



HAL
open science

Still life stereo

Tony Kasvand

► **To cite this version:**

| Tony Kasvand. Still life stereo. [Research Report] RR-1372, INRIA. 1991. inria-00075189

HAL Id: inria-00075189

<https://inria.hal.science/inria-00075189>

Submitted on 24 May 2006

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

IRIA

UNITÉ DE RECHERCHE
IRIA-ROQUENCOURT

Institut National
de Recherche
en Informatique
et en Automatique

Domaine de Voluceau
Roquencourt
BP 105
78153 Le Chesnay Cedex
France
Tél.: (1) 39 63 55 11

Rapports de Recherche

N° 1372

Programme 4
Robotique, Image et Vision

STILL LIFE STEREO

Tony KASVAND

Janvier 1991



★ R R - 1 3 7 2 ★

STILL-LIFE STEREO

A summary of studies carried out during Sept. 1989 to
July 1990 as "Professor Invite" at Project Syntim.

INRIA
Domaine de Voluceau, Roquencourt
78153 Le Chesnay cedex, France

T. Kasvand

July 1990

STEREOVISION INTEMPORELLE

STILL LIFE STEREO

Tony KASVAND
Professeur Invité
Concordia University
Montreal, CANADA

RESUME :

Reconsidérant le problème de la stéréovision dans sa totalité, l'auteur résume ici les études qu'il a menées à cette occasion, lors de son séjour d'1 an à l'INRIA comme professeur invité.

Par l'application judicieuse de techniques existantes de traitement et d'analyse d'images, on espère avoir démontré que le "Problème de correspondance" peut être résolu de diverses manières. Plus particulièrement, la mise en correspondance ne nécessite pas de "connaissances extrinsèques" du type "Modèles d'Objets" ou "Représentation du Monde", et aucune n'a été utilisée pour cette étude.

L'approche présentée est purement expérimentale, dans la mesure où les techniques utilisées sont bien connues. Le but est plutôt de reformuler le(s) problème(s) et les contraintes qui y sont attachées, puis d'analyser le résultat des expérimentations et faire partager le bénéfice des conclusions tirées.

Mots-clés: *analyse de scènes, analyse d'images, approximation analytique, contraintes de similarité, mise en correspondance, reconstruction 3D, relaxation, reconnaissance de formes, segmentation, stéréovision, traitement d'images, vision par ordinateur*

SUMMARY :

The report is a summary of studies carried out while rethinking the entire 3D Stereovision problem, during a one year long stay as Professeur Invité.

By judicious applications of existing image processing and analysis techniques, it is hopefully, partially demonstrated that the so called "Correspondance Problem" can be solved in several ways; especially that, contrary to presently popular opinions, no "external knowledge" such as "Object or World Models" is required to establish correspondance, and none are used in the report.

The present approach is purely experimental, since the techniques are well known. Its aim is rather to reformulate some problem(s) and the related requirements, then analyse and share results of carried out experiments.

Keywords: *analytical approximation, computer vision, image processing, image analysis, matching, pattern recognition, relaxation, segmentation, scene analysis, stereovision, similarity constraints, 3D reconstruction*

ABSTRACT

Measurement of distances to objects in a scene is one of the many practical applications of artificial or computer vision. Distances to objects enable a mobile robot to establish its position, it is necessary for avoiding collisions, and it allows the robot to carry out productive tasks. In other words, vision helps the robot to "see" where it is and what others are doing.

A method of measuring distance is by using stereo vision, exactly like we ourselves do. Two images are required which are taken from slightly different positions (the stereo pair of images), for example, the images from each of our two eyes, from two TV cameras, or by using two photographs or digitized images ("still-life" or "nature morte" stereo). The technical name for this field is (close-range) photogrammetry. Photogrammetry has a long history, for example, in the 1860's Colonel Aime Laussedat used photogrammetry to survey Paris, France. The principle, summarized in one sentence, is: "Given a point in the left image of the stereo pair and a point in the right image, such that these two points correspond to the same point in the scene, then the distance to the point in the scene can be computed from the measurements of the positions of the two points in the images". To find these two points in the images is called the "correspondence problem". Automated machinery exists for establishing correspondence but in close-range photogrammetry the correspondence is mainly determined by a human operator. The automatic methods used are insufficient (cross-correlation and/or Fourier analysis).

By judicious application of existing image processing and analysis techniques it is, hopefully, partially demonstrated in this report that the "correspondence problem" can be solved in several ways. It is mainly a question of our willingness to carry out the necessary work and to design and program suitable computers to make the methods practical. Contrary to presently popular opinions, no "scene understanding", no "object recognition", and no "object or world models" ("external knowledge") are required to establish correspondence and none are used in the present report. The availability of "models" greatly simplifies the problem but they introduce unrealistic constraints if the models are obtained from "external sources". The production of "models" as an integral part of the image analysis procedures is beyond present experimentation and largely beyond contemporary thinking.

The present work is purely experimental since the techniques are well known. Detailed results are shown at the different stages of processing to establish correspondence. The form of the report is: Problem, solution, results; new problem, solution, results; The results are shown in image form but most of the images may carry meaning only for those versed in various aspects of image analysis. Numerous one character per pixel

prints, gray level prints, and stereo pairs are used to illustrate the results. Unfortunately, the gray level prints, especially after copying, look "dismal". For those interested, references are given to computer files from which the images may be displayed or printed. The entire report should only be considered as "introductory reading for dispelling a mystery" for students wanting to study image analysis applied to 3D vision. Research, after all, is a learning experience: Make experiments, see what happens, learn from them in order not to repeat the same mistake twice, and inform others "where not to step".

Chapter 1 briefly discusses some perfectly obvious aspects of stereo images which appear to have been ignored by many or considered as being too trivial. A "similarity principle" is defined, which simply states that "the images of a stereo pair are very similar". The rest of the chapter is devoted to what should be computed from the images and how the correspondence (matching) is to be accomplished. A large number of possibilities are pointed out.

The work starts in Chapter 2 with image processing requirements, followed by grouping of "pixels" to simplify the detection of similarities between the left and right images of the stereo pair. Classical techniques are used, consisting of pixel feature computations, classification of features to create homogeneous regions in the image space, and analytic approximation of the regions followed by "relaxation". (The processes and decision spaces required for images of thin line structures and textures are not included.) Great care is taken to only use features that preserve "stereo fidelity", i.e., images of the features can be seen in stereo with a pair of stereo glasses. It is highly doubtful that features that "ruin the stereo effect" will be of much use in subsequent analysis.

The similar regions found are put into approximate correspondence ("raw matches") in Chapter 3. Two techniques were studied in some depth, called "matched classification" and "direct and-ing" of the regions. Several other methods are mentioned but could not be studied in detail. Since the stereo pair of images used was not of high quality, an inordinate amount of effort was spent on correcting the dissimilarities between the images. This became more of an intellectual challenge than practical necessity since we can see the stereo effect despite the dissimilarities. In practice one would simply readjust the iris on the camera and take another picture. The approximate correspondence is then "adjusted" for better agreement with the information in the two images. In the preliminary experiments a second order polynomial is used for analytic approximation of the gray levels of the matched regions. This is followed by "relaxation". However, better analytic functions and more controlled relaxation processes are required for recovering the "stereo effect" to greater accuracy.

The report ends rather abruptly at this stage with a review and critique. Several reasons existed:

1. It became apparent where the "weak links" in the chain of reasoning were but there is no time to correct these at present.
2. The results "couple to and simplify" the presently popular "edge-based" stereo techniques since the combinatorial "explosion" problem is removed.
3. The procedure is iterative (hierarchical) where the matched facets obtained are simply viewed as "smaller images" to which the same techniques can be applied again, resulting in finer and finer correspondence between the details of the images.

A brief "illustrated summary" is given below. The correctness of the results is verified stereoscopically, i.e., stereo glasses are required to see the 3D effects in the images:

The original image pair is shown in Figures 1a and 1b. Even though it looks like a "room scene" to us, such "understanding" is not available to the machine, nor is it needed. The same stereo pair on a smaller scale is shown in Figure 1ab.

The scene is split into regions of pixels with homogeneous characteristics, see Figures 2a, 2b, and 2ab. (Seven gray levels or "colours" were used to "paint" the regions such that adjacent regions have different "colours".)

The homogeneous regions are matched and as many additional "matched" regions are found as is feasible, see Figures 3a, 3b, and 3ab. These are the "raw matched facets" ("painted" in six "colours"). The "raw" matches, when viewed stereoscopically, create a scene that appears to be partially "exploded". One set of (matched) facets are approximately in their "proper positions" in the scene, while others are "hanging in free space" in front of the scene". Different matching methods tended to produce different visual effects. However, the visually most disturbing effect is caused by minor differences in adjacency relations between the matched regions, i.e., if there is a small space between two facets in one image but not in the other, this space "cannot be matched" by our vision.

The analytic approximation and relaxation steps are an attempt to correct the raw matches. Some preliminary results are shown in Figures 4a, 4b, and 4ab (painted in six colours). As can be seen, the analytic approximation and relaxation techniques described correct the facets only to a limited degree but leave much "fine-grain" noise and may create occasional "gross errors". The small dissimilarities at the boundaries of matching facets produce "visual disturbances" when the results are viewed with

stereo glasses, since our vision is extremely sensitive to such small incompatibilities. In the image space most of these errors only involve a few pixels at a boundary and could be "written off" as "spatial quantization effects".

However, since such small differences are vital to our stereo vision, they are also important for any automated stereo vision system. The differences can be found if the scene is recreated from the analytic approximation. For example, the pixels that could not be assigned to any region are shown as "black zigzags" in Figures 5a, 5b, and 5ab. The pixels at which additional corrections are required are found by comparing the gray levels of the pixels in the original image and in the analytic approximation. For example, if the pixels in the analytic image are replaced by those in the original when the error is greater than 15 units or where the approximation does not exist, the "stereo effect" reappears, see Figures 6a, 6b, and 6ab. Thus, it is known where future corrections are to be made to achieve "full stereo fidelity".

Figure titles:

Figures 1a, 1b, 1ab: The stereo pair of images used for analysis. The left image (1a) and the right image (1b) form a stereo pair but no equipment existed for viewing such images. Hence, a smaller version of the same pair is shown in (1ab).

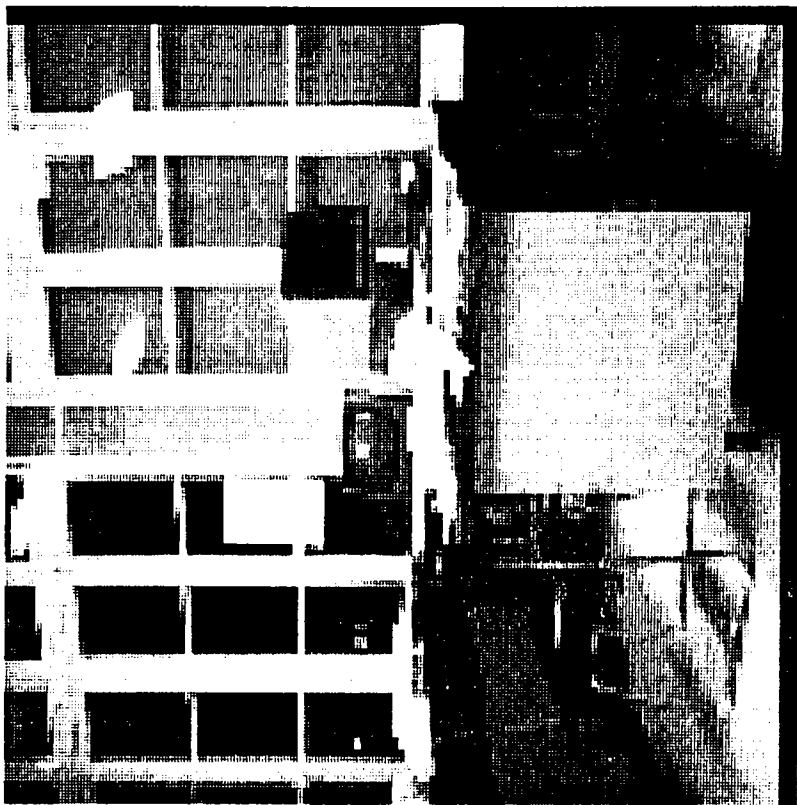
Figures 2a, 2b, 2ab: Regions of pixels with homogeneous characteristics "painted" in seven "colours" such that adjacent regions have different "colours". (2a) is the left image, (2b) the right image, and (2ab) shows a smaller version of the same pair.

Figures 3a, 3b, 3ab: The "raw matched facets" ("painted" in six "colours"). (3a) and (3b) are the left and right images, and (3ab) shows a smaller version of the same pair.

Figures 4a, 4b, 4ab: The matched facets after analytic approximation and relaxation ("painted" in six "colours"). (4a) and (4b) are the left and right images, respectively, and (4ab) is a smaller version of the same pair.

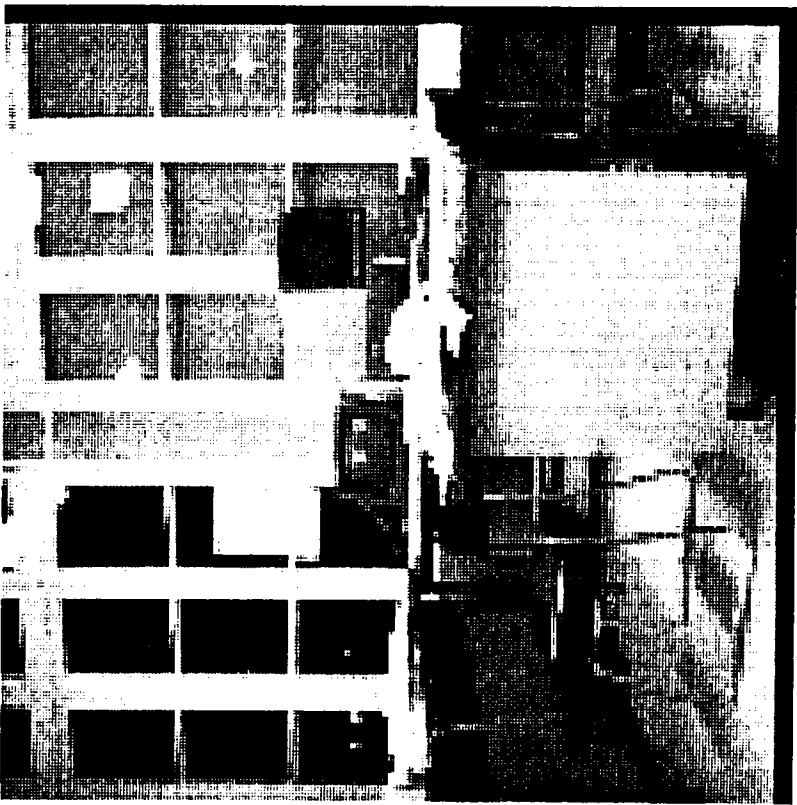
Figures 5a, 5b, 5ab: The scene in Figure 1 recreated from the analytic approximation after relaxation. Pixels where the approximation does not exist are shown as black "zigzags".

Figures 6a, 6b, 6ab: The scene in Figure 1 recreated from the analytic approximation after relaxation. Pixels where the approximation does not exist or where the error between the approximation and the original image is "too large" (>15 units) have been taken from the original image.



if8gray.fsh

Fig.: 1b.



if8gray.lsh

Fig.: 1a.

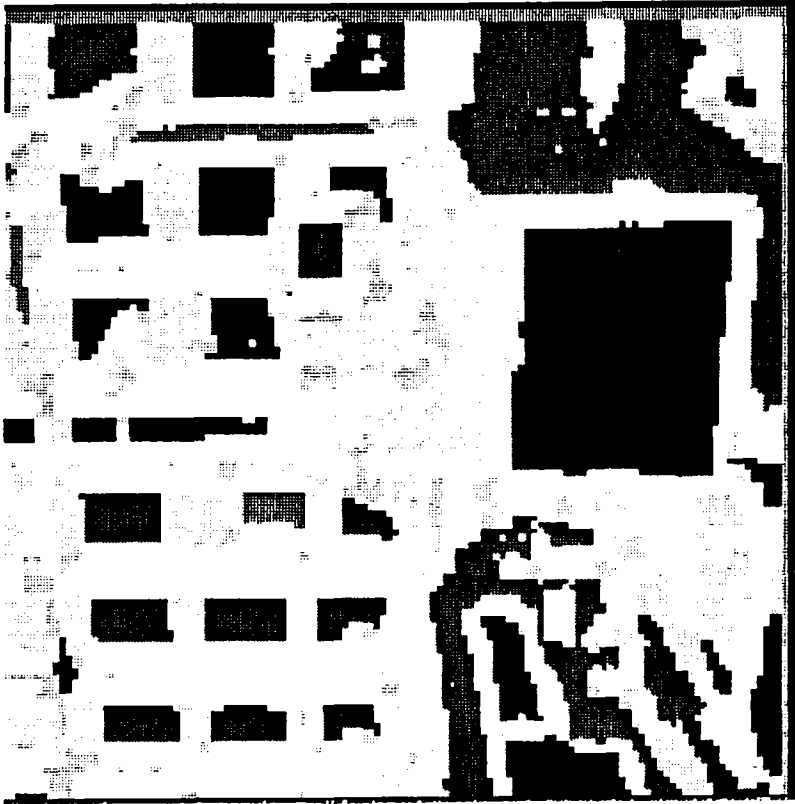


Fig.: 2b.

if8hkpp2.rgr

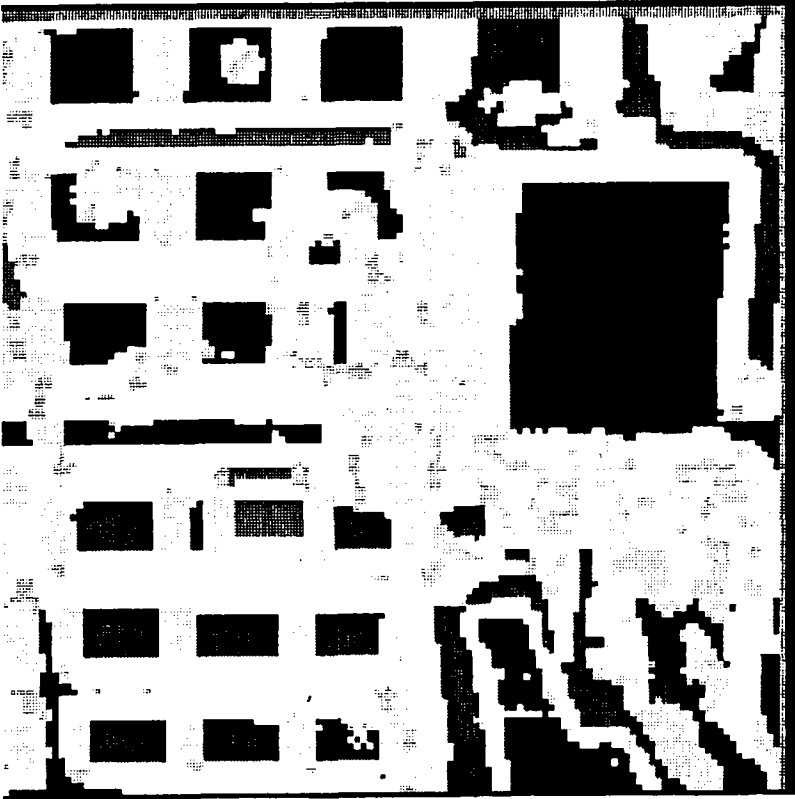


Fig.: 2a.

if8hkpp2.lgr



if8flkan3.rgr

Fig.: 3b.



if8flkan3.lgr

Fig.: 3a.

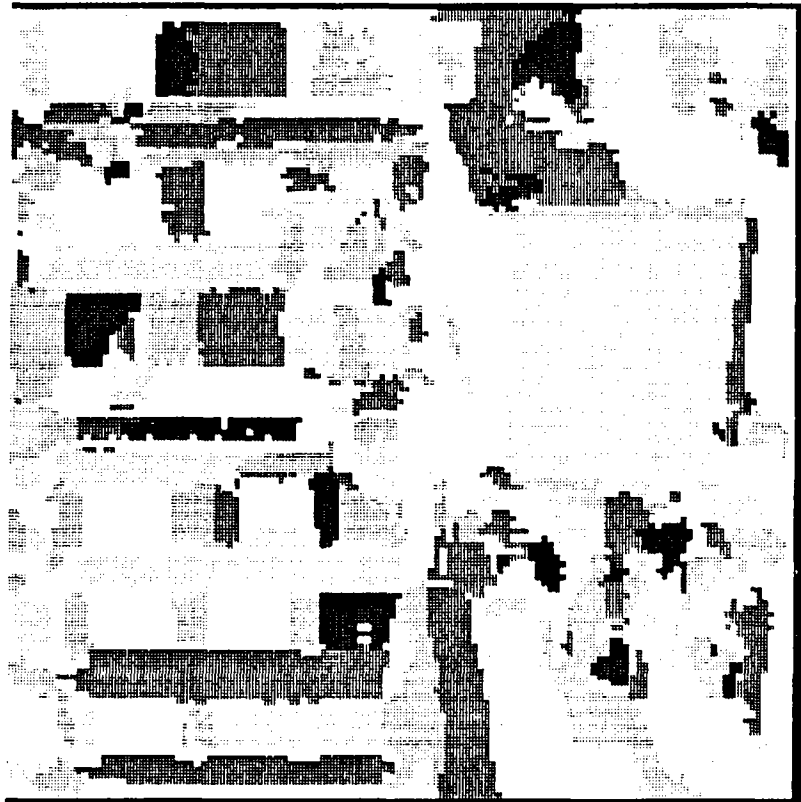


Fig.: 4b.

if8flkan4.rgr

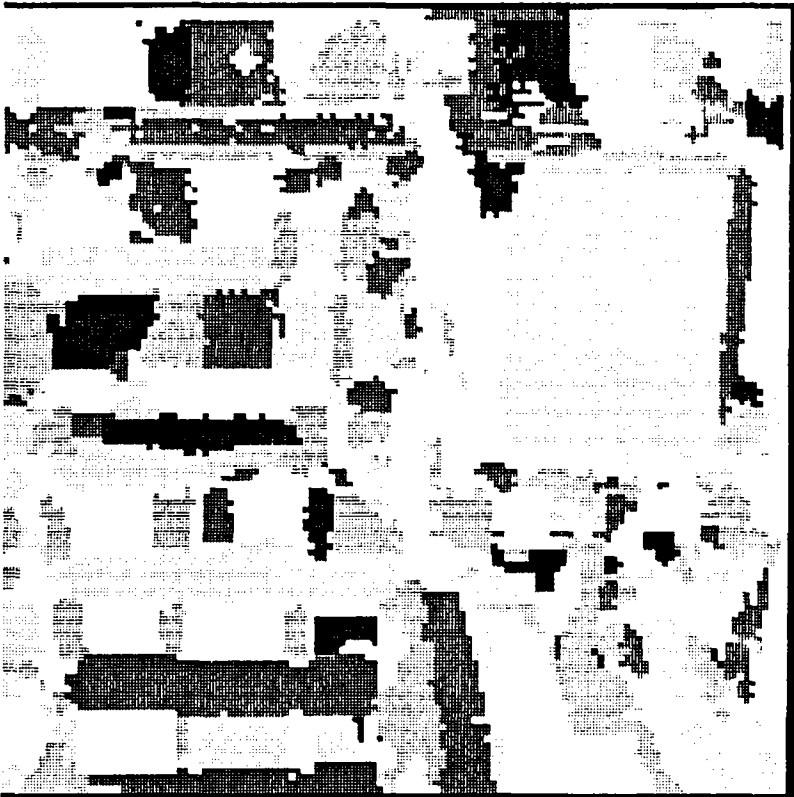


Fig.: 4a.

if8flkan4.lgr

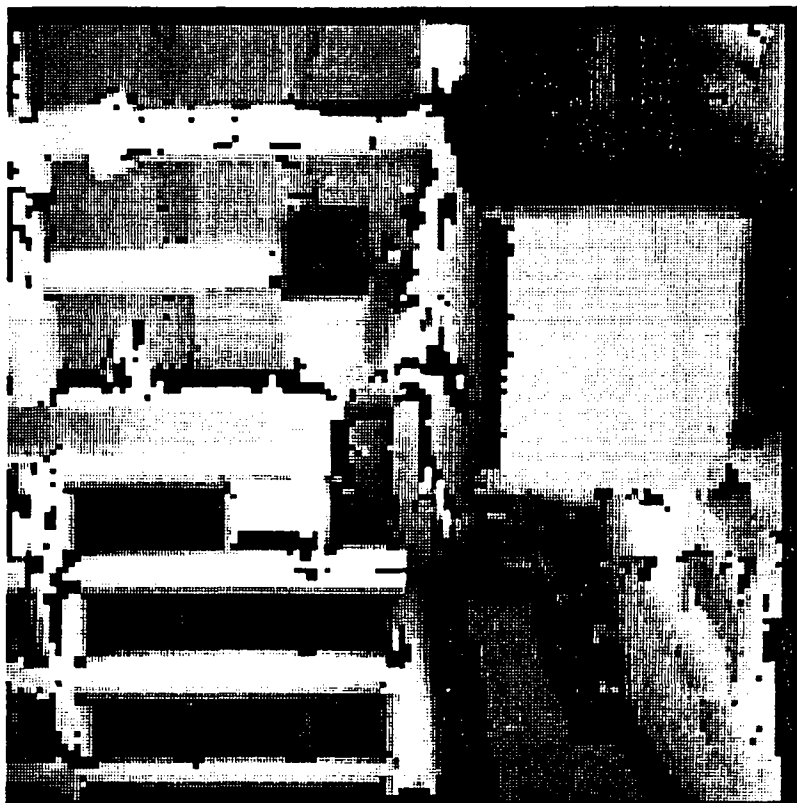


Fig.: 5b.

if8gacan9.rgr

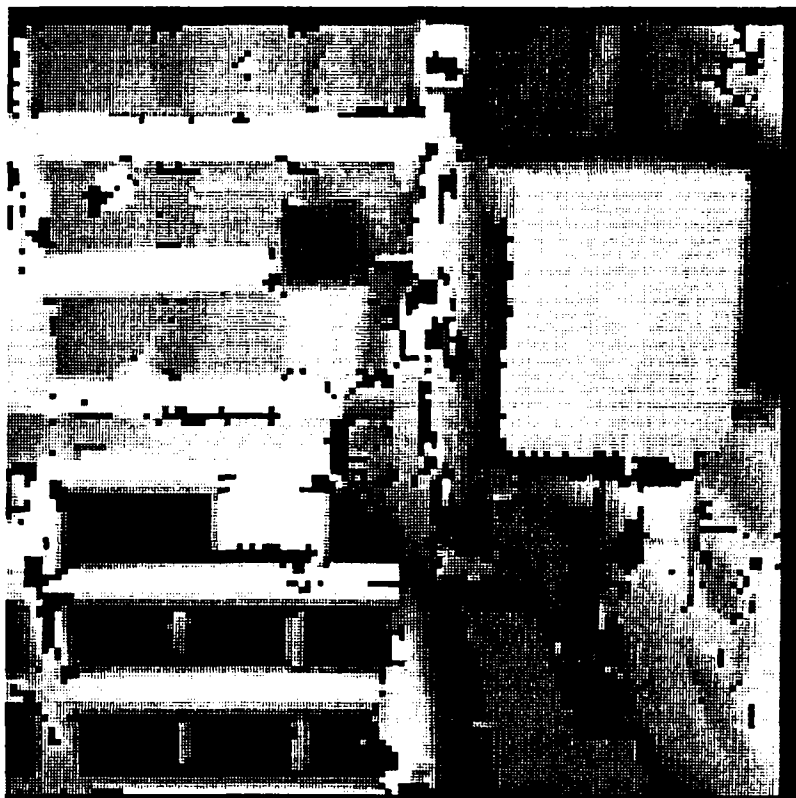
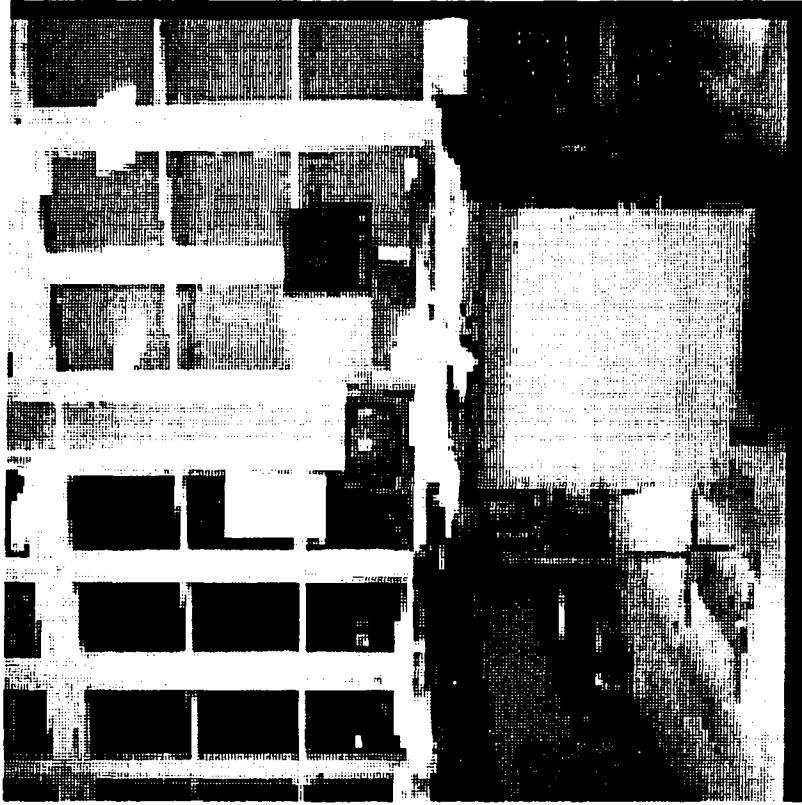
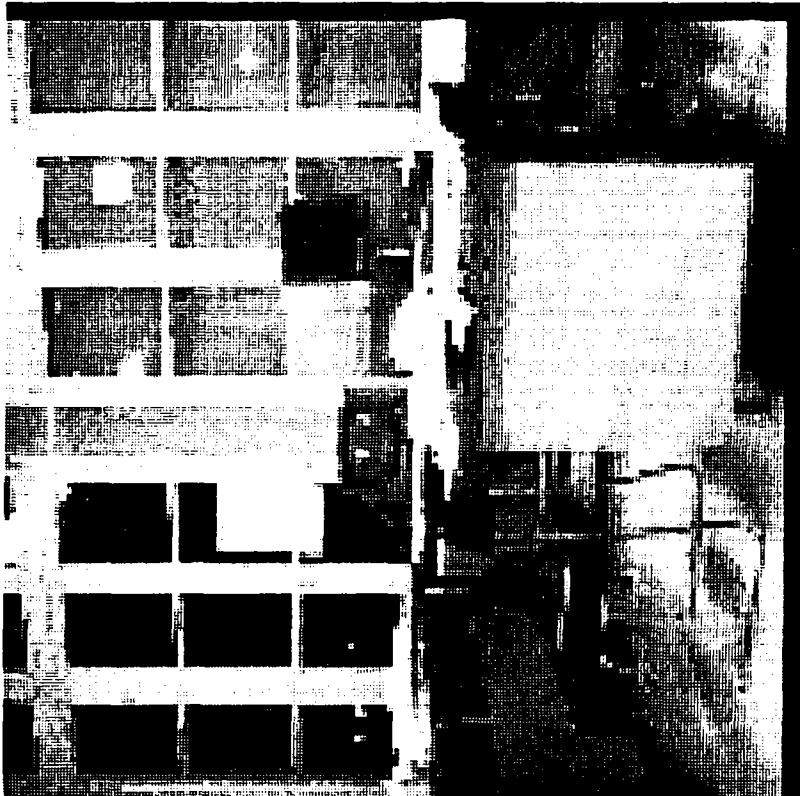


Fig.: 5a.

if8gacan9.lgr



if8qacaud.tgr Fig.: 6b.



if8qacana.lgr Fig.: 6a.



if8gray.lsh

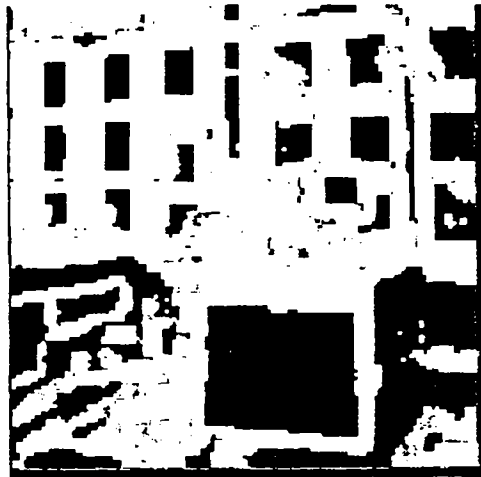


if8gray.rsh

Fig.: 1ab.

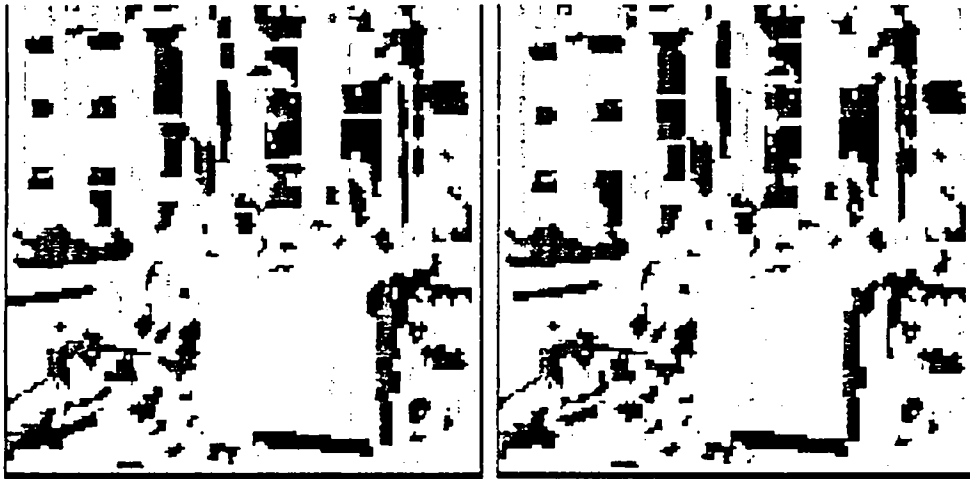


if8hkpp2.lgr



if8hkpp2.rgr

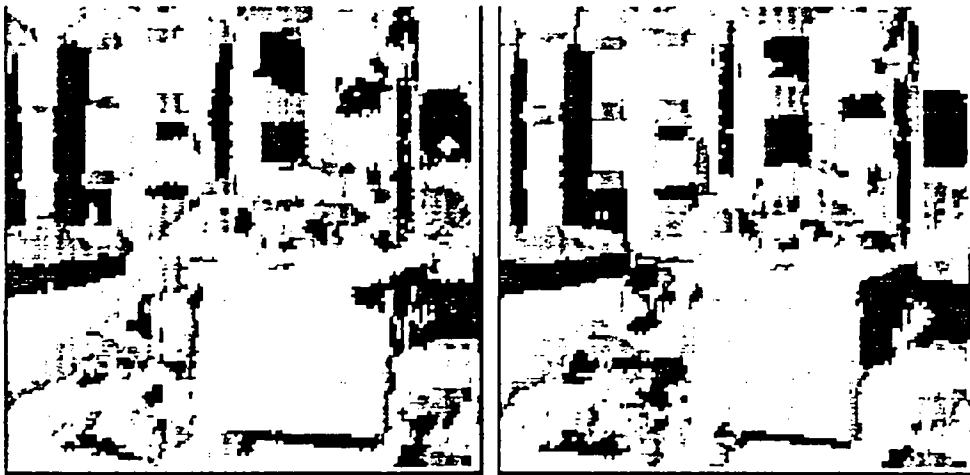
Fig.: 2ab.



if8flkan3.lgr

if8flkan3.rgr

Fig.: 3ab.



if8flkan4.lgr

if8flkan4.rgr

Fig.: 4ab.



if8gacan9.lgr



if8gacan9.rgr

Fig.: 5ab.



if8gacana.lgr



if8gacana.rgr

Fig.: 6ab.

INDEX

Abstract

Index

Acknowledgements

Chapter 1: THE OBVIOUS	page
1.1 Introduction	1
1.2 Seeing distance	4
1.3 The similarity principle	7
1.4 The hierarchical approach	8
1.4.1 The quantities "x" and "X"	10
1.4.2 Matching "X"	11
1.4.3 Comments	13
1.5 Validity requirements	14
1.5.1 The ideal stereo camera	14
1.5.2 The irreducibles	22
1.6 Conclusions	24
Chapter 2: PRELIMINARY COMPUTATIONS	
2.1 Introduction	31
2.2 The image pair	39
2.3 Computable and computed features	41
2.4 Hierarchical feature classification	50
2.5 Analytic approximation and relaxation	73
2.6 Super-regions versus scale-space	97
2.7 Conclusions	93
Chapter 3: FACET MATCHING	
3.1 Introduction	103
3.2 Matching methods	108
3.2.1 Matched classification	109
3.2.2 Direct "and-ing"	119
3.2.3 Correlation matching	130
3.2.4 Facet recognition matching	132
3.2.5 Pixel feature matching	135
3.2.6 Other methods	136
3.2.7 Comments	136
3.3 Modifying the various matches	139
3.3.1 The Lh, L&Rh, and Rh triplet	140
3.3.2 The raw Lm and Rm images	151
3.3.3 Misclassifications	155
3.3.4 Cut	159
3.3.5 Transplant	162
3.3.6 Clip	165
3.3.7 The gaps and the un-matchables	173
3.3.8 Comments	174
3.4 Analytic methods	177
3.4.1 The coupling	178
3.4.2 Constraints on relaxation	181
3.5 Comments	206
REVIEW AND CRITIQUE	207
References	211
Appendix: IMAGE FILES	212

Acknowledgements

The invitation to INRIA, project Syntim, has been highly appreciated for the rare opportunity it offered to rethink the entire 3D stereo problem. Previously, the author had "stayed away" from the stereo problem while "observing" the efforts of others. During this one year the 3D stereo problem was restarted from beginning but in a different direction from the presently popular methods. This is only possible in a place where "liberte, egalite, et fraternite" rules.

Special thanks are due to Dr. Andre Gagalowicz for the invitation, to Mr. Jean-Paul Chieze for general assistance with the "Bora-Cumulus-Archille-Unix-C" environment and for several programs (delret, kasv2im, etc.) without which the work would have been impossible, to Mr. Laurent Vinet for more programs (im_laser, etc.), and to all the other colleagues who have provided valuable advice and assistance, and last but not least, to Mlle. Laurence Bourcier for taking care of the "administrative aspects".

Chapter 1: THE OBVIOUS

1.1 Introduction

Since the history of photogrammetry is rather long, the meaning of the word "photogrammetry" has been redefined several times. According to the Manual of Photogrammetry (1.1), the original meaning was "the science or art of obtaining reliable measurements by means of photographs", later "interpretation of photographs" was included, and by now the definition includes "remote sensing" and more. Possibly automatic image processing, analysis, and "pattern recognition" techniques will be included in the future.

The theoretical as well as the practical aspects of photogrammetry have been well "worked out" and will neither be discussed nor summarized here. Perusal of the more than 1000 pages of the "Manual of Photogrammetry" can give the interested reader a fair overview and numerous references.

As stated in the Abstract, if two corresponding points have been found, one in the "left" and the other in the "right" image of the stereo pair, then the distance to this point in the scene (the 3D world) may be computed from the measurements of point positions in the images. The computational methods include corrections for the way the two photographs are "tilted and rotated" with respect to each other, camera displacements and optics, lens distortions, and so on. The correspondence between the two points is established by an operator (using appropriate equipment), after which the calculations are performed by computers. Methods and equipment exist for establishing correspondence automatically. The methods are based on cross-correlation or (equivalently) on Fourier analysis of small regions in the two images. When the two regions "agree" the (phase) shift corresponds to the height or distance that is to be found. Automatic comparisons give reliable results for (contrasting) "far-away" or rather "flat" scenes for which the two (stereo) images are, essentially, "shifted copies with local variations" of each other.

If the stereo pair is taken at "close range", then the variations between the images are considerably larger and one "eye" may see "things" that the other does not. Of course, automatic comparisons still work, but only in selected parts of the images. Consequently, a human operator is needed to guide and/or assist the machine. Note that the fundamental requirement is to "establish unique correspondence between points in the two images of the stereo pair". When these techniques are applied to, for example, the vision system of a mobile robot then the human operator cannot be involved (all the time). Each image pair has to be analyzed purely automatically in order to establish cor-

respondence between points in the images and to compute the distances.

Since images of real scenes have very complex content, the two standard approaches that have been used are:

1. Ignore the scenes altogether and put some point light sources onto the objects in the scene that are of interest.
2. Ignore the scene and put some corner reflectors onto the objects of interest.

In both cases the images of the scene consist of "uniform black" with some "bright points". The bright points are located with (binary) image analysis methods, put into correspondence, and the distances are computed. The essential difference between the two methods is that in the first (1) the lights can be turned "on and off" thereby making the correspondence problem trivial, while in the second (2) reflectors are placed "far enough apart" to again make the correspondence problem trivial, or the machine is first assisted by an operator who "puts a box around each mark" and establishes correspondence, after which the machine "follows the points in the images" while the object moves, thereby keeping the points in correspondence. Formally, three non-collinear points with known spacing are sufficient to compute the six "degrees of freedom" of an object (marks). Commercial hardware exists (or existed) for both methods.

Numerous variations are possible and have or are being used, such as:

3. Mark the scene but keep the background around the mark "clean" such that the mark can be located with elementary image processing techniques. The marks may be of different shapes and/or in different colours.
4. Use "prominent feature points" in the images of the scene.
5. Use "edges" in the images of the scene.

The last two methods (4,5) are current research topics in many laboratories (1.2). The basic difficulty with these approaches is that they create a "combinatorial explosion" while trying to establish correspondence.

In industrial applications many "special effects" are used to establish distance or position of some "object of interest". In stereo vision the accuracy of the range measurement depends on the distance over which the measurement is made, on the displacement of the two cameras, and on the focal length of the lenses versus the resolution of the image sensor from which the stereo parallax is obtained. Clearly, with closely placed cameras and far-away objects the parallax (or displacement "d") in the two images is practically zero. The accuracy of the distance measures

is proportional to $1/d^2$, i.e., the parallax "d" has to be "significant" for obtaining any accuracy at all with stereo based methods.

However, in addition to all this, there is the entire field of "active vision" for detecting depth. (Unfortunately, the word "active" has now become confused and has two meanings. The "old" meaning used to be "active = emit energy (light)" versus "passive = do not emit energy". The "new" additional meaning is "active = moving or mobile"). The principal active (= emit energy) methods of measuring distance are:

6. The laser range scanner which uses triangulation to measure distance. The scanner gives the distance readings directly and to very high accuracy.
7. The "light strip" method where a flat plane of (laser) light is projected onto a scene. In the images the strip appears as a "crooked" line. The shape of the line is related to the shape of the object and the position of the line is related to the position of the object in the scene.
8. Moiré patterns of interference fringes have been used.
9. Radar, lidar ("light radar"), and sonar ("acoustic radar") measure distance as a function of travel time.

Another categorization of "depth vision systems" is that they are:

- a) Passive, such as stereo, where the observer can be immobile and "emits no radiation". Thus, except for the two "eyes", the observer can remain concealed (hidden) since only the ambient light is used.
- b) Active, such as the laser range finder, which emits a laser beam to triangulate for distance. In the active systems the observer "advertises" his presence by emitting energy in order to "see".

It is hoped that this brief survey places the stereo depth (range) vision system in its proper "perspective", i.e., it is one of many methods. In a practical application all the possibilities should be considered before making a commitment.

1.2 Seeing distance

It may appear from the brief survey in the previous section and from much recent literature that "all aspects of seeing distance (depth, range)" have been covered, but this is not the case at all since the need for seeing the distance has not been defined. The need for seeing distance defines the required accuracy of the distance measurements.

In practice the required accuracy of the distance measurement depends on the application. For close objects distance is a very important "variable", but frequently distance is not important for "far-away" objects. A "relative" measure of distance is often far more important, i.e., is the object "in front of" or "behind" something else. The knowledge that the object is far away is sufficient in itself and it is far more important to know the speed of the object relative to the observer. Even though, mathematically, velocity is but a derivative of distance, the important aspects of "velocity" can be observed without measuring distance. To "see" this (superficial) contradiction between mathematics and reality, consider the driver of a car. Does the driver really measure distances between say 5 to 100 meters to an accuracy of centi- or deci-meters in order to compute the relative velocity vectors of all the other cars to the required accuracy? The human vision cannot measure a distance to the accuracy of a centimetre within a range of even one meter! Clearly, much of the vital information required to "operate" in a 3D world is both different and obtainable by other means.

The human visual system uses about a dozen methods at the same time for detecting depth or distance. Only one among these methods is based on stereo pairs of images and it is claimed that more than 10% of the population cannot see depth based on pure stereo vision, but they may not even be aware of this loss (1.3). There are many cues to depth and it appears that almost all of them are used. The following is a very brief description to indicate the diversity and finessing that one finds in biological information processing systems. For illustrations and better appreciation of the details the original sources should be studied (1.4, Gibson). The ordering of the phenomena used are according to Gibson. Some comments have been included to indicate the nature of the computer vision problem if the use of this effect is contemplated. Thus, directly quoted or slightly modified:

1) The texture perspective or the texture gradient. The sizes and shapes of the texture elements in the image vary as a function of displacement. The detection of texture gradients is a topic in computer vision.

2) Size perspective. Objects vary in size as a function of distance. This requires object recognition and the effect can be used to create visual illusions by producing conflicting situations between the foreground (objects) and the background. Such conflicts, however, seldom occur in nature.

3) Linear perspective or parallel lines "meet" at infinity. From the image processing side this requires line or contour detection, estimation of slopes and curvatures, and clustering of the results. Certain aspects of this effect are used when the scene is illuminated with structured light.

4) Binocular perspective or stereo vision. This uses both the displacement and skew of the image in one eye with respect to the image in the other eye. The necessary comparisons can only be carried out if the image contains details (contours, textures). In order to simplify the comparison problem, for example, the scene may be illuminated with randomly structured light.

5) Motion perspective. The changes that occur in the image due to motion are very complex but also highly informative. The scene expands, contracts, skews, or rotates around the fixation point, i.e., where we fix our gaze. Added complexity is introduced when the objects in the scene are also moving independently of the observer. These effects are easy to observe, for example, when watching the landscape from a moving train or car. The total "visual flow" gives a very compelling impression of one's state (position, orientation, and motion) with respect to the rest of the world. It is also well known that actions can be easily recognized by motion alone even when the objects themselves are invisible, for example, when the moving objects carry some bright spots while the objects themselves are in the dark. Illusions occur when the observer is not aware that he is being moved. Clearly, the analysis of motion is exceedingly well developed in the biological vision systems. Only some simple aspects of motion and visual flow have been studied in computer vision.

6) Aerial perspective or the haziness, blueness, and desaturation of colours as a function of distance. This varies with the illumination of the scene and the scattering properties of the intervening medium. Good artists know how to use this effect in their paintings. These effects have not been used in computer vision, even though the scientific aspects of light scattering are quite well developed and sometimes used to correct satellite or high altitude photographs.

7) The perspective of blur or the variation of the quality of blur as a function of displacement from the centre of clear vision. This effect is hard to observe since we automatically focus on what we observe.

8) Relative upward location in the visual field (on the retina or in the image) if the background is assumed to be a (horizontal) terrain. This effect is easy to measure once the objects in the image have been identified.

9) Sudden shift of texture density or linear spacing between texture elements (usually coincident with changes in brightness or colour). The sudden change in texture density generates a

contour where the more densely packed texture elements appear farther away. This is the texture contour problem in computer vision.

10) Shift in the amount of double imagery occurring at a contour. When a contour is observed with two eyes the images differ near the contour. This effect is one source of problems for automatic stereo mapping.

11) Shift in the rate of motion of the texture elements in the image on one side of a contour with respect to the other side of the contour when the observer's head moves. The shift of a contour in the image is directly related to the distance between the contour and the observer, measured with respect to a (stationary) background. Simpler aspects of this problem are called "change detection" in computer vision. Similar triangulation is used in the laser range finder but, of course, the physical realization is very different.

12) Completeness of an object or the continuity of its outline is an indicator of what is in front of what. Complete objects appear closer since only they can obscure what is behind of them. In computer vision the closest parallel is the analysis of "blocks worlds".

13) Transition between light and shade. Sudden shifts in the brightness of adjacent regions generate contours but also generate the appearance of depth. However, this problem appears to be unresolved. For smooth surfaces the brightness gradient has been used to infer shape in computer vision (1.5).

These thirteen effects are cues to depth or variables computed from the image by our visual system, as described in literature. The stationary one-eyed observer can see texture, size, and linear perspective (1,2,3), aerial and blur perspective (6,7), the relative positions of objects (8), as well as depth at a contour (9), continuity of outlines (12), and transitions in brightness (13). The one-eyed mobile observer can, in addition, also see a form stereo with time delay (4), many aspects of motion perspective (5), and the texture shifts occurring at contours (11). The stationary two-eyed observer can, of course, see everything that the stationary one-eyed observer sees, and also stereo (4) and the doubling of texture at contours (10). The tilt, orientation, and curvatures of a surface are also immediately visible but how these are computed in the human visual system is not certain.

The biological vision system is designed for survival. Consequently, in order not to "advertize" one's position or even existence, the vision system has to be passive. Motion is a very easy "feature" to detect and consequently very dangerous. Hence, the understanding of the visual scene must not require motion on the part of the observer. Stereo vision provides reasonably good distance information (lens focus is used by chameleons). Stereo

vision requires minimally two convergent eyes. In principle, three or more eyes could provide "better coverage" but more than two eyes occur only in very primitive organisms. The biological vision system has to "understand" motion, both its own and that of other moving objects, in order to be able to navigate in the three-dimensional environment. Consequently, motion understanding is very well developed.

As may be observed, stereo vision only plays a minor role and when the ability is missing (in about 10 to 11 per cent of the population) the lack may not even be noticed! Motion based "stereo" is far more informative. However, the present work has been restricted to "still-life" (nature morte) stereo where all motion is excluded.

1.3 The similarity principle

The nature of the stereo problem may be "seen" by simply looking at a reasonably complex scene. Look at the scene with one and the other eye without moving your eyes or head and ask: What did one eye see and what did the other eye see? The answer is, of course, perfectly obvious: We see practically the same scene with one eye or the other! Alternatively, we may look at a pair of stereo photographs. Usually it is hard to see any difference at all between the two photographs without careful inspection of object relationships. After one is convinced that this is true, the rest is just a matter of "working out" the technical details.

The "secret" of the stereo vision problem is thus nothing more than knowing that the two images of the stereo pair are highly similar but, of course, not identical. This is the basis of the "similarity principle" upon which the solution is based. The similarity principle is very powerful since it does not require object or scene "recognition", nor any "scene modelling", no "scene understanding", and the "meaning" of the scene can remain unknown even after stereo reconstruction. In other words, none of the "high level methods" are needed, which are currently very popular. In this sense the present approach is "orthogonal" to the current thinking, since all that the present approach requires is a high degree of similarity between the two images of the stereo pair. The technical requirements for similarity and the differences from photogrammetry are discussed in section 1.5.

Of course, the similarity principle has been used implicitly by all stereo based methods. Automatic cross-correlation or Fourier transformation of small areas in the images finds the "shift" between the similar regions in the two images. Feature point based methods "look for" similar features points, and edge based methods search for matching edges in the two images. However, why stop here when there is a practically unlimited number of additional "similarities" to be found between the two images!

1.4 The hierarchical approach

The "similarity principle" only requires the left and the right images of the stereo pair to be "similar", allows unknown images to be matched (put into correspondence), and allows the 3D scene to be reconstructed without the need to know what the scene represents. Consequently, "similarity" requires no "image understanding", object "recognition", or "models", etc., but if this information is available, it can be easily incorporated.

"Similarity" can be used on many "levels" and, consequently, there are many different approaches to solving the stereo problem. The final goal, however, is the same in all cases, namely, each of the pixels in the two images have to be put into a unique correspondence, except where the two "eyes" see different things. The match cannot be one-to-one on the pixel level since, for example, a "stick" in the scene may be represented by 10 pixels in one image and 12 in the other, but this is only a quantization and resampling problem. The phrase "pixel level match" is used for convenience but it should be understood to include resampling or treated as an approximation given the present quantization.

The basic problems to be solved are, thus:

- A. How to extract some quantities "x" from the two images, and what are these x-s?
- B. How to group these "x-s" into some larger units, groups, or "X-s", if necessary, and what are the conditions that these groups have to satisfy?
- C. How to match the x-s or the X-s according to similarity in the two images?
- D. How to compute the distances after the matching has been accomplished?

The technical requirements for similarity are discussed in section 1.5. Here it may suffice to say that: "The larger the region in the images that is represented by one "X", the simpler it is to match the "X-s" in the two images, the simpler it is to find a unique match, and an increasing number of techniques becomes available for matching".

An hierarchical strategy is proposed as sketched in Figure 1.4-1, out of which only the steps marked with a star (*) have actually been programmed and verified experimentally. These steps are described in detail in Chapters 2 and 3.

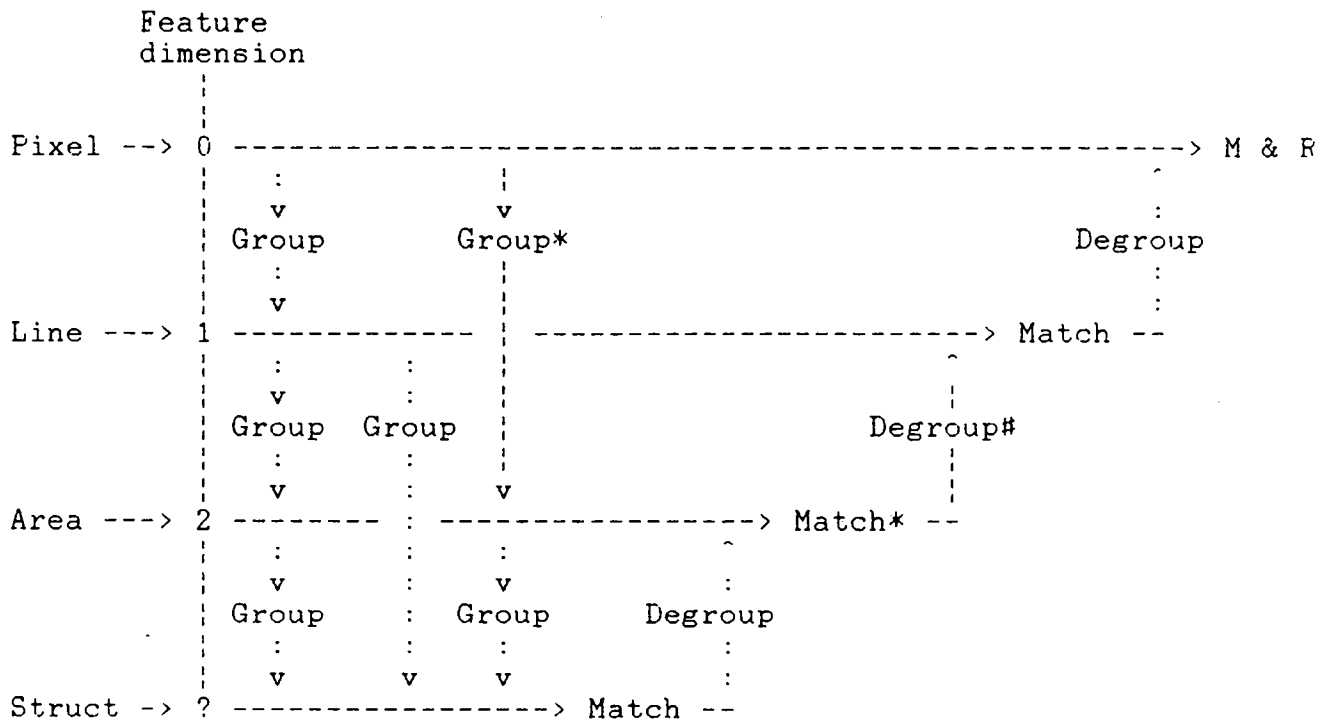


Figure 1.4-1: An hierarchical stereo matching method.
 Struct = structure.
 M & R = match and reconstruct.
 * = Discussed in Chapters 2 and 3.
 # = Edge matching.

1.4.1 The quantities "x" and "X"

The quantities indicated by "x" and "X" have many different names in image processing and analysis literature. In the present case these quantities are categorized according to their dimensionality in the image space:

1. Zero-dimensional. These are features, characteristics, or properties that "belong" to a pixel. In short, these are called the "pixel level" features.
2. One-dimensional. These are single or non-branching one-dimensional regions in the image space normally called "lines" or "edges". In the present case all the features associated with edges, i.e., slope, curvature, strength, etc., are included.
3. Two-dimensional. These are usually called segments, regions, facets, areas, etc.
4. Structures. These are various compositions of the above but in the present case only the composites of (3) are implied.

The zero-dimensional features are of two kinds:

- i) The features measured at pixel level, i.e., the gray level intensity of the pixel or the intensity of its colour components if the image is in colour. Temporal effects (flicker) are also pixel level features if a series of registered images is available.
- ii) The features computed from a local neighborhood around each pixel. The number of such features is practically unlimited.

For the one- and two-dimensional regions in the images and, of course, even more for structures, the numbers and types of features computable is high and difficult to enumerate briefly. From the practical point of view the number and variety of such features may be considered unlimited.

In the usual cases the one-dimensional features are "edges" and the process is called "edge detection", i.e., the pixels are grouped according to selected zero-dimensional features (or functions of such features) into one-dimensional configurations in the image space. Of course there are many methods among which the conventional "edge detection and edge linking" is but one. Clearly, it is highly desirable that the one-dimensional groups (edges) in the left and right images are "true" (very similar) such that when we look at such images, the stereo effect is immediately apparent. However, this may not be an absolute necessity if subsequent processing of the "raw" results will "chop", "cut", and "clip" these regions such that the stereo effect

reappears. This "avenue" was not explored in the present study, partially due to much prior work already completed on straight edges, and the simplifications that can result from using two-dimensional features.

The two-dimensional quantities "X" are pixels grouped according to selected zero-dimensional (but sometimes also one-dimensional) features (or functions of such features) such that these groups represent two-dimensional regions (areas) in the image space. These methods are normally known as "region segmentation" or "region growing". However, much better results are possible by using classical pattern recognition techniques. The classes or clusters obtained are simply mapped or projected back into the image space, for example, as is done in crop classification used in remote sensing. Actually, both the grouping and subsequent matching process is mainly an application of classical pattern recognition methodology, where both much theoretical as well as practical experience is available.

In the present study the grouping was based on the gray levels and the gradients of the gray levels. The resultant regions in the image spaces had poor "stereo appearance" since there was a marked difference in the two images of the stereo pair. However, the "stereo effect" reappeared after the regions were processed.

Structures represented by "X" can have infinite variety, their descriptions can become exceedingly complex but also unique such that the matching problem reduces to a triviality (but the construction of the descriptions that match in the left and right images is by no means trivial). Prior work in these areas is called "structural pattern recognition". In psychology certain types of these structures are known as "gestalt groupings" or simply "gestalt".

1.4.2 Matching "X"

Matching is just another name for uniquely pairing the quantity "X" in the left and right images of the stereo pair. As usual, there are many methods depending on what the quantity "X" represents. For example, cross-correlation and/or Fourier analysis used in existing photogrammetric instruments uses the gray levels (possibly enhanced), i.e., they depend on one zero-dimensional pixel feature distributed over a small area of the image. Other examples are the "feature point methods" which are also essentially based on zero-dimensional features computed from local neighborhoods, and the "edge-based" methods which use one-dimensional features.

Clearly, the matching method depends on the dimension of the quantity "X". In principle, matching "X" in the left and right images is only a variant of the classical pattern recognition

problem. Consequently, much theory and practical experience is already available. In addition there are other methods such as correlation and even straight-forward "and-ing" (a variant of the logical "and" operator).

If there are enough pixel level (zero-dimensional) features available such that each pixel in the left and right image can be uniquely "paired" then the pairing procedure is the same as the "minimum distance classifier". Clearly, a high number of pixel features is required and the classifier has to include the spatial position of the pixel as well as a condition for resampling to achieve "unique pairing". The last two conditions (resampling and unique pairing) are not part of the classical procedure but the required modifications are obvious. Since the size of the decision space will be huge and exceeded available computational facilities, this method could not be experimented with. However, it is very likely that biological vision systems use some variant of this approach (1.6).

A possible modification of direct pixel feature matching is to combine adjacent pixels in pairs, in triplets, in quadruplets, etc., to create pixel feature combinations. Given enough pixels these n-tuplets will become unique and the matching problem is back in the classical domain. This may also be viewed as another method of "texture detection and matching". This approach was not tried.

Since direct pixel feature matching was not feasible and n-tuplets were not used, the approach was modified by grouping the pixel features into "larger units" or groups. The unique pairing will now be between the groups. Pixel feature grouping is standard practice in classical pattern recognition. The details may be found in Chapters 2 and 3. Grouping of the groups into structures was not tried due to lack of time.

1.4.3 Comments

An hierarchical procedure for matching the left and right images of a stereo pair was suggested mainly due to computational feasibility but also since the classical pattern recognition techniques are directly applicable. The computation of the distances (reconstruction) requires pixel level matching.

If direct pixel matches are available then one could proceed directly to 3D reconstruction and create a "dense" or complete distance (range) image. Solutions are directly available from photogrammetry or, if only a "pin-hole camera" model is used, then the derivation of the necessary equations is only a matter of elementary 3D geometry.

If the matching is only for one-dimensional features (edges), then the points on the edges (pixels) will still have to be matched. Only in case of straight line segments is it possible to interpolate after the line ends are matched, since straight lines in the scene project to straight lines in the images. However, except for the edges, the reconstructed image is "very sparse" i.e., a "vacuum" with only strings of pixels present.

If the matching is on a region level (two-dimensional feature matches) then it is necessary to first match the edges of the regions and then the edges on the pixel level. The interior of the region is simply "another smaller version" of an image and is processed as was done for the whole image initially (iterated hierarchical procedure).

1.5 Validity requirements

The "similarity principle" requires the two images of the stereo pair to be "similar". Of course, the images must differ since they represent the scene "seen" from two slightly different positions. This is a fact geometry and the difference (parallax) must exist for the "stereo effect" to be present. However, it may be easier to comprehend and to develop a "feeling" for the stereo problem if we restart from beginning.

1. Sketch the design of a stereo camera system that best suits image processing and subsequent analysis requirements. Elementary optics, some understanding of image processing, the fundamentals of control systems, knowledge of hardware capabilities, and so on, are sufficient basis for a conceptual design.
2. Assuming that the ideal camera system exists as designed in (1), what are the irreducible problems that still remain, being characteristics of the stereo problem rather than artifacts introduced by equipment or by the conditions under which the stereo image pair was obtained.

1.5.1 The ideal stereo camera

It may be easiest to see the whole problem at a "glance" if we "design" the camera system and its controls in order to satisfy the similarity principle rather than accept an existing design and try to adapt the principle to it. A simple "pinhole" camera model is assumed since there is no need to complicate the situation and this is the "usual" camera model assumed in many articles that deal with the stereo problem.

A "pinhole" camera model is shown in Figure 1.5.1-1. The camera is indicated by a sphere called "eye". It has a pinhole for a "lens" and the image is formed at the back of the sphere indicated by "image". A line that goes through the pinhole and the centre of the sphere will be called the "optic axis". Let the origin of a coordinate system be fixed at the pinhole. However, at the moment we are only concerned with the controls we expect to have in a properly designed system.

In order to have sufficient control over the image, we must be able to "swing" the camera horizontally, vertically, rotate it, zoom it, and adjust the "iris" of the pinhole, all under computer control. In terms of movie maker's or aeronautical engineering jargon, we need controls for yaw (pan), pitch (tilt), roll, zoom (focus), and image intensity (gray level). This allows us to point the camera anywhere in the scene and to focus or zoom in on it and to obtain a properly resolved light intensity image. Control of "exposure" is standard practice in cameras but it is better to control it via programs. Automatic focussing is avail-

able but it is only a feedback loop in the focus system that maximizes the sum of the gradient magnitudes over the image or the distance is measured acoustically. The place where to point the optic axis and the zoom (magnification of image) are defined when two cameras are adjusted to satisfy the "similarity principle".

