

# Statistical Modeling and Performance Characterization of a Real-Time Dual Camera Surveillance System

Michael Greiffenhager\*, Visvanathan Ramesh\*, Dorin Comaniciu\*, Heinrich Niemann<sup>†</sup>

\*Imaging & Visualization Department  
Siemens Corporate Research, Inc.  
755 College Road East  
Princeton, NJ 08540, USA  
{michael,rameshv,comanici}@scr.siemens.com

<sup>†</sup>Lehrstuhl für Mustererkennung  
Universität Erlangen-Nürnberg  
Martensstr. 3  
91058 Erlangen, Germany  
niemann@informatik.uni-erlangen.de

## Abstract

*The engineering of computer vision systems that meet application specific computational and accuracy requirements is crucial to the deployment of real-life computer vision systems. This paper illustrates how past work on a systematic engineering methodology for vision systems performance characterization can be used to develop a real-time people detection and zooming system to meet given application requirements. We illustrate that by judiciously choosing the system modules and performing a careful analysis of the influence of various tuning parameters on the system it is possible to: perform proper statistical inference, automatically set control parameters and quantify limits of a dual-camera real-time video surveillance system. The goal of the system is to continuously provide a high resolution zoomed-in image of a persons head at any location of the monitored area. An omni-directional camera video is processed to detect people and to precisely control a high resolution foveal camera, which has pan, tilt and zoom capabilities. The pan and tilt parameters of the foveal camera and its uncertainties are shown to be functions of the underlying geometry, lighting conditions, background color/contrast, relative position of the person with respect to both cameras as well as sensor noise and calibration errors. The uncertainty in the estimates is used to adaptively estimate the zoom parameter that guarantees with a user specified probability,  $\alpha$ , that the detected person's face is contained and zoomed within the image<sup>1</sup>.*

## 1 Introduction

Rapid improvement in computing power, cheap sensing and more flexible algorithms are facilitating increased development of real-time video surveillance and monitoring systems [6]. The deployment of video understanding systems in certain critical applications in the real-world can be done only if performance guarantees can be provided for these systems. This paper illustrates the use of systematic engineering methodology outlined in [13] to design and validate a real-time system with given computational and accuracy constraints. We show that by judicious choice of the inter-

mediate transforms (components of the system) along with a careful analysis of the influence of various parameters in the system, it is possible to perform proper statistical inference, automatically set the control parameters and quantify the limits of a dual-camera real-time video surveillance system.

The following section discusses a review of methodologies for analysis and synthesis of vision systems and outline our approach. Subsequent sections describe surveillance system example, statistical modeling and performance analysis, validation, and experimental results.

## 2 Statistical Methodologies for Vision System Design

Past works have addressed methodological issues and have demonstrated performance analysis of components and systems ([5], [4], [8], [16]). However, it is still an art to engineer systems that meet a given application requirement in terms of computational speed as well as accuracy. The trend in the community is to emphasize statistical learning methods, more appropriately Bayesian methods for solving computer vision problems (See for example [9]). However, there still exists the problem of choice of the right statistical likelihood model and right priors that suit an application. Even if this were possible, it is still computationally infeasible to satisfy real-time application needs.

Sequential decomposition of the total task into manageable sub-tasks (with reasonable computational complexity) and the introduction of pruning thresholds is the common way to tackle the problem. This introduces problems because of the difficulty in approximating the probability distributions of observables at the final step of the system so that Bayesian inference is plausible. This approach to perceptual Bayesian inference has been attempted, (see for example [13], [7]). [13]'s work places more emphasis on performance characterization of a system, while [7] attempted Bayesian inference (using Bayesian networks) for visual recognition. The idea of gradual pruning of candidate hypotheses to tame the computational complexity of the estimation/classification problem has been presented in [1]. Note that none of the works identify how the sub-tasks (e.g. feature extraction steps) can be chosen automatically given an application context. There has been prior

<sup>1</sup>Note, the higher the probability  $\alpha$  the more conservative the zoom factor would be. We set  $\alpha$  to 0.99 in our current system.

work to identify system or module configurations that best perform a given task using contextual models of algorithms.

Our approach involves the following key steps (based on [13]):

*System Configuration choice:* The first step is to choose the modules for the system. This is done by use of context (in other words: application specific prior distributions for object geometry, camera geometry and error models, illumination models). Real-time constraints are imposed by choosing pruning methods or indexing functions that restrict the search space for hypotheses. The choice of the pruning functions are derived from the application context and prior knowledge. The choice of the indexing function is not necessarily critical, except that the only criterion that needs to be used is that it be of the form that simplifies computation of the probability of false hypothesis or the probability of missing a true hypotheses as a function of the tuning constants.

*Statistical modeling and Performance Characterization:* The second step involves the derivation of statistical models for errors at various stages in the chosen vision system configuration, so that one can quantify the indexing step and to tune the parameters to achieve a given probability of miss-detection and false alarm rate. In addition, we perform a validation of theoretical models for correctness (through Monte-Carlo simulations) and closeness to reality (through real experiments). For more details on the methodology see [13].

