# PDF for Documents, XML, and Rich Content

Ed Chase

Worldwide Standards Engineer

Adobe Systems

April 2006

**Adobe**

# **Part 1**
## - The Origins of PDF

Adobe

August 1990

### The Camelot Project
### J. Warnock

This document describes the base technology and ideas behind the project named "Camelot." This project's goal is to solve a fundamental problem that confronts today's companies. The problem is concerned with our ability to communicate visual material between different computer applications and systems. The specific problem is that most programs print to a wide range of printers, but there is no universal way to communicate and view this printed information electronically. The popularity of FAX machines has given us a way to send images around to produce remote paper, but the lack of quality, the high communication bandwidth and the device specific nature of FAX has made the solution less than desirable. What industries badly need is a universal way to communicate documents across a wide variety of machine configurations, operating systems and communication networks. These documents should be viewable on any display and should be printable on any modern printers. If this problem can be solved, then the fundamental way people work will change.

The invention of the PostScript language has gone a long way to solving this problem. PostScript is a device independent page

Adobe

# Adobe & PostScript

- **Original ideas date back to 1976**
  - John Warnock's concepts from a 3D graphics system invented to track ships in NY harbor
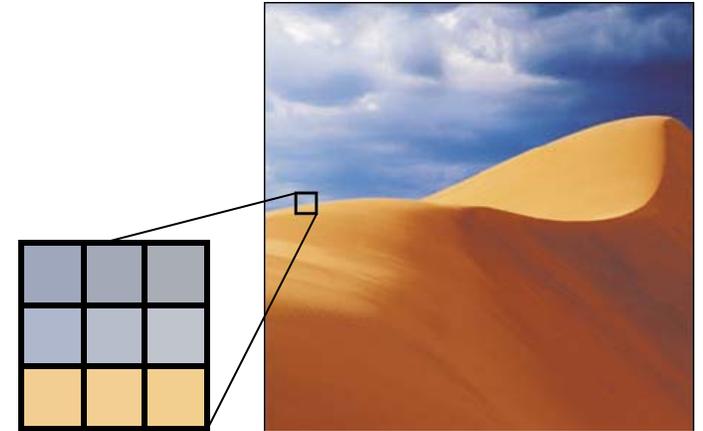  - Adobe founded in 1982 by John Warnock & Chuck Geschke

- **PostScript (1985)**
  - Page description language
  - Replaced ASCII text with structured layout and precise font handling & *hinting*
  - Actually a programming language – running the program interprets the content
  - Includes flow commands & loops
  - Had to be processed by a RIP *(Raster Image Processor)* for display & print
  - Off-loaded processor intensive print jobs to printer
    - *The Apple LaserWriter was a more powerful computer than the Mac desktop PC*
  - Large files, external resources (fonts & images)
  - *EPS (Encapsulated PostScript)* includes a rendered bitmap for easier viewing
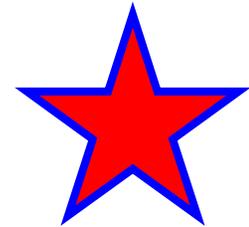
**Adobe**

# Imaging

- ## Bitmap (Raster) Graphics

  - Rows of pixels at different levels of color (or gray)

  - Photos & web images TIFF,GIF, JPEG, PNG

  - Resolution dependant

  - The larger the image, the larger the file

- ## Vector Graphics

  - X & Y coordinates, fixed colors & gradients, fonts, Bezier curves

  - PostScript, WMF, SWF, SVG

  - Resolution independent - decided by display/printer

  - Usually more efficient/smaller than bitmap

<polygon fill="red" stroke="blue" stroke-width="10" points="350,75 379,161 469,161 397,215 423,301 350,250 277,301 303,215 231,161 321,161" />
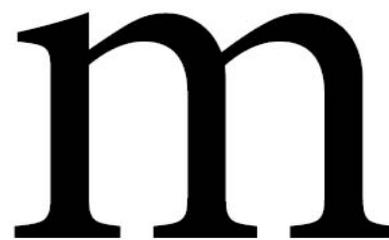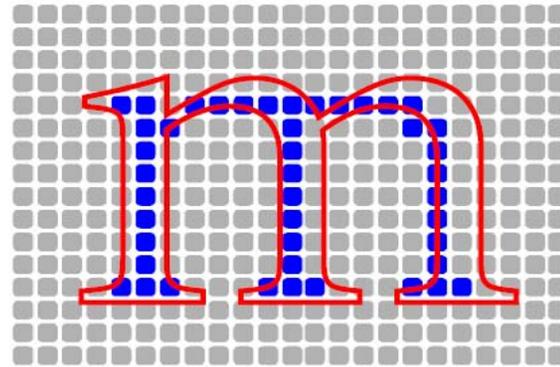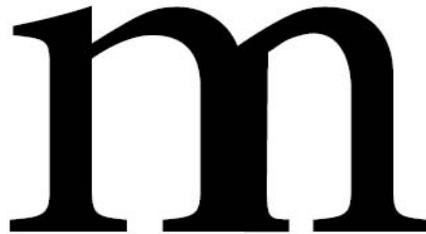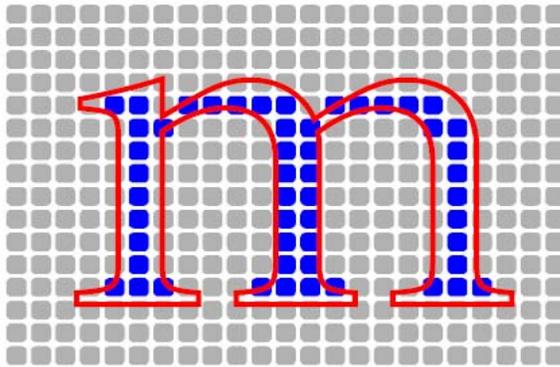
# Precise Typography

