

Characterizing Stereo Matching Problems using Local Spatial Frequency

Mark W. Maimone

May 1996
CMU-CS-96-125

School of Computer Science
Carnegie Mellon University
Pittsburgh, PA 15213-3891 USA

*Submitted in partial fulfillment of the requirements for
the degree of Doctor of Philosophy in Computer Science*

Thesis Committee:

Steven A. Shafer, *chair*
Martial Hebert
Paul Heckbert
Michael Jenkin, York University

©1996 Mark W. Maimone

This research was sponsored in part by the Department of the Army, Army Research Office under grant number DAAH04-94-G-0006, and the NASA Ames Graduate Student Researchers Program NGT 51026. The views and conclusions contained in this document are those of the author and should not be interpreted as necessarily representing official policies or endorsements, either expressed or implied, of the Department of the Army or the United States Government.

Keywords: Computer Vision, Stereo Matching, Local Spatial Frequency, Gabor Filters, Scalogram, Phase-based Stereo, Scale Space, Perspective Foreshortening, Ambiguous (false) Matches, Ground Truth

Abstract

The model of local spatial frequency provides a powerful analytical tool for image analysis. In this thesis we explore the application of this representation to long-standing problems in stereo vision. As the basis for this analysis, we develop a phase-based algorithm for stereo matching that uses an adaptive scale selection process. Our approach demonstrates a novel solution to the phase-wraparound problem that has limited the applicability of other phase-based methods.

The problem of ambiguous matches, or false targets, can greatly reduce the accuracy of a stereo vision system. A common approach to alleviating the problem is the use of a coarse to fine refinement strategy, but we show that this approach imposes some (perhaps overly) strong requirements on the stereo images. Our phase-based method relaxes those requirements, and is therefore able to handle a wider variety of otherwise ambiguous images. But sometimes ambiguity is inherent in the images, so we propose a generalized disparity model to explicitly represent multiple candidates.

Perspective foreshortening, an effect that occurs when a surface is viewed at a sharp angle, can reduce the precision of stereo methods. Many methods tacitly assume that the projection of an object will have the same area in both images, but this condition is violated by perspective foreshortening. We show how to overcome this problem using a local spatial frequency representation. A simple geometric analysis leads to an elegant solution in the frequency domain which, when applied to our phase-based system, increases the system's maximum matchable surface angle from 30 degrees to over 75 degrees.

The analysis of stereo vision algorithms can be greatly enhanced through the use of datasets with ground truth. We outline a taxonomy of datasets with ground truth that use varying degrees of realism to characterize particular aspects of stereo vision systems, and show that each component of this taxonomy can be effectively realized with current technology. We propose that datasets generated in this way be used as the foundation for a suite of statistical analyses to effectively characterize the performance of stereo vision systems.

Acknowledgements

I am indebted to many people for the completion of this thesis.

My thanks go first and foremost to my advisor Steve Shafer. He helped me jump start my research career in vision and image processing, he encouraged me to develop my mathematical debugging skills, and his vision and ability to present ideas clearly have been inspirational. I am especially grateful for his encouraging comments over the years, and for his support during the final wrapup and writing process.

The other members of my committee have done a great deal to improve the content of this thesis. Thanks to Martial Hebert for providing overall direction, Michael Jenkin for our face-to-phase discussions, and Paul Heckbert for teaching me something new at every meeting. This document has benefited greatly from their excellent suggestions, as well as those of Yalin Xiong.

I would also like to express my sincere appreciation for the support, wisdom and friendship of A. Nico Habermann, founding Dean of the School of Computer Science. With his encouragement I learned to look beyond the present and literally reach for the stars.

For financial support, I thank CMU, the Office of Naval Research, and the folks at the NASA Ames Vision Research Lab led by Beau Watson.

The members of the Calibrated Imaging Lab, including Joyoni Dey, John Krumm, Bruce Maxwell, Reg Willson and Yalin Xiong, have been good friends and outstandingly supportive colleagues throughout my development as a vision researcher. Carlo Tomasi helped me overcome my initial fear of Fourier, and the dozens of members of the Vision and Autonomous Systems Center and its director Takeo Kanade have made this an intellectually challenging and personally satisfying environment. I am especially grateful for the support and friendship of Mei Chen, Jennie Kay, Pete Rander, Henry Rowley, Harry Shum, David Simon, Michael Smith and more.

The CMU Computer Science Department has been an incredible place to work. Students like Darrell Kindred, Corey Kosak, Henry Rowley, Bennet Yee, and the amazing folks of SCS Facilities have been a tremendous resource for all things technical. Our liaisons to the world outside of the department do an impossibly good job of hiding the gory details; I cannot give enough thanks to Sharon Burks, Catherine Copetas, Maggie Muller, Karen Olack and the whole SCS Support Staff for their friendship and all of their hard work (but I am certain Aunt Catherine is going to make me try :-).

My musical colleagues helped keep me sane during this whole process. Many thanks to the members and directors of the CMU Jazz Choir and Concert Choir, CS folks Angela Kennedy, Su Yuen Ling, Cliff Mercer, Scott Neal Reilly, Benli Pierce, John Velonis, and especially Geli Zobel for her dedication to singing and keeping me well fed. Special thanks to Robert Page for his energy, inspirational direction, support and personal encouragement of my musical goals.

Thanks too to the many other friends who have made the trip enjoyable in the good times, and bearable in the bad. Classmates Nevin (and Tess) Heintze, Puneet Kumar, Juan Leon, Francesmary Modugno, Joe Tebelskis, and Amy Moormann Zaremski helped start and keep me going. Officemates over the years put up with me and my strange musical tastes: Doug Reece, Raul Valdes-Perez, John Cheng, Keith Gremban, Ari Rapkin and Doug Rohde. And though I'm running out of space and cannot possibly include everyone, Drew Anderson, Eka Ginting, Chris McGlone, Phil Miller and my first advisor Jeannette Wing have all made special contributions. Warm thanks also to Patty, Mary, Chris, Tom and the rest of the late night crew.

Most of all, I thank all of my family for their encouragement, especially my parents and grandparents, without whose love and support none of this would have been possible.

Thanks to Larry Matthies for the image texture used in Figure 1.3 and Chapter 5. Thanks also to Blue Mountain Arts, Inc. for permission to use the image in Figure 4.1 (details follow).

5-DTM created by Stephen Schutz, Ph.D.

Published by Blue Mountain Arts, Inc. (R)

P.O. Box 4549, Boulder, Colorado 80306

Copyright ©1994 Stephen Schutz and Susan Polis Schutz

Reprinted with permission

Notation

Below is a list of notations that are used in this thesis and some of their consequences:

$ x _b$	Remainder	remainder(x, b), or x modulo b , or $x - \lfloor x/b \rfloor b$
$x =_b y$	Modular Equality	remainder(x, b) = remainder(y, b). Note that for all $a \in \mathbb{R}$, $x =_b y$ implies $x + a =_b y + a$, and for all $a \neq 0$, $x =_b y$ implies $ax =_{\frac{b}{a}} ay$
f	Focal Length	Intrinsic camera (lens) parameter associated with the pinhole lens model
$\omega \left(= \frac{1}{\lambda} \right)$	Frequency	Cycles per unit distance. Although this symbol is often reserved for <i>angular</i> frequency (in units of radians), we use ω to avoid confusion with focal length f . Frequency is the reciprocal of wavelength.
$\lambda \left(= \frac{1}{\omega} \right)$	Wavelength	Sinusoid or filter wavelength, often in units of pixels
$f(x) \Longleftrightarrow_{\mathcal{F}} F(\omega)$	Fourier Transform	Function $F(\omega)$ is the frequency-domain equivalent of spatial domain function $f(x)$. In this work $f(x)$ will be real-valued, thus $F(\omega)$ will be complex and odd-symmetric about the origin, so only the positive frequencies will be plotted.
ϕ	Phase	Phase of a complex number (often a Fourier or Gabor coefficient $\rho e^{j\phi}$), or phase of a sinusoid $\sin(x - \phi)$
ρ	Magnitude or Amplitude	Magnitude of a complex number (often a Fourier or Gabor coefficient $\rho e^{j\phi}$)

continued

$[low, high]$	Range (inclusive)	The set of numbers from <i>low</i> to <i>high</i> including the endpoints.
$\lfloor x \rfloor$	Floor	Floor of x , i.e. the greatest integer $i \leq x$.
$*$	Convolution	
$\{x : property\}$	Set Constructor	The set of all x such that <i>property</i> is true.
$X \cap Y$	Set Intersection	If X and Y are sets, the set of elements common to both.
$ S $	Cardinality / Absolute Value	If S is a set, the number of elements in S . Else if S is a scalar, the absolute value of S .

Contents

1	Introduction	1
1.1	Terminology	3
1.2	The Scalogram: A Unified View of Scale Space	7
1.3	Stereo Vision Algorithms	8
1.3.1	Stereo Matching	8
1.3.2	Problems with Traditional Methods	10
1.4	Related Work	12
1.4.1	Evolution of Phase-based stereo	13
1.4.2	Foreshortening	17
1.5	Relationship to Prior Work	18
1.6	Thesis Organization	18
2	A Taxonomy for Stereo Computer Vision Experiments	19
2.1	Overview	20
2.1.1	Representation of Results	22
2.2	The Basics: Mathematical Foundations	23
2.2.1	Example 1: Phase difference as disparity	23
2.3	Noiseless Synthetic Imagery	26
2.3.1	Example 2: Occlusion Masks	27
2.4	Synthetic Imagery with Noise	30
2.4.1	Example 3: Virtual Checkerboard	30
2.5	Controlled Environment	32
2.5.1	Example 4: Calibration Targets	33
2.5.2	Camera Calibration	35

2.5.3	Acquire Dataset Imagery	39
2.6	Measured Environment	40
2.6.1	Camera and Range Sensor Calibration	41
2.6.2	Acquire Imagery and Range Data Concurrently	42
2.6.3	Example 5: Textured Cube	43
2.6.4	Example 6: Model Train Set	43
2.7	Unconstrained Imagery	44
2.8	Implementation	44
2.9	Summary	45
3	Phase-based Stereo	47
3.1	Introduction	47
3.2	Local Spatial Frequency	48
3.2.1	Gabor Filters	51
3.3	Using Phase as Disparity	55
3.3.1	Relating Phase to Feature-based Methods	55
3.4	Phase Difference Measurement Characterization	56
3.4.1	Background on Instantaneous Frequency	57
3.4.2	Measuring Sinusoid Frequency (λ_{ref})	61
3.4.3	Measuring Phase ($\Delta\phi$)	62
3.4.4	Comparing Frequencies (λ)	66
3.5	Choosing Frequencies (<i>for all</i> λ)	66
3.5.1	The Image Scalogram	66
3.6	Combining Filter Estimates	69
3.6.1	The Phase Wraparound Problem	70
3.6.2	Our Filter Combination Method	72
3.7	Our Stereo Method	76
3.7.1	Sample Imagery	78
4	Accuracy: Dealing with Ambiguous Matches	81
4.1	The Problem of Ambiguity	81
4.1.1	Definitions	83
4.1.2	Disparity Space	86

4.2	Effect of Ambiguity on Stereo Methods	89
4.2.1	JISCT Results	89
4.2.2	Jones' Method	92
4.3	Reducing the Ambiguity Factor	93
4.3.1	Coarse to Fine Method	94
4.3.2	Coarse to Fine Results	95
4.3.3	Solution using Phase-based Stereo	96
4.4	Modeling Ambiguity	105
4.4.1	Extended Disparity Representation	105
4.4.2	Demonstrating Improvement on an inherently ambiguous image . . .	109
4.5	Summary	114
5	Precision: Effects of Foreshortening	117
5.1	Relating Disparity to Surface Angle	120
5.2	Expressing the Foreshortening Factor using Image Parameters	121
5.2.1	Verifying the Foreshortening Factor	123
5.3	Applicability	126
5.4	Application	128
5.4.1	Extending Phase-based Stereo Algorithms	128
5.4.2	Results	129
6	Contributions	135
A	Theodolite Error Analysis	139
A.1	Deriving equations for the coordinate axes	140
A.2	Results	142
B	Derivation of Foreshortening Probabilities	145
C	Numbers	147
D	Application of Visual Reconstruction	149

List of Figures

1.1	Illustration of the geometry of a typical stereo vision system.	3
1.2	Evaluation functions are the tools used to find corresponding pixels. Given a pixel in the left image, a set of pixels in the right image is evaluated. The right pixel that exhibits the smallest error in the evaluation function is selected as the correspondent.	5
1.3	Stereo pair illustrating the effects of foreshortening; image compression, differing spatial extents.	6
1.4	Overhead view of a typical Stereo Vision setup. A pair of cameras with focal length f , separated by baseline distance B , are represented by their focus points on the lower left. Object point P has corresponding image coordinates x_{iL} and x_{iR} , and lies at depth Z from the focus points. Mirror image vectors $-x_{iL}$ and $-x_{iR}$ are shown to make the similar triangles used in Equation 1.1 more explicit.	9
1.5	Matlab implementation of the Kuglin-Hines Phase Correlation Image Alignment method (Kuglin & Hines, 1975). Actual computation of the global translation vector occurs in the first two lines of the function, the final six lines extract the indices of the maximum element in a 2D array. The method actually achieves subpixel precision, but this unrefined index-processing code only returns integer results.	14
2.1	Disparity as a function of Phase Difference. Disparity is the horizontal separation between the two signals, and is indicated by the labeled bar just to the left of 0.	25

2.2	Occlusion Mask Generation: <i>Top Row:</i> Stereo pair of images from the Left and Right cameras; actual disparity maps for those images. <i>Middle Row:</i> Pointwise occlusion masks for left and right images (note especially the noise in the right mask around the object borders); disparity maps with pointwise occlusion masks overlaid. <i>Bottom Row:</i> Plane-fitting occlusion masks for left and right images; disparity maps with plane-fitting occlusion masks overlaid.	28
2.3	Finding the best-fit plane in the disparity map. There are four 2x2 windows that contain the center (dark) pixel; the upper right window is highlighted.	29
2.4	Examples of noise easily modeled with synthetic data.	31
2.5	Actual image of the Calibrated Imaging Laboratory (CIL) calibration target with virtual rendering overlaid. The grey background and grid of black dots are part of the original picture, the white dots are rendered dots located at the 3D grid point locations. The dots were rendered as spheres using a virtual camera with the same parameters as those computed from the real image. . .	34
2.6	More sample calibration targets. The left image is the calibration cube from (Shum et al., 1995), right is an image of the MOVI (INRIA) “inverted cube” calibration pattern from (Boyer & Berger, 1995).	35
2.7	Coordinate frames used in dataset acquisition.	37
2.8	Textured cube image and the piecewise-planar patches used by Xiong for error analysis (Xiong, 1995, Figure 3.20). Used with permission.	43
2.9	Image from CIL-0001 dataset with the locations of ground truth measurements.	44
3.1	Local Spatial Frequency illustration (magnitude only). The Original Signal (lower right) is a low frequency sine wave embedded in a high frequency wave. Its Fourier Transform (upper left) has two peaks, one for each frequency. The Ideal Local Spatial Frequency plot (lower left) associates each sample point with its proper frequency. The Spectrogram (upper right) only approximates this ideal plot, but is computed directly from the sample points <i>without</i> explicit knowledge of the original signal’s analytical form. Peaks in the 1D DFT plot correspond to the dark areas in the 2D images.	50
3.2	Gabor filter examples with $\omega = 1/10$ but different m, σ_f on each row. The horizontal axis indicates the number of samples in the filter.	52

3.3	Scanline Pair used in Table 3.1, highlighted in the image from which it was extracted.	54
3.4	Instantaneous Frequency Example: $\omega = 0.06667$. The figure is explained in Section 3.4.1	59
3.5	Instantaneous Frequency Example: $\omega = 0.03333$ to 0.06667 in linear steps, i.e. a <i>chirp</i> signal. The figure is explained in Section 3.4.1	60
3.6	The effect of measurement error vector $\vec{\epsilon}$ on phase angle ϕ	63
3.7	Maximum phase angle error as a function of the length ratio.	64
3.8	Sample Magnitude profile of several Gabor filters, illustrating the peak finding heuristic.	65
3.9	Double sine wave signal and associated scalogram (both magnitude and phase).	67
3.10	Frequency Sampling in the Scalogram and other methods. Grey areas show the central tuning frequencies that are used by each method. Filter widths vary, but typically cover most of the regions between tuning frequencies (see Figure 3.13 for an example).	68
3.11	Ideal phase difference as a function of Disparity, illustrating the phase wraparound problem. Note that coarser scales are on the <i>left</i> in the Frequency plots, but on the <i>right</i> in the Wavelength plots. Each graph has 63 sample points spaced linearly along the X axis. Note that unlike the 2D scalogram plots, here the wavelength axis is <i>horizontal</i>	71
3.12	Filter contributions to the disparity estimate. Pixel intensity represents the value of the error computed by each filter's phase difference comparison, scaled by the magnitude. Darker pixels indicate larger error, and white spots represent areas where the scalograms contained no useful information. The evaluation function (below) is constructed by summing up the values in the columns and dividing by the number of filters used. This plot represents the disparity computation at pixel (128,128) of the stereo pair in Figure 3.14.	74
3.13	Gabor Filter Frequency Response. The Gabor filter magnitude in the frequency domain for filters with wavelength 2, 4, 8, 16, 32, 64, and 128 is presented. The tuning parameters are $\sigma_f = \frac{1}{6}$ and $m = 4$	75

3.14	Synthetic Stereo Pair with Ground Truth. The top two images form the stereo image pair of two frontoplanar surfaces, the middle images are the associated disparity maps with occluded areas marked in black. The lower left image is the disparity map computed by this method, evaluating 301 disparities from 0 to 30. The lower right image is the difference between the stereo disparity and the ground truth. Mean error over the entire figure is 2.41 pixels (variance 23.09); all errors ≥ 2 are in black. (intensities not to scale)	79
3.15	Actual Stereo Pair with Ground Truth. The top two images form the stereo image pair of a flat sheet tilted back slightly, and have been corrected for lens distortion. The middle images are the ground truth data on the left (black areas have unknown ground truth), and disparities computed by our method on the right (161 disparities, 0-40). The bottom image shows the differences between our estimates and the ground truth over the central target, scaled from 0 to 2 (all errors ≥ 2 appear black). Mean error in that region is 1.27701 pixels with $\sigma = 2.8446$. (intensities not to scale)	80
4.1	A 340x340 autostereogram of a wolf. See Figure 4.24 for an elucidation of the embedded structure. Reprinted with permission from Blue Mountain Arts, Inc.	83
4.2	Evaluation Function for pixel (170,100) of Figure 4.1 (with its pre-shifted counterpart), computed using our phase-based method with linear frequency samples. This information is presented in the context of a complete scanline in column 100 of Figure 4.3. Darker pixels denote less error, so the dark stripe from pixel number 75 to 200 at disparity 7 represents the best guess.	84
4.3	Disparity Space for row 170 of Figure 4.1 (with its pre-shifted counterpart), computed using our phase-based method with linear frequency samples. Darker pixels indicate lower matching errors, i.e., the most likely disparity estimates. Column 100 can be seen in an expanded view in Figure 4.2.	87
4.4	Coarse to Fine disparity space and associated scales; fine scales are darker than coarse scales. These results are also from row 170 of Figure 4.1, as are those in Figure 4.3.	89
4.5	Shoe Stereo Pair. These are images SHOE-0 and SHOE-2 from the JISCT Stereo Evaluation study; each is 480x512 pixels ² .	90

4.6	Results of several stereo methods on Figure 4.5 as presented in the JISCT Stereo Evaluation study. Clockwise from the top left are the INRIA 1, INRIA 2, Teleos, and SRI results; each is 59x63 pixels ² . Black pixels were marked as unknown by the algorithms.	90
4.7	Detail view of SRI disparity results and approximate ground truth for the middle row of the shoe stereo pair disparities in Figure 4.6.	91
4.8	Results of Jones stereo method on middle region of speaker images.	92
4.9	Speaker Image.	93
4.10	Evolution of the Coarse to Fine evaluation function at pixel (170,100) in Figure 4.1. The coarsest scale appears on top, the plots below demonstrate the successive refinement. Contrast the bottom plot with Figure 4.2.	97
4.11	Synthetic Speaker Grill. Actual disparity is 18.3445 pixels, but the pattern repeats approximately every 3.8 pixels.	98
4.12	Coarse to fine results on Figure 4.11. Disparity map (upper left) has mean error 14.42 with $\sigma = 23.1308$ pixels, ground truth error image (upper right) maps all errors ≥ 2 to black. Lower plots are the coarse to fine disparity space (left) and scale space (right) for row 240 in the synthetic speaker grill. Ground truth for pixels 200–400 is 23.52 pixels.	99
4.13	Coarse to fine results and images from the Shoe image pair.	100
4.14	Synthetic Grill evaluation functions, illustrating the improvement of the phase-based method over raw correlation (too many minima) and coarse to fine (trapped in local minimum). Actual disparity is 23.52 pixels at pixel (240,300).102	
4.15	Synthetic Grill Phase-based Disparity Space and constraint-filtered Scalogram for line 240. The raw pixel intensities appear above (the same scanline is plotted twice), and illustrate the effect of the grill edges on the plots below. .	103
4.16	Detail view of phase-based disparity results and approximate ground truth for a middle row of the shoe stereo pair. Background disparities are consistent, but incorrect. Figure 4.21 contains the complete disparity space for this scanline.106	
4.17	Image Scalogram for the middle row of the left shoe image.	106
4.18	Phase-based disparity maps for the shoe image pair. Left map is the result from using no constraints, right map is the result from using the heuristic in Section 3.4.3.	107

4.19	A view of the deviousness used in the construction of the Shoe image pair. Although the images were taken of the shoe flat against the texture (left image), the fact that the actual disparity is larger than the period of the checkerboard pattern causes a purely local search to reconstruct the scene with the shoe “floating” above the textured background (right image). Reproduced with permission from Kanade.	107
4.20	Illustration of the peak-finding heuristic from Sections 4.4.1 and 3.4.3 on a sine wave (upper) and the inverted Wolf Evaluation Function from Figure 4.2 (lower). Dashed lines indicate the region in which a Gaussian is fit and then subtracted.	110
4.21	Shoe Disparity Space. The upper plot is row 240 from the Shoe image pair, the lower plot is the disparity space computed by our phase-based method. .	112
4.22	Phase-based Disparity Space with Peaks. Peaks computed using the heuristic have been superimposed on the disparity space from Figure 4.21.	113
4.23	SRI Disparity over Phase-based Disparity Space. The seemingly random results from Figure 4.7 actually fit nicely into the local minima of the phase-based disparity space (the SRI disparity space was not available).	113
4.24	The structure embedded within Figure 4.1. This image was computed by applying our coarse-to-fine stereo method with 5 pixel window and 5 pixel max disparity per level to the original 340x340 image and a copy of the original shifted by 58 pixels.	115
5.1	Overhead view of the foreshortening model. X_S is the distance from the point exactly in front of the left camera (the origin O_S at distance Z_L) to the point (S) on the plate being studied; x_{iL} and x_{iR} are the left and right pixel indices of the image of surface point S ; the cameras are separated by baseline B and the surface tilts away from the cameras at angle θ	118
5.2	The effect of foreshortening on scalogram magnitude. Two views of a flat plate with a sinusoidal texture appear on top, and the scalogram magnitudes for their central scanlines appear below. The responses are similar, but are compressed to higher frequencies in the rotated view.	119
5.3	Overhead view of the foreshortening model. Similar triangles for the left camera geometry are highlighted (see Equation 5.3).	120

-
- 5.4 Overhead view of the foreshortening model. Similar triangles for the right camera geometry are highlighted (see Equation 5.4). 120
- 5.5 Left and right views of a surface tilted 65 degrees. Upper images are the central scanlines, lower images are their corresponding scalograms. You can see similar features in both scalograms: those in the left image are present at higher spatial frequencies because the left image is subject to greater foreshortening effects than the right image. 122
- 5.6 Verifying the Foreshortening Factor - These graphs compare the predicted foreshortening factor (dashed line) against that computed using only image information (solid line). The virtual lab setup (top left) and an example input image with surface angle of 60° (top right) are shown first. Next we have the results derived from a surface angled at 0° (middle left), 30° (middle right), 45° (bottom left), and 60° (bottom right). The virtual surface is 4.0 units from the left camera, both cameras have a field of view of 45° and are separated by a baseline of 0.4 (the surface in the actual images is larger than that shown in the top left rendering). 124
- 5.7 Foreshortening Factor as a function of Depth and Angle. Depth is unitless relative to the baseline, and varies from 3 to 100. Angle varies from zero to 85° . 126
- 5.8 Ground Truth and computed disparity maps for a surface angled at 65° . The top row shows ground truth in perspective on the left, a graph of a representative scanline from all methods on the right. The middle row shows perspective views of the disparity maps computed by the foreshortening-corrected method, Kanade/Okutomi and the uncorrected phase method. The bottom row shows differences between actual disparities and those computed by the foreshortening-corrected method, Kanade/Okutomi and the uncorrected phase method, for pixels that image the plate; darker values denote larger errors. Only differences between 0 and 2 pixels are shown, errors larger than 2 pixels appear as a 2 pixel error. Actual plate disparities range from 25.3 to 39.9 pixels. 131

5.9	Ground truth and disparity (computed by both the uncorrected and foreshortening-corrected phase methods) for the center scanline of the city scene at various rotations. From left to right (and top to bottom): 0, 15, 30, 45, 60, 75, and 80 degrees.	133
A.1	Region of error. The greatest possible error occurs across one of the Target Area diagonals (see Figure A.3 for a close up).	140
A.2	Two dimensional view of the theodolite imaging process.	140
A.3	Target Area Precision (zoom in on Figure A.1): the largest error is the length of one of the diagonals h and v	142
A.4	Error Space for $\delta = 20''$ with unit baseline: depth error $v(\theta, \psi)$ is greater than horizontal error $h(\theta, \psi)$ for all but the nearest points (when the sum of the angles is greater than about 90 degrees)	143

List of Tables

2.1	Stereo Imagery Characteristics and the types of scenarios for which they are available.	20
3.1	Effects of Tuning Parameters on stereo disparity accuracy. Filters with given values for σ_f and m were applied to the scanline stereo pair from Figure 3.3 using the method in Section 3.7. Here the table shows the RMS error between the computed disparity and that known over the portion of the image with visible texture (153 pixels). 401 disparities from 0 to 10 were tested, and the smallest and largest errors are highlighted.	54
3.2	Pseudocode for the basic phase-based stereo algorithm.	77
4.1	Pseudocode for the Coarse to Fine algorithm.	94
5.1	Probability that a surface exhibits $\geq 10\%$ variation between images due to perspective foreshortening. The distribution of surfaces is assumed to be uniform within the range of orientation angles from $-\frac{\pi}{2}$ to $\frac{\pi}{2}$, and depth ratios (distance divided by baseline) are as specified. A sample derivation can be found in Appendix B.	127
5.2	Pseudocode for the foreshortening-corrected algorithm. Column index c must be zero in the center of the image.	129
A.1	GNUplot commands that generated Figure A.4	144
C.1	Default parameter values used to generate images in this thesis.	148

Chapter 1

Introduction

It is a foolish thing to make a long prologue, and to be short in the story itself.

— Second Maccabees 2:32

Stereo matching provides the foundation for many methods in computer vision, and a large number of practical applications: robot navigation, parts inspection, aerial or satellite mapping, and medical imaging to name a few. In all of these, the distance to objects in the scene is computed passively by comparing several images of the world. The distance thus computed might then be used to plan a robot's path, determine the pose of a CAD model, generate landscaping contours, or position a robot arm. Since these tasks may be of critical importance, the distance estimates must be well-characterized and precisely determined. Although one can try to extract distance from a single image (e.g., using shape from shading) or many images (e.g., using multibaseline stereo or optical flow), we will study a minimal configuration for model-independent stereo: two cameras with known position and orientation, observing a scene comprising diffuse objects.

The best way to measure the success of a stereo method is to compare its results against the *ground truth*, or range information measured using means other than stereo. Acquisition of total and accurate ground truth is very expensive, often prohibitively so. But just as there are trade-offs in modeling complexity v. precision in the design of vision algorithms, so too there are trade-offs in the analysis of their performance. Synthetically generated imagery gives total ground truth knowledge, but fails to model the complexities of real-world imaging. Real imagery provides better test data, but greatly reduces the amount

and density of ground truth available for analysis. If computer vision is to become more of an engineering discipline than craftwork (Shaw, 1990), engineers must be able to predict and experimentally characterize the behavior of their systems. Such characterization is only possible when ground truth is available. In Chapter 2 we outline a framework for choosing a reasonable trade-off of ground truth density v. image realism in the analysis of stereo algorithms.

Our approach to the stereo problem will explore the benefits of a *local spatial frequency* representation. There are two fundamental extremes in the study of 2D imagery: the detail-oriented *spatial* view and the wholistic *frequency* view. In the spatial view each individual pixel is paramount; an image is represented by the concatenation of independent pixel values, and each pixel is considered unique and important. In the frequency view only the entire image matters; although the image is broken down mathematically into several frequency components, information in each component relates to the image as a whole. Each view has its benefits: the spatial view can represent discontinuous textures and local (pixel-level) segmentations directly, while the frequency view is a mathematically elegant representation that enables many useful analyses over large regions. The trouble is, the benefits of one view are difficult to attain in the other. Fortunately, a compromise can be reached, one which preserves the localizability of the spatial approach and the analytical benefits of the frequency approach: this combined approach is called *local spatial frequency*, and forms the foundation for the discussion of our stereo method in Chapter 3.

The evaluation of any experimental method requires a characterization of the accuracy and precision of its results. By *precision* we mean the number of significant digits in the numeric result, and by *accuracy* we mean the agreement between the experimentally-derived result and the ground truth. The claim that for all stereo images, all pixels have a disparity of 3.14159 pixel units is very precise, but probably not very accurate. In contrast, the claim that for all stereo images, all pixels have a disparity of 0 ± 512 pixel units is accurate, but not very precise. Chapter 4 will show how to improve accuracy using better models of ambiguous matches and Chapter 5 will show how to improve the precision of stereo results in the presence of perspective foreshortening.

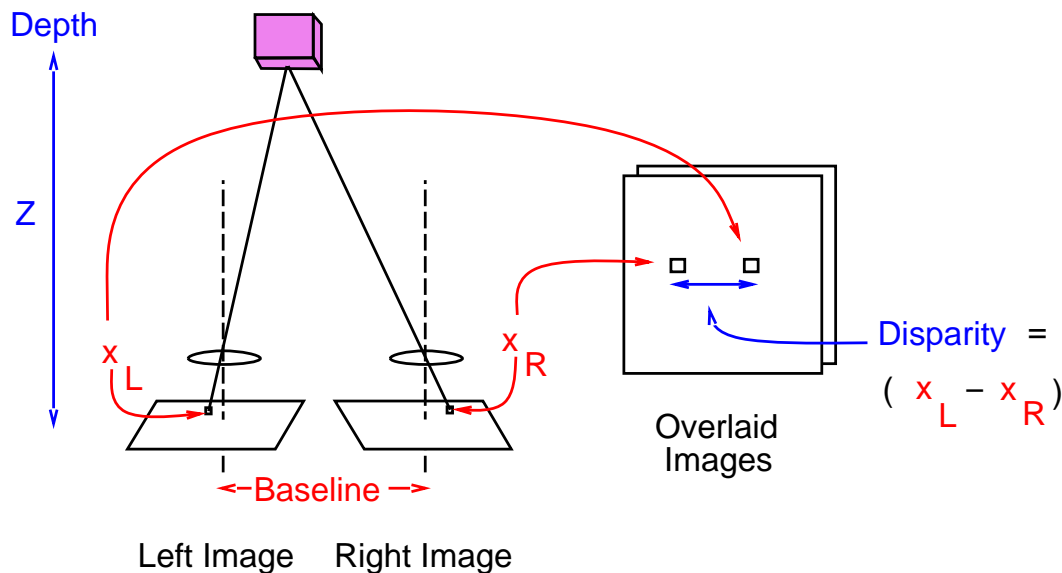


Figure 1.1: Illustration of the geometry of a typical stereo vision system.

1.1 Terminology

Some of the terminology used in describing the behavior of stereo algorithms can be described as follows:

Local Spatial Frequency Two common approaches to image interpretation are frequency analysis (e.g., using Fourier transforms) and direct pixel-based analysis (spatially dividing a picture into a grid of pixels). For years only one of these approaches was used at a time: frequency analysis for global effects, spatial analysis for local effects. But now *local spatial frequency* representations allow both techniques to be used simultaneously, at the cost of greatly increasing the amount of data used to represent the image. A 1D scanline and its Fourier transform can each be expressed in one dimension, but a local spatial frequency version of the same scanline will span two dimensions. Some examples of local spatial frequency representations are wavelets, Wiener filters, scalograms (which are comprised of Gabor filters), and spectrograms (short time Fourier transforms).

Disparity The problem of computing distance from two images can be reduced to that of finding which pixels correspond between the two images. Given a pair of corresponding pixels, the distance between the two cameras and their orientations, it is easy to apply

triangulation to find the distance to the point in world coordinates represented by those pixels. The stereo task is to find the vector offsets between corresponding pixels: this vector is called the *disparity* at a given pixel, and is measured in units of pixels in the image plane. Informally, the disparity tells you how far you must shift a pixel in one image to have it line it up with its correspondent in the other image.

Disparity Image The goal of all stereo methods is the creation of a *disparity image* (also known as a *disparity map*) from a pair of intensity images. Pixels in a disparity image represent not the intensity of incoming light, but rather the amount of disparity assigned to a pixel in the original image pair. Unfortunately the range of disparities can sometimes be very small, so for purposes of visualization all the disparity images in this thesis have been scaled in intensity to bring out as much structure as possible. Smaller disparities indicate objects that are far from the camera, and are rendered with dark grey intensities. Pixels that are occluded (or for which no disparity is known) will be rendered in black. Objects that are nearer have light grey intensities; there will be no purely white or black pixels representing a numeric disparity. Unfortunately our intensity scaling occurs independently in each disparity image, so it is not possible in general to compare absolute intensities between disparity maps. When such a comparison is important, the disparities from a single scanline will be plotted separately so their numeric values can be seen.

Disparity Error Image When the ground truth is known, it is possible to display not only the raw disparities returned by a stereo method, but also the difference between the computed and actual disparities. This is illustrated in a *disparity error image*, in which pixel intensity indicates the absolute difference between the ground truth and measured values. Lighter intensities indicate better matches, and black pixels indicate an absolute error of two or more pixels. Unlike the disparity maps, pixel intensities here span the whole range from white (no error) to black (two or more pixel error). Areas in which ground truth is either not available or not appropriate are rendered in white.

Evaluation Function The function that measures the difference between two small areas from a pair of images is called the *evaluation function*. It gives the pair of areas a

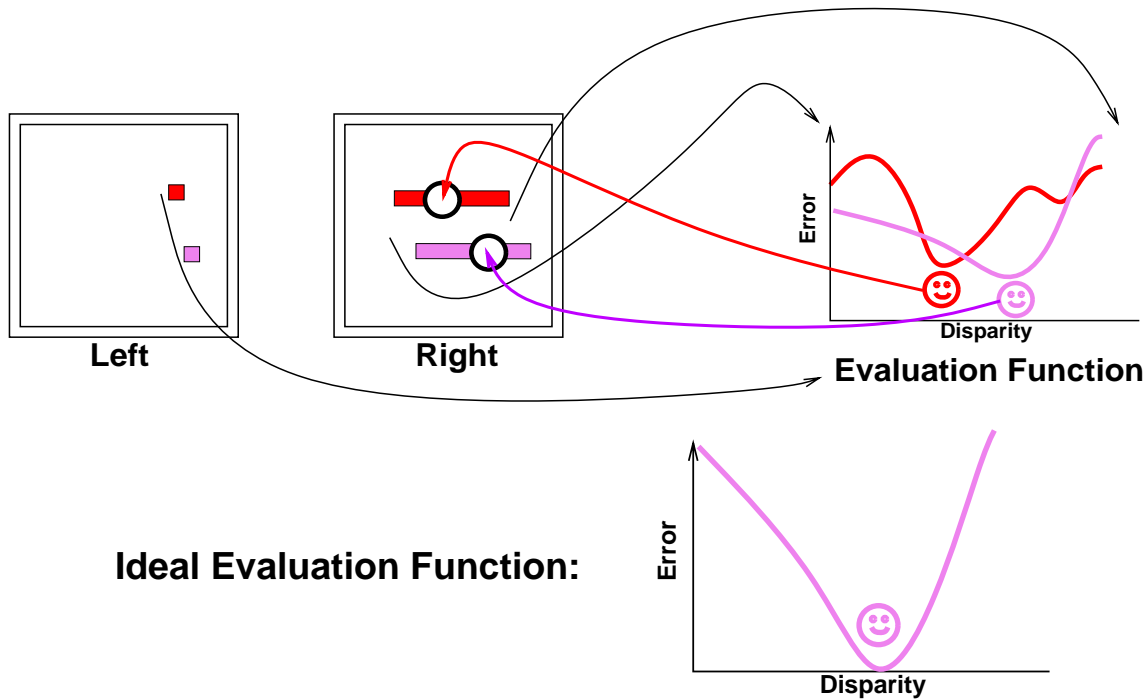


Figure 1.2: Evaluation functions are the tools used to find corresponding pixels. Given a pixel in the left image, a set of pixels in the right image is evaluated. The right pixel that exhibits the smallest error in the evaluation function is selected as the correspondent.

numeric value that indicates the likelihood that the two areas represent different images of the same 3D object. Some common evaluation functions include cross correlation, sum of squared differences (SSD), and sum of absolute differences (SAD). The best match between areas is indicated by larger values for cross correlation, but by smaller error values for SSD and SAD.

Evaluation function profile Since each pixel in one image will be compared with a number of pixels in the other, a visualization of the results of these comparisons can be helpful; see Figure 1.2. The *evaluation function profile* provides this information for a given pixel. The x -axis is Disparity, and the y -axis is the result of applying an evaluation function to the pair of pixel areas (i.e., this pixel and the one that corresponds with disparity x in the other image). The evaluation function profiles that appear in this thesis all use a SAD-like evaluation function, in which smaller values indicate a better match. Thus the y -axis label will be “Error” to reflect the fact that larger values

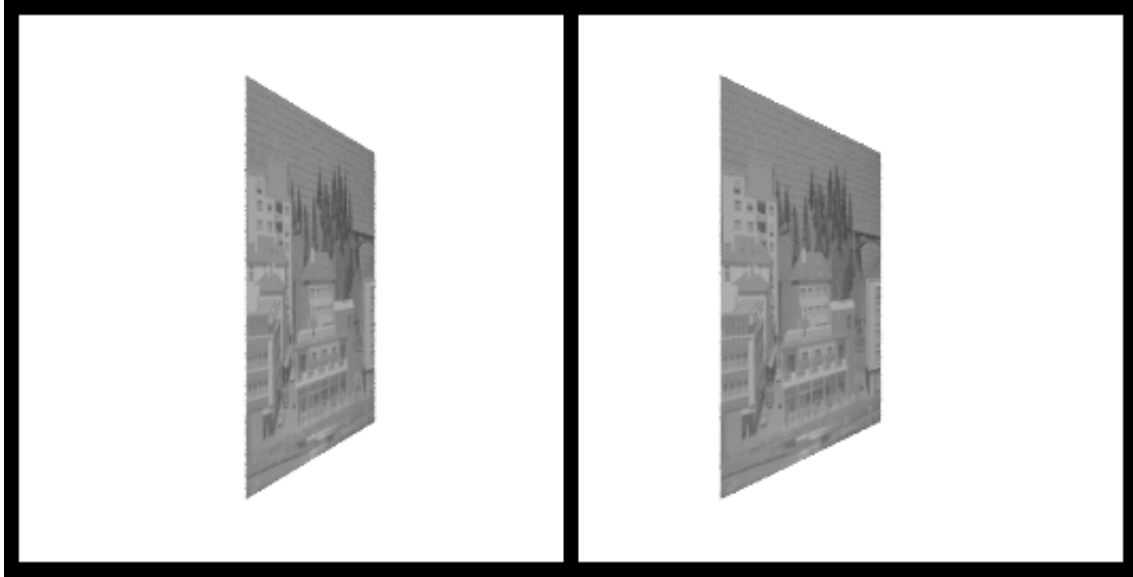


Figure 1.3: Stereo pair illustrating the effects of foreshortening; image compression, differing spatial extents.

indicate worse matches, i.e., matches having greater error.

Window Size The core calculation in all stereo methods is the comparison of a small area in one image with a similar area in the other. The width of this area is known as the *window size* because the comparison occurs only through this small “window” on the original image. While this region might be as small as a pixel or as large as the entire image, often it is a rectangular region several pixels wide. Although windows can have height as well as width, for much of this work we will be concerned with windows that are only 1 pixel tall.

Several factors encourage designers of area-based stereo algorithms to use small search windows. Smaller windows not only reduce the computational cost of a method, but also improve its accuracy by reducing the number of pixels affected by the blurring that inevitably occurs near depth discontinuities, since the search windows are typically run over entire images without regard to an initial segmentation. They cannot be *too* small, however; trying to match single pixels between images is notoriously error prone, since lighting effects and measurement errors often cause corresponding pixels to have different intensities in the stereo images. Even when the corresponding pixels have similar intensities, very often *many* nearby pixels share similar intensities.

Foreshortening

Object surfaces are rarely viewed head-on in both images of a stereo pair. Instead, they may appear more compressed in one image, due to perspective foreshortening, as in Figure 1.3. When a surface has a textured appearance, this effect makes matching its two images very difficult, since its appearance differs so much between the two images. This leads to confusing results from area-based stereo matching techniques, because the visible areas vary so much between the two images. In Chapter 5 we develop a model of perspective foreshortening that enables us to quantitatively predict its effect on stereo image pairs. We present two equivalent forms of a correction factor that allow us to reason about foreshortening effects in both 3D world coordinates and 2D image coordinates. We show how to improve the accuracy of phase-based stereo matching systems using this information, and demonstrate its application to a particular Gabor filter-based stereo system. Applying the correction factor to this system increased its maximum matchable surface angle from 30 degrees to over 75 degrees.

1.2 The Scalogram: A Unified View of Scale Space

We will use a model of Local Spatial Frequency to address long-standing problems in stereo vision. The scalogram is one of several representations that makes use of the local frequency content of an image.

Why Local Frequency? Many of the problems in traditional stereo arise from its limited image representation. Certain imaging phenomena are more succinctly described (and more easily manipulated) in the Fourier domain than in the spatial domain. For example, a sinusoidal pattern at any scale can be fully described by four values in the Fourier domain: amplitude, frequency, phase shift, and a constant (DC) offset. The power of the Fourier transform is its ability to extract this information from any signal in a straightforward and deterministic way. You can think of an image as expressing a function as a sum with delta functions forming the basis, and the Fourier representation as representing the same function but with sinusoidal basis functions. The problem with using the Fourier transform directly is that it extracts frequency information contained *everywhere* in the image; you may find the precise frequencies present in a signal, but you won't know *where* in the signal those frequencies occur. This is unacceptable for image matching.

However, we can compromise by applying the Fourier transform not to the entire image, but rather to a small subset, or *window* of the image. By restricting attention to the parts of the image immediately surrounding a given pixel, we can learn about the *local* frequency context, the patterns present only at that pixel. There is a trade-off, however. By restricting ourselves to a small window on the image, we sacrifice the precision with which we can isolate particular frequencies.

There are many local frequency representations: spectrograms (Short Time Fourier Transforms), Wigner-Ville distributions, wavelets and scalograms to name a few. All are similar in effect, but slightly different in structure. The spectrogram uses a fixed window size at all scales and a logarithmic sampling of wavelengths. This can be useful for texture analysis, where you expect to see the same pattern repeated often at small scales, but seems less useful for image matching. The fixed window size means that high frequency results will not be easily localized, and low frequencies may not have enough support. In contrast, the scalogram uses a variable window size, one which is always a constant number of wavelengths long. This makes high frequencies much more localizable, and provides the necessary support for low frequencies. The scalogram is actually a special case of the very general wavelet functions: the scalogram is a wavelet with a Gabor function as the transfer function. Wigner-Ville is a compromise between spectrograms and scalograms, but often contains many cross terms that complicate automated analysis. A general overview of all these representations can be found in (Rioul & Vetterli, 1991), with a nice overview of their application to computer vision in (Krumm, 1993).

1.3 Stereo Vision Algorithms

1.3.1 Stereo Matching

Stereo matching is a useful tool in a variety of applications: robot vision, parts inspection, and aerial mapping to name a few. It is used to acquire knowledge about distance to objects in the world. It is a passive technology, relying on the interpretation of pixel luminance intensities in two or more images to reconstruct 3D information. Active technologies also exist (e.g., laser rangefinders, sonar, and radar), but we will limit our scope to the passive imaging method of stereo matching.

The minimum system requirements for stereo vision are a pair of cameras positioned with

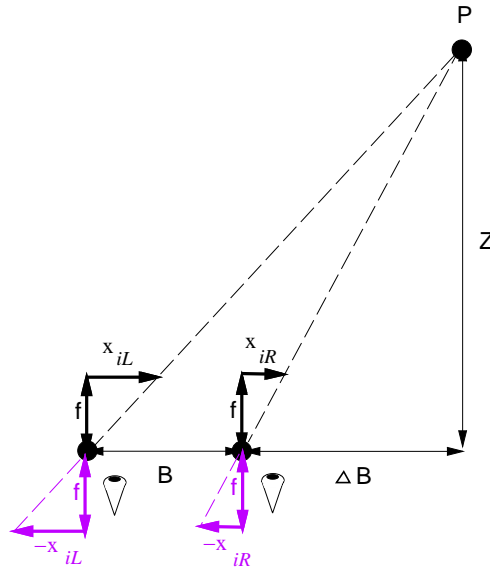


Figure 1.4: Overhead view of a typical Stereo Vision setup. A pair of cameras with focal length f , separated by baseline distance B , are represented by their focus points on the lower left. Object point P has corresponding image coordinates x_{iL} and x_{iR} , and lies at depth Z from the focus points. Mirror image vectors $-x_{iL}$ and $-x_{iR}$ are shown to make the similar triangles used in Equation 1.1 more explicit.

overlapping fields of view. These cameras could be arranged in any number of ways, but to simplify the forthcoming discussion we will restrict our attention to a simple case: both cameras on a horizontal plane with optical and vertical axes parallel, and known *baseline* (distance between them). Thus we explicitly avoid issues dealing with camera vergence (rotation about the vertical axis), torsion (rotation about the optical axis), unknown baseline separation, image rectification, and camera calibration. Such issues can be addressed by other methods (e.g., weak calibration (Robert et al., 1994), Fundamental Matrix recovery, explicit hardware control (Ross, 1993)) and thus can reasonably be avoided here.

A constraint that is often applied to stereo systems is the *epipolar constraint*. The main idea is that given a pixel in the left image, one need not search through the entire right image looking for a correspondent. Instead, attention may be limited to a straight line, the so-called epipolar line. Why is it a line instead of some 2D region? The object represented by a pixel in the left image must lie on a ray that extends in the world from the left focal point through that pixel on the left image plane. The epipolar line is the image on the right image plane of that ray in the world; the projection of a straight line is still a straight

line. A nice introduction to epipolar analysis can be found in (Boufama & Mohr, 1995), and a very detailed description in (Faugeras, 1993). When the optical axes are parallel, as we have assumed here, we enjoy the further property that epipolar lines are guaranteed to be horizontal, and therefore all correspondents will lie on the same-numbered scanline in the other image.

The primary task in stereo matching is to locate pairs of pixels that are images of the same point in space. Once a correspondence has been established, it is a simple matter to determine the distance to that point using triangulation. We can derive the exact relation by considering the overhead view in Figure 1.4. Similar triangles give us two equations relating the pixel indices x_{iL} and x_{iR} to depth Z :

$$\frac{\Delta B}{Z} = \frac{x_{iR}}{f} \qquad \frac{B + \Delta B}{Z} = \frac{x_{iL}}{f} \qquad (1.1)$$

Solving both of these for ΔB and setting them equal, we obtain the canonical expression relating horizontal disparity to depth:

$$Disparity = x_{iL} - x_{iR} = \frac{Bf}{Z} \qquad (1.2)$$

Equation 1.2 gives pointwise disparity only; we will show how to extend this description to surfaces at arbitrary angles in Chapter 5.

