

Fast Road Classification and Orientation Estimation Using Omni-View Images and Neural Networks

Zhigang Zhu, *Member, IEEE*, Shiqiang Yang, Guangyou Xu, *Senior Member, IEEE*, Xueyin Lin, and Dingji Shi

Abstract—This paper presents the results of integrating omnidirectional view image analysis and a set of adaptive backpropagation networks to understand the outdoor road scene by a mobile robot. Both the road orientations used for robot heading and the road categories used for robot localization are determined by the integrated system, the road understanding neural networks (RUNN). Classification is performed before orientation estimation so that the system can deal with road images with different types effectively and efficiently. An omni-view image (OVI) sensor captures images with 360 degree view around the robot in real-time. The rotation-invariant image features are extracted by a series of image transformations, and serve as the inputs of a road classification network (RCN). Each road category has its own road orientation network (RON), and the classification result (the road category) activates the corresponding RON to estimate the road orientation of the input image. Several design issues, including the network model, the selection of input data, the number of the hidden units, and learning problems are studied. The internal representations of the networks are carefully analyzed. Experimental results with real scene images show that the method is fast and robust.

Index Terms—Neural network, omnidirectional vision, road image understanding, rotation invariance, visual navigation.

I. INTRODUCTION

AN AUTONOMOUS mobile robot (vehicle) should have three basic functions in order to move safely in an outdoor road environment: road following, obstacle avoidance, and landmark recognition. All of them need the comprehensive understanding of the natural road scene. In this paper, we deal with the first and part of the last issues in an integrated manner. A robot moves along the road and makes decisions when it reaches some predefined points. It should obtain two kinds of information from the visual sensors: the robot heading (or the road orientation) and its location (in terms of road categories, e.g., straight road, intersection or T-junction). The tasks can be regarded as road classification and orientation identification.

Two problems prevent the existing vision algorithms and systems from being successfully used in this real-world application. First, most previous vision systems of mobile robots can only view objects in front of them due to the narrow view angle of the commonly used TV cameras. As a result, robots may go astray or collide against objects from the

side or behind. Second, most of the vision algorithms for outdoor road understanding only work well in predefined environments. However, whenever the environment changes, they may perform improperly.

To solve the first problem, several researchers have studied the omnidirectional vision system. Elkins and Hall [1] used a fish eye lens for the visual navigation of an outdoor mobile robot. The mobile robot located itself by referring to the known targets in the environment. Yagi *et al.* [2] applied a conic mirror to acquire omnidirectional image for the indoor mobile robot, and vertical edges were used to detect obstacles while the robot carried out a constant linear motion on flat floor. Hong *et al.* [3] used a spherical mirror to capture the omnidirectional images and studied the image-based homing problem in the indoor environment. The one-dimensional (1-D) horizon circle of the omnidirectional image captured at a location was matched with those of a series of predefined “homing” locations and guided the robot to reach the nearest “home.” Stein and Medioni [4] used a rotating camera to acquire the 360-degree panoramic images in the terrain. The omnidirectional curves of the horizon were used to find “drop off” location of the robot. Most of the above approaches share the same characteristics that specific features in the omnidirectional images are used to solve the given problems in predefined environments. Omnidirectional image methods applied to the outdoor road scenes need further investigation.

Artificial neural networks (ANN’s) are a reasonable solution for the second problem due to the following two reasons. First, for the real-world problem of road understanding, there are various aspects that should be taken into consideration, such as the weather, the light, static and dynamic objects on the road, noise, and so on. Therefore, it is difficult for a vision algorithm designed by a human programmer to include all kinds of varieties. Moreover, the need of giving thresholds for feature extraction and parameter estimation often bothers the researchers and engineers in the image analysis of real-world problems. Neural networks, in contrast with being programmed, capture knowledge and skills by training. Second, there are enough data to train the ANN’s. While the robot moves along the road, image sequences are captured at the rate of 25 (or 30) frames/s. For certain road segment with similar surroundings (e.g., straight paved road with trees and bushes on both sides), plenty of representative sample images with similar characteristics can be obtained.

ANN architectures for early vision, especially motion perception, have been proposed based on physiological as well as psychological evidences [5], [6]. In the well-known autonomous land vehicle in a neural net (ALVINN) project, a

Manuscript received October 29, 1996; revised November 10, 1997. This work was supported by the Advanced Research Project of China.

The authors are with State Key Laboratory of Intelligent Technology and Systems, Department of Computer Science and Technology, Tsinghua University, Beijing 100084, China (e-mail: zhuzhg@mail.tsinghua.edu.cn; yangshq@mail.tsinghua.edu.cn; xgy-dcs@mail.tsinghua.edu.cn; lxy-dcs@mail.tsinghua.edu.cn).

Publisher Item Identifier S 1057-7149(98)05314-7.

full-sized self-driving van equipped with video cameras and four onboard workstations has been developed and built at CMU [7], which employ the ANN in real world application. ALVINN is a fully connected three-layered backpropagation network, whose input is 32×32 subsampled road image from a video camera and whose output is the vehicle heading (one out of 45) required to make the van stay on the road. Nine hidden units are used in the system and the output is updated 15 times/s. The ALVINN network is trained using a unique “on the fly” procedure. A road image is processed as the vehicle is driven by a person down a highway. Vehicle headings, while steered by the human driver, provide the feedback necessary for training. Although the ALVINN has successfully driven the Navlab vehicle on various types of the road in various weather conditions, the system is still not perfect. The images taken at different viewing directions by a commonly used TV camera are quite different from each other due to the perspective projection and quite different viewing zones, so the network will be complicated for complex scene. In their recent work, Jochem and Pomerleau [8], [9] proposed the so-called virtual camera method to handle the lane transition problem of highway driving. The basic idea is to find the suitable image subregion and to transform it as if it was captured by a virtual camera from the desired viewing point. As a result the input requirements of the ALVINN are satisfied. Problem may also arise when the vehicle head for directions that are not included in the training range, or when the road scenes vary drastically during the long driving.

Recently we have proposed a method in which omnidirectional imagery and the neural networks are combined in order to reach a better solution for the aforementioned problems [10], [11]. The omnidirectional images provide rotation-invariant features to the neural networks, while the neural networks provide an adaptive way for image classification and identification. Our work, sharing the similar goals with the CMU’s ALVINN, possess four distinctive features.

- 1) The system estimates not only the headings of the robot but also (at the first stage) the road types along the route.
- 2) Omnidirectional view image is used. In this way the system will not be troubled by the limited view angle of the camera that has brought lot of problems in the past works of road following.
- 3) Rotation-invariant features are extracted from rotation-dependence omnidirectional images by a series of image transformations. No image segmentation and explicit feature extraction are needed.
- 4) Since we use the preprocessed image data as the inputs of the network, the number of the input variables is reduced and the networks are concise. The separation of the process of the road classification and the heading decision also greatly reduces the computation complexity.

II. OMNI-VIEW IMAGING SENSOR

To capture the omnidirectional view (omni-view) image of the environment, various imaging methods have been explored, including fish eye lens [1], conic mirror [2], spherical mirror [3], and rotating camera [4]. The time-consuming omni-view imaging by using a rotating camera prevents it from

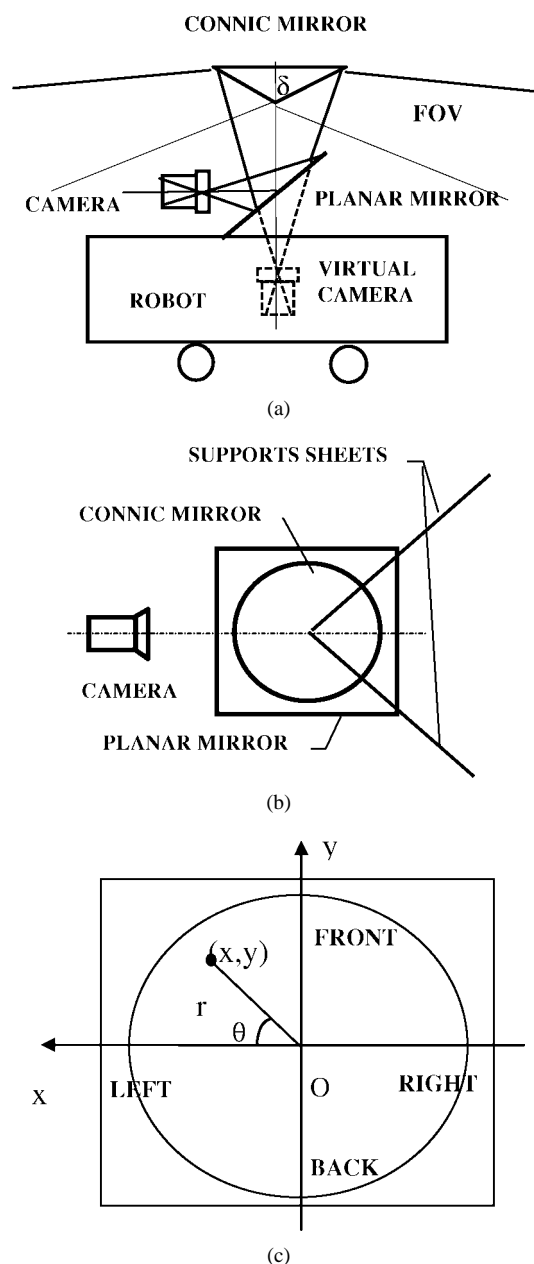


Fig. 1. OVI sensor. (a) Side view (supports not shown). (b) Top view (with supports shown) (c) Image coordinates.

applying to real-time problem. In order to obtain in real-time the omni-view of the road and objects around the mobile robot, the rest of the methods may be used. A fish eye lens yields a wide semispherical view around the camera. However the image of roads and objects near the robot locates along the circular image boundary with poor image resolution. Imagery taken by a spherical mirror provide a similar omnidirectional view of the environment as by a conic mirror, but a large part of the image is occupied by the robot itself and the image along the radius axis is not purely perspective projection but includes a quadratic distortion. So we adopt and modify the conic projection sensor system COPIS proposed by Yagi *et al.* [2], aiming to deal with the situation of the outdoor road scene.

The geometry and configurations of the omni-view image (OVI) sensor are shown in Fig. 1. A conic mirror (with a

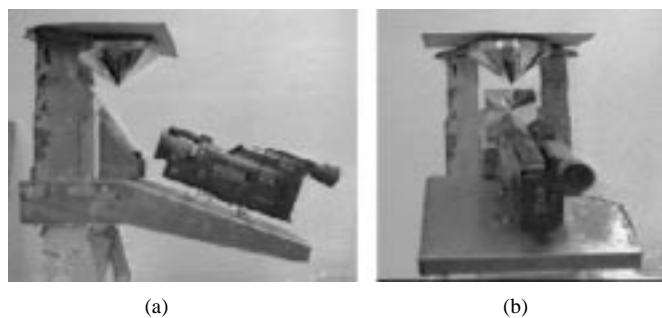


Fig. 2. Prototype of the OVI sensor. (a) Side view (b) Back view.

