# Document Recognition Without Strong Models

## Henry S. Baird

*Computer Science & Engineering*

*Lehigh University*

*(Based on work by & with T. Pavlidis, T. K. Ho, D. Ittner, K. Thompson, G. Nagy, R. Haralick, T. Hong, T. Kanungo, P. Chou, D. Lopresti, G. Kopec, D. Bloomberg, A. Popat, T. Breuel, E. Barney Smith, P. Sarkar, H. Veeramachaneni, J. Nonnemaker, and P. Xiu.)*

# How to Find Good Problems?

*When I was finishing my Ph.D. dissertation,*
*my advisor Ken Steiglitz said to me:*

**"There are a lot of smart people out there who,**
**if you hand them a hard problem,**
**they can solve it.**

**But, <u>picking</u> good problems is a rarer skill."**

*At Bell Labs in 1984, I was free to choose any problem I liked…*

# Document Image Recognition?

I had been interested for years in Computer Vision.
I asked myself:  what seems to be <u>missing</u> ….?

<u>Strategic problem</u>:

**Vision systems were brittle:**

**overspecialized  & hard to engineer.**

Theo Pavlidis & I debated, & decided:
We'd try to invent highly <u>versatile</u> CV systems.
<u>Tactical goal</u>:   **Read *any* page of printed text.**
Open, hard, potentially useful…

But, could this help solve the strategic problem?  (DARPA had doubts…)

# Versatility Goals

- Try to guarantee high accuracy across any given set of:
  - symbols
  - typefaces
  - type sizes
  - image degradations
  - layout geometries
  - languages & writing systems
- First step:   a 100-typeface, full-ASCII classifier
- Automate everything possible:
  - emphasize machine learning (avoid hand-crafted rules)
  - identify good features semi-automatically
  - train classifiers fully automatically
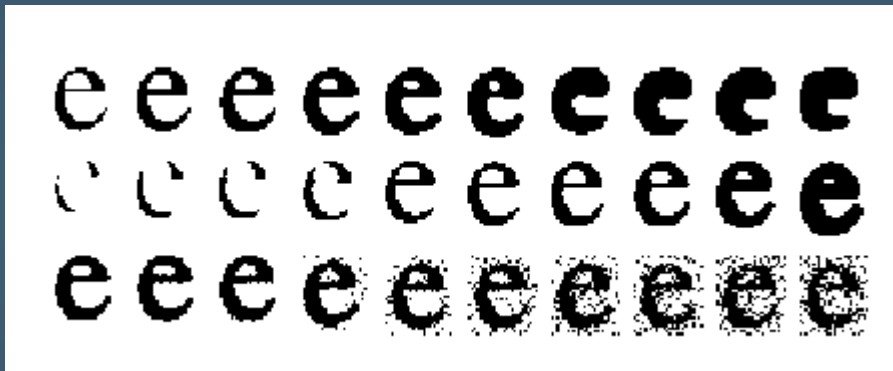  - model image quality, then generate synthetic training data

Pavlidis, Baird, Kahan, & Fossey (1985-1992)
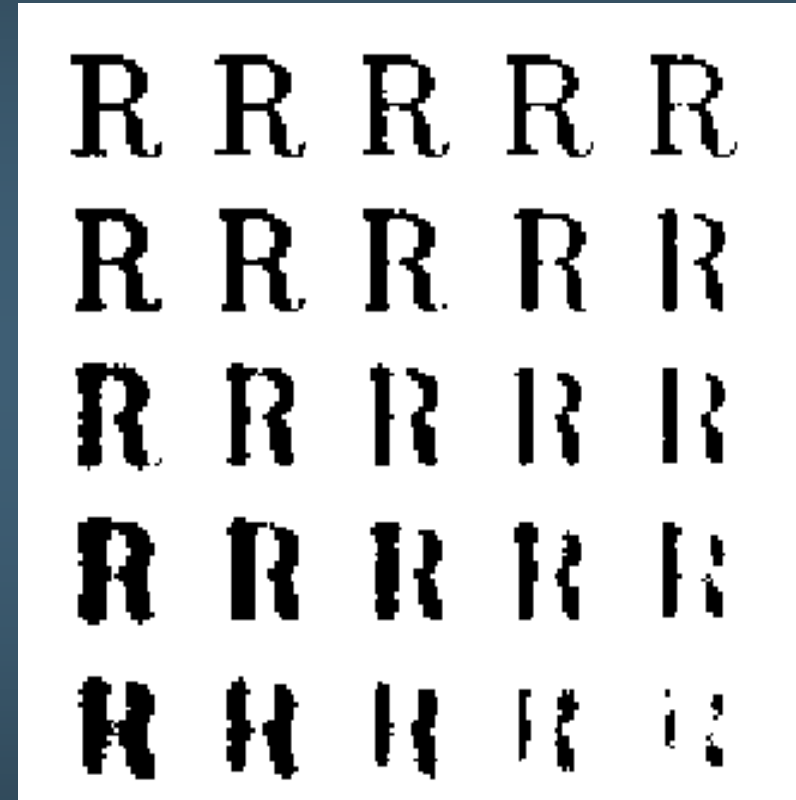
# Image Quality Modeling

thrs x blur

Effects of printing & imaging:

