



ELSEVIER

Available online at www.sciencedirect.com

SCIENCE @ DIRECT®

Computer Vision
and Image
Understanding

Computer Vision and Image Understanding xxx (2004) xxx–xxx

www.elsevier.com/locate/cviu

Dynamic mutual calibration and view planning for cooperative mobile robots with panoramic virtual stereo vision[☆]

Zhigang Zhu,^{a,*} Deepak R. Karuppiah,^b Edward M. Riseman,^b
and Allen R. Hanson^b

^a *Department of Computer Science, the City College, The City University of New York,
New York, NY 10031, USA*

^b *Department of Computer Science, University of Massachusetts, Amherst, MA 01003, USA*

Received 14 May 2001; accepted 24 February 2004

Abstract

This paper presents a panoramic virtual stereo vision approach to the problem of detecting and localizing multiple moving objects (e.g., humans) in an indoor scene. Two panoramic cameras, residing on different mobile platforms, compose a virtual stereo sensor with a flexible baseline. A novel “mutual calibration” algorithm is proposed, where panoramic cameras on two cooperative moving platforms are dynamically calibrated by looking at each other. A detailed numerical analysis of the error characteristics of the panoramic virtual stereo vision (mutual calibration error, stereo matching error, and triangulation error) is given to derive rules for optimal view planning. Experimental results are discussed for detecting and localizing multiple humans in motion using two cooperative robot platforms.

© 2004 Elsevier Inc. All rights reserved.

[☆] This work was supported by DARPA/ITO programs Mobile Autonomous Robot S/W (MARS) (Contract No. DOD DABT63-99-1-0004) and Software for Distributed Robotics (SDR) (Contract No. DOD DABT63-99-1-0022). It is also partially by the Army Research Office under Grant No. DAAD19-99-1-0016, by Air Force Research Lab (Award No. F33615-03-1-63-83) and by the CUNY GRTI program.

* Corresponding author. Fax: 1-212-650-6248.

E-mail address: zhu@cs.cuny.cuny.edu (Z. Zhu).

URL: <http://www-cs.engr.cuny.cuny.edu/~zhu>.

Nomenclature*List of symbols*

χ	Angular resolution of the panoramic image
D_k	Target's distance to camera k ($k = 1, 2$)
$\phi_0, \phi_1,$ and ϕ_2	Three interior angles of the triangle formed by two cameras and the target
θ_i	Bearing angle of the target in image i ($i = 1, 2$)
β_{ij}	Bearing angle of the camera i in image j
B	Baseline length
R	Radius of the cylindrical body of the mobile robot used for mutual calibration
α	Angle subtended by the cylindrical body of the robot
W	Width of the target
w_i	Width of the target in the panoramic image i ($i = 1, 2$)
$T_i^{(k)}$	Feature set of the blob i in camera k
∂X	Error in estimating parameter X ($X = B, \partial\phi_1, \partial\phi_2, \beta_{21}, \theta_1, \theta_2, \alpha, w,$ and D)
∂D_1^B	Distance error due to the baseline error
∂D_1^ϕ	Distance error due to angular error
∂D_1^+	Distance error when $D_1 > B$
∂D_1^0	Distance error when $D_1 = B$
∂D_1^-	Distance error when $D_1 < B$
$\partial D_1^{\text{fix}}$	Distance error when the baseline is fixed
∂D_1^s	Distance error in the size-ratio method.

24 1. Introduction

25 Flexible, reconfigurable vision systems can provide an extremely rich sensing mo-
 26 dality for sophisticated robot platforms. We propose a cooperative and adaptive ap-
 27 proach to the problem of finding and protecting humans in emergency
 28 circumstances, for example, during a fire in an office building. Real-time processing
 29 is essential for the dynamic and unpredictable environments in our application do-
 30 main, and it is important for visual sensing to rapidly focus attention on important
 31 activity in the environment. Any room or corridor should be searched quickly to de-
 32 tect people and fire. Field-of-view issues using standard optics are challenging since
 33 panning a camera takes time, and multiple targets/objectives may require saccades to
 34 attend to important visual cues. Using multiple conventional cameras covering dif-
 35 ferent fields of view could be a solution, but the cost of hardware (cameras, frame
 36 grabbers, and computers) and software (multiple stream data manipulation) would
 37 increase. Thus, we employ panoramic cameras to detect and track multiple objects
 38 (people) in motion, in a full 360-degree view, in real time.

39 We note that there is a fairly large body of work on detection and tracking of hu-
 40 mans [1–5], motivated most recently by the DARPA VSAM effort [6]. On the other
 41 hand, different kinds of omnidirectional (or panoramic) imaging sensors have been

42 designed [7–12], and a systematic theoretical analysis of omnidirectional sensors has
43 been given by Baker and Nayar [7]. Omnidirectional vision has become quite popu-
44 lar with many vision approaches for robot navigation [10,13,14,30,31], 3D recon-
45 struction [15–17,29] and video surveillance [18–20,28]. Research on multiple
46 camera networks with panoramic cameras that are devoted to human and subject
47 tracking and identification can be found in the literature [19–23,30,31]. The most re-
48 lated work is the realtime human tracking system by Sogo et al. [21] using multiple
49 omnidirectional cameras distributed in an indoor environment. The system detects
50 people, measures bearing angles and determine their locations by triangulation. Gen-
51 erally there are two problems in such a system—(1) the correspondence problem
52 among multiple targets, and (2) the measurement accuracy of target locations. The
53 correspondence problem is more difficult in a panoramic stereo than a conventional
54 stereo because both large baseline and low resolution make it hard to establish cor-
55 respondences of the visual features. The second problem arises particularly when a
56 target is (almost) aligned with the stereo pair. In order to solve these problems, Sogo
57 et al. [21] proposed a “N-ocular stereo” approach without visual features that only
58 verifies the correspondences of multiple targets of binocular stereo by a third omni-
59 directional camera. They showed that the uncertainty in estimating 3D locations was
60 reduced by using the best estimations of pairs of four fixed panoramic cameras put in
61 the vertices of a square region. However, the error of localizing a target is still pro-
62 portional to the square of the target’s distance from the cameras with fixed baseline
63 distances; in their simulation, it increases 7-fold when a target moves 3.5 m away
64 from the cameras. Our work differs from theirs in that we deal with panoramic stereo
65 vision on mobile platforms and thus study the issues of dynamic calibration and view
66 planning. We propose a novel concept of mutual calibration and give a detailed error
67 analysis of panoramic stereo that leads to dynamic stereo configurations with adap-
68 tive baselines and viewpoints for best depth estimation. The distinctive feature of our
69 approach is the ability to compose cooperative sensing strategies across the distrib-
70 uted panoramic sensors of a robot team to synthesize optimal “virtual” stereo vision
71 for human detection and tracking.

72 The idea of distributing sensors and cooperation across different robots stems
73 from the requirements of potentially limited (sensor) resources for a large robot team
74 and the need for mobile placement of sensor platforms given the limited resolution in
75 panoramic sensors. Nevertheless, the advantages of cooperative vision arise from
76 more than this compromise. Any fixed-baseline stereo vision system has limited
77 depth resolution because of the physical constraints imposed by the separation of
78 cameras, whereas a system that combines multiple views allows the planning system
79 to take advantage of the current context and goals in selecting viewpoints. In this
80 paper, we focus on the cooperative behavior involving cameras that are aware of
81 each other, residing on different mobile platforms, to compose a virtual stereo sensor
82 with a flexible baseline. In this model, the sensor geometry can be controlled to man-
83 age the precision of the resulting virtual sensor. The cooperative stereo vision strat-
84 egy is particularly effective with a pair of mobile panoramic sensors that have the
85 potential of almost always seeing each other. Once calibrated by “looking” at each
86 other, they can view the environment to estimate the 3D structure of the scene.

87 The organization of the paper is as follows. After the introduction of two depth
88 estimation methods using panoramic sensors in Section 2, we will mainly focus on
89 the following critical issues of a panoramic virtual stereo system:

- 90 (1) Dynamic “mutual calibration” between the two cameras on two separate mobile
91 robots that forms the dynamic “virtual” stereo sensor with a full 360-degree view
92 (Section 3);
93 (2) A detailed numerical analysis of the error characteristics of the panoramic vir-
94 tual stereo in order to derive the rules for optimal view planning of the moving
95 sensor platforms (Sections 4 and 5); and
96 (3) View planning by taking advantage of the current context and goals, based on a
97 thorough error analysis of panoramic virtual stereo (Section 6).
98 Experimental systems and results for multiple human detection and localization
99 will be given in Section 7, and conclusion and future work will be discussed in the
100 last section.

101 2. Panoramic virtual stereo geometry

102 In our experiments we use the panoramic annular lens (PAL) camera system [9] as
103 it can capture its surroundings with a field of view (FOV) of 360-degrees horizontally
104 and $-15 \sim +20$ degrees vertically (Fig. 1). In the application of human tracking and
105 identification by a mobile robot, a vertical viewing angle that spans the horizon is pre-
106 ferred. After image un-warping, distortion rectification, and camera calibration
107 [22,24], we obtain cylindrical images generated by a panoramic “virtual” camera from

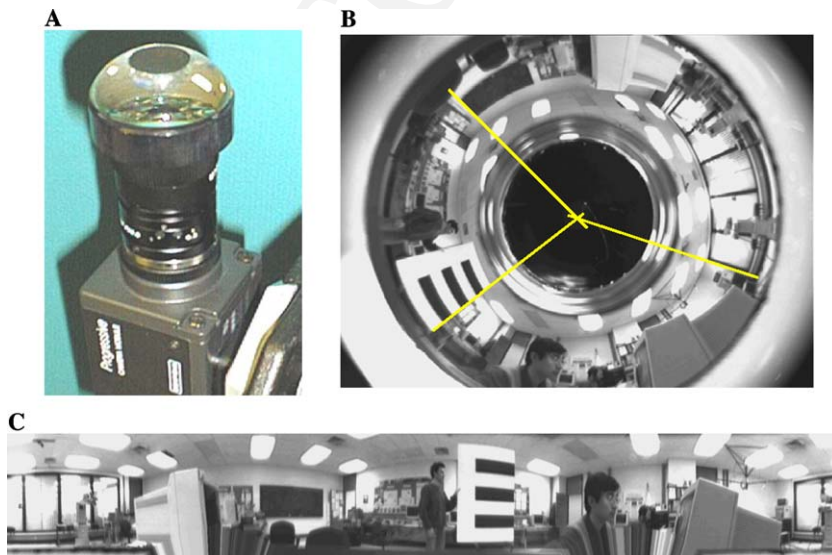


Fig. 1. Panoramic annular lens and images. (A) PAL camera. (B) An original PAL image (768×576). (C) Cylindrical panoramic image.

