

Visual Semantics: Extracting Visual Information from Text Accompanying Pictures *

Rohini K. Srihari and Debra T. Burhans

CEDAR/SUNY at Buffalo
UB Commons, 520 Lee Entrance- Suite 202
Buffalo, NY 14228-2567 USA
rohini@cs.buffalo.edu
burhans@cs.buffalo.edu

Abstract

This research explores the interaction of textual and photographic information in document understanding. The problem of performing general-purpose vision without apriori knowledge is very difficult at best. The use of collateral information in scene understanding has been explored in computer vision systems that use general scene context in the task of object identification. The work described here extends this notion by defining visual semantics, namely, techniques for systematically extracting picture-specific information from text accompanying a photograph. Specifically, this paper discusses the multi-stage processing of textual captions with the following objectives: (i) predicting which objects (implicitly or explicitly mentioned in the caption) are present in the picture, and (ii) generating several types of constraints useful in locating/identifying these objects. The implementation and use of a lexicon specifically designed for the integration of linguistic and visual information is discussed. Finally, the research described here has been successfully incorporated into PICTION, a caption-based face identification system.

Introduction

The general problem being investigated is that of establishing a correspondence between words and the visual depictions they evoke which will be referred to as the *correspondence problem*. This correspondence is not one-to-one in general, but for certain specialized domains it is possible to establish a direct correspondence between words and pictorial elements being referenced by those words (e.g. weather maps: for example, the word ‘pressure’ causes the system to examine isobars in the picture). (Jackendoff 1987) attempts to establish a correspondence between words and 3D models of objects, but the problem is handled primarily at the single-word level (both nouns and verbs) and does not extend to establishing a correspondence be-

tween a sentence/phrase and the complex scene it may evoke.

In the present research, we focus on captioned pictures: specifically, we address the problems of (i) identifying useful information (for picture understanding) in the text and (ii) extracting and representing this information so that it can be used to direct a computer vision system in the task of picture understanding.

We present a new theory of *visual semantics* which applies to situations where pictures are accompanied by descriptive text (caption). The theory includes the following elements, all of which are described in this paper.

- The design of suitable lexicons for integrating linguistic and visual information.
- The representation of visual information as a set of four types of constraints which can be applied to the picture. In computer vision, “visual information” refers to knowledge about objects which is required in order for a computer vision system to detect them in a scene. This includes 3D descriptions of objects in terms of their components and the spatial constraints between them, as well as typical scene information which relates objects in a common context. An example of this is the modelling of a typical neighbourhood scene (comprised of streets, houses, trees, etc.)(Weymouth 1986).
- A systematic procedure for processing a caption whose objective is the generation of visual constraints. Visual information can appear in a variety of syntactic constructs: noun phrases, verbs, relative clauses, subordinate clauses, etc. We show that syntax is often not sufficient for extracting visual information: in such situations a deeper understanding of the semantics of a word or phrase must also be considered.

Visual semantics plays a key role in allowing scene descriptions to be dynamically constructed from descriptive text. These scene descriptions can then be used by an image-processing system to guide knowledge-based interpretation of the associated picture. A system, *PICTION*, based on visual seman-

*This work was supported in part by a grant from ARPA (ARPA 93-F148900-000)

tics is described which uses information obtained from parsing a natural-language caption of a newspaper photograph to label faces in the accompanying photograph. (Strat & Fischler 1991) also discusses the use of context in visual processing but does not cover the generation of constraints.

This research is most relevant in the context of document image understanding. Pictures with captions are ubiquitous in documents, newspapers and magazines. The information contained in both pictures and captions enhances overall understanding of the accompanying text, and often contributes additional information not contained specifically in the text. This information could subsequently be incorporated into an integrated (text and picture) database as part of a larger information retrieval system which permits *content-based retrieval*.