- ## Outline Fonts

  - Vector-based for efficiency at multiple sizes and consistent rendering

  - Scalable & rotate-able

- ## Anomalies

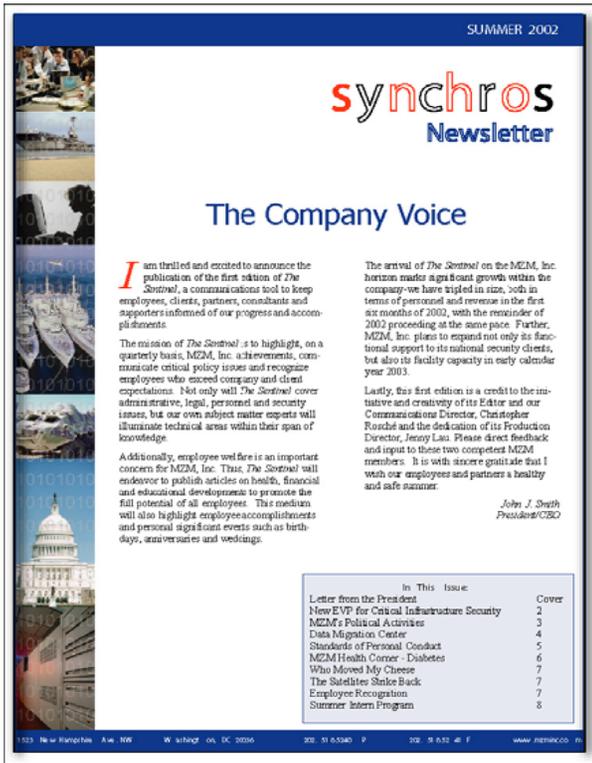  - Uneven stem widths

  - Hints and/or anti-aliasing

Adobe

# The Portable Document Format – PDF (1992)

- View & print rich content from any application across platforms & operating systems

