

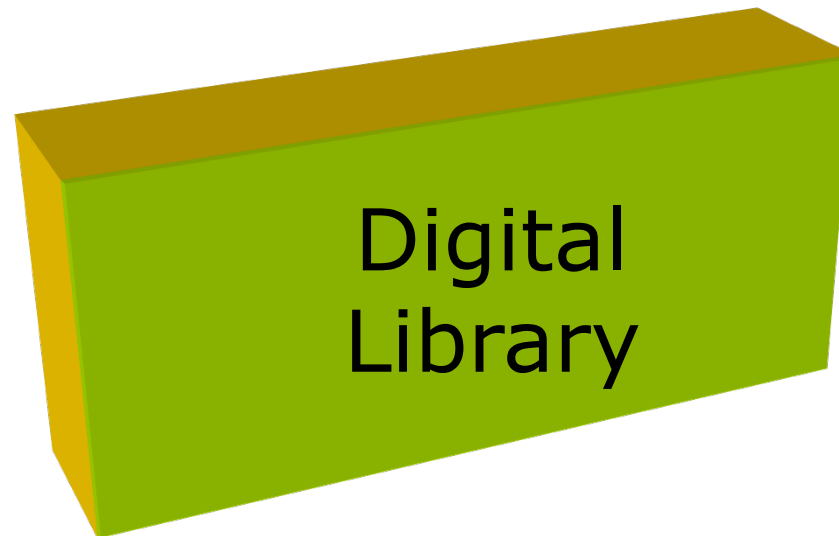


Representation and digitization of multimedia

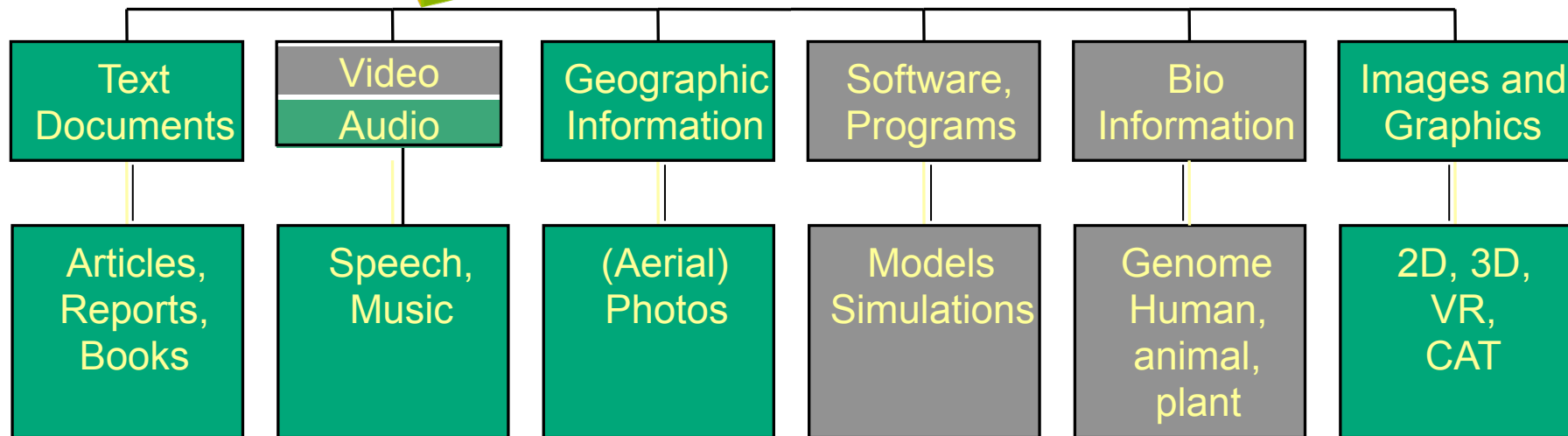
Week 2

Min-Yen KAN

Media types in the DL



Taken from Ed Fox's presentations at VaTech





Distribution of media types in the library

Library Type	LoC Gov't	NUS Acad	U Toronto Acad
Books and manuscripts	19 M	2.2M	9.1 M
Maps	4 M		278 K
Photographs	12 M	22.1 K	622 K
Music	2.7M		186 K
Motion pictures	.9 M		21 K
CD-ROM Databases		1.4K	2.1 K

Question: is the distribution of what we'd like in the digital library the same as in the automated library?

- NUS and LOC figures 2003; U Toronto, 2002
- NUS Libraries multimedia increased over 13% but only 2% for books



Outline

Representation / Digitization

- Textual images
- Images
- Audio
- Coordinated multimedia



Textual images



Cost basis for archives

	Year 1	Year 4	Year 7	Year 10
<i>Depository Library</i>				
Storage cost (per volume)	.24	.27	.30	.34
Access cost (per volume)	3.97	4.46	5.02	5.64
<i>Digital Archive</i>				
Storage cost (per volume)	2.77	1.83	1.21	.80
Access cost (per volume)	6.65	4.76	3.51	2.70

From Lesk (99), pg. 75

Digitization

- Scanning
 - Binding
 - Planetary scanner
- Resolution of scan
 - 300 dpi* for access
 - 600 or higher for archival copy

*** - Depends the smallest point size you need to resolve**



A high speed scanner,
1.2 pps at 200 dpi



Planetary or
"bookeye" scanner



Digitization

○ Purpose:

● Archival

- Quality
- Stability in the long term


● Accessibility

- Delivery
- Editing
- Annotation

1. Initiate the digitalization project
2. Establish start-up costs and secure funding
3. Prepare a detailed project plan include milestones and deliverables
4. Assess and select materials for digitization
5. Digitize materials (prepare source materials, digitize, check quality)
6. Post-process digital materials: edit, OCR, store, catalog and index
7. Deliver and make materials accessible
8. Support and maintenance of materials

-- From Chowdhury and Chowdhury (03)

Document capture costs in USD (ca. 1999)



	Preparation	Scanning	Post-scan Processing	Total (3 years)
Capital	Tables, jogger, \$1,500	Mid-volume scanner plus PC, \$25,000	Two PCs, printer, software, \$12,000	\$47,500 (11%)
Maintenance	None	8% per year \$2,000 per year	8% per year \$1,000 per year	
Labor	Two people \$40,000 per year	One person \$20,000 per year	Two people \$40,000 per year	\$300,000 (71%)
Space	120 square feet \$12,000 per year	40 square feet \$4,000 per year	100 square feet \$10,000 per year	\$78,000 (18%)
Total (3 years)	\$157,500 (37%)	\$103,000 (24%)	\$165,000 (39%)	\$425,500 (100%)

Capacity = $\sim 1,000$ page per hour \times 6.5 hours \times 250 days \times 3 years = 4.8 M.

Cost per page is $\$425,500 / 4,875,000 = \0.09 (8.7 cents)

Images of text

You've scanned in an image like this...

What to do with it?

How would we like to store and access this information?

Model	Queries allowed
Boolean	word, set operations
Vector	words
Probabilistic	words
BBN	words

Table 4.1 Relationship between types of queries and models.

is formed by words and by regular expressions (skipping the ability to allow ...).

