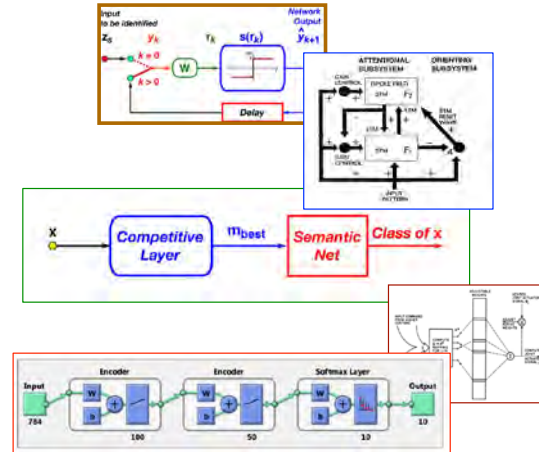


# Neural Networks

Robert Stengel

Robotics and Intelligent Systems, MAE 345,  
Princeton University, 2017

- **Associative/recurrent networks**
  - Hopfield network
  - Adaptive resonance theory network
  - Elman/Jordan networks
- **Unsupervised training**
  - k-means clustering
- **Semi-supervised training**
  - Self-organizing map
- **Cerebellar model articulation controller (CMAC)**
- **Deep learning**
  - Restricted Boltzmann machine
  - Convolutional network
  - Neural Turing Machines



Copyright 2017 by Robert Stengel. All rights reserved. For educational use only.  
<http://www.princeton.edu/~stengel/MAE345.html>

1

## Associative/Recurrent Networks

2

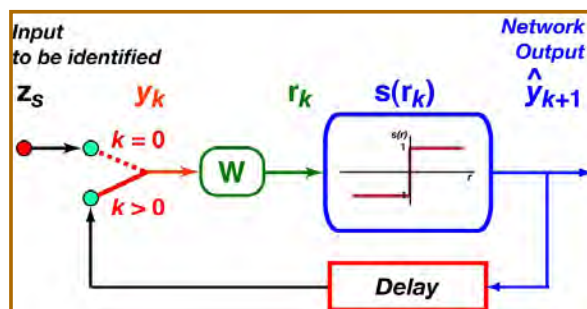
# Associative-Memory Neural Networks

- **Goals**

- Identify symbols from noisy data, given exemplars of possible features
- Retrieve full feature from incomplete samples
  - “To be \_\_\_\_\_”
  - “Snap, crackle, \_\_\_\_\_”
- Build a database from related contextual information, e.g., populate features of one categorical set using features in another



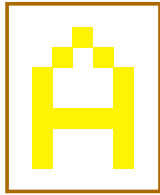
3



## Recurrent Networks

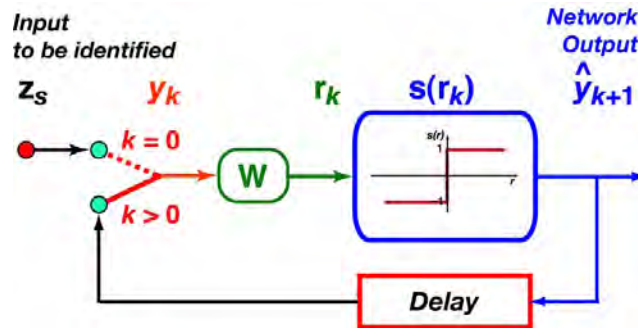
- **Recursion to identify an unknown object**
  - Network is given a single, fixed input, and it iterates to a solution
- **Convergence and stability** of the network are critical issues (*discrete-time dynamic system*)
- Single network may have **many stable states**
  - Classified outputs of the map
  - Pattern recognition with noisy data

4



# Hopfield Network

- **Bipolar (-1,1) inputs and outputs**
  - $\dim(y) = n \times 1$
- **Supervised training with perfect exemplar outputs**
- **Noisy measurement of an exemplar as input to be identified**
- **Network operation**



$$\mathbf{z}_s = \mathbf{y}_s + \mathbf{n}_s$$

$$\hat{\mathbf{y}}_0 = \mathbf{z}_s$$

- **Iterate to convergence**

$$\hat{\mathbf{y}}_{k+1} = \mathbf{s}(\mathbf{r}_k) = \mathbf{s}(\mathbf{W}\hat{\mathbf{y}}_k)$$

$$\hat{y}_{i,k+1} = \begin{cases} 1, & r_{i,k} > 0 \\ \text{Unchanged}, & r_{i,k} = 0 \\ -1, & r_{i,k} < 0 \end{cases}, i = 1 \text{ to } n$$

5

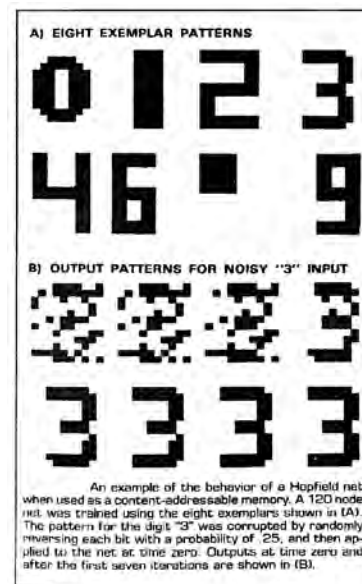
## Training a Hopfield Network

- **Network training**
  - Given  $M$  exemplars,  $\mathbf{y}_s (n \times 1)$ ,  $s = 1, M$
  - Each exemplar is a character represented by  $n$  pixels
  - Batch calculation of weighting matrix

$$\mathbf{W} = \sum_{s=1}^M (\mathbf{y}_s \mathbf{y}_s^T - \mathbf{I}_n)$$

$$= \sum_{s=1}^M \begin{bmatrix} y_1^2 - 1 & y_1 y_2 & \dots \\ y_1 y_2 & y_2^2 - 1 & \dots \\ \dots & \dots & \dots \end{bmatrix}_s$$

- **No iterations to define weights**
- **Large number of weights**
- **Limited number of exemplars ( $< 0.15 n$ )**
- **Similar exemplars pose a problem**



$$n = 120; \quad M = 8$$

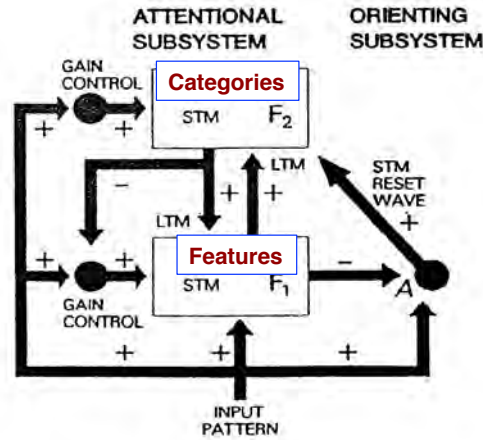
$$\# \text{ weights} = n^2 = 14,400$$

6

# Adaptive Resonance Theory Network

(Grossberg, Carpenter, 1976)

- Self-organizing/stabilizing network for **finding clusters in binary input (ART-1)**
- Broadly based on cerebellar model
  - **Long-Term Memory**
  - **Short-Term Memory**
  - **Stability and plasticity**
  - **Unsupervised and supervised learning**
  - **“Bottom-up” input**
  - **“Top-down” priming**
  - **Pre-cursor to “deep learning”**



7

## ART-1 Network

**Architecture**

**Binary Neurons represent Pattern Pixels**

**Recursive Training Example: adding new templates**

The ART architecture.