- Simplified "object-oriented" PostScript for page layout and vector images

- Supports embedded bitmaps and image compression

- Rich font support with embedding and subsetting

- Single file packaging and structure format with compression

- Less computation-intensive, faster to open - interpreted results of PS code

# Composite Documents
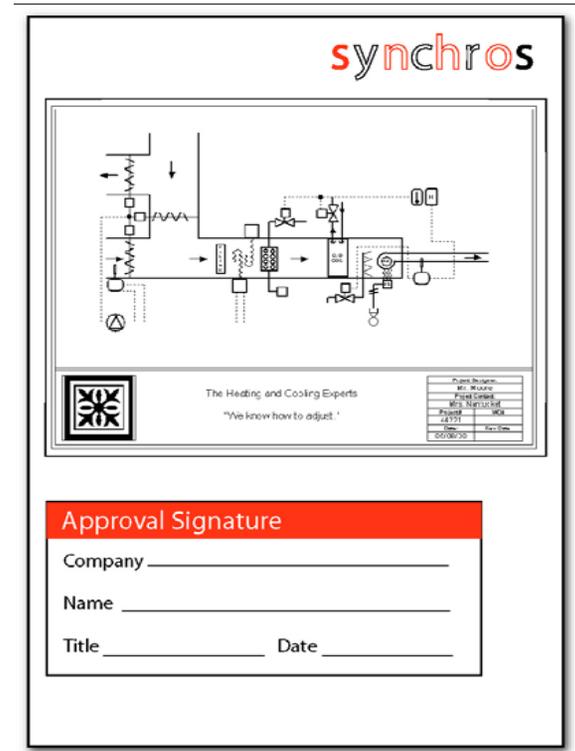


Pages

Images

Graphics

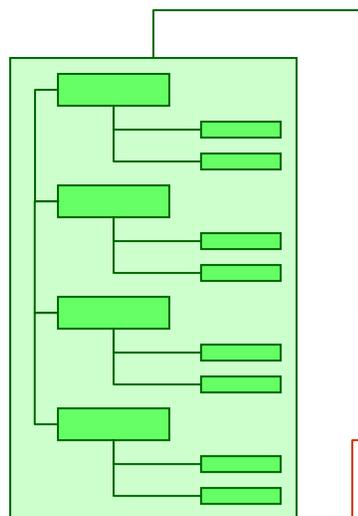Fonts

Colorspaces

Metadata

Annotations

Links
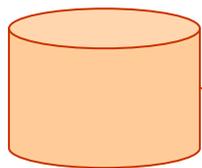
Digital signatures

# Inside of a PDF

**Logical Structure**

**Pages**

**Text, Images, Fonts**

**Metadata**

**XMP**

**PDF** Adobe

**Comments**

**Attachments**

**Signatures**

**Cross-Reference Table**

**Incremental Updates**

**Page 1**

**Page 2**

**Page 3**

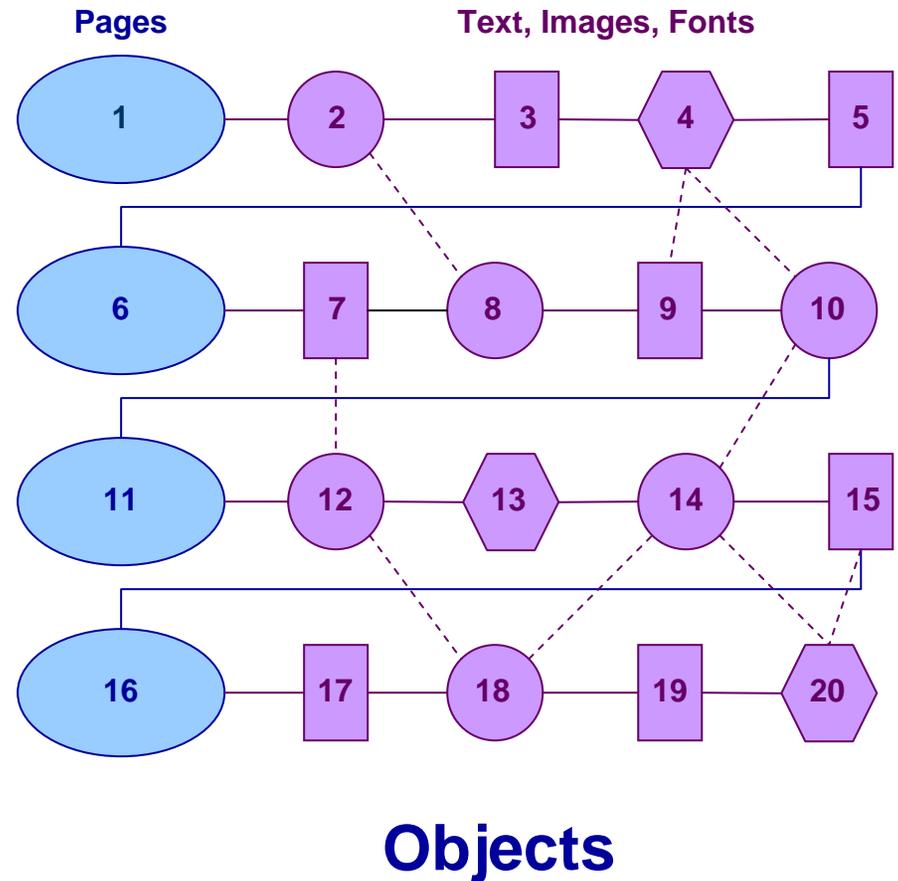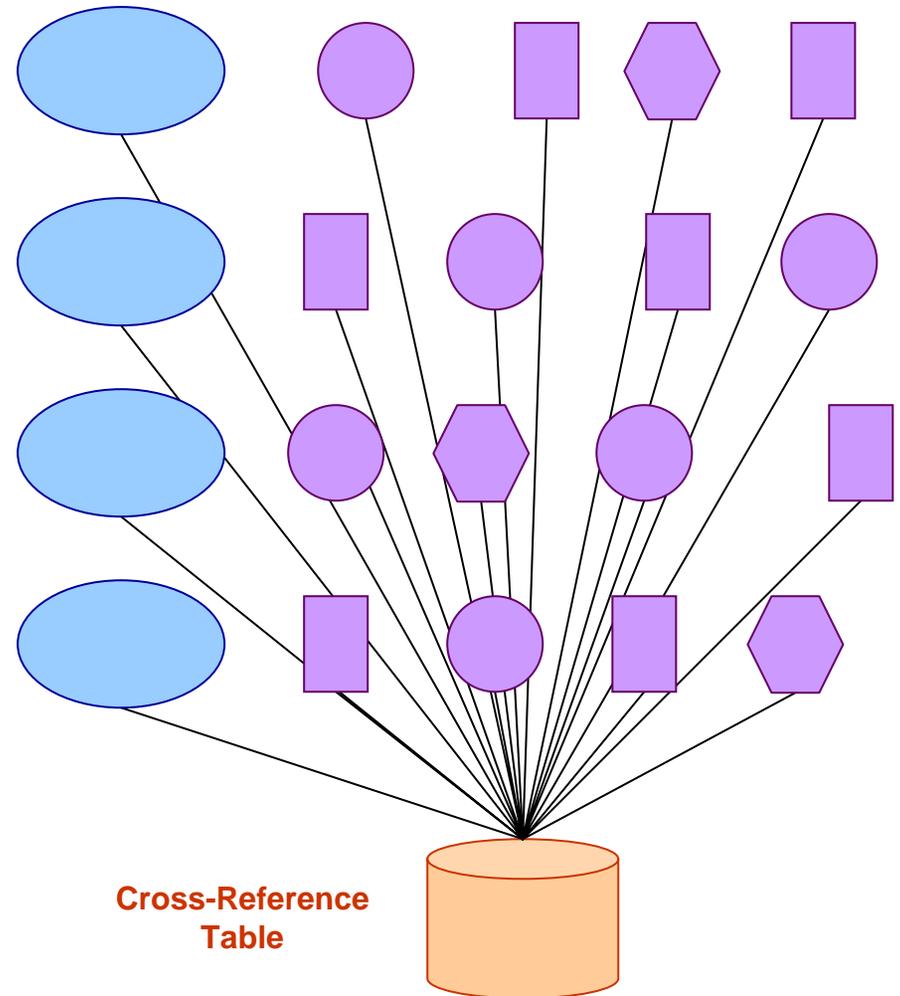**Page 4**

