

Intelligent Missions for MAVs: Visual Contexts for Control, Tracking and Recognition

Sinisa Todorovic, Michael C. Nechyba

Department of Electrical and Computer Engineering, University of Florida, Gainesville, Florida 32611-6200

Email: {sinisha, nechyba}@mil.ufl.edu

Abstract—In this paper, we develop a unified vision system for small-scale aircraft that not only addresses basic flight stability and control, but also enables more intelligent missions, such as ground object recognition and moving-object tracking. The proposed system defines a framework for real-time image feature extraction, horizon detection and sky/ground segmentation, and contextual ground object detection. Multiscale Linear Discriminant Analysis (MLDA) defines the first stage of the vision system, and generates a multiscale description of images, incorporating both color and texture through a dynamic representation of image details. This representation is ideally suited for horizon detection and sky/ground segmentation of images, which we accomplish through the probabilistic representation of tree-structured belief networks (TSBN). Specifically, we propose incomplete meta TSBNs (IMTSBN) to accommodate the properties of our MLDA representation and to enhance the descriptive component of these statistical models. In the last stage of the vision processing, we seamlessly extend this probabilistic framework to perform computationally efficient detection and recognition of objects in the segmented ground region, through the idea of visual contexts. By exploiting visual contexts, we can quickly focus on candidate regions where objects of interest may be found, and then perform additional analysis for those regions only. Throughout, our approach is heavily influenced by real-time constraints and robustness to transient video noise.

I. INTRODUCTION

Over the past several years, researchers at the University of Florida have established an active program in developing *Micro Air Vehicles* (MAVs) and small UAVs with maximum dimensions ranging from 5 to 24 inches [1], [2]. This paper substantially advances our previous work on vision-based flight stabilization for these flight vehicles [2], [3], by considering robust stabilization in less benign settings, and in developing a unified vision-based framework for control, tracking and recognition of ground objects to enable intelligent mission profiles.

Imparting MAVs with autonomous and/or intelligent capabilities is a difficult problem, due to their stringent payload requirements and the relative inadequacy of currently available miniature sensors suitable for these small-scale flight vehicles. The one sensor that is absolutely critical for almost any potential MAV mission (e.g. remote surveillance), however, is an on-board video camera. While

an on-board camera is essential for a remote human supervisor, it also offers a wealth of information that can be exploited in vision-based algorithms for guidance and control of MAVs. As such, this paper presents a unified computer vision framework for MAVs that not only addresses basic flight stability and control, but also enables more intelligent missions, such as ground object recognition and moving-object tracking. In this framework, we first seek to extract relevant features from the flight video that will enable higher level goals. Then, we apply this image representation towards horizon detection and sky/ground segmentation for basic flight stability and control. Finally, we extend this basic framework through the idea of *visual contexts* to perform object recognition in the flight video. Throughout, the most important factors that inform our design choices for the vision system are: (1) real-time constraints, (2) robustness to video noise, and (3) complexity of various object appearances in flight images.

Our paper is organized as follows. In Section II, we briefly review *Multiscale Linear Discriminant Analysis (MLDA)*, motivate it as our principal feature representation, and illustrate how it naturally allows us to perform horizon detection. Then, in Section III, we present *incomplete meta tree-structured belief networks (IMTSBNs)* and discuss their appropriateness for the problem of sky/ground segmentation. Next, the context-based object recognition algorithm is explained in Section IV. Finally, in Section V, we demonstrate the performance of the proposed unified computer vision system.

II. MLDA AND HORIZON DETECTION

A. Multiscale Linear Discriminant Analysis

Our choice of feature selection is largely guided by extensive prior experimentation [4], [5], from which we conclude that a prospective feature space must span both color and texture domains. Recently many new image analysis methods, such as wedgelets, ridgelets, beamlets, etc. [6], [7] have been proposed. Aside from the multiscale and localization properties of wavelets [8], these methods exhibit characteristics that account for concepts beyond the wavelet framework, as does our approach to feature extraction and image representation—namely, *Multiscale*

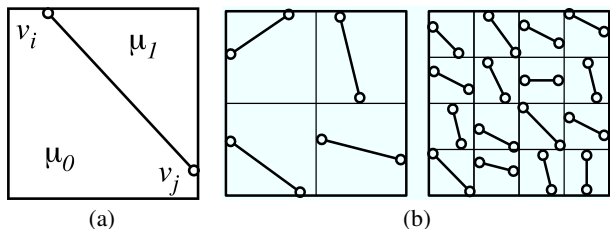


Fig. 1. (a) MLDA atom; (b) MLDA dyadic decomposition.