108 the virtual viewpoint on the axis of the cylindrical image surface (Fig. 1C). Other om-
 109 nidirectional sensors can also be applied; for example, the omnidirectional cameras
 110 proposed by Nayar's research group [7] have been used in our current experiments.

111 Panoramic virtual stereo vision is formed by two panoramic cameras residing on
 112 two separate (possibly mobile) platforms. Let's assume that both of them are subject
 113 to only planar motion on the floor and are at the same height above the floor. Sup-
 114 pose that in Fig. 2A, O_1 and O_2 are the viewpoints of the two cameras and they can
 115 be localized by each other in the panoramic images as P_{12} and P_{21} , respectively. B
 116 is the baseline (i.e., distance O_1O_2) between them. The projection of a target T is rep-
 117 resented by T_1 and T_2 in the two panoramic images. Then a triangle O_1O_2T can be
 118 formed. By defining an arbitrary starting orientation for each cylindrical image,
 119 three angles ϕ_1 , ϕ_2 , and ϕ_0 of the triangle can be calculated from the following four
 120 bearing angles: θ_1 and θ_2 , the bearings of the target in image 1 and image 2, respec-
 121 tively, β_{12} and β_{21} , the bearing angles of camera 1 in image 2, and camera 2 in image
 122 1, respectively. Therefore, the distances from the two cameras to the target can be
 123 calculated by *triangulation* as

$$D_1 = B \frac{\sin \phi_2}{\sin \phi_0} = B \frac{\sin \phi_2}{\sin(\phi_1 + \phi_2)}, \quad D_2 = B \frac{\sin \phi_1}{\sin \phi_0} = B \frac{\sin \phi_1}{\sin(\phi_1 + \phi_2)}. \quad (1)$$

125 With stationary cameras, triangulation error, i.e., the error in estimating D_1 (or D_2)
 126 varies with the target location—larger errors when the target is close to the baseline
 127 and smaller errors when better triangulation is possible. Here, we show first that
 128 panoramic stereo can almost always estimate the distance of the target in the full
 129 360° view. It is commonly known that the triangulation relation in Eq. (1) ap-
 130 proaches singularity as the target moves towards the baseline O_1O_2 . Fortunately,
 131 near colinearity of the sensors and the target can be easily verified, and even then the
 132 3D location of the target can still be estimated by using the size-ratio of the target in
 133 two panoramic images

$$D_1 = B \frac{w_2 \cos \phi_2}{w_1 \cos \phi_1 + w_2 \cos \phi_2} \cos \phi_1, \quad D_2 = B \frac{w_1 \cos \phi_1}{w_1 \cos \phi_1 + w_2 \cos \phi_2} \cos \phi_2, \quad (2)$$

135 where w_1 and w_2 are the widths of the target in the panoramic image pair. Note that
 136 the cosines in the above equations only give signs since the angles are either 0° or

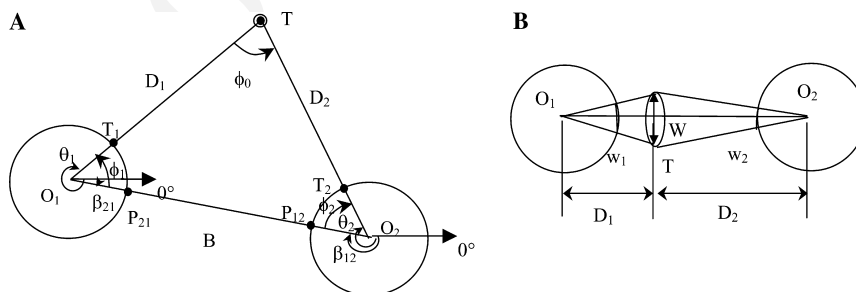


Fig. 2. Two 3D estimation methods. (A) Panoramic triangulation (top view). (B) Panoramic size-ratio method (top view).

137 180°. As an example, if the target lies between O_1 and O_2 (Fig. 2B), the distances to
138 them can be calculated as

$$D_1 = B \frac{w_2}{w_1 + w_2}, \quad D_2 = B \frac{w_1}{w_1 + w_2}. \quad (3)$$

140 In the *size-ratio method*, since the two cameras view the target (e.g., a human)
141 from exactly the opposite direction, the widths of the objects in the two images cor-
142 respond *approximately* to the same width in 3D space (Fig. 2B), which makes the es-
143 timation plausible. As an alternative, we can also use the height information (in the
144 same way as we use width) since the height of an object is more invariant. However,
145 it is only applicable when the top and/or bottom of the figure are visible in both of
146 the panoramic images and can be accurately localized. In contrast, the width infor-
147 mation is easier to extract and more robust since we can integrate the results from
148 different heights of the object. Realizing that the object and the robots may occlude
149 (part of) each other when in a collinear alignment, we will use the width and height
150 information adaptively.

151 In panoramic virtual stereo, where the viewpoint and baseline relation can
152 change, it is interesting to find the best configuration for estimating the distance
153 of a target. For this purpose, first a dynamic mutual calibration approach will be
154 presented in Section 3. Then a detailed numerical analysis of the distance estimation
155 error by the panoramic virtual stereo (with both the triangulation and size-ratio
156 methods) will be given in Section 4, which will lead to useful results for view planning
157 between the two mobile platforms with panoramic cameras.

158 3. Dynamic mutual calibration

159 To estimate the distance of a target, we need to first estimate the baseline and the
160 orientation angles of the two panoramic cameras. In stereo vision, an *epipole* is de-
161 fined as the projection of one camera's center in the other camera's image plane. In a
162 stereo system with normal FOVs, epipoles are usually out of the FOVs in both cam-
163 eras, therefore we must use a third target in the scene for stereo calibration. In con-
164 trast, the panoramic stereo has two "visible epipoles" because the two panoramic
165 cameras can see each other. Here, we propose a special dynamic calibration proce-
166 dure called *mutual calibration* based on the visible epipole property in panoramic ste-
167 reo. Mutual calibration neither needs to setup any additional calibration targets nor
168 requires the use of a third object in the environment. Instead, each of the panoramic
169 cameras can use the other as the calibration target. The advantage of "sensor as the
170 target" in mutual calibration is that the geometric structures and the photometric
171 properties of the sensors as well their platforms can be well designed and are known
172 a priori.

173 Several practical approaches have been proposed for this purpose by using special
174 structures, such as cylinders, vertical lines, and rectangular planar surfaces [24]. The
175 basic idea is to make the detection and calculation robust and fast. One of the ap-
176 proaches is to design the body of each robot as a cylinder with some vivid colors

177 (e.g., white in the intensity images of our current implementation), which can be eas-
 178 ily seen and extracted in the image of the other robot's camera (Fig. 3A). We assume
 179 that the rotation axis of each panoramic camera is coincident with the rotation axis
 180 of the cylindrical body of the corresponding robot, therefore the baseline between
 181 the two panoramic cameras can be estimated using the occluding boundary of either
 182 of the two cylinders, e.g., from the image of camera 2 we have

$$B = R / \sin\left(\frac{\alpha}{2}\right), \quad (4)$$

184 where α is the angle between two occluding projection rays measured in the image of
 185 camera 2, and R is the radius of the 1st cylindrical body (Fig. 3B). The orientation
 186 angle (β_{12}) of the line O_2O_1 is simply the average of the bearings of two occluding
 187 boundary points P_1 and P_2 . We can do the same in the image of camera 1.

188 Fig. 4 shows a calibration result. The cylindrical body of each robot (pointed at
 189 by an arrow in Figs. 4A and B) is detected and measured in the panoramic image of

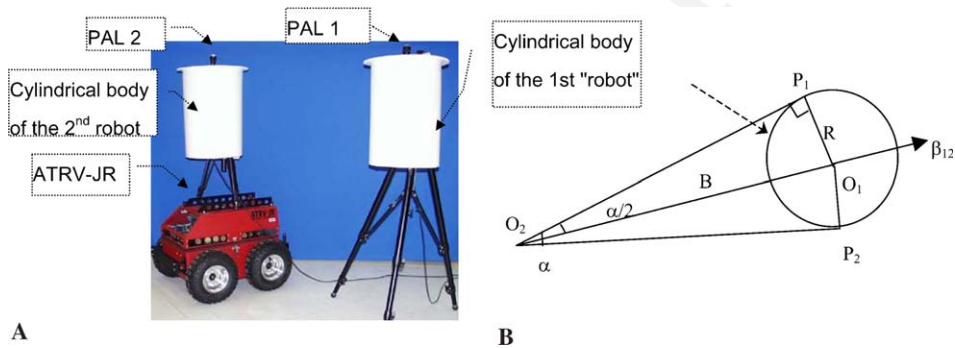


Fig. 3. Finding the orientation and the distance using a cylinder (top view). (A) Setup. (B) Geometry.

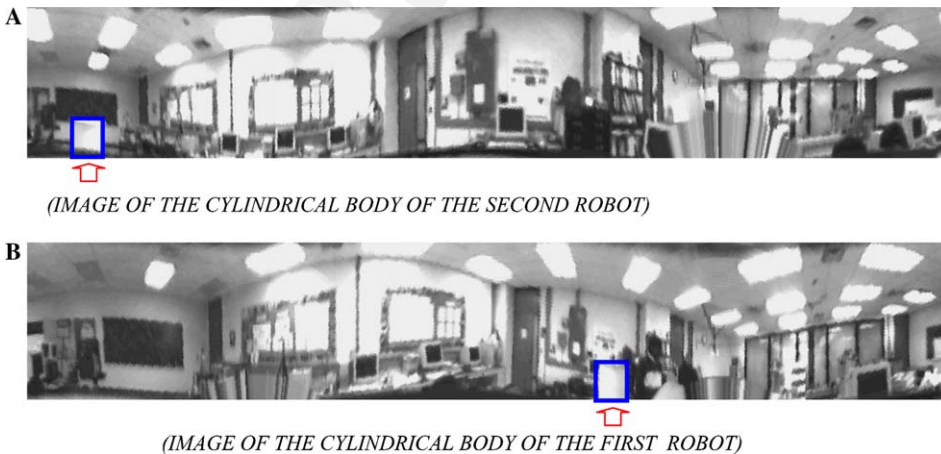


Fig. 4. Dynamic calibration by cylinders (which are pointed by arrows). (A) Pano 1: $F_h = 159.15$ (pixels), $\alpha = 11.52^\circ$ (32 pixels), $\beta_{21} = 23.76^\circ$, $B = 180$ cm. (B) Pano 2: $F_h = 159.15$ (pixels), $\alpha = 11.52^\circ$ (32 pixels), $\beta_{21} = 227.88^\circ$, $B = 180$ cm.