**Adobe**

# Objects

- PDF files are made from "objects"

- Objects are sequential, but can occur in any order in a file

  - Single-pass file generation

- Objects can refer to each other by number

- References can create a cross-linked set of objects (mathematical graph)

- Cross reference table maps object numbers to places within the file

**Pages**   **Text, Images, Fonts**

**Objects**

**Adobe**

# The Cross Reference Table

- No logical order to objects
  - Arbitrary references among objects – fonts, images, colorspaces
  - Processing follows references not sequence

- Cross-reference table provides exact locations of each object

- Allows fast and efficient *random access* to objects without having to parse entire file

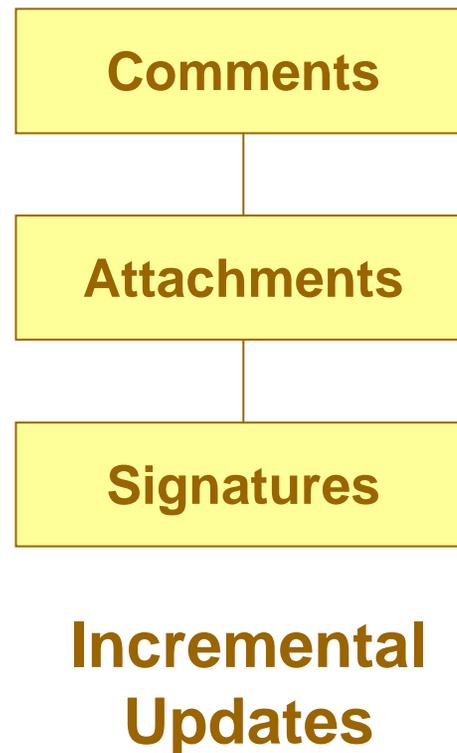- Located at end of file to allows single-pass generation of PDFs

**Cross-Reference Table**

**Adobe**

# Annotations & Incremental Updates

- ## Annotations

  - Comments, notes, highlights, sounds, movies, attachments, stamps, signatures

  - Interactive

  - May use external viewers & codecs

- ## Incremental Updates

  - Modifications are added to the end of the file

  - Cross-reference table is updated

  - Allows for versioning and detecting modifications between signatures

**Comments**

**Attachments**

**Signatures**

**Incremental Updates**

**Adobe**

# Metadata & Logical Structure

- Metadata

  - Descriptive information about a document

    - Subject, author, keywords

  - XMP – eXtensible Metadata Platform

  - Based on XML RDF

    - Dublin Core + Extensible
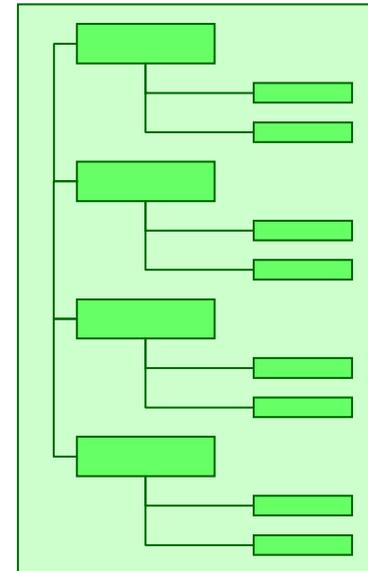
    - Published spec & Open Source tools

- Logical Structure

  - Tagged PDF

  - Presents contents in logical human reading order of a document

  - Re-flow (PDAs & Phones)