blur

thrs
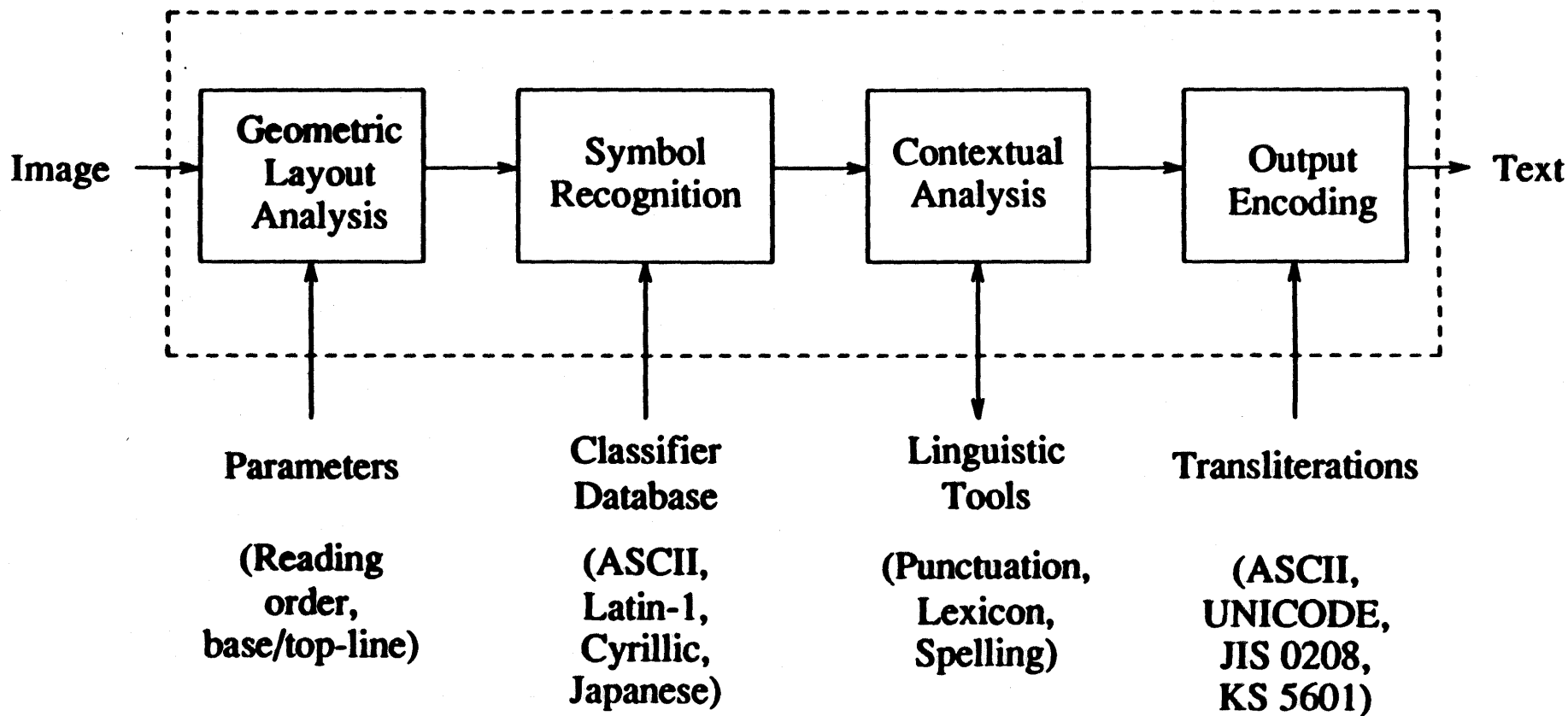
sens

Also, 8 other parameters

Baird & Pavlidis (1985-1992)

# Image Quality Models:
# Fitting to Real Data & Using Safely

- Testing dissimilarity of two sets of images

    a sensitive bootstrap statistic:  indirectly infer parameters

    (Kanungo Ph.D., 1996 ff)

- Estimating parameters directly from sample images

    a few character images are sufficient

    (Barney Smith Ph.D., 1998 ff)

- Ensuring the safety of training on synthetic data

    by interpolation in generator parameter space

    (Nonnemaker Ph.D., 2008)

Many open questions remain (several Ph.D.s' worth?)

# Model-Driven Architecture



Several application-specific models of knowledge:

most can be acquired (trained, hand-crafted, bought) off-line.

Baird & Ittner (1988-1994)

# Accuracy was High, but Not Uniform

Best:

Garamond Roman:  Pack my box with five dozen liquor

*Texttype Italic:  Pack my box with five dozen liquor ju*

*Plantin Light Italic:  Pack my box with five dozen liquor*

Aster:  Pack my box with five dozen liquor jugs.

> 99.97%

Average:

Typewriter Gothic:   Pack my box with five

*Bell Italic:  Pack my box with five dozen liquor jugs.*

Serifa:  Pack my box with five dozen liquor jugs.

Caslon Old Face:  Pack my box with five dozen liquor jugs.

~ 99.7%

Worst:

Benguiat Book:  Pack my box with five dozen liquor jugs

*Gill Sans Italic:  Pack my box with five dozen liquor jugs.*

Weiss:  Pack my box with five dozen liquor jugs.

*Avant Garde Italic:  Pack my box with five dozen li*

< 99%

# Single-font Classifiers are far more accurate on their font

Trained one 100-font classifier:

    tested it on all 100 fonts

    4.2% error rate

Trained 100 single-font classifiers:

    tested each on its own font

    0.81% error rate

Single-font classifiers are *much* better:

    $\times 5.2$ reduction in error (multiplicative factor)

**Generic:**

versatile, but…

not a best fit to many fonts

**Specific:**

brittle, but…

the best fit to some font

If can recognize the input font, can benefit a lot

(but, hard to do)

# "Strong" versus "Weak" Models

*I sometimes find it helpful to distinguish between them.*

- **Strong** models:
  - application-specific,
  - a close fit to the input,
  - often detailed and formal.

- **Weak** models:
  - generic,
  - applicable to other inputs too,
  - often informal or imprecise.

