



## CONTRIBUTED ARTICLE

# On the Match Tracking Anomaly of the ARTMAP Neural Network

GUSZTI BARTFAI

Victoria University of Wellington

(Received 24 March 1994; revised and accepted 10 May 1995)

**Abstract**—This article analyses the match tracking anomaly (MTA) of the ARTMAP neural network. The anomaly arises when an input pattern exactly matches its category prototype that the network has previously learned, and the network generates a prediction (through a previously learned associative link) that contradicts the output category that was selected upon presentation of the corresponding target output. Carpenter *et al.* claimed that such an anomalous situation will never arise if the (binary) input vectors have the same number of 1s (Carpenter *et al.*, 1991, *Neural Networks*, 4, 565–588).

This paper shows that such situations can in fact occur. The timing according to which inputs are presented to the network in each learning trial is crucial: if the target output is presented to the network before the corresponding input pattern, certain pattern sequences will lead the network to the MTA. Two kinds of MTA are distinguished: one that is independent of the choice parameter ( $\beta$ ) of the ART<sub>b</sub> module, and another that is not. Results of experiments that were carried out on a machine learning database demonstrate the existence of the match tracking anomaly as well as support the analytical results presented here.

**Keywords**—Match tracking, ARTMAP, Adaptive resonance theory, Supervised learning, Self-organization, Hierarchical clustering, Machine learning, Zoo database.

## 1. INTRODUCTION

The ARTMAP neural network architecture (Carpenter *et al.*, 1991a) is a supervised learning system capable of self-organising stable recognition categories in response to arbitrary sequences of input patterns. It is built up of a pair of adaptive resonance theory (Carpenter & Grossberg, 1987a) modules (ART<sub>a</sub> and ART<sub>b</sub>) that are connected through an inter-ART associative memory module. The network has an internal mechanism that conjointly maximises predictive generalisation and minimises predictive error through a self-regulating process called *match*

*tracking*. Whenever the network makes a wrong prediction through a previously learned associative link, the *vigilance* parameter  $\rho_a$  of the ART<sub>a</sub> module will be raised by the minimal amount needed to correct the predictive error at the ART<sub>b</sub> module. The ART<sub>a</sub> module will then start searching for another category for the current input until it finds a correct prediction, or creates a new ART<sub>a</sub> category and its associative link to the corresponding ART<sub>b</sub> category. The ARTMAP architecture has been applied to several machine learning benchmark problems (Carpenter *et al.*, 1991a, 1992), and compared favourably with traditional AI learning methods.

There are, however, anomalous situations that can occur during training, in which an ART<sub>a</sub> input matches its category prototype perfectly and the network makes a wrong prediction (Carpenter *et al.*, 1991a). Raising ART<sub>a</sub> vigilance above that (perfect) matching level will prevent the network from finding any other ART<sub>a</sub> category, and the wrong prediction will not be corrected. Carpenter *et al.* claimed that this situation (we shall call it *match tracking anomaly*, or MTA henceforth) will never arise if the (binary) input patterns have the same number of 1s (Carpenter *et al.*, 1991a, pp. 584–585).

---

Acknowledgements: This work was in part supported by Victoria University of Wellington (Research Grant No. RGNT/594/371). I am grateful to Peter Andrae, Department of Computer Science, Victoria University of Wellington, New Zealand, for his critical reading of the drafts, valuable suggestions for improving the style and presentation of this article, and even providing further insights into the similarities and differences of the two MTA conditions discussed here. I wish to thank the reviewers for their constructive comments, too.

Requests for reprints should be sent to G. Bartfai, Department of Computer Science, Victoria University of Wellington, P.O. Box 600, Wellington, New Zealand.

In this article, we show that the MTA *can* in fact arise even if the inputs have the same number of 1s and the training set is non-contradictory. More specifically, we discuss the importance of the timing according to which the input and target patterns are presented to the network in each learning trial. We first show (in Section 4.1) that the MTA will *never* arise if each input (with the same number of 1s) is presented to the network such that it can make a prediction (through an existing associative link) and “prime” itself before the corresponding (non-contradictory) target is presented. This timing proves to be crucial in that if the ART<sub>b</sub> (target) input is presented *before* the ART<sub>a</sub> input, the MTA *can* arise (see Section 4.2). We show this by specifying a set of initial conditions for the network and a sequence of input–target pairs that will lead the network to the MTA. Furthermore, we show that there are two distinct classes of conditions that cause the MTA: one that depends solely on the sequence in which inputs are presented from a training set (see the proof in Section 4.2.1), and another, in which MTA can only occur if the *choice parameter* ( $\beta$ ) of the ART<sub>b</sub> module is above a certain threshold, which depends on the size of the target inputs (see the proof in Section 4.2.2). The combined effect of these two conditions can be observed clearly from the results of the experiments we carried out on a benchmark machine learning database (in Section 5).

We have also extended the classes of problems to which the ARTMAP can be applied by allowing the ART<sub>b</sub> module, too, to develop “true” clusterings. Carpenter et al. introduced the ARTMAP network to be used for pattern classification where the ART<sub>b</sub> (target) inputs are expected *category codes* to which the corresponding inputs should belong. This implies that there is no overlap between any two target patterns. Therefore the ART<sub>b</sub> module will develop the same categories on a given sequence of training patterns at *any* non-zero vigilance level. This makes timing of input presentations irrelevant (see Section 4.3). This issue was, therefore, of no special concern in Carpenter et al. (1991a), and timing was only discussed in Appendix A.2 where all possible combinations of ART<sub>a</sub> and ART<sub>b</sub> input timings were listed with and without predictions. If, however, the ARTMAP network is applied to a more general “mapping” problem where target outputs may overlap, the level of ART<sub>b</sub> vigilance will affect the development of categories. One such problem is to learn a *two-level hierarchy* of input classes from a training set where each input pattern is identical to its target pair (i.e., auto-associative map). We discuss this in more detail in Section 5.

Section 2 introduces the ARTMAP network at a level that is necessary for understanding the main results of this paper. Section 3 defines the match

tracking anomaly (MTA)—which is the centre of this article. Section 4 discusses the circumstances under which the MTA condition cannot and can occur. Relevant theorems and proofs are given in Sections 4.1 and 4.2, respectively. A theorem regarding the important case of ART<sub>b</sub> vigilance equal to 1 is shown in Section 4.3. The experiments are discussed and the results are presented in Section 5, followed by a discussion in Section 6. Finally, conclusions are drawn in Section 7.

## 2. THE ARTMAP NETWORK

*Adaptive resonance theory* (ART) architectures are neural networks that develop stable recognition codes in real time by self-organisation, in response to arbitrary sequences of input patterns (Carpenter & Grossberg, 1987a). They were designed to solve the stability–plasticity dilemma that every intelligent machine learning system has to face: how to keep learning from new events without forgetting previously learned information. ART networks were originally proposed to accept binary input patterns, and were later extended for both continuous and binary inputs (Carpenter & Grossberg, 1987b; Carpenter et al., 1991b).<sup>1</sup>

An ART module has three layers: the *input* layer (F0), the *comparison* layer (F1), and the *recognition* layer (F2), with  $m$ ,  $m$  and  $n$  neurons, respectively (see module ART<sub>a</sub> or ART<sub>b</sub> in Figure 1). The neurons, or nodes, in the F2 layer represent input categories. The F1 and F2 layers interact with each other through weighted bottom-up and top-down connections, which are modified when the network learns. There are additional gain control signals in the network (not shown in Figure 1) that regulate its operation, but those will not be detailed here.

At each presentation of a non-zero binary input pattern  $\mathbf{x}$  ( $x_i \in \{0, 1\}$ ,  $i = 1, 2, \dots, m$ ), the network attempts to classify it into one of its existing categories based on its similarity to the stored prototype of each category node. More precisely, for each node  $j$  in the F2 layer, the bottom-up activation  $T_j = \sum_{i=1}^m x_i Z_{ij}$  is calculated, where  $Z_{ij}$  is the strength of the bottom-up connection between F1 node  $i$  and F2 node  $j$ . Since both the input and the bottom-up weight vectors are binary (with  $Z_{ij}$  being the normalised version of  $z_{ij}$ ),  $T_j$  can also be expressed as

$$T_j = |\mathbf{x} \cap \mathbf{Z}_j| = \frac{|\mathbf{x} \cap \mathbf{z}_j|}{\beta + |\mathbf{z}_j|}, \quad (1)$$

<sup>1</sup> From now on, we shall use ART to refer to the binary ART (called ART1) unless otherwise stated. Also, for easy reference, this paper attempts to use the same notation as Carpenter et al. (1991a), since their claim is the focus of this article.