Linear Discriminant Analysis (MLDA) [9]. Not only does MLDA overcome the shortcomings of wavelets, but it also incorporates color information.

Limited space does not permit a detailed discussion of MLDA; see [9] for a more comprehensive treatment. Here, we illustrate the essentials of MLDA through Figures 1 and 2. In MLDA, we first seek linear discriminants through multiple scales that maximize the Mahalanobis distance J between the RGB color distributions of the two regions on either side of the linear discriminants. Figure 1a illustrates this concept for a square image as a whole, while Figure 1b does so for finer scales. Given this decomposition, individual MLDA atoms (e.g. discriminants) can naturally be represented as a tree \mathcal{T} . However, the distance J will obviously differ for different dyadic squares in the MLDA representation. For example, dyadic squares with relatively uniform color distributions or at finer scales will exhibit smaller distances J between color distributions. Since these discriminants are less useful features, we control the complexity of the MLDA tree \mathcal{T} by pruning those MLDA atoms (e.g. discriminants) or whole subtrees from the complete tree \mathcal{T} . Thus, our final image representation leads to an *incomplete* tree, as depicted in Figure 2.

Examples of MLDA-represented images can be found throughout this paper. Note that MLDA implicitly encodes information about spacial frequencies in the image (i.e. texture) through the process of tree pruning, and that the MLDA tree can easily be examined for spacial interrelationships of its linear discriminants, such as connectedness, collinearity, and other properties of curves in images. Therefore, MLDA is appropriate for computer-vision tasks where both color and texture are critical features.

B. Horizon detection through MLDA

The MLDA framework presents an efficient solution to the horizon-detection problem, incorporating our basic assumptions that the horizon can be interpolated by piecewise linear approximations (i.e., linear discriminants), and that the horizon separates the image into two homogeneous regions. In benign flight video, the horizon-detection problem can be adequately solved by computing only the root atom of the MLDA tree, where the linear discriminant of the root is the optimal solution for the horizon estimate, as illustrated in Figure 3; in fact, the

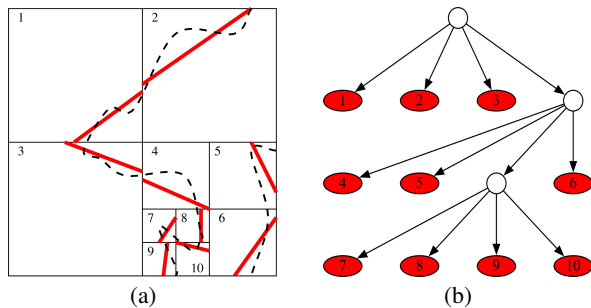


Fig. 2. (a) MLDA graph: the dashed line depicts the actual curve; (b) corresponding MLDA tree: ellipsoid nodes represent the leaf MLDA atoms.

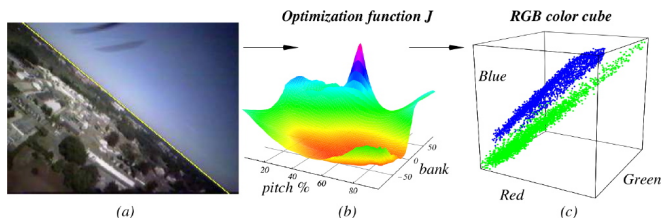


Fig. 3. (a) original image; (b) optimization criterion J as a function of bank angle and pitch percentage, that is, vertices along the perimeter of the image; (c) resulting classification of sky and ground pixels in RGB space.

work presented in [2], [3] is equivalent to this special case. There are, however, unfavorable image conditions (e.g., when the horizon is not a straight line and/or an image is corrupted by video noise), when it is necessary to examine discriminants at finer MLDA scales. For example, in Figure 4, we show that the discriminant of the root MLDA atom does not coincide with the true horizon due to video noise. Expanding the MLDA tree corresponds to image filtering and leads to more accurate positions for the linear discriminants, which then give more accurate evidence of the true horizon’s location in the image. Also, in Figure 5, we show an example, where MLDA detects the true horizon more correctly, as compared to a single-line discriminant. Additional examples and videos can be found at <http://mil.ufl.edu/~nechyba/mav>.