*Expensive to acquire*
*More accurate*
  *(on the right input)*

We often feel forced to choose one over the other

*Cheap to acquire*
*Less accurate*
  *(on average)*

# Strong Models:  e.g.  *Sahovsky Informator*



Chess encyclopaedia in 20+ volumes
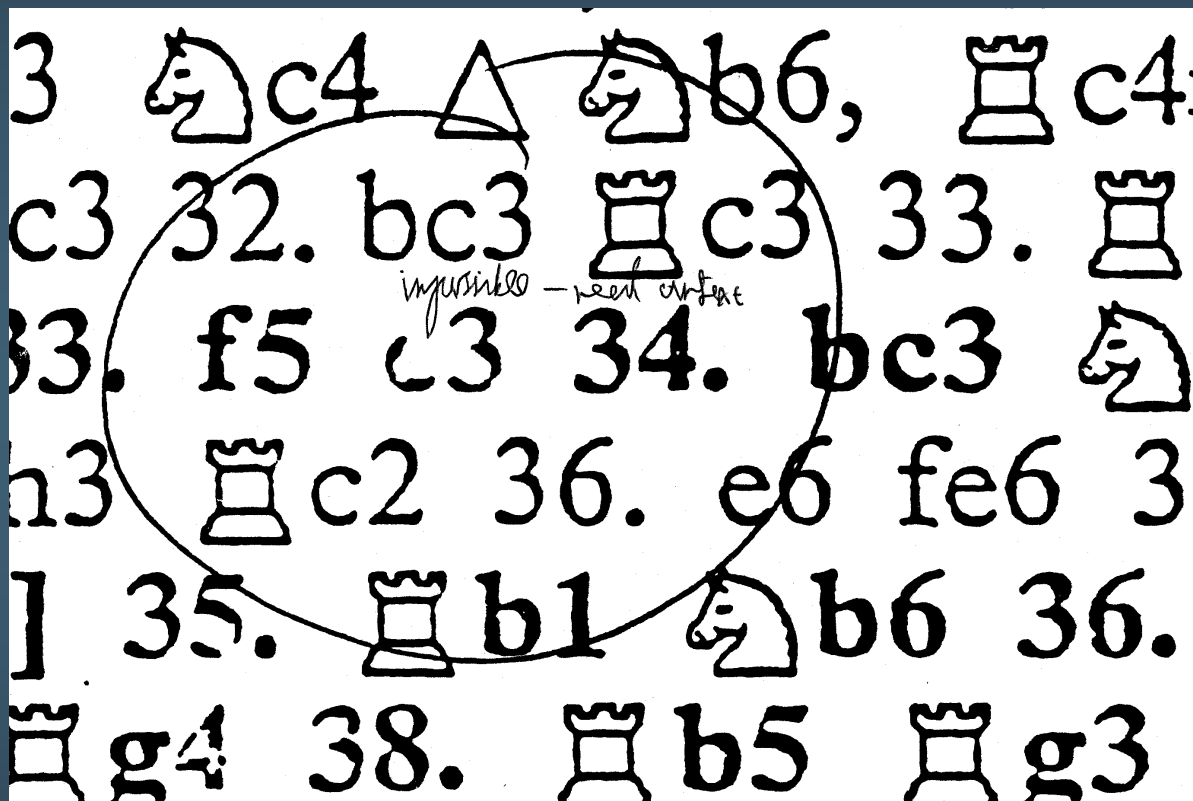
Games of theoretical interest

Ken Thompson wanted to teach these games to *Belle*, his chess machine

Ken coded-up syntax & semantic models

Baird & Thompson (1990)

# Challenging Print Quality



I trained on this special chess font

Near-perfect OCR is impossible on such poor quality

But, unless *entire games* are correctly read, then it's not worth doing…!

# *Informator* Syntax is Computable

**114.**      (R 76/b)    **A 64**

## HULAK — NUNN
### Toluca (izt) 1982

**1. d4** ♘**f6 2. c4 c5 3. d5 e6 4.** ♘**c3 ed5
5. cd5 d6 6.** ♘**f3 g6 7. g3** ♗**g7 8.** ♗**g2
0—0 9. 0—0 a6 10. a4** ♘**bd7 11.** ♘**d2**
♖**e8 12. h3** ♖**b8 13.** ♘**c4** ♘**e5 14.** ♘**a3**
♘**h5 15. e4** ♗**d7?! [15... f5?!; 15...** ♖**f8]
16. a5! [16. g4 b5! 17. ab5 ab5 18.** ♘**ab5**
♗**b5 19.** ♘**b5** ♘**g3! 20.** ♘**d6** ♕**d6 21. fg3
c4** △ ♘**d3∞]** ♕**a5 17. g4** ♘**f6 18. g5! [18.
f4** ♘**eg4 19. hg4** ♘**g4 20.** ♘**c2** ♕**d8∞]**
♘**h5 19. f4** ♘**c4 20.** ♘**c4** ♕**a1   21.** ♘**d6**
♗**d4 22.** ♔**h2** ♖**e7 23.** ♕**f3± [23. e5?!**
♖**e5 24.** ♘**f7** ♖**e7∞] b5 24. e5 b4 25.**

**Game** has the form:

   HEADER MOVE MOVE …

   Find header using layout geometry

   Ignore commentary:   !?   ±   △   [MOVE … (MOVE …) …]

**Move** has the form:

   N. PLY PLY      ( 3. ♖f3 c5 )

   Numbers N must ascend: 1, 2, …

   White ply, then Black ply ("half-moves")

**Ply** has the form:

   PIECE LETTER DIGIT    ( ♖f3 )

   also: LLD LD PLLD PDLD CASTLE

| | |
|---|---|
| PIECE | ♙ ♘ ♗ ♖ ♕ ♔ |
| LETTER | a b c d e f g h |
| DIGIT | 1 2 3 4 5 6 7 8 |
| CASTLE | O—O   O—O—O |

# Chess Semantics is Computable

Apply the rules of chess

Check: **Is the $i$th move legal?**

   prior context:   $1$  ...  $i-1$

Check: **Is the $i$th move suspect?**

   later context:   $i+1$  ...  *end*

Generate: **Which $i$th moves are legal?**

   prior context:   $1$  ...  $i-1$

   list *all* alternatives:  typically 20-50

# Fully Automatic Extraction of Games



Image of page

'Galley-proof' format output from the OCR

Database of games, moves

this game = 83 half-moves

# Semantic Model Astonishingly Helpful

Characters:

| 99.5% | OCR Alone |
| 99.8% | Syntax |
| 99.995% | Semantics |

Games:

| 42% | OCR Alone |
| 76% | Syntax |
| 97% | Semantics |

On Over 2 Million Characters

Syntactic model cuts errors in half

Semantics cuts errors by another <u>factor of 40!</u>

99.5% OCR accuracy implies that
    game accuracy is only 40%

After semantic analysis,
    almost all games are completely correct

# Lessons from Reading Chess

An extreme illustration of strong modeling:
- Syntax & semantics fitted precisely to these books
- Remarkably high performance:   50 errors per million chars

But:  wasn't this a unique event?
- Can we model syntax and semantics of other books?
- Will our users be domain experts w/ software skills?

