

D. J. Burr and R. T. Chien
 Coordinated Science Laboratory
 University of Illinois at Urbana-Champaign
 Urbana, Illinois 60801

The goal of this work is the specification of a system which extracts 3-D features from images by stereo comparison, and which recognizes objects by comparing these features to geometric models of objects. Models are encoded from hand measurements and consist of piecewise-linear wire frame structures. The work is in progress and some partial results are shown.

Bulk correlation methods [4] are used with three images of a scene to insure accurate triangulation for smooth directed edges. The following pictures are obtained: 1. center. 2. scene rotated 20° east. 3. scene rotated 10° north. The program compares center and east pictures for vertical edges and center and north pictures for horizontal edges.

Edge orientation is computed after preprocessing the center picture to find edge chains. An operator developed by Burr [1,2] finds local gradient extrema and tracks them by searching near neighbors (Fig. 1a). These chains are approximated by line segments [3] in Fig. 1b. The edge direction at a node, or line end, is taken as the average of the line direction on each side weighted by its length.

Stereo correlation is implemented as a mean-square difference of image intensity over 9 x 9 pixel windows. Search for the minimum in the corresponding view is restricted to a single line segment [4], and further, to locations where the intensity gradient exceeds a threshold. Triangulation determines a 3-D location for each node of Fig. 1b. A separate program attempts to match this 3-D structure (Fig. 2a) to a model.

Model matching proceeds in two steps as follows: (1) A search is made to find a model edge and perceived edge whose lengths agree. If found, each defines the z-axis of a cylindrical coordinate system (r,θ,z). The rotational, or θ-ambiguity, is relieved by finding an additional edge pair which have nearly equal lengths, center positions (r_c,z_c), and (r,θ,z) direction cosines. The position of this feature relative to each z-axis now defines the y-axis direction of a cartesian reference frame for both model and perceived structure. The x-axes are just y x z.

(2) There may be further implied edge matches due to similarity of x-y-z edge coordinates and directions relative to each reference frame. An attempt is made to measure the confidence of this coordinate transformation by finding these implied edge matches. The criterion for an edge match is — perceived edge length < model edge length,

"This work was supported by the Joint Services Electronics Program (U.S. Army, U.S. Navy, and J.S. Air Force) under Contract DAAB-07-72-C-0259.

perceived edge center lies within a tolerance cylinder about the model edge, and direction cosines agree. The confidence is thus the total length of all matched edges divided by the total length of all scene edges. If this confidence exceeds the last computed value then it becomes the new threshold, and subsequent proposals must either exceed it or be rejected. The match of highest confidence is taken as the correct one, and the model is rotated and projected into the center view for observation (Fig. 2b).

Further work is being done on efficient tracking of 3-D edges and piecewise circular representations for 3-D edges. The ultimate goal is automatic identification of an object in an occluded scene from a model data base.

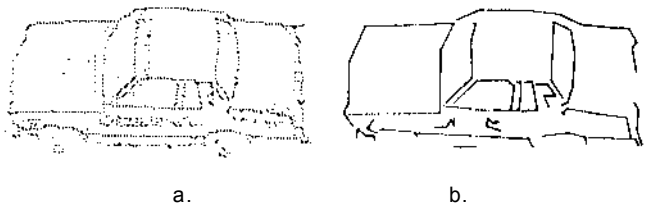


Figure 1. a. Thresholded gradient extrema of car image. b. Linear approximation of edges in 1a.

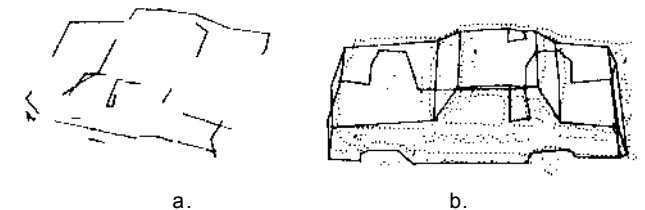


Figure 2. a. Rotated view of Fig. 1b after 3-D computation and removal of highly sloped depth edges. b. Best match of car model to 3-D scene (confidence=.713, PDP-10 run timer=13 sec).

