

Image indexing

Tomasz Neugebauer

The theoretical difficulties in indexing images include: images do not satisfy the requirements of a language, whereas textual materials do; images contain layers of meaning that can only be converted into textual language using human indexing; and images are multidisciplinary in nature, where the terms assigned are the only access points. Theoretical foundations for image indexing consist of distinctions between classes of terms including 'of' and 'about', syntactic and semantic, specific and generic, and answers to the questions who? what? where? and when? Content-based indexing can be used to generate terms for the color, texture, and basic spatial attributes of images. Image searchers use textual descriptor search terms that require human description-based indexing of the semantic attributes of images.

It is my intention to present – through the medium of photography – intuitive observations of the natural world which may have meaning to the spectators.

To the complaint, 'There are no people in these photographs,' I respond, 'There are always two people: the photographer and the viewer.'

Ansel Adams (1902–1984)

It is especially this requirement of differentiation that is not satisfied in the case of images as opposed to text:

Nonlinguistic systems differ from languages, depiction from description, the representational from the verbal, paintings from poems, primarily through the lack of differentiation – indeed through density (and consequent total absence of articulation) – of the symbol system.

(Goodman, 1976: 226)

Introduction

Jacobs argues that the importance of providing improved image access has increased as a result of the democratization of images as legitimate sources of information and education. Managed use of images 'is no longer the fairly elite domain of art historians or specialized archives' (1999: 119). The popularization of the image-rich World Wide Web as a communication and education medium certainly confirms this. Efficient image search tools by Google (<http://images.google.com>), MSN (<http://search.msn.com/images>) and Yahoo! (<http://images.search.yahoo.com/>) which index millions of freely available images online are clear evidence that image indexing and searching is common and widespread in today's visual culture.

Difficulties

The main difference and difficulty in indexing and classifying images as opposed to textual materials is that images do not satisfy the requirements of a language. Nelson Goodman defines the requirements of a language as satisfying at least 'the syntactic requirements of disjointness and differentiation' (1976: 226). Disjointness is the requirement that 'no mark may belong to more than one character' (Goodman, 1976: 133), whereas differentiation is defined as:

For every two characters K and K' and every mark m that does not actually belong to both, determination either that m does not belong to K or that m does not belong to K' is theoretically possible.

(Goodman, 1976: 136)

The disjointness and differentiation of textual language systems allows for simple decomposition into its component symbols: letters, words, paragraphs, and this represents at least a syntactic ease in indexing textual material. The semantic ambiguities inherent in indexing the meaning of a text ensure that indexing remains an art form without absolute answers, but at least the syntax lends itself to analysis. This is not the case with the syntax of images. Thus, as Jacobs observes, images 'are usually absorbed holistically by the viewer' (1999: 120).

Furthermore, as Besser points out, textual material is usually written with a clearly defined purpose that is explicitly stated, summarized and abstracted by the publishers in introductions, prefaces and book covers, whereas images are not (1990: 788). Images are 'decidedly multidisciplinary in nature: they contain a variety of features, each of which may be of potential interest to researchers from somewhat diverse fields of study' (Baxter and Anderson, 1996). For example, see Photo 1.

A photograph such as this may be of interest to a student of sculpture, but it may also be:

useful to historians wanting a snapshot of the times, to architects looking at buildings, to urban planners looking at traffic patterns or building shadows, to cultural historians looking at changes in fashion, to medical researchers looking at female smoking habits, to sociologists looking at class distinctions, or to students looking at the use of certain photographic processes or techniques.

(Besser, 1990: 788)

This requires a high number of access points for each image, depending on the audience, purpose of the database,



Photo 1

and the user characteristics. The cultural historian will not be using attributes such as composition, lighting, and perspective, whereas a student of art history might. These difficult requirements are compounded by the symbolic and allegorical meaning of images, which is highly subjective and interpretive, leading to low levels of inter-indexer consistency (Baxter and Anderson, 1996).

