UniS

# Some Ideas for Modelling Image-Text Combinations

Andrew Salway and Radan Martinec*

School of Electronics and Physical Sciences

Department of Computing

CS-05-02

* University of the Arts, London.

THE QUEEN'S
ANNIVERSARY PRIZES
FOR HIGHER AND FURTHER EDUCATION

2002

UniS

University of Surrey

# Some Ideas for Modelling Image-Text Combinations

## Abstract

The combination of different media types is a defining characteristic of multimedia yet much research has concentrated on understanding the semantics of media types individually. Recently, systems have been developed to process correlated image and text data for tasks including multimedia retrieval, fusion, summarization, adaptation and generation. We argue that the further development and the more general application of such systems require a better computational understanding of image-text combinations. In particular we need to know more about the correspondence between the semantic content of images and the semantic content of texts when they are used together. This paper outlines a new area of multimedia research focused on modeling image-text combinations. Our aim is to develop a general theory to be applied in the development of multimedia systems that process correlated image and text data. Here, we propose a theoretical framework to describe how visual and textual information combine in terms of semantic relations between images and texts. Our classification of image-text relations is grounded in aesthetic and semiotic theory and has been developed with a view to automatic classification.

## 1. Introduction

The combination of different types of multimedia data is a defining characteristic of multimedia systems and applications, yet most research on the understanding and modeling of multimedia semantics has focused on understanding and modeling the semantics of single media types. It has been argued that we must address the issue of how the combination of multimedia elements can be greater than the sum of its parts [9]. On the one hand there are benefits for facilitating more effective and engaging multimedia communication. On the other hand there are benefits for solving problems that are hard in single media, but are more computationally tractable when correlated media are considered.

Combinations of visual and verbal media are commonplace throughout human communication, especially these days on the web. Consider the combination of images and texts in three kinds of websites: newspapers, art galleries and science education. In each case the image and the text are serving a different purpose in communication, and in relation to each other. A news photograph captures one salient aspect of a story, whereas the text tells the whole story in more detail; in this case the image could be said to illustrate the text. In contrast, a painting in an art gallery is the prime object of interest, and its accompanying caption points out and explains salient aspects of it; the text could be said to describe the image. Different again is the way in which images and text are sometimes used in scientific expositions. A diagram may to used to present the same information as the text but in a

different way – they are in some sense equivalent.

We define an image-text combination as an instance of verbal and visual information being presented simultaneously and we are interested in developing a computational understanding of how to 'read' an image-text combination. It is not simply a question of adding the result of text content analysis to the result of image content analysis. It seems that a key aspect of understanding an image-text combination is the way in which the image and the text relate to one another – in terms of relative importance, and in terms of how they function to convey meaning. The preceding examples show variations firstly in the relative importance or interdependence of the image and the text, and secondly, in the different kinds of information, or other value, that one adds to the other. Words like 'illustrate', 'describe' and 'equivalent' capture some of our intuitions about how images and texts relate. However, we believe that a computationally tractable framework for classifying image-text relations is required to facilitate better processing of image-text combinations in a wide range of applications.

Section 2 reviews systems that process correlated images and texts, and argues that having access to descriptions of image-text relations could improve such systems. We also review some theories about images and texts from the fields of semiotics and aesthetics. Section 3 presents the framework that we have developed to classify image-text relations: we envisage these relations as being part of multimedia content descriptions for image-text combinations. Our framework is grounded in ideas from semiotics and it has been developed with automatic classification in mind. Section 4 closes by summarizing our progress to date, and proposing future research.

## 2. Background

There is a growing interest in systems that deal with correlated media, including for multimedia adaptation / summarization / generation, cross-modal information retrieval, multimedia data fusion and hypermedia. Our review of some of these systems suggests that little attention has been paid to the ways in which different media types convey information, and to what happens when they are used in combination. It seems that such issues are either ignored, or dealt with in an ad-hoc manner suitable only to a specific application and a specific genre of data. We conclude our review with a set of questions about image-text combinations that we think are relevant for processing correlated media. We go on to look at some semiotic and aesthetic theories about images and texts to see if they might provide some answers.