The ART1 neural network.

backward LTM from:

input pattern	output 1	output 2	output 3	output 4
C	C	not active	not active	not active
E	E	not active	not active	not active
F	F	not active	not active	not active
F	F	not active	not active	not active

Further Developments:

- Continuous inputs (ART-2)
- Fuzzy Logic (Fuzzy ART)
- Dual-Associative Networks for Pattern Recognition (Lapart, Sandia, 2017)

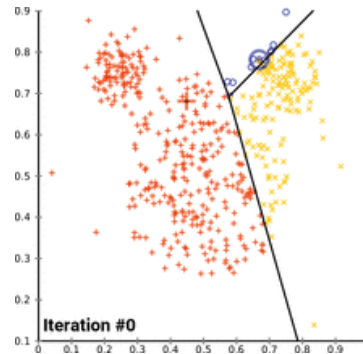
8

## k-Means Clustering

- Least-squares clustering of  $n$  observation sets into  $k$  regions

$$\min_{\mu_i} J = \sum_{i=1}^k \sum_{j=1}^n \|\mathbf{x}_j - \mu_i\|_2$$

- i.e., find **centroids** of each region
- Once centroids are known, boundaries of regions found from **Voronoi diagram**



9

## Self-Organizing Map

(Kohonen, 1981)



- **Competitive, unsupervised learning in 1<sup>st</sup> layer**
- **Premise: input signal patterns that are close produce outputs that are close**
- Ordered inputs produce spatial distribution, i.e., **a map**
- Cells of the map are likened to the cell structure of the **cerebral cortex**

- $x$ : ( $n \times 1$ ) input vector characterizes features (attributes) of a signal
- $m$ : ( $n \times 1$ ) weight vector of a cell that represents an **output class**

10

# Competition in Self-Organizing Map



- Competition is based on minimizing distance from  $x$  to  $m$

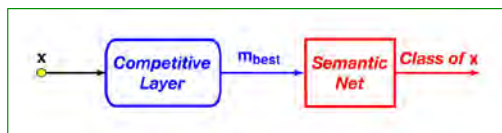
$$Cost = distance = \|x - m_i\|$$

$$\min Cost = \min_{m_i} \|x - m\|$$

- $m$  encodes the output classes
- Supervision:** Semantic Net decodes the output to identify classes

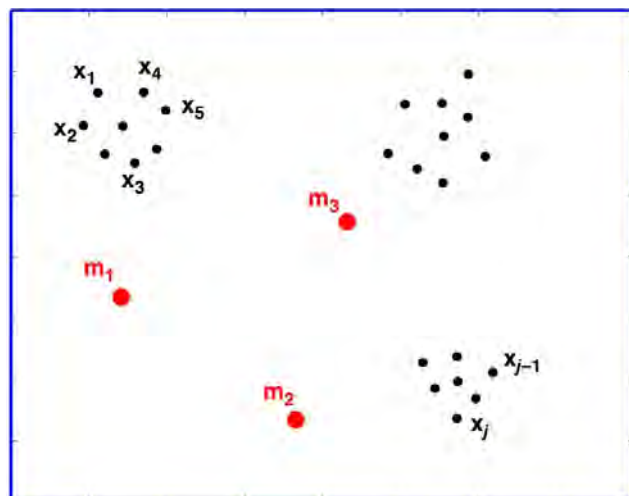
$$m_1 = \begin{bmatrix} 0 \\ 1 \\ 3 \end{bmatrix} \rightarrow \text{Class A}; \quad m_2 = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \rightarrow \text{Class B}$$

11



## Goal of the Self-Organizing Map

- Given:**
  - $l$  output classes
  - Input training set,  $x_j$ ,  $j = 1$  to  $J$
- Find:** Cell weights,  $m_j$ ,  $i = 1$  to  $l$  that best cluster the data (i.e., with minimum norm)
- Initialize the cell weights,  $m_j$ , randomly in the space of  $x$

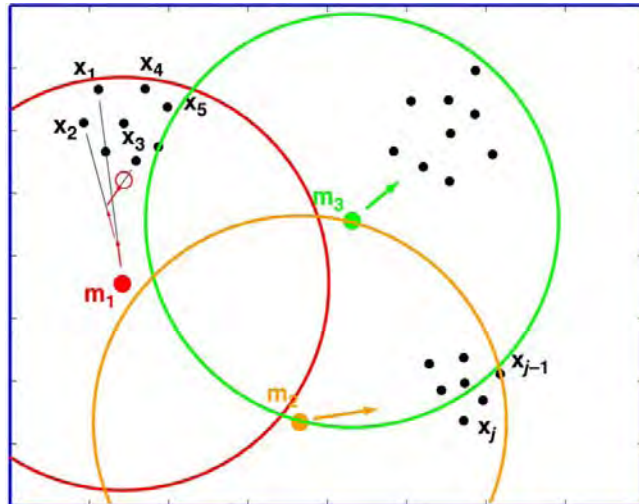


12



## Training the Self-Organizing Map

- Define a neighborhood set within a radius of  $N_c$  around each cell,  $m_i$ 
  - Choose  $N_c$  to overlap with neighboring cells
- Find the **best cell-weight match**,  $m_{\text{best}}$ , (i.e., the closest  $m_i$ ) to the 1<sup>st</sup> training sample,  $x_1$

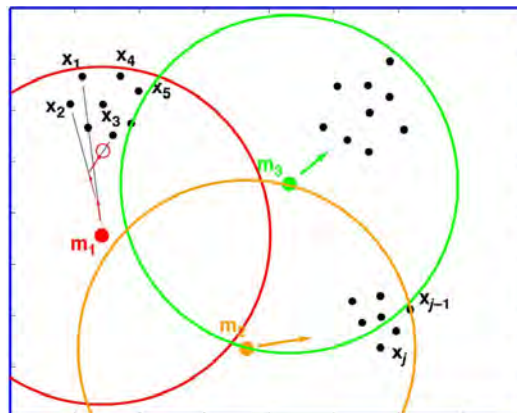


13

## Cell Weight Updates

$$\mathbf{m}_i(k+1) = \begin{cases} \mathbf{m}_i(k) + \alpha_k [\mathbf{x}_1 - \mathbf{m}_i(k)], & \mathbf{m}_i \in N_c \\ \mathbf{m}_i(k), & \mathbf{m}_i \notin N_c \end{cases}$$

