Stefan Auman

# Images on the Net!

The sheer volume of images published on the internet and thus available almost everywhere has reached almost inconceivable dimensions. Within just a few years the internet – or the World Wide Web, to be more precise – has developed into being the most significant medium for publishing images of all kinds. Efficient compression algorithms and the associated low storage requirements for even large corpora have paved the way for this development, as has the increased availability of broadband connections.

If we enter a deliberately non-specific search item – the single letter *s* – into the German Google image search engine[1], it generates around two billion hits[2] (November 2009), with *S-Bahnen* (city trains), *S-Bikes*, the *S-Klasse* (S-Class), the *S-Typ* (S-Type) and the *S-Serie* (S-Series), *S-Kurven* (S-curves) and *S-Budgets*, *S-ATA* and *S-Video*, first names abbreviated to their initial letter, ‹see…› items, and page references (‹page› in German is *Seite*) all appearing high up in the results list. If a newborn human were to spend just one second looking at each result, she or he would have almost reached the current pensionable age by the time they had cast even a fleeting glance at all the images displayed.

Search engines have a purely technical concept of how an *image* as such is to be identified. As a rule, images are referenced in the WWW using the ‹img›-element in hypertext markup language (HTML). In this system, it is not at all important initially whether the files labeled with the required src attribute[3] contain photographs, drawings, maps, digitalized pages of books or archive material, diagrams or any other kinds of illustrations.

The filter options developed continuously by Google over the last few months are made possible by an automated content analysis of the image files. They enable, for example, the display of search results to be limited to the types *face*, *photo*, *clipart*, *lineart* or various color characteristics. In the weighting of results from the image search, of course, other methods are applied. Although the search engine developers keep the details of their weighting algorithms as a trade secret, it is nonetheless possible to identify certain patterns. The spatial proximity of a search term to the ‹img› element in the source text of a page seems to be especially important. The same applies to the name of the image file incorporated via the *src* attribute. The other at-

1   images.google.de (03.12.09).
2   With Microsoft's search engine, Bing, it is 244 million on the same day: www.bing.com/?scope=images (03.12.09).
3   Every image embedded in a WWW page exists in a separate file. The *src* attribute indicates the location where this image file is stored.

tributes of the *‹img›* element – especially the obligatory *alt-* (alternate text) and the optional *title* attribute – are also of considerable relevance: both enable information to be provided about what is contained in the image.[4] Finally, another factor taken into account is whether the search term appears in the title, the address (URL), in a heading or in a link to the page containing the image found.

Despite these efforts to specify the content of images, the conspicuously high proportion of ‹false positives› among the search results immediately leaps to the eye when using Google's image search or that of other providers. The people search, in particular, frequently displays supposed hits that are completely absurd. This disproportionately high volume of ‹noise› is due to the weighting procedures outlined above.

For many years now, however, it has been possible to insert metadata into the image files themselves. In this way, information about the content of the image depicted can be linked and rendered more precisely using specific words and terms. The metadata are embedded directly in the header of the image file and can be accessed there not only by dedicated image processing programs[5] but also by search engines or photo community platforms such as Flickr[6] or Onexposure[7]. The special potential of metadata in the context of academic research lies in the fact that they can be used to create self-documenting and variously annotated images which can be fed into quite diverse applications.

## Metadata in image files

The idea of embedding metadata in image files originated in the early 1990s when the International Press and Telecommunications Council (IPTC) together with the Newspaper Association of America (NAA) adopted its Information Interchange Model (IIM), which came to be known as IPTC-NAA (or IPTC standard for short). The aim of this standard was to improve the electronic transfer of image files between professional photographers and news or photo agencies. In 1995 software manufacturer Adobe developed a proprietary procedure for embedding photographic metadata in the header of image files. What has since been known on the market as IPTC Photo Metadata is in reality a partial implementation of the original standard, as Adobe did not incorporate every field of the IIM. In September 2001 Adobe released its own specification, the Extensible Metadata Platform (XMP), which aims to facilitate a standardized workflow when working with images. Depending on file type, XMP metadata such as a file's IPTC equivalent can be written into the header of an image file or into a ‹sidecar file›.[8] Many image

---

4  The *longdesc* attribute, which enables detailed descriptions of content documented in a separate file to be referenced, is rarely used. On the incorporation of graphics using the *‹img›* element and on its attributes, cf. www.w3.org/TR/REC-html40/struct/objects.html#h-13.2 (03.12.09).