190 the other robot. In the experiment, the perimeter of the cylindrical image is 1000 pix-
 191 els, so the angular resolution in degrees is $360/1000 = 0.36^\circ$ per pixel. We define the
 192 *angular resolution* of the panoramic image as χ in radians for future use, which is
 193 6.28 mrad/pixel in this experiment. The radius of the cylindrical body of each robot
 194 is designed as $R = 18.0$ cm.

195 4. Error analysis

196 Previous work (e.g., [21]) only gave the error distribution of the panoramic stereo
 197 with a fixed stereo geometry. Shum et al. [25] studied the case of an omnidirectional
 198 camera moving within a circular region of the plane and concluded that it was the
 199 best to choose a pair of cameras that are vergent on a point with maximum vergence
 200 angle in order to accurately localize the point. In this paper we will discuss a more
 201 general case where the relations between the two panoramic cameras can change ar-
 202 bitrarily. Our task is to find the distance of a *given* target point from camera 1 by
 203 finding its correspondence in camera 2, so the localization error turns out to be
 204 the distance error. For this reason, we use a different error formulation: for a certain
 205 distance D_1 from camera 1 to the target, what is the error distribution of this distance
 206 with different locations of camera 2, which determines the configurations of baselines
 207 and angles of the panoramic stereo? Can we achieve a better distance estimation for
 208 distant targets with a larger baseline, which is also dynamically determined by the
 209 mutual calibration? Eqs. (1) and (2) show that the accuracy of distance estimation
 210 depends on the accuracy in estimating the baseline and the bearing angles. Here,
 211 we derive an analysis of the error of estimating distance D_1 from the first camera
 212 to the target. First, with the triangulation method, the estimated distance error
 213 can be computed by partial differentials of Eq. (1) as

$$\partial D_1 = \left| \frac{\sin \phi_2}{\sin(\phi_1 + \phi_2)} \right| \partial B + B \left| \frac{\sin \phi_2 \cos(\phi_1 + \phi_2)}{\sin^2(\phi_1 + \phi_2)} \right| \partial \phi_1 + B \left| \frac{\sin \phi_1}{\sin^2(\phi_1 + \phi_2)} \right| \partial \phi_2$$

215 or

$$\partial D_1 = \frac{D_1}{B} \partial B + D_1 |\cot(\phi_1 + \phi_2)| \partial \phi_1 + \frac{D_2}{\sin(\phi_1 + \phi_2)} \partial \phi_2, \quad (5)$$

217 where ∂B is the error in computing the baseline B , and $\partial \phi_1$ and $\partial \phi_2$ are the errors in
 218 estimating the angles ϕ_1 and ϕ_2 from the two panoramic images. Analyzing Eq. (5),
 219 we have found that the distance error comes from three separate error sources:
 220 mutual calibration error, stereo matching error and stereo triangulation error, which
 221 will be discussed below.

222 4.1. Calibration error

223 Dynamic mutual calibration estimates the baseline B , and the bearing angles β_{12}
 224 and β_{21} of the two cameras, all of which are subject to errors in localizing the cali-
 225 bration targets. The error in estimating the baseline by Eq. (4) can be derived as

$$\partial B = \frac{B\sqrt{B^2 - R^2}}{2R} \partial\alpha \leq \frac{B^2}{2R} \partial\alpha, \quad (6)$$

227 where $R \ll B$, and $\partial\alpha$ is the error in estimating the angle α in an image. From Eq. (6)
 228 we can find that the baseline error ∂B is inversely proportional to the dimension of
 229 the cylindrical body for dynamic calibration given the same angle error $\partial\alpha$. On the
 230 other hand, given the radius R and the angle error, the baseline error is roughly
 231 proportional to the square of the baseline itself. The angle error ($\partial\alpha$) is determined by
 232 the errors in localizing the occluding boundaries of the second (or first) cylinder in
 233 the first (or second) panoramic image (Fig. 3). The errors in estimating the bearing
 234 angles β_{21} and β_{12} will introduce errors to the angles ϕ_1 and ϕ_2 of the stereo triangle
 235 (Fig. 2A). Since each bearing angle is the average of the orientations of the two
 236 occluding boundaries, their errors can be roughly modeled as the same as $\partial\alpha$, i.e.,
 237 $\partial\beta_{21} = \partial\beta_{12} = \partial\alpha$. Note that these errors (∂B , $\partial\beta_{21}$, $\partial\beta_{12}$) are derived from the specific
 238 mutual calibration method we are using in this paper. However, the following
 239 general relations hold: larger distance (baseline) between two cameras will introduce
 240 larger errors in estimating the baseline, but the errors in bearing angles are inde-
 241 pendent of the distance as long as the calibration targets can be detected.

242 4.2. Matching error

243 We want to find the distance of a *given* point T_1 in view 1 by finding its corre-
 244 sponding point T_2 in view 2. In this sense, there will be no error in providing the bear-
 245 ing angle θ_1 in view 1, i.e., $\partial\theta_1 = 0$, which implies that the error $\partial\phi_1$ is solely
 246 determined by the error of the angle β_{21} via calibration, i.e., $\partial\phi_1 = \partial\alpha + \partial\theta_1 = \partial\alpha$.
 247 However, the perspective view difference in O_1 and O_2 will introduce a stereo
 248 “matching error” (denoted as $\partial\theta_2$) in θ_2 , the localization of T_1 ’s matching point T_2 ,
 249 which could be a function of the location of the view point O_2 (related to O_1). Thus,
 250 $\partial\phi_2 = \partial\alpha + \partial\theta_2$ is a (complicated) function of the viewpoint location and is generally
 251 larger than $\partial\phi_1$. Generally speaking, the “matching error” is determined by three as-
 252 pects—visibility (the size of a target in the panoramic image), detectability (the con-
 253 trast of the target with the background) and similarity (appearance differences
 254 between images of an object in two widely separated views). In the panoramic virtual
 255 stereo, the sizes and the appearances of a target can suffer from significant perspec-
 256 tive distortion due to widely separated views. The matching error will be directly re-
 257 lated to the primitives we are using for stereo matching.

258 4.3. Triangulation error and overall distance error

259 Now we want to find a numerical result of the following problem: for a certain
 260 distance D_1 from camera 1 to the target, what is the error distribution for different
 261 locations of camera 2, which determines configurations of baselines and angles of
 262 the panoramic stereo? Since it is hard to give a numerical function of the error
 263 $\partial\phi_2$ versus the location O_2 , we will use the same measure error bounds for all the an-
 264 gles, i.e., $\partial\alpha = \partial\phi_1 = \partial\phi_2 \equiv \partial\phi$. We will re-examine this matching error qualitatively

265 after we find the optimal baseline/viewpoints. We decompose the analysis into two
 266 steps. First, by fixing the baseline, we find the optimal angle ϕ_1 . It is equivalent to
 267 finding the optimal position of O_2 on a circle of origin O_1 and radius B (Fig. 5). Sec-
 268 ond, under the optimal angle configuration of all possible baselines, we find the opti-
 269 mal baseline B . An additional consideration is that a human has a size comparable
 270 to the robots, so the distances between a robot and the target should be at least
 271 greater than the dimension of the robot, $2R$. We have the following results by com-
 272 bining the error analysis in the triangulation method (Appendix A) and in the size-
 273 ratio method (Appendix B):

274 **Case 1.** When $B \leq D_1 - 2R$, the best estimation can be achieved when

$$B = 2\sqrt{D_1 R}, \quad \cos \phi_1 = \frac{3BD_1}{2D_1^2 + B^2} \quad (7)$$

276 and the error in the optimal configuration is

$$\partial D_1^+ = \partial D_1^B + \partial D_1^\phi = D_1 \left(\frac{\sqrt{4D_1 R - R^2}}{2R} + \frac{\sqrt{(D_1 - 4R)(D_1 - R)}}{\sqrt{D_1 R}} \right) \partial \phi < 2D_1 \sqrt{\frac{D_1}{R}} \partial \phi, \quad (8)$$

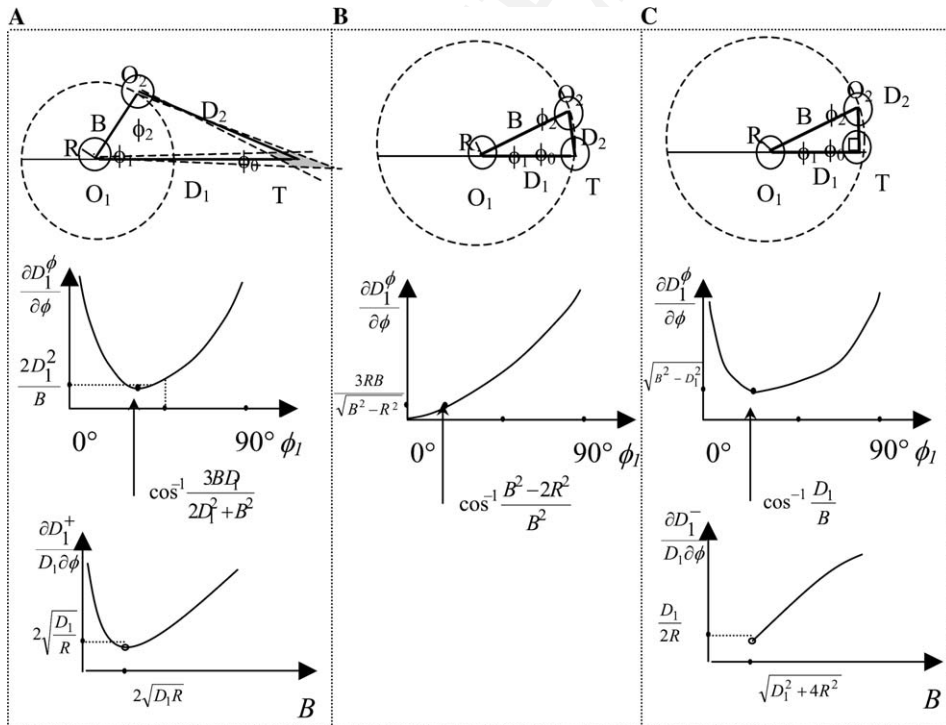


Fig. 5. Best view angles and baselines. (A) $B \leq D_1 - 2R$. (B) $B \in (D_1 - 2R, \sqrt{D_1^2 + 4R^2})$. (C) $B \geq \sqrt{D_1^2 + 4R^2}$.

278 where ∂D_1^B and ∂D_1^ϕ are the distance errors due to the baseline error and angular
 279 errors, respectively. Note that in this case, the minimum error is achieved when
 280 $\phi_1 < 90^\circ$, $\phi_2 > 90^\circ$, and $\phi_0 < 90^\circ$ (see Appendix A). For example, when $R = 0.18$ m,
 281 $D_1 = 4.0$ m, $\partial\phi = 6.28$ mrad (1 pixel), we have the best configuration of $B = 1.70$ m
 282 and $\phi_1 = 54.2^\circ$, and the relative error is $\partial D_1/D_1 = 5.4\%$.

