

Reprinted from

# Robotics and Autonomous Systems

---

Robotics and Autonomous Systems 12(1994) 113–119

## Reliability estimation for neural network based autonomous driving

Dean A. Pomerleau

*School of Computer Science, Carnegie Mellon University, Pittsburgh PA 15213, USA*



# Robotics and Autonomous Systems

## Editors-in-Chief:

Prof. F.C.A. Groen, University of Amsterdam, Faculty of Mathematics and Computer Science, Dept. of Computer Systems, Kruislaan 403, 1098 SJ Amsterdam, The Netherlands; E-mail: groen@fwi.uva.nl

## Founding Editor

Dr. T.M. Knasel, 10324 Lake Avenue, Cleveland, OH 44102-1239, USA; Fax: (216) 651-5136; Tel.: (216) 962-3040

Prof. T.C. Henderson, University of Utah, Dept. of Computer Science, 3160, Merrill Engineering Building, Salt Lake City, UT 84112, USA; E-mail: tch@cs.utah.edu

---

## EDITORIAL BOARD

### USA/Canada

J. ALBUS  
Chair, Robot Systems Division  
Nat. Inst. Standards and Technol.  
Bldg. 220, Room 8124  
Gaithersburg, MD 20899  
USA

CDR H.R. EVERETT Jr.  
Naval Command Control &  
Ocean Surveillance Center  
RDTBE Division  
Code 5303  
San Diego, CA 92152-7381  
USA

G.C. GOLDBOGEN  
Columbia College  
Academic Computing Dept  
600 S. Michigan Ave.  
Chicago, IL 60605  
USA

W. GRUVER  
Simon Fraser University  
School of Engineering Science  
Burnaby, BC  
Canada V5A 1S6

S.A. HAYATI  
Jet Propulsion Laboratory  
4800 Oak Grove Drive, MS198-219  
Pasadena, CA 91109  
USA

AKAK  
Purdue University  
School of Elect. Eng.  
West Lafayette, IN 47907  
USA

L. KANAL  
University of Maryland  
Dept. of Computer Science  
College Park, MD 20742  
USA

D. NITZAN  
Director of Robotics  
Stanford Research Institute  
333 Ravenswood Ave.  
Menlo Park, CA 94025  
USA

### Europe

M. BRADY  
University of Oxford  
Dept. of Engineering Science  
Parks Road  
Oxford OX1 3PJ  
UK

J.L. CROWLEY  
LIFIA (INPG)  
46, Avenue Félix Viallet  
38031 Grenoble  
France

O.D. FAUGERAS  
Computer Vision and Robotics Lab.  
INRIA - SOPHIA  
2004 route des Lucioles  
06902 Valbonne Cedex  
France

A.P. FOTHERGILL  
University of Aberdeen  
King's College  
Dept. of Computer Science  
Old Aberdeen, AB9 2UB  
UK

LO. HERTZBERGER  
University of Amsterdam  
Dept. of Computer systems  
Kruislaan 403  
1098 SJ Amsterdam  
The Netherlands

G. HONDERD  
Technische Universiteit Delft  
Fac. der Electrotechniek  
Mekelweg 4  
2628 CD Delft  
The Netherlands

H.-H. NAGEL  
Fraunhofer-Institut für Informations-  
und Datenverarbeitung (ITB)  
Fraunhoferstr. 1  
7500 Karlsruhe 1  
Germany

C. PETERSSON  
Sensor Control AB  
Pilgatan 8  
572 1 30 Västerås  
Sweden

E.A. PUENTE  
Univ. Politécnica de Madrid  
Dept. de Ing. de Sis. y Autom  
C./J. Gutiérrez Abascal, 2  
Madrid-6  
Spain