1.3.2 Problems with Traditional Methods

In spite of the strong physical constraints available to stereo methods, many problems remain in their implementation. While the mapping from disparity (i.e., pixel correspondence) to depth is well understood, the automatic extraction of disparity is still subject to error. A complete review of all stereo matching problems is beyond the scope of this work, but we will review some of the basic advantages and disadvantages of the standard approaches here.

In feature-based stereo, it is assumed that the objects in the world have distinguishing features that can be easily extracted by a 2D image processing operator. Features might be objects of a particular color, areas with sharp intensity contrast in one direction (edges) or in four directions (the Moravec operator), or any other easily detectable property. One advantage to the use of features in stereo is that the properties of the 3D object represented

by the 2D features are usually very well understood, and can therefore be located with high (subpixel) precision. Another advantage is a computational savings; instead of searching every pixel in the other image, there will generally be only a small number of features detected in the corresponding epipolar line. Thus fewer comparisons are required, and only those image areas that exhibit interesting structure will be checked.

There are many disadvantages to feature-based stereo, however. Since correspondence is only establishing at a small number of pixels, the resulting depth map will be very sparse. Thus any underlying shape may be difficult to extract. There may be large areas in the images that exhibit only small changes in visible texture, and therefore will not have obvious features; no disparity will be computed in such areas. And features visible in one image might not be found in the other; this can cause nonrobust methods to match features incorrectly, resulting in highly inaccurate disparities. Thus while feature-based matching may be useful in applications where the features are known to be visible and only sparse disparities are required, it is not so useful for extracting dense depth maps.

Correlation-based stereo, on the other hand, can provide completely dense depth maps. Subpixel precision may also be obtained, by interpolating the evaluation function's error values in the neighborhood around the minimum disparity. Some of the feature-based methods' problems due to the incorrect matching of features can be alleviated by increasing the size of the correlation window, to provide a disambiguating context around the features. And while this type of stereo requires much computation, there are already many real-time implementations of correlation-based stereo, so computation time is not a bottleneck.

Unfortunately, although correlation-based stereo provides dense disparity maps, there is no guarantee that the results will be accurate. Some of the issues that arise in picking a window size have been discussed in Section 1.1, but the main problem is that the results at depth discontinuities and occlusion boundaries are generally quite unreliable; the extent of an object tends to "spill over" into the background, or vice versa. Also, use of the same-sized matching window in both images means that objects are assumed to be basically planar and viewed from the front. Finally, traditional correlation-based stereo performs matching based on a single view of the world; more robust methods use multiple views, either acquired using additional cameras or generated from the image pair using filters.

1.4 Related Work

Many types of stereo algorithms have been published in the literature. Overviews of the then-strongest techniques can be found in (Barnard & Fischler, 1982) and (Dhond & Aggarwal, 1989), while more recent results on image registration are described in (Brown, 1992). Existing stereo algorithms from the computer vision literature can be loosely classified under one of three headings: traditional correlation-based stereo, feature-based stereo, or frequency-based (often phase-based) stereo.

In correlation-based stereo (henceforth called “traditional” stereo, e.g., (Mori et al., 1973) and the two image case of (Okutomi & Kanade, 1991)), disparity is computed by fixing a small window around a pixel in the left image, then measuring the correlation or Sum-of-Squared-Difference error between intensities in that window and those in similar windows placed at different locations in the right image. The placement that yields the lowest error gives the disparity estimate. This procedure is applied to the image at successively higher resolutions in a coarse-to-fine manner, restricting the search according to the disparity estimate from the previous (lower resolution) image.

In feature-based stereo (e.g., (Matthies, 1989)), a dense image is converted into a spatially sparse set of features (e.g., corners, edges) which are then matched. This results in a sparse disparity map which must be interpolated to yield disparities at every pixel.

Finally, in frequency-based stereo (e.g., (Sanger, 1988)) the original signal is transformed to Fourier space, and some part of the transformed signal is used to compute the disparity. Often the phase of the transformed signal is used, in any of several possible ways (Jenkin & Jepson, 1994). Since our method evolved from those methods which use the phase component of frequency, we present in Section 1.4.1 a short history of phase-based stereo methods and discuss the improvements each method has made over its predecessor.

Several frequency-based methods do not use phase as the primary matching feature, however. Jones and Malik (Jones, 1991) applied 2D oriented derivative-of-Gaussian filters to a stereo pair, and used the magnitude of the filter responses at each pixel as matching features. Their iterative framework allowed them to address scale space selection, half occlusion, and perspective foreshortening, but their solution requires an initial disparity map that is computed without the benefit of models for these problems. Thus their solution is more of a constraint-based post-processing refinement technique. The cepstral method first presented in (Yeshurun & Schwartz, 1989) and improved upon in (Smith & Nandhakumar, 1996) also

works in the frequency domain, but uses cepstral magnitude as the matching feature.

Another post-processing refinement technique is the Kanade/Okutomi variable-window method (Kanade & Okutomi, 1991). This method addresses the occlusion and foreshortening problems by dynamically adjusting the size of the matching window according to constraints on the local variation of both intensity and disparity. In this way they attempt to avoid the boundary problems that arise when the correlation window encompasses two objects at different depths. While the method seems to work very well at depth discontinuities, we will see in Chapter 5 that its ability to handle foreshortening problems is limited (its depth estimates are approximately correct, but imprecise).

Perhaps the most promising spatial domain approach for two image stereo is that of (Belhumeur, 1993). The method uses dynamic programming to segment each scanline into a number of surfaces, and evaluates the surface matches using a Bayesian probability functional; the segmentation with greatest probability yields the final result, which is a completely dense disparity map. Although computationally expensive, the real power of this method is its use of the hypothesized segmentation; it is much easier to compare two objects, whatever their shapes, if you know their boundaries in advance. His examples show excellent results even in the presence of occlusion, surface creases, and foreshortening, but only involve a small number of separate surfaces.

Another successful spatial domain approach is the use of multiple cameras, sometimes known as multibaseline stereo. First described in (Kanade et al., 1992), the approach advocates using a relatively simple SAD stereo matching algorithm over several image pairs. By incorporating multiple views of the world using known camera calibration, many of the shortcomings of the direct yet simple SAD method are eliminated: e.g., specular highlights are ignored, noisy disparity maps become smoother, and some occluded surfaces become visible. Recent extensions have addressed the issues of occlusion (Nakamura et al., 1996) and real-time stereo (Kanade et al., 1995) more directly, but the current approach still suffers from window effects at depth discontinuities and fails to explicitly address foreshortening.

1.4.1 Evolution of Phase-based stereo

One of the first applications to exploit phase as translation was the Kuglin-Hines Phase Correlation Image Alignment method (Kuglin & Hines, 1975). The method assumes that two images — related by pure translation — overlap in some way, and it is conceptually very

```

function y = kuglin (l, r)
% Compute translation between 2D matrices l and r

pdiff = angle(fft2(l))-angle(fft2(r));
kuglin = ifft2 (exp (i*pdiff));

[ rowmax, rowind ] = max(kuglin);
[ allmax, colind ] = max(rowmax);
y = [ rowind(colind) colind ] - 1;

[ row col ] = size (l);
y = y - (y(1) > row/2) .* [row 0] - ...
      (y(2) > col/2) .* [0 col];
end;

```

Figure 1.5: Matlab implementation of the Kuglin-Hines Phase Correlation Image Alignment method (Kuglin & Hines, 1975). Actual computation of the global translation vector occurs in the first two lines of the function, the final six lines extract the indices of the maximum element in a 2D array. The method actually achieves subpixel precision, but this unrefined index-processing code only returns integer results.

simple: take the image FFTs, subtract phases, and apply the inverse FFT (see Figure 1.5). The result is a matrix whose element with maximum magnitude corresponds to the best-fit global translation. In that sense the technique is similar to image correlation, but phase correlation results in a matrix that has better properties: fewer maxima, and a sharper peak around the principal maximum (Horner & Gianino, 1984). Although it returns a single translation vector, it can still be used in the context of stereo vision by comparing subwindows within a larger image. Such a strategy was adopted by Szeliski in (Szeliski, 1994), in which this method provides initial estimates for his image mosaicing technique.

Although the Kuglin-Hines method succeeds at demonstrating the potential of phase-based methods, it has several important problems when considered as a stereo vision solution. Chief among these is the fact that it fails to account for Fourier coefficients whose magnitudes are so small that the phase is unreliable (see Section 3.4.3). Although the paper mentions

that the phase differences can be weighted by magnitude, they provide no details or insight on how to accomplish this. The following method provides one solution to that problem.

The Windowed Fourier Phase (WFP) method by Weng (Weng, 1993) is similar to the multiple-window use of the Kuglin-Hines method, but uses variable-sized windows, and linear interpolation of phases in regions where the magnitude is too small. The justification for interpolation comes from the idea of instantaneous frequency: if the phase derivative is constant, then a linear interpolation between known phases is appropriate. This method uses far fewer filters within each region ($O(\log n)$ where n is the width of the largest window) than Kuglin-Hines ($O(n)$), and allows the window size to vary so that the measurements made by high-frequency filters can be better localized. The filters are assumed to be band pass, and are applied in a coarse-to-fine manner to handle large disparities. Truly ideal filters would respond at a single frequency, but would require input signals of infinite extent, so in practise band pass filters are the best we can expect.

Unfortunately, this method has several problems. The core of the algorithm assumes that phase is linear, or as he calls it “quasi-linear;” that is, in regions where the magnitude is low, the phase is simply interpolated across the scanline to preserve its supposed linearity. In many real images, phase measurements are simply too nonlinear, and this fact cannot be detected reliably simply by using the fixed magnitude threshold he proposes. Also, the use of a pure box filter as the windowing function results in filters that are *not* band pass, but rather have multiple peaks in the frequency domain. Mention is made of the possible use of Gaussian windows (the so-called WFG), but all the results and further theory rely on the box filter. Results obtained with the WFP have been found to be less reliable than those from other phase-based techniques (Jenkin & Jepson, 1994).

To better approximate the band pass behavior desired in the filter set, Sanger proposed a method using Gabor filters (Sanger, 1988).¹ Unlike the coarse-to-fine method of Weng, Sanger’s method combined independent disparity estimates from each filter simultaneously, with a weighting function designed to eliminate outliers and filter out phases with dissimilar magnitudes.

Although this work predates Weng’s, it uses a better filter set and handles filter magnitude

¹The seminal work of Gabor (Gabor, 1946) showed that there are fundamental limits on the resolution attainable by any local-spatial-frequency filter, and the class of filters he described which achieves the maximum possible space/frequency resolution are now known as *Gabor filters*. The real and imaginary components of these filters are sinusoids modulated (multiplied) by a Gaussian window.

more effectively. Still, it still suffers from a naive integration of filter estimates. For example, while Weng’s range of candidate disparities is limited to half the width of the largest filter, Sanger’s is bounded by half the width of the *smallest* filter. The similar-magnitude constraint is also ill-founded; it assumes that each filter’s magnitude will be similar over the entire extent of each image of the same world object. Filter magnitudes rarely appear equal between two images, and are in any case a less important feature than phase (Oppenheim & Lim, 1981) in this context. Finally, this work assumes that the response of a given filter corresponds exactly to the filter’s peak frequency. But since the filter is band pass, it may in fact be measuring responses from any frequency (or many) within a range.

The work of Fleet, Jepson and Jenkin (Fleet et al., 1991) improved this model by allowing the filter’s frequency to vary, providing additional phase measurement constraints, and employing a coarse-to-fine strategy (similar to that later used by Weng) thus eliminating the smallest-filter restriction of Sanger’s method. Instead of assuming that measured frequencies correspond to the filter’s peak frequency, the local phase derivatives from the two images were combined, the assumption being that each filter measured a single frequency at some point in its band pass region. The additional phase constraints also did a better job of filtering out “unstable” phase measurements, thus improving the resulting disparity estimates.

However, this work failed to address the control-strategy problems produced when phase estimates at coarse scales are unavailable. Presumably a simple linear interpolation of the surrounding disparities was used to fill in regions where all the phase information from the $O(\log n)$ filters was not available (i.e., the Gabor filter outputs failed to satisfy Equation 3.10). And while the instantaneous frequency improved the resulting precision somewhat, it was used in a fairly *ad hoc* manner, simply averaging frequency estimates from the two images.

A much more sophisticated solution was developed by Xiong (Xiong, 1995). Instead of assuming that each filter output represents a response at a *single* frequency, he models the response of each Gabor-like filter as a polynomial. Various properties of his so-called *hypergeometric* and *moment* filters make them a much better representation for high-precision matching, since they are able to model directly the effects of finite-width windows and small amounts of perspective foreshortening. Xiong shows how the constraints developed by Fleet *et al* represent a constant approximation to his more general representation. Use of these filters that account for the window effect and local foreshortening yields results that can be

orders of magnitude more precise than those of other methods.

However, like the other methods mentioned above this too is limited to a small range of disparities, being bounded by the filter width. This framework also only allows for a small amount of perspective foreshortening, since it compares phases using the same filters from each image. And finally, like most other stereo methods, this method uses a limited model of disparity that is unable to represent the possibility of an ambiguous match. These topics will all be addressed in the forthcoming chapters.

1.4.2 Foreshortening

Local spatial frequency has already been identified as a valuable tool for modeling surface shape and segmenting multiple textures in a single image (Krumm, 1993; Malik & Perona, 1989). These approaches use filter magnitude in the frequency domain as the feature of interest, and require either that the surface textures exhibit specific properties (e.g., periodicity), or that they be viewed directly head-on.

Local spatial frequency representations have also been successfully applied to optical flow problems (Barron et al., 1994; Xiong, 1995), using phase information as well as magnitude. The stereo problem is more constrained than optical flow, so we should be able to do better by taking advantage of the additional constraints.

There have been a few attempts at modeling foreshortening in the context of stereo matching. Jones and Malik (Jones & Malik, 1991) applied local spatial frequency to the problem, but used an affine transformation matrix in the *spatial* domain without providing a description of its effect in the frequency domain. Belhumeur (Belhumeur, 1993) addressed the problem in the spatial domain, but his method requires an estimate of the disparity derivative, an inherently noisy estimator. The variable window method of Kanade and Okutomi (Kanade & Okutomi, 1990) implicitly addresses foreshortening in the spatial domain by allowing corresponding windows to have different sizes, but is intended to function as a high-precision refinement technique: without proper guidance from other sources it tends to get stuck in local minima and flatten out sloped surfaces. The cepstral method of (Smith & Nandhakumar, 1996) is claimed to be robust in the presence of foreshortening, but only slants up to 43 degrees have been demonstrated.

Several phase-based stereo methods have been described in the literature (Fleet et al., 1991; Sanger, 1988; Weng, 1993), and a review of the more popular variations can be found

in (Jenkin & Jepson, 1994). Although some of these mention foreshortening as an issue, none has explicitly modeled it or corrected for it.

1.5 Relationship to Prior Work

Some of the results in this thesis have appeared (often in compressed form) in other papers. Chapter 2 is a revised version of (Maimone & Shafer, 1996) and includes more information about the CMU Calibrated Imaging Lab (CIL) Stereo Datasets with Ground Truth that have been made available on the Internet.² The specific procedures for using CMU CIL facilities to acquire such data are presented in some detail in (Maimone, 1995). Earlier versions of Chapter 5 and portions of Chapter 3 appeared in (Maimone & Shafer, 1995a) and (Maimone & Shafer, 1995b), but the sample image pair was unreadable, the Internet URL has changed, and the pseudocode in Table 2 of the former reference was missing an important term. Additional information is available on the Internet via the *Calibrated Imaging Lab home page*,³ *Mark Maimone's index page*,⁴ and the *Computer Vision home page*.⁵

1.6 Thesis Organization

Chapter 1 has introduced the problems that inspired this thesis, given background to provide the context for understanding them, and discussed the related work. In Chapter 2 we describe a general framework for the acquisition of high-precision datasets and independent validation of stereo results, and discuss the particular datasets used in this thesis. Chapter 3 presents the concept of local spatial frequency and outlines the phase-based stereo method that provides the foundation for later discussion and analysis. Chapter 4 addresses the problem of ambiguous matches, and Chapter 5 that of perspective foreshortening. Finally, Chapter 6 lists our contributions and summarizes their impact.

²<http://www.cs.cmu.edu/~cil/cil-ster.html>

³<http://www.cs.cmu.edu/~cil/cil.html>

⁴<http://www.cs.cmu.edu/~mwm/>

⁵<http://www.cs.cmu.edu/~cil/vision.html>

Chapter 2

A Taxonomy for Stereo Computer Vision Experiments

*The beauty of truth, as of a picture, is not acknowledged
but at a distance.*

— Joseph Glanvill, *The Vanity of Dogmatizing XV*

Much of computer vision research is an attempt to solve the impossible: to acquire full three-dimensional knowledge given limited two-dimensional data. The state of the art has advanced to a point where there now exists a plethora of partial solutions to computer vision problems. We're getting lots of answers, but just how accurate are they? A few methods provide an estimate of uncertainty with each answer, but those uncertainties do not tell us what we *really* need to know: by how much does the estimated answer differ from the truth?

If computer vision is to become more of an engineering discipline than craftwork (Shaw, 1990), engineers must be able to predict and experimentally characterize the behavior of their systems. Such characterization is only possible when ground truth is available. Synthetically generated imagery gives total ground truth knowledge, but fails to model the complexities of real-world imaging. Real imagery provides better test data, but greatly reduces the amount and density of ground truth available for analysis. In this chapter we outline a framework for choosing a reasonable trade-off of ground truth density v. image realism in the analysis of stereo algorithms.

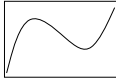
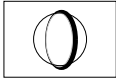

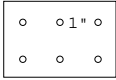
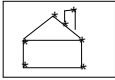
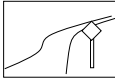
						
	Math	Noiseless Synthetic	Synthetic with Noise	Controlled Environ.	Measured Environ.	No Controls
Fundamental Principles						
Occlusion Maps						
100% Precise Disparity						
Sparse Disparity						
Complex Scenes						
Image Noise						
Real World Noise						
•						
•						
•						

Table 2.1: Stereo Imagery Characteristics and the types of scenarios for which they are available.

In performing the research described in this thesis, especially in Chapter 5, it was determined that the only way to effectively judge its success would be to use stereo imagery with precise ground truth. Unfortunately, very little such data was publicly available, so it became necessary to create our own datasets. We found that there were several distinct scenarios through which ground truth could be acquired, each serving an important function in the evolution of the stereo algorithm. What follows is a description of these various scenarios and the types of situations in which each is most appropriate.

2.1 Overview

The characterization of a vision algorithm's performance is a tedious and difficult task. This is primarily due to the lack of appropriate performance measurement specifications. This chapter will enumerate the scenarios and types of data that can be used, give examples of

interesting properties of stereo vision systems, and explain how best to measure them, in the sense of achieving the most accurate and representative ground truth.

Table 2.1 illustrates some interesting properties of stereo vision data, and matches them with the scenarios in which they can be measured. Because many properties can be measured using any of several scenarios, only a few are likely to be needed for a particular analysis set. The list of rows in Table 2.1 is merely suggestive; other properties can be easily matched with the appropriate scenarios.

A complete statistical framework to take advantage of this data is beyond the scope of this work (see (Matthies, 1989) for an approach assuming unimodal disparity error). We will instead discuss the representation of the data and ways to acquire it, with the understanding that its most effective use will be in the context of a complete analysis package. Such a package would measure interesting properties over several subimages, e.g., the whole image, occluded regions, non-occluded regions, pixels above a known ground plane, etc.

Although comparisons of stereo algorithms have been done before, many have suffered from a lack of available ground truth. One such study, the ARPA JISCT stereo evaluation (Bolles et al., 1993), compared the results of four stereo methods. However, since ground truth was not available, most of their statistics dealt with agreement between the results; not “method A is 80% accurate”, but “methods A and B agree on 80% of the image”. Thus they could neither evaluate stereo methods independently nor quantitatively characterize their performance. The study’s conclusion states in part that “Ground truth is expensive, but there is no substitute for assessing quantitative issues.”

Even those studies that included ground truth have been limited by available technology. An ISPRS study (Gülch, 1991) compared the results of several stereo algorithms using ground truth, but in all but one image pair each true pixel disparity was computed *manually*, without interpolation. Their manual collection of 23,000 total pixel correspondences from eleven 240×240 image pairs was a noble, not to say Herculean, effort, but we argue that the technology for the construction of ground truth images of complex scenes (with much greater density and less required manual intervention) not only exists, but is easy to use.

We demonstrate this claim by outlining a taxonomy of ground truth scenarios (abbreviated in the column headings of Table 2.1), and providing a concrete example for each level. Sections 2.2 through 2.7 present the scenarios, ordered by increasing realism and decreasing amount of ground truth, Section 2.8 describes the tools used to implement these scenarios

and their availability, and Section 2.9 summarizes our contribution.

2.1.1 Representation of Results

Stereo vision is a powerful approach for computing 3D information, but unlike active rangefinding systems stereo works indirectly, by converting pixel correspondences into depth measurements. The quality and density of the resulting depth map depends directly on the character of these correspondences.

Unfortunately, even exact pixel correspondences alone rarely give a complete picture of the range map. This is due to several effects, e.g., nonoverlapping fields of view in parallel and outward-verging cameras, self and half occlusion of objects in the scene, and a lack of intensity variation in areas with bland texture. It is unreasonable to expect a correspondence-based stereo method to calculate completely dense depth maps, since even perfect pixel correspondences can leave gaps in their implied depth maps.

How then should stereo data be represented? There are several possibilities, principal among them disparity maps, depth maps, and object models. Researchers have reported stereo results using all of these representations, but each has its drawbacks. We would prefer to use a representation that allows different stereo methods to be compared on an equal basis. Object models are quite useful, but can only be computed from raw stereo data by making many model-based assumptions. Depth maps would be ideal, but require exact knowledge of the camera geometries (which may not be available), and cannot be completely computed from correspondences alone. To compare results computed from two images without requiring camera calibration information then, disparity maps with occlusion masks are the most general representation.

One need not eliminate metric information completely when providing disparity maps, however. By simply including the parameters of the actual camera system in the dataset (e.g., baseline and camera focal lengths), disparity can be converted easily to depth when needed. Since depth resolution of a stereo system can be increased simply by adjusting the camera separation and/or focal length, the precision of a stereo method is often measured in pixels (units of disparity) rather than Euclidean length (units of depth). But most applications will express requirements in depth units, so this metric information should still be kept available.

The ground truth must also be expressed as a disparity map. This is accomplished by running the 3D ground truth information through the appropriate camera model, to generate

a depth map in the same coordinate system as that of the image being matched. From here it is a simple matter to generate a disparity map, and we will see in Section 2.3.1 how to compute occlusion masks from disparity maps derived in this way.

The model for stereo experimentation is thus to run the stereo algorithm, and compare the computed disparity map with that of the ground truth known from other methods. The rest of this chapter presents several scenarios and examples of test data with ground truth, explains the benefits and limitations of each, and discusses some implementation issues that arose in the course of generating the sample data. Thus we demonstrate that the technology for creating datasets with interesting properties and dense ground truth already exists, and is easy to use.

2.2 The Basics: Mathematical Foundations

The first scenario is that of continuous functions. When describing an algorithm, the first step should be to demonstrate the principle in the simplest possible domain. In stereo vision, the simplest case is typically the comparison of two 1-D functions that represent scanlines. There is often much insight to be gained by focusing attention to a level of detail in which all quantities can be interpreted directly.

This is the level at which the general principles of an algorithm can be demonstrated. At this point it is not necessary for the inputs to have precisely the same qualities as those present in actual discrete imagery. Indeed, using continuous functions as input can often simplify the presentation by allowing the solution to be expressed analytically (in closed form) rather than operationally (e.g., “the result after 10 iterations”) as in (Fleet et al., 1991; Sanger, 1988).

This level of description is also useful for discerning and describing any theoretical limitations of the method, e.g., the points at which its assumptions break down.

2.2.1 Example 1: Phase difference as disparity

The use of phase difference as disparity lies at the heart of many phase-based stereo algorithms (Sanger, 1988; Weng, 1993; Fleet et al., 1991). But unless one is already quite familiar with the frequency domain (and Equation 3.6 in particular), the name itself inspires fear. Suppose one remained uncertain of the underlying principles; how could you convince

yourself that the technique really works? By considering the simplest possible examples according to the representation, and following the processing step by step.

In this example we are interested in studying how the phase of a sine wave relates to stereo disparity. So consider the simple case in which the left and right image scanlines are both sinusoids. A one-dimensional sinusoid is in general completely determined by three parameters: amplitude (A), frequency (ω), and phase (ϕ).

$$\text{Sinusoid: } \boxed{A} \sin \left(2\pi \boxed{\omega} x - \boxed{\phi} \right)$$

For this demonstration we will fix the amplitude A at 1, frequency ω at $\frac{1}{8}$, phase of the left image at 0, and allow the right phase ϕ to vary freely:

$$L(x_L) = \sin \left(\frac{2\pi}{8} x_L \right) \qquad R(x_R) = \sin \left(\frac{2\pi}{8} x_R - \phi \right) \quad (2.1)$$

Stereo disparity is the amount of shift required to make the left and right images appear equal. While in general the disparity in an image will vary at every pixel, in our example all pixel disparities will be equal (this actually happens in a real image whenever a planar surface is viewed head-on, so it is a realistic assumption for this demonstration). Mathematically, disparity is the difference between the left pixel index (x_L) and the right pixel index (x_R). So we find disparity by setting the formulas in Equation 2.1 equal and solving for this difference:

$$\begin{aligned} \sin \left(\frac{2\pi}{8} x_L \right) &= \sin \left(\frac{2\pi}{8} x_R - \phi \right) \\ \frac{2\pi}{8} x_L &=_{2\pi} \frac{2\pi}{8} x_R - \phi \\ \text{Disparity} := x_L - x_R &=_{8} -\frac{8}{2\pi} \phi \end{aligned} \quad (2.2)$$

Thus we see that disparity is indeed related to the difference of the left and right phases (remember the left phase is zero in this example). Figure 2.1 graphically shows the disparities that result from particular values of the right function's phase. You can convince yourself that the mapping from phase to disparity works by plugging the values of ϕ into Equation 2.2 and comparing the answer with the amount of shift visible in the graphs of Figure 2.1.

This example also illustrates an important issue in the design of phase-based stereo techniques: phase-wraparound. In interpreting Figure 2.1, we knew the phase difference was relatively small (i.e., less than 2π), so the disparity could be computed directly. But the

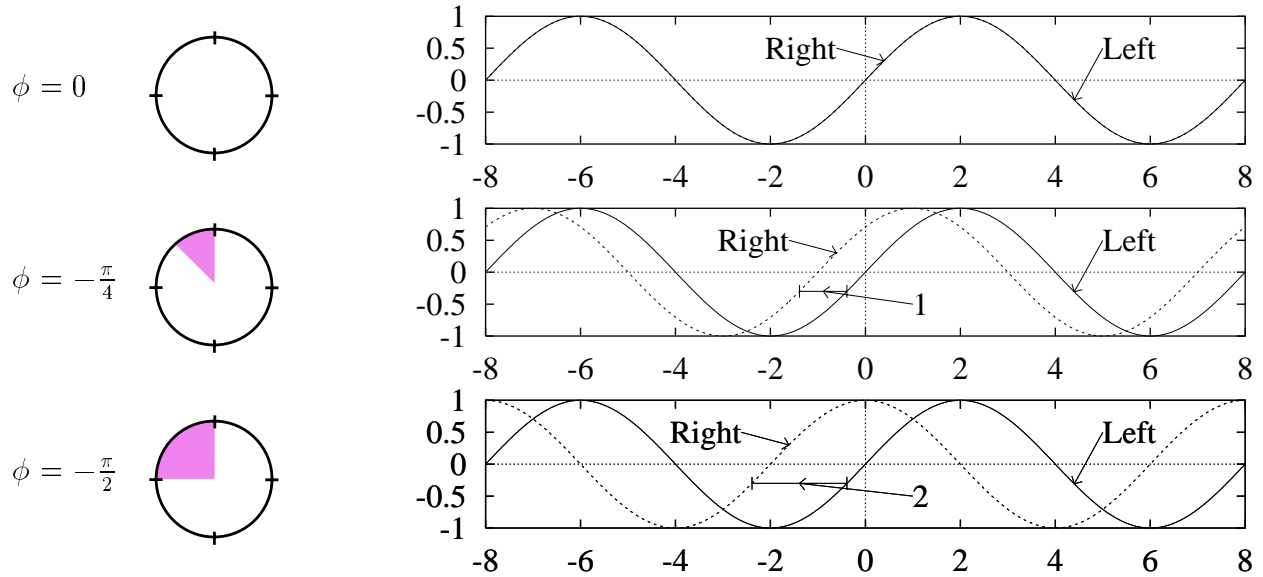


Figure 2.1: Disparity as a function of Phase Difference. Disparity is the horizontal separation between the two signals, and is indicated by the labeled bar just to the left of 0.

disparity formula in Equation 2.2 is only defined *modulo the wavelength of the sine wave*. This means we cannot compute a unique disparity from a single phase value; at this frequency, disparities of 1, 9, 17, ... all appear equivalent. This is an important observation: when evaluating any phase-based stereo method, be sure to consider how it addresses the problem of phase-wraparound. Some authors use a coarse-to-fine approach to alleviate it (Weng, 1993; Fleet et al., 1991), others choose to ignore it (Sanger, 1988) thus restricting themselves to finding only small disparities. Our solution combines phases from several filters at arbitrary frequencies simultaneously, so instead of wrapping around at the wavelength of the smallest or largest filter, our estimates wrap at the least common multiple of all wavelengths that comprise the signal (this is typically larger than the size of the image, and thus tends to yield a unique result). More detail can be found in Section 3.6.1.

This simple, direct analysis has resulted in several important insights: an understanding of the basic technique, and an appreciation for an important property of all phase-based stereo methods.

Limitations: While this level of demonstration is important in communicating the intuition behind an algorithm, it has many limitations. The most obvious is that when images are processed only discrete samples are measured, since images have a fixed resolution. So this general description must be restated in more concrete terms that take the limited

resolution into account. Also, while a purely mathematical scene description is easy to reason about, complex scenes would require such detailed modeling that constructing a continuous version would be too cumbersome.

2.3 Noiseless Synthetic Imagery

The next scenario is that of discretely sampled synthetic imagery. Having established basic principles using continuous functions, the generalization of the method to complete images must be characterized. The broadening of attention from minute pixel-level details to those encompassing entire objects can also yield important insights. For these purposes synthetic data prove most useful. Such data may also be used to verify the implementation of an algorithm on full-sized images.

Synthetic images can be generated by any means, but should initially be created according to a model of the imaging environment in which an algorithm will be deployed. This model should come as close as possible to approximating the real world, though for the moment a noiseless environment should be assumed. In particular, objects being “imaged” should have 3D structure, and should be rendered using the same model as that assumed by the algorithm. For stereo vision this model will typically include full perspective projection, e.g. the pinhole lens model. Use of computationally simpler imaging models such as orthogonal projection or linear (affine) warping, should be avoided except in providing data for debugging an algorithm that makes that assumption. This task is not as difficult as it might seem: the Computer Graphics community has developed many realistic renderers, some of which are freely available and easily modified (Haines, 1993).

The camera model used by the synthetic image generator should be as similar as possible to that used in the actual laboratory camera calibration. Using the same model makes it possible to close the loop; 3D info computed from real imagery can be re-rendered from the same (now virtual) camera position. Figure 2.5 shows how useful this can be: the 3D locations of dots in a real image of a calibration grid are rendered according to the computed camera model, and overlaid on the real image. Such tools can make inspection and validation of 3D reconstructions much easier.

Perfect ground truth is also quite helpful in debugging. While this may seem a trite truism, the difficulty of developing image processing software and the current lack of in-

egrated matrix debugging environments have discouraged vision software developers from adopting this approach of using full-sized image ground truth. Yet freely available software can provide arbitrarily complex test cases that can help the debugging process tremendously. The following example demonstrates how complete ground truth pointed out the flaws in a common technique for generating occlusion masks.

2.3.1 Example 2: Occlusion Masks

Depth maps in complex scenes are typically discontinuous, and therefore difficult to reason about analytically. Also, when opaque objects are imaged from different viewpoints, portions of those objects will be visible only in one image. For these reasons it becomes important to provide occlusion masks with stereo ground truth: correspondence-based stereo algorithms cannot be expected to predict accurate disparity in areas where no correspondence exists.

Occlusion masks are binary images that indicate which pixels are visible in both images of a stereo pair. They are defined relative to a pair of viewpoints: screen pixels in one viewpoint are occluded if the 3D point they represent is not visible from the other viewpoint. How can occlusion masks be generated? One method that works nicely with synthetic data is to simply place a point light source at the focal point of the second viewpoint and re-render the scene (Williams, 1978), taking care to turn off interreflections, translucency, and other light sources. Pixels in shadow are then marked as occluded. However, if the pixel subsampling should differ between the image and occlusion mask renderings, border pixels may be labeled incorrectly.

Another method works directly with the depth map used in the construction of the synthetic image. The top row of Figure 2.2 shows a sample stereo pair of images and their disparity maps. A popular technique in stereo algorithms computes occlusion masks by performing a pointwise comparison between the disparity maps for the left and right images. Any corresponding pixels whose disparities are not equal in magnitude and opposite in sign are marked occluded by this method (Jones, 1991). While this is a useful approximation, it often fails at object boundaries because of steeply-sloped surfaces, as illustrated by the results in the middle row of Figure 2.2. Since this algorithm produces noisy results even with absolutely correct disparity maps, a more robust approach is clearly needed.

There are two main problems with the pointwise method: sharply sloped surfaces may cause the corresponding pixels to point to the same surface but at very different depths, and

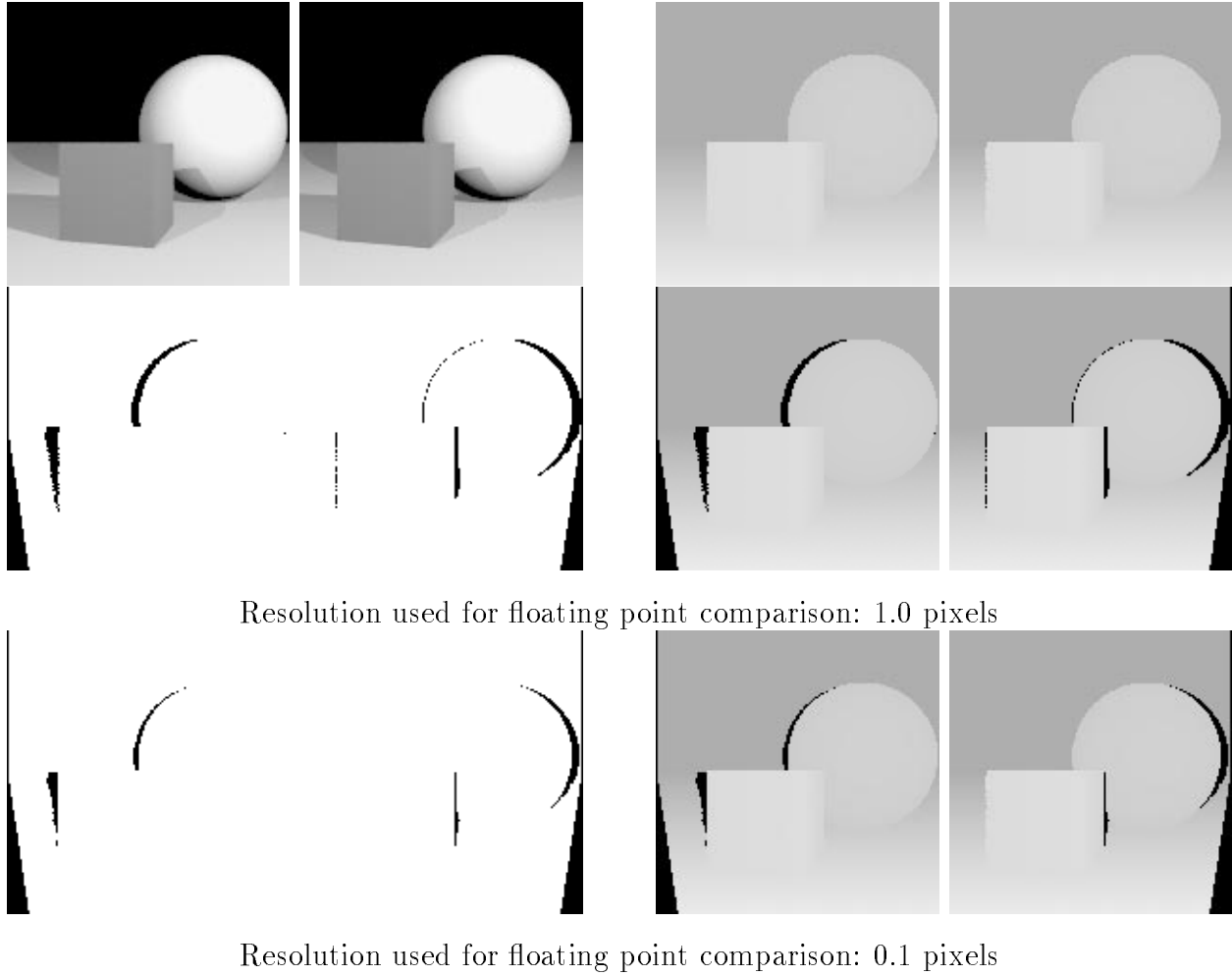


Figure 2.2: Occlusion Mask Generation: *Top Row*: Stereo pair of images from the Left and Right cameras; actual disparity maps for those images. *Middle Row*: Pointwise occlusion masks for left and right images (note especially the noise in the right mask around the object borders); disparity maps with pointwise occlusion masks overlaid. *Bottom Row*: Plane-fitting occlusion masks for left and right images; disparity maps with plane-fitting occlusion masks overlaid.

the pixel subsampling strategy may cause border pixels to point to the wrong object entirely. We can address these problems by fitting a plane to each pixel in the disparity map. If each scene object has an extent of at least two pixels, then it is reasonable to assume that for a given pixel, the 2×2 surrounding window with the least variation in depth will be the appropriate surface patch (see Figure 2.3). We compute each pixel's best plane by finding that 2×2 window which has the least variation in depth (i.e., for which $\max_{2 \times 2}(\text{disparities}) -$

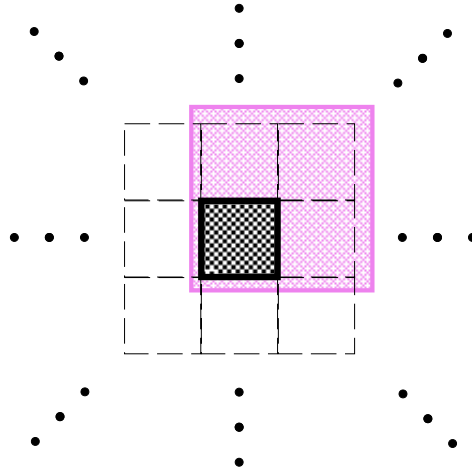


Figure 2.3: Finding the best-fit plane in the disparity map. There are four 2×2 windows that contain the center (dark) pixel; the upper right window is highlighted.

$\min_{2 \times 2}(\text{disparities})$ is minimized). This window selection therefore helps avoid the pointwise method's error of pointing from a slanted surface to the wrong object. (Grant, 1992) provides a good overview of shading techniques in rendering applications. Our approach is most similar to his P2-buffer algorithm, which also fits a plane to adjacent depth values.

The problem of varying depth on a surface is addressed by using the *range* of disparities found in the 2×2 surrounding window selected above. Instead of comparing disparity values directly, the disparity in one image is compared to the range of values contained in the 2×2 window surrounding its correspondent. In practise, we found that extending the measured range by 100% on either side allowed us to increase the floating point resolution of the comparison from 1.0 pixels in the pointwise method to 0.1 pixels; some results are shown in Figure 2.2. That is, two values are considered equal if they differ by at most this amount.

It should be possible to extend this notion to arbitrary viewpoints. The key point is to perform the range check using *normalized inverse depth*. When the optical axes are parallel, this simply means using the disparity. Given two arbitrary viewpoints, the same term Bf/depth can be used, though it will no longer represent the magnitude of the disparity. Converting the depth values from one viewpoint to another will require the application of the complete camera models from both views (not just the linear translation/rotation matrix), but this should pose no problem with synthetic data since both camera models are known exactly and an analytic inverse of the perspective projection is easily coded. This use of arbitrary viewpoints is important for many extensions to the standard stereo model

of parallel optical axes: e.g., verged cameras, multibaseline stereo, and arbitrary motion between frames.

In summary, the use of synthetic data enabled the elucidation of noise-causing effects in a popular stereo method of generating occlusion masks. We also showed how those effects could be mitigated by using a new method of fitting planes locally, and how this method can be extended to a completely general 3D case with arbitrary viewpoints.

This problem is also important to the Computer Graphics community, where occlusion masks are used as shadow maps. However, their goal is to find the amount of incident light over a large area, not to find individual point matches. Most shadow map methods simply smooth over pointwise matches to even out the shading effects (Williams, 1978; Reeves et al., 1987).

Limitations: Though useful for validating software integrity, noiseless synthetic imagery cannot be used to measure robustness to noise, an inevitable problem with real imagery.

2.4 Synthetic Imagery with Noise

A good way to characterize the robustness of a method is to take known data and introduce noise along independent dimensions of the imaging model (e.g., gaussian or focus distance blurring, white noise, lighting intensity, camera misalignment). Using synthetic imagery, the amount of error introduced along each dimension can be quantified precisely, and the degradation of an algorithm according to a particular type of noise model can be determined.

This type of data can be quite useful in experimentally characterizing an algorithm's performance, providing engineers with a quantitative measure of robustness. The ability to track performance loss as a function of noise is an important parameter in engineering design.

2.4.1 Example 3: Virtual Checkerboard

Figure 2.4 illustrates some types of noise that can be easily modeled using synthetic data. The top row contains perfect data, computed as described at the beginning of Section 2.3. As before, we have a perfect disparity map with the occlusion mask in black (in this figure it is computed for the rightmost image). The middle rows illustrate how typical noises affect image formation; here the left image is shown modified according to each particular type

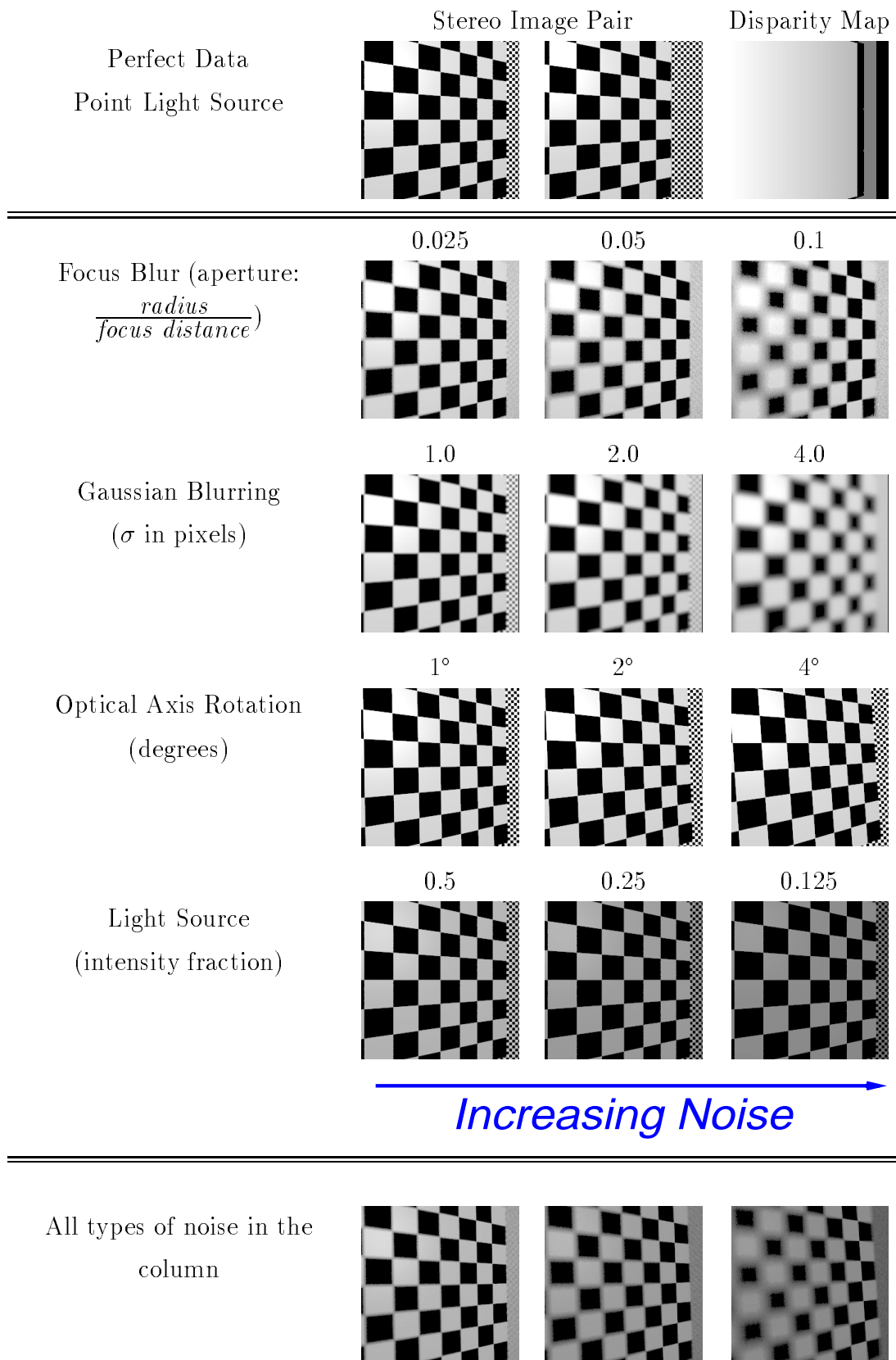


Figure 2.4: Examples of noise easily modeled with synthetic data.

of noise. When actually performing experiments both images will be modified, and the original disparity map used as ground truth. The effect of each type of noise is quantified by comparing the disparities computed from each noisy image pair against the known (constant) ground truth and measuring residual error. In this way the robustness of the algorithm to various types of noise can be estimated, and statistics derived by applying the same noise models to many synthetic datasets and tracking overall error. These noise effects can be combined arbitrarily; the bottom row of Figure 2.4 shows some examples.

Limitations: While these types of experiments are useful for calculating robustness along particular dimensions, they will nevertheless fail to address all the possible effects of real-world imaging. Multiple extended light sources, complex interreflections, nonlambertian objects, complex shapes, and errors introduced in the imaging process such as sensor pixel noise are extremely hard to model precisely, yet they affect every image. Synthetic data are useful for simple validation and characterization, but for algorithm performance verification there is no substitute for actual imagery.

2.5 Controlled Environment

The most useful stereo datasets are those with real imagery and 100% ground truth. Having used synthetic data to establish correctness and characterize robustness along particular model dimensions, one can move on to real images of controlled scenery. This introduces many unmodeled errors in the lighting, camera and optics, but allows them to be characterized by highlighting differences between the disparity map known from the ground truth and that computed by the algorithm.

Several types of noise are introduced here that are typically not modeled in synthetic imagery. The CCD array is subject to preamplifier noise, dark current, shading effects, and photon noise (Photometrics, 1990). Optical effects also become apparent: radial and tangential lens distortions, poor overall focus due to lens manufacturing errors and dust, and misalignment of the lens with the CCD array. Perhaps most importantly, effects such as the interaction of complex light sources with the objects being imaged and the sheer complexity of actual scenes introduce artifacts into real imagery that must be treated as noise by systems that fail to model them.

Some measurement errors can be compensated for in preprocessing. For example, one

common problem with stereo imagery is a difference in gain between the two cameras. The problem is manifest as very different brightness, or distributions of pixel values between the two images. This effect can arise from many causes, e.g., differing apertures on the two lenses and specular highlights on the objects. Yet a simple histogram equalization can bring these distributions closer and make matching easier.

To acquire this type of data the shapes of the objects being imaged must be precisely known. This information ideally will be acquired using methods other than vision, since our objective is to evaluate the quality of a vision-based reconstruction. This can be accomplished by machining an object to precise specifications (as is often done with calibration targets), or by measuring the dimensions of existing objects with known shape (e.g., the sphere in (Jenkin & Jepson, 1994)).

Camera calibration is also a requirement for collecting accurate ground truth. Even if the shapes of objects in the scene are known, their absolute distance and orientation will be unknown. Camera calibration information should be computed even if the stereo algorithm being evaluated does not make use of it, to aid in the computation of dense ground truth. Euclidean camera calibration will enable quantitative analysis of the algorithm's precision.

