

Backpropagation without Multiplication

Patrice Simard, Hans Graf

The back propagation algorithm has been modified to work with out any multiplications and to tolerate computations with a low resolution, which makes it more attractive for a hardware implementation. Numbers are represented in floating point format with 1 bit mantissa and 3 bits in the exponent for the states, and 1 bit mantissa and 5 bit exponent for the gradients, while the weights are 16 bit fixed-point numbers. In this way, all the computations can be executed with shift and add operations. Large networks with over 100,000 weights were trained and demonstrated the same performance as networks computed with full precision. An estimate of a circuit implementation shows that a large network can be placed on a single chip, reaching more than 1 billion weight updates per second. A speedup is also obtained on any machine where a multiplication is slower than a shift operation.

Hoeffding Races: Accelerating Model Selection Search for Classification and Function Approximation

Oded Maron, Andrew Moore

Selecting a good model of a set of input points by cross validation is a computationally intensive process, especially if the number of possible models or the number of training points is high. Techniques such as gradient descent are helpful in searching through the space of models, but problems such as local minima, and more importantly, lack of a distance metric between various models reduce the applicability of these search methods. Hoeffding Races is a technique for finding a good model for the data by quickly discarding bad models, and concentrating the computational effort at differentiating between the better ones. This paper focuses on the special case of leave-one-out cross validation applied to memory based learning algorithms, but we also argue that it is applicable to any class of model selection problems.

Putting It All Together: Methods for Combining Neural Networks

Michael Perrone

Requests for name changes in the electronic proceedings will be accepted with no questions asked. However name changes may cause bibliographic tracking issues.

Authors are asked to consider this carefully and discuss it with their co-authors prior to requesting a name change in the electronic proceedings.

Cross-Validation Estimates IMSE

Mark Plutowski, Shinichi Sakata, Halbert White

Integrated Mean Squared Error (IMSE) is a version of the usual mean squared error criterion, averaged over all possible training sets. If it could be observed, it could be used to determine optimal network complexity or optimal data subsets for efficient training. We show that two common methods of cross-validating average squared error deliver unbiased estimates of IMSE, converging to IMSE with probability one. These estimates thus make possible approximate IMSE-based choice of network work complexity. We also show that two variants of cross validation measure provide unbiased IMSE-based estimates potentially useful for selecting optimal data subsets.

The Power of Amnesia

Dana Ron, Yoram Singer, Naftali Tishby

We propose a learning algorithm for a variable memory length Markov process. Human communication, whether given as text, handwriting, or speech, has multi characteristic time scales. On short scales it is characterized mostly by the dynamics that generate the process, whereas on large scales, more syntactic and semantic information is carried. For that reason the conventionally used fixed memory Markov models cannot capture effectively the complexity of such structures. On

the other hand using long mem(cid:173)ory models uniformly is not practical even for as short memory as four. The algorithm we propose is based on minimizing the sta(cid:173)tistical prediction error by extending the memory, or state length, adaptively, until the total prediction error is sufficiently small. We demonstrate the algorithm by learning the structure of natural En(cid:173)glish text and applying the learned model to the correction of cor(cid:173)rupted text. Using less than 3000 states the model's performance is far superior to that of fixed memory models with similar num(cid:173)ber of states. We also show how the algorithm can be applied to intergenic E. coli DNA base prediction with results comparable to HMM based methods.

On the Non-Existence of a Universal Learning Algorithm for Recurrent Neural Networks

Herbert Wiklicky

We prove that the so called "loading problem" for (recurrent) neural net(cid:173)works is unsolvable. This extends several results which already demon(cid:173)strated that training and related design problems for neural networks are (at least) NP-complete. Our result also implies that it is impossible to find or to formulate a universal training algorithm, which for any neu(cid:173)ral network architecture could determine a correct set of weights. For the simple proof of this, we will just show that the loading problem is equivalent to "Hilbert's tenth problem" which is known to be unsolvable.

Constructive Learning Using Internal Representation Conflicts

Laurens Leerink, Marwan Jabri

We present an algorithm for the training of feedforward and recur(cid:173)rent neural networks. It detects internal representation conflicts and uses these conflicts in a constructive manner to add new neu(cid:173)rons to the network. The advantages are twofold: (1) starting with a small network neurons are only allocated when required; (2) by detecting and resolving internal conflicts at an early stage learning time is reduced. Empirical results on two real-world problems sub(cid:173)stantiate the faster learning speed; when applied to the training of a recurrent network on a well researched sequence recognition task (the Reber grammar), training times are significantly less than previously reported.

A Local Algorithm to Learn Trajectories with Stochastic Neural Networks

Javier Movellan

This paper presents a simple algorithm to learn trajectories with a continuous time, continuous activation version of the Boltzmann machine. The algorithm takes advantage of intrinsic Brownian noise in the network to easily compute gradients using entirely local computations. The algorithm may be ideal for parallel hardware implementations.

Use of Bad Training Data for Better Predictions

Tal Grossman, Alan Lapedes

We show how randomly scrambling the output classes of various fractions of the training data may be used to improve predictive accuracy of a classification algorithm. We present a method for calculating the "noise sensitivity signature" of a learning algorithm which is based on scrambling the output classes. This signature can be used to indicate a good match between the complexity of the classifier and the complexity of the data. Use of noise sensitivity signatures is distinctly different from other schemes to avoid over(cid:173)training, such as cross-validation, which uses only part of the train(cid:173)ing data, or various penalty functions, which are not data-adaptive. Noise sensitivity signature methods use all of the training data and are manifestly data-adaptive and non-parametric. They are well suited for situations with limited training data.

Non-Intrusive Gaze Tracking Using Artificial Neural Networks

Shumeet Baluja, Dean Pomerleau

We have developed an artificial neural network based gaze tracking system which can be customized to individual users. Unlike other gaze trackers, which normally require the user to wear cumbersome headgear, or to use a chin rest to ensure head immobility, our system is entirely non-intrusive. Currently, the best intrusive gaze tracking systems are accurate to approximately 0.75 degrees. In our experiments, we have been able to achieve an accuracy of 1.5 degrees, while allowing head mobility. In this paper we present an empirical analysis of the performance of a large number of artificial neural network architectures for this task.

Learning Curves: Asymptotic Values and Rate of Convergence

Corinna Cortes, L. D. Jackel, Sara Solla, Vladimir Vapnik, John Denker

Training classifiers on large databases is computationally demanding. It is desirable to develop efficient procedures for a reliable prediction of a classifier's suitability for implementing a given task, so that resources can be assigned to the most promising candidates or freed for exploring new classifier candidates. We propose such a practical and principled predictive method. Practical because it avoids the costly procedure of training poor classifiers on the whole training set, and principled because of its theoretical foundation. The effectiveness of the proposed procedure is demonstrated for both single- and multi-layer networks.

How to Describe Neuronal Activity: Spikes, Rates, or Assemblies?

Wulfram Gerstner, J. van Hemmen

What is the 'correct' theoretical description of neuronal activity? The analysis of the dynamics of a globally connected network of spiking neurons (the Spike Response Model) shows that a description by mean firing rates is possible only if active neurons fire incoherently. If firing occurs coherently or with spatio-temporal correlations, the spike structure of the neural code becomes relevant. Alternatively, neurons can be gathered into local or distributed ensembles or 'assemblies'. A description based on the mean ensemble activity is, in principle, possible but the interaction between different assemblies becomes highly nonlinear. A description with spikes should therefore be preferred.

Digital Boltzmann VLSI for constraint satisfaction and learning

Michael Murray, Ming-Tak Leung, Kan Boonayanit, Kong Kritayakirana, James Burg, Gregory Wolff, Tokahiro Watanabe, Edward Schwartz, David Stork, Allen M. Peterson

We built a high-speed, digital mean-field Boltzmann chip and SBus board for general problems in constraint satisfaction and learning. Each chip has 32 neural processors and 4 weight update processors, supporting an arbitrary topology of up to 160 functional neurons. On-chip learning is at a theoretical maximum rate of 3.5×10^8 connection updates/sec; recall is 12000 patterns/sec for typical conditions. The chip's high speed is due to parallel computation of inner products, limited (but adequate) precision for weights and activations (5 bits), fast clock (125 MHz), and several design insights.

Feature Densities are Required for Computing Feature Correspondences

Subutai Ahmad

The feature correspondence problem is a classic hurdle in visual object-recognition concerned with determining the correct mapping between the features measured from the image and the features expected by the model. In this paper we show that determining good correspondences requires information about the joint probability density over the image features. We propose "likelihood based correspondence matching" as a general principle for selecting optimal correspondences. The approach is applicable to non-rigid models, allows nonlinear perspective transformations, and can optimally deal with occlusions and missing features. Experiments with rigid and non-rigid 3D hand gesture r

ecognition support the theory. The likelihood based techniques show almost no decrease in classification performance when compared to performance with perfect correspondence knowledge.

Emergence of Global Structure from Local Associations

Thea Ghiselli-Crippa, Paul Munro

A variant of the encoder architecture, where units at the input and output layers represent nodes on a graph. is applied to the task of mapping locations to sets of neighboring locations. The degree to which the resulting internal (i.e. hidden unit) representations reflect global properties of the environment depends upon several parameters of the learning procedure. Architectural bottlenecks, noise, and incremental learning of landmarks are shown to be important factors in maintaining topographic relationships at a global scale.

Robot Learning: Exploration and Continuous Domains

David Cohn

Requests for name changes in the electronic proceedings will be accepted with no questions asked. However name changes may cause bibliographic tracking issues.

Authors are asked to consider this carefully and discuss it with their co-authors prior to requesting a name change in the electronic proceedings.

Lower Boundaries of Motoneuron Desynchronization via Renshaw Interneurons

Mitchell Maltenfort, Robert Druzinsky, C. Heckman, W. Rymer

Using a quasi-realistic model of the feedback inhibition of motoneurons (MNs) by Renshaw cells, we show that weak inhibition is sufficient to maximally desynchronize MNs, with negligible effects on total MN activity. MN synchrony can produce a 20 - 30 Hz peak in the force power spectrum, which may cause instability in feedback loops.

Catastrophic interference in connectionist networks: Can It Be predicted, can It be prevented?

Robert French

Requests for name changes in the electronic proceedings will be accepted with no questions asked. However name changes may cause bibliographic tracking issues.

Authors are asked to consider this carefully and discuss it with their co-authors prior to requesting a name change in the electronic proceedings.

Signature Verification using a "Siamese" Time Delay Neural Network

Jane Bromley, Isabelle Guyon, Yann LeCun, Eduard Säckinger, Roopak Shah

This paper describes an algorithm for verification of signatures written on a pen-input tablet. The algorithm is based on a novel, artificial neural network, called a "Siamese" neural network. This network consists of two identical sub-networks joined at their outputs. During training the two sub-networks extract features from two signatures, while the joining neuron measures the distance between the two feature vectors. Verification consists of comparing an extracted feature vector with a stored feature vector for the signer. Signatures closer to this stored representation than a chosen threshold are accepted, all other signatures are rejected as forgeries.

Connectionism for Music and Audition

Andreas Weigend

This workshop explored machine learning approaches to 3 topics: (1) finding structure in music (analysis, continuation, and completion of an unfinished piece), (2) modeling perception of time (extrapolation of musical meter, explanation of human data on timing), and (3) interpolation in timbre space.

Central and Pairwise Data Clustering by Competitive Neural Networks

Joachim Buhmann, Thomas Hofmann

Data clustering amounts to a combinatorial optimization problem to reduce the complexity of a data representation and to increase its precision. Central and pairwise data clustering are studied in the maximum entropy framework. For central clustering we derive a set of reestimation equations and a minimization procedure which yields an optimal number of clusters, their centers and their cluster probabilities. A meanfield approximation for pairwise clustering is used to estimate assignment probabilities. A selfconsistent solution to multidimensional scaling and pairwise clustering is derived which yields an optimal embedding and clustering of data points in a d-dimensional Euclidian space.

Processing of Visual and Auditory Space and Its Modification by Experience

Josef Rauschecker, Terrence J. Sejnowski

Requests for name changes in the electronic proceedings will be accepted with no questions asked. However name changes may cause bibliographic tracking issues.

Authors are asked to consider this carefully and discuss it with their co-authors prior to requesting a name change in the electronic proceedings.

