

1996-01

# Distributed Learning, Recognition, and Prediction by ART and ARTMAP Neural Networks

---

<https://hdl.handle.net/2144/2303>

*Downloaded from DSpace Repository, DSpace Institution's institutional repository*

**Distributed learning, recognition, and prediction  
by ART and ARTMAP neural networks**

**Gail Carpenter**

**January, 1996**

**Technical Report CAS/CNS-1996-004**

Permission to copy without fee all or part of this material is granted provided that: 1. The copies are not made or distributed for direct commercial advantage; 2. the report title, author, document number, and release date appear, and notice is given that copying is by permission of the BOSTON UNIVERSITY CENTER FOR ADAPTIVE SYSTEMS AND DEPARTMENT OF COGNITIVE AND NEURAL SYSTEMS. To copy otherwise, or to republish, requires a fee and / or special permission.

Copyright © 1996

Boston University Center for Adaptive Systems  
and  
Department of Cognitive and Neural Systems  
677 Beacon Street  
Boston, MA 02215

**Distributed learning, recognition, and prediction  
by ART and ARTMAP neural networks**

Gail A. Carpenter

Center for Adaptive Systems  
and  
Department of Cognitive and Neural Systems

677 Beacon Street  
Boston University  
Boston, Massachusetts 02215 USA

*Neural Networks*, 1997

Received: January, 1996

Revised: December, 1996

Running title: Distributed ART

Technical Report CAS/CNS TR-96-004  
Boston, MA: Boston University

This research was supported in part by the National Science Foundation (NSF IRI 94-01659) and the Office of Naval Research (ONR N00014-95-1-0409 and ONR N00014-95-0657).

Address for reprint requests: Prof. Gail A. Carpenter, Department of Cognitive and Neural Systems, 677 Beacon Street, Boston University, Boston, MA 02215 USA.

E-mail - [gail@cns.bu.edu](mailto:gail@cns.bu.edu).

**Abstract** -- A class of adaptive resonance theory (ART) models for learning, recognition, and prediction with arbitrarily distributed code representations is introduced. Distributed ART neural networks combine the stable fast learning capabilities of winner-take-all ART systems with the noise tolerance and code compression capabilities of multilayer perceptrons. With a winner-take-all code, the unsupervised model dART reduces to fuzzy ART and the supervised model dARTMAP reduces to fuzzy ARTMAP. With a distributed code, these networks automatically apportion learned changes according to the degree of activation of each coding node, which permits fast as well as slow learning without catastrophic forgetting. Distributed ART models replace the traditional neural network path weight with a dynamic weight equal to the rectified difference between coding node activation and an adaptive threshold. Thresholds increase monotonically during learning according to a principle of atrophy due to disuse. However, monotonic change at the synaptic level manifests itself as bidirectional change at the dynamic level, where the result of adaptation resembles long-term potentiation (LTP) for single-pulse or low-frequency test inputs but can resemble long-term depression (LTD) for higher frequency test inputs. This paradoxical behavior is traced to dual computational properties of phasic and tonic coding signal components. A parallel distributed match-reset-search process also helps stabilize memory. Without the match-reset-search system, dART becomes a type of distributed competitive learning network.

**Keywords** -- Distributed ART, Adaptive Resonance Theory, Distributed coding, Dynamic weight, Fast learning, Competitive learning, ARTMAP, Neural network

## 1. INTRODUCTION: ART, ARTMAP, AND DISTRIBUTED CODES

ART and ARTMAP neural networks feature winner-take-all competitive activation, which permits fast learning and stable coding but which causes category proliferation with noisy inputs. In contrast, multilayer perceptron models feature distributed McCulloch-Pitts activation, which enables good noise tolerance and code compression but which causes catastrophic forgetting with fast learning. This paper introduces a family of neural networks, called distributed ART models, that combine the best of these two worlds: distributed activation provides noise tolerance and code compression while new system dynamics retain stable fast learning capabilities, as follows.

### 1.1 ART and ARTMAP Networks

The theory of adaptive resonance began with an analysis of human cognitive information processing and stable coding in a complex input environment (Grossberg, 1976a, 1980). ART neural network models have added a series of new principles to the original theory and have realized these principles as quantitative systems that can be applied to problems of category learning, recognition, and prediction. Applications of unsupervised ART networks (Carpenter & Grossberg, 1987, 1991; Carpenter, Grossberg, & Rosen, 1991) and supervised ARTMAP networks (Carpenter, Grossberg, Markuzon, Reynolds, & Rosen, 1992; Carpenter, Grossberg, & Reynolds, 1991) include a Boeing parts design retrieval system (Caudell, Smith, Escobedo, & Anderson, 1994), satellite remote sensing (Baraldi & Parmiggiani, 1995; Gopal, Sklarew, & Lambin, 1994), robot sensory-motor control (Bachelder, Waxman, & Seibert, 1993; Baloch & Waxman, 1991; Dubrawski & Crowley, 1994; Srinivasa & Sharma, 1996), robot navigation (Racz & Dubrawski, 1995), machine vision (Caudell & Healy, 1994), 3D object recognition (Seibert & Waxman, 1992), face recognition (Seibert & Waxman, 1993), automatic target recognition (Bernardon and Carrick, 1995; Koch, Moya, Hostetler, & Fogler, 1995; Waxman et al., 1995), medical imaging (Soliz & Donohoe, 1996), electrocardiogram wave recognition (Ham & Han, 1996; Suzuki, 1995), prediction of protein secondary structure (Mehta, Vij, & Rabelo, 1993), strength prediction for concrete mixes (Kasperkiewicz, Racz, & Dubrawski, 1994), signature verification (Murshed, Bortozzi, & Sabourin, 1995), tool failure monitoring (Ly & Choi, 1994; Tarn, Li, & Chen, 1994; Tse & Wang, 1996), chemical analysis from UV and IR spectra (Wienke, 1994), digital circuit design (Kalkunte, Kumar, & Patnaik, 1992), frequency selective surface design for electromagnetic system devices (Christodoulou, Huang, Georgiopoulos, & Liou, 1995), Chinese character recognition (Gan & Lua, 1992; Kim, Jung, Kim, & Kim, 1995), and analysis of musical scores (Gjerdingen, 1990). Despite the growing number of applications, category proliferation from noisy training sets limits the useful domain of fast-learn, winner-take-all (WTA) systems such as ART or ARTMAP. On the other hand, fast learning is often essential for on-line adaptation to rapidly changing circumstances and for encoding rare cases and large databases.

Variants of the basic ART and ARTMAP networks have acquired some of the advantages of distributed coding while maintaining the fast learning capability. For example, ART-EMAP, which uses WTA codes for learning and distributed codes for testing, can significantly improve ARTMAP performance, especially when the size of the training set is small (Carpenter & Ross, 1993, 1995; Rubin, 1995). In medical database prediction problems, which often feature inconsistent training input predictions, ARTMAP-IC improves performance with a combination of distributed prediction, category instance counting, and a new match tracking search algorithm (Carpenter & Markuzon, 1996). A voting strategy further increases predictive accuracy by training the system several times on different orderings of an input set. Voting, instance counting, and distributed representations combine to form confidence estimates for

competing predictions. However, since these and most other ART and ARTMAP variants use WTA coding during learning, they do not solve the primary problem of category proliferation with noisy training sets, unless learning is slow.

The new family of distributed ART models retain stable coding, recognition, and prediction, but allow arbitrarily distributed code representation during learning as well as performance. When the code is winner-take-all, the unsupervised dART model is computationally equivalent to fuzzy ART and the supervised dARTMAP model is equivalent to fuzzy ARTMAP. Distributed ART networks automatically apportion learned changes according to the degree of activation of each coding node. This permits fast as well as slow learning without catastrophic forgetting. Many variations of the basic dART system may be devised but, for clarity, one specific network from the larger class is developed here.

## 1.2 Neural Analogues of dART Network Components

Distributed ART derives primarily from a computational analysis of design principles for constructing a learning system that is fast, stable, and distributed. Nevertheless, many network elements can also be visualized as physical processes with neural interpretations. In distributed ART, the fundamental synaptic memory unit is an adaptive threshold that increases during learning according to a principle of atrophy due to disuse. A dynamic weight that depends on both the coding node activation and the adaptive threshold then replaces the fuzzy ART path weight in the dART algorithm. In contrast, the fundamental synaptic memory unit in nearly all other neural networks is assumed axiomatically to be a multiplicative weight. This view of adaptation is also prevalent in the experimental literature: “Changes in the amplitude of synaptic responses evoked by single-shock extracellular electrical stimulation of presynaptic fibres are usually considered to reflect a change in the gain of synaptic signals, and are the most frequently used measure for evaluating synaptic plasticity.” (Markram & Tsodyks, 1996, p. 807) That is, when long-term potentiation (LTP), or enhanced postsynaptic response to a single test pulse, is observed, the strength, or gain, of a synapse is normally interpreted as having increased. Similarly, long-term depression (LTD) is usually thought of as a weight decrease.

While the multiplicative weight model helps explain classical LTP and LTD experiments, limitations of this hypothesis are beginning to become apparent. Describing their experiments on layer-5 pyramidal neurons in the neocortex, Markram and Tsodyks point out that the enhanced response to single-spike ( $\leq 0.25$  Hz) test probes in an LTP experiment vanishes with 23 Hz test stimuli: “Potentiation of synaptic responses therefore only occurred when the presynaptic frequency was below 20 Hz.” (p. 809) In fact, the Markram and Tsodyks data (Figure 3c, p. 809) actually show depressed postsynaptic responses to higher frequency (30 Hz and 40 Hz) test stimuli. They conclude: “The physiological implications of redistribution of efficacy are also entirely different from unconditional potentiation or depression.” (p. 810)

The dynamic coding behavior of distributed ART model neurons closely resembles this paradoxical “redistribution of efficacy.” In dART, adaptive thresholds increase monotonically during learning, but an increased threshold produces postsynaptic potentiation for lower frequency test inputs and postsynaptic depression for higher frequencies. These bidirectional dynamics are traced to the form of the signal that activates the dART distributed code. This signal is a function of two components with dual computational properties: a *phasic component* that depends on the transmitted input (ligand) and a *tonic component* that is independent of the current input. Both phasic and tonic components depend on the size of the adaptive threshold for the synapse and on the degree of activation of the target node (voltage). Phasic and tonic components can thus be visualized as postsynaptic membrane processes with phasic terms

mediated by voltage-and-ligand-gated receptors and tonic terms mediated by voltage-gated receptors (Nicholls, 1994). At each synapse, phasic and tonic terms dynamically balance one another. During adaptation, phasic terms remain constant while tonic terms may grow. Tonic components then become larger for all inputs, but phasic components become more selective. The net effect is to enhance the total coding signal subsequently sent by input components that are the same as or smaller than the one experienced during training (potentiation) but to reduce the total coding signal sent by input components that are substantially larger than those experienced during training (depression).

Analysis of the Markram and Tsodyks data illustrates how computational modeling of distributed pattern coding by neural network architectures is connected to important current questions concerning the underlying neural mechanisms of learning and memory. Phasic and tonic signals in the dART model, originally derived from a formal analysis of distributed pattern learning, demonstrate how: “Redistribution of synaptic efficacy may therefore serve as a powerful mechanism to alter the dynamics of synaptic transmission in subtle ways and hence to alter the content rather than the gain of signals conveyed between neurons.” (Markram & Tsodyks, 1996, p. 810) The remainder of this paper will henceforth focus primarily on the design of distributed ART.

### 1.3 Outline

Section 2 introduces the dART architecture and formally defines dynamic weights, adaptive thresholds, and phasic and tonic signal components; and characterizes the distributed code that a given input will activate. Section 3 describes a parallel distributed match-reset-search process. Section 4 outlines the distributed outstar used for top-down dART learning and introduces the distributed instar used for bottom-up learning. Dynamics of a distributed competitive learning module are also characterized. Section 5 summarizes a dART algorithm for simulation implementation. With winner-take-all dynamics at the coding field  $F_2$ , the dART algorithm reduces to a fuzzy ART algorithm, and further reduces to an ART 1 algorithm with binary inputs. Section 6 provides a geometric representation of distributed ART and Section 7 includes numerical examples of dART activation, search, and learning. Finally, Section 8 describes distributed ARTMAP.

## 2. DISTRIBUTED ACTIVATION

Over the past decade, an evolving series of neural network models have progressively expanded the domain and function of ART systems. The first model, ART 1 (Carpenter & Grossberg, 1987), is an unsupervised learning system that self-organizes recognition categories for binary input patterns. Fuzzy ART (Carpenter, Grossberg, & Rosen, 1991) generalizes binary ART 1 to the analog input domain, formally replacing set-theoretic intersections with fuzzy set-theoretic intersections (Figure 1a). These and most other ART models use choice, or winner-take-all (WTA), dynamics at the category representation field. Distributed ART (dART) continues the series, generalizing fuzzy ART to permit arbitrarily distributed code representations (Figure 1b). For continuity, dART retains fuzzy ART notation wherever possible.

**Figure 1** (p. 38): Fuzzy ART and distributed ART

### 2.1 dART Network Architecture

Although dART with winner-take-all coding is computationally equivalent to fuzzy ART, the dART architecture differs from the standard ART architecture. Namely, an ART input from a field  $F_0$  passes through a matching field  $F_1$  before activating a coding field  $F_2$ . Activity at  $F_2$  feeds back to  $F_1$ , forming a resonant loop (Figure 1a). ART networks thus encode matched  $F_1$

patterns, rather than the  $F_0$  inputs themselves, a key feature for code stability. With WTA coding, the matched  $F_1$  pattern confirms the original category choice when it feeds back up to  $F_2$ . The critical code confirmation property may not persist in this architecture, however, when  $F_2$  activation is distributed. In contrast, in the distributed ART network, the coding field  $F_2$  receives input directly from  $F_0$ , retaining the bottom-up / top-down matching process at  $F_1$  only to determine whether an active code meets the vigilance matching criterion (Figure 1b). Nevertheless, dART dynamic weights maintain code stability when  $F_2$  coding is winner-take-all. When the matching process is disabled by setting the vigilance parameter to 0, dART becomes a type of feedforward ART network that can also be viewed as a new type of distributed competitive learning architecture.

## 2.2. Activity Vectors

A dART system includes a field of nodes  $F_0$  that represents a current input vector; a field  $F_2$  that represents the active code; and a field  $F_1$  that represents a matched pattern determined by bottom-up input from  $F_0$  and top-down input from  $F_2$ . Vector  $\mathbf{I} \equiv (I_1 \dots I_i \dots I_M)$  denotes  $F_0$  activity,  $\mathbf{x} \equiv (x_1 \dots x_i \dots x_M)$  denotes  $F_1$  activity, and  $\mathbf{y} \equiv (y_1 \dots y_j \dots y_N)$  denotes  $F_2$  activity. Each component of  $\mathbf{I}$ ,  $\mathbf{x}$ , and  $\mathbf{y}$  is contained in the interval  $[0,1]$ . The number of input components ( $M$ ) and the number of coding nodes ( $N$ ) may be arbitrarily large. Although the matched  $F_1$  activity vector  $\mathbf{x}$  does not feed back to  $F_2$  (Figure 1), dART still performs computations that are equivalent to those of fuzzy ART in the special case of category choice at  $F_2$ . The input  $\mathbf{I}$  and the matched pattern  $\mathbf{x}$  may be continuously varying functions of time  $t$ , but the code  $\mathbf{y}$  acts as a content-addressable memory (CAM) that is held constant between resets by strong competition at  $F_2$ .

## 2.3. Dynamic Weights

In fuzzy ART the path from the  $i^{\text{th}}$   $F_1$  node to the  $j^{\text{th}}$   $F_2$  node contains an adaptive weight  $w_{ij}$ , and the path from the  $j^{\text{th}}$   $F_2$  node to the  $i^{\text{th}}$   $F_1$  node contains a weight  $w_{ji}$ . With fast learning,  $w_{ij} \equiv w_{ji}$ . Nearly all neural network models hypothesize such a weight as the unit of long-term memory (LTM). In contrast, in the distributed outstar network (Carpenter, 1994a) the unit of long-term memory is an adaptive threshold  $\tau_{ji}$ . Formally,

$$\tau_{ji} \equiv 1 - w_{ji}. \quad (1)$$

The distributed outstar signal from the  $j^{\text{th}}$   $F_2$  node to the  $i^{\text{th}}$   $F_1$  node is  $[y_j - \tau_{ji}]^+$ , where  $[\dots]^+$  denotes the rectification operator:

$$[\xi]^+ \equiv \max\{\xi, 0\}. \quad (2)$$

This path signal helps avoid catastrophic forgetting because  $[y_j - \tau_{ji}]^+ = [w_{ji} - (1 - y_j)]^+ = 0$  when  $w_{ji}$  is small, unless  $y_j = 1$ . Other types of signals such as the product  $y_j w_{ji}$  remain positive when  $y_j$  is positive, no matter how small the weight has



become, leaving  $w_{ji}$  subject to erosion. When the  $j^{\text{th}}$   $F_2$  node is chosen,  $w_{ji} = (1 - \tau_{ji}) = [y_j - \tau_{ji}]^+$ .

Distributed ART takes this idea one step further, replacing each fuzzy ART weight with a *dynamic weight* that is a joint function of coding node activation and an adaptive threshold. The formal substitution:

$$w_{ji} \rightarrow [y_j - \tau_{ji}]^+ \quad (3)$$

and

$$w_{ij} \rightarrow [y_j - \tau_{ij}]^+ \quad (4)$$

is the key step in converting fuzzy ART to distributed ART. Thresholds  $\tau_{ji}$  in paths from the  $j^{\text{th}}$   $F_2$  node to the  $i^{\text{th}}$   $F_1$  node adapt according to a distributed outstar learning law (Section 4.1), while thresholds  $\tau_{ij}$  in paths from the  $i^{\text{th}}$   $F_0$  node to the  $j^{\text{th}}$   $F_2$  node obey a distributed instar learning law (Section 4.2). Adaptive thresholds remain in the range  $[0,1]$ , starting at or near 0 and increasing monotonically during learning.