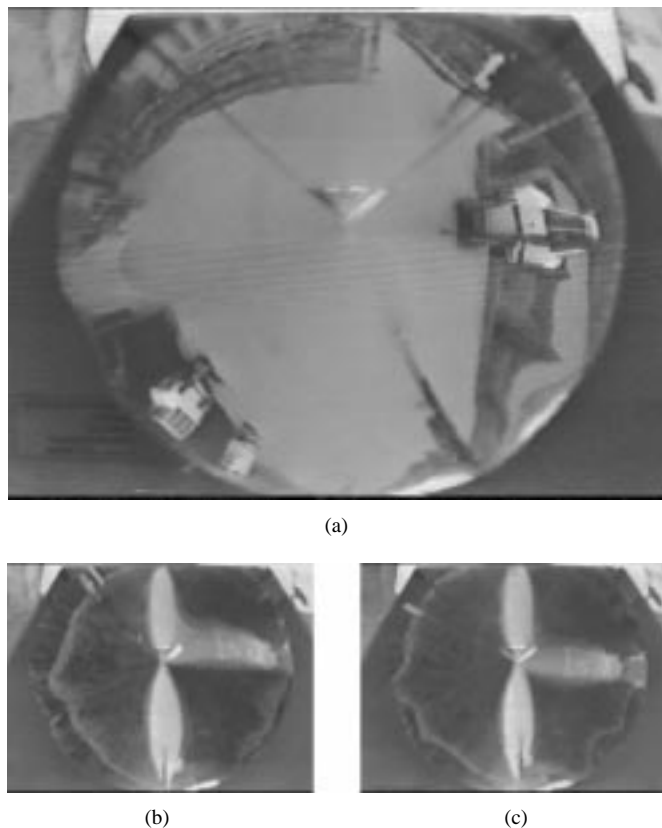


Fig. 3. OVI's. (a) 512×512 OVI image of a square scene (lower-left is a car, right is a truck and upper-right is a person). (b), (c) 256×256 images captured when the robot just entered the T-junction and was in center of the T-junction.

vertex angle $\delta = 55^\circ$) is fixed on the roof of the robot by two very thin sheet metal supports whose thickness is 1 mm. The intersection line of the sheets coincides with the vertical axis of the conic mirror. A planar mirror is placed beneath the conic mirror with a tilt angle $\pi/4$. The camera is mounted horizontally on the roof of the robot. The optical axis of the camera, the vertical axis of the conic mirror, and the normal of the planar mirror lie in the same vertical plane. The distance between the conic mirror and the roof of the robot must be large enough to avoid the occlusion of the field of view by the robot. The position of the planar mirror and the camera should be carefully adjusted to ensure the coincidence of the optical axis of the "virtual camera" inside the planar mirror and the vertical axis of the conic mirror. We use a planar mirror and a horizontally placed camera instead of a camera pointing upward for sake of easy installation and

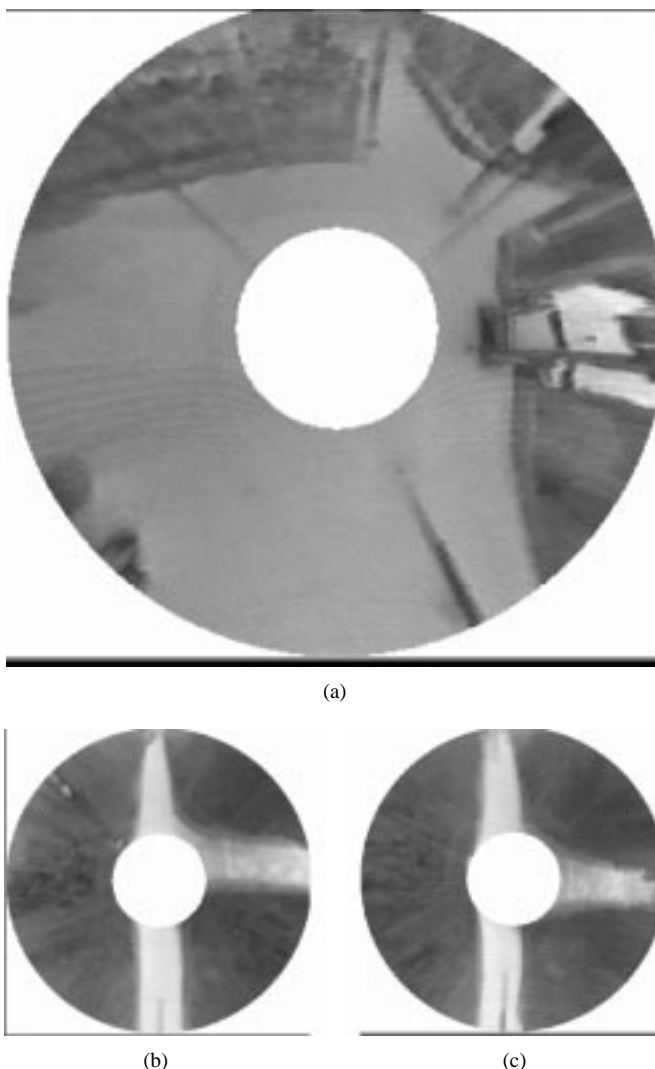


Fig. 4. Ground projections of images in Fig. 3. The robot is located in the center of the white disc in each image.

adjustments, and the protection of the camera lens. Moreover, the OVI after two mirror reflections gives an un-inverse view of the scene, as opposite to the mirror image by the COPIS system. The transparent tube used in the COPIS is replaced by two thin sheet metal supports of the conic mirror because we have found that the commonly available transparent tube is not completely transparent, and the reflection by the tube troubles image analysis, which is more severe in the outdoor environment. The diameter of the conic mirror is 110 mm and the nearest edge of each thin sheet is 100 mm far away from the conic vertical axis, so each sheet occupies less than 0.6° out of the 360° 's field of view. Fig. 2 shows a prototype of the OVI sensor used in our experiments. The camera is placed on a tilted plane (about 15°) instead of a horizontal plane in order to avoid the occlusion of the field of view by the large-size camcorder used in our experiment. The tilt angle of the planar mirror increases to about 60° correspondingly.

The image taken by our OVI sensor represents a 360° view of the scene around the robot, ranging from about 5–30 m on the ground. The omnidirectional image taken by a conic mirror is equivalent to the image taken by a tilted line scan

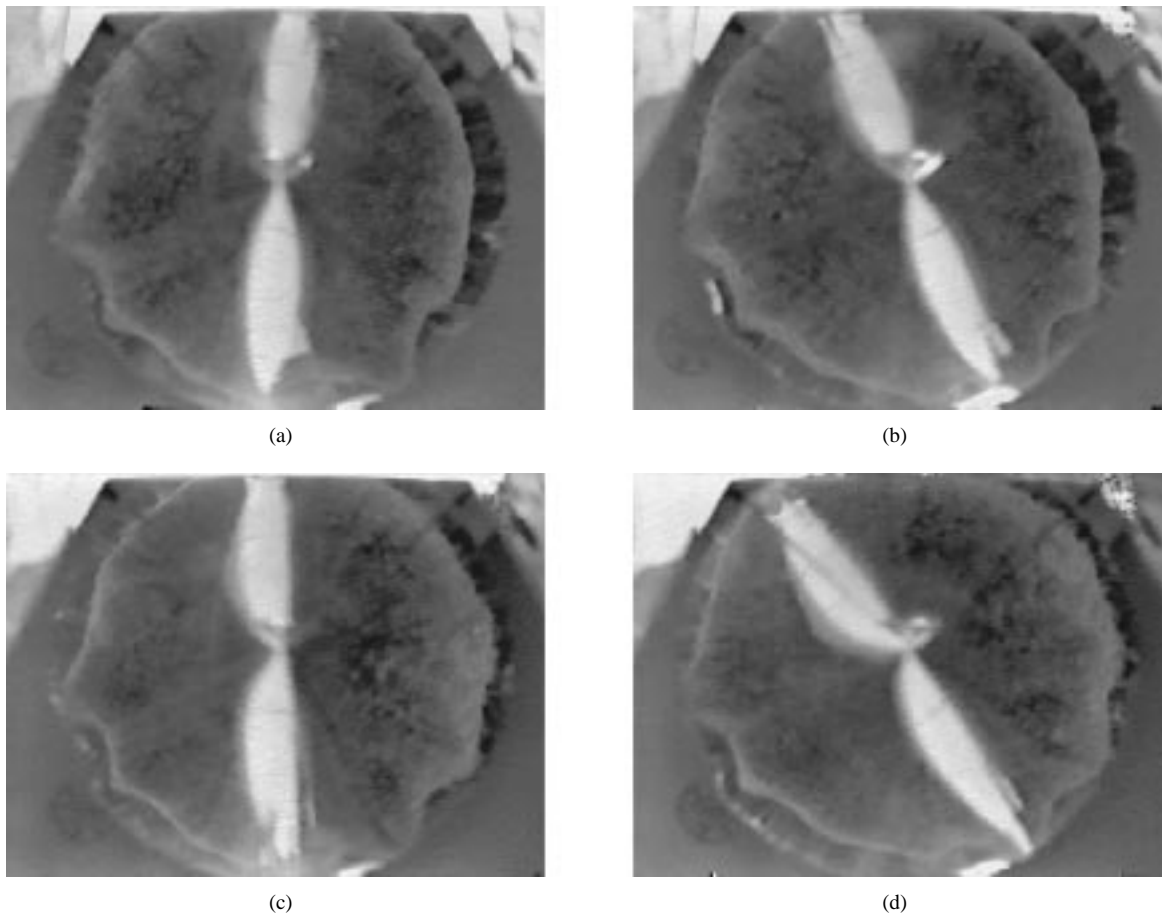


Fig. 5. Images of a paved straight road. The robot (a) headed forward, (b) rotated for an angle, (c) moved to the roadside, and (d) rotated again. The width of the road is about 5 m.

camera rotating around a vertical axis, while its optical center is moving along a circle around the same axis [12]. Fig. 3 shows three OVI's of different scenes: one is captured in the center of a square, the second is near the entry of a T-junction, and the last is in the center of the T, respectively. As the prototype OVI sensor is not accurately adjusted and the conic mirror is not perfect, the sheet supports project two thin gray lines [see Fig. 3(a)] and there are some geometric distortions around the boundaries of the images. Fig. 4 shows the corresponding ground projections of the OVI images in Fig. 3. The ground projection is approximately equivalent to the image taken by a down-looking camera high above the robot [12]. The parallelism of the road has not been completely recovered due to the errors stemmed from coarse system adjustments and a rough calibration procedure.

Although the resolution of the OVI's is relatively low compared with images of a commonly used TV camera, the 360° view image has some distinct advantages when it is used in the road scene understanding:

- 1) It covers all the information in the scene around the robot. As a result, the robot never misses the road.
- 2) The image is of rotation invariance in the sense that the structure of the image and the field of view are not changed at all if the robot rotates around the optical axis of the camera, no matter what kinds of three-dimensional (3-D) structures of the environment are.

- 3) The low-resolution sensing image is quite fit for the qualitative recognition (classification) of road categories. A small amount of lateral offsets of the robot on the road with moderate width, for example, do not bring great changes in the OVI. Appearances in the image remain similar if the robot moves within the same road segment (same category) surrounded by similar scenes. As an example, images with different rotation and lateral offsets for a paved straight road are shown in Fig. 5.