With regards to the adaptation of multimedia, like setting priorities for transmission to low-bandwidth devices and small displays, it is important to recognise which pieces of media are most important for users [22]. It has been noted that web page analysis is required to understand the role of images in this context. Seven functional categories have been defined to describe how images are used on web pages, specifically news websites: story images, preview images, host images, commercial images, headings, formatting, and icons/logos. Images can be

automatically categorized to help prioritize them [8]. Such prioritization is also an issue for generating multimedia summaries, for which it is important to not only select high-priority media but to ensure an appropriate mix of media types. Previous work on multimedia generation has automatically placed images and texts together in a rhetorical structure, to generate instructional manuals [1].

The semantic gap is a well-recognized problem for image retrieval systems and it motivates the integration of information about images from a number of sources, including correlated texts [17, 21]. Indexing terms, and potentially more structured representations of semantic image content, may be extracted from correlated texts. The approach of retrieving one media via another has been called cross-modal information retrieval [14] and seems to be common practice for some web search engines that retrieve images on the basis of keywords found in HTML files, in particular URLs and <alt> tags. Furthermore, the combination of keyword statistics and visual features (colour) improves the retrieval of news web pages [23]. The reliability of keyword selection is improved by combining evidence from image captions and text surrounding images on web pages [4]. In order to be more systematic in the selection of keywords for image indexing it is necessary to recognise when a piece of text is describing the contents of an image, giving information about other aspects of the image, or is not related to the image at all. There seems to have been a tendency to avoid this issue by concentrating only on text very close to the image.

The correlations between visual and textual information are two-way and this fact is exploited by systems that learn to make associations between text features and image features. This learning supports applications such as auto-annotation and auto-illustration [2]. The learning task presupposes a strong degree of correspondence between each image-text pair, or more particularly between image regions and text fragments, but the notion of 'correspondence' has yet to be clearly articulated. Hypermedia is another kind of system that might benefit from better articulations of the correspondences between images and texts relate. Typed hypertext links were proposed to facilitate the navigation of scientific repositories [19]; this idea can perhaps be generalised to typed hypermedia links. An earlier survey of systems that integrate linguistic and visual information categorized the systems and surveyed related disciplines, but did not explicitly address different kinds of image-text combinations [18].

In summary, we suggest that the following questions are relevant for the further development and more general application of systems that process correlated image and text data. With regards to an image and a text in combination:

- How can we tell which is more important for successful communication?
- What correspondence is there between the information conveyed by one and by the other?
- What information, or other value, does one add to the other?
- If we understand the content of one, then what can we infer about the content of the other?

- What conventions are there for combining images and texts in particular genres of communication?

It has been argued that to develop algorithms that extract deeper levels of media content from low-level features, it is important to understand some compositional and aesthetic media principles [5]. We are interested in ideas that help us to describe relations between images and texts, and to recognise these relations automatically from low-level visual and textual features. Note, multimedia data models based on semiotic theories have been proposed, but these have focussed on multimedia data types individually rather than on their combination [6, 15].

Debates about the relative virtues of images and texts go back many thousands of years, with scholars arguing that either visual imagery or language is primary to human thought and communication. Rather than asking questions about their differences, it is perhaps now more important to consider how the two signifying systems interact within a single work. It has been argued that all media are mixed media, and that all representations used by humans to make sense of their world and communicate are heterogeneous [13]. Recent research in semiotics has dealt with thematic systems and cross-modal semiotic relations as part of multimedia discourse analysis [11], the translation between verbal and visual semiotic modes [10] and the combination of moving images and speech in documentary film [20].

