

# Understanding HTM Sequence Memory

Hugo Pristauz, Walter Eder – Oct. 2023 (Draft v0.1)

## Abstract

This paper targets to give aid for people who aim to jump start with Numenta's HTM (Hierarchical Temporal Memory) sequence memory algorithm. Sequence memory is capable of learning sequences, and to predict the next item, when part of a sequence is presented to the sequence memory.

The first three targets of this work are to setup some prerequisites: 1) to set up a formal framework HTM sequence memory and a detailed formulation of the HTM algorithm, which can be efficiently execute in Matlab or Python. 2) Harmonizing the algorithm with the standard representation of nonlinear discrete time state representation, commonly used in system and control theory. 3) Running a toy network comprising a total of 40 neurons with an investigation characteristic input sequences.

The main target is 4) a performance comparison of HTM sequence memory with transformer based neural network alternatives by predicting the continuation of text sequences from "Tiny Shakespeare"

## Introduction

In his revolutionary book *On Intelligence* [1] Jeff Hawkins, with the help of Sandra Blakeslee, describes "how a new understanding of the brain will lead to the creation of truly intelligent machines", suggesting that current "machine intelligence" is still far away from human intelligence. In his later book *A Thousand Brains* [2] Jeff Hawkins suggests a new *definition of intelligence* for beings (humans, animals, machines) on base of the capability to leveraging four abilities, which can be definitely found in the human brain:

- Continuous learning: the being must have the ability of continuous (on-line) learning, which is important for continuously adopting to a continuously changing world. This contrasts with current neural networks, which must be pre-trained before being utilized in an operational mode.
- Creating a model of the world by physical (or mental) movement through the world, with the ability to build an understanding of the world by sensing sequential information of the real word.
- Simultaneous creation of multiple models of the world in order to make accurate predictions. The human brain is mainly a prediction machine, which benefits from the possibility to select the best prediction model of a pool of simultaneously entertained prediction models.
- Utilization of Reference frames for any kind of knowledge storage. Reference frames can be dynamically assembled, which, in consequence allows to "assemble" knowledge which is stored with respect to reference frames.

Capability 2 of this list (world understanding based on sequential sensor information) emphasizes the importance of *dealing with sequential information* (with both spatial and temporal attributes). In this sense *sequence memory*, as it is under focus in this paper, might be a powerful building block for required skill sets, as described above. In addition, capability 3 to predict is not of less importance.

## Biological Background

HTM (Hierarchical Temporal Memory) Sequence memory is a computer model of a biological neuron, which is based on a modern neuro-physical understanding of cortical neuron functionality [3,4].

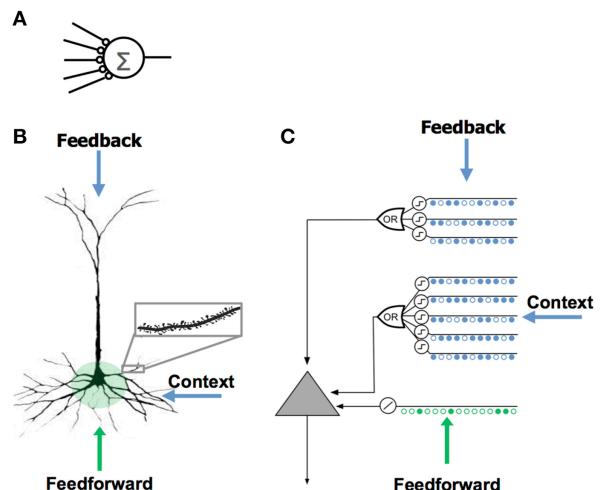


Figure 1: Comparison of neuron models

The neuron model used in most artificial neural networks is known as the perceptron [5], proposed by Frank Rosenblatt in 1957. It has few synapses and no modelling of dendrites (figure 1/A). The perceptron attempts to model the proximal area of the neuron and represents de-facto a mapping of input information arriving at the proximal synapses to the neuron's *output*, as the proximal synapses, those closest to the soma (cell body), have a large effect on the likelihood of a cell generating an *action potential*, which is bursting along the neuron's axon to synapses of other neurons of the network. It is important to realize that the perceptron model is a pure mapping model without any memory.

In contrast, an activation of a distal (non-proximal) synapse has little effect at the soma (cell body). For this reason, it was hard to understand how the thousands of distal synapses can play a role for the cell's responses [4].

Powered by advanced research methods [3], researchers could, however, unlock the secrets, that an activation of neighbored distal synapses within a short time interval leads to a local dendritic NMDA (N-methyl-D-aspartate) spike which causes a depolarization of the soma. Such depolarization subsequently enables the neuron to fire earlier than neighbor neurons with comparable proximal excitation, which gives such neuron the benefit to inhibit neighbored neurons.

Since NMDA spikes have a much longer duration than the action potential spikes of the soma, this mechanism acts like a memory functionality related to the depolarized state, in addition to the (perceptron like) mapping functionality of the neuron. Availability of memory in a neural model offers both the ability of *prediction*, and to implement *sequence state*, as will be shown in further sections.

Thus, compared to the perceptron model of figure 1/A, which generates the output signal according to a mapping of input signals at the proximal synapses (feedforward information) an extended neural model (like the HTM neuron model) deals additionally with context information arriving at distal synapses of several dendritic segments, which causes a state change of the neural model, influencing the neuron's behavior in the near future.

## Synaptic Model

In Neurobiology a synapse forms a connection between a neuron's dendrite and an axon of some other neuron.