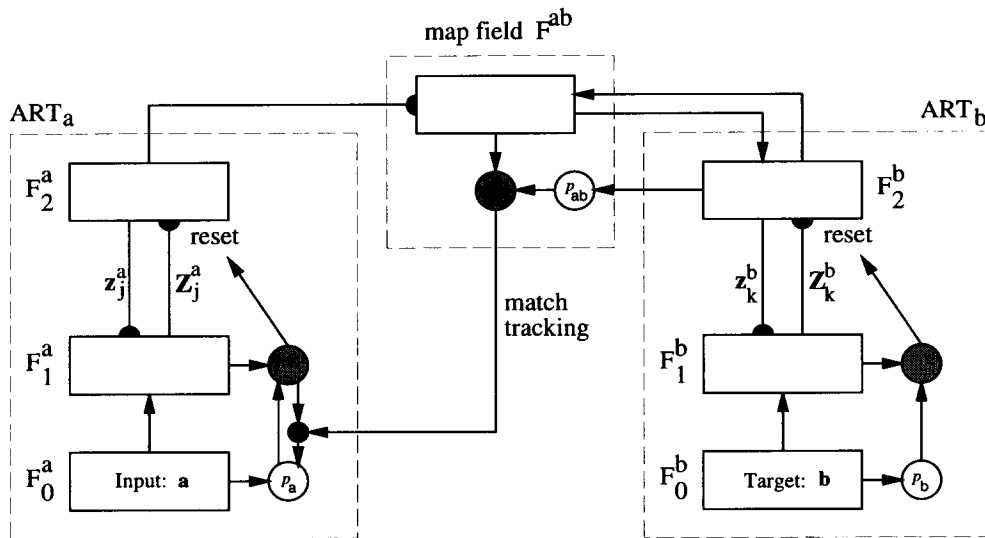


FIGURE 1. Architecture of the ARTMAP network. [This figure is a simplified version of the one that appeared in Carpenter et al. (1992).]

where  $|\cdot|$  is the norm operator ( $|\mathbf{x}| \equiv \sum_{i=1}^m x_i$ ),  $\mathbf{z}_j$  is the (binary) top-down template (or prototype) of category  $j$ , and  $\beta > 0$  is the *choice parameter* (Carpenter et al., 1991a). Then the F2 node  $J$  that has the highest bottom-up activation is selected, i.e.  $T_J = \max\{T_j | j = 1, 2, \dots, n\}$ . The prototype vector of the winning node  $J(\mathbf{z}_J; z_{Ji} \in \{0, 1\}, i = 1, 2, \dots, m)$  is then sent down to the F1 layer through the top-down connections, where it is compared to the current input pattern: the strength of the match is given by

$$\frac{|\mathbf{x} \cap \mathbf{z}_J|}{|\mathbf{x}|},$$

which is compared with a system parameter  $\rho$  called *vigilance* ( $0 < \rho \leq 1$ ). If the input matches sufficiently, i.e., the match strength  $\geq \rho$ , then it is assigned to F2 node  $J$  and both the bottom-up and top-down connections are adjusted for this node:

$$\mathbf{z}_J(\text{new}) = \frac{\mathbf{x} \cap \mathbf{z}_J(\text{old})}{\beta + |\mathbf{x} \cap \mathbf{z}_J(\text{old})|}$$

and

$$\mathbf{z}_J(\text{new}) = \mathbf{x} \cap \mathbf{z}_J(\text{old}).^2 \quad (2)$$

If the stored prototype  $\mathbf{z}_J$  does not match the input sufficiently (match strength  $< \rho$ ), the winning F2 node  $J$  is reset for the period of presentation of the current input. Then another F2 node (or category) will be selected, whose prototype will be matched

against the input. This “hypothesis-testing” cycle is repeated until the network either finds a stored category whose prototype matches the input closely enough, or allocates a new F2 node. Then learning takes place as described above. After an initial period of self-stabilisation, the network will directly (i.e., without search) access the prototype of one of the categories it has found in a given training set. The higher the vigilance level, the larger number of smaller, or more specific, categories will be created. (If  $\rho = 1$ , the network will learn every unique input perfectly with a different category.)

An important feature of the category selection method of the ART network is that if categories  $J_1$  and  $J_2$  are proper subsets of an input pattern  $\mathbf{x}$ , i.e.,  $|\mathbf{x} \cap \mathbf{z}_{J_1}| = |\mathbf{z}_{J_1}|$  and  $|\mathbf{x} \cap \mathbf{z}_{J_2}| = |\mathbf{z}_{J_2}|$ , the category whose prototype is the *larger* subset will be selected first (for the vigilance test). This can be seen from (1), and that the  $f(y) = \frac{y}{\beta + y}$  function monotonically increases for  $\beta > 0$ . This feature will play an important role in proving one of the MTA conditions (see Section 4.2.1).

The architecture of the ARTMAP network can be seen in Figure 1. It consists of two ART modules that are linked together through an “inter-ART” associative memory, called *map field* ( $F^{ab}$ ). Module ART<sub>a</sub> (with a *baseline* vigilance  $\bar{\rho}_a$ ) learns to categorise input patterns presented at layer  $F_0^a$ , while module ART<sub>b</sub> (with vigilance  $\rho_b$ ) develops categories of target patterns presented at layer  $F_0^b$ . Modules  $F_2^a$  and  $F^{ab}$  are fully connected via associative links whose strengths are adjusted through learning. There are one-to-one, two-way, and non-modifiable connections between nodes in the  $F^{ab}$  and  $F_2^b$  layers, i.e., each  $F_2^b$  node is connected to its corresponding  $F^{ab}$  node, and vice versa. A new association between an ART<sub>a</sub> category  $J$  and an ART<sub>b</sub> category  $K$  is learned

<sup>2</sup> This is the so-called fast learning mode of the network, which is typically used in binary ART.

by setting the corresponding  $F_2^a \rightarrow F^{ab}$  link to one and all other links from the same  $ART_a$  node to zero. When an input pattern is presented to the network, the  $F^{ab}$  layer will receive inputs from both the  $ART_a$  module (through the previously learned  $J \rightarrow K$  associative link) and the  $ART_b$  module (from the active  $F_2^b$  category node). If the two  $F^{ab}$  inputs match, i.e., the network's prediction is confirmed by the selected target category, the network will learn by modifying the prototypes of the chosen  $ART_a$  and  $ART_b$  categories according to the ART learning equations shown above. If there is a mismatch at the  $F^{ab}$  layer, a map field reset signal will be generated, and a process called *match tracking* will start, whereby the baseline vigilance level of the  $ART_a$  module will be raised by the minimal amount needed to cause mismatch with the current  $ART_a$  input at the  $F_1^a$  layer. This will subsequently trigger a search for another  $ART_a$  category, whose prediction will be matched against the current  $ART_b$  category at the  $F^{ab}$  layer again. This process continues until the network either finds an  $ART_a$  category that predicts the category of the current target correctly, or creates a new  $F_2^a$  node and a corresponding link in the map field, which will learn the current input/target pair correctly. The  $ART_a$  vigilance is then allowed to return to its resting level ( $\bar{\rho}_a$ ).

After a few presentations of the entire training set, the network will self-stabilise, and will read out the expected output for each input without search.

### 3. THE MATCH TRACKING ANOMALY (MTA)

The purpose of this article is to analyse the behaviour of the match tracking anomaly (MTA), which was introduced in Carpenter et al. (1991a, pp. 584–585) as follows:

In particular, if  $\mathbf{a} \subseteq \mathbf{z}_j^a$  [ $\mathbf{a}$  is an  $ART_a$  input here], match tracking makes  $\rho_a > 1$ , so  $\mathbf{a}$  cannot activate another category in order to learn the new prediction. The following anomalous case can thus arise. Suppose that  $\mathbf{a} = \mathbf{z}_j^a$  but the  $ART_b$  input  $\mathbf{b}$  mismatches the  $ART_b$  expectation  $\mathbf{z}_K^b$  previously associated with  $J$ . Then match tracking will prevent the recoding that would have associated  $\mathbf{a}$  with  $\mathbf{b}$ . That is, the ARTMAP system with fast learning and choice will not learn the prediction of an exemplar that *exactly* matches a learned prototype when the new prediction contradicts the previous predictions of the exemplars that created the prototype. This situation does not arise when all  $ART_a$  inputs  $\mathbf{a}$  have the same number of 1s [...]

(The same number of 1s in the input can be guaranteed by representing input patterns in *complement coding* (Carpenter et al., 1991b), according to which both input  $\mathbf{a}$  and its complement  $\mathbf{a}^c$  (where

every bit in  $\mathbf{a}$  is inverted) are shown to the network. Thus the norm of the new input vector  $|\langle \mathbf{a}, \mathbf{a}^c \rangle|$  will be constant.)

It obviously seems true that the anomalous situation does not arise since the norm of an  $ART_a$  category prototype  $\mathbf{z}_j^a$  can only be reduced through learning due to (2). Once this reduction occurs, no input  $\mathbf{a}$  can be a proper subset of category prototype  $\mathbf{z}_j^a$  (i.e.,  $|\mathbf{a} \cap \mathbf{z}_j^a| = |\mathbf{a}|$ ) if the inputs have the same norm (i.e., number of 1s). The problem can occur only when an input exactly matches a learned prototype that has not been recoded, which then produces a wrong prediction. In the classification problems in Carpenter et al. (1991a), where target patterns were non-overlapping category codes, it would mean that the network had previously learned a different target category for the same input. So the match tracking anomaly could only arise if there were contradictory target categories for the same input in the training set.

In the following, we show that in more general learning tasks, in which target outputs are arbitrary binary vectors (not just simple category codes), the MTA *can* in fact arise even if the inputs to the network are complement coded, and the training set is non-contradictory, i.e., no two identical inputs have different target values.

## 4. CONDITIONS FOR THE MATCH TRACKING ANOMALY

In this section, we present conditions under which the match tracking anomaly can and cannot occur. It will be shown that the timing of the presentation of the  $ART_a$  inputs and  $ART_b$  target outputs as well as the  $ART_b$  vigilance level are of particular importance.

### 4.1. Timing Conditions under which MTA will *never* Arise

The following theorem states the condition under which the match tracking anomaly will *not* occur.