A two-camera system is shown in Figure 1.5.1-2. One of the cameras ("eyes") will be called "left" or L and the other "eye" is called "right" or R. Both cameras have independent controls over yaw, pitch, roll, zoom, and iris. It is now simple to see how to control the cameras to satisfy the similarity principle and to simplify the image processing problems as much as possible.

1. Swing the cameras such that the optic axes intersect at a point called the "fixation point".
2. Zoom the cameras such that the images in the left and right "eyes" are of the same size (have the same scale).
3. Rotate the optic axes such that the images "line up" in the two "eyes" on or near the fixation point.
4. Control the irises such that the images have the same intensities.

If a camera system of this kind and its "drives" are properly "wired" to a computer system then the programming required to achieve the above control structure appears "straight-forward" and no "scene understanding" (and all that) is needed. However, in the absence of such a system to experiment with, these remain conjectures only. Note that the human eye "musculature" (control mechanism) has all these controls and they operate automatically. A simple way to "see" what happens to our stereo vision is to move one of the eye balls (push it slightly with a finger) while looking at a scene.

In photogrammetry a certain amount of control over the above variables (items 1, 2, 3, 4) and others is available but they "take the images as they are" (including lens defects, etc.) and apply appropriate corrections to compute the distances. They need the precision and, consequently, the formulas become rather complicated. Even the "famous" $P_l = A*P_r + T$ formula makes no use of the similarities found in the stereo pair. In $P_l = A*P_r + T$, P_l and P_r are the two corresponding points, one in the left and the other in the right image, A is a general rotation matrix, and T is a translation vector.

These transformations become very simple if the coordinate systems are chosen as shown in Figures 1.5.1-3 and 1.5.1-4. Such a set of coordinates may be called "ego-centric" since they only refer to the observer. In essence, the x-axis goes through the two pin-hole lenses, the y-axis is halfway between the two pin-

holes, and the z-axis goes through the fixation point. The "right-hand" rule is preserved if the directions of the axes are chosen as shown. The eye or camera coordinates are chosen at the pin-hole lenses, see Figure 1.5.1-4. The y_l and y_r axes of the left and right eyes are set parallel to the ego-centric y-axis and, of course, the z_l and z_r axes intersect at the fixation point. If we also assume a set of yaw, pitch, and roll controls for the "head" that carries the two eyes (the "third imaginary eye controls" at the origin of the ego-centric system, which keep the "nose" always pointed in the direction of the fixation point), then the eyes only need to be rotated ("panned") by the same amount (in opposite directions) to "fixate" any point in the scene. The image plane coordinates are chosen as x_l', y_l', z_l' for the left, and x_r', y_r', z_r' for the right eye, where $x_l' \parallel x_l$, $y_l' \parallel y_l$, $z_l' \parallel z_l$, and, $x_r' \parallel x_r$, $y_r' \parallel y_r$, $z_r' \parallel z_r$, with \parallel indicating "is parallel to".

Under the above assumptions, a point P in the scene and its corresponding image points P_l and P_r in the left and right images (which are found by solving the correspondence problem), give:

$$E_l = (x_l, y_l, z_l) = (0, 0, -x_o) = \text{left pin-hole}$$

$$E_r = (x_r, y_r, z_r) = (0, 0, +x_o) = \text{right pin-hole}$$

$$V_l = (-x_{pl}', -y_{pl}', -d) = \text{vector from } P_l \text{ in left eye coordinates.}$$

$$V_r = (-x_{pr}', -y_{pr}', -d) = \text{vector from } P_l \text{ in right eye coordinates.}$$

In the ego-centric coordinates

$$V_{le} = A_l V_l + B_l = (x_l, y_l, z_l)$$

$$V_{re} = A_r V_r + B_r = (x_r, y_r, z_r)$$

$$A_l = \begin{pmatrix} \cos(-T) & 0 & -\sin(-T) \\ 0 & 1 & 0 \\ -\sin(-T) & 0 & \cos(-T) \end{pmatrix}$$

$$A_r = \begin{pmatrix} \cos(+T) & 0 & -\sin(+T) \\ 0 & 1 & 0 \\ -\sin(+T) & 0 & \cos(+T) \end{pmatrix}$$

$$B_l = \begin{pmatrix} -x_o \\ 0 \\ 0 \end{pmatrix}$$

$$B_r = \begin{pmatrix} +x_o \\ 0 \\ 0 \end{pmatrix}$$

The unit vectors V_{leu} and V_{reu} (direction cosines) of V_{le} and V_{re} are

$$V_{leu} = V_{le} / |V_{le}| = (V_{lx}, V_{ly}, V_{lz})$$

$$V_{reu} = V_{re} / |V_{re}| = (V_{rx}, V_{ry}, V_{rz})$$

Approximately, the distances D_l and D_r from the left and right eyes form a triangle with a base of $2x_0$. Let the coordinates of the point P in the scene be (x,y,z) then

$$\begin{aligned}x &= V_{lx}D_l + x_l = V_{rx}D_r + x_r \\y &= V_{ly}D_l + y_l = V_{ry}D_r + y_r \\z &= V_{lz}D_l + z_l = V_{rz}D_r + z_r\end{aligned}$$

which can be solved for D_l and D_r , and (x,y,z) , but the best numerical values should be selected during computation.

If these rays defined by V_{leu} and V_{reu} intersect (at the fixation point) then let $(x$ implies cross product)

$$N_u = V_{leu} \times V_{reu} = (N_{ux}, N_{uy}, N_{uz})$$

which is the normal of a plane in which the vector V_{leu} lies and which is parallel to the V_{reu} vector. The minimum distance Q between the rays V_{leu} and V_{reu} is

$$\begin{aligned}Q &= N_{ux}x_r + N_{uy}y_r + N_{uz}z_r + d' \\0 &= N_{ux}x_l + N_{uy}y_l + N_{uz}z_l + d'\end{aligned}$$

The "squint" Q is a line which has the same direction as N_u and may be used to compute a minor correction to the distance z .

The above brief derivation is only intended to show that with a suitable selection of coordinates, and if there is enough control over the cameras, then the computations can be rather simple (using the shortest distance between two lines in space). The above computations have not been checked experimentally for their numerical stability. If the eyes are independently mobile, i.e., the z -axis does not go through the fixation point, then the computations are slightly more complex but the similarity principle remains the same.

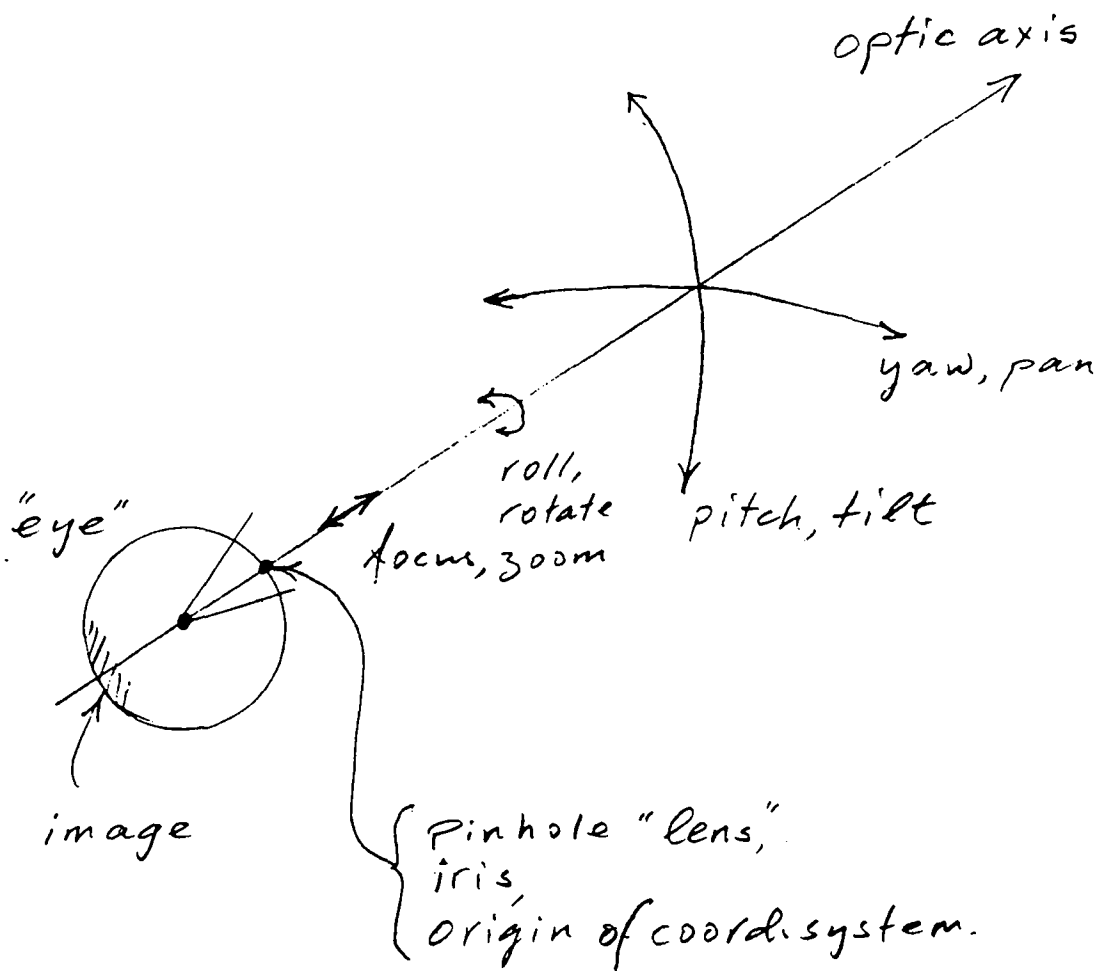


Figure 1.5.1-1: The pin-hole camera model.

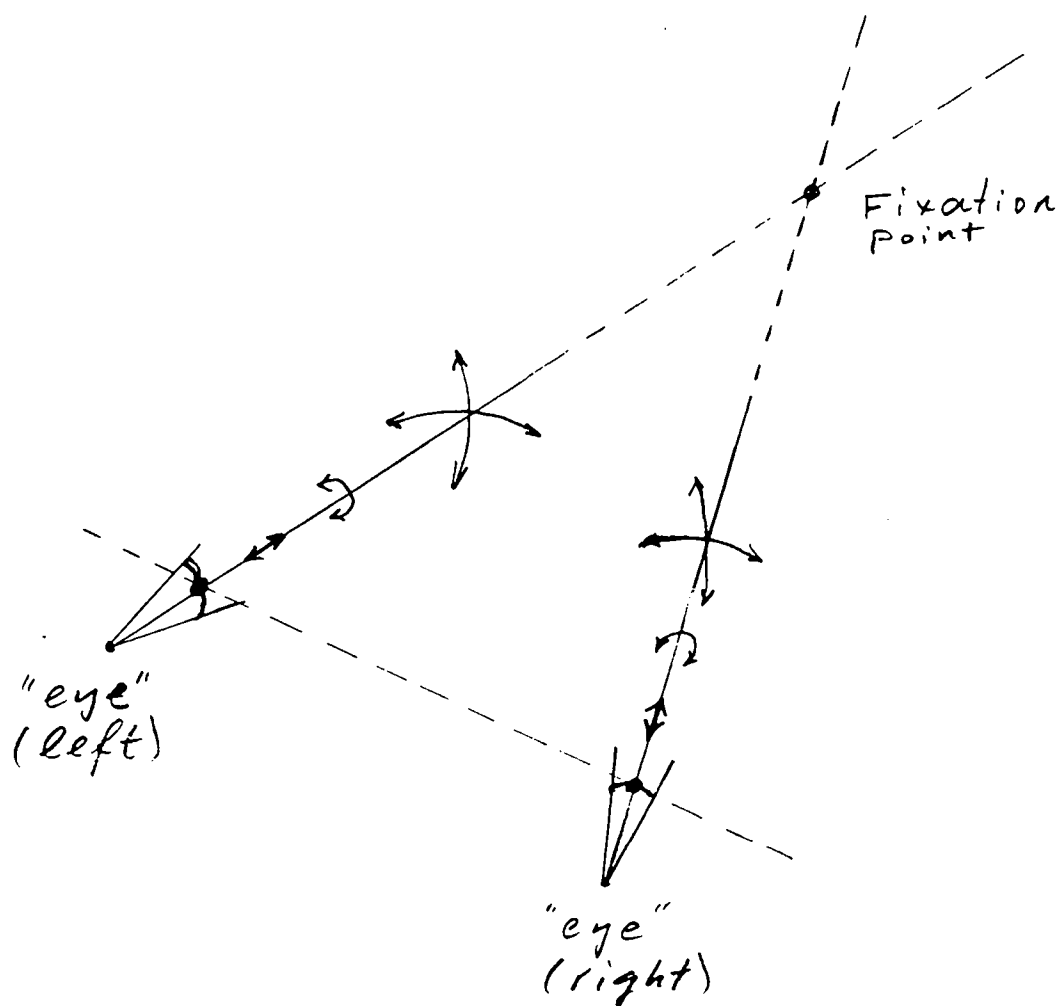


Figure 1.5.1-2: Two pin-hole cameras.

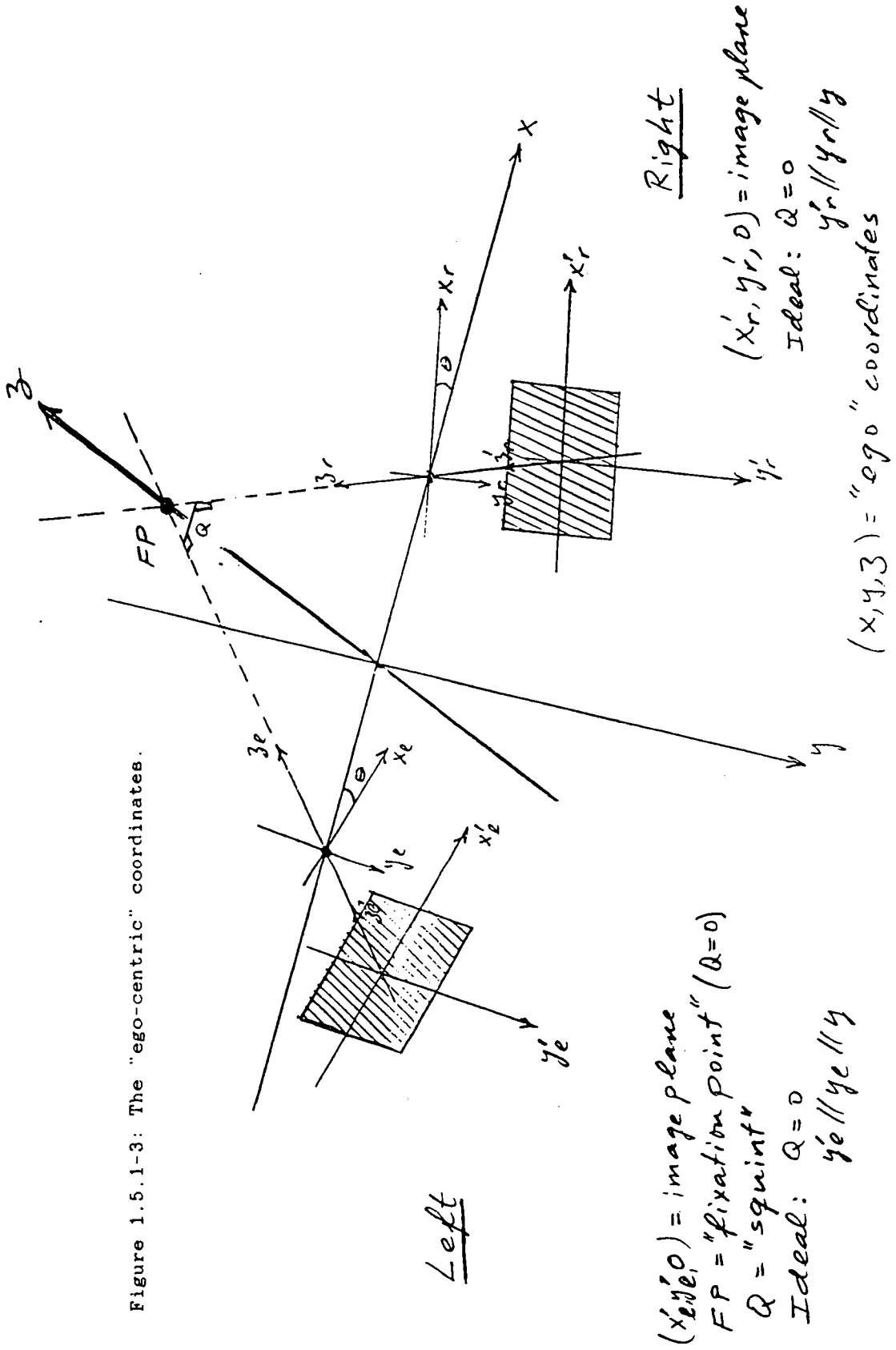


Figure 1.5.1-3: The "ego-centric" coordinates.

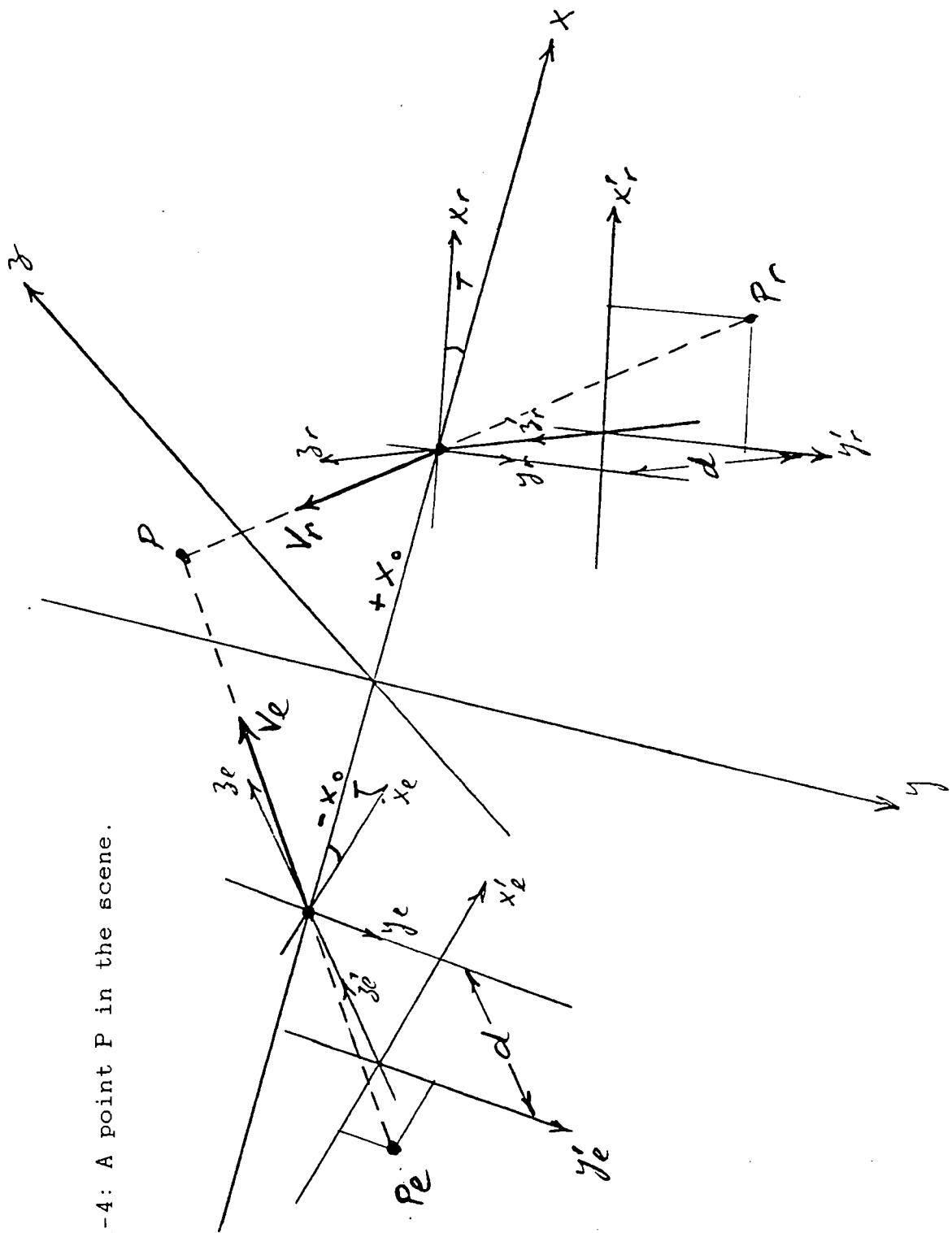


Figure 1.5.1-4: A point P in the scene.

1.5.2 The irreducibles

An object in the 3D space can be anywhere (in any position in the space). An object can have any shape, size, rotation, and orientation. The objects need not have any particular predefined shapes, they may move, rotate, split, merge, and change shape. Objects can be semi-transparent, can obscure each other, and so on. In addition to all this, the illumination adds another set of "variations". Thus, in general, a scene can be very complex, but all these variations do not seriously affect the similarity principle, provided that some elementary conditions are satisfied:

1. The stereo pair is taken at high enough speed and proper focus so that motion blur and defocussing are not noticeable in the images at the resolution required for analysis.
2. The two images have been taken simultaneously or with a single camera of a stationary scene.
3. The cameras are "sufficiently" close such that the disparity (parallax) is not excessive. It should be remembered that there is both a practical and a theoretical limit to the parallax, i.e., one can make the problem ridiculous - point one camera east and the other west and try to reconstruct the 3D scene!
4. The fixation point is on a surface with texture or on an edge and not in "free space".
5. The cameras are "properly aligned" such that:
 - i) The optic centers intersect at the so-called "fixation point".
 - ii) The (TV camera) scan lines are in the same plane.
 - iii) Both cameras have the same light sensitivity (iris setting).
 - iv) Both cameras have the same focus setting.

None of the above conditions are unrealistic. The additional constraint of "still-life" (nature morte) implies that only one pair of images is available, being one "frozen instant" from a live system. It may, however, be interesting to observe some differences between the "still-life" and a "live" stereo system:

- a) In a "live" vision (active) system there is a difference in the analysis if both eyes and head move together or when eye motion and head motion are independent.
- b) In "live" vision there are problems of "how to": Focus the eyes. Line up the images. Position the eyes. Rotate the eyes, etc. There is also the question of "where to look", how to integrate the various views, and so on.
- c) Calibration either has to be carried out at every change of fixation point, or the calibration problem is an artefact dictated by the mathematical method used.

- d) The fixation point can be placed anywhere in the scene. This is an exceedingly important aspect of stereo vision, at least from the processing point of view since, at the fixation point, the parallax is zero. Thus, a "live" system can put any point, edge, or region into correspondence by simply "fixating" on it.
- e) "Still-life" excludes any possibilities of using motion, flicker, exploration (change of fixation point), etc., for helping to establish correspondence. In a "gray still-life" the images only contain gray levels (colours are excluded).

The "still-life" case thus avoids certain problems that a "live" would have to solve but does this complicate or simplify the stereo problem? If the additional computations that the "live" system has to carry out are not considered (they are processes that can be carried out in parallel) the stereo problem tends to "vanish" since one can "look at" any point in the scene that causes problems with stereo correspondence. The "still-life" system does not have this privilege thereby unnecessarily complicating the stereo problem.

However, given a "still-life" pair where the fixation point is in the middle of the images and on a surface, what types of variations or problems can one expect? Most of these are easily understood with the help of a few simple sketches:

1. The absurd case: Even in a "properly taken" stereo pair the two images can be totally different. Put a two-sided thin mirror edge-wise between the eyes and look, see Figure 1.5.2-1. This may also be a thin cardboard with different pictures on either side.
2. The absurd but normal case: Even in a "properly taken" stereo pair parts of the two images can be totally different. Two normal situations are sketched in Figures 1.5.2-2a and 1.5.2-2b. In the first case (Fig. 1.5.2-2a) one has focussed onto a hole or maybe a "fixation target" but the scenes seen through the hole can be different. In the other case (Fig. 1.5.2-2b), which is fairly usual, one eye can see much more than the other of a part of the scene but, of course, this part could also be totally different.
3. A point P in the scene which appears fixed in one image can be anywhere (on a line) in the other image, or not even visible in one, see Figure 1.5.2-3. The trajectory of P may be a "stick" seen "end on" in the other image. If the trajectory is a plane surface (thin) then it appears as a stick in one image and as a surface in the other.
4. The two pin-hole lenses for the left and right eyes are two points P_l and P_r in the ego-centric or a "world" coordinate system, see Figure 1.5.2-4a. Given any point P_o in the

scene, these three points (P_o, P_l, P_r) define a plane P in the space. The left and right image planes I_l and I_r are also planar surfaces. Two plane surfaces that cross each other form a straight line S at the intersection. If the intersection of P and I_l is called S_l and the intersection of P and I_r is called S_r , then the two lines S_l and S_r are called the epipolar lines. Only if P_o coincides with the fixation point F_p then S_l and S_r are in the same plane, otherwise these two lines are "tilted" with respect to each other, see Figure 1.5.2-4b. Clearly, if a point P_x is constrained to move in the plane P , its images are constrained to the lines S_l and S_r .

These simple sketches may be sufficient to indicate how to "see" many of the situations that can be encountered. For example, if there are two indistinguishable points P_x and P_x in each image on S_l and S_r , which of these are to be put into correspondence to obtain the correct reconstruction in the 3D space (cannot decide); do the centers of gravity of the facets in the image space correspond to the centre of gravity of the facet in the 3D scene (no), and so on.

1.6 Conclusions

The brief description in this chapter has, hopefully, put the 3D stereo problem into "perspective". It is also hoped that some of the "mystery" has been removed.

Adequately designed stereo camera systems are available from firms specializing in photogrammetry. However, the problems they address are different from those required in "live" robot vision systems. For the study of still-life (nature morte) stereo a pair of adequately digitized images is sufficient in principle. However, as already mentioned, the solutions to such problems are of more interest to photogrammetry than to the designers of computer vision systems for mobile robots.

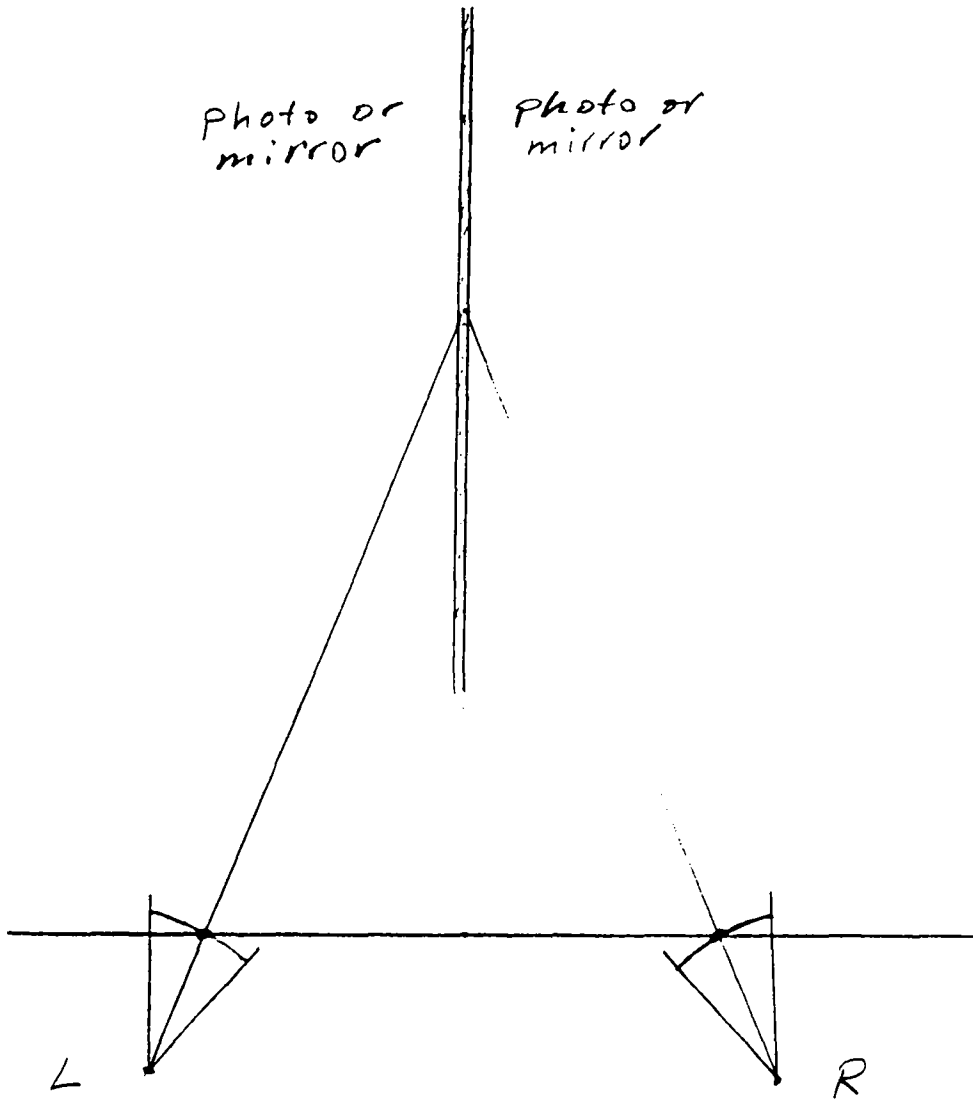


Figure 1.5.2-1: The absurd case.

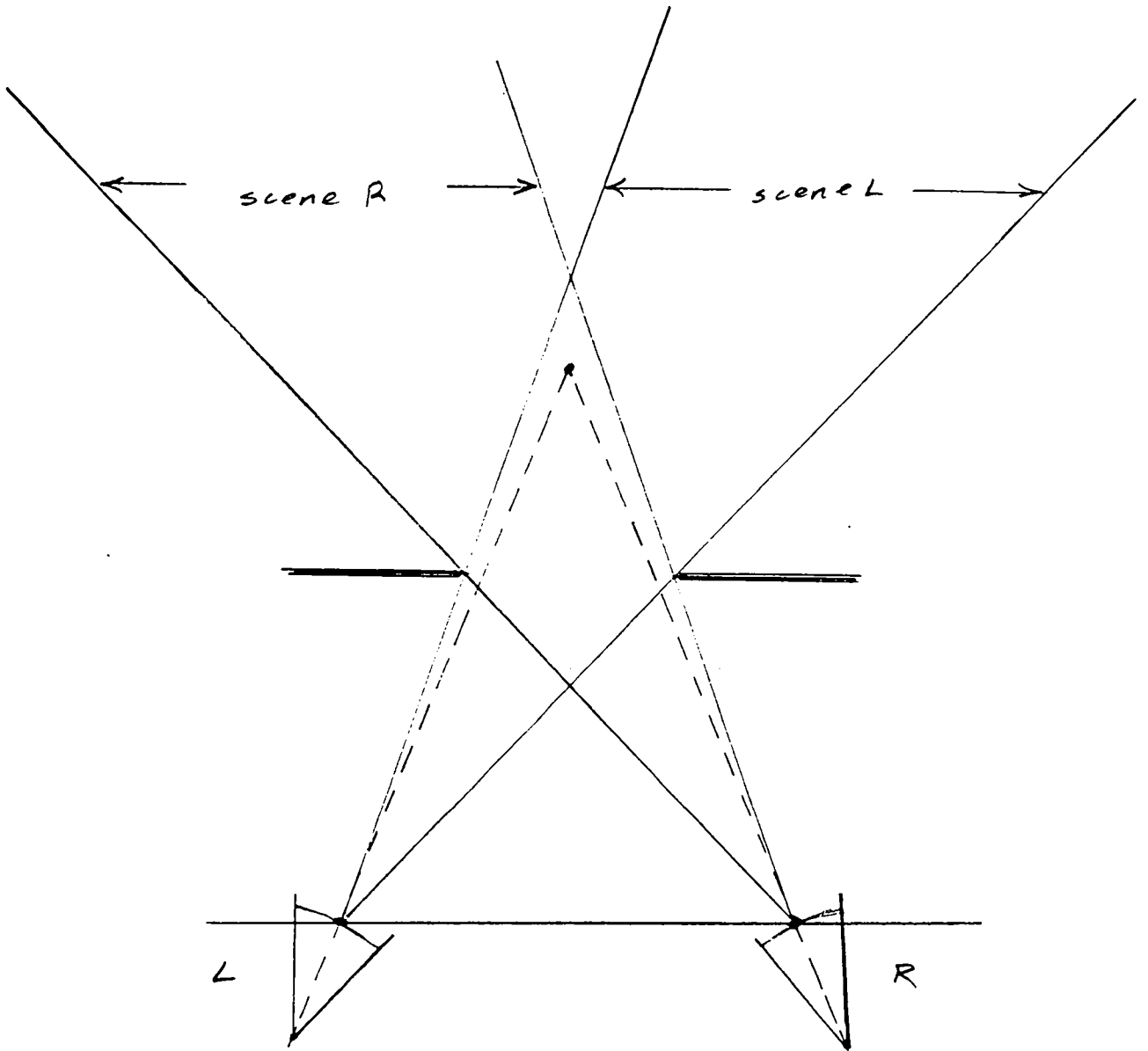


Figure 1.5.2-2a: The absurd but normal case.

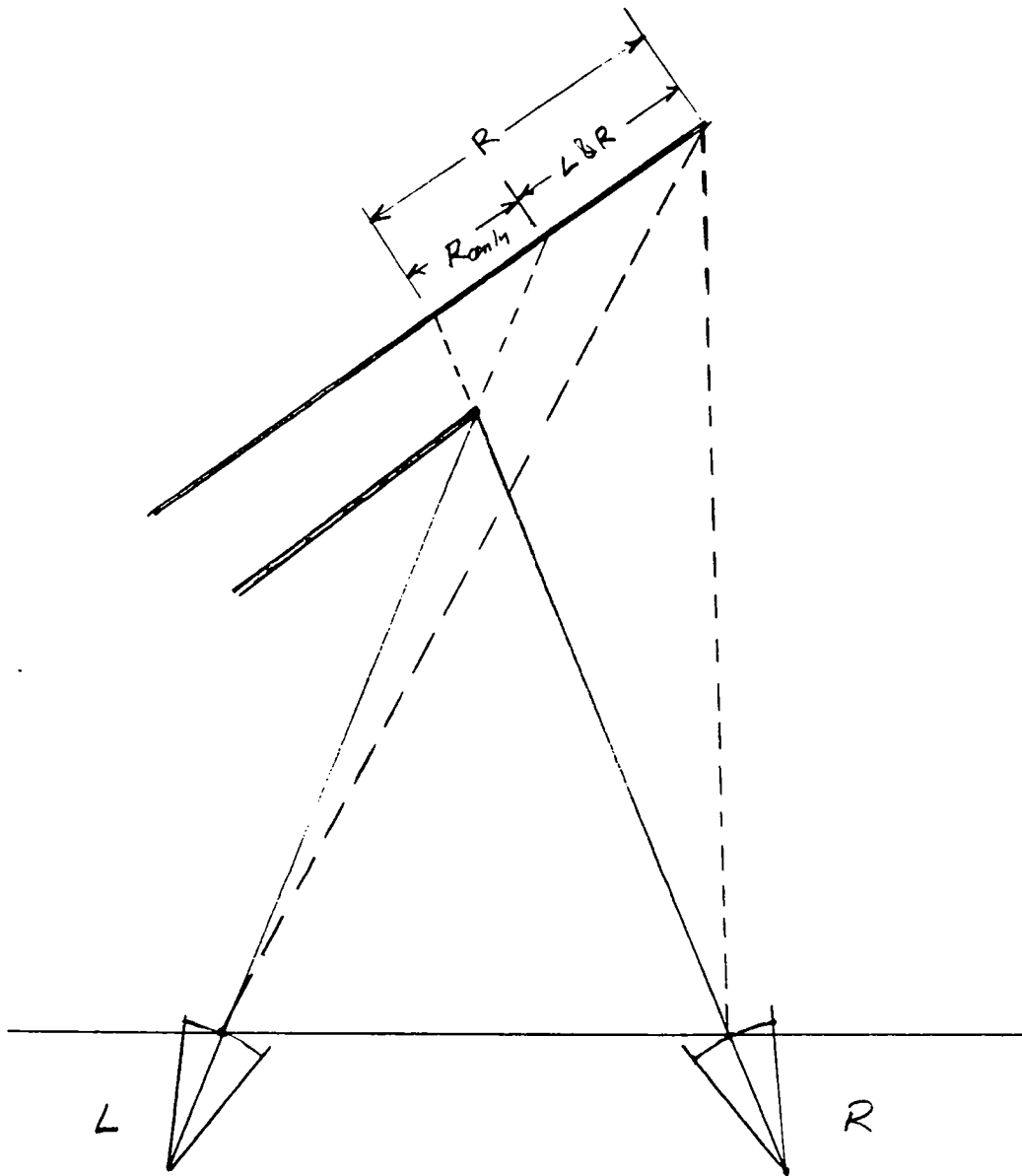


Figure 1.5.2-2b: The usual case.

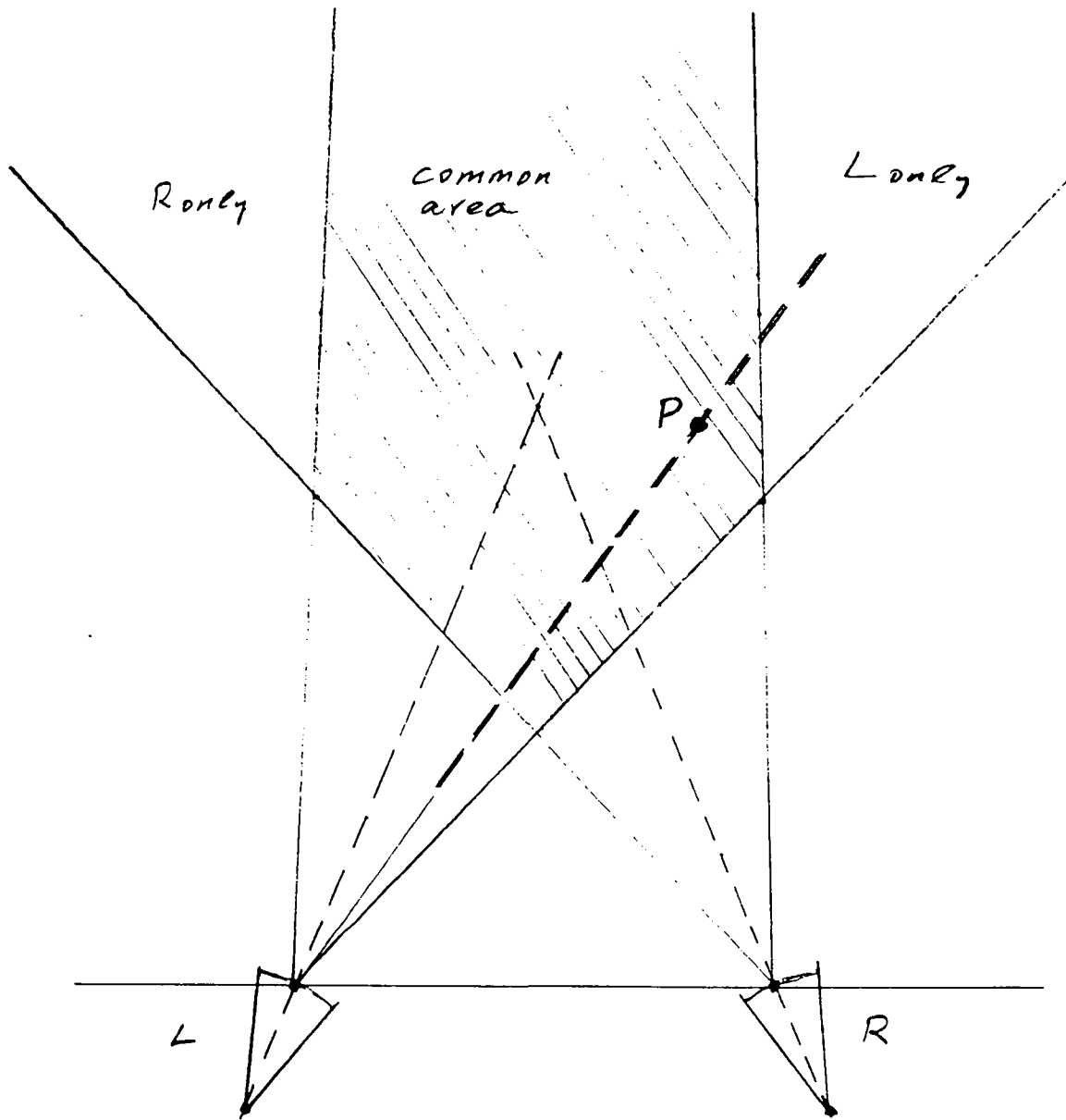


Figure 1.5.2-3: A point P may be anywhere on the epipolar line in one image while it is stationary in the other image.

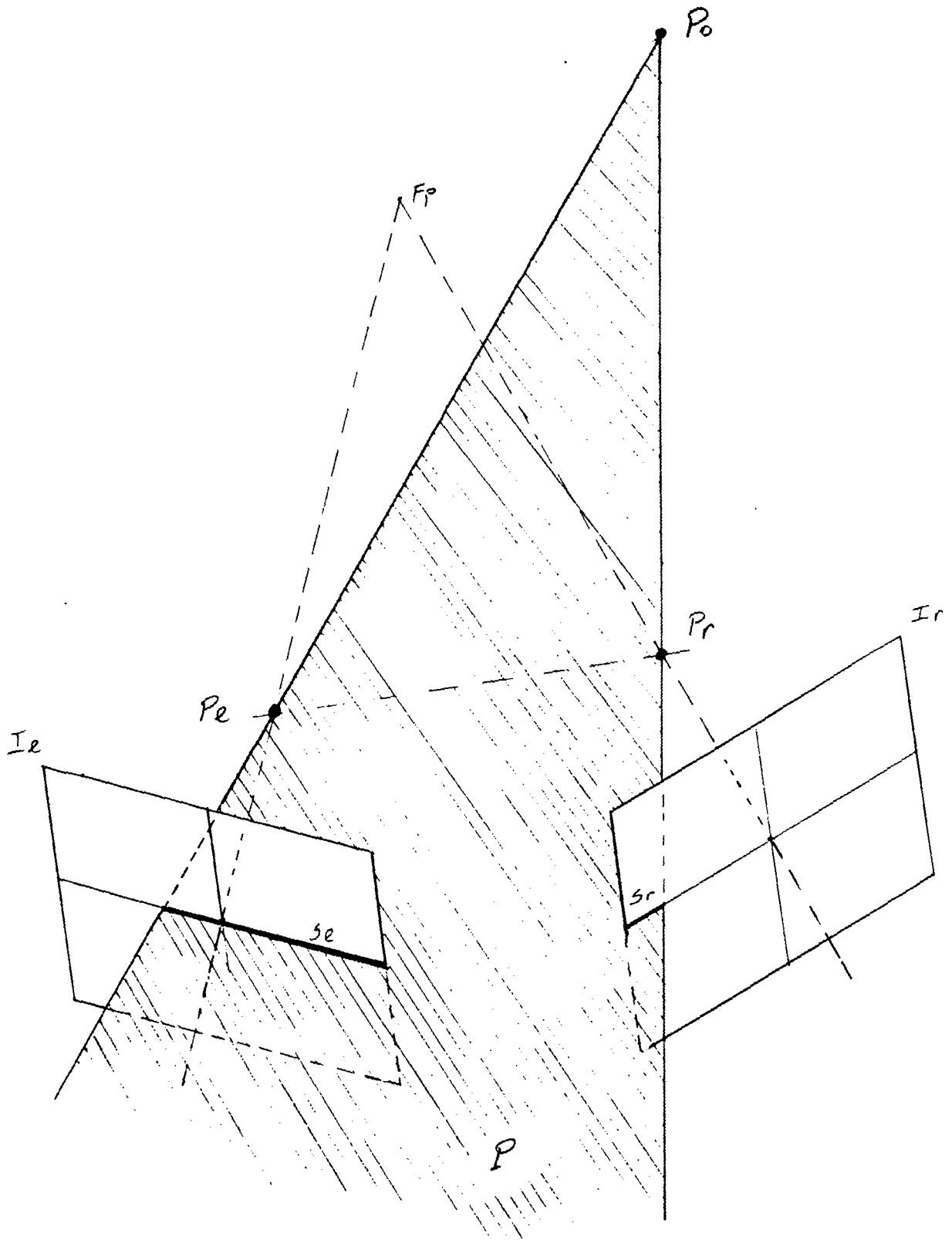


Figure 1.5.2-4a: The epipolar line.

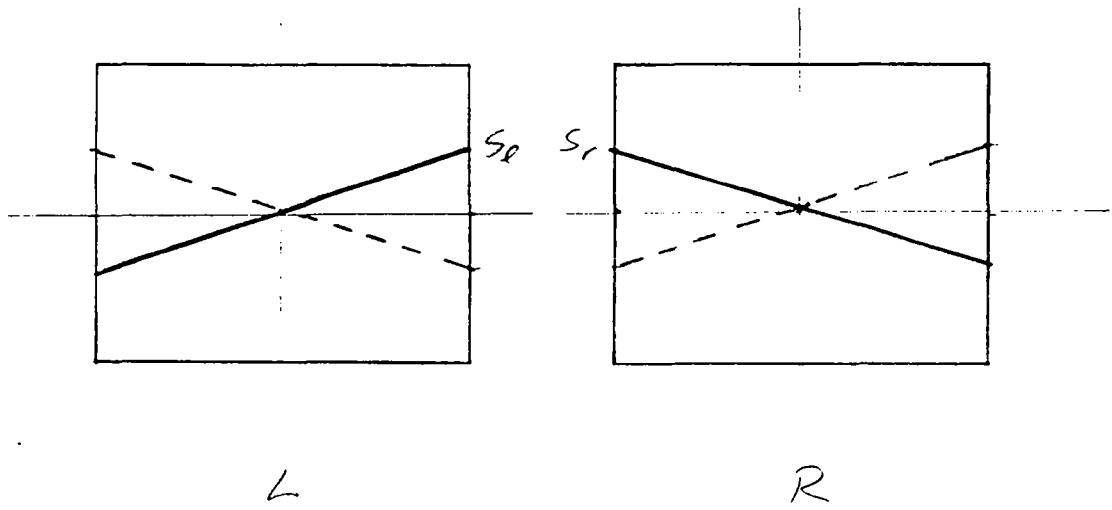


Figure 1.5.2-4b: The epipolar line.

Chapter 2: PRELIMINARY COMPUTATIONS

2.1 Introduction

In order to develop a "feeling" for the actual difficulties with a 3D stereo problem and to "interface" with the existing work at project SYNTIM, the "INRIA canonical image pair #1" (I001g_o and I001d_o) was chosen for experimental studies. However, due to the relatively short visit to INRIA, a rather large personal collection of image processing programs, and the differences in the programming languages and operating systems (Fortran, VMS, and DOS versus C, Unix, and X-windows), the experimental environment consisted of a compromise between private computer equipment brought from Canada and the facilities available at INRIA. In essence, the compromise consisted of the "computational structure" sketched in Figure 2.1-1.

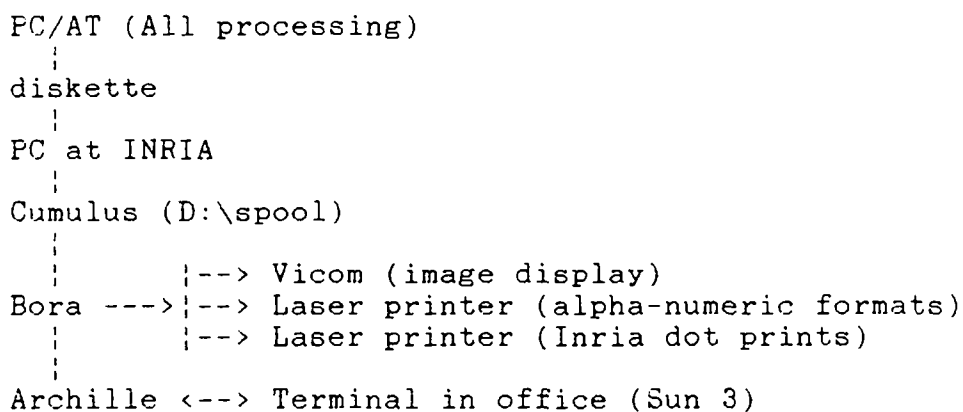


Fig. 2.1-1: The computational structure used for the studies.

The INRIA image pair is shown in Figure 2.1-2 (and called I001.log and I001rog, respectively). These two images appear to form a stereo pair and the depth is visible if the images are looked at through a stereo viewer. Consequently, at least in principle, this image pair should also be adequate for computer vision. The images are arranged side-by-side in Figure 2.1-3a for viewing (using a "Peak 2x viewer", (2.1)) to allow the reader to verify that visual integration of these images presents no difficulties for our vision. However, it was noticed during the experiments that one's stereo vision varies, whatever the rea-

sons. On "good days" the images generated clear impressions of depth, while on the "bad days" stereo integration proved problematic, specially on processed images and even if the originals were relatively far apart. The author is not one of the approximately 11% of the population who, apparently, cannot see "pure stereo" (Julez'es dot stereograms).

In brief, the computational environment was the following:

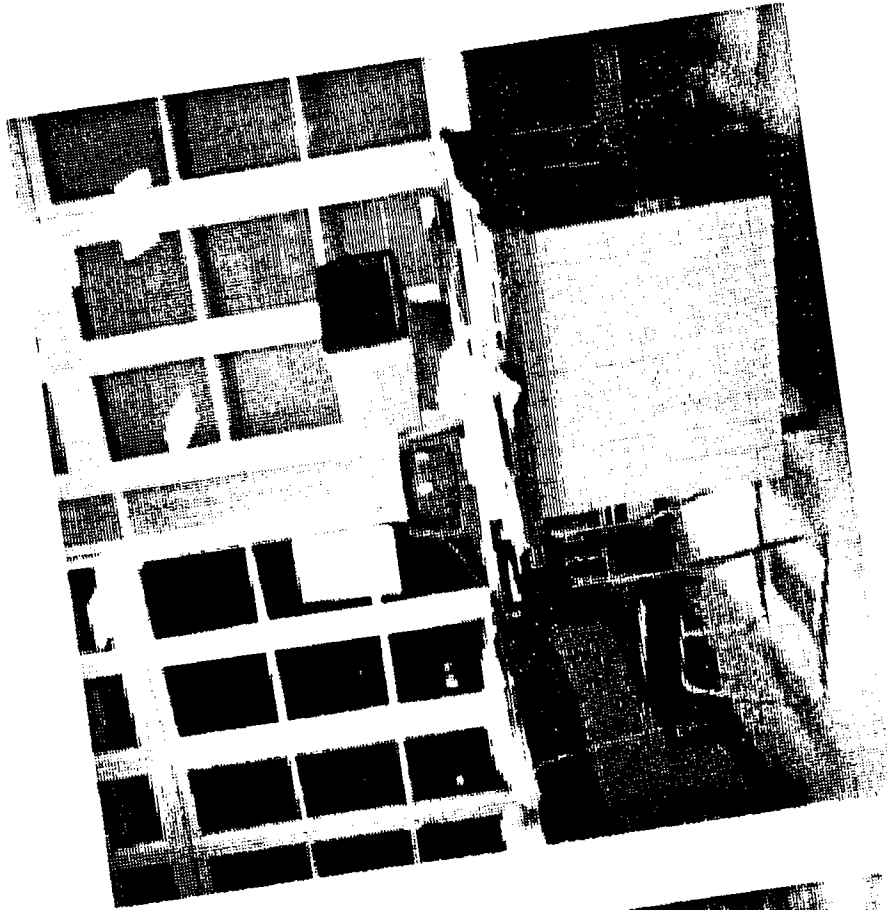
1. Private PC/AT computer using the DOS operating system. This rather small system consists of 512 Kbytes of memory, three hard disks of 32 Mbytes each, a 1.2 Mbyte "floppy", a printer, and a terminal. On this system approximately 1000 previously developed image processing programs (in Fortran) were available from earlier studies of various aspects of image analysis. The computer system was familiar and available "day or night". Consequently, it was preferred over the large systems at INRIA. Initially a telephone interconnection between these systems was considered but these efforts were abandoned due to a variety of reasons and the programming and experimentation was carried out mainly during evenings and week-ends. This compromise placed constraints on some aspects of the research but it represented the "path of least resistance" for quickest results.
2. All the image analysis was carried out on the PC/AT. In order to fit the rather small memory (512 Kbytes) which was accessible to the operating system, the images were reduced to 128x128 bytes by selecting every odd row and every odd pixel on a row. To accommodate the existing "defaults" in the PC, the last row was eliminated, leaving images of 128 by 127 pixels stored in four bytes per pixel. The "integer*4" format is used throughout the analysis. Gray level resolution is preserved by multiplying all applicable data by 100. This stereo pair is shown in Fig. 2.1-3b and a "blown up" version is in Fig. 2.1-4. (These are called "if8fort11.gry" and "if8fort12.gry" in the INRIA Unix system). An 8-level "alphabetic gray" image pair, which can be printed on any alpha-numeric device, is shown in Fig. 2.1-5. Image quality has suffered due to decimation but sufficient remains for processing. Furthermore, since the world is "fractile", it does not matter at what scale the image is digitized, there is always a loss of detail. It was initially intended to develop the programs on the PC/AT using the reduced (decimated) images and then transfer the programs to INRIA facilities to be run at "full resolution". This would have offered a check of the effect of scale on the processing strategies. Unfortunately, the transfer never took place.
3. The "communication" between the PC/AT and the INRIA systems occurred via diskettes using formatted data. In brief, the whole experimental "structure" was as shown in Fig. 2.1-1.

Some difficulties were experienced in obtaining good quality prints rapidly in order to verify the results visually as soon as they were obtained. The program used for printing images on the laser printer, "im_laser -T -s K -inv file_name", had size options from "K=1" to "K=6". Size 6 prints the entire image, while size 5 drops some lines from the image and size 3 drops more lines. The exact cause has not been investigated. A gray level chart was constructed to test the laser printer. The chart has 128, 64, 32, 16, 8, 7, 6, 5, 4, 3, gray levels. The number of distinguishable gray levels is only 4 to 5 on the originals and 4 or less on copies. No texture printing programs existed to bypass this problem. A "4-colour" map painting program was written to try to accommodate the printer but in many cases more than 4 "colours" were needed due to "unknown" regions in the image that required the same colour. (The problem is equivalent to trying to paint the map of the Eurasian continent in 4 colours under the constraint that the seas around it and the lakes within have to be in blue). The size "5" print of the chart is shown in Figure 2.1-6a and two of them fit on one page. Only the size "3" print shown in Figure 2.1-6b could be viewed with the available stereo equipment.

Many of the results are shown as one pixel per character alphabetic prints of images of facets labels (modulu 26). Initially these prints may appear to the reader as "alphabet soup". However, for those who like to see details, these prints are far more informative than gray level or even false colour prints. The texture effects created by groups of similar letters (regions or facets) may suffice for an overall impression. In the vicinity of line 56 these prints may show a "discontinuity" since they were assembled manually from two sheets (X11, xprint, small size).

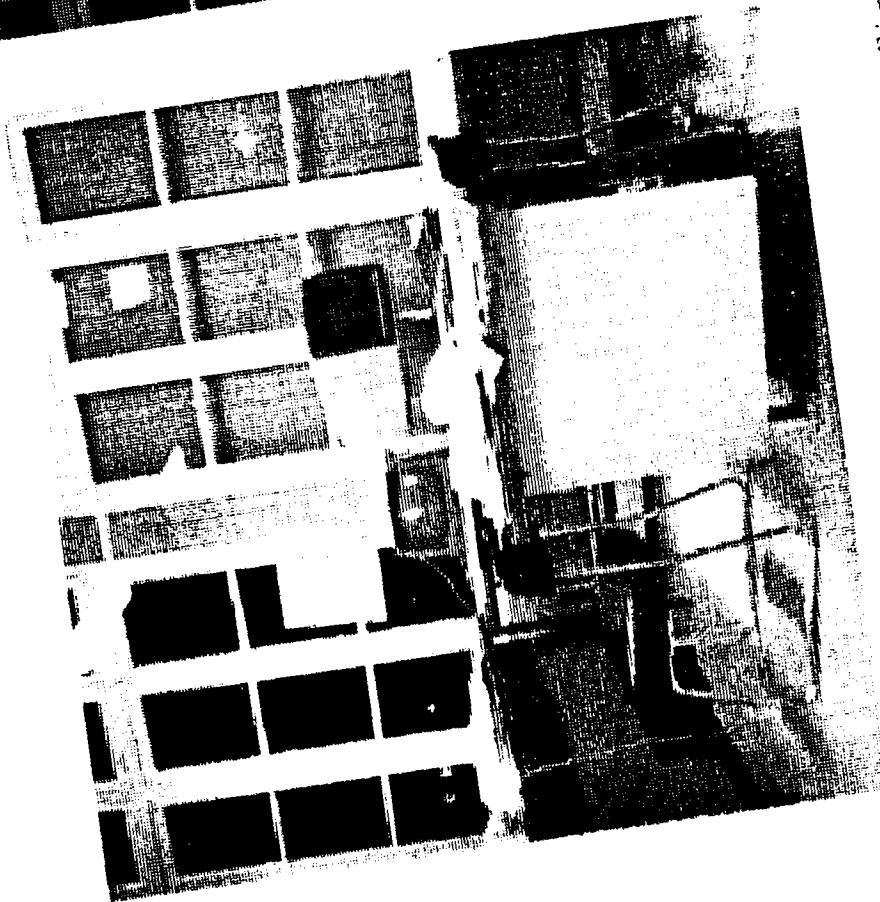
Figure titles:

- Fig. 2.1-2: The "INRIA canonical image pair #1" printed in relatively large size to show detail. I001g_o is the left image and I001d_o is the right image. The digital image matrix size is 256 by 256 bytes. The gray level range = 0 to 255.
- Fig. 2.1-3a and -3b: The "INRIA canonical image pair #1" arranged for stereo viewing (using a "Peak 2x viewer"), and renamed as "i001.log" and "i001rog". (-3a) is the original image while (-3b) shows the decimated image used in computing.
- Fig. 2.1-4: The reduced image pair shown in larger size for illustration of the loss of detail.
- Fig. 2.1-5: The reduced image pair printed in eight gray levels using alphabetic symbols. (Files: Inr021.fig and Inr022.fig)
- Fig. 2.1-6a and -6b: The gray level chart in size 5 (-6a) and size 3 (-6b).



1001.1003

Fig. : 2.1-2.



1001.1003



1001.log

1001.roq

Fig.: 2.1-3a.



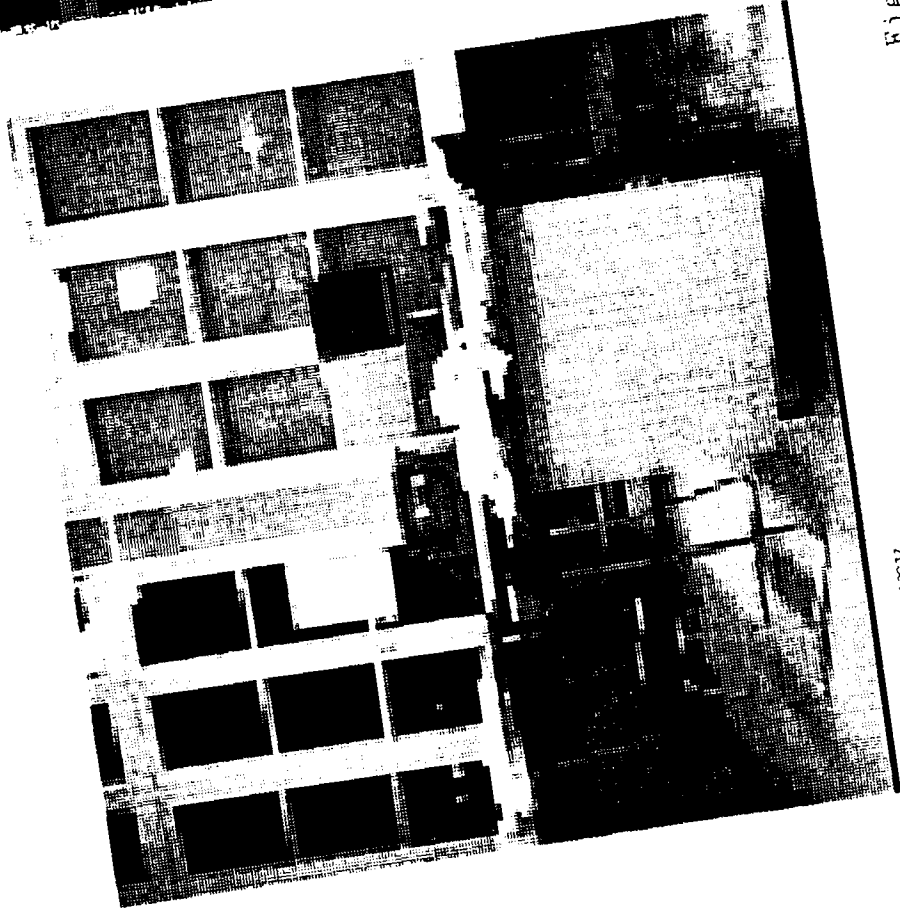
if8fort11.gry

if8fort12.gry

Fig.: 2.1-3b.



if8fort12.gry



if8fort11.gry

Fig. 2.1-4.