283 **Case 2.** When $B \in (D_1 - 2R, \sqrt{D_1^2 + 4R^2})$, the best estimation can be achieved when

$$B = D_1, \quad \cos \phi_1 = \frac{B^2 - 2R^2}{B^2} \quad (9)$$

285 and the error in the optimal configuration is

$$\partial D_1^0 = D_1 \left(\frac{\sqrt{D_1^2 - R^2}}{2R} + \frac{3R}{\sqrt{D_1^2 - R^2}} \right) \partial\phi. \quad (10)$$

287 Note that in this case, $\phi_1 < 90^\circ$ is the minimum angle by physical constraint of the
 288 minimum object distances, and $\phi_2 = \phi_0$.

289 **Case 3.** When $B \geq \sqrt{D_1^2 + 4R^2}$, the best estimation can be achieved when

$$B = \sqrt{D_1^2 + 4R^2}, \quad \cos \phi_1 = \frac{D_1}{B} \quad (11)$$

291 and the error in the optimal configuration is

$$\partial D_1^- = D_1 \left(\frac{\sqrt{D_1^2 + 3R^2}}{2R} + \frac{2R}{D_1} \right) \partial\phi. \quad (12)$$

293 Note that in this case, the minimum error is achieved when $\phi_1 < 90^\circ$, $\phi_2 < 90^\circ$, and
 294 $\phi_0 = 90^\circ$.

295 **Case 4.** In the case of colinearity of sensors and the target, triangulation is invalid.
 296 However we can use the size-ratio method. A similar error analysis (Appendix B)
 297 shows that if the target lies between the two cameras, minimum error is obtained
 298 when the second camera O_2 moves as close as possible to the target, i.e., $D_2 = 2R$,
 299 or

$$B = D_1 + 2R$$

301 and the minimum error can be expressed by

$$\partial D_1^s = D_1 \left(\frac{\sqrt{(D_1 + R)(D_1 + 3R)}}{2R} + \frac{2R}{W} \right) \partial w. \quad (13)$$

303

304 We always have $\partial D_1^s > \partial D_1^-$ given that $B > D_1$, $\partial w = \partial\phi$, and $W \ll D_1$. Similar re-
 305 sults can be obtained when the target lies in one side of both sensors. It can be also
 306 proved that we always have $\partial D_1^- < \partial D_1^0$, which means that it is better to set the

307 baseline slightly greater than the distance D_1 when they have to be approximately
 308 equal. (In addition, the equality condition cannot be satisfied before we have an
 309 accurate estimation of D_1). By some tedious mathematical comparison of Eqs. (8)
 310 and (12) under different D_1 , we arrive at the following observation:

311 **Conclusion 1.** If the distance from camera 1 (the main camera) to the target is
 312 greater than 11.5 times the radius of the robot i.e., $D_1 > 11.5R$, we have
 313 $\partial D_1^+ < \partial D_1^-$, which means that the best configuration is $B = 2\sqrt{D_1 R}$, $\cos \phi_1 = \frac{3BD_1}{2D_1^2 + B^2}$
 314 (Eq. (7)). Otherwise, we have $\partial D_1^+ \geq \partial D_1^-$ i.e., the best configuration¹ is
 315 $B = \sqrt{D_1^2 + 4R^2}$, $\cos \phi_1 = \frac{D_1}{B}$ (Eq. (11)).

316 It is also interesting to compare the panoramic virtual stereo with a fixed baseline
 317 stereo. Assume that in a fixed baseline stereo system on a robot, the two cameras are
 318 mounted as far apart as possible. For a robot with cylindrical body of radius R , the
 319 maximum stereo baseline in that case would be $B = 2R$. Let us assume that there is
 320 no error in stereo camera calibration (i.e., B is accurate). Since we always have
 321 $B < D_1$ in fixed-baseline stereo, we can use Eq. (A.3) in Appendix A to estimate
 322 the distance error in the *best* case, i.e.,

$$\partial D_1^{\text{fix}}|_{B=2R} \approx \frac{D_1^2}{R} \partial \phi. \quad (14)$$

324 Comparing Eq. (14) with (8), we have the following conclusion.

325 **Conclusion 2.** The flexible baseline triangulation method is almost always more
 326 accurate than a fixed baseline stereo. The error in fixed baseline stereo is propor-
 327 tional to D_1^2 , but the error in flexible baseline stereo is proportional to $D_1^{1.5}$. The error
 328 ratio is

$$\partial D_1^+|_{B=2\sqrt{D_1 R}} : \partial D_1^{\text{fix}}|_{B=2R} = 2\sqrt{\frac{R}{D_1}} \quad (15)$$

330 and $\partial D_1^+|_{B=2\sqrt{D_1 R}} < \partial D_1^{\text{fix}}|_{B=2R}$ when $D_1 > 4R$, which is almost always true.

331 The above error analysis results can be used in the optimal view planning. Though
 332 the exact number 11.5R in Conclusion 1 is deduced from the calibration method we
 333 are using, the guidelines apply to general cases. The distance error map under differ-
 334 ent viewpoints of camera O_2 is given in Fig. 6 for $D_1 = 34R = 6$ m to verify the above
 335 conclusion. Minimum error is $\partial D_1/D_1 \partial \phi = 11.2$ when $B = 220$ cm, $\phi_1 = 62.1^\circ$. (We
 336 have two such symmetric locations for O_2 .) The upper bound of the relative error is
 337 $\partial D_1/D_1 = 7.0\%$ when $\partial \phi$ is equivalent to 1 pixel. The selection of optimal viewing
 338 angle and baseline for different distances is shown in Fig. 7. Note that parameters
 339 in Fig. 7 are slightly different from those in Fig. 6 because the curves in Fig. 7 are
 340 drawn using Eqs. (8) and (12) with some approximation and practical consideration.
 341 The error analysis can also be used in the integration of the results from more than
 342 two such stationary sensors.

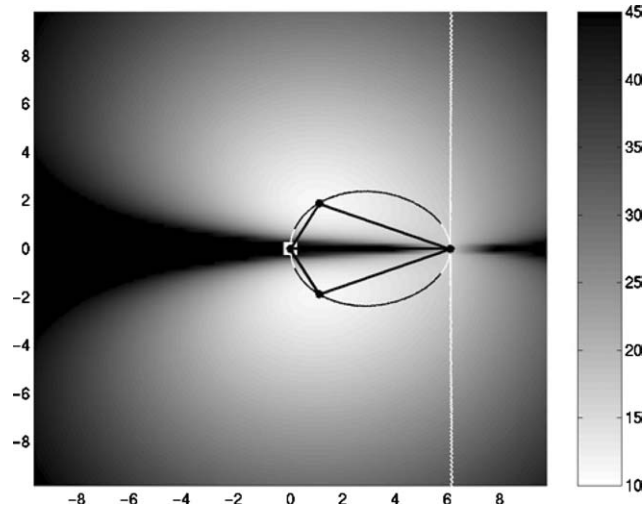


Fig. 6. Error map for distance D_1 when camera O_2 is in different locations of the map by fixing camera O_1 and the target T ($D_1 = 34$, $R = 6$ m, $R = 18$ cm). The labels in the two axes are distances (in meters); the black–white curve shows where the minimum errors can be achieved for viewpoint O_2 on circles with different radii around O_1 (see explanation in the text); the error value ($\partial D_1 / D_1 \partial \phi$) is encoded in intensity: see the corresponding bar.

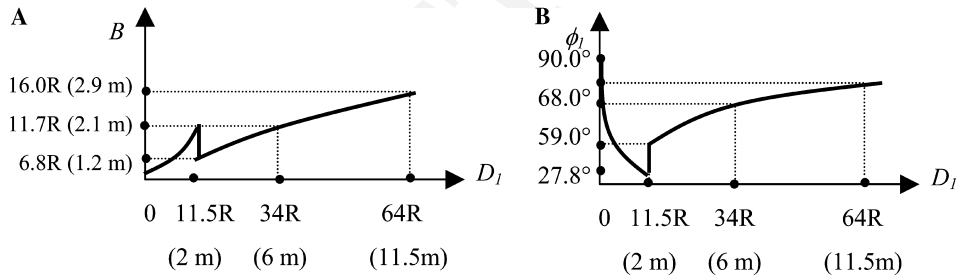


Fig. 7. Best baselines and angles vs. distance curves (the numbers in the parentheses are given when $R = 0.18$ m).

343 5. Matching primitives and matching error revisited

344 Since our primary goal is to detect and to track moving targets (humans) in 3D
 345 space, the primitives of the panoramic virtual stereo are image blobs of human sub-
 346 jects that have already been extracted from the two panoramic images. A fast moving
 347 object extraction and tracking algorithm using motion detection and background
 348 subtraction with a stationary panoramic camera has been developed [24]. Fig. 8 de-
 349 picts the results of our multiple human detection and tracking procedure. Multiple
 350 moving objects (4 people) were detected in real-time while moving around in the scene
 351 in an unconstrained manner; the panoramic sensor is stationary. Each of the four

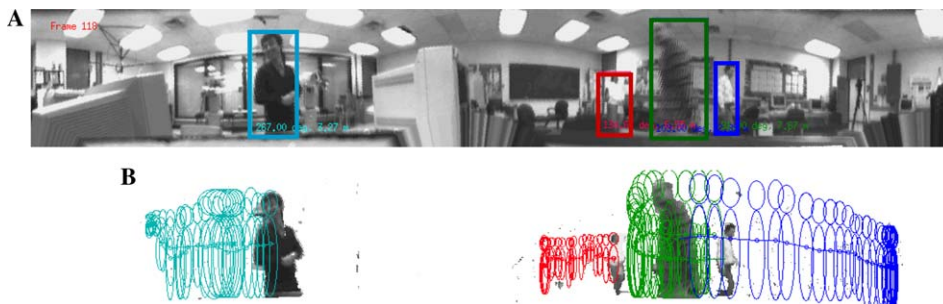


Fig. 8. Tracking multiple moving objects. (A) Cylindrical images with bounding rectangles around moving objects superimposed. (B) Object tracks, each track is for the last 32 frames.

352 people was completely extracted from the complex background, as depicted by the
353 bounding rectangle, direction, and distance of each object. The dynamic track, repre-
354 sented as a small circle, and icon (elliptic head and body) for the last 30 frames of each
355 person is shown in Fig. 8B in different colors. The frame rate for multiple object de-
356 tection and tracking was about 5 Hz in a Pentium 300 MHz PC for 1080×162 pan-
357 oramic images, and thus can be 15–20 Hz with current standard CPUs.