A Connectionist Model of the Owl's Sound Localization System

Daniel Rosen, David Rumelhart, Eric Knudsen

,,,e do not have a good understanding of how theoretical principles of learning are realized in neural systems. To address this problem we built a computational model of development in the owl's sound localization system. The structure of the model is drawn from known experimental data while the learning principles come from recent work in the field of brain style computation. The model accounts for numerous properties of the owl's sound localization system, makes specific and testable predictions for future experiments, and provides a theory of the developmental process.

Identifying Fault-Prone Software Modules Using Feed-Forward Networks: A Case Study

N. Karunanithi

Functional complexity of a software module can be measured in terms of static complexity metrics of the program text. Classifying software modules, based on their static complexity measures, into different fault-prone categories is a difficult problem in software engineering. This research investigates the applicability of neural network classifiers for identifying fault-prone software modules using a data set from a commercial software system. A preliminary empirical comparison is performed between a minimum distance based Gaussian classifier, a perceptron classifier and a multilayer layer feed-forward network classifier constructed using a modified Cascade-Correlation algorithm. The modified version of the Cascade-Correlation algorithm constrains the growth of the network size by incorporating a cross-validation check during the output layer training phase. Our preliminary results suggest that a multilayer feed-forward network can be used as a tool for identifying fault-prone software modules early during the development cycle. Other issues such as representation of software metrics and selection of a proper training samples are also discussed.

Event-Driven Simulation of Networks of Spiking Neurons

Lloyd Watts

A fast event-driven software simulator has been developed for simulating large networks of spiking neurons and synapses. The primitive network elements are designed to exhibit biologically realistic behaviors, such as spiking, refractoriness, adaptation, axonal delays, summation of post-synaptic current pulses, and tonic current inputs. The efficient event-driven representation allows large networks to be simulated in a fraction of the time that would be required for a full compartmental-model simulation. Corresponding analog CMOS VLSI circuit primitives have been designed and characterized, so that large-sca

le circuits may be simulated prior to fabrication.

Grammatical Inference by Attentional Control of Synchronization in an Oscillating Elman Network

Bill Baird, Todd Troyer, Frank Eeckman

We show how an "Elman" network architecture, constructed from recurrently connected oscillatory associative memory network modules, can employ selective "attentional" control of synchronization to direct the flow of communication and computation within the architecture to solve a grammatical inference problem. Previously we have shown how the discrete time "Elman" network algorithm can be implemented in a network completely described by continuous ordinary differential equations. The time steps (machine cycles) of the system are implemented by rhythmic variation (clocking) of a bifurcation parameter. In this architecture, oscillation amplitude codes the information content or activity of a module (unit), whereas phase and frequency are used to "software" the network. Only synchronized modules communicate by exchanging amplitude information; the activity of non-resonating modules contributes incoherent crosstalk noise. Attentional control is modeled as a special subset of the hidden modules with outputs which affect the resonant frequencies of other hidden modules. They control synchrony among the other modules and direct the flow of computation (attention) to effect transitions between two subgraphs of a thirteen state automaton which the system emulates to generate a Reber grammar. The internal crosstalk noise is used to drive the required random transitions of the automaton.

Bayesian Self-Organization

Alan L. Yuille, Stelios Smirnakis, Lei Xu

Recent work by Becker and Hinton (Becker and Hinton, 1992) shows a promising mechanism, based on maximizing mutual information assuming spatial coherence, by which a system can self-organize itself to learn visual abilities such as binocular stereo. We introduce a more general criterion, based on Bayesian probability theory, and thereby demonstrate a connection to Bayesian theories of visual perception and to other organization principles for early vision (Atick and Redlich, 1990). Methods for implementation using variants of stochastic learning are described and, for the special case of linear filtering, we derive an analytic expression for the output.

The "Softmax" Nonlinearity: Derivation Using Statistical Mechanics and Useful Properties as a Multiterminal Analog Circuit Element

I. M. Elfadel, J. L. Wyatt, Jr.

We use mean-field theory methods from Statistical Mechanics to derive the "softmax" nonlinearity from the discontinuous winner-take-all (WTA) mapping.

We give two simple ways of implementing "soft max" as a multiterminal network element. One of these has a number of important network-theoretic properties. It is a reciprocal, passive, incrementally passive, nonlinear, resistive multiterminal element with a content function having the form of information theoretic entropy. These properties should enable one to use this element in nonlinear RC networks with such other reciprocal elements as resistive fuses and constraint boxes to implement very high speed analog optimization algorithms using a minimum of hardware.

Robust Parameter Estimation and Model Selection for Neural Network Regression

Yong Liu

In this paper, it is shown that the conventional back-propagation (BPP) algorithm for neural network regression is robust to leverageages (data with x corrupted), but not to outliers (data with y corrupted). A robust model is to model the error as a mixture of normal distribution. The influence function for this mixture model is calculated and the condition for the model to be robust to outliers is given. EM algorithm [5

] is used to estimate the parameter. The usefulness of model selection criteria is also discussed. Illustrative simulations are performed.

Computational Elements of the Adaptive Controller of the Human Arm

Reza Shadmehr, Ferdinando Mussa-Ivaldi

We consider the problem of how the CNS learns to control dynamics of a mechanical system. By using a paradigm where a subject's hand interacts with a virtual mechanical environment, we show that learning control is via composition of a model of the imposed dynamics. Some properties of the computational elements with which the CNS composes this model are inferred through the generalization capabilities of the subject outside the training data.

Transition Point Dynamic Programming

Kenneth Buckland, Peter Lawrence

Transition point dynamic programming (TPDP) is a memory based, reinforcement learning, direct dynamic programming approach to adaptive optimal control that can reduce the learning time and memory usage required for the control of continuous stochastic dynamic systems. TPDP does so by determining an ideal set of transition points (TPs) which specify only the control action changes necessary for optimal control. TPDP converges to an ideal TP set by using a variation of Q-learning to assess the merits of adding, swapping and removing TPs from states throughout the state space. When applied to a race track problem, TPDP learned the optimal control policy much sooner than conventional Q-learning, and was able to do so using less memory.

Dopaminergic Neuromodulation Brings a Dynamical Plasticity to the Retina

Eric Boussard, Jean-François Vibert

The fovea of a mammal retina was simulated with its detailed biological properties to study the local preprocessing of images. The direct visual pathway (photoreceptors, bipolar and ganglion cells) and the horizontal units, as well as the D-amacrine cells were simulated. The computer program simulated the analog non-spiking transmission between photoreceptor and bipolar cells, and between bipolar and ganglion cells, as well as the gap-junctions between horizontal cells, and the release of dopamine by D-amacrine cells and its diffusion in the extra-cellular space. A 64 x 64 photoreceptors retina, containing 16,448 units, was carried out. This retina displayed contour extraction with a Mach effect, and adaptation to brightness. The simulation showed that the dopaminergic amacrine cells were necessary to ensure adaptation to local brightness.

Monte Carlo Matrix Inversion and Reinforcement Learning

Andrew Barto, Michael Duff

We describe the relationship between certain reinforcement learning (RL) methods based on dynamic programming (DP) and a class of unorthodox Monte Carlo methods for solving systems of linear equations proposed in the 1950's. These methods recast the solution of the linear system as the expected value of a statistic suitably defined over sample paths of a Markov chain. The significance of our observations lies in arguments (Curtiss, 1954) that these Monte Carlo methods scale better with respect to state-space size than do standard, iterative techniques for solving systems of linear equations. This analysis also establishes convergence rate estimates. Because methods used in RL systems for approximating the evaluation function of a fixed control policy also approximate solutions to systems of linear equations, the connection to these Monte Carlo methods establishes that algorithms very similar to TD algorithms (Sutton, 1988) are asymptotically more efficient in a precise sense than other methods for evaluating policies. Further, all DP-based RL methods have some of the properties of these Monte Carlo algorithms, which suggests that although RL is often perceived to be slow, for sufficiently large pro

blems, it may in fact be more efficient than other known classes of methods capable of producing the same results.

Globally Trained Handwritten Word Recognizer using Spatial Representation, Convolutional Neural Networks, and Hidden Markov Models

Yoshua Bengio, Yann LeCun, Donnie Henderson

We introduce a new approach for on-line recognition of handwritten words written in unconstrained mixed style. The preprocessor performs a word-level normalization by fitting a model of the word structure using the EM algorithm. Words are then coded into low resolution "annotated images" where each pixel contains information about trajectory direction and curvature. The recognizer is a convolution network which can be spatially replicated. From the network output, a hidden Markov model produces word scores. The entire system is globally trained to minimize word-level errors.

Robust Reinforcement Learning in Motion Planning

Satinder Singh, Andrew Barto, Roderic Grupen, Christopher Connolly

While exploring to find better solutions, an agent performing on-line reinforcement learning (RL) can perform worse than is acceptable. In some cases, exploration might have unsafe, or even catastrophic, results, often modeled in terms of reaching 'failure' states of the agent's environment. This paper presents a method that uses domain knowledge to reduce the number of failures during exploration. This method formulates the set of actions from which the RL agent composes a control policy to ensure that exploration is conducted in a policy space that excludes most of the unacceptable policies. The resulting action set has a more abstract relationship to the task being solved than is common in many applications of RL. Although the cost of this added safety is that learning may result in a suboptimal solution, we argue that this is an appropriate tradeoff in many problems. We illustrate this method in the domain of motion planning.

Optimal Stopping and Effective Machine Complexity in Learning

Changfeng Wang, Santosh Venkatesh, J. Judd

We study the problem of when to stop learning a class of feedforward networks - networks with linear outputs and fixed input weights - when they are trained with a gradient descent algorithm on a finite number of examples. Under general regularity conditions, it is shown that there are in general three distinct phases in the generalization performance in the learning process, and in particular, the network has better generalization performance when learning is stopped at a certain time before the global minimum of the empirical error is reached. A notion of effective size of a machine is defined.

Exploiting Chaos to Control the Future

Gary Flake, Guo-Zhen Sun, Yee-Chun Lee

Recently, Ott, Grebogi and Yorke (OGY) [6] found an effective method to control chaotic systems to unstable fixed points by using only small control forces; however, OGY's method is based on and limited to a linear theory and requires considerable knowledge of the dynamics of the system to be controlled. In this paper we use two radial basis function networks: one as a model of an unknown plant and the other as the controller. The controller is trained with a recurrent learning algorithm to minimize a novel objective function such that the controller can locate an unstable fixed point and drive the system into the fixed point with no a priori knowledge of the system dynamics. Our results indicate that the neural controller offers many advantages over OGY's technique.

Backpropagation Convergence Via Deterministic Nonmonotone Perturbed Minimization

O. L. Mangasarian, M. V. Solodov

The fundamental backpropagation (BP) algorithm for training artificial neural networks is cast as a deterministic nonmonotone perturbed gradient method. Under certain natural assumptions, such as the series of learning rates diverging while the series of their squares converging, it is established that every accumulation point of the online BP iterates is a stationary point of the BP error function. The results presented cover serial and parallel online BP, modified BP with a momentum term, and BP with weight decay.

Fool's Gold: Extracting Finite State Machines from Recurrent Network Dynamics

John Kolen

Several recurrent networks have been proposed as representations for the task of formal language learning. After training a recurrent network recognize a formal language or predict the next symbol of a sequence, the next logical step is to understand the information processing carried out by the network. Some researchers have begun to extracting finite state machines from the internal state trajectories of their recurrent networks. This paper describes how sensitivity to initial conditions and discrete measurements can trick these extraction methods to return illusory finite state descriptions.

Two Iterative Algorithms for Computing the Singular Value Decomposition from Input/Output Samples

Terence Sanger

The Singular Value Decomposition (SVD) is an important tool for linear algebra and can be used to invert or approximate matrices. Although many authors use "SVD" synonymously with "Eigenvalue vector Decomposition" or "Principal Components Transform", it is important to realize that these other methods apply only to symmetric matrices, while the SVD can be applied to arbitrary nonsquare matrices. This property is important for applications to signal transmission and control. I propose two new algorithms for iterative computation of the SVD given only sample inputs and outputs from a matrix. Although there currently exist many algorithms for Eigenvector Decomposition (Sanger 1989, for example), these are the first true sample-based SVD algorithms.