III. ROTATION-INVARIANT FEATURES

A. Polar Transform and Projection Transform

Suppose the origin of the OVI coordinate system xoy is in the center of the image where the conic vertex is projected, we transform the Cartesian coordinate image $I(x, y)$ into a polar coordinate image (r, Θ) [Fig. 1(c)]

$$r = \sqrt{x^2 + y^2}, \quad \theta = \tan^{-1}(y/x) \quad (1)$$

where r is the radius and θ is the orientation angle ($0 - 2\pi$). Ideally, for a 256×256 original sensor image, the resolution of the angle in the polar image is about 1° , so the dimension of corresponding polar image is 128×360 ($r = 0, \dots, 127$; $\Theta = 0, \dots, 359$). Since the center of the OVI is not in the image center in our actual system, and the near-center zone is too blurry to be useful, the effective range of radius is from

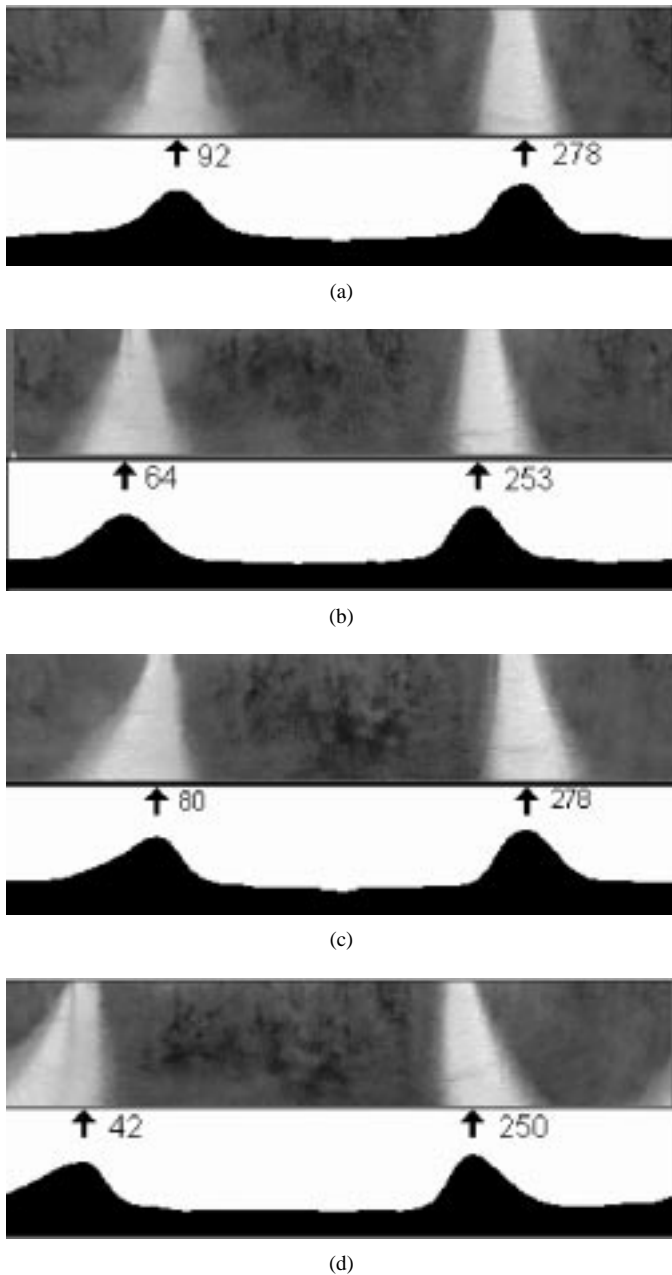


Fig. 6. OVI transformations. (a)–(d) shows the polar transform images (upper row) and the smoothed orientational projections (lower row) for the corresponding images in Fig. 5, respectively. The orientational projections have been smoothed for the sake of peak finding. The arrows between the images of the upper and lower rows indicate the estimated center of the roads, while the number beside them are the corresponding orientation angles in degrees.

30–100. Fig. 6 shows the polar images of the OVI's shown in Fig. 5.

The two-dimensional (2-D) polar image $I(r, \Theta)$ is transformed to a 1-D orientational projection $u(\Theta)$ by accumulating (projecting) the image along radius r (Fig. 6) as

$$u(\theta) = \sum_r I(r, \theta), \quad \theta \in [0, 2\pi). \quad (2)$$

In this paper, we use the Fourier transform of the original projection to extract the rotation-independent features, instead of determining the road orientation in the preprocessing stage

as we did in [10], which may bring errors to the samples. The orientation of the road will be estimated after the road category has been correctly classified.

B. Road Orientation for Data Collecting

In order to prepare the samples for training and testing the BP networks, we should know the ground truth of the road categories and orientations of each image at first. The former can be easily provided by human operators since the image sequence for a certain road category lasts for a long time period. The later, however, should be estimated by the computer automatically due to that road orientations change from frame to frame (see Section V). We use a simple three-step peak finding and tracking algorithm [12] to determine the road orientations if the roads show intuitive peaks in the orientational projection $u(\theta)$. The polar image, orientational projection, and the estimated road orientation for each OVI in Fig. 5 are shown in Fig. 6. Estimated angles between the roads in front of and behind the robot are 186° , 189° , 198° , and 208° for the four images, respectively. Even if the peak finding and tracking method is tolerant of the geometric distortions of the images, it should be pointed out that the method partially relies on the peak finding procedure. For some road categories, the method may be not successful (refer to Fig. 12). So other techniques need to be investigated.

C. Rotation-Invariant and -Dependent Features

Suppose that the orientational projection $u(\theta)$ is sampled to N discrete orientations, and then normalized into $U = \{u(n), n = 0, \dots, N-1\}$ with zero mean and unity standard deviation. The normalized procedure eliminates or at least reduces the influence of any illumination changes of images that are captured at different times. If rotation angle of the robot is

$$\phi = -\frac{2\pi n_0}{N} \quad (3)$$

where $\phi \in [-2\pi, 0]$ then the new orientational projection $v(n)$ will be the circular shift of $u(n)$ by n_0 , and can be denoted as

$$v(n) = u(n - n_0) \quad (4)$$

where $u(n)$ is the orientational projection when the robot heads for the front road (refer to Figs. 5 and 6). The Fourier transform of $u(n)$ is

$$a(k) = \frac{1}{N} \sum_{n=0}^{N-1} u(n) \exp\left(-\frac{j2\pi kn}{N}\right), \quad k = 0, \dots, N-1 \quad (5)$$

and the Fourier transform of $v(n)$ can be expressed as

$$b(k) = a(k) \exp\left(-\frac{j2\pi n_0 k}{N}\right), \quad k = 0, 1, \dots, N-1. \quad (6)$$

By representing $a(k)$ and $b(k)$ in amplitude-phase forms

$$\begin{aligned} a(k) &= a_k \exp(j\psi_k), \\ b(k) &= b_k \exp(j\varphi_k), \quad k = 0, \dots, N-1 \end{aligned} \quad (7)$$

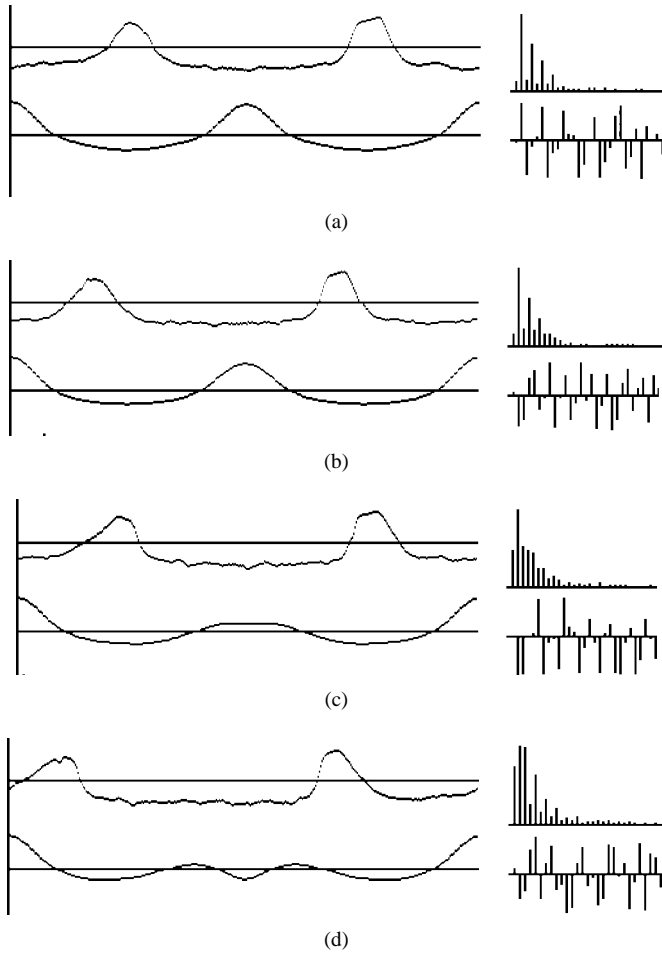


Fig. 7. Each of (a)–(d) shows the normalized projection ($N = 360$, upper left), the first 30 Fourier amplitudes (upper right), the first 30 Fourier phases (lower right) and the ACF's (lower left) of the corresponding polar images in Fig. 6. The ACF's are rotation-invariant, which is clearly shown in (a) and (b). Since there are large geometric distortions between projection in (d) and the others, the ACF's are not identical.

we have the following results:

$$b_k = a_k, \quad k = 0, \dots, N-1 \quad (8)$$

$$e^{j\psi_k} = e^{j(\varphi_k + k\phi)} \quad (9)$$

$$\psi_k = 2\pi m_k + \varphi_k + k\phi, \quad k = 1, \dots, N-1 \quad (10)$$

where m_k is an integer that indicates $2\pi m_k$ additive ambiguous in the k th phase value. Equation (8) says that the Fourier amplitudes are invariant to the rotation of the OVI's. Therefore, they are appropriate features for road scene classifications. Equations (9) and (10) give the basic relation to estimate the orientation difference between two OVI's. For real scene images, the equality can not hold exactly. Analysis and experimental results show that the Fourier phases are very sensitive to the noises, especially to the geometrical distortion of the orientational projections. Fig. 7 shows the normalized projection ($N = 360$), the first 30 Fourier amplitudes, and the phases of the polar images in Fig. 6. It can be seen that the Fourier amplitudes are similar for the four images, but the phases are not stable, especially when the corresponding amplitudes are small.