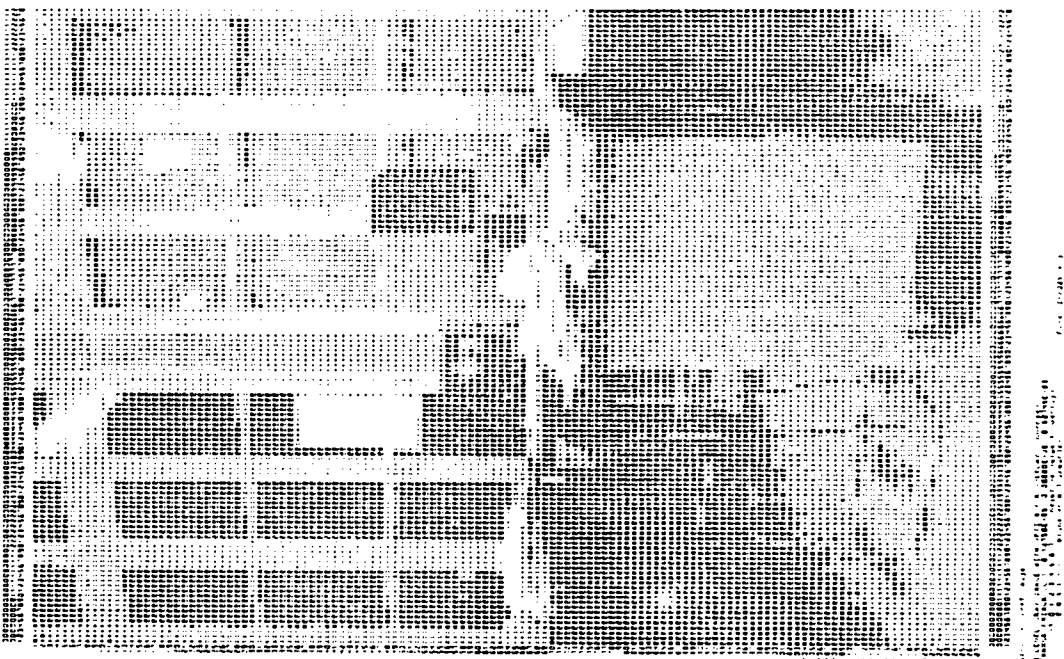
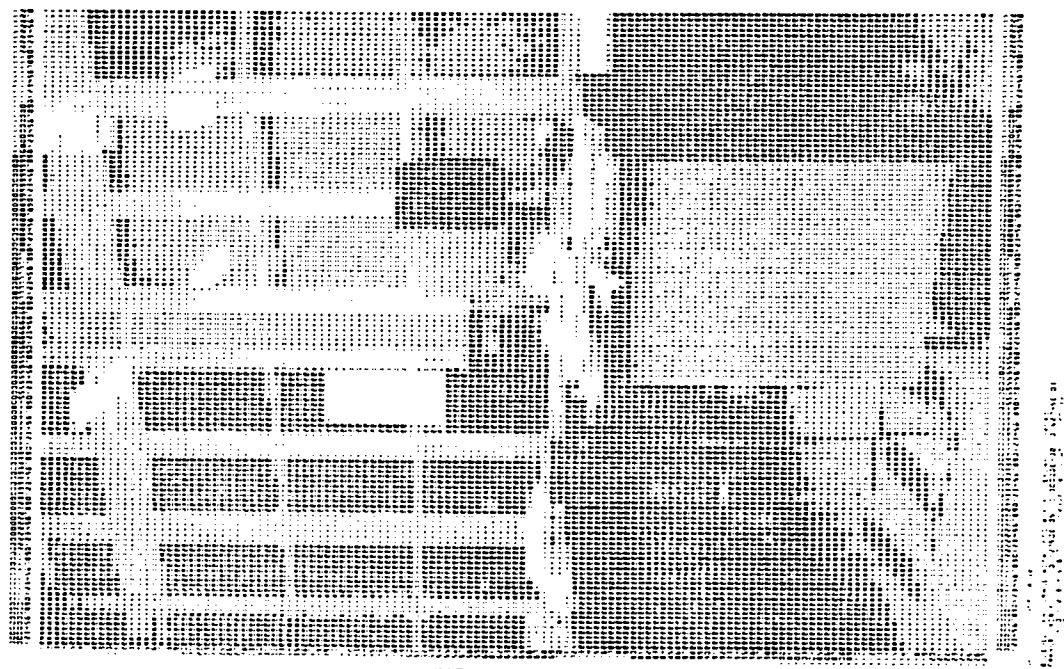
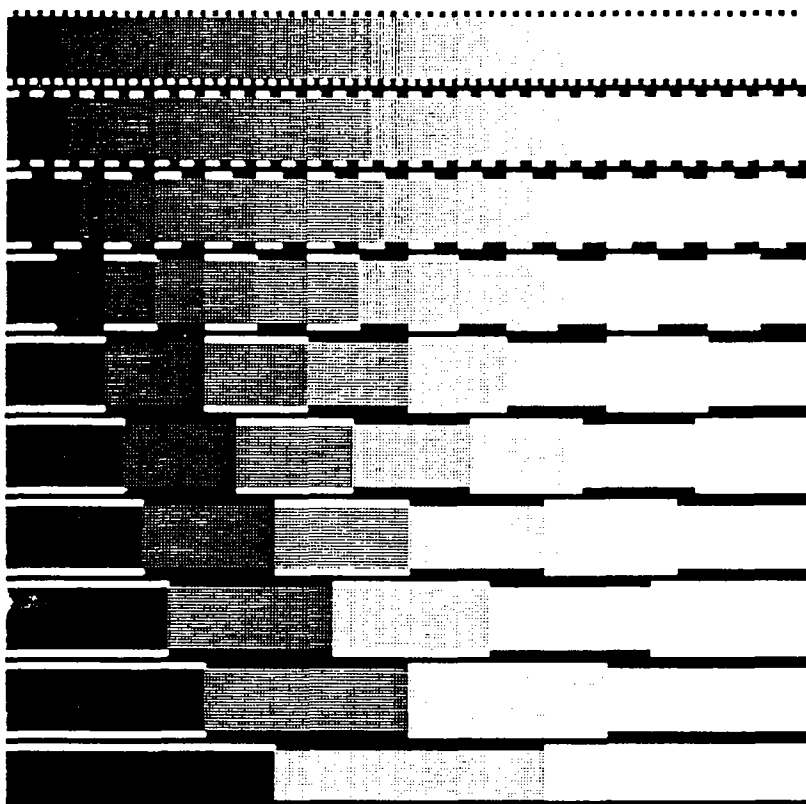
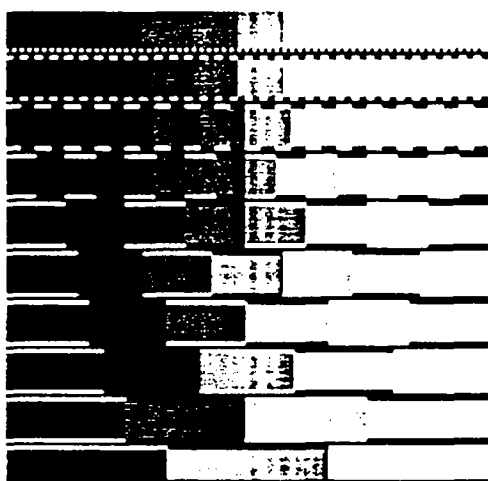


Fig.: 2.1-5.



if8gray.bar

Fig.: 2.1-6a.



if8gray.bar

Fig.: 2.1-6b.

2.2 The image pair

Even though this image pair (I001g&d_o) is adequate for human vision and the "room scene" appears in its correct shape when the images are viewed through stereo glasses, this pair contains numerous aberrations. In brief, the difficulties that became apparent during the analysis were of the following type, where I001g_o (or its decimated version) will be referred to as "left image", "left", or "L", and I001d_o (or its decimated version) will be referred to as "right image", "right", or "R":

1. The gray levels in the two images are rather different. Histograms of the gray levels normalized to a histogram peak of 100 units are shown in Fig. 2.2-1. Normalization to the peak in an histogram tends to "over-dramatize" the difference between the histograms. Non-parametric statistical tests (for ex., Kolmogorov-Smirnov) were considered unnecessary.
2. Closer inspection of the images reveals that, for example, the lower right corner of L contains details which are not visible in the lower right corner of R. A loss of information has occurred in one image with respect to the other. Some of these losses cannot be recovered by filtering.
3. These two 256x256 images appear to be "extracts" from a larger pair of images (>256x256 bytes). The shift between L and R produced in this extraction is unknown and the image centre cannot be relied upon as the optical centre of the two cameras.
4. There does not appear to exist a "fixation point" where the optical axes of the two cameras cross. Consequently, the fixation point, if any, is unknown.
5. The rotation of the cameras around their optical axes is unknown, but appears to be small.
6. The calibration (focal length and camera displacement for the images) is unknown.
7. Reflections of light sources occur in several places, which is normal in ordinary environments.

Initially, some experiments were carried out in order to attempt to "balance" the gray levels in the two images, but the results proved futile due to the near-total loss of detail in some regions of the images. After further thought, however, all the defects in the two images and the missing information about the cameras were viewed as challenges rather than drawbacks. The analysis programs must be able to handle such situations since they occur normally and present no difficulties for our visual system.

2.3 Computable and computed features

It was conjectured in chapter 1 that stereo matching on the pixel level is only feasible if an hierarchical image analysis and image matching structure is used. In three steps, the processing may take the following form:

- A) Segmentation and matching on the "large" scale.
- B) Segmentation and matching of the "subsegments" within each segment obtained in (A).
- C) Matching of edge and pixel features within subsegments of (B).

In other words, hierarchical segmentation and matching is carried out on ever-decreasing regions until there are a sufficient number of pixel-level features to allow unique pixel matching within each sub-region. The only critical requirement is segment and feature similarity between the left and right images, not the meaning of the segments nor the meaning of the pixel features. Consequently, stereo matching or "depth vision" requires no "image understanding" and no "models of the scene" are necessary. This may contradict presently held popular opinions.

The digital image matrix is represented in the computer as a two-dimensional matrix of the form

	1	2	3	4	...	i	...	Idim	---	The i-coordinate
1	*	*	*	*	...	*	...	*		
2	*	*	*	*	...	*	...	*		
3	*	*	*	*	...	*	...	*		
.	*		
j	p		
.	*		
.	*		
Jdim	*	*	*	*	...	*	...	*		

↓
The j-coordinate

where a single pixel (*) will be referred to as (i,j) or "p". The "i" and "j" are image coordinates (the pixel address) of an image matrix Q(i,j). The "Q" may be any alphanumeric combination. For example, in G(i,j) the "G" refers to "gray level". G(p) is "short-hand" for G(i,j). If necessary, the left (gray level) image is represented as Gl(i,j) and the right image as Gr(i,j), where "l" and "r" imply "left" and "right".

In brief, the information available from the images may be classified as follows:

1. Features directly available at single pixel level from the input images. Normally these are intensity, colour, and "flicker", i.e., the gray level, spectral, and temporal

compositions. In the present case only the gray level $G(i,j)$ is available (in intensity range 0 to 255).

2. Features computable from a relatively small local neighborhood N around each pixel p . There is no clearly definable limit to the number of such features. Furthermore, the results depend on the size and shape of N . The more common features are:
 - a) Average gray level $G_{av}(i,j)$ in N .
 - b) The x and y gradients $G_x(i,j)$ and $G_y(i,j)$ in N as well as higher order derivatives.
 - c) "Edge" functions in N , i.e., various "edge detectors" which locate either the extremal value of the gradient in N or compute an "edge amplitude" in N .
 - d) Curvature or change of gradient direction in N .
 - e) Numerous linear or non-linear functions in N , for example, "busyness".
 - f) Co-occurrence matrices and other "structural descriptions" within N .
 - g) Statistical parameters in N , such as, variance, skew, (invariant) moments, etc.

Since there is no clearly definable upper bound to what could, in principle, be computed for a pixel p from the gray levels in N , a logical choice is to restrict the computations only to features that are:

1. Computable in the presence of existing noise.
2. Invariant to aspects in the image that are not wanted.
3. Statistically independent from each other.
4. Able to retain stereo integrity in L versus R .

Noise and improperly resolved details tend to be suppressed by using a relatively large N . Invariance depends on what is wanted and its importance depends on the subsequent processes. In Ref. 2.2 (table 1, page 86) it is shown that the gray level $G(i,j)$ and its first (partial) derivatives $G_x(i,j)$ and $G_y(i,j)$ are independent, and that the first derivative is independent of the second, etc. Verification of "stereo integrity" is carried out by printing the results in image form (whenever feasible) and inspecting the results with stereo glasses. Clearly, if our vision cannot integrate the results into a "3D scene" then subsequent matching by computer is going to be difficult or impossible.

For the present case, the local neighborhood N is chosen relatively large and as isotropic as feasible, see Fig. 2.3-1. The computed features at each pixel p or (i,j) are:

i) The x and y gradient images $G_x(i,j)$ and $G_y(i,j)$, and the corresponding gradient magnitude $G_{pm}(i,j)$ and gradient angle $G_{pa}(i,j)$ images. Gradient angle is measured with respect to the $x=i$ axis. The left and right gradient magnitude images $G_{pml}(i,j)$ and $G_{pmr}(i,j)$ are shown as a stereo pair in Fig. 2.3-2 for visual verification. Both the "black on white" and "white on black" versions are shown in this case (Figures 2.3-2a and -2b) in order to select the "easiest to see" version. The author preferred the "white on black" in most cases. Stereo integrity ("fidelity") is retained in a "ghostly" scene. The corresponding gradient angle images $G_{pal}(i,j)$ and $G_{par}(i,j)$ are shown in Fig. 2.3-3a and -3b. Only the four major directions are printed as $- / \mid \backslash$ ($- = 0 \pm 22.5$; $/ = 45 \pm 22.5$; $\mid = 90 \pm 22.5$; $\backslash = 135 \pm 22.5$ degrees) for gradient magnitude values > 2.5 (in order to avoid the erratic directions at low gradient values).

ii) An "edge amplitude" function $E_f(i,j)$ which indicates "the amount of edge" present at pixel (i,j) .

$$E_f(i,j) = [|G'(i,j)| - W * |G''(i,j)|]_+$$

where, $|G'(i,j)| = G_{pm}(i,j)$ = gradient magnitude, $|G''(i,j)|$ = the magnitude of the second derivative in the direction of the gradient angle $G_{pa}(i,j)$, W = a positive weight, and the $[\cdot]_+$ indicates $\text{If}(E_f(i,j) > 0)$ Then $E_f(i,j) = [\cdot]$ Else $E_f(i,j) = 0$ ($>$ means "greater than"). It should be noted that the weight W can be used to "sharpen" the edge function (until E_f becomes so narrow that it may "vanish" at but not between the pixels). Numerous more formally defined "edge functions" may be found in (2.3). The edge functions $E_{fl}(i,j)$ and $E_{fr}(i,j)$ are shown as a stereo pair in Fig. 2.3-4. Again, stereo fidelity remains in a "ghostly" manner.

iii) A curvature function $C(i,j)$ at values where the edge function $E_f(i,j)$ is nonzero. In essence, $C(i,j)$ represents the change of the direction of the gradient vector within N .

iv) The "automatic volume control" type of contrast enhancement

$$G^{\sim}(i,j) = [G(i,j) - G_{av}(i,j)] / [G_{std}(i,j) + \text{Const}]$$

where $G_{av}(i,j)$ is the average gray in N , $G_{std}(i,j)$ is the standard deviation of $G(i,j)$ in N , and Const is a constant. A stereo pair of $G^{\sim}(i,j)$ is shown in Fig. 2.3-5 ($\text{Const} = 1.0$) and a "larger version" in Fig. 2.3-6 in order to show the "texture" better. (There is a "name" and considerable theoretical analysis of this type of filter, presumably in astronomy, but the author has forgotten the references.)

At the moment the $Ef(i,j)$ and $G\tilde{(i,j)}$ images have only been used to verify stereo consistence ("integrity", "fidelity") for these types of features. The $Gpm(i,j)$ and $Ef(i,j)$ images do present initially a "ghostly" appearance, presumably since we are not used to such images and they have been highly degraded by the low intensity amplitude ("gray") resolution in these prints. $G\tilde{(i,j)}$ is interesting since it tends to confirm the hypothesis that hierarchical segmentation is feasible.

Since these processes are well known, the interest in these experiments centered on stereo fidelity. Clearly, at least some of the features computable from the local neighborhood N retain stereo integrity. As a question to think about: "If stereo integrity were not maintained, is it meaningful to proceed"?

Figure titles:

- Fig. 2.3-1: The neighborhood N used in the present studies. (File: Inr002.fig)
- Fig. 2.3-2a and -2b: The left and right gradient magnitude images $Gpml(i,j)$ and $Gpmr(i,j)$ shown as a stereo pair. (a) Printed as black on white. (b) Printed as white on black.
- Fig. 2.3-3a and -3b: The gradient angle images $Gpal(i,j)$ and $Gpar(i,j)$. Only the four major directions are printed as - / | \ (- = 0 +/-22.5; / = 45 +/-22.5; | = 90 +/-22.5; \ = 135 +/-22.5 degrees) for $Gpm(i,j)$ values > 2.5 (in order to avoid the erratic directions at low gradient values). (Files: Inr041.fig and Inr042.fig)
- Fig. 2.3-4: The edge functions $Efl(i,j)$ and $Efr(i,j)$ are shown as a stereo pair.
- Fig. 2.3-5: The "automatic volume control" functions $G\tilde{l}(i,j)$ and $G\tilde{r}(i,j)$ shown as a stereo pair.
- Fig. 2.3-6: The "automatic volume control" functions $G\tilde{l}(i,j)$ and $G\tilde{r}(i,j)$ shown as a "magnified" stereo pair.

JI=	-10	-9	-8	-7	-6	-5	-4	-3	-2	-1	0	1	2	3	4	5	6	7	8	9	10
10	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
13	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
14	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
16	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
17	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
18	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
19	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
20	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
21	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
22	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
23	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
24	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
25	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
26	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
27	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
28	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

I = Lookup table address.
 dI = I-coordinate increment.
 dJ = J-coordinate increment.
 R = Radius from (0,0) to pixel at dI,dJ.
 A = Angle (degrees) from (0,0) to pixel at dI,dJ.

I	dI	dJ	R	A
(K)	(K)	(K)	(K)	(K)
1	-1	0	0.100000E+01	0.190000E+03
2	0	-1	0.120000E+01	0.320000E+03
3	0	1	0.150000E+01	0.270000E+03
4	1	0	0.180000E+01	0.100000E+03
5	-1	-1	0.1414214E+01	0.132000E+03
6	1	-1	0.1414214E+01	0.225000E+03
7	1	1	0.1414214E+01	0.450000E+03
8	1	1	0.1414214E+01	0.315000E+03
9	-2	0	0.200000E+01	0.180000E+03
10	0	-2	0.200000E+01	0.390000E+03
11	0	2	0.200000E+01	0.270000E+03
12	2	0	0.200000E+01	0.600000E+03
13	-2	-1	0.2236068E+01	0.1534350E+03
14	-2	1	0.2236068E+01	0.2055650E+03
15	-1	-2	0.2236068E+01	0.1135651E+03
16	-1	2	0.2236068E+01	0.2434349E+03
17	1	-2	0.2236068E+01	0.4343437E+03
18	1	2	0.2236068E+01	0.2955651E+03
19	2	-1	0.2236068E+01	0.2655651E+03
20	2	1	0.2236068E+01	0.3334343E+03
21	-2	-2	0.2828427E+01	0.1350000E+03
22	-2	2	0.2828427E+01	0.2250000E+03
23	2	-2	0.2828427E+01	0.4500003E+03
24	2	2	0.2828427E+01	0.3150000E+03
25	-3	0	0.3000000E+01	0.1800000E+03
26	0	-3	0.3000000E+01	0.3999998E+03
27	0	3	0.3000000E+01	0.2700000E+03
28	3	0	0.3000000E+01	0.6000000E+03

File: Inc002.Fig

Fig.: 2.3-1.



if8gradm.lgm

if8gradm.rgm

Fig.: 2.3-2a.



if8gradm.lgm

if8gradm.rgm

Fig.: 2.3-2b.



Figs.: 2.3-3a and -3b.



if8edgefl.lcf

if8edgefl.ref

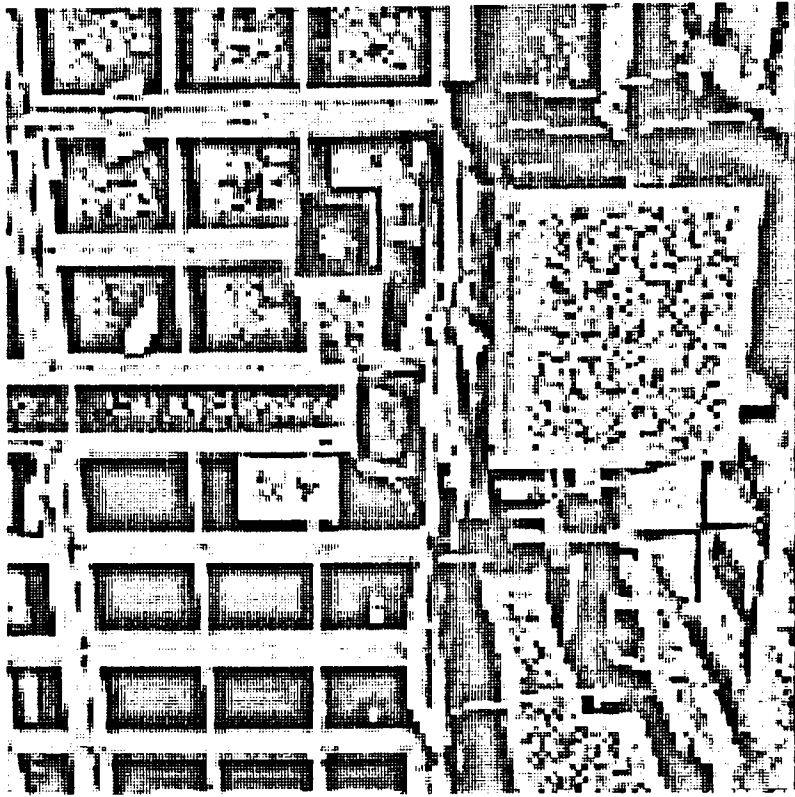
Fig.: 2.3-4.



if8volc2.lgr

if8volc2.rgr

Fig.: 2.3-5.



if8vole2.rqr



if8vole2.lgr

Fig.: 2.3-6.

2.4 Hierarchical feature classification

The next stage in the processing of the left and right images is to detect larger "entities" in the images which should be easier to match than the individual pixel features. The word "entity" is used in the sense of "something larger than a pixel". At the moment the choice is between two kinds of "entities", namely, the one-dimensional entities called "edges", and the two-dimensional entities called "regions". The choice is thus between:

1. "Edge detection", followed by assembling the "edge pieces" into longer edges.
2. "Region detection", consisting of grouping selected pixel features into larger wholes.

In both cases a profusion of techniques is available and, at least in principle, edge and region detection operations are complementary since edges surround regions or regions are surrounded by edges. In practice, however, the contours (edges) around "regions" are automatically guaranteed to be continuous closed (digital) curves. If the "edges" are detected first, then the result is usually a profusion of "edge pieces" which are difficult or impossible to assemble into unique continuous closed curves, except in highly contrasting images. In region detection frequently many small "noise regions" are obtained which partially negate the advantages of a region-based approach. A "contour follower" is a compromise between the two methods but it has its own "problems".

The "region-based" approach was chosen for the following reasons:

- a) Much work has already been done on the edge-based techniques. It is difficult to improve the results and meaningless to repeat these efforts. Furthermore, "edge matching" is difficult in the absence of "additional information".
- b) Regions should be easier to match in the left and right images than edges and the edges are easily obtained after regions have been detected. The regions provide the "additional information" for edge matching without the need of any other sources of "knowledge".
- c) The region-based approach coincides with the endeavors in the SYNTIM project.

Among the many methods available for "segmentation", a modified version of the "potential function" method was adapted (long ago) for low-dimensional decision spaces. The decision space is quantized and tabulated in the computer and then processed to become a look-up table for pixel feature based classification.

This leaves some apparently very important questions unanswered, namely:

1. Which features should be used for segmentation?
2. What is the criterion for judging segmentation results?

Both questions are fundamentally trivial since the left and right images are of the same scene but from slightly different "points of view". Consequently, if the same feature extraction algorithms are applied to the two images, the extracted features are expected to be rather similar (except in some restricted regions in the images where, for example, reflections occur and regions are visible in one but not in the other image). Thus, the only important criterion is that the features form clusters in the decision space and that these clusters represent relatively large connected (or detectable) regions in the image space. Neither the features nor the segments need to have any meaning in the human sense, as long as stereo integrity is maintained in the segmentation.

In the present experiments the pixels were classified according to their unfiltered gray levels $G(i,j)$, the magnitude of the spatial gray level gradient $Gpm(i,j)$, and the direction angle $Gpa(i,j)$ of the gray level gradient with respect to the $x=i$ axis, i.e.:

The left image is characterized by:

$G_l(i,j)$ = Original gray levels in the reduced image of I001g_o.
 $G_{pml}(i,j)$ = Gradient magnitude computed over N from $G_l(i,j)$.
 $G_{pal}(i,j)$ = Gradient angle computed over N from $G_l(i,j)$.
 $G_{xl}(i,j)$ = Gradient in x-direction computed over N from $G_l(i,j)$.
 $G_{yl}(i,j)$ = Gradient in y-direction computed over N from $G_l(i,j)$.

And the right image is characterized by:

$G_r(i,j)$ = Original gray levels in the reduced image of I001d_o.
 $G_{pmr}(i,j)$ = Gradient magnitude computed over N from $G_r(i,j)$.
 $G_{par}(i,j)$ = Gradient angle computed over N from $G_r(i,j)$.
 $G_{xr}(i,j)$ = Gradient in x-direction computed over N from $G_r(i,j)$.
 $G_{yr}(i,j)$ = Gradient in y-direction computed over N from $G_r(i,j)$.

The features computed from any input image are stored in image form and always kept in registration with the input image. The $G_x(i,j)$ and $G_y(i,j)$ images are used later in analytic relaxation.

The original images are assumed to form a stereo pair but image content is assumed unknown. The number of feature classes in the decision space, the shape of the clusters (classes), and their extraction is to be "automatic", i.e., without operator intervention.

The processing consists of the following steps:

1. Individual histogramming of the gray level images $G_l(i,j)$ and $G_r(i,j)$, and gradient magnitude $G_{pml}(i,j)$ and $G_{pmr}(i,j)$ images to establish upper and lower bounds for significant data. The gradient angle $G_{pa}(i,j)$ bounds are 0 and 360 degrees.
2. A two-dimensional decision space (2D histogram) $H(G_{pm},G)$ constructed from both the left and right images scaled to the bounds obtained from step 1. $H(G_{pm},G)$ is a two-dimensional histogram of $G_{pm}(i,j)$ versus $G(i,j)$ quantized to 73 by 73 intervals and tabulated in the computer. The left and right image features are amalgamated in $H(G_{pm},G)$. An one character per "pixel" alpha-numeric print of $H(G_{pm},G)$ is shown in Fig. 2.4-1. Each "pixel" value is represented by 0, 1, ..., 9, A, B, ..., Z, (36 levels, see the scale at the bottom on the figure). In order to show the "less populated" slot values in $H(G_{pm},G)$, the histogram peaks were clipped at 100.
3. Dynamic thresholding of $H(G_{pm},G)$ and conversion to a binary version $H_b(G_{pm},G)$.
4. Labelling of connected clusters in $H_b(G_{pm},G)$.
5. Label "spreading" since dynamic thresholding (in step 3) only "captures" the local peaks in $H(G_{pm},G)$.
6. Application of a constraint after label "spreading". In the present case the constraint consist of only retaining the clusters in $H(G_{pm},G)$ that contain low (and zero) values of the spatial gradient magnitude G_{pm} . After application of the constraint the remaining clusters are relabelled. The labelled decision space $H_l(G_{pm},G)$, which is now a look-up table, is shown in Fig. 2.4-2, where the characters represent cluster labels (C=2, D=3, etc.). Cluster labels are 2, 3, 4, etc.

Comments: The major variables that control the processes from steps 1 to 6 are:

a) The histogram bounds within which the data values are considered to be significant. Usually 1% of the upper and lower histogram areas of $G(i,j)$ and $G_{pm}(i,j)$ are rejected as "outliers".

b) The size of the tabulated decision space $H(G_{pm},G)$. The 73x73 quantization of the histogram is a "default" value from 128x127 range data classification. Clearly, cluster formation depends on the size of the decision space but this dependance is not critical and is also controllable by the parameters in step 3.

c) Dynamic thresholding consists of regularization of $H(Gpm, G)$ with two different kernels, subtraction of the results, and thresholding. The parameters are the "sigmas" of the two Gaussian regularization kernels and the subsequent threshold level. Prior investigations (for range images) indicated that these are not critical and, consequently, the old "default" parameters were retained.

7. Classification of the left and right images based on $Hl(Gpm, G)$ and the gray levels $G(i, j)$ and gradient magnitudes $Gpm(i, j)$. These classified images were only used as "masks" to select the unclassified pixels for the next level of classification (see step 8).
8. A two-dimensional decision space $H(Gpm, Gpa)$ amalgamating the gradient magnitudes $Gpm(i, j)$ and gradient angles $Gpa(i, j)$ for the unclassified pixels in step 7 for both the left and right images. The alpha-numeric print of $H(Gpm, Gpa)$ is shown in Fig. 2.4-3. The 73x73 quantization was a "default".
9. Dynamic thresholding of $H(Gpm, Gpa)$ using the same parameters as in step 3. This gives $Hb(Gpm, Gpa)$.
10. Constraints and labelling adapted to stereo image pairs and subsequent processes.

Comments, continued: Assume that the cameras are horizontally displaced. In order to distinguish between horizontal, vertical, and +/- 45 degree directions (with respect to the plane of the two cameras), artificial clusters were created in the thresholded (binary) decision space $Hb(Gpm, Gpa)$. The horizontal and vertical "directional preferences" could be dictated by gravity, but the others may be purely "cultural effects". Other types of "major directions" are equally easy to impose. The clusters were labelled as in step 4 and the labels were "spread" as in step 5. Another constraint consisted in separating the region in $Hb(Gpm, Gpa)$ corresponding to low values of the gradient magnitude $Gpm(i, j)$ since the corresponding gradient directions $Gpa(i, j)$ are rather erratic and cannot be classified correctly. This region was assigned a "don't care" or "take-me" label (-1). The resultant labelled $Hl(Gpm, Gpa)$ is shown in Fig. 2.4-4, where the "take-me" label is indicated by @. The parameters that control steps 5 and 9 are the same as those for processing $H(Gpm, G)$. The new constraints are:

- d) The number of orientations into which the $Hb(Gpm, Gpa)$ space is split. In the present case the 45 degree interval was chosen.
- e) The subsequent processing was simplified by grouping the 45 degree directions into four major directions, i.e., horizontal, vertical, and +/- 45 degree tilts. For further details, see "Relaxation".

f) The size (width) of the "don't care" or "take-me" region in $H1(Gpm, Gpa)$. A slight improvement in classification was obtained by using the "take-me" labels, see "Post processing" for details.

11. Classification of the remaining pixels in the left and right images using $H1(Gpm, Gpa)$.

12. Combining the two classifications. Since the pixels classified by $H1(Gpm, G)$ and $H1(Gpm, Gpa)$ are mutually exclusive, the classified images may simply be merged (if the labelling is continued from $H1(Gpm, G)$ to $H1(Gpm, Gpa)$). The raw results of classification after absorption of the "don't care" pixels are shown in Figures 2.4-5a and -5b for the left and right images.

The processes described in steps 1 to 12 are summarized in block diagram form in Figures 2.4-6a to -6h.

Comments, continued: It should be noted that the clusters found by using $H(Gpm, G)$ represent relatively uniform gray level regions in the left and right images, while the clusters from $H(Gpm, Gpa)$ represent relatively narrow and elongated regions with "cylindrically distributed" gray level gradient directions. The constraint that unified the 0 and 180, 45 and 225, etc., degree directions was deliberate.

g) The mistakes in the classification are of the following types:

i) A rather sharp reflection of a point light source which creates an approximately spherical light distribution is split into four or more regions. Such a region has known class label distribution but no attempts have been made to collect these into a larger regions.

ii) If a narrow (and long) region is wider than the approximate diameter of the local operator N used in computing the gradients, then the region is split into three narrow strips. One strip represents the centre region (uniform gray with low gradient) and the remaining two are "edge" regions where the gradients are high. However, in "critical" cases the splitting depends on minor variations in the gray levels and is thus erratic.

iii) Very "busy" areas which contain a more or less random distribution of high gray level gradients. Such areas are split into many small regions.

iv) A portion of one image (L or R) contains detectable regions while the corresponding region in the other image (R or L) is saturated (approximately 0 or 255 valued gray levels). Consequently, the regions found will be rather different.

Examples of all of these defects may be found in the classified $Lh(i,j)$ and $Rh(i,h)$ images, see Figures 2.4-5a and -5b.

13. Post processing. In practice the classification stage is followed by a certain amount of "post processing" since the $Lh(i,j)$ and $Rh(i,j)$ images contain:

- i) Unclassified pixels.
- ii) Pixels where the classification is uncertain (indicated by the "take-me" labels).
- iii) Wrongly classified pixels, presumably due to noise and quantization effects.

Clearly, the nature of the post processing will depend on the purpose of the classification, which indicates the nature of additional knowledge that could be introduced at this stage. In the present case the classification is intended to help in matching of the corresponding regions in the left and right images. Consequently, the regions are only made "slightly more uniform" by using "majority logic" in a 3 by 3 local neighborhood $N3$. The steps are essentially as follows:

- 1) Replace the "take-me" or "don't care" labels by the majority label in $N3$.
- 2) Replace the unclassified or 0-valued labels by the majority label in $N3$.
- 3) Replace the "unsupported" or "lonely" and "poorly supported" labels by the majority label in $N3$.

These processes may be iterated but it should be remembered these are "ad hoc" processes and should not be carried "too far". In the present case all the "take-me" labels were replaced if they could be replaced and the remaining ones were zeroed (since some of them may be isolated). This was followed by one iteration for replacing 0-valued labels and one iteration of replacing "lonely" and "poorly supported" labels. The results are shown in Figures 2.4-7a and -7b.

Several attempts were made to try to display the results for visual 3D (stereo) inspection and to represent the results as fairly as possible. The different class boundaries have to be shown but they can also be very disturbing (distracting) to our vision since points and sharp corners are highly noticeable. In the first attempt the gray levels for each class were replaced by their respective average values. Figure 2.4-8a shows the class average gray levels before post processing and Figure 2.4-8b shows them after post processing. The averages differ very little (except for class 1 which is the boundary around the image). The

stereo pair of average gray for each class is shown in Figure 2.4-9. It mostly shows the defects of the printer and no stereo effect is visible. However, on the Vicom film-loop display one could see (or possibly imagine) a stereo effect. The classes were "painted" in a minimal number of "colours". This is the 4-colour problem with the constraint that certain regions have to be assigned common colours. The adjacency matrix for the classes is shown in Figure 2.4-10a and the new assigned colours are in Figure 2.4-10b. Six colours were needed, which is somewhat beyond the capabilities of the laser printer (for printing a proper stereo pair) but the result is much "livelier". The stereo pair of class labelled images "painted" in the new colours is shown in Figure 2.4-11. There may be a hint of stereo in these images but it is very difficult to see. A blown-up version of the stereo pair is shown in Figure 2.4-12.

Comments, continued: The hierarchical classification of the pixel features (gray levels $G(i,j)$, gray level gradient magnitudes $Gpm(i,j)$, and gray level gradient directions $Gpa(i,j)$) produced two types of regions, namely:

- a) Regions of relatively uniform gray values and low gray level gradient magnitudes.
- b) Regions of highly varying gray levels with approximately (within ± 22.5 degree) uniformly directed gray level gradient angles grouped around 0--180, 45--225, 90--270, and 135--315 degree directions.

Clearly, the regions of type (a) may be considered "flat" (uniform) while those of type (b) have "cylindrically distributed" gray levels. The classification has also produced a form of segment recognition based on the illumination distribution over the segment, i.e., it is now known which segments are "flat" and which are "curved". The "curved" segments have been separated into four orientation directions. The post processing step may have introduced some pixels into the various regions which are of either type. The gray level gradients were computed over a relatively large local neighborhood N . Hence, the regions with highly varying gray levels tended to become magnified at the expense of the more uniform regions. This effect is unavoidable since a local operator always needs a local neighborhood N over which to "operate". In the present case this "magnification effect" was also enhanced intentionally (by choosing a large N) in order to assist in matching of the left and right image segments.

From this stage onwards several "avenues" have opened up, i.e., analytic representations of the illumination levels for regions, numerous "shape descriptions" for the segments, direct matching, and so on.

Figure titles:

- Fig. 2.4-1: An one character per "pixel" alpha-numeric print of $H(Gpm,G)$. Each "pixel" value is represented by 0, 1, ..., 9, A, B, ..., Z, see the scale at the bottom of the figure. In order to show the "less populated" slot values in $H(Gpm,G)$, the histogram peaks were clipped at 100.0 (highest value for the scale). (File: Inr009.fig)
- Fig. 2.4-2: The labelled decision space $Hl(Gpm,G)$, which is now a look-up table. The characters represent cluster labels (C=2, D=3, E=4, etc.). (File: Inr010.fig)
- Fig. 2.4-3: An one character per "pixel" alpha-numeric print of $H(Gpm,Gpa)$. (File: Inr013.fig)
- Fig. 2.4-4: The labelled $Hl(Gpm,Gpa)$ where the "don't care" or "take-me" label is indicated by @. (File: Inr014.fig)
- Figures 2.4-5a & -5b: The ("raw") left $Lh(i,j)$ and right $Rh(i,j)$ images of class labels. One character per pixel prints of the raw results of hierarchical classification after absorption of the "don't care" pixels. (Files: Inr017 and Inr018.fig)


```

000000001111111111222222222233333333334444444444555555555566666666667777
1234567890123456789012345678901234567890123456789012345678901234567890123
1 CCCCCCCC.....
2 CCCCCCCC.....
3 CCCCCCCC.....
4 CCCCCCCC.....
5 DDDDDDDDD.....
6 DDD.....
7 DDD.....
8 DDD.....
9 DDD.....
10 DDD.....
11 DDDDD.....
12 DDDDD.....
13 DDDDDDDDD.....
14 DDDDDDDDD.....
15 FFFFFFFF.....
16 FFFFFFFF.....
17 FFFFFFFF.....
18 FFFFFFFF.....
19 FFFFFFFF.....
20 FFFFFFFF.....
21 FFFFFFFF.....
22 FFFFFFFF.....
23 FFFFFFFF.....
24 FFFFFFFF.....
25 FFFFFFFF.....
26 FFFFFFFF.....
27 FFFFFFFF.....
28 FFFFFFFF.....
29 FFFFFFFF.....
30 FFFFFFFF.....
31 FFFFFFFF.....
32 FFFFFFFF.....
33 FFFFFFFF.....
34 FFFFFFFF.....
35 FFFFFFFF.....
36 FFFFFFFF.....
37 FFFFFFFF.....
38 FFFFFFFF.....
39 FFFFFFFF.....
40 FFFFFFFF.....
41 FFFFFFFF.....
42 FFFFFFFF.....
43 FFFFFFFF.....
44 FFFFFFFF.....
45 FFFFFFFF.....
46 FFFFFFFF.....
47 FFFFFFFF.....
48 FFFFFFFF.....
49 FFFFFFFF.....
50 FFFFFFFF.....
51 FFFFFFFF.....
52 FFFFFFFF.....
53 FFFFFFFF.....
54 FFFFFFFF.....
55 FFFFFFFF.....
56 FFFFFFFF.....
57 FFFFFFFF.....
58 FFFFFFFF.....
59 FFFFFFFF.....
60 FFFFFFFF.....
61 FFFFFFFF.....
62 FFFFFFFF.....
63 FFFFFFFF.....
64 FFFFFFFF.....
65 FFFFFFFF.....
66 FFFFFFFF.....
67 FFFFFFFF.....
68 FFFFFFFF.....
69 FFFFFFFF.....
70 FFFFFFFF.....
71 FFFFFFFF.....
72 FFFFFFFF.....
73 FFFFFFFF.....

```

```

000000001111111111222222222233333333334444444444555555555566666666667777
1234567890123456789012345678901234567890123456789012345678901234567890123

```

H1(Gpm,G), processed and labelled histogram of gradient magnitudes and gray levels.

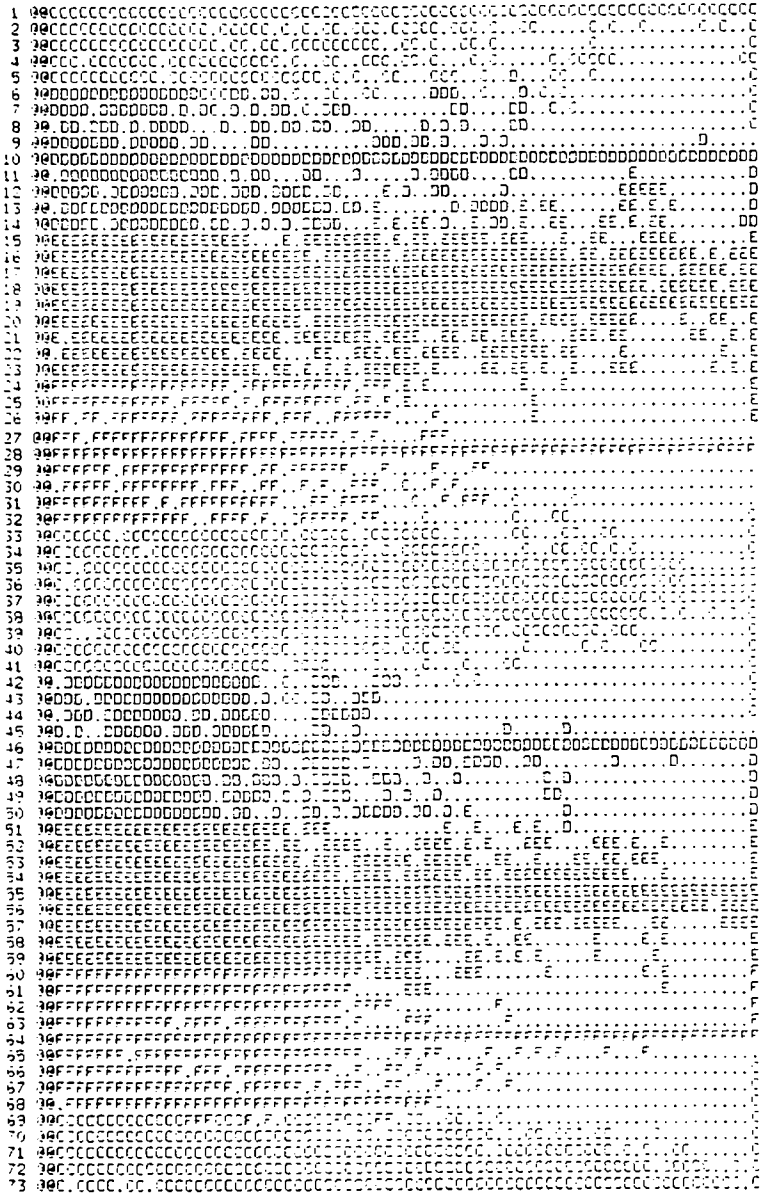
LABEL PRINTING

i= 0	e= 1	C= 2	D= 3	E= 4	F= 5	G= 6	H= 7
i= 8	J= 9	K= 10	L= 11	M= 12	N= 13	O= 14	P= 15
Q= 16	R= 17	S= 18	T= 19	U= 20	V= 21	W= 22	X= 23
Y= 24	Z= 25	[= 26	A= 27	B= 28	C= 29	D= 30	.ETC.
i= 0	j= -1	k= -2	l= -3	m= -4	n= -5	o= -6	p= -7
i= -8	j= -9	k= -10	l= -11	m= -12	n= -13	o= -14	p= -15
q= -16	r= -17	s= -18	t= -19	u= -20	v= -21	w= -22	x= -23
y= -24	z= -25	[= -26	a= -27	b= -28	c= -29	d= -30	.etc.

File: Inr010.fig

Fig.: 2.4-2.

000000001111111122222222333333334444444455555555666666667777
12345678901234567890123456789012345678901234567890123456789012345678901234567890123



000000001111111122222222333333334444444455555555666666667777
12345678901234567890123456789012345678901234567890123456789012345678901234567890123

H1(Gem. Ipa), processed, constrained, and labelled histogram of gradient magnitudes and gradient angles.

LABEL PRINTING

- 0 = 1 C= 2 D= 3 E= 4 F= 5 G= 6 H= 7
- 1= 8 J= 9 K= 10 L= 11 M= 12 N= 13 O= 14 P= 15
- Q= 16 R= 17 S= 18 T= 19 U= 20 V= 21 W= 22 X= 23
- Y= 24 Z= 25 C= 26 A= 27 B= 28 C= 29 D= 30 ,ETC.
- 0 = 1 C= -2 D= -3 E= -4 F= -5 G= -6 H= -7
- 1= -8 J= -9 K= -10 L= -11 M= -12 N= -13 O= -14 P= -15
- Q= -16 R= -17 S= -18 T= -19 U= -20 V= -21 W= -22 X= -23
- Y= -24 Z= -25 C= -26 A= -27 B= -28 C= -29 D= -30 ,etc.

File: Inr014.fig

Fig.: 2.4-4.

$G_l(i,j)$ & $G_r(i,j)$: Left and right gray level images.
 |
 $G_{pml}(i,j)$ & $G_{pmr}(i,j)$: L and R grad magnitude images.
 |
 Histogram
 |
 $H(G_{pm},G)$
 |
 Process, constrain, label
 |
 $H_l(G_{pm},G)$: Look-up table for classification.

Figure 2.4-6a: Process steps 1 to 6, first level of hierarchy.

$H_l(G_{pm},G)$: Look-up table for classification.
 |
 $G_l(i,j)$: Left gray level image.
 |
 $G_{pml}(i,j)$: Left grad magn image.
 |
 Classification
 |
 $Lh_1(i,j)$: Classified left image, level 1.

Figure 2.4-6b: Process step 7, L image, 1'st level of hierarchy.

$H_l(G_{pm},G)$: Look-up table for classification.
 |
 $G_r(i,j)$: Right gray level image.
 |
 $G_{pmr}(i,j)$: Right grad magn image.
 |
 Classification
 |
 $Rh_1(i,j)$: Classified right image, level 1.

Figure 2.4-6c: Process step 7, R image, 1'st level of hierarchy.

$Lh1(i,j)$ & $Rh1(i,j)$: Constraint, L and R classified images.
 |
 $Gpml(i,j)$ & $Gpml(i,j)$: L and R grad magnitude images.
 |
 $Gpal(i,j)$ & $Gpar(i,j)$: L and R grad angle images.
 |
 Histogram
 |
 $H(Gpm,Gpa)$
 |
 Process, constrain, label
 |
 $Hl(Gpm,Gpa)$: Look-up table for classification.

Figure 2.4-6d: Process steps 8 to 10, second level of hierarchy.

$Hl(Gpm,Gpa)$: Look-up table for classification.
 |
 $Gpml(i,j)$: Left grad magnitude image.
 |
 $Gpal(i,j)$: Left grad angle image.
 |
 $Lh1(i,j)$: Classified left image, constraint.
 |
 Classification
 |
 $Lh2(i,j)$: Classified left image, level 2.

Figure 2.4-6e: Process step 11, L image, 2'nd level of hierarchy.

$Hl(Gpm,Gpa)$: Look-up table for classification.
 |
 $Gpml(i,j)$: Right grad magnitude image.
 |
 $Gpar(i,j)$: Right grad angle image.
 |
 $Rh1(i,j)$: Classified right image, constraint.
 |
 Classification
 |
 $Rh2(i,j)$: Classified right image, level 2.

Figure 2.4-6f: Process step 11, R image, 2'nd level of hierarchy.

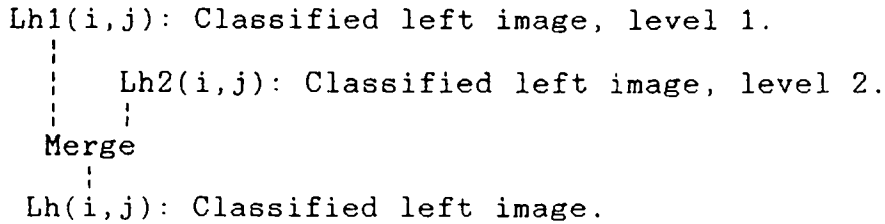


Figure 2.4-6g: Process step 12, merging left image classes.

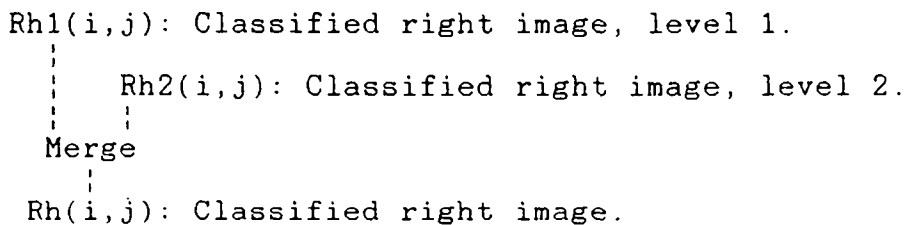


Figure 2.4-6h: Process step 12, merging right image classes.

Figure titles:

Figures 2.4-7a and -7b: The left $Lh(i,j)$ and right $Rh(i,j)$ images after post processing. The "take-me" labels were replaced if possible, followed by one iteration for replacing 0-valued labels and one iteration of replacing "lonely" labels. (Files: Inr019.fig and Inr020.fig)

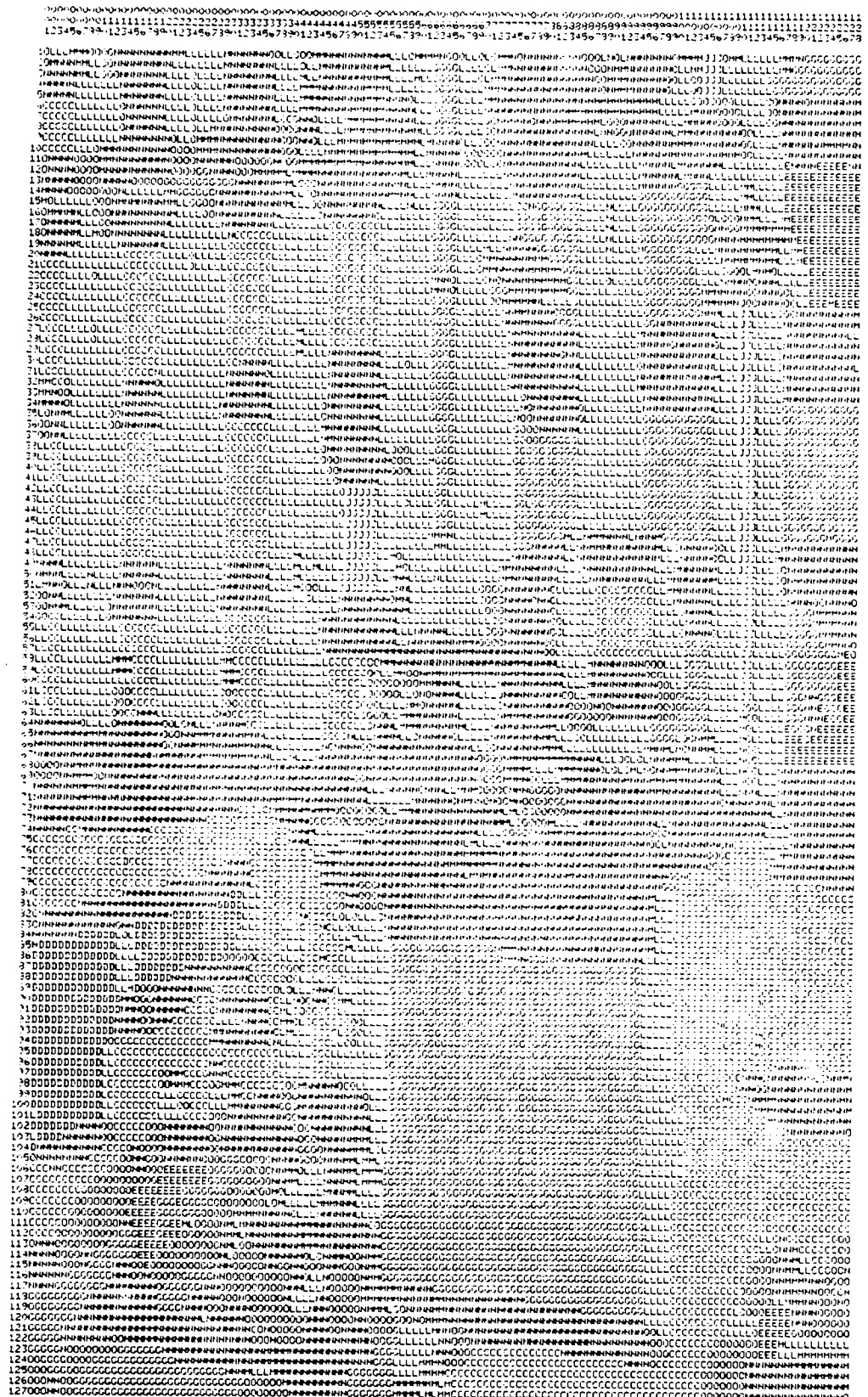
Figures 2.4-8a and -8b: The class average gray levels before post processing (-8a) and after post processing (-8b). (Files: Inr059.fig and Inr053.fig)

Figure 2.4-9: A stereo pair composed of average gray levels for each class.

Figures 2.4-10a and -10b: The adjacency matrix for the classes (-10a) and the new assigned colours (-10b). Six colours were needed. (Files: Inr054.fig and Inr055.fig)

Figure 2.4-11: A stereo pair of class labelled images "painted" in six colours.

Figure 2.4-12: A magnified version of the stereo pair shown in Fig.2.4-11.



Rh(i,j), right, classified pixels, one cycle of "remove -s" and one cycle of "majority label", conn=4, maj=1=1.

LABEL PRINTING

i = 0	j = 1	k = 2	l = 3	m = 4	n = 5	o = 6	p = 7
q = 8	r = 9	s = 10	t = 11	u = 12	v = 13	w = 14	x = 15
y = 16	z = 17	aa = 18	ab = 19	ac = 20	ad = 21	ae = 22	af = 23
ag = 24	ah = 25	ai = 26	aj = 27	ak = 28	al = 29	am = 30	an = 31
ao = 32	ap = 33	aq = 34	ar = 35	as = 36	at = 37	au = 38	av = 39
aw = 40	ax = 41	ay = 42	az = 43	ba = 44	bb = 45	bc = 46	bd = 47
be = 48	bf = 49	bg = 50	bh = 51	bi = 52	bj = 53	bk = 54	bl = 55
bm = 56	bn = 57	bo = 58	bp = 59	bq = 60	br = 61	bs = 62	bt = 63
bu = 64	bv = 65	bw = 66	bx = 67	by = 68	bz = 69	ca = 70	cb = 71
cc = 72	cd = 73	ce = 74	cf = 75	cg = 76	ch = 77	ci = 78	cj = 79
ck = 80	cl = 81	cm = 82	cn = 83	co = 84	cp = 85	cq = 86	cr = 87
cs = 88	ct = 89	cu = 90	cv = 91	cw = 92	cx = 93	cy = 94	cz = 95
ca = 96	cb = 97	cc = 98	cd = 99	ce = 100	cf = 101	cg = 102	ch = 103
ci = 104	cj = 105	ck = 106	cl = 107	cm = 108	cn = 109	co = 110	cp = 111
cq = 112	cr = 113	cs = 114	ct = 115	cu = 116	cv = 117	cw = 118	cx = 119
cy = 120	cz = 121	ca = 122	cb = 123	cc = 124	cd = 125	ce = 126	cf = 127
cg = 128	ch = 129	ci = 130	cj = 131	ck = 132	cl = 133	cm = 134	cn = 135
co = 136	cp = 137	cq = 138	cr = 139	cs = 140	ct = 141	cu = 142	cv = 143
cw = 144	cx = 145	cy = 146	cz = 147	ca = 148	cb = 149	cc = 150	cd = 151
ce = 152	cf = 153	cg = 154	ch = 155	ci = 156	cj = 157	ck = 158	cl = 159
cm = 160	cn = 161	co = 162	cp = 163	cq = 164	cr = 165	cs = 166	ct = 167
cu = 168	cv = 169	cw = 170	cx = 171	cy = 172	cz = 173	ca = 174	cb = 175
cc = 176	cd = 177	ce = 178	cf = 179	cg = 180	ch = 181	ci = 182	cj = 183
ck = 184	cl = 185	cm = 186	cn = 187	co = 188	cp = 189	cq = 190	cr = 191
cs = 192	ct = 193	cu = 194	cv = 195	cw = 196	cx = 197	cy = 198	cz = 199

File: Im020.fig

Fig.: 2.4-7b.

Replace label region by average gray from gray lev image

K	L	R	C	NNFL	NNFR	SSGL	SSGR
1	*	*	1	503	525	51.1	50.7
2	C	C	1	1445	2386	21.0	21.4
3	D	D	1	456	323	18.1	17.7
4	E	E	1	14	256	29.6	29.9
5	F	F	1	13	1	20.1	20.1
6	G	G	1	2303	2373	36.4	36.1
7	H	H	1	29	29	31.4	31.4
9	J	J	1	221	150	104.3	102.0
11	L	L	1	4081	3793	55.6	53.9
12	M	M	1	361	393	48.3	47.8
13	N	N	1	4170	4004	41.3	40.8
14	O	O	1	1570	1413	44.7	42.4

SCALEMAX = 0.20001E+01 0.12789E+03

Average gray levels for classes before cost processing.
KPAW classification results:

Incr017.fig = Class labels, left.

Incr018.fig = Class labels, right.

K = Class label.
L = Class label, left image.
R = Class label, right image.
C = 1 if both labels in left and right images, else 0.
NNFL = Number of pixels per class, left image.
NNFR = Number of pixels per class, right image.
SSGL = Average gray level of class, left image.
SSGR = Average gray level of class, right image.

File: Incr053.fig

Fig.: 2.4-8a.

Replace label region by average gray from gray lev image

K	L	R	C	NNFL	NNFR	SSGL	SSGR
1	*	*	1	503	525	51.1	50.7
2	C	C	1	1445	2438	21.0	21.4
3	D	D	1	457	319	18.1	17.7
4	E	E	1	14	256	29.6	29.9
5	F	F	1	13	1	20.1	20.1
6	G	G	1	2309	2426	36.4	36.1
7	H	H	1	29	29	31.4	31.4
9	J	J	1	221	150	104.3	102.0
11	L	L	1	4376	4001	55.6	53.9
12	M	M	1	372	344	48.3	47.8
13	N	N	1	4684	4414	41.3	40.8
14	O	O	1	1527	1284	44.7	42.4

SCALEMAX = 0.20001E+01 0.12789E+03

Average gray levels for classes after cost processing.

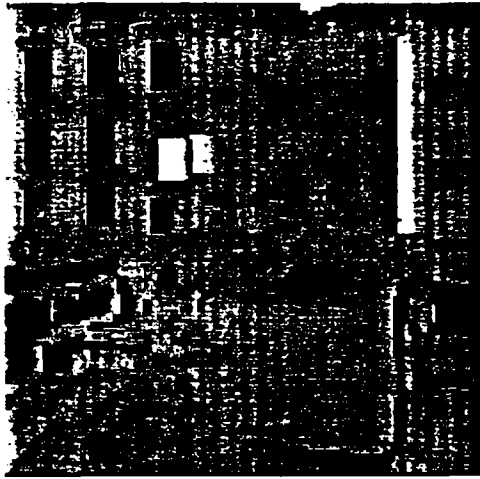
Incr019.fig = Class labels, left.

Incr020.fig = Class labels, right.

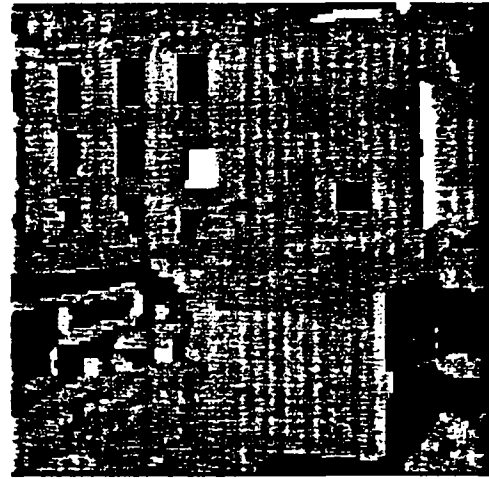
K = Class label.
L = Class label, left image.
R = Class label, right image.
C = 1 if both labels in left and right images, else 0.
NNFL = Number of pixels per class, left image.
NNFR = Number of pixels per class, right image.
SSGL = Average gray level of class, left image.
SSGR = Average gray level of class, right image.

File: Incr053.fig

Fig.: 2.4-8b.



if8hlgryl.lgr



if8hlgryl.rgr

Fig.: 2.4-9.



if8hlgry2.lgr



if8hlgry2.rgr

Fig.: 2.4-11.

```

00000000000000000000000000000000
0000000001111111111122222222223
123456789012345678901234567890

1 00000000011110000000000000000
11 001000000111100000000000000
3 010010000111100000000000000
4 000001000011110000000000000
5 001000000000000000000000000
6 000100000001111000000000000
7 000000000111100000000000000
8 000000000000000000000000000
9 000000000111100000000000000
10 000000000000000000000000000
11 111101101001100000000000000
12 111101101010110000000000000
13 111101101011010000000000000
14 111101101011100000000000000
15 000000000000000000000000000
16 000000000000000000000000000
17 000000000000000000000000000
18 000000000000000000000000000
19 000000000000000000000000000
20 000000000000000000000000000
21 000000000000000000000000000
22 000000000000000000000000000
23 000000000000000000000000000
24 000000000000000000000000000
25 000000000000000000000000000
26 000000000000000000000000000
27 000000000000000000000000000
28 000000000000000000000000000
29 000000000000000000000000000
30 000000000000000000000000000

```

Adjacency matrix for class labelled left and right images.

Inn019.fig = Class labels, left.
Inn020.fig = Class labels, right.

File: Inn054.fig

Fig.: 2.4-10a.

LS	LSL	LSR	LSL	LSR	LSL	LSR	LSL	LSR	LSL	LSR	LSL	LSR
LS = LLS(1) = 11=L												
LS = LLS(2) = 12=M												
LS = LLS(3) = 13=N												
LS = LLS(4) = 14=O												
LS = LLS(5) = 3=D												
LS = LLS(6) = 2=C												
LS = LLS(7) = 6=S												
LS = LLS(8) = 4=E												
LS = LLS(9) = 1=P												
LS = LLS(10) = 7=H												
LS = LLS(11) = 9=J												
LS = LLS(12) = 5=F												

Colour assignments. Six colours were needed.

LSO .. = Class label in images.
NCO .. = Number of contacts with other labels.
NCOOL = New colour assignment.

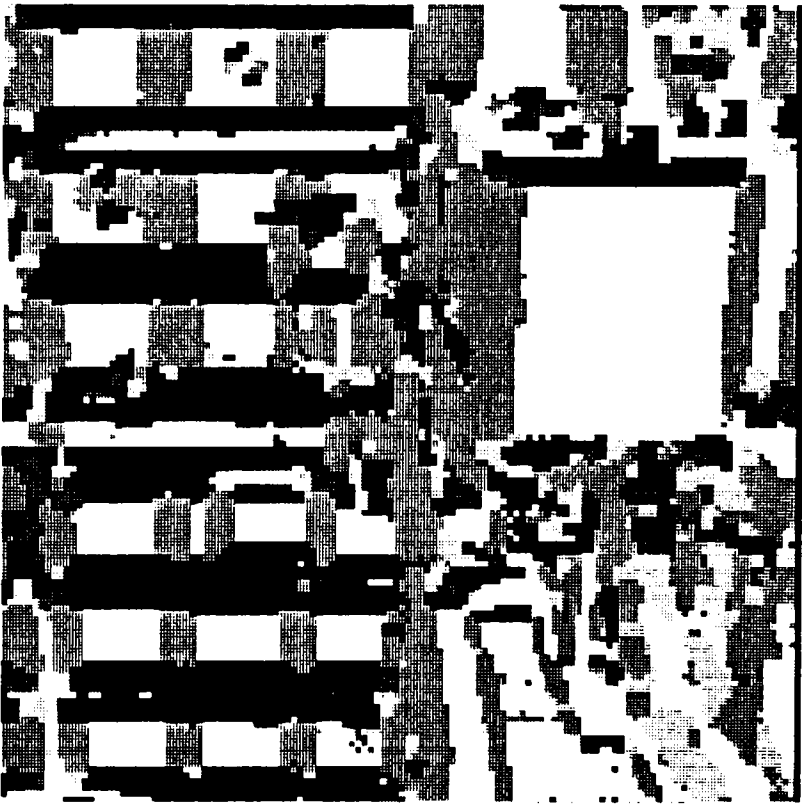
Inn019.fig = Class labels, left.
Inn020.fig = Class labels, right.

File: Inn055.fig

Fig.: 2.4-10b.



if8hlgr2.rgr



if8hlgr2.lgr

Fig.: 2.4-12.

2.5 Analytic approximation and relaxation

The regions or segments produced in the classification stage were of two types, namely, segments with approximately uniform gray levels, and segments with "cylindrically distributed" gray levels. In general, if an analytic approximation is desired, the characteristics upon which the original segmentation was based have to correspond to the form of the analytic function chosen to approximate the gray levels of the segments. In principle, the flat segments may be approximated by first order and the curved segments by a second order polynomials. In the present case this distinction could be made but was not used and all the segments were approximated by a second order polynomial of the form:

$$G_a(x,y) = A + Bx + Cy + Dx^2 + Exy + Fy^2$$

where $G_a(x,y)$ is the analytically computed illumination level of the segment, A, B, C, D, E, F are the coefficients to be computed, $x = i, y = -j$, and (i,j) are the pixel coordinates in the digitized image.

An analytic approximation of each segment (facet) serves the following purposes:

- a) All the pixels belonging to a facet (segment) contribute to the approximation. Thus, the analytic approximation serves as an "umbrella" for representing (the gray levels of) each facet.
- b) The "spreading effect" of the high gradient regions can be reduced with "relaxation".
- c) Functions computed from the coefficients may be used to recognize the facets (according to their gray levels, gradient directions, etc.).
- d) The image may be reconstructed from the analytic representations.