2.5.1 Example 4: Calibration Targets

A common application of totally structured environments is the acquisition of camera calibration parameters, as will be described in Section 2.5.2. Most camera calibration techniques depend on the reliable extraction of 2D feature points from images, and many require precise 3D localization of features in the world as well. Since industrial applications often allow the scene environment to be manipulated during calibration, a common technique is to construct (and position) a target using these constraints:

1. The target must contain many feature points,
2. Those points will be spread throughout most of the camera's field of view when imaged,
3. The points can be easily located by simple image processing techniques, and
4. The 3D locations of the feature points are known to a high degree of precision.

Figure 2.5 shows a sample calibration image used in the Tsai-derived calibration procedure developed by Willson (Willson, 1994), and used in this work as well. The target is a grid

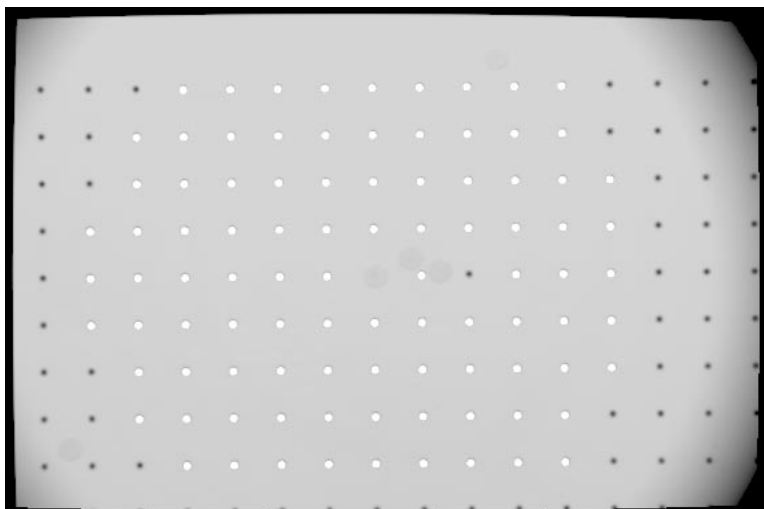


Figure 2.5: Actual image of the Calibrated Imaging Laboratory (CIL) calibration target with virtual rendering overlaid. The grey background and grid of black dots are part of the original picture, the white dots are rendered dots located at the 3D grid point locations. The dots were rendered as spheres using a virtual camera with the same parameters as those computed from the real image.

of black dots on a flat white background, each dot one inch from its horizontal and vertical neighbors. The center dot has been whited out to provide a reference spot for interimage registration. This flat target is mounted on an accurate translation stage, so by imaging it at several locations many 3D calibration points may be acquired. Because the 3D structure is known, the 2D extracted features are easily mapped into a 3D representation. To demonstrate how well the 3D representation fits with the original data, some of the recovered 3D feature points have been rendered as white spheres using the computed camera model (via a computer graphics ray tracer), and overlaid onto the original image in Figure 2.5. This projection back into 2D agrees nicely with the original data.

Figure 2.6 shows some more calibration targets. While CMU's Calibrated Imaging Laboratory (CIL) target from Figure 2.5 must be imaged at several locations to sweep out a 3D volume, these other targets have inherent 3D structure. Shum's cube (Shum et al., 1995) (leftmost in Figure 2.6) is useful because its regular structure means it can be imaged from any angle and still provide a number of features. The disadvantage is that this regularity implies there will always be some ambiguity when matching feature points from different cameras, especially if the camera separation is large and the optical axes all intersect within

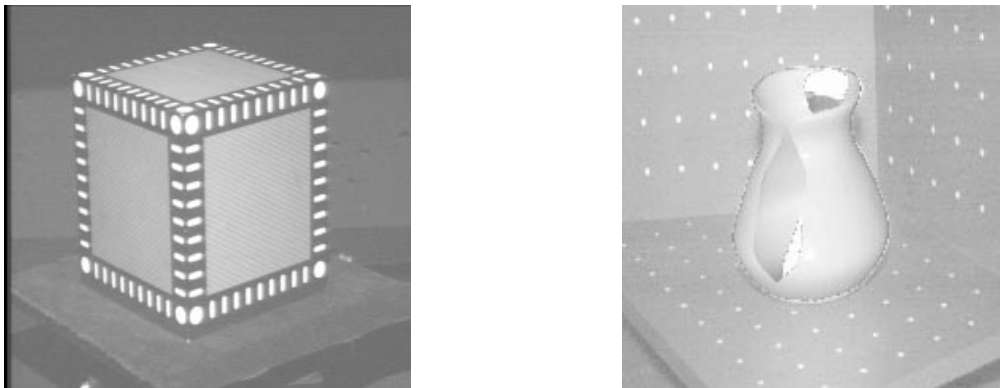


Figure 2.6: More sample calibration targets. The left image is the calibration cube from (Shum et al., 1995), right is an image of the MOVI (INRIA) “inverted cube” calibration pattern from (Boyer & Berger, 1995).

the target. The MOVI “inverted cube” target (Boyer & Berger, 1995) on the right suffers no such ambiguity, because its faces are joined at different angles (the vertical faces are joined at approximately 120 degrees). It is an “inverted” cube because it resembles an office corner: the angles between its visible faces are closer to 90 degrees than to the 270 degree angles found on a cube. A useful property of this target and the CIL target is that either may be left in place during later experiments and can thus provide useful ground truth for the background pixels.

Construction of calibration targets can be a difficult and expensive task. The MOVI target in Figure 2.6 is a set of precisely-machined metal plates, and the CIL target in Figure 2.5 is a PostScript file printed in high resolution on laminated and self-sticking paper, mounted on posterboard and attached to a metal frame on an automated translation stage.

Having seen examples of real objects built to precise specifications, we now outline the procedure required to put them to good use.

2.5.2 Camera Calibration

The task of determining in general how 2D image points correspond to 3D scene points is accomplished by assuming a camera model and performing camera calibration. First the correspondence between a representative set of 2D and 3D points is determined, then the parameters of the chosen camera model that best fit those data are found.

Typical camera models comprise external (or *extrinsic*) and internal (or *intrinsic*) param-

eters. Some common external parameters are X, Y and Z axis translations and rotations, and common internal parameters include image center, thick/thin/pinhole lens focal lengths, lens distortion coefficients, and aspect ratio. In astronomical image processing, lens models must also compensate for specific manufacturing defects by imposing a dense grid of coefficients over the lens surface, to model the point-spread function at each pixel or small group of pixels. Computer vision lens models tend to be much simpler, requiring only a few parameters. This is largely due to the fact that in computer vision, objects can be moved close enough to the camera to compensate for the small-scale lens defects that astronomers are forced to model.

Some calibration techniques require 3D information (Tsai, 1987; Willson, 1994; Bani-Hashemi, 1993), some use correspondences between images to compute the stereo epipolar geometry without complete Euclidean knowledge (Zhang et al., 1995; Faugeras & Toscani, 1986), and others control or restrict camera motion instead (Hartley, 1994; Stevenson & Fleck, 1995; Stein, 1995). We will summarize here the requirements for the first type of procedure (e.g., the one described in (Willson, 1994, Chapter 5)) for acquiring Euclidean geometry without restricting camera motion, and then extend it to include the acquisition of arbitrary ground truth in Section 2.6.

Define and Register Coordinate Frames

The first step is to define and register all coordinate frames. The eventual goal is the acquisition of imagery and co-registered ground truth of well-understood objects. Thus the mapping between images and the world must be modeled to perform the calibration and ground truth registration. This is purely a representational issue, and is used to determine the types (i.e., units) of measurements that will be used in both the calibration and data acquisition steps below.

To relate imagery to the real world, there are four reference frames of interest: 3D object coordinates, 3D world coordinates, 3D viewing coordinates, and 2D screen coordinates. These are illustrated in Figure 2.7 along with a range frame (explained in Section 2.6). The shape of every object in the scene is described in its own coordinate system, each of which is related to a single world coordinate frame. World coordinates are related to viewing coordinates by the external parameters of the camera model (translation and rotation), and viewing coordinates are related to screen coordinates by the remaining internal parameters.

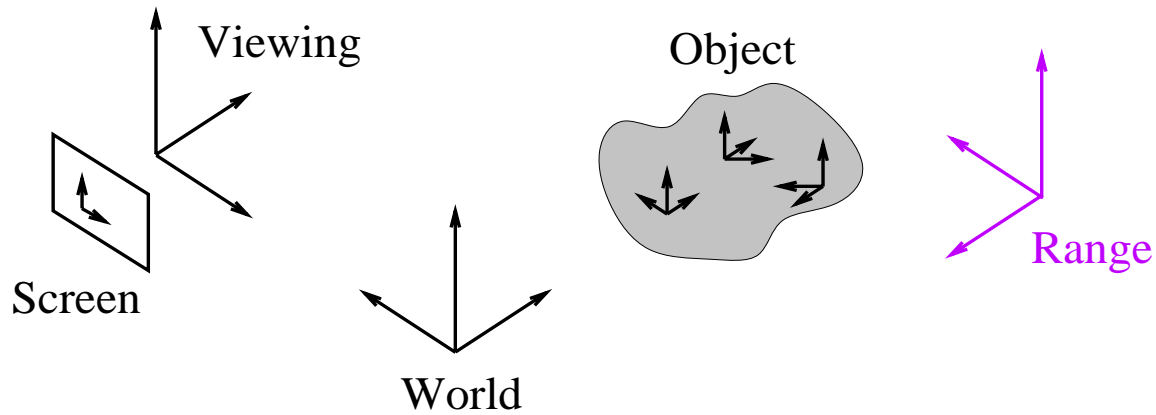


Figure 2.7: Coordinate frames used in dataset acquisition.

In many cases the dataset collection process can be streamlined somewhat by equating the world reference frame with the calibration target's object coordinate frame.

Specification of a coordinate system comprises the directions of the axes, the unit of length along each axis, and the location of the origin. Choosing 2D screen coordinate axes is relatively straightforward: the axes point horizontally and vertically along the CCD array, the units are typically pixels, and the origin is either a corner of the image, the center pixel of the CCD, or is specified by the internal calibration parameters. Specification of the remaining 3D frames will depend upon the application, but the units of the viewing coordinate Z axis will typically be the same as the depth estimates computed by the stereo algorithm, and each set of coordinate axes will usually be orthogonal.

The parameters that relate these coordinate frames are calculated during the camera calibration process.

Establish Calibration Target and Scene Objects

The 3D structure of the calibration target and scene objects must be determined to the best resolution possible, typically on the order of a millimeter or some fraction thereof. This can be accomplished by measuring existing objects or by manufacturing new ones. The calibration target should satisfy the constraints given in Section 2.5.1, but most importantly will ideally sweep out the entire volume of the area to be imaged: not just the volume of the target, but also that of any textured background areas. An easy way to do this is to use a flat or concave calibration target that occupies the entire image, and leave it in the background while imaging the scene objects.

Software that extracts features and builds an internal CAD model of everything in the scene must be developed concurrently with the physical construction of the calibration target and scene objects. These pieces are complementary; e.g., the specifications for the target should be driven by the capabilities of the software. Feature points on the calibration target should be constructed so that their image locations can be accurately measured to subpixel resolution. For example, the targets shown in Section 2.5.1 all use circular feature points for which the centroid can be determined to subpixel precision (Willson, 1994, Appendix D). The software must be able to robustly fit extracted features to the target's known 3D structure, allowing for those distortions that will inevitably occur due to the as-yet unmodeled camera parameters. Although this process is often performed manually (and should always be manually verified), by automating it a substantial bottleneck in dataset acquisition can be avoided.

The costs of establishing a good calibration target and creating recognition software will be easily recovered, in terms of both resources and time, in the form of more accurate measurements and faster dataset acquisition. However, the construction of scenes with complete ground truth is not always possible. For outdoor scenes in particular, the cost of constructing a large enough calibration target and acquiring completely dense ground truth is often prohibitive. The best way to approximate this is to construct models in a laboratory and adjust the lens parameters to simulate outdoor imagery, or compromise the results by using a calibration target smaller than the scene being imaged and acquiring only sparse ground truth such as will be described in Section 2.6.

Acquire Calibration Imagery and Calibrate

An integral part of the calibration process is the acquisition of several images of the calibration target. Before acquiring the calibration imagery, however, the layout must be checked to ensure that the scene images will accurately reflect the desired properties of the scene. The lighting levels must be set, the camera parameters adjusted (e.g., lens aperture, exposure time, focal length), the objects positioned, and some sample imagery taken (both with and without the scene objects). Performing this run-through prior to acquiring the calibration imagery will help ensure that the data collected will measure the properties actually found in the dataset scene imagery.

The calibration data should be acquired under the same conditions as the later scene

imagery. Therefore the camera parameters should be established with both the scene objects and calibration target in mind. For example, when using a target with a white background the aperture must be adjusted so that the calibration imagery is not overexposed.

No matter how many images in the dataset, some calibration target feature points must be visible from all viewpoints. The target should be fixed in one place if possible, to facilitate the inter-image registration that must follow. If the target is moved to accommodate multiple viewpoints, then this transformation must also be known to high precision.

Although the 3D structure of the calibration target and any scene objects are assumed to be known, an independent measurement of the objects (or feature points) can provide a useful sanity check. Tools such as those described in Section 2.6, if available, could be put to good use in such redundant measurements.

Finally, images of the calibration target and the corresponding camera locations (if known) are recorded. These images, the inter-camera transformations, and known structure of the calibration target are all input to the camera calibration procedure.

Camera calibration itself occurs in two steps. First, calibration features are located in the images and mapped to their corresponding 3D world coordinates. This is best accomplished with calibration target-recognition software as described in the previous section. In the second step, these data points (image coordinate and world coordinate vector pairs) are fed into the calibration routine, which uses them to determine the parameters that best fit the camera model.

Once the camera calibration parameters have been computed to a reasonable level of precision, images of the actual scene objects may be taken.

2.5.3 Acquire Dataset Imagery

Finally, the actual stereo datasets can be acquired. The cameras are positioned in the same locations and orientations as were used for the calibration, and with the same settings (e.g., aperture, focal length, and gain). This can be achieved (with difficulty) by building a stable multi-camera head (e.g., as in (Ross, 1993)) or a single camera platform with highly reliable positioning capabilities (e.g., (Willson, 1994, Appendix A)). In addition to imagery, completely dense ground truth will also be acquired in each camera's reference frame.

The best tools for performing the ground truth measurement will vary with the application. One possibility is to use image-based fitting just as was done with the calibration

target. But this technique limits the types of objects that can be imaged to only those with software robust and precise enough to do the fitting. Another is to use measurement tools such as those described in Section 2.6 below to locate key features in 3D world coordinates, and apply the known object shape to fill in the details. In this case the ground truth measurements need only be acquired once since they are expressed in the world coordinate frame; the final result for each camera is computed by simply applying the external parameters for that camera to the world coordinate measurements.

Finally, these completely dense depth maps can be processed as described in Section 2.3.1 to yield occlusion masks for the resulting disparity maps.

Limitations: By controlling the environment, completely dense ground truth may be acquired along with real imagery. However, this approach is very restrictive. The hard work of constructing scenes so that precise ground truth is available over an entire image limits the approach to relatively simple scenes. If the eventual application will include complex scenes as well, imagery of more complex scenes with ground truth must be developed.

2.6 Measured Environment

Adding a range sensor to the laboratory allows images of complex static scenes to be acquired with ground truth. But by relaxing constraints on the objects to be imaged, the acquisition of complete ground truth becomes untenable. Some information is available, e.g., piecewise-planar patches can be measured, but most of the imagery will have unknown ground truth. Even so, this is an important scenario, because it comes the closest to the application of the method outside the laboratory, while still providing some measure of confidence in the results since at least part of the disparity map can be computed precisely. It is generally no longer possible to compute the complete occlusion mask however, because often the missing disparities are exactly those at occlusion boundaries. But at least some of the disparities computed by the stereo method can be verified.

Object shape knowledge and range sensors are used to acquire the ground truth. The locations of feature points in the scene can be measured using pointwise devices such as surveyor's theodolites. Theodolites measure angles rather than distance, but with two theodolites and a simple calibration step, pairs of measured angles can be directly converted to 3D coordinates. By choosing feature points effectively, dense depth maps can be computed

using knowledge of the shape of the objects in the scene. For example, when imaging polyhedral objects, corner point locations can be interpolated to yield dense depth maps over the objects' surfaces.

Why not just use an imaging rangefinder? Rangefinders are indeed useful tools, and could be used to provide some ground truth information, but there are limits (Besl, 1988). Even a perfectly accurate rangefinder would not provide an exact depth map unless it were co-located with the imaging elements. The scientific use of rangefinders is still being studied, and even some popular LIDAR laser rangefinders are subject to outright errors in their measurements, particularly at occlusion boundaries (Hebert & Krotkov, 1992), which makes them unreliable sources of information for evaluating the effect of occlusion on stereo data. Image-based rangefinders that use controlled lighting together with the same CCD array could be very useful in some laboratory situations (see (Tada et al., 1993; Nayar et al., 1995) for two image-based range sensors), but would likely depend on the same camera calibration methods used in the stereo method, resulting in “ground truth” measurements that have some of the same biases as the stereo method.

2.6.1 Camera and Range Sensor Calibration

Data acquisition at this level requires calibration of both the range sensor and the camera. Camera calibration was described in Section 2.5, but the additional requirements and some example datasets are described below.

Calibrate Range Sensor

Ground truth measurements are acquired by introducing a range sensor into the laboratory. Placement and use of the range sensor is complicated by the fact that its use must not interfere with the acquisition of stereo data. Completely dense ground truth measurement will be impossible, for unless the range sensor uses the same CCD elements as the stereo cameras, it cannot be co-located with them, and will therefore be unable to view all the same points.

One choice for range sensor is a pair of surveyor's theodolites. A theodolite is an optical measurement tool consisting of a lens system, stable platform with levels, and instrumentation that measures the angle between an initial direction and that of a point in the world. Angles are converted to distance measurements by combining angles from two theodolites

with the known baseline between the instruments and then triangulating. This system has the desirable properties that it provides range measurements independent of the stereo camera equipment, and makes no restriction on the type of objects that can be measured (except that some features must be visible to each theodolite). The resolution attainable by this system is discussed in Appendix A.

Theodolite calibration occurs in several steps. First the vertical axes are aligned by checking the levels. The horizontal coordinate frame is initialized by aiming each instrument at the other to set a preliminary zero angle, then rotating each by 90 degrees and resetting to zero again. Once both instruments' vertical and horizontal axes have been made parallel, the baseline separating them can be calculated by placing a ruler somewhere in the shared field of view, measuring the angles to two points on the ruler and plugging those angles and the interpoint distance into a known formula.

Determining the best separation for the theodolites is difficult. The best resolution in depth measurements will be obtained with a wide baseline separation, but when the theodolites are too far apart there will be many points near occlusion boundaries that are not visible to both theodolites. This will limit the number of points for which ground truth can be measured, thus reducing the density of the final depth map. Therefore it is important to coordinate the placement of objects, including the camera calibration target, with the positioning of the theodolites.

Define and Register Coordinate Frames

As in Section 2.5.2, the four types of coordinate frames (screen, viewing, world, and object) must be defined and registered. In addition, the coordinate frame of the range sensor must also be determined (see Figure 2.7). This is typically done by locating features on the calibration object using the range sensor, and computing the transformation between the range frame and the calibration object frame.

2.6.2 Acquire Imagery and Range Data Concurrently

Finally the actual datasets can be acquired. Because the recommended range sensor requires a long time to gather measurements, only static scenes may be imaged.



Figure 2.8: Textured cube image and the piecewise-planar patches used by Xiong for error analysis (Xiong, 1995, Figure 3.20). Used with permission.

2.6.3 Example 5: Textured Cube

An example of a dataset from the literature that benefits greatly from even sparse ground truth is the textured cube from (Xiong, 1995). This dataset consists of one stereo image pair, from which an image is reproduced in Figure 2.8. Since each face is known to be planar, only a few points need to be measured to acquire reasonably dense ground truth.

Xiong does in fact use the known planarity to characterize the shape of this object, but instead of taking independent distance measurements he fits planes to the computed disparities on each face. Thus he is able to quantify his algorithm’s ability to recover shape information from stereo disparity, using the residuals from the planar fit as the error measure.

2.6.4 Example 6: Model Train Set

Figure 2.9 illustrates a sample image of a more complex scene with sparse ground truth, taken from the publicly accessible **CIL-0001** dataset.¹ The scene includes many complex surfaces, few of which are absolutely planar, yet some reasonable approximations can be made. For instance, planes can be fit to those roofs where corner locations are available, and to the front faces of the houses and castle towers. The background pixels can also be filled

¹<http://www.cs.cmu.edu/~cil/cil-ster.html>

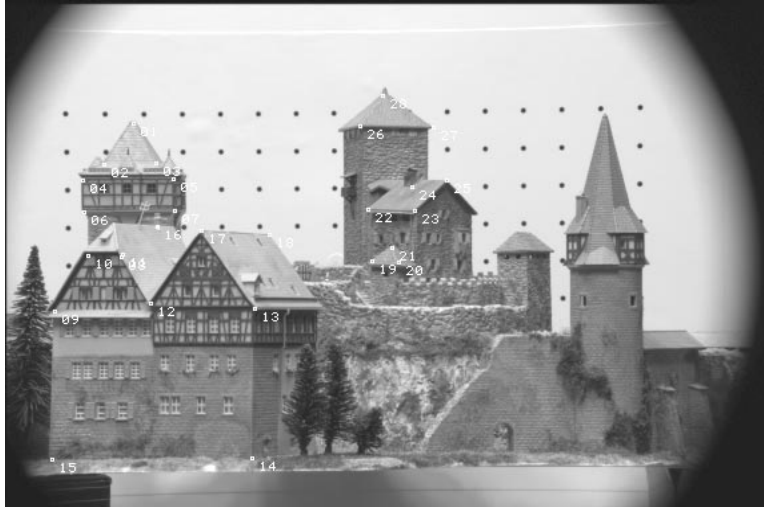


Figure 2.9: Image from CIL-0001 dataset with the locations of ground truth measurements.

in completely, since the shape of the background grid is known from the camera calibration.

The precision of these ground truth measurements is derived in Appendix A.

2.7 Unconstrained Imagery

The final and heretofore most common type of stereo data is that for which no ground truth is made available. The difficulty with such data is that the disparity maps computed by stereo algorithms cannot be assessed metrically, only by human inspection or by ground truth expressed in pixel units.

A common application for this is autonomous navigation. In a typical scenario the acquisition of ground truth is impractical, due to the sheer volume of data being processed. But at least one group has attempted to address this, by simulating road conditions using a static outdoor scene with explicit fiducial marks.²

2.8 Implementation

Our work has benefited greatly from the use of stereo datasets with ground truth. Data from each level of this taxonomy has been used in the research presented throughout this

²The Linköping University Division of Computer Vision provides a small calibrated outdoor dataset at <ftp://isy.liu.se/images/calib.ic/>

thesis. Both synthetic and real imagery were used to develop algorithms, debug their implementations, and characterize their performance. Several examples are given throughout this chapter.

Rayshade, a freeware ray tracing program, provided the foundation for the synthetic datasets used throughout this thesis. Several enhancements were made to adapt this tool to the task of providing stereo datasets with ground truth: automatic depth map extraction, extension of the camera model to include configurable image center and radial lens distortion, addition of back-end tools for depth map manipulation, disparity and occlusion processing, and image format interchange. These extensions and their documentation have been made freely available to the research community.³

Software developed for the Calibrated Imaging Laboratory (CIL) has been successfully used in the collection of stereo datasets with ground truth by several researchers. Camera calibration tools originally developed by Willson (Willson, 1994) have been improved and made more robust, with the result that the time required for dataset acquisition has been reduced from days to minutes. These tools, though somewhat specific to the CIL, are readily available to laboratory members and visitors.

The datasets collected in this laboratory are among the first of their kind to be made available to the general research community: high-quality images with piecewise-dense ground truth.⁰ Publications by other researchers using these datasets have already appeared in peer-reviewed publications, e.g., (Wang & Jepson, 1994).

2.9 Summary

The analysis of stereo vision algorithms can be greatly enhanced through the use of datasets with ground truth. We have outlined a taxonomy of datasets with ground truth that use varying degrees of realism to characterize particular aspects of stereo vision systems, and shown that each component of this taxonomy can be effectively realized with current technology. We proposed that datasets generated in this way be used as the foundation for a suite of statistical analyses to effectively characterize the performance of stereo vision systems.

³See the Computer Vision Source Code web page at <http://www.cs.cmu.edu/~cil/v-source.html>.

Chapter 3

Phase-based Stereo

In theory, a tomographic imaging scanner capable of multiphasic resolution would be able to penetrate this much interference.

— Lt. Commander Data, helping to save the day in
the final Star Trek:TNG episode *All Good Things*

Section 1.4 described several types of stereo vision algorithms. In this chapter we restrict our attention to a particular class of filter-based methods, the so-called *phase-based stereo methods*. We will give a short overview of local spatial frequency, develop our phase-based method, contrasting it with existing ones, and point out issues that must be considered when using any type of phase-based method. The method developed here will be analyzed and extended in later chapters to address issues of foreshortening and ambiguous matches.

3.1 Introduction

There has been a resurgence in interest in filter-based stereo methods. Unlike methods that work by tracking particular features in an image (thus yielding sparse results), filter-based methods allow disparity to be densely computed over entire images. The exponential growth of computer processing power continues to enable new, more powerful and compute-intensive techniques such as multi-camera stereo. Those methods derive their power from the ability to process more data quickly and effectively. Another approach whose appeal is growing is the class of phase-based stereo methods.

In this chapter we will present a phase-based stereo method that requires only a pair of images, but takes advantage of extra processing power to transform its data into representations most appropriate for the task of stereo matching. This presentation of our method takes advantage of the epipolar constraint, and the discussion is therefore restricted to operations on 1D image scanlines. The method can be summarized in one equation, for the disparity at a given pixel:

$$\text{For all } \lambda, \text{ Disparity} = \frac{\Delta\phi(\lambda)}{2\pi} \cdot \lambda_{\text{ref}} \quad (3.1)$$

The parameters of Equation 3.1 will be discussed as follows. The overall framework for the mechanics of the computation is a model of *local spatial frequency* using Gabor filters. *Disparity* will be incorporated as a search parameter in the spatial domain. The *phase* ϕ will be computed by the Gabor filters, and the *phase difference* $\Delta\phi$ will be the principal measured value extracted from the filtered input. The selection of filter wavelengths *for all* λ will be discussed, and the process of wavelength (nee frequency) refinement λ_{ref} will be described.

This gives us the structure for this chapter. We will first review the concept of local spatial frequency, and discuss why Gabor filters seem to be the most appropriate representation for the problem of stereo matching. Next we give some intuition behind the use of the phase of Gabor filter outputs as a measurement of disparity, and the use of instantaneous frequency as a means of refining the value of the signal's wavelength. Having thus outlined the use of a single filter, we discuss the use of multiple filters at different frequencies and the issues that arise in combining their estimates. We then combine all of these into a description of our initial phase-based stereo method, giving several example image pairs and disparity maps. This method will be refined still further in Chapter 5.

3.2 Local Spatial Frequency

There are two fundamental extremes in the study of 2D imagery: the detail-oriented *spatial* view and the wholistic *frequency* view. In the spatial view each individual pixel is paramount; an image is represented by the concatenation of independent pixel values, and each pixel is considered unique and important. In the frequency view only the entire image matters;

although the image is broken down mathematically into several frequency components, information in each component relates to the image as a whole. Each view has its benefits: the spatial view can represent discontinuous textures and local (pixel-level) segmentations directly, while the frequency view is a mathematically elegant representation that enables many useful analyses over large regions. The trouble is, the benefits of one view are difficult to attain in the other. Fortunately, a compromise can be reached, one which preserves the localizability of the spatial approach and the analytical benefits of the frequency approach: this combined approach is called *local spatial frequency*.

In the hope of appealing to those already familiar with Fourier transforms, we will demonstrate this concept using a *spectrogram*, also known as the Short Time Fourier Transform (STFT). If you are already familiar with spectrograms, or do not care to become more familiar, feel free to skip this paragraph. The 2D spectrogram is made up of a sequence of Fourier transforms applied to small, overlapping segments of the signal: each column is the 1D Fourier transform of a window on the signal centered at that same column (i.e., pixel number). The size of this window is fixed (e.g., at 100 pixels in Figure 3.1), but as the window grows adjacent columns will become more and more similar, since most of their windowed input will be the same. It is a highly redundant representation of a signal, an example of which can be found in Figure 3.1. At the limit, if the window size equals the length of the signal, the spectrogram will be only one column wide, and its contents will be the same as the 1D Fourier transform. Now back to the topic at hand.

The general concept of local spatial frequency is illustrated in Figure 3.1. For the moment we consider magnitude only; phase will be introduced later. We begin with a 1D Original Signal, which in this example is a sinusoid with wavelength $\lambda = 20$ pixels (frequency $\omega = 1/20 = 0.05$) embedded within a sinusoid having $\lambda = 7$ pixels ($\omega = 1/7 = 0.14$). This signal is then sampled into 128 pixels from each sinusoid fragment, for a total of 384 data points (we will also call them *image samples*). The frequencies of the original signal can be calculated directly from these samples using the 1D Discrete Fourier Transform (or DFT, see figure), but cannot be localized; we see there are two strong frequencies in the DFT, but we cannot determine where they occur in the original signal. What we *want* to see is a 2D *Ideal* Local Spatial Frequency plot, also shown in the figure. In such a plot, not only would the particular frequencies be known, but also the areas of the original signal in which they occur. In fact such a plot can be approximated, as is demonstrated by the spectrogram in

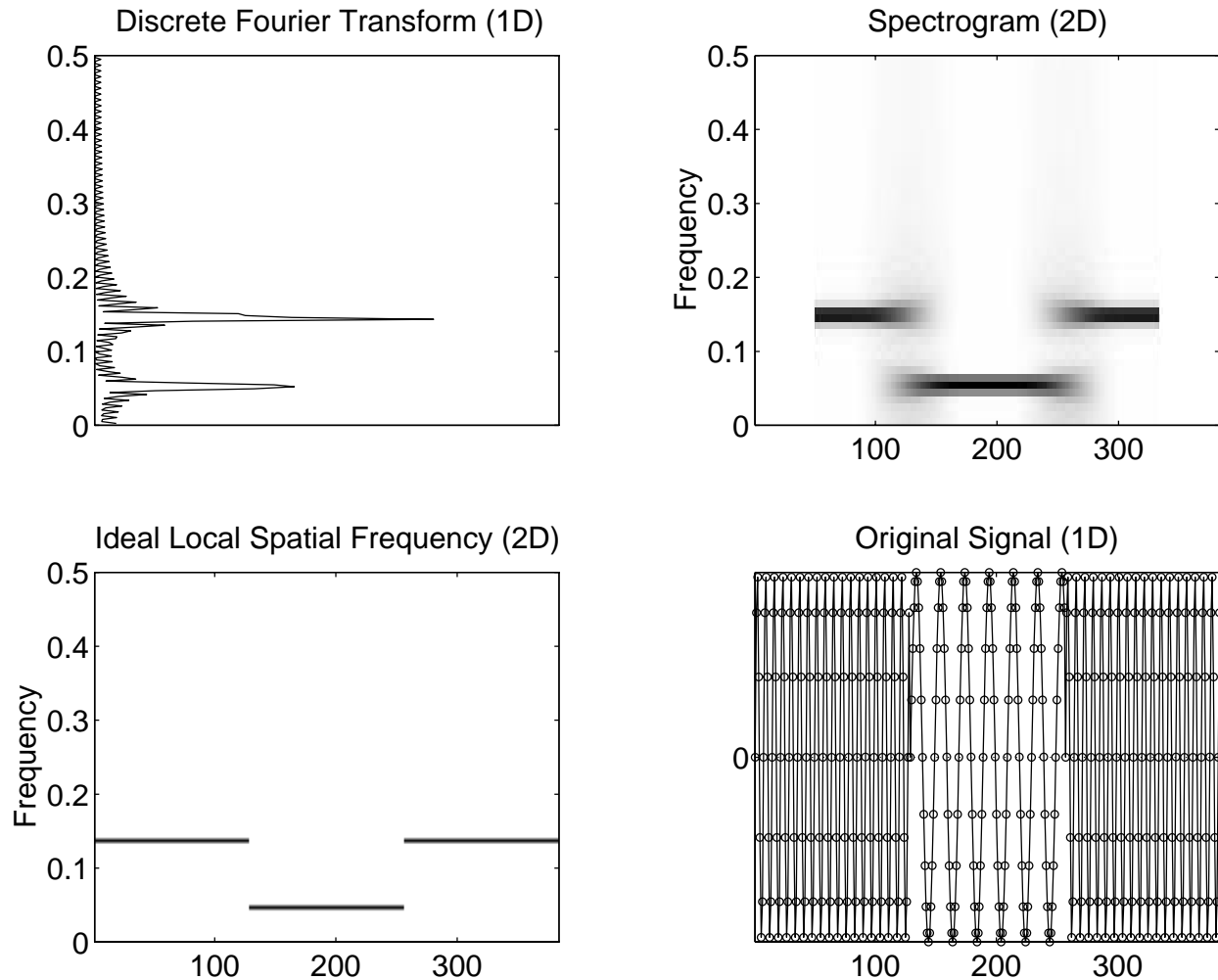


Figure 3.1: Local Spatial Frequency illustration (magnitude only). The Original Signal (lower right) is a low frequency sine wave embedded in a high frequency wave. Its Fourier Transform (upper left) has two peaks, one for each frequency. The Ideal Local Spatial Frequency plot (lower left) associates each sample point with its proper frequency. The Spectrogram (upper right) only approximates this ideal plot, but is computed directly from the sample points *without* explicit knowledge of the original signal’s analytical form. Peaks in the 1D DFT plot correspond to the dark areas in the 2D images.

Figure 3.1, which was computed directly from the image samples *without* prior knowledge of the original continuous signal.

But this representation has some practical and fundamental limitations. On the practical side, the Nyquist interval associated with the Sampling Theorem tells us the analysis window

must be at least twice as wide as the longest wavelength being analyzed; but a wide window will tend to blur the measurements around discontinuities in the original signal (witness the horizontal blurring around the dark lines in the spectrogram of Figure 3.1). More fundamentally, there is a limit to the amount of frequency and localization information that can be extracted. The greater the number of independent frequencies being studied (i.e., the finer the frequency resolution), the coarser the spatial resolution that can be attained, and vice versa. The first mathematical statement of this phenomenon was made by Gabor in (Gabor, 1946), by analogy to the Heisenberg Uncertainty Principle. We can state this “Law of Conservation of Misery” informally as:

$$\Delta x \cdot \Delta \omega \geq \frac{1}{2} \quad (3.2)$$

where Δx and $\Delta \omega$ are the spatial and frequency resolution, respectively. Without going into details (these can be found in (Gabor, 1946, eqs 1.21 – 1.26)), Equation 3.2 says that there is a limit to the amount of local spatial frequency information that can be automatically extracted from sampled data: if you want more frequency information, you have to give up spatial resolution.

3.2.1 Gabor Filters

Having established these fundamental bounds, Gabor went on to describe a family of analysis filters that achieves this optimal resolution. These are now referred to as *Gabor filters*, and they consist of a complex exponential multiplied by a Gaussian window:

$$\text{Gabor}(x, \omega) = u \cdot e^{-i2\pi\omega x} \cdot e^{-0.5\left(\frac{x\omega}{m\sigma_f}\right)^2} \quad \text{for } x \in \left[-\frac{m}{2\omega}, \frac{m}{2\omega}\right] \quad (3.3)$$

$$\text{with Magnitude } \rho(x, \omega) = u \cdot e^{-0.5\left(\frac{x\omega}{m\sigma_f}\right)^2} \quad (3.4)$$

$$\text{and Phase } \phi(x, \omega) = 2\pi\omega x \quad (3.5)$$

where ω is the *tuning frequency* of the filter, $u = \omega/(\sqrt{2\pi}m\sigma_f)$ is a scaling term and m and σ_f are constant tuning parameters; these are discussed below. Gabor filters are closely related to the Fourier transform. In fact, the complex exponential component of the filter is actually *identical* to the kernel of the Fourier transform. Where the Gabor filter differs is in its view of the input data. Whereas the Fourier transform uses the entire image to compute its results, the Gaussian window in the Gabor filter limits attention to a small

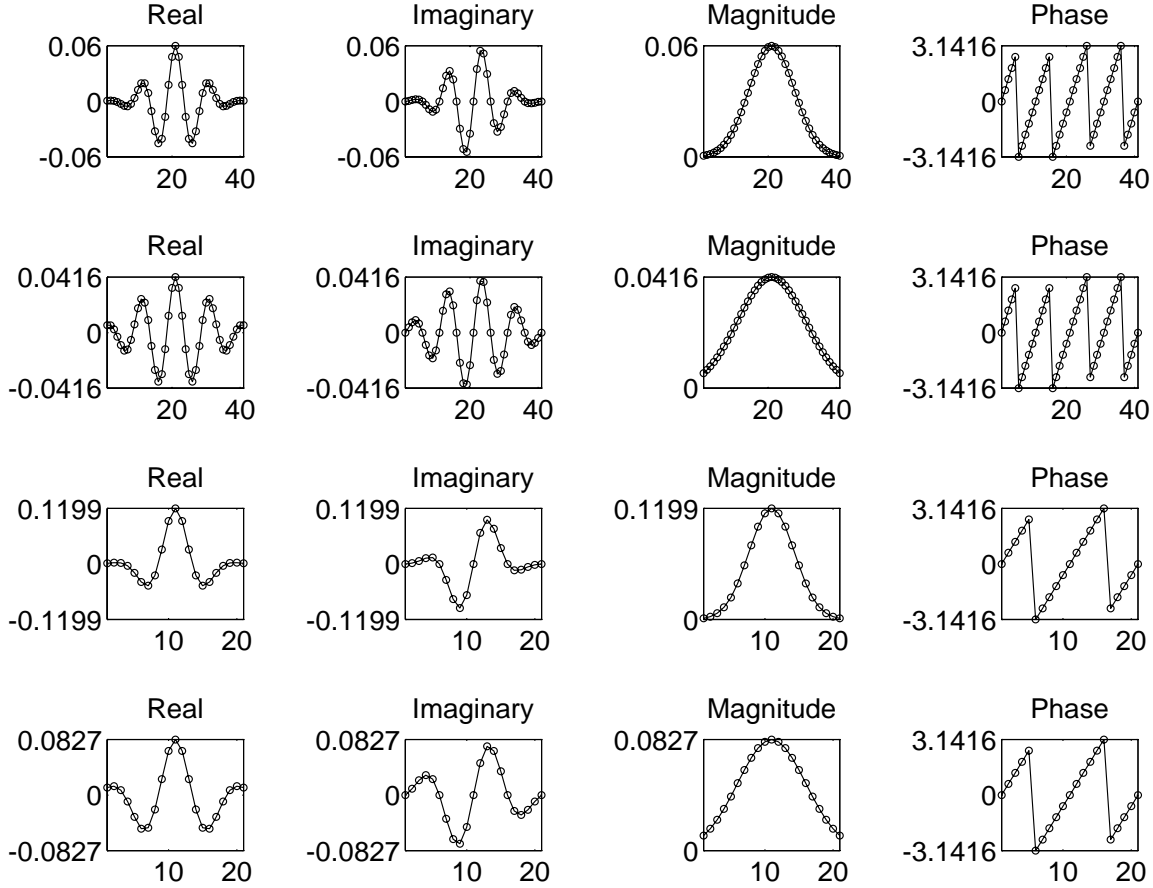


Figure 3.2: Gabor filter examples with $\omega = 1/10$ but different m, σ_f on each row. The horizontal axis indicates the number of samples in the filter.

region in the input image. Also, unlike the Fourier transform which uses a *fixed* window size at each frequency, the Gabor filter uses a window that shrinks and grows as the tuning frequency changes (witness the ω term in the Gaussian component of Equation 3.3). At high frequencies the effective window width will be quite small, thus ensuring that only the nearest pixels will be used to compute filter outputs. At low frequencies the window will be much wider, and will therefore require more data from the original signal.

In addition to the tuning frequency, two additional parameters must be specified to completely define a Gabor filter. The first is the number of wavelengths (m) to include in the window. The Sampling Theorem tells us that for an input signal whose lowest frequency component is ω_L (with corresponding wavelength $\lambda_L = \frac{1}{\omega_L}$), we must include $2\lambda_L$ samples of the signal to reliably extract information about that frequency. More data would be even better, but this number of wavelengths is the minimum. Thus we can choose any value as

long as $m \geq 2$.

The second parameter is how much of the Gaussian to include in the window (σ_f). Parameter σ_f is the fraction of the window size that corresponds to one standard deviation (σ) of the Gaussian, typically $1/6$. Although theoretically the Gaussian has infinite extent, 99.7% of its area is located within three standard deviations from the origin. Thus it is reasonable in practise to limit the extent of a Gaussian window to six standard deviations (three on either side). Since even extreme values of the Gaussian never really equal zero, this causes the tails of the Gaussian to be truncated arbitrarily, but the numerical insignificance of these extreme values makes this a reasonable compromise for digital implementation.

Figure 3.2 illustrates four Gabor filters. Each horizontal row of the figure corresponds to a single filter. Since Gabor filters are complex-valued, two plots are required to completely illustrate their components. The figure presents each Gabor filter in two ways: as a real-imaginary pair, and as a magnitude-phase pair. All of the filters in the figure have constant tuning frequency $\omega = 1/10$, but the extra parameters (m, σ_f) vary from top to bottom: $(4, 1/6)$ on the top, $(4, 1/4)$, $(2, 1/6)$, and $(2, 1/4)$ on the bottom. The plots have been scaled in both directions to make the filters appear as similar as possible, so that the variations in shape and number of samples that result from changes in these parameters can be emphasized.

Parameter Tuning

Tuning of these two additional ‘magic’ parameters is still somewhat of an art. The resolution of the frequency information is greatest when both m and $1/\sigma_f$ are large, allowing lots of data and many wavelengths to be used in the computation. This is fine when the input image consists of entirely one type of signal. Unfortunately, images often contain many objects, with only a small number of pixels representing each one. And in that case we would prefer to use as *little* data as possible (just the relevant pixels) to perform computations; otherwise our results will be corrupted with what is effectively noise from images of the other objects. This adds a new wrinkle to Gabor’s uncertainty principle: not only is the resolution limited, but to accommodate the analysis of multiple objects we should restrict the number of spatial data points even further.

The effects of tweaking these parameters are illustrated in Table 3.1. Ignoring for the moment the details of the stereo algorithm, we can still see that simply changing these filter

$\sigma_f \setminus m$	4	6	8	10	12	14	16
0.2	0.372339	0.204195	0.185642	0.192559	0.187348	0.196636	0.191339
0.25	0.342084	0.192126	0.180964	0.187143	0.166857	0.182314	0.178613
0.3	0.394752	0.232775	0.211052	0.214945	0.166194	0.18913	0.183344

$\sigma_f \setminus m$	18	20	22	24
0.2	0.910356	0.668425	1.18127	1.64487
0.25	0.20824	0.230053	1.10322	1.68861
0.3	0.211346	0.223421	1.09542	1.78614

Table 3.1: Effects of Tuning Parameters on stereo disparity accuracy. Filters with given values for σ_f and m were applied to the scanline stereo pair from Figure 3.3 using the method in Section 3.7. Here the table shows the RMS error between the computed disparity and that known over the portion of the image with visible texture (153 pixels). 401 disparities from 0 to 10 were tested, and the smallest and largest errors are highlighted.

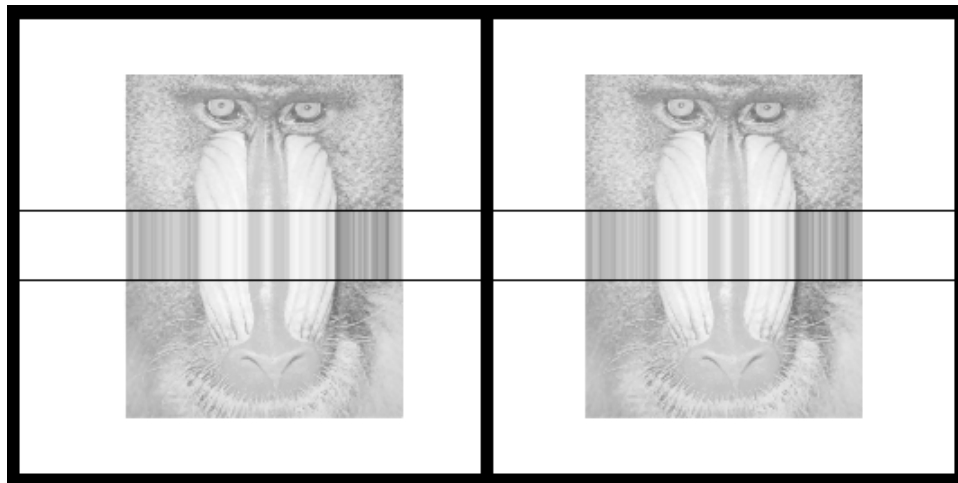


Figure 3.3: Scanline Pair used in Table 3.1, highlighted in the image from which it was extracted.

parameters has a significant effect on the apparent precision of a method using them. In this case, the overall precision varies by an order of magnitude. But predicting such deviations is quite difficult; for example, the results in the table are not strictly unimodal in either dimension, so no general trends can be inferred.

A final note about the magnitude of the Gabor filter envelope as used in this thesis.

Although Gabor filters can be used as wavelets, where shape is preserved across scales, our filters will instead preserve constant area across scales. This is accomplished by using a Gaussian envelope with unit area (scaling with u in Equation 3.3) to normalize the magnitude of the filter outputs across all frequencies. Low frequency filter outputs would otherwise disproportionately dominate the those at high frequencies because the magnitude would grow in proportion to the window size.

3.3 Using Phase as Disparity

The notion that phase corresponds to translation is a familiar one to the signal processing community. The best example of this is the Fourier Shift Theorem, which states that a translation a in space corresponds to a phase shift $2\pi a\omega$ in the frequency domain:

$$\begin{aligned} \text{If } f(x) &\Longleftrightarrow_{\mathcal{F}} F(\omega), \text{ Then} \\ f(x - a) &\Longleftrightarrow_{\mathcal{F}} e^{-j2\pi a\omega} F(\omega) \end{aligned} \quad (3.6)$$

The problem with this theorem is that it is strictly only true for functions with infinite support. Since we want to work with finite images, we have several options: ignore the problem; extend images to infinity using zero-padding, repetition or reflection; or modify the theory to account for the finiteness of function f . Work that extends the theory to handle the finiteness of images (aka the “window effect”) can be found in (Xiong, 1995). Although that work doesn’t completely eliminate the window effect, it reduces it to numerical insignificance through the use of polynomial approximation. But the intuition from this theorem is still good for finite images, as we will see, so our work will call upon the first two options instead.

An illustration of this concept can be found in Section 2.2.1.

3.3.1 Relating Phase to Feature-based Methods

Convolution with a Gabor filter is a procedure similar to finding the local shape of a smoothed version of a signal. More precisely, it is like fitting a sinusoid with varying amplitude and phase to a smoothed input signal. Unlike methods that extract binary features such as edges from their input, by using phase one gets the advantage of a detector of *continuously* varying features. You don’t just get edges, you get “edginess” too. Think of phase as saying what the local signal shape is, and amplitude as the confidence of that interpretation. This view, that

phase values are features to be matched, provides the intuition behind our stereo method which will be presented in Section 3.7.

But you get more than just features from phase measurements. Because the shape of the analysis filter is known (it has a sinusoid as its kernel), you can use the numeric value of the phase measured at one point to predict the phase at another. This is the basis of the so-called *direct phase* methods, which work by assuming that phase will change linearly, so by measuring the local difference you can jump from one pixel directly to its correspondent.

There are problems in interpreting Gabor filter outputs directly, however. The next section discusses the reliability of measurements made using Gabor filters.

3.4 Phase Difference Measurement Characterization

Instead of matching intensities directly, phase methods convolve intensities with continuous filters, then use properties of the filter outputs to determine disparity. This idea was illustrated in Section 2.2.1, and can be derived algebraically using the Fourier Shift Theorem (stated above in Equation 3.6). We will now derive this basic result, already presented in Equation 3.1, but to simplify the derivation we will only consider the effect of translation on a 1D image scanline.