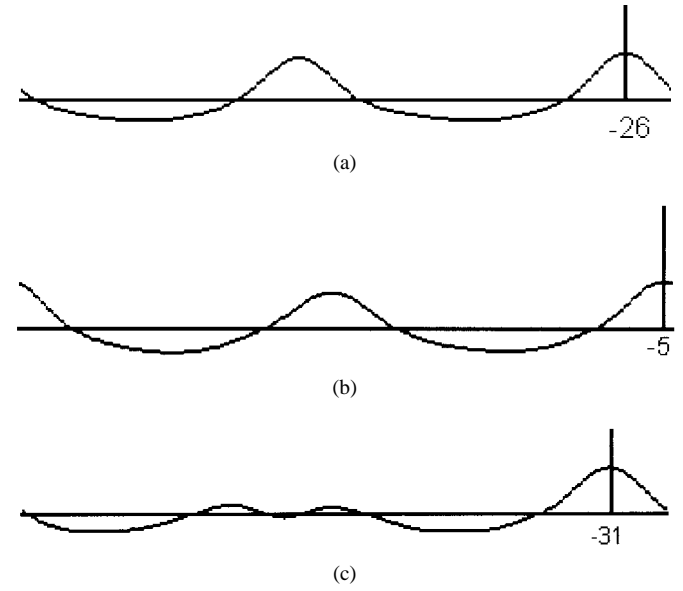


Fig. 8. Orientation differences by using the CCF in the Fourier domain. (a) CCF between projections (a) and (b) in Fig. 7. (b) CCF between projections (a) and (c) in Fig. 7. (c) CCF between projections (c) and (d) in Fig. 7. The vertical lines and the number besides them indicate the position (n_0) where maximum CCF take place (The axis is zero to -360 from right to left).

D. Road Orientation Difference for Data Collecting

Actually, the orientation difference could be estimated by searching the minimum value of the following distance:

$$d(\phi) = \sum_{k=0}^{N-1} (a_k e^{j\psi_k} - b_k e^{j(\varphi_k + k\phi)})^2 \quad (11)$$

for each $\phi(n_0) = -(2\pi n_0/N)$, $n_0 = 0, 1, \dots, N-1$. The computing complex of this procedure is $O(N^2)$. This is equivalent to find the maximum value of the circular cross-correlation function (CCF)

$$C(n_0) = \sum_{n=0}^{N-1} u(n)v((n+n_0) \text{ modulo } N), \quad n_0 = 0, 1, \dots, N-1 \quad (12)$$

which is also $O(N^2)$ computation. Since we have obtained the Fourier transform of $u(n)$ and $v(n)$, we hope to use them for the estimation of orientation difference. The correlation theorem states that the (circular) correlation of two real signal sequence $u(n)$ and $v(n)$ is equal to the inverse discrete Fourier transform (DFT) of the product of conjugation of DFT $a(k)$ and DFT $b(k)$, i.e.,

$$C(n_0) = F^{-1}(a^*(k)b(k)). \quad (13)$$

The alternative approach has only $O(3N \log_2 N)$ -time complexity. The advantage of the global correlation method is that no feature extraction is needed, so the correlation method is more robust to noise and more general in different cases. Fig. 8 shows the experimental results. Compared with the estimations in Fig. 6, the results from the CCF method are the average of the orientation difference between the front roads and the orientation difference between the back roads in the two images. Therefore, the CCF method, together with the

peak finding and tracking method, is used in the data collection and training for the networks since the behavior of the robot can be controlled in the training stage, and the man-machine interaction can be involved. In real operations, the orientation is estimated using the neural networks.

E. Another Rotation-Invariant Feature

From (12) and (13) we can obtain another rotation-invariant sequence, the auto-correlation function (ACF) of $u(n)$

$$R(n_0) = \sum_{n=0}^{N-1} u(n)v((n+n_0) \text{ modulo } N),$$

$$n_0 = 0, 1, \dots, N-1. \quad (14)$$

The ACF of $u(n)$ is equal to the inverse Fourier transform of the energy spectrum of $u(n)$

$$R(n_0) = F^{-1}(a^*(k)a(k)) = F^{-1}(a_k^2). \quad (15)$$

Fig. 7 also shows the ACF's estimated by the inverse Fourier transforms of the energy spectrums of the projections. Compared with the Fourier amplitudes, the similarity of ACF's is more intuitive (refer also to Fig. 13). So the ACF can serve as another rotation-invariant input for road classification.

IV. SYSTEM ARCHITECTURE

In order to successfully work with real-world problems, we must deal with some design issues, including the network model, network size, activation function, learning parameters, and selection of training samples. We will address these issues in this and the following sections, bearing in mind that we face the problems in the outdoor road scene, namely the road classification and orientation estimation.

A. RUNN Architecture

It is commonly accepted that the backpropagation learning procedure has become the most popular method to train multilayer feed-forward networks [13], and the so called backpropagation (BP) networks have been widely used in character recognition, speech recognition, vehicular control and many more cases of applications [7], [10], [13], [14]. There are two main schemes for using ANN's in a pattern classification system [15]. The first one employs an explicit feature extractor (not necessarily a neural network). The extracted features are passed to the input stage of the multilayer BP network. The scheme is very flexible in incorporating a large variety of features. However explicit features, e.g., the boundary of the road, have proved to be very difficult to extract in the outdoor road scene. The other scheme does not explicitly extract features from the raw data. The feature extraction implicitly takes place within the hidden layers of the ANN. A nice property of this scheme is that feature extraction and classification are integrated and trained simultaneously to produce optimal classification results. However, it is not clear whether the types of features that can be extracted by this integrated architecture are the most effective for the given

pattern classification problem. Moreover, this scheme requires a much larger network than the first one. A typical example of the second scheme for visual navigation is the ALVINN [7]–[9].

We take an alternative approach from the two typical schemes. The basic model for road understanding neural networks (RUNN) is an adaptive combination of an image processing module (IPM) and several fully connected BP networks—a single three-layer road classification network (RCN), one two-layer road orientation network (RON) for each road category. Fig. 9(a) shows the system architecture. Raw image data is preprocessed by the IPM before feeding into the neural nets. However, no image segmentation and explicit feature extraction are needed. A composed macro-network, composed of several basic BP networks, is constructed to solve both the road classification and orientation estimation.

B. Configurations of RCN's and RON's

We adopt the convention that a standard Y -layer BP network consists of an input layer, $(Y-2)$ hidden layers and an output layer of units successively connected in a feedforward fashion. The RCN is a fully connected three-layer BP network [Fig. 9(b)]. The inputs of the RCN are P ($\leq N$) rotation-invariant image data (i.e., Fourier amplitudes or ACF), and the outputs are M road categories. The net can be viewed as a nonlinear input-output mapping. The connection between the input and the hidden layers extract special features of input patterns and the connections between the hidden and the output layers recognize specific road categories. Therefore, the hidden units may represent different kinds of features and the number of units in this layer will be decided by an experiment-based approach. A bias unit is connected to the hidden layer and the output layer, respectively.

Each RON is a two-layer BP network without any hidden layer [Fig. 9(c)]. Its function is virtually a correlation operation, and details of this design will be explained in Section VI. The inputs of each RON are Q ($\leq N$) rotation-dependence image data (i.e., the original orientational projection), and the outputs are L road orientations.

The basic model of the neuron unit, or the processing element (PE) for all the networks, is shown in Fig. 9(d). The summation function is

$$I_i = \sum_j w_{ij}x_j + \beta \quad (16)$$

where i stands for the current PE, j stands for a PE that i is connected to, x_j is the output of PE j , w_{ij} is the weight of the connection of i and j , and β is the bias value. The transfer function of each input unit is linear and that of each hidden or output unit is the hyperbolic tangent [TanH, see Fig. 9(e)]

$$\Gamma = I * \text{gain},$$

$$O = \frac{e^\Gamma - e^{-\Gamma}}{e^\Gamma + e^{-\Gamma}} \quad (17)$$

where gain is called the steepness factor. The TanH is quite similar to the sigmoid transfer function. However, its range is -1 to $+1$, as opposed to the sigmoid range of 0 – 1 . Because

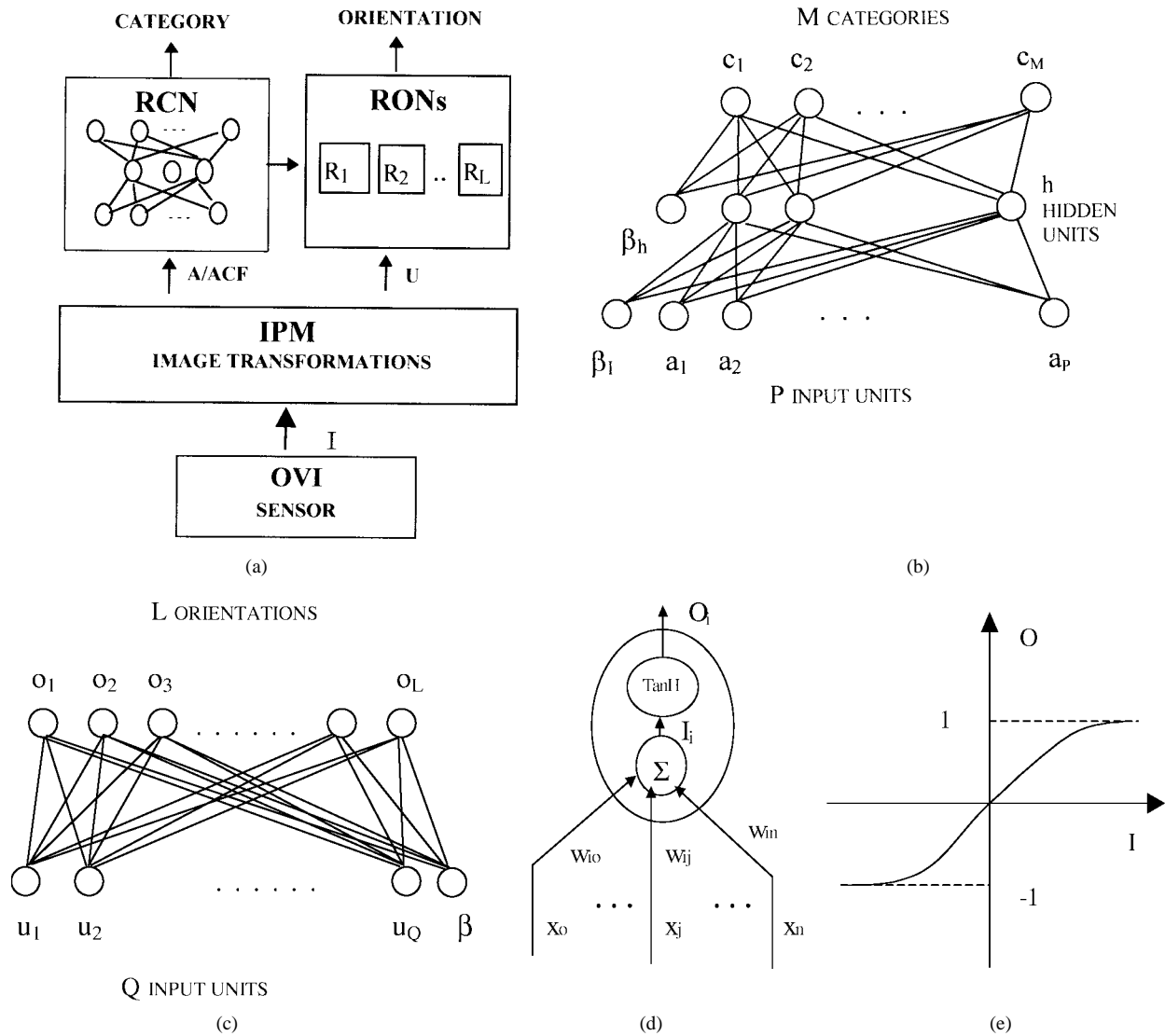


Fig. 9. Architecture of the RUNN. (a) The system is composed of an OVI sensor, an IPM, an RCN, and a set of RON's, i.e., R_1, R_2, \dots, R_L . (b) RCN for road recognition. (c) RON for road orientation. (d) The model of PE. (e) Transfer function TanH.

the output of the transfer function is used as a multiplier in the weight update equations, a range of 0–1 could lead to a bias to learning higher desired output (approaching 1). The TanH gives equal weight to low- and high-end values.

