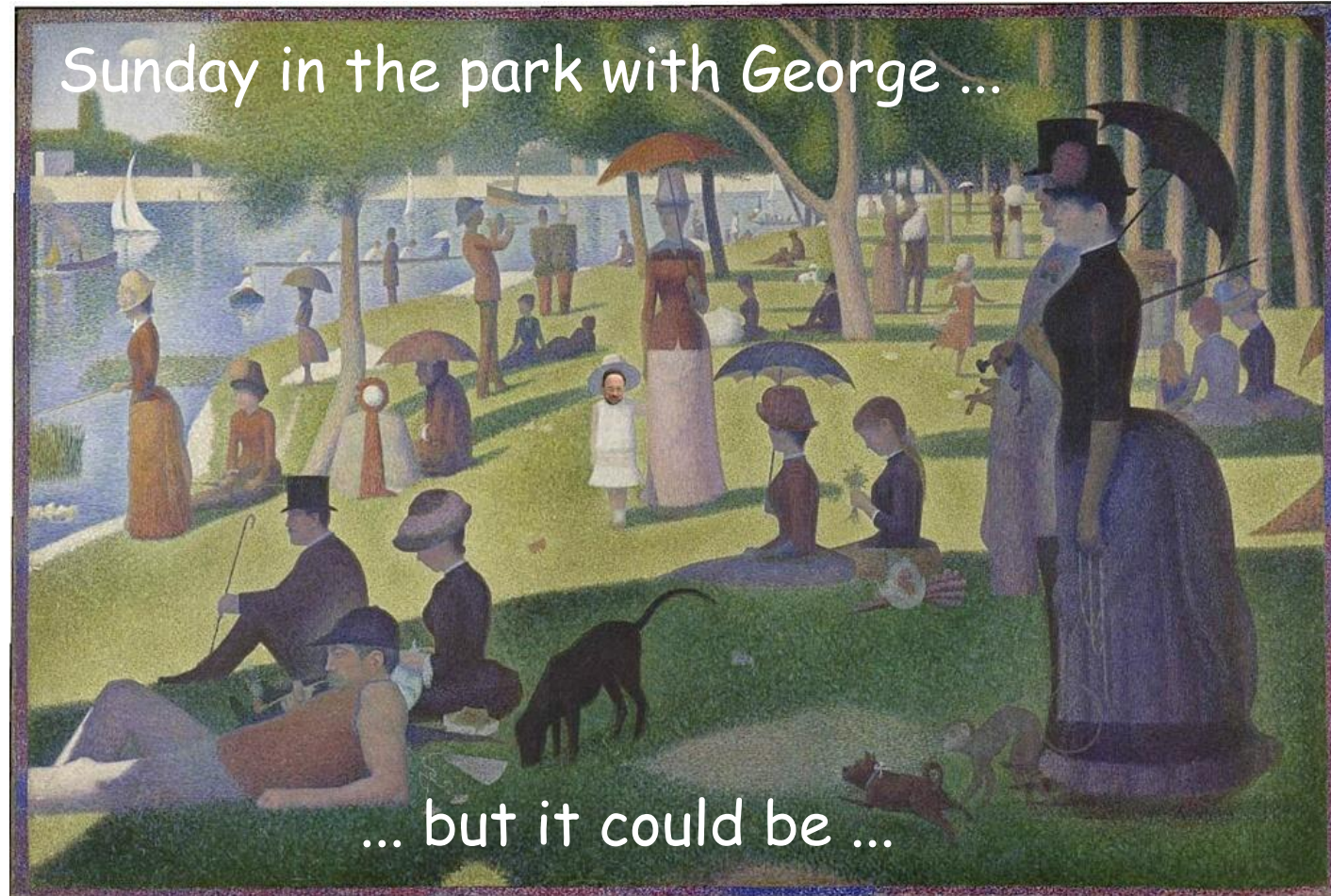


# What this talk is not ...



Adapted from "A Sunday Afternoon on the Island of La Grande Jatte" by Georges Seurat

# What this talk is not ...



*Curious George*

*... but it could be ...*

Adapted from *Curious George* by Margret and H.A. Rey

# What we're really talking about



Troy, NY, December 19, 2008

## Working with George

*Dan Lopresti*

Computer Science & Engineering  
Lehigh University  
Bethlehem, PA, USA

# What's the connection?

- Not a former student.
- Never employed by the same institution.
- Just a lucky bystander?