**THEOREM 1.** *The match tracking anomaly will never arise in the ARTMAP network when complement coded input/target pairs from a non-contradicting training set are presented to the network such that in each learning trial the  $ART_a$  input is presented first, and the network is allowed to make a prediction through the  $F^{ab}$  map field before the  $ART_b$  (target) input is presented.*

If the ARTMAP network is able to make a prediction before the target is presented, the  $ART_b$  module will be “primed” by the  $F_2^b$  layer, which will send down its category template to the  $F_1^b$  layer via the top-down pathways. We assume that the network dynamics are such that this top-down template will be

tested first at the  $F_1^b$  layer as the target vector is registered into the  $F_0^b$  layer. If this top-down template (i.e. the network's prediction) and the target are sufficiently close to each other, no reset wave will be triggered, and the current  $F_2^b$  selection will be confirmed, and subsequently learned.

*Proof.* Let us assume that the network has previously learned to associate ART<sub>a</sub> category  $J$  with ART<sub>b</sub> category  $K$  through presentation of the  $(\mathbf{a}^{(i)}, \mathbf{b}^{(i)})$  input/target pair. Let us assume furthermore that input pattern  $\mathbf{a}^{(i)}$  is coded perfectly by ART<sub>a</sub> category  $J$ , i.e.,  $\mathbf{a}^{(i)} = \mathbf{z}_J^a$ , which is the necessary condition for the MTA to arise if the inputs are complement coded (since  $\mathbf{a}^{(i)} \subseteq \mathbf{z}_J^a$  can only be true if  $\mathbf{a}^{(i)} = \mathbf{z}_J^a$ ). This can always be guaranteed by setting  $\rho_a$  to 1 so only perfect matches will be accepted at the vigilance test.

We now show that the MTA will *not* arise when the input/target pair  $(\mathbf{a}^{(i)}, \mathbf{b}^{(i)})$  is presented again at some later stage during learning. Since there is no other input/target pair involving  $\mathbf{a}^{(i)}$  for which the MTA condition *could* arise (assuming a non-contradictory training set), we shall conclude that the MTA will *never* arise under the circumstances stated in this theorem.

When the input/target pair  $(\mathbf{a}^{(i)}, \mathbf{b}^{(i)})$  is presented again, the MTA condition can arise only if the network makes a wrong prediction first. Since the input set is non-contradictory, this can only be possible if the ART<sub>b</sub> module selects a category other than  $K$ . For that to happen, however, first the ART<sub>b</sub> vigilance test at the  $F_1^b$  layer has to fail with the network's prediction  $\mathbf{z}_K^b$ , i.e.

$$\frac{|\mathbf{b}^{(i)} \cap \mathbf{z}_K^b|}{|\mathbf{b}^{(i)}|} < \rho_b \quad (3)$$

must hold.

In complement coding, the norm of the target patterns is constant (here  $M_b$ ), therefore

$$|\mathbf{z}_K^b| \geq \rho_b M_b \quad (4)$$

for all  $F_2^b$  category prototypes, and it follows from (3) and (4) that

$$|\mathbf{b}^{(i)} \cap \mathbf{z}_K^b| < |\mathbf{z}_K^b| \quad (5)$$

should also be true. However,  $\mathbf{z}_K^b$  is already a subset of  $\mathbf{b}^{(i)}$  since  $F_2^b$  node  $K$  has already coded input  $\mathbf{b}^{(i)}$  before (hence the prediction), and therefore

$$|\mathbf{b}^{(i)} \cap \mathbf{z}_K^b| = |\mathbf{z}_K^b| \quad (6)$$

which directly contradicts (5). Therefore it is impossible that the (first) vigilance test will fail with the current prediction  $\mathbf{z}_K^b$  when target  $\mathbf{b}^{(i)}$  is re-

presented at any later stage. Consequently, the MTA condition will never arise.  $\square$

#### 4.2. Timing Conditions under which MTA *can* Arise

The following sections analyse conditions under which MTA *can* occur. In particular, we show that the MTA can occur when a target pattern is presented to the network *before* the input (i.e., the network is not able to make a prediction), even if all inputs to the network are complement coded. We also distinguish between two MTA conditions: one that is independent of the ART<sub>b</sub> choice parameter ( $\beta_b$ ), and another that is not.

Since in real-time situations, for which the ART systems were designed, there is no guarantee that ART<sub>a</sub> inputs will always arrive first, this problem is a valid one that needs to be investigated.

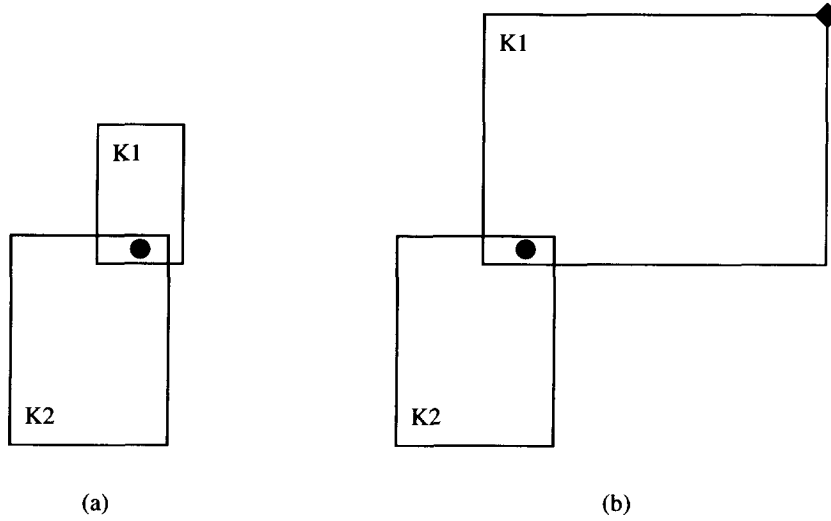
**THEOREM 2.** *The match tracking anomaly can arise in the ARTMAP network when complement coded input/target pairs from a non-contradicting training set are presented to the network such that some ART<sub>b</sub> (target) inputs are presented before their corresponding ART<sub>a</sub> input pair. There are at least two distinct conditions under which the MTA can occur: one that is independent of the ART<sub>b</sub> choice parameter ( $\beta_b$ ), and another that is not.*

We denote the “ $\beta_b$ -independent” MTA condition by MTA-S (for “MTA subset”, as explained below), and the “ $\beta_b$ -dependent” one by MTA- $\beta$ . We shall still use MTA, however, to refer to the match tracking anomaly in general, and shall be more specific only when necessary.

We shall prove this theorem by specifying initial conditions and sequences of training patterns that will lead to the MTA in each of the two cases mentioned in the theorem. We note, however, that there might exist other conditions as well that are different from these two, although the results of experiments we carried out do not suggest this.

Although the two conditions (MTA-S and MTA- $\beta$ ) are different in that one does not depend on  $\beta_b$  and the other does, they have much in common. In both cases,

- the network has previously seen the  $(\mathbf{a}^{(i1)}, \mathbf{b}^{(i1)})$  input/target pair, which causes the MTA,
- the network has already an associative link between ART<sub>a</sub> category  $J1$  and ART<sub>b</sub> category  $K1$ ,
- category  $J1$  perfectly encodes input  $\mathbf{a}^{(i1)}$  (i.e.,  $\mathbf{z}_{J1}^a = \mathbf{a}^{(i1)}$ ),
- category  $K1$  has already been recoded, i.e.,  $|\mathbf{z}_{K1}^b| < M_b$  (where  $M_b$  is the constant size of target patterns), and also includes target  $\mathbf{b}^{(i1)}$  (i.e.,  $\mathbf{z}_{K1}^b \subset \mathbf{b}^{(i1)}$ ),



**FIGURE 2.** Graphical illustration of the MTA-S condition. Target patterns are represented as points on the two-dimensional plane. Categories are represented as rectangles, the area of each rectangle being proportional to the size of the category they represent. Figure 2a shows that category  $K1$  is initially smaller (i.e., more specific) than  $K2$ , and target  $\mathbf{b}^{(i1)}$  (circle) is inside (i.e., proper subset of) both  $K1$  and  $K2$ . If at some later stage, another target input is presented (diamond, in Figure 2b) that causes  $K1$  to generalise to include that input,  $K1$  could grow larger than  $K2$  as a result. When target  $\mathbf{b}^{(i1)}$  is re-presented, this time category  $K2$  will capture it since it has now become more specific than  $K1$ . This will cause the anomaly to emerge as the network still predicts  $K1$  through its previously learned  $J1 \rightarrow K1$  link. See detailed discussion in Section 4.2.1.

- there exists another  $\text{ART}_b$  category ( $K2$ ) “nearby”, which is more specific than  $K1$  (i.e.,  $|\mathbf{z}_{K2}^b| > |\mathbf{z}_{K1}^b|$ ),
- $\text{ART}_a$  vigilance parameter  $\rho_a$  is assumed to be high enough so no other  $\mathbf{a}^{(j)} \neq \mathbf{a}^{(i1)}$  inputs from the training set will be coded by category  $J1$  (this limit depends on the particular training set, but the condition can always be guaranteed by  $\rho_a = 1$ ),
- $\text{ART}_b$  vigilance parameter  $\rho_b$  is assumed to be low enough that target  $\mathbf{b}^{(i1)}$  will be found sufficiently close to both  $\text{ART}_b$  category  $K1$  and  $K2$  (i.e.,  $|\mathbf{b}^{(i1)} \cap \mathbf{z}_{K1}^b| > \rho_b M_b$  and  $|\mathbf{b}^{(i1)} \cap \mathbf{z}_{K2}^b| > \rho_b M_b$ ). Note that this also implies that  $\rho_b$  must be lower than 1.