Note the <u>size of the context</u> is many dozens of moves, all operated on by the semantic analysis.
- Perhaps we can operate on <u>long passages</u> in other ways....
- Would that help…?   (Open question, for years.)

# Beyond Versatility:   George Nagy's
# Adapting Recognizers

Can a recognition system adapt to its input?

Can weak models "self-correct," and so
strengthen themselves fully automatically?

When a 100-font system reads a document in a
single font, can it specialize to it *without:*

- knowing which font it is,
- recognizing the font, or
- using a library of pre-trained single-font classifiers ?

Nagy, Shelton, & Baird (1966 & 1994)

# Toy Example:   a Single-Font Test

{ *0, O, Q, D, G, C* } in *Avant Garde Book Oblique*

10 point size, 300 pixels/inch resolution, 200 sample images each

# The weak (100-font) classifier performs poorly on this….

Confusion matrix of the given polyfont classifier:

|  |  | **t o p - c h o i c e** | | | | | | |
|---|---|---|---|---|---|---|---|---|
|  |  | 0 | O | Q | D | G | C | |
|  | 0 | 96 | 104 | 0 | 0 | 0 | 0 | 52.0 |
| **t** | O | 0 | 200 | 0 | 0 | 0 | 0 | 0.0 |
| **r** | Q | 0 | 1 | 199 | 0 | 0 | 0 | 0.5 |
| **u** | D | 0 | 50 | 4 | 146 | 0 | 0 | 27.0 |
| **e** | G | 0 | 0 | 0 | 0 | 197 | 3 | 1.5 |
|  | C | 0 | 0 | 0 | 0 | 2 | 198 | 1.0 |
|  |  | 0.0 | 43.7 | 2.0 | 0.0 | 1.0 | 1.5 | 13.67 |

Far from perfect:  14% error rate
Especially:   0/O and D/O confusions

# Now, pretending that we believe this classifier, we boldly *retrain*….

Train new classifier on top-choice-labeled images

Confusion matrix of the retrained classifier:

|       |   | 0    | O    | Q   | D   | G   | C   |       |
|-------|---|------|------|-----|-----|-----|-----|-------|
|       |   |      |      | **t o p - c h o i c e** |     |     |     |       |
|       |   | 0    | O    | Q   | D   | G   | C   |       |
|       | 0 | 178  | 22   | 0   | 0   | 0   | 0   | 11.0  |
| **t** | O | 0    | 200  | 0   | 0   | 0   | 0   | 0.0   |
| **r** | Q | 0    | 0    | 200 | 0   | 0   | 0   | 0.0   |
| **u** | D | 0    | 26   | 0   | 174 | 0   | 0   | 13.0  |
| **e** | G | 0    | 0    | 0   | 0   | 200 | 0   | 0.0   |
|       | C | 0    | 0    | 0   | 0   | 0   | 200 | 0.0   |
|       |   | 0.0  | 19.4 | 0.0 | 0.0 | 0.0 | 0.0 | 4.00  |

Error rate drops by a factor of $\times 3.4$ !

The risk of training on (some) mislabeled test data didn't hurt us!
Lucky!! … or is it reliable?

# How lucky can we get…?

Retrain again, using new top-choice labels:

|   |   | **t o p - c h o i c e** | | | | | | |
|---|---|---|---|---|---|---|---|---|
|   |   | O | O | Q | D | G | C |   |
|   | O | 196 | 3 | 0 | 1 | 0 | 0 | 2.0 |
| **t** | O | 0 | 200 | 0 | 0 | 0 | 0 | 0.0 |
| **r** | Q | 0 | 0 | 200 | 0 | 0 | 0 | 0.0 |
| **u** | D | 2 | 16 | 0 | 182 | 0 | 0 | 9.0 |
| **e** | G | 0 | 0 | 0 | 0 | 200 | 0 | 0.0 |
|   | C | 0 | 0 | 0 | 0 | 0 | 200 | 0.0 |
|   |   | 1.0 | 8.7 | 0.0 | 0.5 | 0.0 | 0.0 | 1.83 |

Error rate drops by *another* factor of $\times 3.4$
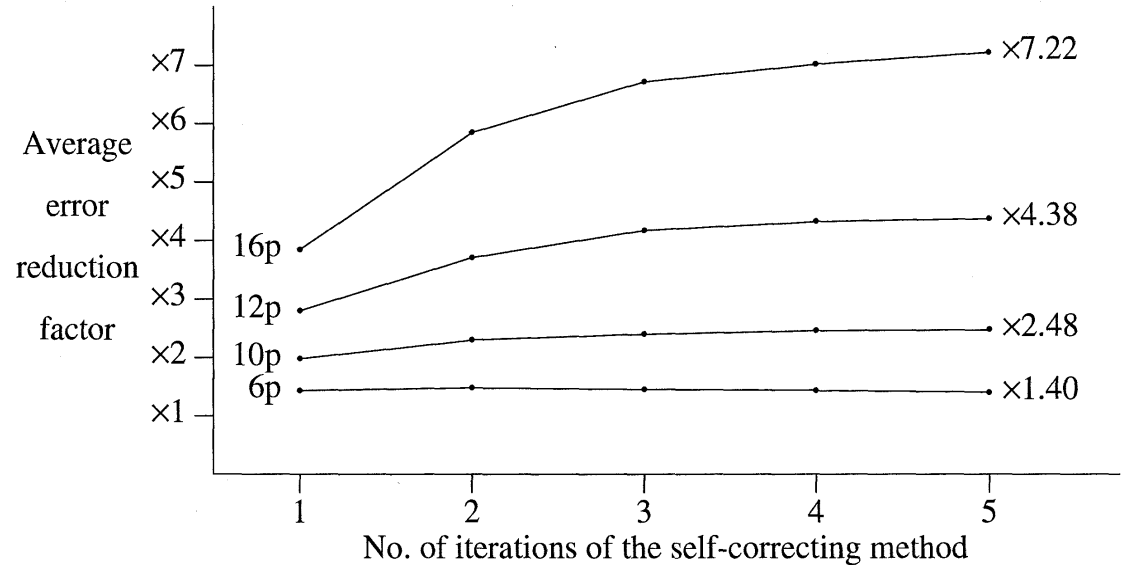
After five iterations:

1.33% error rate

$\times 10.3$ reduction in error rate, overall

# In fact this works reliably *(…but why??)*

Aster Roman
*Aster Italic*
Avant Garde Book Roman (ITC)
*Avant Garde Book Oblique (ITC)*
Bembo Roman
*Bembo Italic*
**Bodoni Roman**
*Bodoni Italic*
Bookman Light Roman [ITC]
*Bookman Light Italic [ITC]*
Breughel Roman
*Breughel Italic*
Caledonia Roman
*Caledonia Italic*
Caslon Old Face #2 Roman
*Caslon Old Face #2 Italic*
Cheltenham Roman
*Cheltenham Italic*
Clearface Regular Roman [ITC]
*Clearface Regular Italic [ITC]*
Cloister Roman
*Cloister Italic*
Corona Roman [Adobe]

Average error-reduction factors over 100 fonts



Some improvement at all four sizes

Three iterations are enough

# Image Quality also is often Constant throughout a Document

Rather like typefaces, a "style" determined by image degradations due to printing, scanning, etc.

# A Theory of Adaptation:    Prateek Sarkar's
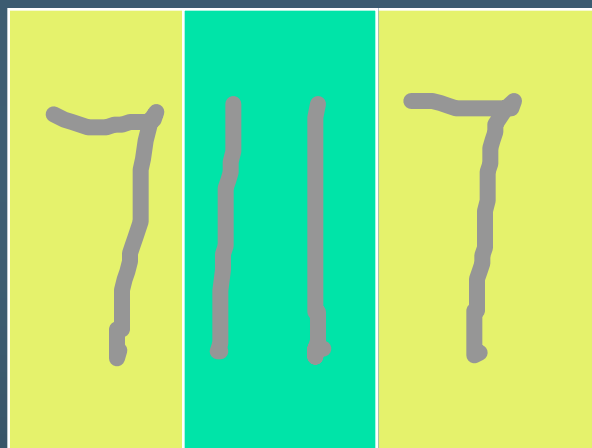# Style-Conscious Recognition

- Many documents possess a consistent *style:*
  - *e.g.* printed in one (or only a few) typefaces
  - or, handwritten by one person
  - or, noisy in a particular way
  - or, using a fixed page layout
  - ….(many examples)

- Broadly applicable idea:   a <u>style</u> is a manner of rendering (or, generating) patterns.

- <u>Isogenous</u>— *i.e.* 'generated from the same source' —documents possess a uniform style

Sarkar, Nagy, Veeramachaneni (2000-2005)

# Style-Consistent Recognition

Sevens, or ones…?  Ambiguous!

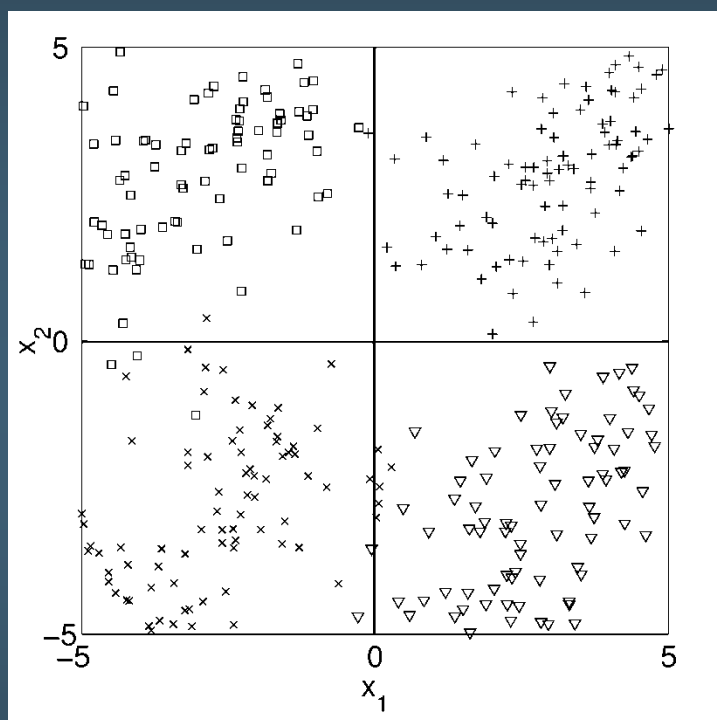**writer 1**                          **writer 2**