358 We have realized that the bearing of the centroid of an entire blob is subject to the
359 effects of the positions of arms and legs, and the errors in body extraction. We have
360 found that the bearing of the head of a human is more accurate than the entire blob
361 of the human subject for three reasons: (1) it is usually visible in the panoramic im-
362 ages; (2) it is almost symmetric from all directions of the robot's viewpoints; and (3)
363 it is easy to extract from the background (see Figs. 8 and 9). The quasi-symmetry
364 property of a head makes it more suitable for matching across two widely separated
365 views. The idea to select invariant features for stereo matching can be further ex-
366 tended by extracting different parts of a human blob for partial match between
367 two views.

368 The head part of a blob is extracted by using the knowledge that it is the top-
369 most part of the blob and it has roughly a fixed height-width ratio (e.g., 3:2) in a
370 panoramic image. Here, the exact height of the head segment is not critical since
371 we only use the bearing angle of the head for triangulation. Fig. 9 shows the ex-
372 tracted human blobs and heads from a pair of panoramic images. Bearing of the
373 head is more suitable for building up correspondence between a pair of human
374 blobs from two widely separated views because of the aforementioned reasons. No-
375 tice that the centroid of each head region gives correct bearing of the head even if
376 the size and view differences are large between two images of the same human sub-
377 ject. The estimated height is not accurate and not consistent across the correspond-
378 ing image pair. For example, the second human subject in the images shows that
379 the bearing of the head is more accurate than the entire blob, which is an inaccur-
380 ate detection of the human body: the left side is “underestimated” due to the sim-
381 ilarity between the shirt and the door, and the right side is “overestimated” due to
382 its shadow.

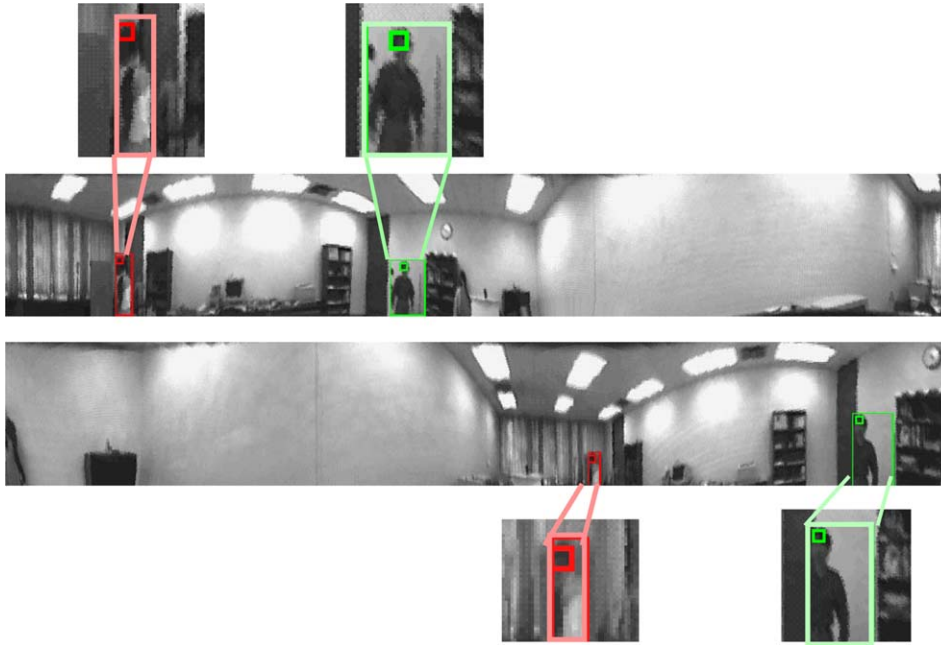


Fig. 9. Head extraction and bearing estimation. The large rectangle around each human subject is the bounding rectangle of the corresponding blob, and the small rectangle inside indicates the centroid of the head.

383 From each panoramic image, a set of objects (blobs) is extracted, which is *anno-*
 384 *tated* by the following parameters

$$\mathbf{T}^{(k)} = \{T_i^{(k)} = (I_i^{(k)}, \theta_i^{(k)}, w_i^{(k)}, h_i^{(k)}), i = 1, \dots, N_k\}, \quad (16)$$

386 where k (1 or 2) is the number of cameras, $I_i^{(k)}, \theta_i^{(k)}, w_i^{(k)}, h_i^{(k)}$ are the photometric
 387 feature, bearing angle of the head of the target i in camera k , the width of the image
 388 blob, and the vertical coordinate of the top of the blob (indicating the height of the
 389 human).

390 The best triangulation configuration is derived when all the angular errors ($\partial\alpha$,
 391 $\partial\phi_1, \partial\phi_2$) are treated as the same, and are assumed to be independent to the view
 392 configuration of the panoramic stereo. However, as we discussed in Section 4.2,
 393 the error $\partial\phi_2$ should be a function of the position of O_2 (given the locations of O_1
 394 and T). A quantitative result can be derived in the same manner as above if the func-
 395 tion is known or can be approximated; but here we only give a qualitative analysis.
 396 The error map in Fig. 6 shows that there is a relatively large region (black part of the
 397 minimum-error curve) with errors that are less than twice the minimum error. The
 398 large errors only occur when angle ϕ_0 is very close to 0° and 180° . Therefore a trade-
 399 off can be made between the matching error (resulting from widely separated views)
 400 and the triangulation error (resulting from small baseline). The target appears sim-
 401 ilar from both the cameras at the best triangulation configuration (in the typical case

402 when $D_1 > 11.5R$), since the distances from the two cameras to the target are com-
 403 parable ($D_2 = \sqrt{D_1^2 - 4D_1R}$). In addition, it is interesting to note that larger view dif-
 404 ference can give a better measurement of the dimension of the 3D object (person),
 405 which is similar to the volume intersection method (Fig. 10).

406 6. Cooperative strategies in the real system

407 In our panoramic virtual stereo vision approach, we face the same problems as in
 408 traditional *motion stereo*: dynamic calibration, feature detection, and matching. In
 409 our scenario, we are also dealing with moving objects before 3D matching, which
 410 seems to add more difficulty. Fortunately, the following cooperative strategies can
 411 be explored between two robots (and their panoramic sensors) to ease these prob-
 412 lems: a “monitor-explore” working mode, mutual awareness, information sharing,
 413 and view planning.

414 6.1. Monitor-explore mode

415 In the two-robot scenario of human searching, one of the robots is assigned as the
 416 “monitor” and the other as the “explorer.” The role of the monitor is to monitor the
 417 movements in the environment, including the motion of the explorer. One of the rea-
 418 sons that we have a monitor is that it is advantageous for it to be stationary while
 419 detecting and extracting moving objects. On the other hand, the role of the explorer
 420 is to follow a moving object of interest and/or find a better viewpoint for construct-
 421 ing the virtual stereo geometry with the camera on the monitor. However, the mo-
 422 tion of the explorer introduces complications in detecting and extracting moving
 423 objects, so we assume that the explorer remains stationary in the beginning of an op-
 424 erational sequence in order to initialize moving objects to be tracked. Then a track-
 425 ing mechanism that can handle ego-motion of the robot continues to track objects of
 426 interest. Such a tracking procedure may integrate the motion, texture and other cues,

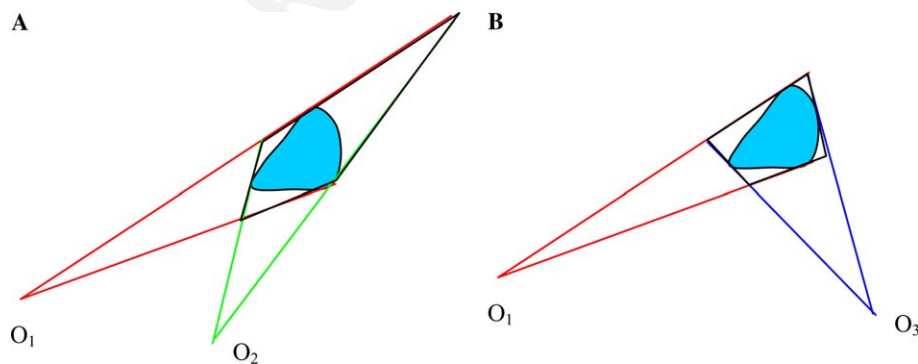


Fig. 10. Viewing differences and distance/dimension estimation. (A) Small viewpoint difference. (B) Large viewpoint difference.

427 which need future work. We also expect that the explorer will remain stationary in an
428 advantageous location after it has found a good viewpoint for 3D estimation. The
429 role of the monitor and the explorer can and will be exchanged during mission exe-
430 cution. The exchange of roles as well as the motion of the explorer may be deter-
431 mined by evaluating expected gain in triangulation accuracy. For example, when
432 the expected improvement is not significant, the robots may just remain in their cur-
433 rent state.

434 6.2. Mutual awareness and information sharing

435 Mutual awareness of the two robots is important for their dynamic calibration of
436 relative orientations and the distance between the two panoramic cameras. In the
437 current implementation, we have designed a cylindrical body with known radius,
438 and color so it is easy for the cooperating robots to detect each other. It is interesting
439 to note that while the motion of the explorer increases the difficulty of tracking other
440 moving objects by itself, tracking information from the monitor is quite useful. It is
441 also possible to use more complicated but known *natural appearances* and geometri-
442 cal models of a pair of robots to implement the mutual awareness and dynamic mu-
443 tual calibration.

444 The two panoramic imaging sensors have almost identical geometric and photo-
445 metric properties. Thus it is possible to share information between them about the
446 targets as well as the robots in the scene. For example, when some number of moving
447 objects are detected and extracted by the stationary monitor, it can pass the informa-
448 tion of the number of objects and their geometric and photometric features of each
449 object to the explorer that may be in motion, thereby increasing robustness of track-
450 ing by the moving explorer. Information sharing is especially useful for the mutual
451 detection of “cooperative” calibration targets since models of the robots are already
452 known a priori. In our simplified case, the cylindrical bodies of both robots always
453 have the same appearances from any viewing angle. Therefore, whenever the monitor
454 has detected the cylindrical body of the moving explorer, it can estimate the bearing,
455 and distance to the explorer. On receiving this information from the monitor, the ex-
456 plorer can try to search for the cylindrical body of the monitor in its image by pre-
457 dicting its size and color under the current configuration and illumination conditions.

458 6.3. View planning for a pair of robots and a single target

459 View planning is applied whenever there are difficulties in object detection and 3D
460 estimation by the virtual stereo system. In our case, we define the view planning as
461 the process of adjusting the viewpoint of the exploring camera so that the best view
462 angle and baseline can be achieved for the monitoring camera to estimate the dis-
463 tance to the target of interest. Occlusion of the human or the robot may occur when
464 an object (either a human or a robot) is between the observing camera and the target,
465 the configuration when triangulation is invalid (we use the size-ratio method in that
466 situation for an initial estimate). The error analysis in Section 4 provides guidelines
467 for “best” viewing planning as follows:

468 (1) *Observation rule*. This rule is applied when the two robots “observe” the target
469 from a distance. If the initial estimated distance from viewpoint O_1 to the target, D_1 ,
470 is greater than $11.5R$, the explorer should move as close as possible to an optimal
471 position that satisfies the minimum distance error conditions, i.e., baseline constraint
472 $B = 2\sqrt{RD_1}$ and the viewing angle constraint $\cos \phi_1 = \frac{3BD_1}{2D_1^2 + B^2}$.