5  Such as Adobe Photoshop, GIMP, Corel Photo-Paint etc.

6  www.flickr.com (03.12.09).

7  1x.com (03.12.09).

8  Sidecar files have the same name as the actual image file, differing only in their extension (*.xmp*).

processing programs synchronize the relevant fields when changes are made to the metadata. In 2005 the IPTC released its ‹IPTC Core Schema for XMP› specification in order to adopt a successor to the IIM that would be more authoritative. Adobe was involved in the preparatory work for IPTC Core. Since July 2009 IPTC Core has been a component of the IPTC Photo Metadata Standards, while the XMP technology developed by Adobe handles the technical implementation aspects.[9]

Professional image archives can no longer be operated without a consistent use of metadata. Nonetheless, conventions affecting the specific design or semantics of the individual fields may well differ from agency to agency or from one provider to another. The example below documents extracts from the metadata associated with a photo of Herta Müller, winner of the 2009 Nobel Prize for Literature, which was distributed via dpa (*Deutsche Presseagentur*)[10]:

| | |
|---|---|
| Keywords: | .Kultur, .Literatur, .Nobelpreise, .Personen |
| Date Created: | 2009:10:07 |
| By-line: | Bernd Weißbrod |
| City: | Stuttgart |
| Province-State: | Baden-Württemberg |
| Country-Primary Location Code: | DEU |
| Headline: | Herta Müller |
| Credit: | picture-alliance/ dpa |
| Source: | Dpa |
| Copyright Notice: | usage worldwide, Verwendung weltweit |
| Caption-Abstract: | Die Berliner Schriftstellerin Herta Müller, aufgenommen vor einer Lesung aus ihrem neuen Werk Atemschaukel im Literaturhaus in Stuttgart am Mittwoch (07.10.2009). Herta Müller ist in den engen Favoritenkreis für den Literaturnobelpreis aufgerückt. Einen Tag vor der diesjährigen Vergabe in Stockholm wurde die in Rumänien geborene Autorin bei allen Spekulationen an vorderster Stelle mitgenannt. Auf den Ladbrokes-Wettlisten ist Müller auf den vierten Platz vorgerückt. Foto: Bernd Weißbrod dpa/lsw +++(c) dpa – Bildfunk+++ |
| Writer-Editor: | bw_dt |

9   For IIM, IPTC, XMP and IPTC Core cf. iptc.org/IPTC4XMP; for XMP see also www.adobe.com/ products/xmp (03.12.09). A helpful introduction to the issues appeared in the German computer magazine c't in 2006: Andrea Trinkwalder: Für die Ewigkeit. Metadatenstandards fürs Bildarchiv (For eternity: Metadata standards for image archives). In: *c't* 16, 2006, pp. 156–158. For the latest developments, cf. iptc.org/cms/site/index.html?channel=CH0099 (03.12.09).

10  Published in Spiegel online: www.spiegel.de/fotostrecke/fotostrecke-47571-5.html (03.12.09). The free extension *Exif Viewer* is available for the Firefox browser; this enables the metadata contained in the files to be displayed directly.

These kinds of conventions and regulatory mechanisms – such as those relating to the design of the keyword field – facilitate a flexible response to specific requirements. A range of procedures used either complementarily or exclusively is available for academic research purposes, including systematic cataloguing codes of a subject classification, a controlled vocabulary, or other procedures for indexing or classificatory subject cataloguing.

### Pilot project: ‹Visual cultures of ecological research›

These were indeed the kind of requirements posed by the project ‹Visual cultures of ecological research›, funded by the Hessen Ministry for Higher Education, Research and the Arts.[11] This project makes use of a mark-up scheme that enables terms allocated via the keywords field to be categorized as subject, personal or geographic keywords. The database tables required for the index search or extended search are also generated on the basis of this categorization. In the course of processing the metadata embedded in the image files, keyword entries are additionally examined to see whether they are part of a controlled vocabulary covered in an additional search option, the thesaurus search.