#### 2.4. Signal Functions

For each input  $\mathbf{I}$  and  $j = 1 \dots N$ , the total signal  $T_j$  from the dART input field  $F_0$  to the  $j^{\text{th}}$   $F_2$  node is a function of the form:

$$T_j = T_j(y_j) = g_j(S_j(y_j), \Theta_j(y_j)). \quad (5)$$

For  $S_j > 0$  and  $\Theta_j > 0$ ,

$$\frac{\partial g_j}{\partial S_j} > \frac{\partial g_j}{\partial \Theta_j} > 0, \quad (6)$$

and

$$g_j(0,0) = 0. \quad (7)$$

The definition of the  $F_0 \rightarrow F_2$  signal  $T_j$  at first appears to be circular:  $T_j$  determines the  $F_2$  code  $\mathbf{y}$  (Figure 1), but  $\mathbf{y}$  in turn determines  $T_j$  (5). However, this circularity does not actually occur in distributed ART dynamics. Because the competitive field  $F_2$  acts as a content-addressable memory, the network holds  $\mathbf{y}$  constant between resets (Section 2.2). Upon reset, a large nonspecific arousal signal breaks the CAM competitive feedback loop, momentarily sending all  $y_j$  values to 1 (Section 2.5). The code  $\mathbf{y}$  at any given time is therefore fully determined by the value of the signals  $T_j(1)$  at the time of the previous reset.  $T_j(y_j)$  represents the synaptic processes that, having survived the competition at reset, determine the

dynamics of search (Section 3.3) and learning (Section 4.2) between resets. Since total  $F_2$  activity is normalized to 1 (Section 2.5), active nodes typically represent a concentrated subset of the field's total capacity ( $N$ ), which can be arbitrarily large. Correspondingly, the signal sum  $T_j(y_j)$  between resets is on average a small fraction of the signal  $T_j(1)$  at the time of reset.

**Figure 2** (p. 39): Distributed instar signal components

In (5) the *phasic* component  $S_j$ , which depends on the input  $\mathbf{I}$ , is a sum:

$$S_j = S_j(y_j) = \sum_{i=1}^M S_{ij}(y_j). \quad (8)$$

A term  $S_j(y_j)$  in the sum (8) may be visualized as a certain fraction of the membrane sites at the  $i^{\text{th}}$  synapse of the  $j^{\text{th}}$   $F_2$  node (Figure 2). After a new input  $\mathbf{I}$  establishes an  $F_2$  code  $\mathbf{y}$  at the time of reset, phasic sites primed by the dynamic weight  $[y_j - \tau_{ij}]^+$  can remain activated by the input  $I_i$ , although a number of these sites ( $\Delta_{ij}$ ) may be refractory, or depleted, due to their recent activation during search (Section 3). Formally,

$$S_{ij}(y_j) = \left[ I_i \wedge [y_j - \tau_{ij}]^+ - \Delta_{ij} \right]^+, \quad (9)$$

where  $\wedge$  represents the fuzzy intersection, or component-wise minimum:

$$(\mathbf{a} \wedge \mathbf{b})_i \equiv (a_i \wedge b_i) \equiv \min(a_i, b_i) \quad (10)$$

(Zadeh, 1965). For  $y_j \in [0, 1]$ ,

$$0 \leq S_j(y_j) \leq \sum_{i=1}^M [y_j - \tau_{ij}]^+ \leq \sum_{i=1}^M y_j = My_j. \quad (11)$$

In (5), the *tonic* component  $\Theta_j$  is a sum:

$$\Theta_j = \Theta_j(y_j) = \sum_{i=1}^M \Theta_{ij}(y_j) \quad (12)$$

where:

$$\Theta_{ij}(y_j) = [y_j \wedge \tau_{ij} - \delta_{ij}]^+. \quad (13)$$

The sum  $\Theta_j(y_j)$ , which is independent of the input  $\mathbf{I}$ , plays the role of a nodal bias term that increases during learning. Once  $\mathbf{y}$  is established following a reset, a fraction of membrane sites

$\tau_{ij}$  are primed by the node's activity ( $y_j$ ), but recently active sites ( $\delta_{ij}$ ) may be refractory during search. Like  $S_j(y_j)$ ,  $\Theta_j(y_j)$  lies in the interval  $[0, My_j]$  since:

$$0 \leq \Theta_j(y_j) \leq \sum_{i=1}^M y_j = My_j. \quad (14)$$

Refractory sites accumulate during a rapid series of resets. On the time scale of learning, the terms  $\Delta_{ij}$  and  $\delta_{ij}$  decay to 0.

By design, the phasic and tonic components of the  $F_2$  input signal  $T_j$  play complementary roles in dART networks. Each phasic term is an increasing function of  $I_i$ . However, when  $\mathbf{I}$  and  $\mathbf{y}$  remain constant during a learning interval, the phasic terms  $S_{ij}(y_j) = I_i \wedge [y_j - \tau_{ij}]^+$  and  $S_{ij}(1) = I_i \wedge (1 - \tau_{ij})$  remain constant (Section 4.2). In contrast, each tonic term is, by definition, independent of  $\mathbf{I}$ . However, the tonic terms  $\Theta_{ij}(y_j) = y_j \wedge \tau_{ij}$  and  $\Theta_{ij}(1) = \tau_{ij}$  increase during learning, when  $y_j$  is large enough. Thus by (5)-(6),  $T_j$  is an increasing function of each component of  $\mathbf{I}$  and  $T_j$  increases during learning.

A distributed version of the fuzzy ART choice-by-difference (CBD) function (Carpenter & Gjaja, 1994) defines one signal rule for  $T_j$  by:

$$T_j = S_j + (1 - \alpha)\Theta_j, \quad (15)$$

with  $0 < \alpha < 1$ . Like  $S_j$  and  $\Theta_j$ , the CBD signal function  $T_j \in [0, My_j]$  since:

$$\begin{aligned} 0 \leq T_j(y_j) &= S_j(y_j) + (1 - \alpha)\Theta_j(y_j) \leq S_j(y_j) + \Theta_j(y_j) \\ &\leq \sum_{i=1}^M I_i \wedge [y_j - \tau_{ij}]^+ + \sum_{i=1}^M y_j \wedge \tau_{ij} \\ &\leq \sum_{i=1}^M \left( [y_j - \tau_{ij}]^+ + y_j \wedge \tau_{ij} \right) \\ &= \sum_{i=1}^M \left( (y_j - y_j \wedge \tau_{ij}) + y_j \wedge \tau_{ij} \right) = \sum_{i=1}^M y_j = My_j. \end{aligned} \quad (16)$$

A distributed version of the Weber law signal function (Carpenter & Grossberg, 1987) defines a different signal rule for  $T_j$  by:

$$T_j = \frac{S_j}{\alpha + My_j - \Theta_j}, \quad (17)$$

with  $\alpha > 0$ . For the Weber law coding function (17),  $T_j \in [0, 1)$  since:

$$\begin{aligned}
0 \leq T_j &= \frac{S_j}{\alpha + My_j - \Theta_j} \\
&\leq \frac{\sum_{i=1}^M I_i \wedge [y_j - \tau_{ij}]^+}{\alpha + My_j - \sum_{i=1}^M y_j \wedge \tau_{ij}} \leq \frac{\sum_{i=1}^M [y_j - \tau_{ij}]^+}{\alpha + \sum_{i=1}^M [y_j - \tau_{ij}]^+} \leq \frac{My_j}{\alpha + My_j} < 1.
\end{aligned} \tag{18}$$

In the case where  $y_j \equiv 1$ ,  $\Delta_{ij} = \delta_{ij} = 0$ , and  $w_{ij} \equiv (1 - \tau_{ij})$ :

$$S_j = S_j(1) = |\mathbf{I} \wedge \mathbf{w}_j| \tag{19}$$

and

$$\Theta_j = \Theta_j(1) = (M - |\mathbf{w}_j|), \tag{20}$$

where  $|\dots|$  represents the city-block norm. In this case the distributed choice-by-difference function (15) reduces to:

$$T_j = T_j(1) = |\mathbf{I} \wedge \mathbf{w}_j| + (1 - \alpha)(M - |\mathbf{w}_j|), \tag{21}$$

which is equivalent to the fuzzy ART choice-by-difference function. The distributed Weber law function (17) reduces to:

$$T_j = T_j(1) = \frac{|\mathbf{I} \wedge \mathbf{w}_j|}{\alpha + |\mathbf{w}_j|}, \tag{22}$$

which is equivalent to the Weber law choice rules originally used in fuzzy ART and, when  $\mathbf{I}$  is binary, ART 1.

## 2.5. Code Representation

In distributed ART networks, activity  $\mathbf{y} = (y_1 \dots y_j \dots y_N)$  at a competitive coding field  $F_2$  is stored as a content-addressable memory. An algorithm that approximates the dynamics of strong competition postulates that external inputs initially determine  $\mathbf{y}$ , but then internal feedback holds  $\mathbf{y}$  constant until  $F_2$  is actively reset. Except during reset,  $\mathbf{y}$  is normalized:

$$|\mathbf{y}| \equiv \sum_{j=1}^N y_j = 1. \tag{23}$$

In ART models,  $F_2$  reset occurs when the bottom-up / top-down matched pattern  $\mathbf{x}$  at  $F_1$  fails to meet a matching criterion defined by a vigilance parameter  $\rho$ . Reset is effected by a large nonspecific arousal signal. In the dART model, reset momentarily sends all  $y_j$  to 1 at a

time  $t=r$ . This allows the values  $T_1(1)|_{t=r} \dots T_N(1)|_{t=r}$  to determine which  $y$  will be established next. Until the next reset,

$$y_j = f_j(T_1(1) \dots T_N(1))|_{t=r}. \quad (24)$$

Realizing  $F_2$  as an on-center off-surround shunting competitive network suggests the hypothesis:

$$\frac{\partial f_j}{\partial T_j} \geq 0. \quad (25)$$

One class of functions that satisfy this hypothesis sets:

$$y_j = \begin{cases} \frac{f(T_j(1))}{\sum_{\lambda \in \Lambda} f(T_\lambda(1))} & \text{if } j \in \Lambda \\ 0 & \text{if } j \notin \Lambda \end{cases} \quad (26)$$

where  $\Lambda$  is a subset of  $\{1 \dots N\}$  such that  $T_J \geq T_j$  for  $J \in \Lambda$  and  $j \notin \Lambda$ ; and where  $f(0) \geq 0$  and  $f'(\xi) \geq 0$  for  $\xi > 0$ . Grossberg (1976b) used a similar class of functions to approximate the dynamics of on-center off-surround shunting competitive networks. The index subset  $\Lambda$  might be the indices of  $T_j$  values that are greater than or equal to the collective average (above-average- $T_j$  rule); or  $\Lambda$  might be the indices of the  $Q$  largest  $T_j$  values ( $Q$ -max rule). Setting  $Q=1$  corresponds to choice, or winner-take-all, coding and setting  $Q=N$  makes all  $y_j$  proportional to  $f(T_j(1))$ . The function  $f$  might realize a power law, with:

$$f(\xi) = \xi^p \quad (27)$$

for  $p > 0$ . Setting  $p=1$  makes  $y_j$  proportional to  $T_j(1)$  for  $j \in \Lambda$ , and increasing the power  $p$  models progressively stronger internal network competition, producing increasingly compressed  $F_2$  codes. In the limit as  $p \rightarrow \infty$ , the system (26)-(27) converges to the choice rule. Other types of coding fields could, for example, represent cooperative or spatially defined interactions as well as competition. Compared to ART and ARTMAP networks, where the coding rule is fixed, applications of dART and dARTMAP networks typically require comparative studies to help choose rules that give the best performance in particular cases.

### 3. DISTRIBUTED SEARCH

The distributed ART match-reset-search process is similar to that of other ART networks. When an  $F_2$  code  $\mathbf{y}$  becomes active, the activity pattern  $\mathbf{x}$  at  $F_1$  represents a match between the current bottom-up input  $\mathbf{I}$  and a top-down input  $\sigma(\mathbf{y})$ . If these inputs fail to meet the vigilance matching criterion, a nonspecific reset signal shuts off the code  $\mathbf{y}$ . Reset also leaves an enduring trace of  $\mathbf{y}$ , or the network would simply reactivate the same code.

The search process plays a variety of roles in ART and ARTMAP systems. Since  $F_2$  is typically a strongly competitive network, active reset of a stored code is needed for each new

input to select a code that is not severely distorted by the previous steady state at  $F_2$ . An *input reset* allows an input to register its own code when it fails to match an active top-down signal  $\sigma(\mathbf{y})$ . Alternatively, a novelty signal can automatically trigger a reset when a new input is presented. Input resets segment a continuously varying input  $\mathbf{I}(t)$  with a discrete series of recognition codes  $\mathbf{y}^{(1)}, \mathbf{y}^{(2)}, \dots$ . While one code remains active, the subset of input features active at  $F_1$  represents a focus of attention. Reset defines the boundary between one attended feature set and the next.

Search also helps to stabilize memory. Immediately after an input activates a code, a *mismatch reset* will quickly shut off  $\mathbf{y}$  if it fails to meet the vigilance matching criterion. Since reset is rapid on the time scale of learning (LTM), an outlier that incorrectly activates a learned code does not disrupt memory. Traces of prior resets should endure on the time scale of short-term memory (STM) and search but should fade on the time scale of learning, since a reset code that was incorrect for one input may be correct for the next. Traces of search are thus a type of medium-term memory (MTM).

Even if  $\mathbf{I}$  and  $\mathbf{y}$  are constant and  $\mathbf{x}$  meets the matching criterion, an increase in the vigilance parameter  $\rho$  can trigger search. Such a *vigilance reset* corresponds to increased “attentiveness” due, for example, to a prediction made by  $\mathbf{y}$  having led to an error. In fact, when an ARTMAP network makes a predictive error during training, the match tracking process raises vigilance until the matching criterion fails, thus triggering a vigilance reset and search. In ARTMAP the vigilance parameter therefore becomes an internally controlled variable that may increase on the MTM time scale but that relaxes to a baseline vigilance level ( $\bar{\rho}$ ) on the LTM time scale. Finally, reset waves might also refresh  $F_2$  periodically, to keep the system from locking into a fixed state even if vigilance is low.

### 3.1. Match Representation

While  $\mathbf{y}$  is fixed between resets, the total input  $\sigma_i$  from  $F_2$  to the  $i^{\text{th}}$   $F_1$  node equals the sum of dynamic weights projecting to that node. That is:

$$\sigma_i = \sigma_i(\mathbf{y}) = \sum_{j=1}^N [y_j - \tau_{ji}]^+, \quad (28)$$

where  $\tau_{ji} \in [0, 1]$  is an adaptive threshold that starts at 0 and may increase during distributed outstar learning (Section 4.1). Since  $\sum_j y_j = 1$ ,  $\sigma_i \in [0, 1]$ . Activity  $\mathbf{x}$  at  $F_1$  then equals the fuzzy intersection of  $\mathbf{I}$  and  $\sigma(\mathbf{y})$ , so:

$$x_i = I_i \wedge \sigma_i(\mathbf{y}) \quad (29)$$

for  $i = 1 \dots M$ . Signals from  $F_2$  thereby prime  $F_1$  in the sense that  $\sigma_i(\mathbf{y})$  imposes an upper bound on inputs  $I_i$  that can be fully represented at the  $i^{\text{th}}$   $F_1$  node.

### 3.2. Resonance or Reset

*Resonance* occurs if the matched pattern  $\mathbf{I} \wedge \sigma(\mathbf{y})$  meets the vigilance criterion:

$$\frac{|\mathbf{x}|}{|\mathbf{I}|} = \frac{|\mathbf{I} \wedge \sigma(\mathbf{y})|}{|\mathbf{I}|} \geq \rho, \quad (30)$$

that is, resonance occurs if

$$|\mathbf{x}| = |\mathbf{I} \wedge \sigma(\mathbf{y})| \geq \rho |\mathbf{I}|. \quad (31)$$

Learning then ensues, as defined below. During a learning interval,  $\mathbf{y}$  remains constant but the input  $\mathbf{I}(t)$  and the vigilance parameter  $\rho$  may vary continuously, as long as the network continues to meet the matching criterion.

*Mismatch reset* occurs if:

$$\frac{|\mathbf{I} \wedge \sigma(\mathbf{y})|}{|\mathbf{I}|} < \rho, \quad (32)$$

that is, if:

$$|\mathbf{x}| = |\mathbf{I} \wedge \sigma(\mathbf{y})| < \rho |\mathbf{I}|. \quad (33)$$

A nonspecific signal to  $F_2$  then momentarily resets all  $y_j$  to 1, until the signal vector  $\mathbf{T}$  establishes a new code  $\mathbf{y}$  (Section 2.5). The search process must be rapid, so that no significant learning can occur with an incorrect code. Mismatch reset must also selectively bias the network against previously active codes or  $\mathbf{T}$ , the same as before, will reactivate the reset code.

### 3.3 Medium-Term Memory

When the  $F_2$  code makes a choice, reset needs simply to deactivate the previously active node  $J$  for the duration of the MTM time scale. When  $\mathbf{y}$  is distributed, a graded bias against the  $j^{\text{th}}$  node needs to reflect how large  $y_j$  has been in previously reset codes, so that highly active nodes can give way to nodes that originally received smaller inputs. Figure 3 shows how such a parallel search process can explore various  $F_2$  code combinations until one is found that satisfies the vigilance criterion. During a rapid series of mismatch-reset events, refractory sites accumulate (Figure 3a-c). During a learning interval, refractory sites recover (Figure 3d).

**Figure 3** (p. 40): Parallel distributed search

Distributed ART realizes the search process by assuming that, when a code  $\mathbf{y}$  is active, sites corresponding to the phasic component  $S_j(y_j)$  (8) and the tonic component  $\Theta_j(y_j)$  (12) become refractory on the MTM time scale. On the time scale of search,

$$\frac{d}{dt} \Delta_{ij} \approx S_{ij}(y_j) = \left[ I_i \wedge [y_j - \tau_{ij}]^+ - \Delta_{ij} \right]^+ \quad (34)$$

and:

$$\frac{d}{dt} \delta_{ij} \approx \Theta_{ij}(y_j) = [y_j \wedge \tau_{ij} - \delta_{ij}]^+. \quad (35)$$

Each term  $S_{ij}(y_j)$  (9) and  $\Theta_{ij}(y_j)$  (13) then quickly converges to 0. When the next reset occurs,  $S_j(1)$  and  $\Theta_j(1)$  are reduced by the previous quantities  $S_j(y_j)$  and  $\Theta_j(y_j)$  (Figure 2). By (6),  $T_j(1)$  is also reduced; with distributed choice-by-difference,  $T_j(1)$  is reduced by the previous quantity  $T_j(y_j)$ . Nodes where  $y_j$  is large tend to have the largest signals and hence the greatest reduction of the subsequent signals after a reset. When  $y_j = 0$ ,  $S_j = \Theta_j = T_j = 0$ , so the signal  $T_j(1)$  at the next reset will be the same as it was before.