  - Accessibility

**Metadata**

**XMP**

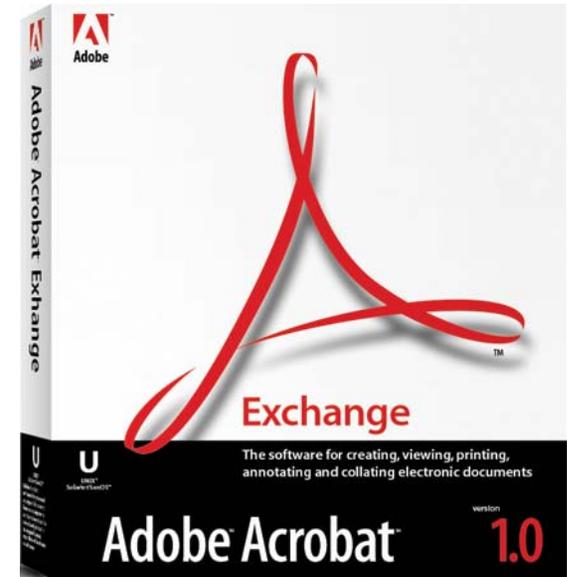**Logical Structure**

**Adobe**

# Part 2
\- 15 Years of PDF

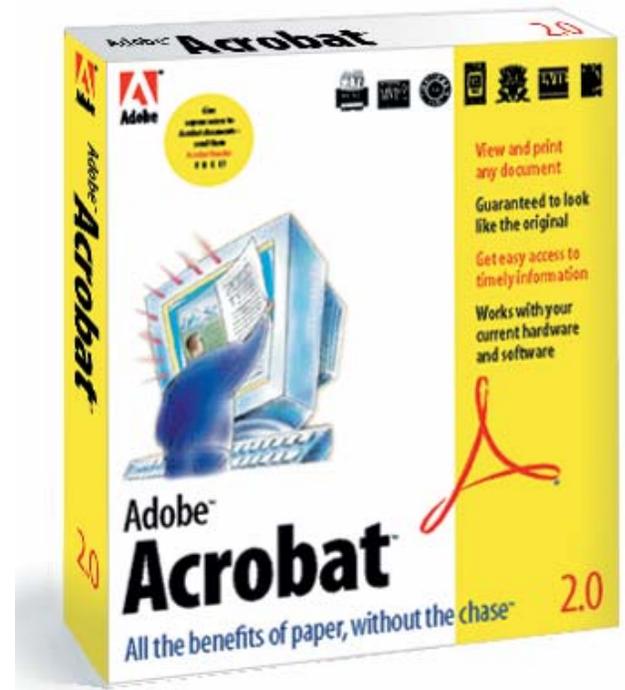**Adobe**

# PDF 1.0 - 1993

- PDF 1.0 Reference  published
  - 214 pages

- View and print anywhere

- Simple find, link, annotations features

- Files were ASCII85 encoded

- Solved the font problem
  - Embedding & subsetting

- Acrobat 1.0 – June, 1993
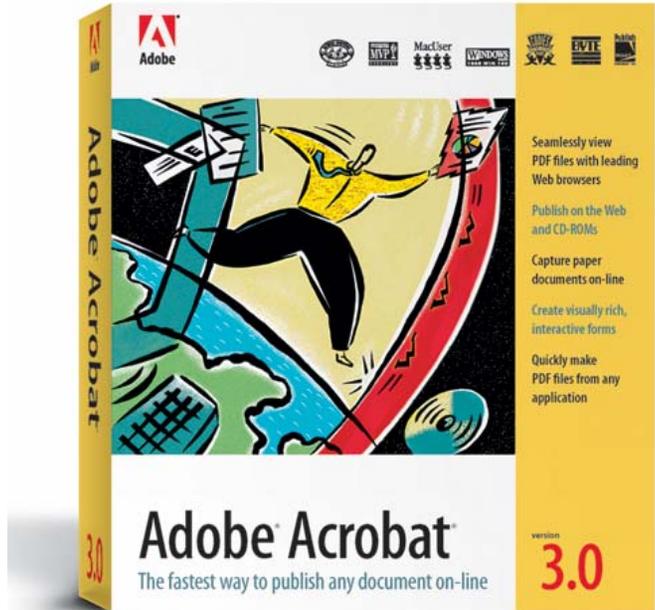  - Reader was $50
  - Creation tools ~ $800