PICTION: A Caption-based Face Identification System

This section discusses the design and implementation of a system called *PICTION* (Srihari 1993; 1991), that identifies human faces in newspaper photographs based on information contained in the associated caption. Since the caption and picture both contribute information to a reader, we refer to these collectively as a *communicative unit*. When given a text file corresponding to a newspaper caption and a digitised version of the associated photograph (referred to as an image), the system is able to locate, label, and give information about objects which are relevant to the communicative unit. *PICTION* was initially tested on a database of 50 pictures and was completely successful at identifying faces in 65

Traditional methods of face recognition employ model-matching techniques: only people for whom pre-stored face models exist can be identified. *PICTION* is noteworthy since it provides a computationally less expensive alternative: faces are identified based solely on visual information conveyed by accompanying text. A key component of *PICTION* is the face locator, which locates (but cannot recognize) human faces in arbitrarily complex news photographs. Spatial constraints obtained from the caption, along with discriminating visual characteristics (e.g., male/female) are sufficient for eliminating false candidates from the face locator and correctly identifying faces.

Figure 1 is an example of a digitised newspaper photograph and accompanying caption that the system successfully processes. The male/female filter is not able to distinguish between Toukie Smith and Robert DeNiro, thus resulting in multiple possible bindings. It should be noted that Sean Penn is identified correctly and uniquely, based on spatial constraints. Since 'Toukie' is an uncommon name, it is not present in a database of proper names. The mention of the word 'sister' allows us to infer that 'Toukie' refers to a fe-

male.

Figure 1 shows the overall control structure of *PICTION*. As the diagram suggests, the three main components of *PICTION*'s architecture are (i) a natural-language processing (NLP) module, (ii) an image understanding (IU) module, and (iii) a language-image interface (LII). The NLP and IU modules interact through the LII which maintains the long-term knowledge base. *PICTION* has been implemented in LOOM (Corporation 1991) which is a high level programming language and environment for constructing knowledge based systems. The visual routines have been implemented in C, and a C/LISP-LOOM interface facilitates communication through the LII. The system runs on a Sun Sparcstation.

The input to the NLP Module is the original newspaper caption; the output is a set of visual constraints. The NLP module (illustrated in Figure 1) is divided into three stages: syntactic parsing, PSI (partial semantic interpretation) and CBCG (caption based constraint generation). The NLP module is called on only once to process the caption and generate visual constraints on the picture. The LII module subsequently converts relevant visual information into a series of directives for the vision module which is then called on to interpret the picture. The vision module has several features which enable it to be guided by the LII module, including: constrained search, ability to characterise objects, crude and refined object-location procedures, ability to change parameters and repeat actions, output compatibility with intermediate representation, and the ability to perform bottom-up interpretation when necessary. Information is consolidated in the LII module by satisfying various types of constraints: this necessitates repeated calls to the vision module.

The remainder of this paper focuses on the NLP module, namely, the processing of the caption.

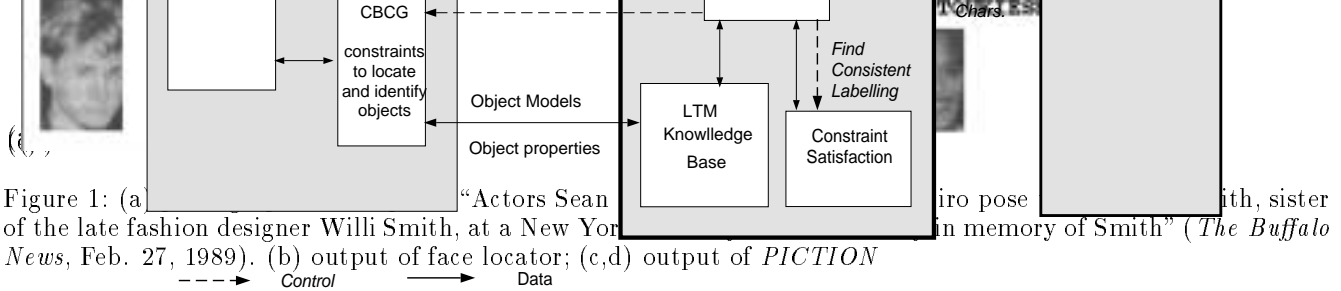


Figure 2: PICTION- System Overview

Lexical Database

PICTION uses a broad coverage syntactic/semantic lexical database which has been constructed from the following sources: (i) Longman Dictionary of Contemporary English (LDOCE), (ii) Oxford Advance Learners Dictionary (OALD), (iii) WordNet(Beckwith *et al.* 1991) (iv) name lists (CLR) and (v) manual augmentation (ontological information and visual semantics). This is similar to work done by Kevin Knight at ISI (Knight & Luk 1994).