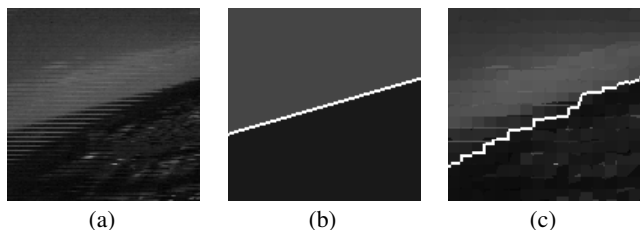


Fig. 4. MLDA efficiently filters video noise: (a) noise-degraded original image; (b) MLDA root atom with the discriminant not equal to the true horizon; (c) MLDA atoms at finer scales as clues for the true horizon position.

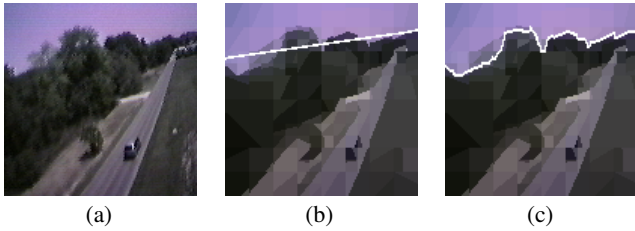


Fig. 5. Non-straight-line horizon: (a) original image; (b) the root MLDA atom with the discriminant not equal to the true horizon; (c) MLDA atoms at finer scales as clues for the true horizon position.

III. SKY/GROUND SEGMENTATION

The horizon detection algorithm determines only the line separating sky and ground regions, which proves sufficient for vision-based flight control under benign flight conditions [2], [3]. However, for complex mission scenarios, correct identification of the sky and ground regions, rather than just the line separating them, takes on increased importance. To accomplish reliable sky/ground image segmentation, we resort to a Bayesian framework, choosing tree-structured belief networks (TSBNs) as the underlying statistical models to describe the sky and ground classes. There are several reasons for this choice of model. Successful sky/ground segmentation implies both accurate classification of large regions, as well as distinct detection of the corresponding boundaries between them. To jointly achieve these competing goals, both large and small scale neighborhoods should be analyzed; this can be achieved through the multiscale structure of TSBNs that naturally arises from our multiscale image representation (i.e., MLDA). Further, it is necessary to employ a powerful statistical model that is able to account for enormous variations in sky and ground appearances, due to video noise, lighting, weather, and landscape variability. Finally, prior work presented in the computer vision literature (e.g. [10]) clearly suggests that TSBNs possess sufficient expressiveness for our goals.

A tree-structured belief network (TSBN) is a generative model comprising hidden and observable random variables (RVs) organized in a tree structure \mathcal{T} . In supervised learning problems, such as our formulation of the sky/ground segmentation problem, an observable RV y represents an

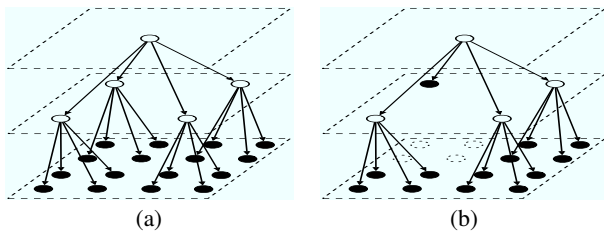


Fig. 6. Differences in TSBN models: (a) “standard” complete TSBN; (b) incomplete TSBN: leaf nodes (depicted black) may occur at coarser resolutions due to MLDA-tree pruning.

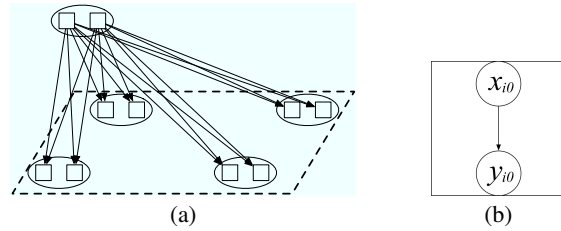


Fig. 7. (a) Zoomed-in look at Fig. 6b: statistical dependencies among parent and children nodes; (b) zoomed-in square from (a): observables are mutually independent given the corresponding hidden RVs.

extracted feature vector (i.e., MLDA atom in our case) and the corresponding hidden RV x identifies the image class, $k \in C$ of y . As is common for TSBNs, the edges between the nodes, representing x 's, describe Markovian dependencies across the scales, whereas y 's are assumed mutually independent given the corresponding x 's, as depicted in Figure 6.