*Hypotheses verification and parameter estimation:* Bayesian estimation is used to evaluate candidate hypotheses and estimate object parameters by using a likelihood model,  $P(\text{measurements}|\text{hypothesis})$ , that takes into account the effects of the pre-processing steps and tuning parameters. In addition, the uncertainty of the estimate is derived in order to predict system performance.

The rest of the paper is organized as follows. Section 3 describes an overview of the surveillance system concept, the mathematics governing the geometry of the people detection and zooming problem, system modules and statistical characterization. Section 4, 5, and 6 describe the validation of the models. Real experiments point out the validity of the models in indoor settings. We conclude in section 7.

### 3 System Description

The task of the two camera surveillance system is to continuously provide zoomed-in high resolution images of the face of a person present in the room. These images represent the input to higher-level vision modules, e.g. face recognition, compaction and event-logging (not discussed in this paper).

#### 3.1 Application Requirements:

The application requirements are as follows: 1) real-time performance on a low-cost PC<sup>2</sup>, 2) person miss-detection rate of  $x_m$ , 3) person false-alarm rate of  $x_f$ , 4) adaptive zooming of person irrespective of background scene structure (with maximal possible zoom based on

<sup>2</sup>Not all system resources in the PC are allocated for visual processing.

uncertainty of person attributes estimated (e.g. location in 3D, height, etc), with performance of the result characterized by face resolution attainable in area of face pixel region (as a function of distance, contrast between background and object, and sensor noise variance and resolution) and bias in the centering of the face. In addition to these requirements, the following assumptions can be made about scene structure: (e.g. the scene illuminant consists of light sources with similar spectrum (e.g. identical light sources in an office area), the number of people to be detected and tracked is bounded, the probability of occlusion of persons (due to other persons) is small.

#### 3.2 System Hardware and Software Configuration:

To continuously monitor the entire scene we use an omnidirectional sensor (OmniCam [11]) mounted below the ceiling. Omni-images are used to detect and estimate the precise location of a given person's foot in the room and this information is used to identify the pan, tilt and zoom settings for a high-resolution foveal camera. Figure 1 shows the system's interface: the overview image, and the high resolution zoomed image of the detected person's face.

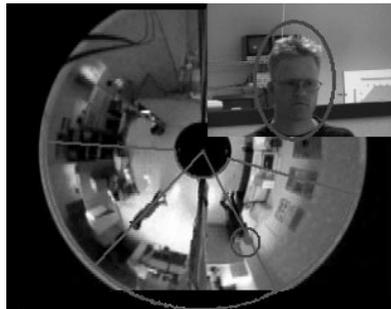


Figure 1: Top: omnidirectional overview image. Red sector: region of interest. Radial lines (green and red) show detected persons. Crosses denote estimated foot/head position. Insert: foveal camera view.

Before we describe the details of how the application requirements translate to the design of individual modules, we discuss the prior distributions (of the 3D scene) reasonable for the given application and identify how these priors induce image priors. The choice of the various estimation steps in the system are motivated from these image priors and real-time requirements. The camera control parameters (pan and tilt) are selected based on the location estimate and its uncertainty (that is derived from statistical analysis of the estimation steps) so as to center the person's location in the image. The zoom parameter is set to maximum value possible so that the camera view still encloses the persons head within the image.

##### 3.2.1 Priors Camera models, Illumination models

The general Bayesian formulation of the person detection and location estimation problem does not suit the real-time constraints imposed by the application. Our approach is to use this formulation only after a pruning step that rules out a majority of false alarms by

designing an indexing step motivated by the 2D image priors (region size, shape, intensity characteristics) induced by the prior distribution in the 3D scene. The prior distributions for person shape parameters: size, height, and his/her 3D location are reasonably simple. These priors on the person model parameters induce 2D spatially variant prior distributions in the projections (e.g. the region parameters for a given person in the image depends on the position in the image) whose form depends on the camera projection model and the 3D object shape.<sup>3</sup> In addition to shape priors, the image intensity/color priors are of importance in our application. Typically we do not assume anything about the object intensity (e.g. homogeneity of object since people can wear variety of clothing and the color spectrum of the light source is not necessarily constrained). However, in the surveillance application, the background is typically assumed to be a static scene (or a slowly time varying scene) with known background statistics (Gaussian mixtures are typically used to approximate these densities). To handle shadowing and illumination changes, these distributions are computed after the calculation of an illumination invariant measure from a local region in an image. The prior distribution of the spectral components of the illuminants in our application are assumed to have same but unknown spectral distribution. Finally, the noise model for the CCD sensor noise is to be specified. This is typically chosen to be i.i.d. zero mean Gaussian noise in each color band.