Let image scanline $l(x)$ be given; this will be the left image. Now let an arbitrary disparity d be given; this is a global disparity that will apply to every pixel in the image. The right image is the same as the left, shifted by the amount d :

$$\begin{aligned} l(x) & \text{ is the left image} \\ r(x) = l(x - d) & \text{ is the right image} \end{aligned}$$

We can compute the Fourier transform of each scanline independently:

$$\begin{aligned} l(x) & \Longleftrightarrow_{\mathcal{F}} L(\omega) \\ r(x) & \Longleftrightarrow_{\mathcal{F}} R(\omega) \end{aligned}$$

The Fourier Shift Theorem tells us how $L(\omega)$ and $R(\omega)$ are related:

$$R(\omega) = e^{-j2\pi d\omega} L(\omega) \tag{3.7}$$

and note that Equation 3.7 holds for any value of ω . So all of the complex Fourier coefficients from the right image are exactly the same as those from the left image, except that their

phases are shifted by some amount. Since the magnitude of the coefficients does not change, let us restrict our attention to just the phase components of Equation 3.7. We will use the notation that $\phi_L = \arg[L(\omega)]$.

$$\begin{aligned}\arg[R(\omega)] &= \arg[e^{-j2\pi d\omega} L(\omega)] \\ \phi_R &= -2\pi d\omega + \phi_L \\ d &= \frac{\phi_L - \phi_R}{2\pi\omega}\end{aligned}$$

Solving for the disparity d results in the formula that was presented back in Equation 3.1, the relation that lies at the heart of all phase-difference or direct phase methods. For convenience we restate it here:

$$disparity = \frac{\Delta\phi}{2\pi frequency} = \frac{\Delta\phi}{2\pi} \cdot \lambda_{\text{ref}} \quad (3.8)$$

This equation can be readily understood by applying it to an example of a continuous sine wave. It simply says that disparity (or shift in x needed to align the signals) is equal to the translation of the sine wave (i.e., phase difference $\Delta\phi$, measured in radians) scaled by a refined estimate of its wavelength, or period λ_{ref} . The same intuition works for *windowed* continuous sine waves as well (e.g., Gabor filters which are sine waves multiplied by Gaussians). Problems arise when applying this to discrete signals, however. How is the actual frequency of the sinusoid measured? Phase plays an important role in the disparity computation, but how accurately can it be measured? Also implicit in Equation 3.8 is the assumption that the left and right phases are measured at the same frequency, i.e., on similar sinusoids; is that assumption reasonable?

To illustrate the sensitivity of these components consider what happens when the signal being studied is a simple sinusoid.

3.4.1 Background on Instantaneous Frequency

The concept of *instantaneous frequency* is useful in many signal processing applications, and to phase-based stereo in particular. The term was introduced in (Carson & Fry, 1937) as part of a description of the mathematics behind frequency modulation transmission (i.e., FM radio). In that application, a low frequency signal (the audio signal) is transmitted on a very high frequency carrier wave (at the radio station's tuning frequency) by adjusting the frequency of the carrier wave over a wide range. Those adjustments, the changes in the

carrier frequency, encode the lower frequency signal of interest, and *instantaneous frequency* is the framework that extracts the original signal from those changes.

A more physical way to think of it is like the measurement of velocity on a car's speedometer. Assuming a 1D signal comprises a possibly changing sinusoid at each moment in time (like a car being driven at different speeds), the frequency of that sinusoid (i.e., the car's velocity) at each moment is the *instantaneous frequency* of the signal. There are several ways to compute it, each of which is applied to a small window of the signal surrounding the point of interest.

One bad way to estimate the instantaneous frequency is to convolve the window with a bank of band pass filters tuned to different frequencies, find the filter with maximum magnitude response, and treat the tuning frequency of that filter as an estimate of the signal's instantaneous frequency. The problem is that if only a small number of filters is used it will not be very accurate. Band pass filters respond to any frequency within their pass band, and since a common width is an octave (Fleet & Jepson, 1993) the resulting estimate may be off by a factor of $\sqrt{2}$. The estimate can be improved somewhat by fitting a curve to the responses of adjacent filters and finding the frequency that corresponds to the peak, but the number of filters required for this approach to be effective makes it less attractive than the following option.

If the input signal is known to have come from a single sinusoid, its instantaneous frequency can be computed directly from the *phase derivative* of a *single* band pass filter. The intuition for this notion does not come easily: frequency and phase are often presented as *independent* components of a sinusoid (for example in Section 2.2.1), so how can they be related?

The relationship is illustrated in Figures 3.4 and 3.5, which demonstrate the same analysis on two different sinusoids. Plot a) shows a sinusoid (with fixed frequency in Figure 3.4 and varying frequency in Figure 3.5) and the magnitude envelope of the Gabor filter that is used to analyze it in Plots d) and e). Plot b) shows the phase of the signal, both as absolute phase (zero at the origin and monotonically increasing) and phase modulo 2π , as it will be estimated below. These plots were computed analytically; the phase is known to change linearly over time for fixed frequencies, and it wraps around at 2π exactly where one period of the sinusoid repeats. Plot c) shows the derivative of this phase, i.e., the slope of the plot above. As you can see, the value of this phase derivative is everywhere equal to the frequency

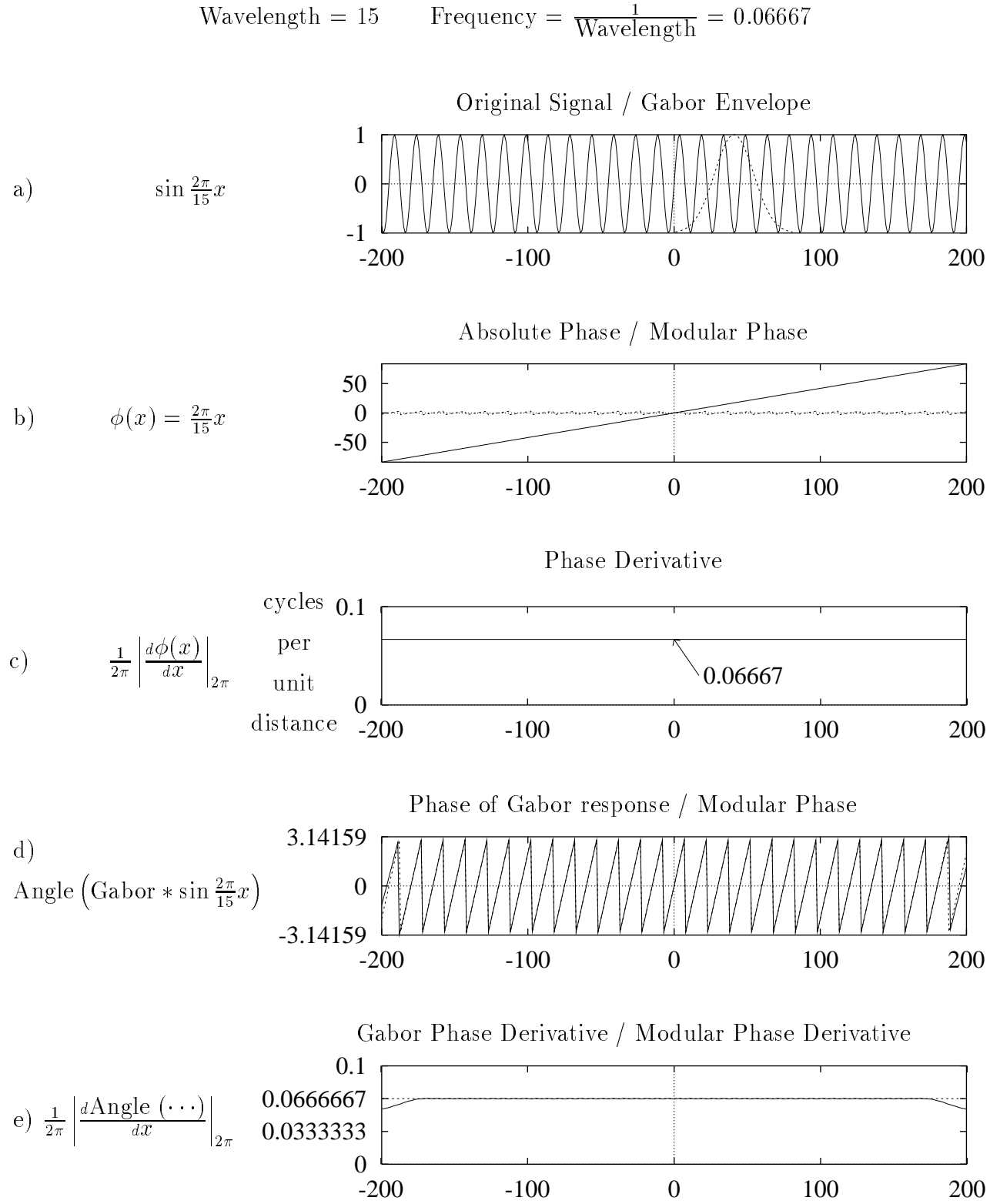


Figure 3.4: Instantaneous Frequency Example: $\omega = 0.06667$. The figure is explained in Section 3.4.1

Wavelength goes from 30 to 15 Frequency = $\frac{1}{\text{Wavelength}}$ = 0.03333 to 0.06667 in linear steps

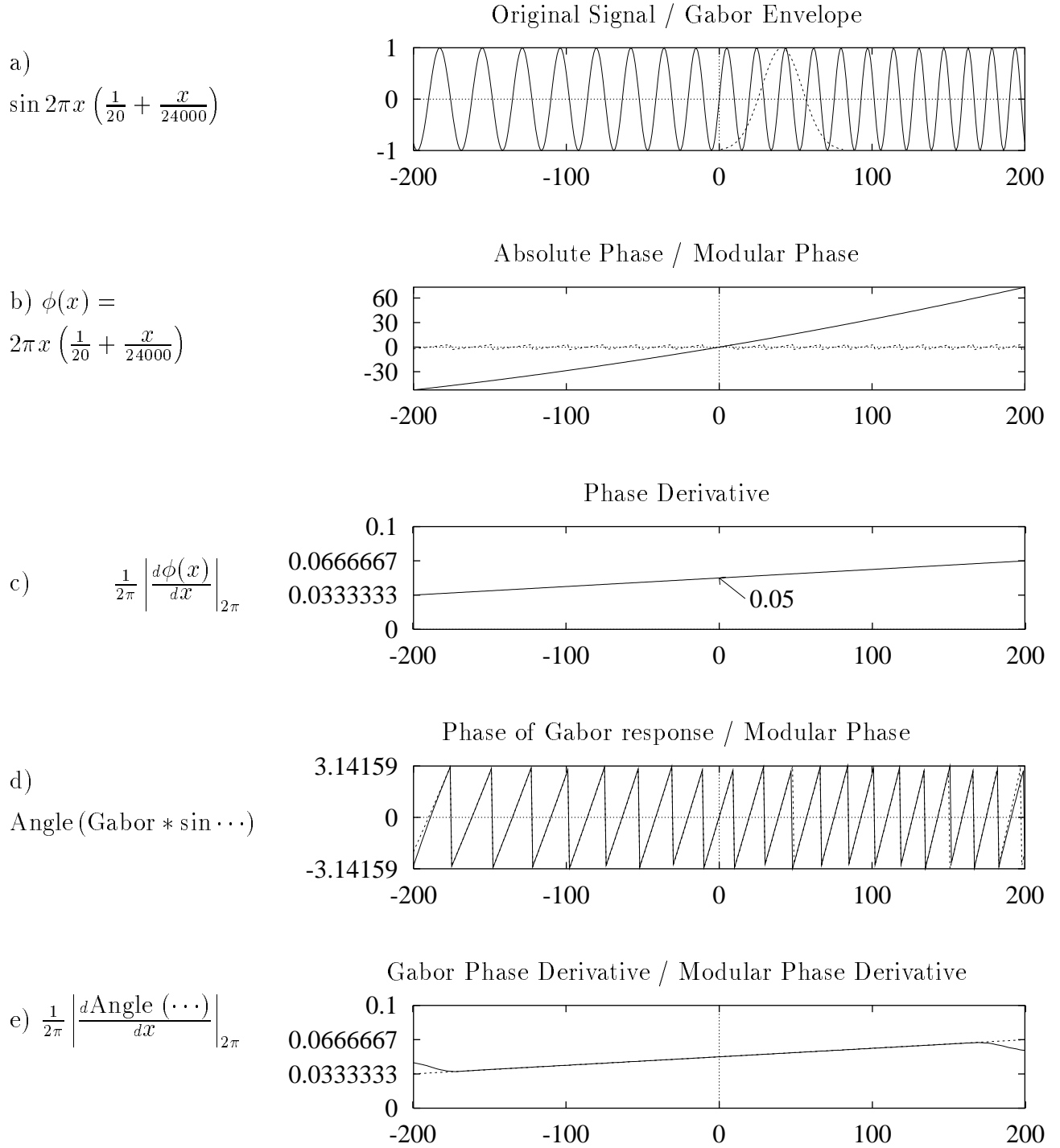


Figure 3.5: Instantaneous Frequency Example: $\omega = 0.03333$ to 0.06667 in linear steps, i.e. a *chirp* signal. The figure is explained in Section 3.4.1

of the original sinusoid. Amazing! To acquire some intuition as to why this must be so, study the figures and compare the original signal's appearance to the X-axis. A sinusoid has a fixed shape, but may be stretched or compressed along the X-axis by changing its frequency. In exactly the same way, the phase of a sinusoid will always change linearly, and its rate of change (i.e., slope) will be related to the original signal's period, or frequency.

The linear change in phase of a sinusoid is not only theoretically interesting, it can also be measured in an arbitrary signal using a band pass filter. The bottom two plots of Figures 3.4 and 3.5 demonstrate this by convolving the original sampled signal with a single Gabor filter whose tuning frequency is $\omega = 1/20$. Even ignoring the details of the computation, one can see in Plot d) that the phase computed by the Gabor filter is approximately equal to the analytically known phase, the dashed line in both Plots b) and d). Finally, the derivative of the Gabor phase is shown in Plot e), again with the analytically known value plotted as a dashed line. Except around the discontinuities at the ends of the signal, the empirically-derived Gabor estimates are extremely close to the analytically known values, even in the chirp signal (Figure 3.5) where the frequency keeps changing. The most interesting property to notice here is that this method of computing instantaneous frequency is reasonably independent of the analysis filter's tuning frequency, since the *same* Gabor filter was used to compute the phase in both figures, and the frequencies in the chirp signal span an entire octave.

As an aside, when the tuning frequency of the filter differs significantly from the frequency of the sinusoid, an *aliased* response may result. In that case the sampled phase derivative may not reflect the instantaneous frequency at all. For a nice explanation and illustration of the effects of aliasing, see (Krumm, 1993).

Now that we have sufficient background in the frequency domain, we must consider the operational issues involved in turning the theory of phase-based stereo matching into a working algorithm. The following sections consider each component of Equation 3.1 in turn, discussing issues of reliability and measurement precision for each parameter in the equation.

3.4.2 Measuring Sinusoid Frequency (λ_{ref})

When filtering a periodic signal there are two common techniques for determining the frequency of the signal: using the frequency of the pass-band filter with highest magnitude response or measuring the phase derivative of the filtered signal.

Using the filter tuning frequency directly is troublesome because any discrete filter will have a blurred response — it combines the responses to all of its passed frequencies. Some systems use the approximation that the filter’s peak tuning frequency is a good enough guess (Sanger, 1988; Weng, 1993), but this assumption may reduce the precision of the results. One “advantage” to using this is that its variance is a fixed quantity; it does not depend on the actual signal content.

The phase derivative (aka *instantaneous frequency*) was shown in Section 3.4.1 to be a good predictor for the example signal. While it provides a more accurate measure of a sinusoid’s frequency in theory, in practise its accuracy depends on the amplitude of the input signal. This has led several researchers to develop constraints that filter out unreasonable phase-derived frequencies (Fleet et al., 1991; Xiong, 1995). The constraint developed by Fleet *et al* filters out those responses whose phase derivatives predict frequencies that lie outside the range of the filter. While this constraint provides a useful “sanity check” on the validity of the computed frequency, it cannot be used to completely eliminate the effects of aliasing or compensate for perspective foreshortening. Still, it will be useful in improving the accuracy of our results, so we discuss it further in the next section.

3.4.3 Measuring Phase ($\Delta\phi$)

The output of a Gabor filter is a complex number. We are primarily concerned with measuring the phase component of this number, but how accurate is a given phase value? Assuming the standard model of broadband white noise (Horn, 1986, p. 129), each complex coefficient in the frequency domain will have a complex error added to it (we represent the error as a 2D vector $\vec{\epsilon}$). $\vec{\epsilon}$ may differ at each frequency, but will always be a perturbation with magnitude of at most ϵ and a random phase orientation. Figure 3.6 gives us some intuition. The vectors represent the Gabor coefficients which are complex numbers; the length of a vector is its magnitude, and the angle between the vector and the horizontal axis is its phase. Adding an error vector $\vec{\epsilon}$ with length ϵ has little effect on the phase of the longer vector, but results in a completely random phase value for the smaller vector. The precision of a particular vector’s phase can be determined exactly by comparing its length to that of the error vector $\vec{\epsilon}$.

Given a vector \vec{v} and an error vector $\vec{\epsilon}$, we can compute the potential error in its phase $\angle\vec{v}$ exactly using Equation 3.9 below, derived as follows. Place the tail of \vec{v} on the origin

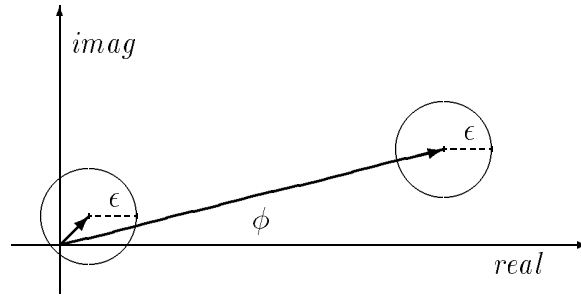


Figure 3.6: The effect of measurement error vector $\vec{\epsilon}$ on phase angle ϕ .

and rotate \vec{v} so that it points to the right on the horizontal axis (i.e., with phase angle zero). Draw a circle of radius ϵ centered at its head. To find the maximum phase error, draw the line tangent to the error circle that also intersects the origin (actually there are two symmetric such lines; draw the one on top). The slope m of that tangent line gives us the maximum error angle: $m = \tan \text{error}$. To find the slope, plug the equation for the line ($y = mx$) into the equation for the circle ($(x - |\vec{v}|)^2 + y^2 = |\vec{\epsilon}|^2$) and solve the resulting quadratic equation for x . At the point of intersection this equation will have two real and equal roots, so we set the discriminant equal to zero and solve for m :

$$\tan \text{error} = m = \sqrt{\frac{1}{\left(\frac{|\vec{v}|}{|\vec{\epsilon}|}\right)^2 - 1}} \quad (3.9)$$

Equation 3.9 tells us that the maximum phase error depends solely on the ratio between the lengths of vectors \vec{v} and $\vec{\epsilon}$. Sample phase angle errors are plotted as a function of this ratio in Figure 3.7.

There are many factors that contribute to the error vector $\vec{\epsilon}$ and therefore reduce the precision of the measured phase value: the precision used in the floating point arithmetic, the blurring of the filter response due to the presence of information at nearby frequencies, the discretization of the filters themselves (especially at high frequencies), and aliasing due to image textures with wavelength less than two pixels to name a few.

A technique for filtering out unreliable phase measurements can be found in (Fleet & Jepson, 1993). They identify two sources of unreliability, and provide a constraint that identifies those regions of space-frequency where measurements should not be trusted. Using the phase and magnitude notation for a Gabor filter with tuning wavelength λ , their constraint

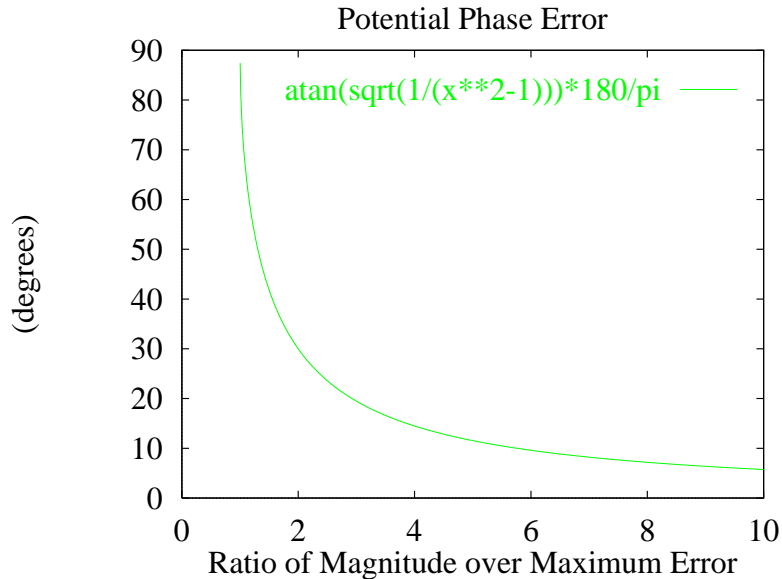


Figure 3.7: Maximum phase angle error as a function of the length ratio.

for finding reasonable measurements can be expressed as:

$$\left| \frac{\delta\phi(x, \lambda)}{\delta x} - \frac{2\pi}{\lambda} + \frac{1}{\rho(x, \lambda)} \frac{\delta\rho(x, \lambda)}{\delta x} \right| < Threshold \quad (3.10)$$

This formula for the constraint was derived from the following two observations:

- The phase derivative (aka instantaneous frequency) of a Gabor filter should be roughly equal to its tuning frequency.
- The amplitude derivative should be small.

These observations (and the constraint in Equation 3.10) are helpful in eliminating unreliable measurements, but they are somewhat simplistic. For example, Chapter 5 will demonstrate that not only *can* the instantaneous frequency differ from the tuning frequency, it *will* differ when the image pair exhibits perspective foreshortening. Also, Xiong has refined this constraint to an arbitrarily higher degree using Moment Filters (Xiong, 1995). But since we will adapt this constraint to handle foreshortening, and since (like Fleet *et al*) we also use Gabor filters, Equation 3.10 will be useful for us as well.

A Peak-finding Heuristic

Another approach at eliminating unreliable estimates that was found to be useful on some images involved a more heuristic process applied to the output of several filters (i.e., one

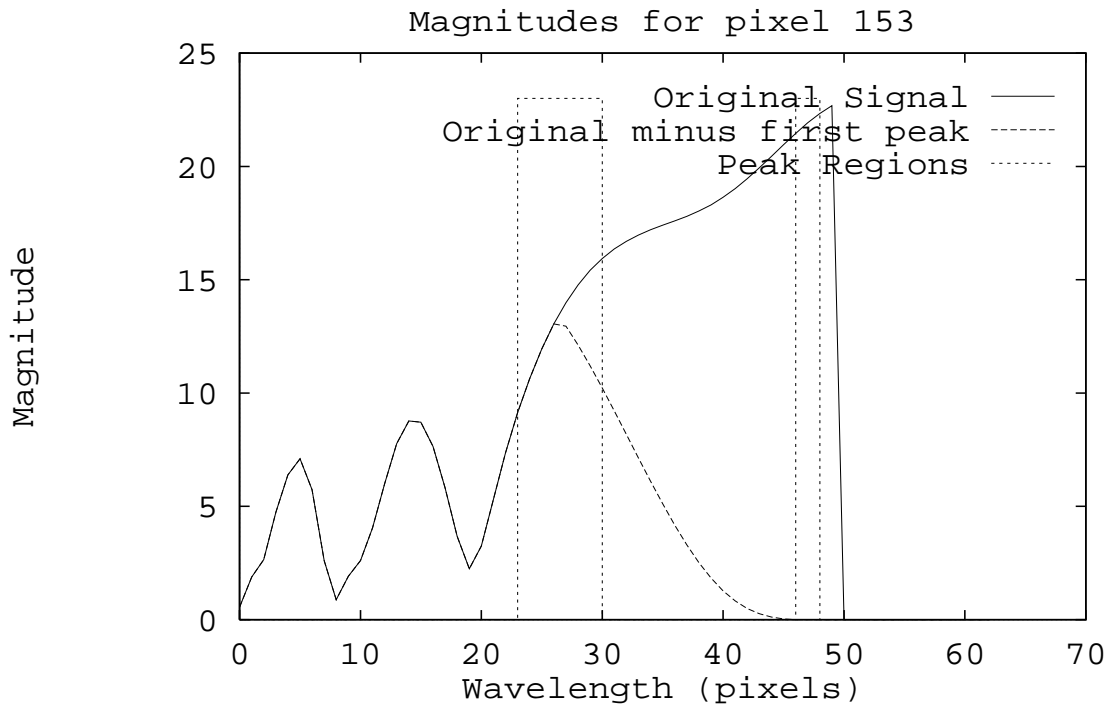


Figure 3.8: Sample Magnitude profile of several Gabor filters, illustrating the peak finding heuristic.

column of the *scalogram*, discussed below). Inspired by the intuition that the most reliable results could be found by concentrating the search on those filters whose magnitude was highest, it isolates those filters whose magnitudes exhibit a peak (or spike) in the overall frequency profile. This kind of scheme works best when the input consists of a sum of band-limited signals, and makes an effort to deal with interference between nearby filters.

The heuristic peak finding scheme works as follows. We attempt to find those regions of the magnitude profile that exhibit explicit peaks by iteratively fitting and subtracting Gaussian curves to the peak regions of the current profile (illustrated in Figure 3.8). The extent of the peak region is determined by a simple heuristic: start with a maximum value, move outward until the second derivative is no longer negative or the end of the function is reached, then back up one or two pixels. The indices of the current peak region are stored at each step, a Gaussian curve is fitted to that region and subtracted from the function, and the process repeats, terminating when the current (subtracted) function's maximum is less than some fraction of the original maximum.

3.4.4 Comparing Frequencies (λ)

In order for phase difference to yield a precise disparity estimate, it must reflect measurements taken at the same frequency. But what if the instantaneous frequencies measured by the same filter on a pair of corresponding image regions differ, as can occur with perspective foreshortening? Others have smoothed over the effects by using the filter tuning frequency, or the average of the two instantaneous frequencies (Fleet & Jepson, 1993), but these are very coarse approximations not based on physical reality. One contribution of this thesis is the development of a theory for finding the proper frequencies based on the physical geometry of the scene, presented in Chapter 5. However, there are difficulties in combining estimates from different filters even when the frequency contents of the image pair are directly comparable, e.g., when there is no foreshortening.

3.5 Choosing Frequencies (*for all* λ)

So far we have discussed primarily the properties of a single filter, applied to signals with a single frequency component. However, unless the input is known to be of a certain form (e.g., an image consisting only of sine waves with restricted frequencies), any useful analysis will require many filters. A nice overview of local spatial frequency models that use multiple frequencies can be found in (Rioul & Vetterli, 1991), with a longer enumeration of specific models in (Hlawatsch & Boudreaux-Bartels, 1992).

3.5.1 The Image Scalogram

For most of this thesis the *image scalogram*'s filter set will be used. The image scalogram is the result of applying a collection of Gabor filters to a 1D image scanline, and therefore consists of a 2D matrix with complex-valued elements; an example is presented in Figure 3.9. Its seemingly unusual triangular shape is due to the interaction between the varying filter window size and the border of the scanline. The high frequency filter outputs found at the top of the scalogram are computed using only a few pixels, so their values can be computed over most of the width of the image. At lower filter frequencies, the number of points that must be sampled increases (as required by the Sampling Theorem), reducing the number of outputs that can be computed from the fixed number of scanline pixel values. Finally, the

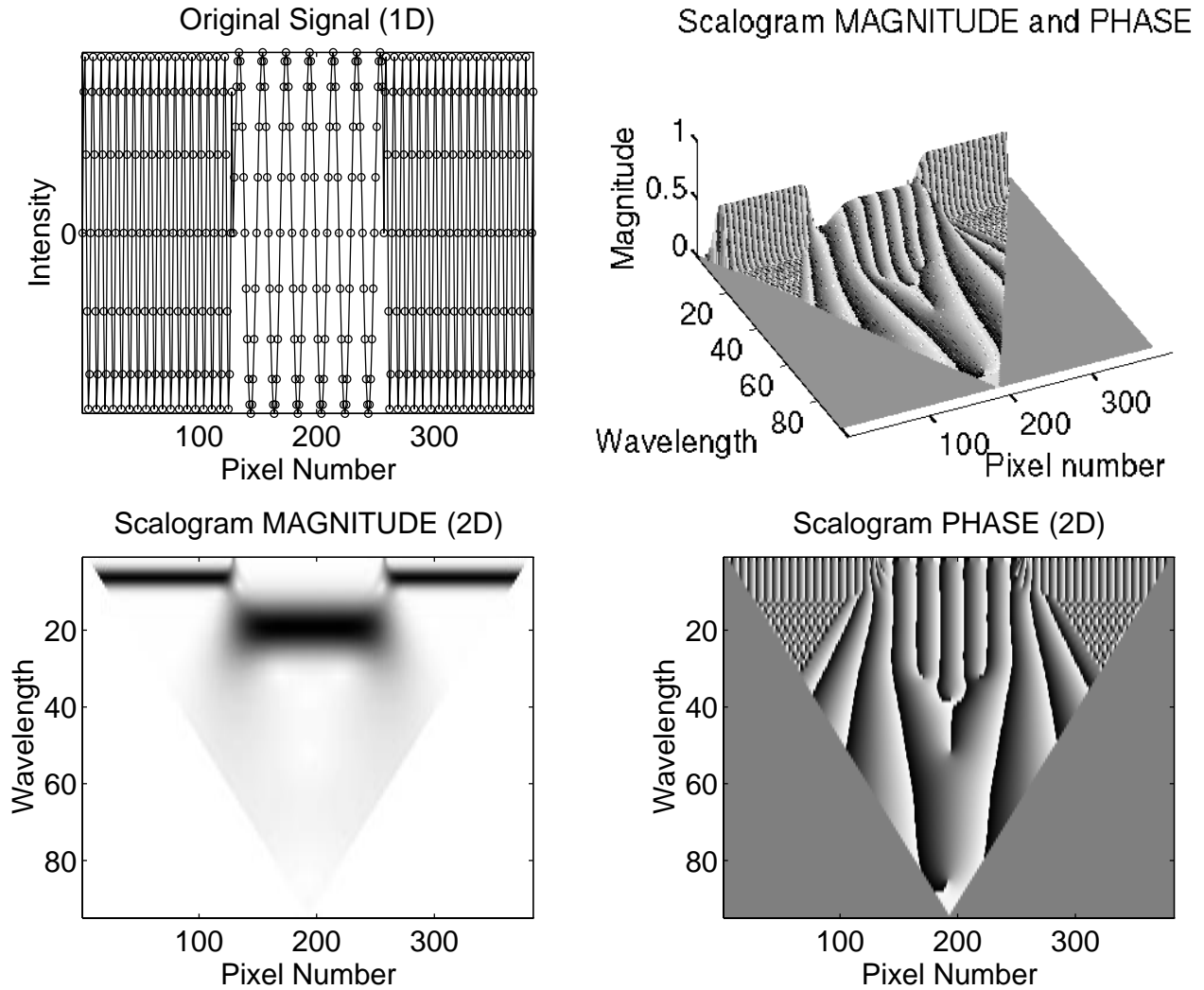


Figure 3.9: Double sine wave signal and associated scalogram (both magnitude and phase).

filter with lowest frequency (located at the bottom of the scalogram) can only be computed at a single point, since it requires the entire scanline as input. Although techniques for filling out the “missing” areas exist, they require either knowledge of the signal beyond the bounds of the input image, or heuristic assumptions about the pattern at the image border. We have chosen to avoid the use of heuristics by considering only those filter outputs for which all the data is available; hence the triangular shape.

The sampling along the (vertical) frequency axis is one of the principal differences between the image scalogram and other local spatial frequency representations. The sampling used in the spectrogram (e.g., Figure 3.1) is the same as that computed by the FFT: a linear sampling of n frequencies from $\omega = 2/n$ to $1/2$ (where n is the width of the sampling window). In

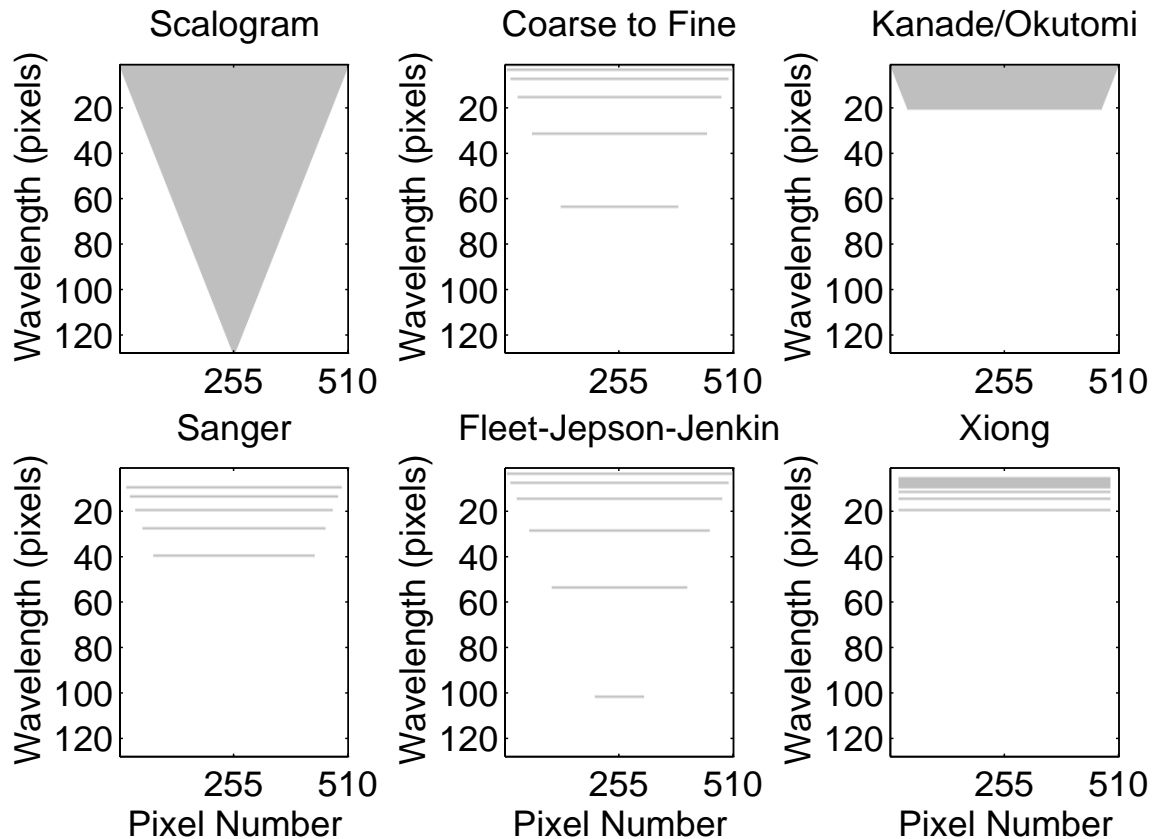


Figure 3.10: Frequency Sampling in the Scalogram and other methods. Grey areas show the central tuning frequencies that are used by each method. Filter widths vary, but typically cover most of the regions between tuning frequencies (see Figure 3.13 for an example).

contrast, the extreme values of the image *scalogram*'s frequency axis are the same, but the sampling is linear in *wavelength*, which is the reciprocal of frequency. Why do we prefer the latter? In part because it was found in practise that weighting the contributions of the lower frequencies more strongly (e.g., by using more low frequency filters) often improved the accuracy of the phase-based stereo method. This dense sampling will also be helpful in compensating for perspective foreshortening effects, as will be explained in Chapter 5.

The scalogram also provides a convenient framework for discussion of other samplings reported in the literature. Because its frequency sampling is based on window width, it is easy to see how the selection of filters in this domain compares with spatial domain approaches like coarse-to-fine (aka image pyramid), and variable window size. Since many of those approaches use far fewer filters than the scalogram, their effectiveness at capturing the

essence of the useful frequency information from the image can be seen by considering how completely they sample the structures evident in the dense scalogram. Some comparisons are illustrated in Figure 3.10 by highlighting the regions of the scalogram that correspond to the filter sets of those other methods.

Some problems with the often-used coarse-to-fine stereo methods become apparent in the context of the image scalogram. A common approach in many spatial domain stereo methods (Matthies, 1992; Nishihara, 1984), coarse-to-fine has also been used in phase-based stereo methods (Weng, 1993; Fleet et al., 1991). In these approaches, information from the lowest frequency (lowest horizontal stripe in each part of Figure 3.10) is used to provide initial coarse disparity estimates, which will be refined by later processing at subsequent higher frequencies. The difficulty is that there are likely to be regions of the scanline over which the filter outputs will be unreliable,¹ and therefore the initial disparity estimates will be unreliable as well. Even when the higher frequencies theoretically provide enough information to compensate for this problem, the final disparities computed by these methods can be highly inaccurate in such regions. This subject will be addressed in detail in Chapter 4, but it points out the need to consider *all* the filter outputs concurrently, rather than in an arbitrary fixed order.

3.6 Combining Filter Estimates

Section 2.2.1 showed how the phase of a sinusoid, and by extension the measurement of phase using a Gabor filter, could be used to compute stereo disparity. In this section we show how estimates from several filters may be combined to yield more accurate results. This technique is mainly useful in conjunction with methods that compare phases between images using the same-numbered pixel (the so-called *direct phase* method), such as those in (Sanger, 1988; Fleet et al., 1991; Weng, 1993), but it will also be useful in our search-based method.

¹The prevalence of unreliable phase measurements is discussed at length in (Fleet & Jepson, 1993). That it is a common occurrence can be inferred from the number of “measles” (white spots) in the scalograms shown throughout this thesis; they represent regions whose Gabor filter outputs exhibit low magnitude, and are found in most scalograms of real imagery.

3.6.1 The Phase Wraparound Problem

Direct phase methods make explicit use of phase difference as disparity, so all such methods have to contend with the problem of phase wraparound. The difficulty lies in the fact that phase can only be measured modulo 2π . Therefore, unless special steps are taken, a given filter will only be able to estimate disparities less than the wavelength specified by its tuning frequency. Mathematically, the problem is that the measured phase difference $\Delta\phi$ might not equal the ideal phase difference $\Delta\phi_{ideal}$:

$$\Delta\phi_{ideal} = 2\pi \frac{Disparity}{\lambda} \quad (3.11)$$

$$\text{But } \Delta\phi = \left| \Delta\phi_{ideal} \right|_{2\pi} \quad (3.12)$$

$$\text{So } \Delta\phi_{ideal} = \Delta\phi + k 2\pi \quad (3.13)$$

$$\text{where } k = \left\lfloor \frac{\Delta\phi_{ideal}}{2\pi} \right\rfloor = \left\lfloor \frac{Disparity}{\lambda} \right\rfloor$$

So $\Delta\phi$ is only part of the answer: the additional factor k must be known if disparity is to be computed exactly from a single filter. Unfortunately, there is no way to measure k without knowing the ideal disparity. While there are some techniques that recover these k values from a set of filter outputs (the so-called *phase unwrapping* techniques of (Tribolet, 1977; Ghiglia & Romero, 1994)), they require accurate phase measurements (and therefore high magnitude response) from *all* filters, a condition rarely met in real images. Because recovery of these k values is so difficult, many methods either assume $k = 0$ or arrange the processing so that k is assumed to be known.

How likely is it that the wraparound issue will arise? Figure 3.11 shows several plots of ideal phase difference as a function of disparity. As the disparity increases, so does the number of times the phase value is expected to wrap around. Methods that rely on many high frequency filters are especially likely to suffer wraparound effects as the disparity grows larger.

This wraparound problem weighs heavily on Sanger's method, which is actually limited to disparities less than the *smallest* wavelength of his filter set. The restriction comes from his technique for combining estimates: averaging the disparity predicted at each filter, assuming $k = 0$ for all phase differences. His method is therefore unable to take advantage of the fine detail present in the highest frequencies of the input images, since the wavelength of

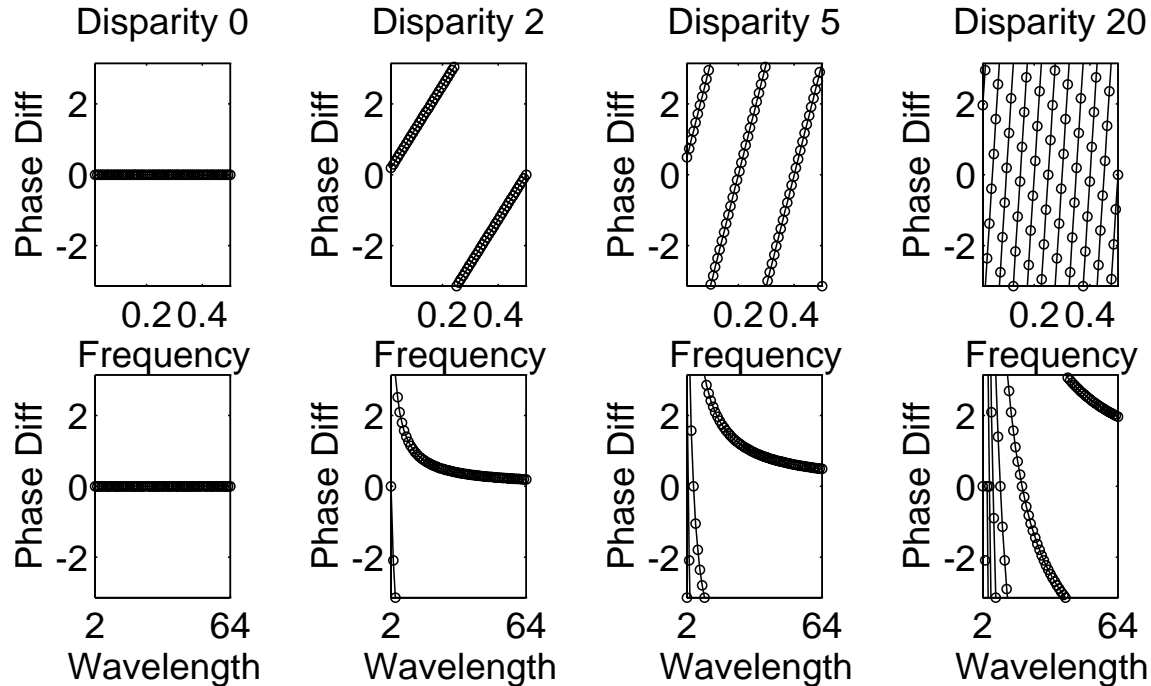


Figure 3.11: Ideal phase difference as a function of Disparity, illustrating the phase wraparound problem. Note that coarser scales are on the *left* in the Frequency plots, but on the *right* in the Wavelength plots. Each graph has 63 sample points spaced linearly along the X axis. Note that unlike the 2D scalogram plots, here the wavelength axis is *horizontal*.

the highest frequency filter must meet or exceed the maximum expected disparity. You can see how his filter set differs from other methods (by lacking the highest frequencies) in Figure 3.10.

Other direct phase methods (Fleet et al., 1991; Weng, 1993) address the phase wraparound problem using a coarse-to-fine approach. Low frequency filters provide initial coarse disparity estimates, which are used only to fix the value of k for the next-higher frequency. This results in two advantages over Sanger's method: disparities as large as the *longest* wavelength in the filter set can be computed, and information from much higher frequencies can also be incorporated. Unfortunately, these approaches make the unreasonable assumption that all filters will have accurate measurements at every pixel. Since there are actually many regions in the scalogram where phase measurements are unreliable (Fleet & Jepson, 1993), the value of k used in the filter with next-higher frequency could be completely wrong, resulting in wildly inaccurate disparity outliers. Weng attempts to address this problem by linearly interpolating phase in the unreliable areas, but his method has still been found to be the least

precise of several phase-based methods (Jenkin & Jepson, 1994). How then can we combine estimates from many filters without knowing in advance which of them will have reliable results?

3.6.2 Our Filter Combination Method

The usual approach in direct phase methods is to compute disparity directly from the phase difference measurements. But we can turn the computation around too, and address the complementary problem: *given* a particular disparity, what are the expected phase differences from all filters? The solution to this problem is straightforward; simply apply Equations 3.11 and 3.12 to the known disparity and filter wavelengths to find the expected or “ideal” phase differences. Figure 3.11 shows the ideal values associated with several disparities. The stereo task then becomes one of finding that disparity whose ideal phase differences best match what was measured. Why bother with this restructuring? By stating the problem in this way, we eliminate the need to compute the appropriate value for k at each filter.

The problem now is how to define the “best” match between measured phase differences and the ideal phase differences associated with a particular disparity. We need an evaluation function that computes a quantitative agreement between these two sets of phase differences. It should combine estimates from any number of filters with arbitrary tuning frequencies, should deemphasize the unreliable measurements, and must account for the potential wraparound effects of comparing phases (phase differences). Empirical tests of various filter combinations (sums of unnormalized v. normalized magnitude), weighting functions (magnitude, squared magnitude, square root of magnitude), and phase comparisons (raw difference, modular difference, with and without absolute value) were performed on several image pairs, and this formula:

$$EvalDirect (disp, c) = \frac{1}{|W|} \sum_{\lambda \in W} \rho(c, \lambda) AbsDiffMod \left(\Delta\phi_{ideal}(disp, \lambda), \Delta\phi(c, \lambda) \right) \quad (3.14)$$

Where $disp$ is the candidate disparity,

c is the current pixel number,

ρ is the magnitude normalized by window size,

$$\Delta\phi(c, \lambda) = |\phi_L(c, \lambda) - \phi_R(c, \lambda)|_{2\pi}, \quad (3.15)$$

and $AbsDiffMod$ is the smallest difference between phases, i.e.

$$\text{AbsDiffMod}(a, b) = \min_{k \in \{-1, 0, 1\}} \{||a|_{2\pi} - |b|_{2\pi} + k2\pi|\} \quad (3.16)$$

was found to provide the most accurate results using the direct phase approach. As mentioned previously, the filters used are typically those of the image scalogram, i.e., linear samples of wavelength from 2 pixels to the one quarter the width of the complete image. Smaller values of Equation 3.14 indicate better matches with the predicted phase differences, so the disparity that yields the minimum value is the best estimate. As an aside, one might be able to replace the $\text{AbsDiffMod}(a, b)$ function with $\cos(a - b)$, then use the *maximum* value to indicate the best match. The benefit would be in using a simpler function, but the effect of its new shape (a curve that is smoothly varying rather than being sharply peaked at zero) on the precision of the method would need to be studied.

Figure 3.12 illustrates how each Gabor filter and each candidate disparity contributes to the final error estimate. Each column in the figure is labeled with a disparity, and represents the goodness-of-fit of all of the measured phase differences to the ideal phase differences predicted by that disparity. Each column is summed over all appropriate wavelengths (witness the \sum in Equation 3.14), and forms the overall estimate for the likelihood that this disparity is correct for this pixel.²

This evaluation function satisfies the requirements given above. It does not depend on a particular set of tuning frequencies; any set of wavelengths W can be used, and they need not be related according to a fixed hierarchical scheme (although nothing prohibits such a sampling). For example, one might use the scalogram's linear wavelength sampling $W = \{2 \dots N/m\}$ pixels, where N is the width of the scanline and m is the Gabor filter tuning parameter for number of wavelengths per window. Or one could use just the sparse sampling of the coarse to fine methods; Figure 3.13 shows the frequency response for such a set of hierarchical Gabor filters). Whatever the sampling, unreliable measurements are deemphasized in Equation 3.14 by eliminating the worst of them from W and by weighting most strongly those outputs from filters with high magnitude (Section 3.4.3 explained the relationship between magnitude and reliability). The wraparound problem is addressed by finding the smallest absolute modular difference between $\Delta\phi$ and $\Delta\phi_{ideal}$.

This method has several advantages. It allows regions with unreliable phase measurements to be ignored safely, without compromising multi-scale processing. Since the mea-

²Figure 3.12 was actually generated using the *error* function from Table 3.2 instead of Equation 3.14, but the summation method is the same for both error functions.