C. Installation and Execution

The RCN and RON's are set up using the Nworks tool [15], which provides a variety of ANN architectures and the flexibility of parameter controls. The IPM's hardware is a pipeline image processing machine named PIPE [16], hosted by a PC 486 with a Nworks environment. The UserIO interface (written in C) of the Nworks connects the image processing module (PIPE) and the neural network environment, working in parallel. Our PIPE system includes a video stage, an input stage, three modular processing stages (MPS's), an iconic-to-symbolic mapping stage (ISMAP) and a output stage. Each of the MPS's includes several frame buffers, two real time convolution processors, ALU's, and most important, a two-valued function (TVF) LUT that greatly facilitated the image geometrical transformations. The polar transformation can be

carried out in an MPS at the video field rate (60 fields/s). The projection transformation is implemented in the ISMAP at the frame rate. The 1-D Fourier transform is carried out by the PC486.

The RUNN works in the following three steps.

- 1) The OVI is captured and transformed by the IPM. The orientational projection U and rotation-invariant Fourier amplitudes are obtained. For comparison, the auto-correlation function (ACF) sequence is also calculated using the inverse Fourier transform of the energy spectrum.
- 2) The rotation-invariant data (A or ACF) is used to decide the road category by the RCN.
- 3) The road category estimation is used to activate the correct RON, and the rotation-dependent data (i.e., the original orientational projection U) is fed into that selected RON to estimate the road orientation.

There are two advantages of the separation of road classification and orientation estimation. First, since rotation-invariant data are used as the input of the RCN instead of the original

image, the distinctiveness of the input units is increased, and therefore the complexity of the network is reduced. If the Fourier amplitude A is used as the input of the RCN, the number of the input units can be reduced to $P \leq (N/2)$. Second, since a separate RON is used to estimate the road orientation for each road category, and the classification result is used to select the corresponding RON, the efficiency of the networks will be improved.

V. DATA COLLECTION FOR TRAINING AND TESTING

The data for the RUNN are collected while the robot is moving on the road. In our experiments, the robot moves along the route around the main building at the campus of Tsinghua University. The OVI sequences are recorded by a video camcorder, and then are played back for processing by the RUNN system.

A. Collecting the Data for the RUNN

At the beginning of data collecting for each category of road segment in a camera shot, the robot heads for the front road (i.e., the road orientation angle is zero), and the desired output of the RCN, representing the road category, is assigned by a human supervisor. The peak finding and tracking algorithm (Section III) with human interaction verifies the orientation of the first frame for each shot. For the current experiments, the input data of the RCN are $N = 32$ subsampled elements of the orientational projection. The road images are classified as six categories ($M = 6$): paved straight road surrounded by bushes and trees (denoted as “||”), T road junction (denoted as “T”), intersection (denoted as “+”), earthy road surrounded by grass and trees (denoted as “D”), narrow curved road passing through the garden in front of the building (denoted as “C”) and the square in front the main building (denoted as “S”). In order to cover most of the situations, the robot zigzags on the road so that captured images can cover most possible directions and various lateral offsets (refer to Fig. 10). For preparing data to train and test the RON, the orientation difference is calculated for the successive image frames within the same road category by find the maximum value of the CCF in (13). The absolute road orientation is obtained by accumulating the orientation differences and is modified by the peak-finding and tracking method. In order to cover most of the rotation (orientation) cases, the sampled orientational projection is shifted by software to simulate all the different road orientations. Both the inputs to the network as well as the desired outputs are mapped into numbers. Fig. 11 shows one typical sample image for each of the six categories. Figs. 12 and 13 show the corresponding polar images, projections, Fourier amplitudes and phases, and the autocorrelation functions of the images in Fig. 11.

B. Selecting and Dividing the Data

It is important to make sure that examples selected for training the network do not have any dubious data fields (e.g., outliers). To this end, we calculate the mean Fourier amplitude vector of all the samples within one category, and the distance between any Fourier amplitude vector of each sample and the

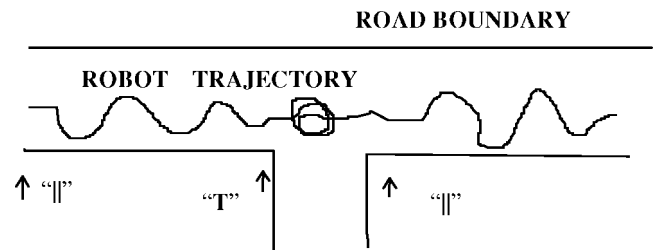


Fig. 10. Collecting the data for training and testing.

TABLE I
SELECTING AND DIVIDING THE DATA FOR THE RCN

Set \ Category		T	+	D	C	S	Total
Original set	1073	371	182	557	441	175	2799
Selected set	930	312	160	445	374	148	2342
Training set	133	133	133	134	133	133	799
Testing set	797	179	26	311	241	15	1569

mean is used to judge whether it is a “good” example. Good examples are chosen from the original raw data set and then are divided into the training set and the testing set. For best results, the selection of training set is based on the following rules: i) every category has roughly same amount of examples; ii) the training set is reasonably representative of each category; iii) it is best to make the testing and training sets completely separate. The actual selection and division results are listed in Table I. The numbers listed in the table are the numbers of real images captured by the OVI sensor, and for the training and testing of the RCN. For the training and testing of the RON’s, each sample has as many as L rotated versions.

C. Scaling the Data

The inputs are already in number forms. The desired outputs are set to either 0 or 1 (e.g., the output vector is “1, 0, 0, 0, 0, 0” for category “||”). Since we use the TanH as transfer function, we will need to scale these values between -1 to $+1$. Fortunately, a so-called minmax table mechanism is provided by the Nworks tool. This preprocessing facility computes the lows and highs of each data field corresponding to each input unit (in the training set or both the training and testing sets) and stores in the minmax table. The Nworks then computes the proper scale and offset for each data field. Real-world values are then scaled to network range (-1 to $+1$) for presenting to the network. Whenever a scaled result is produced, it is descaled to real-world units.

VI. CONSTRUCTION OF THE RUNN

We construct the RUNN during the training and testing process using real image data, and study the following four issues:

- 1) the suitable representation of input data;
- 2) the number of the hidden units;
- 3) the internal representation of the networks;
- 4) the learning problem, for example, the epoch size, the converging speed, etc.

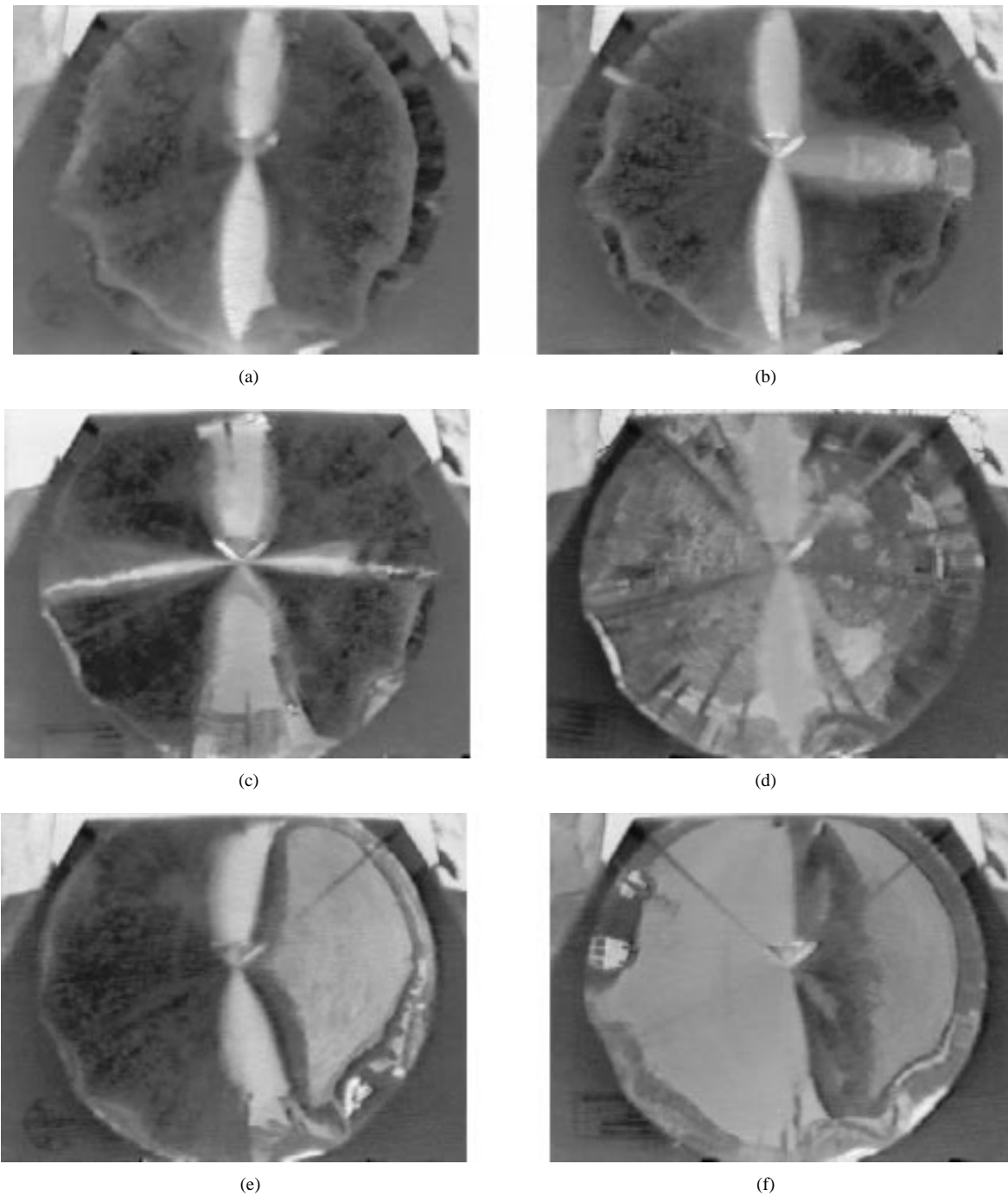


Fig. 11. Sample images of the six road categories (a) Paved straight road (“||”). (b) T-junction (“T”). (c) Intersection (“+”). (d) Earthy road (“D”). (e) Curved road (“C”). (f) Square (“S”).

A. Learning Rule and Schedules

The BP learning strategy is used to training the network with M output units. The error in the output layer is computed as the difference between the desired output ($D = (d_1, d_2, \dots, d_M)$) and the actual output ($Y = (y_1, y_2, \dots, y_M)$). This error E , transformed (scaled) by the derivative of the transfer function TanH, is backpropagated to the *priori* layer where it is accumulated. This backpropagated and transformed error becomes the error term of that *priori* layer. The process of BP continues until the first layer is reached.