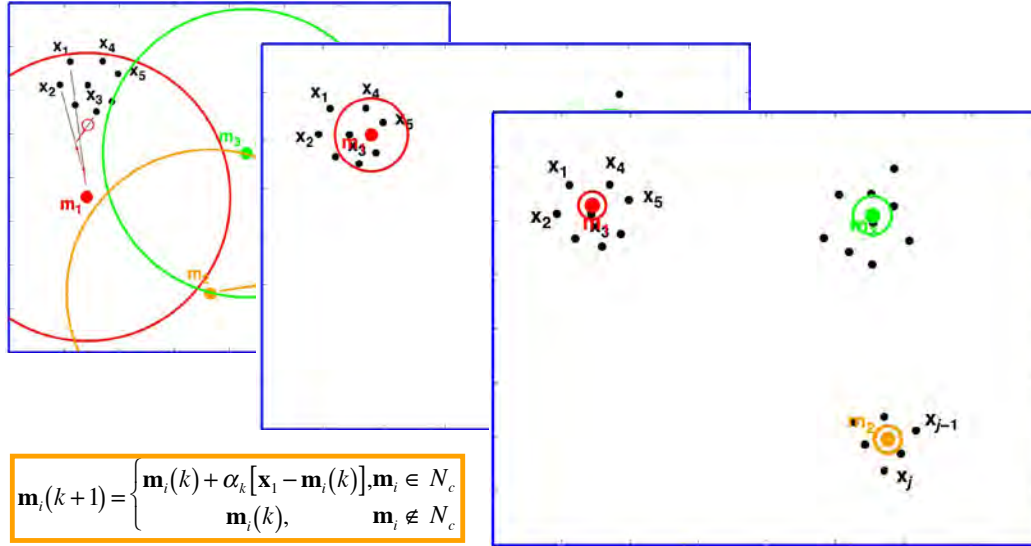
- Update cell weights for all cells in the neighborhood set,  $N_c$ , of  $m_{\text{best}}$ 
  - $\alpha_k$  = adaptation gain or learning rate
- Repeat for
  - $x_2$  to  $x_J$
  - $m_1$  to  $m_I$
- Converse of *particle swarm optimization*



14

# Convergence of Cell Weights

Repeat entire process with decreasing  $N_c$  radius until convergence occurs

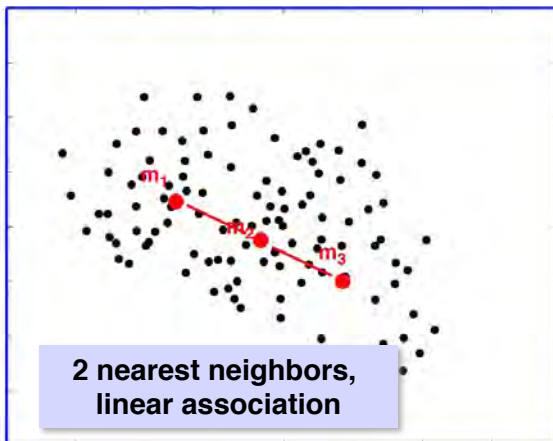


15

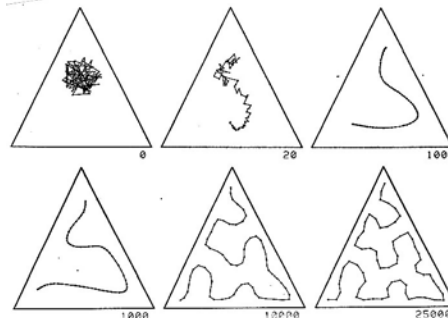
## Semantic Map



- Association of  $\mathbf{m}_{best}$  with categorical information
- Contextual information used to generate map of symbols
- Dimensionality and # of nearest neighbors affects final map
  - Example: linear association of cell weights
  - Points for cell-weight update chosen randomly



Evolution of points on a line that identifies locations of  $\mathbf{m}_i$   
(Uniform random field of data points not shown)

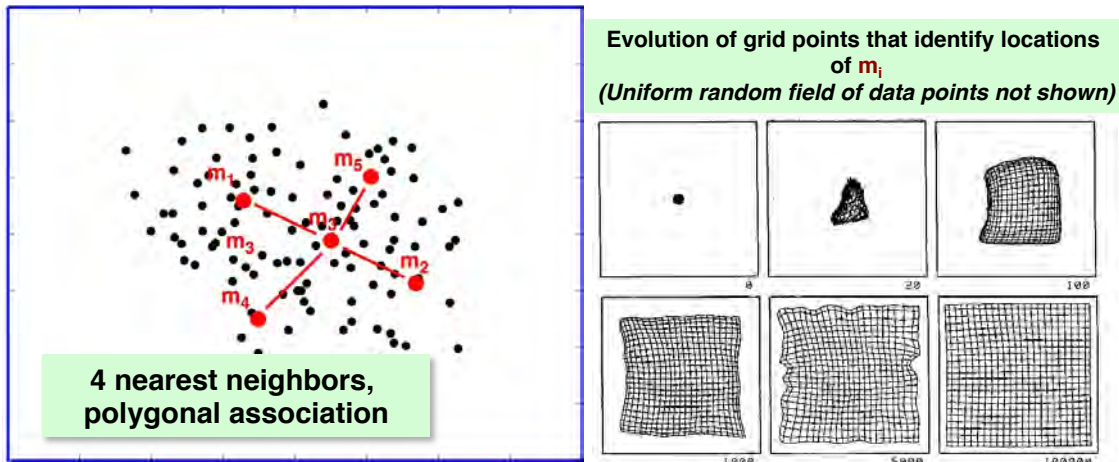


16



# Choice of Neighborhood Architecture

- **Example:** Map is assumed to represent a grid of associated points
- Number of cell weights specified
- Random starting locations for training



17

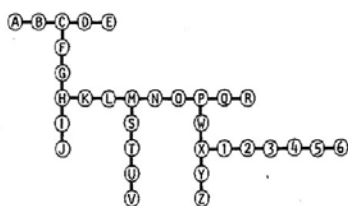
# Minimum Spanning Tree

**Example:** Hexagonal map association identification  
32 points with 5 attributes that may take six values  
(0, 1, 2, 3, 4, 5)