U. REMBOLD  
Universität Karlsruhe  
Institut für Informatik III  
Postfach 6960  
7500 Karlsruhe 1  
Germany

A. STEIGER GARÇAO  
UNL FCT  
Departamento de Informática  
Quinta da Torre  
2825 Monte da Caparica  
Portugal

W. VAN DE VELDE  
Vrije Universiteit Brussels  
AI Laboratory  
Bldg. K Floor 4  
Pleinlaan 2  
1050 Brussels  
Belgium

### Asia

J.H. KIM  
KAIST  
Computer Science Dept.  
3N3-1 Kusong-Dong, Yusong-Gu  
Taejeon 305-N01  
Korea

H. MAKINO  
Yamanashi University  
Faculty of Engineering  
Takeda 4-3-11, Kofu  
Yamanashi  
Japan

Y. SHIRAI  
Osaka University  
Faculty of Engineering  
Dept. of Mechanical Eng. for  
Computer Controlled Machinery  
2-1, Yamadaoka, Suitasi  
Japan

H. YOSHIKAWA  
University of Tokyo  
Dept. of Precision Machinery Eng  
Faculty of Engineering  
7-3-1 Hongo, Bunkyo-ku  
Tokyo  
Japan

# Reliability estimation for neural network based autonomous driving

Dean A. Pomerleau

*School of Computer Science, Carnegie Mellon University, Pittsburgh PA 15213, USA*

---

## Abstract

This paper describes a technique called *Input Reconstruction Reliability Estimation* (IRRE) for determining the response reliability of a restricted class of multi-layer perceptrons (MLPs). The technique uses a network's ability to accurately encode the input pattern in its internal representation as a measure of its reliability. The more accurately a network is able to reconstruct the input pattern from its internal representation, the more reliable the network is considered to be. IRRE provides a good estimate of the reliability of MLPs trained for autonomous driving. Results are presented in which the reliability estimates provided by IRRE are used to select between networks trained for different driving situations.

*Key words:* Neural networks; Reliability estimation; Autoassociator; Autonomous driving

---

## 1. Introduction

In many real-world domains it is important to know the reliability of a network's response since a single network cannot be expected to accurately handle all the possible inputs. Ideally, a network should not only provide a response to a given input pattern, but also an indication of the likelihood that its response is 'correct'. This reliability measure could be used to weight the outputs from multiple networks and to determine when a new network needs to be trained.

This paper describes a technique for estimating a network's reliability called *Input Reconstruction Reliability Estimation* (IRRE). IRRE relies on the fact that the hidden representation developed by an artificial neural network can be considered to be a compressed representation of

important input features. For example, when the network shown in Fig. 1 is trained to produce the correct steering direction from images of the road ahead, the hidden units learn to encode the position and orientation of important features like the road edges and lane markers (see [8] for more details). Because there are many fewer hidden units than input units in the network, the hidden units cannot accurately represent all the details of an arbitrary input pattern. Instead, the hidden units learn to devote their limited representational capabilities to encoding the position and orientation of consistent, frequently-occurring features from the training set. When presented with an atypical input, such as a road with a different number of lanes, the feature detectors developed by the hidden units will not be capable of accurately encode all the actual input features.

Input Reconstruction Reliability Estimation exploits this limitation in representational capacity to estimate a network's reliability. In IRRE, the network's internal representation is used to reconstruct the input pattern being presented. The more closely the reconstructed input matches the actual input, the more familiar the input and hence the more reliable the network's response.

## 2. Reconstructing the input

IRRE utilized an additional set of output units to perform input reconstruction called the encoder output array, as depicted in Fig. 2. This second set of output units has the same dimensionality as the input retina. In the experiments described in this paper, the input layer and encoder output array have 30 rows and 32 columns. The desired activation for each of these additional output units is identical to the activation of the corresponding input unit. In essence, these additional output units turn the network into an autoencoder.

The network is trained using backpropagation both to produce the correct steering response on the steering output units, and to reconstruct the input image as accurately as possible on the en-

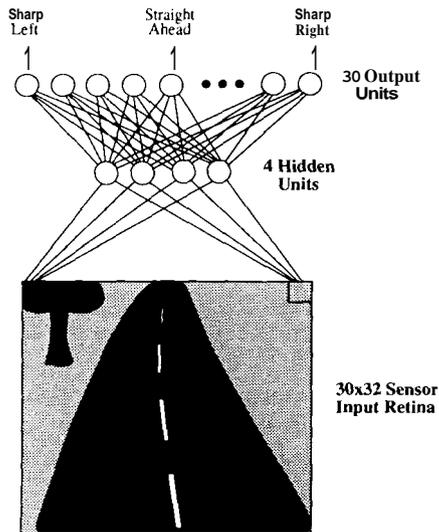


Fig. 1. Original driving network architecture.