During these processing steps much additional information also becomes easily obtainable, for example, the size of each facet, the local thickness of each facet, numerous facet shape descriptions are possible, facet adjacencies (which facet is the neighbor of which other facet), and so on. The approximation and relaxation process is sketched in Fig. 2.5-1. The details are as follows:

For analytic approximation and "relaxation" the following processing steps are used:

1. Label the (4-connected) regions in the left $L_h(i,j)$ and right $R_h(i,j)$ classified images. This gives a unique "tag" or label for each connected region (facet). Clearly, size and

shape information is now immediately computable for each facet separately.

2. Remove facets which are too small to support a given analytic approximation. This process may be "finessed" by selecting approximations that are suitable for each facet or by removing only facets that cannot be approximated uniquely (for example, if all the pixels are on a straight line).
3. Relabel the remaining facets according to size. This is only a convenience feature to allow the facet labels to be used as addresses during computations and parameter tabulations. The (facet) labelled images are called $L_f(i,j)$ and $R_f(i,j)$ for the left L and right R images, respectively, and the "f" indicates a facet.
4. Compute the coefficients of the analytic approximation for each facet. In the present case the L1 norm was used in order not to overemphasize the large errors. The origin ($x=0, y=0$) of each facet may be placed at the image corner ($i=0, j=0$), at the centre of gravity of the facet, or where desired. The procedure consists of: Extract all pixels with a given label k, give these to the approximation program, receive the coefficients and tabulate them as $A(k), B(k)$, etc. It should be noted that this approximation is not an "elastic" approximation (a "spline") of the gray levels. Here the approximation of a facet is independent of any other facets (as well as of its adjacent neighbors).
5. This step is optional but should be carried out to "evaluate by inspection" the accuracy of the approximation, i.e.:
 - i) Compute the gray level values for the left $G_{al}(i,j)$ and right $G_{ar}(i,j)$ images from the analytic approximations and display the images. Here the "a" in $G_a(i,j)$ indicates "analytic".
 - ii) Compute the error images $E(i,j) = |G(i,j) - G_a(i,j)|$ and display the results.

However, it is advisable to carry out these "inspections" after some analytic relaxation of the images to remove minor errors (the regions removed in step 2 leave "holes" in $G_a(i,j)$ which are highly disturbing to our vision).

6. The analytic relaxation procedure is based on comparing the actual gray levels $G(i,j)$ (and their gradients $G_x(i,j)$ and $G_y(i,j)$, if desired) with the analytically computed gray levels $G_a(i,j)$ (and their gradients, if desired). The information on "which pixel (i,j) belongs to which facet k" is available from the corresponding facet labelled image $F(i,j)$ ($F(i,j) = L_f(i,j)$ for the left and $F(i,j) = R_f(i,j)$ for the right image). The parameters are available from the look-up tables $A(k), B(k), C(k), D(k), E(k), F(k)$ for the left and right image, respectively.

In brief, the following "input" information is used for the left and right image, respectively:

$G(i,j)$ = Gray level image.

$G_x(i,j)$ = The x-component of the gradient of $G(i,j)$.

$G_y(i,j)$ = The y-component of the gradient of $G(i,j)$.

$F(i,j)$ = Facet labelled image.

$A(k), B(k), C(k), D(k), E(k), F(k)$ = Parameter look-up tables.

$Gak(x,y) = A(k) + B(k)x + C(k)y + D(k)x^2 + E(k)xy + F(k)y^2$.

$Gak_x(x,y) = B(k) + 2D(k)x + E(k)y = x$ -gradient of $Gak(x,y)$.

$Gak_y(x,y) = C(k) + E(k)x + 2F(k)y = y$ -gradient of $Gak(x,y)$.

i) Self-relaxation: Clearly, whether or not the accuracy of the analytic approximation deviates "too much" from the actual values is easily detectable by computing the error $E(i,j)$ at pixel (i,j) . Thus, in brief:

For all (i,j)

$G(i,j)$ is the original gray level

$k=F(i,j)$ is the label for this pixel

$Gak(x,y) = A(k) + B(k)x + C(k)y + D(k)x^2 + E(k)xy + F(k)y^2$.

$E(i,j) = |G(i,j) - Gak(i,j)|$ is the error

If $(E(i,j) > \text{Threshold})$ then set $F(i,j) = 0$

The label in $F(i,j)$ is zeroed if the error $E(i,j)$ exceeds some threshold and this pixel becomes "unassigned". The unassigned connected pixel regions are also labelled, and the parameters are recomputed for all the facets. In the present context this is one of the processes that can "bring out more details within a facet" after the left and right facets have been matched.

ii) Cross-relaxation: This form of "relaxation" occurs at the boundaries of the facets. A 3 by 3 neighborhood N with the centre pixel p is sufficient. The centre pixel p is on the boundary of one of the facets. Within the 3×3 neighborhood N_3 (4-connected usually) there are pixels with other facet labels.

In brief, the gray level (and its derivatives, if desired) are computed for the centre pixel p from all the pixels (labels) in N that "participate" in the "competition" for assigning its label to the pixel p in $F(i,j)$. The label that gives the best match with respect to $G(i,j)$ (and its derivatives, if desired) is declared the "winner" and it puts its label at location p in $F(i,j)$.

If there are some labels within the neighborhood N_3 in $F(i,j)$, two basic situations occur, namely, the label at pixel $p = (i,j)$ is zero ($0 = \text{Label}(p) = F(i,j)$) and the label at p is non-zero ($k = \text{Label}(p) = F(i,j)$).

If Label(p) is 0 but there are non-zero labels in N3, then it is a question of "conquering unopposed territory" and the label that gives the lowest error term is declared the winner provided that the error is less than some threshold T. If T is large then the "best label wins" and after some iterations all the 0-valued (unassigned) pixels in F(i,j) become filled. This method was used to "fill the gaps" in F(i,j) which were created in step 2 since some of the facets were too small to allow a definite analytic approximation.

If Label(p) is not 0 and there are non-zero labels in N3, then it is a question of "conquering opposed territory". As before, the label that gives the lowest error term is declared the winner provided that the error is less than some threshold T. However, now the behaviour depends on how the error term is defined, on the value of T, and repeated iterations need not terminate with "the number of changed pixels equals zero".

"The number of changed pixels" decreases rapidly first but then "slows down" at small values and may even become slightly "oscillatory". The effect is easily explained by considering a "cylindrical" region which (for whatever reason) became approximated by three facets indicated by A, B, and C. A and C approximate the two sides of the "cylinder" while B approximates both sides of the "cylinder". A part of the contact region between facets A and C (indicated by C') may be best approximated by facet B and B starts to "invade" the contact between A and C, see Fig. 2.5-2.

The error term for a facet k at p=(i,j) is actually defined as a sum of gray level and gradient vector errors:

$$E_k(i,j) = |G(i,j) - G_{ak}(i,j)| + W*|G'(i,j) - G_{ak}'(i,j)|$$

where: $|G(i,j) - G_{ak}(i,j)|$ = Abs value of gray level error.
 $|G'(i,j) - G_{ak}'(i,j)|$ = Gray level gradient error.
 W = Weight to balance the two errors.

With the "best label in N3 wins" strategy, for a small W and large T, the gray levels are approximated as well as possible with the given set analytic coefficients and label distribution but the result may not be best for subsequent matching experiments. In Chapter 3 it is shown how the analytic relaxation process was used after the facets had been matched as well as possible.

However, in order to demonstrate the effects of relaxation, the labelled images in Figures 2.4-7a and -7b were "relaxed" resulting in images shown in Figures 2.5-3a and -3b. The label changes during the relaxation are shown in Fig. 2.5-4 (which is an edited extract of the "history file" during computing). The original parameter lists are given in Figures 2.5-5a and -5b. The parameters have been (in this case) computed with i=0 and j=0 as the origin of the coordinate system. In these lists "K" is the facet label and the "@" in "K=@" indicates the letter used in

printing the facet labelled images (for example, Fig. 2.5-3). The NNEL represents the number of pixels on the facet. NNHL=@ is intended to indicate the corresponding histogram (class) label, but is not used in the present list (set same as K). XX00 and YY00 give the centre of gravity of the facet (XX00=i, YY00=-j). Since the alphabetic code for label printing is modulu 26, the centre of gravity should be used to verify the label code in the image.

In order to show how the image "looks" to the computer the analytic parameters and the labels are used to create an analytic versions $G_{ana_l}(i,j)$ and $G_{ana_r}(i,j)$ of the left and right gray level images. The results are shown in Figure 2.5-6 before analytic relaxation and in Figure 2.5-7 after analytic relaxation. Black dots and "squarish" regions are created where the corresponding pixels have no label or the analytic approximation was "indefinite". Clearly, the relaxation process has "sharpened" the image, but it should be remembered that the initial emphasis was on creating regions that are easier to match rather than more pleasant to look at. Histograms of the gray level differences between the analytic and original images ($|G_{original}(i,j) - G_{analytic}(i,j)|$) are shown in Figure 2.5-8. Images of where the major errors occur are shown in "8-level alphanumeric gray" in Figure 2.5-9.

In order to estimate the effects of segmentation on "stereo fidelity" the analytically created gray level images are shown as stereo pairs in Figures 2.5-10a and -10b. The stereo effect is hardly visible in the raw approximation (Fig. 2.5-10a) but improves somewhat after relaxation (Fig. 2.5-10b). Our vision is greatly disturbed by the "noise points" (missing analytic values) and by the artificial gray level contours created in printing the images.

In order to create more "pleasant-looking" images, the analytically obtained gray level images were "corrected" by replacing all pixels in the analytic images with the corresponding gray level values from the original images if there was no analytic value or if the analytic value differed by more than 15 units. The gray level error of 15 was simply the first choice, being about 12% of the median gray. The pixels that were replaced are shown as "white" (.) in the labelled images (Figures 2.5-11a and -11b). Of course, one could compute new analytic approximations for these regions, if desired, and even correct the class label images $Lh(i,j)$ and $Rh(i,j)$. This may improve "matchability" in subsequent processing steps. The resultant "corrected" images are shown in Figure 2.5-12 on a "magnified" scale and as a stereo pair in Figure 2.5-13a. Except for the false gray level contours, most visible on the "desk end", the stereo effect is recaptured. For comparison, the original gray level image pair is shown in Figure 2.5-13b. However, even though the stereo effect has reappeared, this does not really mean very much since the facet boundaries are no longer visible.

Comments:

The analytic relaxation procedure was iterated on the same set of parameters as shown in Figure 2.5-4. In practice, during relaxation some regions may vanish, others become too small to allow proper analytic approximation, and some others may split into two or more separated connected regions (some of which are usually too small, etc.). Depending on the logic used in the relaxation, regions may also become "vacated" since none of the existing approximations apply. Hence, the regions should be relabelled, the small ones eliminated, and new labels computed according to the sizes of the remaining regions. Now the relaxation process may be reiterated. A block diagram of the process was shown in Figure 2.5-1.

In principle, specially if the analytic approximation uses several forms of analytic functions depending on the nature and size of the region (for example, 0-order for very small facets, 1-st order for "planar" regions, 2-nd order for "curved" regions) then the error can be reduced to very small values. However, in practice an additional form of instability, in addition to that illustrated in Figure 2.5-2 is likely to occur, where pixels on the borders of adjacent facets with very similar parameters will simply continue to "exchange labels". The remedy is to merge such regions first and then continue the iterations.

The analytic approximation and relaxation gives a set of coefficients that characterize the gray level distribution (or any other function computed over the image) for each facet. Since the view-point for a stereo pair of images must be different, the gray level distributions (and the parameters) cannot be exactly the same for the same facet in the left and right image. Consequently, the parameters cannot be compared directly, but several invariants may be computed from the parameters, such as the "curvatures" of the gray level distributions (2.4). Time has neither permitted a closer study of the variations between the parameters nor of the invariants.

A final warning is in order, i.e., if the segmentation (labelling, parameter computations, and relaxation) of the left and right image is done independently, there is no guarantee that the regions in the left and right stereo image are "matchable". In the present case the only constraints to ensure "matchability" were the common $H(G_{pm}, G)$ and $H(G_{pm}, G_{pa})$ decision spaces. However, this is a rather weak constraint to "tie together" the left and right images.

```

G(i,j) = Original gray level image
|
|   F(i,j) = Facet label image <---* (new iteration)
|   |
|   |   Parameter computation
|   |   |
|   |   |   A(.),B(....),F(.) = Parameter lists
|   |   |   |
|   |   |   |   G(i,j) = Original gray levels
|   |   |   |   |
|   |   |   |   |   Gx(i,j) = x-gradient image
|   |   |   |   |   |
|   |   |   |   |   |   Gy(i,j) = y-gradient image
|   |   |   |   |   |   |
|   |   |   |   |   |   |   Analytic relaxation (iterated n times)
|   |   |   |   |   |   |   |
|   |   |   |   |   |   |   |   F'(i,j) = Corrected facet label image
|   |   |   |   |   |   |   |   |
|   |   |   |   |   |   |   |   |   Relabelling (to identify split facets)
|   |   |   |   |   |   |   |   |   |
|   |   |   |   |   |   |   |   |   |   Remove small facets (ana approx impossible)
|   |   |   |   |   |   |   |   |   |   |
|   |   |   |   |   |   |   |   |   |   |   Relabel (according to size of facet)
|   |   |   |   |   |   |   |   |   |   |
|   |   |   |   |   |   |   |   |   |   |   F(i,j) = New facet labels
|   |   |   |   |   |   |   |   |   |   |   |
|   |   |   |   |   |   |   |   |   |   |   |   (Adjacent facet mergers if parameters similar)
|   |   |   |   |   |   |   |   |   |   |   |   |
|   |   |   |   |   |   |   |   |   |   |   |   |   Iterate m times ---->*
|   |   |   |   |   |   |   |   |   |   |   |   |   |
|   |   |   |   |   |   |   |   |   |   |   |   |   |   (Final parameter computation)
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   A(.),B(....),F(.) = Final parameter lists

```

Figure 2.5-1: The analytic approximation and relaxation procedure.

```

A A A A A A A A B B B B B B B   A A A A A A A A A B B B B B B
A A A A A A A A A B B B B B B   A A A A A A A A A B B B B B B
C C C C C C C C C B B B B B B   C C C C C C B B B B B B B B
C C C C C C C C C B B B B B B   C C C C C C C C C B B B B B B
C C C C C C C C C B B B B B B   C C C C C C C C C B B B B B B

    Before relaxation                After relaxation

```

Figure 2.5-2: An example of instability in analytic relaxation.

Figure titles:

Figure 2.5-3a and -3b: The facet labelled images $L_f(i,j)$ and $R_f(i,j)$ after 10 cycles of iteration on the initial parameters. (Files: Inr045.fig and Inr046.fig)

Figure 2.5-4: An edited extract of the "history record" during analytic relaxation. (File: Inr033.fig)

Figures 2.5-5a and -5b: The initial analytic parameters used in the relaxation process. (Files: Inr029.fig and Inr030.fig)

Figure 2.5-6: The analytically recreated gray level images $G_{ana_l}(i,j)$ and $G_{ana_r}(i,j)$ before relaxation.

Figures 2.5-7: The analytically recreated gray level images $G_{ana_l}(i,j)$ and $G_{ana_r}(i,j)$ before relaxation.

Figure 2.5-8: Histograms of gray level differences between the analytic and original images ($|G_{original}(i,j) - G_{analytic}(i,j)|$). (File: Inr049.fig)

Figure 2.5-9: Images of the major errors between original and analytic gray level images, shown in "8-level alphanumeric gray". (Files: Inr039.fig and Inr040.fig)

Figures 2.5-10a and -10b: The analytically created gray level images are shown as stereo pairs. (a) The "raw" approximation after segmentation. (b) Segmentation after analytic relaxation.

Figures 2.5-11a and -11b: The replaced pixels shown as "white" (.) in the labelled images in order to "correct" the analytic approximations for "easier stereo viewing". The gray level error is 15 units or about 12% of the median gray. (Files: Inr051.fig and Inr052.fig)

Figure 2.5-12: The resultant "corrected" analytic images shown on a "magnified" scale.

Figures 2.5-13a and -13b: Stereo pairs. (a) The "corrected" analytic images. (b) The original gray level image pair shown for comparison.

CROSS-CLEAN: Analytic relaxation of left image.

IOXFAP,IOXFAC,IOXG,IOXGX,IOXGY,IOXSCC = 35 33 11 37 39 41
 MON327R: READ INITIAL SURFACE PARAMETER LISTS: IOFP,JDIM = 35 250
 MON120: READ FILE: IVAX,JDIM,JDIM= 35 129 127
 MON120: READ FILE: IVAX,JDIM,JDIM= 11 129 127
 MON120: READ FILE: IVAX,JDIM,JDIM= 37 129 127
 MON120: READ FILE: IVAX,JDIM,JDIM= 39 129 127
 MON307A: ANALYTIC CROSS-CLEANING:
 IDIM,JDIM,KDIM,NCycle,KOBJECT,IDEHDR = 129 127 250 10 4 0
 INDEF,IOPT,IPT,JPT,LPOUT,LPE,NONE,P,LASTHKL = 1 -1 0 0 3 0 0 250
 WEIGHTG,WEIGHTD = 0.10000E+01 0.50000E+00
 EPLIM,CMUL = 0.25000E+02 0.10000E+03

Iteration cycles on initial parameters:
 CYCLE= 1 NUMBER OF CHANGED PIXELS= 0.75430E+04
 CYCLE= 2 NUMBER OF CHANGED PIXELS= 0.31000E+03
 CYCLE= 3 NUMBER OF CHANGED PIXELS= 0.12600E+03
 CYCLE= 4 NUMBER OF CHANGED PIXELS= 0.12500E+03
 CYCLE= 5 NUMBER OF CHANGED PIXELS= 0.35000E+02
 CYCLE= 6 NUMBER OF CHANGED PIXELS= 0.44000E+02
 CYCLE= 7 NUMBER OF CHANGED PIXELS= 0.24000E+02
 AUTO-EXIT: TOT TOTNC1= 0.65380E+02
 MON307A: TOTAL NUMBER OF PIXEL CHANGED = 0.48690E+04
 MON121: WRITE FILE: IVAX,JDIM,JDIM= 41 129 127

CROSS-CLEAN: Analytic relaxation of right image.

IOXFAP,IOXFAC,IOXG,IOXGX,IOXGY,IOXSCC = 36 34 10 33 40 42
 MON327R: READ INITIAL SURFACE PARAMETER LISTS: IOFP,JDIM = 36 250
 MON120: READ FILE: IVAX,JDIM,JDIM= 34 129 127
 MON120: READ FILE: IVAX,JDIM,JDIM= 10 129 127
 MON120: READ FILE: IVAX,JDIM,JDIM= 33 129 127
 MON120: READ FILE: IVAX,JDIM,JDIM= 40 129 127
 MON307A: ANALYTIC CROSS-CLEANING:
 IDIM,JDIM,KDIM,NCycle,KOBJECT,IDEHDR = 129 127 250 1 4 0
 INDEF,IOPT,IPT,JPT,LPOUT,LPE,NONE,P,LASTHKL = 1 -1 0 0 3 0 0 250
 WEIGHTG,WEIGHTD = 0.10000E+01 0.50000E+00
 EPLIM,CMUL = 0.25000E+02 0.10000E+03

Iteration cycles on initial parameters:
 CYCLE= 1 NUMBER OF CHANGED PIXELS= 0.35040E+04
 CYCLE= 2 NUMBER OF CHANGED PIXELS= 0.73500E+03
 CYCLE= 3 NUMBER OF CHANGED PIXELS= 0.21400E+03
 CYCLE= 4 NUMBER OF CHANGED PIXELS= 0.24000E+02
 CYCLE= 5 NUMBER OF CHANGED PIXELS= 0.50000E+02
 CYCLE= 6 NUMBER OF CHANGED PIXELS= 0.33000E+02
 AUTO-EXIT: TOT TOTNC1= 0.34178E+02
 MON307A: TOTAL NUMBER OF PIXEL CHANGED = 0.46430E+04
 MON121: WRITE FILE: IVAX,JDIM,JDIM= 42 129 127

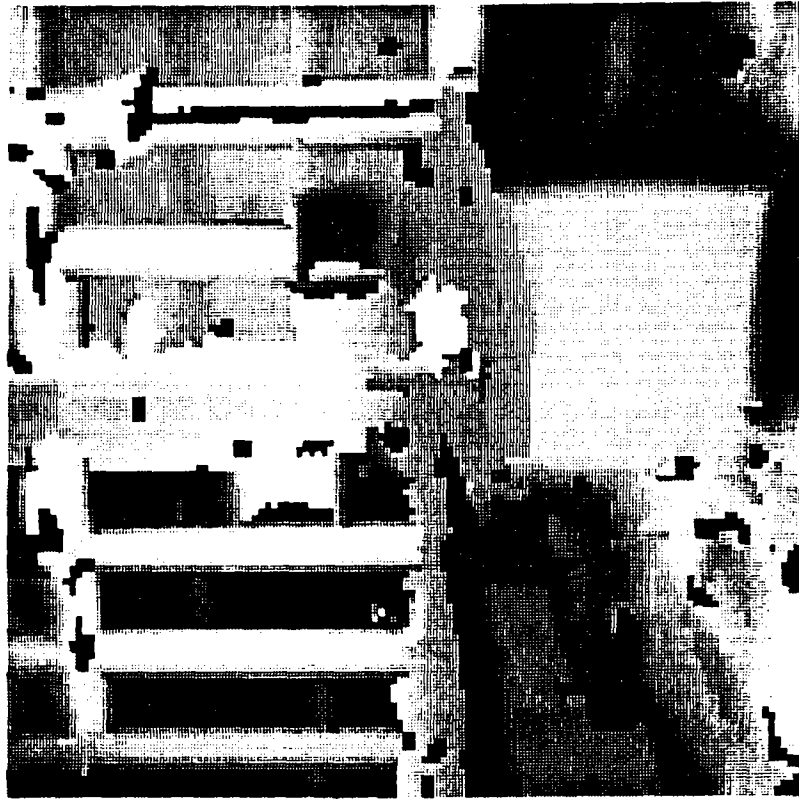
File: Inr033.Flg

Fig.: 2.5-4.

Right image facet parameters. IO:RIP, IO:RAC = 36 34
C(x,y) = 2 + 8*x + 4*y + 2*E*sqrt(x**2 + y**2) + 1.0
C = Facet label, NNEL = Number levels, R00L = Histo label
LM = 1 for valid cases, = 2 for indefinite pairs, = 0 none

Table with columns: K=J, NNEL, N00L, LM, C=O, C=OO, C=000, C=0000, C=00000, C=000000, C=0000000, C=00000000. Rows contain numerical data for various indices from 2=0 to 113=0.

Fig.: 2.5-5b, start.

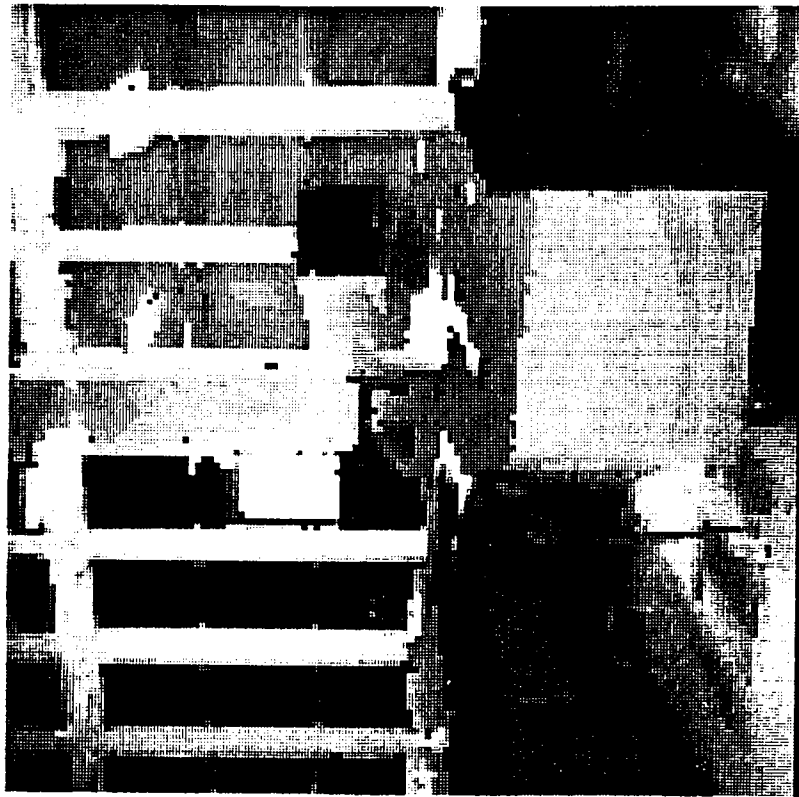


if8ganal.rgr

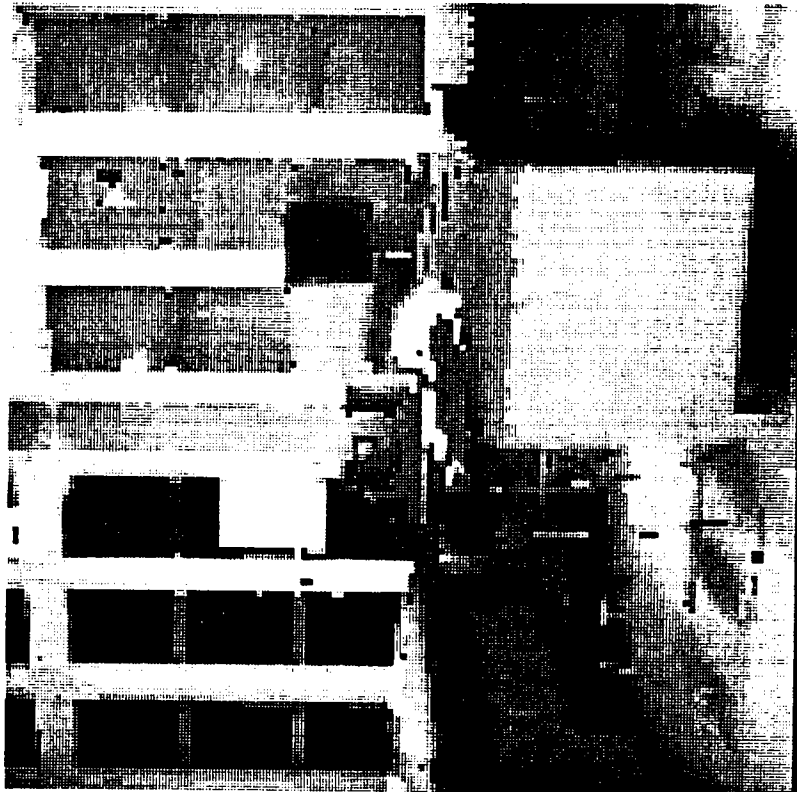


if8ganal.lgr

Fig.: 2.5-6.



if8gana2.rgr



if8qana2.lgr

Fig.: 2.5-7.

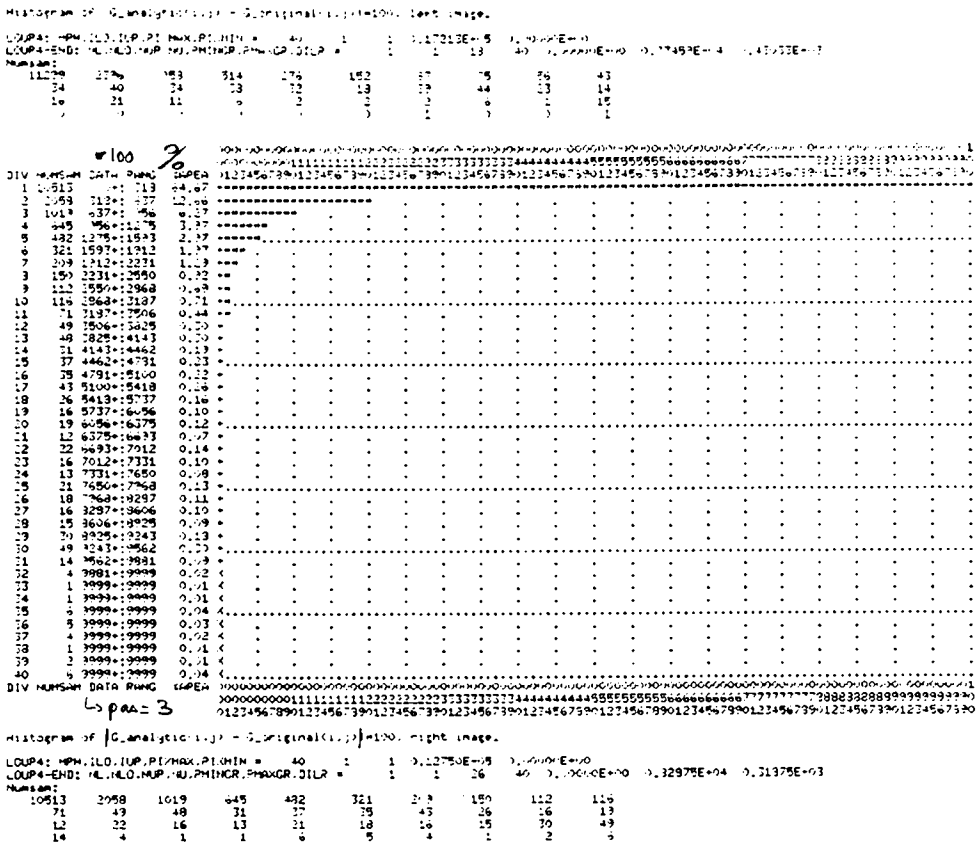
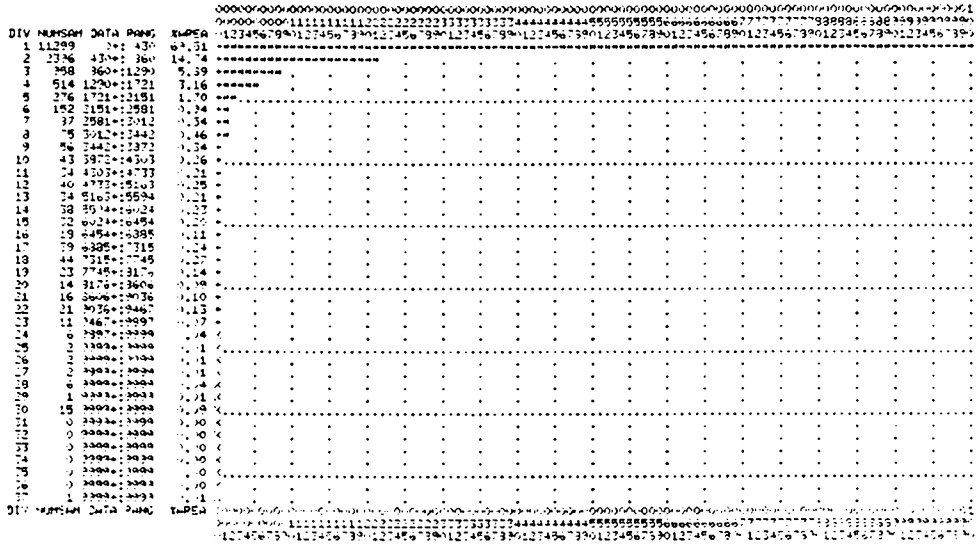


Fig.: 2.5-8.

$|g_{ana} - g_{orig}| * 100$

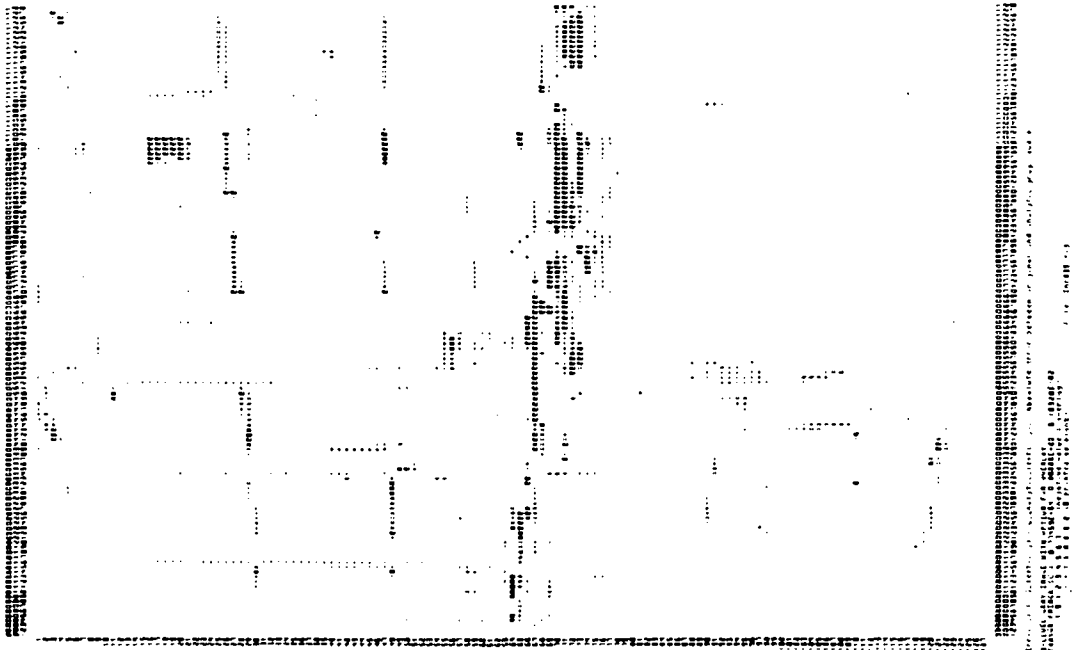
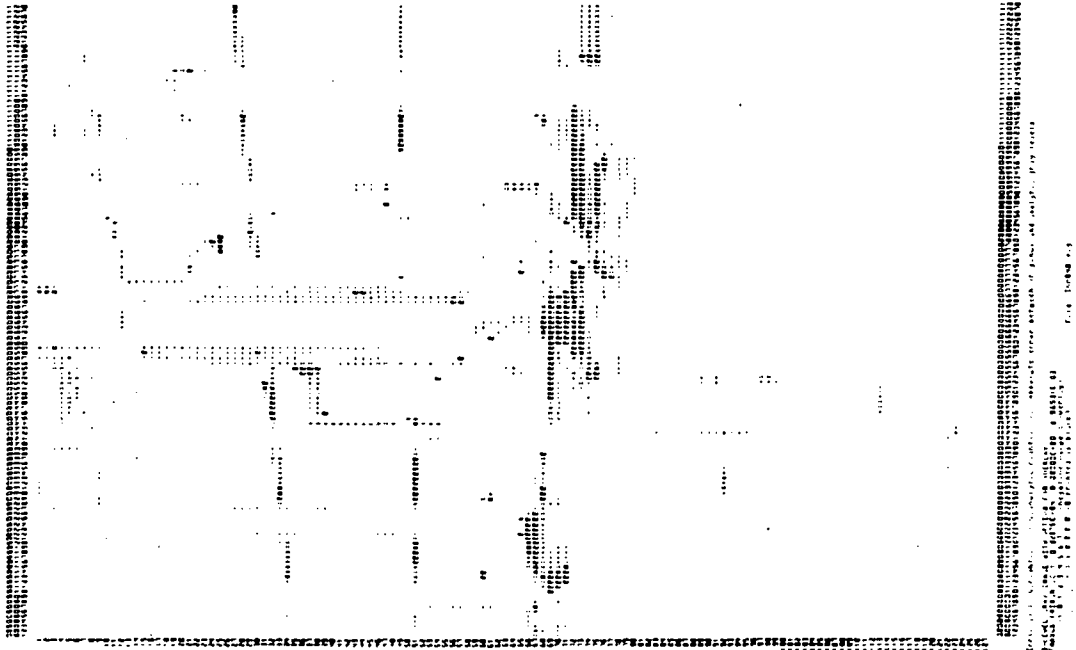


Fig.: 2.5-9.



if8ganal.lgr

if8ganal.rgr

Fig.: 2.5-10a.



if8gana2.lgr

if8gana2.rgr

Fig.: 2.5-10b.

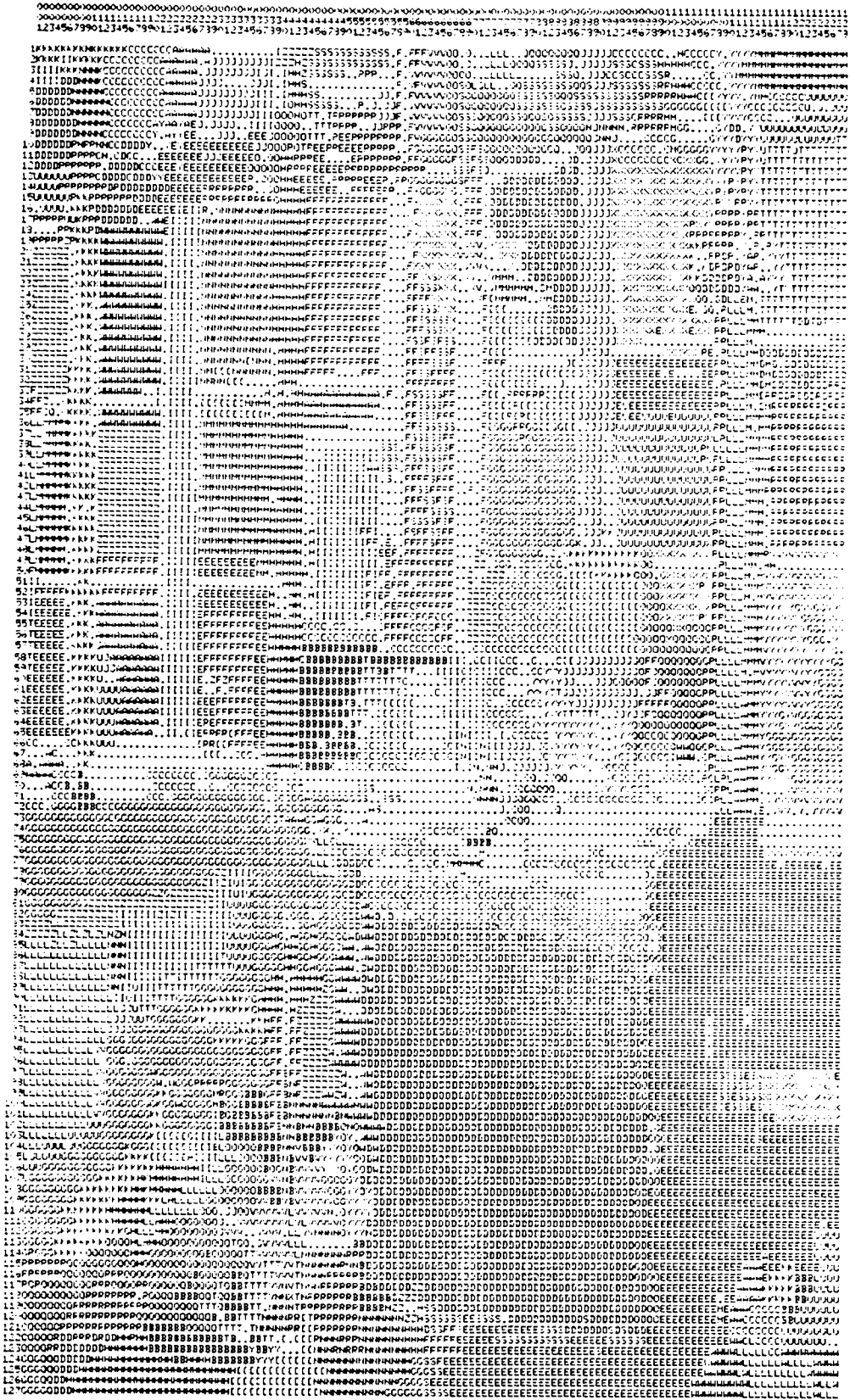
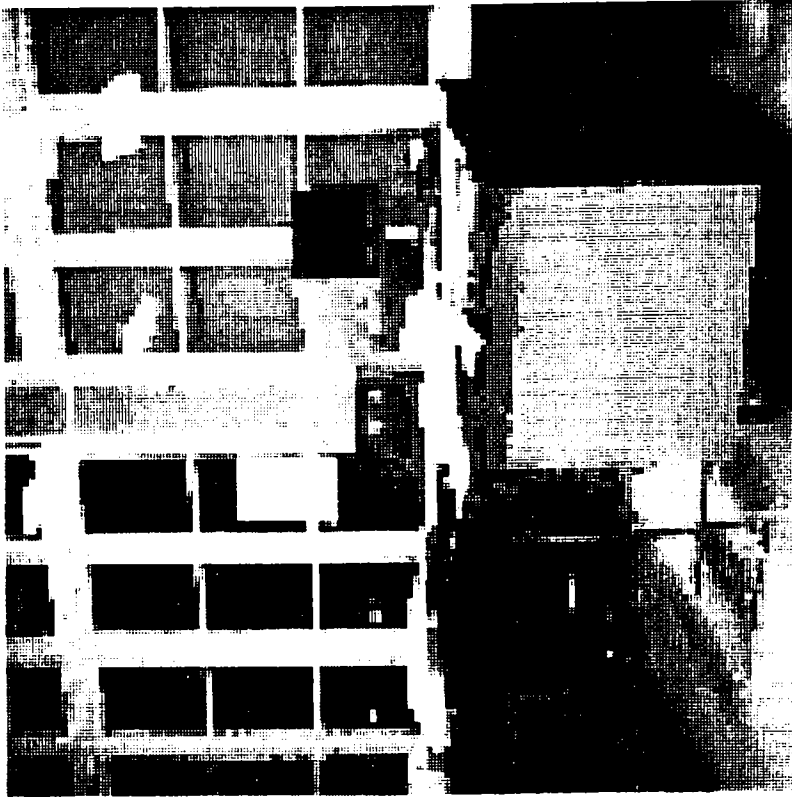
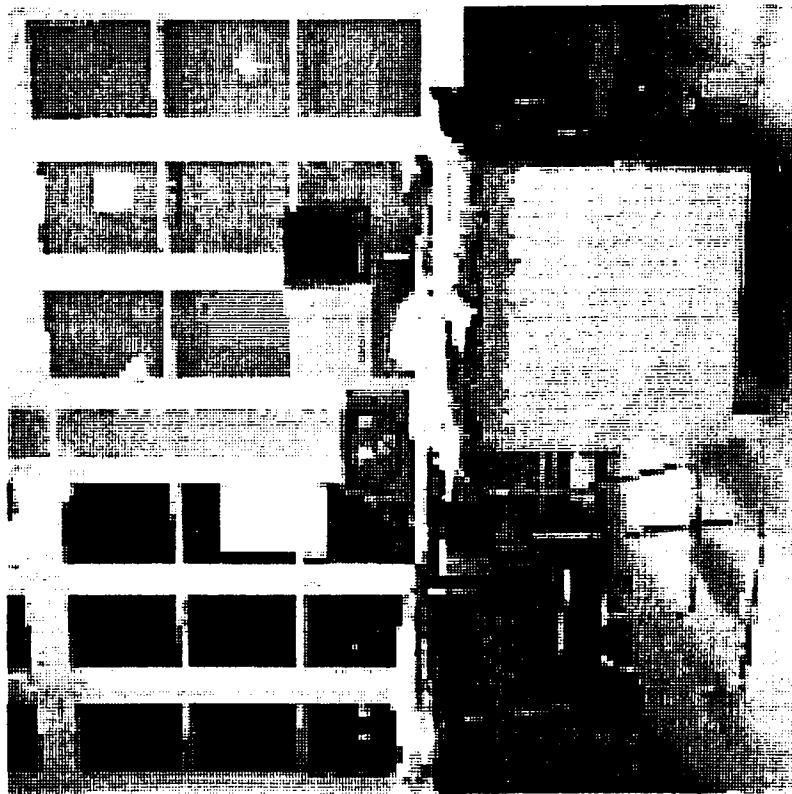


Fig.: 2.5-11b.

File: Inv92.tif



if8gana3.rgr



if8gana3.lgr

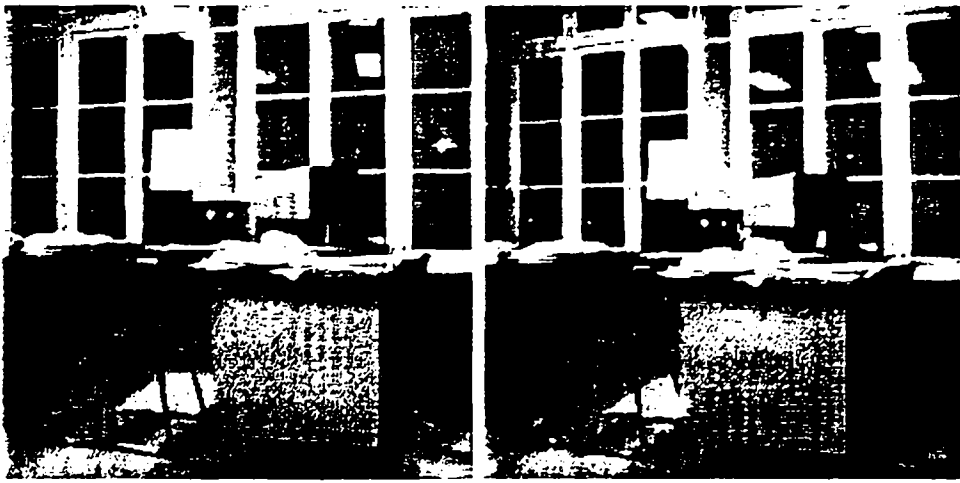
Fig.: 2.6-12.



if8gana3.lgr

if8gana3.rgr

Fig.: 2.5-13a.



i001.log

i001.rog

Fig.: 2.5-13b.

2.6 Super-regions versus scale-space

The initial segmentation described so far is dependent on the contents of the original image (image data), on the processes used to extract additional information (gray level gradients in the present case), on the size of the local operator (N), and on the decision spaces and the clustering procedures used (and to some extent on the "post processing" and on "analytic relaxation", if used). Of course, all the variables are kept constant when the left and right image are processed, but even then the segments produced are only "relatively similar" in the left and the right images of the stereo pair but the similarity is by no means "perfect". The differences occur since in some parts of the images the gray level structure is very different between the left and right image due to illumination changes, reflections, transparencies, or what is visible to one "eye" is not visible to the other due to the geometry of the scene. Thus, a certain amount of "miss-match" between the left and right segmented images is to be tolerated.

At least according to the currently popular opinions, if the size (N) of the local operators is increased then the results are expected to represent increasingly larger regions in the images (the "scale-space" effect). If N is "large enough" then the regions will be "large enough" and matching of the corresponding regions in the left and right images should become increasingly simpler. There is, of course, some truth in this argument, but we should remember that the scale-space approach only amounts to "smearing the information" in the images and is not likely to be the best method.

A more elegant method consists of not destroying the spatial resolution by increasingly larger local operators (N) but rather continuing the hierarchical approach based on the facets already found. There are numerous possibilities which, briefly summarized, are:

- a) Continuing the clustering of non-adjacent similar facets. This actually has to be done if there are structures in the image composed of similar facets since such facets cannot be matched correctly between the left and right images without some additional distinguishing information. A similar problem was once studied in connection with "linear textures" (2.5).
- b) Using adjacent facets to create "new recognition features" since two adjacent facets A and B in one image are also very likely to be adjacent in the other image.

In both cases "super-facets" can be created without destroying the spatial resolution. Time has not permitted experiments with these methods up to the present.

2.7 Conclusions

In this chapter some of the preliminary image processing and analysis steps for stereo vision have been experimented with in order to gain some practical experience with this problem. The image pair used is rather general and, probably, contains a fair example of the problems one may encounter with most gray level images. Exceptions, such as, images consisting only of textures, images of thin line structures, semitransparent scenes, etc., have not been investigated.

As briefly indicated, a very large number of pixel-level features are computable from the stereo pair of images, i.e., it is possible to extract much information from the images upon which to base segmentation. Of course, the same operations, in the same sequences, and with the same parameters, should be used when the two images are processed for pixel features. Hence, it is largely a question of willingness on our part to program the necessary procedures and to produce the computing power (special hardware) to make it practical.

If the original images differ, as is the case with the present pair of images, some preprocessing may be required to "equalize" the images (but this does not mean "histogram equalization"). However, if the information is lost due to saturation, then it is beyond recovery. In the present images this created some interesting problems in matching, see next Chapter. In regions where the pixel level features can be computed, clearly, for stereo vision to be possible, the images of these features must preserve "stereo fidelity", when we look at such a pair of images through a stereo viewer we have to be able to see "depth". Since our stereo vision ability is very powerful, if it exists at all, if an impression of depth cannot be seen then, clearly, there is "something seriously wrong" with the features and it is unlikely that further computing on such features has much meaning for stereo. The features used in the present study could all be "seen in stereo".

It should also be remembered that the segments in themselves do not need to "mean anything" in human terms, even though it is desirable if they do and they often do have a meaning if understood in terms of what the mathematical procedure is doing. In the stereo problem the only important criterion is "similarity" of the segments between the left and right images, and nothing else! Any level of merging is feasible if an hierarchical approach is used. For example, the experiments with the "automatic volume control" filter indicated that much detail can be "brought out" for continued matching on a finer scale. Visual inspection indicated that some of the "texture" brought out preserved texture integrity while, basically, these surfaces had no "intentional texture". In an active ("live") vision system such classification is carried out for each "glance" at the scene.

In the present case, the operations so far carried out and the sequences in which they have been carried out are summarized in the block diagram shown in Figure 2.7-1. The places where the stereo pairs have been created for visual verification of the results are shown in Figure 2.7-2. The abbreviations are explained in Figure 2.7-3.

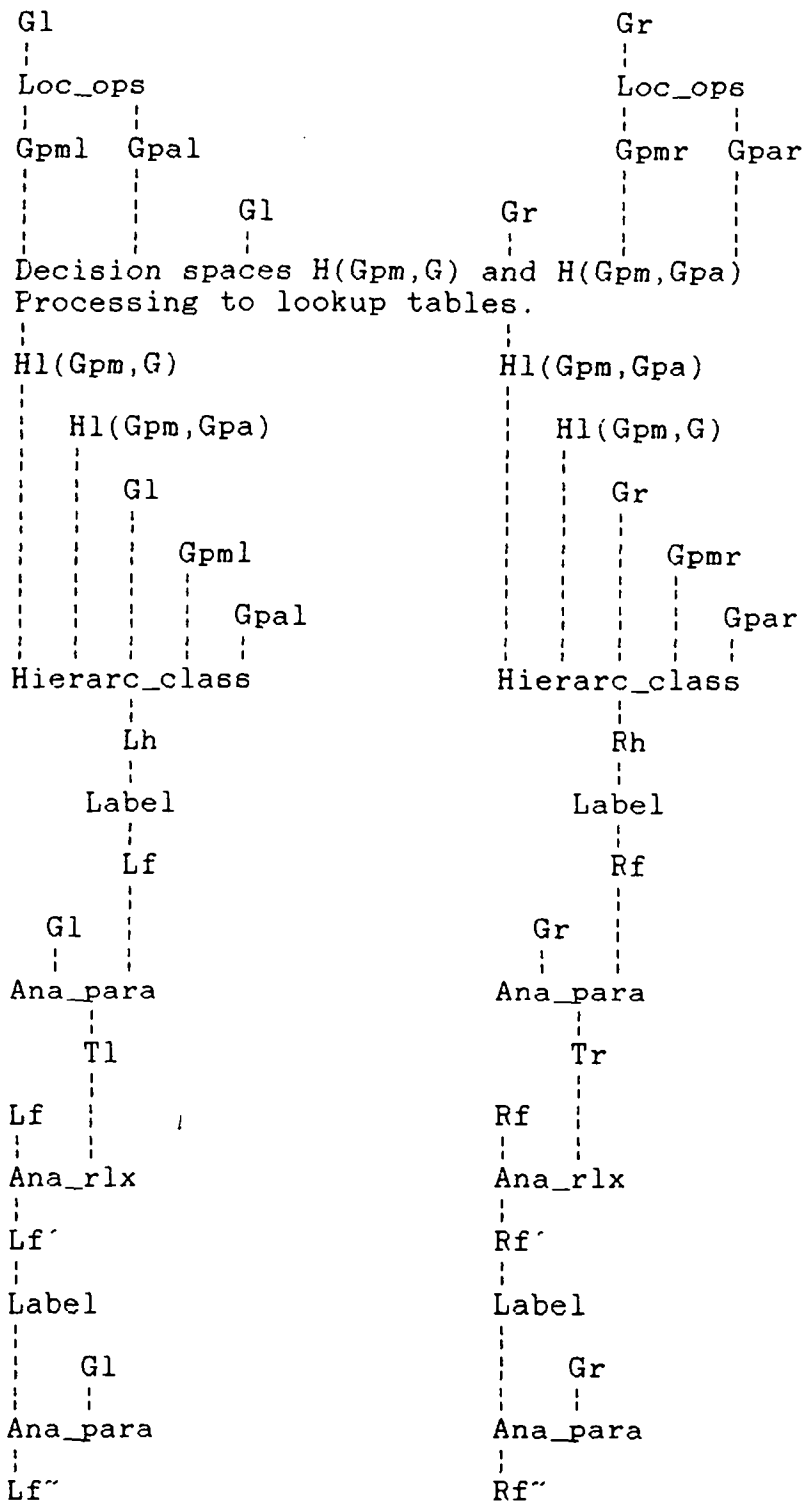


Figure 2.7-1 The basic analysis process. The abbreviations are listed in Fig. 2.7-3.

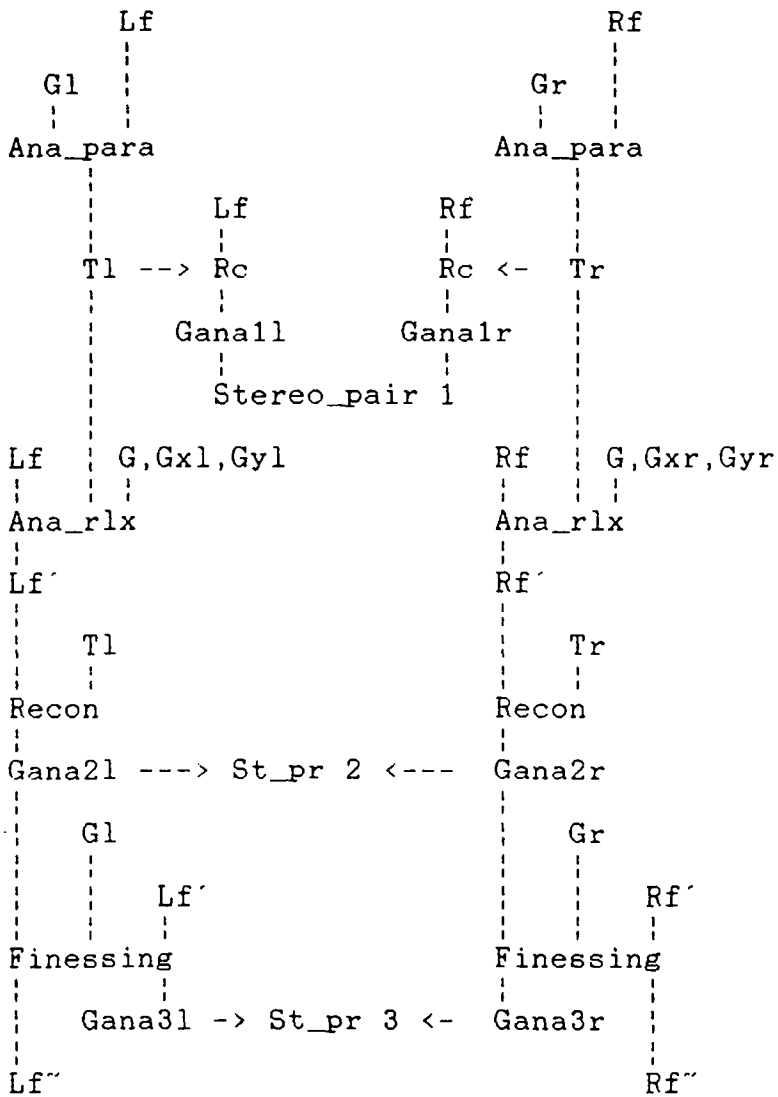


Figure 2.7-2: Stereo pair construction for visual verification.
 The abbreviations are explained in Fig. 2.7-3.

*l = Associated with left image.
 *r = Associated with right image.
 G = Gray level image.
 Gl = Gray level image of stereo pair, left.
 Gr = Gray level image of stereo pair, right.
 Loc_ops = Local operators within neighborhood N.
 Gpm = Gradient magnitude.
 Gpml = Gradient magnitude, left.
 Gpmr = Gradient magnitude, right.
 Gpa = Gradient angle.
 Gpal = Gradient angle, left.
 Gpar = Gradient angle, right.
 H(Gpm,G) = Decision space on Gpm and G for left and right.
 H(Gpm,Gpa) = Decision space on Gpm and Gpa for left and right.
 Hl(Gpm,G) = Lookup table for classifying Gpm and G.
 Hl(Gpm,Gpa) = Lookup table for classifying Gpm and Gpa.
 Hierarc_class = Hierarchical classification.
 *h = Class labelled image.
 Lh = Class labelled image, left.
 Rh = Class labelled image, right.
 Label = Region labelling, 4-connected.
 *f = Facet (region) labelled image.
 Lf = Facet (region) labelled image, left.
 Rf = Facet (region) labelled image, right.
 Gxl,Gyl = Gray level gradients, left.
 Gxr,Gyr = Gray level gradients, right.
 Ana_para = Analytic approximation parameters for facets.
 Tl = Analytic parameter table, left.
 Tr = Analytic parameter table, right.
 Ana_rlx = Analytic relaxation.
 Lf' = Facet (region) labelled image after relaxation, left.
 Rf' = Facet (region) labelled image after relaxation, right.
 Rc = Reconstruction of image from analytic parameters.
 Recon = Reconstruction of image from analytic parameters.
 Gana = Analytically reconstructed image.
 Gana1l = Analytically reconstructed image, 1'st, left.
 Gana1r = Analytically reconstructed image, 1'st, right.
 Stereo_pair = A stereo pair of images shown in text.
 St_pr = A stereo pair of images shown in text.
 Gana2l = Analytically reconstructed image, 2'nd, left.
 Gana2r = Analytically reconstructed image, 2'nd, right.
 Finessing = Replace large errors in Gana by G values.
 Gana3l = Analytically reconstructed image, finessed, left.
 Gana3r = Analytically reconstructed image, finessed, right.
 Lf'' = Facet (region) labelled image, finessed, left.
 Rf'' = Facet (region) labelled image, finessed, right.

Figure 2.7-3: Abbreviations used in Figures 2.7-1 and 2.7-2.

Chapter 3: FACET MATCHING

3.1 Introduction

Successful matching of the left (L) and right (R) image detail of a stereo pair is absolutely essential for "still-life" stereo-based depth vision. The match has to be carried out to pixel level, that is, exact (unique and exclusive) correspondence has to be established for all the pixels in the L and R images which represent regions in the scene visible to both "eyes". Facet matching cannot accomplish this alone, being only a step in an essentially recursive or iterative procedure (from larger facets to increasingly smaller until the facet becomes the size of the pixel and sub-pixel). A vision system that only produces a partial match will leave "unknown" regions in the scene.

For a "live" vision system the fact that a part of the scene is "unknown" or that a facet does not "match", is in itself information (a "centre of interest") to guide the "eyes" to such a location for "taking another look", until the whole 3D scene has been properly reconstructed from multiple "glances" at the scene. Even if the "scene" is only a stereo pair of images viewed through a set of stereo glasses, the "live" system can still "keep searching" until a complete match has been obtained. It is thus conjectured that a "live" vision system that can "move its eyes" (centre of fixation) to any desired place in the scene has a far easier problem to solve than the "still-life" stereo system.

In a "still-life" stereo system only two images are available and the "eyes" as such do not exist. We know from our own experience that such images can be seen "in stereo". However, we can still move our eyes during viewing. The author is not aware of any experiments to verify whether a stereo pair of "fixated images" can be fused by our visual system, but the answer is likely to be affirmative. This question has only been raised to alert the reader to the unlikely possibility that "still-life" stereo vision may not be a very realistic problem, except in photogrammetry.

However, be as it may, this chapter describes several experiments in attempting to match the facets in a left and right "still-life" stereo pair of gray level images. As was demonstrated in the previous chapter, there is no lack of information (individually, i.e., separately) computable from the two images upon which to base the matching experiments. As a brief summary, the information consists of:

- a) Pixel level features belonging to each pixel (gray level, colour, and flicker). In the present case only the gray level is available.
- b) Pixel level features obtained from local operators (over a neighborhood N). There is a profusion of such features.
- c) Homogeneous regions obtained by variously grouping the pixel features from (a) and (b). There is a profusion of such groupings also. The basic criterion is that the facets (regions) in the images become "easier to match" than the original pixels (for example, they have to be "reasonably large" connected regions).
- d) Super-regions of various types obtained by grouping the regions from (c) and/or by using operators that process the whole image as a unit.

Hierarchical matching strategies were suggested. The initial match may only be on the "super-region" level (d) but could also be on any level of detail if the features are unique. Presumably the super-regions are easiest to match since they practically always overlap in the two images (super-region intersection between L and R images is highly likely to be non-zero). This will assist the matching of the "homogeneous regions" (c) especially if they are subsets of super-regions (d). After this matching is accomplished, the local operator may be made more "sensitive" by decreasing the size of local region N (the "scale-space" approach) or the image is enhanced to bring out details within each homogeneous region, or the region boundaries are matched next, etc. In the end the regions should become "small enough" to allow unique pixel level feature matching.

Only the "similarity principle" is used in the matching, i.e., at whatever level the matching is to occur, the items matched are (rather) similar in the L and R images but must have unique characteristics to prevent confusion (combinatorial problems should not exist). The L and R similarity principle extends from pixel features to super-regions. Consequently, object recognition, structure recognition, scene models, etc., are unnecessary and "image understanding" is not needed to "see depth" in stereo images. By definition, the images to be matched must be unknown initially and remain so even after stereo matching.

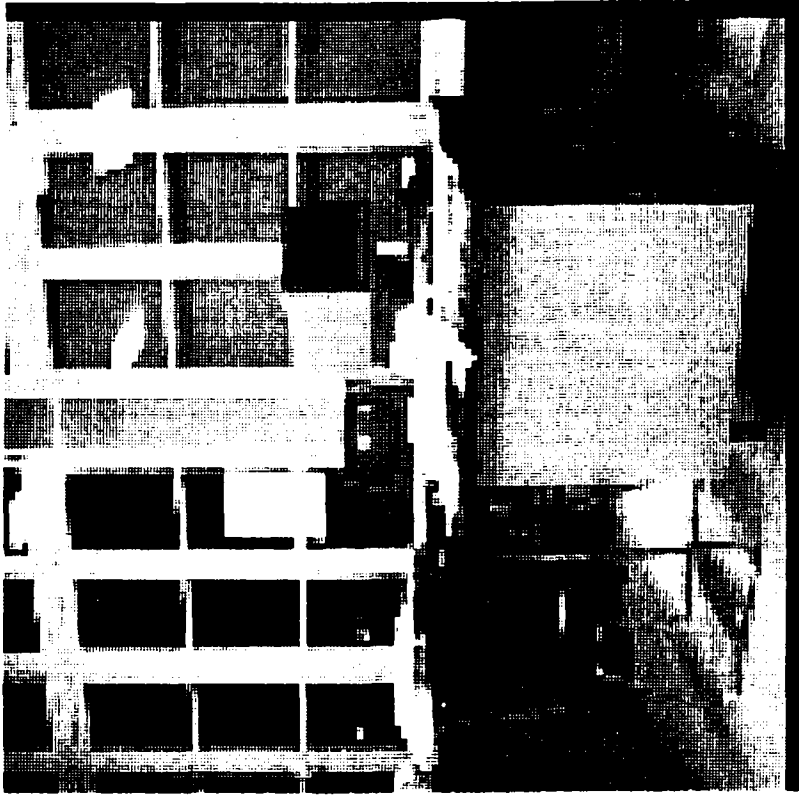
Since no technical information about the two "INRIA images" was available, the right image (R) was shifted slightly with respect to the left image (L) in an attempt to create an apparent "centre of fixation". The R image was shifted by 3 pixels to the left and 2 pixels up with respect to L, see Figure 3.1-1. The corresponding region in L was zeroed. This did not disturb stereo fidelity, as may be ascertained from Figures 3.1-2a and 2b. All feature images needed in the processing were shifted by the same amount (after the processing described in Chapter 2).