In our implementation, the normalized cumulative delta learning rule [16] is used for the RUNN. Cumulative gener-

alize delta rule attempts to alleviate the problem of structured presentation of the training set. The basic idea is to accumulate the weight changes over several training presentations and make the application all at once. The update equations are

$$\begin{aligned}
 m'_{ij} &= m_{ij} + C_1 e_i x_{ij}, \text{ at each iteration;} \\
 \begin{cases} w'_{ij} = w_{ij} + m_{ij} + C_2 a_{ij} \\ a'_{ij} = m_{ij}, \\ m'_{ij} = 0 \end{cases} & \text{if } \text{LCNT} \bmod \text{AUXI} = 0 \quad (18)
 \end{aligned}$$

where C_1 is the learning coefficient (step size), C_2 is the momentum factor, $E = (e_1, \dots, e_n)$ is the error vector

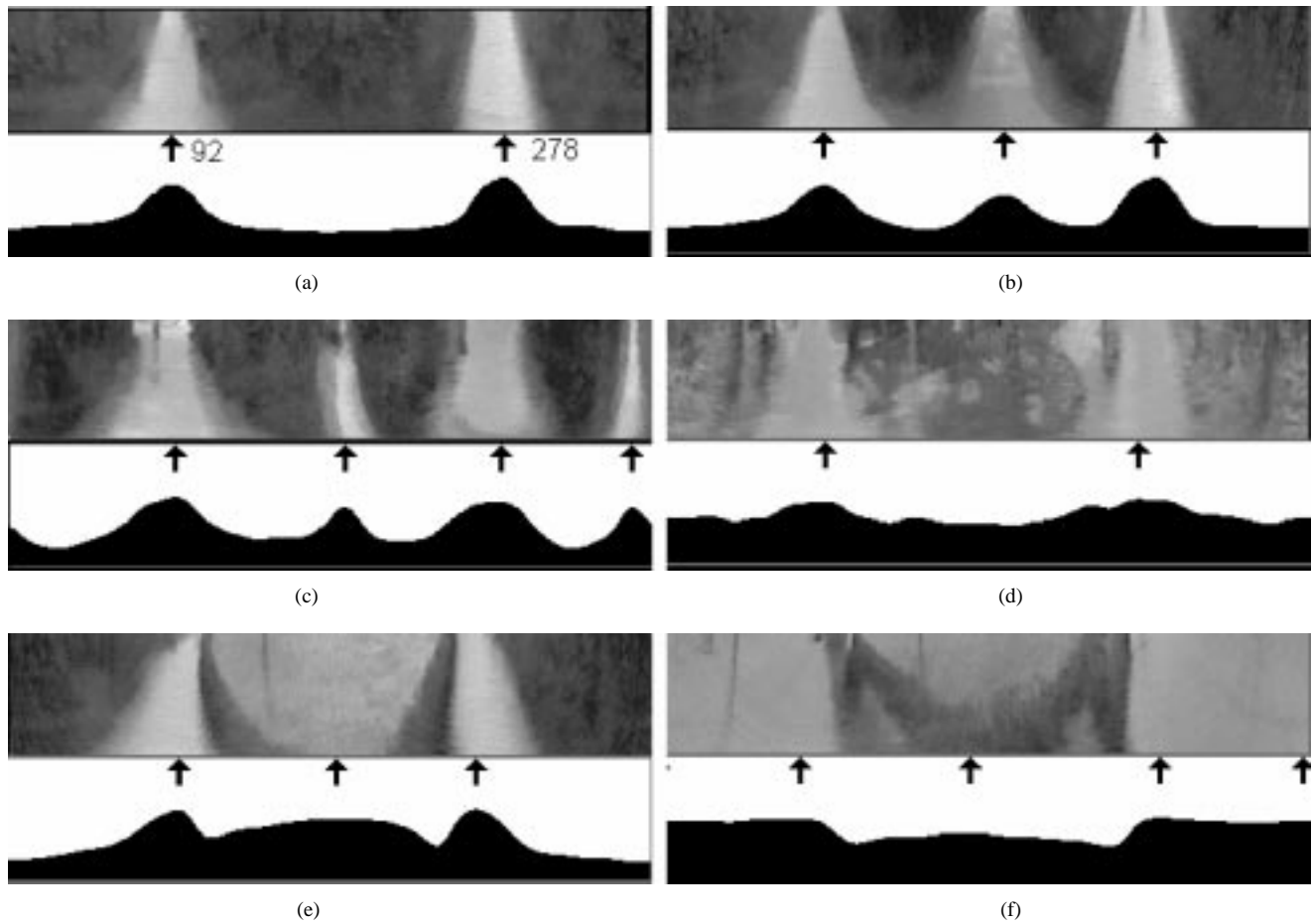


Fig. 12. Polar images and projections. Each of (a)–(e) and (f) shows the polar transform image (upper row) and the smoothed orientational projection (lower row) for the corresponding image in Fig. 12. The arrows between the upper and lower images indicate the center of the roads. Notice that some results of the peak-finding algorithm are not correct for images in (e) and (f).

(as described above), $\mathbf{X}_i = (x_{i0}, \dots, x_{in})$ is the inputs to the i th PE in the current layer, $\mathbf{W}_i = (w_{i0}, \dots, w_{in})$ is the initial weight vector for the i th PE in that layer, and $\mathbf{W}'_i = (w'_{i0}, \dots, w'_{in})$ is the updated weight vector, $\mathbf{M}_i = (m_{i0}, \dots, m_{in})$ is the accumulated weight changes for the i th PE and $\mathbf{A}_i = (a_{i0}, \dots, a_{in})$ is the auxiliary weight field that is used as momentum term. Lcnt is the count of learned samples and AUX1 (epoch size) is the accumulation period. The RMS error (RMSE) is the stopping criterion for training and is defined as

$$\text{RMSE} = \sqrt{\frac{1}{2} \sum_{i=1}^{\text{AUX1}} \|\mathbf{Y}^{(i)} - \mathbf{D}^{(i)}\|^2}. \quad (19)$$

One of the problems with the cumulative delta rule is that the learning coefficient C_1 depends on the epoch size AUX1. As the size AUX1 increases, C_1 should get smaller, otherwise the accumulated weight changes will become too large and cause the learning to diverge. Normalized cumulative delta rule gets around this problem by dividing the accumulated delta weight by the square root of the epoch size before being applied. Moreover, examples in the training set were presented to the network randomly during the training to avoid the “learn one thing but forget others” problem.

During the learning process, different schedules are used for adjusting the learning rates (parameters C_1 and C_2) for the in-

TABLE II
LEARNING SCHEDULE

Schedule for the input layer					
LCNT	50000	0	0	0	0
C_1	0.9000	0.0000	0.0000	0.0000	0.0000
C_2	0.6000	0.0000	0.0000	0.0000	0.0000
Schedule for the hidden layer					
LCNT	10000	30000	70000	150000	310000
C_1	0.3000	0.1500	0.0375	0.0023	0.0000
C_2	0.8000	0.4000	0.1000	0.0063	0.0000
Schedule for the output layer					
LCNT	10000	30000	70000	150000	310000
C_1	0.1500	0.0750	0.0188	0.0012	0.0000
C_2	0.8000	0.4000	0.1000	0.0063	0.0000

put, hidden and output layers, respectively (see Table II). Point fields, e.g., the intervals between transition points increase exponentially, and the coefficient ratio (e.g., 0.5) defines an exponential decay of the C_1 and C_2 for the hidden and output layer, which is sampled at subsequent transition points (e.g., $C_1 = 0.5, 0.5/2, 0.5/4, \dots$). The coefficients for the input layer are not changed with the learning iterations.

B. Road Classification

First the rotation-invariant Fourier amplitudes are used to train the RCN. In this case the RCN has 16 inputs

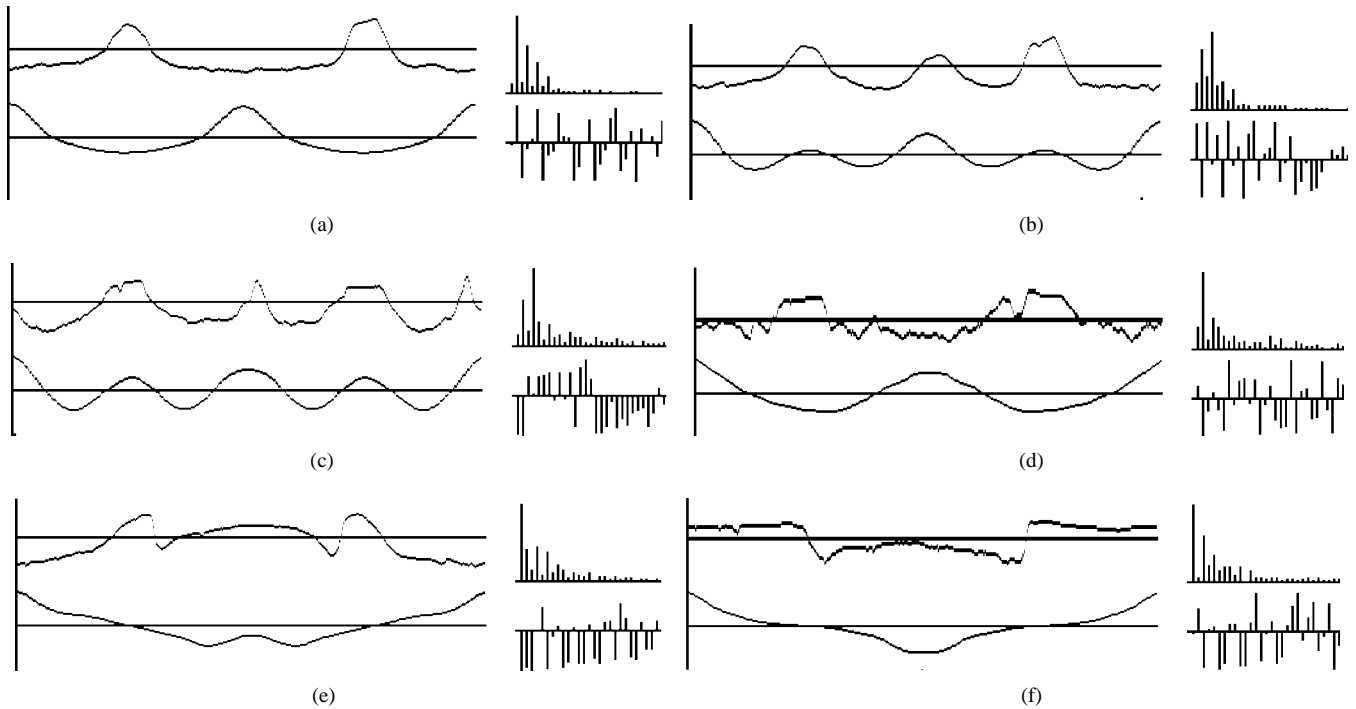


Fig. 13. Normalized projection ($N = 360$, upper left), first 30 Fourier amplitudes (upper right), first 30 Fourier phases (lower right), and ACF (lower left) of the corresponding polar image of each category in Fig. 12. Notice the differences of Fourier amplitudes and ACF curves among the six categories, and the similarities among the sampled images of the same category in Fig. 7.