Efficient Computation of Complex Distance Metrics Using Hierarchical Filtering

Patrice Simard

By their very nature, memory based algorithms such as KNN or Parzen windows require a computationally expensive search of a large database of prototypes. In this paper we optimize the searching process for tangent distance (Simard, LeCun and Denker, 1993) to improve speed performance. The closest prototypes are found by recursively searching included subsets of the database using distances of increasing complexity. This is done by using a hierarchy of tangent distances (increasing the number of tangent vectors from 0 to its maximum) and multiresolution (using wavelets). At each stage, a confidence level of the classification is computed. If the confidence is high enough, the computation of more complex distances is avoided. The resulting algorithm applied to character recognition is close to three orders of magnitude faster than computing the full tangent distance on every prototypes.

Statistics of Natural Images: Scaling in the Woods

Daniel Ruderman, William Bialek

In order to best understand a visual system one should attempt to characterize the natural images it processes. We gather images from the woods and find that these scenes possess an ensemble scale invariance. Further, they are highly non-Gaussian, and this non-Gaussian character cannot be removed through local linear filtering. We find that including a simple "gain control" nonlinearity in the filtering process makes the filter output quite Gaussian, meaning information is maximized at fixed channel variance. Finally, we use the measured power spectrum to place an upper bound on the information conveyed about natural scenes by an array of

receptors.

GDS: Gradient Descent Generation of Symbolic Classification Rules

Reinhard Blasig

Imagine you have designed a neural network that successfully learns a complex classification task. What are the relevant input features the classifier relies on and how are these features combined to produce the classification decisions? There are applications where a deeper insight into the structure of an adaptive system and thus into the underlying classification problem may well be as important as the system's performance characteristics, e.g. in economics or medicine. GDSi is a backpropagation-based training scheme that produces networks transformable into an equivalent and concise set of IF-THEN rules. This is achieved by imposing penalty terms on the network parameters that adapt the network to the expressive power of this class of rules. Thus during training we simultaneously minimize classification and transformation error. Some real-world tasks demonstrate the viability of our approach.

Packet Routing in Dynamically Changing Networks: A Reinforcement Learning Approach

Justin Boyan, Michael Littman

This paper describes the Q-routing algorithm for packet routing, in which a reinforcement learning module is embedded into each node of a switching network. Only local communication is used by each node to keep accurate statistics on which routing decisions lead to minimal delivery times. In simple experiments involving a 36-node, irregularly connected network, Q-routing proves superior to a nonadaptive algorithm based on precomputed shortest paths and is able to route efficiently even when critical aspects of the simulation, such as the network load, are allowed to vary dynamically. The paper concludes with a discussion of the tradeoff between discovering shortcuts and maintaining stable policies.

A Massively-Parallel SIMD Processor for Neural Network and Machine Vision Applications

Michael Glover, W. Miller

This paper describes the MM32k, a massively-parallel SIMD computer which is easy to program, high in performance, low in cost and effective for implementing highly parallel neural network architectures. The MM32k has 32 768 bit serial processing elements, each of which has 512 bits of memory, and all of which are interconnected by a switching network. The entire system resides on a single PC-AT compatible card. It is programmed from the host computer using a C++ language class library which abstracts the parallel processor in terms of fast arithmetic operators for vectors of variable precision integers.

How to Choose an Activation Function

H. N. Mhaskar, C. A. Micchelli

We study the complexity problem in artificial feedforward neural networks designed to approximate real valued functions of several real variables; i.e., we estimate the number of neurons in a network required to ensure a given degree of approximation to every function in a given function class. We indicate how to construct networks with the indicated number of neurons evaluating standard activation functions. Our general theorem shows that the smoother the activation function, the better the rate of approximation.

Efficient Simulation of Biological Neural Networks on Massively Parallel Supercomputers with Hypercube Architecture

Ernst Niebur, Dean Brett

We present a neural network simulation which we implemented on the ma

ssively parallel Connection Machine 2. In contrast to previous work, this simulator is based on biologically realistic neurons with nontrivial single-cell dynamics, high connectivity with a structure modelled in agreement with biological data, and preservation of the temporal dynamics of spike interactions. We simulate neural networks of 16,384 neurons coupled by about 1000 synapses per neuron, and estimate the performance for much larger systems. Communication between neurons is identified as the computationally most demanding task and we present a novel method to overcome this bottleneck. The simulator has already been used to study the primary visual system of the cat.

Convergence of Stochastic Iterative Dynamic Programming Algorithms

Tommi Jaakkola, Michael Jordan, Satinder Singh

Increasing attention has recently been paid to algorithms based on dynamic programming (DP) due to the suitability of DP for learning problems involving control. In stochastic environments where the system being controlled is only incompletely known, however, a unifying theoretical account of these methods has been missing. In this paper we relate DP-based learning algorithms to the powerful techniques of stochastic approximation via a new convergence theorem, enabling us to establish a class of convergent algorithms to which both TD(0) and Q-learning belong.

High Performance Neural Net Simulation on a Multiprocessor System with "Intelligent" Communication

Urs A. Müller, Michael Kocheisen, Anton Gunzinger

The performance requirements in experimental research on artificial neural nets often exceed the capability of workstations and PCs by a great amount. But speed is not the only requirement. Flexibility and implementation time for new algorithms are usually of equal importance. This paper describes the simulation of neural nets on the MUSIC parallel supercomputer, a system that shows a good balance between the three issues and therefore made many research projects possible that were unthinkable before. (MUSIC stands for Multiprocessor System with Intelligent Communication)

An Analog VLSI Model of Central Pattern Generation in the Leech

Micah Siegel

I detail the design and construction of an analog VLSI model of the neural system responsible for swimming behaviors of the leech. Why the leech? The biological network is small and relatively well understood, and the silicon model can therefore span three levels of organization in the leech nervous system (neuron, ganglion, system); it represents one of the first comprehensive models of leech swimming operating in real-time. The circuit employs biophysically motivated analog neurons networked to form multiple biologically inspired silicon ganglia. These ganglia are coupled using known interganglionic connections. Thus the model retains the flavor of its biological counterpart, and though simplified, the output of the silicon circuit is similar to the output of the leech swim central pattern generator. The model operates on the same time- and spatial-scale as the leech nervous system and will provide an excellent platform with which to explore real-time adaptive locomotion in the leech and other "simple" invertebrate nervous systems.

Non-Linear Statistical Analysis and Self-Organizing Hebbian Networks

Jonathan Shapiro, Adam Prügel-Bennett

Neurons learning under an unsupervised Hebbian learning rule can perform a nonlinear generalization of principal component analysis. This relationship between nonlinear PCA and nonlinear neurons is reviewed. The stable fixed points of the neuron learning dynamics correspond to the maxima of the statistically optimized nonlinear PCA. However, in order to predict what the neuron learns, knowledge of the basins of attractions of the n

neuron dynamics is required. Here the correspondence between nonlinear PCA and neural networks breaks down. This is shown for a simple model. Methods of statistical mechanics can be used to find the optima of the objective function of non-linear PCA. This determines what the neurons can learn. In order to find how the solutions are partitioned among the neurons, however, one must solve the dynamics.

Locally Adaptive Nearest Neighbor Algorithms

Dietrich Wettschereck, Thomas Dietterich

Four versions of a k-nearest neighbor algorithm with locally adaptive k are introduced and compared to the basic k-nearest neighbor algorithm (kNN). Locally adaptive kNN algorithms choose the value of k that should be used to classify a query by consulting the results of cross-validation computations in the local neighborhood of the query. Local kNN methods are shown to perform similar to kNN in experiments with twelve commonly used data sets. Encouraging results in three constructed tasks show that local methods can significantly outperform kNN in specific applications. Local methods can be recommended for on-line learning and for applications where different regions of the input space are covered by patterns solving different sub-tasks.

Bayesian Backprop in Action: Pruning, Committees, Error Bars and an Application to Spectroscopy

Hans Thodberg

MacKay's Bayesian framework for backpropagation is conceptually appealing as well as practical. It automatically adjusts the weight decay parameters during training, and computes the evidence for each trained network. The evidence is proportional to our belief in the model. The networks with highest evidence turn out to generalise well. In this paper, the framework is extended to pruned nets, leading to an Occam Factor for "tuning the architecture to the data". A committee of networks, selected by their high evidence, is a natural Bayesian construction. The evidence of a committee is computed. The framework is illustrated on real-world data from a near infrared spectrometer used to determine the fat content in minced meat. Error bars are computed, including the contribution from the dissent of the committee members.

Solvable Models of Artificial Neural Networks

Sumio Watanabe

Solvable models of nonlinear learning machines are proposed, and learning in artificial neural networks is studied based on the theory of ordinary differential equations. A learning algorithm is constructed, by which the optimal parameter can be found without any recursive procedure.

The solvable models enable us to analyze the reason why experimental results by the error backpropagation often contradict the statistical learning theory.

A Comparative Study of a Modified Bumptree Neural Network with Radial Basis Function Networks and the Standard Multi Layer Perceptron

Richard Bostock, Alan Harget

Bumptrees are geometric data structures introduced by Omohundro (1991) to provide efficient access to a collection of functions on a Euclidean space of interest. We describe a modified bumptree structure that has been employed as a neural network classifier, and compare its performance on several classification tasks against that of radial basis function networks and the standard multi-layer perceptron.

Generalization Error and the Expected Network Complexity

Chuanyi Ji

$O(E^{-1}) = O((EK)^{-1} \ln N)$,

Combined Neural Networks for Time Series Analysis

Iris Ginzburg, David Horn

We propose a method for improving the performance of any net(cid:173) work designed to predict the next value of a time series. We advocate analyzing the deviations of the network's predictions from the data in the training set. This can be carried out by a secondary net(cid:173) work trained on the time series of these residuals. The combined system of the two networks is viewed as the new predictor. We demonstrate the simplicity and success of this method, by applying it to the sunspots data. The small corrections of the secondary network can be regarded as resulting from a Taylor expansion of a complex network which includes the combined system. We find that the complex network is more difficult to train and performs worse than the two-step procedure of the combined system.

Illumination-Invariant Face Recognition with a Contrast Sensitive Silicon Retina

Joachim Buhmann, Martin Lades, Frank Eeckman

Changes in lighting conditions strongly effect the performance and reliability of computer vision systems. We report face recognition results under drastically changing lighting conditions for a computer vision system which concurrently uses a contrast sensitive silicon retina and a conventional, gain controlled CCO camera. For both input devices the face recognition system employs an elastic matching algorithm with wavelet based features to classify unknown faces. To assess the effect of analog on-chip preprocessing by the silicon retina the CCO images have been "digitally preprocessed" with a bandpass filter to adjust the power spectrum. The silicon retina with its ability to adjust sensitivity increases the recognition rate up to 50 percent. These comparative experiments demonstrate that preprocessing with an analog VLSI silicon retina generates image data enriched with object-constant features.

Address Block Location with a Neural Net System

Hans Graf, Eric Cosatto

We developed a system for finding address blocks on mail pieces that can process four images per second. Besides locating the address block, our system also determines the writing style, handwritten or machine printed, and moreover, it measures the skew angle of the text lines and cleans noisy images. A layout analysis of all the elements present in the image is performed in order to distinguish drawings and dirt from text and to separate text of advertisement from that of the destination address. A speed of more than four images per second is obtained on a modular hardware platform, containing a board with two of the NET32K neural net chips, a SPARC2 processor board, and a board with 2 digital signal processors. The system has been tested with more than 100,000 images. Its performance depends on the quality of the images, and lies between 85% correct location in very noisy images to over 98% in cleaner images.

Assessing the Quality of Learned Local Models

Stefan Schaal, Christopher Atkeson

An approach is presented to learning high dimensional functions in the case where the learning algorithm can affect the generation of new data. A local modeling algorithm, locally weighted regression, is used to represent the learned function. Architectural parameters of the approach, such as distance metrics, are also localized and become a function of the query point instead of being global. Statistical tests are given for when a local model is good enough and sampling should be moved to a new area. Our methods explicitly deal with the case where prediction accuracy requirements exist during exploration: By gradually shifting a "center of exploration" and controlling the speed of the shift with local prediction accuracy, a goal-directed exploration

on of state space takes place along the fringes of the current data support until the task goal is achieved. We illustrate this approach with simulation results and results from a real robot learning a complex juggling task.