There are, however, substantial differences between our approach and the standard TSBN treatment in the graphical-models literature (e.g., [10], [11]). First, our TSBN model corresponds one-to-one to the MLDA tree. Therefore, due to MLDA-tree pruning, leaf nodes of a TSBN may occur at coarser resolutions as well as at the pixel level, resulting in an incomplete tree structure (see Figure 6b). Second, we enhance the descriptive component of TSBNs by incorporating observable information at all levels of the TSBN model, not only at the lowest pixel level. Thus, to each MLDA atom in the MLDA tree, we assign a pair of observable RVs, (y_{i0}, y_{i1}) , modeling μ_0 and μ_1 RGB pixel color values. Clearly, the corresponding pair of hidden RVs, (x_{i0}, x_{i1}) , models the image classes of μ_0 and μ_1 . Third, Markovian dependencies among x 's across scales model the relation that a parent MLDA atom generates four children MLDA atoms. Given that i is a child of j , it follows that both x_{i0} and x_{i1} have the same parents x_{j0} and x_{j1} , as illustrated in Figure 7a. Therefore, our statistical model is not a tree in the strict sense of the word; however, pairs $(x_{i0}, x_{i1}), \forall i \in \mathcal{T}$, form a tree structure, allowing us to use the belief-network formalism. To emphasize all the aforementioned differences between our model and “standard” TSBNs, we refer to our model as the *incomplete meta tree-structured belief network (IMTSBN)*.

The state of a node, x_{ia} , is conditioned on the state of its parents, x_{jb} , and is given by conditional probability tables, $P_{ijab}^{kl}, \forall i, j \in \mathcal{T}, a, b \in \{0, 1\}, \forall k, l \in C$. It follows that the joint probability of all hidden RVs, $X = \{x_{ia}\}$, can be expressed as

$$P(X) = \prod_{i, j \in \mathcal{T}} \prod_{a, b \in \{0, 1\}} \prod_{k, l \in C} P_{ijab}^{kl} \quad (1)$$

We assume that the distribution of an observable RV, y_{ia} , depends solely on the node state x_{ia} . Consequently, the

joint pdf of $Y = \{y_{ia}\}$ is expressed as

$$P(Y|X) = \prod_{i \in \mathcal{T}} \prod_{a \in \{0,1\}} \prod_{k \in C} p(y_{ia}|x_{ia} = k, \theta_i^k), \quad (2)$$

where $p(y_{ia}|x_{ia} = k, \theta_i^k)$ is modeled as a mixture of M Gaussians whose parameters are grouped in θ_i^k . We usually simplify the notation as $p(y_{ia}|x_{ia} = k, \theta_i^k) = p(y_{ia}|x_{ia})$. For large M , the Gaussian-mixture density can approximate any probability density [12], in particular, the distributions of image classes. In order to avoid the risk of overfitting the models, we assume that θ_i^k 's are equal for all i at the same scale. Finally, an IMTSBN is fully specified by the joint distribution of X and Y given by

$$P(X, Y) = \prod_{i,j \in \mathcal{T}} \prod_{a,b \in \{0,1\}} \prod_{k,l \in C} p(y_{ia}|x_{ia}) P_{ijab}^{kl}. \quad (3)$$

To learn the parameters of an IMTSBN, we iteratively maximize the conditional probability $P(X|Y)$, using a probabilistic inference algorithm for TSBNs, broadly known as Pearl's $\lambda - \pi$ message passing scheme [13], [14]. In the image-processing literature similar algorithms have been proposed [15], [16]. Essentially, all these algorithms perform belief propagation up and down the tree, where after a number of training cycles we obtain all the tree parameters necessary to compute $P(X|Y)$. Note that simultaneously with Pearl's belief propagation, we employ the EM algorithm [17] to learn the parameters of the Gaussian mixture distributions. Since IMTSBNs have observable RVs at all tree levels, the EM algorithm is performed at all scales. Moreover, Pearl's message passing scheme is accommodated to sum λ messages over eight children nodes for each parent when propagating upwards, and to account for π messages of two parents for each child node when propagating downwards (see Figure 7a). Thus, the parameters of the prior models are learned through extensive training, which can be performed off-line¹. Once learned, the parameters fully characterize the likelihoods of image classes, given the pixel values. These likelihoods are then submitted to a Bayes classifier to perform image segmentation.

IV. OBJECT RECOGNITION USING VISUAL CONTEXTS

In our approach to object recognition, we seek to exploit the idea of *visual contexts* [18], and, thus, to incorporate the algorithms for flight stability and control (i.e., the MLDA-based horizon detection and sky/ground segmentation) into a unified framework. For clarity, Figure 8 illustrates the steps discussed below for a sample flight image.