The experiments are far from complete both due to time and equipment limitations. About half a dozen matching experiments are described and even these experiments have been restricted to only the "homogeneous region" level (c). In retrospect, some of the experiments appear to have been rather poorly thought out before the experiment was started but, one learns by doing!

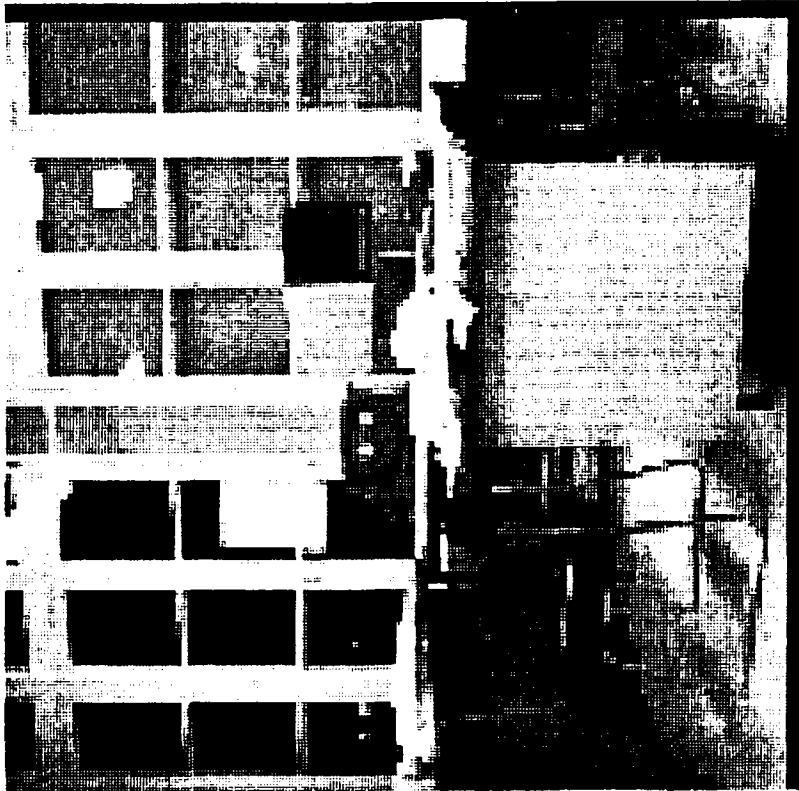
Figure titles:

Figure 3.1-1: The images after shifting in an attempt to create a "centre of fixation".

Figures 3.1-2a and -2b: Stereo pairs of the shifted gray level and the shifted gradient magnitude images.



if8gray.rsh



if8gray.lsh

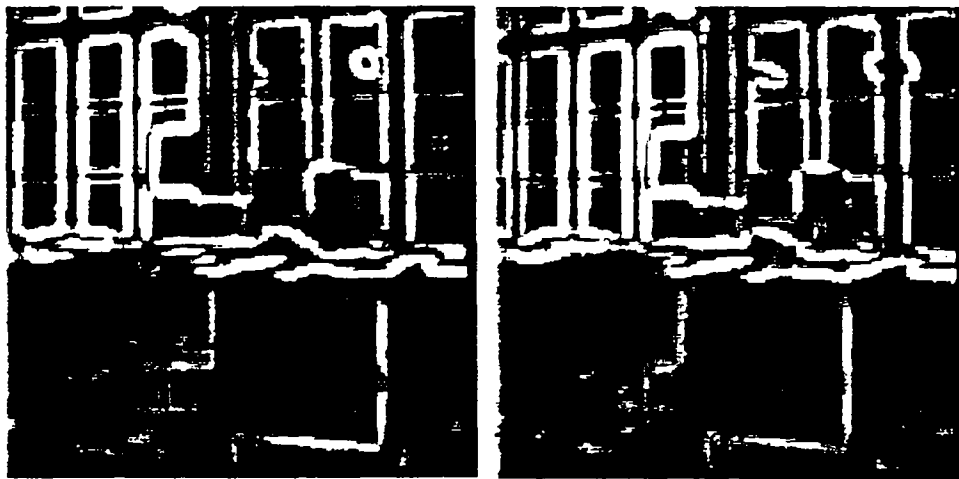
Fig.: 3.1-1.



if8gray.lsh

if9gray.rsh

Fig.: 3.1-2a.



if8gradm.lsh

if9gradm.rsh

Fig.: 3.1-2b.

3.2 Matching methods

Several matching methods were tried. These are briefly described in the sub-sections to follow. Briefly:

1. The first method attempts to use "matched classification" in order to try to classify the pixels in the left and right image at the same time. In the present context the results from this case are used to illustrate the many difficulties one can encounter.
2. The second method is based on simple "and-ing" of class labels which are then post-processed. This method has been somewhat more "refined" to study how some of the problems can be partially "fixed". Except for a few results shown in the Abstract and in section 3.2.2, these results have been left out of the report to reduce its size.
3. The third approach is based on correlation. It is only a slight modification of the techniques used in binary image processing. (Discussed but not completed due to lack of time.)
4. The fourth method is using "mutual recognition" of the facets in the left and right images. At present it is only a slightly modified "minimum distance" classifier.
5. The fifth approach attempted to use pixel feature matching to go directly to a depth image, but became constrained to a very restricted case due to equipment limitations. (Briefly discussed.)

In the first and second methods (1,2) the initial matched image is called $L\&R(i,j)$ and it is accompanied by the corresponding left $L(i,j)$ and right $R(i,j)$ images. The third method (3) gives a list of matching facet labels and also a list of "best shifts" between the facets in L and R . The fourth method (4) only gives a list of matching facet labels in L and R but no shifts. The images of matched facets, $L_m(i,j)$ and $R_m(i,j)$, are constructed from these lists. In all cases, the $L\&R_m(i,j)$ and/or $L_m(i,j)$ and $R_m(i,j)$ images require considerable further processing before the results can be judged. The fifth method attempts to go directly from pixel features to 3D reconstruction.

3.2.1 Matched classification

The two hierarchical labelled decision spaces (look-up tables) used in "homogeneous" region detection ($H1(G_{pm}, G)$ followed by $H1(G_{pm}, G_{pa})$, see Chapter 2) contained about a dozen separate region types but, of course, the actual number of region types is not important. As seen from the results in Chapter 2, independent classification of the L and R images produced "reasonably consistent" regions, except in some areas of the images. The "totally different" regions in L versus R were mainly due to missing gray levels and reflections of light sources.

It was argued that if a pixel p at coordinates (i, j) in the left image L belonged to a given class k , then there should be a pixel p' in the right image R within "an allowable maximal displacement" d_i and d_j from (i, j) which should belong to the same class k if the gray level (G) and gradient magnitude and angle values (G_{pm} and G_{pa}) were "perturbed" by increasing amounts after each search in $d_i * d_j$. This should improve the similarity of the "reasonably consistent" regions but, of course, should not improve the "totally different" regions. Thus, it should be possible to construct a "matched and classified" image $L\&R_h(i, j)$ and "displacement images" $d_i(i, j)$ and $d_j(i, j)$ containing d_i and d_j , where

$$\begin{aligned} i(\text{left image}) \pm d_i &= i(\text{right image}) \\ j(\text{left image}) \pm d_j &= j(\text{right image}). \end{aligned}$$

A block diagram of the processing structure is shown in Figure 3.2.1-1a with explanation of the symbols in Figure 3.2.1-1b.

The "displacement images" $d_i(i, j)$ and $d_j(i, j)$ could not be used since the features were essentially homogeneous within each facet and no pixel matching could take place. The number of features was too limited for unique identification and matching of pixels. The "matched and classified" image $L\&R_h(i, j)$ is shown in Fig. 3.2.1-2. The corresponding $L_h(i, j)$ and $R_h(i, j)$ images are in Figures 3.2.1-3a and -3b. In the present case the "raw" results are shown in order to better illustrate the reasons for many of the processes described in section 3.3, i.e., "post processing" was not used to "homogenize" the $L\&R_h(i, j)$, $L_h(i, j)$, and $R_h(i, j)$ triplet.

As in Chapter 2, attempts were made to display the results "in stereo" for visual verification. When each class was replaced by its average gray level (see Figures 3.2.1-4a and 3.2.1-5) the images became rather "bland" since the gray level gradients are missing (except at class boundaries). Nine "colours" were required to "paint" these regions in different "colours" (see Figures 3.2.1-4b, -4c and -6b). A stereo pair for the average gray levels is shown in Figure 3.2.1-6a and for the "painted" classes in Figure 3.2.1-6b. Mild stereo effects are visible, but the nonmatching "dots" in Figure 3.2.1-6b greatly disturb our vision.

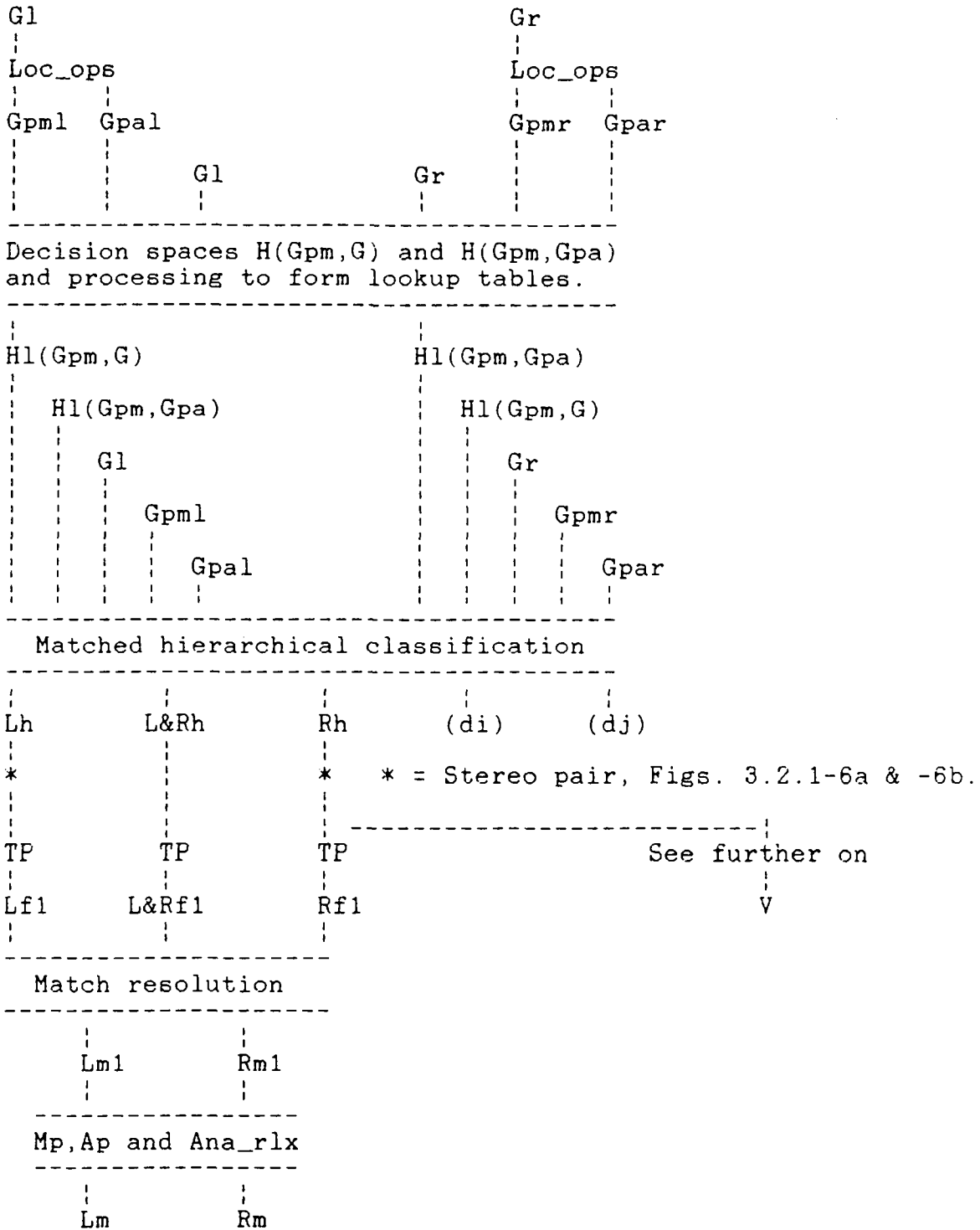


Figure 3.2.1-1a: Matched classification and match resolution. Process shown starting from stereo image pair $G_l(i,j)$ and $G_r(i,j)$. Explanation of symbols is in Fig. 3.2.1-1b.

*l = Associated with left image.
 *r = Associated with right image.
 G = Gray level image.
 Gl = Gray level image of stereo pair, left.
 Gr = Gray level image of stereo pair, right.
 Loc_ops = Local operators within neighborhood N.
 Gpm = Gradient magnitude.
 Gpml = Gradient magnitude, left.
 Gpmr = Gradient magnitude, right.
 Gpa = Gradient angle.
 Gpal = Gradient angle, left.
 Gpar = Gradient angle, right.
 H(Gpm,G) = Decision space on Gpm and G for left and right.
 H(Gpm,Gpa) = Decision space on Gpm and Gpa for left and right.
 Hl(Gpm,G) = Lookup table for classifying Gpm and G.
 Hl(Gpm,Gpa) = Lookup table for classifying Gpm and Gpa.
 *h = Class labelled image.
 Lh = Class labelled image, left.
 Rh = Class labelled image, right.
 L&Rm = Matched class label image.
 di = Displacement image for i-coordinate.
 dj = Displacement image for j-coordinate.
 TP = Lh, L&Rh, and Rh triplet processing
 *f = Facet (region) labelled image.
 Lf1 = Facet (region) labelled image, left.
 Rf1 = Facet (region) labelled image, right.
 L&Rf1 = Facet (region) labelled image of L&Rm.
 Lm1 = Matched image, left.
 Rm1 = Matched image, right.
 Mp,Ap and Ana_rlx = Subsequent processes.
 Mp = Match processing.
 Ap = Analytic approximation.
 Ana_rlx = Analytic relaxation.
 Lm = Matched image, left.
 Rm = Matched image, right.

Figure 3.2.1-1b: Symbols for matched classification, see Fig. 3.2.1-1a.

(Relevant computer files: /user2/kasvand/ima2):

Lh = Inr065.fig	L&Rh = Inr063.fig	Rh = Inr066.fig
Lf1 = Inr069.fig	L&Rf1 = Inr067.fig	Rf1 = Inr070.fig
Match list = Inr071.fig		
Lm1 = Inr073.fig	Rm1(i,j) = Inr074.fig	

The principal difficulties with matched classification are that if nothing is known about the scene then the "allowable maximal displacement" for d_i (along the direction of displacement of the cameras) can be the whole image width, and if the "perturbations" are allowed to be too large then any pixel in L will match any pixel in R and vice versa. However, when the "perturbations" and the " d_i and d_j " tolerances were made small, then the method started to resemble direct "and-ing" of the $L_h(i,j)$ and $R_h(i,j)$ images. Chronologically, this was the second method tried, being a "watered down" version of direct "pixel feature matching", which was beyond the capabilities of the equipment.

Figure titles:

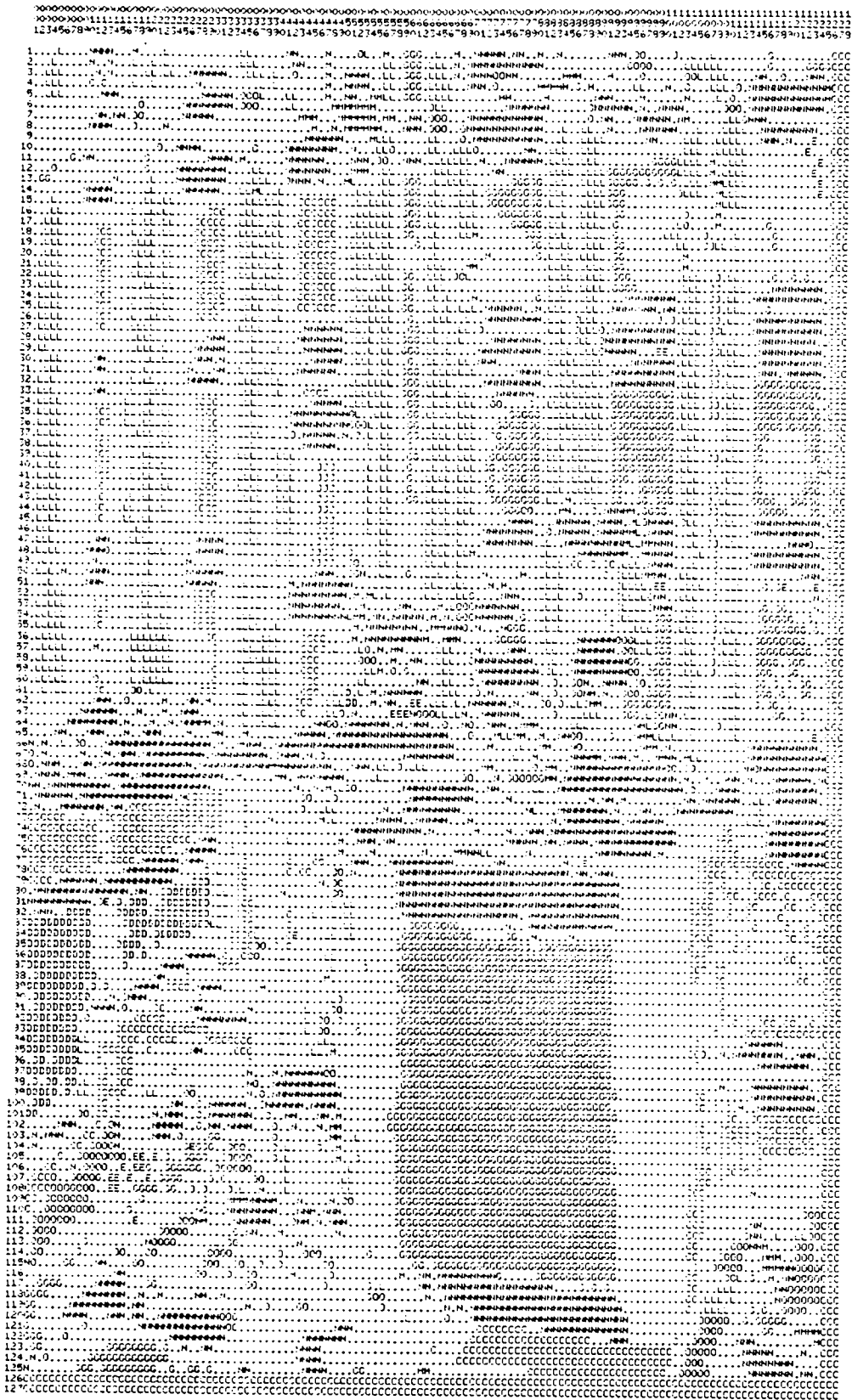
Figure 3.2.1-2: The "matched and classified" image $L \& R_h(i,j)$ showing the class labels from $H_l(Gpm,G)$ and $H_l(Gpm,Gpa)$. The labelling is continued from $H_l(Gpm,G)$ to $H_l(Gpm,Gpa)$. (File: Inr063.fig)

Figures 3.2.1-3a and -3b: The $L_h(i,j)$ and $R_h(i,j)$ class label images. No "post processing" has been used. (Files: Inr065.fig and Inr066.fig)

Figures 3.2.1-4a, -4b and -4c: Average gray levels per class, adjacency matrix for classes, and "colour" assignments for the classes. (Files: Inr095.fig, Inr096.fig, and Inr097.fig)

Figure 3.2.1-5: Each class was replaced by its average gray level.

Figures 3.2.1-6a and -6b: Stereo pairs of the $L_h(i,j)$ and $R_h(i,j)$ images. (a) Using average gray level for each class. (b) Classes "painted" in different "colours" where adjacent classes have different colours.



Labels from matched classification. Top1.2.3 = 1.0.0.

LABEL PRINTING
 * 1 2 3 4 5 6 7
 1 1 2 3 4 5 6 7
 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25
 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45
 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65
 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85
 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100 101 102 103 104 105
 106 107 108 109 110 111 112 113 114 115 116 117 118 119 120 121 122 123 124 125
 126 127 128 129 130 131 132 133 134 135 136 137 138 139 140 141 142 143 144 145 146 147 148 149 150
 151 152 153 154 155 156 157 158 159 160 161 162 163 164 165 166 167 168 169 170 171 172 173 174 175 176 177 178 179 180
 181 182 183 184 185 186 187 188 189 190 191 192 193 194 195 196 197 198 199 200 201 202 203 204 205 206 207 208 209 210
 211 212 213 214 215 216 217 218 219 220 221 222 223 224 225 226 227 228 229 230 231 232 233 234 235 236 237 238 239 240 241 242 243 244 245 246 247 248 249 250
 251 252 253 254 255 256 257 258 259 260 261 262 263 264 265 266 267 268 269 270 271 272 273 274 275 276 277 278 279 280 281 282 283 284 285 286 287 288 289 290
 291 292 293 294 295 296 297 298 299 300 301 302 303 304 305 306 307 308 309 310 311 312 313 314 315 316 317 318 319 320 321 322 323 324 325 326 327 328 329 330
 331 332 333 334 335 336 337 338 339 340 341 342 343 344 345 346 347 348 349 350 351 352 353 354 355 356 357 358 359 360 361 362 363 364 365 366 367 368 369 370
 371 372 373 374 375 376 377 378 379 380 381 382 383 384 385 386 387 388 389 390 391 392 393 394 395 396 397 398 399 400 401 402 403 404 405 406 407 408 409 410
 411 412 413 414 415 416 417 418 419 420 421 422 423 424 425 426 427 428 429 430 431 432 433 434 435 436 437 438 439 440 441 442 443 444 445 446 447 448 449 450
 451 452 453 454 455 456 457 458 459 460 461 462 463 464 465 466 467 468 469 470 471 472 473 474 475 476 477 478 479 480 481 482 483 484 485 486 487 488 489 490
 491 492 493 494 495 496 497 498 499 500 501 502 503 504 505 506 507 508 509 510 511 512 513 514 515 516 517 518 519 520 521 522 523 524 525 526 527 528 529 530
 531 532 533 534 535 536 537 538 539 540 541 542 543 544 545 546 547 548 549 550 551 552 553 554 555 556 557 558 559 560 561 562 563 564 565 566 567 568 569 570
 571 572 573 574 575 576 577 578 579 580 581 582 583 584 585 586 587 588 589 590 591 592 593 594 595 596 597 598 599 600 601 602 603 604 605 606 607 608 609 610
 611 612 613 614 615 616 617 618 619 620 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 636 637 638 639 640 641 642 643 644 645 646 647 648 649 650
 651 652 653 654 655 656 657 658 659 660 661 662 663 664 665 666 667 668 669 670 671 672 673 674 675 676 677 678 679 680 681 682 683 684 685 686 687 688 689 690
 691 692 693 694 695 696 697 698 699 700 701 702 703 704 705 706 707 708 709 710 711 712 713 714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730
 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 754 755 756 757 758 759 760 761 762 763 764 765 766 767 768 769 770
 771 772 773 774 775 776 777 778 779 780 781 782 783 784 785 786 787 788 789 790 791 792 793 794 795 796 797 798 799 800 801 802 803 804 805 806 807 808 809 810
 811 812 813 814 815 816 817 818 819 820 821 822 823 824 825 826 827 828 829 830 831 832 833 834 835 836 837 838 839 840 841 842 843 844 845 846 847 848 849 850
 851 852 853 854 855 856 857 858 859 860 861 862 863 864 865 866 867 868 869 870 871 872 873 874 875 876 877 878 879 880 881 882 883 884 885 886 887 888 889 890
 891 892 893 894 895 896 897 898 899 900 901 902 903 904 905 906 907 908 909 910 911 912 913 914 915 916 917 918 919 920 921 922 923 924 925 926 927 928 929 930
 931 932 933 934 935 936 937 938 939 940 941 942 943 944 945 946 947 948 949 950 951 952 953 954 955 956 957 958 959 960 961 962 963 964 965 966 967 968 969 970
 971 972 973 974 975 976 977 978 979 980 981 982 983 984 985 986 987 988 989 990 991 992 993 994 995 996 997 998 999 1000
 File: Inp063.rtg

Fig.: 3.2.1-2.


```

mjbmmmbb0000000000000000000000
00000000001111111111111111111112
12345678901234567890

1 0111101110111100000000
2 10111100011110000000
3 11011100011110000000
4 11101100011110000000
5 01110000011110000000
6 11110001111100000000
7 11000101101110000000
8 10000110111110000000
9 10000001101110000000
10 11111110011110000000
11 11111110011110000000
12 11111110011110000000
13 11111110011110000000
14 11111110011110000000
15 0000000000000000000000
16 0000000000000000000000
17 0000000000000000000000
18 0000000000000000000000
19 0000000000000000000000
20 0000000000000000000000

ndjssncny mtrcc, for class labelled
left and right images.
Incr017.fig = Class labels, left.
Incr018.fig = Class labels, right.
File: Incr019.fig
  
```

Colour assignments. Nine colours were needed.

Lab ... = Class label in images.
Mean ... = Number of contacts with other labels.
Neicol ... = New colour assignment.

Incr017.fig = Class labels, left.
Incr018.fig = Class labels, right.

File: Incr019.fig

Fig.: 3.2.1-4c.

```

mjbmmmbb0000000000000000000000
00000000001111111111111111111112
12345678901234567890

1 0111101110111100000000
2 10111100011110000000
3 11011100011110000000
4 11101100011110000000
5 01110000011110000000
6 11110001111100000000
7 11000101101110000000
8 10000110111110000000
9 10000001101110000000
10 11111110011110000000
11 11111110011110000000
12 11111110011110000000
13 11111110011110000000
14 11111110011110000000
15 0000000000000000000000
16 0000000000000000000000
17 0000000000000000000000
18 0000000000000000000000
19 0000000000000000000000
20 0000000000000000000000

ndjssncny mtrcc, for class labelled
left and right images.
Incr017.fig = Class labels, left.
Incr018.fig = Class labels, right.
File: Incr019.fig
  
```

Colour assignments. Nine colours were needed.

Lab ... = Class label in images.
Mean ... = Number of contacts with other labels.
Neicol ... = New colour assignment.

Incr017.fig = Class labels, left.
Incr018.fig = Class labels, right.

File: Incr019.fig

Fig.: 3.2.1-4b.

```

K = L R C
MFL  MHF  55GL  55GP
1 = * C 1 485 504 50.5 50.8
2 = C C 1 2129 2953 1.5 0.4
3 = D D 1 509 343 15.1 13.5
4 = E E 1 240 400 26.7 23.2
5 = F F 1 41 12 22.3 23.7
6 = G G 1 2946 3548 55.8 54.2
7 = H H 1 113 37 80.5 81.2
8 = I I 1 41 28 87.1 88.5
9 = J J 1 315 277 105.0 107.3
10 = K K 0 0 5 0.0 96.4
11 = L L 1 3627 3472 55.4 54.0
12 = M M 1 310 738 49.6 47.8
13 = N N 1 3639 3603 43.1 40.6
14 = O O 1 1301 1206 44.6 41.3

SC_GMAX = 0.20001E+01  0.12750E+03
  
```

Average gray levels for classes after matched classification.

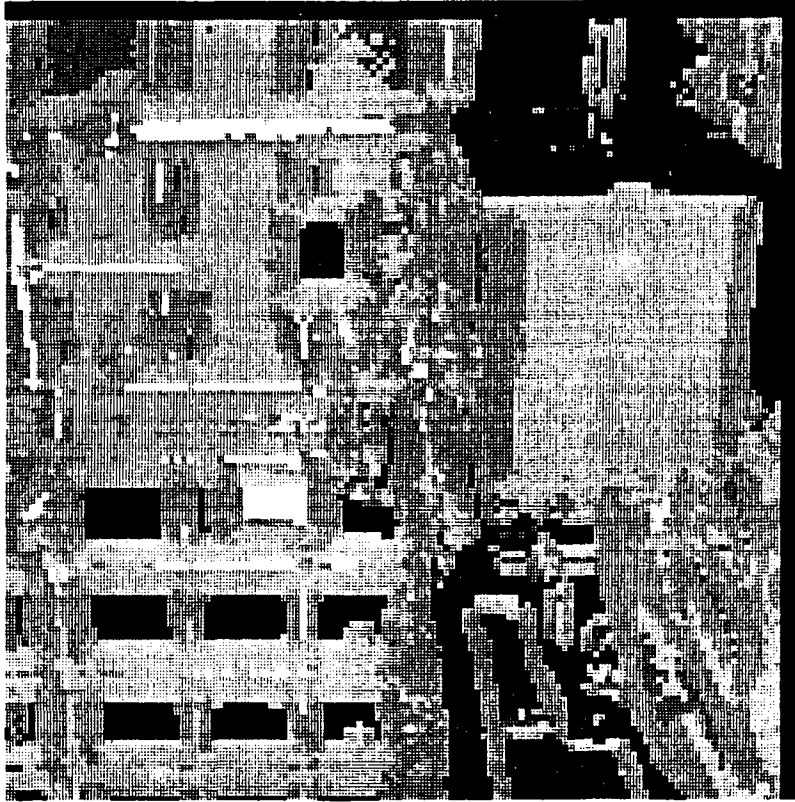
Incr017.fig = Class labels, left.
Incr018.fig = Class labels, right.

K = Class Label.
L = Class Label, left image.
R = Class Label, right image.
C = 1 if both labels in left and right images, else 0.

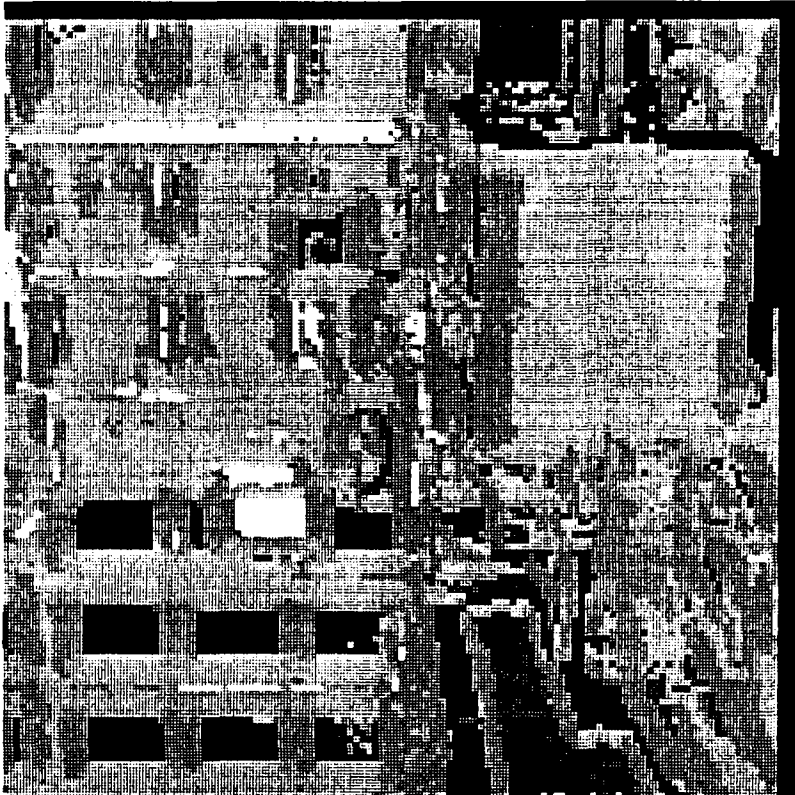
MFL = Number of pixels per class, left image.
MHF = Number of pixels per class, right image.
55GL = Average gray level of class, left image.
55GP = Average gray level of class, right image.

File: Incr019.fig

Fig.: 3.2.1-4a.

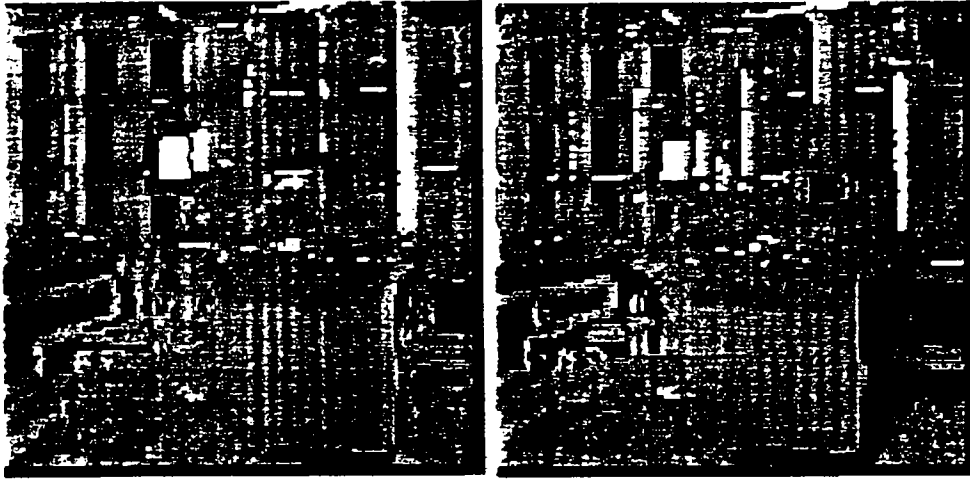


if8hlmcl.rgr



if8hlmcl.lgr

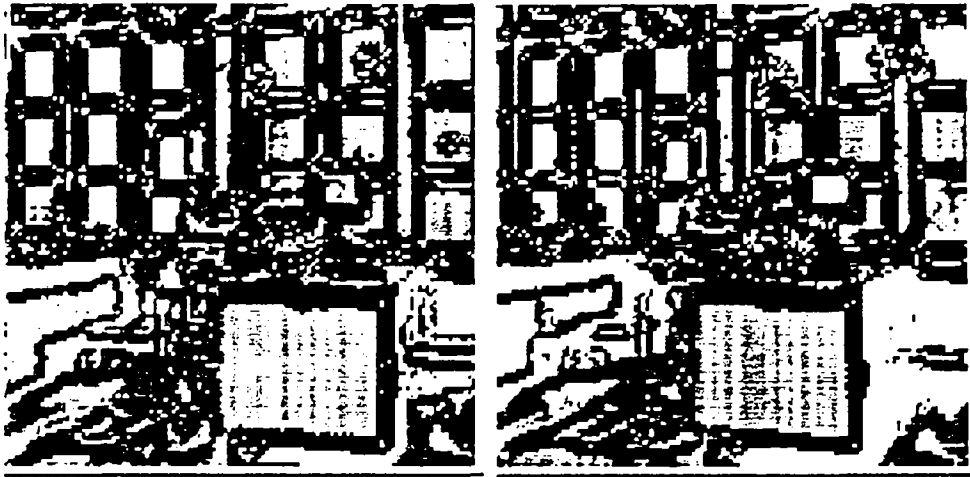
Fig.: 3.2.1-5.



if8hlmcl.lgr

if8hlmcl.rgr

Fig.: 3.2.1-6a.



if8hlmcl2.lgr

if8hlmcl2.rgr

Fig.: 3.2.1-6b.

3.2.2 Direct "and-ing"

The classified left $Lh(i,j)$ and right $Rh(i,j)$ images contain the class labels from $Hl(Gpm,G)$ and $Hl(Gpm,Gpa)$. Direct "and-ing" consists in comparing the labels in $Lh(i,j)$ and $Rh(i,j)$ pixel by pixel (for the same i and j in $Lh(i,j)$ and $Rh(i,j)$) and if the labels are the same then this label is stored in $L\&Rh(i,j)$ else a zero is stored. Thus,

```

For all pixels (i,j)
If( $Lh(i,j).Eq.Rh(i,j)$ ) then  $L\&Rh(i,j) = Lh(i,j)$ 
    else  $L\&Rh(i,j) = 0$ 

```

In block diagram form the process is shown in Figure 3.2.2-1a with explanation of symbols in Figure 3.2.2-1b. There is an obvious danger with this method, specially if the $Lh(i,j)$ and $Rh(i,j)$ images contain (horizontal) repetitive structures. As seen from Fig. 3.2.2-2 for an "one-dimensional" case, if the shift (disparity) between the L and R images is larger than the region size then the matching will be entirely incorrect. Clearly, "grouping" of repetitive structures is needed. Grouping is not a very difficult problem. Some discussion on an one-dimensional version of it may be found in Ref. 2.5.

The $Lh(i,j)$ and $Rh(i,j)$ images, shown in Chapter 2 as Figures 2.4-7a and -7b, are repeated as Figures 3.2.2-3a and -3b. The resultant $L\&Rh(i,j)$ image is shown in Fig. 3.2.2-3c. In this case the Lh and Rh images were subject to post-processing (as compared to "Matched classification") in order to illustrate the improvement (compare with Figs. 3.2.1-3a and -3b). As before, attempts are made to show the results in different ways for verification and evaluation. In Figure 3.2.2-4 the pixel classes are represented such that adjacent regions have different gray levels (seven "colours" were used). Figures 3.2.2-5a and -5b show the same "colours" printed as symbols. Figures 3.2.2-6a and -6b show the stereo pair for Figure 3.2.2-4 and for the original gray level images.

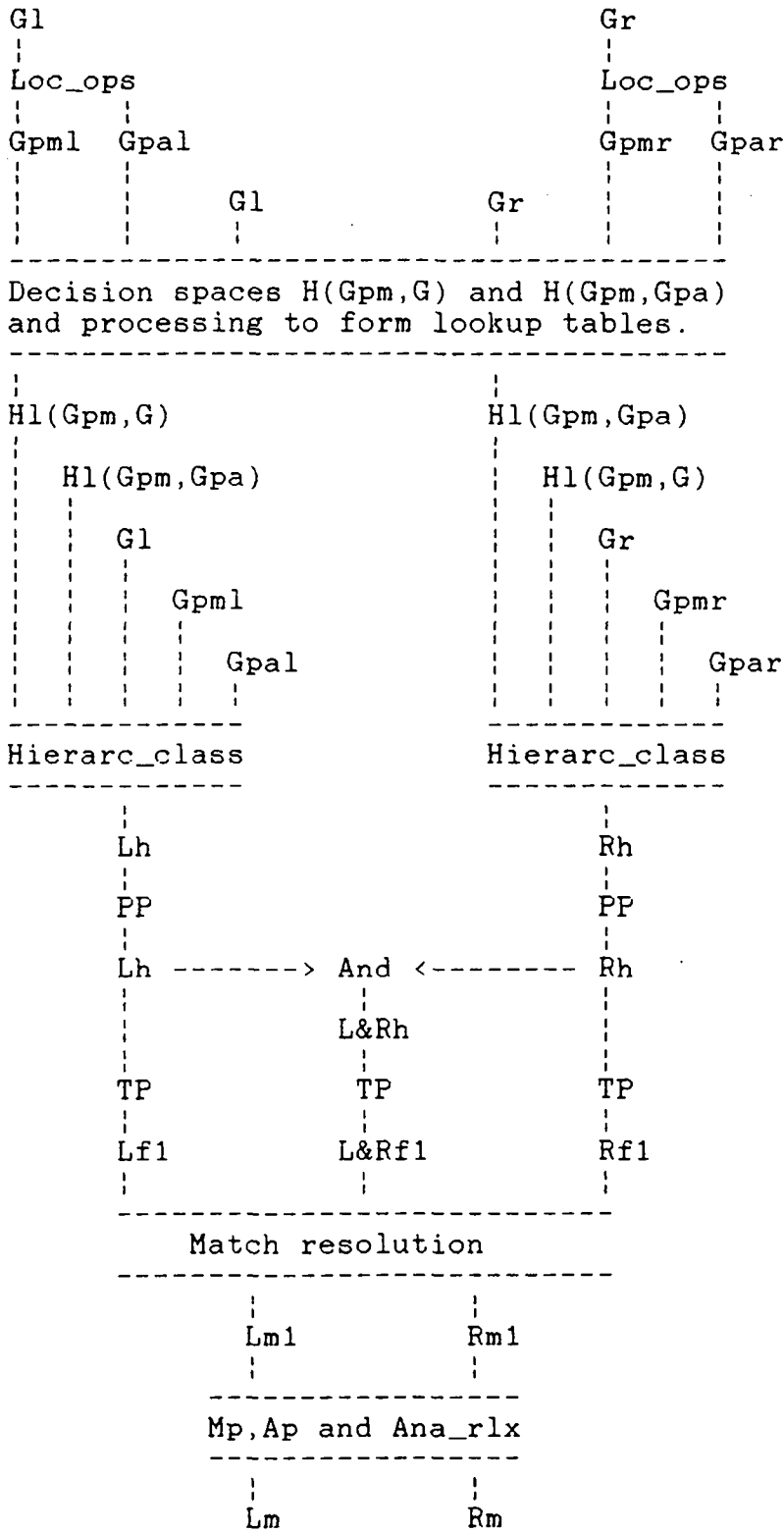


Figure 3.2.2-1a: "And-ing" of class labels and match resolution. The symbols are in Fig. 3.2.2-1b.

*l = Associated with left image.
 *r = Associated with right image.
 G = Gray level image.
 Gl = Gray level image of stereo pair, left.
 Gr = Gray level image of stereo pair, right.
 Loc_ops = Local operators within neighborhood N.
 Gpm = Gradient magnitude.
 Gpml = Gradient magnitude, left.
 Gpmr = Gradient magnitude, right.
 Gpa = Gradient angle.
 Gpal = Gradient angle, left.
 Gpar = Gradient angle, right.
 H(Gpm,G) = Decision space on Gpm and G for left and right.
 H(Gpm,Gpa) = Decision space on Gpm and Gpa for left and right.
 Hl(Gpm,G) = Lookup table for classifying Gpm and G.
 Hl(Gpm,Gpa) = Lookup table for classifying Gpm and Gpa.
 Hierarc_class = Hierarchical classification.
 *h = Class labelled image.
 Lh = Class labelled image, left.
 Rh = Class labelled image, right.
 PP = Lh and Rh pair post-processing.
 And = The "and-ing" process.
 L&Rh = "And-ed" class labels.
 TP = Lh, L&Rh, and Rh triplet processing.
 *f = Facet (region) labelled image.
 Lf1 = Facet (region) labelled image, left.
 Rf1 = Facet (region) labelled image, right.
 L&Rf1 = Facet (region) labelled image of L&Rm.
 Lm1 = Matched image, left.
 Rm1 = Matched image, right.
 Mp,Ap and Ana_rlx = Subsequent processes.
 Mp = Match processing.
 Ap = Analytic approximation.
 Ana_rlx = Analytic relaxation.
 Lm = Matched image, left.
 Rm = Matched image, right.

Figure 3.2.2-1b: Explanation of symbols for matching based on "and-ing", see block diagram in Fig. 3.2.2-1a.

(Relevant computer files: /user2/kasvand/ima2):

L&Rh = Inr100.fig	Lh = Inr101.fig	Rh = Inr102.fig
Lf1 = Inr103.fig	Rf1 = Inr104.fig	
Lm1 = Inr107.fig	Rm1 = Inr108.fig	
Lm = Inr115.fig	Rm = Inr116.fig	
Mp,Ap and ana_rlx:	Inr109.fig	Inr110.fig
	Inr111.fig	Inr112.fig
		Inr113.fig

```

.. AAAA BBBB AAAAA BBBB AAAAA BBBB ..... = Labels in Lh
..... AAAA BBBB AAAAA BBBB AAAAA BBBB = Labels in Rh
..... AAA . BB .. AA ... B ..... = Labels in L&Rh

```

Figure 3.2.2-2: A case of incorrect results from "and-ing". The labels are A for one class and B for the other.

Figure titles:

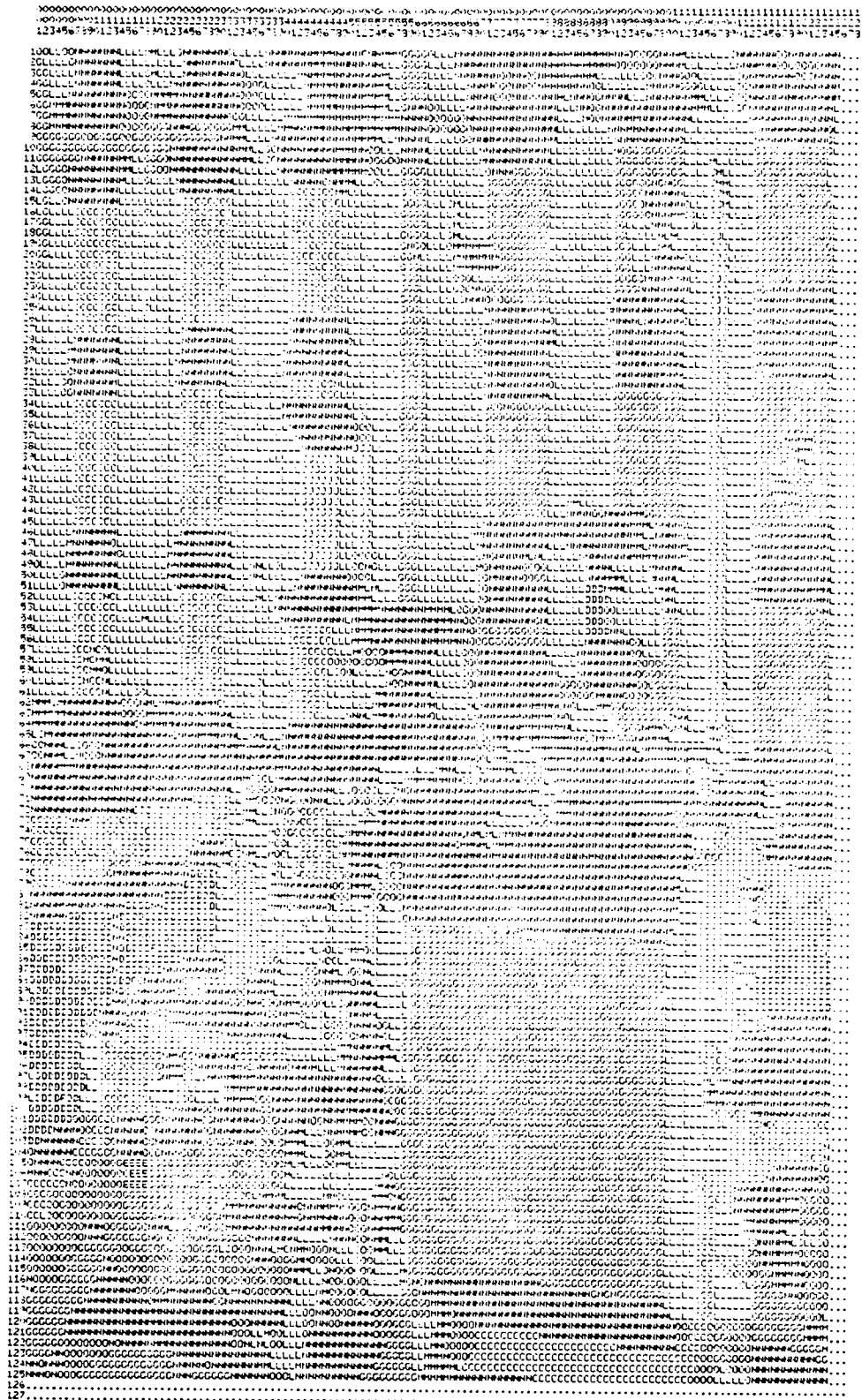
Figures 3.2.2-3a and -3b: The cluster (or class) labels from $Hl(Gpm,G)$ and $Hl(Gpm,Gpa)$ projected into the image space, after some processing. (Files: Inr101.fig, Inr102.fig)

Figure 3.2.2-3c: The "and-ed" $L\&Rh(i,j)$ image showing the class labels from $Hl(Gpm,G)$ and $Hl(Gpm,Gpa)$. (File: Inr100.fig)

Figure 3.2.2-4: The "painted" classes. The pixel classes are represented such that adjacent regions have different gray levels ("colours").

Figures 3.2.2-5a and -5b: The "painted" classes in Figure 3.2.2-4 shown with symbols. (Files: Inr193.fig, Inr194.fig)

Figures 3.2.2-6a and -6b: The stereo pairs for Figure 3.2.2-4 and for the original gray levels.



LN1:J, left: class labelled and post processed image.

LABEL PRINTING

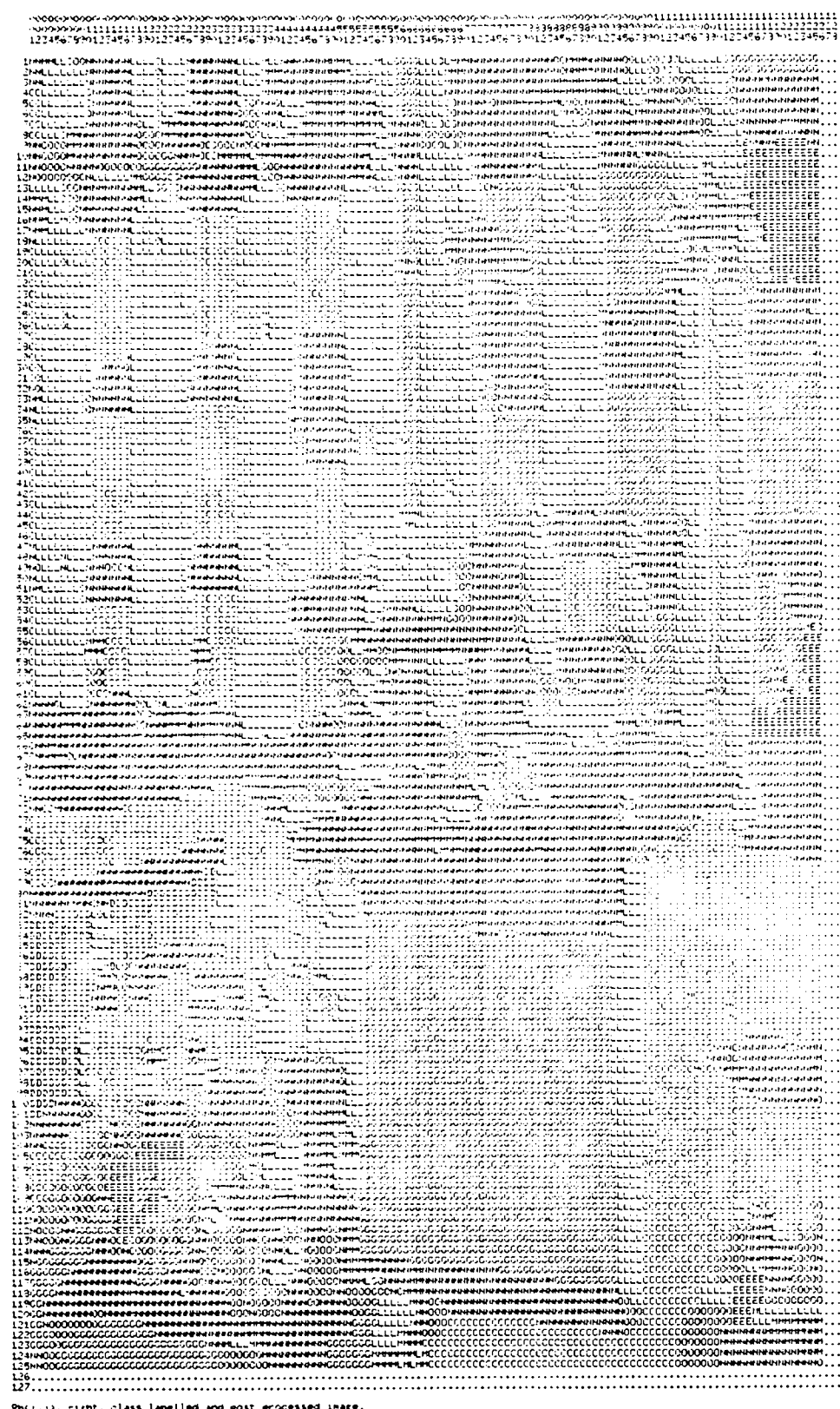
```

# 0 # 1 # 2 # 3 # 4 # 5 # 6 # 7 # 8 # 9 # 10 # 11 # 12 # 13 # 14 # 15 # 16 # 17 # 18 # 19 # 20 # 21 # 22 # 23 # 24 # 25 # 26 # 27 # 28 # 29 # 30 # 31 # 32 # 33 # 34 # 35 # 36 # 37 # 38 # 39 # 40 # 41 # 42 # 43 # 44 # 45 # 46 # 47 # 48 # 49 # 50 # 51 # 52 # 53 # 54 # 55 # 56 # 57 # 58 # 59 # 60 # 61 # 62 # 63 # 64 # 65 # 66 # 67 # 68 # 69 # 70 # 71 # 72 # 73 # 74 # 75 # 76 # 77 # 78 # 79 # 80 # 81 # 82 # 83 # 84 # 85 # 86 # 87 # 88 # 89 # 90 # 91 # 92 # 93 # 94 # 95 # 96 # 97 # 98 # 99 # 100 # 101 # 102 # 103 # 104 # 105 # 106 # 107 # 108 # 109 # 110 # 111 # 112 # 113 # 114 # 115 # 116 # 117 # 118 # 119 # 120 # 121 # 122 # 123 # 124 # 125 # 126 # 127

```

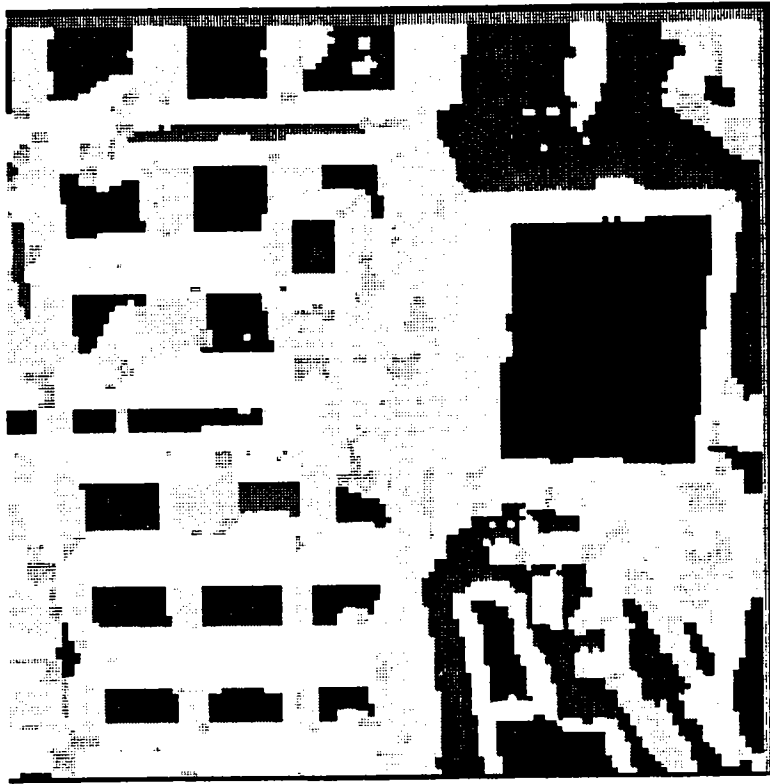
Fig.: 3.2.2-3a.

File: Inr101.Fig

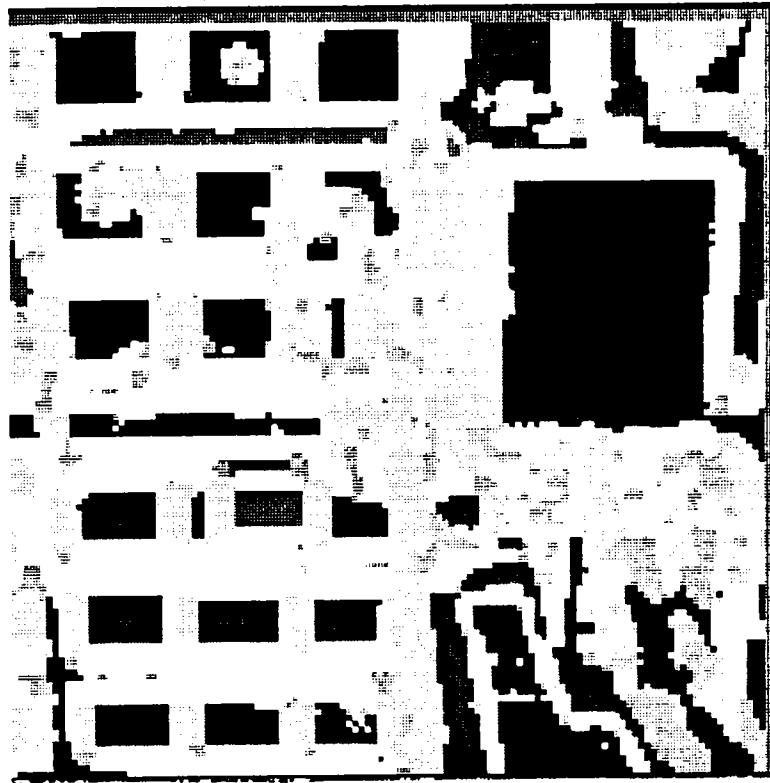


One of right class labelled and most processed image.
L=LABEL PRINTING
a 3 = 3 c= 3 d= 3 e= 4 f= 3 g= 4 h= 7
i= 3 j= 3 k= 10 l= 11 m= 12 n= 14 o= 15
p= 16 q= 17 r= 18 t= 19 u= 21 v= 22 w= 23
x= 24 z= 25 aa= 26 ab= 27 ac= 29 ad= 30 .ETC.
A 0 = 2 B 1 = 2 C 2 = 3 D 3 = 1
E 4 = 9 F 5 = 10 G 6 = 11 H 7 = 14 I 8 = 15
J 9 = 16 K 10 = 17 L 11 = 18 M 12 = 19 N 13 = 14
O 14 = 15 P 16 = 17 Q 17 = 18 R 18 = 19 S 19 = 20 .ETC.
File: im102.fig

Fig. : 3.2.2-3b.



if8hkpp2.rjr



if8hkpp2.lgr

Fig.: 3.2.2-4.

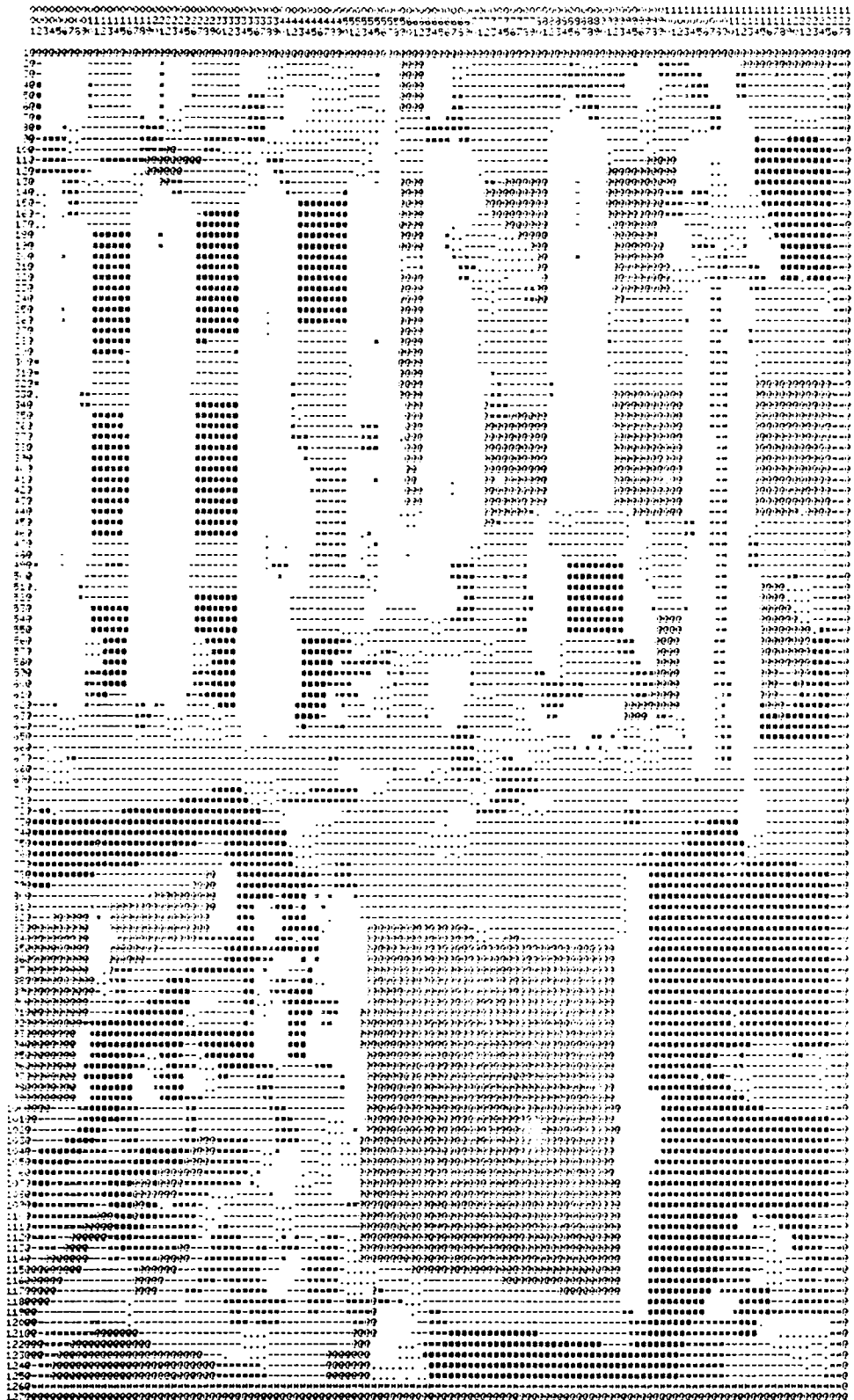


Fig.: 3.2.2-5b.



if8hkpp2.lgr

if8hkpp3.rgr

Fig.: 3.2.2-6a.



if8gray.lsh

if8gray.rsh

Fig.: 3.2.2-6b.

3.2.3 Correlation matching (not completed)

Matching based on correlation has been used for a long time. In photogrammetry the correlation uses fixed size areas (regions) from the left and right images. The (possibly enhanced) gray levels are cross-correlated to find the "shift" for the best match. Alternately, Fourier methods are used for the same purpose. The present case differs in two aspects:

- a) The regions to be correlated have been found by other means, i.e., they are not "arbitrarily" defined.
- b) The size of the regions, their shape, and even combinations of regions may be used as additional "features" for matching.

In this form of matching each region (segment or facet) found in the L and R images is first given "individual identity" by labelling the connected groups of similar cluster labels. The cluster label may serve as a feature describing the facet or it may be ignored. Thus, either facets with similar features are correlated or all facets are correlated irrespective of their features. The shape of regions is "automatically included". The correlation method is, in principle, identical to the one used for matching binary images (where there is only one feature, i.e., 1 on a background of 0's).

In its simplest realization, the result is a huge two-dimensional correspondence matrix $C(F_m, F_n)$ giving the correlation between facets F_m and F_n , $m = 1, 2, 3, \dots, N_f$ and $n = 1, 2, 3, \dots, N_f$, where N_f is the highest facet label number in the $L_f(i, j)$ or $R_f(i, j)$ image. Two additional matrices $I(F_m, F_n)$ and $J(F_m, F_n)$ give the values of the "best shift" in i - and j -coordinates indicating where the best correlation was found. The shift between the $L_f(i, j)$ and $R_f(i, j)$ images should correspond to physical reality (zero at fixation point and constrained to the vicinity of the epipolar line). In order to obtain unique correspondence the matrix $C(F_m, F_n)$ is processed for the "best" correlation between a facet label pair (F_k, F_l) and the corresponding "best" shift (i_k, j_l) is saved, this pair is eliminated and the search is repeated for the "next best", etc. The labels in the original facet label images $L_f(i, j)$ and $R_f(i, j)$ are changed to new matching labels and the best shifts are used for continued processing. In a block diagram form the process is illustrated in Fig. 3.2.3-1. The programming has not been completed at the present moment.