### 3.2.2 System Software Configuration

The software is composed of four functional modules: calibration, illumination-invariant measure computation at each pixel, indexing functions to select sectors of interest for hypothesis generation, statistical estimation of person parameters (e.g. foot location estimation), and foveal camera control parameter estimation. Figure 2 illustrates the step by step transformations applied to the input. The input color image,  $\hat{R}(x, y), \hat{G}(x, y), \hat{B}(x, y)$ , is transformed ( $T: R^3 \rightarrow R^2$ ) typically to compute an illumination invariant measure  $\hat{r}_c(x, y), \hat{g}_c(x, y)$ . The statistical model for the distribution of the invariant measure is influenced by the sensor noise model and the transformation  $T(\cdot)$ . The invariant measure mean ( $B_o(x, y) = (r_b(x, y), g_b(x, y))$ ) and covariance matrix  $\Sigma_{B_o}(x, y)$ , is computed at each pixel  $(x, y)$  from several samples of  $R(x, y), G(x, y), B(x, y)$  for the reference image of the static scene. A change detection measure  $d^2(x, y)$  image is obtained by computing the Mahalanobis distance between the current image data values  $\hat{r}_c(x, y), \hat{g}_c(x, y)$  and the reference image data  $B_o(x, y)$ . This distance image is used as input to two indexing functions  $\psi_1(\cdot)$  and  $\psi_2(\cdot)$ .  $\psi_1(\cdot)$  discards the radial lines  $\theta$  by choosing hysteresis thresholding parameters that satisfy a given combination of probability of false alarm and miss-detection values, while  $\psi_2(\cdot)$  discards segments along the radial lines in the same manner. The result is a set of regions with high probability of significant change. At this point we employ our full

<sup>3</sup>For this application we found that modeling the person as an upright cylinder is a reasonable approximation.

blown statistical estimation technique that uses the 3D model information, camera geometry information, priors on objects, shape, and 3D location to estimate the number of objects and their positions. The last step is to estimate the control parameters for the foveal camera based on the location estimates and uncertainties.

The following sub-sections describe the system, the statistical models, rationale for indexing steps chosen based on application priors, foot position estimation scheme and control parameter selection scheme. First, we outline the omni-directional/pan-tilt camera geometric relations.

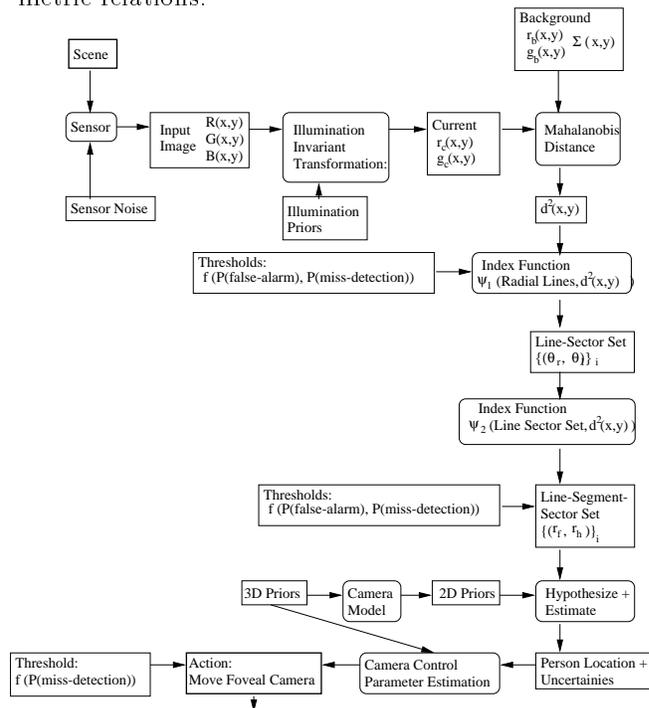


Figure 2: Block diagram: boxes with rounded corners represent transformations while boxes represent data objects.

### 3.2.3 Calibration and Geometry

The projection model for the two cameras are discussed in this section. We denote the geometric model parameters as follows (see Figure 3, 4):

- $H_o$  height of OmniCam above floor (inches)
- $H_f$  height of foveal camera above floor (inches)
- $H_p$  person's height (inches)
- $R_h$  person's head radius (inches)
- $R_f$  person's foot position in world coordinates (inches)
- $D_c$  on floor projected distance between cameras (inches)
- $p(x_c, y_c)$  position of OmniCam center, (in omni-image) (pixel coordinates)
- $r_m$  radius of parabolic mirror (in omni-image) (pixels)
- $r_h$  distance person's head - (in omni-image) (pixels)
- $r_f$  distance person's foot - (in omni-image) (pixels)
- $\vartheta$  - angle between the person and the foveal camera relative to the OmniCam image center (Please see figure 4).



$$\hat{d}^2 = (\hat{\mu}_{\mathbf{b}} - \hat{\mu}_{\mathbf{c}})^T (2\Sigma_{\hat{r}_b, \hat{g}_b})^{-1} (\hat{\mu}_{\mathbf{b}} - \hat{\mu}_{\mathbf{c}}) \quad (7)$$

For background pixels,  $\hat{d}^2$  is approximately  $\chi^2$  distributed with two degrees of freedom. For object pixels  $\hat{d}^2$  happens to be non-central  $\chi^2$  distributed with two degrees of freedom, and non-centrality parameter  $c$ .