Since recovery is slow on the time scale of search, across a rapid series of resets the phasic depletion term  $\Delta_{ij}$  (9) is approximately equal to the largest value  $I_i \wedge [y_j - \tau_{ij}]^+$  has recently attained. The phasic term  $S_{ij}(y_j)$  can then be positive for a new code  $\mathbf{y}$  only if  $y_j$  is larger than it has yet been during the search. Similarly, the tonic depletion term  $\delta_{ij}$  (13) is approximately equal to the largest value  $y_j \wedge \tau_{ij}$  has recently attained. Refractory sites recover on the time scale of learning. For a search where code selection is unbiased by the previous choice, the model assumes that  $\Delta_{ij}$  and  $\delta_{ij}$  converge to 0 during learning. When  $F_2$  makes a choice, with  $y_J = 1$ ,  $S_J(1)$  and  $\Theta_J(1)$  are reduced to 0 as  $\Delta_{iJ} \rightarrow I_i \wedge (1 - \tau_{iJ})$  and  $\delta_{iJ} \rightarrow \tau_{iJ}$  during search. Since  $g_J(0,0) = 0$  (7),  $T_J(1)$  is then also reduced to 0 until it can recover on the time scale of learning.

#### 4. DISTRIBUTED LEARNING

Catastrophic forgetting is a problem faced by all neural networks with distributed activation especially in the fast-learn limit where LTM variables go to asymptote with each input presentation. The instar and outstar learning laws used in previous ART networks would cause catastrophic forgetting if transferred to a network with a distributed code  $\mathbf{y}$ . Stable distributed coding with fast learning requires internal or external control of the learned changes that one input can induce.

The distributed outstar (Carpenter, 1994a) solves the catastrophic forgetting problem for learning in paths that originate from the coding field  $F_2$ . The distributed instar, introduced here, solves the problem in paths that project to  $F_2$ . During distributed outstar learning, the total signal from the coding field to a target node can only decrease, by a principle of atrophy due to disuse. During distributed instar learning, the total signal to a target coding node can only increase, as the tonic component of the signal increases while the phasic component remains constant for a given input. Both learning laws bound the total learned change any one input can impose upon the system.

##### 4.1. Distributed Outstar Learning

Dynamic weights in paths that originate from an  $F_2$  coding node adapt according to a principle of *atrophy due to disuse*. The total top-down priming signal  $\sigma_i(\mathbf{y})$  to the  $i^{\text{th}}$   $F_1$  node equals the sum of dynamic weights projecting to that node (28). During distributed outstar learning, each signal  $\sigma_i(\mathbf{y})$  that exceeds the input  $I_i$  shrinks until it just “covers”  $I_i$ . Each dynamic weight  $[y_j - \tau_{ji}]^+$  falls by an amount that depends upon its contribution to  $\sigma_i(\mathbf{y})$  as the threshold  $\tau_{ji}$  rises according to the equation:



$$\begin{aligned}
\frac{d}{dt} \tau_{ji} &= [y_j - \tau_{ji}]^+ (\sigma_i(\mathbf{y}) - x_i) \\
&= [y_j - \tau_{ji}]^+ (\sigma_i(\mathbf{y}) - I_i \wedge \sigma_i(\mathbf{y})) \\
&= [y_j - \tau_{ji}]^+ [\sigma_i(\mathbf{y}) - I_i]^+,
\end{aligned} \tag{36}$$

by (29). Initially,  $\tau_{ji}(0) = 0$ . By (36), the sum of all thresholds to the  $i^{\text{th}}$   $F_1$  node increases according to:

$$\begin{aligned}
\frac{d}{dt} \sum_{j=1}^N \tau_{ji} &= \sum_{j=1}^N [y_j - \tau_{ji}]^+ (\sigma_i(\mathbf{y}) - x_i) = \sigma_i(\mathbf{y}) (\sigma_i(\mathbf{y}) - x_i) \\
&= \sigma_i(\mathbf{y}) (\sigma_i(\mathbf{y}) - I_i \wedge \sigma_i(\mathbf{y})) = \sigma_i(\mathbf{y}) [\sigma_i(\mathbf{y}) - I_i]^+.
\end{aligned} \tag{37}$$

As long as  $I_i$  remains constant,

$$\sigma_i(\mathbf{y}) \downarrow I_i \wedge \sigma_i(\mathbf{y}) \Big|_{t=r}, \tag{38}$$

where  $t = r$  at the time of the previous reset. That is, either  $\sigma_i(\mathbf{y})$  decreases toward  $I_i$  by atrophy due to disuse; or  $\sigma_i(\mathbf{y})$  is smaller than  $I_i$  to begin with and so remains constant until the next reset. Activity  $\mathbf{x} = \mathbf{I} \wedge \sigma(\mathbf{y})$  at the matching field  $F_1$  thus remains constant during learning, as long as  $\mathbf{I}$  and  $\mathbf{y}$  remain constant.

The distributed outstar equation is simple enough to be solved directly, and its solutions are piecewise linear. If  $\mathbf{I}$  and  $\mathbf{y}$  remain constant during a time interval  $[r, t]$ , then:

$$\tau_{ji}(t) = \tau_{ji}^{old} + \phi(t) \frac{[\sigma_i^{old} - I_i]^+}{\sigma_i^{old}} [y_j - \tau_{ji}^{old}]^+ \tag{39}$$

where  $\tau_{ji}^{old} \equiv \tau_{ji}(r)$  and  $\sigma_i^{old} \equiv \sigma_i(\mathbf{y}) \Big|_{t=r}$ , and where  $\phi(t)$  is an exponential that goes from 0 to 1 as  $t$  goes from  $r$  to  $\infty$  (Carpenter, 1994b). For input presentations of fixed duration, the  $F_2 \rightarrow F_1$  threshold  $\tau_{ji}$  increases during learning from  $\tau_{ji}^{old}$  to  $\tau_{ji}^{new}$ , where:

$$\tau_{ji}^{new} = \tau_{ji}^{old} + \beta \frac{[\sigma_i^{old} - I_i]^+}{\sigma_i^{old}} [y_j - \tau_{ji}^{old}]^+ \tag{40}$$

for a learning rate parameter  $\beta \in [0, 1]$ . Setting  $\beta = 1$  gives the fast-learn limit, where all variables reach asymptote during each input presentation. Equation (39) also provides a formula for the dynamic weight  $[y_j - \tau_{ji}(t)]^+$ , which decreases during learning so that:

$$[y_j - \tau_{ji}(t)]^+ \downarrow [y_j - \tau_{ji}^{old}]^+ \frac{I_i \wedge \sigma_i^{old}}{\sigma_i^{old}}. \tag{41}$$

With category choice at  $F_2$  and fast learning, the distributed outstar reduces to the fuzzy ART outstar, as follows. In the original outstar (Grossberg, 1968, 1970), weights  $w_{ji}$  in paths from a source node with activity  $y_j$  track target node activities  $x_i$ . The specific outstar equation used in fuzzy ART is:

$$\frac{d}{dt}w_{ji} = y_j(x_i - w_{ji}). \quad (42)$$

In fuzzy ART, with  $y_J = 1$  at the chosen  $F_2$  node,  $x_i = I_i \wedge w_{Ji}$  at the  $i^{\text{th}}$   $F_1$  node. Initially, all  $w_{ji}(0) = 1$ , and top-down weights  $w_{ji}$  remain constant for  $j \neq J$ . For  $j = J$ , weights decrease by outstar learning:

$$\begin{aligned} \frac{d}{dt}w_{Ji} &= (x_i - w_{Ji}) = (I_i \wedge w_{Ji} - w_{Ji}) \\ &= -(w_{Ji} - I_i \wedge w_{Ji}) = -[w_{Ji} - I_i]^+. \end{aligned} \quad (43)$$

Correspondingly, in the distributed outstar (36) with the code  $\mathbf{y}$  representing choice at  $F_2$ ,  $\frac{d}{dt}\tau_{ji} = 0$  for  $j \neq J$ . For  $j = J$ ,

$$\begin{aligned} \frac{d}{dt}\tau_{Ji} &= [y_J - \tau_{Ji}]^+ [\sigma_i(\mathbf{y}) - I_i]^+ \\ &= [y_J - \tau_{Ji}]^+ \left[ [y_J - \tau_{Ji}]^+ - I_i \right]^+ = (1 - \tau_{Ji}) [(1 - \tau_{Ji}) - I_i]^+. \end{aligned} \quad (44)$$

Setting  $w_{ji} \equiv (1 - \tau_{ji})$  (1) converts (44) into:

$$\frac{d}{dt}w_{Ji} = -w_{Ji} [w_{Ji} - I_i]^+, \quad (45)$$

with  $\frac{d}{dt}w_{ji} = 0$  for  $j \neq J$ . Thus, except for the convergence rate, the distributed outstar (36) with choice at  $F_2$  reduces to the fuzzy ART outstar (42). With fast learning, the two algorithms are equivalent.

In the fuzzy ART outstar, for fixed  $\mathbf{I}$  and a chosen node  $J$ , the total change in the set of weights from  $F_2$  to the  $i^{\text{th}}$   $F_1$  node is bounded above by  $1 - I_i$ :

$$\sum_{j=1}^N |\Delta w_{ji}| = [w_{Ji}^{\text{old}} - I_i]^+ \leq 1 - I_i, \quad (46)$$

where  $w_{Ji}^{\text{old}} \equiv w_{Ji}(r)$ . In the distributed outstar, the same bound applies, with  $\mathbf{y}$  arbitrarily distributed across  $F_2$ :

$$\begin{aligned}
\sum_{j=1}^N |\Delta \tau_{ji}(t)| &= \sum_{j=1}^N \phi(t) \frac{[\sigma_i^{old} - I_i]^+}{\sigma_i^{old}} [y_j - \tau_{ji}^{old}]^+ \\
&\leq \frac{[\sigma_i^{old} - I_i]^+}{\sigma_i^{old}} \sum_{j=1}^N [y_j - \tau_{ji}^{old}]^+ = [\sigma_i^{old} - I_i]^+ \leq 1 - I_i.
\end{aligned} \tag{47}$$

Thus distributed outstar learning preserves dynamic range and avoids catastrophic forgetting.

## 4.2. Distributed Instar Learning

Distributed instar learning is designed to enhance the competitive advantage of highly active coding nodes with respect to the current input. At the same time, learning makes these nodes more selective, so that different inputs will tend to activate distinct codes. During distributed instar learning a large dynamic weight  $[y_j - \tau_{ij}]^+$  decreases toward a smaller input  $I_i$  according to the equation:

$$\begin{aligned}
\frac{d}{dt} \tau_{ij} &= [y_j - \tau_{ij} - I_i]^+ \\
&= \left[ [y_j - \tau_{ij}]^+ - I_i \right]^+ \\
&= \left( [y_j - \tau_{ij}]^+ - I_i \wedge [y_j - \tau_{ij}^{old}]^+ \right),
\end{aligned} \tag{48}$$

where  $\tau_{ij}^{old} \equiv \tau_{ij}(r)$ , at the time of the previous reset. Initially:

$$\tau_{ij}(0) = \eta_{ij} = 0^+, \tag{49}$$

where the values  $\eta_{ij}$  are small random numbers needed to break the tie among the first  $F_2$  inputs  $T_j(1)$ . As long as  $\mathbf{I}$  and  $\mathbf{y}$  remain constant, the threshold  $\tau_{ij}$  increases:

$$\tau_{ij} \uparrow (y_j - I_i) \vee \tau_{ij}^{old}, \tag{50}$$

where  $\vee$  represents the fuzzy union, or component-wise maximum:

$$(\mathbf{a} \vee \mathbf{b})_i \equiv (a_i \vee b_i) \equiv \max(a_i, b_i) \tag{51}$$

(Zadeh, 1965). As  $\tau_{ij}$  increases, the dynamic weight  $[y_j - \tau_{ij}]^+$  decreases:

$$[y_j - \tau_{ij}]^+ \downarrow I_i \wedge [y_j - \tau_{ij}^{old}]^+. \tag{52}$$

Solving (48) gives:

$$\begin{aligned}\tau_{ij}(t) &= \tau_{ij}^{old} + \phi(t) \left[ y_j - \tau_{ij}^{old} - I_i \right]^+ \\ &= \begin{cases} \tau_{ij}^{old} & \text{if } \tau_{ij}^{old} \geq (y_j - I_i) \\ (1 - \phi(t))\tau_{ij}^{old} + \phi(t)(y_j - I_i) & \text{if } \tau_{ij}^{old} < (y_j - I_i) \end{cases}\end{aligned}\quad (53)$$

where  $\phi(t)$  is an exponential that goes from 0 to 1 as  $t$  goes from  $r$  to  $\infty$ . In addition, the maximum total increase across all the  $MN$  thresholds during one input presentation is bounded above by  $M$ :

$$\begin{aligned}\sum_{j=1}^N |\Delta\tau_{ij}| &= \sum_{i=1}^M \sum_{j=1}^N \left[ y_j - \tau_{ij}^{old} - I_i \right]^+ \\ &\leq \sum_{i=1}^M \sum_{j=1}^N \left[ y_j - I_i \right]^+ \leq \sum_{i=1}^M \sum_{j=1}^N y_j = M.\end{aligned}\quad (54)$$

For input presentations of fixed duration,  $\tau_{ij}$  increases from  $\tau_{ij}^{old}$  to  $\tau_{ij}^{new}$  where:

$$\begin{aligned}\tau_{ij}^{new} &= (1 - \beta)\tau_{ij}^{old} + \beta(y_j - I_i) \vee \tau_{ij}^{old} \\ &= \tau_{ij}^{old} + \beta \left[ y_j - \tau_{ij}^{old} - I_i \right]^+\end{aligned}\quad (55)$$

for the learning rate parameter  $\beta \in [0, 1]$ . In the fast-learn limit,  $\beta = 1$ .

Note that, by (9) and (13),  $S_{ij}(y_j) = I_i \wedge [y_j - \tau_{ij}]^+$  and  $\Theta_{ij}(y_j) = y_j \wedge \tau_{ij}$  during learning, since then  $\Delta_{ij} = \delta_{ij} = 0$ . The distributed instar learning law can thus be written in terms of the phasic and tonic signals, since:

$$\begin{aligned}\frac{d}{dt} \tau_{ij} &= \left[ [y_j - \tau_{ij}]^+ - I_i \right]^+ \\ &= [y_j - \tau_{ij}]^+ - I_i \wedge [y_j - \tau_{ij}]^+ \\ &= y_j - y_j \wedge \tau_{ij} - I_i \wedge [y_j - \tau_{ij}]^+ \\ &= y_j - \Theta_{ij}(y_j) - S_{ij}(y_j).\end{aligned}\quad (56)$$

The term  $S_{ij}(y_j)$  can be thought to represent a set of synaptic sites that are phasically activated by the input  $I_i$ , while  $\Theta_{ij}(y_j)$  represents sites that are tonically activated, independent of  $I_i$ . By hypothesis (6), a phasically active site makes a larger contribution to the overall signal than does a tonically active site. However, the term  $[y_j - S_{ij}(y_j)]$  then represents “disused” sites that are primed by postsynaptic activation  $y_j$  but are not phasically activated by the current input. During learning,  $\tau_{ij}$  remains constant if the  $j^{\text{th}}$  node is relatively inactive

$(y_j \leq \tau_{ij} + I_i)$ . Otherwise,  $\Theta_{ij}(y_j)$  increases toward  $[y_j - S_{ij}(y_j)]$  as disused phasic sites revert to tonic sites.

With category choice at  $F_2$ , the distributed instar reduces to the fuzzy ART instar, as follows. In the original instar (Grossberg, 1972), weights  $w_{ij}$  in paths projecting to an active target node  $j$  track activity  $x_i$  in the incoming paths. The specific instar equation used in fuzzy ART is:

$$\frac{d}{dt} w_{ij} = y_j (x_i - w_{ij}). \quad (57)$$

In fuzzy ART, the path signal from  $F_1$  is  $x_i = I_i \wedge w_{ji}$ , where  $y_J = 1$  at the chosen  $F_2$  node. With fast learning,  $w_{ij} = w_{ji}$  except initially, when  $w_{ij}(0) = 1 - \eta_{ij} = 1^-$ . Bottom-up weights  $w_{ij}$  remain constant for  $j \neq J$ . For  $j = J$ , weights decrease by instar learning:

$$\begin{aligned} \frac{d}{dt} w_{iJ} &= (x_i - w_{iJ}) = (I_i \wedge w_{ji} - w_{iJ}) \\ &= -(w_{iJ} - I_i \wedge w_{iJ}) = -[w_{iJ} - I_i]^+. \end{aligned} \quad (58)$$

Correspondingly, setting  $w_{ij} \equiv (1 - \tau_{ij})$  in the distributed instar (48) with choice at  $F_2$  gives:

$$\begin{aligned} \frac{d}{dt} w_{iJ} &= -\frac{d}{dt} \tau_{iJ} = -[y_J - \tau_{iJ} - I_i]^+ \\ &= -[1 - \tau_{iJ} - I_i]^+ = -[w_{iJ} - I_i]^+ \end{aligned} \quad (59)$$

and  $\frac{d}{dt} w_{ij} = 0$  for  $j \neq J$ . Thus the distributed instar with choice at  $F_2$  and fast learning reduces to the fuzzy ART instar (57).

### 4.3. Distributed Competitive Learning

In a competitive learning network (Grossberg, 1972, 1976b; Malsburg, 1973) inputs  $I_i$  filtered through adaptive pathways produce activations  $y_j$  at nodes of a target competitive field  $F_2$ . At active  $F_2$  nodes, instar learning (57) strengthens the net signal sent by the active input  $\mathbf{I}$ . In general, codes generated by competitive learning networks are unstable, never converging to a consistent representation for certain repeated input sequences (Grossberg, 1976b; Carpenter & Grossberg, 1987). In particular, with fast learning and activation  $\mathbf{y}$  distributed across all  $F_2$  nodes, all the weights  $w_{ij}$  would converge to the same value,  $I_i$ , with each input presentation, a form of catastrophic forgetting.

The  $F_0 \rightarrow F_2$  portion of the dART network, with distributed instar learning (48) and with  $F_2$  signals and activation as described in Section 2, constitutes a distributed competitive learning system. This network is, in fact, equivalent to the dART network with vigilance  $\rho \equiv 0$ , which eliminates the match/reset/search cycle. The  $F_0 \rightarrow F_2$  competitive network is thus a special case of a distributed ART network. The key to this distributed competitive learning

design is the dynamic weight  $[y_j - \tau_{ij}]^+$  that replaces the traditional multiplicative weight  $w_{ij}$ . The distributed instar learning law (48) holds constant all thresholds  $\tau_{ij}$  greater than  $(y_j - I_i)$ . Adaptive increase of a threshold  $\tau_{ij}$  requires a combination of a relatively small starting value  $\tau_{ij}(r)$ , a small path input  $I_i$ , and large coding node activation  $y_j$ . Since  $|y|=1$  but each  $I_i \in [0,1]$ , most thresholds will remain unchanged during learning. When the inequality  $\tau_{ij}(r) < (y_j - I_i)$  permits adaptation,  $\tau_{ij}$  rises only toward the upper limit  $(y_j - I_i)$ , where the dynamic weight  $[y_j - \tau_{ij}]^+$  equals the input  $I_i$ . In contrast, instar learning (57) permits adaptation for all positive  $y_j$ .