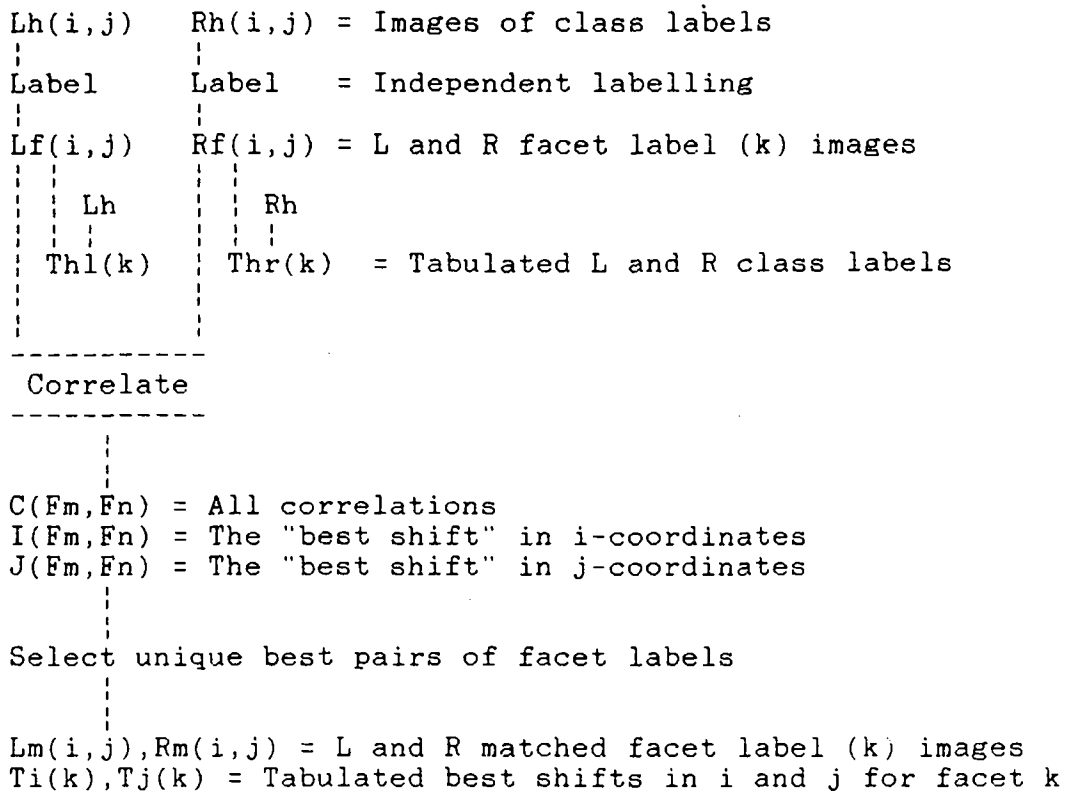


Figure 3.2.3-1: A block diagram of correlation matching.

3.2.4 Facet recognition matching

In matching based on some form of "mutual recognition" only the left versus right image "similarity principle" is used. The problem is to determine similarity between facets in the vicinity of the epipolar line given a set of facet descriptions.

In principle, the entire pattern recognition "paraphernalia" could be employed, and there are many methods available. Since there are many text on "pattern recognition", these will neither be described nor summarized. However, since only the "similarity principle" is needed and not "full recognition" of objects, these methods need to be "seen" from a slightly different point of view. For example, we may consider the left image features (near the epipolar line) as the "learning set" in terms of which the "objects" in the right image (near the epipolar line) are to be "recognized" (and/or vice versa). Thus, there is recognition in the "classical" terms (given a set of objects, acquire or learn the features for cluster or decision surface formulation, compare, etc.) but the objects "keep changing" during the process. However, the deviation from classical procedures is very slight and the needed modifications are obvious.

The number of possible descriptions for the facets is very large but it should be remembered that the segmentation process that produced the facets is prone to errors. Consequently, the descriptions should not be too "elaborate" at this stage. Furthermore, if the image contains repetitive (similar) facets then any "recognition" on the facet level is likely to produce ambiguous results (super-facets are needed). The presently used facet descriptions are:

- a) Facet class label from $Lh(i,j)$ and $Rh(i,j)$.
- b) Number of pixels per facet (area).
- c) Average gray level.
- d) Average x-gradient.
- e) Average y-gradient.
- f) A 4-dimensional shape feature.
- g) Centre of gravity in image plane.

The present "shape feature" consists of the standard deviations of the projections of the facet pixels onto 0, 45, 90, and 135 degree lines (since a "central" standard deviation only requires the centre of gravity for each projection). Most likely the "best" features are those that require no definition of a coordinate system for the facet. Other shape features, such as invariant moments, P^2/A , projected shapes, and so on, are easily obtainable. However, as stated, this is only a "first try" to determine what can be achieved with very unsophisticated methods.

Since the above "features" are a mix of different measures (quantitative, ordinal, and qualitative) which are not mutually comparable (the area of a facet has no relation to its average gray level, the average gray level is at most marginally related

to the average gradient, etc.), the "decision function" was of the following form:

$$D_{mn} = W_f * E_f + W_a * E_a + W_g * E_g + W_s * E_s + W_c * E_c$$

where:

D_{mn} = "Distance" between facets m and n .
 E_f = 0 if facet class labels are the same and
 = 1 if facet class labels differ for facets m and n .
 E_a = Difference in area sizes between m and n , normalized.
 E_g = Magnitude of gradient difference, normalized.
 E_s = Distance in 4-D "shape space", normalized.
 E_c = Distance between centers of gravity.
 W_f, W_a, W_g, W_s, W_c = weights.

$E_a = |A_m - A_n| / (A_m + A_n)$
 A_m = Area (pixel count) of facet m .
 A_n = Area (pixel count) of facet n .
 $E_g = |G'_m - G'_n| / (|G'_m| + |G'_n|)$
 G'_m = Gradient vector for facet m .
 G'_n = Gradient vector for facet n .
 $E_s = \text{Distance}(m,n) / (\text{Distance}(m) + \text{Distance}(n))$.
 $\text{Distance}(m,n)$ = Distance between m and n in sample space.
 $\text{Distance}(m)$ = Distance of sample m to origin.
 $\text{Distance}(n)$ = Distance of sample n to origin.
 $E_c = \text{Distance}((ic, jc)_m, (ic, jc)_n)$.

The match consisted in finding unique pairs of facets (m, n) constrained to the vicinity of the epipolar line for minimum D_{mn} . Thus, the "recognizer" is basically a minimum distance classifier. The weights allowed some "control" over the process.

The gradient directions on "cylindrical" facets are "opposed to each other" on the two sides of the "cylinder" and the directions on flat facets can be rather erratic if the gradient is small. Since the necessary modifications to compensate for these conditions were not included (at present) the weight W_g was set to zero ($W_g=0$) to eliminate the E_g term.

The shape descriptor E_s consisted of the standard deviations of the projections of the facet pixels onto a line. Consequently, this descriptor depends upon the size of the facet and, furthermore, it was not very sensitive to the shape of the facet (as was found from experimental results). A crude normalization consisted of dividing the distance between the sample points m and n ($\text{Distance}(m,n)$) by the sum of their distances to the origin ($\text{Distance}(m) + \text{Distance}(n)$).

The distance E_c between the centers of gravity (ic, jc) was an attempt to force the facets (m and n) to be "not too far from each other" in their respective image planes. This has some validity for "far away" facets but not for near facets in the scene that are far from the centre of fixation.

As stated, the match consisted in finding unique pairs of facets (m,n) constrained to the vicinity of the epipolar line for minimum D_{mn} , where the weights allowed some "control" over the process. The method is crude, ad hoc, and only intended as a "first exploration" of this problem in a reasonably flexible setting. The unique difference between the present problem and the "classical" pattern recognition procedures is that a pair-wise matching is required, i.e., an "object" from the left image is to be matched to exactly one "object" in the right image. This constraint partially exists in classical pattern recognition only if each object is required to be a separate class, but even then, the constraint that there must be only one object recognized per class (or none at all if the facet does not appear in one of the images) is not of the classical type. In addition there is a requirement (which in many cases is true) for an "ordering relation" between the recognized facets. However, these questions will not be considered any further.

The following two methods were used to try to resolve the uniqueness (one-to-one match) and the "mutual" recognition problems at the same time:

1. Find the largest facet "m" in the left (L) or the right (R) image (a facet which has not yet been marked as "done") and find the best matching facet "n" in the other image (R or L) which has not yet been "done". If the (minimal) distance D_{mn} between these facets was less than some threshold T, then declare this pair as "matched" and mark both as "done" to guarantee one-to-one matching. Iterate until all facets that can be matched are matched. The constraint to the epipolar line consisted of a highly rectangular box ($D_i * D_j$). The centers of gravity of the two "candidates" (m and n) had to be within this box for the matching to take place. The problem with this method is that it tends to escalate errors if T is large. If an error is made at a particular facet size then all the smaller facets will have to "choose among the leftovers" for their "mate".
2. The second method was similar to the first (1) except that the "choice for a mate" was not made. Instead all the matches (D_{mn}) less than the threshold T were stored in a "similarity matrix" (actually a list). When no further matches could be found (or when the list was full, the number of entries stored is controlled by T), then the matches were chosen according to "best first" until no more matches were possible.

The results presented no "surprises" since the methods are classical. However, it is best to apply them as "additional help" to other methods or for verification of the results from other methods rather than to use an elementary version of recognition directly.

3.2.5 Pixel feature matching

As indicated in Chapter 2, in addition to super-regions and regions obtained from various forms of clustering to create homogeneous regions, there is a profusion of pixel features obtainable from the local operators and also some features which are strictly on the pixels level (gray level of pixel in this case). Some of these features are statistically independent but mostly the dependence is unknown. In many cases these features are not directly comparable since some may be quantitative, others ordinal, or qualitative. Preferably, the features should be "independent" but, whatever method is to be used, it should allow for various "incompatibilities".

Clearly, the more features there are the more likely it is that each pixel in the left image L and the corresponding pixel in the right image R only have one (and only one) unique combination of features, at least in the epipolar region. The problem of matching is now reduced to finding these unique combinations, whatever they may be. The hierarchical methods will help in defining the regions in the vicinity of the epipolar lines and thereby limit the possible pixels to be matched, but this form of matching will basically produce pixel level correspondences directly.

In order to carry out these experiments in an "interesting manner" (rather than via endless searching of lists) would have required in the order of $M*M*N*F*P$ words of computer memory for an $M*N$ image using F features and P processes to extract the left and right dependencies. The problem is rather similar to processing "stacks" of medical images from a CAT scanner, for example.

In order to see what such an approach might lead to and to accommodate the immediately available computational equipment, a very much scaled down experiment was made. The images L and R were $128*127$ ($M=128$, $N=127$) with depth of 9 ($M*M$ was restricted only to a displacement of 9 pixels), with only one feature ($F=1$) and one process ($P=1$). The feature was the cluster label used in earlier studies. (Chronologically, this was the first matching experiment as soon as the cluster processes in Chapter 2 were completed). In essence, the results of this experiment "pointed out" that the classified regions in the images could either be matched by direct "and-ing" or matched by using correlation. Much more could be said about this approach but these would only be speculations since the experiment over-taxed the computer resources and were discontinued.

3.2.6 Other methods

In the hierarchical method (structure) there are regions of different sizes (from super-regions down to a few pixels or even sub-pixels in final resampling) which are the results of various classifications (clusters) when these are represented in the image space. Preliminary matches can be obtained by "matched classification", "and-ing", correlation, or recognition, either singly, or better, in suitable combinations.

As a first approximation, these matched facets could be directly "projected" into the 3D space by estimating the rotation matrix (A) and translation vector (B), i.e., $F_l = A*F_r + B$, where F_l and F_r represent the pixels, centers of gravity, moments, etc., of the matching facets in the left and right images (3.1). Approximate shapes of the facets (F_l, F_r) are available from the image space, adjacencies are known in image space, etc. In the reconstructed 3D space the adjacencies are likely to remain, and the 3D space should not contain any "cracks" between the reconstructed facets. (The interiors of the facets are "empty" at present). Thus, numerous optimization methods ("energy" functions, "rubber templates", "snakes", and so on), could be formulated to "relax" all the facets at the same time. If, in addition, the facets can "acquire pixels most likely to belong to them" then the methods include segmentation correction based on both the image data and physical reality in the 3D space. A lack of time has at present prevented a "deeper look" at these possibilities.

3.2.7 Comments

The results of the matching experiments that were based on facets may be summarized as follows:

Given large enough facets and a sufficient number of facets to "cover" the images then the matching problem is basically trivial. Exceptions occur for small facets, repetitive structures of similar facets, and in regions where the L and R images differ significantly.

Regions in the left (L) and right (R) images can differ markedly due to:

- a) Reflections which superimpose more or less "transparent" facets onto existing facets. This creates a new type of image processing problem that has not been studied.
- b) Loss of detail where the gray levels have become more or less "saturated" in one but not in the other image. If some significant information remains, these regions could be "equalized in gray" by some form of enhancement. The problem has not been studied in detail.

Matching failed when:

- c) The facets were either too small and did not generate an "intersection" in the L&R(i,j) image or they were too different in shape to allow "mutual recognition" with sufficient reliability.
- d) Due to classification errors where, for example, a small shift in the gray level in one image (L or R) places a well-matching facet into another category in the other image (R or L). This type of error is "fixed" in the next section.
- e) Incorrect facet pairs are formed. These are basically impossible to detect (automatically) at this stage of processing. Only after 3D reconstruction may it be possible to find some inconsistencies, but even this is unlikely since we ourselves are subject to certain visual illusions which can be corrected by an "alternate interpretation" of the scene or by "having another look" usually from a different point of view.
- f) Repetitive similar facets are present. If these facets are "displaced sufficiently" the "and-ing" process will create incorrect matches. If the facets are very similar then the recognition procedure cannot distinguish them. Clearly, such facets have to be first processed into "super-facets" by using similarity criteria before matching is attempted, i.e., the hierarchy in facet based matching needs several levels. This has not been done in the present case.
- g) The same region in the 3D scene resulted in very different facets in the corresponding regions in the L and R images. For example, there is one large facet in one image but many smaller ones in the other image (see the lower right corner of the INRIA image pair). This is a rather "insidious" error and required considerable (and somewhat dubious) processing to correct.

As already implied, some of the matching errors appeared to be "correctable". Other needed refinements became also fairly obvious, which are partially investigated in the next section.

In summary, the major weaknesses of the various methods were as follows:

1. Matched classification on pixel features requires many more features than the few that were attempted. Consequently, the experiment is inconclusive.
2. Direct "and-ing" of class labels fails for small facets and can generate "catastrophic" errors if repetitive structures of (small) facets have not been grouped first.

3. Correlation matching produces ambiguous results on repetitive structures. At best these can be declared as "unmatchable" since the correlation is (nearly) the same. The experiments were not carried out.
4. Pattern recognition methods cannot distinguish between many similar facets in an image. At best these facets can be declared as "unmatchable" since they are very similar.
5. Direct pixel feature matching could not be carried out due to equipment limitations.

This leaves the question: Is there a "best combination" of methods, and if so, what is this combination?

- i) The preliminary processing has to include detection of regular structures of similar facets. These structures become "super-facets". Clearly, the hierarchy needs several levels and not just the one described so far. One has to be able to "mask out" regions in the images during feature computations in order to "reprocess" selected regions. These aspects have been studied in connection with other problems but have not been included in the present case.
- ii) More features and additional decision spaces should be studied.
- iii) Use "and-ing" for large facets since the probability that they overlap in L and R is very high, but a few exceptions can be created if one of the cameras only "sees" the edge of a flat surface while the other sees much more of the same surface. The ultimate extreme is to place the cameras so that one camera sees one side and the other the other side of a thin flat sheet (in 3D space).
- iv) Apply recognition procedures to the leftover facets after direct "and-ing" (iii).
- v) Carry out various "matching correction" procedures to be described in the next section. Once the nature of a problem is known then something can be done about it, unless the weakness in a method is "fatal".

3.3 Modifying the various matches

The "raw" matching results from the previous section are of the following two types:

- a) An "intersection" image $L\&R_h(i,j)$ giving the class labels which intersected in the two images. The intersection could simply be of the "and" type, or a more "refined" one obtained from "matched classification". The $L\&R_h(i,j)$ image is accompanied by the left and right class label images $L_h(i,j)$ and $R_h(i,j)$. Independent labelling of the connected regions followed by "match (intersection) resolution" produces a matched pair of images of the second type (b).
- b) A left and right matched image pair $L_m(i,j)$ and $R_m(i,j)$ where the matching facets have the same label numbers. Non-matching facets are given negative and differing numbers in $L_m(i,j)$ and $R_m(i,j)$ in order not to lose facet information. The $L_m(i,j)$ and $R_m(i,j)$ image pair may or may not be accompanied by i - and j - coordinate "displacement" lists $T_i(k)$ and $T_j(k)$ for the matched facets k , $k = 2, 3, \dots, K_{max}$. The label $k=1$ is not used for "historic" reasons (old programs needed 1 to indicate "raw binary" data).

Several new problems were encountered, which will be considered in some detail in the following sub-sections. However, briefly stated:

1. Sets of connected regions in the $L\&R_h(i,j)$ "intersection" or "coincidence" image frequently "implicated" sets of facets in the $L_h(i,j)$ and $R_h(i,j)$ images. This required "match (intersection) resolution".
2. Depending on which process was used previously and the tolerance settings, two supposedly "matched" facets could be of markedly different shape.
3. Obvious errors occurred where correct matching was "missed" due to differences in original classification.
4. Questions arose about what to do with the unmatched, un-matchable, and undefined (unlabelled) facets.
5. Computation of suitably "balanced" analytic parameters for the facets in order to try "matched" analytic relaxation (in section 3.4).

Methods of checking that the results after processing are at least approximately correct were the same as before, i.e., visual inspection of stereo pairs of results and the "film-loop" in the Vicom image processing and display system.

3.3.1 The Lh, L&Rh, and Rh triplet

The "coincidence" or "intersection" image $L\&Rh(i,j)$ gives the class labels which intersected in the two individually classified $Lh(i,j)$ and $Rh(i,j)$ images. In the more "refined" or "matched classification", the $L\&Rh(i,j)$ image contained "matched" regions which could be considerably larger than those obtained in simple "and-ing". Both types of $L\&Rh(i,j)$ images were accompanied by the left and right class label images $Lh(i,j)$ and $Rh(i,j)$.

The regions in the "more refined" $L\&Rh(i,j)$ image were modified by requiring a match between $L\&Rh(i,j)$, $Lh(i,j)$, and $Rh(i,j)$ class labels or by requiring at least that one of the labels in $Lh(i,j)$ or $Rh(i,j)$ be nonzero. This operation reduced the $L\&Rh(i,j)$ image to the same form as that obtained from direct "and-ing". This simplified the subsequent processing.

The subsequent processing steps consisted of the following, see Figure 3.3.1-1:

1. Independently label the connected regions with similar class labels in the $L\&Rh(i,j)$, $Lh(i,j)$, and $Rh(i,j)$ images. Call these images $L\&Rf0(i,j)$, $Lf0(i,j)$, and $Rf0(i,j)$.
2. Compute the sizes and "thinnesses" of the labelled regions in $L\&Rf0(i,j)$, $Lf0(i,j)$, and $Rf0(i,j)$.
3. Reject all regions in $L\&Rf0(i,j)$, $Lf0(i,j)$, and $Rf0(i,j)$ that were either too small or too thin to allow subsequent analytic approximation.
4. Relabel according to area size the "surviving" facets in $L\&Rf0(i,j)$, $Lf0(i,j)$, and $Rf0(i,j)$. These images will now be referred to as $L\&Rf1(i,j)$, $Lf1(i,j)$, and $Rf1(i,j)$, where the "f" indicates facet labels and the "1" indicates an initial image.
5. Resolve the "set intersection" (match or intersection resolution) problem since sets of connected regions (facets) in the $L\&Rf1(i,j)$ image frequently indicated sets of facets in $Lf1(i,j)$ and $Rf1(i,j)$ images.

The problem (5) could be resolved by creating a huge "coincidence" matrix, or by processing lists, which was used and which became an interesting programming exercise. The problem is resolved by finding unique label pairs in $Lf1(i,j)$ and $Rf1(i,j)$ which have the highest "pixel count" in $L\&Rf1(i,j)$, these are eliminated and the comparison problem is repeated until the pixel count is "too low". Additional requirements are that the matching pair in $Lf1(i,j)$ and $Rf1(i,j)$ have approximately the same areas (sizes) and that the intersection areas are "comparable" to the areas of the two facets. (Additional "pattern recognition" procedures have not yet been introduced). Thus:

If((Amn.Gt.T1).And.(Bmn.Gt.T2).And.(Cmn.Gt.T3)) Then
 assign a new label to the facet pair in Lf1(i,j) and Rf1(i,j).

where:

Amn = Intersection area from L&Rf1(i,j) for facets m and n in
 Lf1(i,j) and Rf1(i,j).

Bmn = Amn/((Am+An)/2) = Normalized intersection area size for
 facets m and n.

Cmn = |Am-An|/(Am+An) = Normalized area difference for facets m
 and n.

Am = Area of facet m.

An = Area of facet n.

T1 = Threshold

T2 = Threshold

T3 = Threshold

The new labels are tabulated in terms of the facet labels in
 Lf1(i,j) and Rf1(i,j) and a new matched set of labelled images
 called Lm1(i,j) and Rm1(i,j) is created. The labels that could
 not be matched are set negative and they differ in Lm1(i,j) and
 Rm1(i,j) (but the label number assignment is continued such that
 the absolute value of a label can be used in addressing).

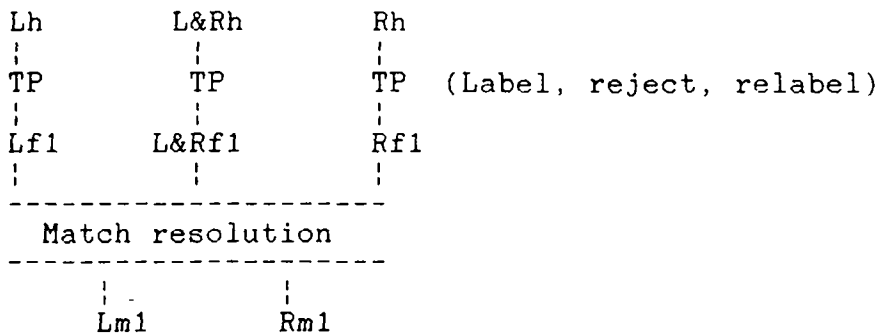
Three "heuristic" parameters, namely T1, T2, and T3, have been
 introduced for the following reasons:

T1 is to eliminate very small intersections in L&Rf1(i,j) which
 are basically unreliable.

T2 is a limit on the intersection area size for the two facets.
 Large regions should intersect more than small ones on a
 relative scale.

T3 is a "mutual recognition" limit. If the two facets differ
 widely in size, clearly they should not be considered as
 matching.

Prior experiments indicated that if the heuristic limits (T1, T2,
 T3) are absent then facets of rather different sizes could be
 associated despite the "best match first" condition in the pair-
 ing.



The process as shown previously.

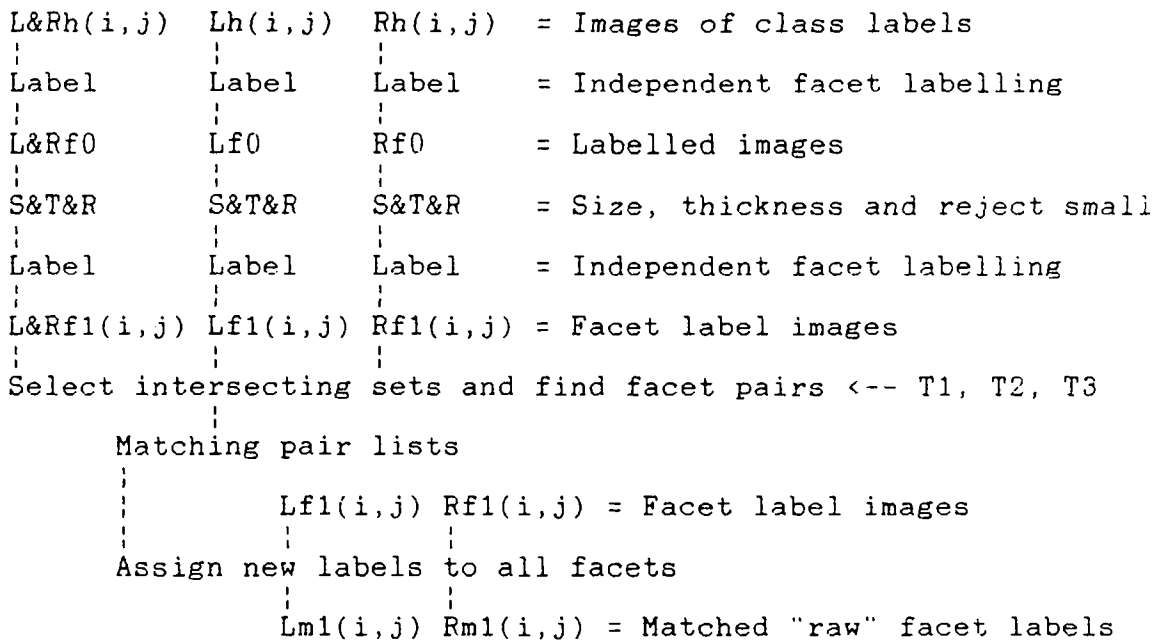


Figure 3.3.1-1: A block diagram of the processing of the L&Rh(i,-j), Lh(i,j), and Rh(i,j) images of class labels.

(Relevant computer files: /user2/kasvand/ima2):
 L&Rh = Inr063.fig Lh = Inr065.fig Rh = Inr066.fig
 Lf1 = Inr069.fig Rf1 = Inr070.fig
 Matching pair lists = Inr071.fig
 Lm1 = Inr073.fig Rm1 = Inr074.fig

A rather extreme output image triplet $L\&Rh(i,j)$, $Lh(i,j)$, and $Rh(i,j)$ was obtained from "Matched Classification" (Figures 3.2.1-2, 3.2.1-3a and -3b) and repeated as Figures 3.3.1-2a, -2b, and -2c. In this case the tolerances were intentionally made very liberal in order to illustrate what happens. After labelling and removal of small and/or very thin facets, the "surviving" facet labels $Lf1(i,j)$ and $Rf1(i,j)$ are shown in Figures 3.3.1-3a and -3b. The matching pair list is not shown since it contains about 220 rows of entries comprehensible only if the details of the programs are known. The resultant "raw" matched facet labels $Lm1(i,j)$ and $Rm1(i,j)$ are shown in Figures 3.3.1-4a and -4b. The matched facets are shown with capital letters while the unmatched ones use lower case letters.

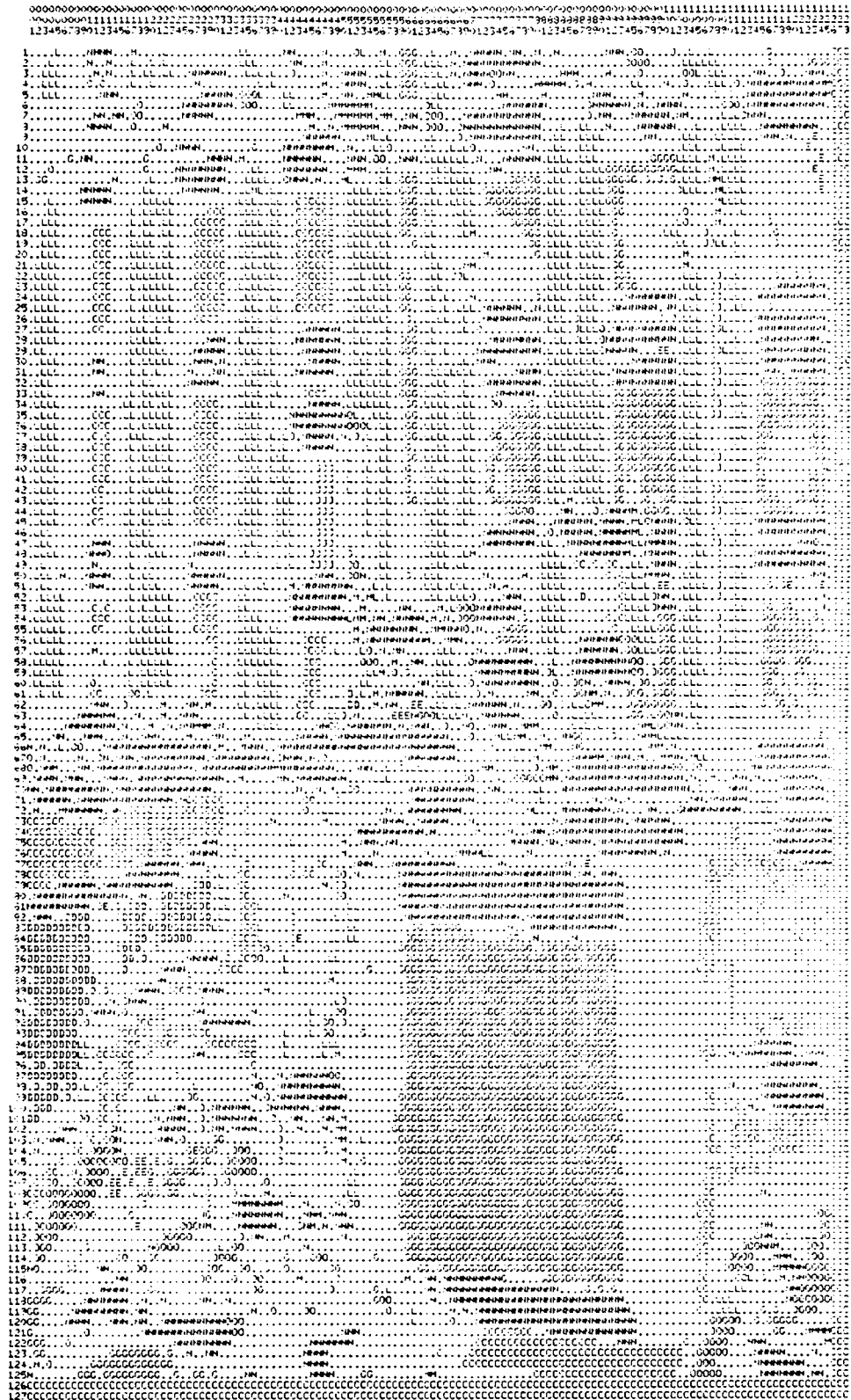
A closer examination of the results ($Lm1(i,j)$ and $Rm1(i,j)$ images) reveals that most of the matches appear to be correct, but there are also some rather comical matches which are not physically feasible in a 3D world. However, the match resolution program knows nothing about the 3D world. This matching should be followed by "recognition matching" in order to pair more facets since most of the unmatched facets are very similar but shifted "beyond the range" of these procedures.

Figure titles:

Figures 3.3.1-2a, -2b, and -2c: Output image triplet $L\&Rh(i,j)$, $Lh(i,j)$, and $Rh(i,j)$ obtained from "Matched Classification" (Figures 3.2.1-2, 3.2.1-3a and -3b).
(Files: Inr063.fig, Inr065.fig, Inr066.fig)

Figures 3.3.1-3a and -3b: The "surviving" facet labels $Lf1(i,j)$ and $Rf1(i,j)$ after removal of small and/or very thin facets.
(Files: Inr069.fig, Inr070.fig)

Figures 3.3.1-4a and -4b: The "raw" matched facet labels $Lm1(i,j)$ and $Rm1(i,j)$, same as Figs. 3.3.2-1a and -1b.
(Files: Inr073.fig, Inr074.fig, Matching list = Inr071.fig)



LRN(i,j) class labels from matched classification. (opt1.2.3 = 1.0.)

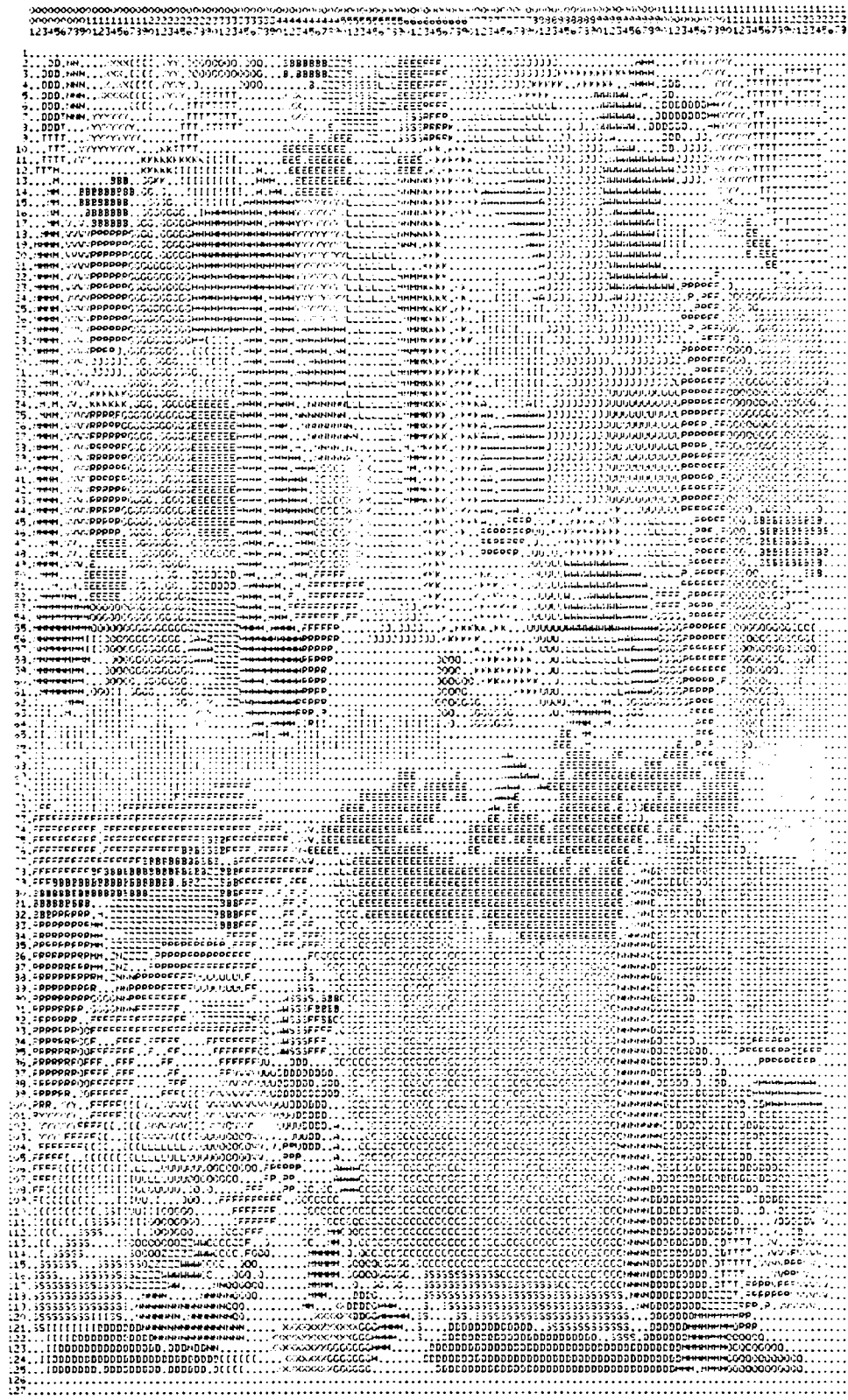
LABEL PRINTING

a	9	=	1	b	=	2	c	=	3	d	=	4	e	=	5	f	=	6	g	=	7		
i	=	8	j	=	9	k	=	10	l	=	11	m	=	12	n	=	13	o	=	14	p	=	15
q	=	16	r	=	17	s	=	18	t	=	19	u	=	20	v	=	21	w	=	22	x	=	23
0	=	24	1	=	25	2	=	26	3	=	27	4	=	28	5	=	29	6	=	30	ETC		

a= 9 b= 2 c= 3 d= 4 e= 5 f= 6 g= 7
 i= 8 j= 9 k= 10 l= 11 m= 12 n= 13 o= 14 p= 15
 q= 16 r= 17 s= 18 t= 19 u= 20 v= 21 w= 22 x= 23
 y= 24 z= 25 aa= 26 ab= 27 ac= 28 ad= 29 ae= 30 etc.

Fig.: 3.3.1-2a.

File: lrn063.fig



right, facet labels.

```

MSEL PRINTING
* 3 10 3 0 3 1 0 4 1 0 11 10 13 10 14 10 15
* 16 17 10 13 17 19 20 18 21 18 23 17
* 24 25 16 26 27 28 29 30 31 30 32 30
* 33 34 31 32 33 34 35 36 37 36 38 36
* 39 40 37 41 42 43 44 45 46 45 47 45
* 48 49 46 50 51 52 53 54 55 54 56 54
* 57 58 55 59 60 61 62 63 64 63 65 63
* 66 67 64 68 69 70 71 72 73 72 74 72
* 75 76 73 77 78 79 80 81 82 81 83 81
* 84 85 82 86 87 88 89 90 91 90 92 90
* 93 94 91 95 96 97 98 99 100 99 101 99
* 102 103 100 104 105 106 107 108 109 108 110 108
* 111 112 109 113 114 115 116 117 118 117 119 117
* 120 121 118 122 123 124 125 126 127 126 128 126
* 129 130 127 131 132 133 134 135 136 135 137 135
* 138 139 136 140 141 142 143 144 145 144 146 144
* 147 148 145 149 150 151 152 153 154 153 155 153
* 156 157 154 158 159 160 161 162 163 162 164 162
* 165 166 163 167 168 169 170 171 172 171 173 171
* 174 175 172 176 177 178 179 180 181 180 182 180
* 183 184 181 185 186 187 188 189 190 189 191 189
* 192 193 190 194 195 196 197 198 199 198 200 198
* 201 202 199 203 204 205 206 207 208 207 209 207
* 210 211 208 212 213 214 215 216 217 216 218 216
* 219 220 217 221 222 223 224 225 226 225 227 225
* 228 229 226 230 231 232 233 234 235 234 236 234
* 237 238 235 239 240 241 242 243 244 243 245 243
* 246 247 244 248 249 250 251 252 253 252 254 252
* 255 256 253 257 258 259 260 261 262 261 263 261
* 264 265 262 266 267 268 269 270 271 270 272 270
* 273 274 271 275 276 277 278 279 280 279 281 279
* 282 283 280 284 285 286 287 288 289 288 290 288
* 291 292 289 293 294 295 296 297 298 297 299 297
* 300 301 298 302 303 304 305 306 307 306 308 306
* 309 310 307 311 312 313 314 315 316 315 317 315
* 318 319 316 320 321 322 323 324 325 324 326 324
* 327 328 325 329 330 331 332 333 334 333 335 333
* 336 337 334 338 339 340 341 342 343 342 344 342
* 345 346 343 347 348 349 350 351 352 351 353 351
* 354 355 352 356 357 358 359 360 361 360 362 360
* 363 364 361 365 366 367 368 369 370 369 371 369
* 372 373 370 374 375 376 377 378 379 378 380 378
* 381 382 379 383 384 385 386 387 388 387 389 387
* 390 391 388 392 393 394 395 396 397 396 398 396
* 399 400 397 401 402 403 404 405 406 405 407 405
* 408 409 406 410 411 412 413 414 415 414 416 414
* 417 418 415 419 420 421 422 423 424 423 425 423
* 426 427 424 428 429 430 431 432 433 432 434 432
* 435 436 433 437 438 439 440 441 442 441 443 441
* 444 445 442 446 447 448 449 450 451 450 452 450
* 453 454 451 455 456 457 458 459 460 459 461 459
* 462 463 460 464 465 466 467 468 469 468 470 468
* 471 472 469 473 474 475 476 477 478 477 479 477
* 480 481 478 482 483 484 485 486 487 486 488 486
* 489 490 487 491 492 493 494 495 496 495 497 495
* 498 499 496 500 501 502 503 504 505 504 506 504
* 507 508 505 509 510 511 512 513 514 513 515 513
* 516 517 514 518 519 520 521 522 523 522 524 522
* 525 526 523 527 528 529 530 531 532 531 533 531
* 534 535 532 536 537 538 539 540 541 540 542 540
* 543 544 541 545 546 547 548 549 550 549 551 549
* 552 553 549 554 555 556 557 558 559 558 560 558
* 561 562 559 563 564 565 566 567 568 567 569 567
* 570 571 568 572 573 574 575 576 577 576 578 576
* 579 580 577 581 582 583 584 585 586 585 587 585
* 588 589 586 590 591 592 593 594 595 594 596 594
* 597 598 595 599 600 601 602 603 604 603 605 603
* 606 607 604 608 609 610 611 612 613 612 614 612
* 615 616 613 617 618 619 620 621 622 621 623 621
* 624 625 622 626 627 628 629 630 631 630 632 630
* 633 634 631 635 636 637 638 639 640 639 641 639
* 642 643 640 644 645 646 647 648 649 648 650 648
* 651 652 649 653 654 655 656 657 658 657 659 657
* 660 661 658 662 663 664 665 666 667 666 668 666
* 669 670 667 671 672 673 674 675 676 675 677 675
* 678 679 676 680 681 682 683 684 685 684 686 684
* 687 688 685 689 690 691 692 693 694 693 695 693
* 696 697 694 698 699 700 701 702 703 702 704 702
* 705 706 703 707 708 709 710 711 712 711 713 711
* 714 715 712 716 717 718 719 720 721 720 722 720
* 723 724 721 725 726 727 728 729 730 729 731 729
* 732 733 730 734 735 736 737 738 739 738 740 738
* 741 742 739 743 744 745 746 747 748 747 749 747
* 750 751 748 752 753 754 755 756 757 756 758 756
* 759 760 757 761 762 763 764 765 766 765 767 765
* 768 769 766 770 771 772 773 774 775 774 776 774
* 777 778 775 779 780 781 782 783 784 783 785 783
* 786 787 784 788 789 790 791 792 793 792 794 792
* 795 796 793 797 798 799 800 801 802 801 803 801
* 804 805 802 806 807 808 809 810 811 810 812 810
* 813 814 811 815 816 817 818 819 820 819 821 819
* 822 823 820 824 825 826 827 828 829 828 830 828
* 831 832 829 833 834 835 836 837 838 837 839 837
* 840 841 838 842 843 844 845 846 847 846 848 846
* 849 850 847 851 852 853 854 855 856 855 857 855
* 858 859 856 860 861 862 863 864 865 864 866 864
* 867 868 865 869 870 871 872 873 874 873 875 873
* 876 877 874 878 879 880 881 882 883 882 884 882
* 885 886 883 887 888 889 890 891 892 891 893 891
* 894 895 892 896 897 898 899 900 901 900 902 900
* 903 904 901 905 906 907 908 909 910 909 911 909
* 912 913 910 914 915 916 917 918 919 918 920 918
* 921 922 919 923 924 925 926 927 928 927 929 927
* 930 931 928 932 933 934 935 936 937 936 938 936
* 939 940 937 941 942 943 944 945 946 945 947 945
* 948 949 946 950 951 952 953 954 955 954 956 954
* 957 958 955 959 960 961 962 963 964 963 965 963
* 966 967 964 968 969 970 971 972 973 972 974 972
* 975 976 973 977 978 979 980 981 982 981 983 981
* 984 985 982 986 987 988 989 990 991 990 992 990
* 993 994 991 995 996 997 998 999 1000 999 1001 999

```

Fig.: 3.3.1-3b.

File: In=070.fig

3.3.2 The raw L_m and R_m images

After the initial or "raw" matching was achieved (L_{m1}, R_{m1}) it was clear that there are numerous rather basic and possibly "philosophical" problems that have to be solved. The raw matched left and right images, i.e., $L_{m1}(i,j)$ and $R_{m1}(i,j)$, contained:

- a) Mainly correctly matched facets but some of these facets could be rather different in size if the tolerances were lowered.
- b) Unmatchable facets without one-to-one correspondence between the facets. For example, in one image a region could be split into two while in the other it remained a single region. The best matching "splinter" was matched leaving the other "splinter" as "unmatchable".
- c) Facets that should have been matched but were not matched due to different classification. For example, the region in one image could be darker than in the other image and, if the gray level was close to a class boundary, the regions became classified as "different".
- d) Unmatched facets that were "too far" from each other and did not generate an intersection in the $L\&R_h(i,j)$ image or the intersection was too small and was rejected.
- e) Unmatched facets that could not be matched since there was no corresponding facet in the other image or its "partner" had already been "grabbed" by another facet (see (b)).
- f) Matches which appeared to be errors since verification of the initial raw results is unreliable.

Most of these problems could be attributed to three causes, namely:

- i) Independent segmentation of the images produces differing sizes of facets in certain regions of the images. The so-called "matched classification" and the use of common decision spaces did not improve the situation markedly. Thus, in "still-life" stereo a certain number of segmentation errors will have to be accepted and dealt with in subsequent processing.
- ii) Facets which should have been "paired" could not be paired since their intersection in the $L\&R_h(i,j)$ image was zero or "too small", the initial classification was wrong, or the "partner" has already been "taken".
- iii) Facets in one image did not exist in the other image or the scene was rather different in certain regions mainly due to reflections and partially due to the geometry of the scene.

Numerous attempts were made to try to deal with these problems. There are no particularly strong "theoretical" justifications to what has been attempted, except to say that "at the time they appeared to be 'reasonable' thing to do". A part of the image pair is shown in Figures 3.3.2-1a and -1b, which are portions of Figures 3.3.1-4a and -4b. Careful inspection reveals that most of errors can be found in these images. However, it should not be forgotten that the printing of the labels is in modulu 26 and it can happen that different but touching labels have the same letter code.

Figure titles:

Figures 3.3.2-1a and -1b: Magnified extracts of the "raw" matched facet label images $Lm1(i,j)$ and $Rm1(i,j)$, from Figures 3.3.1-4a and -4b. (Files: Inr073.fig, Inr074.fig)

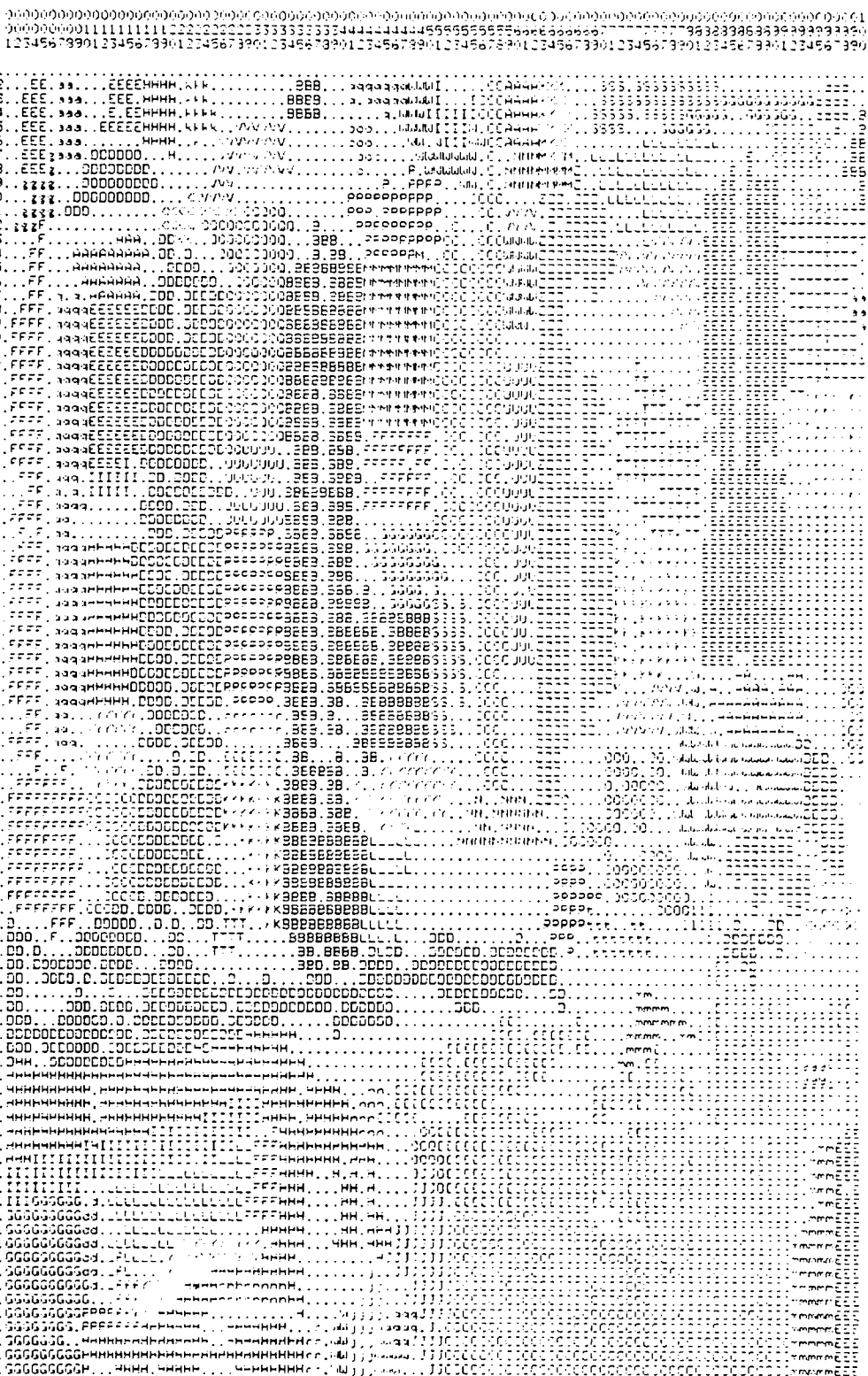


Fig.: 3.3.2-1b.

3.3.3 Misclassifications

Misclassifications were mainly caused by variations in gray levels between the two images ($G_l(i,j)$ versus $G_r(i,j)$) when the gray levels were very close to a class boundary. Since the classes in this case were not "widely separated" this situation will always occur with a segmentation method based on classification (of the whole image in one process).

This problem is easily solved by constructing a coincidence matrix of the unpaired facets and by using the "best match first" principle for assigning the pair a new label. The process is identical to the initial pairing of facets but without requiring a match of class labels. A "finessing" consists of also computing a "distance" between the classes and including the distance in the "best match first" criterion. In the present case the pairing was simply terminated as in the case of initial matching (when the intersection pixel count was less than a threshold).

An edited extract of a part of the "history file" is shown in Figure 3.3.3-1. In this figure the coincidence matrix indicates that the match between labels "-121=n" and "-123=p" has 108 coincident pixels (with temporary addresses (1,1) in the coincidence matrix), located at $(i,j)=(118,16)$ and $(i,j)=(120,14)$, etc. The result after correcting some misclassifications is shown in Figures 3.3.3-2a and -2b.

Figure titles:

Figure 3.3.3-1: Misclassification correction, edited extract from "history file". (File: Inr076.fig, extract from Inr075.fig)

Figures 3.3.3-2a and -2b: The result after correcting some misclassifications in the images shown in Figures 3.3.1-4a and -4b. (Files: Inr077.fig, Inr078.fig)

Corrections for miss-classifications.
 (Find matching negative labels after 'match resolution'.)

Coincidence matrix:

JanIa	1	2	3	4	5	6	7	8	9	10	11	12	13
1	108	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	2	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	2	0	0	0	0	0	0
5	0	0	0	0	0	17	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0	0	0	0	6	0
7	0	0	0	0	0	0	0	0	2	0	0	0	0
8	0	0	0	0	0	0	3	0	0	0	0	0	0
9	0	0	0	0	8	0	0	0	0	0	0	0	0
10	0	0	0	7	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	1	0	0	0	0	0	14	0
12	0	0	0	0	0	0	0	0	0	0	0	0	1
13	13	0	0	0	0	0	0	0	0	0	0	0	0
14	0	0	0	0	0	0	0	0	1	0	0	0	0
15	0	0	0	0	5	0	0	0	0	0	0	0	0
16	0	0	0	0	0	0	0	10	0	0	0	0	0
17	0	0	4	0	0	0	0	0	0	0	0	0	0
18	0	6	0	0	0	0	0	0	0	0	0	0	0

```

NEW LABEL: MAX,KLMAX,KPMAX,LABEL@,LABR@,LABNE@:
108 1 1 -121=0 -127=0 121=N
- LabI@ LabR@ NeI@ NeR@ LI@ LII@ LIR@ II@ JI@ I@ Jr@ Is@ Js@
121 -121=0 0=, 154 0 0 125 0 118 16 0 0 1 0
123 0=, -123=0 0 132 0 0 108 0 0 120 14 0 1
    
```

```

NEW LABEL: MAX,KLMAX,KPMAX,LABEL@,LABR@,LABNE@:
17 6 5 -140=F -150=0 140=F
- LabI@ LabR@ NeI@ NeR@ LI@ LII@ LIR@ II@ JI@ I@ Jr@ Is@ Js@
140 -140=F 0=, 37 0 0 20 0 12 115 0 0 0 0
150 0=, -150=0 0 22 0 0 17 0 0 11 115 0 0
    
```

```

NEW LABEL: MAX,KLMAX,KPMAX,LABEL@,LABR@,LABNE@:
14 11 11 -173=1 -177=0 173=1
- LabI@ LabR@ NeI@ NeR@ LI@ LII@ LIR@ II@ JI@ I@ Jr@ Is@ Js@
173 -173=1 -173=1 0 16 0 0 15 0 0 25 115 0 11
175 -175=1 0=, 16 0 0 14 0 26 115 0 0 11 0
    
```

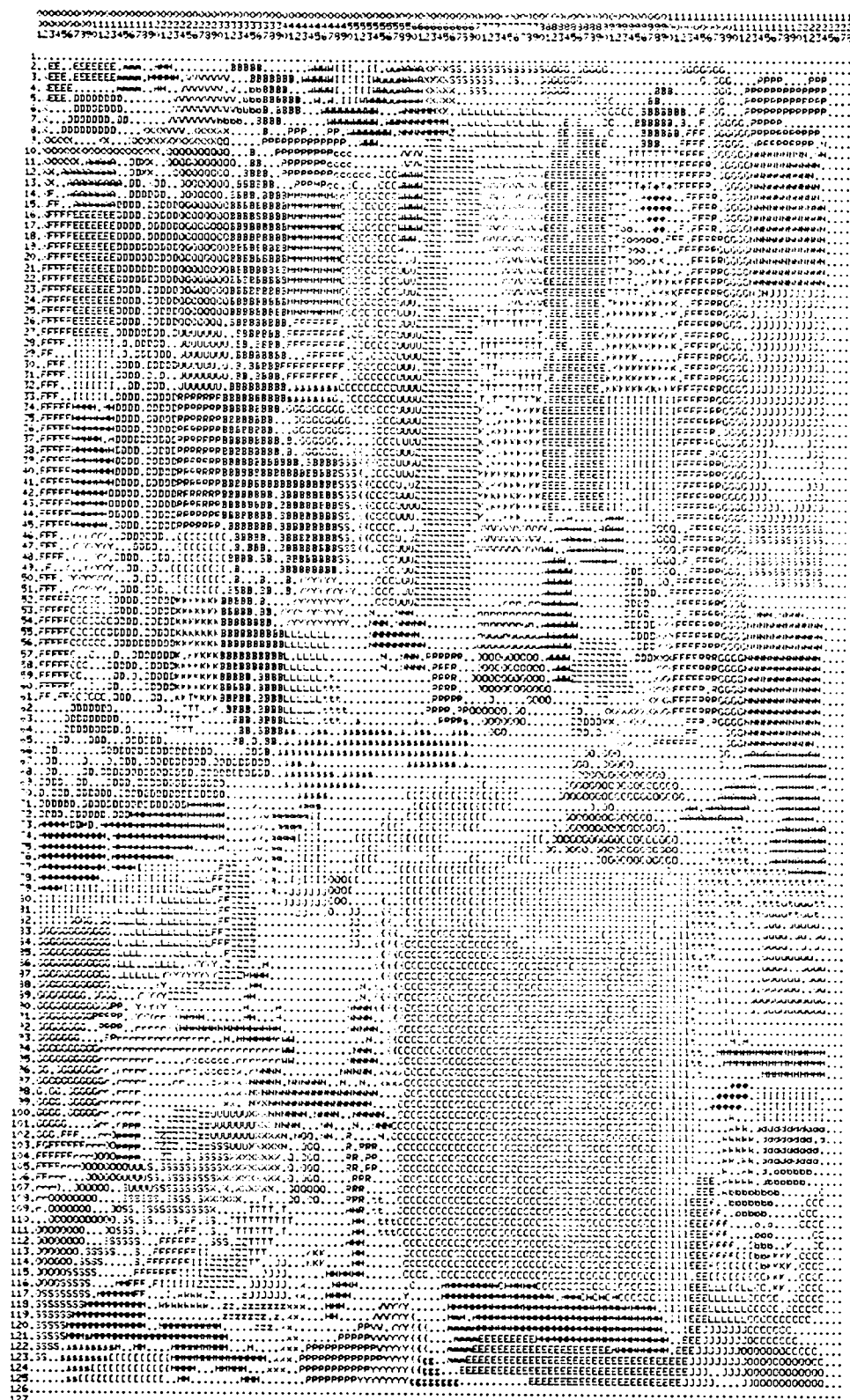
Edited extract of history file:

```

LabI = Left label.
LabR = Right label.
NeI = Number of elements, left.
NeR = Number of elements, right.
LI,JI = i, j coordinates of facets, left.
IR,JR = i, j coordinates of facets, right.
IS,JS = i, j coordinates of coincidence matrix.
    
```

File: Inc076.fig

Fig.: 3.3.3-1.



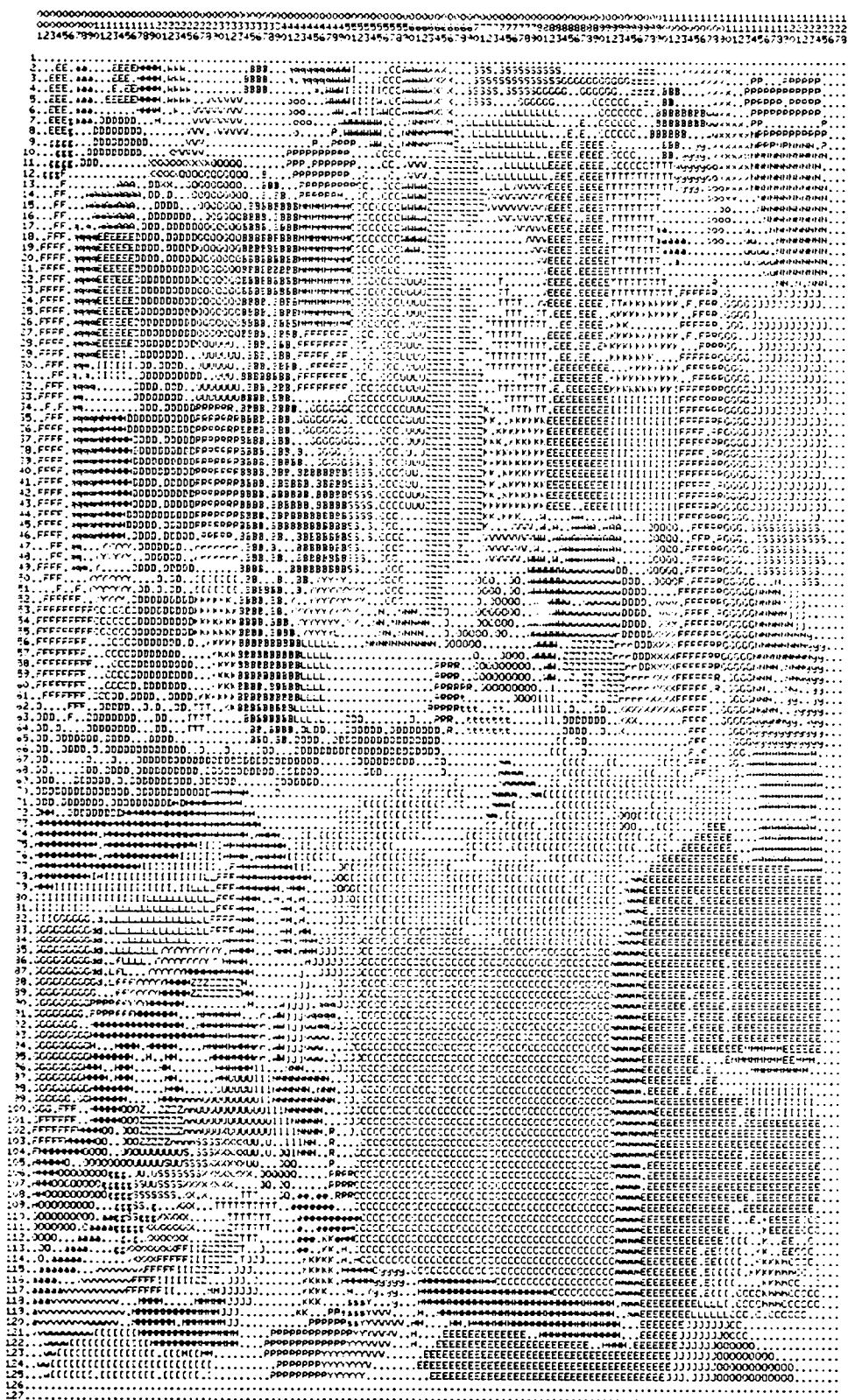
Label (j, left) matched facets after connection for misclassified facets. Area limit = 5.

Label PRINTING

Label	1	2	3	4	5	6	7
1	0	1	0	1	0	1	0
2	1	0	1	0	1	0	1
3	0	1	0	1	0	1	0
4	1	0	1	0	1	0	1
5	0	1	0	1	0	1	0
6	1	0	1	0	1	0	1
7	0	1	0	1	0	1	0
8	1	0	1	0	1	0	1
9	0	1	0	1	0	1	0
10	1	0	1	0	1	0	1
11	0	1	0	1	0	1	0
12	1	0	1	0	1	0	1
13	0	1	0	1	0	1	0
14	1	0	1	0	1	0	1
15	0	1	0	1	0	1	0
16	1	0	1	0	1	0	1
17	0	1	0	1	0	1	0
18	1	0	1	0	1	0	1
19	0	1	0	1	0	1	0
20	1	0	1	0	1	0	1
21	0	1	0	1	0	1	0
22	1	0	1	0	1	0	1
23	0	1	0	1	0	1	0
24	1	0	1	0	1	0	1
25	0	1	0	1	0	1	0
26	1	0	1	0	1	0	1
27	0	1	0	1	0	1	0
28	1	0	1	0	1	0	1
29	0	1	0	1	0	1	0
30	1	0	1	0	1	0	1
31	0	1	0	1	0	1	0
32	1	0	1	0	1	0	1
33	0	1	0	1	0	1	0
34	1	0	1	0	1	0	1
35	0	1	0	1	0	1	0
36	1	0	1	0	1	0	1
37	0	1	0	1	0	1	0
38	1	0	1	0	1	0	1
39	0	1	0	1	0	1	0
40	1	0	1	0	1	0	1
41	0	1	0	1	0	1	0
42	1	0	1	0	1	0	1
43	0	1	0	1	0	1	0
44	1	0	1	0	1	0	1
45	0	1	0	1	0	1	0
46	1	0	1	0	1	0	1
47	0	1	0	1	0	1	0
48	1	0	1	0	1	0	1
49	0	1	0	1	0	1	0
50	1	0	1	0	1	0	1
51	0	1	0	1	0	1	0
52	1	0	1	0	1	0	1
53	0	1	0	1	0	1	0
54	1	0	1	0	1	0	1
55	0	1	0	1	0	1	0
56	1	0	1	0	1	0	1
57	0	1	0	1	0	1	0
58	1	0	1	0	1	0	1
59	0	1	0	1	0	1	0
60	1	0	1	0	1	0	1
61	0	1	0	1	0	1	0
62	1	0	1	0	1	0	1
63	0	1	0	1	0	1	0
64	1	0	1	0	1	0	1
65	0	1	0	1	0	1	0
66	1	0	1	0	1	0	1
67	0	1	0	1	0	1	0
68	1	0	1	0	1	0	1
69	0	1	0	1	0	1	0
70	1	0	1	0	1	0	1
71	0	1	0	1	0	1	0
72	1	0	1	0	1	0	1
73	0	1	0	1	0	1	0
74	1	0	1	0	1	0	1
75	0	1	0	1	0	1	0
76	1	0	1	0	1	0	1
77	0	1	0	1	0	1	0
78	1	0	1	0	1	0	1
79	0	1	0	1	0	1	0
80	1	0	1	0	1	0	1
81	0	1	0	1	0	1	0
82	1	0	1	0	1	0	1
83	0	1	0	1	0	1	0
84	1	0	1	0	1	0	1
85	0	1	0	1	0	1	0
86	1	0	1	0	1	0	1
87	0	1	0	1	0	1	0
88	1	0	1	0	1	0	1
89	0	1	0	1	0	1	0
90	1	0	1	0	1	0	1
91	0	1	0	1	0	1	0
92	1	0	1	0	1	0	1
93	0	1	0	1	0	1	0
94	1	0	1	0	1	0	1
95	0	1	0	1	0	1	0
96	1	0	1	0	1	0	1
97	0	1	0	1	0	1	0
98	1	0	1	0	1	0	1
99	0	1	0	1	0	1	0
100	1	0	1	0	1	0	1
101	0	1	0	1	0	1	0
102	1	0	1	0	1	0	1
103	0	1	0	1	0	1	0
104	1	0	1	0	1	0	1
105	0	1	0	1	0	1	0
106	1	0	1	0	1	0	1
107	0	1	0	1	0	1	0
108	1	0	1	0	1	0	1
109	0	1	0	1	0	1	0
110	1	0	1	0	1	0	1
111	0	1	0	1	0	1	0
112	1	0	1	0	1	0	1
113	0	1	0	1	0	1	0
114	1	0	1	0	1	0	1
115	0	1	0	1	0	1	0
116	1	0	1	0	1	0	1
117	0	1	0	1	0	1	0
118	1	0	1	0	1	0	1
119	0	1	0	1	0	1	0
120	1	0	1	0	1	0	1
121	0	1	0	1	0	1	0
122	1	0	1	0	1	0	1
123	0	1	0	1	0	1	0
124	1	0	1	0	1	0	1
125	0	1	0	1	0	1	0
126	1	0	1	0	1	0	1
127	0	1	0	1	0	1	0

Fig.: 3.3.3-2a.