TABLE III
LEARNING PROCESS OF THE RCN

h	0	3	4	5	6	7	10	12	16
C	53634	51137	86752	73906	65581	75280	71064	86923	50939
e	0.30	0.25	0.20	0.20	0.15	0.15	0.12	0.10	0.14

(h: Number of the hidden units; C: Training iterations; e: RMS error)

TABLE IV
RECOGNITION RATES OF THE RCN USING FOURIER
AMPLITUDES WITH VARIOUS NUMBER OF HIDDEN UNITS

set \ h	0	3	4	5	6	7	10	12	16
training	86.9	91.8	96.3	94.6	96.7	96.6	97.6	98.5	98.5
testing	83.5	86.5	91.6	88.9	92.5	91.6	92.9	94.4	94.2
original	77.2	79.5	86.1	84.1	85.9	85.3	87.6	89.0	88.6

(a_1, \dots, a_{16}) and 6 outputs. Experiments indicated that 16 or 32 is the proper size of the update epoch AUX1. If the epoch is too small (e.g., 8), noise in the data set seems to confuse the network and the network may oscillate; on the other hand, if the epoch size is too large (e.g., 64), proper adjustments may be ignored and the network may not converge. The number of the hidden units, h , is decided by a systematic experimental analysis in order to find the minimum number for the proper classification and best number for the problem. Table III shows the training results for different number of hidden units ($h = 0-16$). “C” is the number of training iterations at which the network becomes stable. The RMS errors “e” are also listed in the table.

Each realization of the RCN was tested by using the training set, testing set, and the original raw data set, which may include noisy data. If the value of one (e.g., k th) of the six network outputs $Y = (y_1, \dots, y_6)$ is greater than 0.5 and the

TABLE V
RECOGNITION RATES OF THE RCN USING ACF

set \ h	3	4	5	6	7	8
training	95.24	94.37	97.37	97.00	80.35	80.85
testing	92.22	90.25	92.99	92.73	90.38	92.67
original	87.17	84.71	87.17	88.85	82.74	84.67

other five values are less than 0.5, then the input road image is thought of being correctly classified and is assigned to k th category. Table IV lists the correct recognition rate (%) for the three data sets under every realization of the RCN.

The training and testing process indicates that four is the minimum number of the hidden units for proper classification, and 12 is the best number. Comparing with the learning process of the networks using rotation-independent orientational projection as inputs in [10], the network RCN converges much slower and the recognition rate is slightly lower, and more hidden units are needed for the best classification. The reason may be that the Fourier amplitudes lose the phase information of the orientational projection and the Fourier transform compacts the energy in the first several terms so that weights between the input and the hidden units are not balanced. However, the rotation-invariance and fewer components of the Fourier amplitude vector makes itself a good choice for the input of the RCN.

We also try to use the ACF, which is calculated as the inverse Fourier transform of the energy spectrum of the orientational projection and is still rotation-invariant, as the inputs of the RCN. Although this transform does not add or lose any information, experiments indicate that this kind of data representation may be more distinctive for the BP network (refer to Fig. 13). The testing results using ACF

TABLE VI
RECOGNITION RATES AND ERRORS OF THE TWO-LAYERED RONS WITH INPUT U

class		T	+	D	C	S
epoch size	4	1	1	16	8	8
training set	100. (0)	99.9 (0)	100. (0)	83.9 (4)*	96.2 (2)	89.4 (1)
testing set	99.0 (0)	100. (0)	100. (0)	83.1 (4)	95.7 (2)	86.7 (1)
original set	96.8 (1)	98.1 (1)	87.9 (5)	70.4 (9)	90.8 (4)	76.5 (6)

*For example 83.9 (4) means the correct rate is 83.9% while the average error for the rest is 4%.

show improvements in convergence speed (half the iteration number as using FA) and recognition rate while the number of hidden units is fewer (Table V). Analysis of the weight patterns shows that the connections between the input and the hidden layers extract more distinctive features from this kind of input patterns, which are more suitable for road classifications. However the improvements of the RCN with input ACF require more computation power for inverse Fourier transform and processing more input units.

C. Road Orientation Estimation

After the road category is determined, the corresponding RON for this category is activated. We compare the results of the networks using original orientational projection U and the Fourier phases Ψ as inputs, and having different number of units in the single hidden layers. The outputs of the RON are L ($= 32$) discrete orientations. The three-layer RON's with different number of hidden units (from zero to 16) do not converge when the input is the phase data. The reason might be that the phase data is sensitive to noise and has additive ambiguous. So we use the original orientational projection U as the inputs of the RCN. Experiments using zero to 16 hidden units in a single hidden layer show that the RON's with no hidden units (i.e., two-layer BP network) perform best for all the road categories. More detailed analysis will be given in the next section. Table VI shows the correct orientation estimation rate, measured by percentage of errorless estimation, and the average orientation error of the rest in percentage ($100 \times \Delta\phi/2\pi$), for each road category. The epoch size of learning process for each road category is also presented in Table VI. The estimation accuracy decreases and the epoch size should be large when the input data become noisy and scattered for a certain category (e.g., "D"). It also indicates that this category should be further divided into several subcategories. The robustness of orientation estimation could be improved by integrating the RON result with the information of the orientation difference of the temporal sequences calculated by (13).

D. Internal Representation

A standard Y -layer feed-forward networks consists of an input layer, ($Y-2$) hidden layers, and an output layer of units successively connected in a feed-forward fashion with no connections between units in the same layer and no feedback connections between layers. The classical single-layer perceptron (two-layer BP network in our context), given two classes of patterns, attempts to find a linear decision boundary separating the two classes. If the two sets of patterns are

$c \setminus h$	1	2	3	4	5	6
	+	+	--	+	+	--
T	--	-	--	--	+	+
+	+	+	--	--	--	-
D	+	--	+	+	+	--
C	+	+	+	+	+	+
S	+	--	+	+	--	--

Fig. 14. Output of the hidden units (+: around 1.0; --: around -1.0; + -: around 0.5; -: around -0.5; c: category; h: no. of the hidden unit).

linearly separable, the perceptron algorithm is guaranteed to find a separating hyperplane in a finite number of steps. It is commonly accepted that a single layer perceptron is inadequate for situations with multiple classes and nonlinear separating boundaries. Hence, the multilayer perceptron network (MLP) was proposed. The MLP net can be viewed as a nonlinear input-output mapping, and the learning process can be seen as fitting a function to the given data set.

A neural network is widely regarded as a black box that reveals little about its predictions. However, analysis of the road classification network RCN with 32 input units and six hidden units reveals some properties of the internal representation, especially the role of the hidden units. We developed the concept in [6] to define the receptive fields of hidden units as the distribution the connection weights from all the input units to each hidden units and the integrating fields as the distribution of the connection weights from all the hidden units to each of the output units. Roughly speaking, each hidden unit extracts some kind of features from the input units through the corresponding receptive field, and the integrating fields integrate several features to conclude the final recognition. It is difficult to describe the mechanism clearly and needs further study. Recent works [18], [19] show that rules can be extracted from ANN's. In our experiments we found that the most of the outputs of the hidden units are saturated (near + or near -1; see Fig. 14). Since we use TanH as the transfer function, it means that some kinds of features are detected (+1) or not detected (-1). The unstable outputs (with absolute values less than 0.5) of the hidden units means that the network is not certain about those features.

Experiments also show that some of the input units are not very important so that pruning them does not affect the recognition results very much. Similarly, some of the hidden units (features), for example the second hidden unit, are not vital for road classification. We found that the RCN, virtually a distributed processing system, still works in case of the

disability or the injury of some parts of the system. The physical meaning of the disability of the input units may be the partial occlusion of the scene by other objects, the partial changes of the environments, or the injury of the “eyes.” Similarly, the disability of the hidden units means the injury of some part of the “brain.” The disability of more important inputs and/or hidden units (e.g., the third one) have more negative effects to the performance of the network, but effects are only obvious to some of the categories (e.g., “S”).

Analysis of the RON’s weights reveals that the operations of the two-layered orientation networks are virtually correlation functions. The Q -dimension weight vector between Q input units and the k th output unit is nearly a circular k -shift version of the representative vector of the input for the given class of the RON. The weight vector corresponding to each output unit is compared with the input vector. The best match indicates the correct orientation of the input vector. In other words, the circular-shifted version of the orientational projection for a given road class are almost linear separable.

The reason we use RON’s to estimate the road orientation instead of direct correlation is that the neural networks can learn the best templates (the representative vector) from the training examples.

E. Practical Considerations

ANN’s are essentially massive parallel computing systems consisting of an extremely large number of simple processors with many interconnections. State-of-the-art computer technology such as VLSI and optics has made this possible. The computation requirements of road scene understanding using the RUNN in the current available serial computer consist of the following four steps.

- 1) *Polar transformation and projection in the PIPE.* The time complexity is $O(RN)$ where R is the radius dimension and N is the orientation dimension of the polar OVI. These two geometrical transformations can be realized in the PIPE in real-time when R and N are 80 and 360, respectively.
- 2) *1-D Fourier transform in PC.* The computation complexity is $O(N \log_2 N)$ using 1-D fast Fourier transform where N is the number of the orientations. This takes about 1.0 s in PC system (CPU 486/66M Hz) when $N = 360$ and takes less than 0.1 s when $N = 32$.
- 3) *Road classification using the RCN.* The time complexity is $O(hP + hM)$ where $P, h,$ and M are the number of units in the input layer, hidden layer and output layer of the RCN, respectively. This number is 424 when $P = 16, h = 12,$ and $M = 6$ and it takes about 0.1 s from the input to the output by using Nworks.
- 4) *Road orientation estimation using the RON.* The time complexity is $O(N^2)$ where N is both the number of units in the input layer and output layer of the RON. This number is 1024 when $N = 32$ and it takes about 0.2 s from the input to the output by using Nworks.

In the current experiments of the RUNN training and testing, the orientational projection is resampled from 360 to 32 since the software Nworks running on a PC486 is

much slow with large number of PE’s in the training process. Correspondingly, the number of the outputs of the RON, the estimated orientations, is sampled to 32 in our principle experimental study. It means the angle resolution is about 11° for the 360° view, which can not meet the practical scene requirements. This problem could be solved by using the $N = 360$ inputs of original orientational projection to the RON’s. Experiments with 360-sampled orientational projections shows that the orientation difference given by the CCF method is the average of the orientation difference between the front roads and the orientation difference between the back roads in the two images (Fig. 8). It means that acceptable results can be produced by the correlation approach. Since the RON acts as a correlation operation, the trained RON with high angular resolution should be better than the fixed CCF method.

Since the energy of Fourier spectrum compacts to the first few items, so the number of the inputs of the RCN, P , is much smaller than N . Experiments show that $P = 30$ is enough to represent each category when N is 360 (refer to Fig. 13), so the scale of the RCN does not increase with N if we used Fourier amplitudes for road classification. Even if the ACF is used, the number of the input data is $N/2$ since ACF is symmetric about $N/2$. However if we expect 1° angular resolution for the orientation estimation, the number of both input and output units, Q and L , should be 360. It means there are N^2 (about 130 K if $N = 360$) connections and weights for each of the RON’s. In the initial stage of the operation or the recovery stage when the robot has missed the road, the system should search for all the 360 directions. Fortunately, in the normal conditions of continuous road following, the system only needs to search for a narrow range of orientations, for example, from -16° to 16° with 1° interval. So only 33 outputs and the corresponding connections are activated, and most of the outputs and the corresponding connections could be disabled. Therefore even the RUNN is simulated in the general von Neumann computer, the real-time computation is possible for real applications if a faster PC system (i.e., Pentium II/266M Hz) is used.

