

Robust Image Acquisition for Vision-Model Coupling by Humanoid Robots

D. Gonzalez-Aguirre, T. Asfour and R. Dillmann

Karlsruhe Institute of Technology, Adenauerring 2, Karlsruhe-Germany.

{david.gonzalez, asfour, dillmann}@kit.edu

Abstract

The visual perception system of a humanoid robot should attain and manage the vision-model coupling. This essential link between the physical world and its modeled abstraction is established by diverse visual tasks including self-localization, object recognition and detection. The efficiency, robustness and precision of these tasks directly depend on their extracted features. Inevitably, the amount, representativeness and repeatability of these features rely upon the quality and stability of the acquired images. Therefore, a novel method for consistent and stable image acquisition based on image fusion is introduced. This method reliably captures the scene's visual manifold by optimal estimation of the sensor irradiance signals. Experimental evaluation with the humanoid robot ARMAR-IIIa corroborates the quality, stability and applicability of the method.

1 Introduction

The structural composition of the humanoid robots allows them to use the existing *made-for-humans* infrastructure. This crucial fact distinguishes them from other robots permitting their full integration and applicability in the society. However, this composition also imposes severe restrictions on their effectors and sensors. Particularly, the natural approach of visual perception through stereoscopic vision is restricted by:

- *Physical constraints*: These conditions restrain the length of the stereoscopic base-line, the size and weight of the cameras and lenses, see Fig.1-a.
- *Complex perturbations*: Inside a humanoid robot head coexist several devices, see Fig.1-b. Their simultaneous operation produces electric, magnetic and thermal perturbations which deteriorate the quality of the sensor signals, see Fig.1-c and d.
- *Circumscribed resources*: Autonomous humanoid robots should dependably realize complex tasks with their limited memory and computational power onboard.
- *Extensive requirements*: The requirements of the visual sensing skills change dynamically e.g., resolution, frame rate and vergence, see [3].

Even under these conditions, the visual perception of a humanoid robot requires consistent and stable sensor measurements in order to attain the vision-model coupling, see diverse approaches [4]-[7].

In order to overcome these conditions, a multi-image fusion method based on density estimation is introduced in Sec.3. The advantages of the method are its stability and robustness to arbitrary multimodal distributions of the irradiance sensor signals. The resulting fusion images improve the global stability and precision of the visual perception processes, see Sec.4.

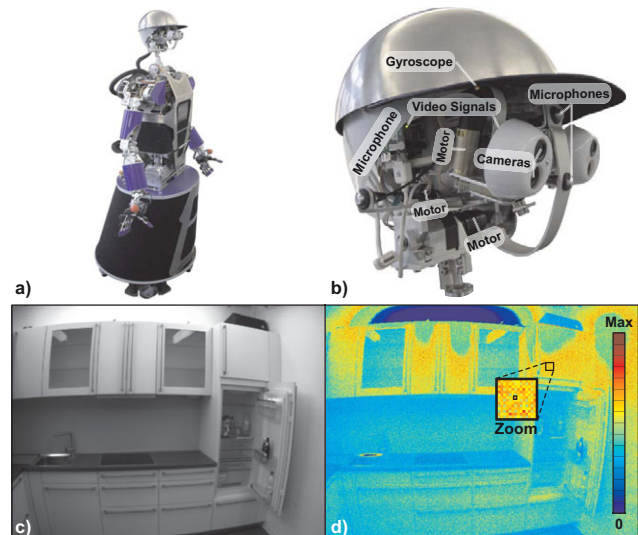


Figure 1. a) The humanoid robot Armar-IIIa [1]. b) The Karlsruhe humanoid head is equipped with a gyroscope, 4 cameras, 6 microphones and 6 motors, see [2]. During task execution, these devices simultaneously work in a *perception-action-cycle*. c) An everyday scene in a made-for-humans environment where the humanoid robot should recognize objects and estimate their poses. d) The pseudo-color deviation image (R_k in Eq.2) shows the sensors instability and detrimental artifacts produced from different noise sources.