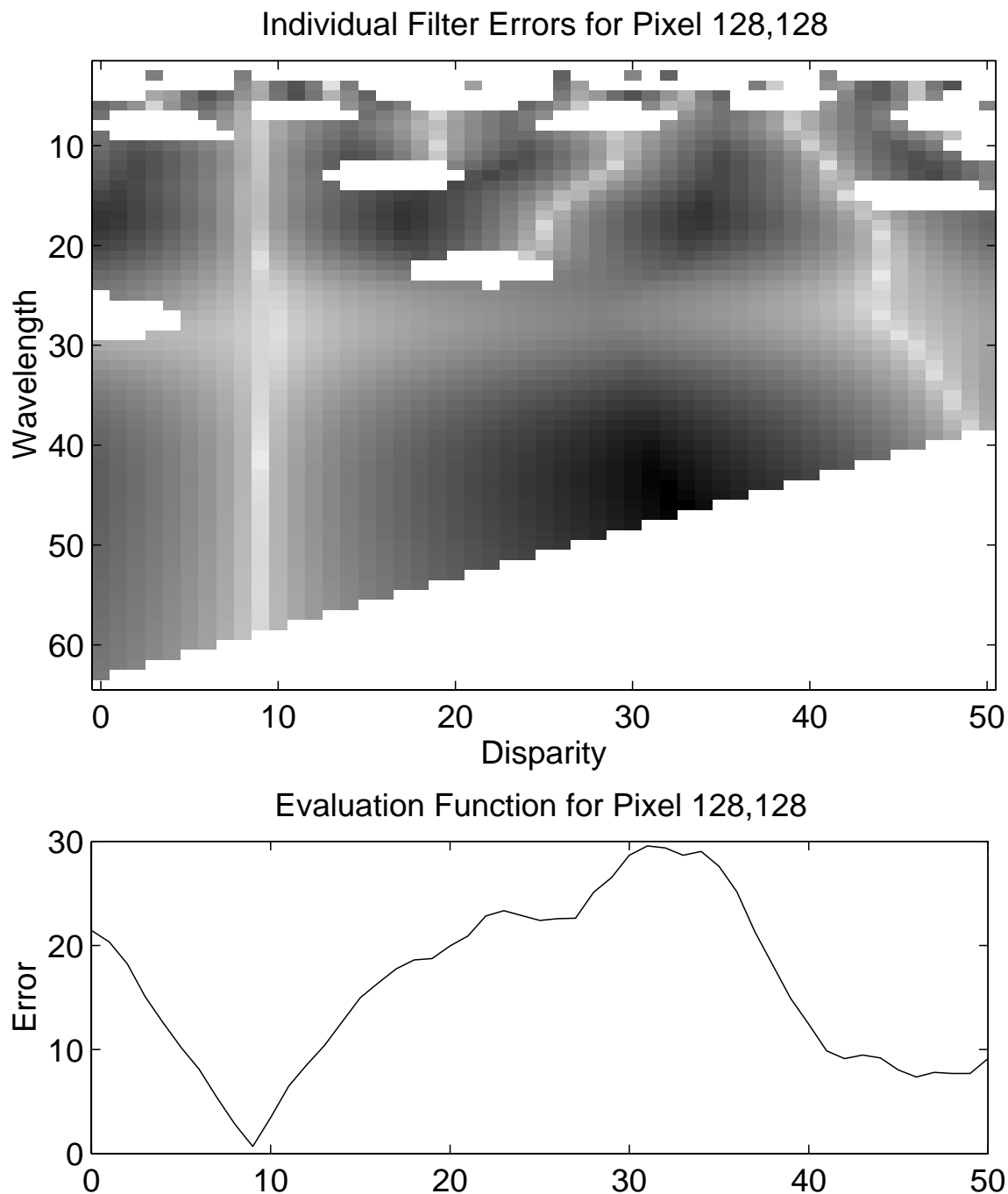


Figure 3.12: Filter contributions to the disparity estimate. Pixel intensity represents the value of the error computed by each filter's phase difference comparison, scaled by the magnitude. Darker pixels indicate larger error, and white spots represent areas where the scalograms contained no useful information. The evaluation function (below) is constructed by summing up the values in the columns and dividing by the number of filters used. This plot represents the disparity computation at pixel (128,128) of the stereo pair in Figure 3.14.

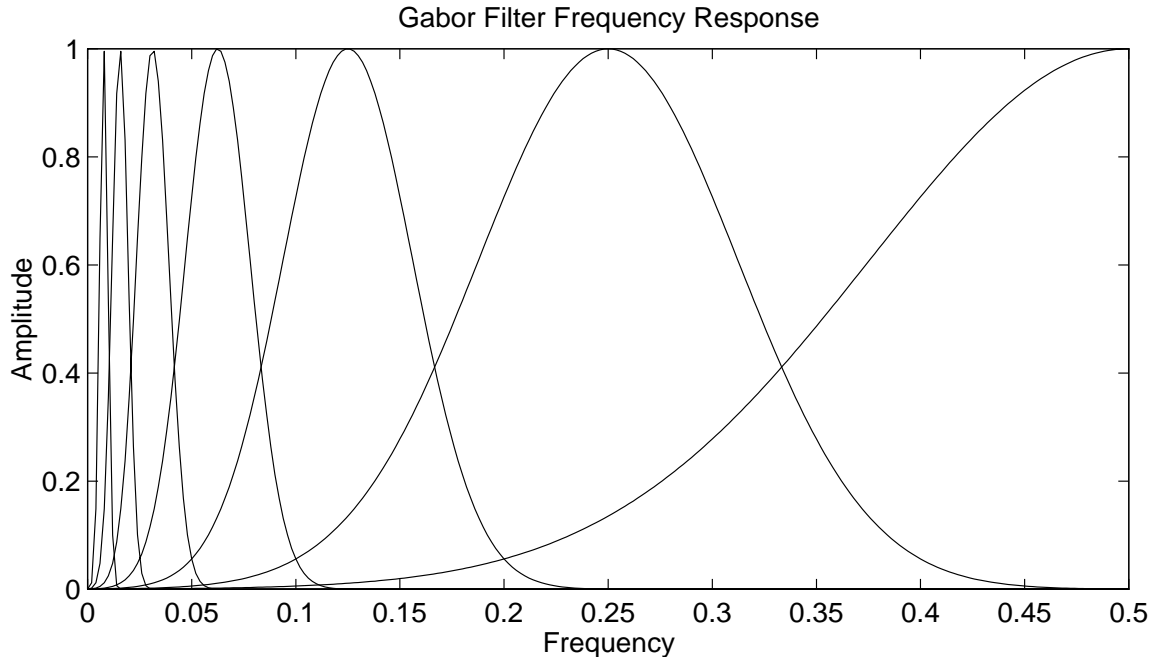


Figure 3.13: Gabor Filter Frequency Response. The Gabor filter magnitude in the frequency domain for filters with wavelength 2, 4, 8, 16, 32, 64, and 128 is presented. The tuning parameters are $\sigma_f = \frac{1}{6}$ and $m = 4$.

surement from each filter (equivalently, each scale) makes a contribution independent of the others, those with unreliable measurements can be safely omitted from W . The method also has the interesting property that it can be used to evaluate the likelihood of disparities much larger than the size of the filters used, when the images being studied have non-interfering periodic texture (i.e., they exhibit specific magnitude peaks in the frequency domain). This is because the evaluation function combines the filter outputs simultaneously, and combines estimates across frequencies based on the content of the signal; the signal determines the scales that are used. Coarse to fine methods, in contrast, use scales that are related in a strictly hierarchical manner and are fixed in advance. For example, given an input signal that is the sum of n sinusoids with wavelengths $\lambda_1, \dots, \lambda_n$, coarse to fine methods can only compute disparities half as large as the largest Gabor filter in the analysis set, whereas our method can evaluate disparities as large as $\text{LCM}(\lambda_1, \dots, \lambda_n)$, a potentially much larger quantity.

However, this version of the direct phase approach has a problem. Because this is a direct phase method, $\Delta\phi$ is calculated by subtracting Gabor phases from the two scanlines

at the *same-numbered pixel* (witness the c in Equation 3.15). From Section 3.2.1 we know that the window size of a Gabor filter varies as a function of wavelength, so filters with small wavelengths (high frequencies) will use relatively small windows for their input. In the presence of a disparity larger than the smallest wavelengths in W , this means that the phases measured by the corresponding filters will be from parts of the image that not only do not overlap, but may not even represent the same object's texture. This problem is faced by all direct phase methods, and provides one motivation for the coarse to fine approach used by others (Westelius, 1995; Weng, 1993; Fleet et al., 1991) and their assumption that the true disparity is less than the filter wavelength. But one cannot depend on responses at the lower frequencies, since the phase measurements there might be unreliable; Chapter 4 will show this in more detail. Fortunately, there is an alternative.

3.7 Our Stereo Method

Our stereo method is constructed out of the pieces described in the earlier sections of this chapter. We compute disparity using many Gabor filters, caching the outputs in 2D image scalogram matrices. Those measurements assumed to be unreliable are flagged and do not enter into the computations. A search over candidate disparities is made at each pixel, and the ideal phase differences associated with each disparity are compared against those measured in the images. A given disparity corresponds to a single, unique phase-difference curve; the curve that best fits the measured phase differences will give our disparity estimate.

Pseudocode for this algorithm can be found in Table 3.2. Matlab code for performing the computation is available on the web from *Mark Maimone's Index Page*.³ This method will be extended in Chapter 5 to address the perspective foreshortening issue, but we outline the basic method here. Let a pair of images be given. We assume that these images have the same dimensions, and were taken with a pair of cameras whose optical axes are parallel (so that the epipolar constraint holds). If the cameras are properly calibrated, this constraint reduces the matching problem from two dimensions to one, because like-numbered scanlines in the two images are guaranteed to correspond to the same plane in the world. We can thus limit the following discussion to the problem of matching a single row in the image pair.

Pick a pair of corresponding rows, each containing n pixel intensities. Compute the

³<http://www.cs.cmu.edu/~mwm/>

Given: A pair of greyscale images, and a list of candidate disparities.

For each row

 Compute Left and Right Scalograms L and R

 For each column c

$$W_L(c) = \{\lambda : \rho_L(c, \lambda) \text{ exists and } \phi_L(c, \lambda) \text{ is reliable}\}$$

$$W_R(c) = \{\lambda : \rho_R(c, \lambda) \text{ exists and } \phi_R(c, \lambda) \text{ is reliable}\}$$

 For each column c

 For each disparity d such that pixel $c + d$ exists

$$W = W_L(c) \cap W_R(c + d)$$

$$error = \frac{1}{|W|} \sum_{\lambda \in W} \rho_L(c, \lambda) \text{AbsDiffMod} \left(\Delta \phi_{ideal}(|d|_1, \lambda), (\phi_L(c, \lambda) - \phi_R(c + d, \lambda)) \right)$$

 Return d that yields minimum $error$

Table 3.2: Pseudocode for the basic phase-based stereo algorithm.

scalograms of these rows, using particular Gabor filter parameters (e.g., $m = 4$ and $\sigma = \frac{1}{6}$). The scalogram for a given row is a two dimensional matrix of complex numbers: for convenience in discussion, we split it into two matrices, magnitude ρ and phase ϕ . So $\rho_L(c, \lambda)$ is the magnitude of the left row's scalogram entry at pixel c with wavelength λ , $\phi_L(c, \lambda)$ is the phase of the left row's entry at that point, and similarly for the right row with ρ_R and ϕ_R .

To compute the disparity for a given pixel number c , we first mark unreliable estimates using Equation 3.10. Having thus ignored potentially unreliable phase values, we use the remaining left and right phases to fit phase differences. First we enumerate a list of possible disparities and compute the “ideal” phase difference curve for each disparity using Equation 3.11. If we let W be the set of wavelengths whose outputs are considered reliable, then the error for a particular disparity is given by the version of Equation 3.14 found in Table 3.2. Finally, we select the disparity that exhibits minimum error as the result for this pixel.

Note that by enumerating candidate disparities, we greatly reduce the phase wraparound problem mentioned in Section 3.6.1. Instead of comparing phase values at the same-numbered column, we use columns that are actually expected to correspond (i.e., we include disparity

d in the phase difference ($\phi_L(c, \lambda) - \phi_R(c + d, \lambda)$) part of Table 3.2's *error* computation). Thus the ideal phase differences need only reflect the remaining *subpixel* disparity (see the $|d|_1$ term in Table 3.2). Since the maximum remaining disparity will be one pixel, and the highest filter frequency will never have a wavelength of less than two pixels, the predicted phase differences will never wrap around. There is no such restriction on the *measured* phase differences, however, so care must still be taken to use the AbsDiffMod function from Equation 3.16 for the comparison.

Some form of phase interpolation is required to achieve the highest precision. However, none is explicitly specified in the implementation given in Table 3.2. The reason is that this interpolation has been encoded into the computation of the ideal phase difference. Rather than interpolate between the measured phases to find the perfect match (with zero disparity and hence uniformly zero ideal phase differences), we tune the ideal difference to compensate for this subpixel shift. So we are using a form of the direct phase method here, but since it is only applied to subpixel disparities we will not incur the wraparound problem.

3.7.1 Sample Imagery

Some example images with stereo results and ground truth are presented in Figures 3.14 and 3.15. The results illustrate some of the limitations of using 1D filters: in particular, errors at depth discontinuities are more jagged than would be expected from a method using 2D filters. And although the average RMS errors seems high (being greater than 1 pixel), it is difficult to compare with other methods since few others have reported ground truth results over complete images. Visual inspection of the area in the center of the images reveals that the error there is usually much less than 1 pixel.

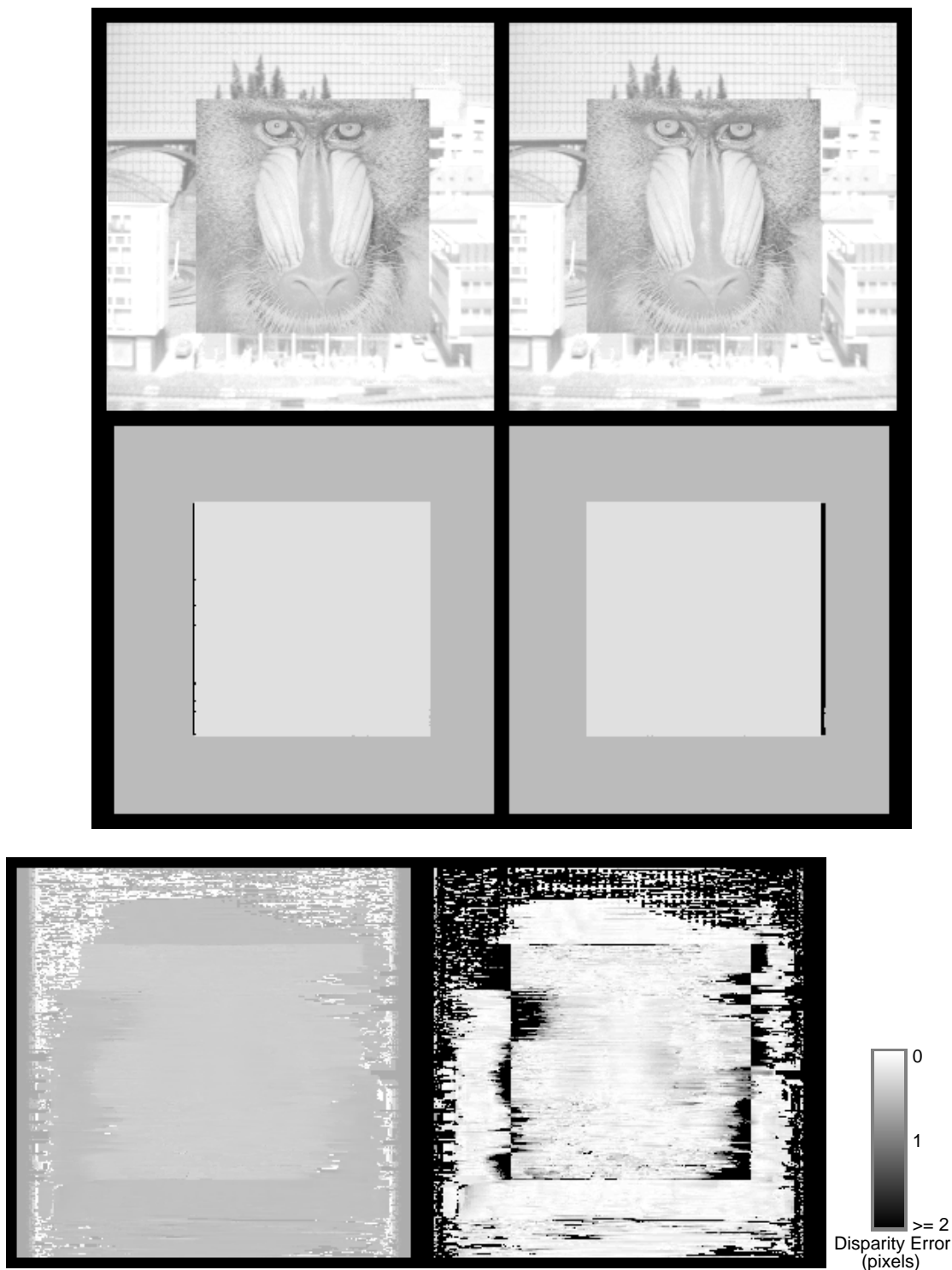


Figure 3.14: Synthetic Stereo Pair with Ground Truth. The top two images form the stereo image pair of two frontoplanar surfaces, the middle images are the associated disparity maps with occluded areas marked in black. The lower left image is the disparity map computed by this method, evaluating 301 disparities from 0 to 30. The lower right image is the difference between the stereo disparity and the ground truth. Mean error over the entire figure is 2.41 pixels (variance 23.09); all errors ≥ 2 are in black. (intensities not to scale)

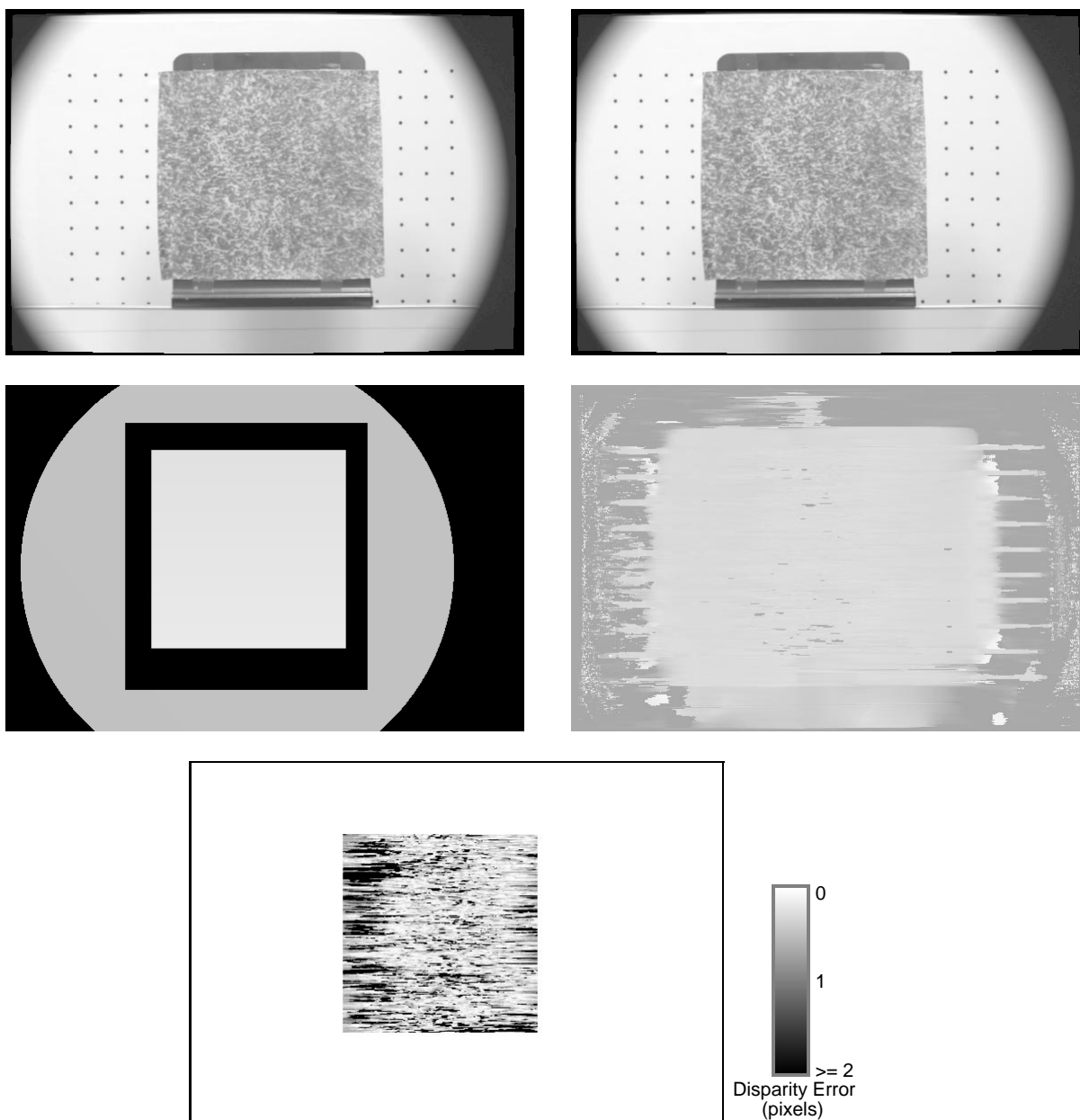


Figure 3.15: Actual Stereo Pair with Ground Truth. The top two images form the stereo image pair of a flat sheet tilted back slightly, and have been corrected for lens distortion. The middle images are the ground truth data on the left (black areas have unknown ground truth), and disparities computed by our method on the right (161 disparities, 0-40). The bottom image shows the differences between our estimates and the ground truth over the central target, scaled from 0 to 2 (all errors ≥ 2 appear black). Mean error in that region is 1.27701 pixels with $\sigma = 2.8446$. (intensities not to scale)

Chapter 4

Accuracy: Dealing with Ambiguous Matches

Time flies like an arrow; fruit flies like a banana.

— Groucho Marx

A problem shared by virtually all stereo methods is that of *ambiguous matches*, or false targets. When two or more parts of an image pair are similar in appearance, as can happen when a repetitive pattern like a checkerboard or picket fence is present, a part of the pattern in one image might seem to match several parts in the other. When this happens, when there are *multiple* potential correspondents for a given pixel, an *ambiguous match* is said to exist. This problem arises because image points are compared in a purely local fashion, and the effects can be seen in disparity estimates that are wildly inaccurate.

In this chapter we will show examples of ambiguous matches, explain how their presence affects several stereo methods, and demonstrate that our phase-based method compensates for the problem more effectively than other methods. The *Disparity Space* provides the framework for this analysis, and motivates a generalized representation of disparity.

4.1 The Problem of Ambiguity

Ambiguous matches are second only to modeling and calibration errors as a source of inaccuracy in stereo vision systems. By modeling errors, we mean the violation of any assumptions

about the image formation process that went into producing a particular stereo pair. Stereo methods by their nature make many implicit assumptions about their input images, e.g., that the units of horizontal and vertical spacing in the real world are identical, the left image was taken by the left camera, surfaces are diffuse, the lens caps have been removed, etc. And by calibration, we mean any computation that will influence (independently of the stereo image data) either the range of disparities being checked or the process by which they are checked in the stereo method. It should be obvious that any part of the image acquisition process that violates modeling assumptions or reduces the accuracy or precision of the calibration results could cause inaccurate results to be generated. Yet even when all of these constraints on the physical system have been met, the correspondence problem remains.

The correspondence problem lies at the heart of all stereo methods. In many methods the determination of correspondence is based solely upon information found in the stereo image pair, without knowledge of the actual 3D scene structure or image segmentation. Yet even though a stereo method returns the “best” answer according to its model, sometimes the “best” correspondence is not the correct one. When an improper correspondence is established, the resulting disparity can be *very* inaccurate.

This problem can be alleviated somewhat by incorporating more data and using better models in the stereo method. There are many such techniques in the literature: using many more than two cameras as in multibaseline stereo (Kanade et al., 1995; Webb, 1993), using more realistic models of image formation (Bhat & Nayar, 1995; Belhumeur, 1995, and Chapter 5 in this text), and interpreting the same data at several scales (the subject of the following sections). Although they differ in their approaches, all of these techniques share the common goal of finding the single best disparity estimate for a given pixel.

But sometimes the “best” answer is not the desired one. Figure 4.1 shows an artificial example of this phenomenon. In the figure a stereo image pair is actually embedded into a single image; this type of image is called an *autostereogram*.¹ In an autostereogram, it would appear that the left and right stereo images are identical. And when stereo processors are given the same left and right stereo images, most (including your eyes) will conclude that the “best” disparity map is a uniform zero shift, i.e., the images contain no interesting 3D structure. However, in this case the interesting stuff happens when you realize that the left and right images of interest are *not* identical; in Figure 4.1, the left image does indeed start

¹Some people have difficulty seeing the depth in images like this. Give it your best shot.

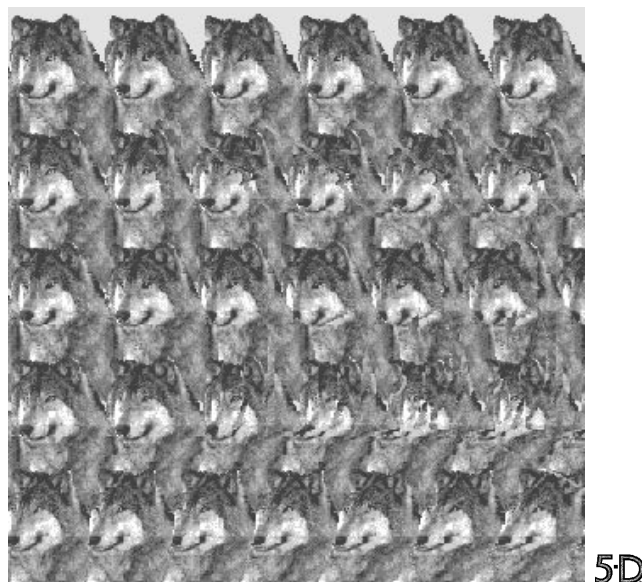


Figure 4.1: A 340x340 autostereogram of a wolf. See Figure 4.24 for an elucidation of the embedded structure. Reprinted with permission from Blue Mountain Arts, Inc.

at the left border, but the right image starts about one sixth of the way in. Running a coarse to fine stereo method on such pre-shifted images brings out the structure that was embedded into the original figure (see Figure 4.24 after trying it yourself), and demonstrates that in this case it is not the “best” disparity that is of interest, but rather the “second best.”

The rest of this section will give us the language and concepts required for our discussion of ambiguity. The sections that follow will discuss both the impact of ambiguity on stereo methods, and extensions to stereo methods that will model, if not eliminate, problems due to ambiguous matches.

4.1.1 Definitions

An *ambiguous match* occurs when a stereo method is unable to determine a unique correspondent for a particular pixel. Ambiguity can be viewed in two ways: as an inherent property of an image, or as an artifact of a particular stereo method. It is certainly possible to construct inherently ambiguous stereo images. Consider two images of a frontoplanar checkerboard, where the checkerboard occupies the entire fields of view: the alignment of the squares cannot be determined from the images alone. Naturally, these types of images will cause problems for stereo matchers, but they are not the only source of ambiguity.

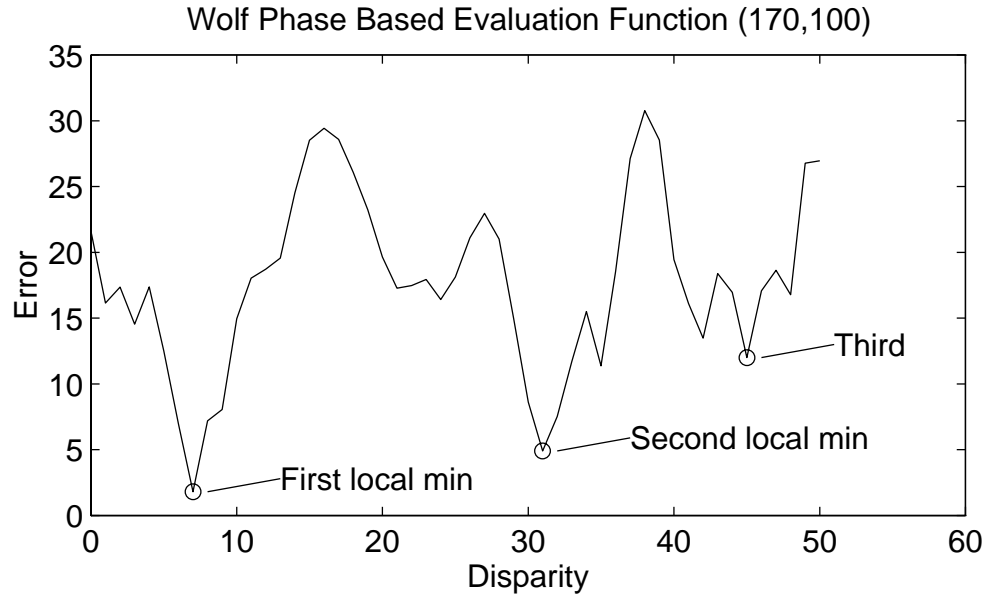


Figure 4.2: Evaluation Function for pixel (170,100) of Figure 4.1 (with its pre-shifted counterpart), computed using our phase-based method with linear frequency samples. This information is presented in the context of a complete scanline in column 100 of Figure 4.3. Darker pixels denote less error, so the dark stripe from pixel number 75 to 200 at disparity 7 represents the best guess.

The limitations of stereo matchers can also cause image regions in stereo imagery to appear ambiguous.

Informally, ambiguous matches will occur when the small area around one pixel “looks like” more than one area in the second image. This is because most stereo algorithms reduce the computational complexity of the matching process by determining correspondence using only small image patches at some stage of their processing. Even the so-called multi-scale coarse to fine methods only make local decisions, based on one or possibly two levels of hierarchy. It is this narrowing of attention, the making of a hard decision based on local information at a single scale, that makes stereo algorithms especially susceptible to ambiguous matches.

The chance of encountering ambiguity looms large in images with repetitive texture, since many image patches will appear to be nearly identical, but ambiguity can also occur even when the images do not exhibit obvious repetition. For many algorithms, all that is necessary is that two small patches appear identical, where “small” and “identical” are defined by the search window size and evaluation function used by the algorithm.

The presence of ambiguous matches can be determined by inspecting a pixel's disparity space, i.e., the result of the evaluation function. If the evaluation function profile is unimodal, i.e., if it has a single minimum and a monotonically increasing first derivative, then a unique disparity exists at that minimum and there is no ambiguity. However, if as is often the case the evaluation function profile is *not* unimodal, then any of the alternate local minima must be considered potential matches as well, and therefore an ambiguous match exists (see Figure 4.2). The degree to which the match at that pixel is ambiguous can be informally determined by considering the ratio of values of the evaluation function at the lowest minima: the closer the ratio is to 1, the more ambiguous is the match.

$$\text{Ambiguity Factor } \alpha = \frac{\text{error (First local minimum)}}{\text{error (Second local minimum)}} \quad (4.1)$$

As a special case, if the function has only one minimum, the error at the second minimum is defined to be infinity (yielding an ambiguity factor of zero).

Computation of the ambiguity factor in actual images is difficult. While it is easy to locate the first local minimum (the global minimum) in a vector, it is a bit harder to isolate the second minimum, because it is not necessarily the vector element with second-lowest value. Rather, it is the element with minimum value where the first derivative changes locally from negative to positive, outside the influence of the original minimum. Of course, in real signals there will be small perturbations in the evaluation function that should not be treated as actual minima, so these will have to be ignored in some way. We have created a useful (but not perfect) heuristic for finding the first n minima, that can be used to aid in the computation of the ambiguity factor. It works automatically by finding the minimum value, locally fitting a Gaussian curve to the values surrounding the peak, eliminating the influence of the Gaussian and iterating on the newly formed signal. The process repeats until the desired number of peaks or an evaluation function threshold is reached. This heuristic is discussed further in Section 4.4.1.

In summary, an ambiguous match occurs when a pixel's evaluation function exhibits multiple local minima. The presence of an ambiguous match does not necessarily mean the disparity estimate will be wrong, just that the potential for a false match exists.

4.1.2 Disparity Space

Having discussed the presense of an ambiguity at a single pixel, we now consider its effect on a scanline.

To understand how ambiguous matches arise we must consider the way stereo methods assign disparities. Recall that the goal of all filter-based stereo algorithms is to associate a unique depth value with each pixel in the input images. Most often a disparity value is assigned to a pixel; other values are possible (e.g., a pixel might be marked as occluded), but for now we only consider disparities. Theoretically speaking, each pixel in a stereo image can potentially be assigned any of a wide range of disparities. Stereo methods work by eliminating improbable disparities until only the most likely ones remain. To understand this process of elimination, we want to know the likelihood that each disparity will be assigned to each pixel. A useful visualization tool for this task is the disparity space.

The concept of *disparity space* is straightforward. Instead of associating a single disparity with each pixel, we associate a vector representing the likelihood of *all* disparities that might be assigned to it. This notion can be applied to a single pixel, a row of pixels, and even a matrix of pixels (i.e., an image). So while the disparity *map* for an image is 2D, the disparity *space* over a whole image will be 3D. To make it easier to visualize, we usually only consider the disparity space for a single pixel (by plotting the disparity evaluation function at that pixel, as in Figure 4.2), or for a single scanline (by creating a 2D image where intensity represents likelihood, as in Figure 4.3). Think of the disparity space as a visualization of the penultimate step in disparity computation; the disparity for a given pixel is determined by finding the disparity space vector element with minimum error.

Figure 4.3 contains an example of the disparity space matrix associated with a pair of image scanlines. The horizontal axis is indexed by Pixel Number and corresponds to the original scanline, and the vertical axis is the list of candidate disparities. Pixel intensities in Figure 4.3 correspond to the match error computed by the evaluation function; column 100 encodes the same information as in Figure 4.2. Dark regions in the image denote areas where the error is large, i.e., they represent those disparities that are least likely to become associated with the pixel numbered below. The disparity space makes the presence of ambiguities on a scanline easy to see. Any column with more than one bright spot represents a pixel likely to be involved in an ambiguous match.

The disparity space is a natural representation for search-based methods like dynamic

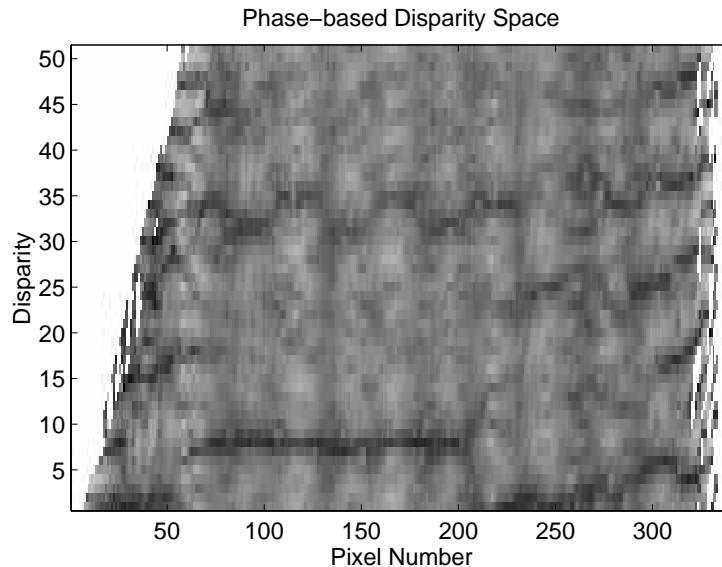


Figure 4.3: Disparity Space for row 170 of Figure 4.1 (with its pre-shifted counterpart), computed using our phase-based method with linear frequency samples. Darker pixels indicate lower matching errors, i.e., the most likely disparity estimates. Column 100 can be seen in an expanded view in Figure 4.2.

programming and our phase-based method from Chapter 3. In fact, the pseudocode for our algorithm presented in Table 3.2 computes exactly this information in its innermost loops. Note that by the time our phase-based method has constructed this image, all of the work in the frequency domain has been completed. The phase measurements have already been combined by the evaluation function to produce the likelihood that a particular disparity is appropriate. So the disparity space framework represents purely spatial phenomena, and can be applied directly to stereo methods that work exclusively in the spatial domain as well. We will elaborate on the benefits of the frequency domain in Section 4.3.3, but for now simply discuss the characteristics of the disparity space.

Once the disparity space is constructed, the usual way to generate a single disparity estimate for each original pixel is to find the minimum entry in each column, and call its associated disparity the best answer. We will often highlight these minimal disparities in disparity space plots by connecting them together across adjacent columns. In this way we can see not only the disparities for a particular scanline, but also the shape of the evaluation functions that led to their extraction. An estimate of precision can be provided as well, by measuring the curvature of the evaluation function around the minimum value (i.e., the

second derivative of the evaluation function with respect to disparity, evaluated at that disparity at which the evaluation function has its minimum). Although one might consider using the raw evaluation function outputs as a measure of confidence and therefore precision, they are likely to be too variable to be of effective use. Simple changes in gain between the two cameras, or changes in object reflectance due to viewpoint change, could cause massive shifts in the amount of error, and therefore it is not in general robust enough to function as an independent precision estimator.

Disparity space images are not limited to search-based methods, they can also model the results of coarse to fine algorithms (e.g., the algorithm presented later in Table 4.1). The results are not quite identical to those of search-based methods, however. The main difference is that coarse to fine methods by their very nature do not compute an evaluation function at every possible disparity. Instead, they navigate through the disparity space in large (“coarse”) steps, refining and evaluating only those subregions that meet some minimum criterion. As a result, the disparity evaluation functions associated with a given pixel are computed in completely different contexts, and therefore cannot be compared directly. So the argument above, that the disparity space was the penultimate step followed by minima extraction, is no longer appropriate. However, coarse to fine disparity spaces can still provide interesting visual information, especially when compared against a complete disparity space generated by a search-based method.

There are two complementary views of the disparity space for coarse to fine algorithms, each illustrated in Figure 4.4. The first view is identical to the one above, where we simply plot the value of the evaluation function between a pixel and a disparity. We cannot extract the minimum value to find the disparity in this representation, but it is still useful for visual inspection of the overall shape of the evaluation functions. The second view is the more appropriate one for coarse to fine methods. In this case the pixel intensities do not indicate the value of the evaluation function, but rather the number of the *highest scale* used to compute the evaluation function for a given pixel/disparity pair, with coarse scales numbered lower than fine scales. This plot has the same advantage for coarse to fine methods that the other plot had for search-based methods: the best disparities can be found by visual inspection, since the disparity returned by the algorithm will always be one of those searched at the highest scale. This plot also provides an explicit road map, demonstrating how the method chose to navigate through the detailed disparity space used by a search-based

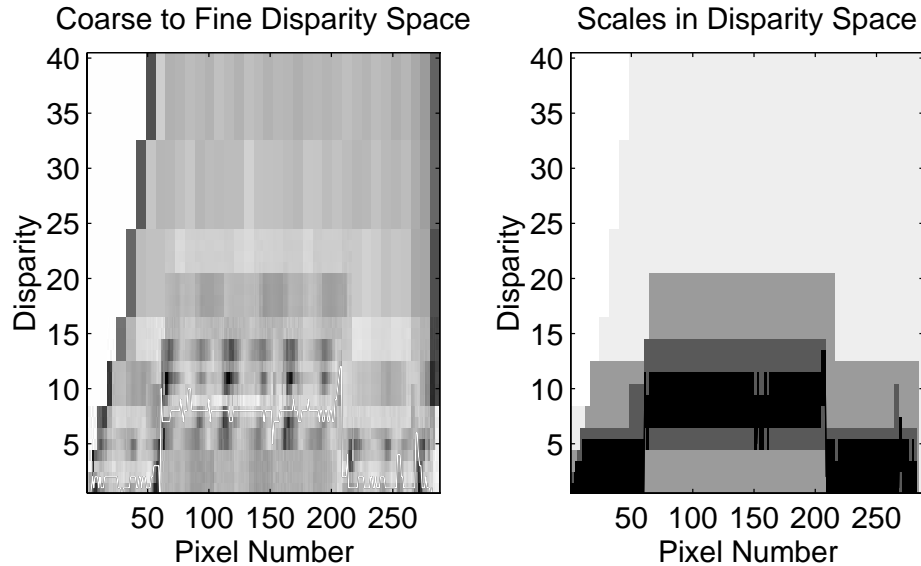


Figure 4.4: Coarse to Fine disparity space and associated scales; fine scales are darker than coarse scales. These results are also from row 170 of Figure 4.1, as are those in Figure 4.3.

method.

The actual content of the disparity space will depend on the evaluation function of each particular stereo method. It provides a way to visualize the search strategy used by that method, but should not be regarded as an inherent property of the image pair. Rather, it summarizes the interpretation of the image pair according to a particular method.

These disparity space representations provide the framework for our discussion of ambiguity.

4.2 Effect of Ambiguity on Stereo Methods

The usual effect of ambiguity on stereo methods is the production of outliers, sharp spikes in the disparity map. In the worst case, such as will be illustrated in Figure 4.7, *none* of the disparities is computed correctly.

4.2.1 JISCT Results

Our first illustration of the problem comes from a recent survey and evaluation of stereo methods, the JISCT Stereo Evaluation study (Bolles et al., 1993). Figure 4.5 is a stereo

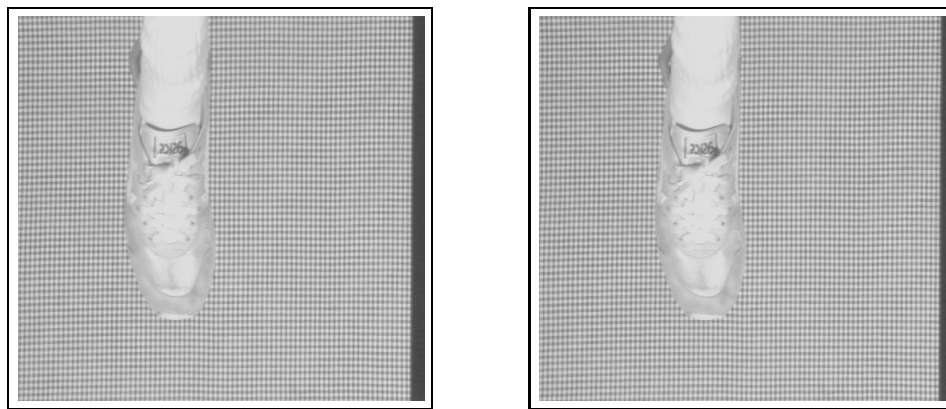


Figure 4.5: Shoe Stereo Pair. These are images SHOE-0 and SHOE-2 from the JISCT Stereo Evaluation study; each is 480×512 pixels².

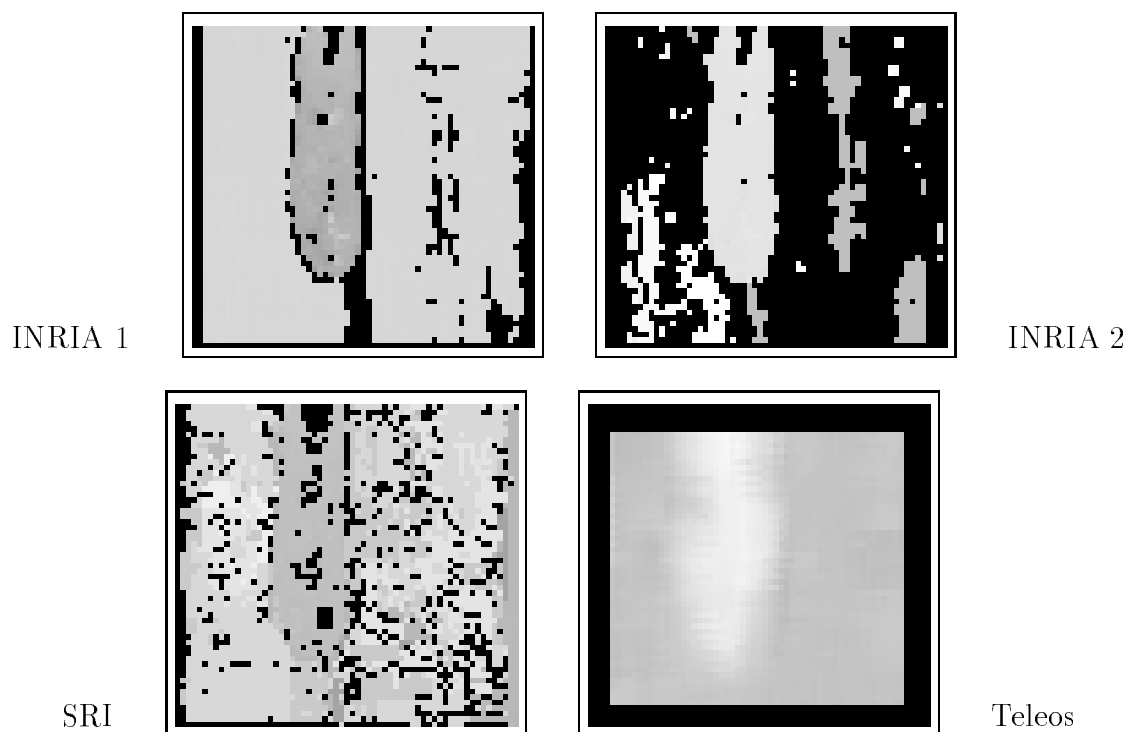


Figure 4.6: Results of several stereo methods on Figure 4.5 as presented in the JISCT Stereo Evaluation study. Clockwise from the top left are the INRIA 1, INRIA 2, Teleos, and SRI results; each is 59×63 pixels². Black pixels were marked as unknown by the algorithms.

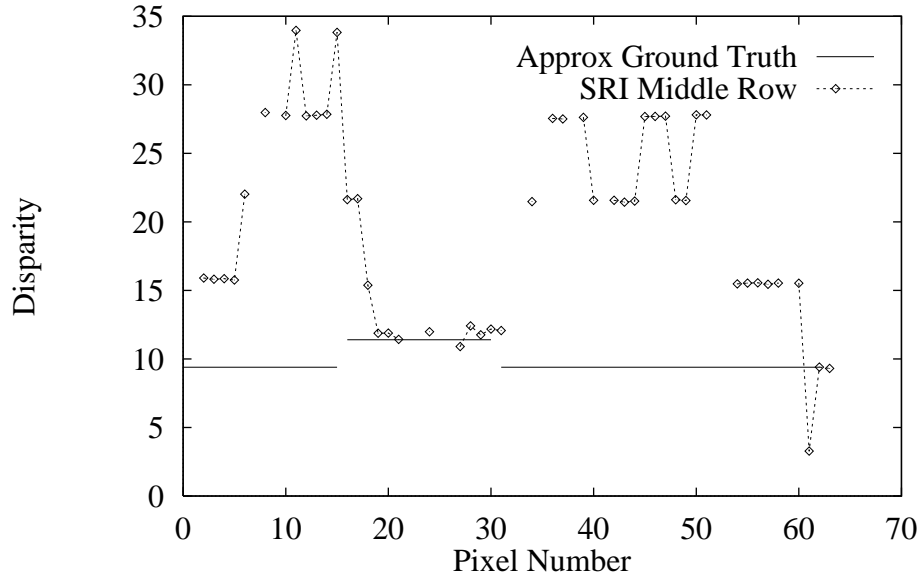


Figure 4.7: Detail view of SRI disparity results and approximate ground truth for the middle row of the shoe stereo pair disparities in Figure 4.6.

pair from that survey of a tennis shoe over a mostly-checked background. Because the background is so regular, almost every pixel in it has multiple potential correspondents within a small area. Figure 4.6 shows the coarse disparity maps computed by several stereo methods in the study, illustrating the problems caused by the checkered pattern. In many of these results, the disparities computed within that pattern either jump around a lot (SRI and INRIA 2) or have significant gaps (all but Teleos).

These problems can be attributed to the repetitive nature of the input image. Each of the methods exhibiting spotty results uses correlation with 11x11 pixel filters as the final step (before interpolation) in assigning correspondence. Because this image is so regular, one match looks very much like another. Thus the background image (a carpet whose 3D shape is actually flat) appears to have many spikes in it. This is in spite of the fact that both methods explicitly try to eliminate such spikes: the INRIA methods employ morphological shrinking of the disparity map to remove outliers, and the SRI method uses multiple passes over the input data to eliminate unreliable matches. Yet you can see just how irregular the results are in Figure 4.7, which plots the SRI disparities for the middle row of the image against the approximate ground truth, generated by manual inspection. Not only are the results erratic, matching any of several candidate disparities, but most of the disparity estimates do not even match the correct piece of the pattern.