**Figure 4** (p. 41): Learning and search at synapse  $i$  of the  $j^{\text{th}}$   $F_2$  node

During distributed instar learning with a given input  $\mathbf{I}$ , the  $F_2$  code  $\mathbf{y}$  remains constant. However, learning may alter the code that this same input will activate later, as follows. Recall that  $\mathbf{y}$  is determined by the size of the  $F_0 \rightarrow F_2$  signal  $T_j(1)$  at the time  $t=r$  of the previous reset. By (24)-(25), each  $y_j$  is an increasing function of  $T_j(1)|_{t=r}$ . During learning, the quantity  $I_i \wedge [y_j - \tau_{ij}(t)]^+$  in the phasic term  $S_{ij}(y_j)$  (9) remains constant, since  $[y_j - \tau_{ij}(t)]^+$  decreases only if it is greater than  $I_i$  (Figure 4). In contrast, the quantity  $y_j \wedge \tau_{ij}(t)$  in the tonic term  $\Theta_{ij}(y_j)$  (13) increases whenever  $\tau_{ij}(t)$  increases, since  $\tau_{ij}(t)$  can increase only if  $y_j > \tau_{ij}(t)$  (48). If  $\mathbf{I}$  is presented again at a later time with no other learned changes having occurred, each term  $S_{ij}(1) = I_i \wedge (1 - \tau_{ij})$  will be the same as it had been when  $\mathbf{I}$  was previously presented and each term  $\Theta_{ij}(1) = \tau_{ij}$  will be the same or larger. Thus by (5)-(6), each increase in a threshold  $\tau_{ij}$  increases the net signal  $T_j(1)$  produced by the same input  $\mathbf{I}$ , all other things being equal. Since  $\mathbf{y}$  is normalized, learning tends to contrast-enhance the  $F_2$  coding pattern activated by a given input: learned changes tend to occur at nodes where  $y_j$  is large, so  $T_j(1)$  becomes larger and the  $j^{\text{th}}$  node will tend to gain an advantage the next time  $\mathbf{I}$  is presented.

Whereas learning can only increase the  $F_0 \rightarrow F_2$  signals  $T_j(1)$  for the active input  $\mathbf{I}$ , subsequent learned changes associated with different inputs could cause either an increase or a decrease in  $T_j(1)$  the next time  $\mathbf{I}$  is presented. Note that  $\tau_{ij}$  increases when an active  $F_2$  node  $j$  is coding an input in which the  $i^{\text{th}}$  component is small (48). The computations below show that, if this happens, the next time input  $\mathbf{I}$  is presented the larger threshold  $\tau_{ij}$  will cause a larger signal  $T_j(1)$  where  $I_i$  is small but will cause a smaller  $T_j(1)$  where  $I_i$  is large. That is, learning has caused node  $j$  to become more responsive to the set of all inputs where the  $i^{\text{th}}$  component is small.

**Figure 5** (p. 42): Effect of learned changes on coding signals.

Suppose that the last time  $\mathbf{I}$  was presented,  $\tau_{ij}$  was equal to  $\tau_{ij}^{old}$  but that  $\tau_{ij}$  has, in the mean time, risen to  $\tau_{ij}^{new}$ . Suppose that  $\mathbf{I}$  is now presented again. If  $I_i$  is small ( $I_i \leq 1 - \tau_{ij}^{new}$ ), then the phasic term  $S_{ij}(1)$  will be the same as before but the tonic term  $\Theta_{ij}(1)$  will have increased from  $\tau_{ij}^{old}$  to  $\tau_{ij}^{new}$  (Figure 5a). Thus when  $I_i$  is small an increased threshold  $\tau_{ij}$  leads to a larger signal  $T_j(1)$ , by (5), (6), and (12). With choice-by-difference (15),  $T_j(1)$  increases by  $(1 - \alpha)(\tau_{ij}^{new} - \tau_{ij}^{old})$ . If  $I_i$  is large ( $I_i \geq 1 - \tau_{ij}^{old}$ ), then  $S_{ij}(1)$  will have decreased from  $(1 - \tau_{ij}^{old})$  to  $(1 - \tau_{ij}^{new})$  while  $\Theta_{ij}(1)$  will have increased by the same amount, from  $\tau_{ij}^{old}$  to  $\tau_{ij}^{new}$  (Figure 5c). Thus (6) implies that, when  $I_i$  is large, an increased threshold  $\tau_{ij}$  leads to a smaller signal  $T_j(1)$ . With choice-by-difference,  $T_j(1)$  decreases by  $\alpha(\tau_{ij}^{new} - \tau_{ij}^{old})$ . If  $I_i$  is in between ( $1 - \tau_{ij}^{new} \leq I_i < 1 - \tau_{ij}^{old}$ ), an increased threshold  $\tau_{ij}$  may lead either to a smaller or a larger signal  $T_j(1)$ , depending on the function  $g_j(S_j, \Theta_j)$  that defines  $T_j$  (5). The phasic term  $S_{ij}(1)$  will have decreased from  $I_i$  to  $(1 - \tau_{ij}^{new})$  while the tonic term  $\Theta_{ij}(1)$  will have increased from  $\tau_{ij}^{old}$  to  $\tau_{ij}^{new}$  (Figure 5b). With choice-by-difference, the change in  $T_j(1)$  is:

$$\begin{aligned} \Delta T_j(1) &= \left(1 - \tau_{ij}^{new} - I_i\right) + (1 - \alpha)\left(\tau_{ij}^{new} - \tau_{ij}^{old}\right) \\ &= \left(1 - \tau_{ij}^{old} - I_i\right) - \alpha\left(\tau_{ij}^{new} - \tau_{ij}^{old}\right). \end{aligned} \quad (60)$$

Thus the increased threshold  $\tau_{ij}$  leads to a larger signal  $T_j(1)$  only if the choice parameter  $\alpha$  is small enough; that is if:

$$(1 - \alpha)\tau_{ij}^{old} + \alpha\tau_{ij}^{new} < (1 - I_i). \quad (61)$$

## 5. A DISTRIBUTED ART ALGORITHM

The algorithm below summarizes distributed ART (Figure 1b) computations with inputs  $\mathbf{I}^{(1)}, \mathbf{I}^{(2)}, \dots, \mathbf{I}^{(n)}, \dots$  presented for equal time intervals. An algorithm to approximate dART dynamics for a continuously varying input  $\mathbf{I}(t)$  would set  $\mathbf{I}^{(n)} = \mathbf{I}(n\Delta t)$ , with the time step  $\Delta t$  and the learning rate parameter  $\beta$  small. Other dART variations are implemented with appropriate substitutions.

(1) Variables:  $i = 1 \dots M, \quad j = 1 \dots N$

STM	MTM	LTM	$F_0 \rightarrow F_2$ signal	$F_2 \rightarrow F_1$ signal
$I_i - F_0$ (input)	$\Delta_{ij} -$ Phasic	$\tau_{ij} - F_0 \rightarrow F_2$	$S_j -$ Phasic	$\sigma_i -$ Total
$x_i - F_1$ (matching)	$\delta_{ij} -$ Tonic	$\tau_{ji} - F_2 \rightarrow F_1$	$\Theta_j -$ Tonic	
$y_j - F_2$ (coding)			$T_j -$ Total	

(2) Signal rule: Define the  $F_0 \rightarrow F_2$  signal function  $T_j = g_j(S_j, \Theta_j)$ , where  $g_j(0,0) = 0$   
and  $\frac{\partial g_j}{\partial S_j} > \frac{\partial g_j}{\partial \Theta_j} > 0$  for  $S_j > 0$  and  $\Theta_j > 0$ .

E.g.,  $T_j = S_j + (1 - \alpha)\Theta_j$  with  $\alpha \in (0,1)$  (choice-by-difference) or  
 $T_j = S_j / (\alpha + My_j - \Theta_j)$  with  $\alpha > 0$  (Weber law).

(3) CAM rule: Define the  $F_2$  steady-state function  $y_j = f_j(T_1 \dots T_N)$ , where  $\partial f_j / \partial T_j \geq 0$ .

E.g., For a power  $p > 0$  (power law)  $y_j = \begin{cases} \frac{(T_j)^p}{\sum_{\lambda \in \Lambda} (T_\lambda)^p} & \text{if } j \in \Lambda \\ 0 & \text{if } j \notin \Lambda \end{cases}$  where

$\Lambda = \{j: T_j \geq \bar{T}\}$  with  $\bar{T} = \frac{1}{N} \sum_{j=1}^N T_j$  (above-average- $T_j$ ); or

$\Lambda =$  the set of  $Q$  indices  $j$  where  $T_j$  is maximal ( $Q$ -max).

#### (4) Parameters

Number of input components -  $i = 1 \dots M$

Number of coding nodes -  $j = 1 \dots N$

Signal rule - E.g.,  $\alpha \in (0,1)$  (choice-by-difference) or  $\alpha > 0$  (Weber law)

CAM rule - E.g.,  $p$  (power law) and  $Q$  ( $Q$ -max), with  $p \rightarrow \infty$  or  $Q = 1$  for choice

Learning rate -  $\beta \in [0,1]$ , with  $\beta = 1$  for fast learning

Vigilance -  $\rho \in [0,1]$

A set of small, positive, random numbers, for initial  $\tau_{ij}$  values -  $\eta_{ij} = 0^+$

(5) First iteration:  $n = 1$

MTM depletion -  $\Delta_{ij} = \delta_{ij} = 0$

$F_0 \rightarrow F_2$  threshold -  $\tau_{ij} = \eta_{ij}$

$F_2 \rightarrow F_1$  threshold -  $\tau_{ji} = 0$

Input -  $I_i = I_i^{(1)}$



(6) Reset: New STM steady state at  $F_2$  and  $F_1$

$F_0 \rightarrow F_2$  signal

$$\text{Phasic - } S_j = \sum_{i=1}^M [I_i \wedge (1 - \tau_{ij}) - \Delta_{ij}]^+$$

$$\text{Tonic - } \Theta_j = \sum_{i=1}^M [\tau_{ij} - \delta_{ij}]^+$$

$$\text{Total - } T_j = g_j(S_j, \Theta_j) \quad ((2) \text{ Signal rule})$$

$$F_2 \text{ activation - } y_j = f_j(T_1 \dots T_N) \quad ((3) \text{ CAM rule})$$

$$F_2 \rightarrow F_1 \text{ signal - } \sigma_i = \sum_{j=1}^N [y_j - \tau_{ji}]^+$$

$$F_1 \text{ activation - } x_i = I_i \wedge \sigma_i$$

(7) MTM depletion:  $F_2$  sites refractory on the time scale of search

$$\text{Phasic - } \Delta_{ij}^{old} = \Delta_{ij}$$

$$\Delta_{ij} = \Delta_{ij}^{old} \vee (I_i \wedge [y_j - \tau_{ij}]^+)$$

$$\text{Tonic - } \delta_{ij}^{old} = \delta_{ij}$$

$$\delta_{ij} = \delta_{ij}^{old} \vee (y_j \wedge \tau_{ij})$$

(8) Reset or resonance: Check the  $F_1$  matching criterion

$$\text{If } \sum_{i=1}^M x_i < \rho \sum_{i=1}^M I_i, \text{ go to } (6) \text{ Reset}$$

$$\text{If } \sum_{i=1}^M x_i \geq \rho \sum_{i=1}^M I_i, \text{ go to } (9) \text{ Resonance}$$

(9) Resonance: New LTM thresholds and MTM recovery on the time scale of learning

$$\text{Old values - } \tau_{ij}^{old} = \tau_{ij}, \quad \tau_{ji}^{old} = \tau_{ji}, \quad \sigma_i^{old} = \sigma_i$$

$$\text{Increase } F_0 \rightarrow F_2 \text{ threshold - } \tau_{ij} = \tau_{ij}^{old} + \beta \left[ y_j - \tau_{ij}^{old} - I_i \right]^+$$

$$\text{Increase } F_2 \rightarrow F_1 \text{ threshold - } \tau_{ji} = \tau_{ji}^{old} + \beta \frac{\left[ \sigma_i^{old} - I_i \right]^+}{\sigma_i^{old}} \left[ y_j - \tau_{ji}^{old} \right]^+$$

$$\text{Decrease } F_2 \rightarrow F_1 \text{ signal - } \sigma_i = \sigma_i^{old} - \beta \left[ \sigma_i^{old} - I_i \right]^+$$

MTM recovery -

$$\Delta_{ij} = \delta_{ij} = 0$$

(10) Next iteration: Increase  $n$  by 1

$$\text{New input - } I_i = I_i^{(n)}$$

$$\text{New } F_1 \text{ activation - } x_i = I_i \wedge \sigma_i$$

Go to (6) Reset

## 6. DISTRIBUTED ART GEOMETRY

A geometric interpretation of fuzzy ART represents categories as boxes in input space that expand during learning (Carpenter, Grossberg, and Rosen, 1991). A generalized version of this geometric representation illustrates dART dynamics, as follows.

### 6.1. Complement Coding

In fuzzy ART, input normalization prevents a type of category proliferation that could otherwise occur when weights erode. Complement coding doubles the dimension of an input vector  $\mathbf{a} \equiv (a_1 \dots a_M)$  by concatenating  $\mathbf{a}$  and its complement  $\mathbf{a}^c$ . The input to a fuzzy ART network is then a  $2M$ -dimensional vector:

$$\mathbf{I} = \mathbf{A} \equiv (\mathbf{a}, \mathbf{a}^c), \quad (62)$$

where

$$\left( \mathbf{a}^c \right)_i \equiv (1 - a_i). \quad (63)$$

Complement coded inputs are normalized because

$$|\mathbf{A}| = \left| \left( \mathbf{a}, \mathbf{a}^c \right) \right| = \sum_{i=1}^M a_i + \sum_{i=1}^M (1 - a_i) = M. \quad (64)$$

If  $\mathbf{a}$  represents input features, then complement coding allows a learned category representation to encode the degree to which each feature is consistently absent as well as the degree to which it is consistently present when that category is active. Because of its computational advantages, complement coding is used in nearly all fuzzy ART and fuzzy ARTMAP applications. Similar advantages can be expected for dART and dARTMAP applications. Except for changing the number of components of  $\mathbf{I}$ ,  $\mathbf{x}$ , and the corresponding LTM vectors from  $M$  to  $2M$ , the description of network dynamics is unchanged since complement coding is only a preprocessing step. A dART algorithm with complement coding can be embedded in an ARTMAP network to form the basis of a dARTMAP algorithm (Section 8).

## 6.2. Fuzzy ART Category Boxes

A geometric interpretation of fuzzy ART associates with each weight vector  $\mathbf{w}_j \equiv (w_{1j} \dots w_{2M,j})$  a box  $R_j$  in  $M$ -dimensional space. In the  $i^{\text{th}}$  dimension ( $i = 1 \dots M$ ), the side of the  $j^{\text{th}}$  box is defined by the interval  $[w_{ij}, w_{i+M,j}^c]$ . That is,  $R_j$  is the set of points  $\mathbf{q}$  for which:

$$w_{ij} \leq q_i \leq (1 - w_{i+M,j}). \quad (65)$$

The size of  $R_j$  is defined as the sum of these intervals:

$$|R_j| = \sum_{i=1}^M \left( (1 - w_{i+M,j}) - w_{ij} \right) = M - \sum_{i=1}^{2M} w_{ij} = M - |\mathbf{w}_j|. \quad (66)$$

When  $M = 2$ ,

$$\mathbf{A} = \left( \mathbf{a}, \mathbf{a}^c \right) \equiv \left( a_1, a_2, a_1^c, a_2^c \right) \quad (67)$$

and category boxes are rectangles in the plane (Figure 6a). Note that, formally, the interval would be “reversed” if  $w_{ij} > (1 - w_{i+M,j})$ . Initially, all  $w_{ij} = 1$  and  $|\mathbf{w}_j| = 2M$ , so initially  $|R_j| = -M$ . During learning,  $R_j$  may grow toward a maximum size  $M$  as weights shrink. Because top-down weights  $w_{ji}$  equal bottom-up weights  $w_{ij}$  in fuzzy ART,  $R_j$  can represent both.

**Figure 6** (p. 43): ART and dART geometry

When a fuzzy ART  $F_2$  node  $j$  is chosen,  $\sigma(\mathbf{y}) = \mathbf{w}_j$  and the matched  $F_1$  pattern  $\mathbf{x} = \mathbf{I} \wedge \mathbf{w}_j$  must satisfy the vigilance criterion (30) for  $j$  to remain active (Figure 1a). This is equivalent to requiring that, for category  $j$  to remain active during a learning interval,

$$|R_j \oplus \mathbf{a}| \leq M(1 - \rho) \quad (68)$$

where  $R_j \oplus \mathbf{a}$  is the smallest box containing both  $R_j$  and  $\mathbf{a}$ . When the  $j^{\text{th}}$   $F_2$  node does remain active, instar learning (57) implies that  $w_{ij}$  may decrease toward  $a_i = A_i$  and  $w_{i+M,j}$  may decrease toward  $a_i^c = A_{i+M}$ . As weights shrink, the size of the interval  $[w_{ij}, w_{i+M,j}^c]$  expands. A wide interval signals that the  $i^{\text{th}}$  feature is uninformative with respect to the  $j^{\text{th}}$  category: since both weights  $w_{ij}$  and  $w_{i+M,j}$  are then small, the corresponding feature has been neither consistently present nor consistently absent when the  $j^{\text{th}}$   $F_2$  node has been active. When node  $j$  remains active during a fast learning interval, box  $R_j$  expands to  $R_j \oplus \mathbf{a}$ . Thus, with category choice and fast learning,  $R_j$  is the smallest box that contains all the training set inputs  $\mathbf{a}$  coded by category  $j$ .

### 6.3. Distributed ART Coding Boxes and Matching Boxes

A geometric representation of distributed ART substitutes dART dynamic weights  $[y_j - \tau_{ij}]^+$  for the fuzzy ART weights  $w_{ij}$ . For each  $j=1\dots N$ , where fuzzy ART weights define a single category box  $R_j$ , dART dynamic weights define a family of *coding boxes*  $R_j(y_j)$ , one for each  $y_j \in [0,1]$ . Fuzzy ART boxes  $R_j$  can represent top-down matching as well as bottom-up category activation since only one  $F_2$  node at a time is active and  $w_{ji} = w_{ij}$ . In dART, however, the  $F_2 \rightarrow F_1$  input vector  $\sigma(\mathbf{y})$  (28) may depend on activities of all  $F_2$  nodes. Top-down dynamic weights  $[y_j - \tau_{ji}]^+$  therefore define a *matching box*  $R(\mathbf{y})$  for each  $F_2$  activity vector  $\mathbf{y}$ .