### 3.2.5 Indexing for Hypothesis generation

To address real-time computational requirements of the application it is crucial to identify sectorized segments in the image that potentially contains people of interest. To perform this indexing step in a computational efficient manner we define two index functions  $\psi_1()$  and  $\psi_2()$  that are applied sequentially as shown in the system block diagram (2). Essentially  $\psi_1()$  and  $\psi_2()$  are projection operations. For instance, define  $\hat{d}^2(R, \theta)$  as the change detection measure image in polar coordinates with coordinate system origin at the omni-image center  $p(x_c, y_c)$ . Then,  $\psi_1()$  is chosen to be the projection along radial lines to obtain  $\hat{M}_\theta$ , the test statistic that can be used to identify changes along a given direction  $\theta$ . This test statistic is justified by the fact that the object projection is approximated by a line-set (approximated as an ellipse) whose major axis passes through the omni-image center with a given length distribution that is a function of the radial foot position coordinates of the person in the omni-image. This section derives the expressions for the probabilities of false alarm and misdetection at this step as a function of the input distributions for  $\hat{d}^2(R, \theta)$ , the prior distribution for the expected fraction of the pixels along a given radial line belonging to the object, and the noncentrality parameter of  $\hat{d}^2(R, \theta)$  in object locations.

Let  $L_\theta^{x_c, y_c}$  be a radial line through  $p(x_c, y_c)$ , parameterized by angle  $\theta$ , and  $\hat{M}(\theta) = \sum_r d_\theta^2(r)$  denote the accumulative measure of  $d^2$  values at image position  $p(\theta, r)$  parameterized by angle  $\theta$  and distance  $r$  in a polar coordinate system at  $p(x_c, y_c)$ . Applying Canny's hysteresis thresholding technique ([13]) on  $\hat{M}(\theta)$ , provides the sectors of significant change bounded by left and right angles  $\theta_l$  respectively  $\theta_r$ . Let  $r_m$  be the total number of pixels along a radial line  $L_\theta^{x_c, y_c}$ , and  $k$  be the expected number of object pixels along this line (The distribution of  $k$  can be derived from the projection model and the 3D prior models for person height, size, and position described previously, see [15]). The distribution of the cumulative measure is:

$$\text{Background} \quad M_\theta \sim \chi_{2r_m}^2(0) \quad (8)$$

$$\text{Object} \quad M_\theta \sim (r_m - k)\chi_{2(r_m - k)}^2(0) + k\chi_{2k}^2(c) \quad (9)$$

with  $c \in [0 \dots \text{inf}]$ .

To guarantee a false-alarm rate for false sectors of equal or less than  $x_f\%$  we can set the lower threshold  $T_l$  so that

$$\int_0^{T_l} \chi_{\hat{M}_\theta}^2(\xi) d\xi = 1 - x_f\% \quad (10)$$

To guarantee a misdetection rate of equal or less than  $x_m\%$ , theoretically, we can solve for an upper threshold  $T_u$  similarly by evaluating the distribution in eqn. (9). Note that  $k$  is a function of  $H_p$ ,  $R_f$ , and  $c$  (see [15]).

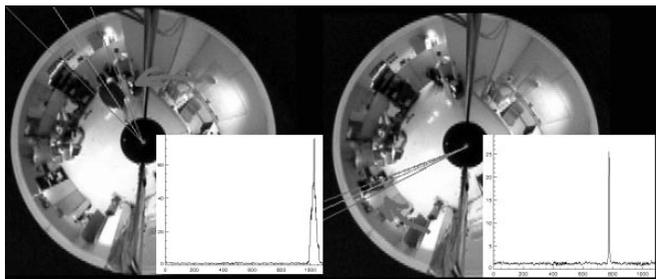


Figure 6: Area of significant change (Left and right lines correspond to  $\theta_l$  and  $\theta_r$ ; Center line denotes the angular position  $\hat{\theta}_f$ .) Inserts show corresponding radial profile  $M_\theta$ .

Therefore we would need to know the distributions of  $H_p$ ,  $R_f$ , and  $c$  to solve for  $T_u$ . Unfortunately, we cannot make any assumptions about the distribution of non-central parameter  $c$ , so we have to resort to the use of a LUT  $T_u(x_m)$  generated by simulations instead.

The second index function  $\psi_2()$  essentially takes as input the domain corresponding to the radial lines of interest and performs an pruning along the radial lines  $R$ . This is done by the computation of  $\hat{d}_{\theta_f}^2(r)$  the integration of the values  $\hat{d}^2()$  along  $\theta_f = \theta + \pi/2$  (within a finite window whose size is determined by the prior density of the minor axis of the ellipse projection), for each point  $r$  on the radial line  $\theta$ . The derivation of the distribution of the test statistic and the choice of the thresholds are exactly similar to the above step.