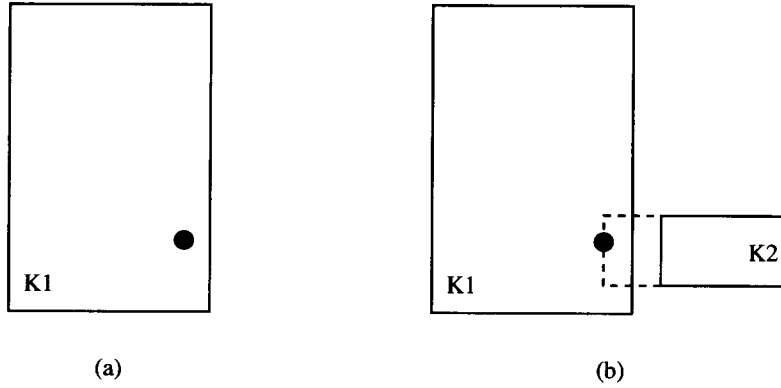
In both cases, the MTA will arise when the  $(\mathbf{a}^{(i1)}, \mathbf{b}^{(i1)})$  input/target pair is presented again, and  $\text{ART}_b$  category  $K2$  “captures” target input  $\mathbf{b}^{(i1)}$ , which contradicts the network’s prediction of  $\text{ART}_b$  category  $K1$  (through its previously learned  $J1 \rightarrow K1$  association). Since input  $\mathbf{a}^{(i1)}$  matches  $\text{ART}_a$  category  $J1$  perfectly, the network will not be able to find a better match through match tracking, therefore the MTA will emerge.

The difference between the two conditions is that in the MTA-S case both the  $K1$  and  $K2$   $\text{ART}_b$  category prototypes ( $\mathbf{z}_{K1}^b$  and  $\mathbf{z}_{K2}^b$ ) are proper subsets of target  $\mathbf{b}^{(i1)}$  (hence the name MTA-S), i.e.,  $|\mathbf{b}^{(i1)} \cap \mathbf{z}_{K1}^b| = |\mathbf{z}_{K1}^b|$  and  $|\mathbf{b}^{(i1)} \cap \mathbf{z}_{K2}^b| = |\mathbf{z}_{K2}^b|$ , therefore “capturing” of  $\mathbf{b}^{(i1)}$  by  $K2$  will not depend on  $\beta_b$  (see the paragraph on the ART network’s category selection in case of proper subsets in Section 2). In the MTA- $\beta$  case,  $\text{ART}_b$  category  $K2$  prototype  $\mathbf{z}_{K2}^b$  is not

a proper subset of target  $\mathbf{b}^{(i1)}$  (while  $K1$  prototype  $\mathbf{z}_{K1}^b$  still is), so the selection of  $\text{ART}_b$  category  $K1$  or  $K2$  will depend on the strength of the bottom-up match between  $\mathbf{z}_{K2}^b$  and target  $\mathbf{b}^{(i1)}$ , which in turn depends on the value of  $\beta_b$  due to (1). Figures 2 and 3 illustrate graphically the MTA-S and MTA- $\beta$  conditions, respectively.

We prove Theorem 2 by showing the existence of MTA in either of the two conditions mentioned in the theorem (defined as MTA-S and MTA- $\beta$  above). We specify initial conditions and training sequences in either case that will lead to the MTA. We must point out, however, that these training sequences are artificial (and simplified) examples that are only used to prove the *existence* of MTA. In real applications, the necessary initial conditions (as a result of prior learning) as well as the “critical training sequence” can be observed through a variety of actual training sequences. Moreover, the two distinct MTA conditions (MTA-S and MTA- $\beta$ ) can even have a combined effect if  $\beta_b$  is above a threshold, which will produce an even richer variety of paths leading to the MTA. These were clearly observed in the experiments we carried out, the results of which are presented in Section 5.

**4.2.1. Proof of the MTA-S Condition.** In this section, we show how the ARTMAP network can be led to the MTA-S condition. First we shall assume that the network is in a particular state after prior learning, and show that presenting the input/target pair sequence  $(\mathbf{a}^{(i1)}, \mathbf{b}^{(i1)})$ ,  $(\mathbf{a}^{(i2)}, \mathbf{b}^{(i2)})$ ,  $(\mathbf{a}^{(i1)}, \mathbf{b}^{(i1)})$  then



**FIGURE 3.** Graphical illustration of the MTA- $\beta$  condition. Target patterns are represented as points on the two-dimensional plane. Categories are represented as rectangles, the area of each rectangle being proportional to the size of the category they represent. Figure 3a shows category  $K1$  and target input  $\mathbf{b}^{(i1)}$  (circle) inside it. As a result of further training, some other target(s) have created category  $K2$  “near”  $K1$  (see Figure 3b).  $K2$  is more specific than  $K1$ , and initially does not include target  $\mathbf{b}^{(i1)}$ , although it could (with low enough  $\rho_b$  vigilance). When target  $\mathbf{b}^{(i1)}$  is re-presented, it will be captured by category  $K2$  if it achieves a higher bottom-up activation than  $K1$ . This depends on the degree of overlap between  $\mathbf{b}^{(i1)}$  and  $\mathbf{z}_{K2}^b$ , and the choice parameter  $\beta_b$  due to (1). See detailed discussion in Section 4.2.2.

will cause the MTA to arise. Then we show that the particular (assumed) state of the network before this sequence *can* result from prior learning.

Let us assume that the network has previously learned to associate ART<sub>a</sub> category  $J1$  with ART<sub>b</sub> category  $K1$ . Also, ART<sub>a</sub> category  $J1$  has perfectly coded input  $\mathbf{a}^{(i1)}$ , and there is no other input in that category, i.e.,

$$\mathbf{z}_{J1}^a = \mathbf{a}^{(i1)}. \quad (7)$$

There is also another input ( $\mathbf{a}^{(i2)}$ ), for which

$$|\mathbf{a}^{(i2)} \cap \mathbf{z}_{J1}^a| < \rho_a M_a \quad (8)$$

where  $M_a$  is the norm of the input patterns.

Furthermore, target inputs  $\mathbf{b}^{(i1)}$  and  $\mathbf{b}^{(i2)}$ , and another ART<sub>b</sub> category,  $K2$ , satisfy the following conditions:

$$\mathbf{b}^{(i1)} \cap \mathbf{z}_{K1}^b = \mathbf{z}_{K1}^b, \quad \mathbf{b}^{(i1)} \cap \mathbf{z}_{K2}^b = \mathbf{z}_{K2}^b, \quad (9)$$

i.e., both category prototypes are proper subsets of target  $\mathbf{b}^{(i1)}$ , and

$$|\mathbf{b}^{(i2)} \cap \mathbf{z}_{K1}^b| > |\mathbf{b}^{(i2)} \cap \mathbf{z}_{K2}^b|, \quad (10)$$

i.e., target  $\mathbf{b}^{(i2)}$  matches  $K1$  category prototype  $\mathbf{z}_{K1}^b$  better than  $\mathbf{z}_{K2}^b$ , and

$$|\mathbf{z}_{K1}^b| > |\mathbf{z}_{K2}^b| > |\mathbf{b}^{(i2)} \cap \mathbf{z}_{K1}^b| \geq \rho_b M_b, \quad (11)$$

i.e., category  $K1$  is more specific than  $K2$ , and matches target  $\mathbf{b}^{(i2)}$  well enough, but not as well as category prototype  $\mathbf{z}_{K2}^b$ .

If we now present the input/target pair ( $\mathbf{a}^{(i1)}$ ,  $\mathbf{b}^{(i1)}$ )

first, ART<sub>b</sub> category  $K1$  will be selected due to (9) and (11), which will be confirmed by the network through the selection of category  $J1$  upon presenting input  $\mathbf{a}^{(i1)}$  due to (7), and the existing  $J1 \rightarrow K1$  association.<sup>3</sup> Due to (7) and (9), category prototypes  $\mathbf{z}_{J1}^a$  and  $\mathbf{z}_{K1}^b$  will not change during learning.

Next, input/target pair ( $\mathbf{a}^{(i2)}$ ,  $\mathbf{b}^{(i2)}$ ) is presented. Here, target input  $\mathbf{b}^{(i2)}$  will select ART<sub>b</sub> category  $K1$  again due to (10) and (11). The network will then either find an existing ART<sub>a</sub> category associated with ART<sub>b</sub> category  $K1$ , or will create a new ART<sub>a</sub> category to learn this association. In either case, ART<sub>a</sub> category  $J1$  will not be affected due to (8). ART<sub>b</sub> category  $K1$ , however, will be modified this time, i.e.,

$$\mathbf{z}_{K1}^b(\text{new}) = \mathbf{b}^{(i2)} \cap \mathbf{z}_{K1}^b(\text{old}).$$

Then input/target pair ( $\mathbf{a}^{(i1)}$ ,  $\mathbf{b}^{(i1)}$ ) is presented again. Here, target input  $\mathbf{b}^{(i1)}$  will read out ART<sub>b</sub> category  $K2$  this time due to (9) and (11). However,  $\mathbf{a}^{(i1)}$  input will again read out ART<sub>a</sub> category  $J1$  since it still matches  $\mathbf{z}_{J1}^a$  perfectly. As a result the  $J1 \rightarrow K1$  prediction will not be confirmed at the  $F^{ab}$  layer, and the match tracking process will be triggered. This will, however, cause MTA since  $\mathbf{a}^{(i1)}$  matches  $\mathbf{z}_{J1}^a$  perfectly (i.e.,  $|\mathbf{a}^{(i1)} \cap \mathbf{z}_{J1}^a| \|\mathbf{a}^{(i1)}\|^{-1} = 1$ ), and  $\rho_a$  cannot be raised above 1. Table 1 illustrates the above sequence and its effect on the network.