Among the other requirements of the information system to be realized as part of the project was the possibility of the images being edited by more than one user in one location. This meant having to devise a workflow that enables not only the project partners in Darmstadt and Marburg but also future partners to undertake all aspects of cataloguing work using their own locally available and established tools.

The project partners make use of Adobe's Photoshop Lightroom[12] software, an integrated workflow solution aimed in part, though not exclusively, at professional photographers. In its library module Lightroom offers an exemplary mode of support for working with cataloguing metadata: Keywords, for example, can be organized hierarchically, allocated synonyms, and exported in many different ways[13]. The allocation of keywords to one or more images is possible, as is a differentiated filtering of the total inventory of images.

### Workflow

Unlike more straightforward electronic image browsers, digital asset management (DAM) programs like Lightroom manage their metadata in their own program database, the catalog. In addition to the metadata generated automatically by digi-

---

11   bildkulturen.online.uni-marburg.de (03.12.09).

12   www.adobe.com/products/photoshoplightroom (03.12.09).

13   For example, with regard to keywords that are part of a hierarchy or a thesaurus, it is possible to specify whether the generic terms, subject headings etc. under which they appear are to be exported automatically as well.

*Fig. 1: The Lightroom library module, here in grid view. For cataloguing the user switches to a full-screen view of the individual data set*

tal cameras[14], the catalog also manages the cataloguing information, development settings for specific images and program-specific data – for example about images' association with user-defined collections[15] (Fig. 1).

The first task in the project workflow is to digitalize the images to be included. Once they have subsequently been imported, all the other data are then recorded in Lightroom. This includes both formal information – such as the original place of publication – as well as the actual content-related cataloguing. Specifically, data are stored regarding the individuals, institutions or objects pictured, the location where the image was created, and issues relating to media technology and methodology. This is followed by keying in a title for the image and a content description.

As soon as the images have finished being edited – the number of images can be determined freely – the process of exporting them for transfer onto the WWW information system is set in motion. The file format (here JPEG) is set within the

14  So-called exif data (*Exchangeable Image File Format*). These document photographic settings such as aperture, exposure, ISO setting, focal length, details about the lens and the image size, date and time taken.

15  For Lightroom's catalog concept, cf. Marc Altmann: Foto-Verwaltung. Katalogkonzept von Lightroom, Teile 1–4 (Organizing Photos. Lightroom's catalog concept, parts 1 – 4). In: *c't special* 02, 2009. pp. 128–139.

*Fig. 2: The map interface of the information system*

export settings. These also determine the storage location as well as the format of the metadata to be embedded. In order to guarantee the greatest degree of self-documentation possible, the metadata are written in parallel in both the IPTC and the XMP area of the image file header. This makes them accessible, in principle, to programs that do not yet support the more modern XMP format.

All that is left to do then is to transfer the exported image files to the WWW server. The procedure is just as straightforward for all the other project partners; whether they do their cataloguing locally using Lightroom or some other application that enables the incorporation of metadata is irrelevant. All that has to be set up are the necessary (once and for all) authorizations to upload files. As soon as the image files have been transferred to the WWW server they are available for fully automated further processing.

Metadata can be extracted from the image files using command line programs such as the free, very efficient *Exiftool*,[16] whose source code is also available. There are also dedicated collections of subroutines (libraries) for a host of programming languages, where the required function can be obtained and added in. The Perl program developed by the project team draws on several such libraries to enable various stages of processing to be combined:

---

16  owl.phy.queensu.ca/~phil/exiftool (03.12.09).

- scanning in all the image files;
- extracting the metadata;
- creating database tables for individuals, places, institutions, objects etc.;
- generating previews for the individual image files.

The automatic processing procedure can be activated as often as users see fit. In many cases it is sufficient if the automated mechanism checks once a day whether new files are contained in the upload folders on the server.

**The search functions**

The information system required as part of the visual cultures project was designed as a heuristic tool. The aim was and is to create for the object of the project – the study of visualization strategies of ecological research – differentiated search options which facilitate quite different strategies for accessing the processed images. These include:

- the *simple search* similar to the function offered by search engines;
- the *extended search*, in which various descriptors can be linked or time filters set;
- the *thesaurus search*: this facilitates a systematic and hierarchical extension or restriction of the overall image inventory with regard to media technology, object, methodology, institution, individual;
- the *index search*, which allows for more exploratory points of access via alphabetical indices; the index search is helpful not least in gaining an overview of the basic cataloguing categories, such as main headings.