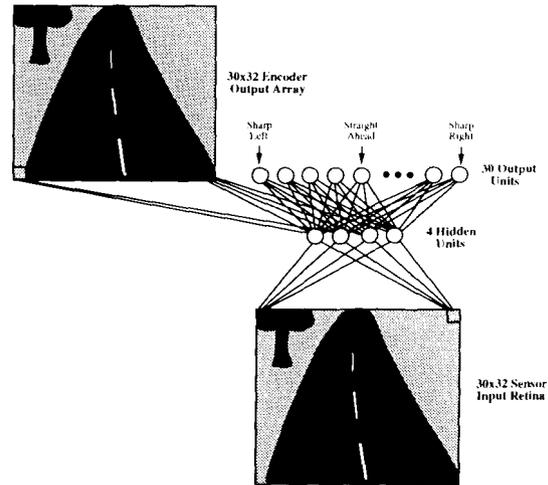


Fig. 2. Network architecture augmented to include an encoder output array.

coder output array. During the training process, the network is presented with several hundred images taken with a camera onboard our test vehicle as a person drives (see [8] for more details). Training typically requires approximately 3 minutes during which the person drives over a 1/4 to 1/2 mile stretch of road.

During testing on a new stretch of road, images are presented to the network and activation is propagated forward through the network to produce a steering response and a reconstructed input image. The reliability of the steering response is estimated by computing the correlation coefficient  $\rho(I, R)$  between the activation levels of units in the actual input image  $I$  and the reconstructed input image  $R$  using the following formula:

$$\rho(I, R) = \frac{\overline{IR} - \bar{I} \cdot \bar{R}}{\sigma_I \sigma_R}$$

where  $\bar{I}$  and  $\bar{R}$  are the mean activation value of the actual and the reconstructed images,  $\overline{IR}$  is the mean of the set formed by the unit-wise product of the two images, and  $\sigma_I$  and  $\sigma_R$  represent the standard deviations of the activation values of each image. The higher the correlation between the two images, the more reliable the network's response is estimated to be. The reason correla-

tion is used to measure the degrees of match between the two images is that, unlike Euclidean distance, the correlation measure is invariant to differences in the mean and variance between the two images. This is important since the mean and variance of the input and the reconstructed images can sometimes vary, even when the input image depicts a familiar situation.

### 3. Results and applications

The degree of correlation between the actual and the reconstructed input images is an extremely good indicator of network response accuracy in the domain of autonomous driving, as shown in Fig. 3. It shows a trained network's steering error and reconstruction error as the vehicle drives down a quarter mile stretch of road that starts out as a single lane path and eventually becomes a two-lane street. The solid line indicates the network's steering error, as measured by the difference in turn curvature between the network's steering response and a person's steering response at that point along the road. The dashed line represents the network's 'reconstruction error', which is defined to be the degree

of static independence between the actual and reconstructed images, or  $1 - \rho(I, R)$ .

The two curves are nearly identical, having a correlation coefficient of 0.92. This close match between the curves demonstrates that when the network is unable to accurately reconstruct the input image, it is also probably suggesting an incorrect steering direction. Visual inspection of the actual and reconstructed input images demonstrates that the degree of resemblance between them is a good indication of the actual input's familiarity, as shown in Fig. 4. It depicts the input image, network steering response, and reconstructed input at the three points along the road, labeled A, B and C in Fig. 3. When presented with the image at point A, which closely resembles patterns from training set, the network's reconstructed image closely resembles the actual input, as shown by the close correspondence between the images labeled 'Input Acts' and 'Reconstructed Input' in the left column of Fig. 4. This close correspondence between the input and reconstructed images suggests that the network can reliably steer in this situation. It in fact can steer accurately on this image, as demonstrated by the close match between the network's