Syntax, Semantics obtained from LDOCE, WordNet

A LOOM knowledge base instance is constructed for each word in the lexicon. Name instances contain information on gender; word instances contain semantic and syntactic information from LDOCE (including subject field codes, semantic restrictions and verb subcategorization), morphological information from OALD and pointers to concepts representing WordNet synsets. WordNet synsets are represented as LOOM concepts organized into a hierarchy which reflects the organization of WordNet. A top level ontology based on the Naive Semantics model (Dahlgren 1988) is incorporated into the model. Synset concepts contain part-of, is-part and (where appropriate) verb classification information from WordNet. Visual and procedural semantic properties have been added manually to a subset of the concepts.

Our approach uses and extends the *Naive Semantics* model which is a theory of associating commonsense knowledge with words. It is necessary to capture fixed, visual properties of objects, such as size and colour, as well as procedural information for certain words and phrases. For example, a recent caption identified one of the people in the corresponding photograph as “the person wearing the hat”. This should generate a call to an object finder with the location “above head”, and the scalar and shape properties of the object (hat).

Examples for the entries for hat and wear:

The entry labelled “WNhat” represents information obtained from the corresponding WordNet entry for the sense of hat defined as lexicon entry “W137”. Synonyms for hat include “chapeau” and “lid”. The “scale”, “shape” and “function” entries for “WNhat” have been manually added and are used in generating visual information. We specify that if the event “wear”

is associated with “hat”, and the subject of “wear” is a human, then a typical location for hat is on the head of the person wearing it. This information is specified procedurally. The CBCG uses this information to generate a locative constraint (described later) which is subsequently used by the LII and the vision module in identifying the person wearing the hat.

Visual Hierarchies

WordNet organizes words into synsets and relates synsets through a hierarchy of superconcepts. However, it is clear that a visual hierarchy is also required, namely, a hierarchy for objects in terms of the visual procedures required to recognise them. We have defined visual hierarchies in terms of *visual superconcepts* which reflect the type (man-made, natural), shape, textural properties, boundary properties etc. of an object. New links (*visual-is-a*, *visual-part-of*) have been defined from existing WordNet synsets (representing concrete objects) to these superconcepts. Figure illustrates visual hierarchies superimposed upon existing WordNet hierarchies. Objects inherit all the properties of their visual superconcepts; however, attributes such as size, colour, expected location, etc. can be specialised. Depending on the visual classification, appropriate recognition tools (such as segmentation tools, edge detectors, surface detectors) can be invoked. At the lowest (i.e., most specific) levels of the hierarchy are specialized object detectors for certain common object classes (e.g., human face, tree, building, car) ¹.

Our visual hierarchy is based on a *blob* theory of object recognition. There are many objects for which it is difficult to construct detailed shape descriptions. In such cases it is sufficient to state some of their properties such as natural/man-made, boundary description, etc. These are useful in delineating these objects from the image. In cases where it is necessary to find these objects, context from the accompanying text can be used to constrain the search. For example, consider a picture of a man holding a trophy accompanied by the caption “Thomas Smith holding the trophy he won at ...”. If it were necessary to identify the trophy, the system would first find the face of Thomas Smith and then search in the appropriate vicinity for an object

¹According to (Biederman 1988), there are about 3000 common entry-level objects which the human perceptual system can detect.

exhibiting the required properties (man-made, small-medium size, etc.)