Synchronization, oscillations, and $1/f$ noise in networks of spiking neurons

Martin Stemmler, Marius Usher, Christof Koch, Zeev Olami

We investigate a model for neural activity that generates long range temporal correlations, $1/f$ noise, and oscillations in global activity. The model consists of a two-dimensional sheet of leaky integrate-and-fire neurons with feedback connectivity consisting of local excitation and surround inhibition. Each neuron is independently driven by homogeneous external noise. Spontaneous symmetry breaking occurs, resulting in the formation of "hotspots" of activity in the network. These localized patterns of excitation appear as clusters that coalesce, disintegrate, or fluctuate in size while simultaneously moving in a random walk constrained by the interaction with other clusters. The emergent cross-correlation functions have a dual structure, with a sharp peak around zero on top of a much broader hill. The power spectrum associated with single units shows a $1/f$ decay for small frequencies and is flat at higher frequencies, while the power spectrum of the spiking activity averaged over many cells-equivalent to the local field potential-shows no $1/f$ decay but a prominent peak around 40 Hz.

What Does the Hippocampus Compute?: A Precursor of the 1993 NIPS Workshop

Mark Gluck

Requests for name changes in the electronic proceedings will be accepted with no questions asked. However name changes may cause bibliographic tracking issues.

Authors are asked to consider this carefully and discuss it with their co-authors prior to requesting a name change in the electronic proceedings.

Unsupervised Learning of Mixtures of Multiple Causes in Binary Data

Eric Saund

This paper presents a formulation for unsupervised learning of clusters reflecting multiple causal structure in binary data. Unlike the standard mixture model, a multiple cause model accounts for observed data by combining assertions from many hidden causes, each of which can pertain to varying degree to any subset of the observable dimensions. A crucial issue is the mixing-function for combining beliefs from different cluster-centers in order to generate data reconstructions whose errors are minimized both during recognition and learning. We demonstrate a weakness inherent to the popular weighted sum followed by sigmoid squashing, and offer an alternative form of the nonlinearity. Results are presented demonstrating the algorithm's ability successfully to discover coherent multiple causal representations of noisy test data and in images of printed characters.

A Hodgkin-Huxley Type Neuron Model That Learns Slow Non-Spike Oscillation

Kenji Doya, Allen Selverston, Peter Rowat

A gradient descent algorithm for parameter estimation which is similar to those used for continuous-time recurrent neural networks was derived for Hodgkin-Huxley type neuron models. Using membrane potential trajectories as targets, the parameters (maximal conductances, thresholds and slopes of activation curves, time constants) were successfully estimated. The algorithm was applied to modeling slow non-spike oscillation of an identified neuron in the lobster stomatogastric ganglion. A model with three ionic currents was trained with experimental data. It revealed a novel role of A-current for slow oscillation below -50 mV.

Memory-Based Methods for Regression and Classification

Thomas Dietterich, Dietrich Wettschereck, Chris G. Atkeson, Andrew Moore

Requests for name changes in the electronic proceedings will be accepted with no questions asked. However name changes may cause bibliographic tracking issues.

Authors are asked to consider this carefully and discuss it with their co-authors prior to requesting a name change in the electronic proceedings.

Neural Network Methods for Optimization Problems

Arun Jagota

Requests for name changes in the electronic proceedings will be accepted with no questions asked. However name changes may cause bibliographic tracking issues.

Authors are asked to consider this carefully and discuss it with their co-authors prior to requesting a name change in the electronic proceedings.

A Comparison of Dynamic Reposing and Tangent Distance for Drug Activity Prediction

Thomas Dietterich, Ajay Jain, Richard Lathrop, Tomás Lozano-Pérez

In drug activity prediction (as in handwritten character recognition), the features extracted to describe a training example depend on the pose (location, orientation, etc.) of the example. In handwritten character recognition, one of the best techniques for addressing this problem is the tangent distance method of Simard, LeCun and Denker (1993). Jain, et al. (1993a; 1993b) introduce a new technique-dynamic reposing-that also addresses this problem. Dynamic reposing iteratively learns a neural network and then reposes the examples in an effort to maximize the predicted output values. New models are trained and new poses computed until models and poses converge. This paper compares dynamic reposing to the tangent distance method on the task of predicting the biological activity of musk compounds. In a 20-fold cross-validation,

Hoo Optimality Criteria for LMS and Backpropagation

Babak Hassibi, Ali H. Sayed, Thomas Kailath

We have recently shown that the widely known LMS algorithm is an HOO optimal estimator. The HOO criterion has been introduced, initially in the control theory literature, as a means to ensure robust performance in the face of model uncertainties and lack of statistical information on the exogenous signals. We extend here our analysis to the nonlinear setting often encountered in neural networks, and show that the backpropagation algorithm is locally HOO optimal. This fact provides a theoretical justification of the widely observed excellent robustness properties of the LMS and backpropagation algorithms. We further discuss some implications of these results.

Classifying Hand Gestures with a View-Based Distributed Representation

Trevor J. Darrell, Alex P. Pentland

We present a method for learning, tracking, and recognizing human hand gestures recorded by a conventional CCD camera without any special gloves or other sensors. A view-based representation is used to model aspects of the hand relevant to the trained gestures, and is found using an unsupervised clustering technique. We use normalized correlation networks, with dynamic time warping in the temporal domain, as a distance function for unsupervised clustering. Views are computed separably for space and time dimensions; the distributed response of the combination of these units characterizes the input data with a low dimensional representation. A supervised classification stage uses labeled outputs of the spatio-temporal units as training data. Our system can correctly classify gestures in real time with a low-cost image processing accelerator.

Optimal Unsupervised Motor Learning Predicts the Internal Representation of Barn Owl Head Movements

Terence Sanger

(Masino and Knudsen 1990) showed some remarkable results which suggest that he

ad motion in the barn owl is controlled by distinct circuits coding for the horizontal and vertical components of movement. This implies the existence of a set of orthogonal internal coordinates that are related to meaningful coordinates of the external world. No coherent computational theory has yet been proposed to explain this finding. I have proposed a simple model which provides a framework for a theory of low-level motor learning. I show that the theory predicts the observed microstimulation results in the barn owl. The model rests on the concept of "Optimal Unsupervised Motor Learning", which provides a set of criteria that predict optimal internal representations. I describe two iterative Neural Network algorithms which find the optimal solution and demonstrate possible mechanisms for the development of internal representations in animals.

Amplifying and Linearizing Apical Synaptic Inputs to Cortical Pyramidal Cells

Öjvind Bernander, Christof Koch, Rodney Douglas

Intradendritic electrophysiological recordings reveal a bewildering repertoire of complex electrical spikes and plateaus that are difficult to reconcile with conventional notions of neuronal function. In this paper we argue that such dendritic events are just an extrinsic expression of a more important mechanism - a proportional current amplifier whose primary task is to offset electrotonic losses. Using the example of functionally important synaptic inputs to the superficial layers of an anatomically and electrophysiologically reconstructed layer 5 pyramidal neuron, we derive and simulate the properties of conductances that linearize and amplify distal synaptic input current in a graded manner. The amplification depends on a potassium conductance in the apical tuft and calcium conductances in the apical trunk.

Decoding Cursive Scripts

Yoram Singer, Naftali Tishby

Online cursive handwriting recognition is currently one of the most intriguing challenges in pattern recognition. This study presents a novel approach to this problem which is composed of two complementary phases. The first is dynamic encoding of the writing trajectory into a compact sequence of discrete motor control symbols. In this compact representation we largely remove the redundancy of the script, while preserving most of its intelligible components. In the second phase these control sequences are used to train adaptive probabilistic acyclic automata (PAA) for the important ingredients of the writing trajectories, e.g. letters. We present a new and efficient learning algorithm for such stochastic automata, and demonstrate its utility for spotting and segmentation of cursive scripts. Our experiments show that over 90% of the letters are correctly spotted and identified, prior to any higher level language model. Moreover, both the training and recognition algorithms are very efficient compared to other modeling methods, and the models are 'on-line' adaptable to other writers and styles.

Foraging in an Uncertain Environment Using Predictive Hebbian Learning

P. Montague, Peter Dayan, Terrence J. Sejnowski

Survival is enhanced by an ability to predict the availability of food, the likelihood of predators, and the presence of mates. We present a concrete model that uses diffuse neurotransmitter systems to implement a predictive version of a Hebb learning rule embedded in a neural architecture based on anatomical and physiological studies on bees. The model captured the strategies seen in the behavior of bees and a number of other animals when foraging in an uncertain environment. The predictive model suggests a unified way in which neuromodulatory influences can be used to bias actions and control synaptic plasticity.

Analysis of Short Term Memories for Neural Networks

Jose C. Principe, Hui-H. Hsu, Jyh-Ming Kuo

Short term memory is indispensable for the processing of time varying information with artificial neural networks. In this paper a model for linear memories is presented, and ways to include memories in connectionist topologies are discussed. A comparison is drawn among different memory types, with indication of what is the salient characteristic of each memory model.

Credit Assignment through Time: Alternatives to Backpropagation

Yoshua Bengio, Paolo Frasconi

Learning to recognize or predict sequences using long-term context has many applications. However, practical and theoretical problems are found in training recurrent neural networks to perform tasks in which input/output dependencies span long intervals. Starting from a mathematical analysis of the problem, we consider and compare alternative algorithms and architectures on tasks for which the span of the input/output dependencies can be controlled. Results on the new algorithms show performance qualitatively superior to that obtained with backpropagation.

Analyzing Cross-Connected Networks

Thomas Shultz, Jeffrey Elman

Jeffrey L. Elman

Recognition-based Segmentation of On-Line Cursive Handwriting

Nicholas Flann

This paper introduces a new recognition-based segmentation approach to recognizing on-line cursive handwriting from a database of 10,000 English words. The original input stream of x, y pen coordinates is encoded as a sequence of uniform stroke descriptions that are processed by six feed-forward neural-networks, each designed to recognize letters of different sizes. Words are then recognized by performing best-first search over the space of all possible segmentations. Results demonstrate that the method is effective at both writer dependent recognition (1.7% to 15.5% error rate) and writer independent recognition (5.2% to 31.1% error rate).

Implementing Intelligence on Silicon Using Neuron-Like Functional MOS Transistors

Tadashi Shibata, Koji Kotani, Takeo Yamashita, Hiroshi Ishii, Hideo Kosaka, Tadahiro Ohmi

We will present the implementation of intelligent electronic circuits realized for the first time using a new functional device called Neuron MOS Transistor (neuMOS or vMOS in short) simulating the behavior of biological neurons at a single transistor level. Search for the most resembling data in the memory cell array, for instance, can be automatically software manipulation. Soft Hardware, which we named, can arbitrarily change its logic function in real time by external control signals without any hardware modification. Implementation of a neural network equipped with an on-chip self-learning capability is also described. Through the studies of vMOS intelligent circuit implementation, we noticed an interesting similarity in the architectures of vMOS logic circuitry and biological systems.

Probabilistic Anomaly Detection in Dynamic Systems

Padhraic Smyth

This paper describes probabilistic methods for novelty detection when using pattern recognition methods for fault monitoring of dynamic systems. The problem of novelty detection is particularly acute when prior knowledge and training data only allow one to construct an incomplete classification model. Allowance must be made in model design so that the classifier will be robust to data generated by classes not

included in the training phase. For diagnosis applications one practical approach is to construct both an input density model and a discriminative class model. Using Bayes' rule and prior estimates of the relative likelihood of data of known and unknown origin the resulting classification equations are straightforward. The paper describes the application of this method in the context of hidden Markov models for online fault monitoring of large ground antennas for spacecraft tracking, with particular application to the detection of transient behaviour of unknown origin.

Learning Complex Boolean Functions: Algorithms and Applications

Arlindo Oliveira, Alberto Sangiovanni-Vincentelli

The most commonly used neural network models are not well suited to direct digital implementations because each node needs to perform a large number of operations between floating point values. Fortunately, the ability to learn from examples and to generalize is not restricted to networks of this type. Indeed, networks where each node implements a simple Boolean function (Boolean networks) can be designed in such a way as to exhibit similar properties. Two algorithms that generate Boolean networks from examples are presented. The results show that these algorithms generalize very well in a class of problems that accept compact Boolean network descriptions. The techniques described are general and can be applied to tasks that are not known to have that characteristic. Two examples of applications are presented: image reconstruction and hand-written character recognition.