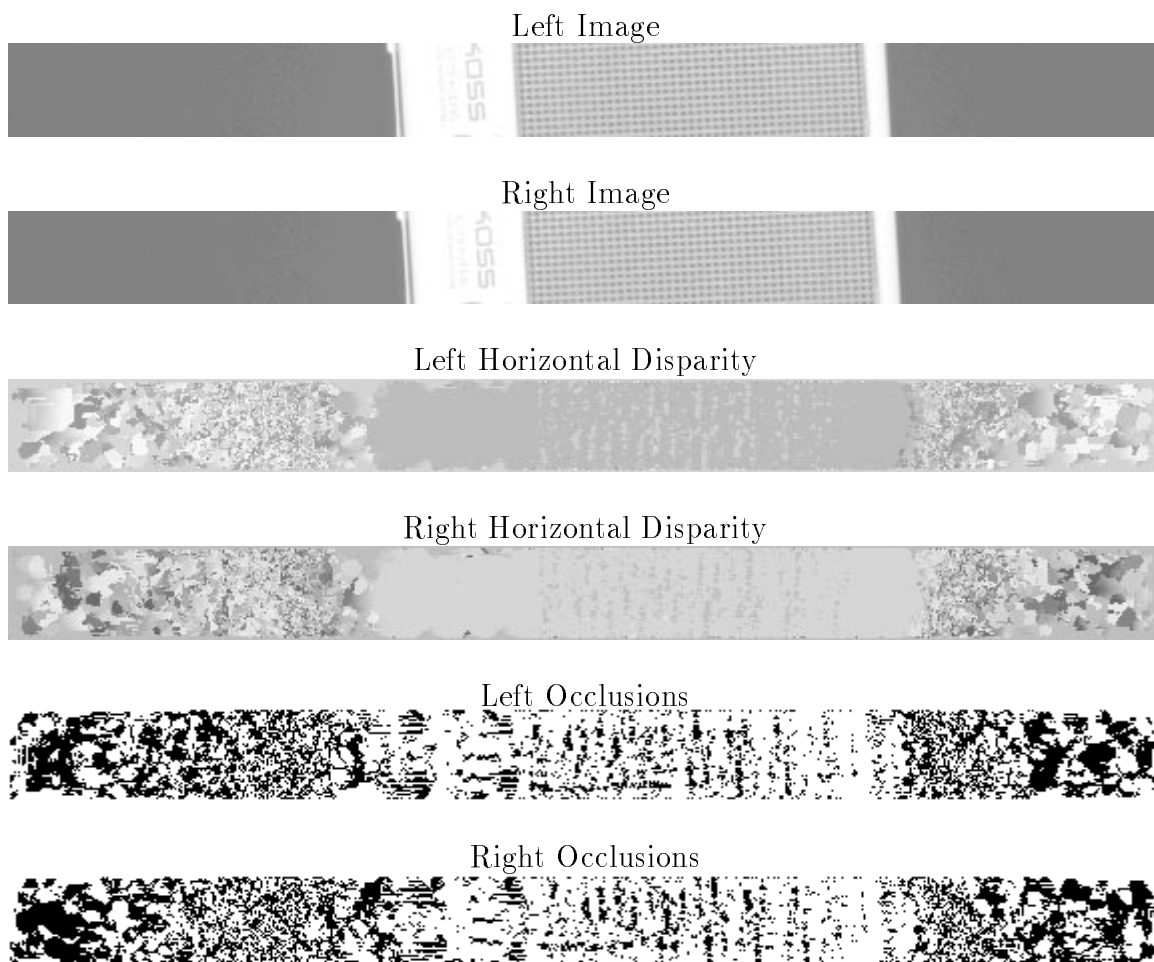


Figure 4.8: Results of Jones stereo method on middle region of speaker images.

Of the four methods in the JISCT study, only the Teleos algorithm was able to match the background successfully. Teleos' use of very large windows (from 25x25 to 90x90 pixels) helped them avoid getting trapped in the abundance of similar features at the finest scales, instead seeming to track stains in the carpet whose influence can only be seen at larger scales (Bolles et al., 1993). But this feature comes at a cost; since the search windows are so large, disparity is not computed at the pixels around the border of the image. Even more importantly, the shape of the tennis shoe (readily apparent in all the other methods) is greatly washed out, blurred by the extra-large search windows.

4.2.2 Jones' Method

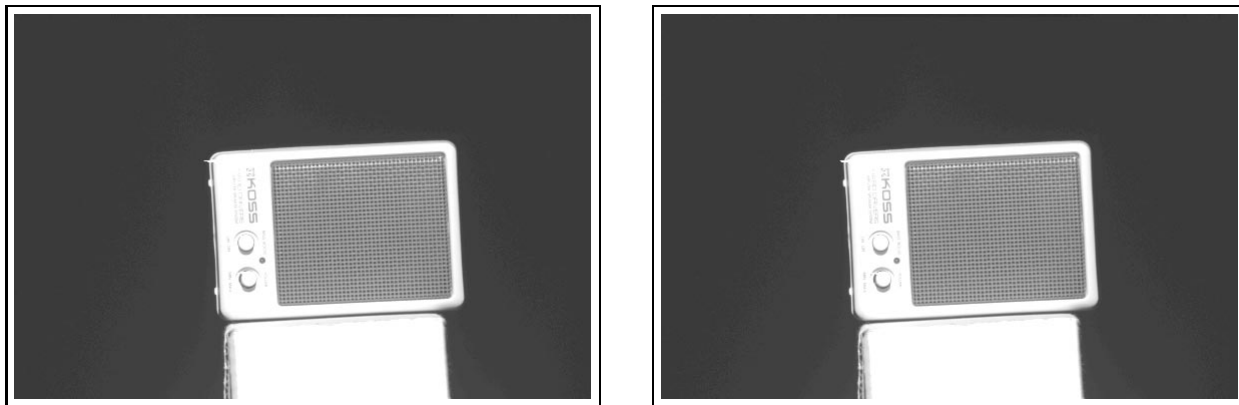


Figure 4.9: Speaker Image.

Another consequence of the disparity spikes caused by ambiguous matches is their effect on postprocessing techniques. For instance, Jones' stereo method (Jones, 1991) includes a special postprocessing step that addresses occlusions. The method compares disparities computed for the left and right images; if corresponding pixels do not have the same disparity, they are marked as occluded.² Given sharp disparity spikes, however, this step results in many spurious occluded pixels, as can be seen in Figure 4.8, especially in the area of the speaker grill. In this implementation of Jones' method, only the first pass and occlusion processing have been used.

4.3 Reducing the Ambiguity Factor

The goal of many stereo methods is to construct or modify the evaluation profile so that the Ambiguity Factor α is as small as possible. In the multibaseline method of (Kanade et al., 1995) this is accomplished by summing the evaluation functions across multiple image pairs. In a similar way, our phase based method sums evaluations functions from multiple *filter outputs*. However, the usual approach is to employ a coarse to fine refinement strategy. Coarse to fine methods smooth over the evaluation function, in an attempt to avoid falling into a local minimum.

In this section we will compare the coarse to fine approach with our phase-based filter combination method. We will see that although coarse to fine works well in some situations, in the presence of ambiguity its naive approach fails to resolve the ambiguity. We will

²A more robust method is outlined in Section 2.3.1 of this text.

Given: A pair of greyscale images L and R , search window size W , max disparity D , smallest resolution R , averaging filter *halve*.

Function *c2f*:

Matrix *prior*, *result*

Disparity results from coarser level, and this level

If (*resolution* / 2 > R)

prior = *c2f* (*halve* (L), *halve* (R), ...)

Else

prior = $\lfloor \frac{D+2}{4} \rfloor \cdot I[\text{resolution}/2, \text{resolution}/2]$

For each row r

For each column c

For disparity $d = -\frac{D}{2}$ to $\frac{D}{2}$

$error(d) = \frac{1}{W^2} \sum \sum_W | \text{left window } (r, c) - \text{right window } (r, 2 \text{ prior}[r/2, c/2] + d) |$

Pick d with minimum *error*, preferring values near $2 \text{ prior}[r/2, c/2]$ in a tie

result (r, c) = d

return *result*

Table 4.1: Pseudocode for the Coarse to Fine algorithm.

demonstrate that our phase-based technique does a better job of compensating for this ambiguity.

4.3.1 Coarse to Fine Method

A common approach to overcoming the problem of ambiguous matches in two-camera stereo is to employ a coarse to fine (or image pyramid) search strategy. The intent is to limit the number of disparities checked at a given resolution; the fewer disparities checked, the lower the likelihood of an ambiguous match. An additional benefit is that a wide range of disparities may be checked for relatively small cost; the coarsest levels have the least data, yet span the widest disparity range in the original image.

To see how such methods compare with our phase-based approach, we have constructed a “typical” coarse to fine stereo algorithm. The method is typical in that it starts with a SAD correlation of the coarsest versions of the images, then iterates using only the results computed at the previous level to restrict the search at the current level. Pseudocode for this method can be found in Table 4.1. In practise we have applied this algorithm using a 3x3 Gaussian smoothing filter ($\sigma = 1$), search windows that are 5 pixels wide at each level, a 5 pixel maximum disparity at each level, and 20 pixels as the smallest resolution.

There is one aspect of this type of method in particular that distinguishes it from our phase-based method. Although both approaches consider the scale space decomposition of the original images, the coarse to fine approach marches through the scale space in a fixed order without regard to the content of the image. Decisions made at the coarse levels become prior assumptions at finer levels, and are taken at face value. Figure 3.10 demonstrates the scales used: the results from longer wavelengths are used to constrain those computed at higher wavelengths. In contrast, our phase-based method considers the entire scale space profile as a unit, selecting only those scales whose data are known to be valid, and combining them without regard to a predefined order. We will demonstrate shortly the benefit our phase-based method derives from this approach.

4.3.2 Coarse to Fine Results

The coarse to fine approach to image analysis is a popular one for two reasons. It is efficient, allowing large image regions to be searched at relatively low computational cost; and often it reduces errors, smoothing over small variations in the evaluation function.

An assumption implicit in the design of coarse to fine methods is that evaluation functions have an overall “bowl” shape, with a unique minimum. Minor variations (i.e., local minima) are tolerated as well, because they will be smoothed over by the processing at the coarsest levels. Another aspect of that assumption, not often stated explicitly, is that the evaluation functions at *all* scales must exhibit that structure; even the coarsest image is assumed to have a well-behaved evaluation function.

Figure 4.10 shows how the evaluation function evolves with each progressive level for a pixel in the wolf stereogram (Figure 4.1). Contrasting the final shape in Figure 4.10 with the evaluation function computed by the phase-based method in Figure 4.2, we see that the coarse to fine method was able to focus attention on the correct region in spite of the

presence of the second and third local minima. So in this case the scale-based smoothing imposed by the coarse to fine method worked well, and allowed the correct local minimum to be extracted. You can view the final disparity map for this image (pair), as computed by the coarse to fine method, in Figure 4.24.

Problem with Coarse to Fine

Unfortunately, in some cases this model can lead to significant errors. The method worked well on the wolf stereogram because that image has a lot of texture at many scales; many details can be seen in the coarse images. But in some cases the coarse versions of images will *not* have sufficient structure. When that occurs, the method will still select a region in which to continue the search at a finer level, but that decision will be based on incomplete information. The coarse to fine strategy does not allow sufficiently for the possibility that a given level's evaluation function may not exhibit an obvious and unique minimum. Even though our implementation attempts to address this problem by preferring disparities near to the one predicted by the coarser level, it does not always succeed.

Consider the synthetic speaker grill in Figure 4.11. Although the original image contains much structure, and the correspondence between the two images seems obvious to a person, all of the interesting structure is washed out at coarser levels. Therefore, the coarse to fine search strategy cannot use the results from coarser levels to effectively constrain the search at finer levels. The resulting disparity image and coarse to fine disparity space for the middle row of the synthetic speaker grill illustrates the effect in Figure 4.12. Although the true disparity lies at a constant 23.52 pixels in the middle of the image, the coarse to fine method is unable to maintain the proper focus. Therefore most of its estimates are incorrect except at the very ends of the grill, where the grill and background colors cause a sharp edge to appear at all scales.

This effect is not limited to synthetic imagery, it occurs in real images as well, as can be seen in the results (and coarse images) from the Shoe stereo pair in Figure 4.13.

4.3.3 Solution using Phase-based Stereo

Our phase based method does a better job at addressing these problems. In this section we demonstrate that our method, as presented in Chapter 3, effectively reduces the ambiguity factor when the image pair contains enough information.

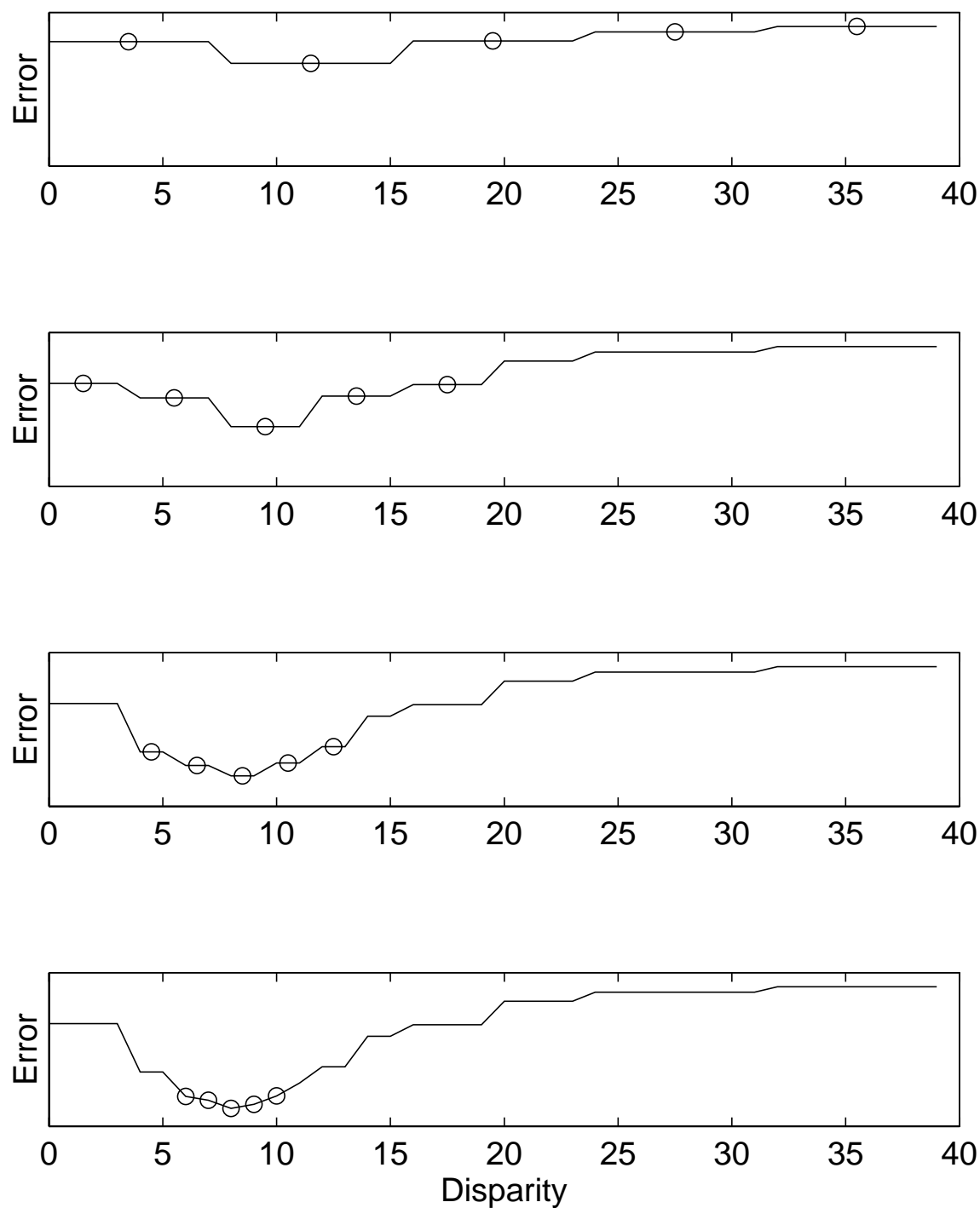


Figure 4.10: Evolution of the Coarse to Fine evaluation function at pixel (170,100) in Figure 4.1. The coarsest scale appears on top, the plots below demonstrate the successive refinement. Contrast the bottom plot with Figure 4.2.

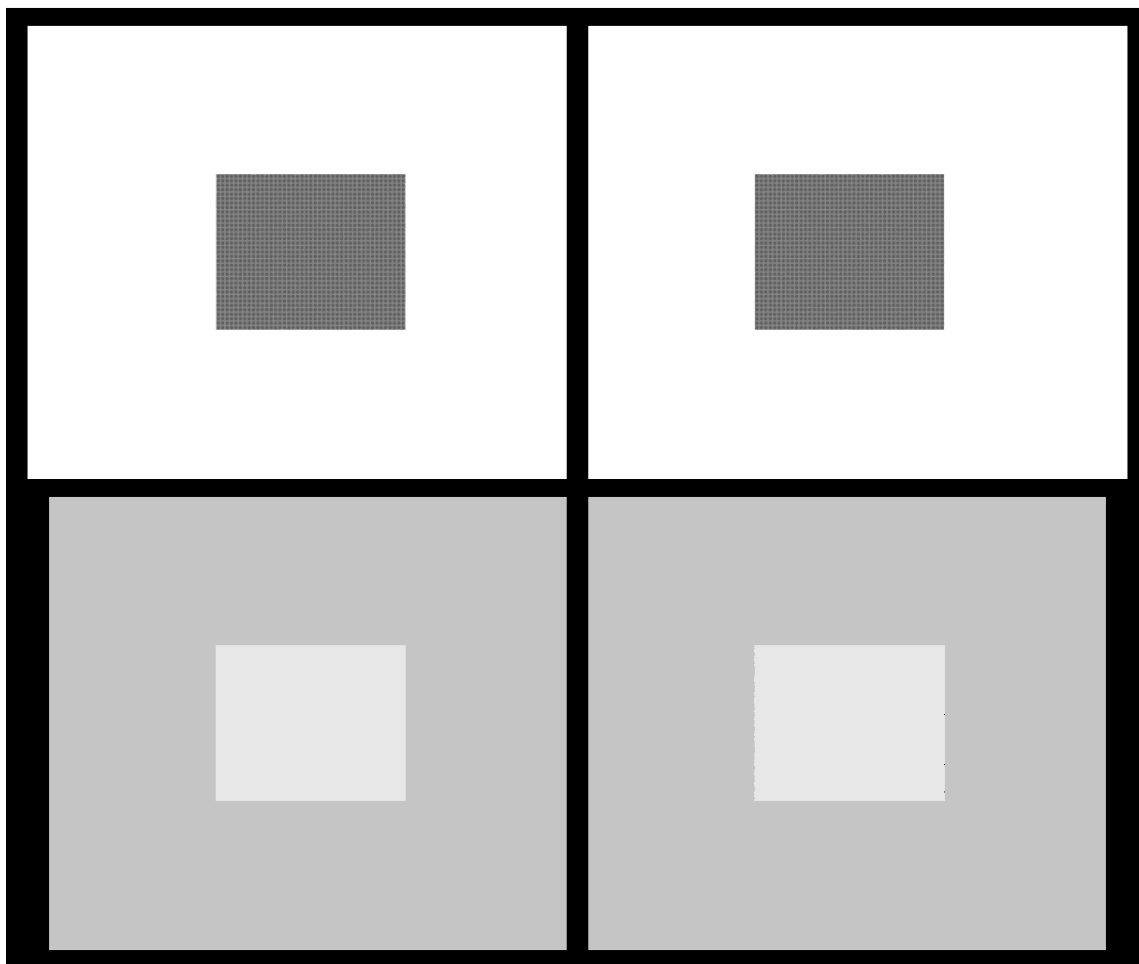


Figure 4.11: Synthetic Speaker Grill. Actual disparity is 18.3445 pixels, but the pattern repeats approximately every 3.8 pixels.

Our phase based method addresses the problem using the same tool as the coarse to fine method: a multiscale representation of the image. However, instead of enforcing an arbitrary ordering on the scales that are used, a more intelligent adaptive search is performed. In addition, results from different scales are combined according to the known relationship between the filters that created them. Because each scale generates a real-valued phase measurement, we have high-precision disparity estimates from *all* scales, not just the finest ones. These factors combine to yield a better behaved evaluation function in general, and in ambiguous areas in particular.

Confidence Estimates One of the advantages to the local spatial frequency framework is that each multiscale measurement includes a confidence estimate. In our phase-based

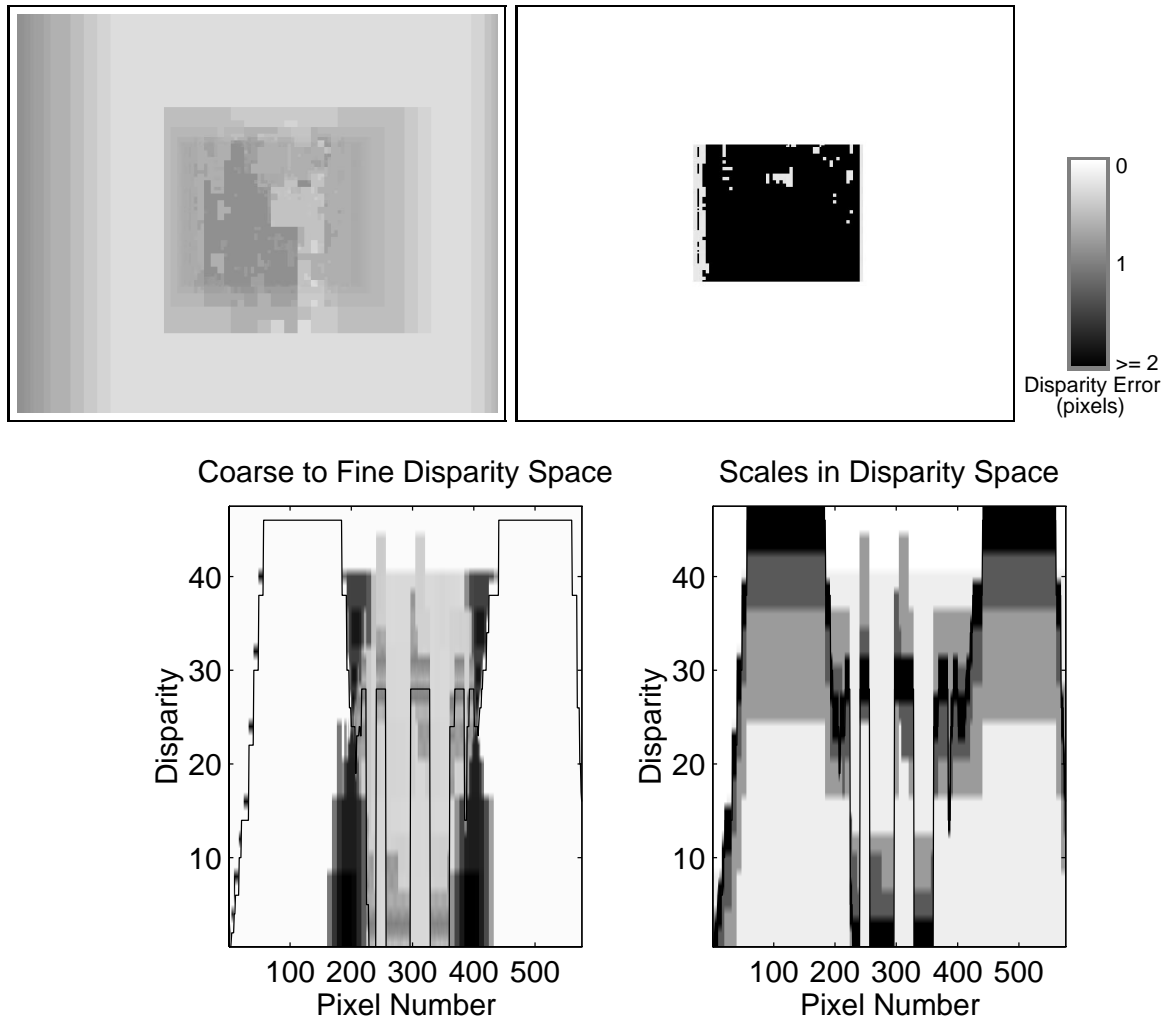


Figure 4.12: Coarse to fine results on Figure 4.11. Disparity map (upper left) has mean error 14.42 with $\sigma = 23.1308$ pixels, ground truth error image (upper right) maps all errors ≥ 2 to black. Lower plots are the coarse to fine disparity space (left) and scale space (right) for row 240 in the synthetic speaker grill. Ground truth for pixels 200–400 is 23.52 pixels.

method, the magnitude of the Gabor filter output is a measure of confidence in the phase value. More precisely, a combination of high magnitude and the constraint in Equation 3.10 provide a measure of confidence. This is an advantage over the coarse to fine approach in the spatial domain, which uses multiscale measurements without the benefit of an independent evaluation of their utility.

There are other phase based methods that use the coarse to fine structure, and therefore enjoy the benefit of a confidence estimate with their measurements. But they lack the ability

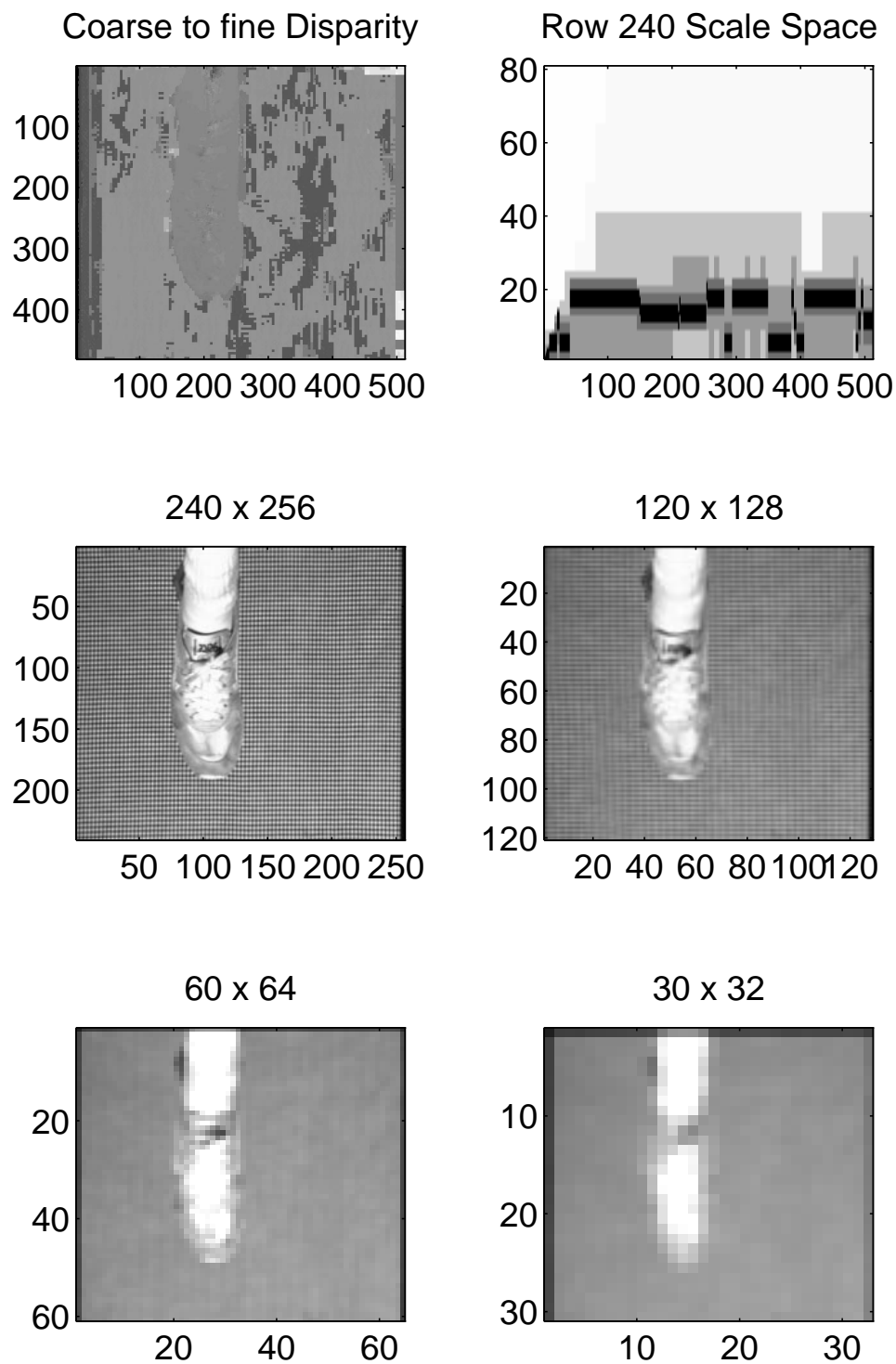


Figure 4.13: Coarse to fine results and images from the Shoe image pair.

to correct for data missing at the coarsest scales.

Scale Selection Thanks to the evaluation function, our method is able to use an *adaptive search through scale space*. To understand how this works, consider the error evaluation function in Equation 3.14. While the heart of the error computation is the difference between the ideal phase difference and that measured in the image, there is also a weighting term, which is the magnitude of the filter output. The effect of that term is to amplify the measurement errors of those filters with high magnitude. In other words, deviations around filters with low magnitude are less important than those with high magnitude, since they contribute less to the error.

Although this is an indirect argument, the point is that the filters with most reliable measurements will contribute the most to the matching error term. Thus the algorithm will choose the proper filters (i.e., scales) for the computation based on the characteristics of the images themselves, not because of an arbitrary requirement imposed by an inflexible search strategy.

Scale Independence Our method not only allows arbitrary scales to be selected, but also considers each scale independently from the others. Because each filter's contribution is merely summed with all the rest, there is no artificial order in which they must be evaluated. This is in contrast to coarse to fine methods, which rely on the coarsest scales to constrain potential mismatches at the finest scales.

High Precision Another advantage of the phase-based approach is the quality of the estimates produced at lower frequencies. Coarse to fine methods typically use lower frequency estimates only to restrict the range of disparity estimates checked at higher frequencies, but our method produces real phase measurements at each scale.

Examples

To understand why our phase-based method outperforms the coarse to fine approaches in instances of ambiguity, we will compare their evaluation functions. And we will use the scalogram to understand the reasons behind the improved shape of the phase-based evaluation functions.

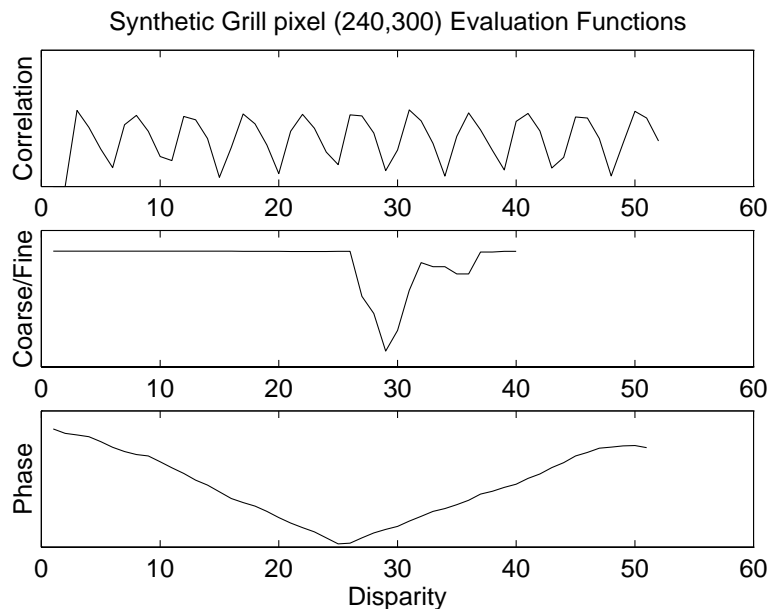


Figure 4.14: Synthetic Grill evaluation functions, illustrating the improvement of the phase-based method over raw correlation (too many minima) and coarse to fine (trapped in local minimum). Actual disparity is 23.52 pixels at pixel (240,300).

Synthetic Speaker Grill

The synthetic speaker grill image provides a nice test case. Its regular structure makes it a prime candidate for ambiguous matches in *any* stereo method, as can be seen in the evaluation functions for raw correlation, coarse to fine, and phase based stereo in Figure 4.14. In the figure, the raw correlation evaluation function profile exhibits ten local minima, all with approximately the same amount of error, resulting in an ambiguity factor of nearly 1. It seems clear that any local method such as this will be prone to ambiguous matches in the grill image pair, but what about a coarse to fine approach? Its evaluation function (also in Figure 4.14) has a better profile, with a unique global minimum and low ambiguity factor, but still has problems. The evaluation function has more than one local minimum, but even worse, the unique global minimum found by the algorithm is *not* the correct disparity. Both of these algorithms were fooled into making a false match by the numerous candidate matches, a problem which their search strategy was unable to model successfully.

Our phase based approach not only finds the correct solution, but also exhibits an evaluation function profile that would be considered ideal by any stereo method: a virtually unimodal function with a single minimum located at the correct disparity. And as the Dis-

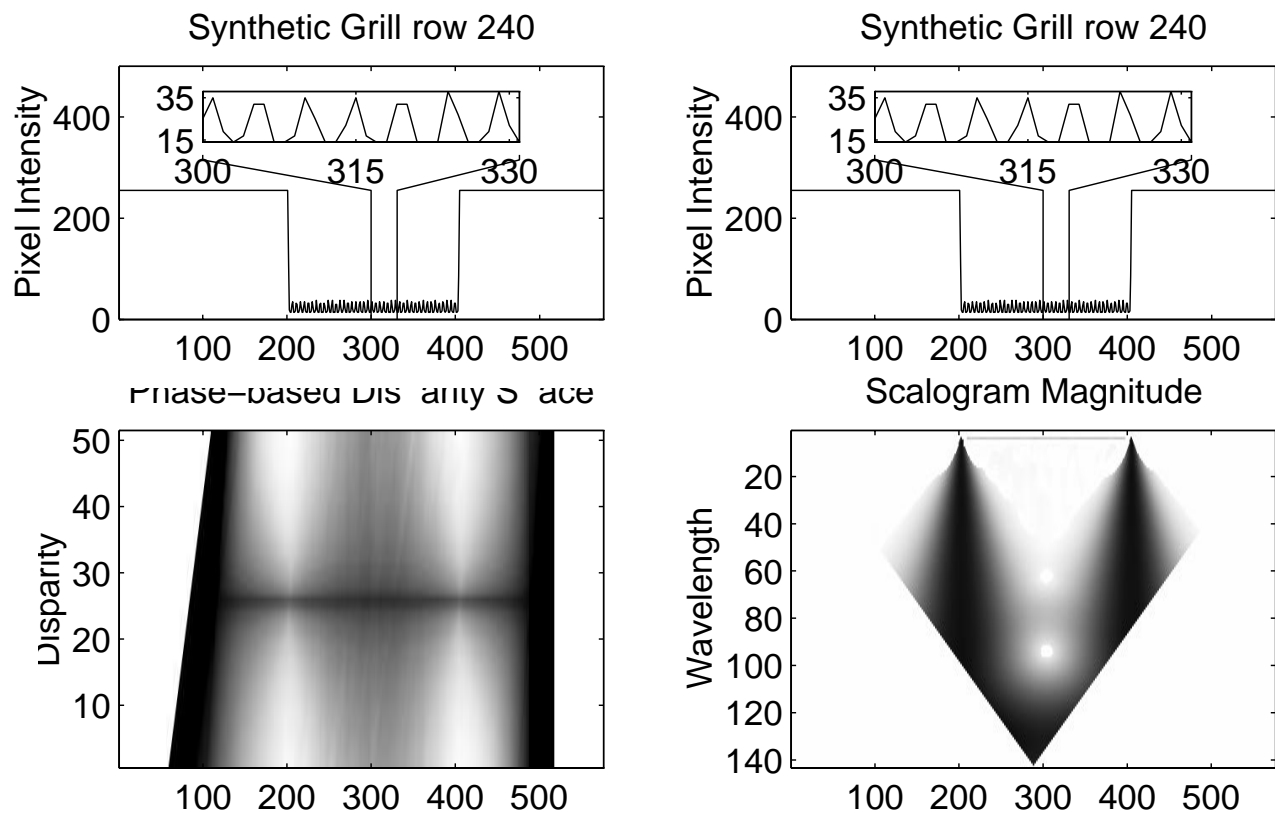


Figure 4.15: Synthetic Grill Phase-based Disparity Space and constraint-filtered Scalogram for line 240. The raw pixel intensities appear above (the same scanline is plotted twice), and illustrate the effect of the grill edges on the plots below.

parity Space in Figure 4.15 shows, *every* pixel in this row has a similarly-shaped evaluation function pointing to the correct disparity. How was our method able to succeed where the others failed?

The explanation can be found in an analysis of the left magnitude scalogram for this row, shown in Figure 4.15. By far the most prominent features are the dark vertical columns that correspond to the edges of the speaker grill. Their presence is to be expected; the Fourier transform of a step edge (e.g., the light/dark transitions in the original image) has energy content at all frequencies. Another expected feature is the horizontal bar at the top of the scalogram. This bar extends across the width of the speaker grill, and its row number corresponds to the period of repetition of the holes in the grill (about 3.8 pixels). The presence of a specific peak in each column like this is exactly what you would expect from an image intensity profile that looks like a sinusoid. But the key to understanding why the evaluation function maintains such a useful profile over the whole row is the dispersion of the dark columns that occurs at the edges of the grill.

Recall that our phase-based stereo method compares columns in the left and right scalograms. This dispersion results in higher magnitudes, and therefore more reliable phase measurements, not only at the edges of the grill but also at low frequencies throughout the central grill region. If only the highest frequencies in that region were considered, then the evaluation functions for the center pixels would look much like the raw correlation profile in Figure 4.14, thanks to the phase wraparound problem mentioned in Section 3.6.1 and illustrated in Section 2.2.1. But in this case the analysis filters with lowest frequencies were wide enough to include the strong influence of the step edge at the borders of the region. So when the columns that represent pixels in the middle of the grill are considered, the analysis also includes nonlocal information about the full extent of the similarly-textured region; the filters in the central columns know about the edges of the grill too. Therefore, our approach to multiscale analysis is able to combine the very high and very low frequency measurements effectively, eliminating the ambiguity faced by other methods.

This analysis also clarifies the reason the coarse to fine method failed to resolve the ambiguity. Even though there is good information at coarser scales (i.e., low frequencies), there is *no* information at medium scales, as the white region in the middle of the scalogram in Figure 4.15 demonstrates. And since the coarse to fine method requires that some decision be made at *each* scale, it is not surprising that decisions made in the absence of useful

information at middle scales would lead to incorrect results.

Shoe

Our method is not foolproof, however. Consider the results of the phase-based method on the Shoe image pair. The results look quite encouraging; the shoe disparities look good, the background color is solid. But the background pattern is a repetitive checkerboard, and the disparity found by our method is *not* the actual disparity; this is shown graphically in Figure 4.16.

In fact this image pair was designed to produce ambiguity, as Figure 4.19 illustrates. The fact that real imagery can exhibit inherent ambiguity like this forces us to redefine our notion of disparity at a pixel.

4.4 Modeling Ambiguity

Sometimes it is simply not possible to reduce the degree of ambiguity present in an image. In such situations it is best to model the ambiguity, rather than attempt to correct the evaluation function by assuming no ambiguity exists. In this section we present a framework for such an analysis, and demonstrate its potential for improving disparity estimates.

4.4.1 Extended Disparity Representation

Our ambiguity model requires us to extend the traditional notion of depth values recovered by a stereo method. The usual model for depth at a pixel is a single disparity estimate, possibly with an indication of the precision, or variance, of the result. Unfortunately, this model does not include any information regarding the accuracy of the estimate. In fact, it presumes an “all or nothing” approach to depth measurement; either the measurement exists and is known to a certain precision, or no depth estimate is known (e.g., the pixel is occluded). Thus the model implicitly assumes that all evaluation functions will have a single unique minimum. As this chapter has demonstrated, e.g., in Figure 4.14, there are times when this simplified model is not sufficient.

The problem is that this model assumes a unique disparity can always be found. While the 3D point imaged at a given pixel will certainly have a unique depth, there may be many plausible disparities under a particular 2D evaluation function. The only way to incorporate

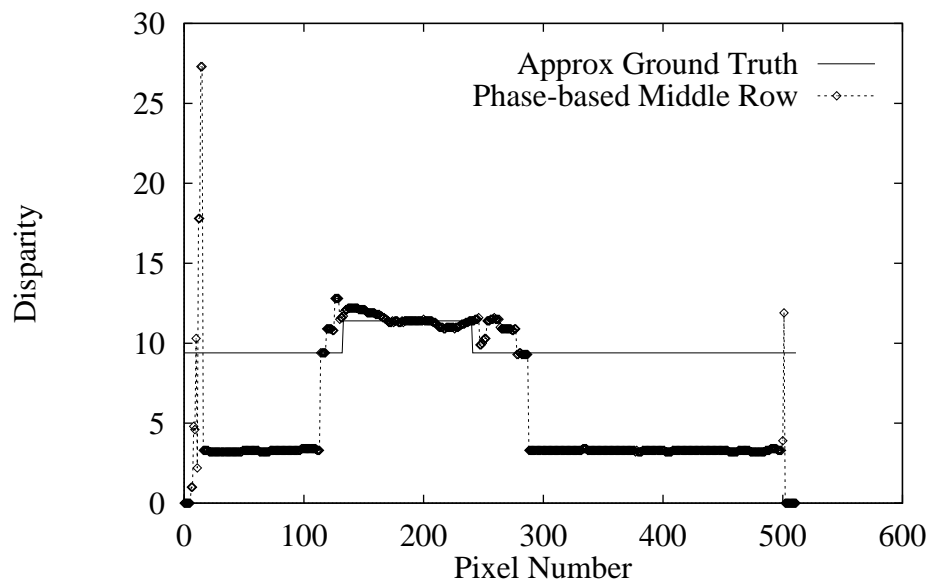


Figure 4.16: Detail view of phase-based disparity results and approximate ground truth for a middle row of the shoe stereo pair. Background disparities are consistent, but incorrect. Figure 4.21 contains the complete disparity space for this scanline.

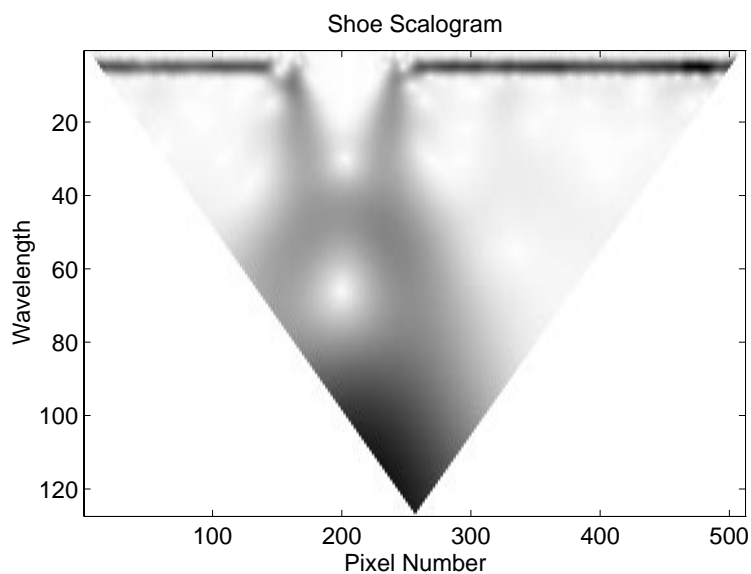


Figure 4.17: Image Scalogram for the middle row of the left shoe image.



Figure 4.18: Phase-based disparity maps for the shoe image pair. Left map is the result from using no constraints, right map is the result from using the heuristic in Section 3.4.3.

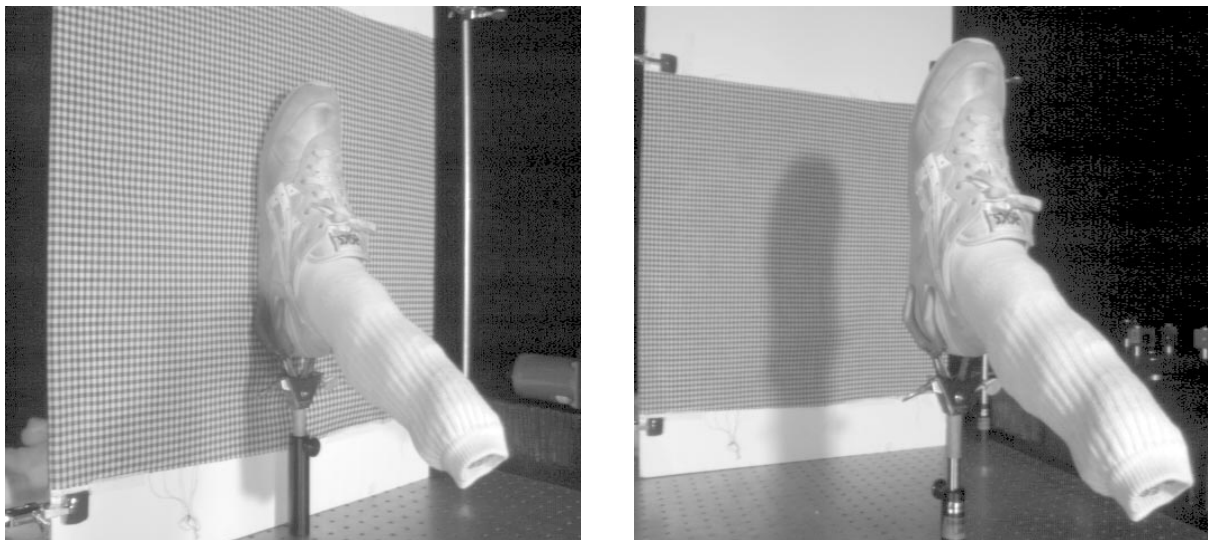


Figure 4.19: A view of the deviousness used in the construction of the Shoe image pair. Although the images were taken of the shoe flat against the texture (left image), the fact that the actual disparity is larger than the period of the checkerboard pattern causes a purely local search to reconstruct the scene with the shoe “floating” above the textured background (right image). Reproduced with permission from Kanade.

multiple candidate disparities under the unique disparity model would be to decrease the precision of the estimate by many pixels so that all candidates are included. That would be a needless waste of information since each estimate is likely to be known rather precisely, often to subpixel precision. A better model would allow a pixel to have *several* disparity estimates, each with its own precision and relative likelihood. More formally,

$$Estimates = (\langle disparity, variance, error \rangle \dots) \quad (4.2)$$

where *disparity* is the best estimate at a local minimum, *variance* is the curvature of that minimum, and *error* is the evaluation function value at that minimum, expressed relative to the other estimates at this pixel using the formula:

$$Error = \frac{error_n}{\max_{i \in PeakIndices} error_i} \quad (4.3)$$

The *PeakIndices* are those disparities whose evaluation functions exhibit local minima. Their automatic extraction is in general very difficult, but we have constructed a heuristic approach that works well enough to demonstrate the principle. By applying the following heuristic to the evaluation function of single pixel, e.g., Figure 4.2, not only the global minimum but *all* useful minima may be extracted. In this way multiple disparity estimates may be generated at a pixel.

Automatic Estimate Extraction

The locations of local minima in a sampled function can be identified using the heuristic peak finder first described in Section 3.4.3. Recall that the method works by locating the global maximum, extracting a window of values around the maximum (up to the point where the sign of the second derivative changes), fitting and subtracting a Gaussian to those points, then iterating on the newly subtracted signal. The intuition behind it comes from the observation that the same process that creates the local maximum will have residual effects nearby, effects that diminish with distance from the peak. This occurs in the frequency domain application of Section 3.4.3 because the frequency response of the Gabor filters has a Gaussian envelope, so nearby overlapping frequencies will indeed see diminishing effects that fall off as a Gaussian. In this new domain of evaluation functions, the assumption is that the evaluation function error will increase monotonically when moving away from the minimum.