# Serendipity

serendipity (sɛrən dɪpɪti)

— n

the faculty of making fortunate discoveries by accident

*1994 = ancient history!*

- "Validation of Simulated OCR Data Sets" G. Nagy, *Proceedings of the Third Annual Symposium on Document Analysis and Information Retrieval*, April 1994, Las Vegas, NV, pp. 127-135.
- "Validation of Document Defect Models for Optical Character Recognition" Y. Li, D. Lopresti, and A. Tomkins, *Proceedings of the Third Annual Symposium on Document Analysis and Information Retrieval*, April 1994, Las Vegas, NV, pp. 137-150.

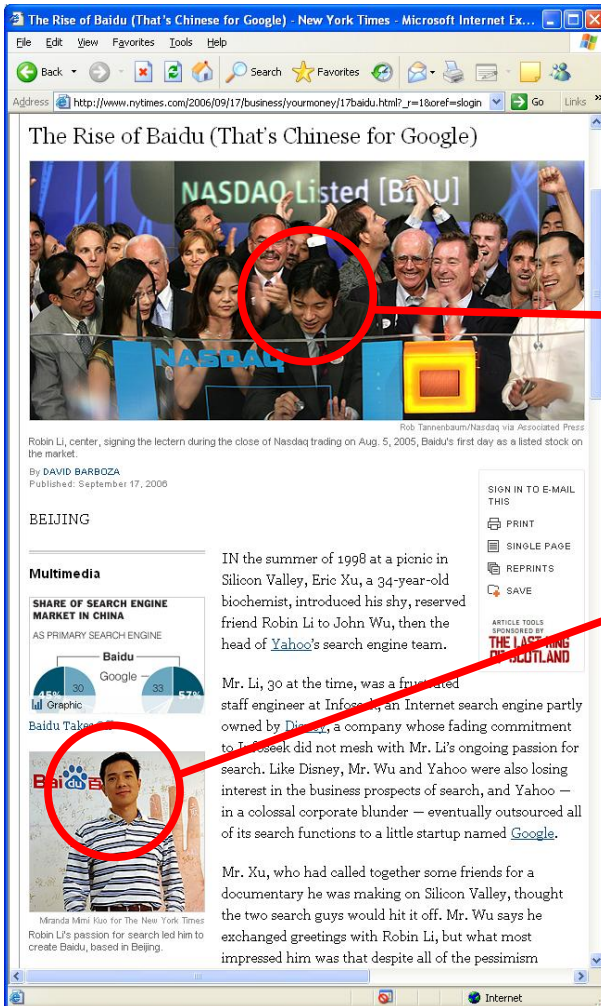
# Papers with George I



- "Spatial Sampling Effects in Optical Character Recognition," D. Lopresti, J. Zhou, G. Nagy, and P. Sarkar, *Proceedings of the Third International Conference on Document Analysis and Recognition*, August 1995, Montréal, Canada, pp. 309-314.
- "Validation of Image Defect Models for Optical Character Recognition," Y. Li, D. Lopresti, G. Nagy, and A. Tomkins, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 18, no. 2, February 1996, pp. 99-108.
- "Spatial Sampling Effects on Scanned 2-D Patterns," J. Zhou, D. Lopresti, P. Sarkar, and G. Nagy, *Advances in Visual Form Analysis*, C. Arcelli, L. P. Cordella, and G. Sanniti di Baja, eds., Singapore: World Scientific, 1997, pp. 666-675.
- "Spatial Sampling of Printed Patterns," P. Sarkar, G. Nagy, J. Zhou, and D. Lopresti, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 20, no. 3, March 1998, pp. 344-351.
- "Automated Table Processing: An (Opinionated) Survey," D. Lopresti and G. Nagy, *Proceedings of the Third IAPR International Workshop on Graphics Recognition*, September 1999, Jaipur, India, pp. 109-134.
- "Issues in Ground-Truthing Graphic Documents," D. Lopresti and G. Nagy, *Proceedings of the Fourth IAPR International Workshop on Graphics Recognition*, September 2001, Kingston, Ontario, Canada, pp. 59-72.