The performance of the method allows to acquire semi-dynamic¹ scenes in real applications.

2 Related Work

Methods which improve the image acquisition can be categorically split into rectification through *image enhancement* and synthesis by means of *image fusion*.

2.1 Image Enhancement

These methods deal with the inverse problem of estimating the ideal “*noiseless*” image from a single noise-contaminated image. For more than four decades, image denoising and image enhancement algorithms have achieved considerable progress for image restoration, see extensive survey in [8].

The neighborhood filters such as the k-nearest neighbors and the non-local means filters [9] and [10] (see comparison in [11]) provide considerable results without artifacts when using small window radius. Their outstanding computational performance (up to 500 fps using a **gpu** as reported in [12]) decreases quadratically depending on the applied window radius. Nevertheless, when the noise affects a region beyond one or two pixel(s) the resulting images present severe artifacts. Thus, the robustness of these methods is not adequate for everyday humanoid robot applications.

A more recent method [13] shows outstanding results with highly contaminated images. However, it is extremely expensive, namely, more than 3 minutes even

¹The scene's content remains static during the sampling.

with a low resolution (256×256 pixels) image. Its performance makes it prohibitive for online applications.

Furthermore, all image enhancement methods can only improve the image up to a certain limit. This occurs due to those diminishing factors which cannot be filtered from a single image. It happens in common situations, e.g., the flickering produced by the artificial lighting in an indoor environment, or particularly in humanoid robots, the electromagnetic perturbations (produced by the head motors) generate noise waves on the image, see Fig.2. In such situations, a single image loses the local information and the noiseless-inverse extraction cannot be properly solved.

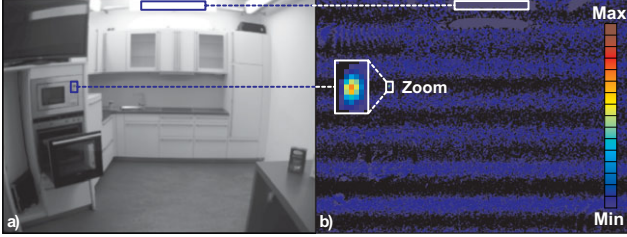


Figure 2. Complex sensing perturbations. a) Everyday scene in a household environment. b) The deviation of a single image (see Sec.3.2) shows three important aspects: i) Noise waves with amplitudes that corrupt (partially vanish or even completely occlude) the underlying saliences such as edges and corners. ii) Overexposed regions (lamp area) show no deviation. iii) The microwave lcd display shows high dynamic deviations. See attached video.

2.2 Image Fusion

These methods exploit the available information in two or more images in order to synthesize an improved image. Image fusion has been widely investigated and successfully applied in various domains [14]. Depending on the data and the desired results, these methods are categorically divided into:

- *Image registration*: Transforms image sets captured from different viewpoints and sensors into a consistent coordinate system, see [15].
- *Super resolution*: Enhances the resolution of the images assuming either structural regularity [16] or statistical similarity of the image [17].
- *High-dynamic-range imaging*: Increases the dynamic range of luminance beyond the sensor capabilities by fusing multi-exposed images, see [18].
- *Multi-focus imaging*: Expands the depth of field by fusing multi-focused images, see [19].
- *Poisson blending*: Composes gradient domain images for the seamless image content insertion [20].
- *Image based rendering*: Uses multi-images to generate novel scene viewpoints, see [21].
- *Stitching*: Combines images with overlapping field of view in order produce panoramas, see [22].

In contrast to these methods, the aim of the method proposed in this article is to improve the quality and stability of the images acquired in each viewpoint using the available sensor resolution. Additionally, the following proposed method holds these considerations:

- No assumptions regarding the image content.
- No requirement of long static scene or robot pose.