Input Data Matrix

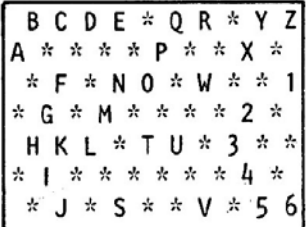
Attribute	Item																																
	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	1	2	3	4	5	6	
$a_1$	1	2	3	4	5	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3
$a_2$	0	0	0	0	0	1	2	3	4	5	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3
$a_3$	0	0	0	0	0	0	0	0	0	0	1	2	3	4	5	6	7	8	3	3	3	3	6	6	6	6	6	6	6	6	6	6	6
$a_4$	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	2	3	4	1	2	3	4	2	2	2	2	2	2	
$a_5$	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	2	3	4	5	6

**Minimum spanning tree:**  
smallest total edge length



Minimal spanning tree corresponding to Table

**Hexagonal lattice of grid points**  
that identify locations of  $m_i$



Self-organized map of the data matrix of Table

18

# Semantic Identification

## Example of semantic identification

Each item for training has symbolic expression and context  
 Categories: noun, verb, adverb

Bob/Jim/Mary	1	<b>Sentence Patterns:</b>		Mary likes meat	
horse/dog/cat	2	1-5-12	1-9-2	2-5-14	Jim speaks well
beer/water	3	1-5-13	1-9-3	2-9-1	Mary likes Jim
meat/bread	4	1-5-14	1-9-4	2-9-2	Jim eats often
runs/walks	5	1-6-12	1-10-3	2-9-3	Mary buys meat
works/speaks	6	1-6-13	1-11-4	2-9-4	dog drinks fast
visits/phones	7	1-6-14	1-10-12	2-10-3	horse hates meat
buys/sells	8	1-6-15	1-10-13	2-10-12	Jim eats seldom
likes/hates	9	1-7-14	1-10-14	2-10-13	Bob buys meat
drinks/eats	10/11	1-8-12	1-11-12	2-10-14	cat walks slowly
much/little	12	1-8-2	1-11-13	2-11-4	Jim eats bread
fast/slowly	13	1-8-3	1-11-14	2-11-12	cat hates Jim
often/seldom	14	1-8-4	2-5-12	2-11-13	Bob sells beer
well/poorly	15	1-9-1	2-5-13	2-11-14	(etc.)

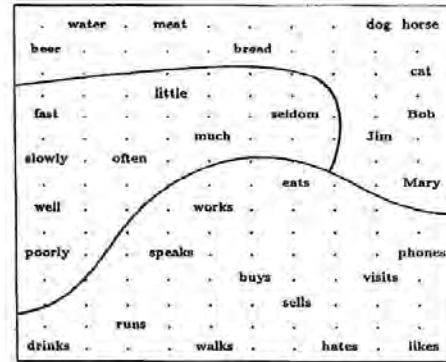
(a)

(b)

(c)

Outline of vocabulary used in this experiment. (a) List of used words (nouns, verbs, and adverbs), (b) sentence patterns, and (c) some examples of generated three-word-sentences.

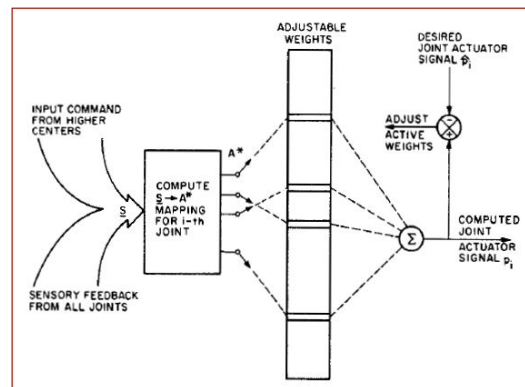
Ritter, Kohonen, 1989



"Semantic map" obtained on a network of 10 × 15 cells after 2000 presentations of word-context-pairs derived from 10 000 random sentences of the kind shown in Fig. 10(c). Nouns, verbs, and adverbs are segregated into different domains. Within each domain a further grouping according to aspects of meaning is discernible.

# Cerebellar Model Articulation Controller (CMAC)

- Another precursor to deep learning
- Inspired by models of human cerebellum
- **CMAC**: Two-stage mapping of a vector input to a scalar output
- **First mapping**: Input space to **association space**
  - $s$  is fixed
  - $a$  is binary
- **Second mapping**: Association space to **output space**
  - $g$  contains learned weights



$$S : X \rightarrow a$$

Input  $\rightarrow$  Selector vector

$$g : a \rightarrow y$$

Selector vector  $\rightarrow$  Output

Albus, 1975

# Example of Single-Input CMAC Association Space

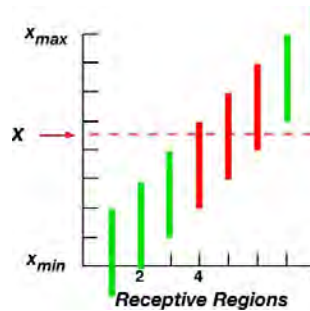
$$s : x \rightarrow \mathbf{a}$$

*Input*  $\rightarrow$  *Selector vector*

- $x$  is in  $(x_{min}, x_{max})$
- Selector vector,  $\mathbf{a}$ , is binary and has  $N$  elements
- Input quantization =  $(x_{max} - x_{min}) / N$
- Receptive regions of association space map  $x$  to  $\mathbf{a}$ 
  - Analogous to neurons that “fire” in response to stimulus
- $N_A =$  Number of receptive regions

$$N_A = N + C - 1 = \dim(\mathbf{a})$$

- $C =$  Generalization parameter = # of overlapping regions



$$\mathbf{a} = \begin{bmatrix} 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 \end{bmatrix}^T$$

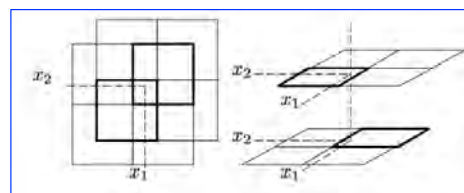
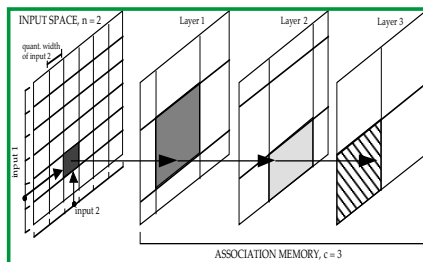
$C = 3$

21

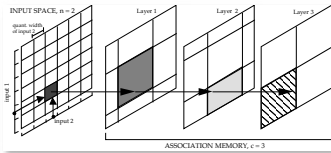
## CMAC Output and Training

- In higher dimensions, association space is  $\dim(\mathbf{x})$ , a plane, cube, or hypercube
- Potentially large memory requirements
- Granularity (quantization) of output
- Variable generalization and granularity

