# Handwriting Recognition Technology in the Newton's Second Generation "Print Recognizer" (The One That Worked)

Larry Yaeger
Professor of Informatics, Indiana University
Distinguished Scientist, Apple Computer

World Wide Newton Conference
September 4-5, 2004

WWNC 2004

# Handwriting Recognition Team

## Core Team

Larry Yaeger   (ATG)        Brandyn Webb   (Contractor)
Dick Lyon      (ATG)        Les Vogel        (Contractor)
Bill Stafford  (ATG)

## Other Contributors

Rus Maxham      Kara Hayes        Gene Ciccarelli    Stuart Crawford
Chris Hamlin    George Mills      Dan Azuma          Boris Aleksandrovsky
Josh Gold       Michael Kaplan    Ernie Beernink     Giulia Pagallo

## Testers

Polina Fukshansky    Glen Raphael      Julie Wilson    Emmanuel Euren
Ron Dotson           Denny Mahdik

# Recognizer History

- ʋ '92 ATG "Rosetta" project demos well at Stewart Alsop's "Demo '92" (blows socks off Nathan Myhrvold's MS demo) and WWDC
- ʋ '93 Head of ATG suggests abandoning handwriting recognition for interactive TV project
- ʋ '93-'94 Rosetta nearly ships in "PenLite" pen-based Mac product
- ʋ Jan '94 Port to Newton started
- ʋ '94 Brief interest in Rosetta for abortive "Nautilus" Mac product
- ʋ … testing with tethered Newtons, *much* accuracy improvement…
- ʋ 18 Nov '94 Provided handful of untethered Newtons for testing
- ʋ 1 Feb '95 Beta 1 build (Merry Xmas!)
- ʋ '95 Rosetta ships as "Print Recognizer" in Newton (120?)
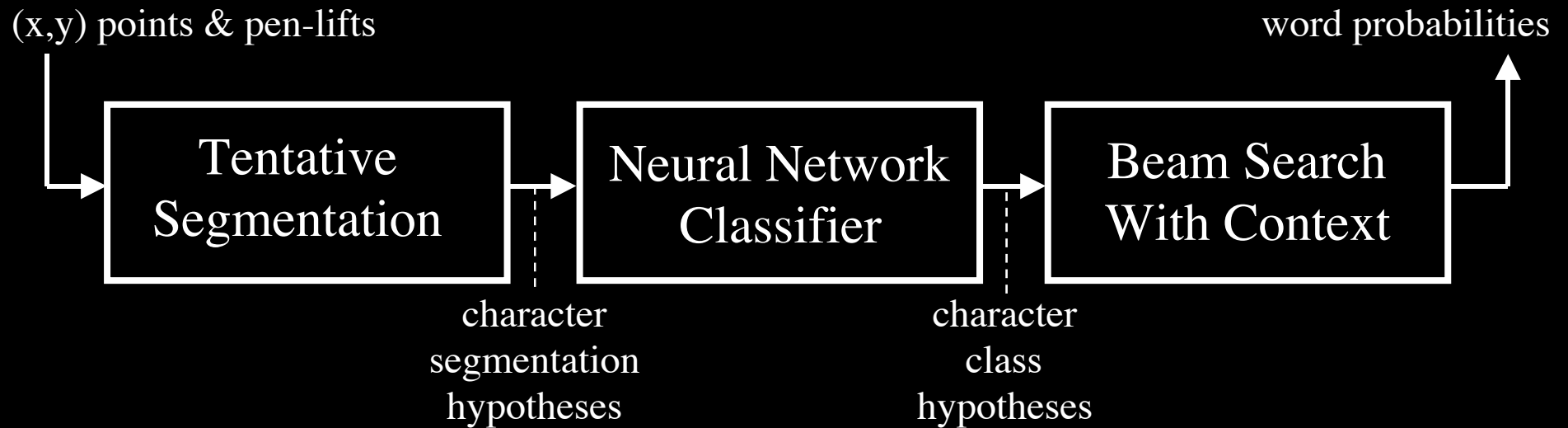- ʋ '95 Rosetta widely acknowledged as world's first usable handwriting recognizer

# Recognizer History

υ 13 Nov '95 John Markoff writes about Rosetta in NY Times

υ Nov or Dec '95 receive cease-and-desist demand for use of "Rosetta" name (Mac-based SmallTalk platform)

υ Jan '96 team picks "Mondello" codename, "Neuropen" product name

υ '96 Short-lived "Hollywood" pen-based Mac project

υ Mar '97 cursive almost working

υ 18 Mar '97 ATG laid off

υ May '00 "Inkwell" for Mac OS 9 declares beta

υ May '00 Marketing declares "no new features on 9", OS X work begins

υ Jul '02 Inkwell for Mac OS X declares GM (10.2 / Jaguar)

υ Sep '03 Inkwell APIs and additional languages declare GM (10.3 / Panther)

υ Apr '04 Motion announced with gestural interface, including tablet and in-air ink-on-demand
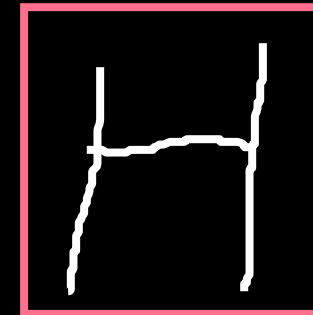
# Recognizer Overview

