

LandScan: A Natural Language and Computer Vision System for Analyzing Aerial Images

Ruzena Bajcsy
Aravind Joshi
Eric Krotkov
Amy Zwarico

CIS Dept/D2
University of Pennsylvania
Philadelphia, Pennsylvania 19104

Abstract

LandScan (LANguage Driven SCene ANalysis) is presented as an integrated vision system which covers most levels of both vision and natural language processing. Computations are both data-driven and query-driven. In the report we focus on the design of the vision and control modules. Future work will investigate in more detail the design of the natural language interface. The data-driven system employs active control of stereo cameras for image acquisition, and dynamically constructs a surface model from multiple aerial views of an urban scene. The query-driven system allows the user's natural language queries to focus analysis to pertinent regions of the scene. This is different than many image understanding systems which present a symbolic description of the entire scene regardless of what portions of that picture are actually of interest.

1. Introduction

The aim of our research on LandScan (LANguage Driven SCene ANalysis) is to develop a system capable of dynamically updating and maintaining a model of an urban world over multiple aerial views. The system will have a natural language front end through which users can query the system about a scene, and interactively assist the vision processing by restricting the analysis to those areas of the scene which are of current interest. A unique contribution of the work is that processing is both data-driven (bottom up, determined by sensor data) and query-driven (top down, determined by user queries). The integration of both methods into one system can help overcome the shortcomings of each method employed independently. For example, if data-driven processing were able to segment a graph of edges derived from the image into *several* different connected components, query-driven information about what the system *should* be looking for can help impose structure, and a unique segmentation, upon the otherwise ambiguous data.

The data-driven processing starts with stereo aerial images and reconstructs the surfaces in the scene. The query-driven processing constructs a logical representation of the scene using the queries to guide analysis. High-level scene analysis is performed using an Augmented Transition Network (ATN).

This research was supported by the following grants: ARO DAAG-29-84-k-0061 AFOSR 82-NM-2W NSF MCS-8219190-CER NSF MCS 82-07294 AVRO DAAB07-84-K-F077 NIH 1-ROI-HL-29085-01.

As an example, suppose the user asks, "Is there a car on the street?*" The output from this query would be: the objects to be recognized, car and street; the relation ON which must hold between them; and an indication that this query is responded to by a yes/no answer with some explanation. The vision system would then be called to find a car and a street in the relation ON. The car and street would then be added to the Scene Model (if not there already) and the system would reply with an affirmative response.

This paper will describe some related research, the implementation of the data-driven and query-driven portions of the LandScan system, and our plans for future work. A later paper will detail how natural language queries will interface with LandScan to guide the scene analysis.

2. Related Research

A large corpus of research on aerial image understanding *per se* exists, and many general vision techniques are applicable to the aerial domain. Large aerial projects have been undertaken at USC [Nevatia 83], CMU [Herman 83] and SRI [Fischler 83]. However, very few integrated systems have been successfully implemented, and the best system architecture is still an open question. In particular, the problem of providing high-level feedback to the vision system has not been adequately addressed.

ATN's have been used primarily in the domain of natural language [Bates 81], [Winograd 83]. A notable exception is the system designed by Tropf and Walter (Tropf 83) which uses an ATN model for the recognition of 3D objects with known geometries.

The work of Talmy [Talmy 83] and Herskovits [Herskovits 82] influenced the design of both the topological relations in the models and the choice of linguistic attributes which must be associated with objects in order to ensure a robust and reliable natural language interface.

3. Data-driven System Implementation and Results

This section will describe the data-driven vision module?, which must be effective in an urban world, seen from above. Urban scenes are characterized by an abundance of straight lines and planar surfaces. Under

these constraints, the scene may usefully be approximated as polyhedra.

We have tested the modules on real, highly complex aerial images; it is very difficult to present these results. We show results derived from imaging a subset of the scale model depicted in Figure 6-2: a "mock up" of an urban scene. One advantage of using the scale model is the clarity of the results, and the ease of verifying their compatibility with reality.

3.1. Vision Modules

A stereo pair of images is acquired. The gradient ∇I of the images blurred at multiple resolutions is computed, and the Canny operator is used to locally suppress non-maxima in the gradient magnitude $\|\nabla I\| = \text{SQRT}(I_x^2 + I_y^2)$. We call the surviving pixels "edgels". To find corners, the variance σ^2 in the gradient direction $\theta = \tan^{-1}(I_y/I_x)$ is computed over a local neighborhood. The corneriness of an edgel is proportional to the product $\sigma^2 \|\nabla I\|$.