- Geometric consistency is well-kept for the Euclidean metric extraction from stereo images.
- Fixed focal length is preserved according to the intrinsic camera calibrations.
- The stability and quality of the acquired images are improved in terms of image processing results, not in human perceptual metrics.

3 Methodology

The first step for robust image acquisition by multi-image fusion is to analyze the *sensor deviation behavior*. Subsequently, based on this behavior, an efficient and robust *fusion strategy* is introduced to overcome both the sensor deficiencies and the unsuitable environmental circumstances. Afterwards, the *convergence analysis* provides a deeper insight into the selection and effects of the *sampling horizon*. Finally, two fundamental feature extraction operations are used to evaluate the stability and precision improvements.

3.1 Sensor Deviation Behavior

The sensor deviation is the amount of intensity variation assuming a semi-dynamic scene. In order to anticipate possible sampling artifacts during the multi-image capture, the maximal available frame rate is used and the images are directly stored for offline processing. Notice that this *indirect mode* is adequate only to establish the sensor deviation behavior, whereas the fusion strategy for online applications (see Sec.3.2) is partially done within the inter-frame interval, see Sec.4

3.1.1 Irradiance Signals

The pixel location $\mathbf{x} \in \mathbb{N}^2$ within the image area Ω is limited by the width w and height h . The Ω set is the domain of the time varying image function $I_t : \mathbb{N}^2 \mapsto \mathbb{N}$. The value associated with a location \mathbf{x} is a random variable independent and identically distributed over the intensity set Θ , namely

$$\mathbf{x} \in \Omega := \{ x \mid (1, 1) \leq x \leq (w, h) \} \subset \mathbb{N}^2,$$

$$I_t(\mathbf{x}) \in \Theta := \{ i \mid 0 \leq i \leq (2^m - 1) \} \subset \mathbb{N},$$

where the temporal subindex t stands for the time stamp and m denotes the bits per pixel.

3.1.2 k-Temporal-Scope

The observation time scope including $k > 1$ images involves the descriptive statistics: maximum \mathbf{U}_k , minimum \mathbf{L}_k , range \mathbf{R}_k and mean \mathbf{M}_k (see Fig.3), formally expressed as

$$\mathbf{U}_k(\mathbf{x}) := \max \left[I_t(\mathbf{x}) \right]_{t=1}^k \quad (1)$$

$$\mathbf{L}_k(\mathbf{x}) := \min \left[I_t(\mathbf{x}) \right]_{t=1}^k \quad (2)$$

$$\mathbf{R}_k(\mathbf{x}) := \mathbf{U}_k(\mathbf{x}) - \mathbf{L}_k(\mathbf{x}) \quad (3)$$

$$\mathbf{M}_k(\mathbf{x}) := \frac{1}{k} \sum_{t=1}^k I_t(\mathbf{x}). \quad (4)$$

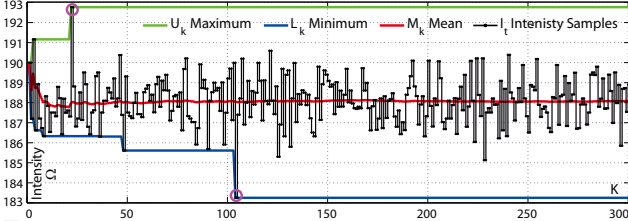


Figure 3. A single pixel's sequential intensity values within the k -temporal-scope show the descriptive statistics. Notice the magenta circles marking the upper and lower intensity outliers. The pixel source location is marked in the center of the zoom in Fig.1-d.

3.2 Fusion Strategy

In order to soundly fuse the observed images, the probability density function (pdf) of each pixel is used to determine the representativeness of the samples. The pdf is attained by kernel density estimation (kde) [23]. This technique provides many advantages compared to the mean or median fusion, see Fig.4.