### 3.2.6 Hypothesis and Estimation Step

We have derived the distributions of the  $\hat{d}^2$  image measurements, and have narrowed our hypotheses for people location and attributes. The next step is to perform the Bayes estimation of person locations and attributes. This step uses the likelihood models  $L(\hat{d}^2|background)$  and  $L(\hat{d}^2|object)$  along with 2D prior models for person attributes induced by 3D object priors  $P(R_p)$ ,  $P(H)$ ,  $P(\theta)$  and  $P(S)$ . In our current application we make use of the fact that the probability of occlusion by persons is small to assert that the probability of a sector containing multiple people is rather small.<sup>4</sup> The center angle  $\theta_f$  of a given sector would in this instance give us the estimate of the major axis of the ellipse corresponding to the person. It is then sufficient to estimate the foot location of person along the radial line corresponding to  $\theta_f$ . The center angle  $\theta_f$  of the sector defines the estimate for the **angular component** of the foot position, see Figure 6. We approximate  $\theta_f$  to be normal distributed with unknown  $\theta_f$  and variance  $\sigma_{\theta_f}$ .  $\theta_f$ 's are estimated as the center positions of the angular sectors given by  $\psi_1()$ . The standard deviation of a given estimate is currently obtained by assuming that the width of the angular sector gives the 99 percentile confidence interval.<sup>5</sup>

Given the line  $\theta_f$  it is necessary to estimate the foot position of the person along this radial line. To find

<sup>4</sup>Estimation when occlusions occur in the OmniImage within a sector is a subject of future research.

<sup>5</sup>We observe that this estimate is an upper bound of the true standard deviation and our experiments have indicated that true variance is a small quantity for a wide range of person positions and background to person contrasts.

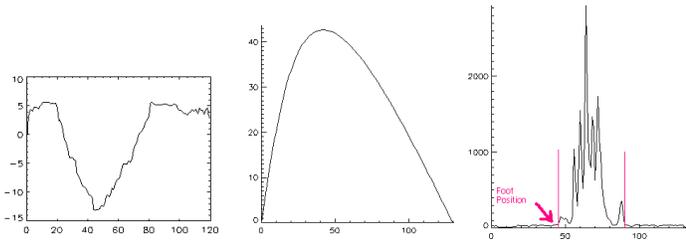


Figure 7: Top left: Bayes error as function of hypothesized foot position  $r_f$ , here: most probable foot position at position  $r_f = 47$ . Bottom left: Projected person length  $k$  as function of  $r_f$ . Note:  $k(r_f = 47) = 43$ . Right: Profile  $\bar{d}_{\theta_f}^2(r)$ : by minimizing Bayes error, responses in interval [43...90] are classified as object responses.

this estimate and variance of the **radial foot position**  $r_f$  we choose the best hypothesis for the foot position that minimizes the Bayes error.<sup>6</sup> Let  $P(h_i|m)$  denote the posterior probability to be maximized, where  $h_i$  denotes the  $i$ th out of multiple foot position hypotheses and  $m$  the measurements ( $\bar{d}_{\theta_f}^2(r)$ ), that are statistically independent; hyper-script  $b$  or  $o$  denotes *background* respectively *object*:

$$\begin{aligned} & P(h_i|m) \\ &= P(h_i^b|m^b)P(h_i^o|m^o) = P(h_i^b|m^b) (1 - P(\bar{h}_i^o|m^o)) \\ &= \frac{p(m^b|h_i^b)P(h_i^b)}{p(m^b)} \frac{p(m^o) - p(m^o|\bar{h}_i^o)P(\bar{h}_i^o)}{p(m^o)} \end{aligned} \quad (11)$$

where  $p$  denotes the density function.  $P(h_i|m)$  becomes maximal for maximal  $p(m^b|h_i^b)$  and minimal  $p(m^o|\bar{h}_i^o)$ , so that

$$r_f = \operatorname{argmax}_{r_f} \log \left( \frac{p(m^b|h_i^b)}{p(m^o|\bar{h}_i^o)} \right) = \quad (12)$$

$$\operatorname{argmax}_{r_f} \left( \sum_{r=0}^{r_f-1} \bar{d}_{\theta_f}^2(r) + \sum_{r=r_h(r_f)}^{r_m} \bar{d}_{\theta_f}^2(r) - \sum_{r=r_f}^{r_h(r_f)-1} \bar{d}_{\theta_f}^2(r) \right)$$

Finally, we estimate the uncertainty in the foot position  $r_f$ . We do not have a close form for that yet, though, our approach provides us with the pdf's up to the latest step in the algorithm. At this point it is affordable to simulate the distribution of  $r_f$  and generate  $\sigma_{r_f}^2$  via perturbation analysis, since only few estimates with known distributions are involved in few operations. Experiments show, that one can approximate  $\hat{r}_f$  as Gaussian distributed with unknown mean  $r_f$ , and variance  $\sigma_{r_f}^2$ .

<sup>6</sup>The prior distribution of person heights and sizes are assumed to be Gaussian. While we actually need to estimate the person height as well as width on the projection using the Bayesian formulation we use the assumption that the variance of the height and size is small and just fix  $H_p$  and  $S_p$  as constants. The geometric transformations are still taken into account to identify 2D projection lengths as a function of radial position along the radial line.

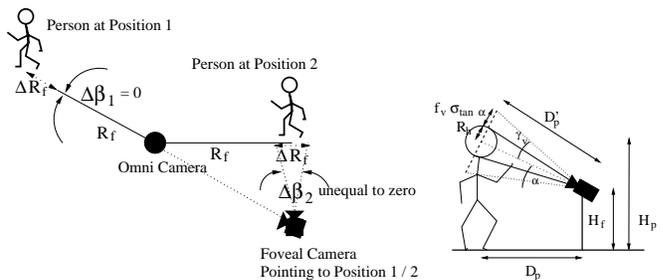


Figure 8: Left: local dependency - same uncertainty in  $R_f$ , different  $\Delta\beta$ . Right: geometric relations for vertical angle of view calculation (view from the side).