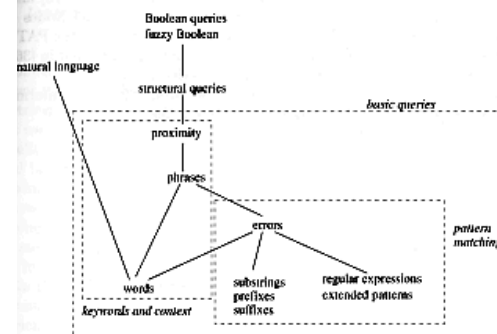


Figure 4.4 The types of queries covered and how they are structured.

The area of query languages for text databases is definitely moving towards greater flexibility. While text models are moving towards the goal of achieving better understanding of the user needs (by providing relevance feedback, for example), the query languages are allowing more and more power in the specification of the query. While extended patterns and searching allowing errors permit finding patterns without complete knowledge of what is wanted, querying on structure of the text (and not only on its content) provides greater expressiveness and increased functionality.

Another important research topic is visual query languages. Visual meta-queries can help non-experienced users to pose complex Boolean queries. Also, a visual query language can include the structure of the document. This topic is related to user interfaces and visualization and is covered in Chapter 10.



Storing a textual image

- Mostly bi-level (two-tone) until recently
 1. CCITT Fax III and IV
 - Bi-level transmission and storage standard
 - Optimized for Roman alphabet
 2. Textual image compression
 - Codebook of marks
 - A level for access and one for preservation

Horizontal run length encoding

- Always starts with white run (possibly of length zero)
- Huffman code stores a *terminating code* (TC) of all lengths shorter than 64 pixels
- Longer length encoded by 64, 128, 256 = 2^k make up code + TC

run length	color of run	
	white	black
0	00110101	0000110111
1	000111	010
2	0111	11
3	1000	10
4	1011	011
5	1100	0011
6	1110	0010
7	1111	00011
8	10011	000101
9	10100	000100
...

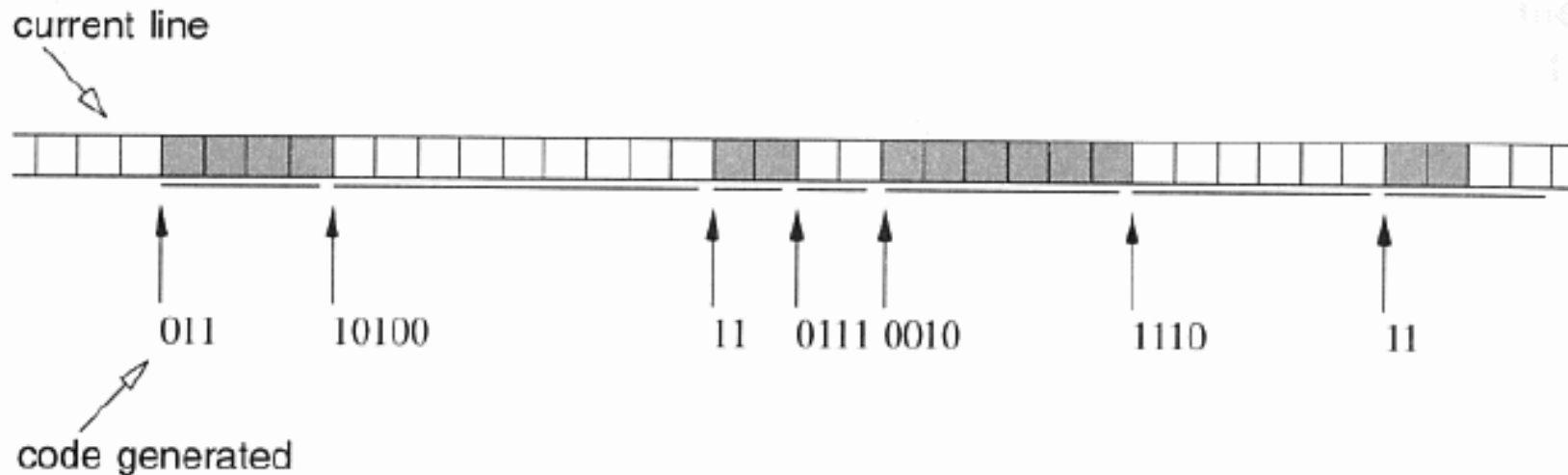
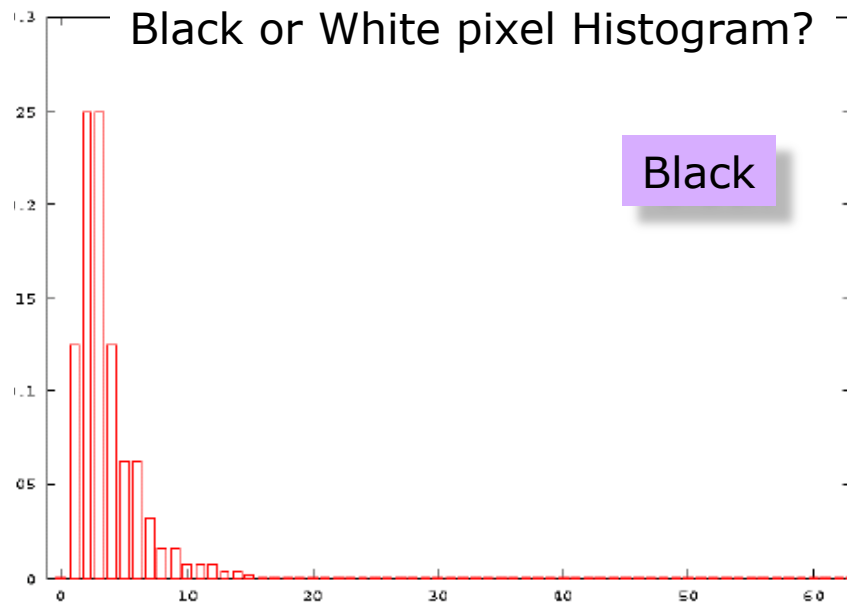
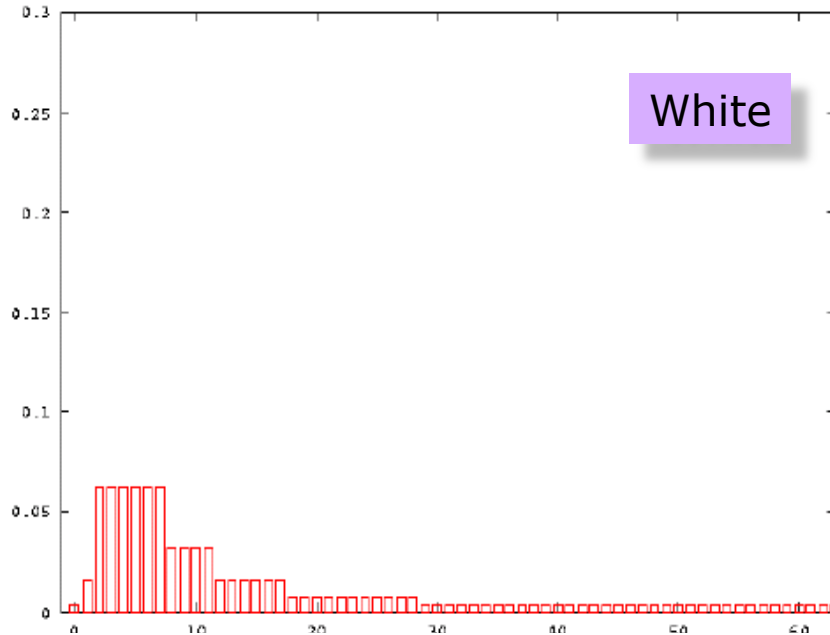


Figure 6.2 Example of one-dimensional coding.