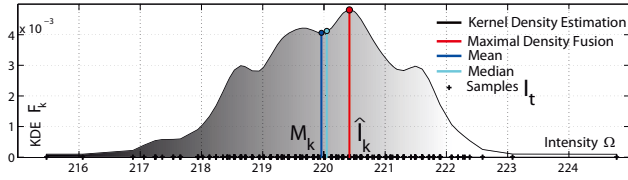


Figure 4. The probabilistic distribution function of the irradiance values illustrates why the straightforward fusion methods such as mean or median produce lower density values. This occurs if the data is not symmetrically spread, skewed or multimodal distributed.

The kde is expressed as

$$f_k(\mathbf{x}, i) := \sum_{t=1}^k K_\lambda \left(i - I_t(\mathbf{x}) \right), \quad (5)$$

where the K_λ denotes the smoothing kernel with bandwidth λ . The Epanechnikov kernel [23] was selected for kde due to its performance, theoretical advantages and the experimental quality of its results. Since the bandwidth assessment is an issue itself [24], the method in [25] was used for its selection. In contrast to ϵ -truncation or kernel dependency as in [24], the kde is efficiently approximated by the next adaptive method.

The most likely fusion value is the maximal density intensity $\hat{I}_k: \mathbb{N}^2 \mapsto \mathbb{R}$ (the top most red value in Fig.4)

$$\hat{I}_k(\mathbf{x}) := \underset{i \in \mathbb{R}}{\operatorname{argmax}} \left[f_k(\mathbf{x}, i) \right]_{i=L_k(\mathbf{x})}^{U_k(\mathbf{x})}. \quad (6)$$

It is semicontinuously computed by a two stage interval analysis: In the first stage, the kde is coarsely sampled with evenly distributed $\alpha < 1 \in \mathbb{R}^+$ increments

$$\hat{I}_k^\alpha(\mathbf{x}) := \underset{i \in \mathbb{R}}{\operatorname{argmax}} \left[f_k(\mathbf{x}, i) \right]_{i=\alpha \cdot j + L_k(\mathbf{x}), j \in \mathbb{N}_0}^{U_k(\mathbf{x})} \quad (7)$$

subsequently, a refinement $\beta < \alpha$ is performed as

$$\hat{I}_k^\beta(\mathbf{x}) := \underset{i \in \mathbb{R}}{\operatorname{argmax}} \left[f_k(\mathbf{x}, i) \right]_{i=\beta \cdot j + \hat{I}_k^\alpha(\mathbf{x}) - \alpha, j \in \mathbb{N}^+}^{\hat{I}_k^\alpha(\mathbf{x}) + \alpha - \beta}. \quad (8)$$

In this manner, the number of iterations is implicitly adjusted according to the observed range. This auto adjustment sagaciously adapts the required computational power while obtaining high accuracy.

3.3 Convergence Behavior

Based on these descriptive statistics, it is clearly noticeable that the stability of the sensor values is reached when the range increment is neglectable. This behavior is globally and smoothly depicted by the mean range expansion

$$\mathbf{E}_k := \frac{1}{wh} \sum_{\mathbf{x} \in \Omega} \mathbf{R}_k(\mathbf{x}), \quad \text{and its rate } \psi_k := \frac{\delta \mathbf{E}_k}{\delta k}. \quad (9)$$

Usually, the range expansion requires a large k -temporal-scope in order to converge, see Fig.5. This occurs due to the outliers in the noisy lower and upper quantiles of the cumulative distribution function (cdf), e.g., the samples at time stamps 20 and 104 in Fig.3.

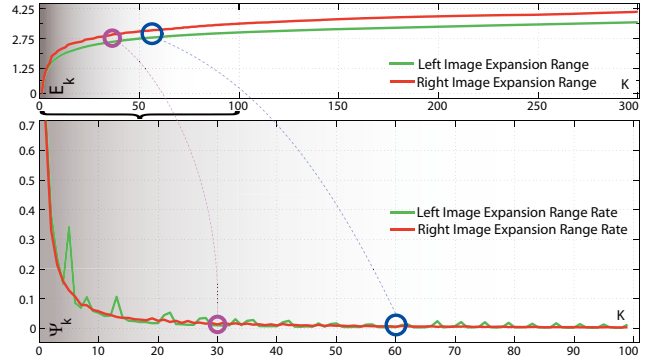


Figure 5. The range expansion E_k (upper plot) and its rate Ψ_k (lower plot) provide an intuitive illustration of the sensor deviation convergence. The range expansion stability is reached at the time position where the curve shows a quasilinear tendency, it occurs approximately after $k > 60$ (blue circles) whereas the rate Ψ_k converges in a shorter (roughly) $k < 30$ temporal-scope (magenta circle).