Our representation allows objects to be modelled at various resolutions. For example, at the most general level, a tree is a natural object with a fractal boundary representation. A more detailed visual model of tree defines the visual parts of a tree as well as the spatial relationship between these parts. These parts are classified according to the chosen set of visual super-concepts. At the most specific level, the description of a tree includes a specialised recognition module for trees.

The has-part information in visual object models may differ from the has-part information used in existing WordNet entries. There are 2 reasons for this: (i) the has-part information in WordNet may be too fine-grained to ever be exploited by the vision system (e.g., shoes can have laces), and (ii) the names of the parts may differ (e.g., the WordNet entry for 'tree' includes crown, trunk, branches and stump, however the visual description for tree will have a treetop and a trunk). This is illustrated in Figure .

A study in human cognition regarding the task of identifying objects reveals that there is justification for having different abstractions for words and pictures. (Linde 1982) postulates the existence of two separate semantic-memory representations. (Jolicoeur, Gluck, & Kosslyn 1984) states that both words and pictures may use the same semantic-memory representation; they differ however in the *entry point*, namely the particular level of abstraction in the hierarchy at which the association is made.

Parsing

Pre-Processing Input

There are two objectives to this phase. The first is the elimination of *directive* phrases such as "left-of", "front row, left to right", etc. These are not essential to syntactic parsing and are passed directly to the CBCG in the parser output structure. This is accomplished using a set of rules which indicate the noun phrase to which the directive information should be attached.

The second objective of pre-processing is the detection and classification of proper noun sequences. Input texts usually contain unknown words (words not in the lexicon). Since the parser requires all input words to be known, the detection of unknown words, most of which are proper nouns, is critical. In the following example taken from a recent edition of *The Buffalo News*, the proper noun sequences are bracketed:

[Vladimir Horowitz] at [Steinway and Sons],
[New York],...

In addition, "Vladimir Horowitz" is classified as a name, "Steinway and Sons" is classified as an organization name and "New York" is looked up and classified as a location.

A hidden Markov model is used to detect proper noun sequences, with a portion of the Penn Treebank

corpus (excerpts from *The Wall Street Journal*) used as the training set. We are currently able to detect proper noun sequences with about 90% accuracy. This figure can be improved by increasing the size of the training set, which is being done. Proper noun lists containing appropriate classification information such as gender, location, organization, etc. as well as heuristic rules (involving typical suffixes such as "Inc.") are employed in the categorization of proper nouns. Correct classification of these proper noun sequences still poses a significant challenge (Mani *et al.* 1993).

LFG Parser

An LFG which covers basic features of English grammar has been implemented, requiring approximately 1000 lines of code.. Morphological analysis, which is a part of syntactic analysis, is embedded in the lexicon searching functions. We employ an LR parser augmented by pseudo/full unification packages (Tomita 1987). The parser is made efficient through compilation of the grammar into an LR Parsing Table. The parser outputs a recursive structure of pairs of features and values. It utilizes some basic semantic features (e.g., animate/inanimate) in order to verify semantic restrictions during parsing. Such information is purely ontological and is obtained through lexicon lookup. These semantic features are incorporated into the final output structure.

2

Partial Semantic Interpretation

The purpose of this phase is to disambiguate among multiple syntactic structures output by the parser. Specifically, this phase includes (i) prepositional phrase disambiguation, (ii) word sense disambiguation and (iii) anaphora resolution. Currently, we are handling only cases that can be disambiguated by ontological information present in the lexical database. We are experimenting with statistical techniques for handling (i) and (ii). Anaphoric references are being resolved by the use of weights assigned to various referents depending on their role in the sentence.

Caption-based Constraint Generator (CBCG)

Input to the CBCG is the disambiguated parse from the PSI stage; output is the *conceptualized graph*. The CBCG makes use of the long-term memory database (LTM) in the LII for retrieval of information about well-known people, and the semantic lexicon for assigning visual semantics to the parse.