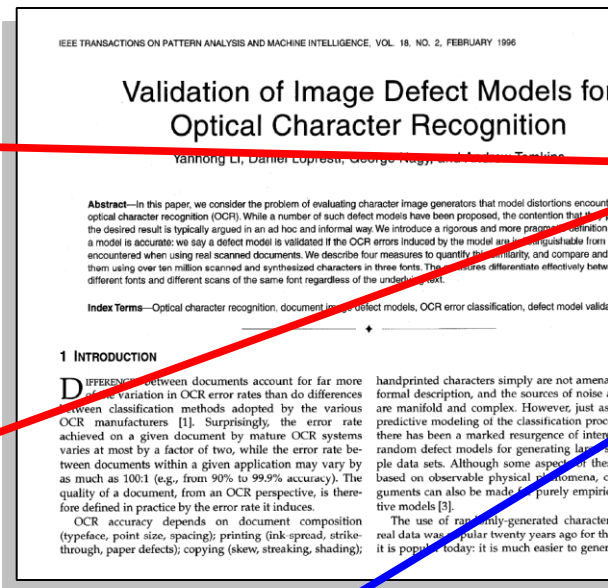
# Brush with fame and fortune

From the New York Times, Sunday Sept. 17, 2006.



Robin Li, a student who worked with us, is now billionaire founder of Baidu.

Robin



**Robin Li** received the BS degree in information science from Peking University, China in 1991 and the MS degree in computer science from the State University of New York at Buffalo in 1993. Since 1994, he has been a senior software engineer at GARI Software, IDD Information Services. From 1992 to 1994, he was a research assistant at the Center of Excellence for Document Analysis and Recognition (CEDAR) of SUNY Buffalo. In the summer of 1993, he did an internship at Matsushita Information Technology Laboratory in Princeton, N.J. His research interests include document analysis, information retrieval, text compression, and financial information routing.

**Daniel Lopresti** received the AB degree in mathematics and engineering from Dartmouth College in 1982 and the PhD degree in computer science from Princeton University in 1987. From 1986 until 1991, he was on the faculty of the Department of Computer Science at Brown University. In 1991 he joined the newly formed Matsushita Information Technology Laboratory as a senior scientist and leader of the Carbon Project. His research interests include pattern matching and recognition, parallel VLSI architectures, and computational aspects of molecular biology.

**George Nagy** received the BEng. and MEng. degrees from McGill University and the PhD in electrical engineering from Cornell University in 1962 (on neural networks). For the next ten years he conducted research on various aspects of pattern recognition and OCR at the IBM T.J. Watson Research Center in Yorktown Heights. From 1972 to 1985 he was a professor of computer science at the University of Nebraska-Lincoln and worked on remote sensing applications, geographic information systems, computational geometry, and human-computer interfaces. Since 1985 he has been a professor of computer engineering at Rensselaer Polytechnic Institute. He has held visiting appointments at the Stanford Research Institute, Cornell, the University of Montreal, the National Scientific Research Institute of Quebec, the University of Genoa and the Italian National Research Council in Naples and Genoa, AT&T Bell Laboratories, IBM Almaden, McGill University, and the Institute for Information Science Research at