The range expansion index rate ψ_k from Eq.9 provides an upper limit to estimate the number of images needed for optimal fusion. However, after fusing a certain number of images the resulting synthesized image does not improve substantially. In order to determine the minimal number of images required for the convergence, the following technique is performed: First, a large image set is captured, its cardinality is called n -horizon. Using this set, the fusion strategy is performed through the Eq.7-8. Now, the n -horizon fusion image $\check{I}_n(\mathbf{x}) \in \mathbb{R}$ is regarded as the ground truth reference in order to analyze both the convergence trade-off over the k -temporal-scope and the abnormality.

3.3.1 Convergence Trade-off

It expresses the global convergence versus the k -temporal-scope (see Fig.6) as

$$\chi_k := \left(\frac{1}{wh} \left[\sum_{\mathbf{x} \in \Omega} \left(\hat{I}_t(\mathbf{x}) - \check{I}_n(\mathbf{x}) \right)^2 \right]_{t=1}^k \right)^{\frac{1}{2}}. \quad (10)$$

3.3.2 Abnormality Distribution

The abnormality distribution describes the comprehensive spatiotemporal deviation of the intensity values. It provides an insight into the sensor anomalous distribution. For instance, the intensity value $l \in \Theta$

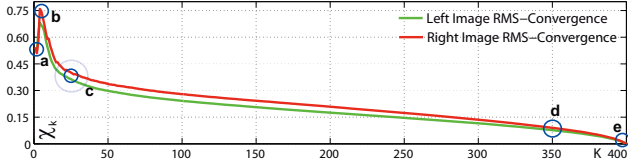


Figure 6. The convergence behavior χ_k from Eq.10. Notice the four regions delimited by circles; a) The χ_2 initial deviation. b) After $k > 4$ samples the maximal deviation is reached. c) Within less than one sampling second ($15 \leq k \leq 25$, @ 30fps) the convergence slows down into a linear behavior. d) Not until a very long $k > 350$ temporal scope (depending on the hardware and scene factors) the deviation convergence behaves nonlinearly. e) Plenary convergence at the sampling $n = 400$ -horizon. Notice that the convergence behavior within c) and d) is quasilinear with a small negative slope.

has the upper abnormality $A_u(l)$ depicting the maximal value found in the whole k -temporal scope and spatial Ω domain which actually corresponds to l in the n -horizon fusion image $\check{I}_n(\mathbf{x})$. Likewise, the lower abnormality A_l , the abnormality range A_r and the rms-abnormality $A_\zeta(l)$ provide the complementary description of the abnormality distribution expressed as

$$A_u(l) := \max \left[I_t(\mathbf{x}) - l \right]_{t=1}^n \quad (11)$$

$$A_l(l) := \min \left[I_t(\mathbf{x}) - l \right]_{t=1}^n \quad (12)$$

$$A_r(l) := A_U(l) - A_L(l) \quad (13)$$

$$A_\zeta(l) := \left[\frac{1}{n} \sum_{t=1}^n \left(I_t(\mathbf{x}) - l \right)^2 \right]^{\frac{1}{2}} \quad (14)$$

all subject to $l = \text{round}(\check{I}_n(\mathbf{x})) : \forall \mathbf{x} \in \Omega$, see Fig.7.