In semiotics, pioneering work was done by Barthes who proposed three image-text relations: *illustration, anchorage* and *relay* [3]. The definitions of these relations combine the idea of status (the relative dominance of either image or text), and the idea of function. In *illustration* the image is said to be parasitic on the text, i.e. the text is dominant, and the image serves the function of elucidating the text – as for example with a photograph accompanying a newspaper story. In *anchorage* the situation is reversed so that the image is the main source of information, and the text only serves to elucidate what is already there. When image and text have equal status, and provide information not present in the other, then they share the relation of *relay*.

The idea of status was also applied by Halliday when he analyzed the relations between clauses in language. Two clauses are in an equal, or paratactic, relationship if both can be meaningfully understood on their own; otherwise they are in an unequal, or hypotactic relationship [7]. Halliday also analyzed function in terms of the logico-semantic relations between clauses, such as *elaboration*, *enhancement* and *extension*. One clause elaborates on the meaning of another by giving more specific detailed description of it. A clause enhances another by qualifying its meaning circumstantially in terms of time or location. Finally, one clause extends the meaning of another by adding further, but related information. These ideas are part of a semantically-oriented grammar that explains how the elements of language are combined to express meanings. The emphasis on meanings in this approach means that it is possible to envisage a multimedia grammar in which different media are integrated by one set of media-independent relations.

## 3. Image-Text Relations

We think that recognizing how an image and a text are related is a crucial step when understanding an image-text combination. We want to develop a framework so that for any given image-text combination we can classify the relations that hold between the image and the text, or between their parts. It is intended that these relations could be classified automatically, or entered as part of the multimedia authoring process, and that they would be used by systems that process image-text combinations. Our approach is to synthesise ideas from Barthes and from Halliday. We hope to maintain generality across all kinds of image-text combinations, though we present the framework here with a focus on common examples from web pages. Recently, a taxonomy of 49 relationships between images and texts was proposed, however, it was intended for human-human communication, e.g. between writers and illustrators, rather than for automatic classification [12].

### 3.1. Example Image-Text Combinations

This section discusses some relatively simple image-text combinations from web pages. Note that in our analyses we concentrate on the main images and texts and ignore peripheral images and text fragments, like banners, buttons and the text of hyperlinks.

### 3.1.1. Online News Stories
This section is based on web pages from three prominent news websites: news.bbc.co.uk, www.cnn.com and www.nytimes.com.

Almost every online news story comes with at least one image that is normally a photograph, and sometimes with two or three. The text on each webpage comprises: headline, image caption and story. The text of the story can be divided into the first paragraph, which gives an overview of the whole story and is often in bold font/type, and the remaining paragraphs that serve to give more details. Regarding page layout, the first image is positioned to the right of the opening paragraph and at the same level under the headline, as shown in Figure 1.
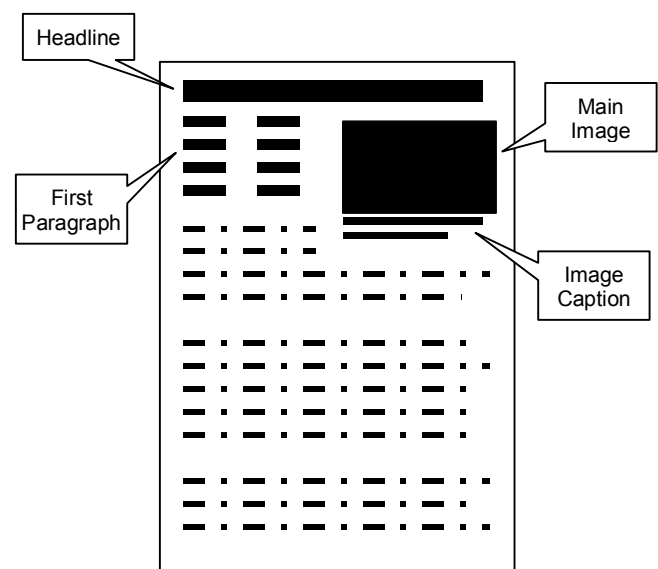


Fig. 1 Example of the layout of image and text on a news webpage.