A Unified Gradient-Descent/Clustering Architecture for Finite State Machine Induction

Sreerupa Das, Michael C. Mozer

Although recurrent neural nets have been moderately successful in learning to emulate finite-state machines (FSMs), the continuous internal state dynamics of a neural net are not well matched to the discrete behavior of an FSM. We describe an architecture, called DOLCE, that allows discrete states to evolve in a net as learning progresses. DOLCE consists of a standard recurrent neural net trained by gradient descent and an adaptive clustering technique that quantizes the state space. DOLCE is based on the assumption that a finite set of discrete internal states is required for the task, and that the actual network state belongs to this set but has been corrupted by noise due to inaccuracy in the weights. DOLCE learns to recover the discrete state with maximum a posteriori probability from the noisy state. Simulations show that DOLCE leads to a significant improvement in generalization performance over earlier neural net approaches to FSM induction.

Convergence of Indirect Adaptive Asynchronous Value Iteration Algorithms

Vijaykumar Gullapalli, Andrew Barto

Reinforcement Learning methods based on approximating dynamic programming (DP) are receiving increased attention due to their utility in forming reactive control policies for systems embedded in dynamic environments. Environments are usually modeled as controlled Markov processes, but when the environment model is not known a priori, adaptive methods are necessary. Adaptive control methods are often classified as being direct or indirect. Direct methods directly adapt the control policy from experience, whereas indirect methods adapt a model of the controlled process and compute control policies based on the latest model. Our focus is on indirect adaptive DP-based methods in this paper. We present a convergence result for indirect adaptive asynchronous value iteration algorithms for the case in which a look-up table is used to store the value function. Our result implies convergence of several existing reinforcement learning algorithms such as adaptive real-time dynamic programming (ARTDP) (Barto, Bradtke, & Singh, 1993) and prioritized sweeping (Moor

e & Atkeson, 1993). Although the emphasis of researchers studying DP-based reinforcement learning has been on direct adaptive methods such as Q-Learning (Watkins, 1989) and methods using TD algorithms (Sutton, 1988), it is not clear that these direct methods are preferable in practice to indirect methods such as those analyzed in this paper.

Fast Pruning Using Principal Components

Asriel Levin, Todd Leen, John Moody

We present a new algorithm for eliminating excess parameters and improving network generalization after supervised training. The method, "Principal Components Pruning (PCP)", is based on principal component analysis of the node activations of successive layers of the network. It is simple, cheap to implement, and effective. It requires no network retraining, and does not involve calculating the full Hessian of the cost function. Only the weight and the node activity correlation matrices for each layer of nodes are required. We demonstrate the efficacy of the method on a regression problem using polynomial basis functions, and on an economic time series prediction problem using a two-layer, feedforward network.

Hidden Markov Models for Human Genes

Pierre Baldi, Søren Brunak, Yves Chauvin, Jacob Engelbrecht, Anders Krogh

Human genes are not continuous but rather consist of short coding regions (exons) interspersed with highly variable non-coding regions (introns). We apply HMMs to the problem of modeling exons, introns and detecting splice sites in the human genome. Our most interesting result so far is the detection of particular oscillatory patterns, with a minimal period of roughly 10 nucleotides, that seem to be characteristic of exon regions and may have significant biological implications.

Postal Address Block Location Using a Convolutional Locator Network

Ralph Wolf, John Platt

This paper describes the use of a convolutional neural network to perform address block location on machine-printed mail pieces. Locating the address block is a difficult object recognition problem because there is often a large amount of extraneous printing on a mail piece and because address blocks vary dramatically in size and shape. We used a convolutional locator network with four outputs, each trained to find a different corner of the address block. A simple set of rules was used to generate ABL candidates from the network output. The system performs very well: when allowed five guesses, the network will tightly bound the address delivery information in 98.2% of the cases.

An Analog VLSI Saccadic Eye Movement System

Timothy Horiuchi, Brooks Bishofberger, Christof Koch

In an effort to understand saccadic eye movements and their relation to visual attention and other forms of eye movements, we - in collaboration with a number of other laboratories - are carrying out a large-scale effort to design and build a complete primate oculomotor system using analog CMOS VLSI technology. Using this technology, a low power, compact, multi-chip system has been built which works in real-time using real-world visual inputs. We describe in this paper the performance of an early version of such a system including a 1-D array of photoreceptors mimicking the retina, a circuit computing the mean location of activity representing the superior colliculus, a saccadic burst generator, and a one degree-of-freedom rotational platform which models the dynamic properties of the primate oculomotor plant.

An Optimization Method of Layered Neural Networks based on the Modified Information Criterion

Sumio Watanabe

This paper proposes a practical optimization method for layered neural networks, by which the optimal model and parameter can be found simultaneously. We modify the conventional information criterion into a differentiable function of parameters, and then, minimize it, while controlling it back to the ordinary form. Effectiveness of this method is discussed theoretically and experimentally.

A Computational Model for Cursive Handwriting Based on the Minimization Principle

Yasuhiro Wada, Yasuharu Koike, Eric Vatikiotis-Bateson, Mitsuo Kawato

We propose a trajectory planning and control theory for continuous movements such as connected cursive handwriting and continuous natural speech. Its hardware is based on our previously proposed forward-inverse-relaxation neural network (Wada & Kawato, 1993). Computationally, its optimization principle is the minimum torque change criterion. Regarding the representation level, hard constraints satisfied by a trajectory are represented as a set of via-points extracted from a handwritten character. Accordingly, we propose a via-point estimation algorithm that estimates via-points by repeating the trajectory formation of a character and the via-point extraction from the character. In experiments, good quantitative agreement is found between human handwriting data and the trajectories generated by the theory. Finally, we propose a recognition schema based on the movement generation. We show a result in which the recognition schema is applied to the handwritten character recognition and can be extended to the phoneme timing estimation of natural speech.

Discontinuous Generalization in Large Committee Machines

H. Schwarze, J. Hertz

The problem of learning from examples in multilayer networks is studied within the framework of statistical mechanics. Using the replica formalism we calculate the average generalization error of a fully connected committee machine in the limit of a large number of hidden units. If the number of training examples is proportional to the number of inputs in the network, the generalization error as a function of the training set size approaches a finite value. If the number of training examples is proportional to the number of weights in the network we find first-order phase transitions with a discontinuous drop in the generalization error for both binary and continuous weights.

Inverse Dynamics of Speech Motor Control

Makoto Hirayama, Eric Vatikiotis-Bateson, Mitsuo Kawato

Progress has been made in computational implementation of speech production based on physiological data. An inverse dynamics model of the speech articulator's musculo-skeletal system, which is the mapping from articulator trajectories to electromyographic (EMG) signals, was modeled using the acquired forward dynamics model and temporal (smoothness of EMG activation) and range constraints. This inverse dynamics model allows the use of a faster speech motor control scheme, which can be applied to phoneme-to-speech synthesis via musculo-skeletal system dynamics, or to future use in speech recognition. The forward acoustic model, which is the mapping from articulator trajectories to the acoustic parameters, was improved by adding velocity and voicing information inputs to distinguish acoustic parameter differences caused by changes in source characteristics.

Learning Temporal Dependencies in Connectionist Speech Recognition

Steve Renals, Mike Hochberg, Tony Robinson

Hybrid connectionist/HMM systems model time both using a Markov chain and through properties of a connectionist network. In this paper, we discuss the nature of the time dependence currently employed in our systems using recurrent networks (RNs) and feed-forward multi-layer perceptrons (MLPs). In particular, we int

duce local recurrences into a MLP to produce an enhanced input representation. This is in the form of an adaptive gamma filter and incorporates an automatic approach for learning temporal dependencies. We have experimented on a speaker (cid:l73) independent phone recognition task using the TIMIT database. Results using the gamma filtered input representation have shown improvement over the baseline MLP system. Improvements have also been obtained through merging the baseline and gamma filter models.

Odor Processing in the Bee: A Preliminary Study of the Role of Central Input to the Antennal Lobe

Christiane Linster, David Marsan, Claudine Masson, Michel Kerszberg

Based on precise anatomical data of the bee's olfactory system, we propose an investigation of the possible mechanisms of modulation and control between the two levels of olfactory information processing: the antennal lobe glomeruli and the mushroom bodies. We use simplified neurons, but realistic architecture. As a first conclusion, we postulate that the feature extraction performed by the antennal lobe (glomeruli and interneurons) necessitates central input from the mushroom bodies for fine tuning. The central input thus facilitates the evolution from fuzzy olfactory images in the glomerular layer towards more focussed images upon odor presentation.

Surface Learning with Applications to Lipreading

Christoph Bregler, Stephen Omohundro

Most connectionist research has focused on learning mappings from one space to another (eg. classification and regression). This paper introduces the more general task of learning constraint surfaces. It describes a simple but powerful architecture for learning and manipulating nonlinear surfaces from data. We demonstrate the technique on low dimensional synthetic surfaces and compare it to nearest neighbor approaches. We then show its utility in learning the space of lip images in a system for improving speech recognition by lip reading. This learned surface is used to improve the visual tracking performance during recognition.

Figure of Merit Training for Detection and Spotting

Eric Chang, Richard P. Lippmann

Spotting tasks require detection of target patterns from a background of richly varied non-target inputs. The performance measure of interest for these tasks, called the figure of merit (FOM), is the detection rate for target patterns when the false alarm rate is in an acceptable range. A new approach to training spotters is presented which computes the FOM gradient for each input pattern and then directly maximizes the FOM using backpropagation. This eliminates the need for thresholds during training. It also uses network resources to model Bayesian a posteriori probability functions accurately only for patterns which have a significant effect on the detection accuracy over the false alarm rate of interest. FOM training increased detection accuracy by 5 percentage points for a hybrid radial basis function (RBF) - hidden Markov model (HMM) wordspotter on the credit-card speech corpus.

Generation of Internal Representation by α -Transformation

Ryotaro Kamimura

Requests for name changes in the electronic proceedings will be accepted with no questions asked. However name changes may cause bibliographic tracking issues.

Authors are asked to consider this carefully and discuss it with their co-authors prior to requesting a name change in the electronic proceedings.

Autoencoders, Minimum Description Length and Helmholtz Free Energy

Geoffrey E. Hinton, Richard Zemel

An autoencoder network uses a set of recognition weights to convert an input vector into a code vector. It then uses a set of generative weights to convert t

he code vector into an approximate reconstruction of the input vector. We derive an objective function for training autoencoders based on the Minimum Description Length (MDL) principle. The aim is to minimize the information required to describe both the code vector and the reconstruction error. We show that this information is minimized by choosing code vectors stochastically according to a Boltzmann distribution, where the generative weights define the energy of each possible code vector given the input vector. Unfortunately, if the code vectors use distributed representations, it is exponentially expensive to compute this Boltzmann distribution because it involves all possible code vectors. We show that the recognition weights of an autoencoder can be used to compute an approximation to the Boltzmann distribution and that this approximation gives an upper bound on the description length. Even when this bound is poor, it can be used as a Lyapunov function for learning both the generative and the recognition weights. We demonstrate that this approach can be used to learn factorial codes.

Dual Mechanisms for Neural Binding and Segmentation

Paul Sajda, Leif Finkel

We propose that the binding and segmentation of visual features is mediated by two complementary mechanisms; a low resolution, spatial-based, resource-free process and a high resolution, temporal-based, resource-limited process. In the visual cortex, the former depends upon the orderly topographic organization in striate and extrastriate areas while the latter may be related to observed temporal relationships between neuronal activities. Computer simulations illustrate the role the two mechanisms play in figure/ground discrimination, depth-from-occlusion, and the vividness of perceptual completion.

The Statistical Mechanics of k-Satisfaction

Scott Kirkpatrick, Géza Györgyi, Naftali Tishby, Lidror Troyansky

The satisfiability of random CNF formulae with precisely k variables per clause ("k-SAT") is a popular testbed for the performance of search algorithms. Formulae have M clauses from N variables, randomly negated, keeping the ratio $a = M / N$ fixed. For $k = 2$, this model has been proven to have a sharp threshold at $a = 1$ between formulae which are almost always satisfiable and formulae which are almost never satisfiable as $N \rightarrow \infty$. Computer experiments for $k = 2, 3, 4, 5$ and 6, (carried out in collaboration with B. Selman of ATT Bell Labs), show similar threshold behavior for each value of k . Finite-size scaling, a theory of the critical point phenomena used in statistical physics, is shown to characterize the size dependence near the threshold. Annealed and replica-based mean field theories give a good account of the results.