A distributed ART coding box  $R_j(y_j)$  depends on the  $F_2$  activity level  $y_j$  as well as on the  $F_0 \rightarrow F_2$  thresholds  $\tau_{ij}$  ( $i=1\dots 2M$ ). For each  $y_j \in [0,1]$ ,  $R_j(y_j)$  is the set of points  $\mathbf{q}$  for which:

$$[y_j - \tau_{ij}]^+ \leq q_i \leq \left(1 - [y_j - \tau_{i+M,j}]^+\right) \quad (69)$$

(Figure 6b). With  $w_{ij} \equiv 1 - \tau_{ij}$  and  $y_j = 1$ , the dART coding box  $R_j(y_j)$  is the same as the fuzzy ART category box  $R_j$ . As  $y_j$  decreases from 1 to 0, the box  $R_j(y_j)$  grows, filling the entire unit box when  $y_j$  smaller than all the thresholds  $\tau_{ij}$  ( $i=1\dots 2M$ ), i.e., when all the dynamic weights  $[y_j - \tau_{ji}]^+$  equal 0. Ignoring MTM aftereffects (i.e., with  $\Delta_{ij} = \delta_{ij} = 0$ ), the size of  $R_j(y_j)$  is:

$$\begin{aligned}
|R_j(y_j)| &= \sum_{i=1}^M \left( \left( 1 - [y_j - \tau_{i+M,j}]^+ \right) - [y_j - \tau_{ij}]^+ \right) \\
&= M - \sum_{i=1}^{2M} [y_j - \tau_{ij}]^+ = M - \sum_{i=1}^{2M} (y_j - y_j \wedge \tau_{ij}) \\
&= M - 2My_j + \sum_{i=1}^{2M} y_j \wedge \tau_{ij} = M(1 - 2y_j) + \Theta_j(y_j).
\end{aligned} \tag{70}$$

Thus  $R_j(y_j)$  represents the tonic component  $\Theta_j(y_j)$  (12)-(13) of the  $F_0 \rightarrow F_2$  signal  $T_j(y_j)$  (5). The expanded box  $R_j(y_j) \oplus \mathbf{a}$  represents the phasic component  $S_j(y_j)$  (8)-(9), as follows.

For a given input  $\mathbf{a} \equiv (a_1 \dots a_M)$ ,  $R_j(y_j) \oplus \mathbf{a}$  is the set of points  $\mathbf{q}$  where:

$$\left\{ [y_j - \tau_{ij}]^+ \wedge a_i \right\} \leq q_i \leq \left\{ \left( 1 - [y_j - \tau_{i+M,j}]^+ \right) \vee a_i \right\} \tag{71}$$

(Figure 6c). For  $i = 1 \dots M$ ,

$$[y_j - \tau_{ij}]^+ \wedge a_i = [y_j - \tau_{ij}]^+ \wedge A_i \tag{72}$$

and:

$$\begin{aligned}
\left( 1 - [y_j - \tau_{i+M,j}]^+ \right) \vee a_i &= \left( 1 - [y_j - \tau_{i+M,j}]^+ \right) \vee (1 - (1 - a_i)) \\
&= 1 - [y_j - \tau_{i+M,j}]^+ \wedge (1 - a_i) = 1 - [y_j - \tau_{i+M,j}]^+ \wedge A_{i+M}.
\end{aligned} \tag{73}$$

Thus,

$$\begin{aligned}
|R_j(y_j) \oplus \mathbf{a}| &= \sum_{i=1}^M \left( \left\{ \left( 1 - [y_j - \tau_{i+M,j}]^+ \right) \vee a_i \right\} - \left\{ [y_j - \tau_{ij}]^+ \wedge a_i \right\} \right) \\
&= \sum_{i=1}^M \left( 1 - [y_j - \tau_{i+M,j}]^+ \wedge A_{i+M} - [y_j - \tau_{ij}]^+ \wedge A_i \right) \\
&= M - \sum_{i=1}^{2M} [y_j - \tau_{ij}]^+ \wedge A_i = M - S_j(y_j).
\end{aligned} \tag{74}$$

Therefore the expanded box  $R_j(y_j) \oplus \mathbf{a}$  represents the phasic component  $S_j(y_j)$  of the  $F_0 \rightarrow F_2$  signal  $T_j(y_j)$ .

The boxes  $R_j(y_j)$  and  $R_j(y_j) \oplus \mathbf{a}$  also provide a geometric representation of the distributed choice-by-difference signal rule. Defining the distance  $d(R_j, \mathbf{a})$  from  $R_j$  to  $\mathbf{a}$  by:

$$d(R_j, \mathbf{a}) \equiv |R_j \oplus \mathbf{a}| - |R_j|, \quad (75)$$

(70) and (74) imply that the choice-by-difference signal function (15) can be written as:

$$\begin{aligned} T_j(y_j) &= S_j(y_j) + (1 - \alpha)\Theta_j(y_j) \\ &= \left(M - |R_j(y_j) \oplus \mathbf{a}|\right) + (1 - \alpha)\left(|R_j(y_j)| - M(1 - 2y_j)\right) \\ &= M\left(1 - (1 - \alpha)(1 - 2y_j)\right) - d(R_j(y_j), \mathbf{a}) - \alpha|R_j(y_j)|. \end{aligned} \quad (76)$$

Recall that the values  $y_1 \dots y_N$  will assume following a reset are determined by  $T_1(1) \dots T_N(1)$  (24). By (76),

$$\begin{aligned} T_j(1) &= \left(M - |R_j(1) \oplus \mathbf{a}|\right) + (1 - \alpha)\left(|R_j(1)| + M\right) \\ &= M(2 - \alpha) - d(R_j(1), \mathbf{a}) - \alpha|R_j(1)|. \end{aligned} \quad (77)$$

Geometric interpretation of distributed choice-by-difference thus shows that, except for MTM aftereffects during search, an input  $\mathbf{a}$  will most strongly activate an  $F_2$  node  $j$  when  $\mathbf{a}$  is in or near  $R_j(1)$  and when  $R_j(1)$  is small. The relative importance of distance vs. size depends on the choice parameter  $\alpha$ . When  $\alpha$  is close to 0, a maximal  $T_j(1)$  is one that minimizes the distance from  $R_j(1)$  to  $\mathbf{a}$ , with the size of  $R_j(1)$  used only to break ties if  $\mathbf{a}$  is contained in more than one coding box. When  $\alpha$  is close to 1, a maximal  $T_j(1)$  is one that minimizes the size of the expanded box  $R_j(1) \oplus \mathbf{a}$ .

During distributed instar learning, while  $\mathbf{a}$  and  $\mathbf{y}$  remain constant,  $R_j(y_j)$  expands toward  $R_j(y_j) \oplus \mathbf{a}$  (Figure 6c) as the tonic terms  $y_j \wedge \tau_{ij}$  in  $\Theta_j(y_j)$  grow and the phasic terms  $A_i \wedge [y_j - \tau_{ij}]^+$  in  $S_j(y_j)$  remains constant. No threshold  $\tau_{ij}$  will change during learning if  $\mathbf{a}$  is contained in  $R_j(y_j)$ , even if  $\mathbf{a}$  is not contained in  $R_j(1)$ . With fast learning following a reset at time  $t = r$ , thresholds grow from  $\tau_{ij}^{old} \equiv \tau_{ij}(r)$  to  $\tau_{ij}^{new}$  and:

$$\begin{aligned}
\sum_{i=1}^{2M} |\Delta \tau_{ij}| &\equiv \sum_{i=1}^{2M} (\tau_{ij}^{new} - \tau_{ij}^{old}) = \sum_{i=1}^{2M} [y_j - A_i - \tau_{ij}^{old}]^+ \\
&= \sum_{i=1}^{2M} \left( [y_j - \tau_{ij}^{old}]^+ - [y_j - \tau_{ij}^{old}]^+ \wedge A_i \right) \\
&= \sum_{i=1}^{2M} [y_j - \tau_{ij}^{old}]^+ - \sum_{i=1}^{2M} [y_j - \tau_{ij}^{old}]^+ \wedge A_i \\
&= \left( M - |R_j^{old}(y_j)| \right) - \left( M - |R_j^{old}(y_j) \oplus \mathbf{a}| \right) \\
&= d(R_j(y_j), \mathbf{a}) \Big|_{t=r}.
\end{aligned} \tag{78}$$

That is, the total threshold change at the  $j^{th}$   $F_2$  node equals the distance from  $R_j(y_j)$  to  $\mathbf{a}$  at time  $t = r$ , the start of the learning interval. Thus, setting  $\alpha = 0^+$ , which favors  $F_2$  nodes for which  $R_j(1)$  is closest to  $\mathbf{a}$  in (77), also favors nodes that will minimize the total  $F_0 \rightarrow F_2$  threshold change during learning. In fuzzy ART, the parameter limit where  $\alpha$  is close to 0 was called the *conservative limit*, since category choice then favors weight conservation wherever possible.

Once a dART code  $\mathbf{y}$  becomes active, the signal  $\sigma_i(\mathbf{y})$  from  $F_2$  to the  $i^{th}$   $F_1$  node equals the sum of the top-down dynamic weights  $[y_j - \tau_{ji}]^+$ . The signals  $\sigma_i(\mathbf{y})$  define a matching box  $R(\mathbf{y})$  as the set of points  $\mathbf{q}$  where:

$$\sigma_i(\mathbf{y}) \leq q_i \leq 1 - \sigma_{i+M}(\mathbf{y}) \tag{79}$$

for  $i = 1 \dots M$  (Figure 6d). The expanded box  $R(\mathbf{y}) \oplus \mathbf{a}$  is the set of points  $\mathbf{q}$  where:

$$\sigma_i(\mathbf{y}) \wedge a_i \leq q_i \leq (1 - \sigma_{i+M}(\mathbf{y})) \vee a_i. \tag{80}$$

As in (72)-(73),

$$\sigma_i(\mathbf{y}) \wedge a_i = \sigma_i(\mathbf{y}) \wedge A_i \tag{81}$$

and:

$$(1 - \sigma_{i+M}(\mathbf{y})) \vee a_i = 1 - \sigma_{i+M}(\mathbf{y}) \wedge A_{i+M} \tag{82}$$

for  $i = 1 \dots M$ . Thus, by (29), (64), and (80),

$$\begin{aligned}
|R(\mathbf{y}) \oplus \mathbf{a}| &= \sum_{i=1}^M \left( \{ (1 - \sigma_{i+M}(\mathbf{y})) \vee a_i \} - \{ \sigma_i(\mathbf{y}) \wedge a_i \} \right) \\
&= M - \sum_{i=1}^{2M} \sigma_i(\mathbf{y}) \wedge A_i = |\mathbf{A}| - |\sigma(\mathbf{y}) \wedge \mathbf{A}| = |\mathbf{A}| - |\mathbf{x}|.
\end{aligned} \tag{83}$$

Therefore, the active dART code  $\mathbf{y}$  meets the matching criterion (31) when:

$$|R(\mathbf{y}) \oplus \mathbf{a}| \leq M(1 - \rho), \quad (84)$$

as in fuzzy ART (68). Geometrically, resonance requires that the expanded box  $R(\mathbf{y}) \oplus \mathbf{a}$  not be too big, by (84). When the matching criterion is met,  $R(\mathbf{y})$  expands toward  $R(\mathbf{y}) \oplus \mathbf{a}$  during learning. No top-down learned changes occur if  $\mathbf{a}$  is already contained in  $R(\mathbf{y})$ .

## 7. DISTRIBUTED ART COMPUTATION

The dART algorithm (Section 5) summarizes a general solution to the distributed ART system of equations. Once a specific network and parameters are selected for a particular application, computational analysis is usually required to trace network coding in response to a given input sequence. When network dimensions are small, as in the examples below, explicit system solutions are simple enough to permit direct calculation, without use of a computer. Each example uses a choice-by-difference signal rule, a power law CAM rule, fast learning, and 2 or 3 coding nodes at  $F_2$  to illustrate dART activation, search, and learning. As in Figure 6, 2-dimensional inputs are complement-coded, so  $\mathbf{a} = (a_1, a_2)$  and  $\mathbf{I} = \mathbf{A} = (\mathbf{a}, \mathbf{a}^c)$ .

### 7.1. dART Learning

Figure 7 illustrates distributed ART learning in a system with dimensions  $M = N = 2$  and input  $\mathbf{a} = (0.7, 0.8)$ . The index set  $\Lambda = \{1, 2\} = \{1 \dots N\}$  in the power law CAM rule, so  $y_j = T_j^p(1) / (T_1^p(1) + T_2^p(1))$  for  $j=1, 2$ . For  $j=1$ , initial threshold values  $(\tau_{1j} \dots \tau_{4j}) = (0.9, 0.3, 0.4, 0.9)$  are represented by the coding box  $R_1(1)$ , with  $d(R_1(1), \mathbf{a}) = 0.3$  and  $|R_1(1)| = 0.5$ . For  $j=2$ , initial threshold values  $(\tau_{1j} \dots \tau_{4j}) = (0.9, 0.9, 0.9, 0.2)$  are represented by the coding box  $R_2(1)$ , with  $d(R_2(1), \mathbf{a}) = 0.6$  and  $|R_2(1)| = 0.9$ . By (77),

$$T_j(1) = 2(2 - \alpha) - d(R_j(1), \mathbf{a}) - \alpha |R_j(1)|. \quad (85)$$

Thus when  $\alpha = 0.2$ ,  $(T_1(1), T_2(1)) = (3.20, 2.82)$  and, when  $p=1$ ,  $(y_1, y_2) = (0.532, 0.468)$  (Figure 7a). Then  $\mathbf{a} \in R_1(y_1)$  but  $\mathbf{a} \notin R_2(y_2)$ :  $d(R_1(y_1), \mathbf{a}) = 0$  and  $d(R_2(y_2), \mathbf{a}) = 0.068$ . During learning,  $R_2(y_2)$  expands to include  $\mathbf{a}$  as  $\tau_{42}$  increases from 0.2 to 0.268. If  $\mathbf{a}$  is repeatedly presented and no other learned changes take place,  $\tau_{42}$  will continue to increase toward 0.274, the point where  $(y_1, y_2) = (0.526, 0.474)$ , where  $R_2(y_2)$  would just include  $\mathbf{a}$ .

If the power  $p$  increases to 5, with the network otherwise the same,  $(y_1, y_2) = (0.653, 0.347)$  (Figure 7b). Compared to the case where  $p=1$ , the higher power stores a more contrast-enhanced representation of the signal  $(T_1(1), T_2(1))$  in the CAM system at  $F_2$ . In this case  $\mathbf{a} \in R_1(y_1)$  and  $\mathbf{a} \in R_2(y_2)$ , so no changes occur during learning.

If the choice parameter  $\alpha$  increases to 0.8 but the network is otherwise the same as in Figure 7b, the signal  $(T_1(1), T_2(1)) = (1.70, 1.08)$ . Then  $\mathbf{y}$  is further contrast-enhanced, with



$(y_1, y_2) = (0.906, 0.094)$  (Figure 7c). In this case,  $\mathbf{a} \in R_2(y_2)$  but  $\mathbf{a} \notin R_1(y_1)$ :  $d(R_2(y_2), \mathbf{a}) = 0$  and  $d(R_1(y_1), \mathbf{a}) = 0.206$ . During learning,  $R_1(y_1)$  expands to include  $\mathbf{a}$  as  $\tau_{31}$  increases from 0.4 to 0.606.

If  $\mathbf{a}$  is presented again later with no other learned changes having taken place in the mean time (and no MTM distortion),  $(y_1, y_2) = (0.916, 0.084)$  (Figure 7d). Learning has thus contrast-enhanced the code  $\mathbf{y}$ . If  $\mathbf{a}$  is repeatedly presented and no other learned changes take place,  $\tau_{31}$  will continue to increase toward 0.616, the point where  $(y_1, y_2) = (0.916, 0.084)$ , where  $R_1(y_1)$  would just include  $\mathbf{a}$ .

**Figure 7** (p. 44): Distributed ART activation and learning

**Table 1** (p. 37): Distributed ART activation and learning

Table 1 shows steady-state  $\mathbf{y}$  values of the system described above (Figure 7) as the power  $p$  increases from 1 to 5 and as the choice parameter  $\alpha$  increases from 0.01 to 0.99. During learning,  $\tau_{ij} = \tau_{31}$  increases for  $y_1 > 0.7$ , when  $\mathbf{a} \notin R_1(y_1)$ , since  $d(R_1(1), \mathbf{a}) = 0.3$ ; and  $\tau_{ij} = \tau_{42}$  increases for  $y_2 > 0.4$ , when  $\mathbf{a} \notin R_2(y_2)$ , since  $d(R_2(1), \mathbf{a}) = 0.6$  (boldface values of  $y_j$ ). In all other cases, no changes occur during learning. If the same input  $\mathbf{a}$  is presented again and no other learned changes have meanwhile occurred, a larger  $\tau_{ij} = \tau_{31}$  value implies a larger rectangle  $R_1(1)$ , a smaller distance  $d(R_1(1), \mathbf{a})$ , and larger  $T_1(1)$  and  $y_1$  values. The code  $(y_1, y_2)$  is thus contrast-enhanced by the learning process. On the other hand, a larger  $\tau_{ij} = \tau_{42}$  value, which implies a larger rectangle  $R_2(1)$ , a smaller distance  $d(R_2(1), \mathbf{a})$ , and larger  $T_2(1)$  and  $y_2$  values, would make the code  $(y_1, y_2)$  more uniform.