Which histogram is which?

- CCITT Fax group III uses Huffman encoding to decide close to optimal encoding
- We show a black pixel histogram and white pixel histogram here. Which is which?
- Here's a hint:

run length	color of run	
	white	black
0	00110101	0000110111
1	000111	010
2	0111	11
3	1000	10
4	1011	011
5	1100	0011
6	1110	0010
7	1111	00011
8	10011	000101
9	10100	000100
...

CCITT fax group IV

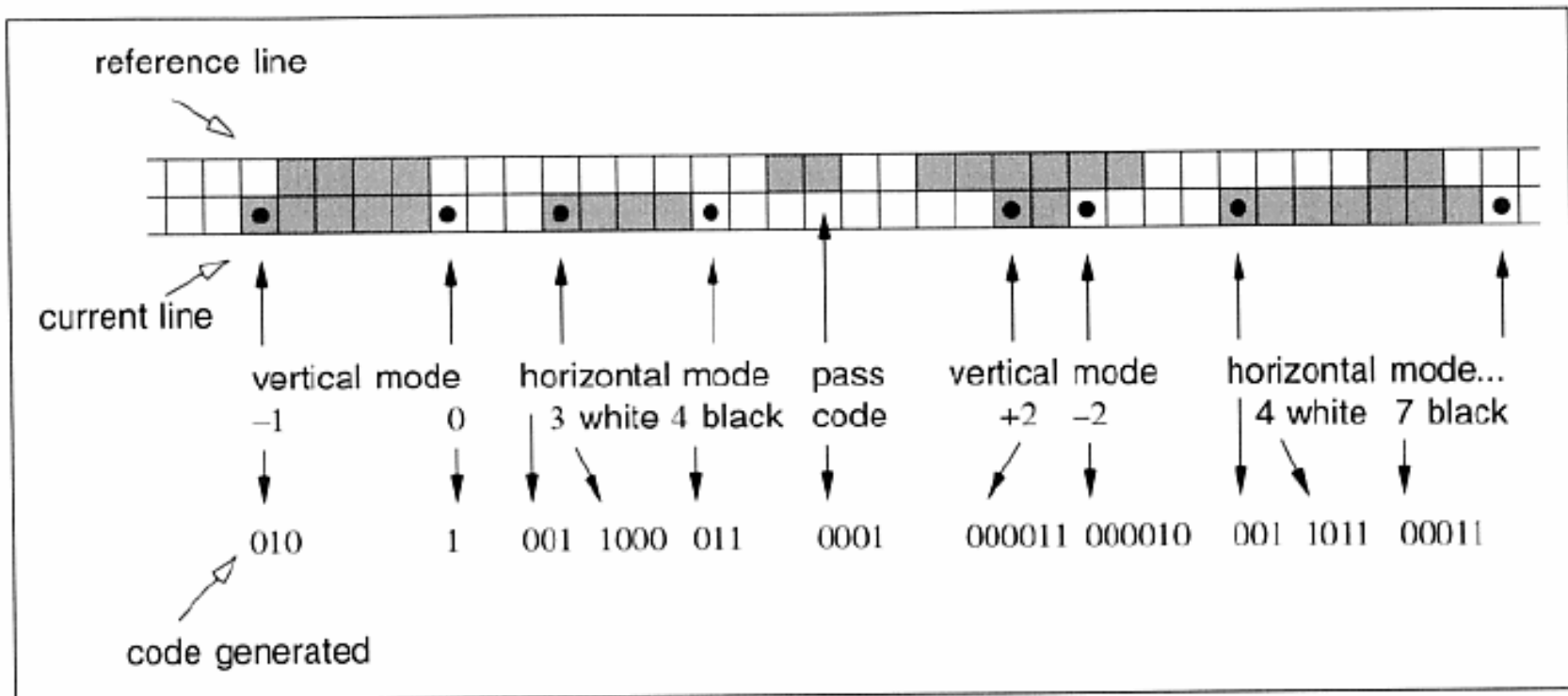


Figure 6.3 Example of two-dimensional coding.

- Takes vertical redundancy into account
- Three methods of encoding: vertical, horizontal and pass



Textual image compression

1. Find and isolate ***marks*** (connected group of black pixels)
2. Construct library of symbols
3. Identify the symbol closes to each mark and get coordinates
4. Store information
5. *Store additional information to reconstruct original image

Library

- Resolutie van de staten generael der Vereenighde Nederlanden, dienende tot antwoord op de memo-
e by de ambassadeurs van sijne majesteit van
Vranckrijck.
's Graven-hage, 1678. 4°. Fag. H. 2. 80. N°. 20.
Fa . H. 2. 85. N°. 17. Fag. H. 3. 42. N°. 4.
- ractaet van vrede gemaect tot Nimwegen op
den 10 Augusty, 1678, tusschen de ambassadeurs
van [V.] ende de ambassadeurs vande
staten generael der Vereenighde Nederlanden.

(symbol, x-offset, y-offset)

(1,50,13) (28,73,121)

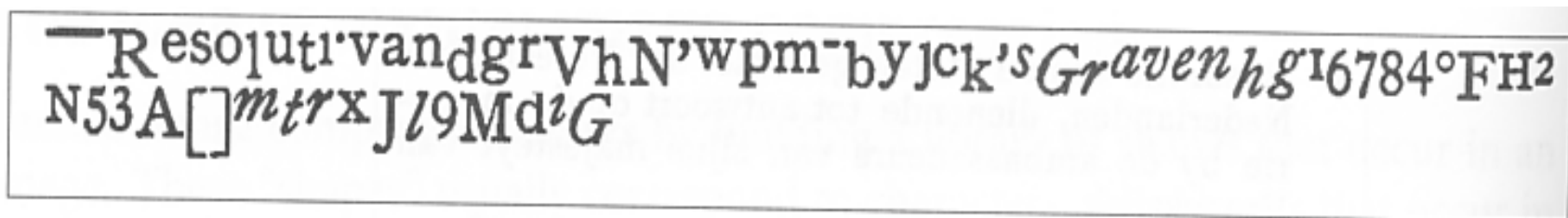


Figure 7.2 Library of symbols created from the example image.