Presumably when indexing images for a particular database, we have a notion of who the end-users are so that we can try to anticipate the most appropriate access points for each image. However, image collections are especially

multidisciplinary, and targeting the indexing for a subset of the end-users truly excludes the others. The image is not textually expressive in itself, and the assignment of terms is a purely interpretive exercise by the indexer, which ends up being the most important access point for retrieval by the user. In the case of textual information, the user can always switch to natural language full-text searching if frustrated by exclusive descriptors, whereas this option is simply not available in image searching.

Table 1 Panofsky's levels of meaning

Level of meaning	Examples	Panofsky's term		Shatford's term
First	Factual: man, woman, tree Expressional: grief, peacefulness, happiness, anger	Pre-iconography (description) – 'primary or natural subject matter' requiring 'every-day familiarity with objects and events.'		Generic description of the objects and actions represented. The factual are descriptions <i>Of</i> , expressional represent the <i>About</i> .
		Factual	Expressional	
Second	Fat jolly man sitting in lotus position is <i>of</i> the Buddha <i>about</i> compassion.	Iconography (analysis) – 'secondary or conventional matter,' requiring knowledge of culture.		Specific objective meaning descriptions <i>Of</i> and mythical abstract or symbolic <i>About</i>
Third	Zanussi's <i>Illumination</i> (film) is an instance of cinema-verité, and about science, ethics, and enlightenment. An image of a grey brick wall is <i>about</i> socialism.	Iconology (interpretation) – 'intrinsic meaning of content,' requiring synthesis, knowledge of artistic, social and cultural setting		Not intended to be indexed except in highly specialized domains since agreement on meaning at this level is too difficult to achieve.

Source: Shatford (1986).



Photo 2

Theory

Shatford extends Panofsky's theory (1962) of three levels of meaning in a work of art to general subject analysis of pictorial work (Shatford, 1986: 43). Only the first and second levels of Panofsky's theory are to be indexed, since the third seems exemplified by the content of a subjective critical review. Jacobs interprets Panofsky's third level as requiring extensive interpretation and cultural knowledge, which is reserved for art history and similar domains (1999: 120). Panofsky's first two levels are the basis for Shatford's generic and specific levels of description. (See Table 1.)

Shatford uses Frege's distinction between sense and reference to show that 'images may be defined as referents for the sense of the words used to describe them' (Shatford, 1986: 46). A single image can be the referent for many different senses of words, or conversely and more conventionally: a picture is worth a thousand words.

The kinds of subject to be indexed in pictures correspond to answers to the four questions who? what? when? and where?, and 'each of these basic facets may then be subdivided into aspects based on *Of* in the specific sense, *Of* in the generic sense and *About*' (Shatford, 1986: 48) To illustrate let us take Photo 2.

The first level represents the literal general meaning, which requires little or no cultural knowledge: it is a photograph of staircases and about a city street in winter. The second level of meaning represents the symbolic and specific, requiring some cultural background to discern: it is

a photograph of St Urbain Street in Montreal, and about a disappearing winter beauty of iron craftsmanship.

Shatford's who, what, when and where is the basis of a faceted classification of images, divided into:

- 'who or what, beings and objects, is this picture *Of*?'
- are these symbols, representations, abstractions or personifications, in other words, 'what are they *About*?' (Shatford, 1986: 50).

The 'who' *Of* is divided into the generic and specific. Berinstein gives the examples shown in Table 2.

Table 2 Berinstein's examples

Who of specifically	Who of generally	Who about
Nancy Garman	Woman	Editor
Starship Enterprise	Spaceship	Exploring space, the future

Source: Berinstein (1999: 86).

The 'what' facet is broken down into *Of* (events and actions) and *About* (conditions and emotions). For example, Photo 3 is *Of* (general) a woman holding her hand, *Of* (specific) Rossitsa Valkanova, and *About* a film producer, support, and perseverance.