To complete the proof, we also need to show that such initial assumption for the ART<sub>b</sub> category prototypes  $\mathbf{z}_{K1}^b$  and  $\mathbf{z}_{K2}^b$  and target pairs  $\mathbf{b}^{(i1)}$  and

<sup>3</sup> Alternatively, ART<sub>a</sub> category  $J1$  can be a newly created category, which will learn to associate ART<sub>b</sub> category  $K1$  during this trial, but it does not affect the proof.

**TABLE 1**  
**A Sequence of Input/Target Patterns that Causes the MTA under Conditions Discussed in the Text**

Step	Input $\mathbf{a}^{(j)}$	ART <sub>a</sub> Cat. No.	$F_2^a \rightarrow F_2^b$ Mapping		ART <sub>b</sub> Categories	Target $\mathbf{b}^{(j)}$
...	...	...	...		$\mathbf{z}_{K1}^b$	...
$k$	$\mathbf{a}^{(i1)}$	$J1$	$J \rightarrow K1$	$\triangleright$	$\mathbf{z}_{K1}^b$	$\mathbf{b}^{(i1)}$
$k+1$	$\mathbf{a}^{(i2)}$	$\neq J1$	$\neq J1 \rightarrow K1$	?	$\mathbf{b}^{(i2)} \cap \mathbf{z}_{K1}^b$	$\mathbf{b}^{(i2)}$
$k+2$	$\mathbf{a}^{(i1)}$	$J1$	$J \rightarrow 1?2$	$\triangleright$	$\mathbf{b}^{(i2)} \cap \mathbf{z}_{K1}^b$	$\mathbf{b}^{(i1)}$

The ART<sub>a</sub> categories are now shown here, only the index of the winning node in each step. Two of the developed ART<sub>b</sub> categories, which are referred to in the text, are shown. The symbol " $\triangleright$ " appears beside the ART<sub>b</sub> category that the network predicts through the  $F_2^a \rightarrow F_2^b$  associations in each trial ("?" means it does not matter which category is predicted). The symbol " $\triangleleft$ " denotes the category that is selected by the ART<sub>b</sub> module before the corresponding input is presented. So in step  $k$ , the network's prediction is confirmed (i.e., " $\triangleright$ " and " $\triangleleft$ " appear at the same category), while there is a mismatch in step  $k+2$ , which eventually causes the MTA.

$\mathbf{b}^{(i2)}$  is "realistic", i.e., they can result from previous training of the network.

It is, however, not trivial to see. We cannot do it simply by first creating the proper  $J1 \rightarrow K1$  association and  $\mathbf{z}_{J1}^a$ ,  $\mathbf{z}_{K1}^b$  prototypes, and then another ART<sub>b</sub> category ( $K2$  with a prototype  $\mathbf{z}_{K2}^b$ ), which is also a proper subset of another input [here  $\mathbf{b}^{(i1)}$ , and see (9) as the assumption] such that  $|\mathbf{z}_{K1}^b| > |\mathbf{z}_{K2}^b|$ , as was assumed in (11). The network would not have selected category  $K2$  when  $K1$  is clearly a better choice. This situation can only come about if the required  $J1 \rightarrow K1$  association and  $\mathbf{z}_{J1}^a$ ,  $\mathbf{z}_{K1}^b$  prototypes are created first, and then a particular sequence of input/target patterns is presented, through which a new ART<sub>b</sub> prototype ( $\mathbf{z}_{K2}^b$ ) is first created, which is then gradually "eroded" (i.e., its norm decreased) in small enough pieces until condition (11) is reached, while leaving  $\mathbf{z}_{K1}^b$  unchanged in each trial. For simplicity, here we assume the smallest possible size for a "piece" by which  $|\mathbf{z}_{K2}^b|$  is decreased, which is one.

More precisely, we present a sequence of input/target pairs ( $\mathbf{a}^{(j)}$ ,  $\mathbf{b}^{(j)}$ ,  $j = 1, 2, \dots$ ) such that each input and the category prototypes in question will satisfy

$$|\mathbf{b}^{(j)} \cap \mathbf{z}_{K1}^b| = 1, \quad (12)$$

and

$$|\mathbf{z}_{K2}^b(\text{new})| = |\mathbf{b}^{(j)} \cap \mathbf{z}_{K2}^b(\text{old})| = |\mathbf{z}_{K2}^b(\text{old})| - 1 \quad (13)$$

until  $|\mathbf{z}_{K1}^b| > |\mathbf{z}_{K2}^b|$  (see condition (11)).

Now we need to show that in each step

$$\frac{1}{\beta + |\mathbf{z}_{K1}^b|} < \frac{|\mathbf{z}_{K2}^b| - 1}{\beta + |\mathbf{z}_{K2}^b|}. \quad (14)$$

In other words, ART<sub>b</sub> category  $K2$  will be selected (over  $K1$ ) at the  $F_2^b$  layer and allowed to change, while  $K1$  remains unchanged.

This can be shown by considering the following two cases:

Case 1.  $2 < |\mathbf{z}_{K2}^b| \leq |\mathbf{z}_{K1}^b|$

Here, (14) will hold since the numerator of the right hand side is greater than 1, and the denominator is smaller than that of the left hand side.

Case 2.  $2 < |\mathbf{z}_{K1}^b| < |\mathbf{z}_{K2}^b|$

To see that (14) will hold in this case as well, we need to see first that

$$\frac{1}{\beta + |\mathbf{z}_{K1}^b|} < \frac{1}{\beta + |\mathbf{z}_{K1}^b|} + \frac{|\mathbf{z}_{K1}^b| - 2}{\beta + |\mathbf{z}_{K1}^b|} = \frac{|\mathbf{z}_{K1}^b| - 1}{\beta + |\mathbf{z}_{K1}^b|}. \quad (15)$$

Now (14) holds again, since the

$$f(x) = \frac{x - 1}{\beta + x}$$

function increases monotonically when  $x > 2$  (here  $|\mathbf{z}_{K1}^b| > 2$ ) with arbitrary  $\beta > 0$ .

Therefore, the assumptions about the network's initial state [see inequalities (9) to (11)] were realistic. This completes the proof of the MTA-S condition in theorem 2.  $\square$

A concrete example with a complete sequence leading to the MTA-S is given in Table 2.

**4.2.2. Proof of the MTA- $\beta$  Condition.** This section shows that there are situations where the MTA can occur if  $\beta_b$  is above a certain threshold.<sup>4</sup> As in the previous section, we provide initial conditions and a training sequence that will lead to the MTA- $\beta$ .

Let us take three input/target pairs ( $\mathbf{a}^{(i1)}$ ,  $\mathbf{b}^{(i1)}$ ), ( $\mathbf{a}^{(i2)}$ ,  $\mathbf{b}^{(i2)}$ ) and ( $\mathbf{a}^{(i3)}$ ,  $\mathbf{b}^{(i3)}$ ). The patterns are complement coded (as was assumed), thus the norm of the input and target vectors ( $\mathbf{a}^{(i)}$  and  $\mathbf{b}^{(i)}$ ) is constant ( $M_a$  and  $M_b$ , respectively).

<sup>4</sup> For clarity, we shall refer to  $\beta_b$  as  $\beta$  from here on unless the distinction from  $\beta_a$  needs to be made clear.



TABLE 2  
An Example Sequence of Input/Target Pairs that Lead the ARTMAP Network to the MTA-S

Step	Input $\mathbf{a}^{(i)}$	ART <sub>a</sub> Cat. No.	$F_2^a \rightarrow F_2^b$ Mapping	ART <sub>b</sub> Categories		Target $\mathbf{b}^{(i)}$
				$\mathbf{z}_1^b$	$\mathbf{z}_2^b$	
1	1 1 1 1 1	1	1 → 1	1 1 1 1 1	Unused	1 1 1 1 1
2	1 1 1 0 0	2	2 → 1	1 1 1 x x	Unused	1 1 1 0 0
3	0 0 0 1 1	3	3 → 2	1 1 1 x x	0 0 0 1 1	0 0 0 1 1
4	1 0 0 1 1	4	4 → 2	1 1 1 x x	x 0 0 1 1	1 0 0 1 1
5	0 1 0 1 1	5	5 → 2	1 1 1 x x	x x 0 1 1	0 1 0 1 1
6	0 0 1 1 1	6	6 → 2	1 1 1 x x	x x x 1 1	0 0 1 1 1
7	1 1 1 1 1	1	1 → 1	▷ 1 1 1 x x	x x x 1 1	1 1 1 1 1
8	1 0 0 0 0	7	7 → 1	1 x x x x	x x x 1 1	1 0 0 0 0
9	1 1 1 1 1	1	1 → 1?2	▷ 1 x x x x	x x x 1 1	1 1 1 1 1

Here, each input pattern is identical to its corresponding target (i.e., auto-associative learning), and the norm of all patterns is 5. For clarity, the complement patterns are not shown. In category prototypes, a "x" means "don't care", which corresponds to a "00" bit pair in complement coding. For the ART<sub>a</sub> module, only the winning nodes are shown. We assume that the ART<sub>a</sub> vigilance is high enough (e.g.,  $\rho_a = 0.9$ ) so the 11111 pattern will be the only one in its own category. ART<sub>b</sub> vigilance is less than 1/5 (e.g.,  $\rho_b = 0.1$ ). The use of the symbols "▷" and "◁" is the same as in Table 1. (The absence of the "▷" symbol in rows 1–6 indicates that the network did not have a prediction in these trials, and the corresponding  $J \rightarrow K$  links were created in those trials.) The  $F_2^a \rightarrow F_2^b$  mappings that the network created or used in each trial are also shown. Steps 1–6 show the "erosion" process of ART<sub>b</sub> category  $\mathbf{z}_2^b$ , while steps 7 and 8 "prepare" the MTA condition that occurs in step 9. (Here, concrete category numbers are used, since these are the only patterns presented to the network).