# PDF 1.1 - 1994

- **Password Security**

- **Article threads**

- **Actions (inc. external links)**

- **No longer required ASCII85**

- **Binary bits at start of file**
  - Transport issues with ASCII

- **Acrobat 2.0 – November, 1994**
  - MS-DOS, UNIS, OS2 versions
  - Plug-ins for developers
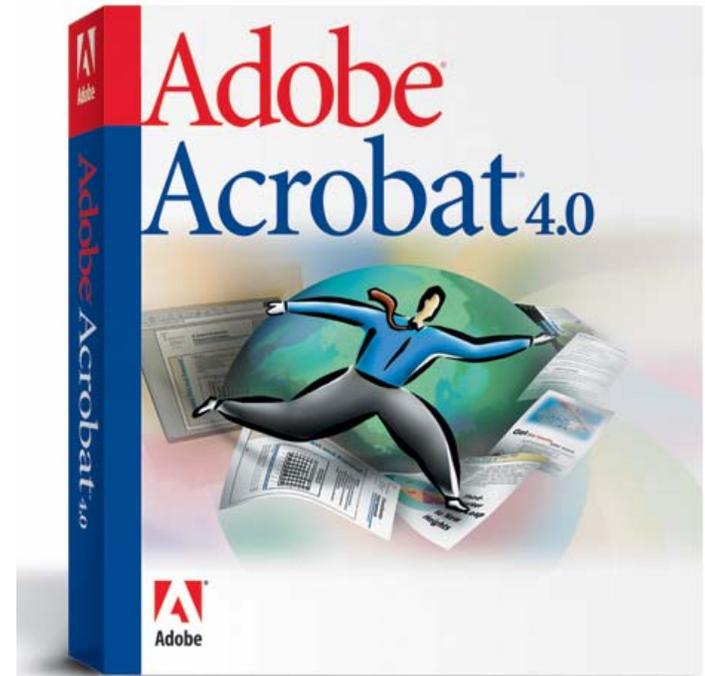  - Reader is free (still fits on a floppy disk)

# PDF 1.2 - 1996

- The Acrobat "Internet" release
  - Browser window, byte serving, fast first page
- Support for CJK in March, 1997 (v. 3.5)
- Prof. publisher features (PS Level II)
  - Spot colors, halftones, OPI
  - Undercolor removal, overprint
- Simple fill-in forms
- Start of an event model
  - On open, buttons, mouse over, etc.
- Support for external movie files

- Acrobat 3.0 – November, 1996
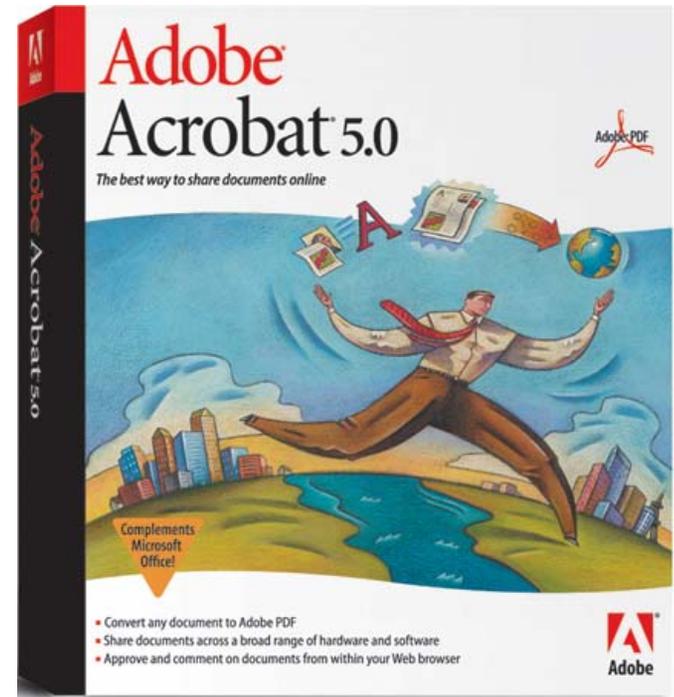  - First set of specialized features for the "Professional Publisher"

# PDF 1.3 - 1999

- Collaboration

  - Sticky notes, highlights, round-trip comments

- File attachments

- PostScript III parity

- Digital signatures

- Programmable (JavaScript) forms

- Logical structure

- First eBooks published


- Acrobat 4.0 – April, 1999

  - Dropped the name "Exchange"
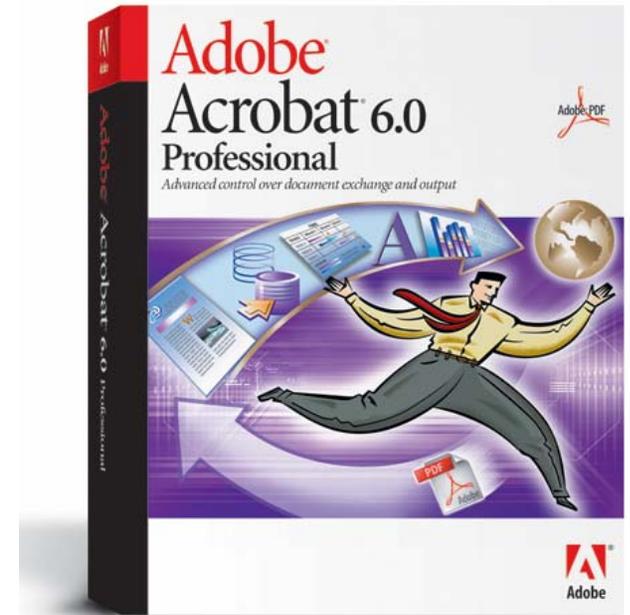
  - Self-sign signatures