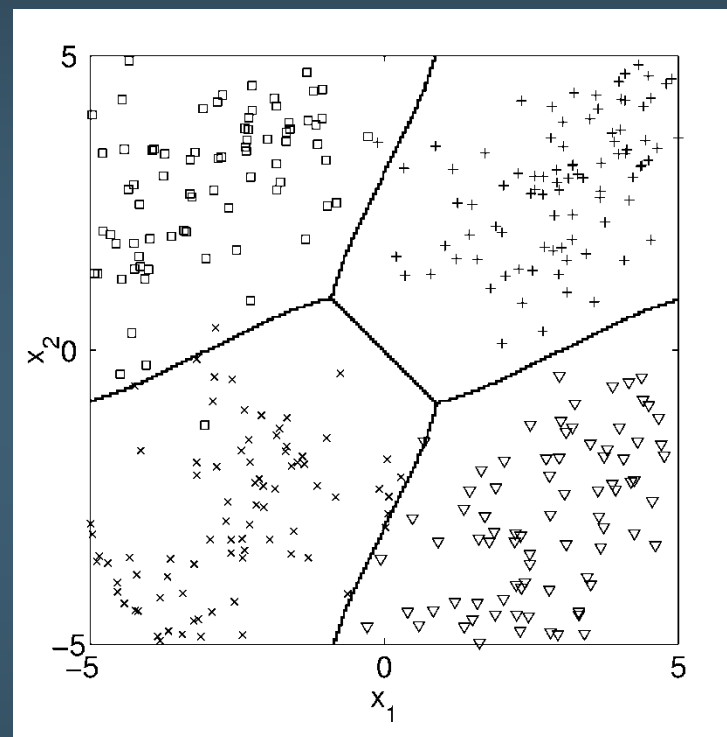Ambiguity is resolved by style-consistency.

# Style-Conscious Methodology

- Modeling style-consistency improves classification on isogenous input

- Improvement is higher on longer input passages

- Styles and style parameters can be estimated without style labels

- Style models complement, and do not impede, other recognition models (*e.g.* linguistic)

- *Lesson: <u>weak models can become stronger</u> when operating on long isogenous passages*

# Refined Classifier Decision Regions
## (for a passage with two symbols)



Style-unconscious: suboptimal      Style-conscious: optimal

# PARC's Document Image Decoding

Text image

.Kn,eHjll 4MLXOapykw,QjJ0YO
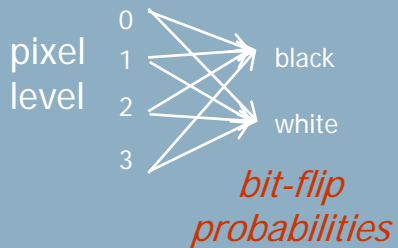
## Multi-level image degradation model

Character templates

pixel level

0
1
2
3

black

white

bit-flip probabilities

Side-bearing model

Recognition

Text content

.Kn,eHjl1 4MLXOapykw,QjJ0YO

Kopec, Chou, Kam, & Lomelin (1994-1997)

# DID can learn Strong Models
# of even extreme image degradations

Error rate <1%

.Kn.cHjll 4MLXOapyl·w .QjJ0Y()

l.q.e97 4VyuZ mj

7 l·vqa9PDT NA qp

5 uqeNl·R 9 Zrrql·

KiZHPWU3aQM 49F1

- Works over a wide range of image qualities

- A system can adapt to any of a large set of pre-trained qualities

Sarkar, Baird, & Zhang (2003)

# DID is Model-Intensive

- Explicit formal stochastic models of
  - **text generation**:   language, format
  - **image rendering**:   typefaces, layout
  - **image quality**:   asymmetric bit-flip
  ( combined in a single finite-state Markov network )

- Search algorithms find best 'decoding'

  - provably optimal (under MAP criterion)

  - Viterbi and Iterated Complete Path:   often fast

- Joint over many models, some weak:
  - linguistic:  char *N*-gram & imperfect lexica
  - quality:  simplistic bit-flip model

Kopec, Chou, Minka, Popat, Bloomberg (1994-2001)

# Weak Language Models Can Help Overcome Severe Image Noise

Degraded, subsampled, greyscale image



DID recognition without a language model

WHITR.KITTIVI HAO BEEN HAVING IT.,.RACE,WASHEI4.BX THB.UI,D CAT FOR

DID w/ n-gram char model, Iterated Complete Path search algorithm

WHITE KITTEN HAD BEEN HAVING ITS FACE WASHED BY THE OLD CAT FOR

Kopec, Popat, Bloomberg, Greene (2000-2002)

# Lessons from DID

- Combining several models, even if some are weak, can yield high accuracy

- Joint recognition over many models---iconic, linguistic, quality, layout---can be performed provably optimally, and fast

- Recognizing entire text-lines at a time helps


- *Weak models can provide the basis for high performance recognition systems*

# Extremely Long Passages: "Whole-Book" Recognition

Operate on the <u>complete set</u> of a book's page images, using automatic unsupervised adaptation to improve accuracy.

*Given:* *(1) images of an entire book,*

*(2) an initial transcription (generally erroneous), &*

*(3) a dictionary (generally imperfect),*

*Try to:* *improve recognition accuracy fully automatically,*

*guided only by evidence within the images.*

Xiu & Baird (2008-2011)

# Start with Two Weak Models
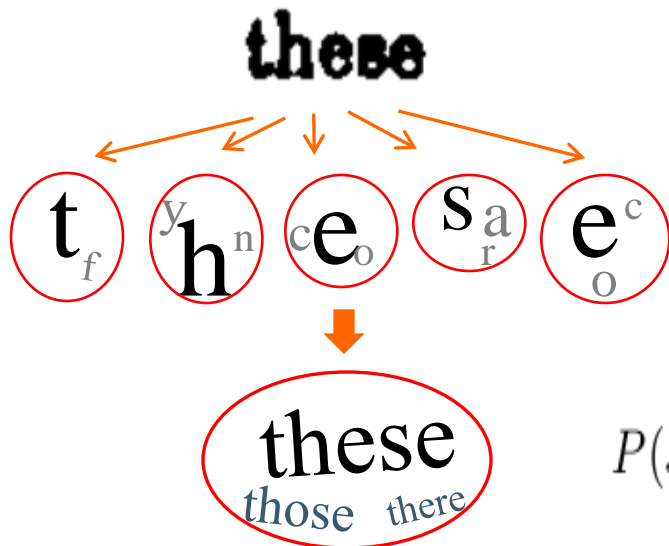
- **Iconic model:**
  - *Describes image formation and determines the behaviour of a character-image classifier*
  - *For example, the prototypes in a template-matching character classifier.*
  - *Weak: inferred from buggy OCR transcription*