Lipreading by neural networks: Visual preprocessing, learning, and sensory integration

Gregory Wolff, K. Prasad, David Stork, Marcus Hennecke

We have developed visual preprocessing algorithms for extracting phonologically relevant features from the grayscale video image of a speaker, to provide speaker-independent inputs for an automatic lipreading ("speechreading") system. Visual features such as mouth open/closed, tongue visible/not-visible, teeth visible/not-visible, and several shape descriptors of the mouth and its motion are all rapidly computable in a manner quite insensitive to lighting conditions. We formed a hybrid speechreading system consisting of two time delay neural networks (video and acoustic) and integrated their responses by means of independent opinion pooling - the Bayesian optimal method given conditional independence, which seems to hold for our data. This hybrid system had an error rate 25% lower than that of the acoustic subsystem alone on a five-utterance speaker-independent task, indicating that video can be used to improve speech recognition.

Dynamic Modulation of Neurons and Networks

Eve Marder

Biological neurons have a variety of intrinsic properties because of the large number of voltage dependent currents that control their activity. Neuromodulatory substances modify both the balance of conductances that determine intrinsic properties and the strength of synapses. These mechanisms alter circuit dynamics, and suggest that functional circuits exist only in the modulatory environment in which they operate.

Training Neural Networks with Deficient Data

Volker Tresp, Subutai Ahmad, Ralph Neuneier

We analyze how data with uncertain or missing input features can be incorporated into the training of a neural network. The general solution requires a weighted integration over the unknown or uncertain input although computationally cheaper closed-form solutions can be found for certain Gaussian Basis Function (GBF) networks. We also discuss cases in which heuristical solutions such as substituting the mean of an unknown input can be harmful.

When will a Genetic Algorithm Outperform Hill Climbing

Melanie Mitchell, John Holland, Stephanie Forrest

We analyze a simple hill-climbing algorithm (RMHC) that was previously shown to outperform a genetic algorithm (GA) on a simple "Royal Road" function. We then analyze an "idealized" genetic algorithm (IGA) that is significantly faster than RMHC and that gives a lower bound for GA speed. We identify the features of the IGA that give rise to this speedup, and discuss how these features can be incorporated into a real GA.

Unsupervised Parallel Feature Extraction from First Principles

Mats Österberg, Reiner Lenz

We describe a number of learning rules that can be used to train unsupervised parallel feature extraction systems. The learning rules are derived using gradient ascent of a quality function. We consider a number of quality functions that are rational functions of higher order moments of the extracted feature values. We show that one system learns the principle components of the correlation matrix. Principal component analysis systems are usually not optimal feature extractors for classification. Therefore we design quality functions which produce feature vectors that support unsupervised classification. The properties of the different systems are compared with the help of different artificially designed datasets and a database consisting of all Munsell color spectra.

Classification of Electroencephalogram using Artificial Neural Networks

A C Tsoi, D S C So, A Sergejew

In this paper, we will consider the problem of classifying electroencephalogram (EEG) signals of normal subjects, and subjects suffering from psychiatric disorder, e.g., obsessive compulsive disorder, schizophrenia, using a class of artificial neural networks, viz., multi-layer perceptron. It is shown that the multilayer perceptron is capable of classifying unseen test EEG signals to a high degree of accuracy.

Comparison Training for a Rescheduling Problem in Neural Networks

Didier Keymeulen, Martine de Gerlache

Airline companies usually schedule their flights and crews well in advance to optimize their crew pool activities. Many events such as flight delays or the absence of a member require the crew pool rescheduling team to change the initial schedule (rescheduling). In this paper, we show that the neural network comparison paradigm applied to the backgammon game by Tesauro (Tesauro and Sejnowski, 1989) can also be applied

to the rescheduling problem of an aircrew pool. Indeed both problems correspond to choosing the best solution from a set of possible ones without ranking them (called here best choice problem). The paper explains from a mathematical point of view the architecture and the learning strategy of the backpropagation neural network used for the best choice problem. We also show how the learning phase of the network can be accelerated. Finally we apply the neural network model to some real rescheduling problems for the Belgian Airline (Sabena).

Bounds on the complexity of recurrent neural network implementations of finite state machines

Bill Horne, Don Hush

In this paper the efficiency of recurrent neural network implementations of m -state finite state machines will be explored. Specifically, it will be shown that the node complexity for the unrestricted case can be bounded above by $O(m \log m)$. It will also be shown that the node complexity is $O(m \log m)$ when the weights and thresholds are restricted to the set $\{-1, 1\}$, and $O(m)$ when the fan-in is restricted to two. Matching lower bounds will be provided for each of these upper bounds assuming that the state of the FSM can be encoded in a subset of the nodes of size $\log m$.

Classification of Multi-Spectral Pixels by the Binary Diamond Neural Network
Yehuda Salu

A new neural network, the Binary Diamond, is presented and its use as a classifier is demonstrated and evaluated. The network is of the feed-forward type. It learns from examples in the 'one shot' mode, and recruits new neurons as needed. It was tested on the problem of pixel classification, and performed well. Possible applications of the network in associative memories are outlined.

Optimal Brain Surgeon: Extensions and performance comparisons

Babak Hassibi, David Stork, Gregory Wolff

a second-order

Agnostic PAC-Learning of Functions on Analog Neural Nets

Wolfgang Maass

There exist a number of negative results ([J], [BR], [KV]) about learning on neural nets in Valiant's model [V] for probably approximately correct learning ("PAC-learning"). These negative results are based on an asymptotic analysis where one lets the number of nodes in the neural net go to infinity. Hence this analysis is less adequate for the investigation of learning on a small fixed neural net.

with relatively few analog inputs (e.g. the principal components of some sensory data). The latter type of learning problem gives rise to a different kind of asymptotic question: Can the true error of the neural net be brought arbitrarily close to that of a neural net with "optimal" weights through sufficiently long training? In this paper we employ some new arguments in order to give a positive answer to this question in Haussler's rather realistic refinement of Valiant's model for PAC-learning ([H], [KSS]). In this more realistic model no a-priori assumptions are required about the "learning target", noise is permitted in the training data, and the inputs and outputs are not restricted to boolean values. As a special case our result implies one of the first positive results about learning on multi-layer neural nets in Valiant's original PAC-learning model. At the end of this paper we will describe an efficient parallel implementation of this new learning algorithm.

Mixtures of Controllers for Jump Linear and Non-Linear Plants

Timothy Cacciatore, Steven Nowlan

We describe an extension to the Mixture of Experts architecture for modelling and controlling dynamical systems which exhibit multiple modes of behavior. This extension is based on a Markov process model, and suggests a recurrent network for gating a set of linear or non-linear controllers. The new architecture is demonstrated to be capable of learning effective control strategies for jump linear and non-linear plants with multiple modes of behavior.

A Hybrid Radial Basis Function Neurocomputer and Its Applications

Steven Watkins, Paul Chau, Raoul Tawel, Bjorn Lambrigtsen, Mark Plutowski

A neurocomputer was implemented using radial basis functions and a combination of analog and digital VLSI circuits. The hybrid system uses custom analog circuits for the input layer and a digital signal processing board for the hidden and output layers. The system combines the advantages of both analog and digital circuits, featuring low power consumption while minimizing overall system error. The analog circuits have been fabricated and tested, the system has been built, and several applications have been executed on the system. One application provides significantly better results for a remote sensing problem than have been previously obtained using conventional methods.

Resolving motion ambiguities

K. I. Diamantaras, D. Geiger

We address the problem of optical flow reconstruction and in particular the problem of resolving ambiguities near edges. They occur due to (i) the aperture problem and (ii) the occlusion problem, where pixels on both sides of an intensity edge are assigned the same velocity estimates (and confidence). However, these measurements are correct for just one side of the edge (the non occluded one). Our approach is to introduce an uncertainty field with respect to the estimates and confidence measures. We note that the confidence measures are large at intensity edges and larger at the convex sides of the edges, i.e. inside corners, than at the concave side. We resolve the ambiguities through local interactions via coupled Markov random fields (MRF). The result is the detection of motion for regions of images with large global convexity.

Developing Population Codes by Minimizing Description Length

Richard Zemel, Geoffrey E. Hinton

The Minimum Description Length principle (MDL) can be used to train the hidden units of a neural network to extract a representation that is cheap to describe but nonetheless allows the input to be reconstructed accurately. We show how MDL can be used to develop highly redundant population codes. Each hidden unit has a location in a low-dimensional implicit space. If the hidden unit activities form a bump of a standard shape in this space, they can be cheaply encoded by the center of this bump. So the weights from the input units to the hidden units in an autoencoder are trained to make the activities form a standard bump. The coordinates of the hidden units in the implicit space are also learned, thus allowing flexibility, as the network develops a discontinuous topography when presented with different input classes. Population-coding in a space other than the input enables a network to extract nonlinear higher-order properties of the inputs.

Learning in Computer Vision and Image Understanding

Hayit Greenspan

Requests for name changes in the electronic proceedings will be accepted with no questions asked. However name changes may cause bibliographic tracking issues.

Authors are asked to consider this carefully and discuss it with their co-authors prior to requesting a name change in the electronic proceedings.

Neural Network Definitions of Highly Predictable Protein Secondary Structure Classes

Alan Lapedes, Evan Steeg, Robert Farber

We use two co-evolving neural networks to determine new classes of protein secondary structure which are significantly more predictable from local amino sequence than the conventional secondary structure classification. Accurate prediction of the conventional secondary structure classes: alpha helix, beta strand, and coil, from primary sequence has long been an important problem in computational molecular biology. Neural networks have been a popular method to attempt to predict these conventional secondary structure classes. Accuracy has been disappointingly low. The algorithm presented here uses neural networks to simultaneously examine both sequence and structure data, and to evolve new classes of secondary structure that can be predicted from sequence with significantly higher accuracy than the conventional classes. These new classes have both similarities to, and differences with the conventional alpha helix, beta strand and coil.

Coupled Dynamics of Fast Neurons and Slow Interactions

A.C.C. Coolen, R. Penney, D. Sherrington

A simple model of coupled dynamics of fast neurons and slow interactions, modelling self-organization in recurrent neural networks, leads naturally to an effective statistical mechanics characterized by a partition function which is an average over a replicated system. This is reminiscent of the replica trick used to study spin-glasses, but with the difference that the number of replicas has a physical meaning as the ratio of two temperatures and can be varied throughout the whole range of real values. The model has interesting phase consequences as a function of varying this ratio and external stimuli, and can be extended to a range of other models.

Using Local Trajectory Optimizers to Speed Up Global Optimization in Dynamic Programming

Christopher Atkeson

Dynamic programming provides a methodology to develop planners and controllers for nonlinear systems. However, general dynamic programming is computationally intractable. We have developed procedures that allow more complex planning and control problems to be solved. We use second order local trajectory optimization to generate locally optimal plans and local models of the value function and its derivatives. We maintain global consistency of the local models of the value function, guaranteeing that our locally optimal plans are actually globally optimal, up to the resolution of our search procedures.

Connectionist Modeling and Parallel Architectures

Joachim Diederich, Ah Tsoi

The introduction of specialized hardware platforms for connectionist modeling

("connectionist supercomputer") has created a number of research topics. Some of

these issues are controversial, e.g. the efficient implementation of incremental learning techniques,

the need for the dynamic reconfiguration of networks and possible programming environments for these machines.

Structural and Behavioral Evolution of Recurrent Networks

Gregory Saunders, Peter Angeline, Jordan Pollack

This paper introduces GNARL, an evolutionary program which induces recurrent neural networks that are structurally unconstrained. In contrast to constructive and destructive algorithms, GNARL employs a population of networks and uses a fitness function's unsupervised feedback to guide search through

network space. Annealing is used in generating both gaussian weight changes and structural modifications. Applying GNARL to a complex search and collection task demonstrates that the system is capable of inducing networks with complex internal dynamics.

Two-Dimensional Object Localization by Coarse-to-Fine Correlation Matching

Chien-Ping Lu, Eric Mjolsness

We present a Mean Field Theory method for locating two-dimensional objects that have undergone rigid transformations. The resulting algorithm is a form of coarse-to-fine correlation matching. We first consider problems of matching synthetic point data, and derive a point matching objective function. A tractable line segment matching objective function is derived by considering each line segment as a dense collection of points, and approximating it by a sum of Gaussians. The algorithm is tested on real images from which line segments are extracted and matched.