473 (2) *Approaching rule*. This rule is applied when both the two robots are close to the
474 target and the explorer is trying to “approach” the target. If the estimated distance is
475 smaller than $11.5R$, the explorer should approach to the target to satisfy the baseline
476 constraint $B = \sqrt{D_1^2 + 4R^2}$ and the viewing angle constraint $\cos \phi_1 = \frac{D_1}{B}$.

477 (3) *Mutual-awareness rule*. When two panoramic cameras are aware the existence
478 of each other, the maximum distance of the baseline is $B = 2R/w\chi$, given the angular
479 resolution of the panoramic image, χ , the size of the cylindrical robot body, R , and
480 minimum number of detectable pixels of the robots, w .

481 For example, assume that $w = 10$ pixels is the minimum detectable width, then the
482 maximum baseline is $B = 2.8$ m given $R = 0.18$ m and $\chi = 6.28$ mrad/pixel. This con-
483 straint on the baseline still allows the optimal configuration of the panoramic virtual
484 stereo to provide an effective estimation of the distance of a target 10 m away
485 (Fig. 7).

486 (4) *Navigation rule*. The view planning strategy should also consider the cost of
487 moving in finding a navigable path to the selected position. This cost is also a func-
488 tion of distance, smoothness of the path and time to travel.

489 Note that the explorer is always trying to find a best position in the presence of a
490 target’s motion. It is a more difficult problem than localizing a stationary target. For a
491 stationary target, as the robot is moving, the system could collect many stereo esti-
492 mates (of the target) along the way, and integrating them to form a much higher qual-
493 ity estimate of the actual position. However, for a continuously moving object, the
494 integration (if possible) requires formation of the dynamic track of the moving object.
495 This integration could use our derived error model with a Kalman filter methodology.

496 6.4. View planning for multiple robots and multiple objects

497 These strategies can be extended to more than two cooperative robots, and in fact
498 more than two robots will make the work much easier. For example, we can keep
499 two of the three robots in a team stationary so that they can easily detect the moving
500 objects in the scene, including the third robot in motion. Thus, the locations of all the
501 moving objects can be estimated from the pair of stationary panoramic cameras.
502 Then, for a target of interest, we can find (dynamically) the best viewpoint for the
503 third robot in order to estimate the target’s distance from either of the two stationary
504 robots. By using the knowledge of the (dynamic) locations of the target, other mov-
505 ing objects and the three robots, a navigable path for the third robot can be planned
506 to the desirable goal. These measurements can also facilitate the detection of the tar-
507 get and the two stationary robots by the mobile robots, for example, by tracking the
508 objects with visual features inherited from the other two robots. Thus, the stereo tri-
509 angulation relation can be constructed between the moving and the stationary plat-
510 forms.

511 On the other hand, the view planning rules for a single moving object can also be
 512 extended to deal with multiple moving objects. There are three interesting cases.

513 (1) In general, $N + 1$ robots can construct optimal configurations for N moving
 514 objects ($N > 2$), i.e., a main robot can cooperate with each of the N robots for the
 515 detection and localization of each of the N objects (Fig. 11A). However, this method
 516 is inefficient and needs to move the N robots.

517 (2) As a special case (Fig. 11B), two moving robots with panoramic cameras (O_1
 518 and O_2) can construct optimal configurations for estimating the distances of two
 519 moving objects ($T^{(1)}$ and $T^{(2)}$), by the alignment of the two viewpoints of the cameras
 520 to mirror each other.

521 (3) As an approximation method, two moving robots with panoramic cameras
 522 can construct near optimal configurations for estimating the distances of multiple
 523 moving objects. This can be done by clustering the targets into two groups, and
 524 the two cameras then configure two best triangulations for the centers of the two
 525 groups (Fig. 11B). It should be apparent that more than two robots usually can
 526 do a better job in view planning.

527 7. Experimental system and results

528 In our experimental system, we mounted one panoramic annual lens (PAL) cam-
 529 era on an RWI ATRV-Jr. robot (the explorer), and the other PAL camera on a tri-
 530 pod (the monitor) (Fig. 3A). Two Matrox-Meteor frame grabbers, each connected to
 531 a PAL camera were installed on the ATRV-JR and a desktop PC, respectively, at the
 532 time both had 333 MHz PII processors. The communication between two platforms
 533 is through sockets over an Ethernet link (wireless Ethernet communication will be
 534 used in the future system). The 3D moving object detection and estimation programs
 535 run separately on the two machines at about 5 Hz. Only camera and object param-
 536 eter data (i.e., baseline, bearing angles, sizes, and photometric features in Eq. (16))
 537 were transmitted between two platforms so the delay in communication can be ig-
 538 nored at the current processing rate (5 Hz). In the implementation of examples

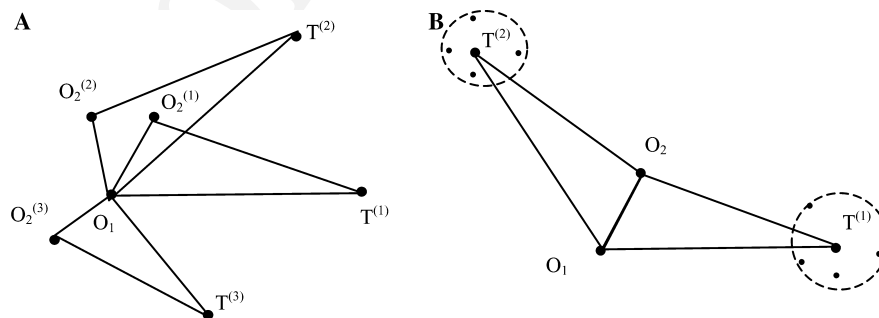


Fig. 11. View planning for multiple robots and multiple objects. (A) N objects, $N + 1$ robots. (B) 2 (groups of) objects, 2 robots.

539 shown in this paper, we assume that the most recent results from both platforms cor-
 540 respond to the events at same time instant. Synchronized image capture is currently
 541 using network time protocol (NTP) with temporal interpolation [26], and it not an
 542 issue with higher speed CPUs.

543 Fig. 12 shows the result from an experiment to evaluate the panoramic stereo's
 544 performance of tracking a single person walking along a known rectangular path
 545 when the two cameras were stationary. Each red dot (dark in B/W) represents a lo-
 546 cation of the person. The dense clusters of dots show the six locations where the per-

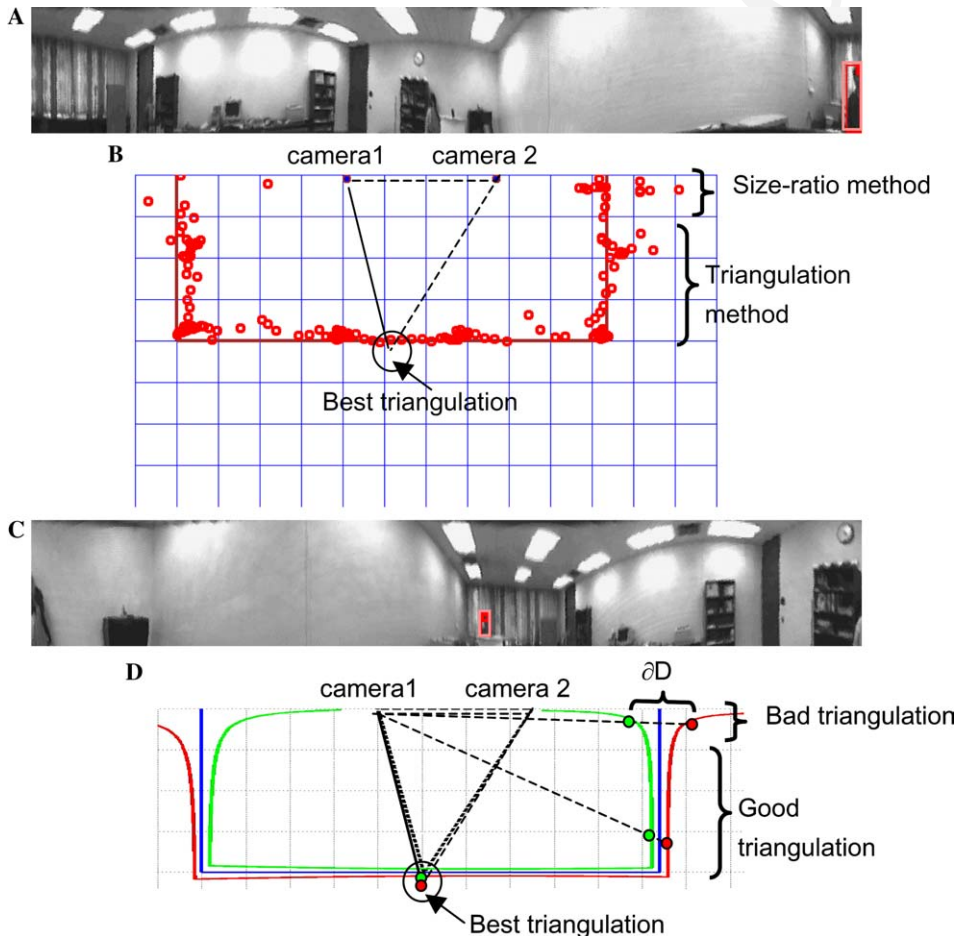


Fig. 12. Panoramic stereo tracking result. The images in (A) and (C) are the panoramic image pair from two panoramic cameras. Each image is actually the corresponding background with the superimposed blob images and their annotations of the blobs. In (C) the real localization results are plotted in the top view of the room where each grid is $50 \times 50 \text{ cm}^2$. Each small circle (red in color version) represents a location of the person in walking. In (D) the theoretical distance error bounds from camera 1 are shown for the same track. The real estimates for best triangulation results validate the theoretical analysis.

547 son made turns during the walking. We used two methods to localize the moving
548 subject—the triangulation method when a good triangle of the target and the two
549 cameras can be formed, and the size-ratio method when the target was near the lo-
550 cations of colinearity. The theoretical error bounds superimposed on the 2D map in
551 Fig. 12D, were computed assuming that all the angular errors in Eqs. (5) and (6)
552 were equivalent to 1 pixel. The target (T) position where the theoretical best triangu-
553 lation on this track can be expected is shown in (Figs. 12B and D), which is consis-
554 tent with the real experimental results. Even if the localization errors in images are
555 larger than the assumed 1 pixel error in our previous analysis, the average error of
556 the estimated track is ± 10 cm, which is comparable to the theoretical error bounds.
557 The bad localizations occurred when the extraction of the human blobs is not cor-
558 rect. For example, in the beginning and the end of the walk, where the size-ratio
559 method is used, large errors occurred not because the size-ratio method is not stable
560 computationally, but because the detection of the human figure is not correct when it
561 is small and occluded by one sensor platform.