Having previously identified the overall type of scene, we can then proceed to recognize specific objects/structures within the scene. Thus, objects (e.g., cars,

buildings), the locations where objects are detected (e.g., road, meadow), and the category of locations (e.g., sky, ground) form a taxonomic hierarchy. There are two main advantages to this type of approach. First, contextual information helps disambiguate the identity of objects despite the poverty of scene detail in flight images. Second, visual contexts significantly reduce the search required to locate specific objects of interest, obviating the need for an exhaustive search for objects over various scales and locations in the image.

For each image in a video sequence, we first compute its MLDA representation and perform categorization—that is, sky/ground image segmentation. Then, we proceed with localization—that is, recognition of global objects/structures (e.g., road, forest) in the ground² region. Here, we extend the set of classes C from sky/ground, generalizing the meaning of classes to any global object that may appear in flight video. Thus, the results from Section III are readily applicable, except that, since we label only ground regions, we build an IMTSBN of smaller dimensions than the image size and, thus, achieve significant computational savings.

While MLDA proves sufficient for reliable identification of global objects, it yields less reliable results for the recognition of small objects. Therefore, in the next step of the algorithm, our intention is not to actually recognize and differentiate between particular objects, but rather to detect candidate locations where objects of interest might appear in the image. For this task, it is necessary to resort to a more descriptive image representation than merely μ_0 and μ_1 of MLDA atoms. Hence, we examine spacial interrelationships of linear discriminants among neighboring MLDA atoms, such as connectedness and collinearity, since these types of geometric regularities should be most evident in artificial objects, which we seek to detect. Scanning the identified ground region in the image, we analyze 3×3 boxes of leaf-level MLDA atoms, checking whether neighboring linear discriminants are parallel, perpendicular or collinear, though the size of the analyzed window and testing procedure could easily be tailored to other application requirements. The examined image areas with the spacial frequency of the aforementioned geometric properties higher than an application-dependent threshold are extracted as candidate locations for object appearances.

To further examine the candidate object locations from the previous step, it becomes necessary to consider image analysis tools other than MLDA. The results of our extensive experimentation [4], [5], where various feature extraction methods have been investigated, suggest that the

¹In this case, there are no real-time constraints for training.

²Recognition of objects in the sky region is of lesser importance to us; however, if required, it can be easily incorporated into the algorithm.

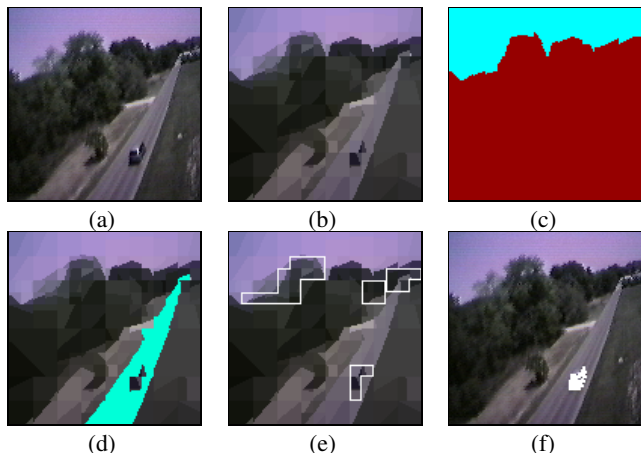


Fig. 8. The hierarchy of visual contexts conditions gradual image interpretation: (a) a 64×64 MAV-flight image; (b) MLDA image representation; (c) categorization: sky/ground segmentation; (d) localization: recognition of global objects; (e) detection of candidate object locations; (f) car recognition using CWT and HSI features.

Complex Wavelet Transform (CWT) [19] together with HSI color features yield robust and reliable object recognition. Thus, having detected candidate object locations, we proceed with new feature extraction of those particular image regions, only. CWT and HSI coefficients form complete quad-tree structures for which we then build TSBNs, similar to the procedure explained in Section III. The off-line training of TSBNs is performed for a predefined set of image classes (over CWT and HSI features). For instance, for remote traffic monitoring, a global object *road* may induce a set of objects $\{car, cyclist, traffic\ sign\}$. Consequently, for object recognition, we, in fact, consider only a small finite number of image classes, which improves recognition results.

Following inference down the hierarchy of visual contexts, we gradually perform categorization, localization and object recognition, as demonstrated in Figure 8. For space reasons, we present only the recognition of a global object *road*, then, the detection of candidate object locations in the ground region, and final recognition of an artificial object *car*, though the algorithm performs simultaneous searches over a larger set of objects. Obviously, the proposed algorithm performs well despite video noise and the poverty of image detail in the flight image.