It seems that the textual components of the webpage normally convey most, if not all, of the important information about the event being reported. The image depicts just one aspect of the event and seems to be used to grab attention (shocking scenes) or make a story seem more personal (pictures of people involved). The image is not essential in order to understand the story told by the text, but the image on its own without the text would often fail to communicate

anything meaningful. This aspect of the image-text relationship is reflected in the newspaper production process - a photo editor often adds a picture once the reporter has finished writing the story. The image caption often summarises the news story, like a more literal headline, rather than describe what is depicted in the image.

Some different ways in which the content of images relates to the main text of news stories can be seen by considering some specific examples; here we assume that the main story is summarised in the first paragraph of the text, and that the main image-text relation is between the image and the first paragraph. Sometimes the image is a close-up shot of one or two people at the centre of the story – see Figure 3a in Appendix A. In this case the name of the person, e.g. 'The Queen', appears in the headline, the image caption and the first paragraph of the story. Rather than a specific person, an image will sometimes depict an unnamed person as an example of a group of people who are the subject of a story, for example an image of an unnamed Aids sufferer that accompanies a story about the worldwide Aids situation. When reporting on an unexpected event then the only available image may be of the resulting state, like the devastation left after an explosion, Figure 3b. An image may also depict a cause for the event being reported, as in a story about more troops being sent to a war in response to increased resistance which is accompanied by an image depicting some resistors.

### 3.1.2. Online Art Galleries

This section is based on web pages from four prominent online art gallery websites – specifically pages showing individual paintings from their collections: www.metmuseum.org, www.sfmoma.org, www.louvre.fr, www.tate.org.uk.

Paintings displayed in online art galleries, like those displayed in physical galleries, are accompanied by some text comprising catalogue details (title, artist, date) and a caption which is part description of the painting's content and part background information about the painting. Compared to the news web pages there is less regularity in the layout of art gallery web pages. Nevertheless, there appears to be a tendency for the image to be placed at the top of the page, with catalogue details to one side of the painting at the same level – either to the left or to the right. The caption, often running to hundreds of words, is located either beneath the painting, or to the right, Figure 2.
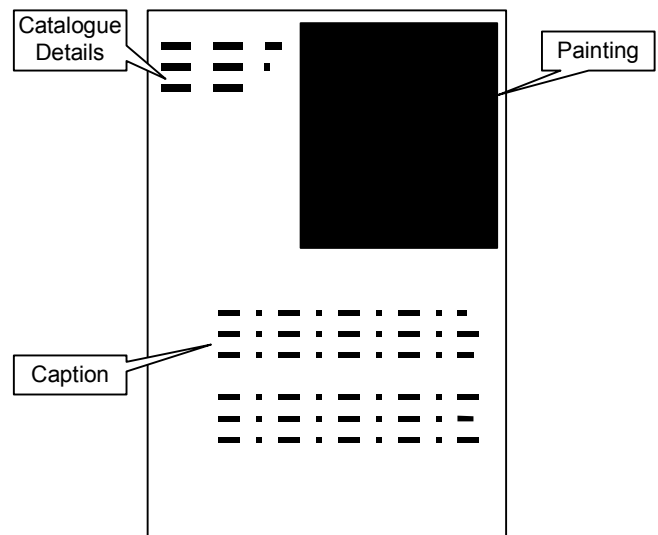


Fig. 2 Example of the layout of image and text on an art webpage.