A few adjustments do need to be made to accommodate the new domain. Since the original method finds local maxima, and we are now interested in minima, the evaluation function values must first be negated and translated to have a nonnegative minimum value before calling the peak-finder. Also, while the magnitude profile of overlapping filters tends to vary smoothly, in this case the evaluation function profile may have more abrupt transitions, and generally will be subject to more noise. Therefore another threshold is introduced, to enable the method to ignore small perturbations in its input. Whereas the previous method expanded the window around the maximum until the second derivative became greater than zero (i.e., changed sign), for this application we will set the threshold slightly above zero, e.g. 0.5. The effect of this threshold is to allow the support window to grow larger in the presence of noisy data, which should enable a better Gaussian fit. Finally, in order to avoid spurious results from the smaller peaks that result from the subtraction step, we mark all values that contribute to the Gaussian fit as unusable peaks. Should one of those pixels be found to be the maximum, the iteration of the peak-finding procedure will terminate.

Figure 4.20 illustrates this heuristic procedure on two examples. Regions of the function that contributed to the Gaussian fit are marked with dashes in the figure, and each extracted maximum is highlighted with a dark circle.

Output from this heuristic provides exactly the multiple estimates demanded by Equation 4.2. The location of the peak indicates the *disparity*, the curvature of the evaluation function provides the *variance*, and given a complete set of peaks the *error* can be computed using Equation 4.3. In the next section we apply this heuristic to an image pair to illustrate the benefits of the representation.

4.4.2 Demonstrating Improvement on an inherently ambiguous image

Given the locations of the local minima in a pixel's evaluation function, multiple precise disparity estimates may be associated with that pixel. Although multiple disparities can model the results of stereo methods more accurately, in practise it is extremely difficult to put them to good use. The problem lies in the generalization from pixels to surfaces: if all possible combinations of even as few as two estimates were considered, the number of possible surfaces in a single *scanline* would be 2^{512} , i.e., beyond astronomical. We will not propose a general solution to this problem, but rather will demonstrate that useful results

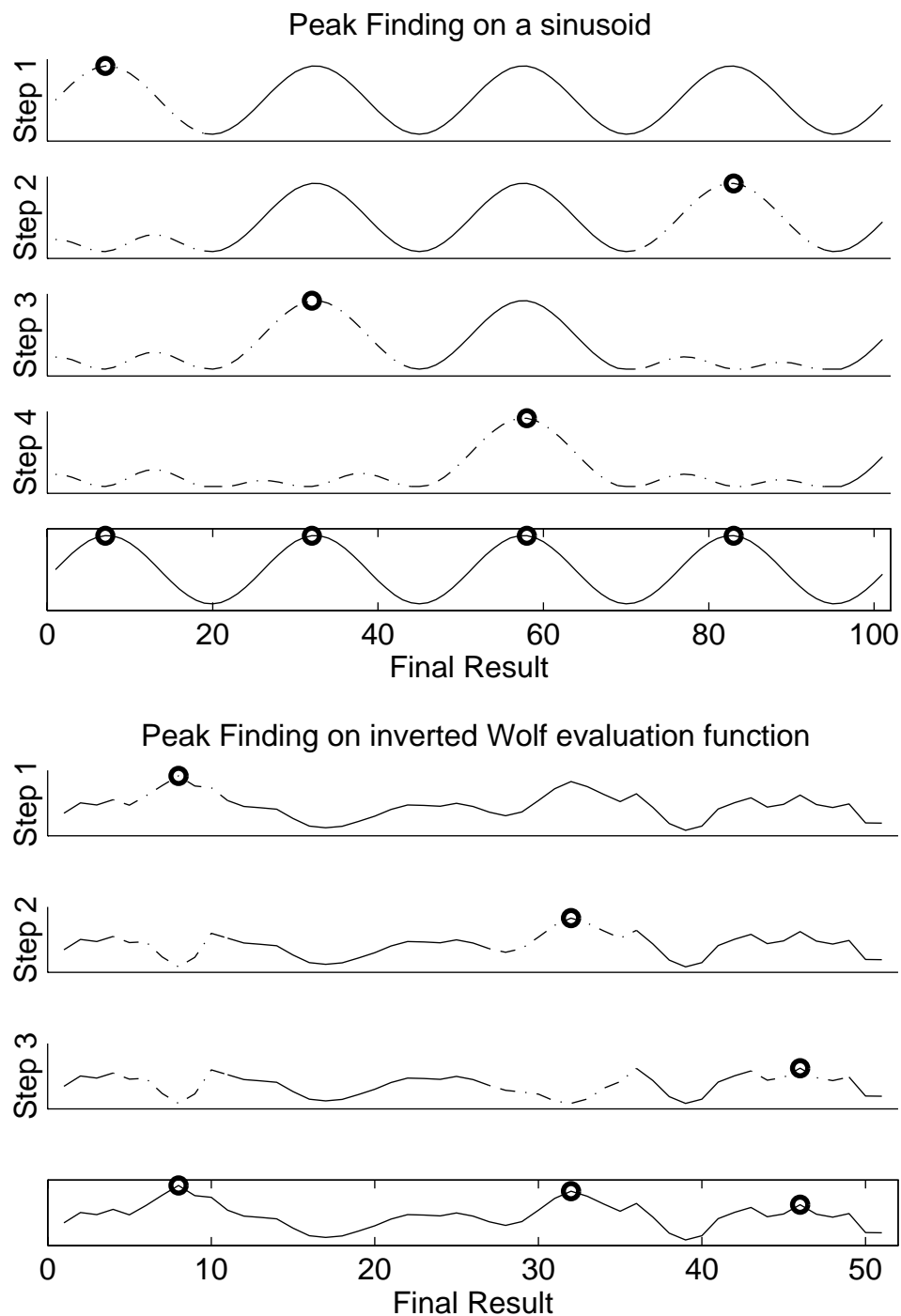


Figure 4.20: Illustration of the peak-finding heuristic from Sections 4.4.1 and 3.4.3 on a sine wave (upper) and the inverted Wolf Evaluation Function from Figure 4.2 (lower). Dashed lines indicate the region in which a Gaussian is fit and then subtracted.

may still be obtained.

We will use the Shoe image pair for this demonstration, since even our phase-based method was fooled by its inherent ambiguity. Figure 4.21 illustrates the phase based disparity space for the central scanline of the image pair, with the scanline plotted above for comparison. In addition, local minima extracted by the heuristic (using a peak threshold of 70% of the maximum) have been identified and plotted directly over the disparity space image in Figure 4.22. Only the peak locations have been plotted, their variance and error components are implicit in the disparity space image. These local minima make any inherent ambiguities plain: when a column exhibits more than one local minimum, the corresponding pixels will have a non-zero ambiguity factor.

The nature of the overall ambiguity becomes clear by inspecting the connectivity between adjacent columns. All columns in the area representing the shoe (pixels 150–250) have a single minimum at the correct disparity, and are thus clearly unambiguous. The area of the checkerboard pattern has lots of peaks, though; individual columns have as few as two or as many as nine potential matches within the given disparity range (0–50 pixels). The majority of those peaks fall into disparities about 6 pixels apart: this corresponds nicely to the period of the checkerboard in the original image, which agrees with our intuition that the self-similar checkerboard texture makes pixels in the background prime candidates for ambiguous matches.

The utility of this representation becomes clear when you contrast this combined minima plot (Figure 4.22) with the disparity results of the phase based method (Figure 4.16) and of SRI’s method (Figures 4.7 and 4.23). From manual inspection of the original image pair in Figure 4.5 (especially the black stripe on the right hand side that indicates the end of the texture), we know that most of the actual image disparities lie around 10 ± 2 pixels. The incorrect results from both methods are easily explained as alternative local minima in the disparity space that were chosen because the evaluation function rated them more favorably. By plotting all of the local minima, we are more likely to find the actual disparity along with other potential match points.

Once all the candidate matches have been found, we can accurately describe the match quality over the entire scanline. Columns with a single match can be interpreted as before, where a single disparity value represents the most accurate reconstruction possible. In areas with multiple candidates, techniques other than finding the raw evaluation function

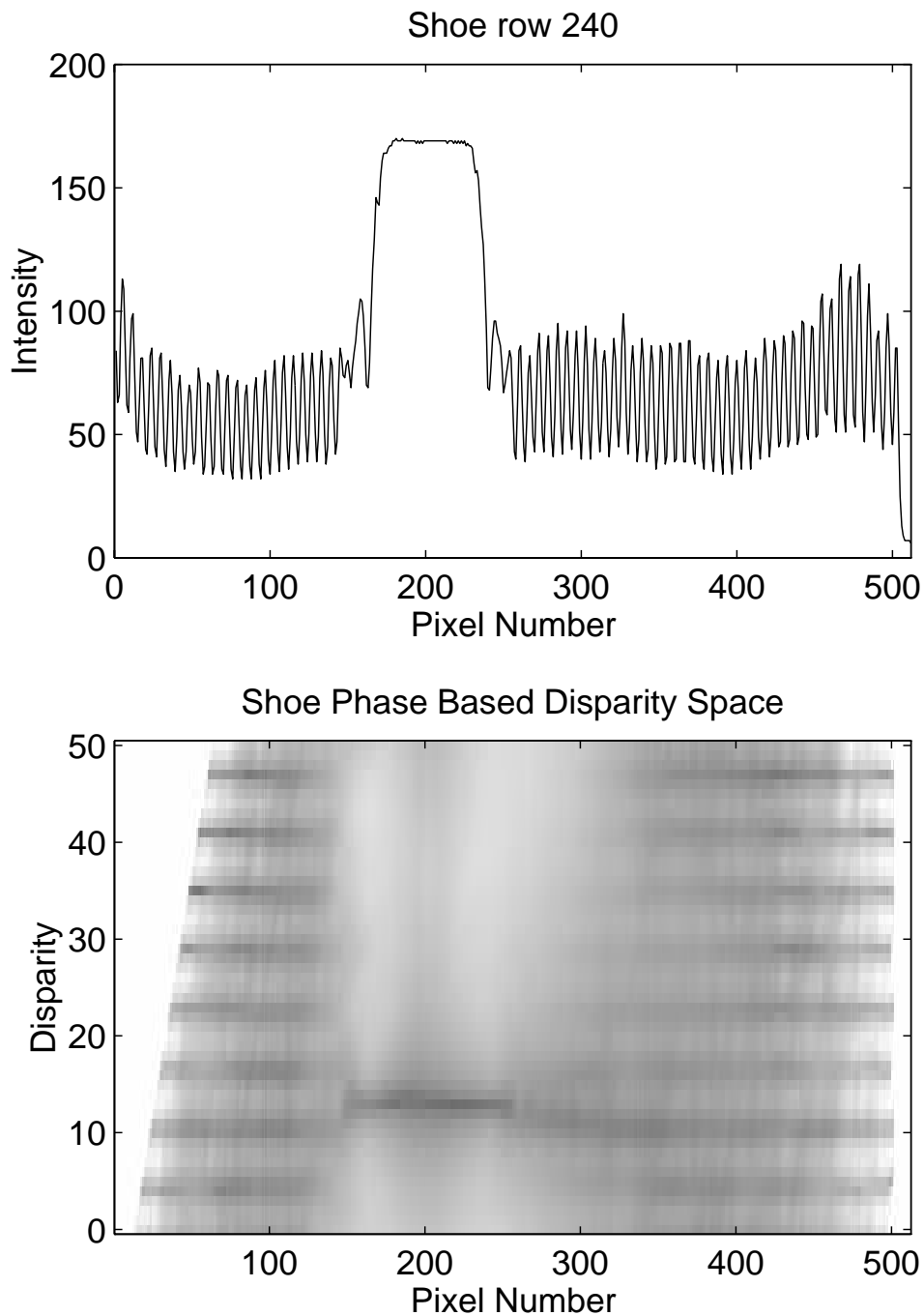


Figure 4.21: Shoe Disparity Space. The upper plot is row 240 from the Shoe image pair, the lower plot is the disparity space computed by our phase-based method.

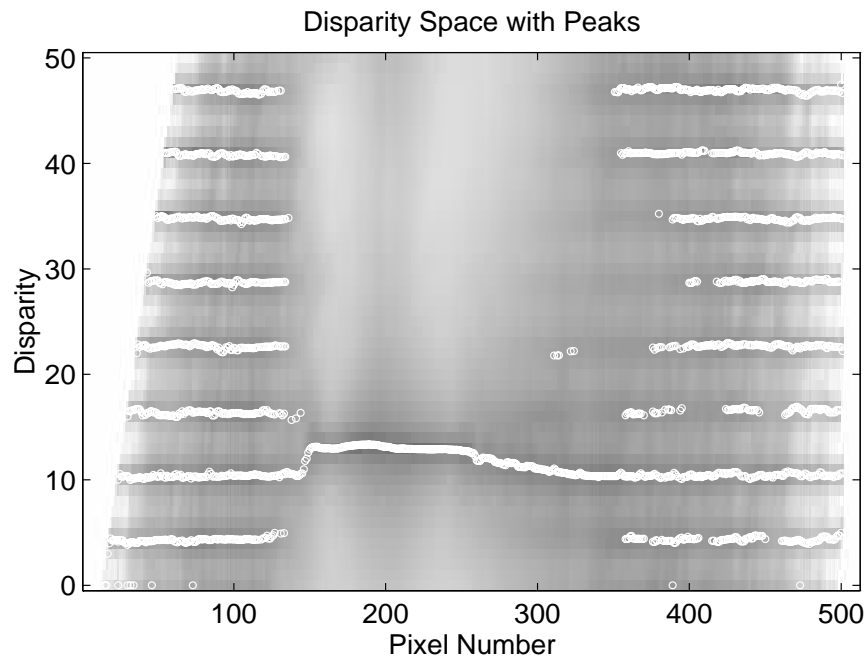


Figure 4.22: Phase-based Disparity Space with Peaks. Peaks computed using the heuristic have been superimposed on the disparity space from Figure 4.21.

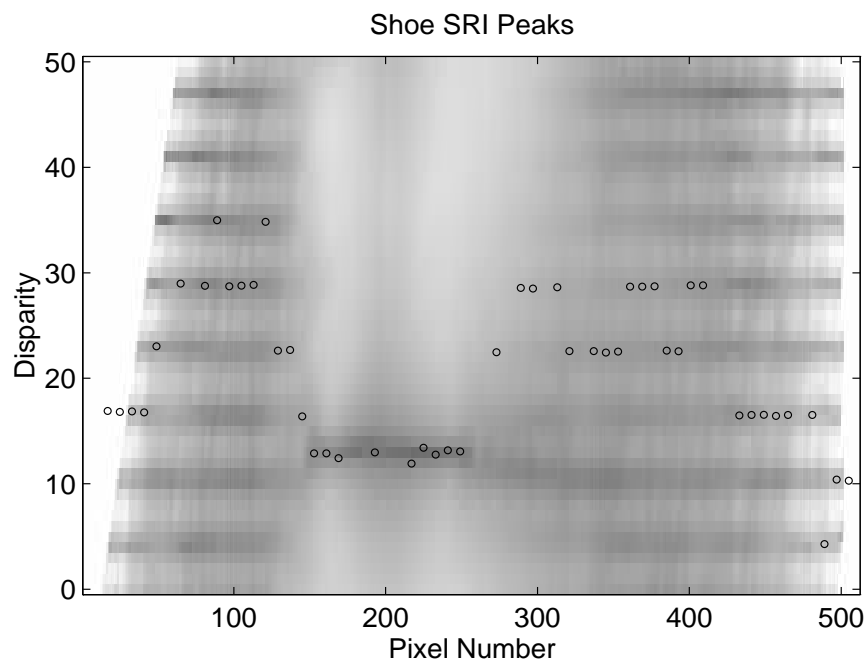


Figure 4.23: SRI Disparity over Phase-based Disparity Space. The seemingly random results from Figure 4.7 actually fit nicely into the local minima of the phase-based disparity space (the SRI disparity space was not available).

minimum can help eliminate false targets. For example, morphological operations on a binary “peaks-only” image could be used to group similar disparity estimates together and eliminate outliers. In the shoe example, such a procedure might find that the background region has eight possible interpretations that can be grouped into three categories: a plane that lies much closer to the camera than the shoe, nearly the same distance as the shoe, or much further than the shoe. Such a procedure might average the *error* terms of the candidate matches to present those multiple interpretations in a particular order. By grouping the multiple candidate matches together at adjacent pixels, we reduce the complexity of the model from an astronomical number of independent pixel disparities down to a more manageable eight potential surfaces.

A similar approach has already been demonstrated in (Zitnick & Webb, 1995). The core of their stereo algorithm groups match candidates into surfaces, then chooses the surface with the greatest number of points as the correct one. However, their method currently requires that a pattern be projected into the scene, so it is not a passive stereo method. Further details of their method are beyond the scope of this work, but the interesting point is that nearly all prior work on stereo has focused on improving the shape of the evaluation function, while Zitnick and Webb eliminate the evaluation function entirely, relying instead on pixel to pixel matches (i.e., a 1x1 pixel evaluation window). This chapter has presented a compromise; improve the evaluation function, *and* work with the extracted potential mismatches.

4.5 Summary

The problem of ambiguous matches, or false targets, can greatly reduce the accuracy of a stereo vision system. We have shown that the usual approach to alleviating the problem, a coarse to fine refinement strategy, imposes some (perhaps overly) strong requirements on the stereo images. Our phase-based method relaxes those requirements, and is therefore able to handle a wider variety of otherwise ambiguous images. We have also proposed a generalized disparity model that explicitly represents multiple candidates. This model allows higher level functions to reason more accurately about the structure of the image.

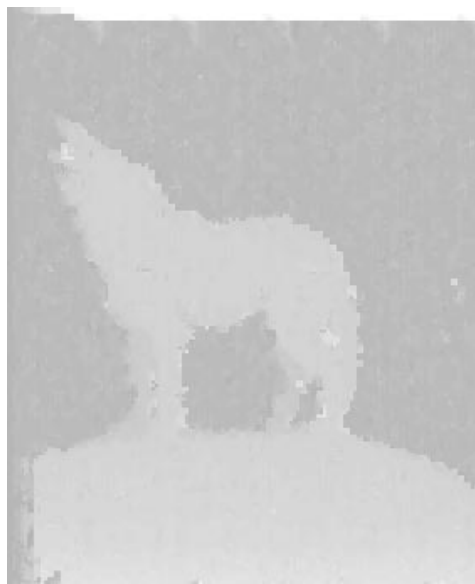


Figure 4.24: The structure embedded within Figure 4.1. This image was computed by applying our coarse-to-fine stereo method with 5 pixel window and 5 pixel max disparity per level to the original 340x340 image and a copy of the original shifted by 58 pixels.

Chapter 5

Precision: Effects of Foreshortening

Changes in latitude, changes in attitude.

— Jimmy Buffett

In this chapter we show how perspective foreshortening is manifest in the local spatial frequency representation of stereo images. Ours will be a forward-reasoning analysis, beginning with complete knowledge of the three-dimensional geometry of the scene and ending with its two-dimensional projection in the image plane. The primary result is the presentation (in Equation 5.10) of the Foreshortening Factor that allows us to compensate for arbitrary foreshortening effects without explicitly warping the images. This result makes no restrictions on the surface texture, and will not require the use of disparity derivatives. The complementary technique (starting with the projections to determine three-dimensional geometry) will be presented in Section 5.4.

To simplify the analysis, we assume the only object in the world is a textured flat plate that is either parallel to the image plane, or rotated about the vertical axis by some angle θ . We further assume that the stereo cameras have parallel optical (depth) and vertical (height) axes. Note that we can restrict our attention to the effects of foreshortening in one-dimensional image *scanlines*, rather than complete two-dimensional images, since all disparities will be horizontal under this assumption. Our world model will likewise be a two-dimensional slice through the three-dimensional scene. Figure 5.1 shows an overhead schematic of a horizontal slice through the world. We adopt the convention that parameters measuring distances in the world will be capitalized (e.g., X_S , Z_L), and those measuring

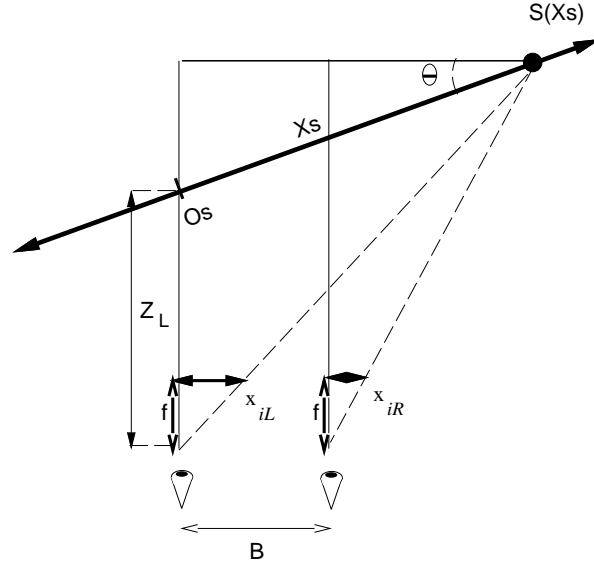


Figure 5.1: Overhead view of the foreshortening model. X_S is the distance from the point exactly in front of the left camera (the origin O_S at distance Z_L) to the point (S) on the plate being studied; x_{iL} and x_{iR} are the left and right pixel indices of the image of surface point S ; the cameras are separated by baseline B and the surface tilts away from the cameras at angle θ .

pixel or camera distances will be lower case (e.g., x_{iL} , f).

Figure 5.2 illustrates the effect of this foreshortening in the frequency domain. To simplify the demonstration, a flat plate that has a surface texture with a single frequency component is used: a sine wave. The figure has two images of the plate on top, and the corresponding scalogram magnitude plots below. The head-on view of the plate on the left side of the figure has the expected scalogram; a horizontal line centered at the frequency of the sine wave, with some extra energy (dark regions) at the edges of the plate where it borders the plain white background. The rotated view also has a straight line in the scalogram, but it appears at higher frequencies and is no longer horizontal. This transformation will be quantified precisely in the closed-form *foreshortening factor* developed later in this chapter.

Although our ultimate goal is to find the disparity between two stereo images, we must first determine how the appearance of the object's surface texture will differ between them. Specifically, we want to know how the sampling rate varies between the two images. This is a geometric formulation; what matters is how much of the surface is being mapped to each pixel, not the actual surface texture (i.e., color intensity). So for each location X_S on the

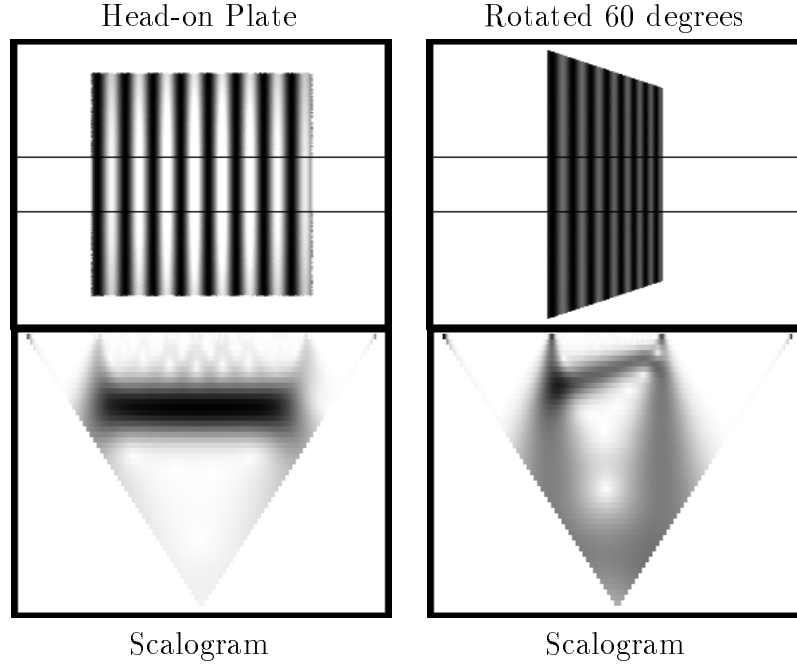


Figure 5.2: The effect of foreshortening on scalogram magnitude. Two views of a flat plate with a sinusoidal texture appear on top, and the scalogram magnitudes for their central scanlines appear below. The responses are similar, but are compressed to higher frequencies in the rotated view.

surface, we want to compare the pixel areas in the left and right images. Mathematically, we want to compare the left sampling rate $\frac{\delta X_S}{\delta x_{iL}}$ to the right sampling rate $\frac{\delta X_S}{\delta x_{iR}}$:

$$\text{Sampling ratio} = \frac{\frac{\delta X_S}{\delta x_{iL}}}{\frac{\delta X_S}{\delta x_{iR}}} = \frac{\delta x_{iR}}{\delta x_{iL}} \quad (5.1)$$

Simplifying the ratio in this way proves most useful. The resulting formula tells us we can compute the sampling ratio (which will be called *foreshortening factor* later) in *image space*, without having to explicitly model the distance X_S along the object. Unfortunately, it also implies that we need the disparity derivative (recall $\delta \text{disparity} / \delta x_{iL}$ is simply $\delta(x_{iL} - x_{iR}) / \delta x_{iL} = 1 - \text{Sampling Ratio}$). Since our ultimate goal is to estimate disparity, it would be best if we could avoid using both disparity and its derivative in our calculations (the derivative of a noisy signal will be even noisier). The remainder of this section will show how we can express this ratio with terms that do *not* require disparity derivatives.

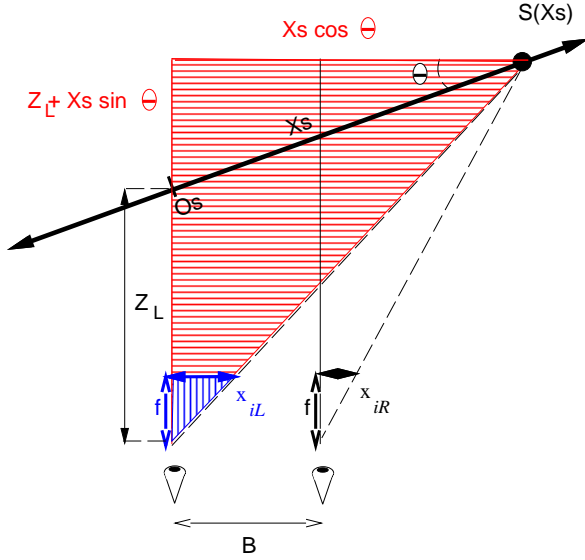


Figure 5.3: Overhead view of the foreshortening model. Similar triangles for the left camera geometry are highlighted (see Equation 5.3).

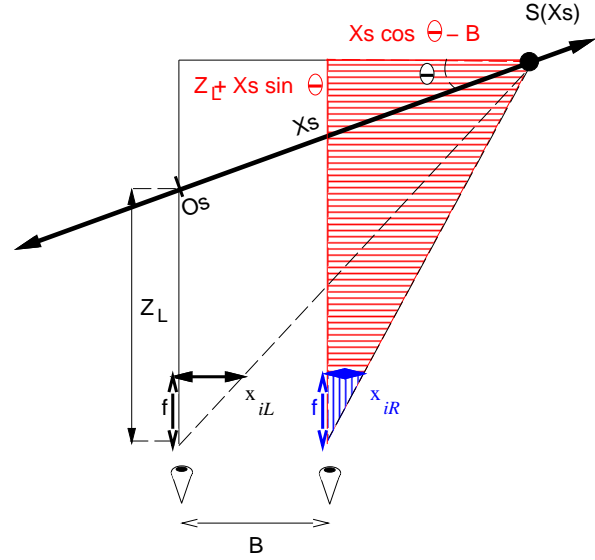


Figure 5.4: Overhead view of the foreshortening model. Similar triangles for the right camera geometry are highlighted (see Equation 5.4).

5.1 Relating Disparity to Surface Angle

How is disparity related to the surface angle? Equation 1.2 gives the disparity for an individual point, but we will now show how it varies across a surface. We will focus our attention on the distance from the left camera to the surface point immediately in front of it, expressing other depths in terms of this value Z_L .

Recall that disparity is the difference of the left and right pixel indices. So let's see how each of the left and right indices (x_{iL} and x_{iR}) relates to the surface angle θ . A quick look at Figure 5.1 shows us the general answer using similar triangles:

$$\frac{\text{pixel index}}{\text{focal length}} = \frac{X \text{ World Coordinate}}{Z \text{ World Coordinate}} \quad (5.2)$$

Figures 5.3 and 5.4 highlight the similar triangles for the left and right scene geometries. Applying Equation 5.2 to those figures we obtain expressions for x_{iL} and x_{iR} :

$$\frac{x_{iL}}{f} = \frac{X_S \cos \theta}{Z_L + X_S \sin \theta} \quad (5.3)$$

$$\frac{x_{iR}}{f} = \frac{X_S \cos \theta - B}{Z_L + X_S \sin \theta} \quad (5.4)$$

Equations 5.3 and 5.4 give us expressions for x_{iL} and x_{iR} in terms of the focal length f , baseline B , distance in front of the left camera Z_L , surface angle θ , and location on the surface X_S . These equations represent projections of the same surface point X_S into two image planes, and we can find the relationship between them by solving Equations 5.3 and 5.4 for X_S and setting them equal.

$$\frac{x_{iL}Z_L}{f \cos \theta - x_{iL} \sin \theta} = \frac{x_{iR}Z_L + Bf}{f \cos \theta - x_{iR} \sin \theta} \quad (5.5)$$

Solving Equation 5.5 for the right pixel index gives us:

$$x_{iR} = x_{iL} \left(1 + \frac{B}{Z_L} \tan \theta \right) - \frac{Bf}{Z_L} \quad (5.6)$$

And finally, recalling that disparity is the difference of the two indices:

$$disparity = x_{iL} - x_{iR} = \frac{Bf}{Z_L} - x_{iL} \frac{B}{Z_L} \tan \theta \quad (5.7)$$

Equation 5.7 is nearly the answer we want. It relates disparity to the scene parameters, and does not depend on knowing the actual surface location. It does require knowledge of Z_L (distance to the surface point in front of the left camera), unfortunately, but we will eliminate this restriction below.

Equation 5.7 has some interesting interpretations. When the surface is frontoplanar (i.e., $\theta = 0$ and thus $\tan \theta = 0$) it reduces to the familiar expression relating disparity to depth from Equation 1.2; this is correct since all surface points would lie at the same depth Z_L . And for an arbitrary fixed angle θ the disparity *derivative* is constant, i.e., the disparity varies linearly with respect to the image location x_{iL} . While we won't take advantage of this property of the derivative, it could prove useful to shape-recovery techniques.

5.2 Expressing the Foreshortening Factor using Image Parameters

Now that we know how the disparity and pixel locations relate to surface angle, let us return to the Foreshortening Factor (Equation 5.1) and eliminate the derivative by substituting for x_{iR} :

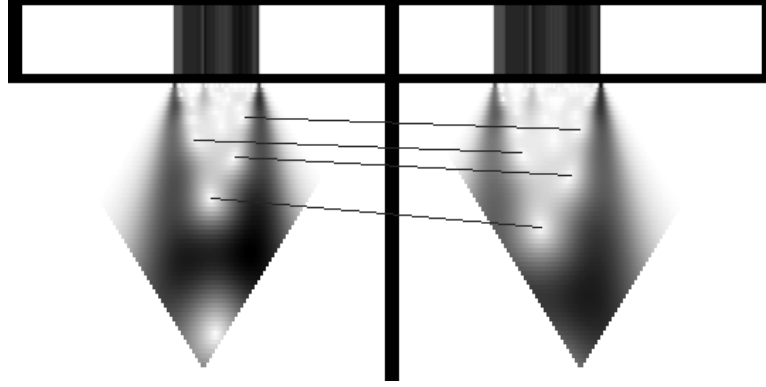


Figure 5.5: Left and right views of a surface tilted 65 degrees. Upper images are the central scanlines, lower images are their corresponding scalograms. You can see similar features in both scalograms: those in the left image are present at higher spatial frequencies because the left image is subject to greater foreshortening effects than the right image.

$$\begin{aligned} \text{Foreshortening Factor} &= \frac{\delta x_{iR}}{\delta x_{iL}} = \frac{\delta \left(x_{iL} \left(1 + \frac{B}{Z_L} \tan \theta \right) - \frac{Bf}{Z_L} \right)}{\delta x_{iL}} && \text{from Equation 5.6} \\ \text{Geometric Form} &= 1 + \frac{B}{Z_L} \tan \theta && (5.8) \end{aligned}$$

This expression is very interesting. It tells us that for a given flat surface, the Foreshortening Factor is *constant* over both images of the surface. In other words, the local spatial frequencies of the left and right images are related by a simple constant scale factor. You can get a feel for this by visually tracking the low magnitude phase singularities (white spots) between the two image scalograms in Figure 5.5.

The fact that foreshortening causes frequency shifts has been noted in the literature (Fleet & Jepson, 1993), but no explicit model was given to explain it in the context of stereo vision (but see (Krumm, 1993; Malik & Perona, 1989) for single image texture-based models). Instead, the instantaneous frequency was recovered using a heuristic averaging technique. This technique yielded somewhat better results than the use of direct frequency, but did not take advantage of the scene geometry to compute the precise shift. This averaging technique also failed whenever the frequency shift caused the instantaneous frequency to fall outside the range of the filter in either of the images. Our model overcomes these problems by making use of all available frequency bands, rather than limiting attention to a small number.

The result in Equation 5.8 is useful for describing the form of the foreshortening effect

(that of a constant scale factor), but it would be useless in a stereo matcher since it requires knowledge of the depth Z_L . A program that computed depth given depth would not be very impressive. So how can we eliminate the need to know Z_L ? Consider the ratio $\frac{B}{Z_L}$. We can rewrite Equation 5.7 as:

$$\frac{B}{Z_L} = \frac{\text{disparity}}{f - x_{iL} \tan \theta} \quad (5.9)$$

and replace that in Equation 5.8, giving us this final expression for the projected form of the Foreshortening Factor:

$$\text{Projected form} = 1 + \frac{\text{disparity} \tan \theta}{f - x_{iL} \tan \theta} \quad (5.10)$$

This is what we want! Equation 5.10 relates parameters in the image plane to the surface slope θ , but does not require prior knowledge of the distance to the object or an estimate of the disparity derivative. It does require use of some known parameters (focal length f , image location x_{iL}) and variables being estimated (disparity, surface angle θ), but we will see how to manage these algorithmically in Section 5.4.

In this section we described the effect of perspective foreshortening in terms of local spatial frequency. We developed this theory in steps to demonstrate several properties: the frequency shift (aka *foreshortening factor*) between images of an oriented flat surface is constant, it is independent of the surface texture, and it can be expressed using only disparity and surface angle (without disparity derivatives). Section 5.4 will show how these results can be applied to a stereo matching system.

5.2.1 Verifying the Foreshortening Factor

Before continuing, we will verify the geometric form of this Foreshortening Factor using a simple example: a flat surface with a sinusoidal texture. If the model is correct, the surface's apparent spatial frequencies will be shifted between the two images by the amount given in Equation 5.8. Note that we're not solving the stereo problem yet, in fact this demonstration will use the *known* disparity to compare the left and right *image* frequencies at the same *surface* locations. What this will show is that Equation 5.8 accurately predicts the frequency

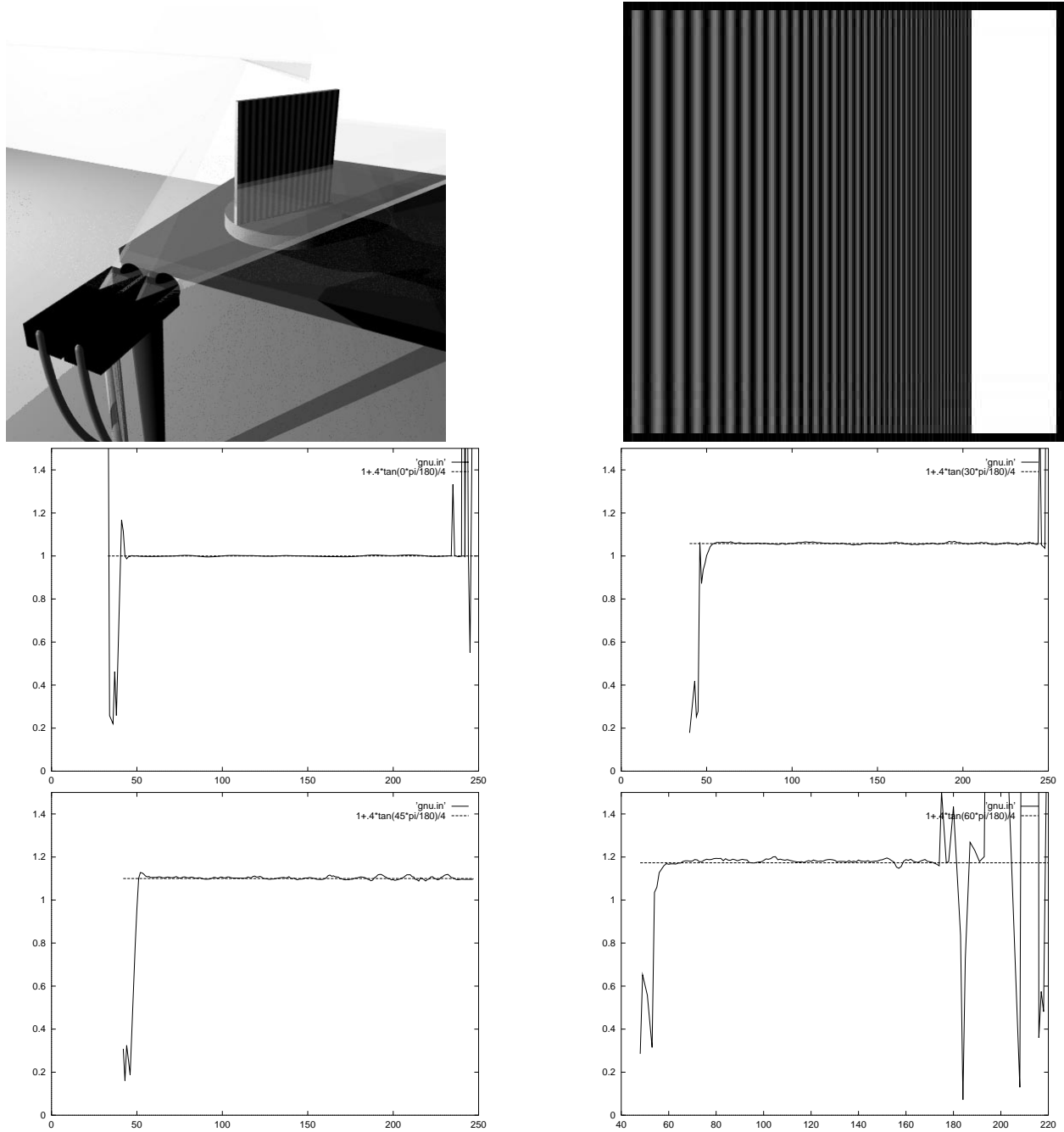


Figure 5.6: Verifying the Foreshortening Factor - These graphs compare the predicted foreshortening factor (dashed line) against that computed using only image information (solid line). The virtual lab setup (top left) and an example input image with surface angle of 60° (top right) are shown first. Next we have the results derived from a surface angled at 0° (middle left), 30° (middle right), 45° (bottom left), and 60° (bottom right). The virtual surface is 4.0 units from the left camera, both cameras have a field of view of 45° and are separated by a baseline of 0.4 (the surface in the actual images is larger than that shown in the top left rendering).

shift of a simple signal. We will use synthetic data so that our ground truth can be as precise as possible.

Recall the geometric form of the Foreshortening Factor from Equation 5.8:

$$\text{Foreshortening Factor} = 1 + \frac{B}{Z_L} \tan \theta$$

Just what is this foreshortening factor? It describes the relationship between the spatial frequencies at two image pixels representing the *same* surface point. How can we measure such frequencies, and how do we know they correspond to the same surface point?

Finding the frequency is easy, but imprecise; we will use an artificial surface texture that contains a single peak in the positive frequency domain, i.e., a sine wave. Its apparent frequency can be found simply by locating the filter output with highest magnitude.¹ As a further refinement, we will use the instantaneous frequency (phase derivative) of that filter output as our frequency estimate. Under the scalogram representation this corresponds to picking the maximum magnitude value in each column.

The procedure for finding corresponding points is somewhat complex, but simply stated involves using knowledge of the ground truth to give the disparity at each pixel (disparity is inversely related to depth, which is known from the 3D model). Remember, we are not trying to solve the stereo problem at this point, we are simply trying to verify a property of corresponding pixels.

Having established the correspondence in the 2D images, we extract the apparent frequency at each pixel using the method described above, linearly interpolating the instantaneous frequency measurements from the right image. Finally, we graph the ratio of the computed image frequencies values against the predicted ratio in Figure 5.6, for several surface angles. The computed ratio is quite accurate but gets progressively less precise as the angle increases. The loss of precision occurs from several factors, e.g., our use of simple linear interpolation to compute the frequencies, and our filter set which only samples the highest frequencies very sparsely.

¹In practise our windowing scheme provides only high frequency info at the image borders, so our computed Foreshortening Factor will become inaccurate at the ends of the graph since the actual spatial frequency is lower than the lowest measured by the filters at that pixel.

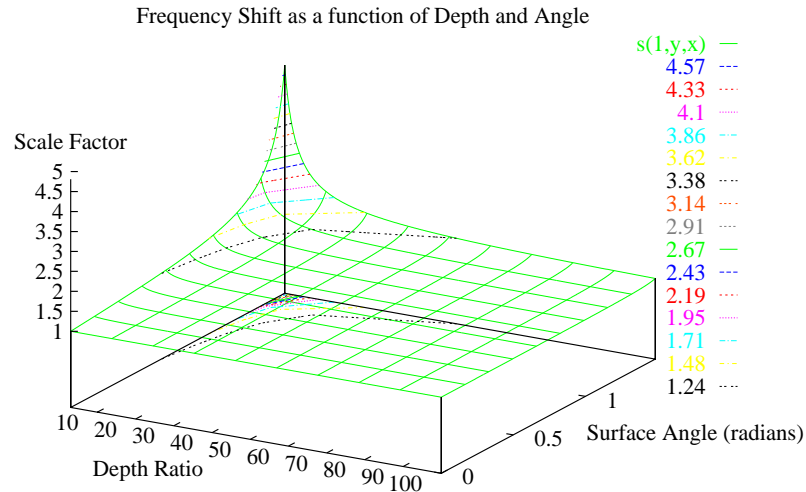


Figure 5.7: Foreshortening Factor as a function of Depth and Angle. Depth is unitless relative to the baseline, and varies from 3 to 100. Angle varies from zero to 85° .

5.3 Applicability

How important is this foreshortening analysis? More specifically, how often do situations arise in which the assumption that a surface is frontoplanar can cause problems for stereo systems? Intuitively the analysis would seem to be needed any time a surface is slanted at a sharp angle; but what if the surface is so far away the slant can't be measured? One might also think it only necessary for surfaces at the sharpest angles; but close up images can exaggerate even small angles. We will use the Foreshortening Factor to quantify these effects in the spatial domain.

Since we want to consider the scenery being imaged rather than the images themselves, we will use the geometric formulation of the Foreshortening Factor from Equation 5.8. Although this expression is a function of three variables, we can reduce it to two if we consider the ratio of depth over baseline $\frac{Z}{B}$ to be a single variable. In the rest of this chapter the word *depth* will denote this unitless version of depth, expressed relative to the camera baseline. For example, the distance between a person's eyes would be 1, the distance to their computer monitor 4-6, and the distance to the far wall in a typically small three-person graduate student office about 100. Figure 5.7 plots the near-complete Foreshortening Factor space for a person looking at objects in such an office.

Depth Range	P($\geq 10\%$ effect)	Example Domain
0-100	0.210355	<i>Human in office</i>
5-20	0.354404	<i>Robot Vehicle</i>
30-100	0.0808227	<i>Inspection Robot</i>

Table 5.1: Probability that a surface exhibits $\geq 10\%$ variation between images due to perspective foreshortening. The distribution of surfaces is assumed to be uniform within the range of orientation angles from $-\frac{\pi}{2}$ to $\frac{\pi}{2}$, and depth ratios (distance divided by baseline) are as specified. A sample derivation can be found in Appendix B.

Figure 5.7 shows the Foreshortening Factor computed from many combinations of depth and orientation (except for the extreme values near the point at which it approaches infinity). The graph makes it clear that the Foreshortening Factor has its greatest impact when objects are sharply slanted and/or located near the cameras. We can quantify its influence using the contour lines that separate regions of large and smaller foreshortening effects. Suppose we assume that surface depth and orientation are uniformly distributed throughout a scene. Then we can compute the probability that a surface will require at least a 10% correction term by finding the area under the 1.1 Foreshortening Factor contour curve. The derivation follows in Appendix B, but the result is that given a uniform distribution of angles from 0° to 90° and depths from 0 to 100, the probability that a surface will require at least a 10% correction is 0.210355. Try it out; if you're sitting in an office, see if you can find one sharply foreshortened surface for each set of four nearly head-on surfaces in your immediate vicinity.

Of course the probability of finding foreshortened surfaces depends very much on the domain being studied. Robot vehicles like Carnegie Mellon's NAVLAB often use a very wide baseline, on the order of one meter. With the nearest visible ground point being about five meters away, depth ratios of 5 to 20 are common in this domain. In that range, under the same assumptions of uniform distribution, the probability of finding a foreshortened surface jumps to better than one in three (see Table 5.1). Inspection robots typically use much smaller baselines, with corresponding depth ratios from 30 to 100. Even in that range, the probability of finding a 10% foreshortened surface is significant (nearly one in twelve). These results suggest that a wide variety of stereo vision systems could benefit from an analysis that considers the effects of foreshortening.

5.4 Application

The analysis in Section 5 is not only theoretically interesting, it can also improve the performance of real stereo algorithms. Phase-based methods such as (Fleet et al., 1991; Sanger, 1988; Weng, 1993) as well as our method can benefit from this analysis. In this section we explain how to apply the Frequency Shift to these phase-based stereo matching algorithms and demonstrate how its application to our system increased the maximum matchable surface angle from 30 degrees to over 75 degrees.

5.4.1 Extending Phase-based Stereo Algorithms

Some have argued that a small number of Gabor filters are sufficient for stereo matching. (Fleet & Jepson, 1993; Weng, 1993) The idea is that although the phase may vary slightly across nearby frequencies, the amount of variation is small enough that the error introduced in measuring it at what might be the wrong frequency will be insignificant. But the assumption is made that the same filters can be applied to both images, i.e., that both images can be sparsely sampled at the same set of spatial frequencies. As was shown in the preceding section, that assumption is not true when perspective foreshortening occurs in the images; indeed, we have seen that frequency shifts of even 10% can occur often. Instead of introducing error by sampling at the wrong frequency, we would like to turn these perturbations to our advantage by using them to confirm hypotheses of surface slant.

We will need a dense sampling of the phase space to get the most accurate results. We will also interpolate phase values between adjacent frequencies when possible. The image scalogram provides a useful framework for such computations, and will be used as the basis for our foreshortening-corrected stereo algorithm.

The method outlined in Section 3.7 uses a global minimization strategy to find the best disparity from a list of candidates. This framework makes it easy to include a foreshortening correction term: in addition to searching disparity space, we also search over surface angle. Recall from Figure 5.2 that a foreshortened object will generate a response at *different* frequencies (i.e., scalogram rows) in the two images. Searching over surface angle allows us to predict the corresponding frequency directly. Pseudocode for this revised algorithm is given in Table 5.2. The only difference between this and the original algorithm is the presence of the correction term on the right image phase measurements. This simple presentation of the

Given: A pair of greyscale images, lists of potential disparities and surface angles, focal length f .

For each row

 Compute Left and Right Scalograms L and R

 For each column c

 For each disparity d

 For each angle a

$$correction = 1 + \frac{d \tan a}{c \tan a - f}$$

$$error = \sum_{\lambda: \rho(\lambda) > threshold} \rho_L(c, \lambda) \cdot$$

$$|\Delta\phi_{ideal}(d, \lambda) - (\phi_L(c, \lambda) - \phi_R(c + d, \lambda \cdot correction))|_{2\pi}$$

 Return d (and a) that yield minimum *error*

Table 5.2: Pseudocode for the foreshortening-corrected algorithm. Column index c must be zero in the center of the image.

algorithm is only made feasible because of the large number of filters used in the scalogram. The large filter set gives us a dense set of phases at many scales from which to compute the appropriate subsampled phases.

There are several implementation details that arise from this simple correction factor. It depends on three variables: the currently hypothesized disparity, surface angle, and the current location within the image. Because these values vary at each pixel on the image scanline, it must be recomputed for each hypothesis. And as was mentioned above, the corrected frequency will almost certainly not be one of those already present in the scalogram; some method of interpolation will be required. These are not serious problems, but imply that their implementation will be very compute-intensive.