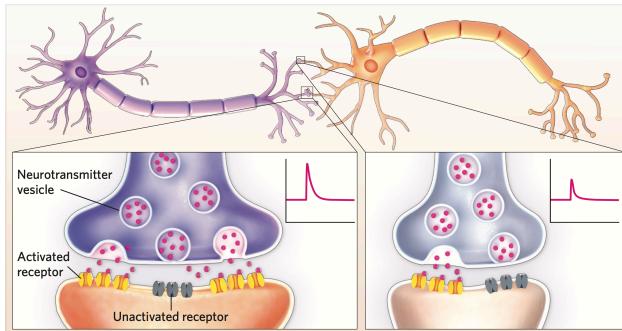


Figure 2: axon/dendrite connections via synapses

In formal neuron models synapses are usually assigned with a weight, influencing the transmission efficacy of the synapse, which is changed during the learning process. In the HTM approach there are two fundamental differences compared to the learning rules of most neural models.

- First, learning occurs by growing and removing synapses from a pool of "potential" synapses [7].

- Second, Hebbian learning and synaptic change occur at the level of a dendritic segment, not the entire Neuron [8].

Learning in an HTM neuron is modeled by the growth of new synapses from a set of potential synapses. A "permanence" value is assigned to each potential synapse and represents its growth. Learning occurs by incrementing or decrementing the permanence values. The synaptic weight is a binary value set to 1 if the permanence is above a threshold (figure 3).

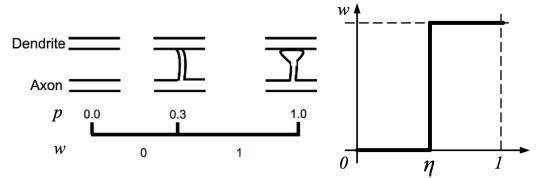


Figure 3: permanence  $p$  and weight of a synapse

Since the pre-synaptic signal is also binary, the synapse acts as a logical gate. For sufficient high permanence  $p \geq \eta$ , the synapse is called *connected* and allows a pre-synaptic signal to cause a post-synaptic effect ( $q=z$ ), while otherwise the synapse behaves as *unconnected* with no post-synaptic effect, regardless of the value of the pre-synaptic signal ( $q=0$ ).

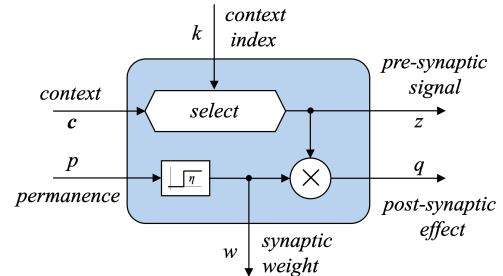


Figure 4: Block diagram for a single HTM neuron

This is shown in figure 4, where the (binary) pre-synaptic signal  $z$  is gated by the binary synaptic weight  $w = (p \geq \eta)$ , the Boolean result of the comparison of permanence  $p \in [0,1]$  with synaptic threshold  $\eta$  (a typical setting is  $\eta = 0.5$ ). The other block input is context  $c \in \{0,1\}^L$ , typically a high dimensional binary vector collecting neuron outputs from a large neuron context set. The *select* function block in figure 5 means to pick the  $k$ -th element of context  $c$  for the pre-synaptic signal  $z$ . In this sense index  $k$  contributes to the definition of the neural network topology. The formal model for an HTM-synapse is as follows:

$$w = (p \geq \eta) \quad w \in \{0,1\}, \text{synaptic weight} \quad (1a)$$

$$z = i_k^T \cdot c = c_k \quad z \in \{0,1\}, \text{pre-synaptic signal} \quad (1b)$$

$$q = w \cdot z \quad q \in \{0,1\}, \text{post-synaptic effect} \quad (1c)$$

with  $i_k^T$  denoting the  $k$ -th unit vector of  $\{0,1\}^L$ , and  $^T$  used for transpose.

## Synaptic Regions

So far (1a-1c) tells us whether a given context  $c$  presented to a synapse will cause a post-synaptic effect  $q$  with respect to the current synaptic permanence  $p$ , which is a sort of synaptic state.

A formal learning rule, however, is still missing, and we postulated for our model, that *learning* and *synaptic change* shall occur on the level of dendritic segments and not the entire neuron. This leads us to the concept of *synaptic regions* where groups of synapses collaborate with each other in order to generate some group effect.

As an example, the synapses of all proximal dendrites (receiving feedforward signals - see figure 1) collaborate with each other in order to cause occasionally the soma to generate an action potential. In this case we speak about the *proximal synaptic region*.

In contrast we have several *distal synaptic regions* formed by a collection of all synapses of distal dendritic segments. If those synapses generate sufficient collective post-synaptic effect, the distal dendritic segment can generate an NMDA spike, which can put the neuron into a so called *predictive state* by depolarization.

While we do not cover the proximal synaptic region in this paper, we will deal with the formal HTM-model for the distal synaptic regions related to distal dendritic segments. We will start with a generic formal model for a synaptic region (figure 6).

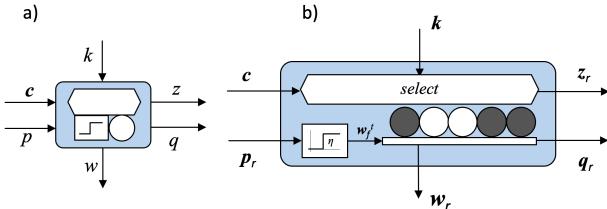


Figure 5: block symbol for a) single synapse, b) synaptic region

While in figure 5/a input/output signals of a single neuron's model are scalars (except context  $c$ ), we extend the single-synapse model in figure 5/b to model a synaptic region with a total number of  $s$  synapses. The individual synapses are still considered to work independent from each other according to (1a-1c), except the  $s$  instances of signals  $p$ ,  $w$ ,  $q$ ,  $z$ , and  $k$  are now grouped in terms of vectors  $\mathbf{p}_r$ ,  $\mathbf{w}_r$ ,  $\mathbf{q}_r$ ,  $\mathbf{z}_r$  augmented with index  $r$  to refer to synaptic region  $r \in \{1, \dots, d\}$ ,  $d$  denoting the total number of *synaptic regions*.