## 7.2. dART Search

Figure 8 illustrates distributed ART search in a system that is much like the one in Figure 7a except that  $F_2$  has three coding nodes ( $N = 3$ ). A distributed choice-by-difference signal rule sets the choice parameter  $\alpha = 0.2$ , a power law CAM rule sets  $p = 1$ , and input  $\mathbf{a} = (0.7, 0.8)$ . For  $j = 1, 2$ , thresholds are the same as the initial  $\tau_{ij}$  values in Section 7.1. With  $\Lambda$  equal to the index set of above-average  $T_j$ , activity  $y_j$  is proportional to  $T_j(1)$  when  $T_j(1)$  is greater than or equal to the average ( $\bar{T}$ ); otherwise  $y_j = 0$ . By hypothesis,  $T_3(1) \leq 2.44$  so that initially, with all  $\Delta_{ij} = \delta_{ij} = 0$ ,  $T_1(1) = 3.20 > T_2(1) = 2.82 \geq \bar{T} > T_3(1)$ . Thus  $y_1 = 0.532$  and  $y_2 = 0.468$ , as in Figure 7a; and  $y_3 = 0$ . While this code  $\mathbf{y}$  is active, the MTM depletion terms  $\Delta_{ij}$  and  $\delta_{ij}$  quickly go to equilibrium, sending the phasic terms  $S_{ij}(y_j)$  (9) and the tonic terms  $\Theta_{ij}(y_j)$  (13) to 0. Then for  $j = 1$ ,  $(\Delta_{1j} \dots \Delta_{4j}) = (0.0, 0.232, 0.132, 0.0)$ , which will reduce  $S_1(1)$  by 0.364 at reset; and  $(\delta_{1j} \dots \delta_{4j}) = (0.532, 0.3, 0.4, 0.532)$ , which will reduce  $\Theta_1(1)$  by 1.764 at reset. Thus (15) implies that, with  $\alpha = 0.2$ ,  $T_1(1)$  will be reduced by 1.78 at reset. For  $j = 2$ ,  $(\Delta_{1j} \dots \Delta_{4j}) = (0.0, 0.0, 0.0, 0.2)$ , which will reduce  $S_2(1)$  by 0.2 at reset; and  $(\delta_{1j} \dots \delta_{4j}) = (0.468, 0.468, 0.468, 0.2)$ , which will reduce  $\Theta_2(1)$  by 1.604 at reset. Thus,

with  $\alpha = 0.2$ ,  $T_2(1)$  will be reduced by 1.48 at reset. Since  $y_3 = 0$ ,  $\Delta_{ij}$  and  $\delta_{ij}$  remain equal to 0 for  $j = 3$  and  $i = 1 \dots 4$ .

**Figure 8** (p. 45): Distributed ART search

**Table 2** (p. 37): Distributed ART search

A reset with input **a** still active would then leave  $T_3(1)$  unchanged but would reduce  $T_1(1)$  from 3.20 to 1.42 and would reduce  $T_2(1)$  from 2.82 to 1.34. What the next code **y** will be depends on the size of  $T_3(1)$  (Table 2). When  $T_3(1)$  is large ( $1.5 < T_3(1) \leq 2.44$ ), node  $j = 3$  is the only one active following reset, since  $T_1(1)$  and  $T_2(1)$  are then below average. With smaller  $T_3(1)$  values ( $1.38 \leq T_3(1) \leq 1.5$ ), nodes  $j = 1$  and  $j = 3$  share activation following reset. With even smaller  $T_3(1)$  values ( $1.26 < T_3(1) < 1.38$ ),  $T_2(1)$  and  $T_3(1)$  are below average, so node  $j = 1$  is the only one active following reset. Finally, when  $T_3(1)$  is very small ( $0 \leq T_3(1) \leq 1.26$ ), nodes  $j = 1$  and  $j = 2$  share activation following reset, as they did before. However, the code **y** is now more uniform, with  $y_1$  smaller and  $y_2$  larger before the reset.

## 8. DISTRIBUTED ARTMAP

ARTMAP networks for supervised learning self-organize mappings from input vectors, representing features such as patient history and test results, to output vectors, representing predictions such as the likelihood of an adverse outcome following an operation. The original binary ARTMAP (Carpenter, Grossberg, & Reynolds, 1991) incorporates two ART 1 modules,  $ART_a$  and  $ART_b$ , that are linked by a *map field*  $F^{ab}$ . At the map field the network forms associations between categories via outstar learning and triggers search, via the ARTMAP match tracking rule, when a training set input fails to make a correct prediction. Match tracking increases the  $ART_a$  vigilance parameter  $\rho_a$  in response to a predictive error at  $ART_b$ . Fuzzy ARTMAP (Carpenter, Grossberg, Markuzon, Reynolds, & Rosen, 1992) substitutes fuzzy ART for ART 1. Distributed ARTMAP (dARTMAP) substitutes dART for fuzzy ART and distributed outstar learning for outstar learning at the map field (Figure 9a).

**Figure 9** (p. 46): Distributed ARTMAP

Many applications of supervised learning systems such as ARTMAP are classification problems, where the trained system tries to predict a correct category given a test set input vector. A prediction might be a single category or distributed as a set of scores or probabilities. For this class of problems, the dARTMAP architecture illustrated in Figure 9b does not require the full  $dART_b$  architecture. Even in this case, however, dARTMAP implementation requires a number of judicious design choices, in contrast to the few choices required for fuzzy ARTMAP implementation. Recent benchmark simulation studies have demonstrated that, with fast learning and noisy training data, dARTMAP maintains the predictive accuracy of fuzzy ARTMAP while dramatically improving code compression. Ongoing research seeks to characterize how a distributed learning system such as dARTMAP can combine speed, performance, generalization, and code compression in a variety of new applications.

## REFERENCES

- Bachelder, I.A., Waxman, A.M., & Seibert, M. (1993). A neural system for mobile robot visual place learning and recognition. In *Proceedings of the World Congress on Neural Networks (WCNN'93)* (pp. I-512-517). Hillsdale, NJ: Lawrence Erlbaum Associates.
- Baloch, A.A., & Waxman, A.M. (1991). Visual learning, adaptive expectations, and behavioral conditioning of the mobile robot MAVIN. *Neural Networks*, **4**, 271-302.
- Baraldi, A., & Parmiggiani, F. (1995). A neural network for unsupervised categorization of multivalued input patterns: An application of satellite image clustering. *IEEE Transactions on Geoscience and Remote Sensing*, **33**, 305-316.
- Bernardon, A.M., & Carrick, J.E. (1995). A neural system for automatic target learning and recognition applied to bare and camouflaged SAR targets. *Neural Networks*, **8**, 1103-1108.
- Carpenter, G.A. (1994a). A distributed outstar network for spatial pattern learning. *Neural Networks*, **7**, 159-168.
- Carpenter, G.A. (1994b). Distributed recognition codes and catastrophic forgetting. In *Proceedings of the World Congress on Neural Networks (WCNN'94)* (pp. IV-133-142). Hillsdale, NJ: Lawrence Erlbaum Associates.
- Carpenter, G.A., & Gjaja, M.N. (1994). Fuzzy ART choice functions. In *Proceedings of the World Congress on Neural Networks (WCNN'94)* (pp. I-713-722). Hillsdale, NJ: Lawrence Erlbaum Associates.
- Carpenter, G.A., & Grossberg, S. (1987). A massively parallel architecture for a self-organizing neural pattern recognition machine. *Computer Vision, Graphics, and Image Processing*, **37**, 54-115.
- Carpenter, G.A., & Grossberg, S. (1991). Pattern recognition by self-organizing neural networks. Cambridge, MA: MIT Press.
- Carpenter, G.A., Grossberg, S., Markuzon, N., Reynolds, J.H., & Rosen, D.B. (1992). Fuzzy ARTMAP: A neural network architecture for incremental supervised learning of analog multidimensional maps. *IEEE Transactions on Neural Networks*, **3**, 698-713.
- Carpenter, G.A., Grossberg, S., & Reynolds, J.H. (1991). ARTMAP: Supervised real-time learning and classification of nonstationary data by a self-organizing neural network. *Neural Networks*, **4**, 565-588.
- Carpenter, G.A., Grossberg, S., & Rosen, D.B. (1991). Fuzzy ART: Fast stable learning and categorization of analog patterns by an adaptive resonance system. *Neural Networks*, **4**, 759-771.
- Carpenter, G.A., & Markuzon, N. (1996). ARTMAP-IC and medical diagnosis: Instance counting and inconsistent cases. CAS/CNS Technical Report CAS/CNS-96-017, Boston, MA: Boston University.

- Carpenter, G.A., & Ross, W.D. (1993). ART-EMAP: A neural network architecture for learning and prediction by evidence accumulation. In *Proceedings of the World Congress on Neural Networks (WCNN'94)* (pp. III - 649-656). Hillsdale, NJ: Lawrence Erlbaum Associates.
- Carpenter, G.A., & Ross, W.D. (1995). ART-EMAP: A neural network architecture for object recognition by evidence accumulation. *IEEE Transactions on Neural Networks*, **6**, 805-818.
- Caudell, T.P., & Healy, M.J. (1994). Adaptive Resonance Theory networks in the Encephalon autonomous vision system. In *Proceedings of the 1994 IEEE International Conference on Neural Networks (ICNN'94)* (pp. II-1235-1240). Piscataway, NJ: IEEE.
- Caudell, T.P., Smith, S.D.G., Escobedo, R., & Anderson, M. (1994). NIRS: Large scale ART-1 neural architectures for engineering design retrieval. *Neural Networks*, **7**, 1339-1350.
- Christodoulou, C.G., Huang, J., Georgiopoulos, M., & Liou, J.J. (1995). Design of gratings and frequency selective surfaces using Fuzzy ARTMAP neural networks. *Journal of Electromagnetic Waves and Applications*, **9**, 17-36.
- Dubrawski, A., & Crowley, J.L. (1994). Learning locomotion reflexes: A self-supervised neural system for a mobile robot. *Robotics and Autonomous Systems*, **12**, 133-142.
- Gan, K.W., & Lua, K.T. (1992). Chinese character classification using an Adaptive Resonance network. *Pattern Recognition*, **25**, 877-88.
- Gjerdingen, R.O. (1990). Categorization of musical patterns by self-organizing neuronlike networks. *Music Perception*, **7**, 339-370.
- Gopal, S., Sklarew, D.M., & Lambin, E. (1994). Fuzzy-neural networks in multi-temporal classification of landcover change in the Sahel. In *Proceedings of the DOSES Workshop on New Tools for Spatial Analysis*. Lisbon, Portugal, DOSES, EUROSTAT. ECSC-EC-EAEC: Brussels, Luxembourg, pp. 55-68.
- Grossberg, S. (1968). A prediction theory for some nonlinear functional-differential equations, I: Learning of lists. *Journal of Mathematical Analysis and Applications*, **21**, 643-694.
- Grossberg, S. (1970). Some networks that can learn, remember, and reproduce any number of complicated space-time patterns, II. *Studies in Applied Mathematics*, **49**, 135-166.
- Grossberg, S. (1972). Neural expectation: Cerebellar and retinal analogs of cells fired by learnable or unlearned pattern classes. *Kybernetik*, **10**, 49-57.
- Grossberg, S. (1976a). Adaptive pattern classification and universal recoding, II: Feedback, expectation, olfaction, and illusions. *Biological Cybernetics*, **23**, 187-202.
- Grossberg, S. (1976b). Adaptive pattern classification and universal recoding, I: Parallel development and coding of neural feature detectors. *Biological Cybernetics*, **23**, 121-134.
- Grossberg, S. (1980). How does a brain build a cognitive code? *Psychological Review*, **87**, 151.
- Ham, F.M., & Han, S. (1996). Classification of cardiac arrhythmias using fuzzy ARTMAP. *IEEE Transactions on Biomedical Engineering*, **43**, 425-430.

- Kalkunte, S.S., Kumar, J.M., & Patnaik, L.M. (1992). A neural network approach for high resolution fault diagnosis in digital circuits. *Proceedings of the International Joint Conference on Neural Networks*, I, (pp. 83-88). Piscataway, NJ: IEEE.
- Kasperkiewicz, J., Racz, J., & Dubrawski, A. (1995). HPC strength prediction using artificial neural network. *Journal of Computing in Civil Engineering*, 9, 279-284.
- Kim, J.W., Jung, K.C., Kim, S.K., & Kim, H.J. (1995). Shape classification of on-line Chinese character strokes using ART 1 neural network. *Proceedings of the World Congress on Neural Networks (WCNN'95)* (pp. II-191-194). Hillsdale, NJ: Lawrence Erlbaum Associates.
- Koch, M.W., Moya, M.M., Hostetler, L.D., and Fogler, R.J. (1995). Cueing, feature discovery, and one-class learning for synthetic aperture radar automatic target recognition. *Neural Networks*, 8, 1081-1102.
- Ly, S., & Choi, J.J. (1994). Drill condition monitoring using ART-1. In *Proceedings of the IEEE International Conference on Neural Networks (ICNN'94)* (pp. II-1226-1229). Piscataway, NJ: IEEE.
- Malsburg, C. von der (1973). Self-organization of orientation sensitive cells in the striate cortex. *Kybernetik*, 14, 85-100.
- Markram, H., & Tsodyks, M. (1996). Redistribution of synaptic efficacy between neocortical pyramidal neurons. *Nature*, 382, 807-810.
- Mehta, B.V., Vij, L., & Rabelo, L.C. (1993). Prediction of secondary structures of proteins using fuzzy ARTMAP. In *Proceedings of the World Congress on Neural Networks (WCNN'93)* (pp. I-228-232). Hillsdale, NJ: Lawrence Erlbaum Associates.
- Murshed, N.A., Bortozzi, F., & Sabourin, R. (1995). Off-line signature verification, without *a priori* knowledge of class  $\omega_2$ . A new approach. In *Proceedings of the Third International Conference on Document Analysis and Recognition (ICDAR'95)*.
- Nicholls, D. G. (1994). Proteins, transmitters and synapses. Oxford: Blackwell Science Ltd.
- Racz, J., & Dubrawski, A. (1995). Artificial neural network for mobile robot topological localization. *Robotics and Autonomous Systems*, 16, 73-80.
- Rubin, M.A. (1995). Application of fuzzy ARTMAP and ART-EMAP to automatic target recognition using radar range profiles. *Neural Networks*, 8, 1109-1116.
- Seibert, M., & Waxman, A.M. (1992). Adaptive 3D object recognition from multiple views. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 14, 107-124.
- Seibert, M., & Waxman, A.M. (1993). An approach to face recognition using saliency maps and caricatures. In *Proceedings of the World Congress on Neural Networks (WCNN'93)* (pp. III-661-664). Hillsdale, NJ: Lawrence Erlbaum Associates.
- Soliz, P., & Donohoe, G.W. (1996). Adaptive resonance theory neural network for fundus image segmentation. *Proceedings of the World Congress on Neural Networks (WCNN'96)* (pp. 1180-1183). Hillsdale, NJ: Lawrence Erlbaum Associates.

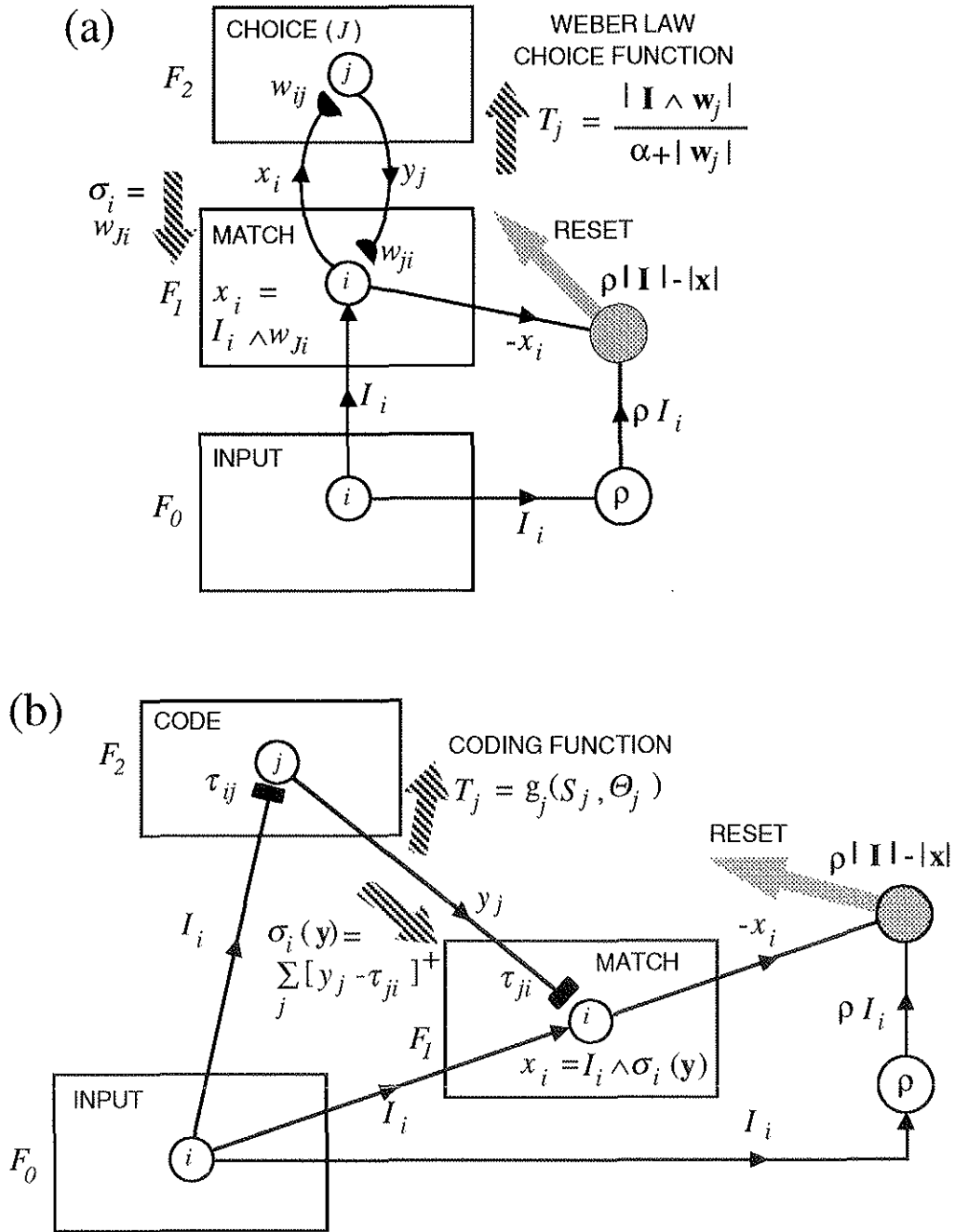
- Srinivasa, N., & Sharma, R. (1996). A self-organizing invertible map for active vision applications. *Proceedings of the World Congress on Neural Networks (WCNN'96)* (pp. 121-124). Hillsdale, NJ: Lawrence Erlbaum Associates.
- Suzuki, Y. (1995). Self-organizing QRS-wave recognition in ECG using neural networks. *IEEE Transactions on Neural Networks*, **6**, 1469-1477.
- Tarng, Y.S., Li, T.C., & Chen, M.C. (1994) Tool failure monitoring for drilling processes. In *Proceedings of the 3rd International Conference on Fuzzy Logic, Neural Nets and Soft Computing* (pp. 109-111), Iizuka, Japan.
- Tse, P., & Wang, D.D. (1996). A hybrid neural networks based machine condition forecaster and classifier by using multiple vibration parameters. *Proceedings of the 1994 IEEE International Conference on Neural Networks*, **IV**, (pp. 2096-2100). Piscataway, NJ: IEEE.
- Waxman, A.M., Seibert, M.C., Gove, A., Fay, D.A., Bernardon, A.M., Lazott, C., Steele, W.R., & Cunningham, R.K. (1995). Neural processing of targets in visible, multispectral IR and SAR imagery. *Neural Networks*, **8**, 1029-1051.
- Wienke, D. (1994). Neural resonance and adaption - Towards nature's principles in artificial pattern recognition. In L. Buydens and W. Melssen (Eds.), *Chemometrics: Exploring and exploiting chemical information*. Nijmegen, NL: University Press.
- Zadeh, L. (1965). Fuzzy sets. *Information and Control*, **8**, 338-353.