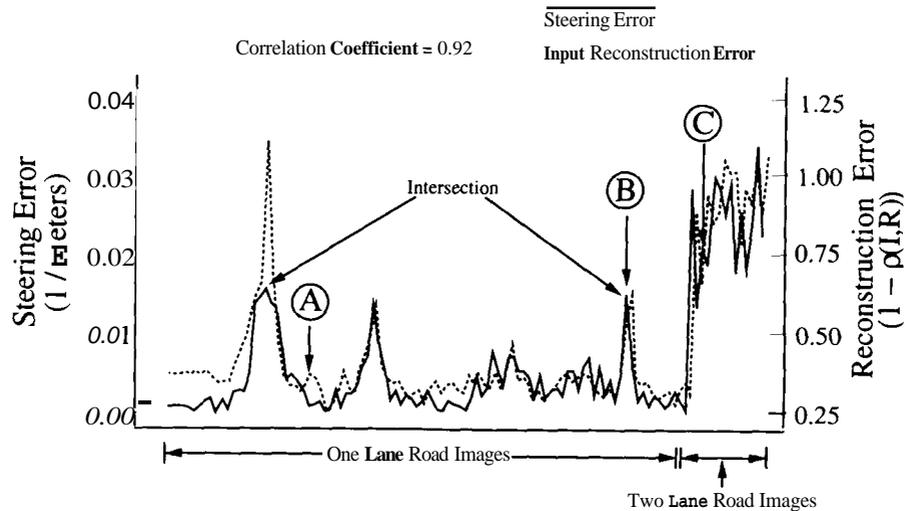


Fig. 3. Reconstruction error obtained using autoencoder reconstruction versus network steering error over a stretch of one-lane and two-lane road.

steering response labeled 'Output Acts' and the desired steering response labeled 'Target Acts' in the upper left corner of Fig. 4.

When presented with a situation the network did not encounter during training, such as the

fork image and the two-lane road image shown in the other two columns of Fig. 4, the reconstructed image bears much less resemblance to the original input. This suggests that the network is confused. This confusion results in an incorrect