The formal model is straight forward,

$$\mathbf{w}_r = (\mathbf{p}_r \geq \eta) \quad \mathbf{w}_r \in \{0, 1\}^d, \text{ synaptic weights} \quad (1a)$$

$$\mathbf{z}_r = \mathbf{c}(\mathbf{k}_r) \quad \mathbf{z}_r \in \{0, 1\}^d, \text{ pre-synaptic signals} \quad (1b)$$

$$\mathbf{q}_r = \mathbf{w}_r \circ \mathbf{z}_r \quad \mathbf{q}_r \in \{0, 1\}^d, \text{ post-synaptic effects} \quad (1c)$$

with  $\circ$  denoting elementwise product, and  $\mathbf{c}(\mathbf{k}_r)$  meaning to pick  $d$  elements from  $\mathbf{c}$  relating to the indices given by  $\mathbf{k}_r$ .

## Mathematical Model of HTM Neuron

In HTM terminology a single HTM neuron is called a *cell*, which can be treated in the first approach as a black box with the input/output model shown in figure 3.

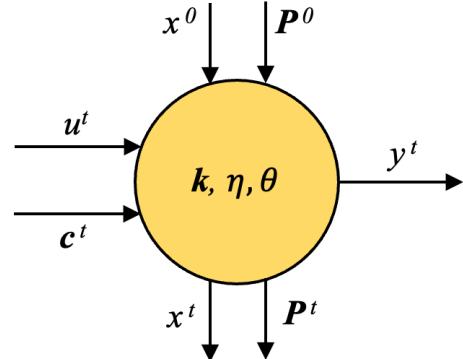


Figure 6: input/output model of a cell (HTM neuron)

All quantities shown in figure 3 are related to a discrete time index  $t$  ( $t = 0, 1, 2, \dots$ ). The HTM neuron presented in this paper has two kinds of inputs:

- the (scalar) binary *feedforward input*  $u^t \in \{0, 1\}$
- the (vectorial) binary *context input*  $\mathbf{c}^t \in \{0, 1\}^L$ .

Further has the HTM neuron a state comprising

- the (scalar) binary cell state  $x^t \in \{0, 1\}$
- the permanence matrix  $\mathbf{P}^t \in [0, 1]^{s \times d}$

with  $[0, 1]$  denoting the closed interval  $\{x \in \mathbb{R} \mid 0 \leq x \leq 1\}$ . Since the tuple  $(x^t, \mathbf{P}^t)$  represents a state of the cell, which will be sequentially updated at each time step  $t$ , we need to provide

- initial cell state  $x^0 \in \{0, 1\}$
- initial permanence matrix  $\mathbf{P}^0 \in [0, 1]^M$ ,

which is sketched in figure 3 by the two input signals at the top. Finally, we have

- the (scalar) binary cell output  $y^t \in \{0, 1\}$ .

In addition to inputs and outputs we have some cell parameters  $\mathbf{k} \in (1, 2, \dots, N)^M$  and threshold values  $\eta \in [0, 1]$  and  $\Theta \in \{0, 1, \dots, N\}$ , which will be explained in more detail later.

Depending on the binary values of their input, state and output we will call a cell in the following as

- *excited*, when its input is activated ( $u^t = 1$ )
- *predictive*, when its state is activated ( $x^t = 1$ )
- *active*, when its output is activated ( $y^t = 1$ )

## HTM Neuron Layer

To implement HTM sequence memory we need many *cells*, and we arrange them in terms of an  $mn$  matrix with  $n$  columns, each containing  $m$  cells (figure 4), resulting in a layer of total  $N := mn$  cells.

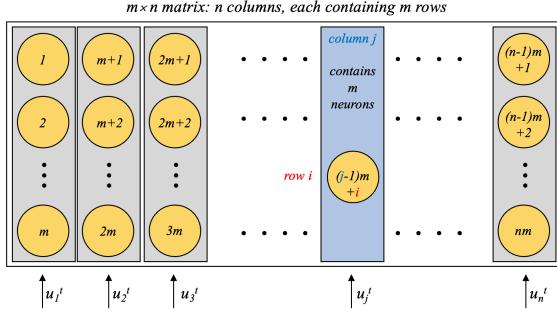


Figure 4: HTM layer with a matrix of  $m n$  neurons

In HTM terminology each column of the layer matrix is called a *minicolumn*, and each of them is attached with a scalar binary input  $u_j^t \in \{0,1\}$ . In total the HTM layer receives  $n$  scalar inputs  $u_j^t$ , which are combined as the (binary) *input vector*

$$\mathbf{u}^t = [u_1^t, u_2^t, \dots, u_n^t]^T \in \{0,1\}^n. \quad (1)$$

As figure 4 suggests, each cell of a particular minicolumn  $j$  receives the *same* binary input  $u_j^t$ . In a similar way we provide quantities  $x^t, y^t$  and  $\mathbf{d}^t$  of figure 3 with cell index  $k \in \{1, \dots, N\}$  and combine them to (binary) vectors  $\mathbf{x}^t, \mathbf{y}^t$  and (non-binary) matrix  $\mathbf{D}^t$ .

$$\mathbf{x}^t = [x_1^t, x_2^t, \dots, x_N^t]^T \in \{0,1\}^N \quad (2a)$$

$$\mathbf{y}^t = [y_1^t, y_2^t, \dots, y_N^t]^T \in \{0,1\}^N \quad (2b)$$

$$\mathbf{D}^t = [\mathbf{d}_1^t, \mathbf{d}_2^t, \dots, \mathbf{d}_N^t] \in [0,1]^{MXN} \quad (2c)$$

Note that the binary context input  $\mathbf{v}^t$  collects all cell outputs  $\mathbf{y}^t$  and is applied as a common context input to all cells (thus needs no augmentation). Formally we write  $\mathbf{v}^t = \mathbf{y}^t$ , which introduces a contextual feedback loop to the cell layer, as it is shown in figure 5.