- **Linguistic model:**
  - *Describes word-occurrence probabilities*
  - *For example, a dictionary*
  - *Weak: not a perfect lexicon: too small (or too large)*

*Word recognition, driven by (1) iconic model alone, and (2) both iconic and linguistic models (jointly), may get <u>different results</u>, indicating "disagreements" between the models.*

# Disagreements can be Detected Statistically



**Char recognition** (apply iconic model alone):

$$P(s_1|x_1) \cdots P(s_T|x_T)$$

**Word recognition** (iconic & linguistic jointly):

$$P(s_1 \cdots s_T | x_1 \cdots x_T) = \frac{1}{\alpha} P(s_1|x_1) \cdots P(s_T|x_T) C(s_1 \cdots s_T)$$

***Character** disagreement:*

$$\epsilon(i|X) \equiv - \sum_{s \in \Sigma} P(s_i = s | X) \cdot \log P(s|x_i)$$

(**cross entropy** on a char)

***Word** disagreement:*

$$\epsilon(X) = \sum_{i=1}^{T} \epsilon(i|X)$$

(…w/in a word)

***Passage** disagreement:*

$$\epsilon(\mathcal{P}) = \sum_{X \in \mathcal{P}} \epsilon(X)$$

(…w/in the whole book)

# Disagreement-Driven Model Adaptation Algorithm

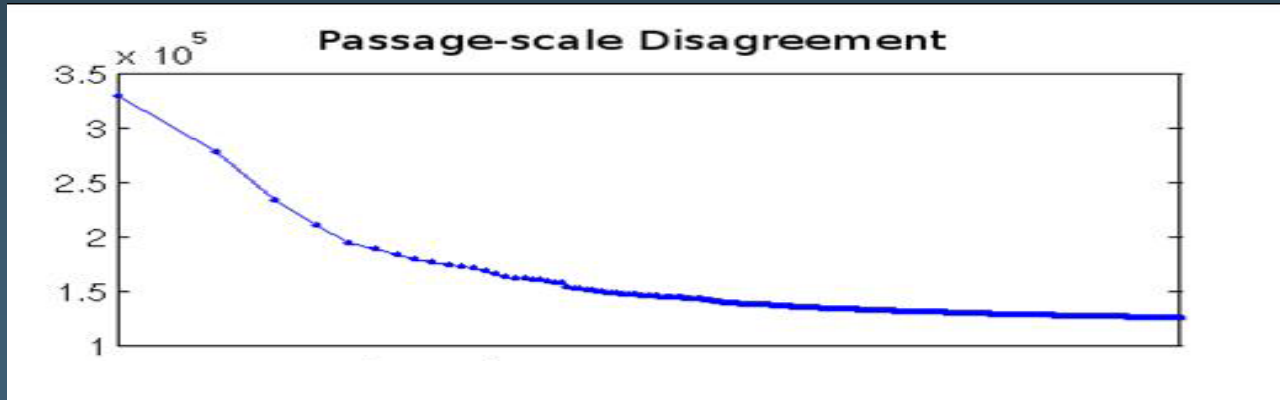Iterate many times….

- Compute all character, word & passage disagreements
- Identify words and characters where the two models most disagree.
- Propose adaptations to the models to reconcile them.
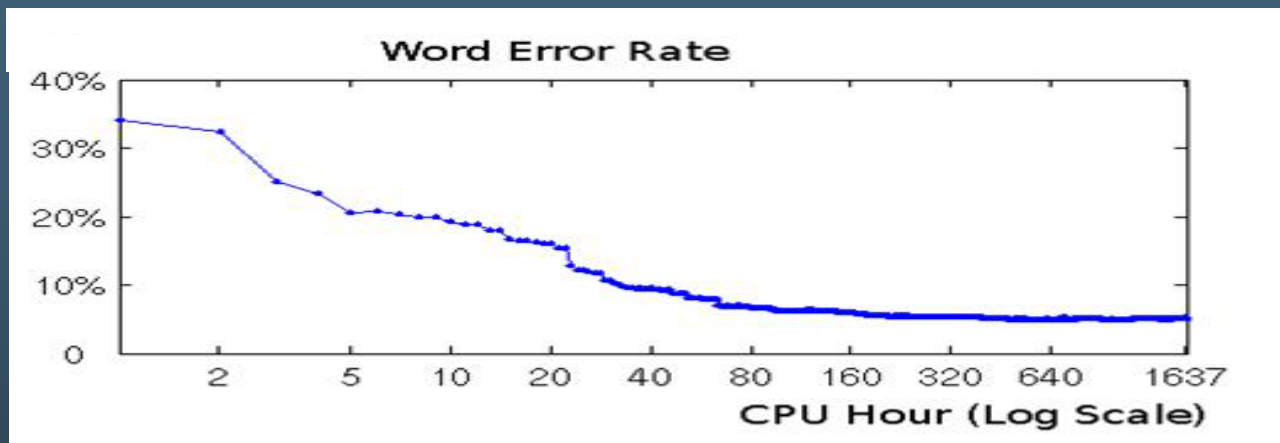- Check that each proposed adaptation reduces passage disagreement:  if so, accept the adaptation.