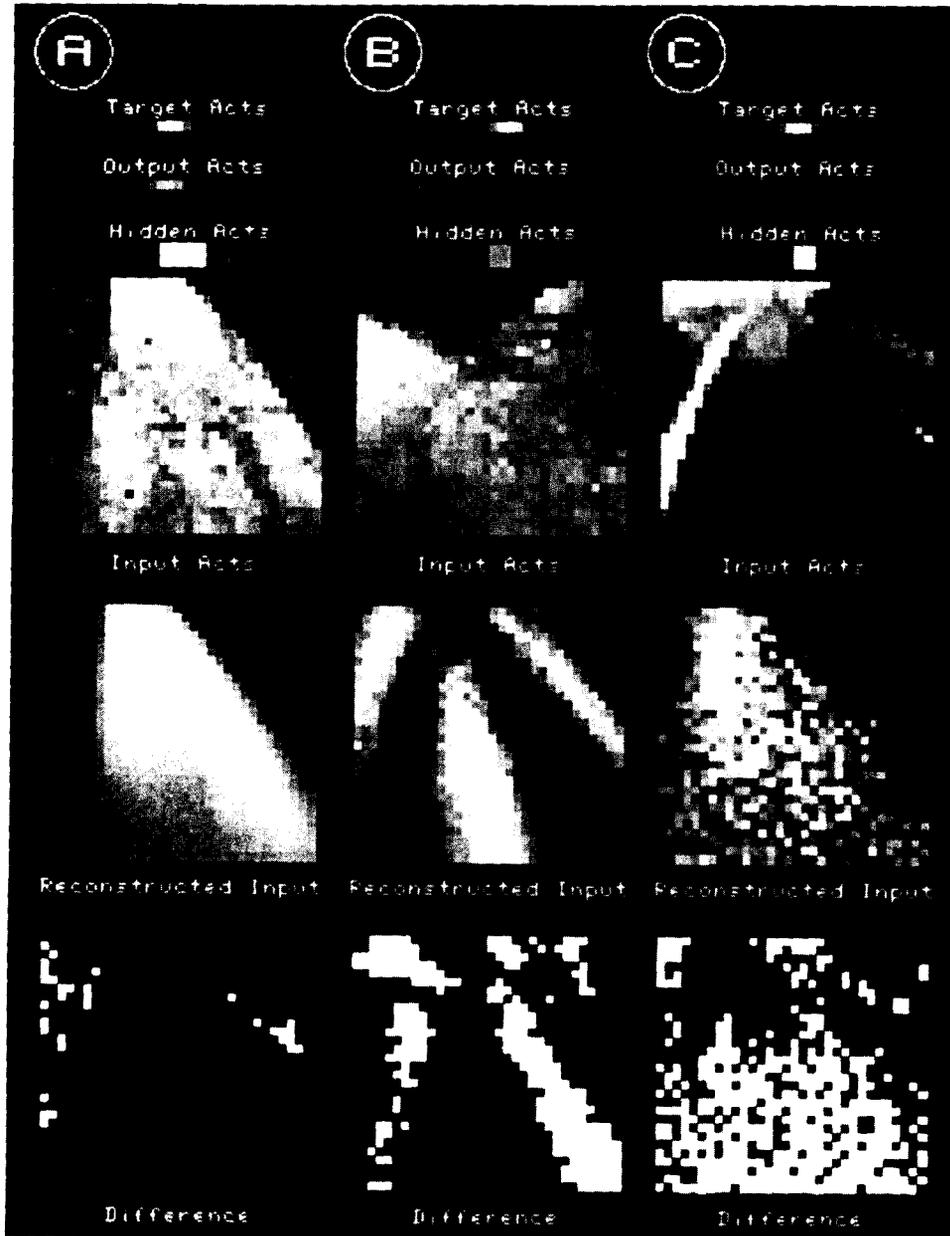


Fig. 4. The actual input, the reconstructed input and the point-wise absolute difference between them on a road image similar to those in the training set (labeled **A**), and on two atypical images (labeled **B** and **C**).

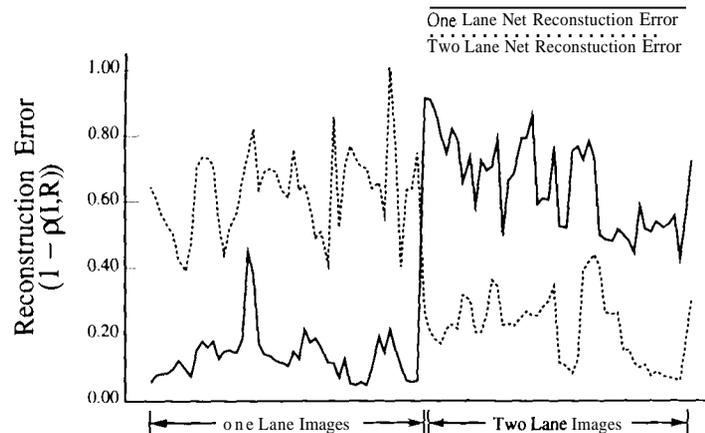


Fig. 5. Reconstruction error of networks trained for one-lane road driving (solid line) and two-lane road driving (dashed line)

steering response, illustrated in the discrepancy between the network's steering response and the target steering response for the two atypical images.

The reliability prediction provided by IRRE has been used to improve the performance of the neural network based autonomous driving system in a number of ways. The simplest is to use IRRE to control vehicle speed. The more accurate the input reconstruction, the more confident the network, and hence the faster the system drives the vehicle. A second use the system makes of the reliability estimate provided by IRRE is to update the vehicle's position on a rough map of the area being driven over. When the map indicates there should be an intersection or other confusing situation up ahead, a subsequent sharp rise in reconstruction error is a good indication that the vehicle has actually reached that location, allowing the system to pinpoint the vehicle's position. Knowing the vehicle's location on a map is useful for integrating symbolic processing, such as planning, into the neural network driving system (for more details, see [9]).