The 'when' facet includes specific linear and general cyclical time: 'specific dates and periods (June 1885,

Renaissance) and recurring time (Spring, Night)' (Shatford, 1986: 53). The *about* aspect of the 'when' facet will rarely be used: it answers the question 'Is the element of time represented in the picture a manifestation of an abstract idea? A symbol?' (Shatford, 1986: 53). Berinstein interprets this as the 'linking of ideas associated with time, such as "the end of the world"' (1999: 86). Similarly, the *About* aspect of the 'where' facet will 'only be present if the place symbolizes another place or an abstract idea (Heaven, Hell, Paradise)' (Berinstein, 1999: 86). The 'where' facet includes 'geographic, cosmographic and architectural space,' and is divided in the specific *Of* aspect that includes 'locale, site, or place represented' and generic *Of* including 'kinds of places . . . landscape, cityscape, interior, planet, jungle' (Shatford, 1986: 53).

Two separate approaches

Jørgensen identifies some broad research areas in the domain of image access closely related to each other: indexing and classification, image users and uses, and image parsing (1999: 294). In the area of indexing and classification, there is a fundamental question of the appropriateness of textual indexing for a non-textual medium.

Studies of the research literature in image indexing and retrieval confirm that there are two distinct approaches to the problem: content-based and description-based (Chu, 2001: 1,017). The content-based approach is the auto-

mated extraction of textual and image properties, whereas description-based refers to human indexing using captions and keywords for artist and image data (Chu, 2001: 1,011). Content-based techniques are primarily the domain of computer scientists developing technologies for 'direct retrieval of visual image content such as shape, texture, color, and spatial relationship' (Jørgensen, 1999: 300). In content-based retrieval systems, instead of entering search terms, for example, the user selects an area of an image in order to narrow down the search and retrieve related items (Jørgensen, 1999: 300). Alternatively, the system fetches similar items based on quantifiable aspects of the image or its parts. Although content-based research is popular, it has not produced systems that satisfy end-user information needs. This is because properties of an image such as shape, texture, and color contribute to our understanding of an image, but do not define it. Text-based search techniques remain the most efficient and accurate methods for image retrieval (Jørgensen, 1999: 302).

Indexing for retrieval

In image database retrieval systems, the most common method of indexing and retrieval is the establishment of 'fields containing bibliographic data (artist, photographer, title, date, etc.) together with a field containing some descriptive text to support the image field' (Baxter and Anderson, 1996: 68). Shatford argues that although different image



Photo 3

attributes need to be indexed depending on the type of collection and users, the general categories of biographical, subject, exemplified and relationship attributes can be used as a checklist (1994: 584). The biographical attributes include those about the history of the image's creation, such as creator attributes, creation date, as well as its 'travels . . . Where it is now, where it has been, who has owned it, how much it costs or has cost, whether it has been altered in any way' (Shatford, 1994: 584). Exemplified attributes are distinct from the subject: they describe the visual object as an instance-of: an etching, photograph, poster or an 8-bit GIF. Relationship attributes are those between different images (for example, a digital image of a class hierarchy expressed in Universal Modeling Language and an image of the Graphical User Interface that is an implementation of that hierarchy), or between images and texts (such as the list of functional requirements that were used to create the class hierarchy).

Jørgensen summarizes research in cognitive psychology, suggesting that users search for objects in images at a 'Basic Level', which is 'neither the most specific nor the most abstract level but is rather an intermediate level,' and carries out studies that confirm this (1999: 305). The layers and levels of indexing of images inevitably vary with each collection, its audience and purpose. Schroeder describes the General Motors Media Archive (GMMA) project dealing with over 3 million still photographs, and using a layered indexing approach consisting of the object (trees, sky, dirt road, etc.), style (glamour, documentary, engineering) and implication (purpose, unique qualities) (Schroeder, 1998). The implication layer is the greatest challenge to the indexer and of highest value to the searcher, as it provides the differentiating unique quality descriptions of the images, enabling an increase in precision.