*2-dimensional association space*  
*Rectangular receptive regions*



22



# CMAC Output and Training

- CMAC output,  $y$ , (i.e., control command) from activated cells of  $c$  Associative Memory layers

$$y_{CMAC} = \mathbf{w}^T \mathbf{a} = \sum_{i=j}^{j+C-1} w_{i_{activated}}$$

$j$  = index of first activated region

- Least-squares training of CMAC weights,  $w$ 
  - Analogous to synapses between neurons

$$w_{j_{new}} = w_{j_{old}} + \frac{\beta}{c} \left( y_{desired} - \sum_{i=1}^c w_{i_{old}} \right)$$

$\beta$  is the learning rate and  $w_j$  is an activated cell weight

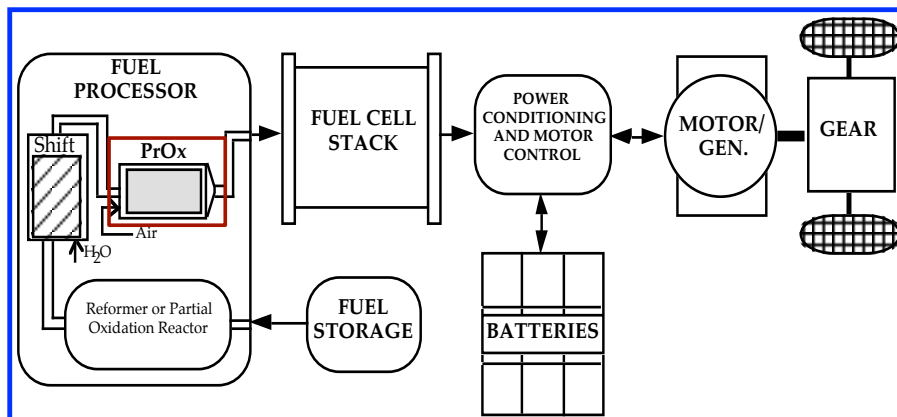
- Localized generalization and training

23

## CMAC Control of a Fuel-Cell Pre-Processor

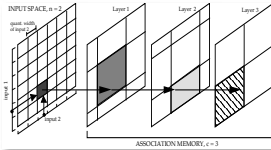
(Iwan and Stengel)

Fuel cell produces electricity for electric motor

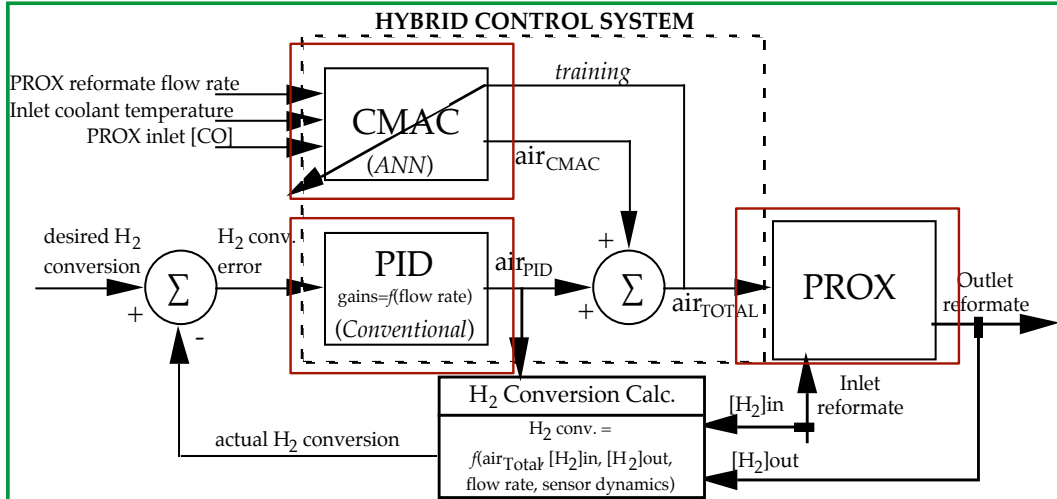


Pre-processor produces hydrogen for the fuel cell and carbon monoxide, which “poisons” the fuel cell catalyst

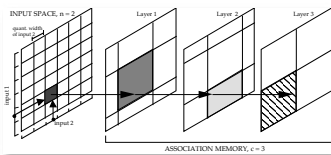
24



# CMAC/PID Control System for Preferential Oxidizer



25



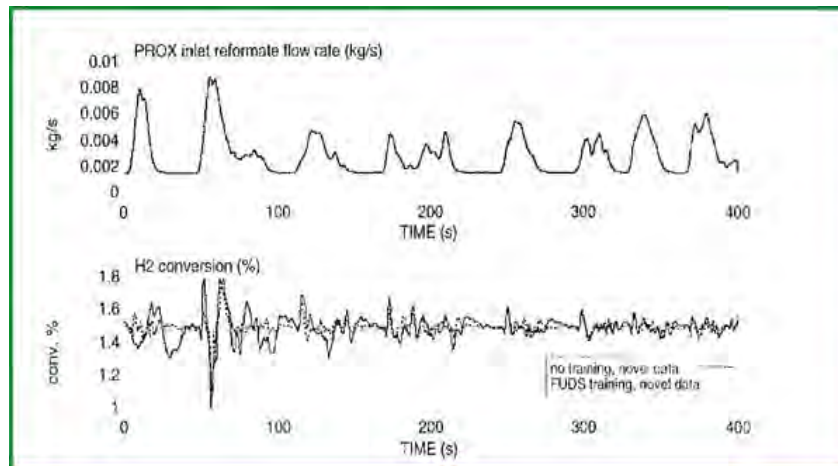
# Summary of 3-D CMAC Characteristics

- **Inputs and Number of Divisions for receptor cubes:**
  - PrOx inlet reformat flow rate (95)
  - PrOx inlet cooling temperature (80)
  - PrOx inlet CO concentration (100)
- **Output: PrOx air injection rate**
- **Associative Layers, C: 24**
- **Number of Associative Memory Cells/Weights and Layer Offsets: 1,276 and [1,5,7]**
- **Learning Rate,  $\alpha$  : ~0.01**
- **Sampling Interval: 100 ms**

26

# Flow Rate and Hydrogen Conversion of CMAC/PID Controller

- H<sub>2</sub> conversion command (across PrOx only): 1.5%
- Novel data, with (---) and without pre-training (—)
- Federal Urban Driving Cycle (= FUDS)



27

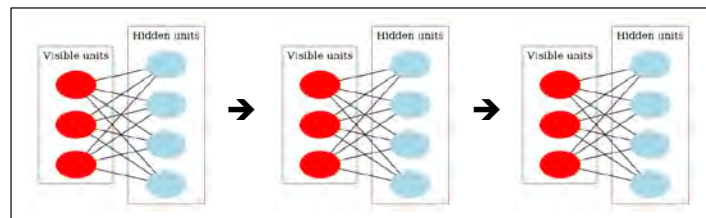
## Comparison of PrOx Controllers on Federal Urban Driving Cycle

	mean H <sub>2</sub> error		maximum H <sub>2</sub> error		mean CO out	
	%	%	ppm	ppm	max. CO out	%
• Fixed-Air	0.68	0.87	6.3	28	57.2	
• Table Look-up	0.13	1.43	6.5	26	57.8	
• PID	0.05	0.51	7.7	30	58.1	
• CMAC/PID	0.02	0.16	7.3	26	58.1	
						net H <sub>2</sub> output