**The two models are "criticizing" one another, & correcting one another—although both are imperfect!**

# Disagreements Identify Errors

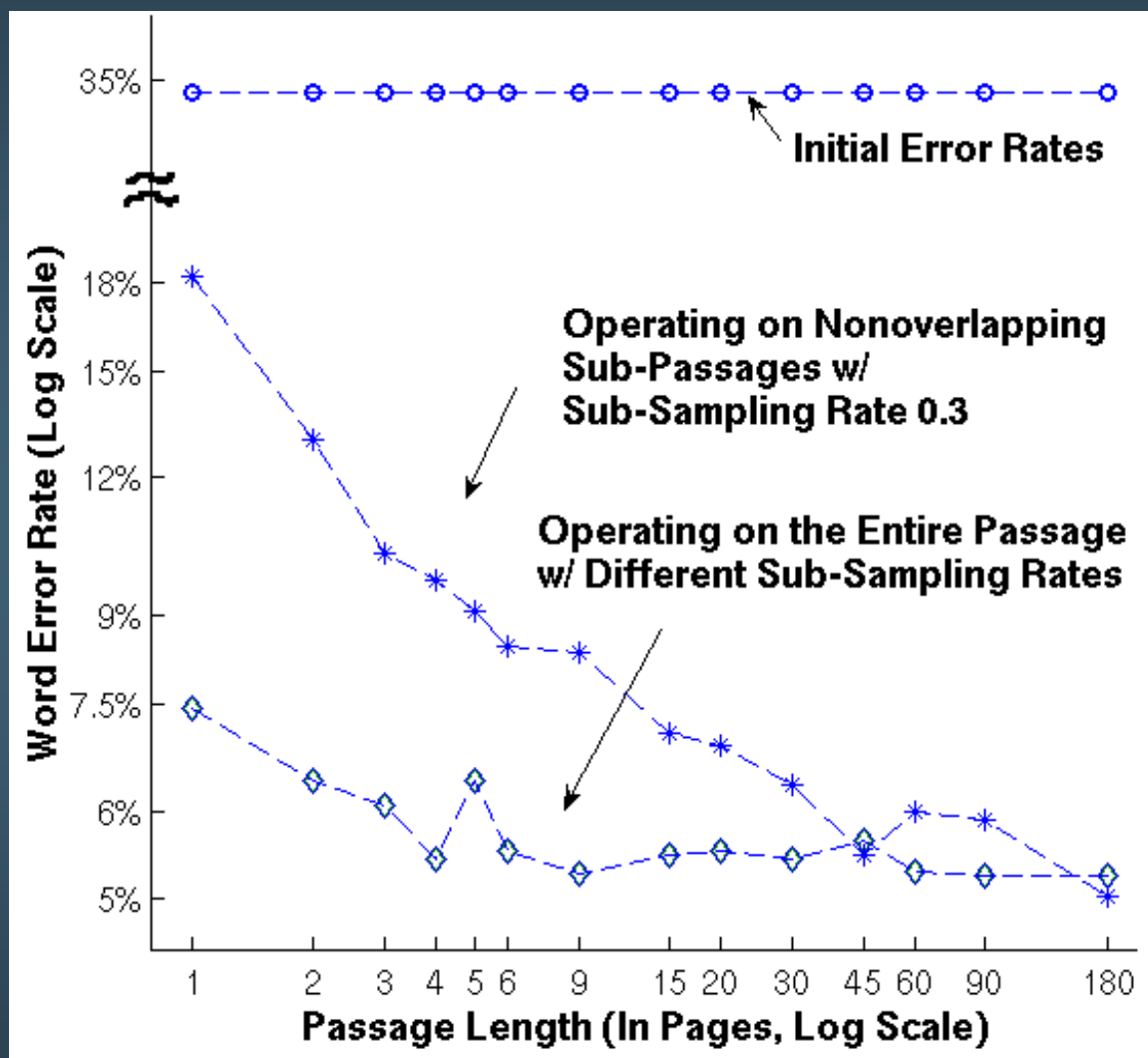The algorithm drives disagreements down.....



...and, disagreements are correlated with errors....



...so, the algorithm drives down errors.

# Longer Passages Improve More

**Benefits of isogeny:**

*Longer passages are driven to lower error rates ultimately.*

# Pingping Xiu's
# Whole-Book Recognition

- The larger the input passage is, the better the algorithm performs: the lower the final error rate.

- The algorithm can be sped up by two orders of magnitude using randomization and caching.

- Rigorous sufficient conditions for the algorithm to succeed have been proven.

- Two weak models, although both are imperfect, can criticize and correct one another, both becoming stronger.

# Enables 'Anytime' Recognition

- Recognizers which run 'forever'
  - safe, since accuracy improves nearly monotonically
  - trade runtime for (eventual) accuracy

- Can be interrupted at any time to see the best interpretation found so far
  - system is always operating on the *entire* document

- A good fit to 'personal recognition' needs
  - users are unskilled:   can't engineer; won't correct
  - no tight deadline:   soak up idle cycles

# Twenty-five Years of DAR Research: Model-Intensive Recognition

- **Specify the domain precisely:**

  *define quantitative generative models*

  *of document images to be recognized*

- **Learn models from examples:**

  *synthetic training data can be safe;*

  *affordable weak models may be good enough*

- **Strive for provable performance guarantees:**

  *invent joint recognition algorithms which are*

  *formally optimal w.r.t. to the models*

- **Adapt weak models, strengthen automatically:**

  *on short passages, apply style-conscious adaptation;*

  *on long passages, mutual criticism and correction*

# Advantages of Model-Intensive Recognition

- When the models are strong (closely fit the input),
  <u>results are the best possible</u>

- When models can be trained nearly automatically,
  <u>effort required for best results is minimized</u>

- When training is known to work across a wide range,
  <u>confidence in high performance is high</u>

- If the system isn't yet good enough:
  improve the models:  adaptively perhaps
  ---but <u>not</u> the recognition algorithms!

# Focusing on a peculiar distinction: 'Strong' *versus* 'Weak' Models

Shifts our attention away from end-results:
accuracy, speed, and costs of engineering
—and towards this question:

*How well do our models fit the
particular input which our system
is trying to recognize?*

The answer to this can determine accuracy, engineering costs, even speed….

*By working this way, we may enjoy the best of both:
affordable engineering costs, plus high accuracy!*

# Thanks!

*And thanks to all those who inspired me, especially:*

*Theo Pavlidis, Tin Kam Ho, David Ittner, Ken Thompson, George Nagy, Robert Haralick, Tao Hong, Tapas Kanungo, Phil Chou, Dan Lopresti, Gary Kopec, Dan Bloomberg, Ashok Popat, Tom Breuel, Elisa Barney Smith, Prateek Sarkar, Harsha Veeramachaneni, Jean Nonnemaker, and Pingping Xiu.*