The painting is the primary object of interest and is intended to convey meaning in its own right – the text is added later by someone other than the artist. The text may help the

viewer's appreciation of the painting, but is not essential. The text however makes little sense on its own and its dependence on the image is sometimes made apparent with phrases like 'This painting depicts…', 'this landscape…' and '…the characters on the left…'. Some of the caption, but often only a small part, describes the content of the painting directly. Note that sentences describing image content seem to always be in the present tense, e.g. 'Discord chooses the apple…' – Figure 3c. This is one event in a sequence, the consequences of which are described in the text but not shown in the painting. Other parts of a caption may give background explanations about the socio-historical context of the painting, or about the artist's motivations and inspirations. Further details about who commissioned the painting and where it has been displayed may also be given.

### 3.1.3. Online Science Textbooks
This section is based on scientific diagrams from: www.accessexcellence.org/RC/.

Diagrams are used extensively and in various ways throughout scientific literature. Authors take advantage of the fact that a diagram can be a clearer or more precise way to express their ideas. Sometimes there is a kind of redundancy when the text and an accompanying image convey essentially the same information, albeit in different ways. There is often extensive cross-reference between the text and images – particularly with labelled diagrams.

Figure 3d shows an example in which the image and the text refer to the same state of affairs, i.e. the structure of the DNA molecule, and both seem to communicate almost equivalent information. That is to say, the viewer of the web page could probably do without one or the other and still get the essential message. The equality of the image and the text is perhaps realized in the way that both are centered on the page. In other examples the relatedness of images and texts is realized by references in the text like 'Figure X shows…' or '…shown in Figure Z', and references to labeled parts of the diagram.

## 3.2. Our Classification Scheme

The examples discussed in Section 3.1 begin to hint at some of the ways in which images and texts can relate to one another. In our classification of image-text relations we distinguish two kinds of relations that we take to be mutually independent. Status relations are to do with the relative importance of the text and the image, or the dependence of one on the other. Logico-semantic relations are to do with the functions that images and texts serve for one another. We need to recognise that different relations may hold between different parts of images and texts, i.e. between image regions and text fragments.

### 3.2.1. Status Relations
The relation between an image and a text is *equal* when:

Either
• both the image and the text are required for successful communication, in which case they are *equal-complementary;*
Or
• both the image and the text can be understood individually, in which case they are *equal-independent*.

The relation between an image and a text is *unequal* when either the image or the text can be understood individually. That which cannot be understood individually is *subordinate* to the other.

Consider the examples discussed in Section 3.1. Images tend to be subordinate to the main text on news web pages whereas the text tends to be subordinate to the image on art gallery web pages. Image and text often share an equal relationship in scientific textbooks. The relationship is equal-independent when both convey the same information in different ways, and when there is cross-reference between the image and the text. Some technical images may require text to identify them, in which case the image and the text (probably a caption) are equal-complementary.

### 3.2.2. Logico-Semantic Relations

A text *elaborates* the meaning of an image, and vice versa, by further specifying or describing it. Photographs on news web pages frequently elaborate the text, specifically the first paragraph. For example, the image in Figure 3a specifies what the Queen looks like – this is information missing from the text. Part of the painting caption in Figure 3c elaborates the image through the description of image content. It is possible for an image to elaborate a text and a text to elaborate an image at the same time, as in Figure 3d.

A text *extends* the meaning of an image, and vice versa, by adding new information. Consider the image captions in Figures 3a and 3b. As well as naming the content of the image, i.e. 'The Queen' and 'The hotel', each caption gives new information not provided by the image. Another example would be car adverts comprising images of cars, and text giving information about price and performance.

A text *enhances* the meaning of an image, and vice versa, by qualifying it with reference to time, place and/or cause-effect. In Figure 3b there is an enhancement relation between the image and the first paragraph of the text because the image depicts the effect of the explosion reported in the text.