Fig. 5 illustrates how IRRE can be used to integrate the outputs from multiple expert networks. The two lines in this graph illustrate the reconstruction error of two networks, one trained to steer on one-lane roads (solid line), and the other trained to steer on two-lane roads (dashed line). The reconstruction error for the one-lane

network is low on the one-lane road images, and high on the two-lane road images. The opposite is true for the network trained for two-lane road driving. The reliability estimate provided by IRRE allows the system to determine which network is most reliable for driving in the current situation. By simulating multiple networks in parallel, and then selecting the one with the highest reliability, the system have been able to drive the Navlab test vehicle on one, two and four lane roads automatically at speeds of up to 55 miles per hour.

#### 4. Discussion

The effectiveness of input reconstruction reliable estimation stems from the fact that the network has a small number of hidden units and is only trained in a narrow range of situations. These constraints prevent the network from faithfully encoding arbitrary input patterns. Instead, the hidden units learn to encode features in the training images that are most important for the task. Baldi and Hornik [1] have shown that if an autoencoder network with a single layer of  $N$  linear hidden units is trained with back-propagation, the activation levels of the hidden units will represent the first  $N$  principal components of the training set. Since the units in the driving network are non-linear, this assertion does not strictly

hold in this case. However, Cottrell and Munro [2] have found empirically that autoencoder networks with a sigmoidal activation function develop hidden units that span the principal subspace of the training images, with some noise on the first principal component due to network non-linearity. Because the principal components represent the dimensions along which the training examples varies most, it can be shown that using linear combinations of the principal components to represent the individual training patterns optimally preserves the information contained in the training set [7].

However, the compressed representation developed by a linear autoencoder network is only optimal for encoding images from the same distribution as the training set. When presented with images very different from those in the training set, the image reconstructed from the internal representation is not as accurate. The results presented in this paper demonstrate that this reconstruction error can be employed to estimate the likelihood and magnitude of error in MLPs trained for autonomous driving.

However, the input reconstruction technique presented here has a serious potential shortcoming, namely that it forces the network's hidden units to encode all input features, including potentially irrelevant ones. While this increased representation load on the hidden units has the potential to degrade network performance, this effect has not been observed in the tests conducted so far. In support of this finding, Cluck [3] has found that forcing a network to autoencode its input frequently improves its generalization. In [4], Gluck and Myers use the representation developed by an autoencoder network as a model for simple types of learning in biological systems. The model suggests that the hippocampus acts as an autoencoder, developing internal representations that are then used to perform other tasks.

But if interference from the autoencoder task proves to be a problem, one way to eliminate it would be to have separate groups of hidden units connected exclusively to one group of outputs or the other. Having a separate set of hidden units for the autoencoder task would ensure that the representation developed for the input recon-

struction does not interfere with representation developed for the 'normal' task. It remains to be seen if this decoupling of internal representations will adversely affect IRRE's ability to predict network errors.

As a technique for multi-network integration, IRRE has several advantages over existing connectionist arbitration methods, such as Hampshire and Waibel's Meta-Pi architecture [5] and the Adaptive Mixture of Experts Model of Jacobs et al. [6]. It is a more modular approach, since each expert can be trained entirely in isolation and then later combined with other experts without any additional training by simply selecting the most reliable network for the current input. Since IRRE provides an absolute measure of a single network's reliability, and not just a measure of how appropriate a network is relative to others, IRRE can be also used to determine when none of the experts is capable of coping with the current situation.

A potentially interesting extension to IRRE is the development of techniques for reasoning about the difference between the actual input and the reconstructed input. For instance, it should be possible to recognize when the vehicle has reached a fork in the road by the characteristic mistakes the network makes in reconstructing the input image. Another important component of future work is to test the ability of IRRE to estimate network reliability in domains other than autonomous driving.