υ Powerful state-of-the-art technology

    υ Neural network character classifier

    υ Maximum-likelihood search over letter segmentation, letter class, word, and word segmentation hypotheses

    υ Flexible, loosely applied language model with very broad coverage

υ Now part of "Inkwell" in Mac OS X

υ Also provides gesture recognition

    υ System

    υ Application (Motion)

# Recognition Block Diagram

(x,y) points & pen-lifts                                                word probabilities

| Tentative Segmentation | Neural Network Classifier | Beam Search With Context |
|---|---|---|

character
segmentation
hypotheses

character
class
hypotheses

# Character Segmentation

- υ Which strokes comprise which characters?
- υ Constraints
  - υ All strokes must be used
  - υ No strokes may be used twice
- υ Efficient pre-segmentation
  - υ Avoid trying all possible permutations
  - υ Based on order, overlap, crossings, aspect ratio...
- υ Integrated with recognition
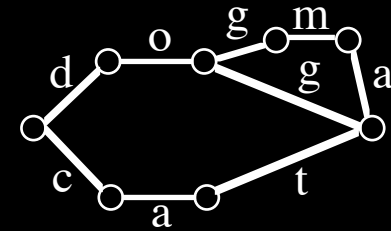  - υ Forward & reverse "delays" implement implicit graph of hypotheses

# Neural Network Character Classifier

υ Inherently data-driven

υ Learn from examples

υ Non-linear decision boundaries
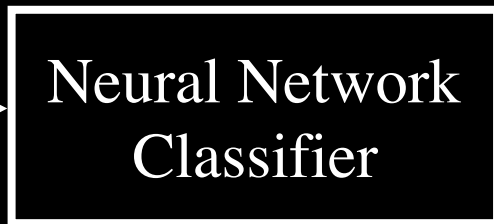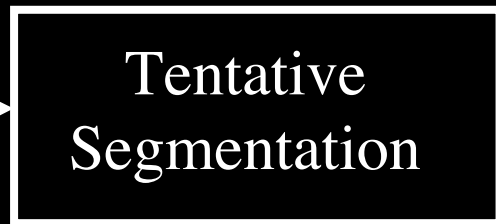
υ Effective generalization

# Context Is Essential

ʋ Humans achieve 90% accuracy on characters in isolation (our database)

   ʋ Word accuracy would then be ~ 60%  $(.9^5)$

ʋ Variety of context models are possible

   ʋ N-grams

   ʋ Variable (Memory) Length Markov Model

   ʋ Word lists

   ʋ Regular expression graphs

ʋ "Out of dictionary" writing also required

   ʋ "xyzzy", unix pathnames, technical/medical terms, etc.
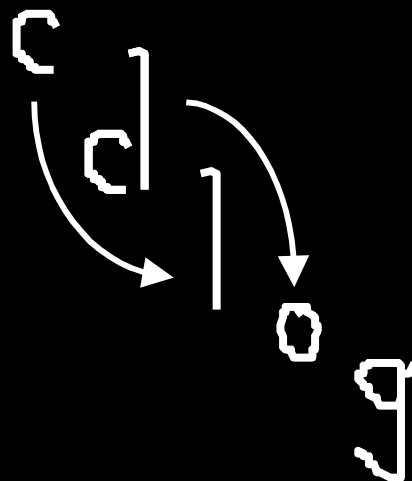
# Recognition Technology

(x,y) points & pen-lifts

word probabilities

| Tentative Segmentation | Neural Network Classifier | Beam Search With Context |
|---|---|---|

character segmentation hypotheses

character class hypotheses

| | | | | | |
|---|---|---|---|---|---|
| a | .1 | .0 | .0 | .0 | .0 |
| b | .0 | .1 | .0 | .0 | .0 |
| c | .7 | .0 | .0 | .1 | .0 |
| d | .0 | .7 | .0 | .0 | .0 |
| e | .1 | .0 | .0 | .1 | .0 |
| f | .0 | . | .0 | .0 | .0 |
| g | .0 | .0 | . | .0 | .7 |
| ... | ... | ... | ... | ... | ... |
| l | .0 | .1 | 1. | .0 | .0 |
| ... | ... | ... | ... | ... | ... |
| o | .1 | .0 | .0 | .8 | .0 |
| ... | ... | ... | ... | ... | ... |

# Character Segmentation

| Ink | Segment Number | Segment | Stroke Count | Forward Delay | Reverse Delay |
|-----|----------------|---------|--------------|---------------|---------------|
| | 1 | c | 1 | 3 | 0 |
| | 2 | d | 2 | 4 | 1 |
| clog | 3 | do | 3 | 4 | 2 |
| | 4 | l | 1 | 2 | 0 |
| | 5 | lo | 2 | 2 | 1 |
| | 6 | o | 1 | 1 | 0 |
| | 7 | g | 1 | 0 | 0 |