5.4.2 Results

We added the correction term to the algorithm presented in Section 3.7 using linear interpolation between adjacent phases. In this section we present the results of our method on real images that have been synthetically mapped onto planar surfaces. The use of synthetic

data allows us to quantify its precision using perfect knowledge of the ground truth.

Consider the stereo pair in Figure 1.3. It shows a synthetic stereo image pair of a flat plate rotated 65 degrees from the image planes, with the image of a city scene texture-mapped onto the plate. The actual disparity map (known from the 3D world model) and differences between the ground truth and disparity computed by three stereo methods are presented in Figure 5.8. The figure shows disparity maps rendered as perspective surfaces; only the area known to have texture is shown since the plain white background makes depth recovery impossible in those areas.

For this demonstration of the foreshortening-corrected algorithm, a set of 501 potential disparities were considered (0 to 50 in steps of 0.1), and the angle was fixed at 65 degrees. The RMS error of this result was 0.38 pixels over the entire plate, with $\sigma = 0.63$. The bulk of this error can be attributed to two causes: the dark spots and a subtle systematic error over the surface. The spots most likely arise from an artifact of the rendering process which caused a few nearby pixels in one image to map to the same intensity. The more subtle effect is that the disparity error, while within measurement bounds at the ends and center of the plate, varies by as much as 0.5 pixels between the center and end of the plate (see Figure 5.8, upper right).

The Kanade-Okutomi variable-window refinement method (Kanade & Okutomi, 1990) uses a statistical analysis to grow the window from 3x3 to some maximum, stopping when an error criterion (based on local changes in intensity and disparity) is exceeded. For this test we let disparity vary between 0 and 50 pixels (as in our method), let the window size vary from 3 to 21 pixels, and ran the method for 10 iterations. It approximated the surface shape well, but produced many more outliers and quantized the flat tilted surface into several stair-step frontoplanar patches (see Figure 5.8, upper right). The RMS error of this method was 0.99 pixels over the entire plate, with $\sigma = 2.36$.

The uncorrected phase method results are also shown in Figure 5.8. The same 501 potential disparities were considered, but foreshortening correction was not applied. The RMS error of this result was 3.77 pixels over the plate, with $\sigma = 6.23$. The main source of error is a general flattening trend over the entire plate, most likely due to the larger windows used at lower frequencies. Like most traditional stereo matchers, the uncorrected method has a strong bias toward frontoplanar surfaces, but unlike Kanade/Okutomi this uncorrected phase method is unable to restrict its attention to the smallest-sized windows.

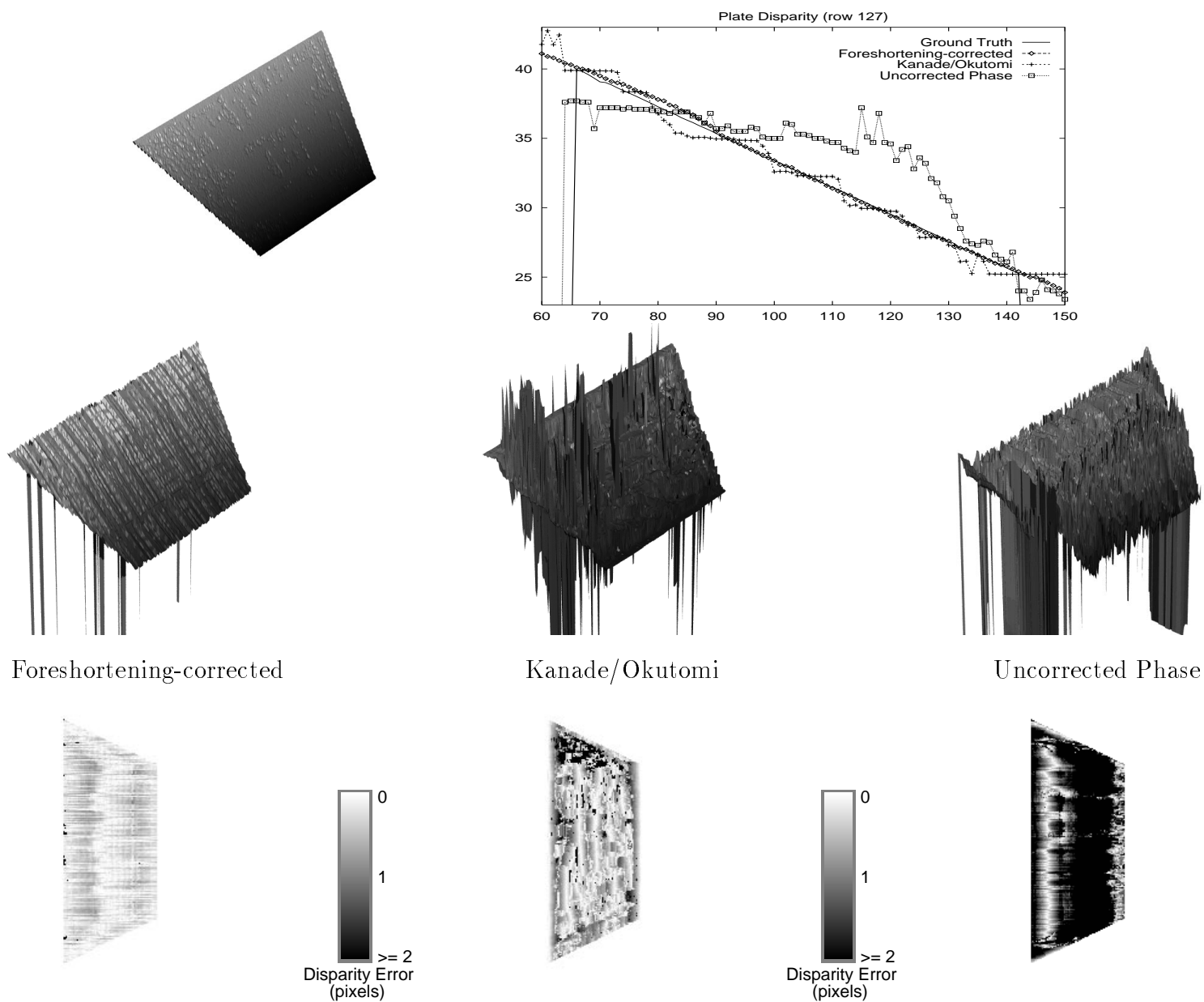


Figure 5.8: Ground Truth and computed disparity maps for a surface angled at 65° . The top row shows ground truth in perspective on the left, a graph of a representative scanline from all methods on the right. The middle row shows perspective views of the disparity maps computed by the foreshortening-corrected method, Kanade/Okutomi and the uncorrected phase method. The bottom row shows differences between actual disparities and those computed by the foreshortening-corrected method, Kanade/Okutomi and the uncorrected phase method, for pixels that image the plate; darker values denote larger errors. Only differences between 0 and 2 pixels are shown, errors larger than 2 pixels appear as a 2 pixel error. Actual plate disparities range from 25.3 to 39.9 pixels.

Other Rotation Angles A cross-section of results for different angles of rotation is presented in Figure 5.9. For these results only a representative scanline is shown, to demonstrate how closely the computed disparity matches the actual ground truth. Only the disparities on the plate itself are correct because the region behind it is a plain white background, and there is no way to distinguish the correct disparity of a featureless surface.

The uncorrected method does reasonably well with small angles, but at slants greater than 30° its performance degrades by several pixels. In contrast, the foreshortening-corrected method performs well even at 75° , though at 80° the systematic error becomes more apparent.

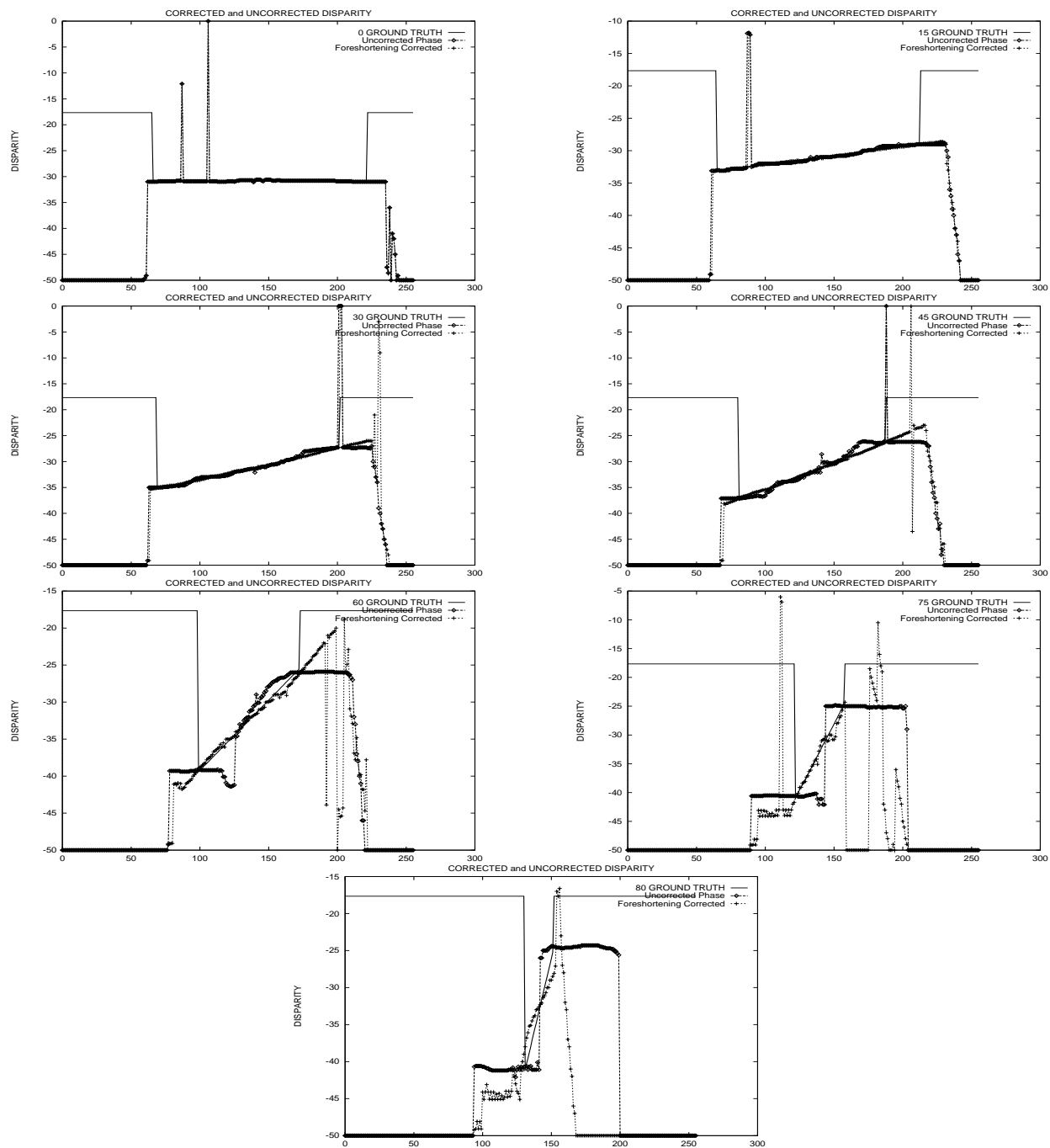


Figure 5.9: Ground truth and disparity (computed by both the uncorrected and foreshortening-corrected phase methods) for the center scanline of the city scene at various rotations. From left to right (and top to bottom): 0, 15, 30, 45, 60, 75, and 80 degrees.

Chapter 6

Contributions

And at the laste he took conclusioun.

— Chaucer, *the Knight's Tale*

This thesis has addressed several long-standing problems in stereo vision: perspective foreshortening, ambiguous matches, and the quantitative evaluation of stereo results. We also demonstrate by example the utility of the local spatial frequency representation in the context of stereo vision. Some particular contributions include:

Perspective Foreshortening We addressed the long-standing problem of perspective foreshortening in stereo vision. Ours is the first work to provide an analytical closed-form expression for the effect of perspective foreshortening on stereo matching in the frequency domain, and demonstrate that the model improves the results of a real stereo method. Results relating foreshortening to the frequency domain have appeared in the shape-from-texture literature (Krumm, 1993; Malik & Perona, 1989), but without a description of the use of disparity; instead, the mapping between two areas was left as a general affine matrix. Similarly, (Jones & Malik, 1991) uses an affine matrix to relate two corresponding image patches in the context of stereo, but does not provide the analytic relationship. Ours is the first presentation to unify these three themes (stereo, foreshortening, frequency domain) and the first to demonstrate good matching even at slants up to 75 degrees (the spatial-domain system of (Belhumeur, 1993) might perform as well, but this has not been demonstrated).

Adaptive Scale Selection We demonstrated the improvement of our scale-adaptive algorithm over traditional multiscale (e.g., coarse to fine) algorithms. Ours is the first multiscale phase-based method that is not confused by missing information at intermediate scales. Most of the prior stereo systems work either at a single scale or using a coarse-to-fine approach in a fixed order. Ours is the first system to invoke automatic scale selection (via filter magnitude weighting) in a nonrestrictive manner. (Sanger, 1988) also used magnitude weighting, but imposed a restrictive limit on candidate disparities. (Jones, 1991) also used some scale-space processing, but only to eliminate the coarsest scales at depth discontinuities. Ours is the first system that can handle missing information at *any* scales.

Phase Wraparound We eliminated the restriction of previous phase-based stereo algorithms on the maximum disparity range, and described a new stereo algorithm that eliminates the problem of phase wraparound. Ours is the first phase-based stereo method to overcome these limitations.

New Disparity Model We proposed a new error model for disparity that more accurately represents the inaccuracies that result from ambiguous matches. This new model characterizes the results of stereo processing more appropriately than the traditional two parameter model (disparity value with variance).

Datasets We proposed a new taxonomy for stereo vision experiments, and provided some of the first public datasets with piecewise dense ground truth. These datasets and the easy-to-use tools for creating them should encourage the community to use (and demand) more quantitative evaluations of stereo systems.

To summarize, our method uses an adaptive search through scale space. We combine estimates from the most reliable scales in a framework that can be used to evaluate the likelihood of *arbitrary* disparities at each pixel; we are not limited by the wavelength of any single filter, as are all previous phase-based methods. The method eliminates the need to perform explicit phase unwrapping, thereby improving accuracy, but at the cost of an additional search over candidate disparities.

Potential future extensions to this work include:

-
1. Addition of a precision parameter (variance) to the representation of ground truth in Chapter 2 to account for measurement errors and variation within a pixel.
 2. Development of a synthetic data generator with better lighting models (e.g., using radiosity rather than ray tracing).
 3. Addition of sequences of stereo images to the taxonomy in Chapter 2.
 4. Experimentation with the shape of the evaluation function (as in Equation 3.14); a replacement for AbsDiffMod such as a cosine might yield smoother results, or be faster to implement.
 5. Speed up the processing by incorporating the fast wavelet transform with appropriate interpolation in place of the complete scalogram computation, using a smaller set of filters, and a smaller set of foreshortening angle candidates.
 6. Extend the occlusion and disparity models to account for multiple depths within a pixel.
 7. Explore the use of nonsymmetric filters (in place of the Gaussian envelope used by Gabor filters) to better address the disparity spillover that occurs at depth discontinuities.
 8. Develop search strategies over the binary “peaks-only” images of Chapter 4 to merge pixel disparities into potential surfaces.

Appendix A

Theodolite Error Analysis

Both stereo disparity and ground truth measurements have finite precision which should be made explicit. As a first step toward extending our notion of ground truth to include this precision, we present in this section an analysis of the resolution obtainable using surveyor's theodolites in their present configuration as part of the Calibrated Imaging Laboratory (CIL).

The Calibrated Imaging Laboratory theodolites (Sokkisha, 1984) can repeatably measure angles to within about 20 seconds of arc. That is, during a single test run, an individual can repeatedly aim the theodolite site at a target, unlock it, then aim again and be confident that the difference between successive measurements will never be more than 20 seconds. This is in spite of the fact that the instrument readout is apparently measured to the nearest tens of seconds (*tens* not *tenths*).

We would like to know how accurate subsequent **X-Y-Z** computations can be, under this limitation. A simple two-dimensional (**X-Z**) analysis will give us a rough idea of the magnitude of the precision in the horizontal plane. Figure A.1 shows the overall model: depending on the angles measured, the computed depth D might lie anywhere within the shaded region. Since that region is polygonal, we know that the largest possible error (i.e., the maximum distance between any two points in the region) will occur between the endpoints of one of its two diagonals. Just how long are the diagonals? To determine that, we need to derive equations for the horizontal and depth coordinates. We address the horizontal first.

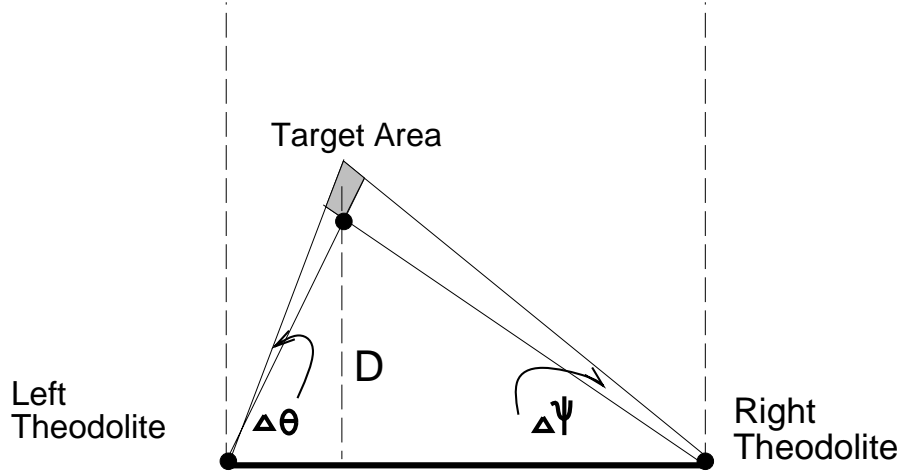


Figure A.1: Region of error. The greatest possible error occurs across one of the Target Area diagonals (see Figure A.3 for a close up).

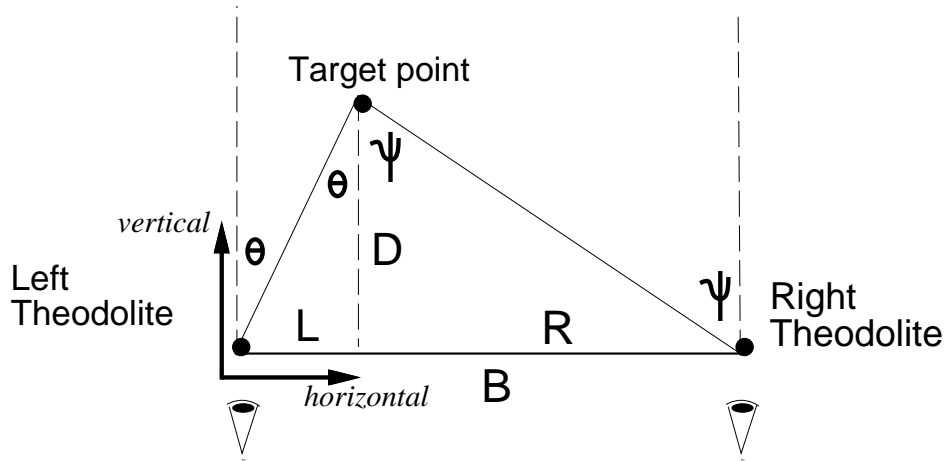


Figure A.2: Two dimensional view of the theodolite imaging process.

A.1 Deriving equations for the coordinate axes

Figure A.2 shows the geometry of the scene. We will treat the two theodolites (as well as the target) as points. θ is the angle measured by the left theodolite, ψ is that measured by the right. B is the length of the baseline between the two theodolites; the baseline is split in two at the projection of the target point: $B = L + R$. D is the distance from the baseline to the target point. If we define the left theodolite to be the origin of a coordinate system with horizontal axis along the baseline, we have D as the vertical *depth* coordinate, and L

as the horizontal coordinate.

Now we need to express D and L as functions of the two angles and baseline alone. The left and right angles bear a simple relationship to the two baseline parts:

$$\tan \theta = \frac{L}{D} \quad \text{and} \quad \tan \psi = \frac{R}{D}$$

Solving for D and setting them equal, we have:

$$\frac{L}{R} = \frac{\tan(\theta)}{\tan(\psi)}$$

Recalling that L and R sum to the baseline B , we have by substitution:

$$L = \frac{B}{1 + \frac{\tan(\psi)}{\tan(\theta)}} \quad \text{and} \quad R = \frac{B}{1 + \frac{\tan(\theta)}{\tan(\psi)}} \quad (\text{A.1})$$

This solves our horizontal coordinate problem: we have an expression for L that depends only on angles θ , ψ and the baseline. Moving on, these equations give us two expressions for the distance D between the baseline and target point:

$$D = \frac{L}{\tan \theta} = \frac{R}{\tan \psi}$$

both of which reduce to the same symmetric result:

$$D = \frac{B}{\tan \theta + \tan \psi} \quad (\text{A.2})$$

So now we have the depth coordinate D as well. Without loss of generality, we assume the baseline is a constant factor and write $D = d(\theta, \psi)$ and $L = l(\theta, \psi)$. Now we can compute the lengths of the two diagonals in the Target Area; the larger one will give us the maximum possible error.

To compute the lengths of the diagonals h and v (shown in Figure A.3), we find the Euclidean distance between their endpoints. Call the theodolite measurement error δ : for the CIL theodolites δ is 20". Then the length of the “horizontal” diagonal h (it’s not really horizontal) is:

$$h(\theta, \psi) = \sqrt{(d(\theta, \psi - \delta) - d(\theta - \delta, \psi))^2 + (l(\theta, \psi - \delta) - l(\theta - \delta, \psi))^2} \quad (\text{A.3})$$

The vertical diagonal v (it’s not really vertical) is computed in the same way:

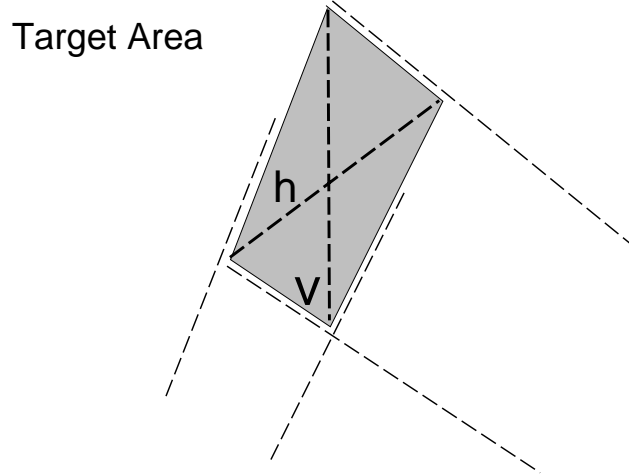


Figure A.3: Target Area Precision (zoom in on Figure A.1): the largest error is the length of one of the diagonals h and v .

$$v(\theta, \psi) = \sqrt{(d(\theta - \delta, \psi - \delta) - d(\theta, \psi))^2 + (l(\theta - \delta, \psi - \delta) - l(\theta, \psi))^2} \quad (\text{A.4})$$

Now we're done; the length of the horizontal diagonal is the maximum error in the horizontal direction, the vertical diagonal is the maximum error in depth. All that remains is to plug in the measurement error $\delta = 20''$.

A.2 Results

In the current laboratory configuration, the theodolite measurements will vary from 10 to 70 degrees, if all objects of interest lie on or above the object optical table. Figure A.4 shows the shape of the error surface for this configuration, with $\delta = 20''$ and assuming a unit baseline. How do we interpret this?

The largest error in Figure A.4 is about 0.0015, when both theodolites have angles of 10° . What does this really mean? Assuming a baseline of 86.1 inches (219cm), it means an object 6.21m away can only be measured to within 3.3mm.¹ However, the far end of the optical table is only about 3m away from the baseline, and according to the model the center point on the far edge gives angles of about 20° for each theodolite. The error for those angles is 0.0004, which means the best precision for the far end of the table is 0.876mm = $0.0004 \cdot 219\text{cm}$. Thus

¹6.21m = $\frac{2.19\text{m}}{\tan 10^\circ + \tan 10^\circ}$ and 3.3mm = $219\text{cm} \cdot 0.0015$

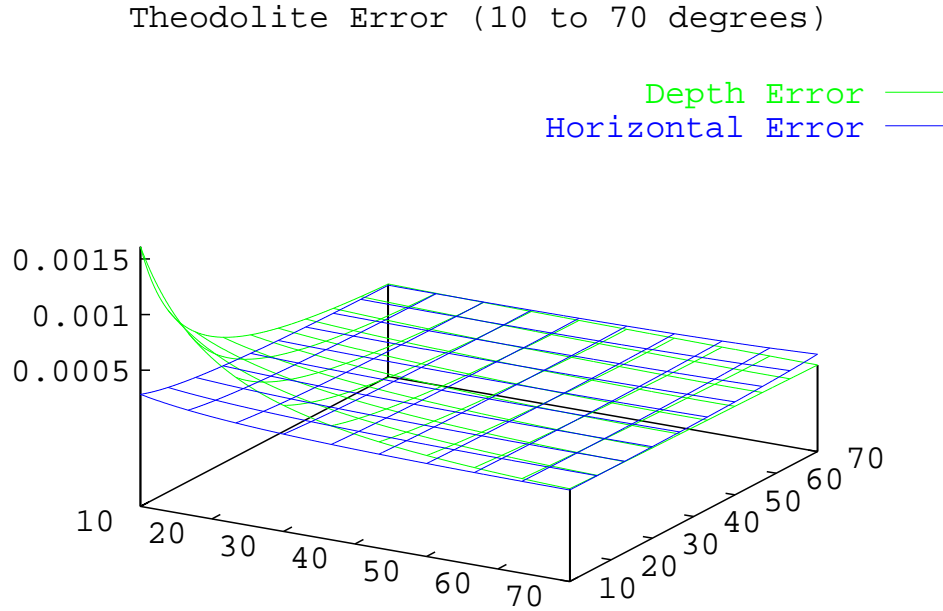


Figure A.4: Error Space for $\delta = 20''$ with unit baseline: depth error $v(\theta, \psi)$ is greater than horizontal error $h(\theta, \psi)$ for all but the nearest points (when the sum of the angles is greater than about 90 degrees)

it's safe to say our theodolite measurements over the optical table are in general accurate to within a millimeter.

Nearer measurements have better accuracy. For example, in the **CIL-0001** Stereo Dataset the left horizontal angles range from 41 to 50 degrees, and right angles from 31 to 42 degrees. The maximum error in that range is $0.296mm = 0.000135 \cdot 219cm$.

It is in fact possible to compute the precision at each point. Simply plug the angles measured into the error terms defined above. Or to estimate the precision for a whole region, run the code in Table A.1 through the GNUplot package with appropriate limits (instead of 10° to 70°), and visually pick out the largest error.

Limitations: The acquisition of ground truth requires that static objects be imaged in a laboratory environment. However there are many applications for which dynamic imagery is required, and stereo systems must be tested using comparable data. The acquisition of such data implies that no ground truth will be available.

```
d2r = pi / 180.0
hor(x,y) = 1/(1+tan(x * d2r)/tan(y * d2r))
ver(x,y) = 1/(tan(x * d2r) + tan(y * d2r))
d = 20.0/3600
splot [10:70] [10:70] \
    sqrt((hor(x-d,y-d)-hor(x,y))**2 + \
        (ver(x-d,y-d)-ver(x,y))**2) title "Depth Error", \
    sqrt((hor(x-d,y)-hor(x,y-d))**2 + \
        (ver(x-d,y)-ver(x,y-d))**2) title "Horizontal Error"
```

Table A.1: GNUplot commands that generated Figure A.4

Appendix B

Derivation of Foreshortening Probabilities

We present here the derivation of one of the probabilities from Table 5.1. Figure 5.7 showed the Foreshortening Factor computed from many combinations of depth and orientation (except for the extreme values near the point at which it approaches infinity). The graph makes it clear that the Foreshortening Factor has its greatest impact when objects are sharply slanted and/or located near the cameras. We can quantify its influence using the contour lines that separate regions of large and smaller foreshortening effects. Suppose we assume that surface depth and orientation are uniformly distributed throughout a scene. Then we can compute the probability that a surface will require at least a 10% correction term by finding the area under the 1.1 Foreshortening Factor contour curve. This derivation assumes the depth range begins at zero (the more general results require a little more work).

We want to find:

$$P(\text{Foreshortening Factor} \geq 1.1 \text{ or } \leq 0.9) = P\left(\left|\frac{\tan \theta}{d}\right| \geq 0.1\right)$$

for $d = \frac{Z_L}{B} \in [0 : 100]$ and $\theta \in [-\frac{\pi}{2} : \frac{\pi}{2}]$, each uniformly distributed. Since \tan is symmetric we can eliminate the absolute value by restricting the angle θ to $[0 : \frac{\pi}{2}]$. Continuing:

$$\begin{aligned} P\left(\left|\frac{\tan \theta}{d}\right| \geq 0.1\right) &= P\left(\frac{\tan \theta}{d} \geq 0.1\right) = \frac{\int_0^{\frac{\pi}{2}} \min\left(\frac{\tan \theta}{\text{Foreshortening Factor}-1}, 100\right) d\theta}{\int_0^{\frac{\pi}{2}} \int_0^{100} dd d\theta} \\ &= \frac{\int_0^{\frac{\pi}{2}} \min\left(\frac{\tan \theta}{0.1}, 100\right) d\theta}{50\pi} \end{aligned}$$

To eliminate the min from the integral we must find the minimum angle requiring a 10% correction at distance 100:

$$\theta_{\min} = \arctan 100(\text{Foreshortening Factor} - 1) = \arctan 10 = 84.2894^\circ$$

Now we can split up the integral into two parts and evaluate it:

$$\begin{aligned} \int_0^{\frac{\pi}{2}} \min\left(\frac{\tan \theta}{0.1}, 100\right) d\theta &= \int_0^{\theta_{\min}} \frac{\tan \theta}{0.1} d\theta + \int_{\theta_{\min}}^{\frac{\pi}{2}} 100 d\theta \\ &= \left[\frac{\ln \sec \theta}{0.1} \right]_0^{\theta_{\min}} + 9.96688 \\ &= 23.0756 + 9.96688 \end{aligned}$$

This brings us to the final result:

$$P\left(\left|\frac{\tan \theta}{d}\right| \geq 0.1\right) = \frac{33.0425}{50\pi} = 0.210355$$

So under the assumption of uniform distribution on depth ratio from 0 to 100 and angle from -90° to 90° , the probability of a surface exhibiting at least a 10% foreshortening effect is 0.210355.

Appendix C

Numbers

Table C.1 shows the default parameter values used to generate the images in this thesis. Additional details can be found in the Matlab source code used to generate many of the actual images, especially those in Chapters 3 and 4. This code is available on the web from *Mark Maimone's Index Page*.¹

¹<http://www.cs.cmu.edu/~mwm/>

Parameter	Value	Comments
Disparity Range d	$0, 1, \dots, 50$	The range of candidate disparities can usually be inferred by figure axes, and the step size by the visible quantization in disparity space images. Often the step size will be less than 1, which indicates subpixel precision.
Wavelengths W	$2, 3, \dots, n/4$	Unless otherwise stated, we typically use the image scalogram sampling which is linear in wavelength.
Wavelengths per window m	4	The Gaussian envelope of a Gabor filter will be truncated outside this many wavelengths. This determines the window size.
Sigma fraction σ_f	$1/6$	Given a fixed window size (m), the σ parameter of the Gaussian will be this fraction of the window.
Fleet threshold	.05	Cutoff for Equation 3.10.

Table C.1: Default parameter values used to generate images in this thesis.

Appendix D

Application of Visual Reconstruction



XXXX XXXXXXXX XXXXX

Pittsburgh, PA XXXXX

(412) XXX - XXXX

work: (412) 268 - 7698

email: mwm@cs.cmu.edu

October 3, 1995

Prof. Peter Schickele

XXXXX XXXXXXXXXXXX / NY Campus of U of SND at H

XXX XXXX XXXX XXXXXX

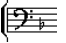
New York, NY XXXXX

Dear Colleague:

Greetings! I am writing to inform you of a most remarkable discovery: the visual appearance of Herr P.D.Q. Bach's headstone within the mausoleum on the outskirts of Baden-Baden-Baden. Although said mausoleum was destroyed in the 1840s ([1] pp. 112-113), computer analysis of the data collected by your esteemed research team ([1] pp. 131-135) has made possible visualization of this artifact whose historical significance cannot be un-

derestimated.

As you are no doubt aware, my own area of expertise lies in techniques for reconstruction of three-dimensional shapes from oblique imagery [2]. Application of these techniques to those passages of text referring to Herr Bach's exploits at Madame Höllender's ([1] pp. 100-101) resulted in images of a variety of organic experience [3] unsuitable for publication in any *reputable* medium of scientific discourse. I remained *firmly committed* to the project, however (anybody who says I'm not is a liar), and following a, hunch, applied the principles of super-resolution [4] to your original images graciously made available to me by a random house-keeper for a remarkably small fee.

I shan't bore you with the details of the reconstruction, but merely outline the process. Having accepted our dough, a dear (a female dear) house-keeper started us off with the data. Following the principles of laser holography [5], a specially-modified computer filtered the light ray (a drop of gold on Sun interface was required) passed by me through a Headpiece to the Staff of Fa ([6] ) over a distance of some kilometers, a long long way to run. So like a needle pulling thread¹, the beam was fed through the IID/NM T-junction [7] while I retired for a drink with jam and bread (that will bring us back to dough). Having scaled these obstacles – naturally filtering out any non-diatons we found floating around – we achieved our final result, an image of which is enclosed herewith.

Should you find yourself with access to the World Wide Web I encourage you to examine:

<http://www.cs.cmu.edu/~mwm/pdq/smix.html>

for background.

Anxiously awaiting your reply I remain...

Your Honor's and my most especially Highly
Honored Sir's most obedient servant,

Mark W. Maimone

¹La

References

- [1] P. P. Schickele. *The Definitive Biography of P.D.Q. Bach*. Random House, 1976.
- [2] M. W. Maimone and S. A. Shafer. Modeling foreshortening in stereo vision using local spatial frequency. In *Intelligent Robotics and Systems*, pages 519–524, 1995.
- [3] P. Schickele. Schickele Mix no. 40. Public Radio International (formerly APR).
- [4] P. Cheeseman, B. Kanefsky, R. Kraft, J. Stutz, and R. Hanson. Super-resolved surface reconstruction from multiple images. Technical Report FIA-94-12, NASA Ames AI Research, December 1994.
- [5] S. A. Benton. Photographic holography. In *SPIE Proc. Optics in Entertainment*, volume 391, pages 2–9, January 1983.
- [6] I. Jones, M. Brody, and Sallah. Scaling the tannis map room. [recently declassified Top Secret Army Intelligence Report 9906753], 1936.
- [7] D. Adams. *The Restaurant at the End of the Universe*, pages 11–13. Random House, 1980.



Bibliography

- Bani-Hashemi, A. (1993). A Fourier Approach to Camera Orientation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 15(11):1197–1202.
- Barnard, S. T. & Fischler, M. A. (1982). Computational stereo. *Computing Surveys*, 14(4):553–572.
- Barron, J. L., Fleet, D. J., & Beauchemin, S. S. (1994). Performance of optical flow techniques. *International Journal of Computer Vision*, 12(1):43–77. <ftp://csd.uwo.ca/pub/vision>.
- Belhumeur, P. N. (1993). A binocular stereo algorithm for reconstructing sloping, creased, and broken surfaces in the presence of half-occlusion. In *International Conference on Computer Vision*, pages 431–438.
- Belhumeur, P. N. (1995). A bayesian approach to binocular stereopsis. To appear in IJCV.
- Besl, P. J. (1988). Active, optical range imaging sensors. *Machine Vision and Applications*, 1:127–152.
- Bhat, D. N. & Nayar, S. K. (1995). Stereo in the presence of specular reflection. In *International Conference on Computer Vision*, pages 1086–1092.
- Bolles, R. C., Baker, H. H., & Hannah, M. J. (1993). The JISCT Stereo Evaluation. In *ARPA Image Understanding Workshop*, pages 263–274.
- Boufama, B. & Mohr, R. (1995). Epipole and fundamental matrix estimation using virtual parallax. In *International Conference on Computer Vision*, pages 1030–1036.
- Boyer, E. & Berger, M. O. (1995). 3D Surface Reconstruction Using Occluding Contours. In *International Conference on Computer Analysis of Images and Patterns*.

- Brown, L. G. (1992). A survey of image registration techniques. *ACM Computing Surveys*, 24(4):325–376.
- Carson, J. R. & Fry, T. C. (1937). Variable-frequency electric circuit theory. *Bell System Technical Journal*, 16:513–540.
- Dhond, U. R. & Aggarwal, J. K. (1989). Structure from stereo — a review. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1489–1510.
- Faugeras, O. D. (1993). *Three-Dimensional Computer Vision*. MIT Press.
- Faugeras, O. D. & Toscani, G. (1986). The calibration problem for stereo. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 15–20.
- Fleet, D. & Jepson, A. (1993). Stability of phase information. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 15(12):1253–1268.
- Fleet, D. J., Jepson, A. D., & Jenkin, M. R. M. (1991). Phase-based disparity measurement. *CVGIP: Image Understanding*, 53(2):198–210.
- Gabor, D. (1946). Theory of communication. *Journal of the Institution of Electrical Engineers*, pages 429–457.
- Ghiglia, D. C. & Romero, L. A. (1994). Robust two-dimensional weighted and unweighted phase unwrapping that uses fast transforms and iterative methods. *Journal of the Optical Society of America A: Optics, ..., Vision*, 11(1):107–117.
- Grant, C. W. (1992). *Visibility Algorithms in Image Synthesis*. PhD thesis, University of California Davis Computer Science.
- Gülch, E. (1991). Results of test on image matching of ISPRS WG III/4. *ISPRS Journal of Photogrammetry and Remote Sensing*, 46:1–18.
- Haines, E. (1993). Free ray tracer summary. *Ray Tracing News*, 6(3):article 4. <ftp://ftp-graphics.stanford.edu/pub/Graphics/RTNews/html/rtnv6n3.html#art4>.
- Hartley, R. I. (1994). Self-calibration from multiple views with a rotating camera. In *European Conference on Computer Vision*, pages 471–478.

- Hebert, M. & Krotkov, E. (1992). 3D measurements from imaging laser radars: How good are they? *Intl. Journal of Image and Vision Computing*, 10(3):170–178.
- Hlawatsch, F. & Boudreaux-Bartels, G. F. (1992). Linear and quadratic time-frequency signal representations. *IEEE Signal Processing Magazine*, pages 21–67.
- Horn, B. K. (1986). *Robot Vision*. MIT Press.
- Horner, J. L. & Gianino, P. D. (1984). Phase-only matched filtering. *Applied Optics*, 23(6):812–816.
- Jenkin, M. R. M. & Jepson, A. D. (1994). Recovering local surface structure through local phase difference measurements. *Computer Vision, Graphics and Image Processing: Image Understanding*, 59(1):72–93.
- Jones, D. G. (1991). *Computational Models of Binocular Vision*. PhD thesis, Stanford University.
- Jones, D. G. & Malik, J. (1991). Determining three-dimensional shape from orientation and spatial frequency disparities part ii - using corresponding image patches. Technical Report UCB/CSD 91/657, University of California Berkeley Computer Science Department.
- Kanade, T., Kano, H., Kimura, S., Yoshida, A., & Oda, K. (1995). Development of a video-rate stereo machine. In *International Robotics and Systems Conference (IROS)*, volume 3, pages 95–100.
- Kanade, T. & Okutomi, M. (1990). A stereo matching algorithm with an adaptive window: Theory and experiment. In *DARPA Image Understanding Workshop Proceedings*, pages 383–398.
- Kanade, T. & Okutomi, M. (1991). A stereo matching algorithm with an adaptive window: Theory and experiment. In *Intl. Conference on Robotics and Automation*, pages 1088–1095.
- Kanade, T., Okutomi, M., & Nakahara, T. (1992). A multiple-baseline stereo method. In *ARPA Image Understanding Workshop*, pages 409–426.

- Krumm, J. (1993). *Shape from Texture and Segmentation using Local Spatial Frequency*. PhD thesis, Carnegie Mellon Robotics Institute.
- Kuglin, C. D. & Hines, D. C. (1975). The phase correlation image alignment method. In *Proceedings of the IEEE Int. Conference on Cybernetics and Society*, pages 163–165.
- Maimone, M. W. (1995). Watch the Birdie: A Guide to Imaging in the Calibrated Imaging Laboratory. Unpublished CMU Calibrated Imaging Laboratory (CIL) Tech Report.
- Maimone, M. W. & Shafer, S. A. (1995a). Modeling foreshortening in stereo vision using local spatial frequency. Technical Report CMU-CS-95-104, Carnegie Mellon University Computer Science Department.
- Maimone, M. W. & Shafer, S. A. (1995b). Modeling foreshortening in stereo vision using local spatial frequency. In *International Robotics and Systems Conference (IROS)*, pages 519–524. IEEE Computer Society Press. <http://www.ius.cs.cmu.edu/project/cil/fore/tr.html>.
- Maimone, M. W. & Shafer, S. A. (1996). A taxonomy for stereo computer vision experiments. *ECCV'96 Workshop on Performance Characteristics of Vision Algorithms* <http://www.ius.cs.cmu.edu/project/cil/tax/>.
- Malik, J. & Perona, P. (1989). A computational model of texture segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 326–332.
- Matthies, L. (1989). *Dynamic Stereo Vision*. PhD thesis, Carnegie Mellon University Computer Science Department.
- Matthies, L. H. (1992). Stereo vision for planetary rovers: stochastic modeling to near real-time implementation. *International Journal of Computer Vision*, 8(1):71–91.
- Mori, K.-I., Kidode, M., & Asada, H. (1973). An iterative prediction and correction method for automatic stereocomparison. *Computer Graphics and Image Processing*, 2(3/4):393–401.
- Nakamura, Y., Matsuura, T., Satoh, K., & Ohta, Y. (1996). Occlusion detectable stereo — occlusion patterns in camera matrix. In *IEEE Conference on Computer Vision and Pattern Recognition*.

- Nayar, S. K., Watanabe, M., & Noguchi, M. (1995). Real-time focus range sensor. In *International Conference on Computer Vision*, pages 995–1001.
- Nishihara, H. K. (1984). Prism: A practical real-time imaging stereo matcher. Technical Report AI Memo 780, Massachusetts Institute of Technology AI Laboratory.
- Okutomi, M. & Kanade, T. (1991). A multiple-baseline stereo. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 63–69.
- Oppenheim, A. V. & Lim, J. S. (1981). The importance of phase in signals. *Proceedings of the IEEE*, 69(5):529–541.
- Photometrics (1990). *Charge-Coupled Devices for Quantitative Electronic Imaging*. Photometrics, Ltd., 3440 E. Britannia Drive, Tucson, AZ 85706.
- Reeves, W. T., Salesin, D. H., & Cook, R. L. (1987). Rendering antialiased shadows with depth maps. In *ACM Computer Graphics SIGGRAPH*, pages 283–291.
- Rioul, O. & Vetterli, M. (1991). Wavelets and signal processing. *IEEE Signal Processing Magazine*, pages 14–38.
- Robert, L., Buffa, M., & Hebert, M. (1994). Weakly-calibrated stereo perception for rover navigation. In *ARPA Image Understanding Workshop*.
- Ross, B. (1993). A practical stereo vision system. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 148–153.
- Sanger, T. D. (1988). Stereo disparity computation using gabor filters. *Biological Cybernetics*, 59:405–418.
- Shaw, M. (1990). Prospects for an engineering discipline of software. Technical Report CMU-CS-90-165, Carnegie Mellon University Computer Science Department. Also IEEE Software, Nov 1990.
- Shum, H.-Y., Hebert, M., Ikeuchi, K., & Reddy, R. (1995). An integral approach to free-form object modeling. Technical Report CMU-CS-95-135, Carnegie Mellon University Computer Science Department.

- Smith, P. W. & Nandhakumar, N. (1996). An improved power cepstrum based stereo correspondence method for textured scenes. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 18(3):338–348.
- Sokkisha (1984). *Electronic Digital Theodolite DT20E Operation Manual*. Sokkisha Co., Ltd., Keio Yoyogi Building 5th Floor, No. 1, 1, 1-chome, Tomigaya, Shibuta-ku, Tokyo, 151 Japan.
- Stein, G. P. (1995). Accurate internal camera calibration using rotation, with analysis of sources of error. In *International Conference on Computer Vision*, pages 230–236.
- Stevenson, D. E. & Fleck, M. M. (1995). Robot aerobics: Four easy steps to a more flexible calibration. In *International Conference on Computer Vision*, pages 34–39.
- Szeliski, R. (1994). Image mosaicing for tele-reality applications. Technical Report 94/2, DEC Cambridge Research Lab.
- Tada, S., Gruss, A., & Kanade, T. (1993). CMU very fast range-imaging system. Technical Report CMU-CS-93-179, Carnegie Mellon University Computer Science Department.
- Tribolet, J. (1977). A new phase unwrapping algorithm. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 25(2):170–177.
- Tsai, R. Y. (1987). A versatile camera calibration technique for high-accuracy 3D machine vision metrology using off-the-shelf tv cameras and lenses. *IEEE Journal of Robotics and Automation*, RA-3(4):323–344. <http://www.cs.cmu.edu/~rgw/TsaiCode.html>.
- Wang, Z. & Jepson, A. (1994). A new closed-form solution for absolute orientation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 129–134.
- Webb, J. (1993). Implementation and performance of fast parallel multi-baseline stereo vision. In *Computer Architectures for Machine Perception*. New Orleans, LA. Also appeared in ARPA Image Understanding Workshop, Washington, D.C.
- Weng, J. (1993). Image matching using the windowed Fourier phase. *International Journal of Computer Vision*, 11(3):211–236.
- Westelius, C.-J. (1995). *Focus of Attention and Gaze Control for Robot Vision*. PhD thesis, Linköping University. Dept. of Electrical Engineering Dissertation 379.

- Williams, L. (1978). Casting curved shadows on curved surfaces. In *ACM Computer Graphics SIGGRAPH*, pages 270–274.
- Willson, R. G. (1994). *Modeling and Calibration of Automated Zoom Lenses*. PhD thesis, Carnegie Mellon Electrical and Computer Engineering.
- Xiong, Y. (1995). *High Precision Image Matching and Shape Recovery*. PhD thesis, Carnegie Mellon Robotics Institute.
- Yeshurun, Y. & Schwartz, E. L. (1989). Cepstral filtering on a columnar image architecture: A fast algorithm for binocular stereo segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 11(7):759–767.
- Zhang, Z., Faugeras, O., & Deriche, R. (1995). Calibrating a binocular stereo through projective reconstruction using both a calibration object and the environment. In *Europe-China Workshop on Geometrical modelling and Invariance for Computer Vision*, pages 253–260.
- Zitnick, C. L. & Webb, J. A. (1995). The reformulation of stereo vision. Unpublished Tech Report by Visual Interface, Inc.

Index

- accuracy, 2
- Ambiguity
 - Factor α , 85
- autostereogram, 82
- banana, 81
- calibration
 - errors, 81
- camera
 - calibration, 38
 - models, 35
 - motion, 36
- checkerboard
 - virtual, 31
- correspondence problem, 82
- cross correlation, 5
- Data, Lt. Commander, 47
- disparity, 3
- disparity
 - error image, 4
 - image, 4
 - map, 4
- disparity space, 86
- epipolar
 - constraint, 9
 - line, 9
- evaluation function, 4
- evaluation function
 - profile, 5
- false target, 81
- Fourier Shift Theorem, 55
- fruit flies, 81
- Gabor filters, 51
- histogram equalization, 33
- image pyramid, 94
- JISCT, 89
- local frequency, 8
- modeling errors, 81
- MOVI target, 35
- multibaseline stereo, 1, 13, 30, 82, 93
- Nyquist interval, 50
- occlusion, 93
- occlusion masks, 27
- phase
 - unwrapping, 70, 136
 - wraparound, 24, 70, 72, 73, 77, 104, 136
- precision, 2

SAD

(sum of absolute differences), 5, 95

Sampling Theorem, 50, 52, 66

scalogram, 7, 8, 66–69, 73, 76, 77, 101, 103,
104, 106, 119, 122, 128, 129, 137

SSD

(sum of squared differences), 5

(sum of squared distance), 12

synthetic noise, 30

window effect, 55

window size, 6