First, let  $\rho_a$  be high enough to cause ART<sub>a</sub> to allocate an uncommitted  $F_2^a$  node for each new, previously unseen pattern. This can always be guaranteed if  $\rho_a = 1$ .

The target patterns satisfy the following conditions:

$$|\mathbf{b}^{(i1)} \cap \mathbf{b}^{(i2)}| = L_{12} \quad (L_{12} > 0), \quad (16)$$

$$|\mathbf{b}^{(i1)} \cap \mathbf{b}^{(i3)}| = L_{13} \quad (L_{13} > 0), \quad (17)$$

such that

$$L_{13} > L_{12}, \quad (18)$$

and

$$|\mathbf{b}^{(i2)} \cap \mathbf{b}^{(i3)}| = 0. \quad (19)$$

From (16), (17) and (19), it follows that<sup>5</sup>

$$L_{12} + L_{13} \leq M. \quad (20)$$

Also, let

$$\rho_b < \frac{L_{12}}{M}. \quad (21)$$

We now show that under these conditions, presenting the above patterns in the order of  $(\mathbf{a}^{(i1)}, \mathbf{b}^{(i1)})$ ,  $(\mathbf{a}^{(i2)}, \mathbf{b}^{(i2)})$ ,  $(\mathbf{a}^{(i3)}, \mathbf{b}^{(i3)})$ ,  $(\mathbf{a}^{(i1)}, \mathbf{b}^{(i1)})$  will cause MTA to emerge if  $\beta$  is above a certain threshold.

In this case, we start with an untrained network, and look at the behaviour of the ART<sub>b</sub> module. (We also assume that category nodes become committed in the order of 1, 2, 3, ..., and thus these numbers will be used instead of  $J1$ ,  $K1$ , etc., used in the previous sections.)

In the first step, when the  $(\mathbf{a}^{(i1)}, \mathbf{b}^{(i1)})$  pair is presented,  $\mathbf{b}^{(i1)}$  will be learned by category 1, and the corresponding  $1 \rightarrow 1$  link will also be created (since input  $\mathbf{a}^{(i1)}$  will be learned by ART<sub>a</sub> category 1).

In the second step, target  $\mathbf{b}^{(i2)}$  will be captured by category 1, as it will achieve the maximal bottom-up activation (no other  $F_2^b$  nodes are active as yet, and the initial bottom-up weights are chosen to be small enough for the uncommitted nodes), and

$$|\mathbf{b}^{(i2)} \cap \mathbf{z}_1^b| = L_{12} > \rho_b M,$$

due to (21). Prototype  $\mathbf{z}_1^b$  is then modified to learn  $\mathbf{b}^{(i1)} \cap \mathbf{b}^{(i2)}$ .

In the third step, target  $\mathbf{b}^{(i3)}$  will create a new ART<sub>b</sub> category  $\mathbf{z}_2^b$  as there will be a total mismatch between  $\mathbf{b}^{(i3)}$  and  $\mathbf{z}_1^b$  due to (19).

The interesting case is step four, when pattern  $(\mathbf{a}^{(i1)}, \mathbf{b}^{(i1)})$  is presented again. Here both existing ART<sub>b</sub> categories will match the input sufficiently due to (18) and (21), so the selection will be based entirely on the bottom-up activation levels, which can be calculated in the following way.

The bottom-up weights for ART<sub>b</sub> categories 1 ( $\mathbf{z}_1^b$ ) and 2 ( $\mathbf{z}_2^b$ ) will be

$$\frac{1}{\beta + L_{12}} \text{ and } \frac{1}{\beta + M},$$

<sup>5</sup> For clarity, we refer to  $M_b$  simply as  $M$  since  $M_a$  is irrelevant to the proof here.

**TABLE 3**  
**The Sequence of Input/Target Pairs ( $\mathbf{a}^{(1)}, \mathbf{b}^{(1)}$ ), ( $\mathbf{a}^{(2)}, \mathbf{b}^{(2)}$ ), ( $\mathbf{a}^{(3)}, \mathbf{b}^{(3)}$ ), ( $\mathbf{a}^{(1)}, \mathbf{b}^{(1)}$ ), and the Developed ARTMAP Categories and  $F_2^a \rightarrow F_2^b$  Associations**

Step	Input	ART <sub>a</sub> Categories			$F_2^a \rightarrow F_2^b$ Mapping	ART <sub>b</sub> Categories		Target
		$\mathbf{z}_1^a$	$\mathbf{z}_2^a$	$\mathbf{z}_3^a$		$\mathbf{z}_1^b$	$\mathbf{z}_2^b$	
1	$\mathbf{a}^{(1)}$	$\mathbf{a}^{(1)}$	Unused	Unused	1 → 1	$\mathbf{b}^{(1)}$	Unused	$\mathbf{b}^{(1)}$
2	$\mathbf{a}^{(2)}$	$\mathbf{a}^{(1)}$	$\mathbf{a}^{(2)}$	Unused	2 → 1	$\mathbf{b}^{(1)} \cap \mathbf{b}^{(2)}$	Unused	$\mathbf{b}^{(2)}$
3	$\mathbf{a}^{(3)}$	$\mathbf{a}^{(1)}$	$\mathbf{a}^{(2)}$	$\mathbf{a}^{(3)}$	3 → 2	$\mathbf{b}^{(1)} \cap \mathbf{b}^{(2)}$	$\mathbf{b}^{(3)}$	$\mathbf{b}^{(3)}$
4	$\mathbf{a}^{(1)}$	$\mathbf{a}^{(1)}$	$\mathbf{a}^{(2)}$	$\mathbf{a}^{(3)}$	1 → 1?2	$\mathbf{b}^{(1)} \cap \mathbf{b}^{(2)}$	$\mathbf{b}^{(3)}$	$\mathbf{b}^{(1)}$

Since  $\rho_a$  is high enough, each new input creates a new ART<sub>a</sub> category in steps 1, 2 and 3. The winning ART<sub>b</sub> categories in each step are the ones that appear on the right-hand side of the arrow in the " $F_2^a \rightarrow F_2^b$  mapping" column. The re-presentation of the first pair in step 4 will cause the MTA if the bottom-up activation of ART<sub>b</sub> node 2 is higher than that of node 1. (Here, concrete category numbers are used, since these are the only patterns presented to the network).

respectively. The bottom-up activations will therefore be

$$\frac{L_{12}}{\beta + L_{12}} \text{ and } \frac{L_{13}}{\beta + M},$$

respectively, due to (1). For fixed  $L_{12}$ ,  $L_{13}$  and  $M$ ,  $\beta$  can always be set such that the  $\mathbf{b}^{(i)}$  target will select either category 1 or 2. This can be seen by examining the

$$\frac{L_{12}}{\beta + L_{12}} < \frac{L_{13}}{\beta + M} \tag{22}$$

inequality. If (22) is satisfied, then ART<sub>b</sub> category 2 will be selected, and thus the MTA will emerge.

Rearranging (22), we get

$$\beta > \frac{L_{12}(M - L_{13})}{L_{13} - L_{12}}. \tag{23}$$

For fixed  $L_{12}$ ,  $L_{13}$  and  $M$  that satisfies (20), the occurrence of the MTA condition will therefore be determined by the choice of  $\beta$ . Since the recommended setting of  $\beta$  is to be small (Carpenter et al., 1991a, p. 579), it is worthwhile looking at how the lower limit of  $\beta$  in (23) is affected by the choice of  $L_{12}$ ,  $L_{13}$  and  $M$  (which are determined by the problem at hand).

To get the minimum value of the right-hand side of (23), we need to minimise  $L_{12}(M - L_{13})$ , and maximise  $L_{13} - L_{12}$ . The minimum values for  $L_{12}$  and  $M - L_{13}$  are 1 and  $L_{12}$  [from (20)], respectively. The resulting choice for  $L_{13}(= M - L_{12})$  will maximise  $L_{13} - L_{12}$ . So the lower limit for  $\beta$  is

$$\beta_{\min} = \frac{1}{M - 2}. \tag{24}$$

Therefore, for a given  $M_b$ ,  $\beta_b$  should always be chosen to be smaller than  $\frac{1}{M_b - 2}$  to avoid the MTA- $\beta$  for patterns that satisfy conditions (16)–(21).

Table 3 shows the developed ART<sub>a</sub> and ART<sub>b</sub> categories as well as  $F_2^a \rightarrow F_2^b$  links in response to the input/target pairs discussed here.

This completes the proof of the MTA- $\beta$  condition in theorem 2 as well as theorem 2 as a whole.  $\square$

A concrete example is given in Table 4, with  $M = 4$ ,  $L_{13} = 3$ , and  $L_{12} = 1$ . Here, if  $\beta > \frac{1}{2}$ , the MTA condition will arise.