Me

George

# Papers with George II



- "Why Table Ground-Truthing is Hard," J. Hu, R. Kashi, D. Lopresti, G. Nagy, and G. Wilfong, *Proceedings of the Sixth International Conference on Document Analysis and Recognition*, September 2001, Seattle, WA, pp. 129-133.
- "A Nonparametric Classifier for Unsegmented Text," G. Nagy, A. Joshi, M. Krishnamoorthy, Y. Lin, D. Lopresti, S. Mehta, and S. Seth, *Proceedings of Document Recognition and Retrieval XI (IS&T/SPIE International Symposium on Electronic Imaging)*, January 2004, Santa Jose, CA, pp. 102-108.
- "Chipless ID for Paper Documents," D. Lopresti and G. Nagy, *Proceedings of Document Recognition and Retrieval XII (IS&T/SPIE International Symposium on Electronic Imaging)*, January 2005, San Jose, CA, pp. 208-215.
- "Mobile Interactive Support System for Time-Critical Document Exploitation," G. Nagy and D. Lopresti, *Symposium on Document Image Understanding Technology*, November 2005, College Park, MD, pp. 111-119.
- "Match Graph Generation for Symbolic Indirect Correlation," D. Lopresti, G. Nagy, and A. Joshi, *Document Recognition and Retrieval XIII (IS&T/SPIE International Symposium on Electronic Imaging)*, January 2006, San Jose, CA, pages 606706.1-606706.9.



# Papers with George III



- "Notes on Contemporary Table Recognition," D. Embley, D. Lopresti, and G. Nagy, *Proceedings of the Seventh IAPR International Workshop on Document Analysis Systems*, H. Bunke and A. L. Spitz, eds., Berlin: Springer-Verlag, 2006, pp. 164-175.
- "Interactive Document Processing and Digital Libraries," G. Nagy and D. Lopresti, *Proceedings of the Second International Conference on Document Image Analysis for Libraries*, April 2006, Lyon, France, pp. 2-11.
- "Table Processing Paradigms: A Research Survey," D. Embley, M. Hurst, D. Lopresti, and G. Nagy, *International Journal on Document Analysis and Recognition*, vol 8, no. 2-3, June 2006, pp. 66-86.
- "A Maximum-Likelihood Approach to Symbolic Indirect Correlation," A. Joshi, G. Nagy, D. Lopresti, and S. Seth, *Proceedings of the Eighteenth International Conference on Pattern Recognition*, August 2006, Hong Kong, pp. 99-103.
- "Multi-Character Field Recognition for Arabic and Chinese Handwriting," D. Lopresti, G. Nagy, S. Seth, and X. Zhang, *Proceedings of the Summit on Arabic and Chinese Handwriting Recognition*, September 2006, College Park, MD, pp. 93-100.

# Papers with George IV



- "A Document Analysis System for Supporting Electronic Voting Research," D. Lopresti, G. Nagy, and E. Barney Smith, *Proceedings of the Eighth IAPR International Workshop on Document Analysis Systems*, IEEE Computer Society Press, September 2008, Nara, Japan, pp. 167-174.
- "Ballot Mark Detection," E. Barney Smith, D. Lopresti, and G. Nagy, *Proceedings of the Nineteenth International Conference on Pattern Recognition*, December 2008, Tampa, FL, pages 4 (CD-ROM).
- "Mark Detection from Scanned Ballots," E. Barney Smith, D. Lopresti, and G. Nagy, *Proceedings of Document Recognition and Retrieval XVI (IS&T/SPIE International Symposium on Electronic Imaging)*, January 2009, San Jose, CA, pages 7247-26.01-7247-26.10.
- "Tools for Monitoring, Visualizing, and Refining Collections of Noisy Documents," D. Lopresti and G. Nagy, *Proceedings of the Third Workshop on Analytics for Noisy Unstructured Text Data*, July 2009, Barcelona, Spain, pp. 9-16.
- "Document Photography in Vitro," G. Nagy, B. Clifford, A. Berg, G. Saunders, E. Barney Smith, and D. Lopresti, *Proceedings of the Third International Workshop on Camera-Based Document Analysis and Recognition*, July 2009, Barcelona, Spain, pp. 26-33.