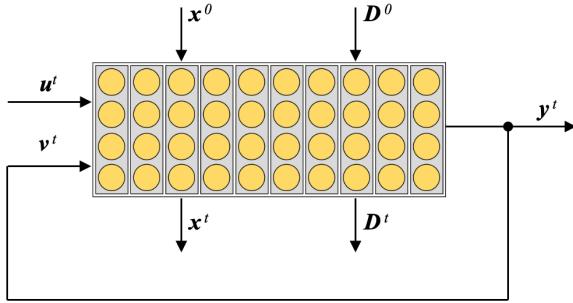


Figure 5: input/output model of cell layer with contextual feedback loop

## Input Sequences

Let us now talk in more detail about how input sequences are presented to the cell layer of figure 5. Let us consider two sentences “Mary likes to sing.” and “John likes to

dance.” To formally deal with these sequences, we define a vocabulary

$$V = \{\text{'Mary'}, \text{'John'}, \text{'likes'}, \text{'to'}, \text{'sing'}, \text{'dance'}, \text{'.}\}$$

and related sequences

$$\begin{aligned} S_1 &= (\text{'Mary'}, \text{'likes'}, \text{'to'}, \text{'sing'}, \text{'.}) \\ S_2 &= (\text{'John'}, \text{'likes'}, \text{'to'}, \text{'dance'}, \text{'.}). \end{aligned}$$

Since each input at time instance  $t$  to our cell layer must be a binary vector  $\mathbf{u}^t \in \{0,1\}^N$ , we must map each symbol of sequence  $S_1$  and  $S_2$  to the binary vector space  $\{0,1\}^n$ , and we call this process tokenizing, and talk about the binary vectors  $\mathbf{u}^0, \mathbf{u}^1, \mathbf{u}^2, \dots$  as *tokens*, which are sequentially presented as input vectors  $\mathbf{u}^t$  to the cell layer at time instances  $t = 0, 1, 2, \dots$ .

In this sense we can define a tokenizer as a function  $\Phi: V \rightarrow \{1,0\}^n$  mapping from a vocabulary  $V$  to the binary input space  $\{1,0\}^n$ .

## Example 1: 4 x 10 Toy Layer with 5 x 2 synapses per cell

Let us define a 410 cell layer ( $m=4, n=10$ ) where each cell supports 2 dendritic segments ( $d=2$ ) with 5 synapses per segment ( $s=5$ ). Thus, the layer has  $n = 10$  minicolumns with a total of  $N = m \cdot n = 40$  cells, where each cell entertains  $M = s \cdot d = 10$  synapses in total. Referring to vocabulary

$$V = \{\text{'Mary'}, \text{'John'}, \text{'likes'}, \text{'to'}, \text{'sing'}, \text{'dance'}, \text{'.}\}$$

we could define a tokenizer  $\Phi: V \rightarrow \{1,0\}^{10}$  with the mapping

$$\begin{aligned} \text{'Mary'} &\mapsto [0 0 1 0 0 1 1 1 0 0]^T \\ \text{'John'} &\mapsto [1 0 0 1 0 1 1 0 0 0]^T \\ \text{'likes'} &\mapsto [0 1 0 0 0 0 0 1 1 1]^T \\ \text{'to'} &\mapsto [0 1 0 1 1 0 0 0 0 1]^T \\ \text{'sing'} &\mapsto [1 1 0 0 1 0 1 0 0 0]^T \\ \text{'dance'} &\mapsto [1 0 0 1 1 0 0 0 0 1]^T \end{aligned}$$

In this sense the tokenizer maps the sequence

$$S_1 = (\text{'Mary'}, \text{'likes'}, \text{'to'}, \text{'sing'}, \text{'.})$$

to the input token sequence

$$U_1 = \{\mathbf{u}^0, \mathbf{u}^1, \mathbf{u}^2, \mathbf{u}^3, \mathbf{u}^4\} = \left\{ \begin{matrix} \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, & \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, & \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}, & \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 1 \end{bmatrix}, & \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \\ \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, & \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}, & \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, & \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}, & \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \\ \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, & \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 1 \end{bmatrix}, & \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 1 \end{bmatrix}, & \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 1 \end{bmatrix}, & \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}, \\ \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}, & \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}, & \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, & \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, & \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \\ \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}, & \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}, & \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, & \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, & \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \end{matrix} \right\}, \end{aligned}$$

and sequence

$$S_2 = (\text{'John'}, \text{'likes'}, \text{'to'}, \text{'dance'}, \text{'.})$$

To the input token sequence

$$U_2 = \{\mathbf{u}^0, \mathbf{u}^1, \mathbf{u}^2, \mathbf{u}^3, \mathbf{u}^4\} = \left\{ \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 1 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \right\}$$

## Distal Dendritic Segments

Let us next investigate the distal dendritic segments of a cell and consult figure 5, which shows a neuron with two ( $d=2$ ) distal dendritic segments, each of them entertaining five synapses ( $s=5$ ).