Residue

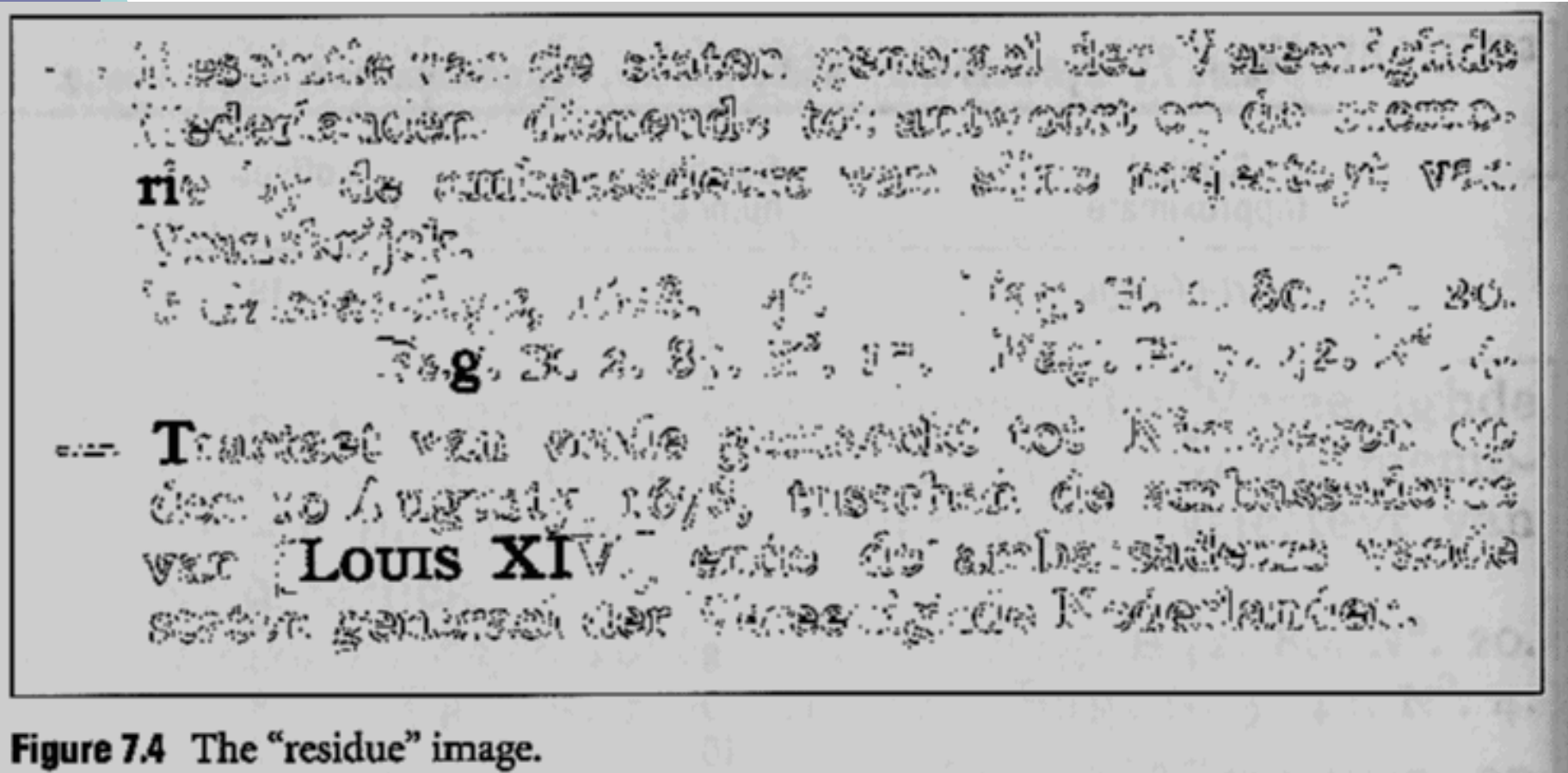


Figure 7.4 The "residue" image.



Text image outline

- Storage ✓
 - CCITT Fax Group III and IV ✓
 - Textual image compression ✓
- Access
 - De-skew
 - Segmentation
 - Media detection

De-Skew

- Projection profile
 1. Accumulate Y-axis pixel histogram
 2. Rotate to find most crisp histogram

- One of three common algorithms

Abstract
We present results of two methods for assessing the event profile of news articles as a function of verb type. The unique contribution of this research is the focus on the role of verbs, rather than nouns. Two algorithms are presented and evaluated, one of which is shown to accurately discriminate documents by type and semantic properties, i.e. the event profile. The initial method, using WordNet (Miller et al. 1990), produced multiple cross-classification of article, primarily due to the locally nature of the verb tree coupled with the sense disambiguation problem. Our second approach using English Verb Classes and Alternations (EVCAs) Levin (1993) showed that transitive categorization of the frequent verbs in WSI made it possible to identify discriminative documents. For example, our results show that articles in which communication verbs predominate tend to be opinion pieces, whereas articles with a high percentage of agreement verbs tend to be about mergers or high cases. An evaluation is performed on the results using Kendall's τ . We present convincing evidence for using verb semantic classes as a discriminant in document classification.

1 Motivation
We present techniques to characterize document type and event by using semantic classification of verbs. The intuition motivating our research is illustrated by an examination of the role of verbs in the following example:
The window was broken by the implementation by James Blair, and very valuable resources from Vantage International, Kathleen McKown and Stan the child. Partial funding for the project was provided by NSF award #9430979-RTMILITE. Consulting Contract and Synchrotron Facilities (CNSC) and Kiewit, and by the Columbia University Chase for Research on Information Access.

1.1 Focus on the Noun
Many natural language analysis systems on nouns and noun phrases in order to identify information on who, what, and where. For example, in summarization, Barzilay and Eliezer (1997) and Lin and Hovy (1997) focus on n word noun phrases. For information extraction tasks, such as the TACRED-sponsored Meeting Understanding Conference (1992), only project use verb phrases (Verney), Rappaport et al. (1993), Lin (1993). In contrast named entity tasks, which identify nouns from noun phrases, has generated numerous projects on information access.

Abstract
We present results of two methods for assessing the event profile of news articles as a function of preposition, using the framework of entailment (1983). Each category permits the formation of a wh-question, e.g. for *FROM* "did you buy?" can be answered by the "what". The wh-questions for *ACTION* (EVENT) can only be answered by verbal structures, e.g. in the question "what did you do?", where the response must be a verb *you, write, fell, etc.*

<i>FROM</i>	<i>ACTION</i>	<i>FROM</i>
<i>PLACE</i>	<i>MOVING</i>	<i>EVENT</i>
<i>AMOUNT</i>		

The distinction in the ontological category of nouns and verbs is reflected in information extraction systems. For example, given the phrases *James and US Air* that occur with particular articles, the reader will know who story is about, i.e. *James* and *US Air*. How the reader will not know the *EVENT*, i.e. happened to the *James* or to *US Air*. Did a piece rise, fall or stabilize? These are the most typically applicable to pieces, and embody the event.