V. RESULTS

Below, we report experimental results on the proposed framework on a set of sample flight and other natural-scene images. First, the MLDA feature extraction and our horizon-detection algorithm has been demonstrated to run at 30Hz on an Athlon 2.4GHz PC with a down-sampled image resolution of 64×64 . For the test images, representative samples of which are illustrated in Figure 3–5, our

horizon-detection algorithm correctly identifies the horizon in over 99.9% of cases.

Second, for training IMTSBNs for the sky and ground image classes, we used two sets of 500 sky and ground training images. We carefully chose the training sets to account for the enormous variability within the two classes. For each image, first the MLDA representation was found and then the corresponding IMTSBN model was trained. The training time for 1000 images of both classes takes less than 40s on an 2.4GHz x86 processor. Correct segmentation takes only 0.03s–0.07s, depending on the number of levels in the MLDA tree, for 64×64 image resolution. Clearly, the algorithm runs faster for a small number of MLDA terminal nodes, but at the cost of increased segmentation error. Note that while we are very close to meeting our 30Hz processing goal, it is not crucial for sky/ground segmentation, as long as this segmentation is updated sufficiently often. In between segmentation updates, horizon detection suffices for flight stability and control.

Having trained our sky and ground statistical models, we tested the segmentation algorithm on 300 flight images. For accuracy, we separated the test images into three categories of 100 images each: (I) easy-to-classify sky/ground appearances (e.g. clear blue sky over dark-colored land), (II) challenging images due to landscape variability, and (III) noise-degraded images. In Figure 9, we illustrate classification performance on representative test images from each category. To measure misclassification, we marked the “correct” position of the horizon for each image by visual inspection and then computed the percentage of erroneously classified pixels. In Table I, we summarize our segmentation results. Averaging results seemed inappropriate, because only a small number of test images in each category generated the larger error percentages.

Third, object recognition tests have been carried out for numerous types of objects, of which, for space reasons, we present only results for the following image classes: *road, car, cyclist, traffic sign*. For training TSBNs, we carefully chose 100 MAV-flight images for each image class. After experimenting with different image resolutions, we found that reliable object recognition was achievable for resolutions as coarse as 64×64 pixels. With this resolution, the processing time of the object-recognition procedure (i.e., sky/ground segmentation, road localization and car/cyclist/traffic-sign recognition), ranges from 1s

TABLE I
PERCENTAGE OF SKY/GROUND MISCLASSIFIED PIXELS

category I	category II	category III
2% - 6%	2% - 8%	5% - 14%

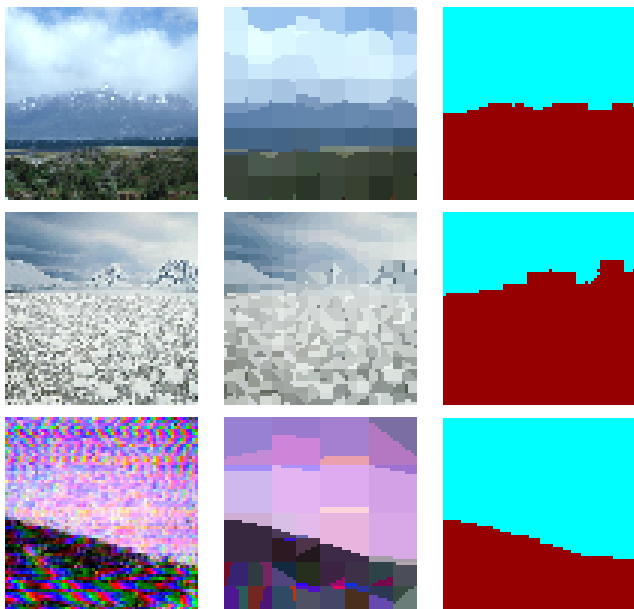


Fig. 9. Sky/ground segmentation of the three categories of test images: (top) mountain view (category I); (center) water surface covered with broken patches of ice similar in appearance to the sky (category II); (bottom) noise-degraded MAV flight image (category III).

to 4s on an Athlon 2.4GHz PC, depending on image complexity. Computing additional CWT and HSI features for candidate object locations is the most time consuming; therefore, we implement an optimization procedure that adaptively selects the optimal feature subset, as discussed in [5]. Thus, for simple images, where only one candidate location is detected, and where only one feature (e.g., wavelet coefficients oriented in the direction of 15°) is used for pixel labeling, we achieve processing times of about 0.5s, which is sufficient for the purposes of moving-car or cyclist tracking. Moreover, for a sequence of images in video, the categorization and localization steps could be performed only for images that occur at specified time intervals.