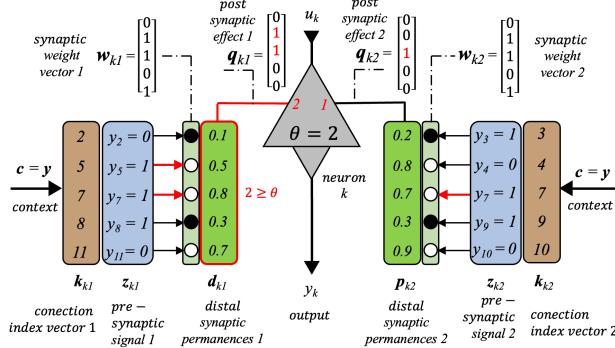


Figure 5: neuron with 2 distal dendritic segments and 5 synapses per segment ( $s = 5, d = 2$ , time index has been omitted)

According to figure 3 the model for cell  $k$  (HTM neuron  $k$ ) provides both a *parameter vector*  $\mathbf{k}_k$  as well as a *distal permanence vector*  $\mathbf{p}_k^t$  for the total set of distal segments. Broken down to a single distal dendritic segment these vectors can be split up into partial parameter vectors  $\mathbf{k}_{kl}$ ,  $\mathbf{k}_{k2}$ , ...,  $\mathbf{k}_{kd}$  and partial distal permanence vectors  $\mathbf{p}_{kl}^t$ ,  $\mathbf{p}_{k2}^t$ , ...,  $\mathbf{p}_{kd}^t$  according to

$$\mathbf{k}_k = [\mathbf{k}_{k1}^T, \mathbf{k}_{k2}^T, \dots, \mathbf{k}_{kd}^T]^T \in \{i \mid 0 \leq i \leq N \wedge i \neq k\}^N \quad (3a)$$

$$\mathbf{p}_k^t = [\mathbf{p}_{kl}^{tT}, \mathbf{p}_{k2}^{tT}, \dots, \mathbf{p}_{kd}^{tT}]^T \in \{i \mid 0 \leq i \leq N \wedge i \neq k\}^N \quad (3b)$$

Referring to figure 5 we have for example

$$\mathbf{k}_k = [2 \ 5 \ 7 \ 8 \ 11 \mid 3 \ 4 \ 7 \ 9 \ 10]^T$$

$$\mathbf{p}_k^t = [0.1 \ 0.5 \ 0.8 \ 0.3 \ 0.7 \mid 0.2 \ 0.8 \ 0.7 \ 0.3 \ 0.9]^T$$

with partial vectors

$$\mathbf{k}_{kl} = [2 \ 5 \ 7 \ 8 \ 11]^T, \mathbf{k}_{k2} = [3 \ 4 \ 7 \ 9 \ 10]^T$$

$$\mathbf{p}_{kl}^t = [0.1 \ 0.5 \ 0.8 \ 0.3 \ 0.7]^T, \mathbf{p}_{k2}^t = [0.2 \ 0.8 \ 0.7 \ 0.3 \ 0.9]^T$$

The connection index vectors  $\mathbf{k}_{k\mu}$  are telling us which other *cell outputs*  $y_i$  are connected to synapses of the *distal dendritic segment*  $\mu$ . As figure 5 demonstrates, segment 1 has 5 synapses which connect to cell outputs  $y_2, y_5, y_7, y_8$ ,

and  $y_{11}$ , according to the index values contained in  $\mathbf{k}_{kl}$ . In a similar way the 5 synapses of segment 2 connect to cell outputs  $y_3, y_4, y_7, y_9$ , and  $y_{10}$ , related to the index values contained in  $\mathbf{k}_{k2}$ .

Formally the connection vector  $\mathbf{k}_{k\mu}$  defines a projection  $P_{k\mu}$ :  $\{0,1\}^N \rightarrow \{0,1\}^s$  given  $\mathbf{v}^t \mapsto \mathbf{v}^t(\mathbf{k}_{k\mu}) =: \mathbf{z}_{k\mu}^t$ . In other words, from the total context  $\mathbf{v}^t$  a distal dendritic segment  $\mu$  of cell  $k$  “sees” only the  $s$  values of the projected vector  $\mathbf{z}_{k\mu}^t$ , which are picked according to the  $s$  index values defined by  $\mathbf{k}_{k\mu}$ . Since the context vector  $\mathbf{v}^t$  is identical with layer output vector  $\mathbf{y}^t$ , we have also:

$$\mathbf{z}_{k\mu}^t = \mathbf{y}^t(\mathbf{k}_{k\mu}) \quad \text{for all segments } \mu \text{ in cell } k$$

If we continue the example shown in figure 5 by assuming

$$\mathbf{c}^t = [0 \ 0 \ 1 \ 0 \ 1 \ 0 \ 1 \ 1 \ 1 \ 0 \ 0 \dots]$$

then we get

$$\mathbf{z}_{k1}^t = \mathbf{c}^t(\mathbf{k}_{kl}) = \mathbf{c}^t([2, 5, 7, 8, 11]) = [0 \ 1 \ 1 \ 1 \ 0]^T$$

$$\mathbf{z}_{k2}^t = \mathbf{c}^t(\mathbf{k}_{k2}) = \mathbf{c}^t([3, 4, 7, 9, 10]) = [1 \ 0 \ 1 \ 1 \ 0]^T$$

for the projected context. Whether these projected cell outputs will all contribute to the excitation of the distal dendritic segment, depends on the following *synaptic rule*:

- A synapse contributes to the excitation of a distal dendritic segment, if, and only if

- (i) the synapse is active (a binary value), which is given if permanence value of the synapse is greater than or equal the synaptic threshold  $\eta$
- (ii) the synapse is connected to an active cell output, i.e.,  $y^t = 1$

Formally this can be modelled by introducing the binary synaptic weight vectors

$$\mathbf{w}_{k\mu}^t = (\mathbf{p}_{k\mu}^t \geq \eta) \quad \mathbf{s}_{k\mu}^t \in \{0,1\}^s \quad (4)$$

which has components equal to 1, if the relating component of permanence vector  $\mathbf{p}_{k\mu}$  has a value greater or equal to synaptic threshold  $\eta$ , and zero otherwise.