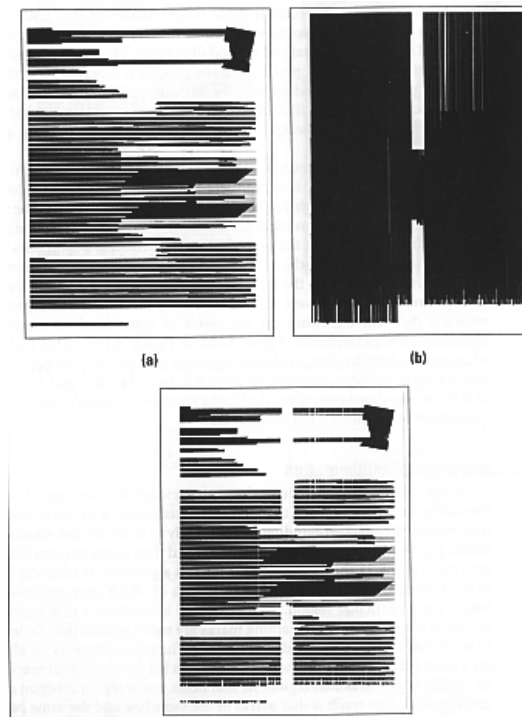
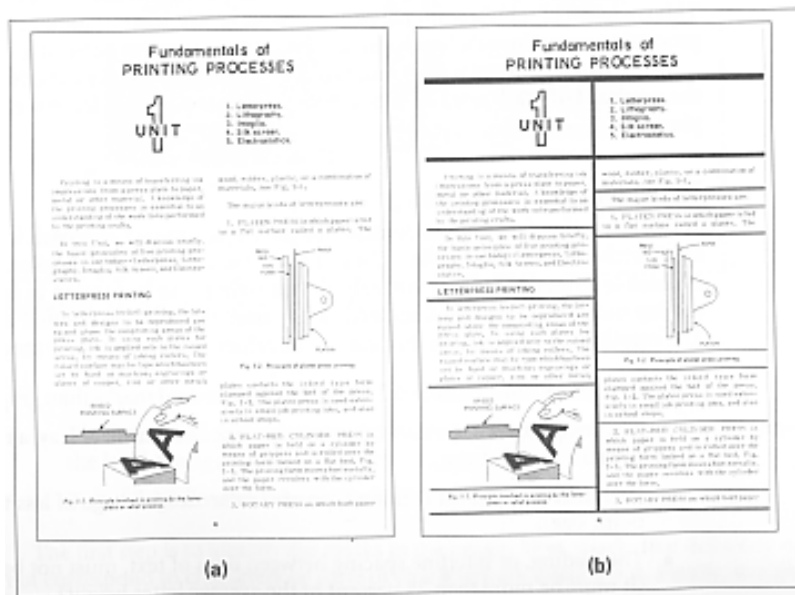
1 Motivation
We present techniques to characterize the document type and event by using semantic classification of verbs. The intuition motivating our research is illustrated by an examination of the role of verbs in the following example:
The window was broken by the implementation by James Blair, and very valuable resources from Vantage International, Kathleen McKown and Stan the child. Partial funding for the project was provided by NSF award #9430979-RTMILITE. Consulting Contract and Synchrotron Facilities (CNSC) and Kiewit, and by the Columbia University Chase for Research on Information Access.

1.1 Focus on the Noun
Many natural language analysis systems on nouns and noun phrases in order to identify information on who, what, and where. For example, in summarization, Barzilay and Eliezer (1997) and Lin and Hovy (1997) focus on n word noun phrases. For information extraction tasks, such as the TACRED-sponsored Meeting Understanding Conference (1992), only project use verb phrases (Verney), Rappaport et al. (1993), Lin (1993). In contrast named entity tasks, which identify nouns from noun phrases, has generated numerous projects on information access.

Segmentation

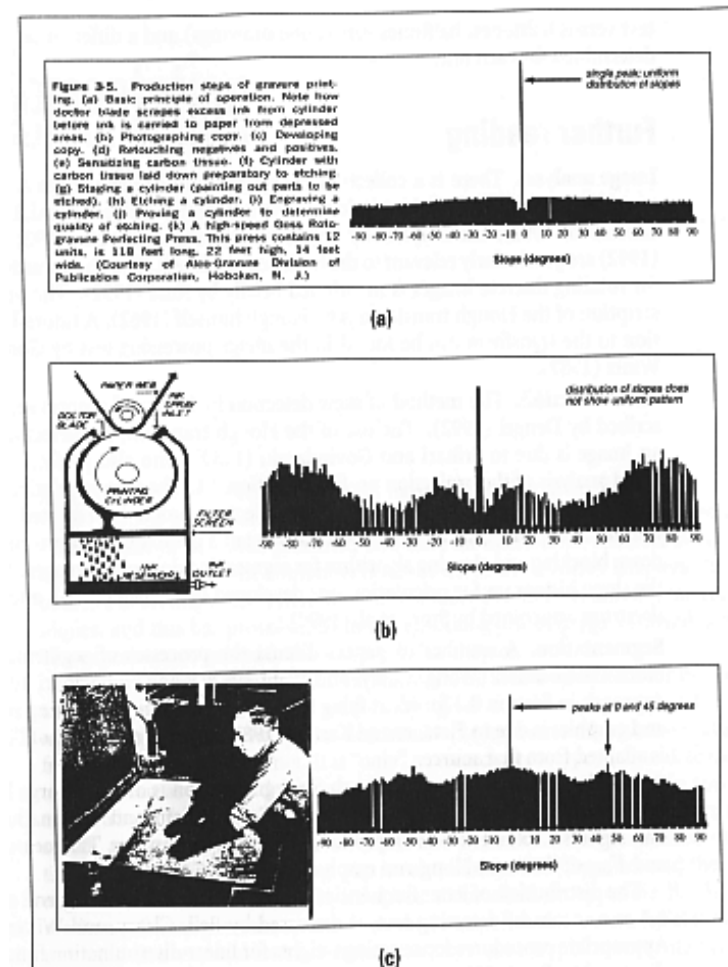
○ Top-down
(e.g., X-Y cut)

○ Bottom-up
(e.g. smearing)



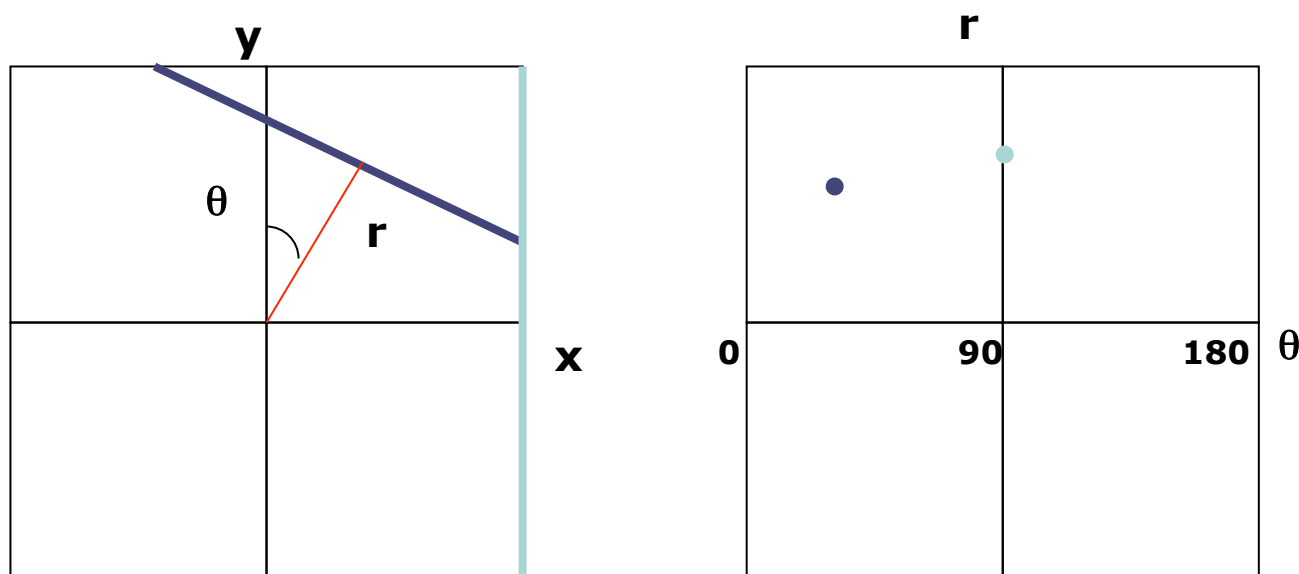
Classification

- Separate:
 - Images
 - Text
 - Line art
 - Equations
 - Tables
- One technique:
 - Slope Histogram (Hough transform)