References

1. Burr, D. J. and R. T. Chien, "The Minimal Spanning Tree in Visual Data Segmentation," Proc. 3rd IJCP, San Diego, California, Nov., 1976,
2. Chien, R. T. and L. J. Peterson, "Image Compression by Feature Extraction and Reconstruction," Proc. Workshop on Picture Data Description and Management, Chicago, 111., April 21-22, 1977.
3. Duda, R. O. and P. E. Hart, Pattern Classification and Scene Analysis, P, 338, John Wiley and Sons, New York, 1973-
4. Hannah, M. J., "Computer Matching of Areas in Stereo Images," Stanford AI Memo AIM-239, Stanford University, July, 1974.

This report describes ongoing research on a working system which drives a vehicle through cluttered environments under computer control, guided by images perceived through an onboard tv camera. The emphasis is on reliable and fast low level visual techniques which determine the existence and location of objects in the world, but do not identify them. Included are an interest operator for choosing distinctive regions in images, a correlator for finding matching regions in similar images, a camera solver which determines camera displacement and distance to objects from stereo information (Gennery, D.B., this Proceedings) and an automatic geometric distortion corrector for camera nonlinearities. Many of these use pictures reduced in linear dimension by powers of 2 by summation of pixels. Other operators are a high pass filter, a point noise remover, a contrast normalizes a vertical roll corrector, a picture comparator and an operator for reducing pictures by other than powers of two.

Our hardware includes an electric vehicle, called the cart, remotely controlled over a CB radio link by a PDP-KL10. It carries a b/w tv camera whose picture is broadcast over a UHF channel, and occasionally digitized by the computer. It has motors for the wheels, steering and camera pan. Each can be made to run forward or backward. There are potentiometers on the steering and pan which enables them to be commanded to point straight ahead.

Budgetary and personnel limitations have resulted in crude mechanical arrangements. The motor speeds are poorly regulated, and video is the only feedback to the computer. Dead reckoning errors are about 30%. Our small resources have been spent gaining software experience before undertaking serious hardware work. In my opinion our major hardware limitation is one shared by all other vision work, and AI in general, namely a critical shortage of raw processing power. For instance it would take about 100,000 efficiently programmed PDP-10's to match the human visual system.

Results

Early versions of the routines described below were used in a program which drove the vehicle in straight lines or uniform arcs. It acquired and tracked distant features, using their motion from frame to frame to build up a model of vehicle response, and to servo on the desired path. With the cart outdoors on a dirty road, it worked well. Ten runs of about 60 steps were completed. The runs were usually terminated by serious hiccups of the radio control link. Each step took the cart 2 feet forward, and used 30 compute seconds. The radio link has since been much improved.

Several runs involving the distortion corrector, camera solver and new versions of the interest operator and correlator have been completed. The new program tries to determine the distance to the features by applying the camera solver after tracking them through several images. The performance is poor. The camera solver results are erratic, seemingly due to the degenerate nature of the solution. Objects lying near the camera axis (most of the scene) provide no depth information.

Next

We are trying a new approach, replacing the camera pan mechanism with one which provides 21 inches of side to side motion, in three 7 in. steps. This should provide adequate parallax, and also close spacing to make the correlations easy. Since the camera motion parameters will be known the correlation searches become one dimensional, and an absolute scale factor is known. The camera solving is also easy. The idea is to locate nearby features in 30 at each vehicle stop. The vehicle motion can be found from the apparent feature motions between stops. The location of the ground can be deduced from the camera height and orientation.

Interest Operator

This routine is used to acquire new features in a scene. It selects a relatively uniform scattering, to minimize the probability of missing important obstacles, and chooses distinctive areas for unambiguous correlation. This is achieved by returning regions which are local maxima of a directional variance measure. Featureless areas and simple edges (which have no variance in the direction of the edge) are thus avoided.

Directional variance is measured over small square overlapping windows of specified size (typ. 4x4 to 8x8). Sums of squares of differences of pixels adjacent in each of four directions (horizontal, vertical and two diagonals) over the window are obtained. The variance of the window is the minimum of these four sums.