Based on the concept of the pre-synaptic signal  $\mathbf{z}_{k\mu}$  and the synaptic vector  $\mathbf{s}_{k\mu}$  we can define now the vector of post-synaptic effects for distal dendritic segment  $\mu$

$$\mathbf{q}_{k\mu}^t = \mathbf{z}_{k\mu}^t \circ \mathbf{w}_{k\mu}^t, \quad \mathbf{q}_{k\mu}^t \in \{0,1\}^s \quad (5)$$

with  $\circ$  denoting elementwise multiplication. If an element of the post-synaptic effect is 1, we know that both the permanence value of the related synapse is greater than or equal to synaptic threshold  $\eta$ , and we also know that the specific context variable relating to the synapse is also active (having value 1). Otherwise, one of the two conditions is violated.

In the next section, where we present the full HTM algorithm for the cell layer in figure 5, we will see, that the

segment activation vectors  $\{\mathbf{q}_{k1}, \mathbf{q}_{k2}, \dots, \mathbf{q}_{kd}\}$ , which can be summarized as the *segment coincidence matrix* of cell  $k$ .

$$\mathbf{Q}_k^t = [\mathbf{q}_{k1}^t, \mathbf{q}_{k2}^t, \dots, \mathbf{q}_{kd}^t], \quad \mathbf{Q}_k^t \in \{0,1\}^{s \times d} \quad (6)$$

It will play a key role for both state updates of a cell and reinforced learning, i.e., adaption of the synaptic permanence values (figure 6). For a compact formulation of the HTM algorithm we also define the *layer coincidence matrix*  $\mathbf{Q}^t$  as

$$\mathbf{Q}^t = [\mathbf{Q}_1^t(:, \cdot), \mathbf{Q}_2^t(:, \cdot), \dots, \mathbf{Q}_d^t(:, \cdot)], \quad \mathbf{Q}^t \in \{0,1\}^{M \times N} \quad (7)$$

Before we tackle the HTM algorithm, let us complete the sample description of figure 5. First we calculate the synaptic vectors for each segment  $\mu$ , assuming the common setting of  $\eta = 0.5$ .

$$\begin{aligned} \mathbf{s}_{k1}^t &= (\mathbf{d}_{k1}^t \geq \eta) = ([0.1 \ 0.5 \ 0.8 \ 0.3 \ 0.7] \geq 0.5) = [0 \ 1 \ 1 \ 0 \ 1] \\ \mathbf{s}_{k2}^t &= (\mathbf{d}_{k2}^t \geq \eta) = ([0.2 \ 0.8 \ 0.7 \ 0.3 \ 0.9] \geq 0.5) = [0 \ 1 \ 1 \ 0 \ 1] \end{aligned}$$

Knowing now which synapses are currently in action, we can calculate the segment excitation.

$$\begin{aligned} \mathbf{q}_{k1}^t &= \mathbf{z}_{k1}^t \circ \mathbf{s}_{k1}^t = [0 \ 1 \ 1 \ 1 \ 0]^T \circ [0 \ 1 \ 1 \ 0 \ 1]^T = [0 \ 1 \ 1 \ 0 \ 0]^T \\ \mathbf{q}_{k2}^t &= \mathbf{z}_{k2}^t \circ \mathbf{s}_{k2}^t = [1 \ 0 \ 1 \ 1 \ 0]^T \circ [0 \ 1 \ 1 \ 0 \ 1]^T = [0 \ 0 \ 1 \ 0 \ 0]^T \end{aligned}$$

$$\mathbf{Q}_k^t = [\mathbf{q}_{k1}^t, \mathbf{q}_{k2}^t] = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 1 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}$$

## The HTM Algorithm

We are ready now to formulate the coarse form of the HTM algorithm for sequence memory. The original form is of this algorithm has been presented in [6]. The form which is used here makes more use of vector/matrix operations and less use of if-branches, which ends up in an algorithm with sequential, straight forward flow.

The computational scheme of the HTM algorithm for a sequence memory layer (as shown in figure 5) is outlined in figure 6.

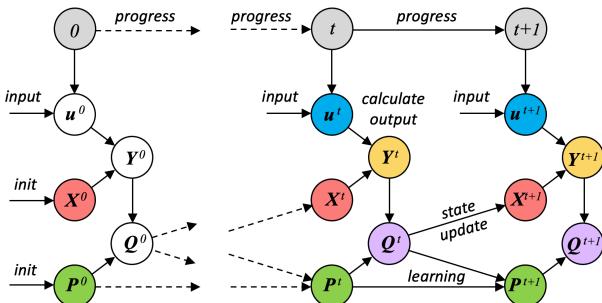


Figure 6: computational scheme for the HTM algorithm  
Besides of the quantities  $\mathbf{u}^t, \mathbf{x}^t, \mathbf{D}^t$  and  $\mathbf{y}^t$ , which are shown in figure 5, the *coincidence matrix*  $\mathbf{Q}^t$  introduced in the

previous section plays an important computational role as an intermediate quantity. Figure 6 can be interpreted as follows: After initializing the layer state  $(\mathbf{x}^0, \mathbf{D}^0)$  at time  $t=0$  a sequence of updating iterations is performed repeatedly for the discrete time sequence  $t = 0, 1, 2, \dots$

```
initialize  $\mathbf{x}^0, \mathbf{D}^0$ 
for  $t = 0, 1, 2, \dots$  :
    input  $\mathbf{u}^t$ 
    calculate  $\mathbf{y}^t$  from  $\mathbf{x}^t, \mathbf{u}^t$  (output activation)
    calculate  $\mathbf{Q}^t$  from  $\mathbf{D}^t, \mathbf{y}^t$  (coincidence matrix)
    update  $\mathbf{x}^{t+1}$  from  $\mathbf{Q}^t$  (cell state update)
    update  $\mathbf{D}^{t+1}$  from  $\mathbf{D}^t, \mathbf{Q}^t$  (cell state update)
```