# Papers with George V



- "Camera-based Ballot Counter," G. Nagy, B. Clifford, A. Berg, G. Saunders, D. Lopresti, and E. Barney Smith, *Proceedings of the Tenth International Conference on Document Analysis and Recognition*, July 2009, Barcelona, Spain, pp. 151-155.
- "Style-Based Ballot Mark Recognition," P. Xiu, D. Lopresti, H. Baird, G. Nagy, and E. Barney Smith, *Proceedings of the Tenth International Conference on Document Analysis and Recognition*, July 2009, Barcelona, Spain, pp. 216-220.
- "Document Analysis Issues in Reading Optical Scan Ballots," D. Lopresti, G. Nagy, and E. Barney Smith, *Proceedings of the Ninth IAPR International Workshop on Document Analysis Systems*, June 2010, Boston, MA, pp. 105-112.
- "Characterizing Challenged 2008 Minnesota Ballots," G. Nagy, D. Lopresti, E. H. Barney Smith, and Z. Wu, *Document Recognition and Retrieval XVIII (IS&T/SPIE International Symposium on Electronic Imaging)*, January 2011, San Francisco, CA.
- "When is a Problem Solved?," D. Lopresti and G. Nagy, to be presented at the *Eleventh International Conference on Document Analysis and Recognition (ICDAR 2011)*, September 2011, Beijing, China.
- "Towards Improved Paper-based Election Technology," E. Barney Smith, D. Lopresti, G. Nagy, and Z. Wu, to be presented at the *Eleventh International Conference on Document Analysis and Recognition (ICDAR 2011)*, September 2011, Beijing, China.

# The Debate

"Defect Models are Important to Advance the State-of-the-Art of Optical Character Recognition"

April 1996

Las Vegas, NV

For:

- Henry Baird
- Bob Haralick

Against:

- Dan Lopresti
- George Nagy



# Decades of Influence



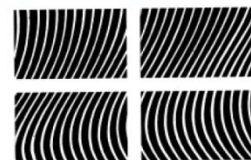
George Nagy graduate Physics (fencing and solving Euler's Second awarded the PhD at Cornell Rosenblatt build Tobe network for speech re character recognition claims credit for IBM reverse sabbatical at t trains from cats' medi Department of Comp where he dabbled in c 1985 he has been Prof

Troy, NY. Nagy's credits in document analysis incl Casey, "self-corrective" character recognition with years later with Henry Baird), character recognition

## Keynote talk on datasets by George Nagy at DAS 2010 ...

### Data Sets advertised in IEEE Computer

#### January 1972 (6 data sets)



PATTERN RECOGNITION DATA BASES AVAILABLE

In order to encourage research in the field of pattern recognition, the IEEE Computer Society's Technical Committee on Pattern Recognition has begun collecting data bases from a variety of sources. These data bases, including substantial back-up documentation, may be ordered by using the form at the bottom of the page.

**1.1.1 Machine Imprinted Alphanumeric Characters** - Dr. H. F. Ryan, Cornell Aeronautical Laboratory, Ito/U.S. Postal Service  
An alphanumeric character data base of 100,000 samples of 48 character sizes. (Also known as the CAL-U.S. Postal Service Alphanumeric Character Data Base). Thresholded binary images of slanted, centered, mixed-font, machine-imprinted characters. Resolution is 24 x 24. Magnetic tape, 9 track, 2 reels, 1600 BPI.  
Price: \$112.50 (\$80.75 with furnished tape)  
Member's discount price: \$90, (\$55, with furnished tape)

4ICOMPUTER/JANUARY/FEBRUARY 1972

Discounts off the data base list prices are available to IEEE members and members of the American Federation of Information Processing Societies' constituent societies.

When ordering, you may elect to send us your own blank tapes, if you do, be sure they are in good condition and have no other data recorded on them.

**1.1.2 Handwritten Numeric Characters** - Dr. A. L. Kroll, Honeywell Information Systems, Data Systems Division  
The data base consists of 50 samples of each numeric character generated by 9 different authors. Simple printing rules were specified but not always followed. The samples were selected from those distributed. The images are binary with a resolution of 25 x 21. Punched cards.  
Price: \$27.50  
Member's discount price: \$30.

#### January 1973

### PATTERN RECOGNITION DATA BASES AVAILABLE



but not always followed. The samples were selected from those distributed. The images are binary with a resolution of 25 x 21. Punched cards.  
Price: \$41.25  
Member's discount price: \$33.