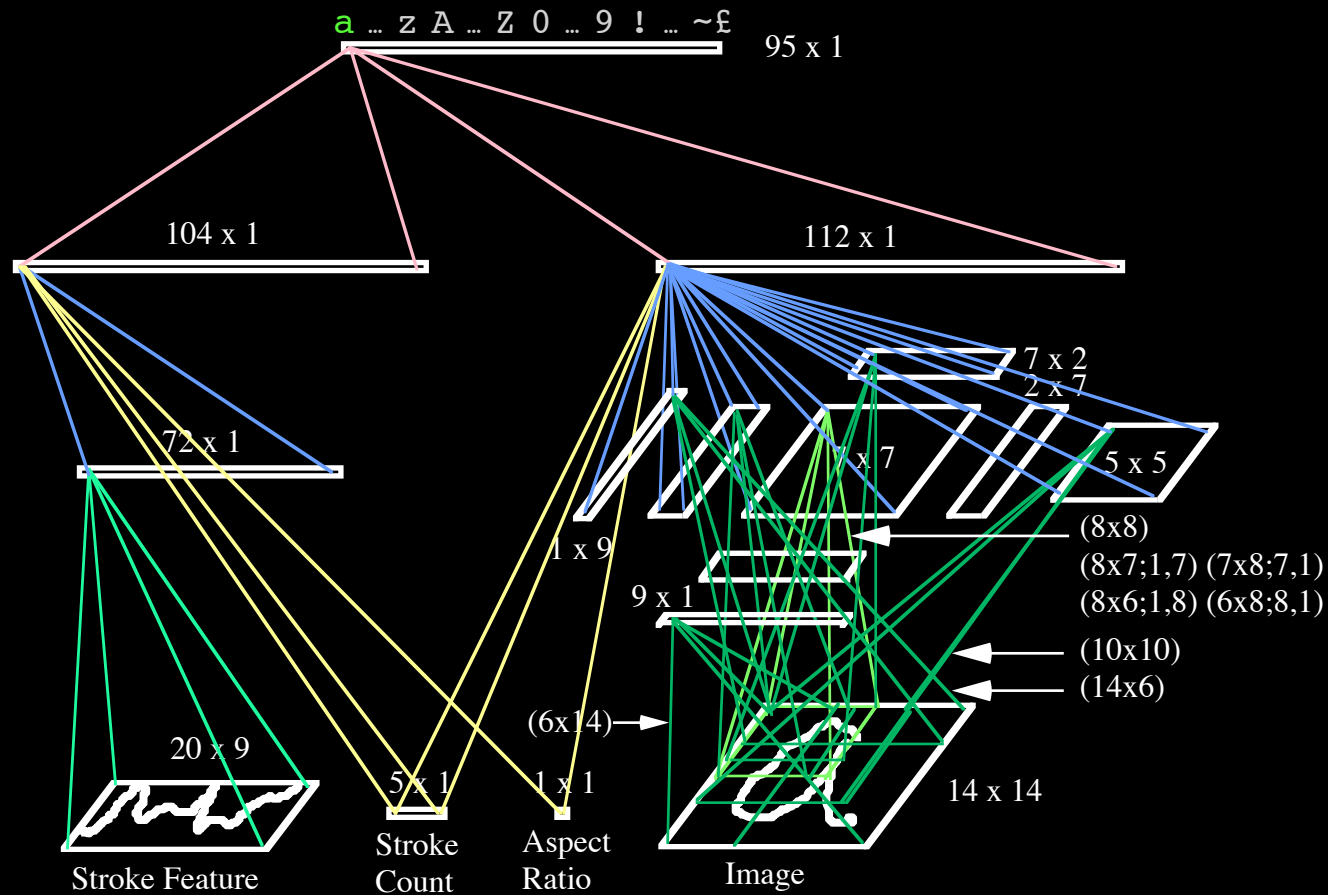$i \rightarrow j$ is legal *iff* $FD_i + RD_j = j - i$

# Network Design

- ʋ Variety of architectures tried
  - ʋ Single hidden layer, fully-connected
  - ʋ Multi-hidden layer, with receptive fields
  - ʋ Shared weights (LeCun)
  - ʋ Parallel classifiers combined at output layer
- ʋ Representation as important as architecture
  - ʋ Anti-aliased images
  - ʋ Baseline-driven with ascenders and descenders
  - ʋ Stroke features

# Network Architectures

a … z A … Z 0 … 9 ! … ~

a … z A … Z 0 … 9 ! … ~

a … z A … Z 0 … 9 ! … ~

# Neural Network Classifier

a … z A … Z 0 … 9 ! … ~£

95 x 1

104 x 1

112 x 1

7 x 2
2 x 7

72 x 1

x 7

5 x 5

(8x8)
(8x7;1,7) (7x8;7,1)
(8x6;1,8) (6x8;8,1)

1 x 9

9 x 1

(10x10)
(14x6)