Before diving into the algorithmic details let us recap that we call a cell  $k$  with activated input *excited* ( $u_k^t = 1$ ). In case of activated output ( $x_k^t = 1$ ) we call the *active*, and we talk about a *predictive* cell when its cell state is set ( $x_k^t = 1$ ). Also keep in mind that all cells belonging to the same minicolumn  $j$  share the same input  $u_j^t$ . Thus, an excitation of one cell in a minicolumn implies an excitation of all cells in the minicolumn, which suggests calling the whole minicolumn *excited*, whenever one cell (which implies: all cells) of the minicolumn is excited.

In addition, we call a minicolumn *predictive*, if it contains at least one predictive cell, otherwise the minicolumn is called *non-predictive*. An excited minicolumn with no predictive cells is called a *bursting minicolumn*.

Finally keep in mind that *cell excitation* (referring to input  $u_k^t$ ) and *segment excitation* (referring to distal dendritic segment excitation  $\mathbf{Q}_k^t$ ) are completely different concepts, and neither *cell excitation* implies necessarily *segment excitation*, or vice versa.

### Step 1) Output Activation

The rules for output activation, where calculate  $\mathbf{y}^t$  is calculated from  $\mathbf{x}^t, \mathbf{u}^t$ , are as follows:

- excited predictive cells get active
- cells of bursting minicolumns get active
- any other cell gets inactive

The first rule is biologically inspired by the capability of pyramidal cells to transition into a mid-term ( $> 100\text{ms}$ ) depolarized state caused by distal dendritic NMDA spikes [3,4]. Such depolarized neurons, which play a predictive role, are easier to activate and will start to fire earlier than non-depolarized neurons [6].

The second rule covers just the default, where there are no depolarized neurons in a minicolumn. In this case no neuron has an activation advantage against other neurons in the minicolumn, and all cells of the bursting minicolumn (which is by definition *excited*) become *active* [6].

The calculations are as follows:

$$\mathbf{B}^t = \mathbf{I} - \mathbf{I}^m \cdot \max(\mathbf{X}^t) \quad \text{burst matrix } \mathbf{B}^t \in \{0,1\}^{mxn} \quad (8a)$$

$$\mathbf{C}^t = \mathbf{X}^t + \mathbf{B}^t \quad \text{candidates } \mathbf{C}^t \in \{0,1\}^{mxn} \quad (8b)$$

$$\mathbf{Y}^t = \mathbf{C}^t \circ (\mathbf{I}^m \cdot \mathbf{u}^{tT}) \quad \text{output activation} \quad (8c)$$

Even this kind of formulation looks a bit unusual, the chosen form has the advantage of generating a straightforward sequence of intermediate matrix results

$$\mathbf{x}^t \rightarrow \mathbf{X}^t \rightarrow \mathbf{B}^t \rightarrow \mathbf{C}^t \rightarrow \mathbf{y}^t$$

with helpful interpretations, as will be demonstrated in the subsequent example. In the consolidated form (9) we clearly see the dependency of  $\mathbf{y}^t$  from  $\mathbf{x}^t$  and  $\mathbf{u}^t$ .

$$\mathbf{Y}^t = (\mathbf{X}^t + (1 - \mathbf{I}^m \cdot \max(\mathbf{X}^t)) \circ (\mathbf{I}^m \cdot \mathbf{u}^{tT})) \quad (9)$$

## Step 2) Calculate Coincidence Matrix

In this step the *layer's coincidence matrix*  $\mathbf{Q}^t$  is calculated from the permanence state  $\mathbf{D}^t$  and the layer output  $\mathbf{y}^t$ . The rationale behind is that each cell entertains  $d$  distal dendritic segments which act as coincidence detectors. Any element of the binary matrix  $\mathbf{Q}^t$ , which equals 1, tells about a coincidence of a related *enabled* distal synapse with a connected active cell. In contrast, a value of 0 means that either the synapse is disabled, or the connected cell is inactive, or both.

The details for the calculation of the coincidence matrix have been explained previously in detail. Thus, we only have to summarize the required calculations of this step.

$$\mathbf{S}^t = \sigma(\mathbf{D}^t - \eta) \quad \text{synaptic matrix} \quad (10a)$$

$$\mathbf{Z}^t = \mathbf{y}^t(\mathbf{K}) \quad \text{projected context} \quad (10b)$$

$$\mathbf{Q}^t = \mathbf{S}^t \circ \mathbf{Z}^t \quad \text{coincidence matrix} \quad (10c)$$

Again (10a-c) can be consolidated to a single operation:

$$\mathbf{Q}^t = \sigma(\mathbf{D}^t - \eta) \circ \mathbf{y}^t(\mathbf{K}) \quad (11)$$

Each Boolean element of the *synaptic matrix*  $\mathbf{S}^t$  in (10a) tells us, whether the related synapse is enabled or disabled, depending on whether its permanence value exceeds threshold  $\eta$  or not. In (10b) the layer output  $\mathbf{y}^t$  is projected into a subspace to yield  $\mathbf{Z}^t$  according to the index values of  $\mathbf{K}$  in order to achieve compatible matching with the synaptic matrix  $\mathbf{S}^t$ . Finally  $\mathbf{S}^t$  and  $\mathbf{Z}^t$  are element-wise multiplied to establish the coincidence values.