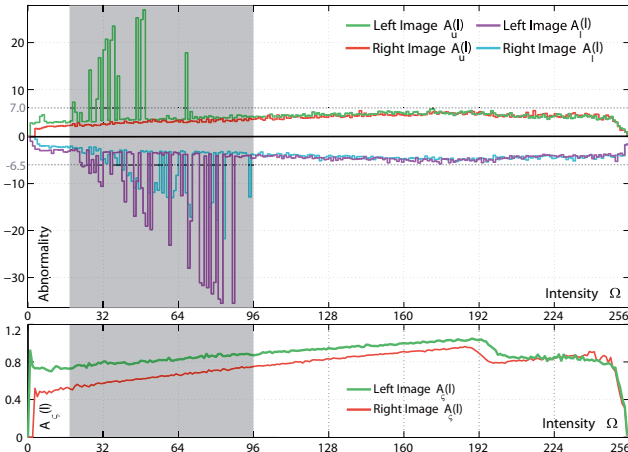


Figure 7. Plots showing the relation between noise distribution and the intensity values. In the upper plot, the fusion abnormality extrema of Fig.1. In the gray region, prominent outliers of the upper $A_u(l)$ and lower $A_l(l)$ abnormality produce detrimental effects for feature extraction and image segmentation. The lower plot shows the rms-abnormality $A_\zeta(l)$ distribution. The global effects of outliers are not globally relevant but locally pernicious for feature extraction.

4 Experimental Evaluation

In order to simultaneously support our claim and evaluate the improvement effects, two important feature extraction tasks were performed.

Edge Stability: Due to the importance of the edge-cue, an evaluation of the stability improvement produced by the multi-image fusion is presented, see Fig.8.

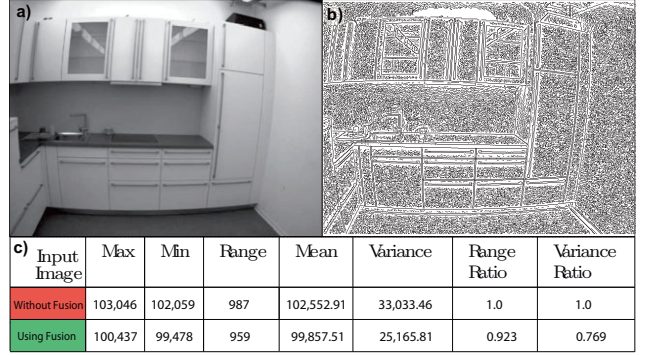


Figure 8. a) Semi-dynamic scene. b) A trial result of the edge detector based on [27]. c) Results show **23.81%** reduction in the variation range of the amount of edge pixels and **7.7%** less variance when using fused images attained by the proposed method.

Segmentation Stability: The noise artifacts from the directly captured image (such as those in the gray region in Fig.7) have negative side effects, e.g., discontinuous segments are connected and vice versa. Hence, when using directly captured images, the segmentation is not consistent in terms of the amount and sizes of the resulting segments. These issues were overcome by the proposed method see results in Fig.9.

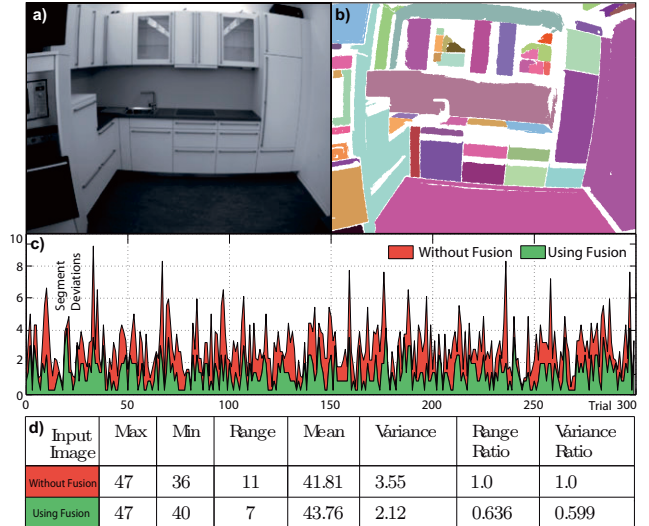


Figure 9. a) Semi-dynamic scene. b) A trial result of the segmentation based on adaptive region growing similar to [26]. c) The absolute variation of the total segments relative to the mean from all trials. d) Results show **36.36%** reduction in the variation range of the amount of segments and **40.05%** less variance using fused images attained by the proposed method.