(6x14)→

20 x 9

14 x 14

5 x 1

1 x 1

Stroke Feature

Stroke
Count

Aspect
Ratio

Image

# Normalizing Output Error

- υ Normalize "pressure towards zero"
- υ Based on recognition of the fact that most training signals are zero
- υ Training vector for letter "x"
  ```
  a ... w x y z A ... Z 0 ... 9 ! ... ~
  0 ... 0 1 0 0 0 ... 0 0 ... 0 0 ... 0
  ```
- υ Forces net to attempt to make unambiguous classifications
- υ Makes it difficult to obtain meaningful 2nd and 3rd choice probabilities
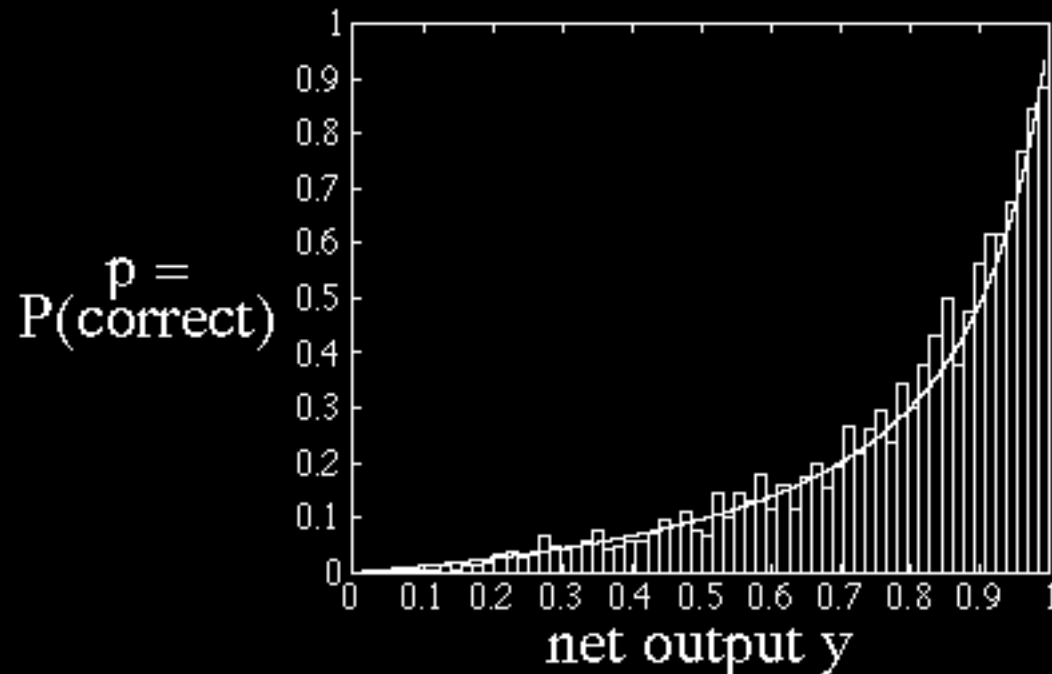
# Normalized Output Error

ᴜ We reduce the BP error for non-target classes relative to the target class

ᴜ By a factor that "normalizes" the non-target error relative to the target error

ᴜ Based on the number of non-target vs. target classes

ᴜ For non-target output nodes

$$e' = e\,A$$

where $A = 1 / d\,(N_{outputs} - 1)$

ᴜ Allocates network resources to modeling of low-probability regime

# Normalized Output Error

ʊ Converges to MMSE estimate of
$$f( P(class|input), A )$$

ʊ We derived that function:
$$\langle \hat{e}^2 \rangle = p \, (1-y)^2 + A \, (1-p) \, y^2$$
where $\quad p = P(class|input),$
$\qquad y = $ output unit activation

ʊ Output y for particular class is then:
$$y = p \, / \, (A - A \, p + p)$$

ʊ Inverting for p:
$$p = y \, A \, / \, (y \, A - y + 1)$$

# Normalized Output Error



Empirical p vs. y histogram for a net trained with
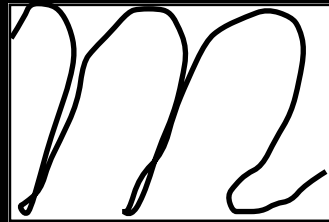A=0.11 (d=0.1), with corresponding theoretical curve
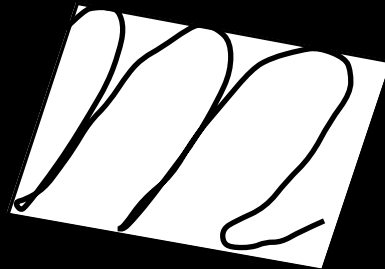
# Normalized Output Error

# Stroke Warping

ʋ Produce random variations in stroke data during training

ʋ Small changes in skew, rotation, x and y linear and quadratic scaling

ʋ Consistent with stylistic variations

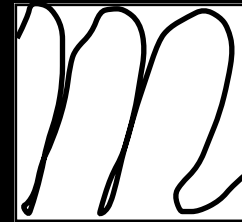ʋ Improves generalization by effectively adding extra data samples
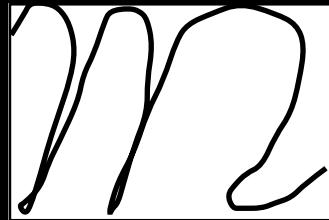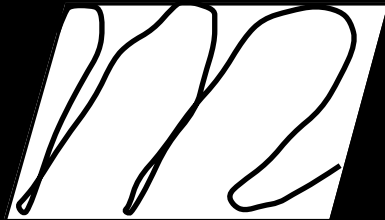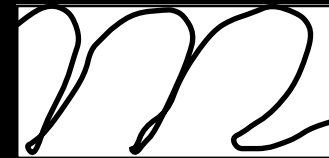
# Stroke Warping