The operator is applied to a reduced version of the picture, where the specified window size shrinks to 2 or 3 pixels. Noise sensitivity is reduced and speed increased. Partly hand coded, the routine takes 75 ms for a 260x240 image, with 8x8 windows.

Given a feature in one picture, the correlator attempts to find the matching region in another image. It takes the position in the first picture, a rectangular search area (often the whole image) in the second picture, and a feature window size n.

The search uses a coarse to fine strategy, which begins in reduced versions of the pictures. The order of reduction is chosen to shrink the smaller dimension of the search rectangle to between n and 2n pixels. An n by n window in the shrunken source image, centered on the desired feature, is considered. It covers about 25% of this tiny version of the picture. A correlation coefficient is calculated for each possible placement of this window on the search area. For a search area exactly 2n by 2n, there are $(n+1)^2$ positions. The one with the highest coefficient becomes the search area for the next level of refinement.

This is repeated with pictures reduced one step less, i.e. linearly twice as large. An n by n window is again centered around the location of the feature, and is searched for in the best matching window from the previous search, which expands to 2n by 2n at the new reduction. This goes on in successively larger versions of the pictures until an n by n window is matched in the unreduced images. There are about $\log_2(w/n)$ searches in all, where w is the smaller dimension of the search rectangle in the unreduced picture.

This approach has advantages over a simple pass of a correlation coefficient. It needs only 1/150 the number of pixel comparisons to find an 8x8 window in a 256x256 picture (smaller advantage for smaller searches). The simple method comparisons are without context, and a match may be found in totally unrelated parts of the image. In our technique coarse structure guides the higher resolution comparisons, and further speedup is possible because smaller windows work. The searches at coarse levels rarely fail, possibly because noise and distortions are reduced by reduction.

The correlation measure used, designed to have limited contrast sensitivity, was obtained by multiplying the normalized correlation coefficient by twice the cosine of the angle with the line a-b. It is: $2\frac{ab(\Sigma a^2 + \Sigma b^2)}{\Sigma a^2 \Sigma b^2}$. Normalized correlation is the sum of the pairwise products of a and b divided by the geometric mean of the sum of their squares. The new measure, referred to as pseudo-normalized correlation, is the sum of the products divided by the arithmetic mean of the sums of the squares.

By in-line coding the source window and using a table of squares the bulk of the correlation is done in 3 instructions per pixel comparison. An 8x8 window is found in a 260x280 area in 75 ms. The error rate is 10% on interest operator selected features. Typical image pairs are taken two feet apart with a 60 degree lens.

Scale Changes

As the vehicle moves the image it sees changes. The major element of this transformation is an enlargement of nearby objects. We have tried correlating across images reduced by different geometric scale factors by generating pictures $2^{2/3}$ as large as each of the binary steps. We obtain effective scale changes of 1, $2^{1/3}$, $2^{2/3}$ and 2 by comparing various combinations of reductions of the first and second images. The results are disappointing. The method often introduces as many new errors as it corrects. Experiments in applying it more selectively are planned.

Camera Distortion Correction

Electron optics tend to have geometric distortions undesirable when using a camera as a measuring instrument. We have written a camera calibration program which is given an image of a square array of black spots on a white background, and told the array to lens center/spot spacing distance ratio. It computes a polynomial for transforming feature image positions accurately to angle in space.

It tolerates a wide range of image sizes (3 to 12 spots across) and illumination and arbitrary rotation. After intense fiddling with a training set of 20 images, it has worked without error on 80 widely differing new images. Our test pattern is a ten foot square painted on a wall, with two inch spots at one foot intervals.

The algorithm gets an image of such an array, and finds four major peaks in the magnitude of the fourier transform of a reduced version of it, to find its rotation and spacing. The Interest operator is used to find a starting spot, and a special operator, which does local thresholding and finds centroids and moments of black areas, pinpoints all the spots, guided by the rotation/spacing information. A fourth degree least squares polynomial in two variables relating the actual to the ideal position of the spots is then generated.

Acknowledgement: This work was supported in part by contracts and grants from ARPA, NASA and NSF.