562 Fig. 13 shows the results of detecting and tracking two people who walked from
563 the opposite directions along the same known rectangular path. In this example, a
564 simple “greedy” match algorithm [22] was used where the similarities of the match-
565 ing primitives in intensity and the consistencies in 3D measurements are calculated.
566 In the 2D map of the room (center of each picture in Fig. 13), the red (which is dar-
567 ker in B/W print) dot sequence shows the path of one person, and the green (which is
568 lighter in B/W print) dot sequence shows that of the other. The proposed 3D match,
569 localization and tracking algorithms produced rather good results with consistent
570 3D localization for both people. The average localization error is about 20 cm. There
571 are about 5% mis-matches in this set of experiment, which happened in two places.
572 One place is when the shadow of a person was projected on the wall and was detected
573 and mis-matched by the system. The second place of error is when the two people
574 met. Further improvements and experiments on stereo match, view planning and
575 evaluation are needed.

576 8. Concluding remarks

577 This paper has presented a panoramic virtual stereo approach for two (or more)
578 cooperative mobile platforms. There are three main contributions in our approach:
579 (1) a simple but effective dynamic mutual calibration between two panoramic sen-
580 sors; (2) a thorough error analysis for the panoramic virtual stereo vision system;
581 and (3) viewing planning based on optimal stereo configurations. The integration
582 of omnidirectional vision with mutual awareness and dynamic calibration strategies
583 allows intelligent cooperation between visual agents, which provides a nice way to
584 solve problems of limited resources, view planning, occlusions and motion detection
585 of mobile robot platforms. Experiments have shown that this approach is encourag-
586 ing. At the system level, the panoramic virtual stereo is one of the important modules
587 to localize multiple moving human subjects in the distributed sensor network archi-
588 tecture proposed in [26]. In particular, the error modeling and the view planning

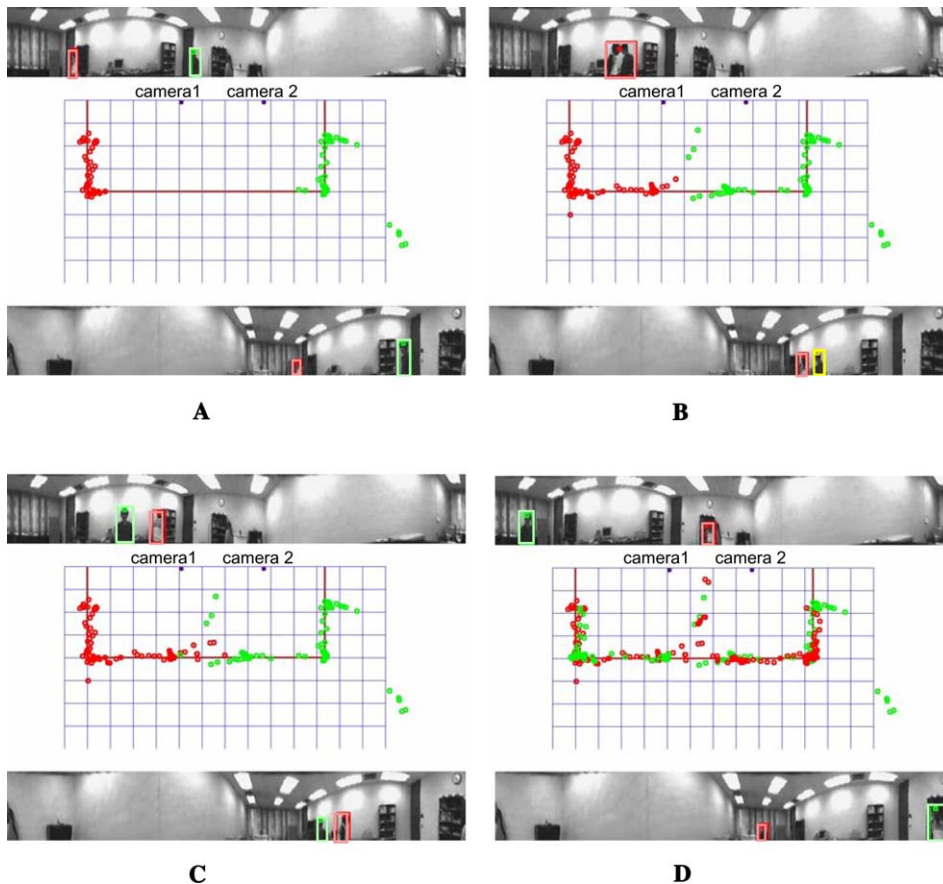


Fig. 13. Panoramic stereo tracking two people. The four pictures show localizing and tracking results (A) before they met, (B) when they met, (C) after they departed, and (D) when they arrived at their goals, out of 214×2 localization results. Each picture of (A)–(D) has the same layout as in Fig. 12. Each small circle in darker tone (red in color version) or in lighter tone (green in color version) represents a location of the corresponding person marked by a bounding rectangle with the same tone (color) in images.

589 strategy developed in this paper are applied. Interesting future work include the fol-
590 lowing:

- 591 (1) *Improvement of the mutual calibration accuracy.* By integrating a panoramic cam-
592 era with a pan/tilt/zoom camera, the system can increase the capability in both
593 viewing angle and image resolution to detect the cooperating robots as well as
594 the targets. Robust and accurate dynamic mutual calibration is one of the key
595 issues in cooperative stereo vision.
- 596 (2) *Improvement of 3D matching.* By using the image contours of objects and more
597 sophisticated features (color, texture, etc), more accurate results can be expected.
598 This is another significant factor that affects the robustness and accuracy of 3D
599 estimation.

600 (3) *Tracking of 3D moving objects.* More sophisticated algorithms for tracking mov-
 601 ing objects should be incorporated, in the presence of occlusion, and by moving
 602 cameras as well as stationary cameras.

603 9. Uncited reference

604 [27].

605 Appendix A. Best baseline and viewpoint when $B < D_1$

606 In the first step, we are trying to find the minimum value of the error due to the
 607 second and third terms of Eq. (5), i.e.,

$$\partial D_1^\phi = D_1 |\cot(\phi_1 + \phi_2)| \partial \phi_1 + \frac{D_2}{\sin(\phi_1 + \phi_2)} \partial \phi_2. \quad (\text{A.1})$$

609 It is equivalent to find the optimal position of O_2 on a circle of origin O_1 and
 610 radius R . We first consider the case, where $B < D_1$. In this case, $(\phi_1 + \phi_2) > 90^\circ$,
 611 so Eq. (A.1) can be re-written as a function of ϕ_1 by using the sine and cosine
 612 laws

$$\partial D_1^\phi = \frac{B^2 + 2D_1^2 - 3BD_1 \cos \phi_1}{B \sin \phi_1} \partial \phi, \quad (\text{A.2})$$

614 where we assume that the same measure errors in angles, i.e., $\partial \phi_1 = \partial \phi_2 = \partial \phi$. By
 615 some mathematical deductions, we can find that the minimum error can be achieved
 616 when $\cos \phi_1 = \frac{3BD_1}{2D_1^2 + B^2}$. The minimum error under the *best configuration* is

$$\partial D_1^\phi |_{\min} = \frac{\sqrt{(D_1^2 - B^2)(4D_1^2 - B^2)}}{B} \partial \phi < \frac{2D_1^2}{B} \partial \phi. \quad (\text{A.3})$$

618 The error in Eq. (A.2) increases from the minimum value to ∞ when the angle ϕ_1
 619 changes from the optimal value to 0° and 180° , respectively (Fig. 6A). Note that in
 620 this case, the minimum error is achieved when $\phi_1 < 90^\circ$, $\phi_2 > 90^\circ$, and $\phi_0 < 90^\circ$.
 621 Here, we compare this result with three special cases (Fig. 14):

622 (1) *Max-vergent configuration.* Two rays O_1T and O_2T have the maximum vergent
 623 angle given the fixed baseline B . In this case $\phi_2 = 90^\circ$, the distance error due to an-
 624 gular errors is

$$\partial D_1^\phi |_{\phi_2=90^\circ} = \frac{2D_1 \sqrt{D_1^2 - B^2}}{B} \partial \phi > \partial D_1^\phi |_{\min}. \quad (\text{A.4})$$

626 (2) *Symmetric configuration.* Two rays O_1T and O_2T have the same length given
 627 the fixed baseline B . In this case $\phi_1 = \phi_2$, the distance error due to angular errors is

$$\partial D_1^\phi |_{\phi_1=\phi_2} = \frac{D_1 \sqrt{4D_1^2 - B^2}}{B} \partial \phi > \partial D_1^\phi |_{\min}. \quad (\text{A.5})$$

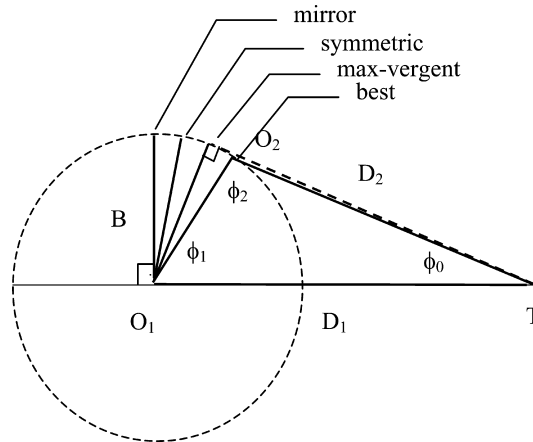


Fig. 14. Best viewpoints given the baseline distance.

629 (3) *Mirror configuration*. The ray O_1T is perpendicular to the baseline B . In this
 630 case $\phi_1 = 90^\circ$, which seems to be a “mirror” of case 1, however, the distance error
 631 is larger

$$\partial D_1^\phi \Big|_{\phi_1=90^\circ} = \frac{2D_1^2 + B^2}{B} \partial\phi > \frac{2D_1^2}{B} \partial\phi. \quad (\text{A.6})$$

633 By a simple comparison we have

$$\partial D_1^\phi \Big|_{\phi_1=90^\circ} > \frac{2D_1^2}{B} \partial\phi > \partial D_1^\phi \Big|_{\phi_1=\phi_2} > \partial D_1^\phi \Big|_{\phi_2=90^\circ} > \partial D_1^\phi \Big|_{\min} \quad (\text{A.7})$$

635 which implies three conclusions: (1) Given the fixed baseline, the distance error in the
 636 max-vergent configuration or the symmetric configuration is slightly larger than the
 637 best configuration. (2) The errors in all the three configurations (max-vergent,
 638 symmetric, and the best) are smaller than $\frac{2D_1^2}{B} \partial\phi$, which is smaller than the mirror
 639 configuration error. (3) Given the fixed baseline, the max-vergent configuration is the
 640 closest to the best configuration.

