made fuzzy, so that a final decision can be made after word recognition.

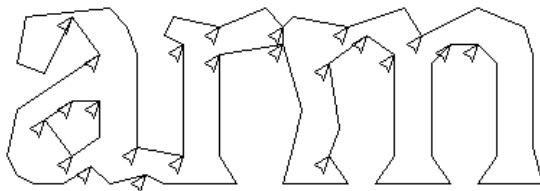
first pass:

- ▶ attempt to recognize each word
- ▶ 'good' words are used as adaptation data

second pass:

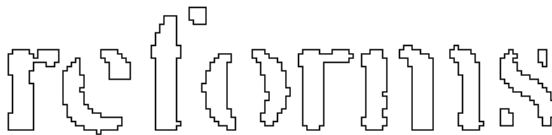
- ▶ recognize words that were not recognized well again
- ▶ use adaptation data

- ▶ does not need de-skewing
- ▶ filter large and small blobs out
- ▶ fit remaining blobs to parallel text line model
- ▶ re-assign left out blobs
- ▶ fit baselines as splines



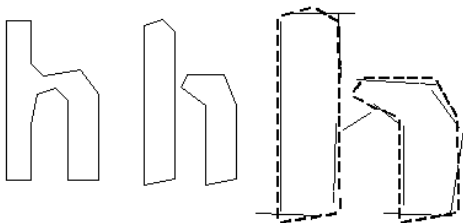
While the result from a word is unsatisfactory, Tesseract attempts to improve the result by chopping the blob with worst confidence.

Candidate chop points are found from concave vertices of a polygonal approximation of the outline, and may have either another concave vertex opposite, or a line.



A* search of the segmentation graph of possible combinations of the maximally chopped blobs into candidate characters without actually building the segmentation graph, but instead maintaining a hash table of visited states

character classifier can recognize broken characters directly



- ▶ match outlines of test many-to-one
- ▶ test features 3-dimensional, (x, y position, angle), typically 50-100
- ▶ prototype features are 4-dimensional (x, y, position, angle, length), typically 10-20 features in a prototype
- ▶ hierarchical classification and use of lookup-tables for speed-up



baseline and moment normalization used in
adaptive and static classifier

- ▶ no need for broken characters in training
- ▶ 20 samples of 94 characters from 8 fonts in a single size, but with 4 attributes (normal, bold, italic, bold italic), making a total of 60,160 training samples
- ▶ other classifiers: often more than 1,000,000 training samples

very basic linguistic analysis:
look-up in different dictionaries

OCR - Introduction

OCR fonts

Tesseract

Sources of OCR Errors

Imaging Defects

- Heavy Print
- Light Print
- Heavy and Light Print
- Stray Marks
- Curved Baselines

Similar Symbols

- Similar Vertical Symbols
- Other Similar Symbols

Punctuation

- Commas and Periods
- Quotation Marks
- Special Symbols

Typography

- Italics and Spacing
- Underlining
- Shaded Backgrounds
- Reverse Video
- Unusual Typefaces
- Very Large Print
- Very Small Print

(Nagy⁺ 2000)

RECOMMENDATIONS

RecDmENSATIONS
RECO-HENATIDNS
REC-ENDATIOffS

Fig. 1a. The Siamese **M**'s baffle all three systems.
Small pieces are missing from the second **M**, the first **N**, and the **D**.

sites

site.
sites
sitea

Fig. 1b. A few pixels can make all the difference. Each system recognizes the first **s**,
but only one can identify the second. Perhaps comparing the two **s**'s would help.

1/4-lb.

I /4 lb.

1/4-lb.

1/4-lb.

Fig. lb.

Fig. lb.

Fig. lb.

Fig. lb.

Proc. 1st Cong.

Proc. I st Cong.

Proc. I st Cong.

Proc. 1st Cong.

4,500 3,500
4.500 3,500
4.500 3,500
4.500 3,500

4,300 residents
4.300
4.300
4.300

Fig 3. All three devices conspire for a thousand-fold reduction of 4,500 and 4,300, but a few extra pixels save 3,500.

Office of the Dean

OfficeoftheD=n
O]7IceoftLyeDe6m
Office#the Dean

Congressional

Conaresslional
mmmmnimhmgm
Conaresslional

Reporters

J?epoztez'
cReIbozeu
CRepoziezi

Fig. 4. Kerned italics cause a variety of problems, include the omission of interword blanks. Underscoring is also difficult. All the devices are bewildered by the enormous flourish on the capital **R** (truly a 'majuscule'), and by the z-like appearance of the lower-case **r**.

(Nagy⁺, 2000)

- ▶ improved image processing
- ▶ adaptation of the classifier to the current document
- ▶ multi-character recognition
- ▶ increased use of context
- ▶ [combination of multiple OCR results]

- ▶ OCR systems generally optimized for average performance over large sets of characters of different
 - ▶ fonts
 - ▶ sizes
 - ▶ scan qualities
- ▶ improve performance of character recognizer by using style information

17 77

- ▶ modeling the sample distribution [Breuel 2001]
 - ▶ each document contains only a single style
 - ▶ modeling sample distribution as a mixture of Gaussians
- ▶ hierarchical Bayesian approach [Mathis et al. 2002]
 - ▶ each document contains a small number of fonts
 - ▶ estimation of prior distributions of a style variable
- ▶ style constrained classifiers [Sarkar et al. 2005]
 - ▶ style consistency constraint is hidden variable
 - ▶ combination of class and style represented by Gaussian mixture model
 - ▶ each mixture component trained by estimating it directly from set of samples from a specific class and style
- ▶ maximum likelihood linear regression [Senior et al. 1997]
 - ▶ based on work on speaker adaptation [Legetter 1995]
 - ▶ adaptation using linear transformations