File: tw-077.fig



Printed, right matched facets after correction for misclassified facets. Area limit = 5.

LABEL PRINTING

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100
---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	-----

Fig.: 3.3.3-2b.

File: Ino78.fig

3.3.4 Cut

The matched facets that were markedly different in size could be "cut to size" since usually a facet in one image should not differ very much in size from the corresponding facet in the other image. Of course, exceptions exist which may cause new problems later. A block diagram of the "cut" process is shown in Figure 3.3.4-1.

```
L R = Lm2(i,j) Rm2(i,j) = Matched facet label images
| |
| |
| | Extract a matching facet pair
| |
| | \ /
| | Cut
| | / \
L R = Lm3(i,j) Rm3(i,j) = Matched "cut" label images
```

Figure 3.3.4-1: The "cutting" process.

The "cutting" process was based on defining a "maximal allowable" size difference rectangle ($D_i * D_j$), where D_i is the allowed size variation along the i-axis and D_j is the one along the j-axis. In the experiments $D_i = 2$ and $D_j = 7$ pixels. The basic program logic for any matched facet F is as follows, where F1 is a facet label in L and F2 is the corresponding facet label in R ($F1=F2$). (The method assumes that the matching pair F of facets in L and R have been extracted and stored in a separate pair of images):

```
If, for any pixel at (i,j) in L with label F1
the label F2 at ((i +/- Di),(j +/- Dj) in R is not equal to F1
Then zero the pixel (i,j) in R
Else leave pixel unchanged. (Similar logic for R versus L)
```

Thus, for any pair of matched facets F in the left (L) and right (R) images, the pixels that were outside the "box" defined by D_i and D_j are zeroed. These pixel regions are labelled with 0 and are called the "unknown" regions. Of course, the unmatched facets cannot be "cut" since they have no "partners".

Some results showing the effects of "cutting" are shown in Figures 3.3.4-2a and -2b. The most noticeable change occurred to the facet labelled "E" in Figure 3.3.1-4b, lower right corner, which has been "cut to size" in Figure 3.3.4-2b leaving a large "undefined" area. Experimental results tended to show that "cutting" was justified but the facets after "cutting" were still rather different in size since D_i and D_j were quite "liberal".

Figure titles:

Figures 3.3.4-2a and -2b: Results showing the effects of "cutting". (Files: Inr081.fig, Inr082.fig)

3.3.5 Transplant

The "transplanting" process was one attempt to try to fill the unknown regions in the L and R images. It was simply argued that, if a pixel (i,j) in R has no label (label=0) while there is a non-zero unmatched (negative) label in the same location (i,j) in L, then, why not simply "transplant" this label from L into R at location (i,j), and vice versa for R versus L. Thus, where .Lt. means "less than" and .Eq. means "equal to":

```
If((L(i,j).Lt.0).And.(R(i,j).Eq.0)) Then R(i,j)=L(i,j)
If((R(i,j).Lt.0).And.(L(i,j).Eq.0)) Then L(i,j)=R(i,j)
```

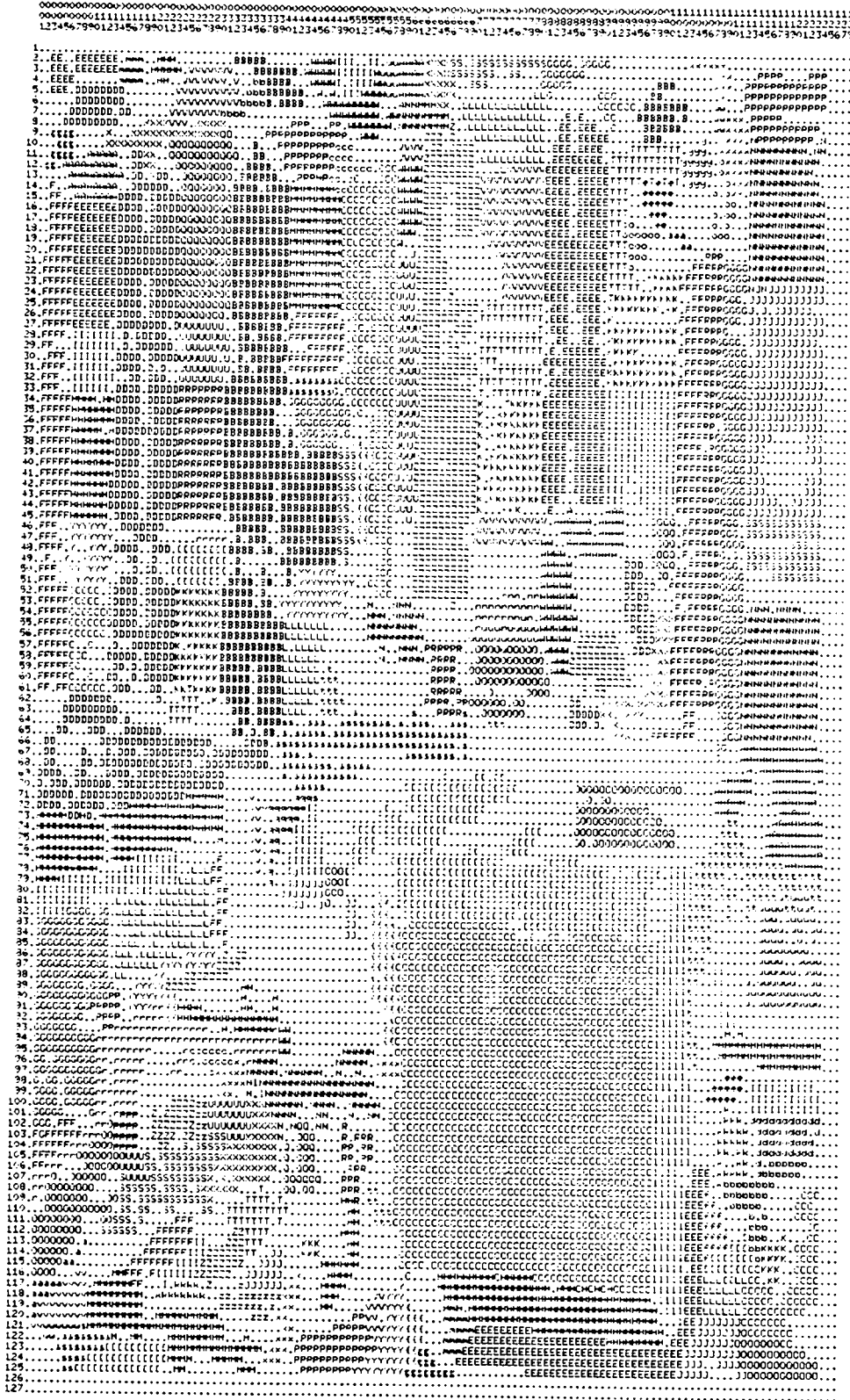
Of course, this is approximately valid only for "distant" facets in the scene but will generate drastic errors for "close" facets which are not near the centre of fixation. However, close facets should have large intersection areas and should not create this problem in the first place. A block diagram is shown in Figure 3.3.5-1 and some experimental results of "transplanting" are in Figures 3.3.5-2a and -2b. The most noticeable change again occurred in the lower right corner of the images where the original gray levels were rather different.

```
L R = Lm3(i,j) Rm3(i,j) = Matched "cut" label images
 \ /
Transplant
 / \
L R = Lm4(i,j) Rm4(i,j) = Matched facet label images
```

Figure 3.3.5-1: The label "transplanting" process.

Figure titles:

Figures 3.3.5-2a and -2b: Some results of "transplanting" the labels. (Files: Inr083.fig, Inr084.fig)



Left: matched labels after "out" and "transplant". Transplanted and unmatched labels are negative. D17 D1=2.

LABEL PRINTING

Q=	1	2	3	4	5	6	7
1	2	3	4	5	6	7	8
9	10	11	12	13	14	15	16
17	18	19	20	21	22	23	24
25	26	27	28	29	30	31	32
33	34	35	36	37	38	39	40
41	42	43	44	45	46	47	48
49	50	51	52	53	54	55	56
57	58	59	60	61	62	63	64
65	66	67	68	69	70	71	72
73	74	75	76	77	78	79	80
81	82	83	84	85	86	87	88
89	90	91	92	93	94	95	96
97	98	99	100	101	102	103	104
105	106	107	108	109	110	111	112
113	114	115	116	117	118	119	120
121	122	123	124	125	126	127	128

File: In=03.Fig

Fig.: 3.3.5-2a.

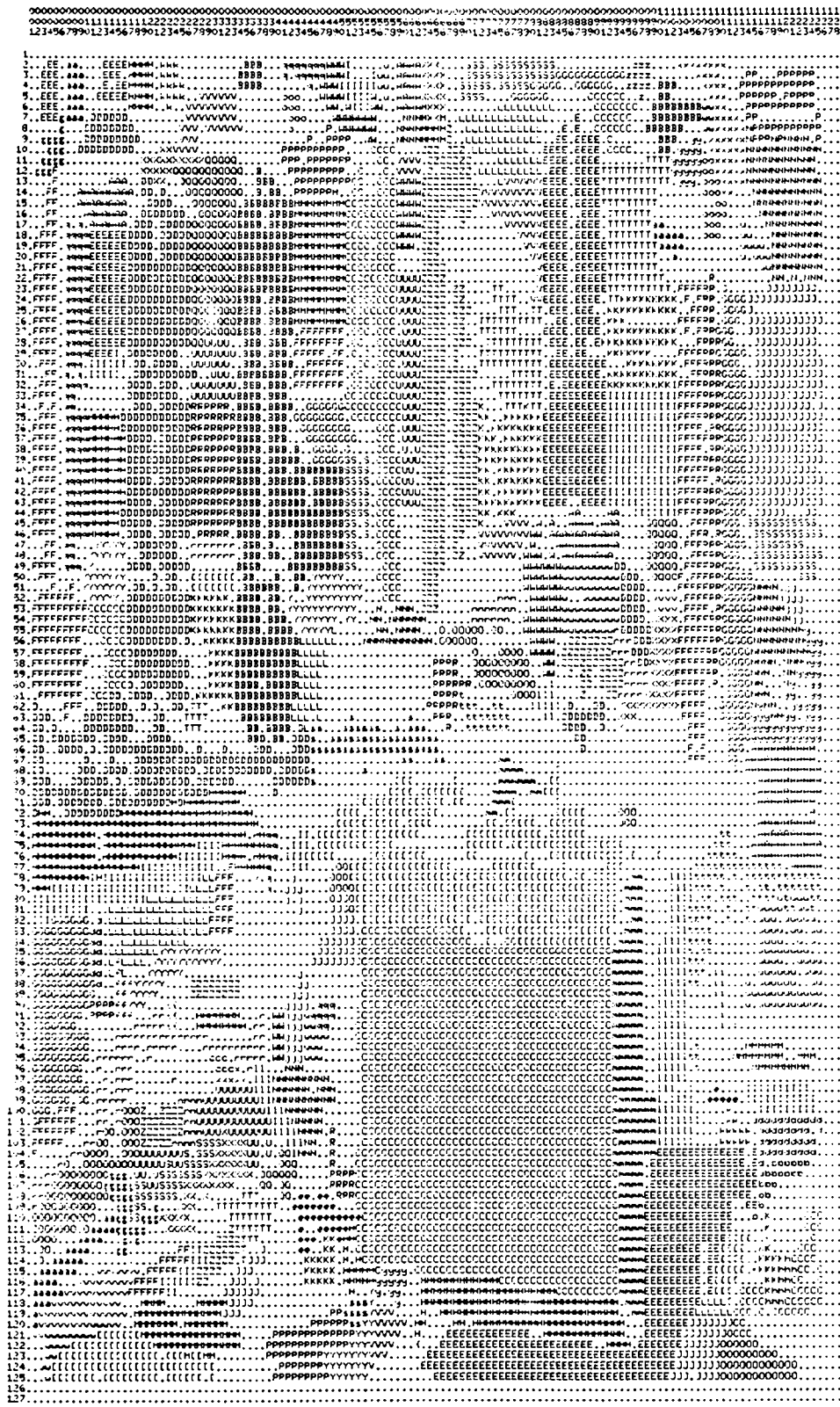


FIG. 3. right, matched labels after "cut and transplant", transplanted and unmatched labels are negative. D=7 D=2.

LABEL PRINTING

0	1	2	3	4	5	6	7	8	9
10	11	12	13	14	15	16	17	18	19
20	21	22	23	24	25	26	27	28	29
30	31	32	33	34	35	36	37	38	39
40	41	42	43	44	45	46	47	48	49
50	51	52	53	54	55	56	57	58	59
60	61	62	63	64	65	66	67	68	69
70	71	72	73	74	75	76	77	78	79
80	81	82	83	84	85	86	87	88	89
90	91	92	93	94	95	96	97	98	99

File: In=84.Fig

Fig.: 3.3.5-2b.

3.3.6 Clip

The purpose of the so-called "clipping" process is to try to give an "even start" in the analytic relaxation process to each pair of matched facets. Two alternatives were considered:

A) Centre of gravity based clipping.

- a) For any pair of matched facets in $Lm1(i,j)$ and $Rm1(i,j)$ compute their individual centers of gravity.
- b) Shift one of the facets (of a matched pair) so that the centers of gravity in $Lm1(i,j)$ and $Rm1(i,j)$ overlap.
- c) "And" the shifted facets. This reduces the facet pair to the same size in the new matched facet images $Lm2(i,j)$ and $Rm2(i,j)$.

B) Correlation based clipping (not completed).

- a) Correlate the pair of matched facets in $Lm1(i,j)$ and $Rm1(i,j)$ to find the shift d_i and d_j for the best match. (The d_i and d_j lists are directly available from the correlation based matching.)
- b) Shift one of the facets (of the matched pair) as given by d_i and d_j for best overlap.
- c) "And" the shifted facets (as in (A)). This reduces the facet pair to the same size in the new matched facet images $Lm2(i,j)$ and $Rm2(i,j)$.

A block diagram of the process is in Figure 3.3.6-1.

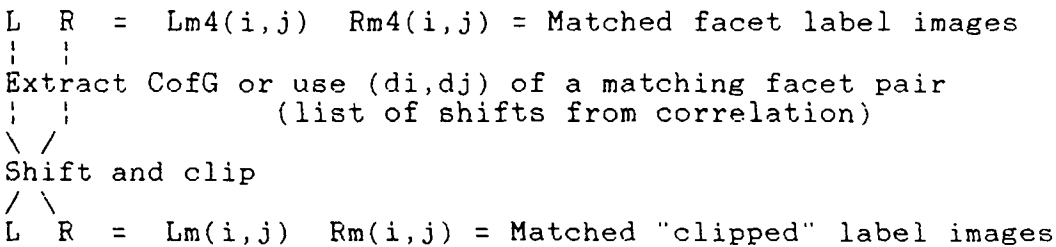


Figure 3.3.6-1: The label "clipping" process.

The "allowable displacements" between the two facets were limited by a "box" $D_i * D_j$, as in previous cases. Some results are shown in Figures 3.3.6-2a and -2b. The matching facet pairs are reduced to the same sizes in the $L_m(i,j)$ and $R_m(i,j)$ images. This is thought to be "the best that can be done" for "balancing" the subsequent relaxation process, see section 3.4. Note that the clipping process (as it is at present) also created holes within facets. For example, if a region F_l in the left image has a hole but the corresponding region F_r in the right image does not, then the "and-ing" process will make a hole in F_r (a region of zero labels). If the shift between F_l and F_r is less than the "diameter" of the hole then some of the 0-valued regions will coincide, see further on under "gaps".

The matching pairs of facets are of the same size in the $L_m(i,j)$ and $R_m(i,j)$ images. Geometrically such images are impossible. However, when viewed stereoscopically, a rather strange scene is seen (but it may take some initial effort to "integrate" these images visually). Consequently, several versions of stereo pairs were printed.

Figures 3.3.6-3a and -3b show the matched facets (F_l, F_r) printed with averaged grays ($Gray = (Gray(F_l) + Gray(F_r)) / 2$). The unmatched facets were printed in average gray over both images. In Figures 3.3.6-4a and -4b the facets were "painted" such that adjacent facets have different gray levels. Five gray levels ("colours") were needed.

Some rather interesting stereo effects were visible in the original prints but the copying procedures may have reduced all these efforts to nought. Thus, in words, the shifted matched facets create a scene that generally appears correct but there are gaps between the facets. This creates an "exploded" view of the 3D scene. In addition, there are "freely hanging" facets which are "in front of" the room scene. Only the unmatched facets create visual disturbances. This creates interesting implications for the subsequent reconstruction:

1. The initially matched (and clipped) facets are likely to be immediately reconstructible since they "appear correct" in the 3D scene. However, since these facets are (mostly) "exploded and hanging in free space" the reconstruction will have gaps between the facets and, of course, the interiors of the facets are empty.
2. The facets that "hang in front of the 3D scene" do not integrate with the rest of the scene. They act to "obscure" parts of the background.
3. The unmatched facets create visual disturbances since they "want to integrate" with the scene.

The full implications of the above observations are not entirely comprehended at the moment. It is unlikely that any visual illusions are involved, but the author has seen "too much" of this particular stereo pair to be a reliable observer. Experiments on colleagues have indicated that they are mostly unfamiliar with stereo viewing. After some "training" they see some stereo effect to a certain degree in some images but "independent confirmation" has been hard to obtain with the existing image quality and equipment.

The obvious implications are:

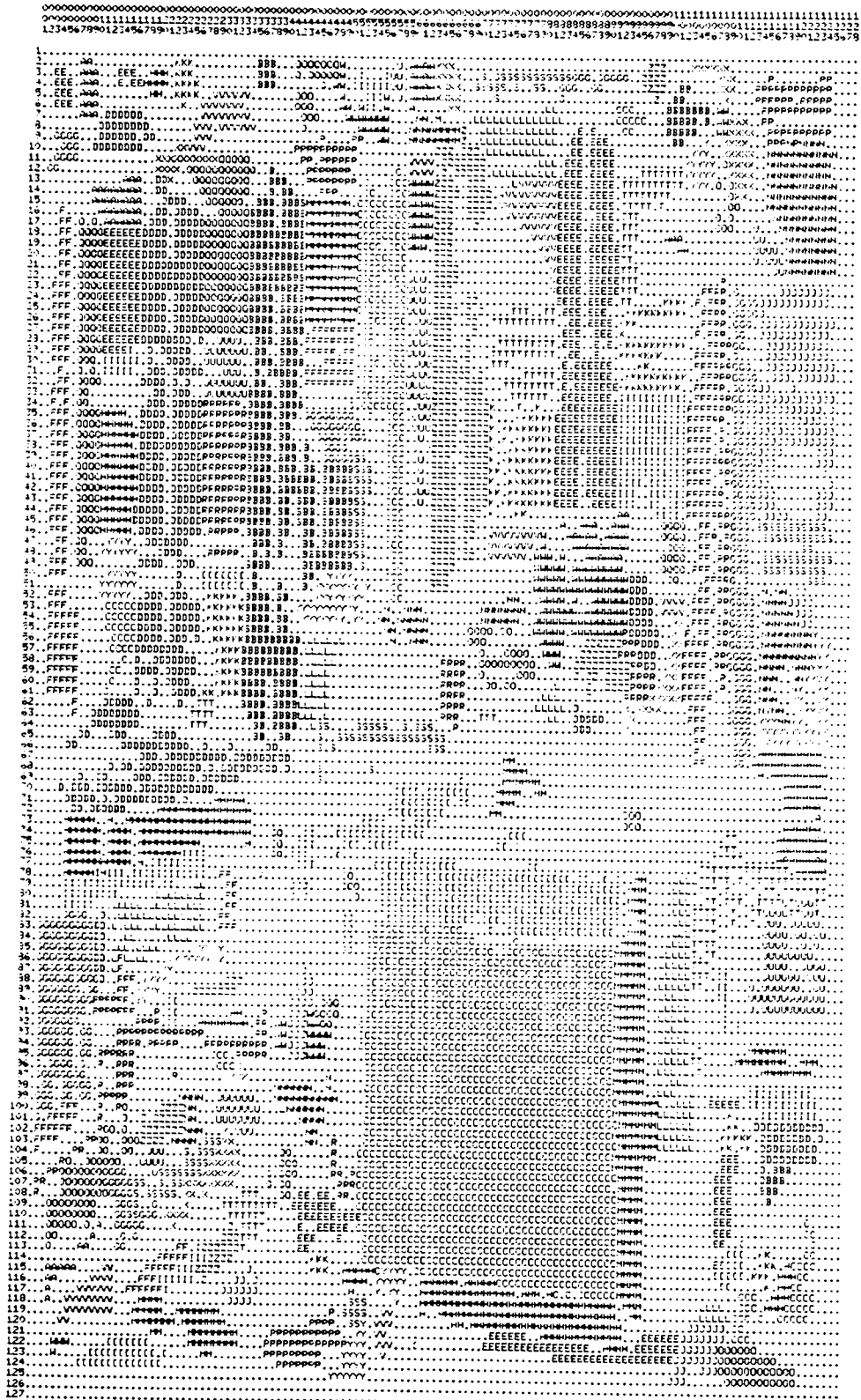
- a) The matched facets should be "expanded" by some form of "acclomerative pixel acquisition method" (acclomerative merging, snakes, etc.). In the present case this acclomeration is done via analytic approximation and relaxation, see section 3.4.
- b) The facets that "hang in front of the scene" (case 2) should also "acquire a parallax" after the processing indicated in (a).
- c) The unmatched facets should either be eliminated or not allowed to participate in the process (a). The unmatched facets are most likely to be errors but, of course, anything seen with one "eye" only will be unmatched.

Figure titles:

Figures 3.3.6-2a and -2b: The matched facets after clipping.
(Files: Inr087.fig and Inr088.fig)

Figures 3.3.6-3a and -3b: The matched facets (F1,Fr) printed with averaged grays ($\text{Gray} = (\text{Gray}(F1) + \text{Gray}(Fr)) / 2$). The unmatched facets are printed in average gray over both images.

Figures 3.3.6-4a and -4b: The facets "painted" in five "colours" such that adjacent facets have different gray levels.

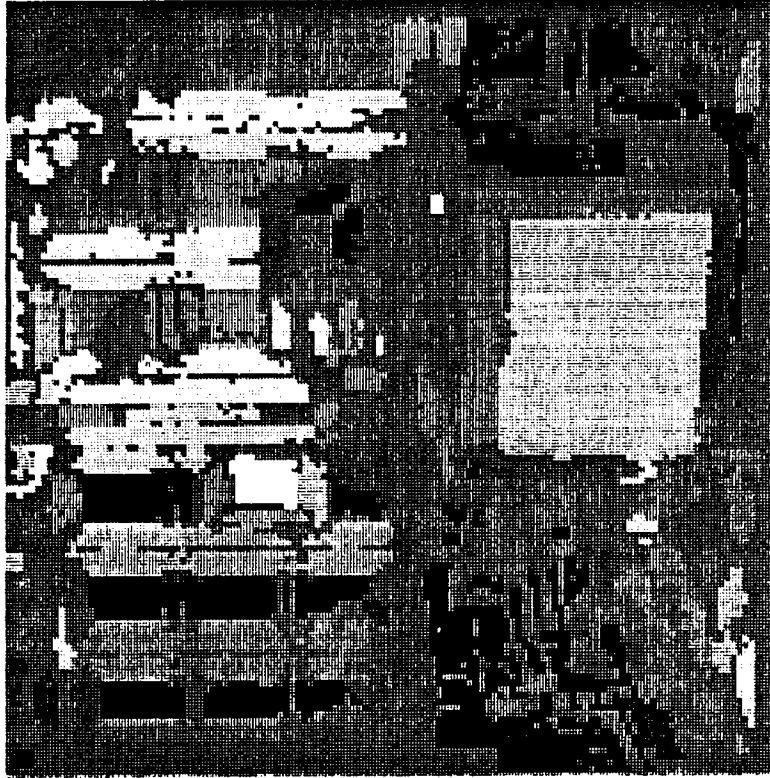


REL(I,J), right, matches, cut, translated, and flipped facet labels, Di=1, Dj=2.

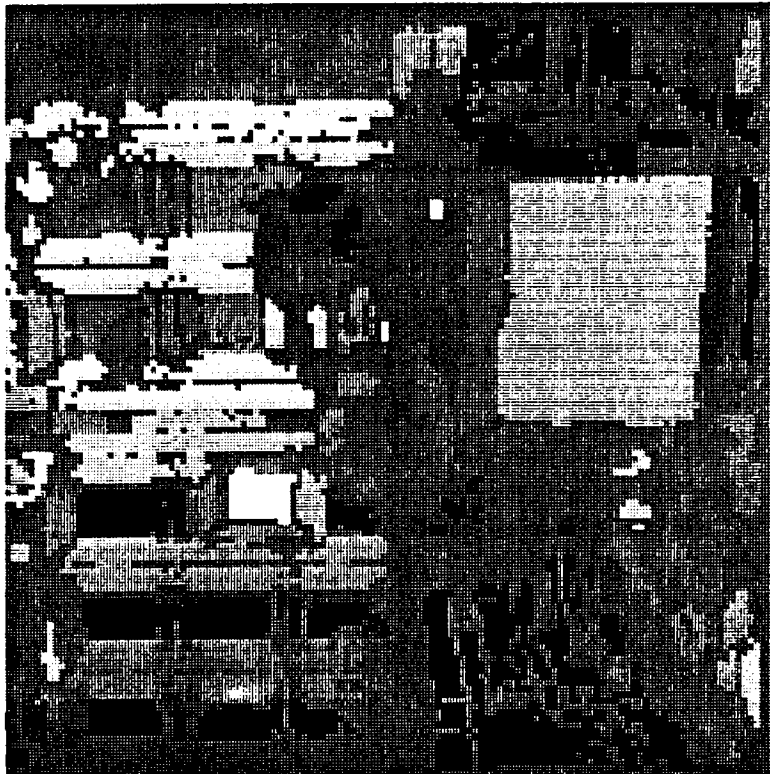
LABEL PRINTING

1	2	3	4	5	6	7
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29	30	31	32	33	34	35
36	37	38	39	40	41	42
43	44	45	46	47	48	49
50	51	52	53	54	55	56
57	58	59	60	61	62	63
64	65	66	67	68	69	70
71	72	73	74	75	76	77
78	79	80	81	82	83	84
85	86	87	88	89	90	91
92	93	94	95	96	97	98
99	100	101	102	103	104	105
106	107	108	109	110	111	112
113	114	115	116	117	118	119
120	121	122	123	124	125	126
127						

File: Inv-088.fig

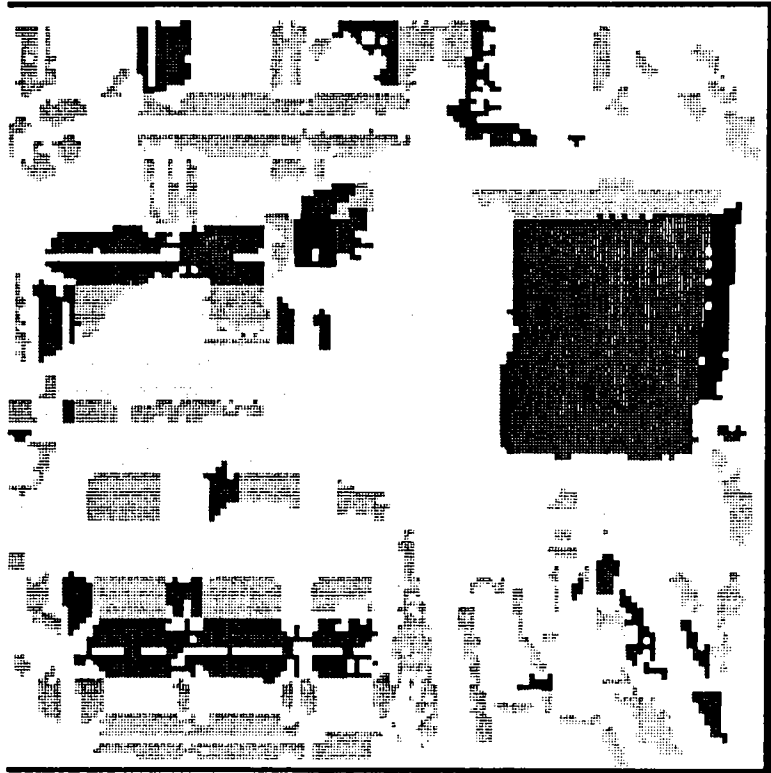


if8flgmc6.rgr

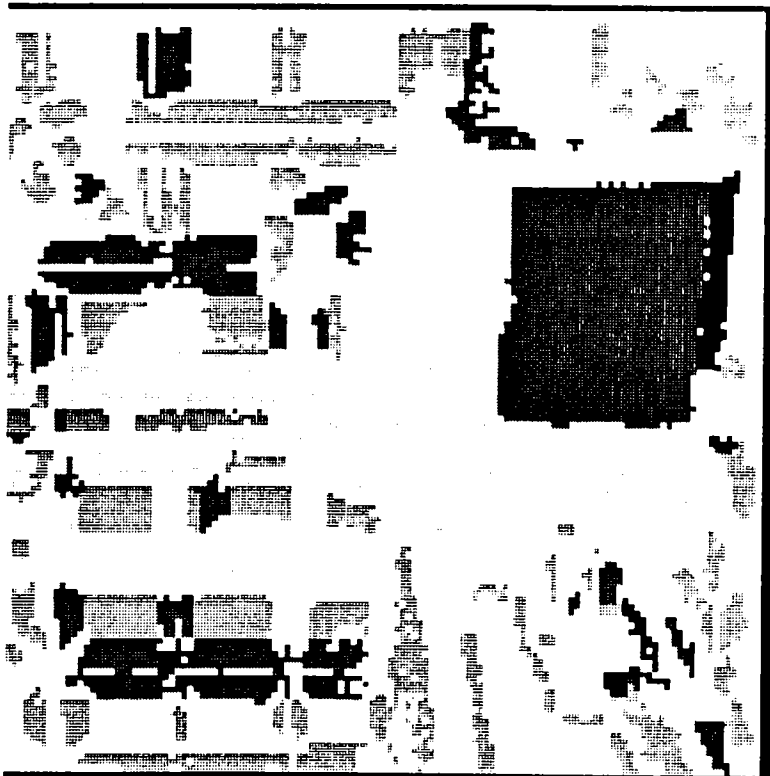


if8flgmc6.lgr

Fig.: 3.3.6-3a.

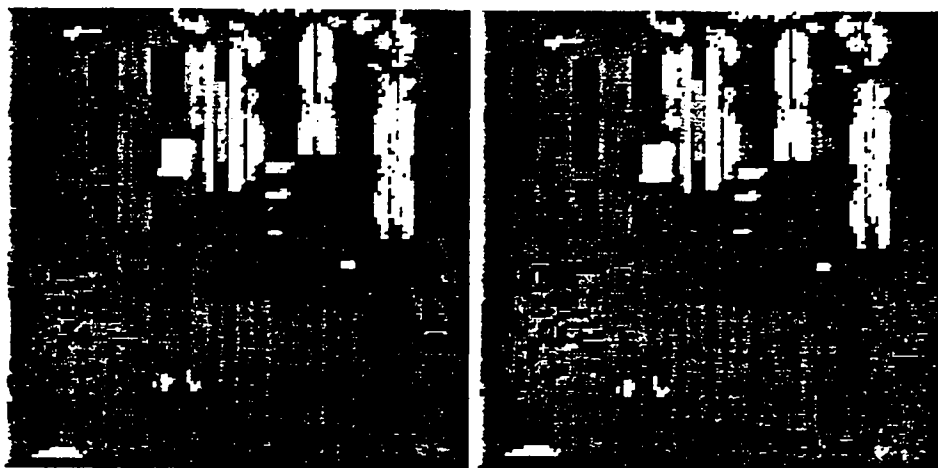


if8flgmc5.rgr



if8flgmc5.lgr

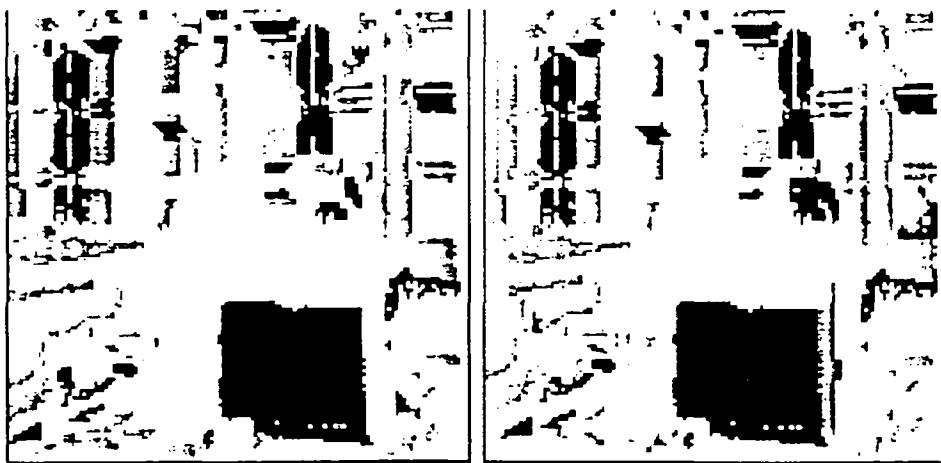
Fig.: 3.3.6-4a.



if8flgmc6.lgr

if8flgmc6.rgr

Fig.: 3.3.6-3b.



if8flgmc5.lgr

if8flgmc5.rgr

Fig.: 3.3.6-4b.

3.3.7 The gaps and the un-matchables

There are three kinds of facets that are "left over" from the previous processes, i.e.:

1. The "facets" (regions) where the labels are zero. There is "nothing" at these pixels.
2. The "unmatched" facets where there is a facet in one image of the stereo pair but not in the other. Even after the "transplant" operation there remain some unmatched facets.
3. The "unmatched" facets that could not be matched even though a "partner" may exist in the other image. Reason for the inability to match are that the two facets are "far apart" or one of the facets has become "fragmented".

In case of the unlabelled regions (case 1), if there is "nothing" in both images at the same pixel (i,j) then the "common characteristic" (similarity) is "nothing" (which is information about similarity) and this can be used to create a region. The process is as follows:

i) Create regions of pixels indicating "nothing". i.e.,

```
If((L(i,j).Eq.0).And.(R(i,j).Eq.0)) Then B(i,j) = 1
Else B(i,j) = 0
Endif
```

where $B(i,j)$ is a binary image, ".Eq." means "is equal to", and ".And." means logical "And". Since the previous processes may have reduced the sizes of the existing regions, this operation can create a binary image $B(i,j)$ consisting of "blobs connected by strings". This brings the process directly into the domain of binary image processing.

ii) Erode $B(i,j)$ to "disconnect" the "blobs".

iii) Label the "blobs" but continue the label numbers from highest label in $Lm(i,j)$ and/or $Rm(i,j)$.

iv) Put these regions "back" into the $Lm(i,j)$ and $Rm(i,j)$ images.

At least now the "larger" unlabelled regions have become matched and labelled. However, they may not perform too well in the subsequent analytic approximation and relaxation processes if the regions cannot be properly approximated by the simple functions used for the matched facets.

The unmatched facets (negatively labelled regions) that still remain after all these operations depend on the sequence in which the processes were carried out. In general several interpretations are possible:

- a) These regions are only visible to one "eye".
- b) The regions are caused by reflections (for example, distant bright objects) and are thus at different locations in the 3D scene than the region in which they appear in the image, or they are from transparent regions (a "hole" in the foreground scene). Frequently, the reflected region is greatly modified (broken up) by the reflecting surface.

3.3.8 Comments

The modifications that the matched facet labelled images $L_m(i,j)$ and $R_m(i,j)$ have been "subjected to" are summarized in Figure 3.3.8-1 and listed below. The purpose has been to illustrate some of the processing steps that could be used to "fix" the majority of errors. Hopefully, fewer errors have been introduced than have been fixed.

1. A certain number of misclassifications are removed by finding the "intersection" between unmatched facets.
2. Matched facets that are "too different in size" are "cut" to "reasonable size".
3. Since "cutting" leaves blank regions these are "filled in", if possible, by transplanting labels from the other image.
4. All the facets are then "clipped" to the same size in order to give different facets an "even chance" during the subsequent operations (analytic approximation and relaxation, see section 3.4).
5. The remaining "unknown" regions are then also given labels in order to allow them to "defend their territories" during subsequent processes.
6. The facet regions are next approximated analytically and the regions are corrected by analytic relaxation, see section 3.4.

The basic question is, how many of these processes are logically justifiable, but most of all, how do they perform in practice?

The correction of misclassifications may be looked upon as and extension or a correction to the label matching process and justified in the same terms as the initial matching. The cutting process can only be justified in probabilistic terms, i.e., it is more likely than not that the corresponding facets in the L and R images are approximately of the same size. Some exceptions were mentioned earlier. The transplantation process may also be justified in probabilistic terms since, if an unmatched region is

"large enough", it is very likely to contain facets that are similar to those visible in the other image, provided there are any. The subsequent clipping process is not logically justifiable except in terms of processes that follow. In the present case the subsequent process is analytic relaxation that is expected to "grow" the facets back to their proper sizes in L and R. Hence, the clipping process is only justified as "a method of giving an 'equal chance' to all the facets" in the subsequent analytic relaxation. The filling of the unknown regions in both images by "artificial facets" may also be viewed as a method of giving an "equal chance" to such regions in the analytic relaxation process. An alternative would be to reclassify the image for only these regions (see Chapter 2).

The most dubious processes are cutting and transplanting. Besides their ad hoc nature, the dilemma with the "cutting" and "transplanting" processes is:

- i) Cutting is too "rough", it either leaves or removes a large area of a facet and leaves nothing in its place, depending on the D_i and D_j . The logical "upper limit" for D_i is the whole i -dimension of the image (in the vicinity of the epipolar line).
- ii) Transplanting puts "something" into the vacated area but this "something" is only a part or a copy of the facets that were in the other image. Consequently, unless the facets are from "far away places", the "transplants" are in the wrong positions.

The "cutting" and "transplanting" processes were attempting to deal with unequal segmentation in the left (L) and right (R) images. An alternate approach may be envisaged, which is best explained in connection with a simple example. Consider a region in the image which in L has been split into two facets F_{11} and F_{12} , while the region is only one facet F_r in R. Of course, the intersection of $F_{11}+F_{12}$ with F_r is not precise (+ indicates "union"). Both F_{11} and F_{12} "lay claims" on F_r . The amount of these "claims" may be the intersection areas $F_{11}*F_r$ and $F_{12}*F_r$ (* indicates "intersection"). At present the facet with majority wins. More subtle processes are possible but their implications have not been thought out in detail.

In the final analysis the important aspect is how these operations behave in practice on images for which such operations are appropriate. Experimental results tended to show that "cutting" and "transplanting" were largely justified but there were some facets that remained the same in L and R after "matched" analytic relaxation. The facets after "cutting" were still rather different in size since D_i and D_j were rather "liberal". Consequently, the analytic parameters for the two "matching" facets were not really representing the situation in a satisfactory manner and the meaning of the matched analytic relaxation become dubious. Both of these problems were "alleviated" by including

the clipping process. The experiments so far are insufficient to either confirm or reject these as viable methods.

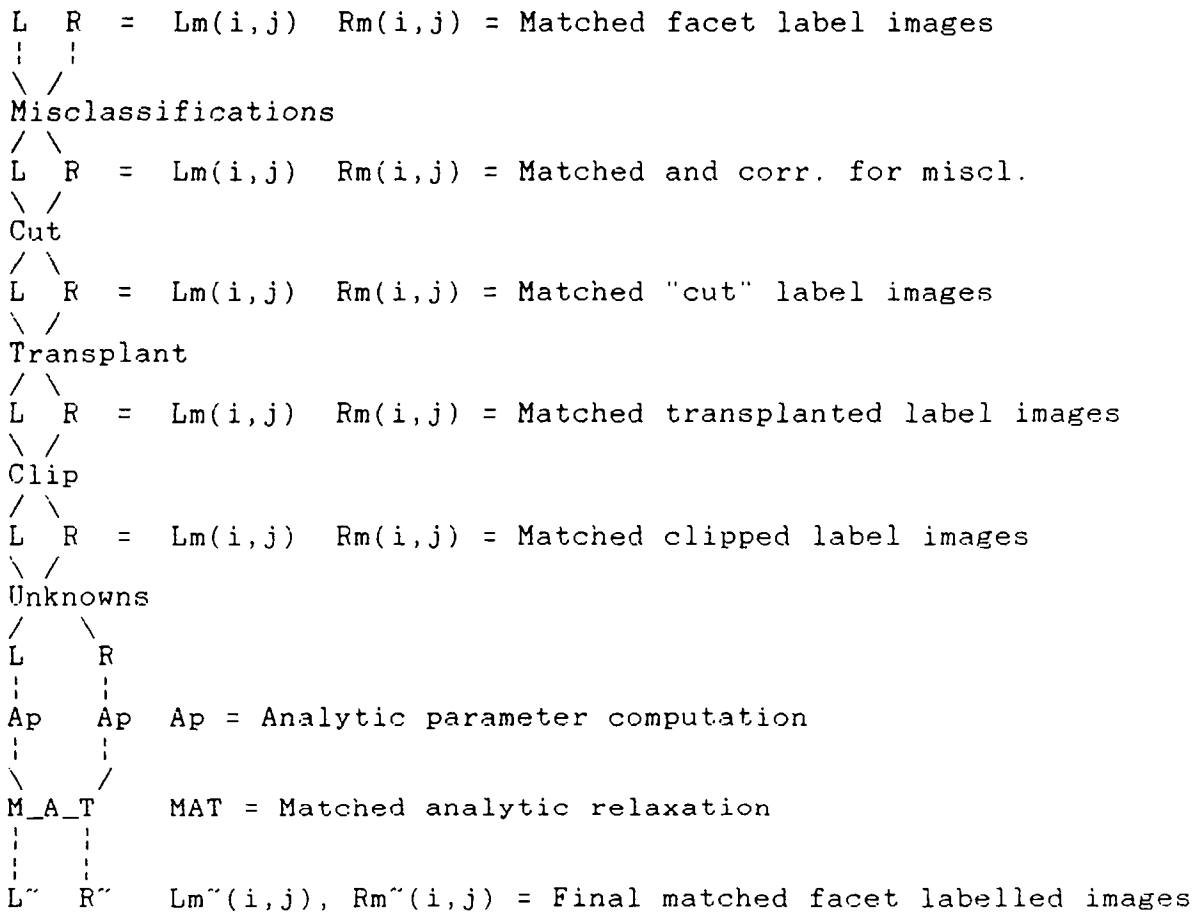


Figure 3.3.8-1: A sketch of the "fixing" processes.

3.4 Analytic methods

The analytic approximation of the gray levels of each facet is intended to serve two purposes which have been expressed here in ordinary words and concepts:

1. The approximation is an "umbrella" or "cover" for the facet which has been "defined" by all the presently known pixels that belong to the facet. The presently chosen second order polynomial approximation is a very crude model of light distribution on a facet.
2. This "umbrella" serves to "defend the interests" of the pixels belonging to the facet and also serves as the tool for acquiring new "like-minded" pixels during the analytic relaxation.

Hence, the approximation and subsequent analytic relaxation may be likened to a military government of a country in its role as the instrument of the country's coherence and during times of conflict the defender of the country's territorial integrity and the tool for its expansionist interests. The reader may translate this into mathematical symbols!

The process was intentionally expressed in "clear" language in order to explain some of the problems that are complicating the analytic approximation and relaxation strategy.

- a) In order to have a balanced relaxation (a "fair fight" between the facets) the initial approximations for all the facets should be as "fair" as possible. In the absence of any better method, the supposedly paired facets were made of equal size, since it was noticed in prior experiments that unequal starting conditions were undesirable.
- b) All regions in the images should have analytic expressions in order to properly participate in the analytic relaxation. If a facet was bordering an unlabelled region (an "undefended" region without an analytic approximation for it) then the labelled facets tended to expand into it during relaxation since the only constraint was the limit on the overall error term between the analytically computed and actual data.
- c) The unmatched but labelled facets created a dilemma during relaxation.

If the unmatched facet was a "true" facet representing a region that was only visible in one image then, of course, this facet should "defend its territory" during relaxation. However, if the unmatched facet was false (a part of a matched region split off in one image but not in the other image), then the relaxation should cause the false facet to be absorbed by the region from which it was split off originally. In general, this does not happen since the analytic approximation of a partial

facet is usually rather different from that of a whole facet. Consequently, the "false" facet was able to "defend itself" and the desired correction did not take place. An alternative is to simply ignore such facets, i.e., the unmatched facets are treated as "unlabelled", which contradicts (b). Hence, considerable effort was spent in order to try to "match everything" before the analytic process was started, see previous section.

- d) There is no such thing as "two falsely matched facets". Expressed in other terms, if a match has been found between two facets then, at the present stage of processing, there is no other information available to detect whether or not this match is correct. After the reconstruction stage it may be possible to detect some contradictions in the 3D scene, for example, "two visible facets obscure each other", but these aspects of the problem have not been studied.
- e) If a region was "incorrectly" segmented (from human point of view) but the "incorrect" segments were matched correctly, then they are correct segments (despite human opinion). According to the similarity principle "recognition" or "understanding" of the facets is neither required nor needed.

3.4.1 The coupling

The greatest conceptual problem is created by the need to have some form of "coupling" or link between each of the left and right matched facets (Flk,Frk) during relaxation.

Let the k 'th left facet be indicated by Flk (where "l" is for "left" image), let the k 'th right facet be Frk ("r" indicates "right" image), and let F symbolize a facet. For a matched facet pair (Flk,Frk) the index k is the same in the labelled images $L_m(i,j)$ and $R_m(i,j)$, and the analytic parameters are obtained from tables addressed by k (the left image parameter tables are $T_l(k,n)$ and the right image ones are $T_r(k,n)$, for $n = 1,2,\dots$ for the degree of the analytic approximation). Thus, the data management in the relaxation program is very simple. A block diagram of the process is shown in Figure 3.4.1-1. Some details of the relaxation process were given in Chapter 2.

The basic problem with the "coupling" between the left and right matched facets (Flk,Frk) during relaxation is that the procedure is analytic (and the author is not aware of any other form of mathematics which is both precise and general at the same time)! Even though the approximated values can be computed at any pixel (i_l, j_l) in the left image L and at any pixel (i_r, j_r) in the right image R, for a "coupling" to exist between L and R the pixels (i_l, j_l) and (i_r, j_r) have to be matched in the sense of stereo matching. Thus, ideally, the L and R images will have to be matched on the pixel level before the relaxation technique can

be correctly applied! [A similar problem in pattern recognition versus segmentation is well-known: An image has to be segmented correctly before the objects can be recognized but the objects have to be recognized before the image can be segmented correctly! In ordinary terms this is known as the "chicken and egg" puzzle, i.e., which came first, the chicken or the egg?]

Several strategies are possible in order to try to bypass this problem without resorting to "models". Models restrict the solution to a very small domain and create "combinatorial explosions" during matching.

1. Improve the segmentation independently in the two images. Attempts to simply try to "improve" the segmentation using the existing set of features may help but become futile "in the limit" since the images do differ. The alternative is to compute more features for each pixel and to improve both the classification and the matching techniques. In the limit this would resemble direct pixel feature matching and eliminates most of the intervening steps.
2. Try to solve the problem with approximations and iterations. As with any iteration technique, if the starting point is too far from the real solution then either no solution is obtained and the iterations do not converge, or some form of false convergence is obtained. Active or "live" vision is able to "look at a funny facet" and this process immediately produces a better match and corrections to the iterations since the centre of fixation is on the facet. In still-life vision the wrong solution may be just another way of seeing the 3d scene. In the absence of additional "knowledge" how is one to verify that this is indeed the wrong solution? Our own visual illusions tend to confirm this opinion.

The approximate matches from (2) should be very helpful in pixel feature matching (1) since the problem is now confined to the vicinity of each matched facet. As a final step pixel level matching is required anyway, whatever the intermediate methods in the hierarchy may have been.

3. Approximate reconstruction. If the 3D scene (depth map or image) is reconstructed from the presently matched information (facets, edges, pixels) then, of course, the results depend on the matches. In other words, the matching information can only be obtained from the left and right images and a partially reconstructed 3D scene is no better than the matches (in L and R) that were used. It is not clear at the moments whether a partially reconstructed 3D scene would be of any help in resolving the paradox without using additional information (more features, better analytic procedures, or even "extraneous" knowledge, i.e., models).
4. Apply constraints that are appropriate to the accuracy at which the matching is carried out, see next section.

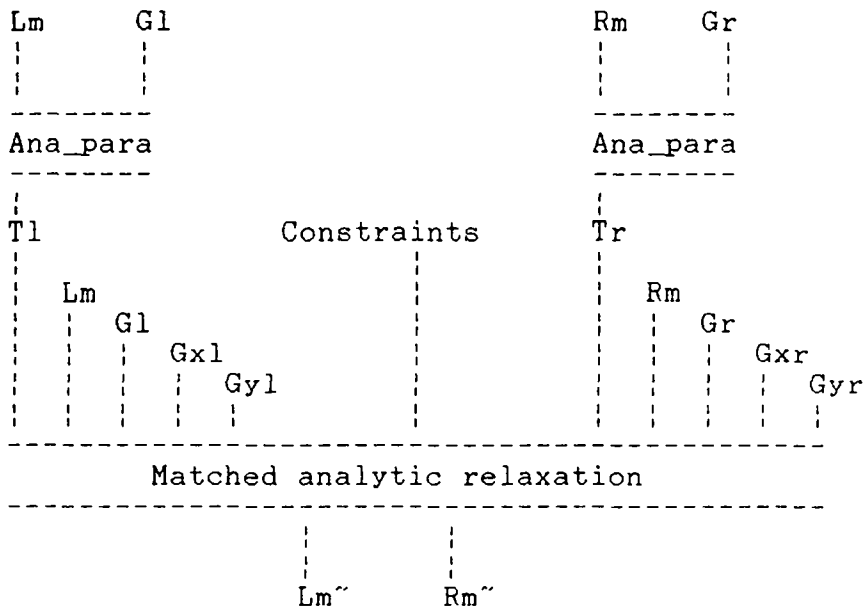


Figure 3.4.1-1: A sketch of the analytic process.

Lm, Rm = Left and right matched facet label images.
 G1, Gr = Left and right gray level images.
 Ana_para = Analytic parameters for facets.
 Tl, Tr = Tabulations of analytic parameters.
 Gx1, Gy1 = X and y components of gradient of G1.
 Gxr, Gyr = X and y components of gradient of Gr.
 Lm'', Rm'' = Analytically corrected matched facet label images.

3.4.2 Constraints on relaxation

There is no lack of approximate constraints that could be placed on the relaxation process in order to "couple" the two images of the stereo pair. Given any two facets A and B in the L and R images, some of the constraints may be the following (which have been listed without any particular order of preference):

1. Edge within a "tolerance box" $D_i * D_j$.
2. Presence of an edge for facet A.
3. Presence of an edge between facets A and B.
4. Adjacency of A to B.
5. Edge with a certain slope tolerance.
6. Edge with a certain gradient tolerance.
7. Edge with a certain curvature tolerance.
8. Local maximum of edge function at edge.
9. Combinations of above including local extrema.

- a) The presence of an edge (border) for facet A in both L and R, while ignoring the other facet (B).

This simply requires the facet to be matched (A is both in the left 'L' and right 'R' images). The relaxation is allowed to progress independently in L and R provided the edge of A stays within a tolerance box $D_i * D_j$ centered on the epipolar line. The constraint on the "progress" of relaxation is obtained independently from the L and R images. Pixel-level matching between L and R is not required. The method is similar to the "independent relaxation" in Chapter 2, except that only the matched facets are allowed to "relax" provided they stay within $D_i * D_j$.

- b) Improvements on the first method (a) consist of requiring increasingly precise matches between the borders of facet A. For example:

- i) Gray level and gray level gradient were introduced separately in Chapter 2. Here one may require a certain amount of similarity between these in both images at the same time.

- ii) The local maximum of the "edge function" E_f (see Chapter 2) should be reached in both images, provided $E_f >$ some tolerance. $E_f(i, j)$ was defined as a function with "amplitude" rather than "on-off" for exactly this reason since it allows "hill climbing". A "smeared" on-off type edge function can serve in the same way but most likely with poorer results.

- iii) The slope of the edge at the border of A in L and R should be approximately the same. For more or less straight edges this constraint does not require much pixel-level correspondence between L and R.

- iv) The curvature of the edge at the border of A in L and R should be "similar". This places more stringent pixel-level

match requirements between L and R where the curvatures are high (corners).

It may be worth noticing that these methods of constraining the relaxation include a certain degree of "edge matching" between for a facet A in the left (L) and right (R) images. This information may be saved for subsequent use. Edge matching will not be discussed further due to lack of time.

- c) Various more or less "artificial" constraints may be used, such as:
- i) There has to be a contact between the matched facets A and B in both images.
 - ii) The displacement of the boundary where the relaxation takes place in L and R is to stay within a "reasonable tolerance limit" $D_i * D_j$.
 - iii) Many combinations of methods (a), (b), and (c).

In the present experiments the constraints only consisted of the following:

1. Gray level and gray level gradient constraints applied individually as in independent relaxation.
2. The requirement of a contact between two matched facets (c-i) was optional but abandoned due to many unlabelled pixels.
3. The boundaries of a matched facet in L and R had to stay within a "box" (c-ii) around the epipolar line.

The reasons were twofold. Firstly, these were initial experiments before a "deeper comprehension" of what was involved and, secondly, the programs started to overwhelm the computational resources.

The facet labels from "matched classification", after various corrections, are shown in Figures 3.3.6-2a and -2b, together with several stereo pairs. As was described earlier, there were two kinds of facets, i.e.:

- a) Those that were approximately in the "right position" when viewed stereoscopically (but there were "cracks" between them in the 3D scene).
- b) The facets that were "hanging in free space" in front of the 3D scene.

The resulting facet labels after relaxation are shown in Figures 3.4.2-1a and -1b. In these figures the "acquired pixels" have been indicated by lower case letters to show the "conquests"

made by the various labelled regions. For example, the pixels conquered by facet E are indicated by e. Pixels that had no label and could not be "conquered" are indicated by dots (.). The iteration was run for ten cycles using the same set of "raw" or initial facet parameters. The number of pixels conquered in each iteration is shown in Figure 3.4.2-1c. As seen, the relaxation process is stable in a "global" sense but, as will be seen soon, the small details in the images do not obey "stereo fidelity" requirements.

If the gray level image is recreated (computed) from the analytic values, with the gray levels for the "missing" (0-valued) labels replaced by the gray values from the original image, then the result looks "reasonable", see Figures 3.4.2-2a and -2b. In this figure "gross errors" were also replaced, i.e., if $|G_{\text{analytic}} - G_{\text{original}}| > 15$ then the original gray level pixel was used instead of the analytic value. The "surviving labels" are shown in Figures 3.4.2-3a and -3b. The stereo pair is shown in Figure 3.4.2-4a together with a stereo pair of the original gray level image for comparison, see Figure 3.4.2-4b.

The replacement of "missing pixels" (zero-valued labels and where the approximation is "too poor") is a valid operation, being just another form of "correction". However, the facets to which these "missing pixels" should be assigned, or the new facets created by this operation, have not yet been determined. Unfortunately, these results are deceptive since the "stereo effect" is rather subtle.

To investigate further, the relaxed labels in Figure 3.4.2-1a and -1b were used for analytically recreating different gray level images. The gray levels of the 0-valued labels was replaced by average gray for all the 0-valued labels in the images. The results deteriorated dramatically, see Figures 3.4.2-5a and -5b. In order to see the errors the gray levels for the "missing labels" are taken from the original images and the pixels with the labels are shown as "average gray", see Figures 3.4.2-6a and -6b. The corresponding stereo pairs are shown in Figures 3.4.2-7a and -7b.

The same experiment was repeated with the "corrected" labels, see Figures 3.4.2-3a and -3b. Figures 3.4.2-8a and -8b show the analytic gray levels with average gray for the unlabelled pixels. Figures 3.4.2-9a and -9b show the errors with average grays for the labelled pixels. The corresponding stereo pairs are in Figures 3.4.2-10a and -10b. The results did not change significantly.

Next, the minimum number of "colours" was assigned to each facet in the two images such that adjacent facets in the two images have different "colours" (Figures 3.4.2-1a and -1b with conquered labels set positive). Six colours were needed. The "coloured" facets are shown in Figures 3.4.2-11a and -11b with the top left corner of the adjacency matrix in -11c. A print is

shown in Figures 3.4.2-12a and -12b, and the stereo pair in Figure 3.4.2-14a. The facets were also shown with average gray per facet, see Figures 3.4.2-13a and -13b, with the stereo pair in Figure 3.4.2-14b.

The profusion of displays in various combinations represents an attempt to "pin down" the critical aspects of a stereo image. The analytic results were or could be corrected and/or displayed in several ways:

1. Whenever the analytic value exists, i.e., the pixel has a label then, if the error " $|G_{ana} - G_{orig}| > Limit$ " then eliminate this pixel (set the label to zero). This operation is equivalent to using another constraint on the output of analytic relaxation (a form of post-processing of analytic results).
2. A method that is similar to the first (1) but the labels are corrected only at the edges of the labelled regions. This has not yet been done.
3. A method that is similar to the first (1) or the second (2) but the analytic gray values for the pixels that are unlabelled or that "fail the test in 1" are replaced by the gray levels from the original image. The resultant gray level images were shown in Figures 3.4.2-2a and -2b, and as stereo pairs in Figures 3.4.2-4a. However, this method of presenting the results does not show the facets and prevent visual verification of facet shapes too see if the boundaries between the facets are approximately correct.
4. The adjacent facets are "painted" in different "colours". This was largely an attempt to bypass the weaknesses of the laser printer, but as seen, with little success.

All the results "point" in the same direction, i.e., the stereo appearance of these results is "disturbed" by the mismatching random "dots", by the differences in the details of the contours of the supposedly matching facets, and differences in adjacency relations between facets in the two images. Our stereo vision is greatly disturbed by relatively minor deviations between the facets and by inconsistent adjacency relations. The full implications remain to be investigated, especially for smooth contours without visible quantization effects, but the immediate conclusion is that pixel level correspondence is required (for highly spatially quantized contours and facets), and that adjacency relations cannot be violated (in a gross manner).

Figure titles:

Figures 3.4.2-1a, -1b, -1c: (1a) and (1b) show $Lm\tilde{(i,j)}$ and $Rm\tilde{(i,j)}$, the resulting facet labels after relaxation. The "acquired pixels" have been indicated by lower case letters to show the "conquests" made by the various labelled regions, i.e., the pixels conquered by facet E are indicated by e, etc. (1c) shows the number of pixels changed during relaxation (extract from history file). (Files: Inr093.fig, Inr094.fig, Inr092.fig).

Figures 3.4.2-2a, -2b: The analytically created gray level images after missing pixels and gross errors have been replaced.

Figures 3.4.2-3a, -3b: The "surviving" labels after analytic reconstruction with error limit of 15. (Files: Inr187.fig, Inr188.fig).

Figures 3.4.2-4a, -4b: (4a) A stereo pair of Figures 3.4.2-2a and -2b together with a stereo pair of the original gray level image for comparison (4b).

Figures 3.4.2-5a, -5b: The gray levels of the 0-valued labels replaced by average gray for all the 0-valued labels in the images.

Figures 3.4.2-6a, -6b: The errors when the gray levels for the "missing labels" are taken from the original images and the pixels with the labels are shown as "average gray". "Inverse" display compared to Fig. 3.4.2-5.

Figures 3.4.2-7a, -7b. The stereo pairs corresponding to Figs. 3.4.2-5 and -6.

Figures 3.4.2-8a, -8b: The analytic gray levels with average gray for the unlabelled pixels using the "corrected" labels in Figures 3.4.2-3a and -3b.

Figures 3.4.2-9a, -9b: The errors with average grays for the labelled pixels. The "inverse" to Fig. 3.4.2-8.

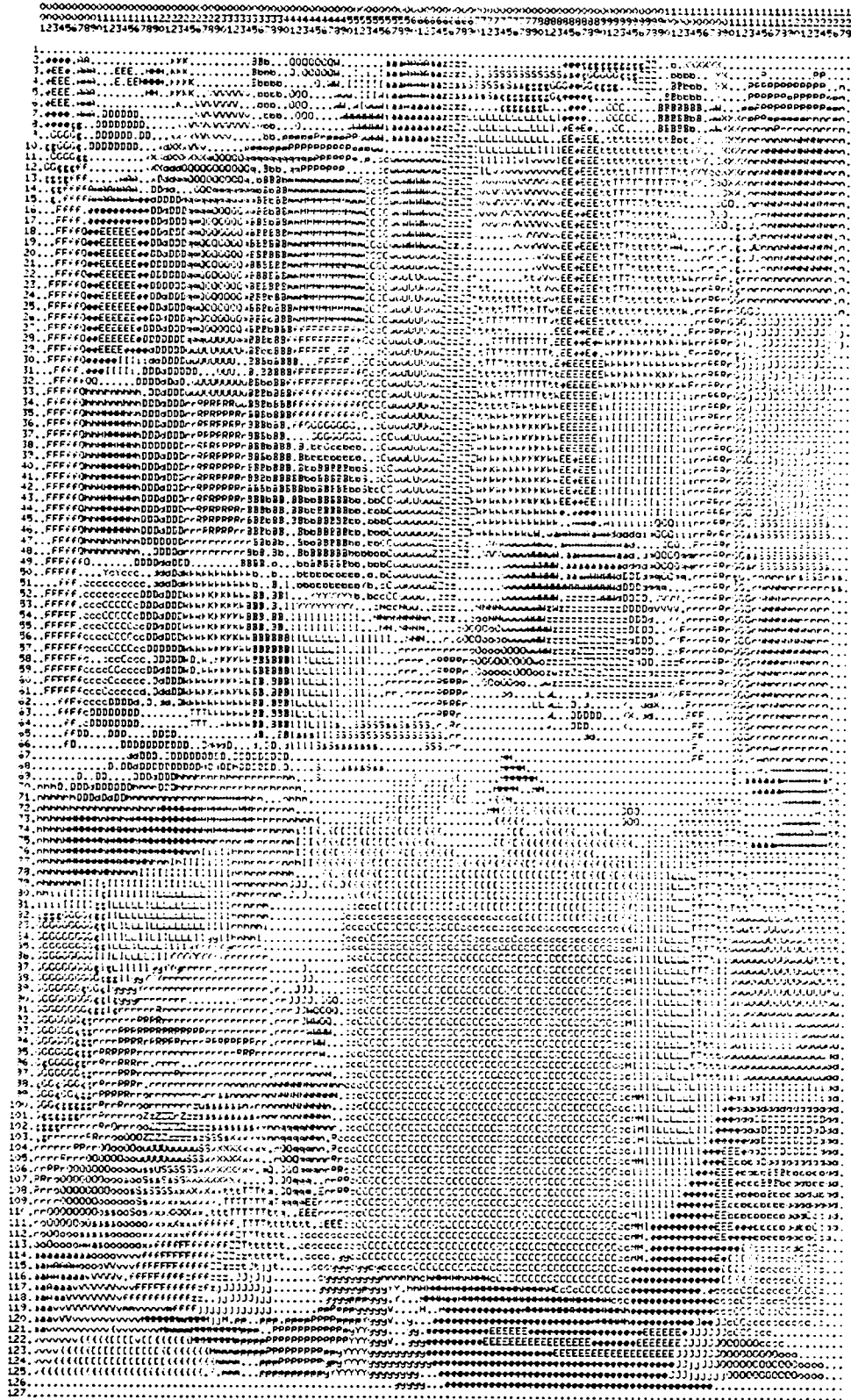
Figures 3.4.2-10a, -10b. The stereo pairs corresponding to Figs 3.4.2-8 and -9.