Constraint Types

In PICTION, visual information derived from text generates a set of objects (items predicted to be in the

²These lists were obtained from the Consortium for Lexical Research at New Mexican State University.

picture) and constraints on those objects. The set of constraints are divided into four types:

- **Spatial Constraints** are geometric or topological constraints, such as left-of, above, inside, etc. They are generally binary and describe inter-object relationships. Complex spatial constraints such as “surround” are broken down into a set of constraints based on spatial primitives.
- **Locative constraints** express information about the location of objects in the picture with respect to a particular frame of reference. The information conveyed is procedural in nature, for example, if you are told there is a chair in the corner, it results in the following high-level procedure construct: *locate_in_vicinity(chair, region(corner(entire_image)))*.
- **Characteristic Constraints** are unary constraints which describe properties of objects. Examples include gender and hair color.
- **Contextual Constraints** are those which describe the setting of the picture, and the objects which are expected to appear. For example, the people present (explicitly mentioned in the caption), whether it is an outdoor scene, and general scene context (apartment, airport, etc).

Some visual constraints are expressed explicitly as assertions, for example *left_of(person_a person_b)*. Locative and characteristic constraints are implicit in the object model. Contextual constraints consist of the object instances which have been instantiated, and a general scene context which is also asserted, hence these are explicit in nature.

As an example, consider the photograph and caption of Figure which was processed successfully by the face-identification system *PICTION*. The output of the CBCG for this example would be the following:

Automatic Generation of Visual Constraints

The CBCG has been written as a rule-based system which uses LOOM concepts to drive the semantic “parsing”. LOOM concepts representing different syntactic parse structures which can be encountered, at both the top and inner levels of the PSI output, have been created. LOOM methods and rules which drive the CBCG have corresponding actions which are invoked when the input matches a particular concept.

There are three main categories of rules.

- **Word-based:**
Spatial and characteristic constraints are frequently indicated by single words. Examples include left, right, above and below, as well as characteristics such as hair colour and titles like President.
- **Phrase-based:**
Locative and characteristic constraints are often indicated by directive phrases. Examples include “be-

tween the two buildings” and “wearing the striped shirt”.

- **Sentence-based:**

Contextual constraints can generally be inferred at the sentence level, taking into account the various objects mentioned and their relations and properties. An example of this is the “SVOPP” rule which states that if the sentence is of the form subject-verb-object-prepositional_phrase, and, if both the subject and object represent humans, and, if the PP represents a time and/or location, then, propose that both the subject and the object are in the picture.

The top-level rules are based on syntactic structure and attempt to predict which people (or objects) are in the picture. Verifying the antecedents of these rules (e.g., is the person deceased?) causes other rules to be fired. The final action is to generate identifying information for every person/object predicted to be in the picture.

As an example, consider the caption “Actors Sean Penn, left, and Robert DeNiro pose with Toukie Smith, sister of the late fashion designer Willi Smith, at a New York celebrity auction Sunday”. All three people mentioned will be predicted to be in the picture as the sentence has the form <subject-list> <verb> <object> <adverbial place> <adverbial time>.

The picture is predicted to be of an indoor scene when the rule is-indoor is fired, which happens when the inside predicate concept is satisfied. “auction” has the semantic feature \$social-event stored in the semantic lexicon, and there is nothing mentioned which has the semantic feature of being outdoors.

Characteristics such as gender are generally indicated by single words. In this example, the fact that Toukie Smith is female is inferred by the presence of the word “sister” in the directional phrase associated with her name: the synset containing “sister” has a \$female semantic feature, and a rule for directional phrases states that the first or last word will generally refer to the person named.

The following rule which is a phrase based constraint infers that a person is deceased if the phrase which describes him or her contains the adjective “late” followed by a person word. In this case a boolean flag is set which will prevent the CBCG inferring the presence of this person in the picture:

As in the above examples, the \$person semantic feature is stored in the lexicon.

Finally, the production rules for spatial constraints will fire thereby generating the spatial constraints shown in the previous section.

Summary

This paper has presented a new theory of visual semantics which concerns the use of descriptive text in interpreting accompanying photographs. Although the examples are from captioned newspaper photographs,

and the application is knowledge-based vision, this work can be extended to any domain where both language and pictures are used to communicate information. A highlight of this work is the development of a semantic lexicon which includes visual hierarchies, as well as a systematic procedure for generating visual constraints from text accompanying a picture.