Original       Rotation       X Linear

X Quadratic       X Skew       Y Linear

# Class Frequency Balancing

υ Skip and repeat patterns

υ Instead of dividing by the class priors

    υ Eliminates noisy estimate of low freq. classes

    υ Eliminates need for renormalization

    υ Forces net to better model low freq. classes

υ Compute normalized frequency, relative to average frequency $\quad F_i = S_i / \overline{S}$

$$\overline{S} = 1 / C \sum_{i=1}^{C} S_i$$

# Class Frequency Balancing

υ Compute repetition factor

$$R_i = ( a / F_i )^b$$

υ Where  `a (0.2 to 0.8)`  controls amount of skipping vs. repeating

υ And  `b (0.5 to 0.9)`  controls amount of balancing

# Stroke-Count Frequency Balancing

υ Compute frequencies for stroke-counts in each class

υ Modulate repetition factors by stroke-count sub-class frequencies

$$R_{ij} = R_i [(S_i/J)/S_{ij}]^b$$

# Adding Noise to Stroke-Count

υ Small percentage of samples use randomly selected stroke-count (as input to the net)

υ Improves generalization by reducing bias towards observed stroke-counts

υ Even improves accuracy on data drawn from training set

# Negative Training

υ Inherent ambiguities force segmentation code to generate false segmentations

υ Ink can be interpreted in various ways...

clog

υ "dog", "clog", "cbg", "%g"

υ Train network to compute low probabilities for false segmentations

# Negative Training

υ Modulate negative training two ways…

  υ Negative error factor (0.2 to 0.5)

    υ Like A in normalized output error

  υ Negative training probability (0.05 to 0.3)

    υ Also speeds training

υ Too much negative training

  υ Suppresses net outputs for characters that look like elements of multi-stroke characters
  `(I, 1, l, |, o, O, 0)`

υ Slight reduction in character accuracy, large gain in word accuracy

# Error Emphasis

υ Probabilistically skip training for correctly classified patterns

υ Never skip incorrectly classified patterns

υ Just one form of error emphasis

  υ Can reduce learning rate/error for correctly classified patterns

  υ And/or increase learning rate/error for incorrectly classified patterns

  υ Maintain pool of samples from which correctly classified patterns are flushed each epoch

# Training Probabilities and Error Factors

| Segment | Type | Prob. of Usage | | Error Factor | |
|---|---|---|---|---|---|
| | | Correct | Incorrect | Target Class | Other Classes |
|  | POS | 0.5 | 1.0 | 1.0 | 0.1 |
|  | NEG | 0.18 | | NA | 0.3 |

# Annealing & Scheduling

υ Start with large learning rate, then decay

  υ When training set's total squared error increases

υ Start with high error emphasis, then decay

υ Start with minimal negative training, then increase
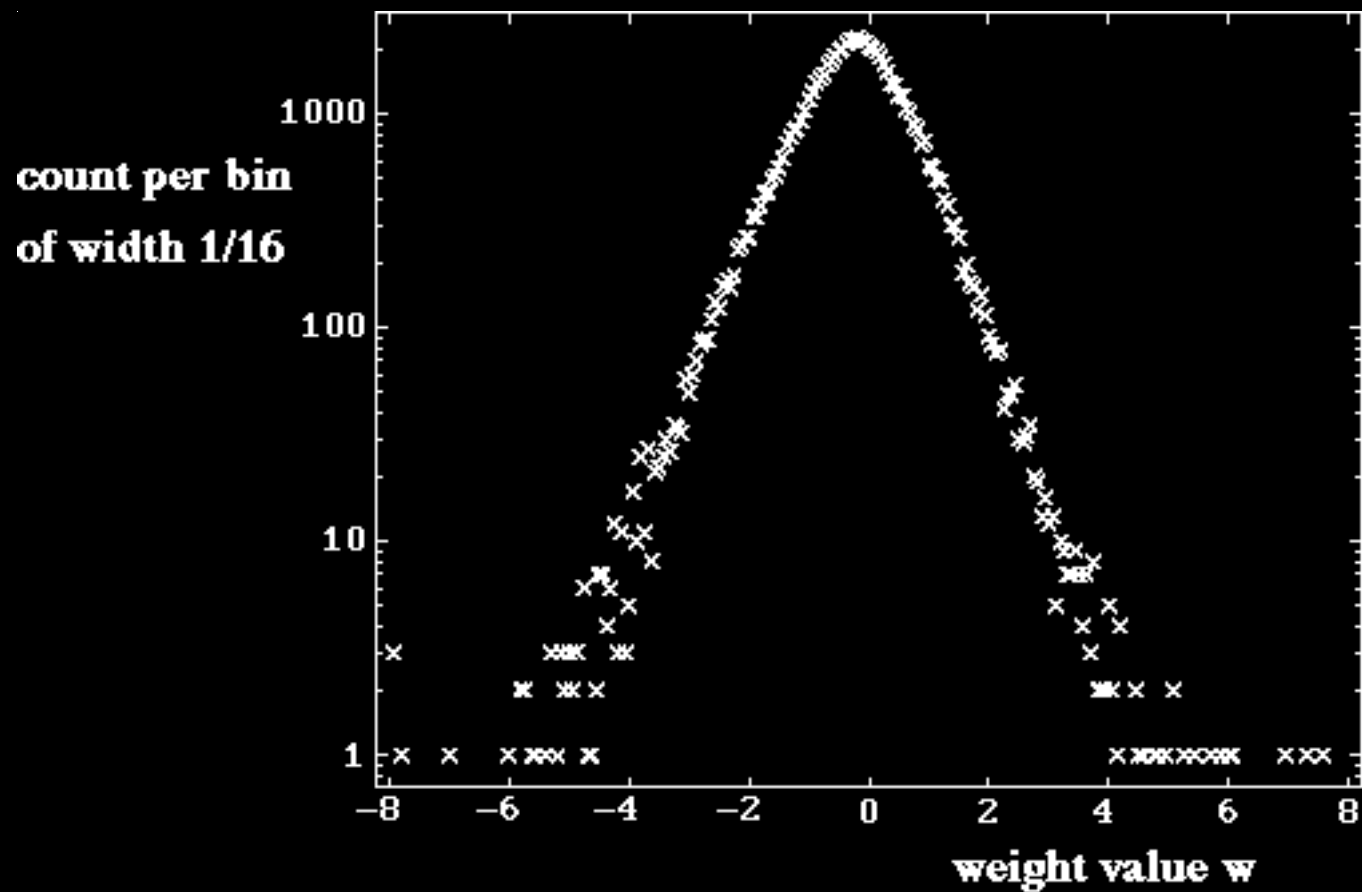
  υ Mostly for pragmatic reasons