Similar to the sky/ground segmentation testing, we carried out object recognition experiments on three categories of test sets of 20 images each, containing appearances of cars, cyclists and traffic signs, as illustrated in Figure 10 and Figure 11. By visual inspection, we hand-labeled the pixels belonging to objects for each test image. In Table II, we report the recognition results of 20, 38, and 20 objects in the first, second and third category of images, respectively. The outcomes are organized into three groups: correctly recognized objects (correct recognition – CR), detected object appearances with erroneous identification (swapped identity – SI), and undetected object appearances (complete miss – CM). For the CR group of outcomes, in Table III, we present the percentage of correctly classified pixels for the three categories of images. When considering

TABLE II
OBJECT RECOGNITION RESULTS

category I (20 objects)			category II (38 objects)			category III (20 objects)		
CR	SI	CM	CR	SI	CM	CR	SI	CM
16	2	2	31	3	4	11	2	7

TABLE III
PERCENTAGE OF CORRECTLY CLASSIFIED PIXELS FOR CR OUTCOMES

category I	category II	category III
94%	91%	87%

the overall quality of recognition results, it is important to emphasize that we are dealing with relatively poor-quality images, not high-resolution, high-quality images. For example, the relatively larger error associated with category III images is obviously due to video-noise degradation.

VI. CONCLUSION

In this paper, we presented a unified computer-vision system for enabling intelligent MAV missions. We first reviewed MLDA as a feature extraction tool and demonstrated its applicability to real-time horizon detection. Then, we explained the use of IMTSBN statistical models in segmenting flight images into sky and ground regions. We extended this basic framework through the idea of visual contexts to detect and recognize object/structure appearances. Finally, we presented results of our extensive testing, which appear to validate our vision-system design choices. While testing to date has largely been done on recorded flight images and video, we are currently moving towards further closed-loop flight tests of the proposed contextual vision system. Together with other researchers, we are also working on recovering 3D-scene structure from flight video. Moreover, work is on-going in miniature-sensor integration (e.g. GPS, INS), flight vehicle modeling and characterization, and advanced control algorithm development. We expect that the totality of our efforts will eventually enable MAVs to not only fly in unobstructed environments, but also in complex urban 3D environments.

REFERENCES

- [1] P. G. Ifju, S. Ettinger, D. Jenkins, and L. Martinez, "Composite materials for Micro Air Vehicles," *SAMPE Journal*, vol. 37, no. 4, pp. 7–12, 2001.
- [2] S. M. Ettinger, M. C. Nechyba, P. G. Ifju, and M. Waszak, "Vision-guided flight stability and control for Micro Air Vehicles," *Advanced Robotics*, vol. 17, no. 7, pp. 617–40, 2003.
- [3] —, "Vision-guided flight stability and control for Micro Air Vehicles," in *Proc. IEEE Int'l Conf. Intelligent Robots and Systems (IROS)*, vol. 3, 2002, pp. 2134–40.
- [4] S. Todorovic, M. C. Nechyba, and P. G. Ifju, "Sky/ground modeling for autonomous MAVs," in *Proc. IEEE Int'l Conf. Robotics and Automation (ICRA)*, vol. 1, 2003, pp. 1422–7.
- [5] S. Todorovic and M. C. Nechyba, "Towards intelligent mission profiles of Micro Air Vehicles: object recognition in flight images," to appear in *Proc. Eight European Conf. Comp. Vision*, May 2004.

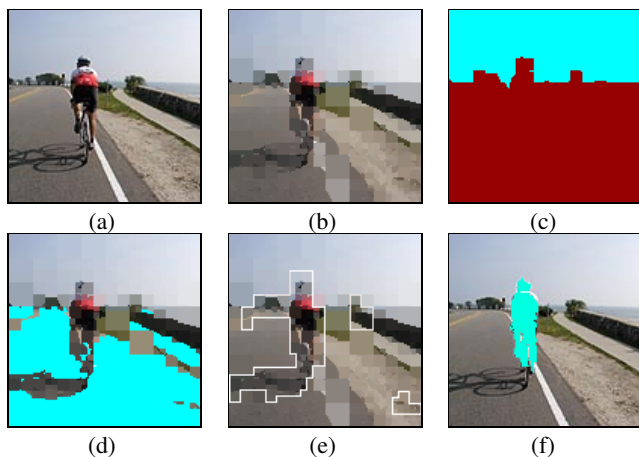


Fig. 10. The hierarchy of visual contexts conditions gradual image interpretation: (a) a category I image; (b) MLDA image representation; (c) categorization: sky/ground segmentation; (d) localization: recognition of global objects; (e) detection of candidate object locations; (f) cyclist recognition using CWT and HSI features.