28

# Deep Learning with Restricted Boltzmann Machine

- Multiple layers of RBMs
- **Semi-supervised learning**
  - Clustering (visible) units
  - Sigmoid (hidden) units
- Pre-train each layer separately and contextually (*unsupervised*)
- Fine-tune with backpropagation (*supervised*)
- Restrict connections between layers
- Goal is to overcome “*vanishing or exploding gradient problem*” in multi-layer back-propagation

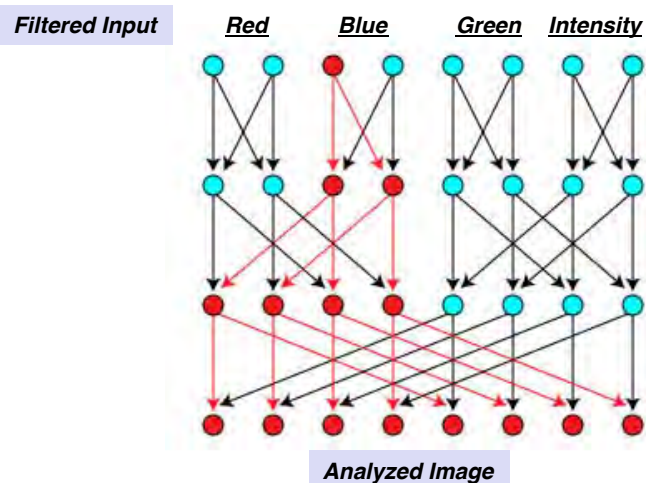


Hinton et al, 2006

29

## Sparse Deep Network

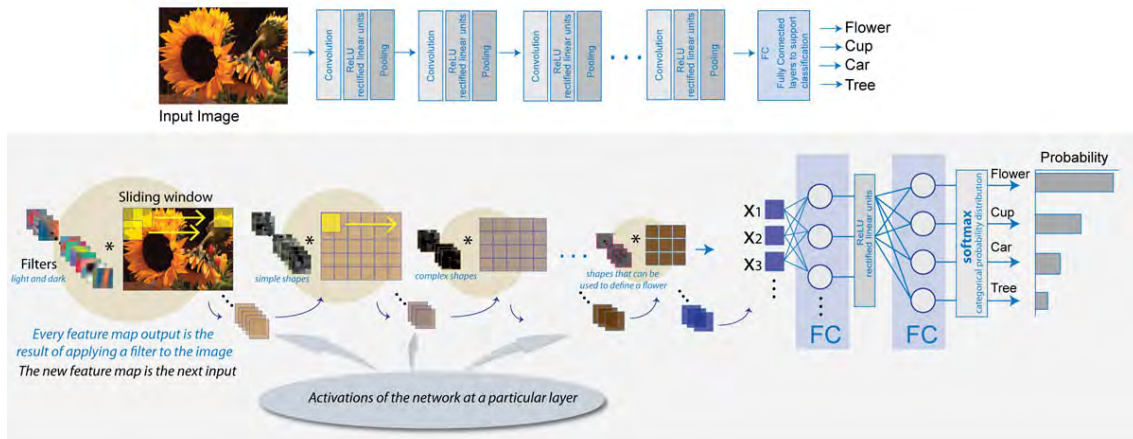
- Partitioned input space
- Expanding network connections



- Fully connected final layer

30

# Convolutional Neural Network

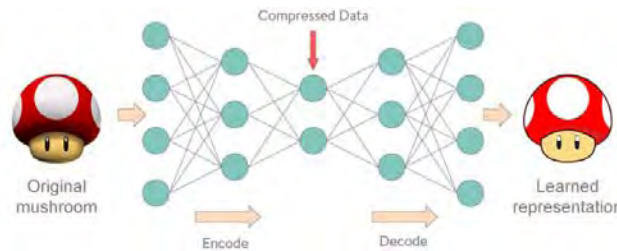


- **Decomposition of image**
  - Sliding window of receptive fields
  - Pooling (dimension reduction)
  - Simply connected networks
- **Repeated sequence of operations**
  - Convolution (cross-correlation)
  - Rectification neurons (ReLU)
  - Fully connected networks

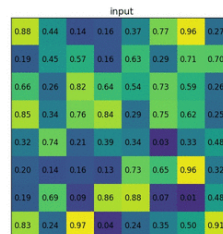
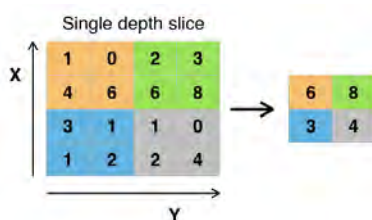
31

# Autoencoding and Pooling

- **Autoencoding:** Same number of inputs and outputs
- **Compression and decompression layers identify important attributes**



- **Max pooling:** selection of important attributes
- **Dimension reduction of features**
- **Enhanced invariance in characterization of a feature in different perspectives**



32



# Convolution

- Cross-correlation of outputs from previous layers
- Apply to partitioned receptive fields

“Heat Map”

$$C_S = D^T D$$

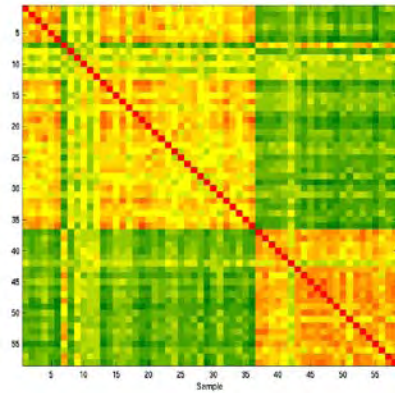
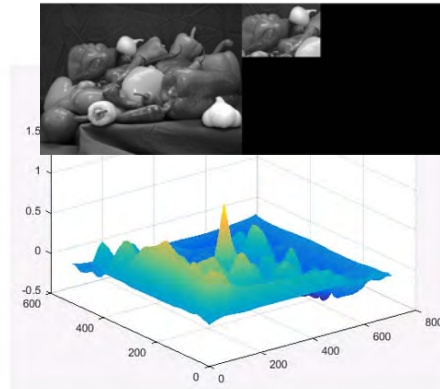