In addition to this, a *map interface* was created using OpenStreetMap technologies[17] (Fig. 2).

One characteristic feature of the information system is the variable presentation of results. The grid view provides a rapid guide to the results set (Fig. 3); additional cataloguing data can be superimposed in list view. While the detailed view shows the data set in its entirety, the map view makes it possible to view the results set in such a way that all the locations identified in the hits are represented in a dynamically generated map.

As information about individuals, locations and subjects – the latter include a list of the thesaurus terms allocated to the image in question – are represented as links in the detailed view, new results sets can be formed in an ad hoc way, thus enabling relationships between images to be rendered visible or indeed established (Fig. 4).

---

17  OpenStreetMap is a wiki project aimed at producing a free map of the world, cf. www.openstreetmap.org (03.12.09).

*Fig. 3: The results set from a simple search in grid view*



*Fig. 4: Representation of a hit in detailed view; all metadata originally recorded as keywords are represented as links*

All the options provided by the information system for searching for images and displaying results are ultimately based on the metadata embedded in the images, thus demonstrating the power of this simple approach.[18] All that has been added to them are basic data concerning geographical units (locations, landscapes), as these coordinates need to be available for the map interface and map view. Finally, a literature database has been integrated into the information system, through which bibliographic details can be incorporated into the detailed view. The primary keys for the linkage are again part of the image metadata.

The information system itself was programmed using a web application framework. Such frameworks support the development of dynamic WWW applications by providing components for database access, roles and rights management, internationalization (I18N), localization (L10N) and much more. The visual cultures project makes use of the Zend Framework[19], based on the PHP programming language. It makes available very efficient and fully developed libraries for full text search (Apache Lucene Technology[20]), PDF generation, input validation, internationalization, authentification and authorization, and mail etc.

A further strength of the Zend Framework is that it supports search engine-friendly addresses (*URLs*) very well. This, along with the way the results are presented and the metadata embedded in the files themselves, makes it highly likely that the images will be found even if the information entered into search engines is not very specific – and this in the context of the information system through which all further search options are opened up. Good support for technologies in the Web 2.0 environment additionally provides ideal conditions for the future expansion of the information system.

## Works cited

Altmann, Marc: Foto-Verwaltung. Katalogkonzept von Lightroom, Teile 1–4 (Organizing Photos. Lightroom's catalog concept, parts 1–4). In: c't special 02, 2009. pp. 128–139.
Trinkwalder, Andrea: Für die Ewigkeit. Metadatenstandards fürs Bildarchiv (For eternity: Metadata standards for image archives). In: ct 16, 2006, pp. 156–158.

1x.com (03.12.09).
bildkulturen.online.uni-marburg.de (03.12.09).
framework.zend.com (03.12.09).
images.google.de (03.12.09).
iptc.org/cms/site/index.html?channel=CH0099 (03.12.09).
iptc.org/IPTC4XMP (03.12.09).

18   A minor limitation is given simply due to the fact that changes necessitate the renewed exporting of the images concerned. Although it would easily be possible to make corrections, additions etc. via the WWW (the changes could even be written back into the metadata of the images), the associated risk of discrepancies between the cataloguing system and the information system is averted by ensuring that all writing operations take place in principle in the cataloguing system.

19   framework.zend.com (03.12.09).

20   lucene.apache.org/java/docs (03.12.09).

lucene.apache.org/java/docs(03.12.09).
owl.phy.queensu.ca/~phil/exiftool (03.12.09).
www.adobe.com/products/photoshoplightroom (03.12.09).
www.adobe.com/products/xmp (03.12.09).
www.bing.com/?scope=images (03.12.09).
www.flickr.com (03.12.09).
www.openstreetmap.org (03.12.09).
www.spiegel.de/fotostrecke/fotostrecke-47571-5.html (03.12.09).
www.w3.org/TR/REC-html40/struct/objects.html#h-13.2 (03.12.09).