WATTLE: A Trainable Gain Analogue VLSI Neural Network

Richard Coggins, Marwan Jabri

This paper describes a low power analogue VLSI neural network called Wattle. Wattle is a 10:6:4 three layer perceptron with multiplying DAC synapses and on chip switched capacitor neurons fabricated in 1.2um CMOS. The on chip neurons facilitate variable gain per neuron and lower energy/connection than for previous designs. The intended application of this chip is Intracardiac Electrogram classification as part of an implantable pacemaker / defibrillator system. Measurements of the chip indicate that 10pJ per connection is achievable as part of an integrated system. Wattle has been successfully trained in loop on parity 4 and ICEG morphology classification problems.

Neurobiology, Psychophysics, and Computational Models of Visual Attention

Ernst Niebur, Bruno Olshausen

Requests for name changes in the electronic proceedings will be accepted with no questions asked. However name changes may cause bibliographic tracking issues.

Authors are asked to consider this carefully and discuss it with their co-authors prior to requesting a name change in the electronic proceedings.

Bayesian Backpropagation Over I-O Functions Rather Than Weights

David H. Wolpert

The conventional Bayesian justification of backprop is that it finds the MAP weight vector. As this paper shows, to find the MAP i-o function instead one must add a correction term to backprop. That term biases one towards i-o functions with small description lengths, and in particular favors (some kinds of) feature-selection, pruning, and weight-sharing.

Fast Non-Linear Dimension Reduction

Nanda Kambhathla, Todd Leen

We present a fast algorithm for non-linear dimension reduction. The algorithm builds a local linear model of the data by merging PCA with clustering based on a new distortion measure. Experiments with speech and image data indicate that the local linear algorithm produces encodings with lower distortion than those built by five layer auto-associative networks. The local linear algorithm is also more than an order of magnitude faster to train.

Observability of Neural Network Behavior

Max Garzon, Fernanda Botelho

We prove that except possibly for small exceptional sets, discrete time analog neural nets are globally observable, i.e. all their corrupted pseudo-orbits on computer simulations actually reflect the true dynamical behavior of the network. Locally finite discrete (boolean) neural networks are obser

vable without exception.

Optimal Signalling in Attractor Neural Networks

Isaac Meilijson, Eytan Ruppin

In [Meilijson and Ruppin, 1993] we presented a methodological framework describing the two-iteration performance of Hopfield(cid:173) like attractor neural networks with history-dependent, Bayesian dynamics. We now extend this analysis in a number of directions: input patterns applied to small subsets of neurons, general con(cid:173)nectivity architectures and more efficient use of history. We show that the optimal signal (activation) function has a slanted sigmoidal shape, and provide an intuitive account of activation functions with a non-monotone shape. This function endows the model with some properties characteristic of cortical neurons' firing.

Neural Network Exploration Using Optimal Experiment Design

David Cohn

Consider the problem of learning input/output mappings through exploration, e.g. learning the kinematics or dynamics of a robotic manipulator. If actions are expensive and computation is cheap, then we should explore by selecting a trajectory through the in(cid:173)put space which gives us the most amount of information in the fewest number of steps. I discuss how results from the field of opti(cid:173)mal experiment design may be used to guide such exploration, and demonstrate its use on a simple kinematics problem.

Directional Hearing by the Mauthner System

Audrey Guzik, Robert Eaton

We provide a computational description of the function of the Mauthner system. This is the brainstem circuit which initiates fast start escapes in teleost fish in response to sounds. Our simulations, using back propagation in a realistically constrained feedforward network, have generated hypotheses which are directly interpretable in terms of the activity of the auditory nerve fibers, the principle cells of the system and their associated inhibitory neurons.

Asynchronous Dynamics of Continuous Time Neural Networks

Xin Wang, Qingnan Li, Edward Blum

Motivated by mathematical modeling, analog implementation and distributed simulation of neural networks, we present a definition of asynchronous dynamics of general CT dynamical systems defined by ordinary differential equations, based on notions of local times and communication times. We provide some preliminary results on globally asymptotical convergence of asynchronous dynamics for contractive and monotone CT dynamical systems. When applying the results to neural networks, we obtain some conditions that ensure additive-type neural networks to be asynchronizable.

Development of Orientation and Ocular Dominance Columns in Infant Macaques

Klaus Obermayer, Lynne Kiorpes, Gary Blasdel

Maps of orientation preference and ocular dominance were recorded optically from the cortices of 5 infant macaque monkeys, ranging in age from 3.5 to 14 weeks. In agreement with previous observations, we found that basic features of orientation and ocular dominance maps, as well as correlations between them, are present and robust by 3.5 weeks of age. We did observe changes in the strength of ocular dominance signals, as well as in the spacing of ocular dominance bands, both of which increased steadily between 3.5 and 14 weeks of age. The latter finding suggests that the adult spacing of ocular dominance bands depends on cortical growth in neonatal animals. Since we found no corresponding increase in the spacing of orientation preferences, however, there is a possibility that the orientation prefer

ences of some cells change as the cortical surface expands. Since correlations between the patterns of orientation selectivity and ocular dominance are present at an age, when the visual system is still immature, it seems more likely that their development may be an innate process and may not require extensive visual experience.

Counting function theorem for multi-layer networks

Adam Kowalczyk

We show that a randomly selected N -tuple x of points of R^n with probability > 0 is such that any multi-layer perceptron with the first hidden layer composed of h_1 threshold logic units can implement exactly 2^{h_1} different dichotomies of x . If $N > h_1$ then such a perceptron must have all units of the first hidden layer fully connected to inputs. This implies the maximal capacities (in the sense of Cover) of 2^n input patterns per hidden unit and 2 input patterns per synaptic weight of such networks (both capacities are achieved by networks with single hidden layer and are the same as for a single neuron). Comparing these results with recent estimates of VC-dimension we find that in contrast to the single neuron case, for sufficiently large n and h_1 , the VC-dimension exceeds Cover's capacity.

Supervised Learning with Growing Cell Structures

Bernd Fritzke

We present a new incremental radial basis function network suitable for classification and regression problems. Center positions are continuously updated through soft competitive learning. The width of the radial basis functions is derived from the distance to topological neighbors. During the training the observed error is accumulated locally and used to determine where to insert the next unit. This leads (in case of classification problems) to the placement of units near class borders rather than near frequency peaks as is done by most existing methods. The resulting networks need few training epochs and seem to generalize very well. This is demonstrated by examples.

A Learning Analog Neural Network Chip with Continuous-Time Recurrent Dynamics

Gert Cauwenberghs

We present experimental results on supervised learning of dynamical features in an analog VLSI neural network chip. The recurrent network, containing six continuous-time analog neurons and 42 free parameters (connection strengths and thresholds), is trained to generate time-varying outputs approximating given periodic signals presented to the network. The chip implements a stochastic perturbative algorithm, which observes the error gradient along random directions in the parameter space for error-descent learning. In addition to the integrated learning functions and the generation of pseudo-random perturbations, the chip provides for teacher forcing and long-term storage of the volatile parameters. The network learns a 1 kHz circular trajectory in 100 sec. The chip occupies $2\text{mm} \times 2\text{mm}$ in a $2\mu\text{m}$ CMOS process, and dissipates 1.2 mW.

Functional Models of Selective Attention and Context Dependency

Thomas Hildebrandt

Requests for name changes in the electronic proceedings will be accepted with no questions asked. However name changes may cause bibliographic tracking issues.

Authors are asked to consider this carefully and discuss it with their co-authors prior to requesting a name change in the electronic proceedings.

Learning Classification with Unlabeled Data

Virginia de

One of the advantages of supervised learning is that the final error metric is available during training. For classifiers, the algorithm can directly reduce the number of misclassifications on the training set. Unfortunately

cid:173) nately, when modeling human learning or constructing classifiers for au (cid:173) tonomous robots, supervisory labels are often not available or too ex(cid:173) pensive. In this paper we show that we can substitute for t he labels by making use of structure between the pattern distributions to diffe rent sen(cid:173) sory modalities. We show that minimizing the disagreement bet ween the outputs of networks processing patterns from these different modalitie s is a sensible approximation to minimizing the number of misclassifications i n each modality, and leads to similar results. Using the Peterson-Barney vowel dataset we show that the algorithm performs well in finding ap(cid:173) pr opriate placement for the codebook vectors particularly when the con(cid:173) fu seable classes are different for the two modalities.

Temporal Difference Learning of Position Evaluation in the Game of Go

Nicol Schraudolph, Peter Dayan, Terrence J. Sejnowski

The game of Go has a high branching factor that defeats the tree search appro ach used in computer chess, and long-range spa(cid:173) tiotemporal inter actions that make position evaluation extremely difficult. Development of conv entional Go programs is hampered by their knowledge-intensive nature. We de monstrate a viable alternative by training networks to evaluate Go positions v ia tem(cid:173) poral difference (TD) learning. Our approach is based on ne twork architectures that reflect the spatial organization of both input and reinforcement signals on the Go board, and training protocols that pro vide exposure to competent (though unlabelled) play. These techniques yiel d far better performance than undifferentiated networks trained by self(cid:173) 3) play alone. A network with less than 500 weights learned within 3,000 games of 9x9 Go a position evaluation function that enables a primitive one-ply sear ch to defeat a commercial Go program at a low playing level.

Recovering a Feed-Forward Net From Its Output

Charles Fefferman, Scott Markel

We study feed-forward nets with arbitrarily many layers, using the stan(cid:173) 3) dard sigmoid, $\tanh x$. Aside from technicalities, our theorems are:
1. Complete knowledge of the output of a neural net for arbitrary inputs uniq uely specifies the architecture, weights and thresholds; and 2. There are only finitely many critical points on the error surface for a generic training problem.

Optimal Stochastic Search and Adaptive Momentum

Todd Leen, Genevieve Orr

Stochastic optimization algorithms typically use learning rate schedules that behave asymptotically as $J.t(t) = J.to/t$. The ensem(cid:173) ble dyna mics (Leen and Moody, 1993) for such algorithms provides an easy path to resu lts on mean squared weight error and asymp(cid:173) totic normality. We ap ply this approach to stochastic gradient algorithms with momentum. We sh ow that at late times, learning is governed by an effective learning rate $J.tejJ = J.to/(1 - f3)$ where $f3$ is the momentum parameter. We descri be the behavior of the asymptotic weight error and give conditions on $J.tejJ$ that insure optimal convergence speed. Finally, we use the resu lts to develop an adaptive form of momentum that achieves optimal convergence speed independent of $J.to$.

The Parti-Game Algorithm for Variable Resolution Reinforcement Learning in Multi dimensional State-Spaces

Andrew Moore

Parti-game is a new algorithm for learning from delayed rewards in high dimensional real-valued state-spaces. In high dimensions it is essenti al that learning does not explore or plan over state space uniformly. Part i-game maintains a decision-tree partitioning of state-space and app lies game-theory and computational geom(cid:173) etry techniques to efficie ntly and reactively concentrate high reso(cid:173) lution only on critical

areas. Many simulated problems have been tested, ranging from 2-dimensional to 9-dimensional state-spaces, including mazes, path planning, non-linear dynamics, and uncurl(cid:173)ing snake robots in restricted spaces. In all cases, a good solution is found in less than twenty trials and a few minutes.

Complexity Issues in Neural Computation and Learning

V. P. Roychowdhury, K.-Y. Siu

Requests for name changes in the electronic proceedings will be accepted with no questions asked. However name changes may cause bibliographic tracking issues.

Authors are asked to consider this carefully and discuss it with their co-authors prior to requesting a name change in the electronic proceedings.

Stability and Observability

Max Garzon, Fernanda Botelho

The theme was the effect of perturbations of the defining parameters of a neural net(cid:173)

work due to: 1) measurement (particularly with analog networks); 2) discretization

due to a) digital implementation of analog nets; b) bounded-precision implementation(cid:173)

tion of digital networks; or c) inaccurate evaluation of the transfer function(s); 3)

noise in or incomplete input and/or output of the net or individual cells (particu(cid:173)

larly with analog networks).

The Role of MT Neuron Receptive Field Surrounds in Computing Object Shape from Velocity Fields

G. Buracas, T. Albright

The goal of this work was to investigate the role of primate MT neurons in solving the structure from motion (SFM) problem. Three types of receptive field (RF) surrounds found in area MT neurons (K.Tanaka et al.,1986; Allman et al.,1985) correspond, as our analysis suggests, to the 0th, 1st and 2nd order fuzzy space-differential operators. The large surround/center both differentiate smooth velocity fields and discontinuity detection at boundaries of objects. The model is in agreement with recent psychophysical data on surface interpolation involvement in SFM. We suggest that area MT partially segregates information about object shape from information about spatial relations necessary for navigation and manipulation.