VII. CONCLUSION AND DISCUSSION

In this paper, we present the experimental results of training and testing the BP network for the outdoor road scene understanding using OVI’s. Both the road orientations used for robot heading and the road categories used for robot localization are determined through the integration of invariant image analysis and adaptive neural networks. Several design issues, including the network model, the selection of input data, the number of the hidden units and the learning problems are studied. The internal representations of the networks are carefully analyzed, which could guide further research and applications. Experimental results with real scene images are promising. In order to actually apply the neural networks to real-world autonomous robot navigation in the outdoor natural scene environment, the following aspects are in consideration.

A. Using 2-D Image Patterns

Keeping the requirement of extracting rotation invariant image features in mind, and at the same time reducing the

dimension of the original 2-D OVI's, we are investigating the possibility of using the principal component analysis—the Karhunen–Loeve transform (KLT) along the radius axis r of the omni-view polar image $I(r, \theta)$ to obtain eigenfeatures for a given θ

$$\mathbf{U}(\theta) = (u_0(\theta), u_1(\theta), \dots, u_S(\theta)) \quad (20)$$

where S is expected to be much smaller than the original dimension of r . Preliminary experiments show that the first three components of KL coefficients of the 1-D radius image can properly represent the original 2-D OVI with r from zero to 127. The sampled eigenfeature sequence along the orientation direction can be expressed as

$$U = \{\mathbf{U}(n) = (u_0(n), u_1(n), \dots, u_S(n)), n = 0, \dots, N - 1\}. \quad (21)$$

Sequence $\{u_0(n)\}$, the first component of $\{\mathbf{U}(n)\}$, is approximately the orientational projection define in this paper. The rotation-invariant vector sequence can be obtained by applying the 1-D DFT to each component sequence. The computation complexity for 1-D KLT-DFT method is only $O(SN(R + \log 2N))$ where S is the dimension of the eigenfeature, N is the number of orientations and R is the radius of the polar image.

B. Using Temporal Coherence

A spatio-temporal pattern recognition (SPR) network was proposed by Grossberg [20], [21] to explain certain cognitive process for recognizing sequence of events. The primary application of SPR networks appears to be in the area of recognizing repetitive audio signals. It is straightforward to apply the SPR network for recognizing image sequence of outdoor road scene. Since the same road category will last for a period of time, SPR network should not be sensitive to the occasional image events and could give a robust recognition.

C. Self-Organization and Unsupervised Learning

The natural extension of the work is to use the unsupervised self-organization neural network [22]. When the robot has enough ability to travel around the known world, we can expect that it can also explore the unknown world by itself.

D. Multiple Sensor Fusion and Integration

Visual navigation of a mobile robot in the natural environment is a difficult and comprehensive subject, which is related to almost every aspects of computer vision researches. The fundamental tasks of visual navigation are composed of global localization, road following and obstacle detection. Environment modeling is the foundation of visual navigation. A task-oriented, multiscale and full-view scene modeling strategy is proposed for visual navigation in natural environment [12]. It combines the panoramic vision for scene modeling, omnidirectional vision for road understanding and binocular vision for obstacle detection into an integrated system. This approach overcomes the drawbacks of traditional visual navigation methods that mainly depended on local and/or single view visual information.

ACKNOWLEDGMENT

The authors are grateful to the anonymous reviewers of this paper for their valuable comments and suggestions.

REFERENCES

- [1] R. T. Elkins and E. L. Hall, "Three dimensional line following using omnidirectional vision," *Proc. SPIE*, vol. 2354, pp. 130–144, 1994.
- [2] Y. Yagi, S. Kawato, and S. Tsuji, "Real-time omnidirectional image sensor (COPIS) for vision-guided navigation," *IEEE Trans. Robot. Automat.*, vol. 1, pp. 11–27, Feb. 1994.
- [3] J. Hong, *et al.*, "Image-based homing," in *Proc. Int. Conf. Robotics and Automation*, Apr. 1991, pp. 620–625.
- [4] F. Stein and G. Medioni, "Map-based localization using the panoramic horizon," in *Proc. Int. Conf. Robotics and Automation*, 1992.
- [5] Y. T. Zhou and R. A. Chellappa, "A network for motion perception," in *Proc. IJCNN*, 1990, pp. II841–851.
- [6] E. Atsumi *et al.*, "Internal representation of a neural network that detects local motion," in *Proc. IJCNN*, 1993, pp. 198–201.
- [7] D. A. Pomerleau, "Neural network based autonomous navigation," in *Vision and Navigation: The CMU Navlab*. Boston, MA: Kluwer, 1990.
- [8] T. Jochem, "Using virtual active tools to improve autonomous driving tasks," Tech. Rep., Carnegie Mellon Univ., Pittsburgh, PA, , Nov. 1993.
- [9] ———, D. A. Pomerleau, and C. Thorpe, "Vision guided lane transition," *IEEE Symp. on Intelligent Vehicles*, Detroit, MI, Sept. 25–26, 1995.
- [10] Z. G. Zhu and G. Y. Xu, "Neural networks for omni-view road image understanding," *J. Comput. Sci. Technol.*, vol. 11, no. 4, July 1996.
- [11] Z. G. Zhu, H. J. Xi, and G. Y. Xu, "Combining rotation-invariance images and neural networks for road scene understanding," *Proc. IEEE ICNN*, 1996, pp. 1732–1737.
- [12] Z. G. Zhu, "Environment modeling for visual navigation," Ph.D. dissertation, Tsinghua Univ., Beijing, China, May 1997.
- [13] D. E. Rumelhart, B. Widrow, and M. A. Lehr, "The basic ideas in neural networks," *Commun. ACM*, vol. 37, pp. 87–91, Mar. 1994.
- [14] R. Linggard *et al.*, Ed., *Neural Networks for Vision, Speech and Natural Language*. London, U.K.: Chapman & Hall, 1992.
- [15] A. K. Jain, J. Mao, and K. M. Mohiuddin, "Artificial neural networks: A tutorial," *IEEE Trans. Comput.*, vol. 29, pp. 31–44, Mar. 1996.
- [16] Neuralware Inc., *Neural Computing: Using Nworks of Professional II/Plus*, 1991.
- [17] Aspx Corporation, *The PIPE Reference Manual*, 1989.
- [18] G. G. Towell and J. W. Shariik, "Extracting refined rules from knowledge-based neural networks," *Mach. Learn.*, vol. 13, pp. 71–101, Oct. 1993.
- [19] R. Setiono and H. Liu, "Symbolic representation of neural networks," *IEEE Trans. Comput.*, vol. 29, pp. 71–78, Mar. 1996.
- [20] S. Grossberg, "Some networks that can learn, remember and reproduce any number of complicated space-time pattern, I," *J. Math. Mechan.*, vol. 19, pp. 53–91, 1969.
- [21] ———, "Some networks that can learn, remember and reproduce any number of complicated space-time pattern, II," *Stud. Appl. Math.*, vol. 49, pp. 135–166, 1970.
- [22] T. Kohonen, T., *Self-Organization and Associative Memory*. Berlin, Germany: Springer-Verlag, 1984.



Zhigang Zhu (S'93–M'98) was born in Shanxi, China, in 1964. He received the B.S., M.S., and Ph.D. degrees in computer science and technology from Tsinghua University, Beijing, China, in 1988, 1991, and 1997 respectively.

In 1991, he joined the faculty of the Department of Computer Science and Technology, Tsinghua University, where he is currently Associate Professor. Since 1997, he has been Director of the Information Processing and Application Division.

His research interests include computer vision, image and video analysis, intelligent mobile robots, and virtual reality.

Dr. Zhu received the Science and Technology Achievement Award (second-class) from the Ministry of Electronic Industry, China, in 1996, the Outstanding Young Teacher Award from Beijing Administration in 1997, and the C. C. Lin Applied Mathematics Award (first prize winner) from Tsinghua University in 1997.



Shiqiang Yang received the M.S. degree in 1983.

He is currently Associate Professor in the Department of Computer Science and Technology, Tsinghua University, Beijing, China. His research interests include neural networks, image coding, multimedia system, human-computer interaction, and parallel processing in computer vision.



Guangyou Xu (SM'91) was born in Shanghai, China. He graduated from the Department of Automatic Control Engineering, Tsinghua University, Beijing, China, in 1963.

He was an Assistant Professor with the Department of Automatic Control at Tsinghua University, from 1963 to 1978. He was a visiting scholar at the School of Electrical Engineering, Purdue University, West Lafayette, IN, from December 1982 to December 1984, working with Prof. K. S. Fu. From 1978 to 1986, he was a Lecturer with the Department of Electronic Engineering, Tsinghua University. He was an Associate Professor with the Department of Computer Science and Technology, Tsinghua University, from 1986 to 1989. Since 1989, he has been a Professor with the Department of Computer Science and Technology at Tsinghua University. He served as Director of Information Processing and Application Division there from 1986 to 1997. He was a Visiting Professor at the Beckman Institute, University of Illinois at Urbana-Champaign, from December 1993 to June 1994. His current research interests include artificial intelligence, computer vision, image processing, and multimedia information systems. He has published numerous technical papers and books. He is the author of *Floppy Disc Driver* (Railroad: 1986, in Chinese), *Artificial Intelligence and Its Application* (with K. S. Fu and Z. X. Cai, Tsinghua University Press, 1987), *Multimedia Technology and Systems* (Railroad: 1993), and *Personal Multimedia Computer* (Post: 1995).

Prof. Xu was appointed as a member of Machine Vision Expert Group in the field of Automation by the National High-Tech Project Committee in 1990. He is a member of Robot Measurement Committee TC-17 of Imeko and a standing member of the Council of China Image and Graphic Association. He has received a number of science and technology awards, including the second-class award of the National Science and Technology Progress Prize of China (1986), and four second-class award of the Research Subject Prize issued by the National High-Tech Project Committee in the field of Automation, State Education Committee, Ministry of Electronic Industry China, respectively, from 1992–1996, and Guanghua Science and Technology Achievement Award (1997).



Xueyin Lin received the B.S. degree in automatic control from Tsinghua University, Beijing, China, in 1962.

In 1962, he joined the Department of Automatic Control, Tsinghua University. Since then he was a Assistant Professor, Lecturer and Associate Professor. He is currently Professor in the Department of Computer Science and Technology, Tsinghua University. He was a visiting scholar at the Department of Electrical Engineering, University of Cincinnati, Cincinnati, OH, from April 1983 to July 1985 and from March 1993 to June 1994. His research interests include image processing, computer vision, and pattern recognition. He has published numerous technical papers and books.



Dingji Shi was born in Nanjing, China, in 1933. He graduated from the Department of Electrical Engineering, Tsinghua University, Beijing, China, in 1954.

He has been teaching at Tsinghua University since that time, and is currently a Professor with the Department of Computer Science and Technology. His current scientific activity is in the area of computer vision and its application.