Corresponding edgels in the two images are matched using 2-sided correlation at multiple scales. In images of parts of the scene in Figure 5-1, 83% of the vertically oriented edgels are matched, and the disparity at each is computed. The resultant sparse depth map is refined by linear interpolation (acceptable under the constraints of this domain), first across columns and then across rows. Figure 5-2 depicts the interpolated depth map (corresponding to only the more distant objects in Figure 5-1).



Figure 5-1: Scale model.

3.2. Surface Model

A graph is constructed to serve as the surface model [Krotkov 84]. The construction algorithm converts a set of contours into a set of closed contours represented as a graph (a linked list of vertices, edges, and faces) by traversing edges and at trihedral junctions choosing the path making the most acute angle with respect to the present path.

Surface attributes and relations are computed in the

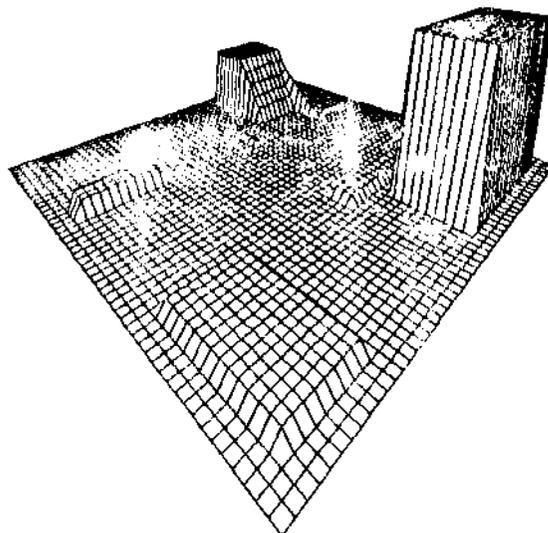


Figure 6-2: Depth map of partial view of 5-

SurfsUP [Radack, et al 84] geometrical modelling system. Attribute values for each face in the surface graph are computed: compactness, centroid, normal, area, type (building, sidewalk, field, street, and unknown), and number of sides. These values are computed once and stored on an attribute list. Computed topological relations are *above*, *adjacent* (touching), *contiguous* (sharing an edge), *contains* (proper inclusion), *looksadjacent*, *lookscontiguous* (respectively adjacent and contiguous under perspective transformations). Relations (and indirectly their complements) are computed once and stored as Boolean arrays.

4. Query-driven System Implementation and Results

This section describes the design and implementation of the query-driven processes. These include object recognition and scene modelling [Zwarico 84], high level reasoning processes, and query handling.

4.1. Object Recognition

The ATN formalism has been chosen as the paradigm for object recognition in LandScan. It is composed of three parts: the grammar, a dictionary, and an interpreter. The grammar represents the *a priori* or world knowledge that the system must have in order to recognize objects. The dictionary represents the actual data: a list of all of the faces which have been segmented by vision and the relations between them. The third component is a Lisp program which provides the control structure for the process. Figure 5-3 shows the results of running the recognizer on the scene in Figure 5-1.

4.2. The Scene Model

The Scene Model is composed of two components: a list of objects currently known to be in the scene and a set of matrices representing the primitive relations. The objects on the object list have already been recognized. Each object has associated with it a list of surfaces, its location, and a subtype. The relations which are the same as in the surface model, are represented by their adjacency matrices because the adjacency matrix is easily updated and makes composition of relations a simple matter of Boolean matrix multiplication.

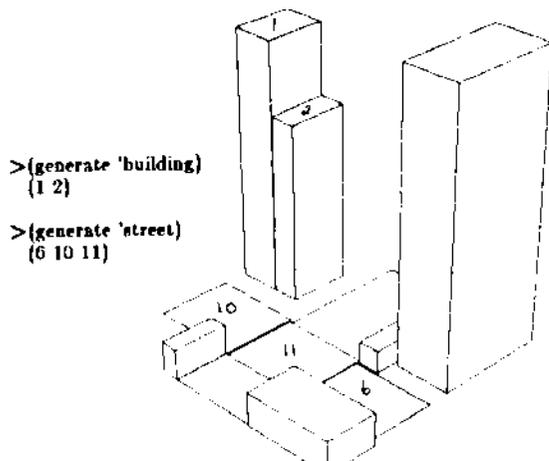


Figure 5-8: Reconstructed planar surfaces.