The large number of terms necessary in image indexing has led to a large number of specialized thesauri-based systems. The Art and Architecture Thesaurus uses hierarchically arranged 'terminology describing physical attributes, styles and periods, agents, activities, materials and objects' (Baxter and Anderson, 1996). ICONCLASS is another classification system used extensively to index images in the domain of art history (Baxter and Anderson, 1996). Thesauri-based systems raise issues of cost-effectiveness: it is time-consuming to translate the multitude of terms into structured vocabularies that are not standardized for the entire domain of images. The subjective elements of 'aboutness' of an image, formed by the viewer during the sensory visual experience of understanding an image, require too many terms to be exhaustively indexed. This is the reason why content-based algorithmic techniques continue to be the focus of research and development.

The misconception about algorithmic techniques is that they do not require human indexers. Modeling indexing activities for a computer requires the kind of sound theoretical foundation in indexing that only professional indexers and information professionals can provide. The separation of computer scientists working towards content-based indexing from traditional indexing theorists is not ideal. It seems unlikely that mathematical analysis of images as bit-

depth digit maps alone will result in usable systems. It is the hybrid approaches that are most promising.

Besser's proposed hybrid solution of 'text-based cataloguing and indexing sufficient to allow the user to narrow a retrieval set to a reasonable size, coupled with some kind of procedure for browsing through the retrieved set of images' (1990: 790) anticipates the popular image search interfaces by Google, Yahoo! and MSN. A division into a pyramid of syntactic and semantic levels for visual descriptor terms has produced 'consistent and positive results' for image representation and retrieval (Jørgensen et al., 2001: 945). This approach satisfies Small's Aristotelian 'Principle Number One' for image indexing:

'Do not make your datum more accurate than it is.'
This principle may be rephrased as 'Preserve the Mess.' Preservation of ambiguity, however, does not mean a lack of either organization or controls.
(Small, 1991: 52)

The reasonable approach is to combine algorithmic approaches for syntactic attributes that can be extracted from images automatically, and human indexing for semantic attributes that require world knowledge. Jørgensen et al's pyramid contains the following syntactic elements that lend themselves to automatic (computer-generated) indexing: type/technique (e.g. black and white, color), spectral sensitivity (color), frequency sensitivity (texture), and image components such as dot, line, tone, spatial layout of elements (Jørgensen et al, 2001: 940). Automated assignment of keywords saves time and money, so it should be done as much as possible. However, humans 'mainly use higher level attributes to describe, classify and search for visual material' (Jørgensen et al, 2001: 940), and these semantic attributes cannot at this time be algorithmically assigned. The semantic attributes include generic objects (e.g. table, telephone, chair), generic scenes (e.g. city, landscape, indoor, outdoor, portrait), specific objects (e.g. Notre Dame Basilica, Bruce Lee), specific scenes (e.g. Warsaw, Plains of Abraham), and abstract objects (e.g. guilt, remorse) (Jørgensen et al, 2001: 940).

The popularity of the World Wide Web has emphasized the need to index images in a multimodal environment which lends itself especially to hybrid approaches. The MARIE-3 system, for example, uses web page layout and word syntax to extract captions of photographs (accompanying text) from the rest of the text on the page (Rowe and Frew, 1998). Extracted captions from multimodal documents can be used as an automatic content-based indexing strategy for improving search precision (Srihari, Zhang, and Aibing, 2000). With the addition of human semantic-level indexing to the system using a structured set of categories for visual descriptors (see the Jørgensen et al, 2001 pyramid), we have the kind of hybrid approach that results in usable image search applications for the Web.

Conclusion

This paper presents visual materials as distinct from the textual in that only the latter satisfy Nelson Goodman's

syntactic requirements for a language. When indexing images with textual descriptors, indexers are in fact creating linguistic interpretations of their subjective experience of the images. The end-users are similarly creative when searching for images using written language. The result of this is the need for a general theoretical basis for crosswalks from these image-experiences to linguistic descriptions, because as Paula Berinstein points out, 'If your inquiring mind and that of the cataloger don't meet, you won't find the pictures you need' (1999: 85).