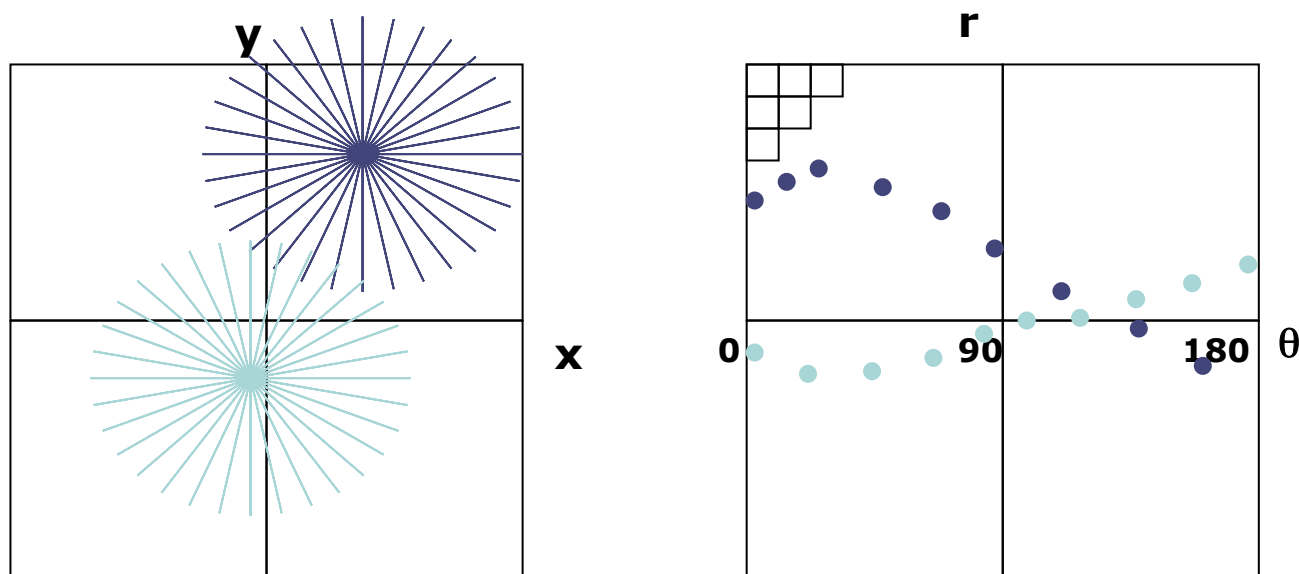
Hough Transform

- A line-to-point transform
- In practice, used to find lines in an image (e.g., set of pixels on a line)



Hough Transform

- Create virtual lines for each point
- Accumulate counts for bin in Hough space



Effective as not doing pairwise comparison



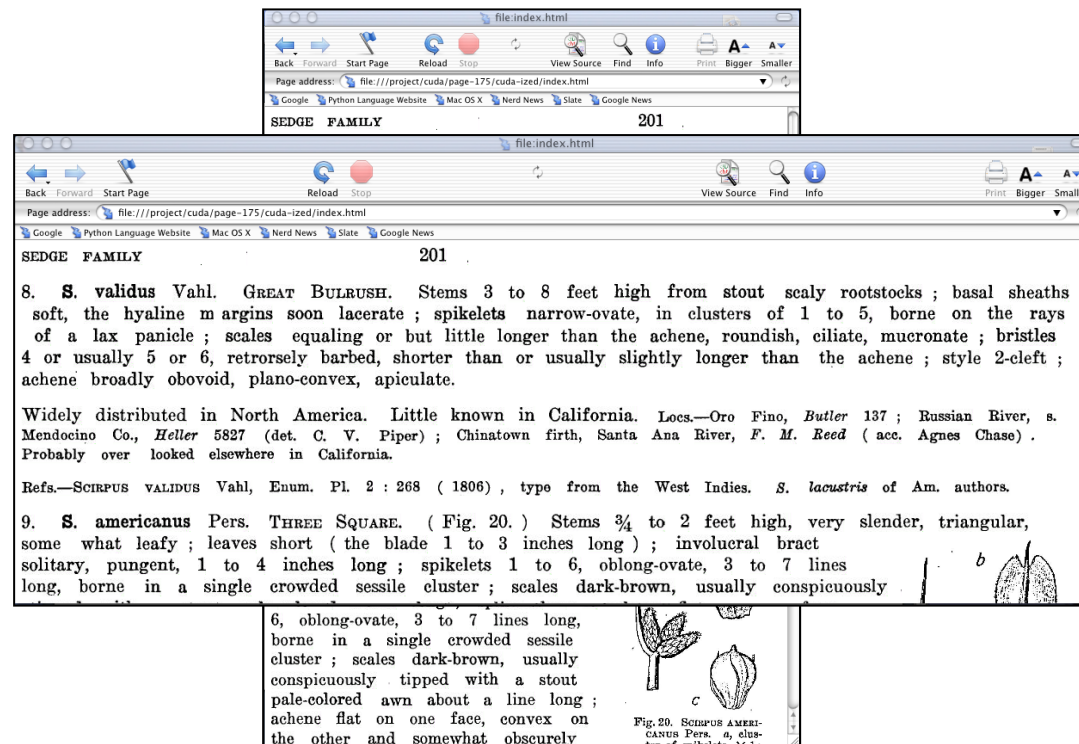
Robust Document Understanding

- OCR and document understanding are (currently) fragile technologies
 - Full scan \Rightarrow OCR \Rightarrow store pipeline makes many assumptions
 - What are some?
 - Type faces, h/w styles
 - Image qualities
 - Layout geometries
 - Writing systems
 - Languages
 - Domains of discourse

Scholarly and historical DL are much harder!

A solution (one of many)

- Courtesy Henry Baird's ICDAR 03 slides.



<http://www.cse.lehigh.edu/~baird/Talks/icdar03.ppt#21>



To think about

- How does the Hough transform save on pairwise comparisons? Can you tune the Hough transform for accuracy vs. efficiency?
- We have been discussing two-tone (b or w) images. Is dealing with color easier or harder? How do you deal with dithered images?
- Other questions about CCITT on the forum

- References:
 - Self Study module on Huffman Coding (available from Syllabus page in website)
 - Lesk (1997), Chapter 3, Images of Pages.
 - Lesk (1997), Chapter 4, Multimedia Storage and Access.
 - Witten, Moffat and Bell (1999), Chapter 6.1 - 6.2



Image data

- Raster graphics
 - As an array of pixels
- Vector graphics
 - As a collection of vectors
- Which format appropriate for which images?
 - Maps
 - Photographs
 - Line art
- For which use?
 - Fidelity?
 - Re-scaling?
 - Compression?