The Scene Model is dynamic because information can be added to it as further image analysis occurs. A new object is added to the head of the list. The relations are updated by calculating the relations between the new entity and the current object list.

4.3. Linguistic Analyzer

Given a query, the Linguistic Analyzer will symbolically represent this utterance so that it can be used by the reasoning process to analyze the image. The Linguistic Analyzer will parse the query, determine the query type, and categorize all implicit subqueries in the actual utterance. Its output will contain a list of the objects to be found, the relations which must hold between these objects, and the query type (so that an appropriate response can be generated). This is not yet implemented.

4.4. Reasoning

The Reasoner analyzes the query, determines the strategy for obtaining an answer to the query, and provides feedback to the vision system. In order to obtain the information necessary for the generation of the response the Reasoner must have both runtime data (the current Scene Model and the query) and global knowledge (the World Model and the Object Model).

The World Model describes the features and relations of the objects in the domain, buildings, streets, sidewalks, etc. The world is represented by a labelled directed multigraph in which the nodes are the objects in the domain and the arcs are labelled with the relation which can hold between them. The Object Model represents the expected physical features (subparts) and linguistic properties (features which affect the usage and interpretation of a spatial construct) of the objects in the domain.

If the reasoning processes fail to produce a positive response (the query fails to have an answer), the Reasoner performs two types of query failure analysis. The first type of query failure involves a query violating the global knowledge. In this case, the system will respond with a message indicating why the query is conceptually ill-formed in this domain. The other type of failure involves not finding the information requested in the scene model. In this case, rather than simply responding "not present",

the system may ask the user whether a new view of the scene should be analyzed.

5. Discussion

This paper has presented LandScan, a prototype integrated system under development that covers most of the different levels of vision and natural language processing. It may be used both to guide the low-level vision processing, and to provide communication of visual information to a user. While LandScan is not complete in the sense that all of it is successfully implemented, it provides a computational model for a vision system guided by natural language.

In summary, the data-driven subsystem of LandScan takes stereo images and builds a surface graph representing three-dimensional geometric and topological attributes. The query-driven modules recognize objects and build a Scene Model which represents the user's interest in the image.

The natural language interface which uses the scene representation still has to be designed. It must be able to apply locative linguistic constructs to some representation of visual data and reason about this data. When this is operative, the scene analysis will be truly query-driven and the goals of the system will have been reached.

References

- [Bates 51] Bates, Madekine. *The Theory and Practice of Augmented Transition Network Grammars*. In Leonard Bole (editor), *Natural Language Communication with Computers*. Springer-Verlag, 1081.
- [Fischler 83] Martin Fischler. *Image Understanding Research and its Application to Cartography and Computer-Based Analysis of Aerial Imagery*. Technical Report, SRI International, September, 1983.
- [Herman 83] Herman, Martin, Takeo Kanade, Shigeru Kuroe. *The 3D MOSAIC Scene Understanding System*. In *Proceeding of the 8th International Joint Conference on Artificial Intelligence*. 1983.
- [Herskevsts 82] Herskovits, Annette. *Space and the Prepositions in English: Regularities and Irregularities in a Complex Domain*. PhD thesis, Department of Linguistics, Stanford University, 1982.
- [Krotkov 81] Krotkov, Eric. *Construction of a Three Dimensions Surface Model*. Technical Report MS-CIS-81-JO, CIS Department, University of Pennsylvania, 1984.
- [Talmy 83] Talmy, Leonard. *How language Structures Space*. Technical Report 4, Berkeley Cognitive Science Report, January, 1983.
- [Tropf 83] Tropf and Walters. *An ATN for 3-D Recognition of Solids in Single Images*. In *Proceedings of the 5th International Joint Conference on Artificial Intelligence*. 1983.
- [Winograd 83] Winograd, Terry. *Language as a Cognitive Process*. Addison-Wesley Publishing Co., 1983.
- [Zwarico 84] Amy Zwarico. *The Recognition and Representation of 3D Images for A Natural language Driven Scene Analyzer*. Technical Report MS-CIS-84-29, University of Pennsylvania, 1984.

driver involves looking closely at the end of the blade; the relatively localized *context* of the business end of the blade is established by the grosser levels of the hierarchy, where it is recognized (for example) that the tool is not a hammer or wrench. In this way, the Marr-Nishihara proposal tends (heuristically) to relate large scale geometric structure to gross functional use.