### 4.3. ART<sub>b</sub> vigilance level and MTA

This section deals with an important special case of the ARTMAP network where the ART<sub>b</sub> vigilance level ( $\rho_b$ ) is 1. This setting is commonly used in

**TABLE 4**  
**An Example Input/Target Pattern Sequence that Demonstrates the MTA- $\beta$  Condition**

Step	Input	ART <sub>a</sub> Categories			$F_2^a \rightarrow F_2^b$ Mapping	ART <sub>b</sub> Categories		Target
		$\mathbf{z}_1^a$	$\mathbf{z}_2^a$	$\mathbf{z}_3^a$		$\mathbf{z}_1^b$	$\mathbf{z}_2^b$	
1	1 1 1 1	1 1 1 1	Unused	Unused	1 → 1	1 1 1 1	Unused	1 1 1 1
2	1 0 0 0	1 1 1 1	1 0 0 0	Unused	2 → 1	1 x x x	Unused	1 0 0 0
3	0 1 1 1	1 1 1 1	1 0 0 0	0 1 1 1	3 → 2	1 x x x	0 1 1 1	0 1 1 1
4	1 1 1 1	1 1 1 1	1 0 0 0	0 1 1 1	1 → 1?2	1 x x x	0 1 1 1	1 1 1 1

Here, each input pattern is identical to its corresponding target, and the norm of all patterns is 4. The ARTMAP categories and links that were developed in response to this sequence are also shown. For clarity, the complement patterns are not shown here. In category prototypes, a "x" means "don't care", which corresponds to a "00" pair in complement coding. The ART<sub>a</sub> and ART<sub>b</sub> vigilance parameters ( $\rho_a$ ,  $\rho_b$ ) are set to 0.8 and 0.3, respectively. If  $\beta > \frac{1}{2}$ , the MTA will occur in step 4 when input/target pair (1111, 1111) is re-presented, because ART<sub>b</sub> category 2 will be selected. This contradicts the network's prediction of category 1 through the link that was created in step 1. The MTA will thus emerge.

pattern classification applications where outputs represent target classes. The following theorem reassures us that the MTA will *not* arise when  $\rho_b = 1$ .

**THEOREM 3.** *The match tracking anomaly will never arise in the ARTMAP network if  $\rho_b = 1$  regardless of the timing used to present complement coded input/target pairs from a non-contradicting training set.*

*Proof.* If  $\rho_b = 1$ , then no ART<sub>b</sub> category will be recoded during training, i.e.,  $|z_k^b| = M_b$  for all ART<sub>b</sub> category prototypes. Therefore, upon presentation of input/target pair  $(\mathbf{a}^{(i)}, \mathbf{b}^{(i)})$  that the network has seen before, either a previously learned  $J1 \rightarrow K$  association will be confirmed, in which case  $\mathbf{z}_{j1}^a = \mathbf{a}^{(i)}$ , or a new  $J2 \rightarrow K$  association will be created through match tracking, in which case  $|z_{j2}^a| < M_a$  and therefore raising  $\rho_a$  will not cause MTA. In neither of these cases will ART<sub>b</sub> category  $K$  be recoded so the MTA will not arise.  $\square$

## 5. EXPERIMENTAL RESULTS

In the previous sections, we proved the existence of the match tracking anomaly by specifying initial conditions and constructing training sequences that lead to the MTA. These were relatively simple and “clean” examples that let us concentrate on the underlying problem. In reality, however, there are a number of different ways the network can be driven into the MTA, which are dependent on the training set and the presentation order of input/target pairs during training. For example, before any step of the above example training sequences, there can be a number of input/target presentations that are irrelevant to the occurrence of a given MTA. This section presents experimental results that will demonstrate the existence of the two MTA conditions analysed here.

We carried out experiments on a machine learning benchmark database. Unlike in Carpenter et al. (1991a), where target outputs were codes for correct output categories, here the ARTMAP network was used to develop two-level class hierarchies from the training set (Bartfai, 1994). In particular, each input was identical to its target pair, and both the ART<sub>a</sub> and ART<sub>b</sub> modules developed clusterings of patterns from the same training set at different levels of vigilance.

For the experiments, we used the “zoo” machine learning benchmark database (Murphy & Aha, 1992). It contains 101 instances of animals described with 18 attributes. Out of these attributes, we used the 15 boolean ones that indicate the presence or absence of certain features like “hair”, “aquatic”, “domestic”, and so on. We also used the “number of legs” attribute, which is a set of six integers. The “type” (or

class) attribute was ignored.<sup>6</sup> The patterns were presented to the ARTMAP network in *complement coding*. However, complement coding was defined for binary patterns only in Carpenter et al. (1991a). Since one of the features here (see “legs” above) is not binary, we had to encode that feature as well in a way that was compatible with complement coding. For this, we have defined a coding scheme that we call *generalised complement coding*, which can be applied to arbitrary-sized sets.

According to generalised complement coding, a value  $v_i$  of a set  $\mathcal{S} = \{v_1, \dots, v_s\}$  will be mapped to a binary vector  $(\mathbf{c})$  of length  $s$  such that

$$c_j = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{otherwise} \end{cases} \quad j = 1, \dots, s.$$

So, for example, if an attribute takes on the value 4 from the set of  $\{0, 2, 4, 5, 6, 8\}$ , then the corresponding binary vector will be 001000. Note that if  $s = 2$ , this scheme will be equivalent to the original form of complement coding.

This way, we can keep the size of the (binary) input pattern to the network constant, thereby implementing a form of “binary normalisation” of the input. In the case of the database used in the above way, the input and target patterns were 36-element binary vectors, in which the number of 1s was the same (16).

The simulations were carried out using a public domain neural network simulator program (PlaNet, see Miyata, 1991) running under Unix and X-windows.

The network parameters were chosen to be similar to those presented in Carpenter et al. (1991a). In particular, the initial values of the bottom-up weights in both ART modules were chosen such that  $F_2$  nodes become active in the order  $j = 1, 2, \dots$  and were small enough so the network selected an uncommitted ART<sub>a</sub> or ART<sub>b</sub> node only if  $|\mathbf{a}^{(i)} \cap \mathbf{z}_j^a| = 0$  or  $|\mathbf{b}^{(i)} \cap \mathbf{z}_k^b| = 0$ , respectively, for all committed nodes. Also, the  $\beta_a$  and  $\beta_b$  parameters were sufficiently small that, among committed nodes, bottom-up activations

$$\frac{|\mathbf{a}^{(i)} \cap \mathbf{z}_j^a|}{\beta_a + |\mathbf{z}_j^a|} \text{ and } \frac{|\mathbf{b}^{(i)} \cap \mathbf{z}_k^b|}{\beta_b + |\mathbf{z}_k^b|}$$

were determined by the sizes of  $|\mathbf{a}^{(i)} \cap \mathbf{z}_j^a|$  and  $|\mathbf{b}^{(i)} \cap \mathbf{z}_k^b|$  relative to  $|\mathbf{z}_j^a|$  and  $|\mathbf{z}_k^b|$ , respectively. This also satisfied the  $\beta_b < \beta_{\min}$  condition in Section 4.2.2. The network was used in *fast learning* mode, i.e., during learning, the connection weights were allowed

<sup>6</sup> The 18th attribute (“animal name”) was simply used as a label for the individual instances.

**TABLE 5**  
**Frequencies of MTA Occurrence at Different Levels of the  $\beta$  Choice Parameter when the ARTMAP Network was Trained on the "Zoo" Database (see details in text)**

$\beta$	Total Epochs	MTA in epochs		MTA in trials	
		Total	%	Total	%
0.001	215	30	14.0	78	0.4
0.01	215	30	14.0	78	0.4
0.1	215	30	14.0	78	0.4
1.0	252	49	19.4	119	0.5
10.0	370	228	61.6	1270	3.4
100.0	491	250	50.9	1662	3.4

Each row shows the results of the same set of 100 training sessions. In each training session, the network was trained completely on the training set, i.e., training stopped when no further weight changes occurred in the network for a full training cycle (or epoch). The  $ART_a$  and  $ART_b$  vigilance parameters ( $\rho_a$ ,  $\rho_b$ ) were 0.4 and 0.1, respectively. The "total epochs" column shows the total number of presentations of the entire training set during the 100 training sessions. The "MTA in epochs" columns show the number of epochs in which MTA occurred at least once (both in absolute and percentage). The "MTA in trials" columns show the total number of learning trials (i.e., input/target presentations) in which MTA occurred during the 100 training sessions.

to reach their asymptotic values while the current input/target pair was presented.

Training proceeded by showing input/target pairs repeatedly, until the network stabilised itself, i.e., each input read out its prototype directly, so no further learning took place.

Table 5 shows the results of measurements we carried out into the frequency of MTA when the network was trained on the "zoo" database. It can be seen that with small  $\beta$  values, the MTA condition occurred in about 14% of the training epochs (i.e., one full presentation of the entire training set). This accounts for only 0.4% of the total number of input/target presentations. It demonstrates that the MTA condition can arise even when  $\beta$  is small.<sup>7</sup> These are the MTA-S cases that were discussed in Section 4.2.1. However, we can see a sharp rise in the frequency of MTA as  $\beta$  increases beyond a certain threshold. This demonstrates that increasing  $\beta$  does affect the MTA (See MTA- $\beta$  discussed in Section 4.2.2). The fact that this threshold appears in the range of [0.1, 1], which is slightly higher than the limit

$$\frac{1}{M-2} = 14^{-1} \approx 0.07$$

that was derived in Section 4.2.2, suggests that MTA-

<sup>7</sup> We tried even smaller values for  $\beta$ . The numbers did not change until rounding took effect due to the finite accuracy of representable numbers in computer simulations. Very often, the two effects (MTA and rounding) cancelled each other out suggesting falsely that the MTA condition will disappear if  $\beta$  becomes very small. This may serve as a warning about the limits of computer simulations.