Figures 3.4.2-11a, -11b, -11c: (11a and 11b) "Coloured" facets from Figures 3.4.2-1a and -1b. The "colours" were assigned to each facet in the two images such that adjacent facets in the two images have different "colours". Six colours were needed. (11c) The top left corner of the adjacency matrix. (Files: Inr093b.fig, Inr094b.fig, Inr125a.fig)

Figures 3.4.2-12a, -12b: A print of the "colored" facets shown in Figures 3.4.2-11a and -11b.

Figures 3.4.2-13a, -13b: The facets shown with average gray per facet.

Figures 3.4.2-14a, -14b: Stereo pairs, (14a) for Figure 3.4.2-12 and (14b) for Figure 3.4.2-13.



Plot 1, J, right. Facet labels after 10 iterations in initial parameters.

```

LABEL PRINTING
  = 3   = 2   D = 3   E = 4   F = 5   G = 6   H = 7
  = 3   J = 3   K = 10   L = 11   M = 12   N = 13   O = 14   P = 15
  O = 16   R = 17   S = 18   T = 19   U = 20   V = 21   W = 22   X = 23
  Y = 24   Z = 25   a = 26   b = 27   c = 28   d = 29   e = 30   f = 31
  g = 32   h = 33   i = 34   j = 35   k = 36   l = 37   m = 38   n = 39
  o = 40   p = 41   q = 42   r = 43   s = 44   t = 45   u = 46   v = 47
  w = 48   x = 49   y = 50   z = 51   aa = 52   ab = 53   ac = 54   ad = 55
  ae = 56   af = 57   ag = 58   ah = 59   ai = 60   aj = 61   ak = 62   al = 63
  am = 64   an = 65   ao = 66   ap = 67   aq = 68   ar = 69   as = 70   at = 71
  au = 72   av = 73   aw = 74   ax = 75   ay = 76   az = 77   ba = 78   bb = 79
  bc = 80   bd = 81   be = 82   bf = 83   bg = 84   bh = 85   bi = 86   bj = 87
  bk = 88   bl = 89   bm = 90   bn = 91   bo = 92   bp = 93   bq = 94   br = 95
  bs = 96   bt = 97   bu = 98   bv = 99   bw = 100   bx = 101   by = 102   bz = 103
  ca = 104   cb = 105   cc = 106   cd = 107   ce = 108   cf = 109   cg = 110   ch = 111
  ci = 112   cj = 113   ck = 114   cl = 115   cm = 116   cn = 117   co = 118   cp = 119
  cq = 120   cr = 121   cs = 122   ct = 123   cu = 124   cv = 125   cw = 126   cx = 127
  cy = 128   cz = 129   da = 130   db = 131   dc = 132   dd = 133   de = 134   df = 135
  dg = 136   dh = 137   di = 138   dj = 139   dk = 140   dl = 141   dm = 142   dn = 143
  do = 144   dp = 145   dq = 146   dr = 147   ds = 148   dt = 149   du = 150   dv = 151
  dv = 152   dv = 153   dv = 154   dv = 155   dv = 156   dv = 157   dv = 158   dv = 159
  dv = 160   dv = 161   dv = 162   dv = 163   dv = 164   dv = 165   dv = 166   dv = 167
  dv = 168   dv = 169   dv = 170   dv = 171   dv = 172   dv = 173   dv = 174   dv = 175
  dv = 176   dv = 177   dv = 178   dv = 179   dv = 180   dv = 181   dv = 182   dv = 183
  dv = 184   dv = 185   dv = 186   dv = 187   dv = 188   dv = 189   dv = 190   dv = 191
  dv = 192   dv = 193   dv = 194   dv = 195   dv = 196   dv = 197   dv = 198   dv = 199
  dv = 200
  
```

Fig.: 3.4.2-1b.

File: In-094.fic

Analytic relax of matched facets.
 MN397: ANA RELAX OF MATCHED PAIRS
 I=AKL,IVAKR, IDIM,JDIM,KDIM = 55 56 129 127 590
 IDIM1,JDIM1 = 1 3
 JSTART, JEND, JSTP, IFEDBK, = 1 127 1 0
 IOPT, ICYCLE, IJECT, NSTEP1, NSTEP2, INDEF, AORDER = -1 10 4 7 2 0 2
 ITEST1, ITEST2, ITEST3, LBTEST = 0 0 0 0 0 0
 IPI, IPI1, IPI2 = 0 0 1
 EPLU, EPLU1, EPLU2 = 0.1000E+01 0.1000E+01
 WEIGHTC, WEIGHTD = 0.1000E+01 0.1000E+01
 EPLIM, EPLIMD = 0.25000E+02 0.15000E+02

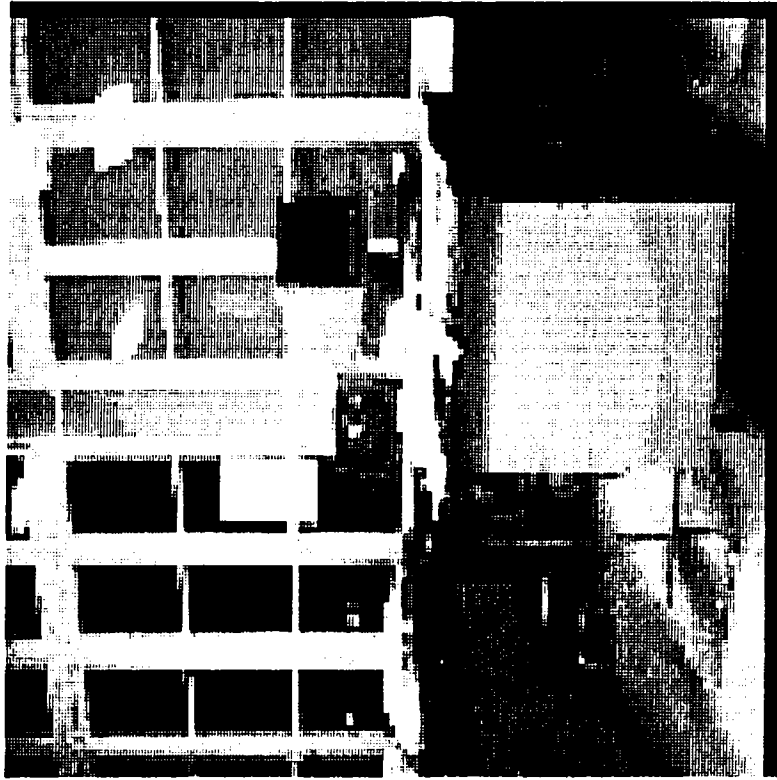
	NC	TOTL	TOTR	TOT
CYCCHNG,PIX: NC,TOTL,TOTR,TOT = 1	1754.0	2716.0	2449.0	
CYCCHNG,PIX: NC,TOTL,TOTR,TOT = 2	1672.0	1634.0	2204.0	
CYCCHNG,PIX: NC,TOTL,TOTR,TOT = 3	931.0	842.0	1362.0	
CYCCHNG,PIX: NC,TOTL,TOTR,TOT = 4	602.0	558.0	1160.0	
CYCCHNG,PIX: NC,TOTL,TOTR,TOT = 5	415.0	357.0	772.0	
CYCCHNG,PIX: NC,TOTL,TOTR,TOT = 6	277.0	270.0	547.0	
CYCCHNG,PIX: NC,TOTL,TOTR,TOT = 7	203.0	213.0	416.0	
CYCCHNG,PIX: NC,TOTL,TOTR,TOT = 8	141.0	139.0	281.0	
CYCCHNG,PIX: NC,TOTL,TOTR,TOT = 9	120.0	140.0	260.0	
CYCCHNG,PIX: NC,TOTL,TOTR,TOT = 10	37.0	119.0	207.0	

Extract of printout during analytic relaxation.
 Edited extract from history file: INNO31.FIG.

NC = Iteration count.
 TOTL = Number of pixels changed per iteration in left image.
 TOTR = Number of pixels changed per iteration in right image.
 TOT = Number of pixels changed per iteration in both images.
 WeightC = Weight for error between analytic and actual gray values.
 WeightD = Weight for error between analytic and actual gray value
 gradients.
 Eplim = Overall error limit. The formula is in Chapter 2.
 Istep1 = Di of the 'box'.
 Istep2 = Dj of the 'box'.
 WEIGHTC, WEIGHTD, EPLIM = 0.1000E+01 0.1000E+01 0.2500E+02
 NSTEP1, NSTEP2 = 2

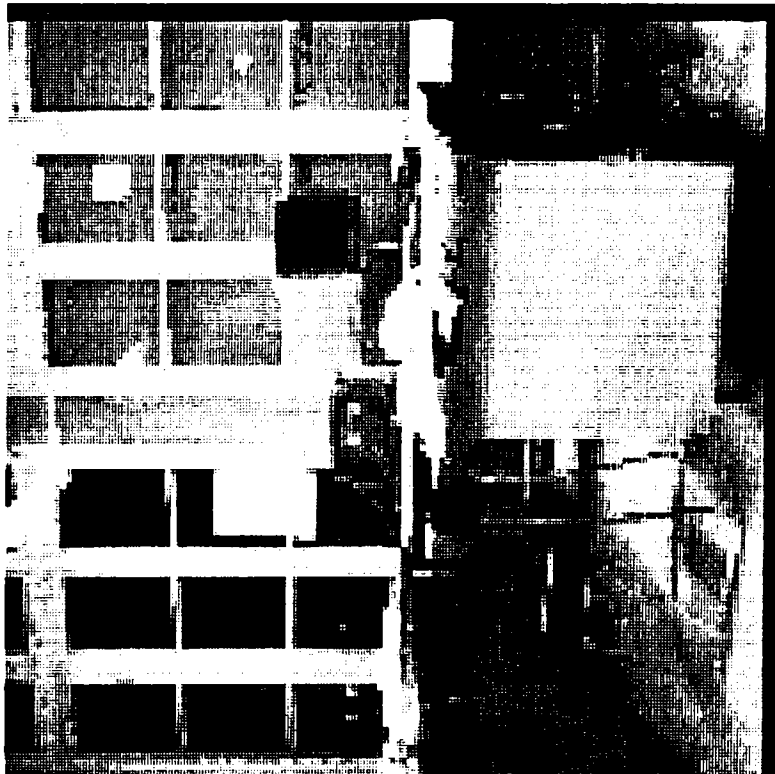
File: INNO31.FIG

Fig.: 3.4.2-1c.



if8gacmc9.rgr

Fig.: 3.4.2-2b.



if8gacmc9.lgr

Fig.: 3.4.2-2a.

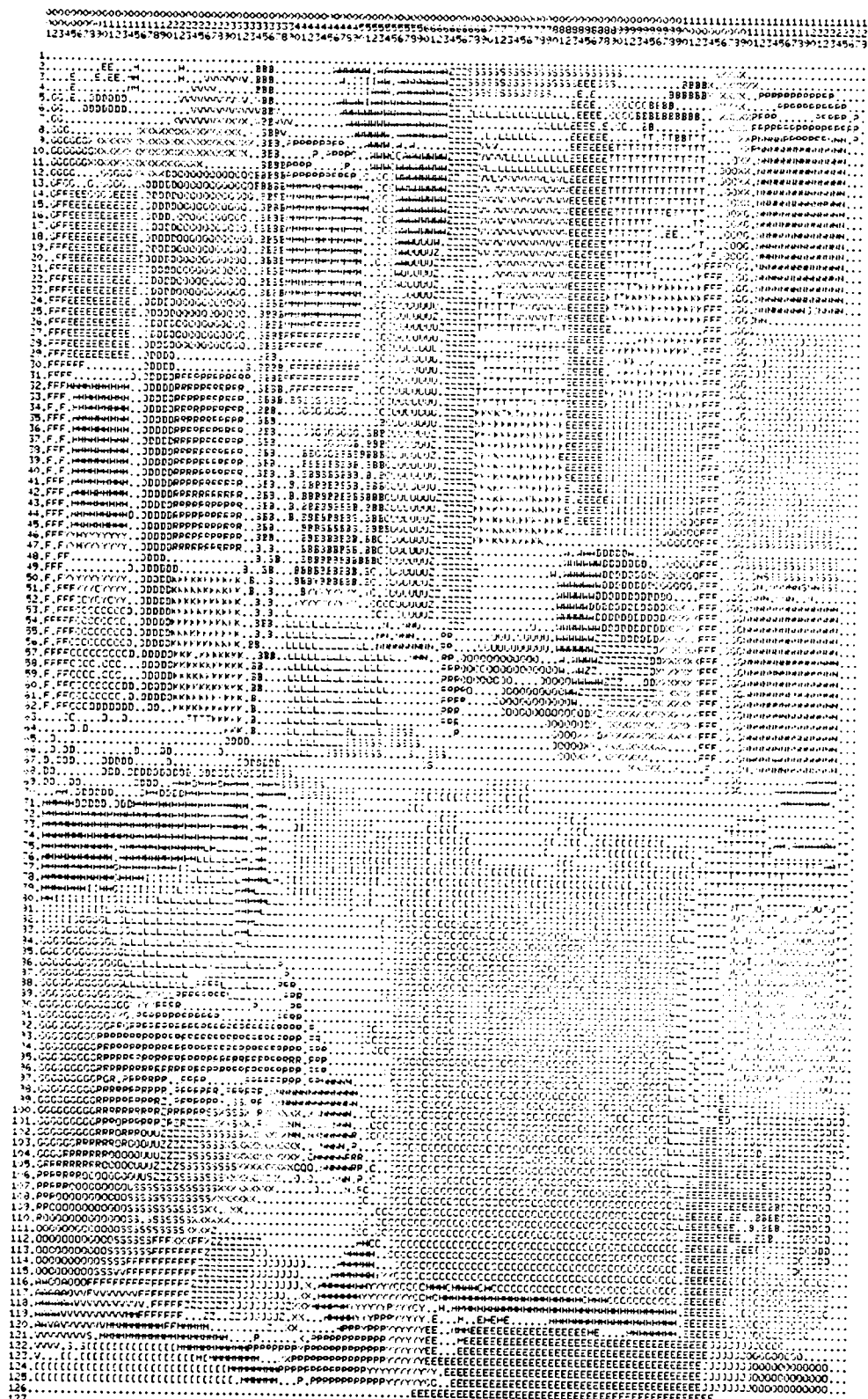


Fig.: 3.4.2-3a.

File: in107.fig



if8gacmc9.lgr



if9gacmc9.rgr

Fig.: 3.4.2-4a.

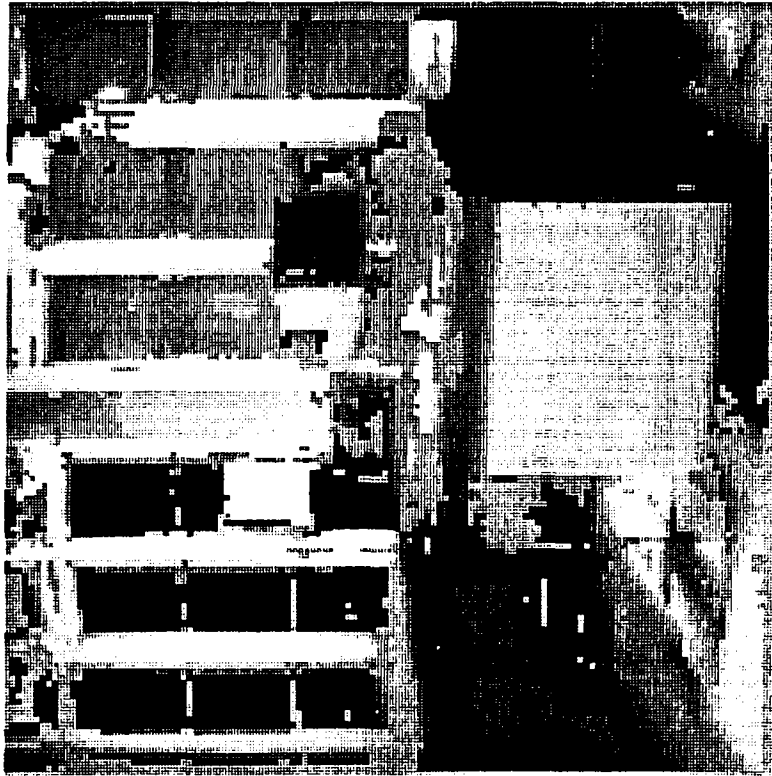


if8gray.lsh



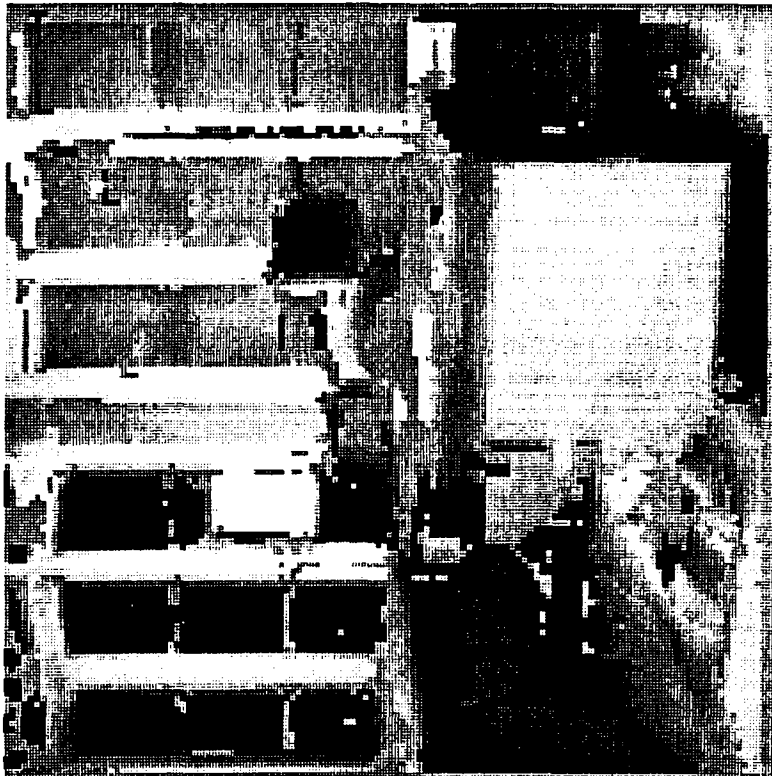
if8gray.rsh

Fig.: 3.4.2-4b.



if8ganmc.rgr

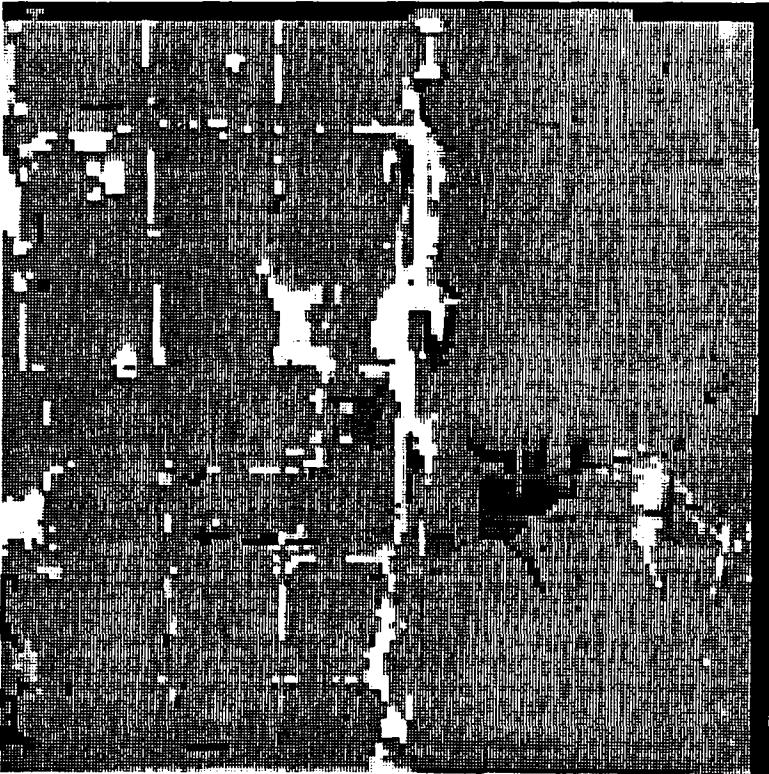
Fig.: 3.4.2-5b.



if8ganmc.lgr

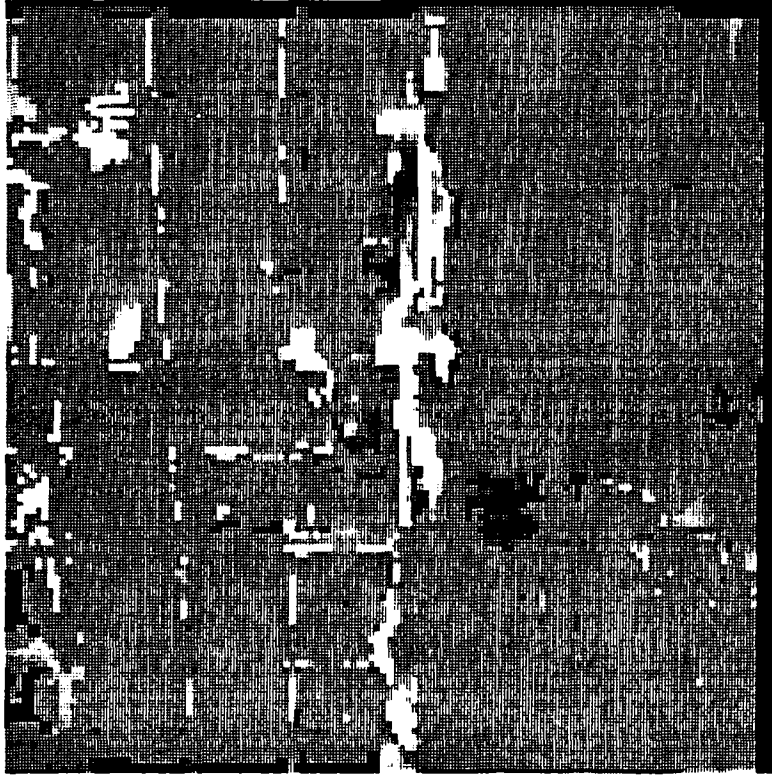
Fig.: 3.4.2-5a.

2 415,20



i18germc.lgr

Fig.: 3.4.2-6a.



i18germc.lgr

Fig.: 3.4.2-6b.

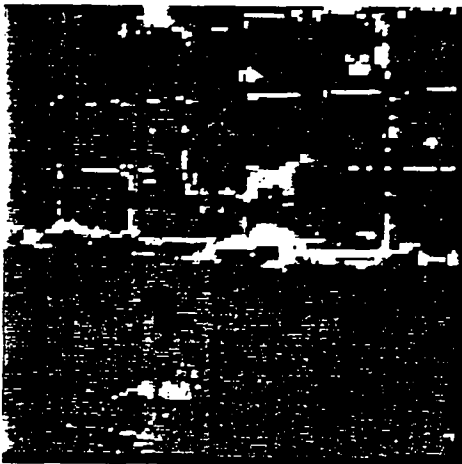


if8ganmc.lgr

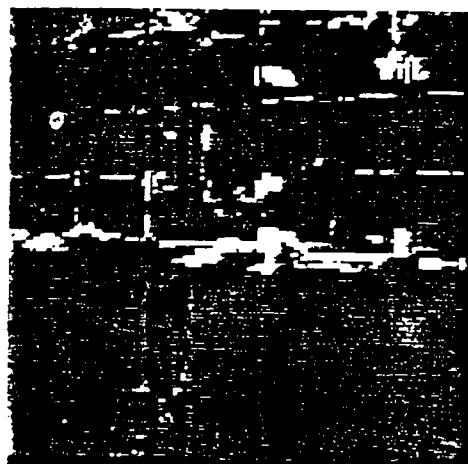


if8ganmc.rgr

Fig.: 3.4.2-7a.



if8germc.lgr



if8germc.rgr

Fig.: 3.4.2-7b.

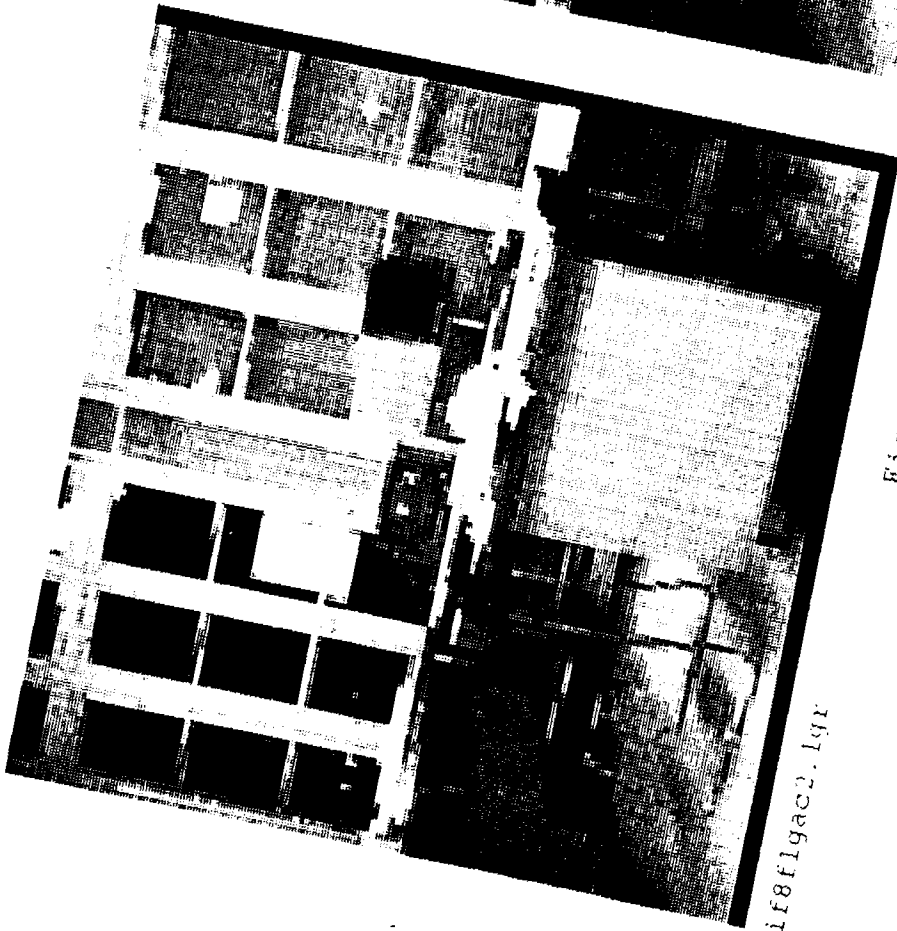


Fig.: 3.4.2-8a.

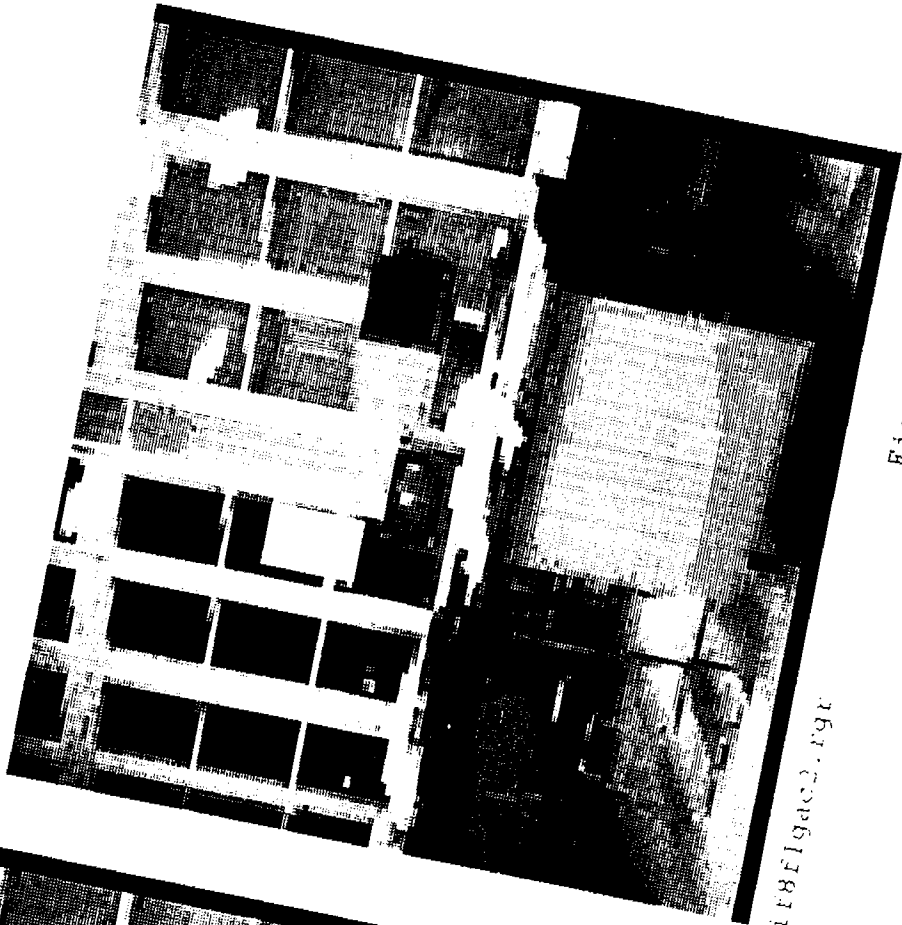
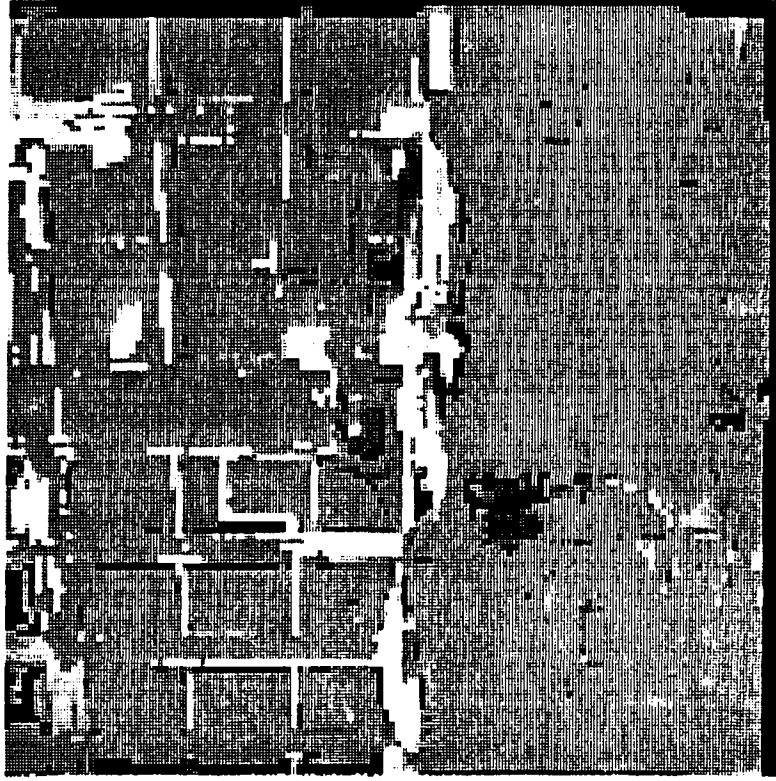
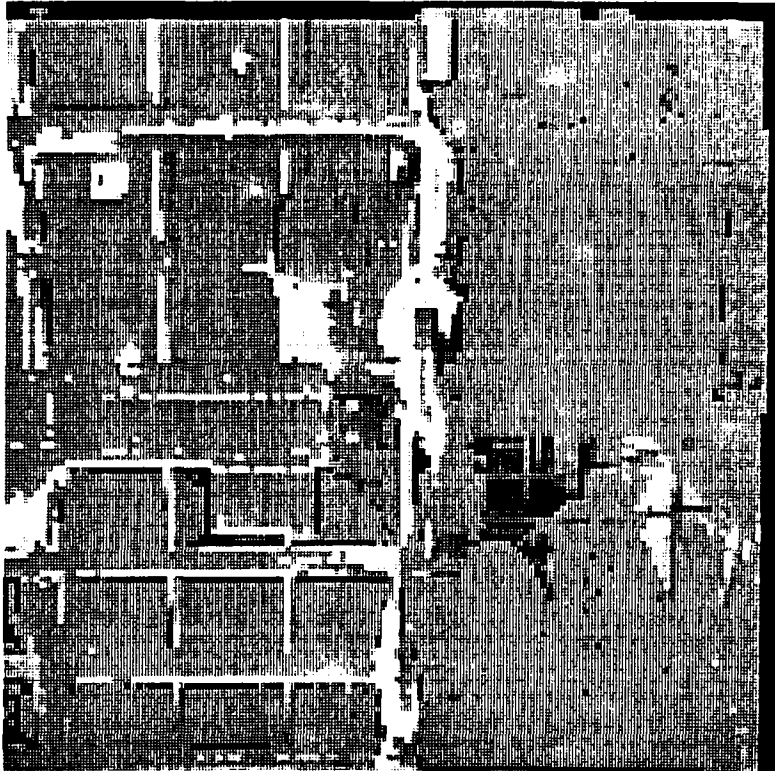


Fig.: 3.4.2-8b.



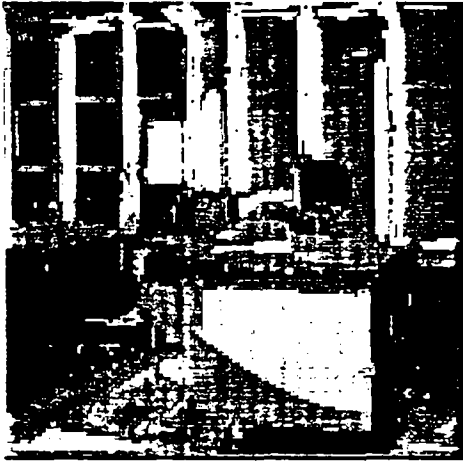
if8germc2.lgr

Fig.: 3.4.2-9b.



if8germc2.lgr

Fig.: 3.4.2-9a.

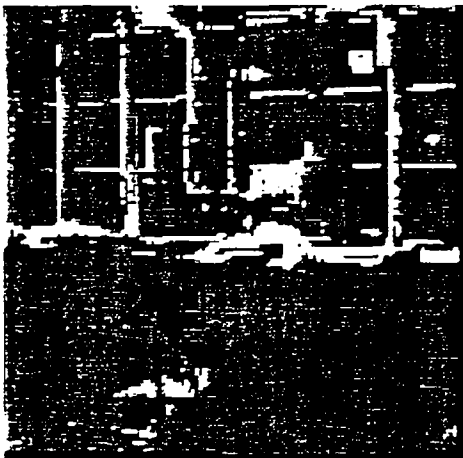


if8ganmc2.lgr



if8ganmc2.rgr

Fig.: 3.4.2-10a.

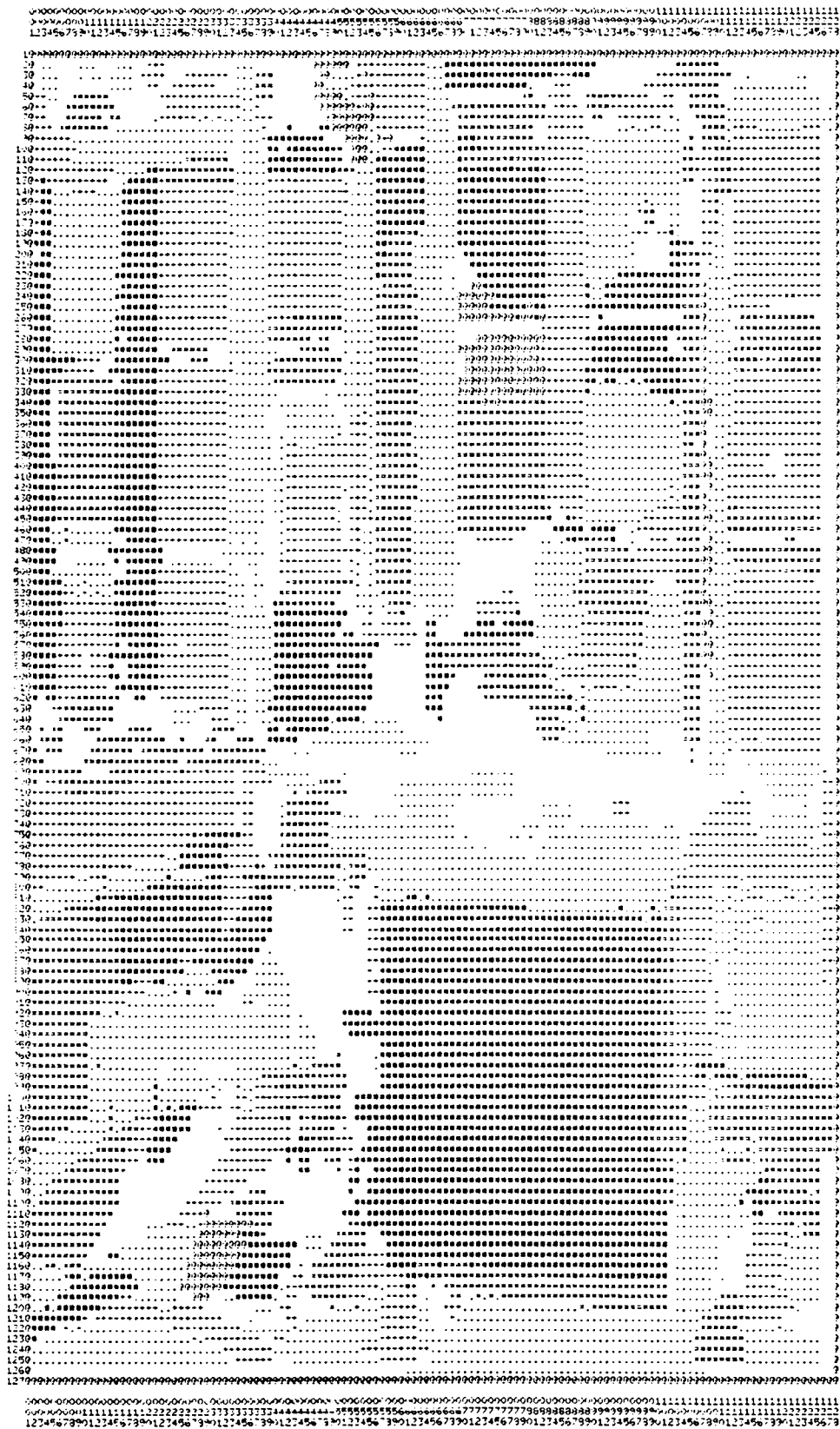


if8germc2.lgr



if8germc2.rgr

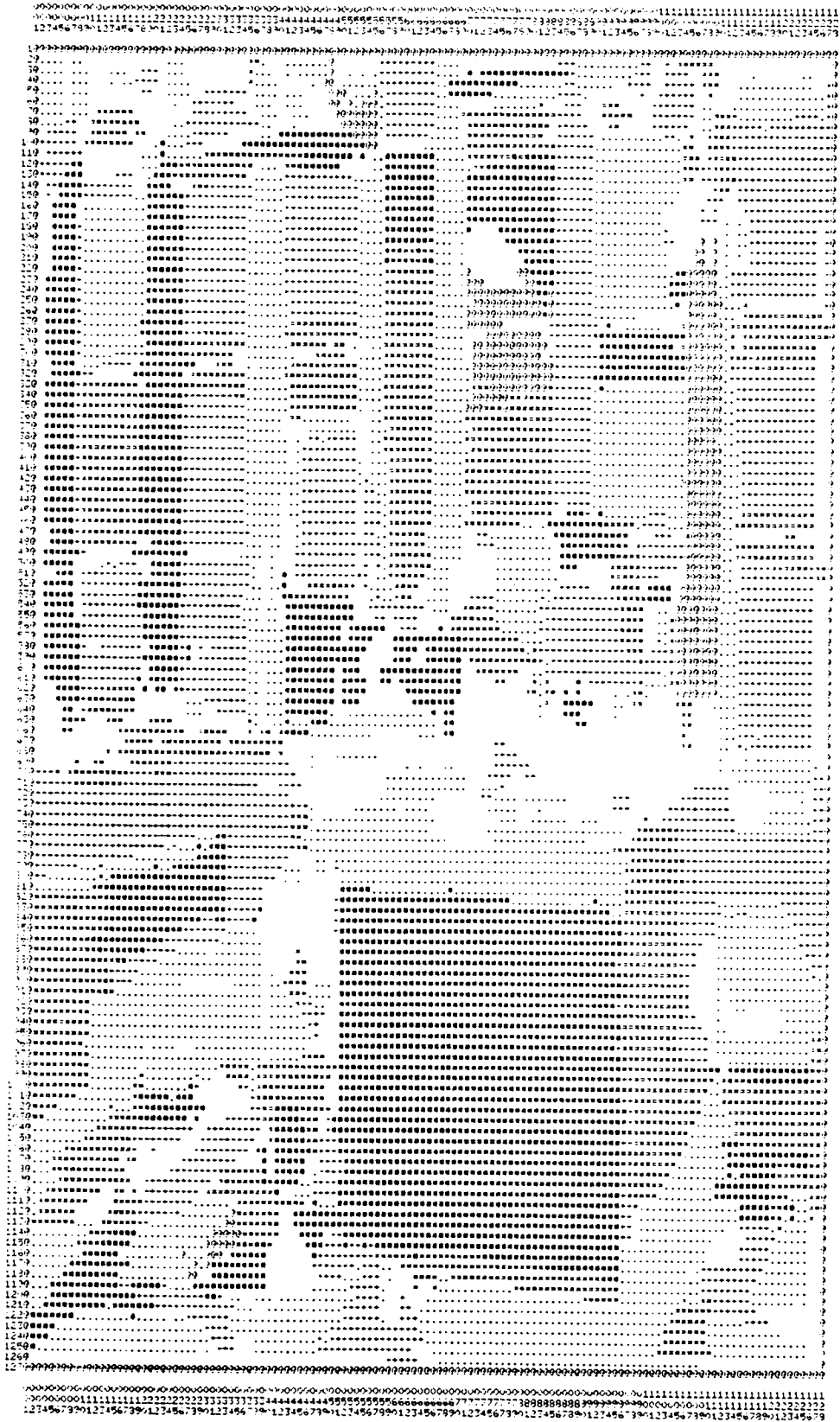
Fig.: 3.4.2-10b.



LEVEL GRAY IMAGE WITH OPTION FOR OVERLAY INVERSE
..... 1 2 3 4 5 6 7 negative value = overlay
..... 8 9 0 * * * * * blank from * onwards

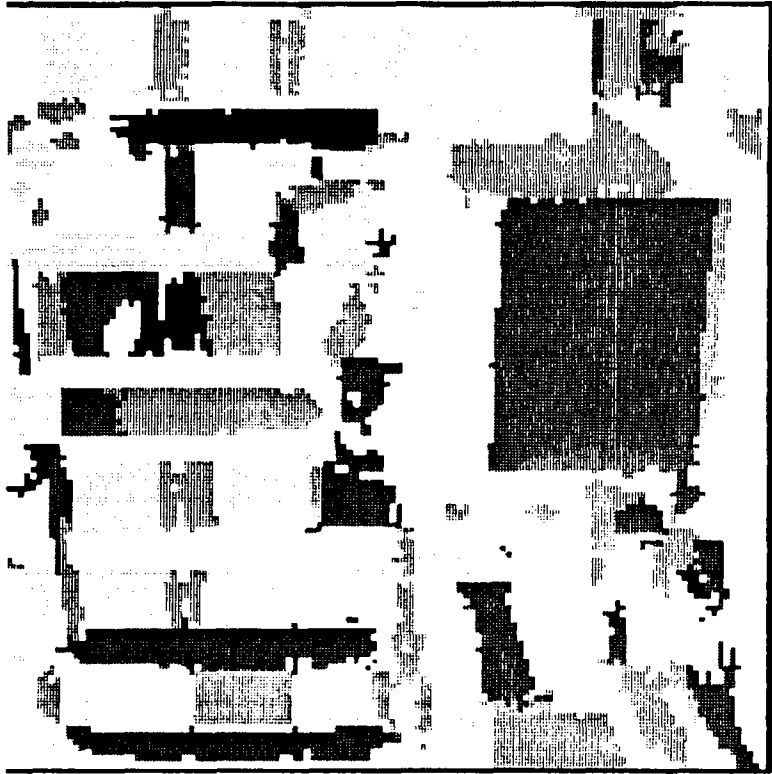
File: Inu93b.fig

Fig.: 3.4.2-11a.



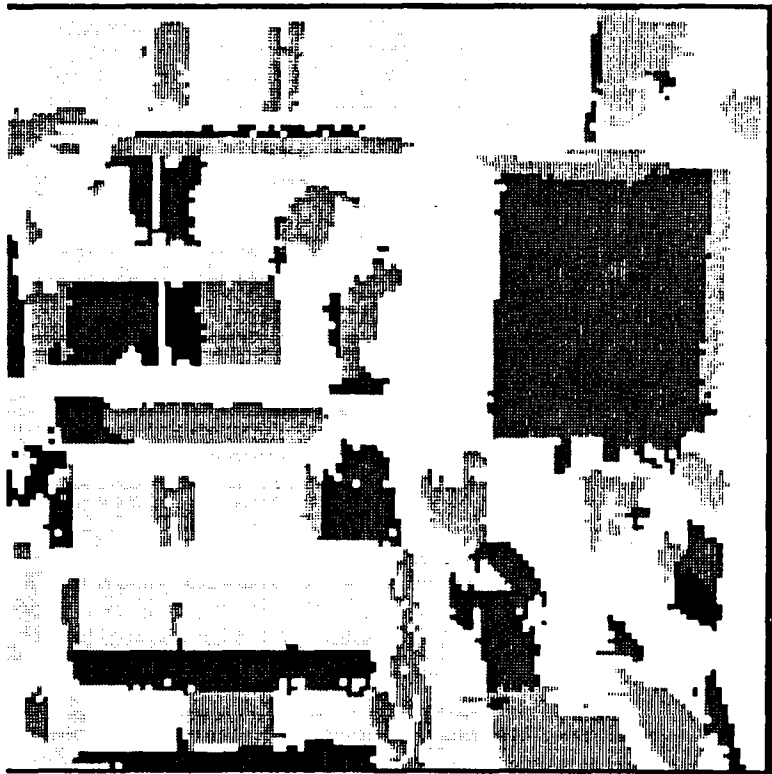
MAXOR FINMR IC = 0.2956E+1 3.3036E+0 3.1173E+0
 3-LEVEL IRAY IMAGE WITH OPTION FOR 3-D-LEVEL overlay
 1 2 3 4 5 6 7 negative value = overlay
 0 = (black from 7 colours)
 inv:94b.fig = R1(1).j. right, matched classified and related facets, printed in 6 colours.
 file: inv:94b.fig

Fig.: 3.4.2-11b.



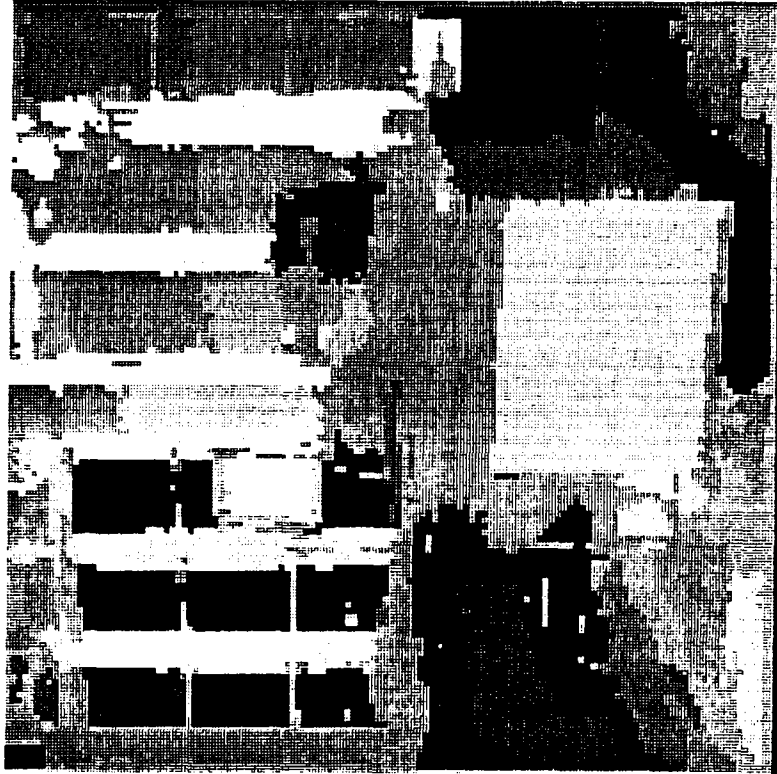
if8flkmc8.lgr

Fig.: 3.4.2-12b.



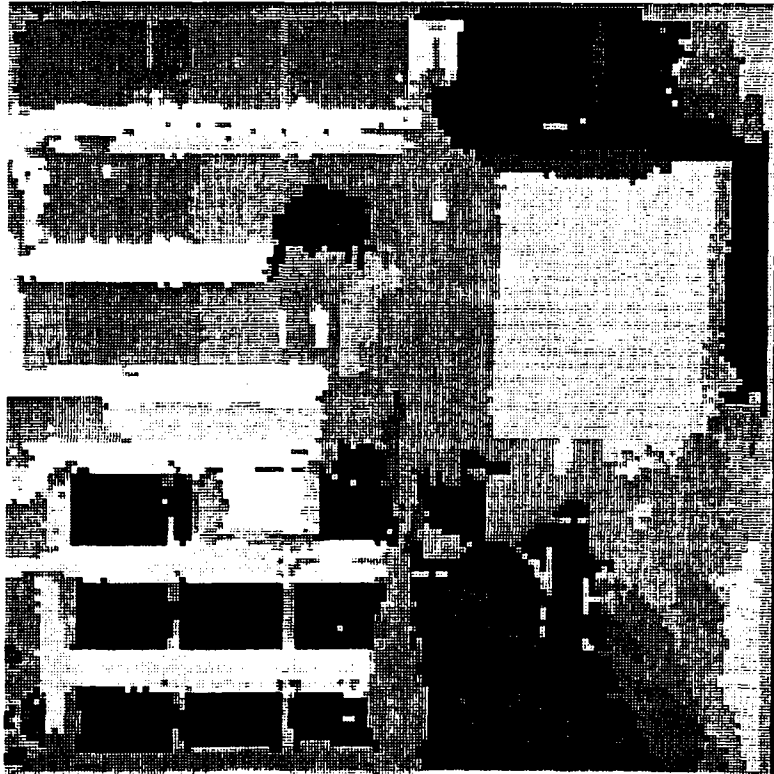
if8flkmc8.lgr

Fig.: 3.4.2-12a.



if8flqmc7.rqr

Fig.: 3.4.2-13b.



if8flqmc7.lqr

Fig.: 3.4.2-13a.



if8flkmc8.lgr

if8flkmc8.rgr

Fig.: 3.4.2-14a.



if8flgmc7.lgr

if8flgmc7.rgr

Fig.: 3.4.2-14b.

3.5 Comments

The processing steps following the matching are all well known, i.e., these consist of:

- (a) Improving the edges since they are rather "ragged" and have violated the "similarity principle". The presently popular method is called "snakes". The better established methods are based on regularization.
- (b) Segmentation of the edges and matching the edge pieces in the left and right images. If a piece-wise linear approximation is used then the results "interface" immediately to the edge based techniques. However, "edge segmentation" also has to obey the "similarity principle".
- (c) A problem that requires "further thought" is how to guarantee valid adjacency relations and in what conditions are these relations "imperative".
- (d) The interior of each facet is "just a smaller image". Consequently, the techniques in Chapter 2 apply again, but now on a smaller scale. However, there are also other methods, for example, the "rubber sheet" or "flexible template" approaches, as they used to be known.

However, instead of continuing, it is better to review and criticize the results since several weaknesses have become apparent.

REVIEW AND CRITIQUE

As described in the preceding chapters, the "similarity principle" is the "fundamental secret" for solving the 3D stereo problem since similarity requires no "higher level knowledge". For a considerable time the so-called "low level processing" of images has been considered either "well known requiring no further research" or "too low to be worthy of study". Immense efforts have been spent on "using high level knowledge" and on the application of "elegant mathematics" but without demonstratable results based on realistic images. Consequently, at least in the opinion of the present author, application of computers to image analysis, among other topics which are "natural and easy for ourselves", has been in a state of "Brownian motion without apparent drift" at least for the last decade or two.

In the present case all the "trends and fashions" were ignored and approximately half a year was spent on "thinking and experimenting". Brute computing power was available via the "INRIA net", ranging from "Suns to Connection machines", but a PC/AT was considered sufficient for preliminary experiments, especially since an image processing program library existed on the PC/AT. Considerable programming was done but most of the work consisted in modifying existing programs and assembling them into new running sequences. The experiments so far required 13 "machine loads" of programs (see "Inria*.for" in /scrt2/kasvand/ima3).

In brief, the processing steps may be summarized as follows, where the comments or observations have been included within the paragraph or inserted as "retrospective comments":

1. The computer is given two matrices (G_l and G_r) of integer numbers which are called the "left" (l or L) and the "right" (r or R) images of the stereo pair. No other information is neither available nor wanted since no "understanding" or other human concepts need to be evoked. The similarity principle is sufficient. Anything that is to be found is to be found from these two matrices of numbers.
2. These two matrices (G_l, G_r) are processed for gradients, which creates more matrices of numbers. G_l gives G_{pml} and G_{pal} , and G_r gives G_{pmr} and G_{par} ("p" indicates gradient operator, "m" indicates magnitude, "a" indicates angle, and "l" and "r" are for "left" and "right" images). Notice that only two of the practically unlimited number of features have been used.
3. The pixels of these matrices (G_l, G_{pml}, G_{pal} and G_r, G_{pml}, G_{par}) are clustered according to more or less classical pattern

recognition procedures. The regions in the images represented by the clusters are slightly "cleaned".

4. The connected regions of similar pixel classes are given "individual identity" (labelled) and the more "robust" regions are kept. The labelled regions are now called "facets".
5. The facets that produced an "intersection" (i.e., by "and-ing" the left and right images) are given similar labels in both images. Facets that could not be matched were also given new labels (negative), but these labels differed in the left and right images. This is a straight-forward application of well-known techniques and is but one of several methods for matching the left and right image regions.

Retrospective comments: The "new and unknown territory" begins at the end of step (3), in (4), and in (5). The "newness" is not in the techniques but in the modifications, improvements, and combinations required to make them "work better". The mistakes are, basically:

- (a) The classical techniques are too "myopic", i.e., the 3x3 4-connected neighborhood is unsuitable (too small).
- (b) More techniques should be used together to match the facets.

Now follows a long string of operations (steps 6 to 13) which were mainly studied for "curiosity" rather than necessity. In better quality images most of these operations are superfluous.

6. Some obvious misclassifications are corrected by not requiring class label coherence in the left and right images during matching.
7. The highly dissimilar but matched facets are "cut" to "reasonable" size by comparing them against each other. This step is justifiable in probabilistic terms but alternate methods should also be tried, for example, regions where the facets are very dissimilar should be reclassified.
8. Facet pairs which remained unmatched in "and-ing" even though they were very similar (since their intersection was zero or too small) are matched by recognition techniques. However, elementary pattern recognition methods can also produce many erroneous matches if many of the facets do not differ much from each other.
9. The facets are "clipped to the same size". This step was thought to be necessary to give a "fair start" to all the facets in the subsequent analytic relaxation procedure but otherwise has no strong logical justification. The use of

the centre of gravity of the "raw" facets, matching the centers of gravity, and "clipping" produces "good looking" facets but the adjacency relations between the facets are not properly preserved. Our stereo vision becomes "disturbed" by the "cracks" between the facets if the cracks do not match in the left and right images. Clearly, adjacency has to be preserved or recovered.

10. The facets that have remained unmatched are "transplanted" from one image into the other if the other image contains no labelled pixels at the same coordinate locations. This procedure, also, has only a probabilistic justification. The resultant facets, when viewed in stereo, "hang in free space in front of the scene".
11. The unmatched labels (negative) that still remain may be retained or eliminated by setting them to zero. Zeroing these labels has poor justification and should not be done but zeroing enlarges some of the "unknown" (zero-label) regions. This problem is at present unresolved.
12. The remaining unlabelled regions in the two images, if they "match", are extracted, "rounded" by binary processing, and labelled. This allows the "unknown facets" to "participate" in the analytic procedures but otherwise the procedure can only be justified in terms of probability.
13. The "unknown facets" found in the previous step (12) are put back into the left and right images. The left and right images are now rather well populated by "matched" facets.

The preliminary or "raw" matching of the facets in the left and right images is now complete. Most of the matches are correct and the matching method can be adjusted to the type of facet present, i.e., large facets can be "and-ed" for match, smaller facets require recognition techniques, including the adjacency relations to already matched facets, etc. In other words, this was a preliminary study to see what rather crude techniques can produce.

The next set of experiments are attempts to improve the shapes of the "raw" matched facets.

14. The gray level of each of the facets is approximated by an analytic function. A second order polynomial was used (since the programs were available from prior work).
15. The facets are now "relaxed" to try to assign as many pixels as possible to the various facets while preserving or recovering "stereo fidelity". The present relaxation program "ties together" the left and right images only by requiring that $E_r < E_{rlim}$ in

$E_r = W_g * |G_{org} - G_{ana}| + W_d * |G'_{org} - G'_{ana}|$, where

G_{org} = Original gray level value,
 G_{ana} = Analytic gray level value,
 G'_{org} = Original gray level gradient value,
 G'_{ana} = Analytic gray level gradient value,
 W_g, W_d = Weights,

provided that a pixel at (i, j) in the left image had an "homologue" which was within a box ($i = +/-d_i, j = +/-d_j$; $d_i, d_j = 2, 7$) and the homologue was within the error tolerance (Erlim). This constraint was far too weak.

Retrospective comments: The list of mistakes made is rather long. Some of these "mistakes" were deliberate in order to see if it is possible to "get away with" some crude method. Briefly summarized:

- a) The analytic approximations of the gray levels of the facets did not correspond to functions for surface brightness.
- b) No analytic approximation was used for the edges of the facets and thus the shapes of the edges could not be used.
- c) During relaxation decisions were made at pixel level without including sufficient information from the neighborhood around each pixel.
- d) The relaxation procedure was not sufficiently constrained and, consequently, produced rather "noisy" results.
- e) Most importantly of all, the relaxation procedure did not properly incorporate the similarity principle.

The present study ends at this stage since:

- i) The processes in Chapter 3 require modifications in light of the experience acquired and the remaining time is insufficient to complete the changes for this report.
- ii) The results can serve as starting points for numerous edge based techniques. The combinatorial explosion problem has been eliminated since only the edges of each of the matching facets have to be considered rather than all the edges in the image at the same time (1.2, 4.1).
- iii) Each matched facet may be looked upon as a small image and processed in the same way. Thus, the process may be iterated for each facet in order to achieve the next level of correspondence.

References

(The author apologizes. Time did not allow a careful search and review of the available material. For additional refs., see 1.1, 1.2, 4.1, and books entitled "... pattern recognition ...".)

- 1.1) Manual of Photogrammetry, Fourth Edition, C. C. Slama, C. Theurer, and S. W. Henriksen, eds. 1980.
- 1.2) Faugeras, O., A few steps towards artificial 3D vision, INRIA, research report No. 790, February, 1988.
- 1.3) Julesz, B., Foundations of Cyclopean Perception, University of Chicago Press, 1971.
- 1.4) Gibson, J.J. The Perception of the Visual World, Houghton Mifflin Co., 1950.
- 1.5) Horn, B.K.B., Obtaining shape from shading information. In P.H. Winston (Ed.), The Psychology of Computer Vision, McGraw Hill, 1975, pp. 115-155.
- 1.6) Grossberg, S., and Marshall, J.A., Stereo boundary fusion by cortical complex cells: A system of maps, filters, and feedback networks for multiplexing distributed data, Neural Networks, Vol. 2, 1989, pp. 29-51.
- 2.1) Peak viewer 2x, No.1994-2, made in Japan, approximate cost is \$40.- Canadian.
- 2.2) Kass, M., Computing visual correspondence, in From Pixels to Predicates, A.P. Pentland, ed., Ablex Pub. Co., 1986, pp. 78-92.
- 2.3) Otsu, N., and Kasvand, T., A new family of nonlinear edge detectors for noisy images. Paper #435-02. SPIE.
- 2.4) Forsynth, D., Mundy, J.L., Zissermann, A., and Brown, C.M., Projective invariant representations using implicit algebraic curves, Computer Vision - ECCV90, pp. 427-436.
- 2.5) Kasvand, T., Extraction of lines of unspecified texture from unconstrained line drawings, Science of Form: Proc. 1st Internat. Symp. on Science of Form, Y. Kato, General editor, KTK Scientific Pub., Tokyo, 1986, pp. 421-430.
- 3.1) Peyret, M., and Gagalowicz, A., Reconstruction 3D basee sur une analyse en regions d'un couple d'images stereo, Inria, Aug. 1989.
- 4.1) Remion, Y., Stereovision par zones, outils et structures d'un systeme expert, these docteur ingenieur, l'Ecole Nationale Supérieure des Telecommunications, 8 Juin 1988.

Appendix: IMAGE FILES

All the results shown in this report, including results not discussed, are available as computer files in the "Bora" system in /user2/kasvand/ima2/ for images and /scrt2/kasvand/ima3/ for some of the programs and the text of this report.

1. Text of report. The text has been written using PC/AT Wordperfect 4.2 but stored in Bora without the control codes (#1 foreign file output mode). The control codes have to be inserted for other types of word processors.

ich0p0n.prn = Beginning of report.
 ich1p0n.prn = Chapter 1.
 ich2p0n.prn = Chapter 2.
 ich3p0n.prn = Chapter 3.
 ich4p0n.prn = Review to end of report.

2. Images in one pixel per character alphabetic or alphanumeric form and "history files" of the runs for this report. Normal printing method is via X11 window using "small" print and "cat file_name". The images produce two pages of print which have to be "glued together" to produce the results shown.

inr000.fig = Index to inr*.fig and to image files.
 inr*.fig = Image or history files, see index (inr000.fig).
 * = 001 to approximately 200.

3. Display files, 256x256 one byte. These may be displayed on the Vicom image display system or printed on the laser printer using ~/im_laser -T -s n (-inv -xcen -ycen) file_name. The size of the image is defined by "n", n = 1, 2, ..., 6. Size 5 (n=5) is used for the "big" prints and size 3 (n=3) for the stereo pairs. Unfortunately, the gray level quality is very low and sinks with increasing n.

if8*.l* = Left image of the pair.
 if8*.r* = Right image of the pair.
 The if8*. part is the same for the left and right image.
 The .l* and .r* portions may sometimes be inconsistent.
 ifort8*.l* = Some older file, left image of pair.
 ifort8*.r* = Some older file, right image of pair.
 n*. * = Normalized display file (Inrimage "norma" used).

4. Composite display files (512x512 by one byte) for the Vicom system. These were intended for showing integrated results but were discontinued due to lack of disk space.

pvic*.l51 = Left image, composite of four results.
 pvic*.r51 = Right image, composite of four results.

ISSN 0249-6399