- A-kind-of hierarchies

Family hierarchies are ubiquitous, and apply as much to visual shape representations as to the more cognitive situations in which they were developed in Artificial Intelligence. *ACRONYM* represents the fact that the sets of B747-SPs, B747s, wide-bodied jets, jets, and aircraft, are ordered by subset inclusion. Similarly, a claw hammer is a-kind-of framing hammer, which is a-kind-of hammer. In general, a subset hierarchy is a partially-ordered set, but not a tree. From the domain of tools, for example, a shingle ax is both a-kind-of ax, and a-kind-of hammer.

3. Learning

The commonest form of inductive generalization used to learn concepts from positive examples is the *drop condition* heuristic (Dietterich and Michalski 1981, Winston 1984, page 398). This is the method used in our program. Through careful design of the representation the method has been extended to allow generalizations of intervals and structural graphs.

The idea behind the heuristic is that if two things belong to the same class then the differences between them must be irrelevant. Accordingly, when we have a partial model of a concept and receive a new example, we modify the model by deleting all the differences between it and the example. This can be seen by comparing Figure 2b with Figure 3b. Notice that the network in Figure 3 puts very little constraint on the size or shape of the head. This is because the shapes of the heads in the examples vary widely. For instance, the heads of the first and third hammer are straight while the head of the second hammer is curved. Note also that the manner in which the handle joins the head is only loosely specified. This is because the handle is joined to the side of the head in the first two examples but to the end of the head in the third example.

This is a simplified explanation of the learning algorithm. The matching involved is not graph isomorphism nor is it, merely counting the number of required features an object has. Rather it is a complex local matching scheme. Consider using the semantic net shown in Figure 1 as the model for the *airplane* concept. For an object to match this model, at the top level it must have three pieces which look similar to the three in the model. A piece of the example is similar to the wing model if, first of all, it has the shape specified in the network and, second, it has two things which look like engines attached to it. Suppose that a certain piece has the right shape for a wing but has only one engine attached to it. At the level

of the wing model the program notices that there is a discrepancy yet judges that the piece is still close enough to the description to be called a wing. When the top level of the matcher asks if the piece in question looks like a wing the answer is "yes". No mention is made of the fact that the wing is missing an engine. The difference only matters locally and is isolated from the higher levels of matching.

Another important concern is limiting the scope of generalizations made. Imagine that the program is shown a positive example that is substantially different from its current model. Altering the model by the usual induction heuristics typically leads to gross over-generalization. This, in turn, runs counter to what Winston [1984, page 401] has dubbed *Martin's law*, namely: learning should proceed in small steps. Therefore our program creates a new, separate model based on the new example, splitting the concept being taught into a disjunction.

In some cases, the disjunction will be replaced by a single model as positive examples are taught that are intermediate to the disjuncts. For example, suppose that the first example of a hammer shown to the program is a claw hammer, and that the second is a sledge hammer. The program will create a disjunction as its concept of hammer, but it will be consolidated into a single model once it has seen such examples as a mallet and ballpeen hammer.

Even though the program only generalizes a concept using an example that is structurally similar, it is sometimes deceived and must recover from over-generalization. We follow Winston [1984] and provide censors that override the offending rule. Censors can be generalized and there can be disjunctive censors; in fact this is the usual case. Since censors can be generalized they also have the possibility of being over-generalized. This is countered by putting censors on the censors. In general, a concept is not represented by a single model but by a group of models. There can be several positive models corresponding to the disjuncts as well as several negative non-models summarizing the exceptions to the other models.

4. Current Work

The goals of our research are not limited to learning. The work reported here forms part of the *Mechanic's Mate* project [Brady, Agre, Braunegg, and Conneil 1984], which is intended to assist a handyman in generic assembly and construction tasks. The primary goal of that project is to understand the interplay between reasoning that involves tools and fasteners and representations of their shape.

For example, instead of learning that a certain geometric structure is called a hammer, we learn that something which has a graspable portion and a striking surface can be used as a hammer. These two functional concepts are then defined geometrically in terms of the shape representation. Reasoning from function as well as from form

allows more flexibility. For instance, faced with a hammering task, but no hammer, one might try mapping the hammer structure onto that of any available tool. A screw driver provides a good match, identifying the blade of a screw driver with the handle of the hammer, and the (assumed flat) side of the screw driver handle with the striking surface of the head of the hammer. In this way, the Mechanic's Mate can suggest improvisations, like using a screw driver as a hammer.