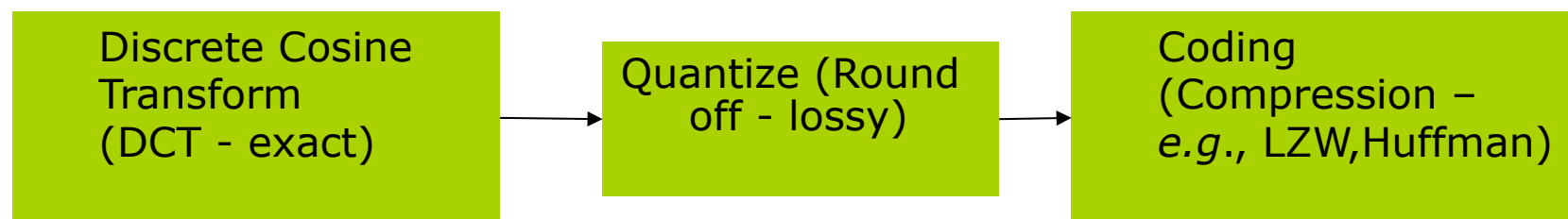
GIF / PNG

- **GIF** (‘jiff’, Graphics Interchange Format)
 - Stable, lossless color format
 - Compression achieved by:
 - 8-bit format (256 colors)
 - LZW encoding (**Unisys patent**)
 - Good for large areas of like-colored pixels.
 - Interlacing options for low-bandwidth accessibility
- **PNG** (‘ping’, Portable Network Graphics)
 - Uses public-domain variant of LZW, *gzip*
 - Up to 48 bits of color (compared to 8 in GIF)
 - Support for alpha channels (transparency) and gamma correction (white balancing)



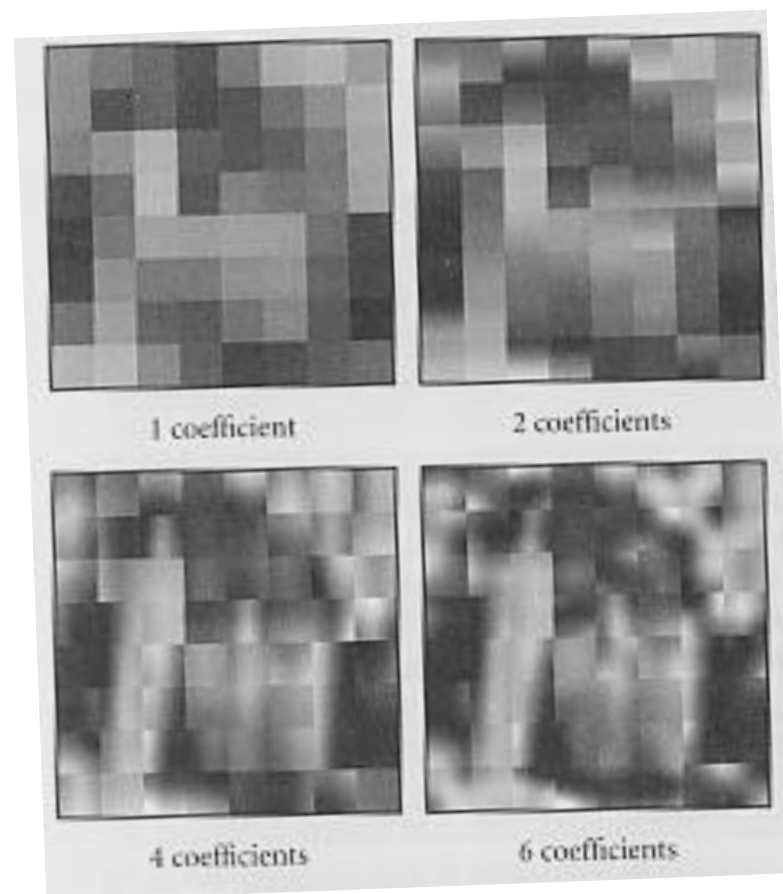
Joint Photography Experts Group

- Breaks image into 8×8 pixel blocks, each pixel 24 bits (YUV channels = 3×8 bits each)
- Compresses each block separately, without reference to neighbors



JPEG, continued

- Transform yields coefficients
- Ordered from low frequency (gradual change) to high frequency
- Gradual changes well represented
 - Good for scenery, natural images
- JPEG 2000 incorporates wavelet compression
 - Better for sharp edges





Postscript

- A programming language whose operators draw graphics on the page.
 - Text is deemed a type of graphic
 - To “draw” a page, you construct a paths used to create the image.
- A stack based, usually interpreted language
- Uses reverse polish notation



A simple Postscript example

A method to place some text down the left margin of the a page.

- You can use this after the marker for the beginning of a page.

```
gsave                                % save graphics state on stack
90 rotate                             % rotate 90 degrees
100 .55 -72 mul moveto                % go to coords 100, (.55*-72)
/Times-Roman findfont                % Get the font (set of operators) Times-Roman
10 scalefont                          % set the font size
setfont                               % Use the specified font
0.3 setgray                           % Change the color to gray
(PUT NOTE HERE) show                 % call the individual operators P,U,T ...
                                      % to draw letters
grestore                              % restore the graphics state
```




Portable Document Format

- An object database
 - Subset of Postscript, makes it faster to process
 - Can use several different compression techniques (*e.g.*, LZW and Huffman)
 - **Proprietary**
 - Has capabilities for hyperlinks



Geospatial Datasets

- Which image format is best for maps?
Hmm, let's think about it. What goes into a map?
1. Geographic information,
which provides the position and shapes of specific geographic features.
 2. Attribute information,
which provides additional non-graphic information about each feature.
 3. Display information,
which describes how the features will appear on the screen.

-- Excerpted from Geo Community, 04

Pop Quiz: Some digital maps do not contain all three types of information.

**Raster maps often do not store Attribute information,
but vector maps often do not store Display information.**

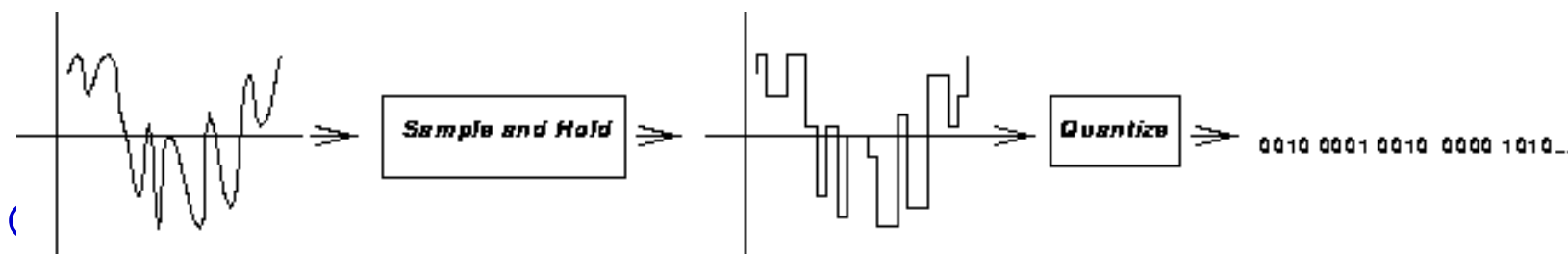


Audio

- Limit representation to what people can hear
 - Humans: \sim 20 Hz to 20 KHz
- Highest frequency (pitch) determines storage size.
 - Speech: limited range: up to 3 KHz
 - Music: full dynamic range, 20 KHz
 - Can be referred to as its *bandwidth*

Sampling

- Take continuous signal and discretize
- Higher sampling rate = better fidelity



sampling rate = $2 \times$ bandwidth

- Music: full dynamic range: $\sim 22\text{K} \times 2 = 44\text{K}$
- Speech: $4\text{K} \times 2 = 8\text{K}$