**Table 1:** Distributed ART activation and learning.

	CBD signal	$p=1$	$p=2$	$p=5$	$p=5$
	$(T_1(1), T_2(1))$	$(y_1, y_2)$	$(y_1, y_2)$	$(y_1, y_2)$	$(y_1, y_2)$
				$\xrightarrow{\text{LEARNING}}$	
$\alpha = 0.01$	(3.67, 3.37)	(0.521, <b>0.479</b> )	(0.543, <b>0.457</b> )	(0.605, 0.395)	(0.605, 0.395)
$\alpha = 0.20$	(3.20, 2.82)	(0.532, <b>0.468</b> )	(0.563, <b>0.437</b> )	(0.653, 0.347)	(0.653, 0.347)
$\alpha = 0.50$	(2.45, 1.95)	(0.557, <b>0.443</b> )	(0.612, 0.388)	( <b>0.758</b> , 0.242)	( <b>0.769</b> , 0.231)
$\alpha = 0.80$	(1.70, 1.08)	(0.612, 0.388)	( <b>0.713</b> , 0.287)	( <b>0.906</b> , 0.094)	( <b>0.916</b> , 0.084)
$\alpha = 0.99$	(1.23, 0.53)	(0.699, 0.301)	( <b>0.843</b> , 0.157)	( <b>0.985</b> , 0.015)	( <b>0.985</b> , 0.015)

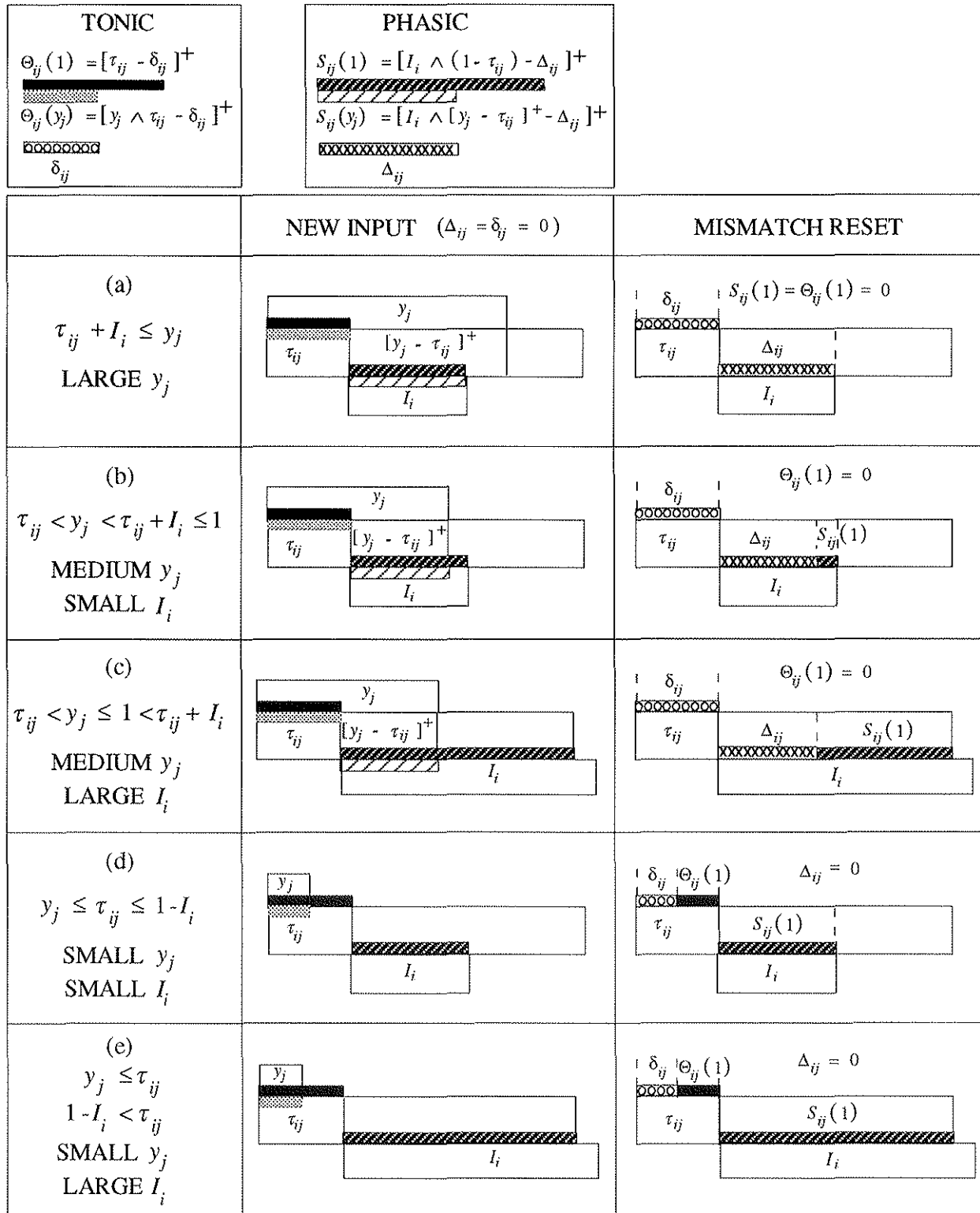
**Table 2:** Distributed ART search in response to an input  $\mathbf{a} = (0.7, 0.8)$ , with complement coding, a power law CAM rule ( $p=1$ ) for above-average  $T_j(1)$ , a choice-by-difference signal rule ( $\alpha=0.2$ ), and  $N=3$ .

<b>Before reset:</b> $T_1(1) = 3.20, T_2(1) = 2.82$				
$T_3(1)$	$y_1$	$y_2$	$y_3$	$\bar{T}$
$0 \leq T_3 \leq 2.44$	0.532	0.468	0	$2.01 \leq \bar{T} \leq 2.82$
<hr/>				
<b>After reset:</b> $T_1(1) = 1.42, T_2(1) = 1.34$				
$T_3(1)$	$y_1$	$y_2$	$y_3$	$\bar{T}$
(a) $1.5 < T_3 \leq 2.44$	0	0	1	$1.42 < \bar{T} \leq 1.73$
(b) $1.38 \leq T_3 \leq 1.5$	(*)	0	(**)	$1.38 \leq \bar{T} \leq 1.42$
(c) $1.26 < T_3 < 1.38$	1	0	0	$1.34 < \bar{T} < 1.38$
(d) $0 \leq T_3 \leq 1.26$	0.514	0.486	0	$0.92 \leq \bar{T} \leq 1.34$
(*) $0.507 \geq \frac{1.42}{1.42 + T_3} \geq 0.486$		(**) $0.493 \leq \frac{T_3}{1.42 + T_3} \leq 0.514$		

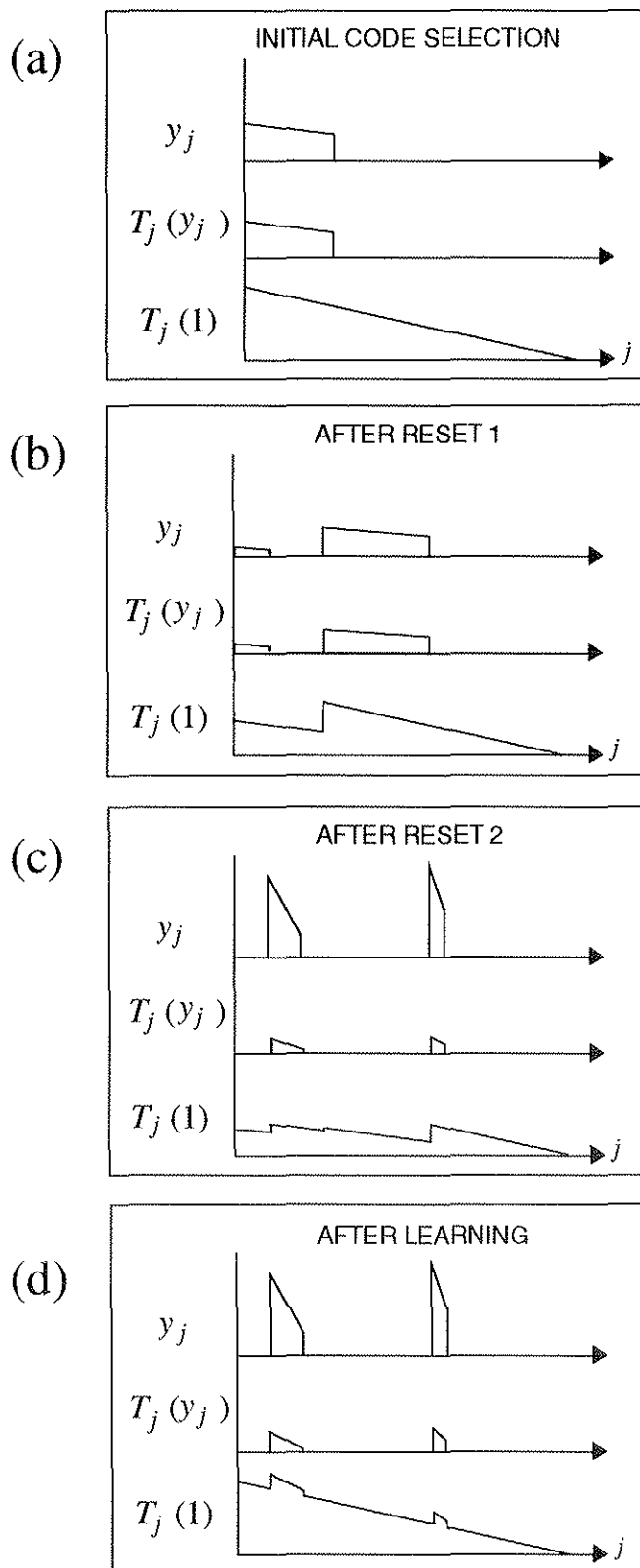


**Figure 1:** Fuzzy ART and distributed ART. (a) In fuzzy ART, the  $F_2$  the node ( $j = J$ ) that receives the largest input  $T_j$  from  $F_1$  becomes active. Activity  $x$  at the field  $F_1$  reflects the match between the bottom-up input  $I$  and the top-down input  $\sigma$ , which is equal to the weight vector  $w_j$ . When  $x$  fails to meet the vigilance matching criterion, reset leaves node  $J$  refractory on the time scale of search. Refractory nodes recover on the time scale of learning. (b) Like fuzzy ART, distributed ART computes a matched pattern  $x$  at  $F_1$  and resets  $F_2$  if  $x$  fails to meet the vigilance matching criterion. In dART, however,  $F_2$  receives input directly from  $F_0$ . The code  $y$ , which is a function of phasic components  $S_j$  and tonic components  $\Theta_j$ , may be arbitrarily distributed. The  $i^{th}$   $F_1$  node receives a positive signal from each  $F_2$  node at which activity  $y_j$  exceeds an  $F_2 \rightarrow F_1$  adaptive threshold  $\tau_{ji}$ . With choice at  $F_2$  and fast learning, distributed ART is computationally equivalent to fuzzy ART.



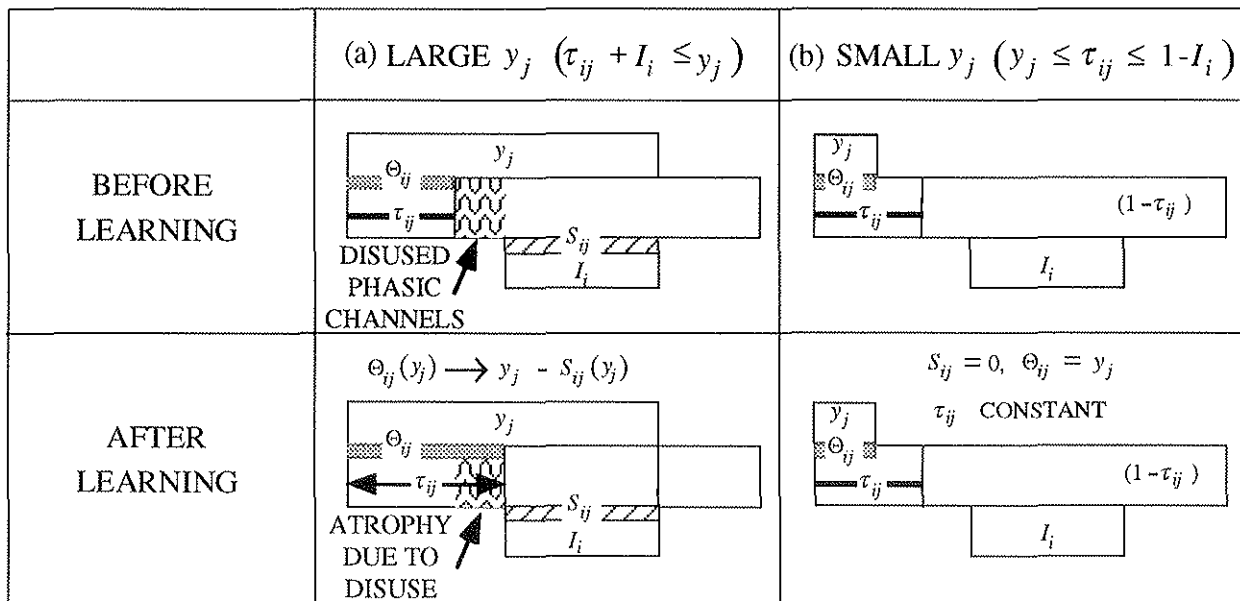


**Figure 2:** Visual representation of distributed instar signal components as a fraction of total membrane sites. The phasic term  $S_{ij}(y_j)$  and the tonic term  $\Theta_{ij}(y_j)$  depend on the adaptive threshold  $\tau_{ij}$  at the  $i^{\text{th}}$  synapse of the  $j^{\text{th}}$   $F_2$  node. At reset, nonspecific arousal momentarily sends all  $y_j \rightarrow 1$ . The terms  $S_{ij}(1)$  and  $\Theta_{ij}(1)$  at the time of reset then determine the next code  $y$ . A given  $y_j$  value gates membrane sites, so that  $S_{ij}(y_j)$  and  $\Theta_{ij}(y_j)$  may be large for large  $y_j$  but must be small for small  $y_j$ . Phasic and tonic terms thus correspond to membrane processes that are gated by postsynaptic voltage ( $y_j$ ), and the phasic term  $S_{ij}$  is also gated by the released presynaptic transmitter, or ligand ( $I_i$ ). After reset mismatch, previously active sites  $\Delta_{ij}$  (phasic) and  $\delta_{ij}$  (tonic) are depleted, or refractory, and remain so on an MTM time scale. During a search, phasic and tonic terms  $S_{ij}(1)$  and  $\Theta_{ij}(1)$  can be large only if  $y_j$  has recently remained small.



**Figure 3:** Parallel distributed search, with the  $F_2$  code  $y_j$  proportional to  $T_j(1)$  for  $j \in \Lambda \subseteq \{1 \dots N\}$  and a choice-by-difference signal rule. (a)  $T_j(y_j) = 0$  for  $j \notin \Lambda$ . (b) After reset,  $T_j(1)$  is diminished by the previous value of  $T_j(y_j)$ . A new set  $\Lambda$  of  $F_2$  nodes where  $T_j(1)$  is maximal leads to a new active code  $\mathbf{y}$ . (c) Following another reset on the time scale of search,  $T_j(1)$  is further reduced by the previous value of  $T_j(y_j)$ . (d) Refractory sites recover on the time scale of learning, so  $T_j(1)$  reverts to its original value at inactive sites while  $T_j(1)$  may increase where  $y_j > 0$ . These values of  $T_j(1)$  would determine the next code  $\mathbf{y}$  if another reset should then occur with the same I due, say, to a sudden increase in vigilance.

$\frac{d}{dt} \tau_{ij} = y_j - S_{ij}(y_j) - \Theta_{ij}(y_j)$	<b>TONIC:</b> $\Theta_{ij}(y_j) = y_j \wedge \tau_{ij}$
<b>NO MTM:</b> $\Delta_{ij} = \delta_{ij} = 0$	<b>PHASIC:</b> $S_{ij}(y_j) = I_i \wedge [y_j - \tau_{ij}]^+$

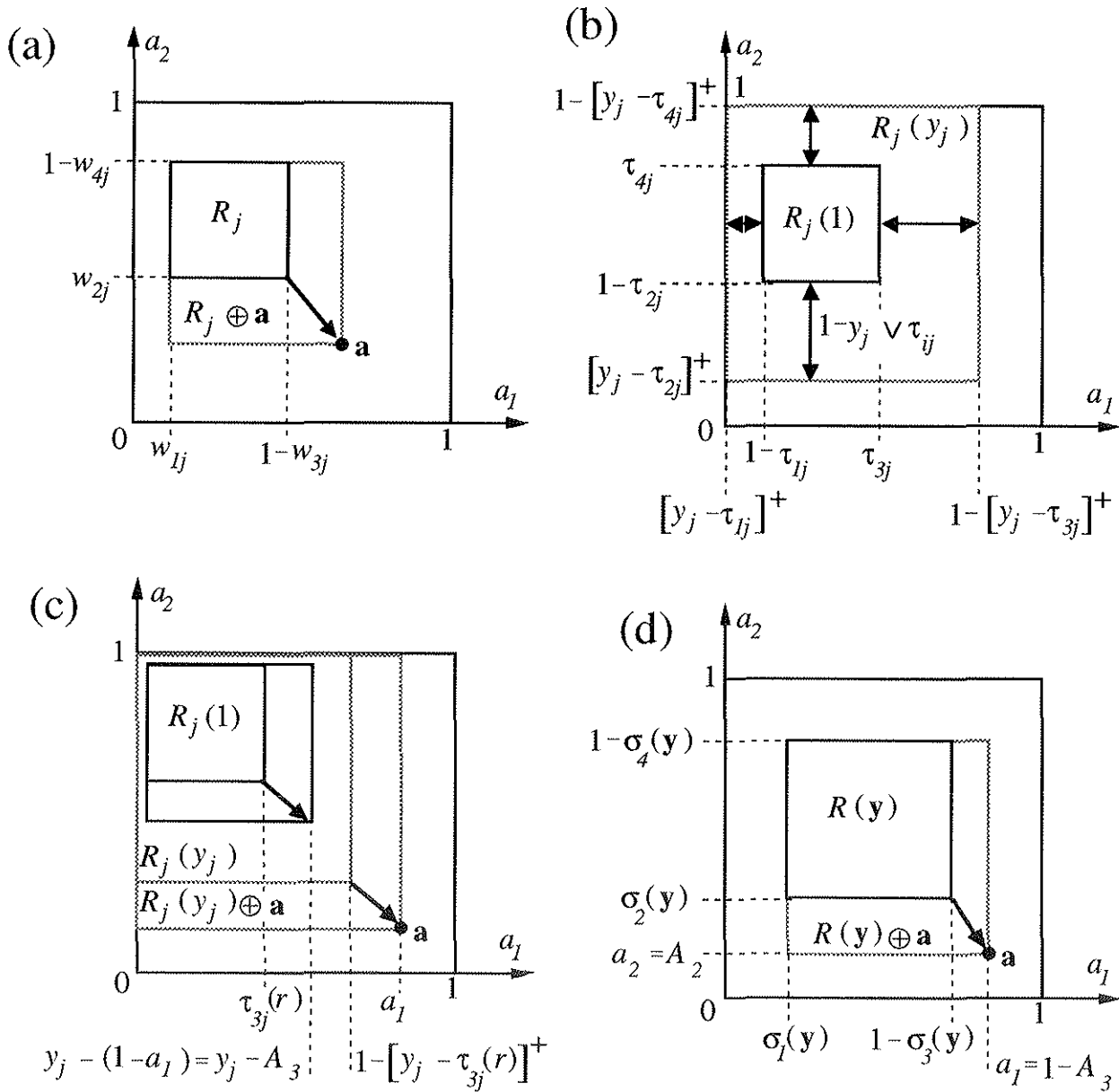


**Figure 4:** Distributed instar learning at synapse  $i$  of the  $j^{\text{th}}$   $F_2$  node: disused phasic channels (pattern) that are primed by  $y_j$  but not occupied by  $I_i$  revert to tonic channels. (a) A large  $y_j$  may permit the threshold  $\tau_{ij}$  to increase during learning. When  $\tau_{ij}$  is increasing, the tonic terms increase because then  $\Theta_{ij}(y_j) = \Theta_{ij}(1) = \tau_{ij}$  while the phasic terms remain constant because then  $S_{ij}(y_j) = S_{ij}(1) = I_i$ . (b) A small  $y_j$  tends to leave  $\tau_{ij}$  constant during learning because then  $\Theta_{ij}(y_j) = y_j$  and  $S_{ij}(y_j) = 0$ .