### 3.2.7 Foveal Camera Control Parameter Estimation

Once the foot position  $p(\theta_f, r_f)$  is known, we can apply formula 1- 4 to estimate 3D distances  $R_p, D_p$ , and foveal camera control parameter tilt  $\alpha$ , pan  $\beta$  and zoom factor  $z$ .

Figure 8 (left) illustrates how uncertainties in 3D radial distance  $R_p$  influence the foveal camera control parameters.

We have seen that the foot position estimate error can be approximated as a zero mean Gaussian random variate. For the following error propagation steps we will assume that  $\hat{r}_m, \hat{r}_p, \hat{H}_o, \hat{H}_p, \hat{H}_f$ , and  $\hat{D}_c$  are Gaussian random variables with true unknown means  $r_m, r_p, H_o, H_p, R_h, H_f$ , and  $D_c$ , and variances  $\sigma_{r_m}^2, \sigma_{r_p}^2, \sigma_{H_o}^2, \sigma_{H_p}^2, \sigma_{R_h}^2, \sigma_{H_f}^2$ , and  $\sigma_{D_c}^2$  respectively (all estimated in the calibration phase). By applying linearization in the geometric transformations, and making independence assumptions on variables where applicable (see sections 4, 5) it is easy to show (see [15]) how the estimates and its uncertainties propagate through the geometric transformations outlined in section 3. For limited space reason we only print the final results for the uncertainties in tilt  $\alpha$ , and pan  $\beta$ , which were used to calculate the zoom parameter  $z$  in section 3.2.7. (for more details, and derivations of  $\sigma_{R_p}^2, \sigma_{D_p}^2$  see [15]):

$$\begin{aligned} \sigma_{\tan \hat{\alpha}}^2 &= \frac{\sigma_{D_p}^2}{D_p^4} \left( (H_p - R_h - H_f)^2 + \sigma_{H_p}^2 + \sigma_{R_h}^2 + \sigma_{H_f}^2 \right) \\ &+ \frac{\sigma_{H_p}^2 + \sigma_{R_h}^2 + \sigma_{H_f}^2}{D_p^2} \end{aligned} \quad (13)$$

$$\begin{aligned} \sigma_{\sin \hat{\beta}}^2 &= \frac{R_p^2 \sigma_{\hat{\beta}}^2 \cos^2 \vartheta}{D_p^2} + (\sin^2 \vartheta + \sigma_{\hat{\beta}}^2 \cos^2 \vartheta) * \\ &* \left( \frac{R_p^2 \sigma_{D_p}^2}{D_p^4} + \frac{\sigma_{R_p}^2}{D_p^2} + \frac{\sigma_{R_p}^2 \sigma_{D_p}^2}{D_p^4} \right) \end{aligned} \quad (14)$$

Given the uncertainties in the estimates, we can derive the horizontal and vertical angle of view for the foveal camera,  $\gamma_h$  respectively  $\gamma_v$ , which map directly to the zoom parameter  $z$ . Figure 8 (right) shows the geometric relationships for the vertical case. Following

equation provides the vertical angle of view.

$$\gamma_v = 2 \operatorname{atan} \left( \frac{\hat{R}_h + f_v \sigma_{\tan \hat{\alpha}} \hat{D}'_p}{\sqrt{\hat{R}_h^2 + \hat{D}'_p{}^2}} \right) \text{ with } \hat{D}'_p = \frac{\hat{D}_p}{\cos \alpha} \quad (15)$$

where factor  $f_v$  solves for  $\int_0^{\frac{f_v}{2}} N(0, 1) d\xi = \frac{x_z}{2} \%$  given user specified confidence percentile  $x_z$  that the head is display in the foveal frame. Similar derivations apply for the horizontal case.

#### 4 Validation of Assumptions

We verified the correctness of our theoretical expressions and approximations through extensive simulations. Due to lack of space, we only show plots validating expressions for illumination normalization (eqn. 5, 5, Figure 9), and for foveal camera control parameters (eqn. 13, 14, Figure 10). This validation assumes correctness of the underlying statistical models. Validation of the models on real data is done in section 5. In the following we include plots showing theoretically predicted answers and the differences between these predictions and simulated results, based on 10000 samples of normal distributed parameters. For demonstration purpose, parameters that represent the worst-case system behavior were chosen (see [15] for details). The figures show results obtained by using the following settings for the standard deviations:  $\sigma_{\hat{H}_o} = 1\text{cm}$ ,  $\sigma_{\hat{H}_p} = 5\text{cm}$ ,  $\sigma_{\hat{H}_f} = 1\text{cm}$ ,  $\sigma_{\vartheta} = 1^\circ$ ,  $\sigma_{D_c} = 10\text{cm}$ ,  $\sigma = 1$  graylevel. For validation of the distribution of the normalized color values, we used a value for  $B = 50$  while varying  $R$  and  $G$  values in the range 0 through 255 (see Figure 9).

In reality uncertainties are calculated on-line from the current data and are functions of the object, background and location of the object as well as the sensor noise.