References

- Beckwith, R.; Fellbaum, C.; Gross, D.; and Miller, G. A. 1991. WordNet: A Lexical Database Organized on Psycholinguistic Principles. In *Lexicons: Using On-line Resources to Build a Lexicon*. Lawrence Erlbaum.
- Biederman, I. 1988. Aspects and extensions of a theory of human image understanding. In Pylyshyn, Z., ed., *Computational Processes in Human Vision: An interdisciplinary perspective*. Ablex.
- Corporation, I. 1991. *LOOM Users Guide , Version 1.4*.
- Dahlgren, K. 1988. *Naive Semantics for Natural Language Understanding*. Boston: Kluwer Academic Press.
- Jackendoff, R. 1987. On Beyond Zebra: The Relation of Linguistic and Visual Information. *Cognition* 26(2):89-114.
- Jolicoeur, P.; Gluck, M. A.; and Kosslyn, S. M. 1984. Pictures and Names: Making the Connection. *Cognitive Psychology* 16:243-275.
- Knight, K., and Luk, S. 1994. Building a Large Scale Knowledge Base for Machine Translation. Forthcoming. In *Proceedings of AAAI-94*.
- Linde, D. J. 1982. Picture-word differences in decision latency. *Journal of Experimental Psychology: Learning, Memory and Cognition* 8:584-598.
- Mani, I.; MacMillan, T. R.; Luperfoy, S.; Lusher, E. P.; and Laskowski, S. J. 1993. Identifying Unknown Proper Names in Newswire Text. In *Proceedings of the Workshop on Acquisition of Lexical Knowledge from Text*, 44-54.
- Srihari, R. K. 1991. PICTION: A System that Uses Captions to Label Human Faces in Newspaper Photographs. In *Proceedings of the 9th National Conference on Artificial Intelligence (AAAI-91, Anaheim)*, 80-85. AAAI Press.
- Srihari, R. K. 1993. Intelligent Document Understanding: Understanding Photos with Captions. In *Proceedings of the International Conference on Document Analysis and Recognition (ICDAR-93)*, 664-667.
- Strat, T. M., and Fischler, M. A. 1991. Context-Based Vision: Recognizing Objects Using Information from Both 2-D and 3-D Imagery. *IEEE PAMI* 13(10):1050-1065.
- Tomita, M. 1987. An Efficient Augmented-Context-Free Parsing Algorithm. *Computational Linguistics* 13(1-2):31-46.
- Weymouth, T. 1986. *Using Object Descriptions in a Schema Network for Machine Vision*. Ph.D. Dissertation, University of Massachusetts at Amherst.

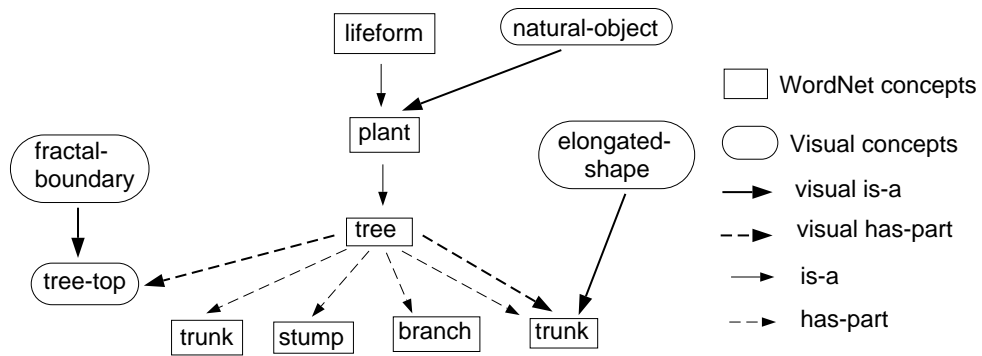


Figure 3: Visual hierarchy superimposed on WordNet concept hierarchy