The theoretical distinction between the *Of* and *About* of a picture has been used to create structured layers of visual descriptors. Syntactic attributes of images can be indexed algorithmically; however, human information seekers use more abstract terms for searching that require the indexing and image recognition abilities of the human mind. Computer scientists are continually improving algorithms for automatic content-based image analysis and indexing. The use of image captions in multimodal environments as a source of automatically generated index terms has proved to be particularly successful. However, this is nothing more than the extraction of human-generated indexing and description from multi-modal environments. End-users search for images using the kind of abstract concepts that require human processing and classification of the images. It is inevitable that in the absence of existing textual descriptions for images, the indexer will have to create these in order to provide access. A solid theoretical background in the types and classes of descriptors for the domain of images will improve inter-indexer consistency, but the interpretation of images will always be a creative process.

Acknowledgments

This is a slightly edited version of an article that originally appeared in *Photography Media Journal*, March 2005.

All the photos are by the author (Tomasz Neugebauer).

References

- Adams, A. (1988) *Ansel Adams: the most notable quotations, 1950–1988*, comp. J. B. Simpson. Boston, Mass.: Houghton Mifflin. Available at: www.bartleby.com/63/3/5803.html (accessed 18 March 2010).
- Baxter, G. and Anderson, D. (1996) Image indexing and retrieval: some problems and proposed solutions. *Internet Research* 6(4).
- Berinstein, P. (1999) Do you see what I see? Image indexing principles for the rest of us. *Online* March/April, 85–8.
- Besser, H. (1990) Visual access to visual images: the UC Berkley Image Database Project. *Library Trends* 38(4), 787–98.
- Chu, H. (2001) Research in image indexing and retrieval as reflected in the literature. *Journal of the American Society for Information Science and Technology* 52(12), 1011–18.
- Goodman, Nelson (1976) *Languages of art: an approach to a theory of symbols*. Indianapolis: Hackett.
- Jacobs, C. (1999) If a picture is worth a thousand words, then . . . *The Indexer* 21(3), 119–21.
- Jørgensen, C. (1999) Access to pictorial material: a review of current research and future prospects. *Computers and Humanities* 33, 293–318.
- Jørgensen, C., Jaimes, A., Benitez, A. B. and Chang, S. (2001) a

conceptual framework and empirical research for classifying visual descriptors. *Journal of the American Society for Information Science and Technology* 52(11), 938–47.

- Panofsky, E. (1962) *Studies in iconology: humanistic themes in the art of the Renaissance*. New York: Harper & Row.
- Rowe, N. C. and Frew, B. (1998) Automatic caption localization for photographs on World Wide Web. *Information Processing and Management* 34(1), 95–107.
- Schroeder, K. (1998) Layered indexing of images. *The Indexer* 21(1), 11–15.
- Shatford, S. (1986) Analyzing the subject of a picture: a theoretical approach. *Cataloging and Classification Quarterly* 6(3), 39–62.
- Shatford Layne, S. (1994) Some issues in the indexing of images. *Journal of the American Society for Information Science* 45(8), 583–8.
- Small, J., P. (1991) Retrieving images verbally: no more key words and other heresies. *Library Hi Tech* 9(1), 51–60.
- Srihari, R. K., Zhang, Z. and Aibing, R. (2000) Intelligent indexing and semantic retrieval of multimodal documents. *Information Retrieval* 2, 245–75.

Tomasz Neugebauer is contactable at: email: photomedia@gmail.com

indexing specialists

Our team of experts have many years' experience in undertaking and designing effective indexes

- Books, journals, manuals ●
- Multi-volume works ●
- From 100,000 to 100 lines ●
- Unfamiliar subjects ●
- Design of style and layout ●
- Short deadlines met ●

Phone +44 (0) 1273 416777
 email indexers@indexing.co.uk
 website www.indexing.co.uk

INDEXING SPECIALISTS (UK) LTD
306A Portland Road, Hove
East Sussex, BN3 5LP, UK