Performance: During the capturing, the value of each pixel sample is stored and compared to its upper and lower bounds. Later, the fusion from Eq.7-8 is performed. Finally, the data structures are cleaned for the next acquisition, see Tab.1.

Phase (ms)	Max	Min	Mean	Deviation
Capturing	2.25	2.06	2.11	0.04
Fusing (30 Frames)	891.89	712.05	762.94	55.92
Resetting	16.60	12.72	13.31	0.44

Table 1. Running performance of the proposed fusion method. These results were obtained with a CPU Intel(R) Core(TM)2 Quad @ 2.33GHz in a non-optimized single thread implementation.

5 Conclusions

The contribution of this article is the multi-image fusion method for robust visual manifold acquisition for complex robot systems including humanoid robots and mobile service robots. The proposed method conveniently overcomes the hardware limitations and simultaneously aids against the unsuitable environmental conditions. Categorically, the method is a temporal image registration and optimization which clearly improves all visual sensing tasks by providing stable and superior quality images.

The experimental evaluation clearly supports our claim by providing up to 40.05% segmentation improvement in terms of stability compared to the result by using directly captured images. Moreover, the edge extraction improves up to 23.81% in the same manner. This method enhances the image for more representativeness and repeatability of the extracted features allowing a wider applicability of the existing vision methods.

6 Acknowledgments

The research leading to these results has received funding the German Research Foundation (DFG: Deutsche Forschungsgemeinschaft) under the SFB 588 and from the European Union Sixths and Seventh Framework Programme FP7/2007-2013 under grant agreement gree270273.