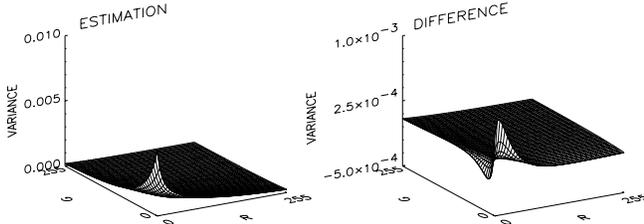


Figure 9: Color normalization: Variance  $\sigma_r^2 + \sigma_g^2 + \sigma_b^2$ . Theoretical values (left) and difference between simulation and theory (right).

Plots show the correctness of the derivations and approximations, give insights of the system limitations depending on user defined tolerances, and show, where the assumptions hold. By examining parametric expressions for uncertainties (eqn. (13),(14), see also [15]) the differences between simulation, and derived prediction can be explained by the linearization error, where the assumption of low signal to noise ratio breaks.

#### 5 Experimental Validation of Models

The correctness of our models is verified by comparing ground truth values against module estimates for mean and variance of the running system. First, we marked eight positions  $P1 - P8$  of different radial distances and pan angles (Figure 11 upper left image).

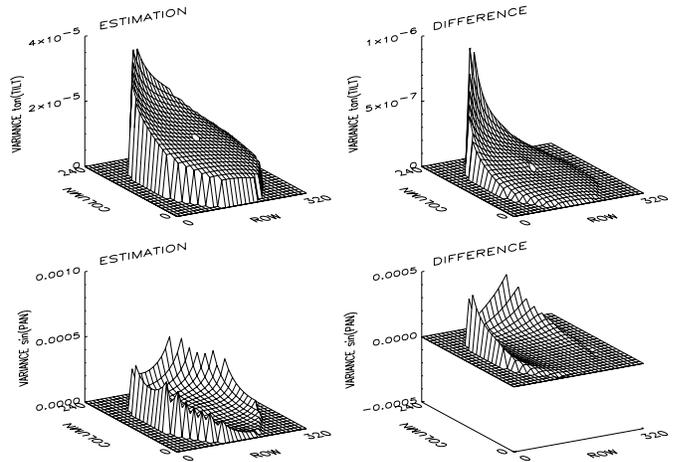


Figure 10: Variances of  $\sin(\hat{\beta})$  and  $\tan(\hat{\alpha})$  plotted as a function of person foot position in omni-image coordinates. Left: Theoretical. Right: Difference between simulation and theory.

Positions, and test persons were chosen to simulate different positions, illumination, and contrast. For limited space reasons we only show the table for the final foveal camera control parameters for one person. Ground truth values for the mean values were taken by measuring tilt angle  $\alpha$ , and pan angle  $\beta$  by hand, and are compared against the corresponding mean of system measurements estimated from 100 trials per position and person. The variances calculated from the system estimates for pan and tilt angle are compared against the average of the corresponding variance-estimates calculated based on the analysis. The comparison between system output and ground truth demonstrates the correctness of the model assumptions in the statistical modeling process (see Table 1).

Table 1: Validation. First two lines show the predicted and experimental variances for the tilt angle, respectively. Next two lines correspond to pan angle.

$\times 10^{-5}$	P1	P2	P3	P4	P5	P6	P7	P8
$\hat{\sigma}_{\tan \hat{\alpha}}^2$	2.10	2.12	1.57	1.40	1.35	1.31	1.31	1.32
$\tilde{\sigma}_{\tan \hat{\alpha}}^2$	2.05	2.04	1.60	1.34	1.36	1.32	1.40	1.31
$\hat{\sigma}_{\sin \hat{\beta}}^2$	28.9	26.1	21.3	17.9	15.3	15.2	18.4	20.1
$\tilde{\sigma}_{\sin \hat{\beta}}^2$	25.9	24.1	19.5	15.1	14.9	15.0	18.1	19.3

#### 6 System Performance

In this section we demonstrate the performance of the running system. Figure 11 demonstrates how the system can precisely track<sup>7</sup> a person's face and zoom, while guaranteeing that the face is in the frame. The output of the foveal camera proved sufficient as input for face recognition algorithms (not discussed in this paper.)

We now illustrate, how the statistical analysis is used to optimize the camera setup. The formulas 13, and 14 suggest that the configuration that minimizes these uncertainties is the one with large inter-camera distance

<sup>7</sup>The tracker uses a simple nearest neighbor prediction.



Figure 11: Top left: omni-image; tracked path, positions at which snapshots were taken. Left right, top down: corresponding foveal image.

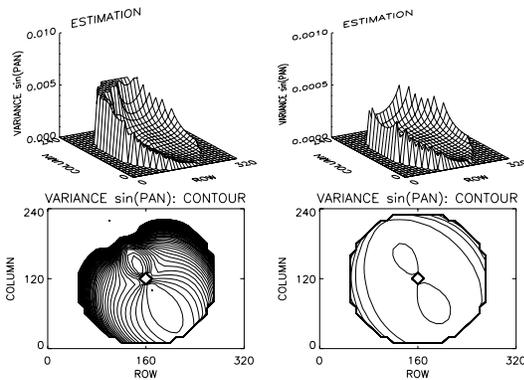


Figure 12: Influence of camera positioning on global and local performance. Top: variance  $\sigma_{\hat{\sin} \beta}$ . Bottom: corresponding contours. Left: close distance  $D_c$  between foveal camera and OmniCam. Right: larger distance, better performance.