641 In the second step, we will find the optimal baseline in the case of optimal angle.
 642 Inserting Eqs. (6) and (A.3) into (5) and assuming that the angle error $\partial\alpha$ in Eq. (6)
 643 also equals to $\partial\phi$, we have

$$\partial D_1 = \left(\frac{D_1 \sqrt{B^2 - R^2}}{2R} + \frac{\sqrt{(D_1^2 - B^2)(4D_1^2 - B^2)}}{B} \right) \partial\phi < D_1 \left(\frac{B}{2R} + \frac{2D_1}{B} \right) \partial\phi. \quad (\text{A.8})$$

645 It is intuitive that the larger is the baseline, the better the triangulation will be
 646 (term 2 in Eq. (A.8)), however the estimated error in the baseline is also larger (term
 647 1). The minimum value can be achieved when $B \approx 2\sqrt{D_1 R}$, which means that (1)
 648 more accurate baseline estimation can be obtained given a larger cooperative robotic
 649 target (i.e., R), hence the optimal baseline for estimating distance D_1 can be larger,
 650 and (2) the farther the target is, the larger the baseline should be.

651 Assuming that the human object has a size comparable to the robots, the dis-
 652 tances between a robot and the target should be at least greater than the dimension
 653 of the robot, $2R$. So Eq. (3)–(9) is only valid when $D_2 \geq 2R$, hence we should have
 654 $D_1 \geq B + 2R$. Similarly, we can find the optimal solutions when $B = D_1$ and $B > D_1$.

655 Appendix B. Comparison between triangulation and size-ratio approach

656 The error for the size-ratio method can be calculated in a similar way, For exam-
 657 ple, the distance error for Eq. (3) is

$$\partial D_1 = \frac{D_1}{B} \partial B + \frac{D_1}{w_1 + w_2} \partial w_1 + \frac{B - D_1}{w_1 + w_2} \partial w_2. \quad (\text{B.1})$$

659 Using the mutual calibration error Eq. (6) we have

$$\partial D_1 = D_1 \left(\frac{\sqrt{B^2 - R^2}}{2R} + \frac{B - D_1}{W} \right) \partial w, \quad (\text{B.2})$$

661 where W is the width of the target, and we have $w_1 + w_2 = W \left(\frac{1}{D_1} + \frac{1}{D_2} \right)$. We assume
 662 that $\partial w_1 = \partial w_2 = \partial \alpha = \partial w$, where w is measured in radians. Obviously, we have
 663 $B > D_1$, $D_1 > 2R$ and $D_2 > 2R$. Eq. (B.2) implies that a larger target means better
 664 distance estimation. The minimum error is obtained when the baseline is as large as
 665 possible ($B = D_1 + 2R$), i.e., the second camera O_2 moves as close as possible to the
 666 target ($D_2 = 2R$). So the minimum error can be expressed by

$$\partial D_1^s = D_1 \left(\frac{\sqrt{(D_1 + R)(D_1 + 3R)}}{2R} + \frac{2R}{W} \right) \partial w. \quad (\text{B.3})$$

668 We always have $\partial D_1^s > \partial D_1^-$ given that $B > D_1$, $\partial w = \partial \phi$ and $W \ll D_1$.

669 References

- 670 [1] I. Haritaoglu, D. Harwood and L. Davis, W4S: a real-time system for detection and tracking people
 671 in 2.5D, in: Proc. ECCV, 1998.
 672 [2] A.J. Lipton, H. Fujiyoshi, R.S. Patil, Moving target classification and tracking from real-time video,
 673 in: Proc. DARPA Image Understanding Workshop, vol. 1, November 1998, pp. 129–136.
 674 [3] C. Papageorgiou, T. Evgeniou, T. Poggio, A trainable object detection system, in: Proc. DARPA
 675 Image Understanding Workshop, vol. 2, November 1998, pp. 1019–1024.
 676 [4] A. Pentland, A. Azarbayjani, N. Oliver, M. Brand, Real-time 3-D tracking and classification of
 677 human behavior, in: Proc. DARPA Image Understanding Workshop, vol. 1, May 1997, pp. 193–200.
 678 [5] F.Z. Brill, T.J. Olson, C. Tserng, Event recognition and reliability improvements for the autonomous
 679 video surveillance systems, in: Proc. DARPA Image Understanding Workshop, vol. 1, November
 680 1998, pp. 267–284.
 681 [6] DARPA Image Understanding Workshop Proceedings, VSAM—Video Surveillance and Monitoring
 682 Session, Monterey, November 1998.
 683 [7] S. Baker, S.K. Nayar, A theory of catadioptric image formation, in: Proc. 6th Internat. Conf. on
 684 Computer Vision, IEEE, India, 1998.
 685 [8] V. Nalwa, A true omnidirectional viewer, Technical Report, Bell Lab, Holmdel, NJ, February 1996.

- 686 [9] P. Greguss, Panoramic imaging block for three-dimensional space, U.S. Patent 4,566,763 (28 January
687 1986).
- 688 [10] Y. Yagi, S. Kawato, Panoramic scene analysis with conic projection, in: Proc. IROS, 1990.
- 689 [11] Yamazawa, K., Y. Yagi and M. Yachida, Omnidirectional imaging with hyperboloidal projections, in:
690 Proc. IROS, 1993.
- 691 [12] I. Powell, Panoramic lens, *Appl. Opt.* 33 (31) (1994) 7356–7361.
- 692 [13] Z. Zhu, S. Yang, G. Xu, X. Lin, D. Shi, Fast road classification and orientation estimation using
693 omni-view images and neural networks, *IEEE Trans Image Process.* 7 (8) (1998) 182–197, August.
- 694 [14] Hong J, Tan X, Pinette B, R. Weiss, E.M. Riseman, Image-based homing, in: Proc. Internat. Conf. on
695 Robotics and Automation, April 1991, pp 620–625.
- 696 [15] H. Ishiguro, M. Yamamoto, S. Tsuji, Omnidirectional Stereo, *IEEE Trans. PAMI*, 14(2) 1992, 257–
697 262.
- 698 [16] K.G. Konolige, R.C. Bolles, Extra set of eyes, in: Proc. DARPA Image Understanding Workshop,
699 vol. 1 November 1998, pp. 25–32.
- 700 [17] Kawasaki H, Ikeuchi K. Sakauchi M, Spatio-temporal analysis of omni image, in: CVPR'00, 2000,
701 pp. 577–584.
- 702 [18] T. Boulton, E., R. Micheals, X. Gao, P. Lewis, C. Power, W. Yin, A. Erkan, Frame-rate omnidirectional
703 surveillance and tracking of camouflaged and occluded targets, in: Proceedings of the Second IEEE
704 Workshop on Visual Surveillance, June 1999, pp. 48–58.
- 705 [19] K.C. Ng, H. Ishiguro, M. Trivedi, T. Sogo, Monitoring dynamically changing environments by
706 ubiquitous vision system, in: Proceedings of the Second IEEE Workshop on Visual Surveillance, June
707 1999, pp. 67–73.
- 708 [20] D. Gutchess and A.K. Jain, Automatic Surveillance Using Omnidirectional and Active Cameras, in:
709 Proc. 4th Asian Conf. Comput. Vis., Taipei, January 2000.
- 710 [21] T. Sogo, H. Ishiguro, M.M. Trivedi, N-ocular stereo for real-time human tracking, in: R. Benosman, S.B.
711 Kang (Eds.), *Panoramic Vision: Sensors, Theory and Applications*, Springer-Verlag, New York, 2000.
- 712 [22] Z. Zhu, K.D. Rajasekar, E. Riseman, A. Hanson, Panoramic Virtual Stereo Vision of Cooperative
713 Mobile Robots for localizing 3D Moving Objects, in: Proc. IEEE Workshop on Omnidirectional
714 Vision—OMNIVIS'00, Hilton Head Island, June 2000, pp. 29–36.
- 715 [23] M. Trivedi, K. Huang, I. Mikic, Intelligent Environments and Active Camera Networks, *IEEE*
716 *Systems, Man and Cybernetics*, October 2000.
- 717 [24] Z. Zhu, E.M. Riseman, A.R. Hanson, Geometrical modeling and real-time vision applications of
718 panoramic annular lens (PAL) camera, Technical Report TR #99-11, Computer Science Department,
719 University of Massachusetts Amherst, February, 1999.
- 720 [25] H.-Y. Shum, A. Kalai, S.M. Seitz, Omnivergent stereo, in: Proc. IEEE Seventh Internat. Conference
721 on Comput. Vis., September 1999, pp 22–29.
- 722 [26] D.R. Karuppiiah, Z. Zhu, P. Shenoy, E.M. Riseman, A fault-tolerant distributed vision system
723 architecture for object tracking in a smart room, in: IEEE Second Internat. Workshop on Computer
724 Vision Systems, in: B. Schiele, G. Sagerer (Eds.), Springer Lecture Notes in Computer Science 2095,
725 Vancouver, Canada, July 2001, pp. 201–219.
- 726 [27] Z. Zhu, D.R. Karuppiiah, E.M. Riseman, A.R. Hanson, Adaptive Panoramic Stereo Vision for
727 Human Tracking and Localization with Cooperative Mobile Robots, Accepted by Robotics and
728 Automation Magazine, special issue on panoramic robots.
- 729 [28] G. Cielniak, M. Miladinovic, D. Hammarin, L. Göransson, A. Lilienthal, T. Duckett, Appearance-
730 based tracking of persons with an omnidirectional vision sensor, IEEE Workshop on Omnidirectional
731 Vision (in conjunction with CVPR), June 2003.
- 732 [29] O. Shakernia, R. Vidal and S. Sastry, Structure from small baseline motion with central panoramic
733 cameras, IEEE Workshop on Omnidirectional Vision (in conjunction with CVPR), June 2003.
- 734 [30] G. Adorni, S. Cagnoni, M. Mordonini, A. Sgorbissa, Omnidirectional stereo systems for robot
735 navigation. IEEE Workshop on Omnidirectional Vision (in conjunction with CVPR), June 2003.
- 736 [31] E. Menegatti, A. Scarpa, D. Massarin, E. Ros, E. Pagello, Omnidirectional distributed vision system
737 for a team of heterogeneous robots, IEEE Workshop on Omnidirectional Vision (in conjunction with
738 CVPR), June 2003.