References

- [1] Asfour T., Regenstern K., Azad P., Schroder J., Bierbaum A., Vahrenkamp N., Dillmann R., "ARMAR-III: An Integrated Humanoid for Sensory-Motor Control", Humanoids 2006. 6th IEEE-RAS Conference on, pp.169-175.
- [2] Asfour T., Welke K., Azad P., Ude A., Dillmann R., "The Karlsruhe Humanoid Head", Humanoids 2008. 8th IEEE-RAS Conference on, pp.447-453.
- [3] Welke K., Asfour T., Dillmann R., "Bayesian Visual Feature Integration with Saccadic Eye Movements", Humanoids 2009. 9th IEEE-RAS Conference on, pp.256-262.
- [4] Okada K., Kojima M., Sagawa Y., Ichino T., Sato K., Inaba M., "Vision Based Behavior Verification System of Humanoid Robot for Daily Environment Tasks", Humanoid 2006. 6th IEEE-RAS Conference on, pp.7-12, 2006.
- [5] Okada K., Kojima M., Tokutsu S., Mori Y., Maki T., Inaba M., "Task Guided Attention Control and Visual Verification in Tea Serving by the Daily Assistive Humanoid HRP2JSK", Intelligent Robots and Systems, IEEE/RSJ Conference on, 2008.
- [6] Gonzalez-Aguirre D., Asfour T., Bayro-Corrochano E., Dillmann R., "Model-Based Visual Self-Localization using Geometry and Graphs", ICPR 2008, pp.1-5, 2008.
- [7] Gonzalez-Aguirre D., Asfour T., R. Dillmann, "Towards Stratified Model-Based Environmental Visual Perception for Humanoid Robots", Pattern Recognition Letters, ISSN 0167-8655. 2010.
- [8] Rao D., Panduranga P, "A Survey on Image Enhancement Techniques: Classical Spatial Filter, Neural Network, Cellular Neural Network, and Fuzzy Filter", Industrial Technology, IEEE Conference on, 2821-2826, 2006.
- [9] Mahmoudi M. Sapiro G., "Fast Image and Video Denoising via Nonlocal Means of Similar Neighborhoods", Signal Processing Letters, IEEE, pp.839-842, 2005.
- [10] Tasdizen, T., "Principal Components for Non-local Means Image Denoising", Image Processing, 15th IEEE Conference on, pp.1728-1731, 2008.
- [11] Buades A., Bartomeu C., Morel J., "On Image Denoising Methods", Technical Note, Centre de Mathematiques et de Leurs Applications, 2004.
- [12] Kharlamov A., Podlozhnyuk V., "Image Denoising", Technical report, NVIDIA Corporation, Santa Clara, CA, 2007.
- [13] Kostadin D., Ro F., Katkovnik V., Egiazarian, K., "BM3D Image Denoising with Shape Adaptive Principal Component Analysis", Workshop on Signal Processing with Adaptive Sparse Structured Representations, 2009.
- [14] Irshad H., Kamran M., Siddiqui A., Hussain A., "Image Fusion Using Computational Intelligence: A Survey", Environmental and Computer Science, 2nd Conference on, pp.128-132, 2009
- [15] Battiato S., Bruna A., Puglisi G., "A Robust Block-Based Image/Video Registration Approach for Mobile Imaging Devices", Multimedia, IEEE Trans. on, pp.622-635, 2010.
- [16] Glasner D., Bagon S., Irami M., "Super-Resolution from a Single Image", Computer Vision, 2009 IEEE 12th International Conference on, pp.349-356, 2009.
- [17] Sung P., Min P., Moon K., "Super-Resolution Image Reconstruction: a Technical Overview", Signal Processing Magazine, IEEE, pp. 21-36, 2003.
- [18] Gonzalez-Aguirre D., Asfour T., Dillmann R., "Eccentricity Edge-Graphs from HDR Images for Object Recognition by Humanoid Robots", Humanoid Robots 2010.
- [19] Cheng-I C., Phen-Lan L., Po-Whei H., "A New Fusion Scheme for Multi-Focus Images Based on Dynamic Salient Weights on Discriminative Edge Points", Machine Learning and Cybernetics, Conference on, pp.351-356, 2010.
- [20] Szeliski R., Uyttendaele M., Steedly D., "Fast Poisson Blending using Multi-Splines", TechReport MSR-TR-2008-58.
- [21] Pessoa S., Moura G., Lima J., Teichrieb V., Kelner J., "Photorealistic rendering for Augmented Reality: A Global Illumination and BRDF Solution", Virtual Reality Conference IEEE, pp.3-10, 2010.
- [22] Linqiang C., Wang X., Liang X., "An Effective Video Stitching Method", Computer Design and Applications, Conference on, pp.297-301, 2010.
- [23] Duda R., Hart P., Stork D., "Pattern Classification", ISBN 978-0-471-05669-0, New York 2001.
- [24] Elgammal A., Duraiswami R., Davis L, "Efficient Kernel Density Estimation using the Fast Gauss Transform with Applications to Color Modeling and Tracking", Pattern Analysis and Machine Intelligence, IEEE Trans. on, pp.1499-1504, 2003.
- [25] Raykar V., Duraiswami R., "Fast Optimal Bandwidth Selection for Kernel Density Estimation", Proc. of the sixth SIAM Conference on Data Mining, pp.524-528, 2006.
- [26] Zhu Q., Hong-e R., "A New Idea for Color Annual Ring Image Segmentation Adaptive Region Growing Algorithm", Information Engineering and Science, Conference on, pp.1-3, 2009.
- [27] Grigorescu C., Petkov N., and Westenberg M., "Contour Detection Based on Nonclassical Receptive Field Inhibition", Image IEEE Transactions on In Image Processing, vol. 12, no. 7, pp. 729739, 2003.