### 3.3. Towards Automatic Classification

A preliminary analysis has suggested some image features and text features that might be used together to classify image-text relations automatically. Further analysis is required to see if these features are present in a wider range of examples, and to determine which combination of features best classifies each image-text relation. The classification will be easier if the genre of the image-text combination is known because they may be preferred image-text relations with genre-specific realisations. Features of interest to us include:

- **Page layout and formatting:** the relative size and position of the image and the text; font type and size; image border
- **Lexical references in text:** for example, 'This picture shows…'; 'See Figure 1'; 'on the left'; 'is shown by'
- **Grammatical characteristics of the text:** tense – past / present; quantification – single / many; full sentences or short phrases
- **Modality of images:** a scale from realistic to abstract, or from

photographic to graphic – a function of depth, colour saturation, colour differentiation, colour modulation, contextualisation, pictorial detail, illumination and degree of brightness [10] – may correlate with use of GIF / JPEG

- **Framing of images:** for example, one centred subject, or no particular subject

### 3.3.1. Features to classify Status relations?

Two kinds of features that seem most indicative of status relations are page layout and lexical references. English is a language read left-to-right, and top-to-bottom, and there is an expectation that important information is positioned to be read first, and given most space. Extending this to the image-text scenario then perhaps we should expect a similar principle to guide web page layout, such that the subordinate media type should appear to the right / below, and take less space. On news websites, where images are subordinate to text, the photographs accompanying news stories are positioned to the right of the main story and occupy approximately 5-10% of the page space compared to the main text. In contrast, on most of the major art gallery websites we looked at, where texts are subordinate to images, the painting was positioned to the left and given about 50-100% the space of the text. Lexical references in text such as 'this painting…' or 'Figure X shows…' are strong indications that the text is about, and therefore subordinate to, the image – especially when occurring at the beginning of the text. However, some lexical references like 'this is shown in Figure X' may suggest equal status – especially when near the end of the text.

### 3.3.2. Features to classify Logico-Semantic relations?

Determining logico-semantic relations involves a comparison between what is depicted in the image and what is referred to by the text. If exactly the same people, objects and events are depicted and referred to, then there is elaboration. If completely new things are depicted / referred to then there is extension. If related temporal, spatial or causal information is provided then there is enhancement. The question is how may such comparisons be computed? For text, information extraction techniques can recognise proper nouns and work out who is the subject of a story, and determine what kind of event or state is being referred to in a text. For images, image processing techniques can detect faces, indoor vs. outdoor scenes, and framing – all of which may give clues about the main subject being depicted, e.g. a portrait, an anonymous person, the kind of event, etc. Working out whether a text refers to the same number of people as depicted in an image (one or many) would involve analysing quantifiers in the text, and detecting numbers of faces in the image. In the case of a single face it would be important to analyse whether it was a portrait of a specific person (centred, main focus), or a more anonymous character (possibly non-centred, turned from camera or out of focus) – this relates to the modality of the image. It might also be interesting to compare the complexity of the image and of the text: image complexity could perhaps be measured as a function of the number of edges / regions or graphic elements. Measures of text complexity relate to sentence length, average word length and use of embedded clauses.

When a text elaborates an image we have noted that often present tense is used, or short phrases rather than complete sentences. When an image elaborates a text in the news domain the image is of a realistic modality and typically depicts a person who is framed to fill the photograph – often head and shoulders; the image does not depict any significant action. The text tends to repeat the name of the person depicted. The enhancement relation of cause-effect is realised when the image depicts a process, and the text refers to a state, or vice versa. The image seems to normally be a general scene, rather than a closely cropped photograph with one main subject.

## 4. Closing Remarks

We are trying to establish a general theory of how images and texts combine which may be seen perhaps as a kind of multimedia grammar, applicable to different kinds of multimedia systems and applications, and to different genres of image-text combinations. Our contribution in this report, based primarily on theoretical analysis, comprises two parts. First, we have tried to open up what we believe is an important topic for multimedia research and have identified a set of general questions about image-text combinations. Second, we have outlined a framework to classify image-text relations, as a step towards understanding image-text combinations. In the framework we have synthesized ideas from semiotic theories about images and texts. Using the framework, we have analyzed some relatively simple examples of image-text combinations from web pages and noted some

preferred combinations in particular genres. This analysis led us to consider some ways in which image-text relations might be classified automatically. We have also tried to argue that machine-executable representations of image-text relations, whether classified automatically or entered as part of an authoring process, could help multimedia retrieval, browsing, adaptation, generation, etc.