# Training Schedule

| Phase | Epochs | Learning Rate | Correct Train Prob | Negative Train Prob |
|---|---|---|---|---|
| 1 | 25 | 1.0 - 0.5 | 0.1 | 0.05 |
| 2 | 25 | 0.5 - 0.1 | 0.25 | 0.1 |
| 3 | 50 | 0.1 - 0.01 | 0.5 | 0.18 |
| 4 | 30 | 0.01 - 0.001 | 1.0 | 0.3 |

# Quantized Weights

υ Forward/classification pass requires less precision than backward/learning pass

 υ Use one-byte weights for classification

  υ Saves both space and time

  υ ±3.4   (-8 to +8 with 1/16 Steps)

 υ Use three-byte weights for learning

  υ ±3.20

 υ First Newton version

  υ ~200KB ROM (~85KB for weights)

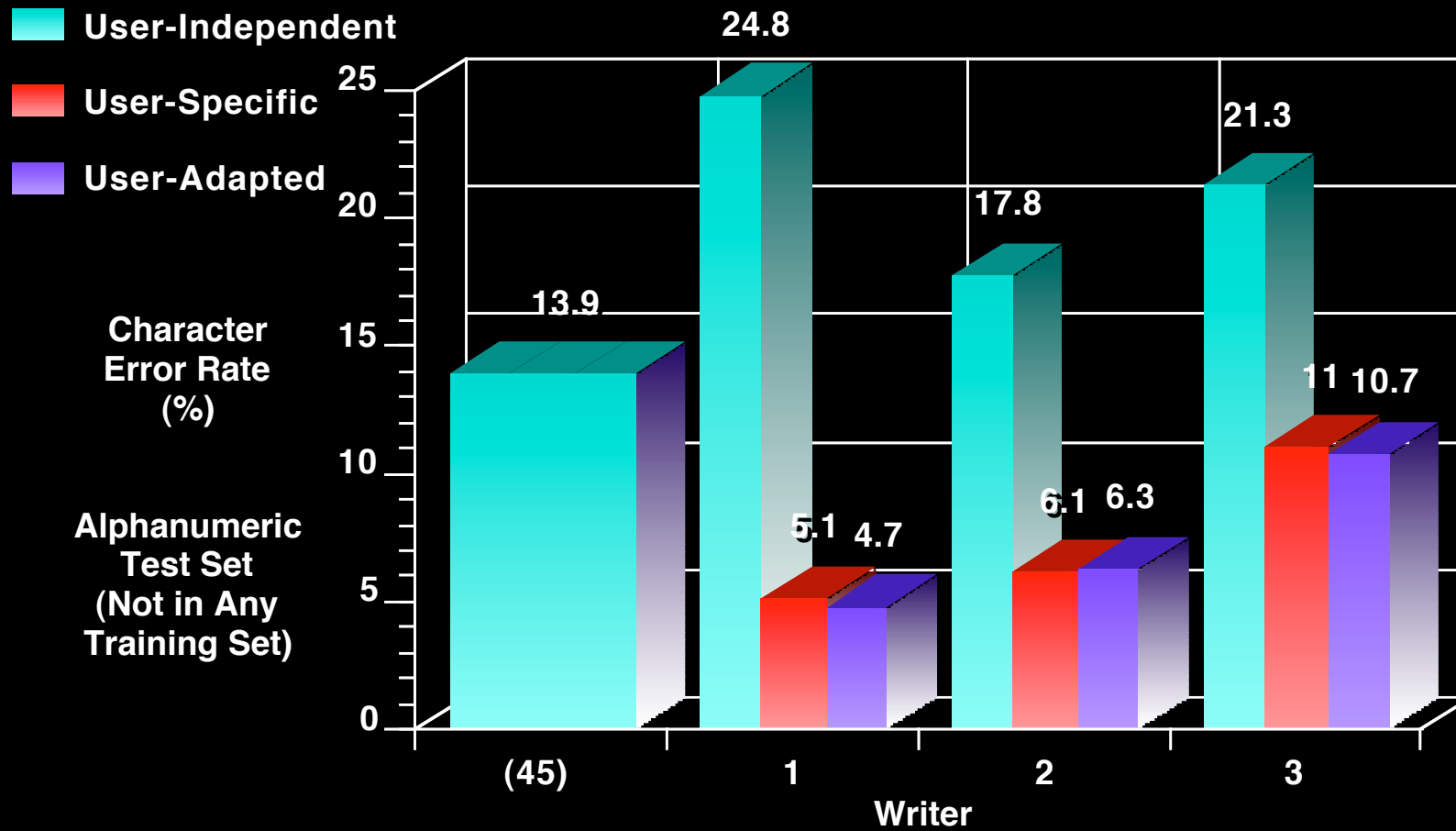  υ ~5KB-100KB RAM

  υ ~3.8 char/second

# Quantized Weights

# User Adaptation

υ Neural net classifer based on an inherently learning technology

υ Learning not used in Newton due to memory constraints

υ Learning not (yet) used in Mac OS X due to limited human resources

υ Can reduce error rates by factor of 2 to 5, yet user-independent "walk-up" performance is maintained!