Our initial goal was to learn shape models cast in the representation described previously. Eventually, the *Mechanic's Mate* will have to learn about the non-geometric properties of objects: weight, material type, and the processes that use them. Currently we are using Katz's English interface [Katz and Winston 1983] to tell our program such things. This is not satisfactory. Instead, we hope to teach dynamic information using a robot arm and hand.

Another area of interest is inducing structural subclasses from examples. Since the subclasses that form the a-kind-of hierarchy are an important part of the shape representation, they should be learnable. However, in learning subclasses there is a danger of combinatorial explosion. Learning subclasses requires a suitable similarity metric. Feature-based pattern recognition systems learn subclasses as clusters in feature space, and clusters are sets that are dense with respect to the Euclidean metric. Part of our research in learning shape descriptions has been to determine what makes objects look similar. This suggests using the metric employed in the learning procedure to form subclasses through a process analogous to feature space clustering. This is the focus of our current work.

5. Acknowledgements

This report describes research done at the Artificial Intelligence Laboratory of the Massachusetts Institute of Technology. Support for the Laboratory's Artificial Intelligence research is provided in part by the the System Development Foundation, the Advanced Research Projects Agency of the Department of Defense under Office of Naval Research contract N00014-80-O-0505, and the Office of Naval Research under contract number N00014-77-C-0389. We thank the people who have commented on the ideas presented in this paper, particularly Phil Agre, Steve Bagley, Hob Berwick, Hen DuBoulay, Alan Bundy, Margaret Fleck, Scott Heide, Boris Katz, Tomas Lozano-Perez, John Mallery, Tom Mitchell, Sharon Salveter, Dan Weld, and Patrick Winston.

References

- Brady, Michael, and Asada, Haruo, [1984], "Smoothed Local Symmetries and their Implementation," *Int. J. Robotics Research*, 3 (3) .
- Brady, Michael, Agre, Philip, Braunegg, David J., and Connell, Jonathan H., [1984], *The Mechanic's Mate*, ECAI 84: Advances in Artificial Intelligence, O'Shea, T. (ed.), Elsevier Science Publishers B.V. (North-Holland).
- Brooks, Rodney A., [1981], "Symbolic Reasoning Among 3-D Models and 2-D Images," *Artif. Intell.*, 17, 285 - 348 .
- Brooks, Rodney A., and Binford, Thomas O., [1980], Representing and Reasoning About Partially Specified Scenes, *Proceedings, DARPA Image Understanding Workshop*, Baumann, Lee S. (ed.), Science Applications Inc., 150 156.
- Canny, John Francis, [1983], Finding Edges and Lines in Images, MIT Artificial Intelligence Laboratory, Cambridge Mass., AI-TR-720.
- Connell, Jonathan H., [1985], forthcoming SM thesis, MIT Department of Electrical Engineering.
- Dietterich, T. G., and Michalski, R. S., [1981], "Inductive Learning of Structural Descriptions," *Artif. Intell.*, 16 .
- Doyle, Richard J., and Katz, B., [1985], Exploring the Boundary Between Natural Language and Knowledge Representation, MIT Artificial Intelligence Laboratory, Forthcoming AI Memo.
- Heide, S., [1984], A Hierarchical Representation of Shape, SM thesis, MIT Department of Mechanical Engineering.
- Katz, Boris, and Winston, Patrick H., [1983], A Two-way Natural Language Interface, *Integrated Interactive Computing Systems* P. Degano and Erik Sandewall (eds.), North-Holland, Amsterdam.
- Marr, D., and Nishihara, H. K., [1978], "Representation and Recognition of the Spatial Organisation of Three Dimensional Shapes," *Proc. Roy. Soc. Lond.* ii, 200 , 269 - 294.
- Winston, Patrick H., [1980], "Learning and Reasoning by Analogy," *Comm. ACM*, 23 , 689 - 703.
- Winston, Patrick H., [1981], "Learning New Principles from Precedents and Exercises," *Artif. Intell.*, 19 , 321 - 350.
- Winston, Patrick H., [1982], Learning by Augmenting Rules and Accumulating Censors, MIT Artificial Intelligence Laboratory, AIM-678.
- Winston, Patrick H., [1984], *Artificial Intelligence*, 2nd. Ed., Addison-Wesley, Reading, Ma..
- Winston, Patrick H., Binford, Thomas O. , Katz, B., and Lowry, M., [1984], Learning Physical Descriptions from Functional Definitions, examples, and precedents, Robotics Research, Michael Brady and Richard Paul (eds.), MIT Press, Cambridge, 117 - 135.