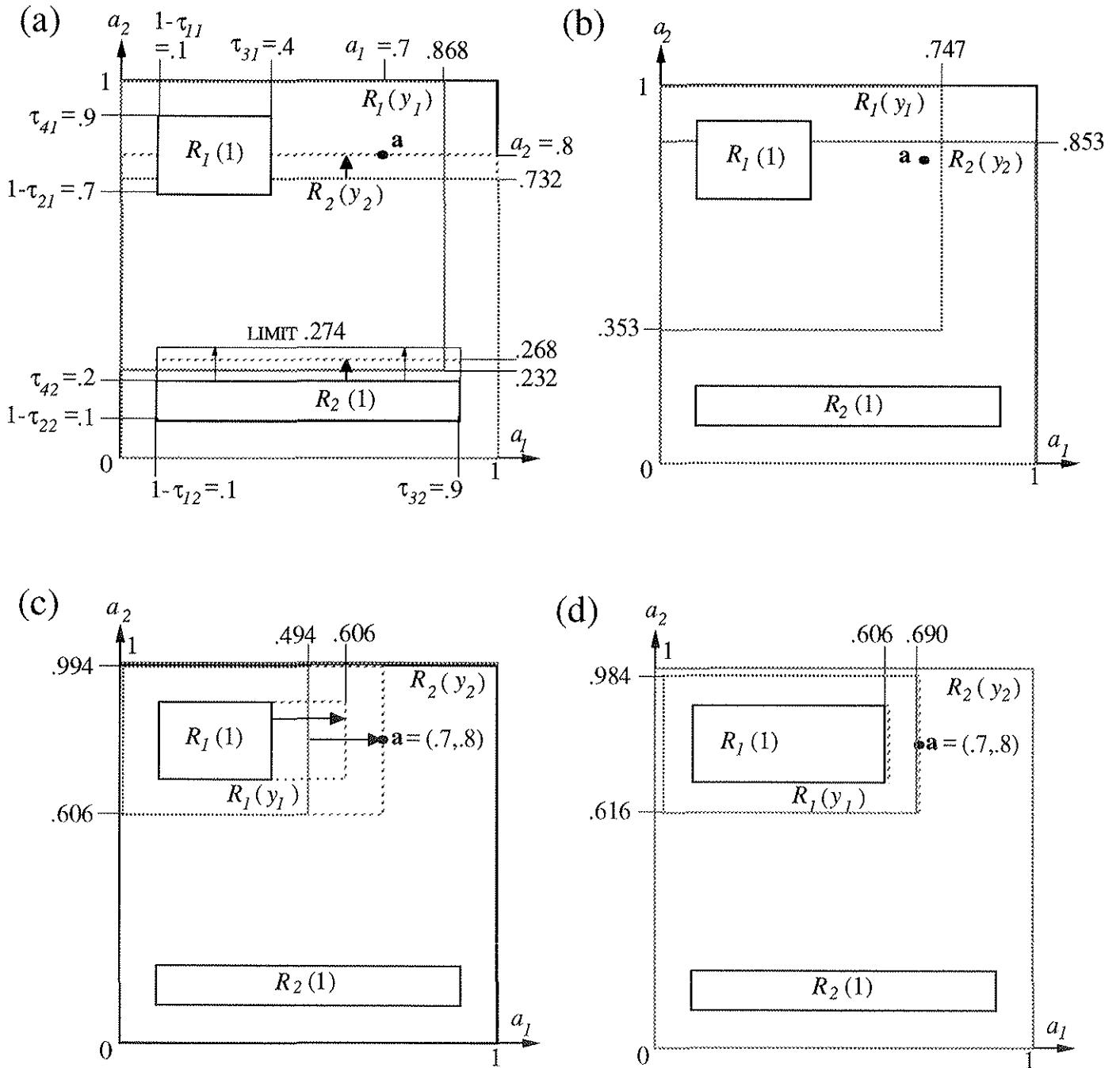
AT RESET: $y_j = 1$	TONIC: $\Theta_{ij}(1) = \tau_{ij}$
NO MTM: $\Delta_{ij} = \delta_{ij} = 0$	PHASIC: $S_{ij}(1) = I_i \wedge (1 - \tau_{ij})$

	$\tau_{ij} = \tau_{ij}^{old}$	$\tau_{ij} = \tau_{ij}^{new}$
(a) SMALL $I_i$ $I_i + \tau_{ij}^{old} < 1$ $I_i + \tau_{ij}^{new} \leq 1$		POSTSYNAPTIC SIGNAL INCREASES 
(b) MEDIUM $I_i$ $I_i + \tau_{ij}^{old} < 1$ $\leq I_i + \tau_{ij}^{new}$		
(c) LARGE $I_i$ $1 \leq I_i + \tau_{ij}^{old}$ $< I_i + \tau_{ij}^{new}$		POSTSYNAPTIC SIGNAL DECREASES 

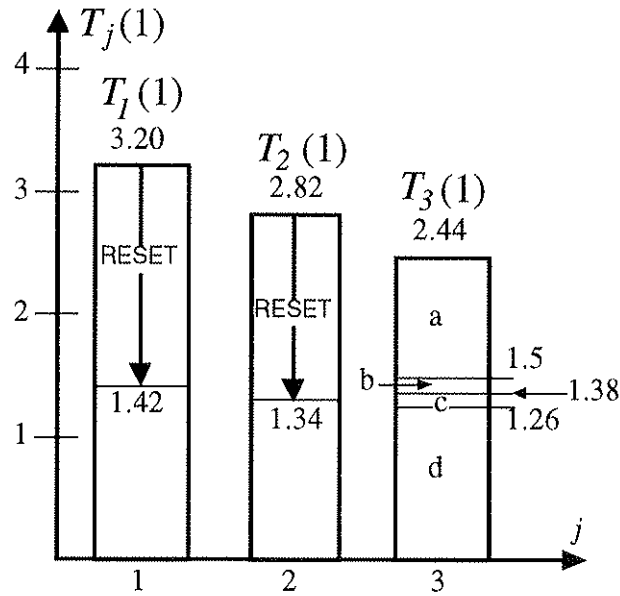
**Figure 5:** Effect of learned changes on coding signals: an increase in the threshold  $\tau_{ij}$  between presentations of an input  $\mathbf{I}$  may make  $T_j(1)$  larger or smaller the next time  $\mathbf{I}$  is presented, depending on the size of  $I_i$ . That is, although learning causes a monotonic change in the LTM representation at the level of receptors ( $\tau_{ij}$ ), this change can resemble either LTP (for a single test pulse or small  $I_i$ ) or LTD (for larger  $I_i$ ) at the level of the postsynaptic potential ( $y_j$ ). (a) When  $I_i$  is small, a higher threshold  $\tau_{ij}$  makes the tonic term  $\Theta_{ij}$  larger while the phasic term  $S_{ij}$  stays the same, so  $T_j(1)$  is larger. (b) When  $I_i$  is neither large nor small, a higher threshold makes  $\Theta_{ij}$  larger and  $S_{ij}$  smaller, and  $T_j(1)$  may be larger or smaller, depending on the signal rule that defines it. (c) When  $I_i$  is large, a higher threshold increases  $\Theta_{ij}$  and decreases  $S_{ij}$  by equal amounts, making  $T_j(1)$  smaller.



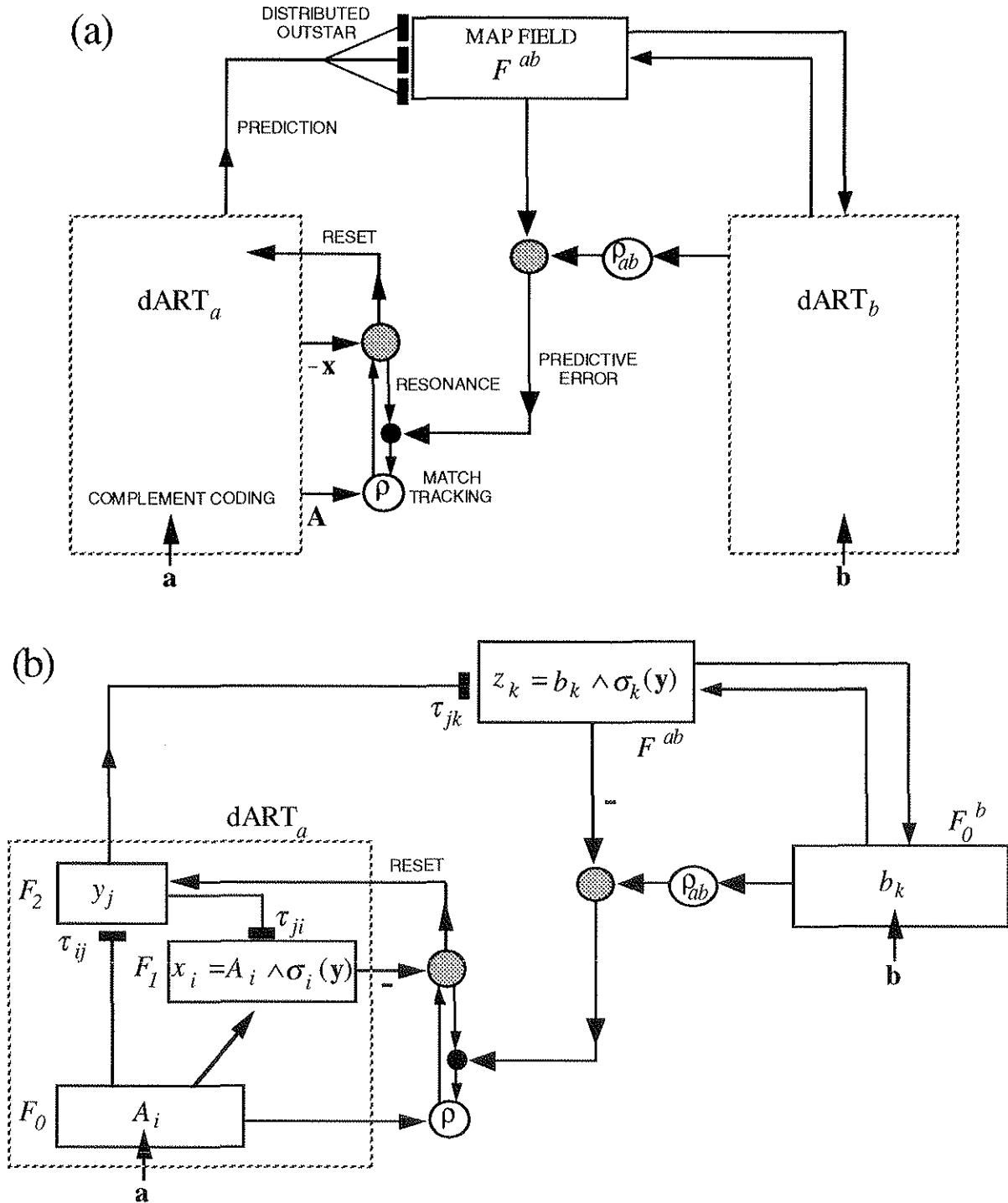
**Figure 6:** ART and dART geometry. (a) The fuzzy ART category box  $R_j$  provides a geometric representation of each weight vector  $w_j$ . Since bottom-up and top-down weight vectors are equal, category boxes can represent the dynamics of choice, search, and learning at  $F_1$  and  $F_2$ . During learning, the chosen box  $R_j$  expands toward  $R_j \oplus \mathbf{a}$ . Before this can occur, however, the search process resets node  $j$  and sends  $T_j$  to 0 if the size of the expanded box  $R_j \oplus \mathbf{a}$  would be greater than  $M(1-\rho)$ . (b) Distributed ART replaces the bottom-up fuzzy ART weights  $w_{ij}$  with a family of dynamic weights  $[y_j - \tau_{ij}]^+$  and replaces the category box  $R_j$  with a corresponding nested family of coding boxes  $R_j(y_j)$ . The dART box  $R_j(1)$  corresponds to the fuzzy ART box  $R_j$ . (c) During distributed instar learning, with activity  $y_j$  at the  $j^{\text{th}}$  node, the box  $R_j(y_j)$  expands toward  $R_j(y_j) \oplus \mathbf{a}$  ( $j=1 \dots N$ ) as some adaptive thresholds  $\tau_{ij}$  increase ( $i=1 \dots 2M$ ). Since  $R_j(0)$  fills the square, no thresholds  $\tau_{ij}$  change when  $y_j = 0$ . The boxes  $R_j(1)$  that will determine the next code  $\mathbf{y}$  expand as much as the larger boxes  $R_j(y_j)$ . However,  $R_j(1)$  will reach  $\mathbf{a}$  only if  $y_j = 1$  or if  $\mathbf{a}$  was already contained in  $R_j(1)$  at the time of the previous reset. (d) When the code  $\mathbf{y}$  is active, a matching box  $R(\mathbf{y})$  represents the  $F_2 \rightarrow F_1$  inputs  $\sigma_i(\mathbf{y})$ . With choice at  $F_2$ ,  $\sigma_i(\mathbf{y}) = (1 - \tau_{ji}) \equiv w_{ji}$ , and  $R(\mathbf{y})$  corresponds to the fuzzy ART box  $R_j$ . The code  $\mathbf{y}$  will be reset if  $|R(\mathbf{y}) \oplus \mathbf{a}| > M(1-\rho)$ . If  $|R(\mathbf{y}) \oplus \mathbf{a}| \leq M(1-\rho)$ ,  $R(\mathbf{y})$  expands toward  $R(\mathbf{y}) \oplus \mathbf{a}$  during distributed outstar learning, as thresholds  $\tau_{ji}$  increase and  $\sigma_i(\mathbf{y})$  converges toward  $\sigma_i(\mathbf{y}) \wedge A_i$ .



**Figure 7:** Distributed ART activation and learning in response to an input  $\mathbf{a} = (0.7, 0.8)$ , with complement coding, a power law CAM rule, a choice-by-difference signal rule, and two coding nodes ( $N=2$ ). (a) When  $p=1$  and  $\alpha=0.2$ ,  $(y_1, y_2) = (0.532, 0.468)$ . During learning,  $R_2(y_2)$  expands to include  $\mathbf{a}$  as  $\tau_{42}$  increases from 0.2 to 0.268. If  $\mathbf{a}$  is repeatedly presented and no other learned changes take place,  $\tau_{42}$  will continue to increase toward 0.274, the point where  $R_2(y_2)$  would just include  $\mathbf{a}$ . (b) When  $p=5$  and  $\alpha=0.2$ ,  $(y_1, y_2) = (0.653, 0.347)$  and no changes occur during learning. (c) When  $p=5$  and  $\alpha=0.8$ ,  $(y_1, y_2) = (0.906, 0.094)$ . During learning,  $R_1(y_1)$  expands to include  $\mathbf{a}$  as  $\tau_{31}$  increases from 0.4 to 0.606. (d) If  $\mathbf{a}$  is presented again,  $(y_1, y_2) = (0.916, 0.084)$ . If  $\mathbf{a}$  is repeatedly presented and no other learned changes take place,  $\tau_{31}$  will continue to increase toward 0.616.



**Figure 8:** Distributed ART search in response to an input  $\mathbf{a} = (0.7, 0.8)$ , with complement coding, a power law CAM rule ( $p = 1$ ) for above-average  $T_j(1)$ , a choice-by-difference signal rule ( $\alpha = 0.2$ ), and  $N = 3$ . Initially,  $T_1(1) = 3.20$  and  $T_2(1) = 2.82$ . When  $T_3(1) \leq 2.44$ ,  $T_3(1) < \bar{T} \leq T_2(1) < T_1(1)$ , so  $y_1 = 0.532$  and  $y_2 = 0.468$ , as in Figure 7a; and  $y_3 = 0$ . For this code  $\mathbf{y}$ ,  $T_1(y_1) = 1.78$ ,  $T_2(y_2) = 1.48$ , and  $T_3(y_3) = 0$ . A reset would therefore leave  $T_3(1)$  unchanged but would reduce  $T_1(1)$  to 1.42 and  $T_2(1)$  to 1.34. The next code  $\mathbf{y}$  then depends on the size of  $T_3(1)$  (Table 2).



**Figure 9:** (a) Distributed ARTMAP substitutes dART modules for the ART modules in ARTMAP, and substitutes distributed outstar learning from  $ART_a$  to the map field  $F^{ab}$  for outstar learning. (b) A simplified dARTMAP network computes classification probabilities, with  $|b|=1$  at an output field  $F_0^b$ .