Image Convolution

$$z_k(x_k, y_k) = \sum_{i=0}^K c_k [x_k y_{k+1}]$$

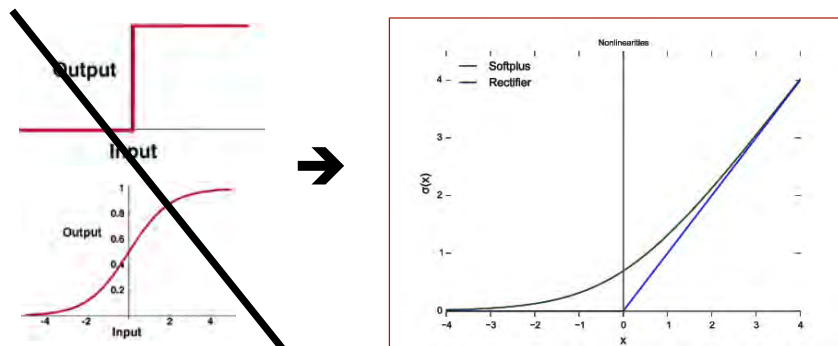


33

## Rectified Linear Unit (ReLU)

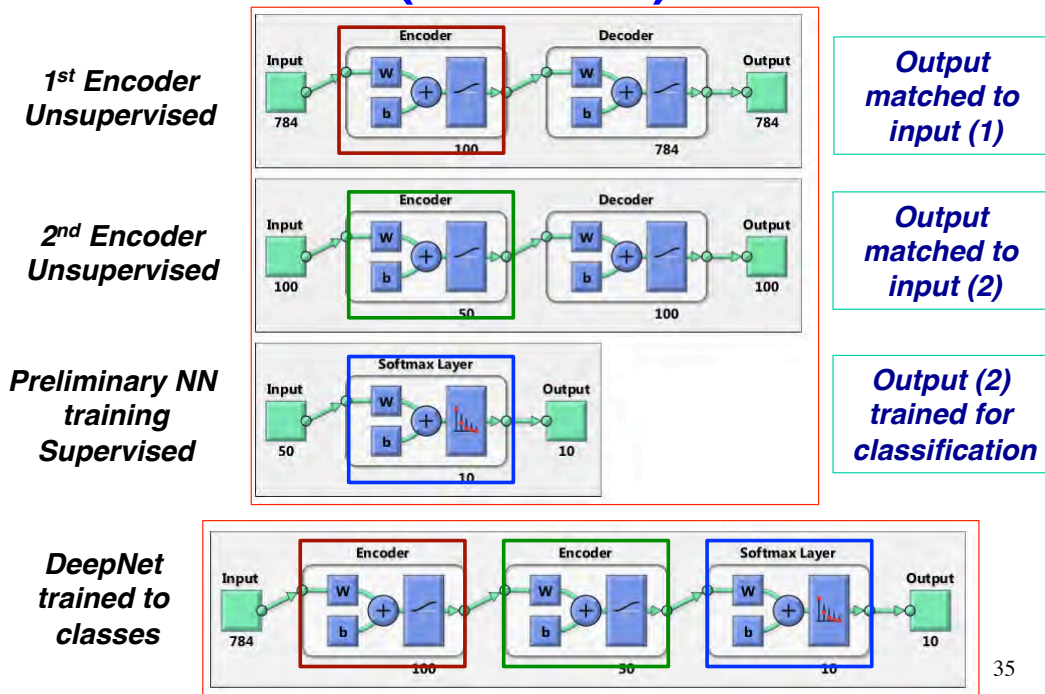
- Simple alternative to hardlim, sigmoid nodes
- Faster, more accurate classification in some applications

$$y = \max(0, y)$$



34

# Convolution Neural Network (ConvNet)



35

## More on Recurrent Neural Networks

- Feedback added to a feed-forward neural network (*discrete-time dynamic system*)
- One-step memory introduced to network

### Elman Network

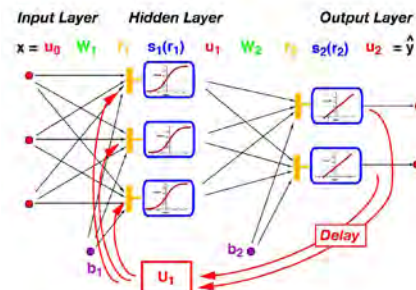
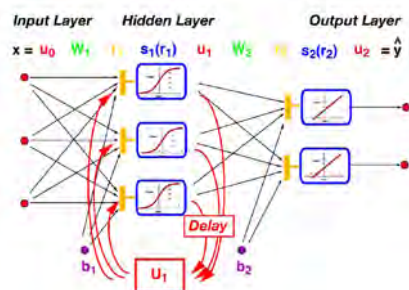
$$\mathbf{u}_1(k) = \mathbf{s}_1 [\mathbf{W}_1 \mathbf{x}(k) + \mathbf{U}_1 \mathbf{u}_1(k-1) + \mathbf{b}_1]$$

$$\mathbf{u}_2(k) = \mathbf{s}_2 [\mathbf{W}_2 \mathbf{u}_1(k) + \mathbf{b}_2]$$

### Jordan Network

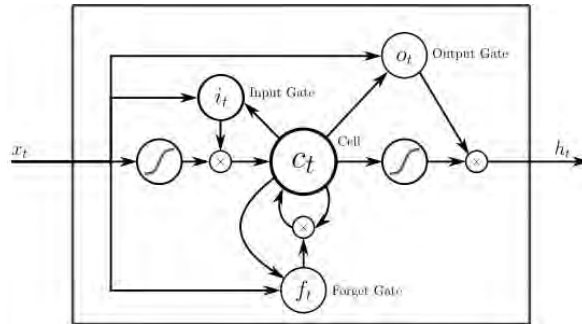
$$\mathbf{u}_1(k) = \mathbf{s}_1 [\mathbf{W}_1 \mathbf{x}(k) + \mathbf{U}_1 \mathbf{u}_2(k-1) + \mathbf{b}_1]$$

$$\mathbf{u}_2(k) = \mathbf{s}_2 [\mathbf{W}_2 \mathbf{u}_1(k) + \mathbf{b}_2]$$



# Long Short-Term Memory

- Memory held until new value overwrites
- Cell, input gate, output gate, forget gate



$$\mathbf{u}_f(k) = \mathbf{s}_g[\mathbf{W}_f \mathbf{x}(k) + \mathbf{U}_f \mathbf{u}_h(k-1) + \mathbf{b}_f]$$

$$\mathbf{u}_i(k) = \mathbf{s}_g[\mathbf{W}_i \mathbf{x}(k) + \mathbf{U}_i \mathbf{u}_h(k-1) + \mathbf{b}_i]$$