# User Adaptation

υ User training scenario

  υ 15-20 min. of data entry

    υ Less for problem characters alone

  υ Possibly < 1 min. network learning

    υ One-shot learning may suffice (single epoch)

    υ May learn during data entry

    υ Maximum of a few minutes (~10-12 Epochs)

υ Learn on the fly

  υ Can continuously adapt

  υ Need system hooks

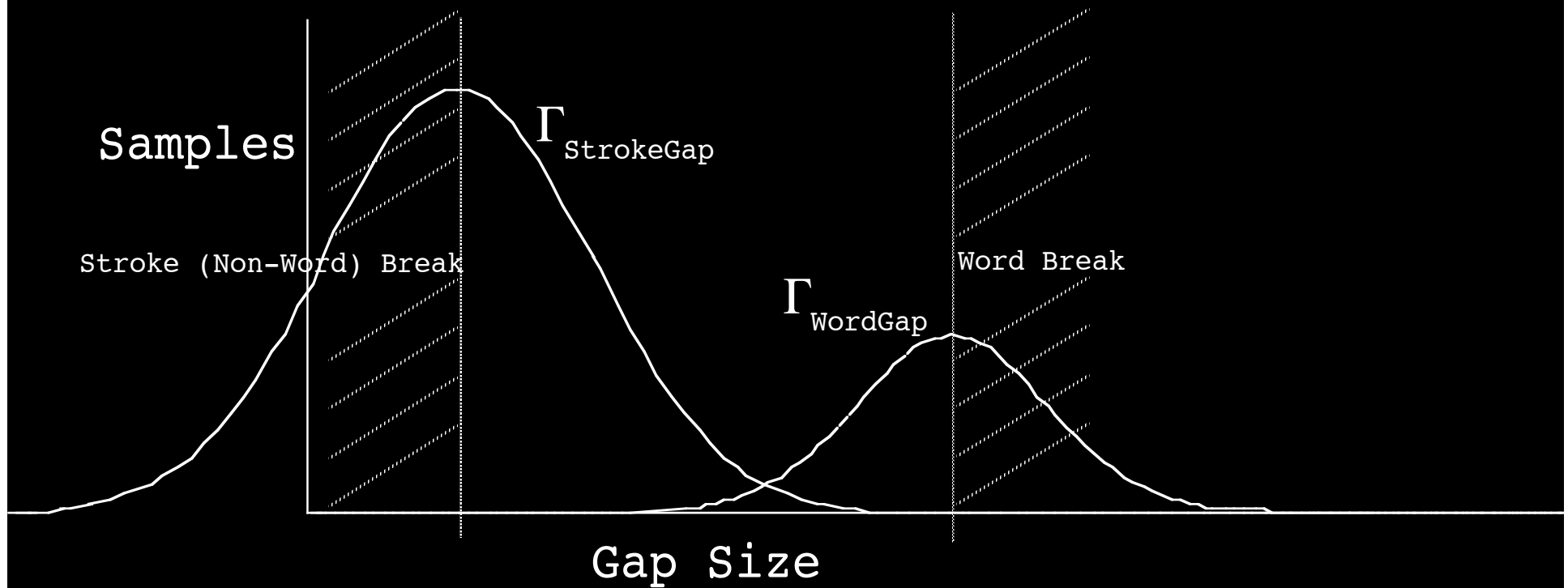  υ Choosing what to train on is key system issue

# User Adaptation

# Integration with Character Segmentation

ʊ Search takes place over segmentation hypotheses (as well as character hypotheses)

ʊ Stroke recombinations are presented in regular, predictable order

ʊ Forward and reverse "delay" parameters suffice to indicate legal time-step transitions

# Integration with Word Segmentation

υ Search also takes place over word segmentation hypotheses

υ Word-space becomes an optional segment/character

υ Weighted by probability ("SpaceProb") derived from statistical model of gap sizes and stroke centroid spacing