### Acknowledgements

An earlier version of this paper appears in *Advances in Neural Information Processing Systems 5*, C.L. Giles, S.J. Hanson and J.D. Cowan (eds.), Morgan Kaufmann Publishers, 1993.

I thank Dave Touretzky, Chuck Thorpe and the entire Unmanned Ground Vehicle Group at CMU for their support and suggestions. Principle support for this research has come from ARPA, under contracts 'Perception for Outdoor Navigation' (contract number DACA76-89-C-0014, monitored by the US Army Topographic Engineering Center) and 'Unmanned Ground Vehicle System'

(contract number DAAE07-90-C-R059, monitored by TACOM).

## References

- [1] P. Baldi and K. Hornik. Neural networks and principal component analysis: Learning from examples without local minima. *Neural Networks* 2 (1989) 53–58.
- [2] G.W. Cottrell and P. Munro. Principal components analysis of images via back-propagation. *Proc. Soc. of Photo-Optical Instr. Eng.*, Cambridge, MA (1988).
- [3] M.A. Gluck. Personal communications (Rutgers Univ., Newark, NJ, 1992).
- [4] M.A. Gluck and C.E. Myers. Hippocampal function in representation and generalization: A computational theory. *Proc. Cogn. Sci. Soc. Conf.* (Erlbaum, Hillsdale, NJ, 1992).
- [5] J.R. Hampshire and A.H. Waibel. The Meta-Pi network: Building distributed knowledge representations for robust pattern recognition, *IEEE Trans. on Pattern Analysis and Machine Intelligence* (1989).
- [6] R.A. Jacobs. M.I. Jordan. S.J. Nowlan and G.E. Hinton, Adaptive mixtures of local experts. *Neural Computation* 3 (1) (1991) 79–87.
- [7] R. Linsker. Designing a sensory processing system: What can be learned from principal component analysis? (IBM Technical Report RC14983 (No. 66896), 1989).
- [8] D.A. Pomerleau, Efficient training of artificial neural networks for autonomous navigation. *Neural Computation* 3 (1) (1991) 88–97.
- [9] D.A. Pomerleau, J. Gowdy and C.E. Thorpe. Combining artificial neural networks and symbolic processing for autonomous robot guidance. *Engineering Applications of Artificial Intelligence* 4 (4) (1991) 279–285.



**Dean Pomerleau** received a BA degree from Williams College in 1987 and a Ph D (Computer Science) from Carnegie Mellon University in 1992. Since 1992 he has been a Research Scientist at Carnegie Mellon University, with a joint appointment in the Robotics Institute and the School of Computer Science. His research focuses on connectionist vision techniques, particularly for mobile robot guidance and human computer interaction. He is the director of the Run-

Off-Road Accident Prevention Project, a US Department of Transportation sponsored consortium of universities and industrial partners developing warning systems to prevent highway accidents



# Robotics and Autonomous Systems

## Aims and Scope

*Robotics and Autonomous Systems* will carry articles describing fundamental developments in the field of robotics, with special emphasis on autonomous systems. An important goal of this journal is to extend the state of the art in both symbolic and sensory based robot control and learning in the context of autonomous systems.

*Robotics and Autonomous Systems* will carry articles on the theoretical, computational and experimental aspects of autonomous systems, or modules of such systems.

Application environments of interest include industrial, outdoor, and outer space where advanced robotic techniques are required for autonomous systems to accomplish goals without human intervention; this includes robotics for hazardous and hostile environments.

*Robotics and Autonomous Systems* will carry brief reports on international meetings in the field, as well as an occasional multi-author debate on current topics of interest. Forthcoming meetings of importance will be listed.

In more detail, the journal will cover the following topics:

- symbol mediated robot behavior control
- sensory mediated robot behavior control
- active sensory processing and control
- industrial applications of autonomous systems
- sensor modeling and data interpretation; e.g., models and software for sensor data integration, 3D scene analysis, environment description and modeling, pattern recognition
- robust techniques in AI and sensing; e.g., uncertainty modeling, graceful degradation of systems
- robot programming; e.g., on-line and off-line programming, discrete event dynamical systems, fuzzy logic
- CAD-based robotics; e.g., CAD-based vision, reverse engineering
- robot simulation and visualization tools
- tele-autonomous systems
- micro electromechanical robots
- robot control architectures
- robot planning, adaptation and learning.