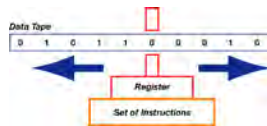
$$\mathbf{u}_o(k) = \mathbf{s}_g[\mathbf{W}_o \mathbf{x}(k) + \mathbf{U}_o \mathbf{u}_h(k-1) + \mathbf{b}_o]$$

$$\mathbf{u}_c(k) = \mathbf{u}_f(k) \cdot \mathbf{u}_c(k-1) + \mathbf{u}_i(k) \cdot \mathbf{s}_c[\mathbf{W}_c \mathbf{x}(k) + \mathbf{U}_c \mathbf{u}_h(k-1) + \mathbf{b}_c]$$

$$\mathbf{u}_h(k) = \mathbf{u}_o(k) \cdot \mathbf{s}_h[\mathbf{u}_c(k)]$$

Klaus, 2015

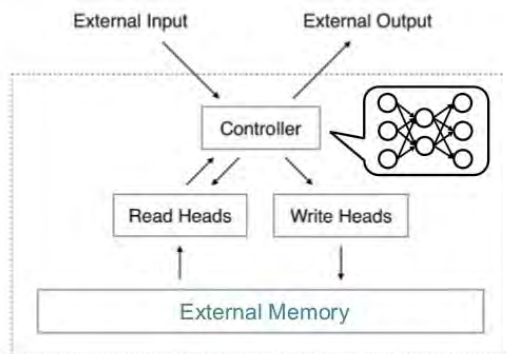
37



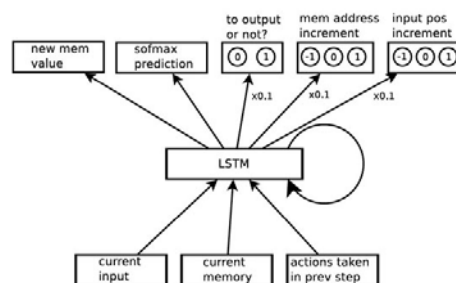
# Neural Turing Machines

*Sigmoid networks can be "Turing complete", Siegelmann & Sontag, 1991, 1995*

- Trainable read/write access to memory
- Controller/program implemented by neural networks
- Reinforcement-Learning NTM
- uses either feed-forward or LSTM neurons
- Improves on LSTM neurons



Graves, 2014



Zaremba, 2015

38

***Next Time:  
Communication, Information,  
and Machine Learning***

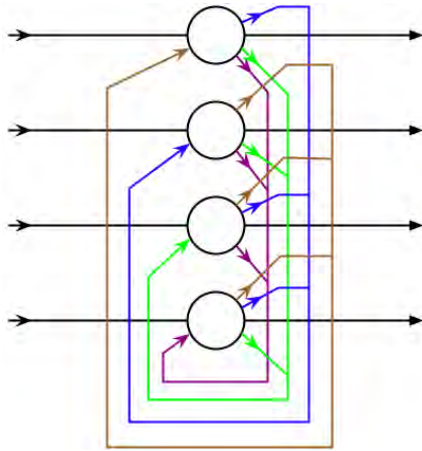
39

***SUPPLEMENTAL  
MATERIAL***

40

# Hopfield Network

Alternative plot of 4-node network



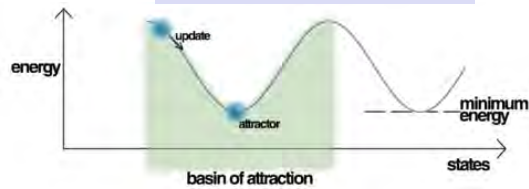
Novel Image



Exemplar



"Energy Landscape"



41



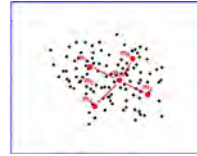
## Linear Vector Quantization

- Incorporation of supervised learning in Semantic Net
- Classification of groups of outputs
- Type 1
  - Addition of codebook vectors,  $m_c$ , with known meaning

$$\mathbf{m}_c(k+1) = \begin{cases} \mathbf{m}_c(k) + \alpha_k [\mathbf{x}_k - \mathbf{m}_c(k)], & \text{if classified correctly} \\ \mathbf{m}_c(k) - \alpha_k [\mathbf{x}_k - \mathbf{m}_c(k)], & \text{if classified incorrectly} \end{cases}$$

42

# Linear Vector Quantization



- **Type 2**
  - **Inhibition of nearest neighbor whose class is known to be different, e.g.,**
    - **$\mathbf{x}$  belongs to class of  $\mathbf{m}_j$  but is closer to  $\mathbf{m}_i$**

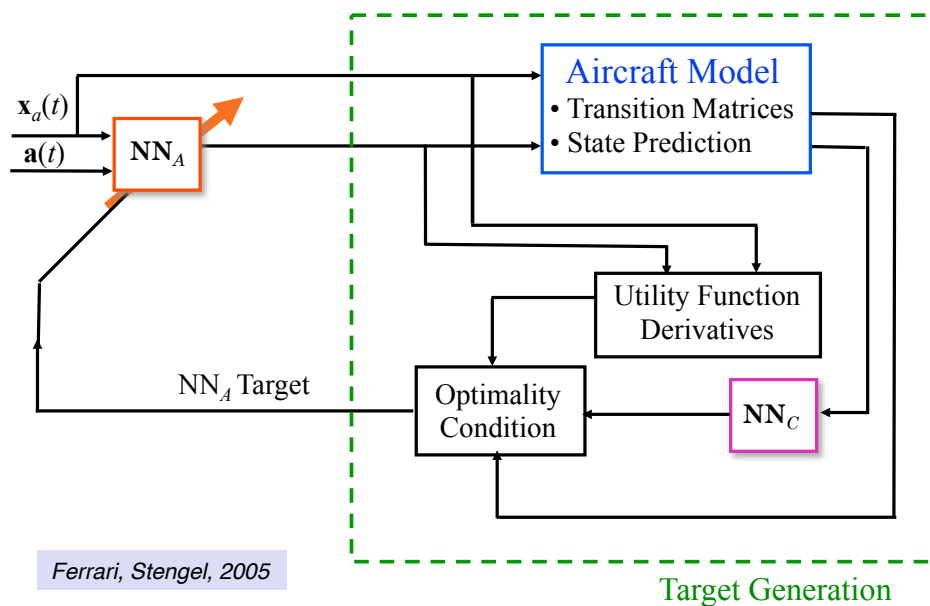
$$\mathbf{m}_i(k+1) = \mathbf{m}_i(k) - \alpha_k [\mathbf{x}_k - \mathbf{m}_i(k)]$$

$$\mathbf{m}_j(k+1) = \mathbf{m}_j(k) + \alpha_k [\mathbf{x}_k - \mathbf{m}_j(k)]$$

43

## Adaptive Critic Proportional-Integral Neural Network Controller

### Adaptation of Control Network



44

# Adaptive Critic Proportional-Integral Neural Network Controller

## Adaptation of Critic Network