**1.1.2 Handprinted FORTRAN Alphanumeric Characters** - Dr. John H. Mueser, Stanford Research Institute  
The data base consists of two parts, with each part on a reel. The first part contains 3 alphabets of 48 characters, corresponding to the nondecimal character set of the basic FORTRAN language, handwritten to each of 25 random characters a total of 3 x 48 x 25 = 4,500 pictures.  
The second part has 2,999 characters printed by a typewriter. There are 928 characters in each of 28 alphabets of 48 characters each; the remaining 2,070 characters are taken from fragments of word ending words. The tapes are binary with a 24 x 24 resolution. Magnetic tape, 7 track, 2 reels, 254 BPI.  
Price: \$116.75 (\$75.75 with furnished

documentation). The samples were selected from those distributed. The images are binary with a resolution of 25 x 21. Punched cards.  
Price: \$41.25  
Member's discount price: \$33.

**1.1.2 Handprinted Numeric Characters** - Phoenix Graphix, Trinity Business Printing Co., Ltd/Tachida Research and Development  
The data base consists of 10,000 hand-printed numeric characters collected from 100 individuals as well as approximately 100,000 characters typed on 1000 with a resolution of 16 x 16. A single percent sample of 50 words. Magnetic tape, 7 track, 2 reels, 500 BPI, 4000 pictures per reel.  
Price: \$116.21 (\$84.76 with furnished tape)  
Member's discount price: \$93, \$54, with furnished tape

**1.1.1 Courier Script** - Dr. L. D. Hanson, Bell Telephone Laboratories  
The data base consists of 10,000 characters with a resolution of 12 x 12. Punched cards.  
Price: \$116.21 (\$84.76 with furnished tape)  
Member's discount price: \$93, \$54, with furnished tape

### Pattern Recognition DATA BASES



**1.1.1 Machine Imprinted Alphanumeric Characters** - Dr. H. F. Ryan, Cornell Aeronautical Laboratory, Ito/U.S. Postal Service  
An alphanumeric character data base of 100,000 samples of 48 character sizes. (Also known as the CAL-U.S. Postal Service Alphanumeric Character Data Base). Thresholded binary images of slanted, centered, mixed-font, machine-imprinted characters. Resolution is 24 x 24. Magnetic tape, 9 track, 2 reels, 1600 BPI.  
Price: \$112.50 (\$80.75 with furnished tape)  
Member's discount price: \$90, (\$55, with furnished tape)

Discounts off the data base list prices are available to IEEE members and members of the American Federation of Information Processing Societies' constituent societies.

When ordering, you may elect to send us your own blank tapes, if you do, be sure they are in good condition and have no other data recorded on them.

**1.1.1 Machine Imprinted Alphanumeric Characters** - Dr. H. F. Ryan, Cornell Aeronautical Laboratory, Ito/U.S. Postal Service  
An alphanumeric character data base of 100,000 samples of 48 character sizes. (Also known as the CAL-U.S. Postal Service Alphanumeric Character Data Base). Thresholded binary images of slanted, centered, mixed-font, machine-imprinted characters. Resolution is 24 x 24. Magnetic tape, 9 track, 2 reels, 1600 BPI.  
Price: \$112.50 (\$80.75 with furnished tape)  
Member's discount price: \$90, (\$55, with furnished tape)

\* Slides available on the DAS 2010 website.

#### AS 2010 (GN) April 1976 (10 data sets)

# A really good bad idea?

## Training Humans to Read Like Machines (or, The Lazy Researcher's Approach to Perfect OCR)

*Dan Lopresti and George Nagy (if he agrees)*

- For decades, document analysis researchers have labored with tremendous effort and unbridled enthusiasm in desperate attempts to raise accuracy rates for optical character recognition to 100%.
- Success has proved to be elusive for all but the cleanest of documents typeset using standard fonts, i.e., boring cases that present absolutely no challenge and that even a moderately-talented trained monkey could handle with one paw tied behind its back.

# A really bad good idea?

- Tired of seeing the field perpetuate this exercise in futility, in this work we propose a novel, radical, earth-shaking, ground-breaking, revolutionary, radical idea.
- We posit that if it is too hard to solve the problem, it is always possible to change the problem and thereby make it easier to solve.
- Our thesis is that if we train humans to read like machines - to make all of the same mistakes that our current computer algorithms make when processing a typical page image - we will instantly achieve 100% OCR accuracy with no additional research effort required.