- [6] D.L. Donoho, "Wedgelets: Nearly-minimax estimation of edges," *Annals of Statistics*, vol. 27, no. 3, 1999.
- [7] A. G. Flesia, H. Hel-Or, E. J. C. A. Averbuch, R. R. Coifman, and D. L. Donoho, "Digital implementation of ridgelet packets," in *Beyond Wavelets*, J. Stoeckler and G. V. Welland, Eds. Academic Press, 2002.
- [8] J. K. Romberg, H. Choi, and R. G. Baraniuk, "Bayesian tree-structured image modeling using wavelet-domain Hidden Markov Models," *IEEE Trans. Image Processing*, vol. 10, no. 7, 2001.
- [9] S. Todorovic and M. C. Nechyba, "Multiresolution linear discriminant analysis: efficient extraction of geometrical structures in images," in *Proc. IEEE Int'l Conf. Image Processing*, vol. 1, 2003, pp. 1029-32.
- [10] S. Kumar and M. Hebert, "Man-made structure detection in natural images using a causal multiscale random field," in *Proc. IEEE Conf. Comp. Vision Pattern Rec.*, 2003.
- [11] X. Feng, C. K. I. Williams, and S. N. Felderhof, "Combining belief networks and neural networks for scene segmentation," *IEEE Trans. Pattern Anal. Machine Intelligence*, vol. 24, 2002.
- [12] M. Aitkin and D. B. Rubin, "Estimation and hypothesis testing in finite mixture models," *J. Royal Stat. Soc.*, vol. B-47, no. 1, 1985.
- [13] J. Pearl, *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. San Mateo: Morgan Kaufmann, 1988.
- [14] B. J. Frey, *Graphical Models for Machine Learning and Digital Communication*. Cambridge, MA: The MIT Press, 1998.
- [15] H. Cheng and C. A. Bouman, "Multiscale Bayesian segmentation using a trainable context model," *IEEE Trans. Image Processing*, vol. 10, no. 4, 2001.
- [16] M. S. Crouse, R. D. Nowak, and R. G. Baraniuk, "Wavelet-based statistical signal processing using Hidden Markov Models," *IEEE Trans. Sig. Processing*, vol. 46, 1998.
- [17] G. J. McLachlan and K. T. Thriyambakam, *The EM algorithm and extensions*. John Wiley & Sons, 1996.
- [18] A. Torralba, K. P. Murphy, W. T. Freeman, and M. A. Rubin, "Context-based vision system for place and object recognition," AI Laboratory - MIT, Tech. Rep., March 2003.
- [19] N. Kingsbury, "Image processing with complex wavelets," *Phil. Trans. Royal Soc. London*, vol. 357, 1999.

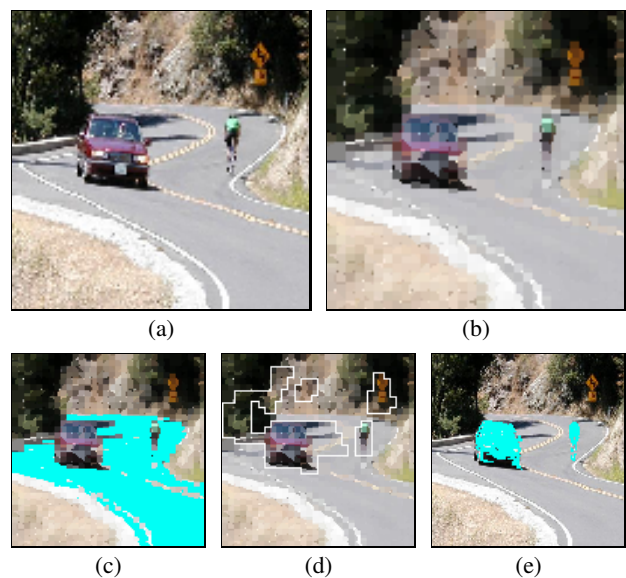


Fig. 11. Simultaneous recognition of objects: (a) a category II image; (b) MLDA image representation; (c) localization: recognition of global objects; (d) detection of candidate object locations; (e) car and cyclist recognition using CWT and HSI features. Note that while the traffic sign was detected as a candidate region, it was not recognized as such due to lack of detail.