## 5. References

1. André, E. The Generation of Multimedia Documents. In *A Handbook of Natural Language Processing*, Dale, Moisl and Somers (eds.), Marcel Dekker Inc., 2000.
2. Barnard, K., Duygulu P. and Forsyth D. Exploiting Text and Image Feature Co-occurrence Statistics in Large Datasets. Chapter to appear in *Trends and Advances in Content-Based Image and Video Retrieval*, 2004.
3. Barthes, R. *Image-Music-Text*. Fontana, London, 1977.
4. Coelho, T.A.S., Calado, P.P., Souza, L.V., Ribeiro-Neto, B. and Muntz, R. Image Retrieval Using Multiple Evidence Ranking. *IEEE Trans. Knowledge and Data Engineering* 16, 4 (April 2004), 408-417.
5. Dorai, C. and Venkatesh, S. Computational Media Aesthetics: Finding Meaning Beautiful. *IEEE Multimedia* 8, 4 (Oct-Dec 2001), 10-12.
6. Gonzalez, R. Hypermedia Data Modeling, Coding and Semiotics. *Procs. of the IEEE* 85, 7 (July 1997), 1111-1140.
7. Halliday, M. A. K. *An Introduction to Functional Grammar*. Edward Arnold 2nd edition, London, 1994.
8. Hu, J. and Bagga, A. Categorizing Images in Web Documents. *IEEE Multimedia* 11, 1 (Jan-Feb 2004), 22-30.
9. Jain, R. Are We Doing Multimedia? *IEEE Multimedia* 10, 4 (Oct-Dec. 2003), 110-111.
10. Kress, G. and van Leeuwen, T. *Reading Images: the grammar of visual design*. Routledge, London and New York, 1996.
11. Lemke, J. L. Travels in Hypermodality. *Visual Communication*, 1, 3 (2002), 299-325.
12. Marsh, E. E. and White, M. D. A Taxonomy of Relationships between Images and Text. *J. of Documentation* 59, 6 (2003), 647-672.
13. Mitchell, W. J. T. *Picture Theory: essays on verbal and visual representation*. University of Chicago Press, Chicago and London, 1994.
14. Owen, C.B. and Makedon, F. Cross-Modal Information Retrieval. In Furht (ed.). *Handbook of Multimedia Computing*. CRC Press, Florida, 1998.

15. Purchase, H. Defining Multimedia. *IEEE Multimedia* 5, 1 (Jan.-March 1998), 8-15.

16. Rowe, L. A. and Jain, R. ACM SIGMM Retreat Report on Future Directions in Multimedia Research. http://www.acm.org/sigmm/

17. Smeulders, A.W.M., Worring, M., Santini, S., Gupta, A. and Jain, R. Content-Based Image Retrieval at the End of the Early Years. *IEEE Trans. Pattern Analysis and Machine Intelligence* 22, 12 (2000), 1349-1380.

18. Srihari, R. K. Computational Models for Integrating Linguistic and Visual Information: A Survey. *Artificial Intelligence Review* 8, 5-6 (1995), 349-369.

19. Trigg, R. *A Networked-Based Approach to Text Handling for the Online Scientific Community*. PhD Dissertation, Uni. of Maryland, 1983. http://www.workpractice.com/trigg/

20. van Leeuwen, T. Conjunctive structure in documentary film and television." Continuum 5:1 (1991), 76-113.