# Support to back up our claims

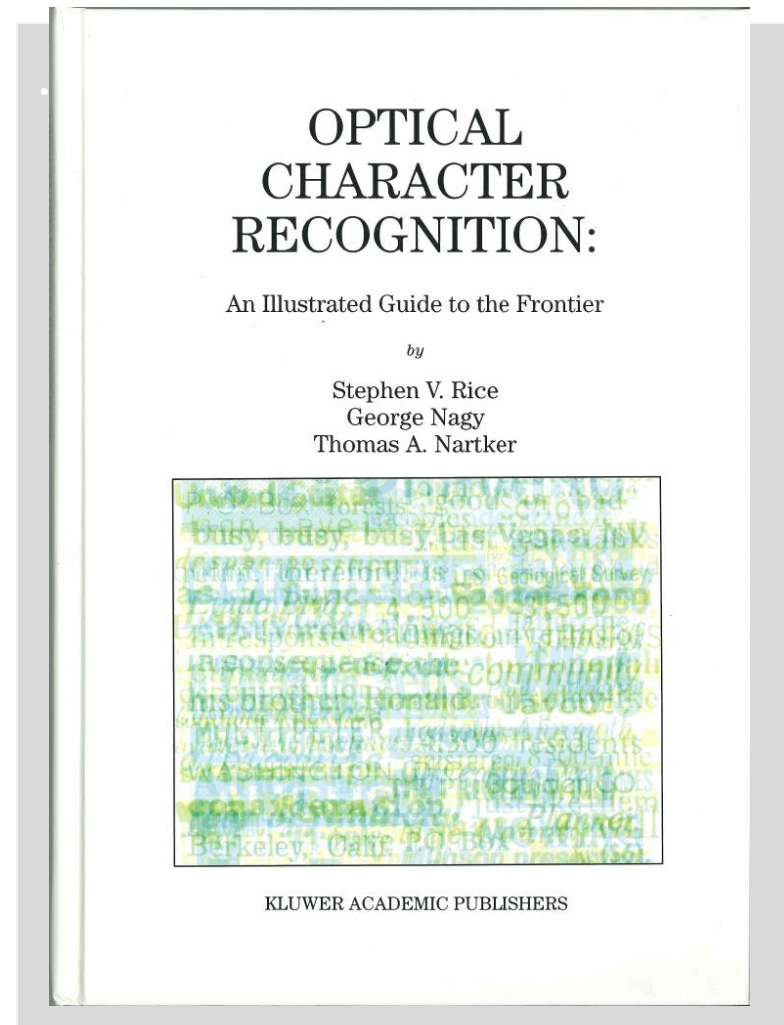
- We are quite certain this is feasible because humans are typically very smart and infinitely adaptable - making mistakes comes naturally to our species.