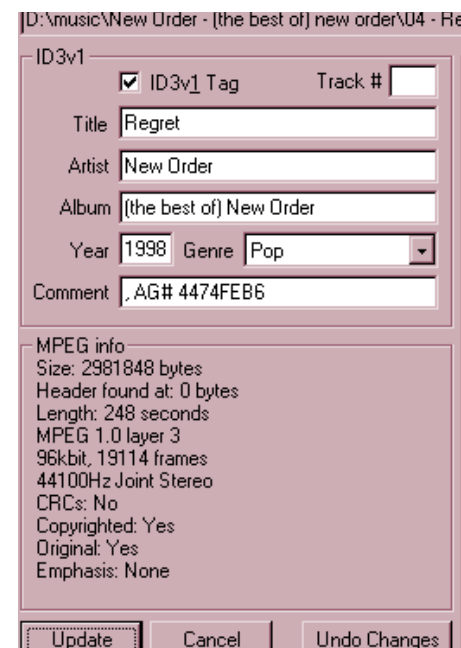


Amplitude and Channels

- Sampling at these time intervals to get *amplitude* of signal
 - a total of ~ 30 -60 dB in loudness
 - Human ear more sensitive to soft sounds
 - *Compend* amplitude (use log scale to more precisely represent low volumes)
 - 1 or 2 bytes
- For each time interval, may have to sample one or more channels
 - Differential coding (joint stereo)
 - Dolby AC 3 = 5 + 1 channels
 - Stereo = 2 channels

Storage Requirements (bitrate)

- Digital Music:
 - $44 \text{ K samples/sec} \times 16 \text{ bits/sample} \times 2 \text{ channels} = \sim 1.4 \text{ M bits/sec}$
- Digital Voice:
 - $8 \text{ K samples/sec} \times 8 \text{ bits/sample} \times 1 \text{ channel} = \sim 64 \text{ K bits/sec}$
- Analog
 - FM stereo: $40 \text{ K samples/sec} \times 8 \text{ bits/sample} \times 3 \text{ channels} = \sim 900 \text{ K bits/sec}$
 - Telephony: $\sim 6 \text{ K samples/sec} \times 2 \text{ bits/sample} \times 1 \text{ channel} = \sim 12 \text{ K bits/sec}$
- Formats
 - MP3: $\sim 1/10 \text{ compression} = 128 \text{ K bits}$
 - GSM: $\sim 1/5 \text{ compression} = 15 \text{ K bits}$

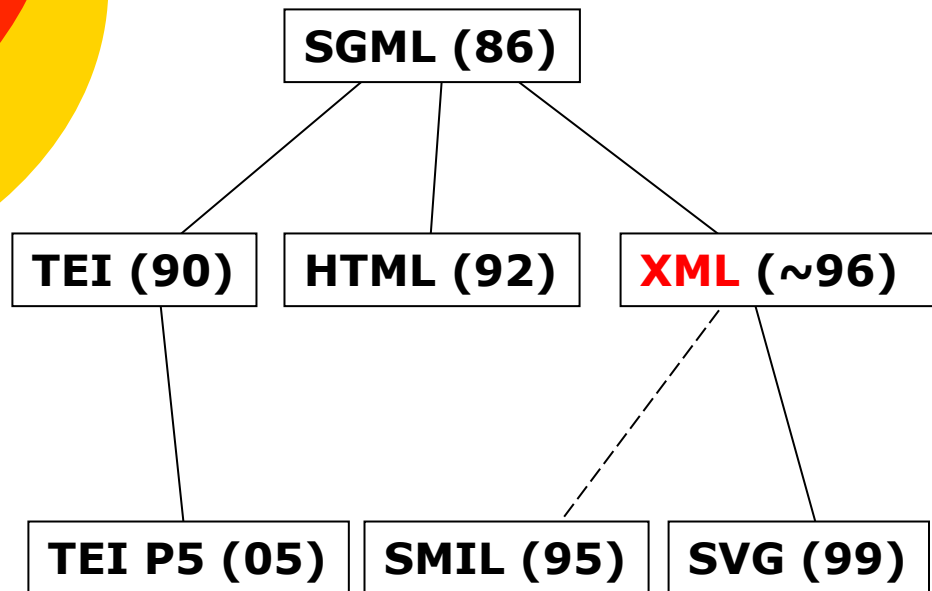




Putting media together

Have multimedia, will travel...

XML



XML says: “My family tree!”

- A basis for many other technologies
- No semantics (eXtensible, not rigid), just allows for hierarchical containment
- A meta markup language



XML, continued

- Features:
 - Separation of content from presentation
 - Content: Document Type Definition (DTD), optional
 - Presentation: CSS, XSLT
 - Enhanced hyperlinking capabilities
 - Bidirectional linking
 - Finer grained linking (XPointer)

Text Encoding Initiative

To encode knowledge “of literary and linguistic texts for online research and teaching”



- better interchange and integration of scholarly data
- support for all texts, in all languages, from all periods
- guidance for the perplexed: **what** to encode --- hence, a user-driven codification of existing best practice
- assistance for the specialist: **how** to encode --- hence, a loose framework into which unpredictable extensions can be fitted

- From the TEI Pizza talk

The “**beef**” in XML. All the semantics and none of the filling.
It’s quite filling, weighing in at 600 K words! (Think 8 kg of books)



Synchronized Multimedia Integration Language :-)

- A script for orchestrating a presentation
 - Think TV news
- Basics:
 - Define a root window
 - Layers
- Timing
 - `<par>` parallel playback
 - `<seq>` sequential playback
 - Media clips have `begin` and `end` attributes

To think about: what's the alternative format to SMIL?
How does it enhance presentation?



Summary

- Representation of knowledge
 - The more you know about the media, the faster, smaller you can transmit and store it
 - Different formats for different purposes, difference isn't superficial
- Multimedia representation
 - Trend toward accessibility, not compressibility
 - Separation of compression from format



References

- More on SMIL:
W3C's SMIL page <http://www.w3.org/AudioVideo/>
- SMIL demos:
<http://www.ludicrum.org/demos/SMILTimingForTheWeb-Demos.html>
- <http://www.geocomm.com/> and <http://www.usgs.gov> are good spots for GIS information.
- Genomic DL indexing and retrieval:
<http://goanna.cs.rmit.edu.au/~jz/fulltext/ieeekade02.pdf>
- JPEG: Pennebaker and Mitchell (93), *The JPEG Still Image Data Compression Standard*
- TEI Consortium:
<http://www.tei-c.org/>



To Think About

- The dichotomy of raster vs. vector for graphic images is carried out in other multimedia domains as well. What are their corresponding formats be in the audio and video domains?
- What about other media and presentation devices that we have missed out on? For example, small mobile devices and non-visual information devices?
- References
 - Lesk (1997), Chapter 4, Multimedia Storage and Access.
 - Witten, Moffat and Bell (1999), Chapter 6.5 (GIF/PNG section only), 6.6 (JPEG section)
 - Witten, Moffat and Bell (1999), Chapter 7, Section 1
 - Witten, Moffat and Bell (1999), Chapter 8.
 - Bainbridge, Nevill-Manning, Witten, Smith and McNab, (1999) Towards a Digital Library of Popular Music.