$D_c$  and foveal camera height  $H_f$  equal to the mean person eye-level height  $H_p$ . Figure 12 illustrates a comparison of the uncertainties in the pan angle for this setup (right plots), versus a camera position setup with lower distance  $D_c$  (left plots). Similar results are obtained for the tilt angle (omitted for limited space reasons).

The system proved reliable in terms of detection and zooming over longtime experiments within the operational limits denoted by the outer line of the upper right contour plot in Figure 11. Outside, the zooming became imprecise, and did not match the user defined tolerances in terms of zoom precision. This is anticipated, since the assumptions do not hold for these regions as one can see from the plots in section 4.

Figure 12 illustrates how the setup of the system (here placement of foveal camera) influences precision globally and locally. Note preferred directions of low uncertainties (top right to lower left in upper plots). This can be used to adapt the system to user defined accuracy constraints in certain areas of the room.

## 7 Conclusion and Future Work

This paper demonstrated how by careful statistical modeling it is possible to develop and quantify a system to perform a visual surveillance task. The essence of the message is that by careful decomposition of the global task into sub-pieces, statistical characterization of the system, and incorporation of application specific priors in various stages of the system, it is possible to build computationally efficient, but yet statistically well motivated systems. To present this essence it was necessary to make certain simplifying priors and illustrate a working system in a constrained environment<sup>8</sup>. Extensive amount of real and synthetic data experiments were used to validate the models derived. Although we mainly discussed person detection and location estimation alone, the actual video surveillance system has tracking algorithms implemented. We hope to address systematic characterization of the tracking algorithm in the next paper.

## References

- [1] Y. Amit and D. Geman, "A computational model for visual selection", *Neural Computation*, 1999.
- [2] P. Allen and R. Bajcsy, "Two sensors are better than one: example of vision and touch", *Proceedings of 3rd International Symposium on Robotics Research*, pp. 48-55, Gouvieux, France, 1986.
- [3] B. Efron, R. Tibshirani "An Introduction to the Bootstrap", Chapman & Hall, New York, 1993.
- [4] Machine Vision & Applications, International Journal, Special Issue on Performance Evaluation; ed. by W. Forstner, Vol. 9, nos. 5/6, 1997.
- [5] R. Haralick, "Overview: Computer Vision Performance Characterization", *Proceedings of the DARPA Image Understanding Workshop*, Vol. 1, pp.663-665, 1994.
- [6] T. Kanade et al, "Advances in Cooperative Multi-Sensor Video Surveillance", *Proceedings of the DARPA Image Understanding Workshop*, Vol. 1, pp. 3-24, 1998.
- [7] W. Mann and T. Binford, "Probabilities for Bayesian Networks in Vision.", *Proceedings of the ARPA IU Workshop*, 1994, Vol. 1, pp. 633-643.
- [8] K. Cho, P. Meer, J. Cabrera: "Performance assessment through bootstrap", *IEEE Trans. Pattern Anal. Machine Intell.*, 19, 1185-1198, 1997.
- [9] D. Mumford, *Pattern theory: a unifying perspective*, in "Perception as Bayesian Inference", edited by D. Knill and W. Richards, Cambridge Univ. Press, 1996.
- [10] S. Nayar and T. Boulton, "Omnidirectional Vision systems: 1998 PI Report", *Proceedings of the DARPA Image Understanding Workshop*, Vol. 1, pp. 93-100, 1998.
- [11] S. Nayar, "Omnidirectional Video Camera", *Proceedings of the DARPA Image Understanding Workshop*, Vol. 1, pp. 235-242, 1997.
- [12] J. Nielsen "Characterization of Vision Algorithms: An experimental Approach", *ECVnet Workshop on Benchmarking*, 1995.
- [13] V. Ramesh et al, "Computer Vision Performance Characterization," *RADIUS: Image Understanding for Imagery Intelligence*, edited by O. Firschein and T. Strat, Morgan Kaufmann Publishers, San Francisco, 1997.
- [14] D. Stoyan, W.S. Kendall, J. Mecke, *Stochastic Geometry and its Applications*, John Wiley and Sons, 1987.
- [15] M. Greiffenhagen, V. Ramesh, "Auto-Camera-Man: Multi-Sensor Based Real-Time People Detection and Tracking System", Technical Report, Siemens Corporate Research, Princeton, NJ, USA, Nov. 1999.
- [16] S. Wang and T. Binford, "Generic, Model-Based Estimation and Detection of Discontinuities in Image Surfaces," *Proceedings of the ARPA IU Workshop*, 1994, Vol. 2, pp. 1443-1450.
- [17] G. Wysecki, W.S. Stiles "Color Science: Concepts and Methods, Quantitative Data and Formulae", John Wiley & Son, 1982.

<sup>8</sup>Although this is reasonable for the application chosen the challenge of accomplishing system characterization with more complicated priors remains.