# PDF 1.4 - 2001

- Transparency

- Tagged PDF

- Accessibility

- JBIG2

- XMP metadata

- Reader Extensions

- Forms connected to backend databases
    - FDF

- *ISO PDF/X based on PDF 1.3 approved in 2001*


- Acrobat 5.0 – April, 2001
    - Easy PDF Creation with PDFMakers
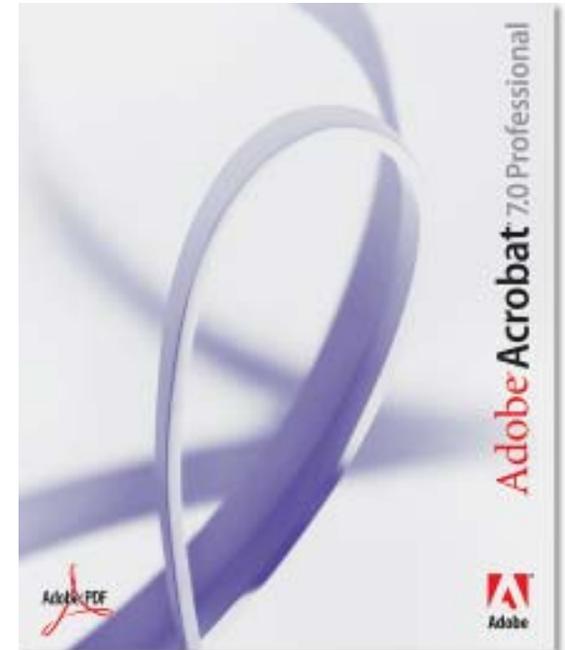    - Palm, PocketPC, and Nokia phone PDF Readers

# PDF 1.5 - 2003

- XML (XFA) Forms

- Embedded multimedia

- Better file compression (object streams)

- Layers (optional content)

- JPEG2000

- Larger than 2 Gb files

- Certificate security (signatures & encryption)



- Acrobat 6.0 –May, 2003

    - Acrobat is split into "Standard" and "Professional", and "Adobe Reader"

    - First set of features for AEC

    - eBook Reader

# PDF 1.6 - 2004

- Embedded 3D (U3D)

- Larger than 200" pages

- OpenType font support

- Various features for AEC

  - Measure, Zoom, Object Data (CAD, Visio)

- PDF 1.6 Reference Manual

  - 1213 pages

- *ISO PDF/A based on PDF 1.4 approved in 2005*

- Acrobat 7 – December, 2004

  - Launch time

  - XML Form Designer

  - Organizer

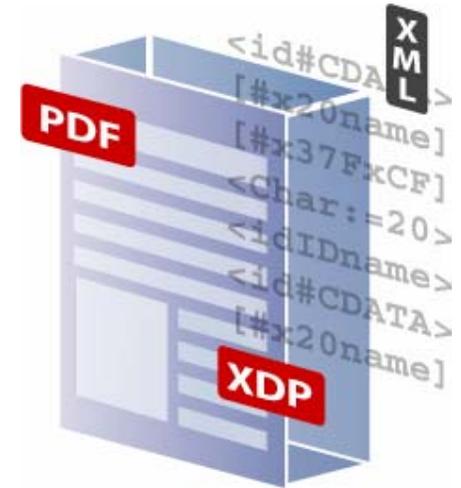  - More AEC and Pro Publish tools

  - Acrobat 3D



Adobe Acrobat 7.0 Professional
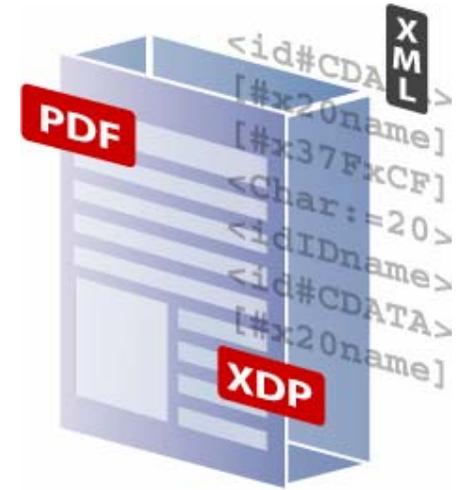
# Part 3
- PDF and XML

# PDF & XML

- ## XML Use with PDF

  - XML + XSL-FO for PDF generation

  - XMP: RDF XML

  - Export structure (tags) as XML

  - XFA & XDP – Adobe Intelligent Documents

- ## PDF or XML for Documents?