υ Non-space hypotheses are weighted by 1-SpaceProb

# Word Segmentation Statistical Model

Samples

$\Gamma_{\text{StrokeGap}}$

Stroke (Non-Word) Break

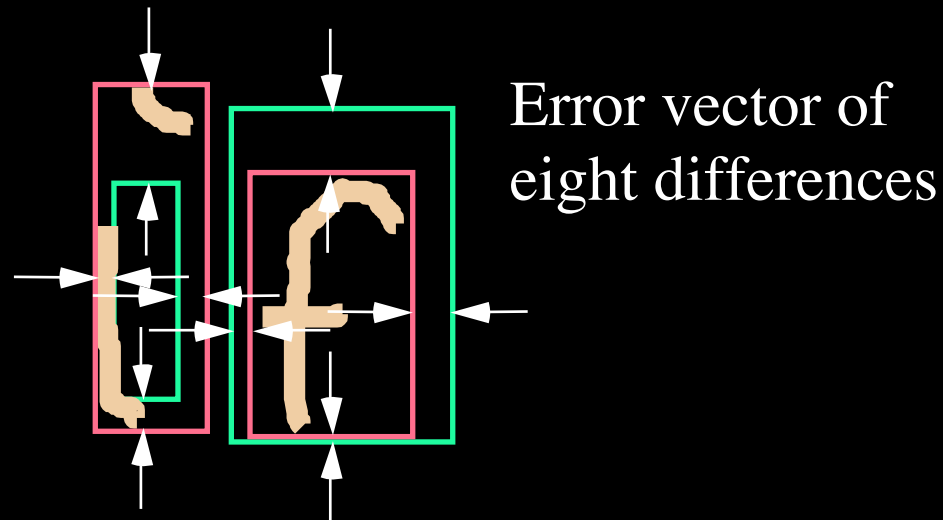$\Gamma_{\text{WordGap}}$

Word Break

Gap Size

$$P_{\text{WordBreak}} = \Gamma_{\text{WordGap}} / (\Gamma_{\text{StrokeGap}} + \Gamma_{\text{WordGap}})$$

# Recognition Ambiguity

scuba

# Geometric Context

"if" from User vs Table

Error vector of
eight differences

(User data scaled and offset to
minimize error magnitude)

# Language Model

- υ Dictionaries
  - υ Word lists
  - υ Regular expression grammars
- υ BiGrammars - combinations of dictionaries
  - υ Probabilistically weighted
  - υ Flexible starts, stops, and transitions

# Regular Expression Grammars

υ Telephone numbers example:

```
dig    = [0123456789]
digm01 =   [23456789]

acodenums = (digm01 [01] dig)

acode  = { ("1-"?    acodenums "-"):40 ,
           ("1"? "(" acodenums ")"):60 }

phone = (acode? digm01 dig dig "-" dig dig dig dig)
```

# Bigrammars

- υ Limited context telephone example:

```
BiGrammar2 Phone

[phone.lang 1. 1. 1.]
```

# BiGrammars

υ (Fairly) general context example:

```
BiGrammar2 FairlyGeneral
[EndPunct.lang  0.  .9  .5  EndPunct.lang .1]
(.8
  (.6
    [WordList.dict .5  .8  1. EndPunct.lang .2]
    [User.dict     .5  .8  1. EndPunct.lang .2]
  )
  (.4
    [Phone.lang    .5  .8  1. EndPunct.lang .2]
    [Date.lang     .5  .8  1. EndPunct.lang .2]
  )
)
(.2
  [OpenPunct.lang  1.  0.  .5
    (.6
      WordList.dict .5
      User.dict     .5
    )
    (.4
      Phone.lang    .5
      Date.lang     .5
    )
  ]
)
```
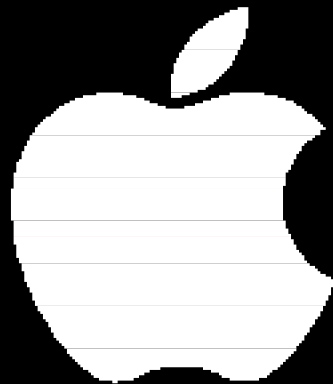
# Old Newton Writing Example

when Year-old Arabian retire tipped off the Christmas wrapping, No square with delights  Santa brought the Attacking hit too dathe would  Problem was, Joe talked Bobbie.  His doll stones at the really in its army Antiques I machine gun and hand decades At its side.  But it says things like 3 "Want togo shopping"  The Pro has claimed responsibility  that's Bobbie Liberation Organization.  Make up of more than 50 Concerned parents 3 Machinist 5 and oth er activists 5 the Pro claims to hsve crop if Housed switched the voice boxes on 300 hit, Joe and Bobbie foils across the United States this holiday Season  we have operations All over the country" said one pro member 5 who wished to remain autonomous.  "Our goal is to cereal and correct Thu problem of exposed stereo in editorials toys."

# Mondello Writing Example

When 7-year-old Zachariah Zelin ripped off the Christmas wrapping, he squealed with delight.  Santa brought the talking G.I. Joe doll he wanted.  Problem was, Joe talked like Barbie.  His doll stands at the ready in its Armyfatigues, machine gun and hand grenades at its side.  But it says things like, "Want to go shopping?"  The BLO has claimed responsibility.  That's Barbie Liberation Organization.  Made up of more than 50 concerned parents, feminists and other activists, the BLO claims to have surreptitiously switched the voice boxes on 300 G.I. Joe and Barbie dolls across the United States this holiday season.  "We have operatives all over the country," said one BLO member, who wished to remain anonymous.  "Our goal is to reveal and Correct the problem of gender-based stereotyping in children's toys."

# Apple-Newton
# Handwriting Recognition

The Power +o he your 6est