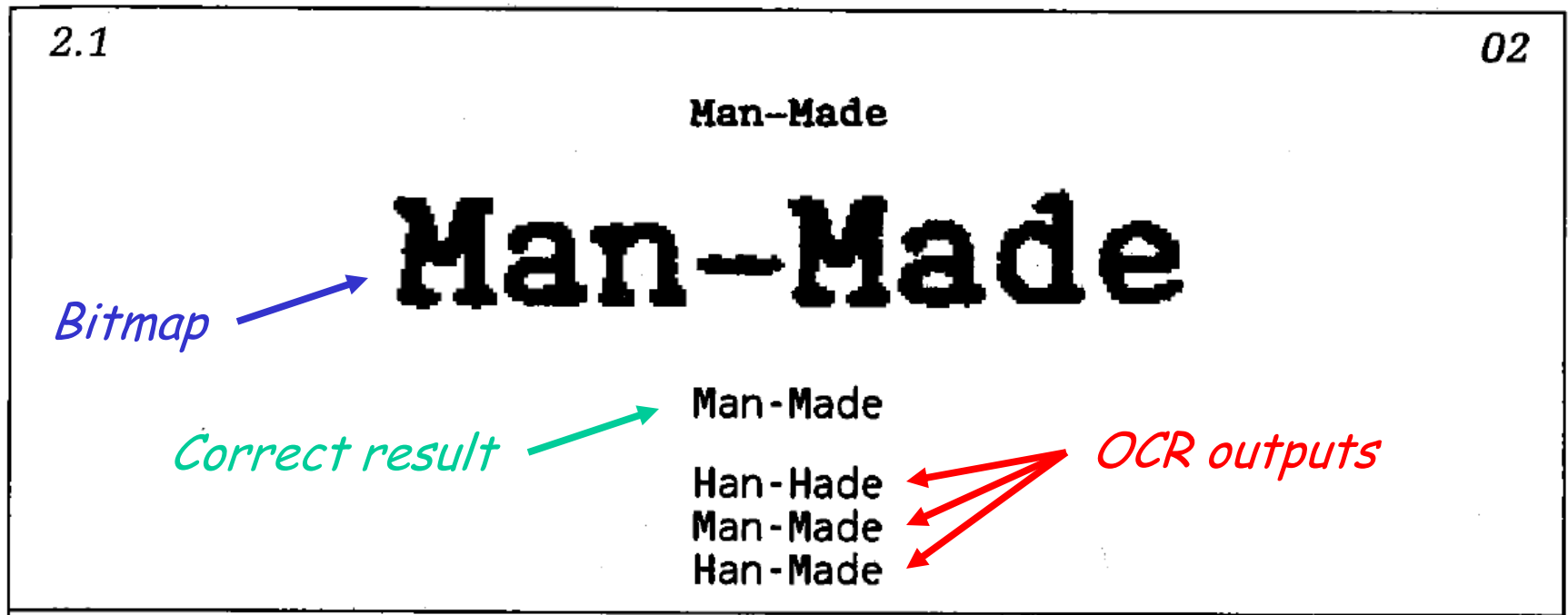


# Proof? We don't need no proof.

- Instead, we shall conduct a brief but revealing demonstration using this classic tome:



# The task: read like a machine



# The task: read like a machine



Unnamed Mineral

What would a machine output for this bitmap?

- (A) Unnamed Mineral
- (B) I think that I shall never see ...
- (C) Unnamed Ninerall ← *Correct answer*

# The task: read like a machine

Holt (1984).

What would a machine output for this bitmap?

- (A) Bolt (1984). ← *Correct answer*
- (B) Holt (1984).
- (C) ... a poem as lovely as a tree,

# The task: read like a machine



**McGovern**

What would a machine output for this bitmap?

- (A) A penny saved is a penny earned.
- (B) McGovern
- (C) McGovem ← *Correct answer*

The task: read like a machine

Imaging Experts

What would a machine output for this bitmap?

(A) Imaging Experts

(B) A stitch in time saves nine.

(C) Iiii~in~ L\1)(4~ ← *Correct answer*

# The task: read like a machine

**Great Lakes**

What would a machine output for this bitmap?

- (A) May you live in interesting times.
- (B) Grear Lakes ← *Correct answer*
- (C) Great Lakes

# The task: read like a machine

4,300 residents

What would a machine output for this bitmap?

- (A) 4.300 residents ← *Correct answer*
- (B) 4,300 residents
- (C) A miss is as good as a mile.



# The task: read like a machine

I have learned

What would a machine output for this bitmap?

- (A) I have learned ← *Correct answer*
- (B) I have learned
- (C) What, me worry?

# The task: read like a machine

just like them

What would a machine output for this bitmap?

(A) just like them

(B) justlike them ← *Correct answer*

(C) A rolling stone gathers no moss.

# The task: read like a machine

Do you think you could learn to make the same mistakes a machine would make?

- (A) Yes, I already make those same mistakes.
- (B) Yes, I'm smarter than a dumb machine.
- (C) Yes, anything you say, just stop talking!

Voilà!

Perfect OCR!!!  
(or "Perfect OCR!!!")

# A final Haiku

Farewell RPI  
George Nagy is retired  
He is all ours now

Congratulations and best wishes, George and Jill, for a long, healthy, enjoyable, fulfilling retirement!