## Step 3) Cell State Transition

The aim of this step is to update the *cell state* by calculating  $\mathbf{x}^{t+1}$  from the *coincidence matrix*  $\mathbf{Q}^t$ .

$$\mathbf{p}^t = [\|\mathbf{Q}_1^t\|_1, \|\mathbf{Q}_2^t\|_1, \dots, \|\mathbf{Q}_n^t\|_1]^T \in \{0 \dots s\}^n \quad (12a)$$

$$\mathbf{x}^{t+1} = (\mathbf{p}^t \geq \theta) \quad (12b)$$

Equation (12) means that the given *coincidence matrix*  $\mathbf{Q}^t$  has to be decomposed according to (7) into its parts  $\mathbf{Q}_k^t$ . Then, according to (6), for each  $\mathbf{Q}_k^t$  we calculate the coincidence number

$$p_k^t = \|\mathbf{Q}_k^t\|_1 = \max(\|\mathbf{q}_{k1}^t\|_1, \|\mathbf{q}_{k2}^t\|_1, \dots, \|\mathbf{q}_{kd}^t\|_1)$$

as the largest column sum, and form vector

$$\mathbf{p}^t = [p_1^t, p_2^t, \dots, p_n^t]^T \in \{0 \dots s\}^n$$

In (12b), finally, we set cell state  $x_k^{t+1} = 1$  when the related coincidence number exceeds a threshold ( $p_k^t \geq \theta$ ), otherwise we set  $x_k^{t+1} = 0$ . Again (12a,12b) can be replaced by the consolidated form (13), which expresses the dependency of  $\mathbf{x}^{t+1}$  from  $\mathbf{Q}^t$

$$\mathbf{X}^{t+1} = ([\|\mathbf{Q}_1^t\|_1, \|\mathbf{Q}_2^t\|_1, \dots, \|\mathbf{Q}_n^t\|_1]^T \geq \theta) (m \times n) \quad (13)$$

## Step 4) Learning

In the last state the permanence values of the distal dendritic segments are updated, which is called learning. The following rules are applied:

- synapses of active cells with coincidence increase their permanence values
- synapses of active cells with no coincidence decrease their permanence values
- inactive cells do not change their permanence values

Formally we have

$$\Delta \mathbf{D}^t = (\delta^+ \mathbf{Q}^t - \delta^-) \circ \mathbf{Y}^t \quad (14a)$$

$$\mathbf{D}^{t+1} = \mathbf{D}^t + \Delta \mathbf{D}^t \quad (14b)$$

with  $\delta^+ > \delta^-$ . The consolidated form is:

$$\mathbf{D}^{t+1} = \mathbf{D}^t + (\delta^+ \mathbf{Q}^t - \delta^-) \circ \mathbf{Y}^t \quad (15)$$

## Summary: HTM Algorithm

Taking (9,) the HTM algorithm looks in its compact form:

```

configure  $\mathbf{K}$ , init  $\mathbf{X}^0, \mathbf{P}^0$ 
repeat for  $t = 0, 1, 2, 3$ 
    input  $\mathbf{u}^t$ 
     $\mathbf{Y}^t = (\mathbf{X}^t + (1 - \mathbf{I}^m \cdot \max(\mathbf{X}^t)) \circ (\mathbf{I}^m \cdot \mathbf{u}^{tT}))$ 
     $\mathbf{Q}^t = \sigma(\mathbf{D}^t - \eta) \circ \mathbf{y}^t(\mathbf{K})$ 
     $\mathbf{X}^{t+1} = ([\|\mathbf{Q}_1^t\|_1, \|\mathbf{Q}_2^t\|_1, \dots, \|\mathbf{Q}_n^t\|_1]^T \geq \theta) (m \times n)$ 
     $\mathbf{D}^{t+1} = \mathbf{D}^t + (\delta^+ \mathbf{Q}^t - \delta^-) \circ \mathbf{Y}^t$ 

```

## References

- [1] Hawkins J., Blakeslee S.: "On Intelligence"; Owl Books (2005).
- [2] Hawkins J.: "A Thousand Brains"; Basic Books, New York (2022).
- [3] Antic, S. D., et al.: "The decade of the dendritic NMDA spike"; *J. Neurosci. Res.* 88, 2991–3001. doi: 10.1002/jnr.22444 (2010).
- [4] Major, G., Larkum, M. E., and Schiller, J.: "Active properties of neocortical pyramidal neuron dendrites"; *Annu. Rev. Neurosci.* 36, 1–24. doi: 10.1146/annurev-neuro-062111-150343 (2013)
- [5] Rosenblatt, F.: „The perceptron. A probabilistic model for information storage and organization in the brain“; *Psychological Reviews*, 65 (1958): S. 386–408.
- [6] Hawkins J., Subutai A.: „Why Neurons Have Thousands of Synapses, a Theory of Sequence Memory in Neocortex“; *frontiers in Neural Circuits*, 2016.
- [7] Hawkins J., Subutai A., Cui Y.: „A Theory of How Columns in the Neocortex Enable Learning the Structure of the World“; *frontiers in Neural Circuits*, 2017.
- [7] Chklovskii et al.: (2004)
- [8] Stuart, Häusser: (2001)
- [9] Berlyand L., Jabin P.E.: „Mathematics of Deep Learning; An Introduction“; *de Gruyter Textbook*, 2023.