# From Paper *to* Web

How to
Make
Information
Instantly
Accessible

**Tony McKinley**

ADOBE
PRESS

# table of contents

## Part 1: Creating Digital Content

# Chapter 6: Decisions in Indexing 117

# Part 3: Searching Digital Content

## Chapter 7: Acrobat Search 131

# Chapter 8: Enhanced PDF Collections on the Web 159

# Chapter 9: Advanced Navigation for Superior Information Access 179

# Part 4: Using Digital Content

# Part 5: Advanced Resource Guide

# index

# Acknowledgments

# Introduction

In this book I've tried to share my excitement and enthusiasm for the great new technology that is providing easy and unprecedented access to information. It might seem like a strange obsession, but I've been dedicated to turning paper into digital form for over 15 years and finally — everything works!

Dating back to the earliest Optical Character Recognition machines, I've worked at finding better ways than typing to get information into computers. Of course, in the early days we mostly wanted to just get information off paper to simply edit it or recompose it and then print it right back out on paper. But, even that was exciting. Look, a machine that reads and types!

By the early '80s, more interesting things were being done with electronic information, and I had the opportunity to work with such online database pioneers as Mead Data Central, Chemical Abstracts and BIOSIS. Now that digital information itself was being consumed over online terminals, OCR seemed to have found a higher calling and I was lucky enough to work for Kurzweil Computer Products.

The information that we wanted to convert to digital form was now in typeset form, unlike the far simpler typewritten documents that early OCR could handle. We could "train" our Kurzweil scanners to recognize magazines for Mead's NEXIS database, and for the giraffe-high stack of Chemical Abstracts published in hard cover. Compared to OCR's earlier limitations of being able to read only about a dozen typestyles, Kurzweil Intelligent Scanning Systems were a real breakthrough.

In those days, there were precious few prospects for our scanners. One limitation may have been the $50,000+ price tag, but the more severe limitation was that there just wasn't much you could do with digital documents. By the late '80s local area networks and personal computers had become more powerful and popular. But still, to equip a PC to run a moderately fast scanner, we had to add special hardware boards that often cost more than the computers they were installed in. By 1990, I took great pride in the fact that we could put together a decent one-seat scanning system for about $25,000.

But every step of the way was still very expensive, and everything had to be beefed up for imaging applications, including storage, networks, displays and printers. The high speed OCR systems still cost $20,000-$30,000 or more. Then, as Moore's law kicked in over a few generations of PCs the computers finally grew into the demanding requirements of imaging and recognition. But still, transmitting images beyond

the LAN was demanding and expensive, requiring dedicated communications lines. And CD-ROM readers cost in the thousands, and CD recorders were beyond the reach of all but the hardiest companies boldly investing in electronic publishing.

It was a lot of fun, but document imaging never had the great breakout that we all expected for so long. We were still limited to niche applications, not far removed from the earliest industrial and professional imaging applications. Then, the Web came along and woke up the slumbering giant of the Internet. Suddenly we had something to do with all the digital documents we could create!

The fact that the PCs, CD drives, storage, monitors, modems and every other piece continually got faster and cheaper certainly helped move things along too.

But my personal expectations remained unfulfilled. Optical character recognition was still expensive because it was difficult to reproduce the appearance of paper documents and it was costly to edit and clean up the results of OCR.

Then one day a couple of years ago, I was walking by the Adobe booth at the AIIM show, the biggest annual event in the imaging industry. It was curious that Adobe was even there, but what I was seeing on the big screen monitor was even more curious. They had just scanned a page, done recognition, and were displaying the results. Unlike the results of OCR, which contains tildes for unrecognizable characters in the text, there were small images of the suspect words blended into the document! And the output looked just like the original! To someone who had spent 15 years toiling in the field of OCR, spanning hundreds of installations and millions of dollars worth of systems, this was absolutely incredible.

As soon as it became available, I got my hands on a copy of Adobe Acrobat Capture 1.0 and put it to the test in our Online OCR Lab. I was stunned. I literally sat in front of my computer in wide-eyed astonishment. Now, I don't expect everyone to feel that way, but then not everyone has personally done more than 2,000 OCR demos like I have. If I'd had Capture over the years I could have met the needs of hundreds and hundreds of customers that conventional OCR simply could not satisfy.

I have one more confession to make, and that is that my interest in converting paper to digital documents goes back beyond my years in OCR. My inspiration for this field comes from Buckminster Fuller, in his 1962 book Education Automation. In that typically freewheeling talk Bucky proposed a universally accessible digital library that would enable anyone, anywhere to study, learn and grow. Bucky figured that this intellectual freedom of the masses would bring humanity's best ideas to reality.

Now, that may sound like a lofty goal, but if we're going to get there, the path is clear before us. The theme of this book is that we can now provide superior access to information. It's not just that we can digitally miniaturize books and get them out of the library without leaving home. The key is that we can provide instant access to the information within the documents.

In my study of Acrobat Capture, I explored the efficiencies of Adobe's Portable Document Format, which offers a range of features from cross-platform viewing and printing to a built-in set of management and search capabilities. Not only does Capture do what my 2,000 OCR demo customers ask for, but PDF seems to meet the vaguely defined but clearly required needs of digital documents and digital libraries.

I've tried to address the entire experience of moving from paper to digital documents on the Web. After finding the most efficient way to create the new documents, we move on to the even more important question of how to organize the information within the digital library. Finally, the issues and techniques of Information Retrieval are explored to give both publishers and users some useful tips and tools for finding what they seek on the Web.

While the lessons learned in this book came from many years in business and academic applications of this technology, I still aspire to the ideals espoused by Bucky Fuller. Now that the technology to make the world a better place that he predicted is here, we're on our way.

— T.M. - 1/20/97

# Dedication

I lovingly dedicate this book to my wife, Patty
and my children, Laura, Tony and Hugh.

# About the Author

Tony McKinley is a industry analyst, writer and consultant, and is a principal in Intelligent Imaging. Mr. McKinley's dedication to document imaging and recognition began in 1978, and includes five years at Kurzweil Computer Products, and six years as President and Founder of an early imaging systems integration company. Mr. McKinley's Online OCR Lab has performed testing for some of the largest vendors in the field, and recent analysis and writing clients have included Adobe, Xerox, Caere, Fujitsu, Canon, Open Text, Excalibur Technologies and ZyLab. He is a Contributing Editor for *ImagingWorld Magazine.* He lives near Valley Forge, Pennsylvania.

""