Estimating analogical similarity by dot-products of Holographic Reduced Representations

Tony A. Plate

Models of analog retrieval require a computationally cheap method of estimating similarity between a probe and the candidates in a large pool of memory items.

The vector dot-product operation would be ideal for this purpose if it were possible to encode complex structures as vector representations in such a way that the superficial similarity of vector representations reflected underlying structural similarity. This paper describes how such an encoding is provided by Holographic Reduced Representations (HRRs), which are a method for encoding nested relational structures as fixed-width distributed representations. The conditions under which structural similarity is reflected in the dot-product rankings of HRRs are discussed.

Segmental Neural Net Optimization for Continuous Speech Recognition

Ying Zhao, Richard Schwartz, John Makhoul, George Zavaliagkos

Previously, we had developed the concept of a Segmental Neural Net (SNN) for phonetic modeling in continuous speech recognition (CSR). This kind of neu(cid:173)

ral network technology advanced the state-of-the-art of large-vocabulary CSR, which employs Hidden Markov Models (HMM), for the ARPA 1000-word Re(cid:173) source Management corpus. More Recently, we started porting the neural net system to a larger, more challenging corpus - the ARPA 20,000-word Wall Street Journal (WSJ) corpus. During the porting, we explored the following research directions to refine the system: i) training context-dependent models with a regularization method; ii) training SNN with projection pursuit; and iii) combining different models into a hybrid system. When tested on both a development set and an independent test set, the resulting neural net system alone yielded a performance at the level of the HMM system, and the hybrid SNN/HMM system achieved a consistent 10-15% word error reduction over the HMM system. This paper describes our hybrid system, with emphasis on the optimization methods employed.

Adaptive knot Placement for Nonparametric Regression

Hossein L. Najafi, Vladimir Cherkassky

We show how an "Elman" network architecture, constructed from recurrently connected oscillatory associative memory network modules, can employ selective "attentional" control of synchronization to direct the flow of communication and computation within the architecture to solve a grammatical inference problem. Previously we have shown how the discrete time "Elman" network algorithm can be implemented in a network completely described by continuous ordinary differential equations. The time steps (machine cycles) of the system are implemented by rhythmic variation (clocking) of a bifurcation parameter. In this architecture, oscillation amplitude codes the information content or activity of a module (unit), whereas phase and frequency are used to "software" the network. Only synchronized modules communicate by exchanging amplitude information; the activity of non-resonating modules contributes incoherent crosstalk noise. Attentional control is modeled as a special subset of the hidden modules with outputs which affect the resonant frequencies of other hidden modules. They control synchrony among the other modules and direct the flow of computation (attention) to effect transitions between two subgraphs of a thirteen state automaton which the system emulates to generate a Reber grammar. The internal crosstalk noise is used to drive the required random transitions of the automaton.

Bayesian Modeling and Classification of Neural Signals

Michael Lewicki

Signal processing and classification algorithms often have limited applicability resulting from an inaccurate model of the signal's underlying structure. We present here an efficient, Bayesian algorithm for modeling a signal composed of the superposition of brief, Poisson-distributed functions. This methodology is applied to the specific problem of modeling and classifying extracellular neural waveforms which are composed of a superposition of an unknown number of action potentials (APs). Previous approaches have had limited success due largely to the problems of determining the spike shapes, deciding how many are shapes distinct, and decomposing overlapping APs. A Bayesian solution to each of these problems is obtained by inferring a probabilistic model of the waveform. This approach quantifies the uncertainty of the form and number of the inferred AP shapes and is used to obtain an efficient method for decomposing complex overlaps. This algorithm can extract many times more information than previous methods and facilitates the extracellular investigation of neuronal classes and of interactions within neuronal circuits.

Supervised learning from incomplete data via an EM approach

Zoubin Ghahramani, Michael Jordan

Real-world learning tasks may involve high-dimensional data sets with arbitrary patterns of missing data. In this paper we present a framework based on maximum likelihood density estimation for learning from such data

set.s. We use mixture models for the den(cid:173) sity estimates and make two distinct appeals to the Expectation(cid:173) Maximization (EM) principle (Dempster et al., 1977) in deriving a learning algorithm-EM is used both for the estimation of mix(cid:173) ture components and for coping with missing data. The result(cid:173) ing algorithm is applicable to a wide range of supervised as well as unsupervised learning problems. Results from a classification benchmark-the iris data set-are presented.

Correlation Functions in a Large Stochastic Neural Network

Iris Ginzburg, Haim Sompolinsky

Most theoretical investigations of large recurrent networks focus on the properties of the macroscopic order parameters such as population averaged activities or average overlaps with memories. However, the statistics of the fluctuations in the local activities may be an important testing ground for comparison between models and observed cortical dynamics. We evaluated the neuronal correlation functions in a stochastic network comprising of excitatory and inhibitory populations. We show that when the network is in a stationary state, the cross-correlations are relatively weak, i.e., their amplitude relative to that of the auto-correlations are of order $1/N$, N being the size of the interacting population. This holds except in the neighborhoods of bifurcations to nonstationary states. As a bifurcation point is approached the amplitude of the cross-correlations grows and becomes of order 1 and the decay time constant diverges. This behavior is analogous to the phenomenon of critical slowing down in systems at thermal equilibrium near a critical point. Near a Hopf bifurcation the cross-correlations exhibit damped oscillations.

Clustering with a Domain-Specific Distance Measure

Steven Gold, Eric Mjolsness, Anand Rangarajan

With a point matching distance measure which is invariant under translation, rotation and permutation, we learn 2-D point-set objects, by clustering noisy point-set images. Unlike traditional clustering methods which use distance measures that operate on feature vectors - a representation common to most problem domains - this object-based clustering technique employs a distance measure specific to a type of object within a problem domain. Formulating the clustering problem as two nested objective functions, we derive optimization dynamics similar to the Expectation-Maximization algorithm used in mixture models.

Speaker Recognition Using Neural Tree Networks

Kevin Farrell, Richard Mammone

A new classifier is presented for text-independent speaker recognition. The new classifier is called the modified neural tree network (MNTN). The NTN is a hierarchical classifier that combines the properties of decision trees and feed-forward neural networks. The MNTN differs from the standard NTN in that a new learning rule based on discriminant learning is used, which minimizes the classification error as opposed to a norm of the approximation error. The MNTN also uses leaf probability measures in addition to the class labels. The MNTN is evaluated for several speaker identification experiments and is compared to multilayer perceptrons (MLPs), decision trees, and vector quantization (VQ) classifiers. The VQ classifier and MNTN demonstrate comparable performance and perform significantly better than the other classifiers for this task. Additionally, the MNTN provides a logarithmic saving in retrieval time over that of the VQ classifier. The MNTN and VQ classifiers are also compared for several speaker verification experiments where the MNTN is found to outperform the VQ classifier.

Tonal Music as a Componential Code: Learning Temporal Relationships Between and Within Pitch and Timing Components

Catherine Stevens, Janet Wiles

This study explores the extent to which a network that learns the temporal relationships within and between the component features of Western tonal music can account for music theoretic and psychological phenomena such as the tonal hierarchy and rhythmic expectancies. Predicted and generated sequences were recorded as the representation of a 153-note waltz melody was learnt by a predictive, recurrent network. The network learned transitions and relations between within pitch and timing components: accent and duration values interacted in the development of rhythmic and metric structures and, with training, the network developed chordal expectancies in response to the activation of individual tones. Analysis of the hidden unit representation revealed that musical sequences are represented as transitions between states in hidden unit space.

Connectionist Models for Auditory Scene Analysis

Richard Duda

Although the visual and auditory systems share the same basic tasks of informing an organism about its environment, most connectionist work on hearing to date has been devoted to the very different problem of speech recognition. We believe that the most fundamental task of the auditory system is the analysis of acoustic signals into components corresponding to individual sound sources, which Bregman has called auditory scene analysis. Computational and connectionist work on auditory scene analysis is reviewed, and the outline of a general model that includes these approaches is described.

Learning in Compositional Hierarchies: Inducing the Structure of Objects from Data

Joachim Utans

I propose a learning algorithm for learning hierarchical models for object recognition. The model architecture is a compositional hierarchy that represents part-whole relationships: parts are described in the local context of substructures of the object. The focus of this report is inducing the structure of learning hierarchical models from data, i.e. model prototypes from observed exemplars of an object. At each node in the hierarchy, a probability distribution governing its parameters must be learned. The connections between nodes reflects the structure of the object. The formulation of substructures is encouraged such that their parts become conditionally independent. The resulting model can be interpreted as a Bayesian Belief Network and also is in many respects similar to the stochastic visual grammar described by Mjolsness.

VLSI Phase Locking Architectures for Feature Linking in Multiple Target Tracking Systems

Andreas Andreou, Thomas Edwards

Recent physiological research has shown that synchronization of oscillatory responses in striate cortex may code for relationships between visual features of objects. A VLSI circuit has been designed to provide rapid phase-locking synchronization of multiple oscillators to allow for further exploration of this neural mechanism. By exploiting the intrinsic random transistor mismatch of devices operated in subthreshold, large groups of phase-locked oscillators can be readily partitioned into smaller phase-locked groups. A multiple target tracker for binary images is described utilizing this phase-locking architecture. A VLSI chip has been fabricated and tested to verify the architecture. The chip employs Pulse Amplitude Modulation (PAM) to encode the output at the periphery of the system.

A Network Mechanism for the Determination of Shape-From-Texture

Kô Sakai, Leif Finkel

We propose a computational model for how the cortex discriminates shape and depth

th from texture. The model consists of four stages: (1) extraction of local spatial frequency, (2) frequency characterization, (3) detection of texture compression by normalization, and (4) integration of the normalized frequency over space. The model accounts for a number of psychophysical observations including experiments based on novel random textures. These textures are generated from white noise and manipulated in Fourier domain in order to produce specific frequency spectra. Simulations with a range of stimuli, including real images, show qualitative and quantitative agreement with human perception.

Encoding Labeled Graphs by Labeling RAAM

Alessandro Sperduti

In this paper we propose an extension to the RAAM by Pollack. This extension, the Labeling RAAM (LRAAM), can encode labeled graphs with cycles by representing pointers explicitly. Data encoded in an LRAAM can be accessed by pointer as well as by content. Direct access by content can be achieved by transforming the encoder network of the LRAAM into an analog Hopfield network with hidden units. Different access procedures can be defined depending on the access key. Sufficient conditions on the asymptotical stability of the associated Hopfield network are briefly introduced.

Learning Mackey-Glass from 25 examples, Plus or Minus 2

Mark Plutowski, Garrison Cottrell, Halbert White

We apply active exemplar selection (Plutowski & White, 1991; 1993) to predicting a chaotic time series. Given a fixed set of examples, the method chooses a concise subset for training. Fitting these exemplars results in the entire set being fit as well as desired. The algorithm incorporates a method for regulating network complexity, automatically adding exemplars and hidden units as needed. Fitting examples generated from the Mackey-Glass equation with fractal dimension 2.1 to an rmse of 0.01 required about 25 exemplars and 3 to 6 hidden units. The method requires an order of magnitude fewer floating point operations than training on the entire set of examples, is significantly cheaper than two contenting exemplar selection techniques, and suggests a simpler active selection technique that performs comparably.

Structured Machine Learning for 'Soft' Classification with Smoothing Spline ANOVA and Stacked Tuning, Testing and Evaluation

Grace Wahba, Yuedong Wang, Chong Gu, Ronald Klein, MD, Barbara Klein, MD

We describe the use of smoothing spline analysis of variance (SSANOVA) in the penalized log likelihood context, for learning (estimating) the probability p of a '1' outcome, given a training set with attribute vectors \mathbf{a} and outcomes. p is of the form $\text{pet} = eJ(\mathbf{t}) / (1 + eJ(\mathbf{t}))$, where, if \mathbf{t} is a vector of attributes, f is learned as a sum of smooth functions of one attribute plus a sum of smooth functions of two attributes, etc. The smoothing parameters governing f are obtained by an iterative unbiased risk or iterative GCV method. Confidence intervals for these estimates are available.