## Information for Authors

A detailed Information for Authors brochure is available free of charge or obligation from the Publisher or the Editors-in-Chief. Contributions in final version can be submitted to one of the Editors-in-Chief; for addresses see second cover page.

## Author's Benefits

1. 30% discount on all book publications of North-Holland.
2. 50 reprints are provided free of charge to the principal author of each paper published.

## Advertising Information

Advertising rates and technical requirements are available from the Advertising Manager or from a National Advertising Representative. An additional charge is asked for color printing and for prominent positions.

**Advertising Manager:** Willeke van Cattenburch, Elsevier Science Publishers B.V., P.O. Box 211, 1000 AE Amsterdam, The Netherlands; tel.: (20) 515 3220, telex: 164 79 elsvi nl., fax: (20) 6833041.

Advertising representative for the USA & Canada: Weston Media Associates. Att: Daniel S. Lipner, P.O. Box 1110, Greens Farms, CT 06436-1110, Tel.: (203) 261.2500, Fax: (203) 261.0101.

## Publication information

Robotics and Autonomous Systems (ISSN 0921-8890). For 1994 volumes 13-14 are scheduled for publication. Subscription prices are available upon request from the publisher. Subscriptions are accepted on a prepaid basis only and are entered on a calendar year basis. Issues are sent by surface mail except to the following countries where air delivery via SAL is ensured: Argentina, Australia, Brazil, Canada, Hong Kong, India, Israel, Japan, Malaysia, Mexico, New Zealand, Pakistan, PR China, Singapore, South Africa, South Korea, Taiwan, Thailand, USA. For all other countries airmail rates are available upon request. Claims for missing issues must be made within six months of our publication (mailing) date. Please address all your requests regarding orders and subscription queries to: Elsevier Science Publishers, Journal Department, P.O. Box 211, 1000 AE Amsterdam, The Netherlands. Tel.: +31-20-5803642, fax: +31-20-5803598.

© 1994, Elsevier Science B.V. All rights reserved.

No part of this publication may be reproduced, stored in a retrieval system or transmitted in any form or by any means, electronic, mechanical, photocopying, recording or otherwise, without the prior written permission of the publisher, Elsevier Science B.V., Copyright & Permissions Dept., P.O. Box 521, 1000 AM Amsterdam, The Netherlands.

Special regulations for authors. Upon acceptance of an article by the journal, the author(s) will be asked to transfer copyright of the article to the publisher. The transfer will ensure the widest possible dissemination of information.

Special regulations for readers in the USA - This journal has been registered with the Copyright Clearance Center, Inc. Consent is given for copying of articles for personal or internal use, or for the personal use of specific clients. This consent is given on the condition that the copier pays through the Center the per-copy fee stated in the code on the first page of each article for copying beyond that permitted by Sections 107 or 108 of the US Copyright Law. The appropriate fee should be forwarded with a copy of the first page of the article to the Copyright Clearance Center, Inc., 27 Congress Street, Salem, MA 01970,

USA. If no code appears in an article, the author has not given broad consent to copy and permission to copy must be obtained directly from the author. The fee indicated on the first page of an article in this issue will apply retroactively to all articles published in the journal, regardless of the year of publication. This consent does not extend to other kinds of copying, such as for general distribution, resale, advertising and promotion purposes, or for creating new collective works. Special written permission must be obtained from the publisher for such copying.

No responsibility is assumed by the Publisher for any injury and/or damage to persons or property as a matter of products liability, negligence or otherwise, or from any use or operation of any methods, products, instructions or ideas contained in the material herein.

Although all advertising material is expected to conform to ethical standards, inclusion in this publication does not constitute a guarantee or endorsement of the quality or value of such product or of the claims made of it by its manufacturer.