21. Wang, X.-J., Ma, W.-Y. and Li, X. Data-Driven Approach for Bridging the Cognitive Gap in Image Retrieval. *Procs. IEEE ICME 2004, Taipei.*

22. Xie, X., Ma, W.-Y. and Zhang, H.-J. Maximizing Information Throughput for Multimedia Browsing on Small Displays. *Procs. IEEE ICME 2004, Taipei.*

23. Zhao, R. and Grosky, W. I. Narrowing the Semantic Gap – Improved Text-Based Web Document Retrieval Using Visual Features. *IEEE Trans. Multimedia* 4, 2 (2002), 189-200.

# Appendix A

Figure 3a

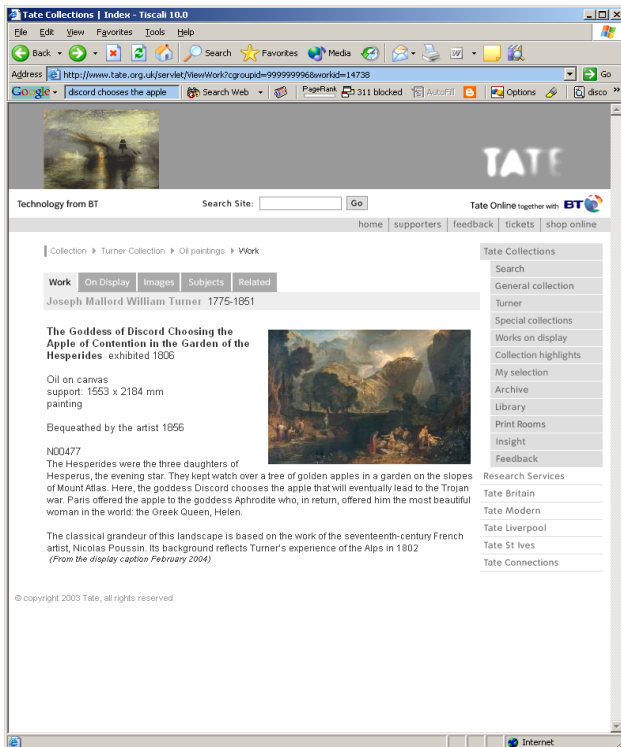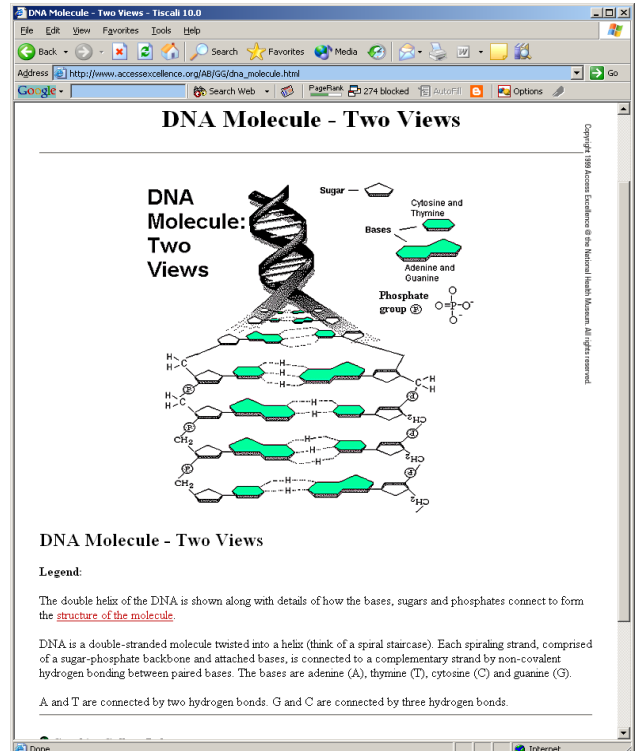Figure 3b

Figure 3c

Figure 3d

# Department of Computing

University of Surrey
Guildford, Surrey
GU2 7XH UK

Tel:          +44 1483 683133
Fax:         +44 1483 686051
E-mail:    a.salway@surrey.ac.uk
www.surrey.ac.uk