  - Difficult to produce the same flexibility and precision for compound documents with current generation of XML formats

  - Efficiency, file size & generation

  - Tagging, structure, metadata in PDF provide XML equivalents and export

  - Random access

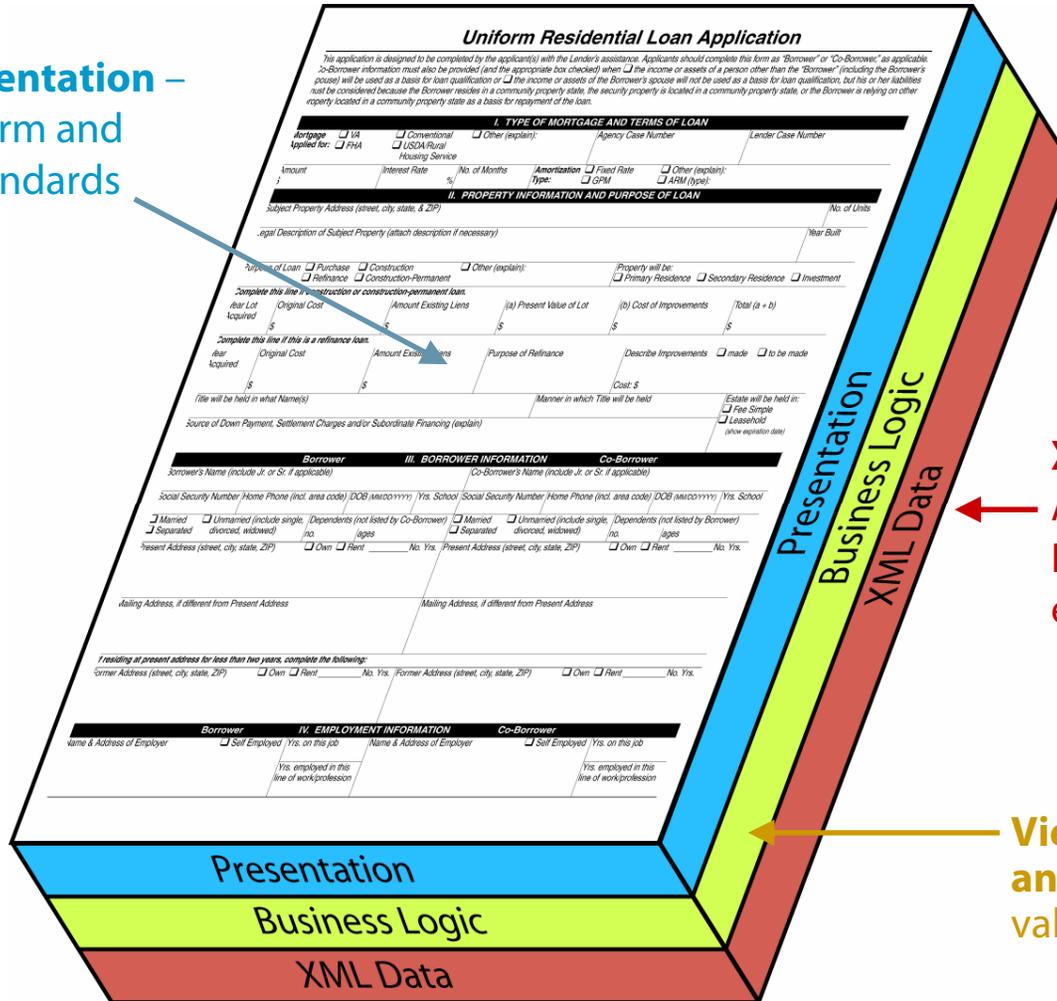  - Standards – PDF/X, PDF/A, PDF/E, PDF/UA

  - Ubiquity

# PDF & XML Part 2 – Intelligent Documents

- *Human-readable PDF + Machine-readable XML data*

- Combines PDF with XML industry standards for interoperability

- Simplifies development and deployment of standards-based applications

- PDF as a platform for industry XML standards

  - PDF can be the bridge across industry and LOB domains

- Secure container that can also leverage industry security and authenticity standards

- *Not an XML version of PDF, but uses XML as the template, data, & container*

# PDF as a Platform for XML Standards



**PDF Presentation** – precise form and layout standards

**XML Data Standards** – ACORD, MISMO, OAGIS, PISCES, RosettaNET, UBL, eDocs, XBRL

**View and Data Mapping and Validation** – interactive validation and rendering

# PDF as a Platform for XML Standards - Examples

- MISMO, PISCES, PRIA, REPI
  - Electronic mortgages & property standards

- ACORD
  - Insurance standards

- UBL, OAGIS, UN eDocs
  - ebXML CCTS-based b2b standards

- RosettaNet
  - Manufacturing

- XBRL
  - Financial reporting

- Paperless electronic mortgages - MISMO eMortgages
  - Diverse mortgage banking ecosystem with many different participants
  - Legally significant documents – notes, disclosures, credit reports, investors
  - Digital identities and signatures
  - MISMO XML data standards + PDF for mortgage documentation

- International trade documentation – UN eDocs
  - Cross-border trade documents – negotiable bills of lading, purchase agreements
  - Wide range of participants – from paper to full automation
  - PDF is recommended as a platform for UN eDocs XML standards

# Part 4
- Demos

# Thank You!

*Ed Chase*
Standards Engineer
Adobe Systems Worldwide Standards
chase@adobe.com

**Adobe**