$\beta$  does not take immediate effect, probably because here it is superimposed on the MTA-S, or that the patterns did not satisfy all of the conditions from (16) to (21), (19) in particular. In the latter case, even the theoretical limit would possibly be lower than what was derived in Section 4.2.2.

We also observed that, on average, two out of every three MTA occurrences were due to race conditions at the bottom-up activations (i.e., two or more  $F_2^b$  nodes achieved the same maximal bottom-up activation for the current input). In the simulations, we broke these ties by declaring the node with the smallest index the winner. As a result, some true MTA conditions may have remained unnoticed as the above tie breaking method chose the correct category (whereas another one would not have). We believe, however, that this did not affect the results significantly. (Besides, the aim of these experiments were to demonstrate that these anomalous situations do exist, and not to derive empirical formulae on just how often they occur.)

We also carried out a separate set of experiments where the  $ART_a$  and  $ART_b$  vigilance levels were varied. The results for various combinations of  $\rho_a$  and  $\rho_b$  can be seen in Table 6.

First, it supports Theorem 3 since no MTA occurred when  $\rho_b$  was 1. Also, an interesting pattern can be observed in the table: the MTA will not occur if  $\rho_a \leq \rho_b$ . This is because in the ARTMAP network, the  $ART_b$  module will always "force" its categorisation onto  $ART_a$  through the internal feedback mechanism. In these "auto-associative" experiments, the ART modules were presented with identical input patterns. As a result, they developed identical categories for  $\rho_a \leq \rho_b$ . Since no recoding of

**TABLE 6**  
**MTA Occurrences at Different  $ART_a$  and  $ART_b$  Vigilance Levels ( $\rho_a, \rho_b \in \{0.1, 0.4, 0.7, 1\}$ )**

$\rho_b$	Vigilance		MTA % of	
	$\rho_a$		Epochs	Trials
0.1	0.1		0	0
	0.4		14.0	0.4
	0.7		36.8	1.0
	1.0		80.2	7.7
0.4	0.1		0	0
	0.4		0	0
	0.7		50.0	1.8
	1.0		94.9	9.9
0.7	0.1		0	0
	0.4		0	0
	0.7		0	0
	1.0		54.1	1.5
1.0	0.1		0	0
	0.4		0	0
	0.7		0	0
	1.0		0	0

The choice parameter ( $\beta$ ) was fixed at 0.001 here. All other settings are the same as in Table 5.

an  $ART_b$  category could occur without recoding its corresponding  $ART_a$  category too, the MTA did not arise in these cases. For the remaining combinations (i.e.,  $\rho_a > \rho_b$ ), we can see an increasing effect of  $\rho_a$  on the MTA frequencies. This is simply because as  $\rho_a$  increases, more and more inputs will be encoded perfectly by  $ART_a$  category nodes, which is a necessary condition for MTA (see Section 3).

Finally, if the timing of the input and target presentations were according to Theorem 1 (i.e.,  $ART_a$  input presented before  $ART_b$  target), no MTA occurrences were found in any of the training sessions we carried out. This supports the claim of Theorem 1, and further underlies the importance of the timing according to which inputs are presented to the ARTMAP network.

## 6. DISCUSSION

In the previous sections, we proved the existence of the match tracking anomaly in the ARTMAP network. The question now arises whether the match tracking anomaly can be eliminated even if the timing of input presentations is “undesired”. We note first that in pattern recognition problems, this can be done simply by setting  $\rho_b = 1$  (see Theorem 3 in Section 4.3). In general, however, the ARTMAP network needs some protection mechanism against MTA.

In the experiments, for example, where the ARTMAP network was used to develop two-level class hierarchies (Bartfai, 1994), and the target was presented before the input in each trial (i.e., the “undesired” timing), we modified the standard ARTMAP learning method as follows:

Whenever the MTA condition arose, the network found an uncommitted  $ART_a$  node and learned the input and its association with the current target (in the original algorithm no learning occurred in these cases). It also “undid” the existing (wrong) association by freeing up the corresponding  $ART_a$  node (i.e., making it uncommitted again) and resetting its  $F_2^a \rightarrow F^{ab}$  link. (This way, we ensured all  $ART_a$  category prototypes remained unique, and every input could access its category directly after self-stabilisation, both of which are important features of ART networks.) In MTA-free learning trials, the standard ARTMAP learning method was applied.

This modification was an adequate solution to eliminate the MTA condition since the training set was not contradictory here. In general, however, the occurrence of MTA may also be the result of a contradictory training set, in which case, learning the new association and forgetting the old one “blindly” may not be the best strategy. Therefore, in a real-

world, autonomous learning environment, the network should be equipped with additional mechanisms that ensure correct operation. In particular, we recommend that:

1. The choice parameter be chosen to be sufficiently small to avoid the “ $\beta$ -dependent” MTA condition (experimental results suggest that the limit may be lower than what was derived in Section 4.2.2);
2. the network be extended with an “internal self checking” process that is activated whenever the MTA occurs. This process should first disable further learning temporarily, and then let the  $ART_a$  module make an internal prediction (while keeping the target input registered in the  $F_0^b$  layer) to get the correct timing, and check for the MTA condition again. If MTA persists, it will be an indication of a contradicting input/target pair, in which case the network should not learn on the current inputs. If MTA disappears, which indicates that the target input was presented first and that was the cause of MTA, the network should learn the current input/target pair with the correct (“ $ART_a$  first, then  $ART_b$ ”) timing. It could be the subject of further research to validate the correctness of this procedure, and to see if it can be implemented as an extension to the current match tracking process.

We also note that although the “anomalous situation” that can arise in the ARTMAP network was reported *in relation to* the match tracking process (Carpenter et al., 1991a), the problem itself lies *not* with the match tracking process *itself*. Whatever method is used to correct a wrong prediction, if the  $ART_a$  category *matches the input perfectly*, there is no way the network can find another category that will capture the same input (thereby produce the correct output) on future trials. In fact, we could even replace the match tracking process with some other method that can correct a wrong prediction of the network, for instance the one that was used in Tan (1992), and would still face the same problem. So regardless of the mechanism used to correct a wrong prediction, it is important that the network be able to “prime” itself with its “expectation” before the target is presented, and subsequently learned.

## 7. CONCLUSION

This paper analysed the match training anomaly of the ARTMAP network. We showed that it can occur even when the inputs from a non-contradictory training set are presented to the network in complement coding. We conclude that the timing according to which the input and target patterns are presented to the network in each learning trial is critical if the  $ART_b$  vigilance level is less than 1.

To avoid the match tracking anomaly when input patterns are complement coded, the ART<sub>a</sub> input has to be presented first, letting the network “prime” the  $F_1^b$  layer by its own “expectation” (through previously learned  $F_2^a \rightarrow F_2^b$  associations), before the target input is presented to the ART<sub>b</sub> module.

If we cannot guarantee this timing, as in most real-world, autonomous learning tasks, we have to equip the network with additional mechanisms, like the one recommended in the previous section which ensure its correct operation.

We expect the MTA condition will also be present in other ARTMAP architectures, like the Fuzzy ARTMAP network (Carpenter et al., 1992), for example. It would be interesting to see how the theorems presented here should change so they incorporate continuous inputs as well as slow learning mode. This again can be the subject of future research.

#### REFERENCES

- Bartfai, G. (1994). Hierarchical clustering with ART neural networks. In *Proceedings of the IEEE International Conference on Neural Networks*, (Vol. 2, pp. 940–944). IEEE Press.
- Carpenter, G., & Grossberg, S. (1987a). A massively parallel architecture for a self-organizing neural pattern recognition machine. *Computer Vision, Graphics, and Image Processing*, 37, 54–115.
- Carpenter, G., & Grossberg, S. (1987b). ART2: self-organizing of stable category recognition codes for analog input patterns. *Applied Optics*, 26(23), 4919–4930.
- Carpenter, G., Grossberg, S., & Reynolds, J. (1991a). ARTMAP: Supervised real-time learning and classification of nonstationary data by a self-organizing neural network. *Neural Networks*, 4, 565–588.
- Carpenter, G., Grossberg, S., & Rosen, D. (1991b). Fuzzy ART: Fast stable learning and categorization of analog patterns by an adaptive resonance system. *Neural Networks*, 4, 759–771.
- Carpenter, G., Grossberg, S., Markuzon, N., Reynolds, J., & Rosen, D. (1992). Fuzzy ARTMAP: A neural network architecture for incremental supervised learning of analog multidimensional maps. *IEEE Transactions on Neural Networks*, 3(5), 698–713.
- Miyata, Y. (1991). *PlaNet, a Tool for Constructing, Running, and Looking into a PDP Network*.
- Murphy, P. M., & Aha, D. W. (1992). UCI repository of machine learning databases [Machine-readable data repository]. Technical report, Department of Information and Computer Science, University of California, Irvine, CA.
- Tan, A.-H. (1992). Adaptive resonance associative map: a hierarchical ART system for fast stable associative learning. In *Proceedings of IJCNN* (Vol. 1, pp. 860–865).