# University of Southampton Research Repository

# UNIVERSITY OF SOUTHAMPTON

Faculty of Engineering and Physical Sciences
School of Electronics and Computer Science

# Memory Consolidation in Memristive Systems

*by*

## Christos Giotis

ORCiD: 0000-0002-1825-0097

*A thesis for the degree of*
*Doctor of Philosophy*

February 2023

University of Southampton

Abstract

**Memory Consolidation in Memristive Systems**

by Christos Giotis

This thesis investigates the challenge of memory consolidation and learning in artificial synapses. The adoption and evolution of artificial intelligence (AI) also byproducts a frequently overlooked exponentially increasing need for information processing and data storage. This issue is either met with the physical expansion of storage facilities or with the inevitable forgetting of old information in favour of new; both of which seriously hinder the performance of embedded AI systems. This work presents a novel approach in emulating the complex biochemical mechanisms which allow neuronal synapses to store multiple memories on top of the other and at different timescales, like a palimpsest, and which give rise to the incredible learning capacity of biological intelligence.

This work mainly focused on exploiting the intrinsic time dependent volatility in emerging memristive nanotechnologies to showcase palimpsest consolidation. Memristive volatility was studied using a data-driven approach and device-agnostic characterisation and mathematical modelling methods were developed to uncover the main properties of the mechanism. It was found that volatility can exist bidirectionally in $TiO_2$ memristors and that its time constants can be manipulated via the invasiveness and/or frequency of device stimulation. Importantly, within a given observation time window, volatility was shown to operate at two timescales; a fast decay of large magnitude followed by a saturating steady state and a small non-volatile residue. By operating memristive devices as binary synapses, spiking plasticity events were able to store long-term memories in the non-volatile residue, while expressing the opposing state in the short-term. Palimpsest consolidation was examined in simulated memory networks which were able to protect long-term memories while expressing up to hundreds of uncorrelated short-term memories. It was also found that these networks bear close resemblance to the visual working memory of mammalian brains. The same plasticity dynamics were finally extended towards the context of neuronal activity detection, where memristive sensors were able to 'learn' during high spiking frequencies and 'forget' during less active timeframes.

The results presented in this thesis verify the candidacy of volatile memristors as natural facilitators of learning in AI. The ability to learn continuously without catastrophically forgetting old memories, can create new possibilities in the way AI can be used to undertake more generalised tasks. Moreover, the same artificial synapses have shown immense potential in neural interfacing. This can potentially reshape the ways AI is currently interpreted and lead to novel research which aims to integrate both biological and artificial intelligence.

# Contents

# List of Figures

# List of Tables

# Declaration of Authorship

I declare that this thesis and the work presented in it is my own and has been generated by me as the result of my own original research.

I confirm that:

This work was done wholly or mainly while in candidature for a research degree at this University;

Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated;

Where I have consulted the published work of others, this is always clearly attributed;

Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work; I have acknowledged all main sources of help;

Where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself;

Parts of this work have been published as:

C. Giotis, A. Serb, S. Stathopoulos, L. Michalas, A. Khiat and T. Prodromakis, "Bidirectional Volatile Signatures of Metal–Oxide Memristors—Part I: Characterization," *IEEE Transactions on Electron Devices*, vol. 67, no. 11, pp. 5158-5165, Nov. 2020, doi: 10.1109/TED.2020.3014854.


C. Giotis, A. Serb, S. Stathopoulos and T. Prodromakis, "Bidirectional Volatile Signatures of Metal-Oxide Memristors—Part II: Modelling," *IEEE Transactions on Electron Devices*, vol. 67, no. 11, pp. 5166-5173, Nov. 2020, doi: 10.1109/TED.2020.3022343.


T. Abbey, C. Giotis, A. Serb, S. Stathopoulos and T. Prodromakis, "Thermal Effects on Initial Volatile Response and Relaxation Dynamics of Resistive RAM Devices" *IEEE Electron Device Letters*, vol. 43, no. 3, pp. 386-389, March 2022, doi: 10.1109/LED.2022.3145620.


C. Giotis, A. Serb, V. Manouras, S. Stathopoulos, and T. Prodromakis, "Palimpsest Memories Stored in Memristive Synapses," *Science Advances*, vol. 8, no. 25, Jun. 2021, doi: 10.1126/sciadv.abn7920.


Signed:.......................................................................        Date:.................

# Acknowledgements

I wouldn't be the person I am without my family. *Mom, Dad*, you've been together through thick and thin, you've given me and my sister the world and always with a smile. I couldn't have asked for a better place to grow up in. You are the perfect example - you have never given up, always hopeful and have always worked towards sharing a better tomorrow with us. You taught us how to love and be loved and for that we will always be complete. My sister *Eirini*, my number 1. You are my best friend, always eager to understand me and selflessly support me. I only hope I am your equal in those regards. I cannot wait to share the rest of my life with you friend! To all of you, thank you. Those who live without love see the world in black and white, play music with half the notes. Our life together is a vibrant colour explosion, the loudest, most playful concert for the ages. It's us four for life.

My work, my aspirations, my life would mean nothing if I couldn't share it with true friends. *Nikola* and *Dorry*, we met as children and became men together. We dreamt wild dreams and did stupid things - I wouldn't change anything. You are my brothers, my life with you has been a trip, thank you for everything! *Pano*, we climbed mountains together, saw Athens for what it is and what it could be and became family. May we keep climbing. *Georgina* you are the fire I wish will always keep burning in my life. Be who you are, passionate and feisty and let's roam free through everything! I will always be here for you! *Petro* my free spirit, your life is an adventure and I am so happy to be a part of it. Thanks for watching floating cities in the skies with me, Pegasus will always be there for us. *Niko*, *Mike*, *Iasona*, you have celebrated with me at my best, you've picked me up at my worst. I can only do the same for you. Our band is a beautiful incoherence and I love every minute of it. *Ruth*, *Stelio*, *Niko*, *Filippe*, I know you forever and you have never stopped bringing joy and laughter in my life. Always looking out for each other, dreaming, dancing and tasting together. I don't want to think of a life without you, I want to be here for you and hope we will share our futures together! I wish all of you nothing but fortune in your lives. I wish you the energy to materialise all your dreams and live free in the way you choose to! Thank you all for being part of me - the best is yet to come!

*And no one showed us to the land*
*And no one knows the where's or why's*
*But something stirs and something tries*
*And starts to climb toward the light.*
*- Pink Floyd*

Life, like science, is art. Art that is passionate and audacious. We seek to find hidden truths and translate them to our language. Seek on, seek free.

Love, Freedom, Passion, Good Vibes and a Smile,

Christos Giotis (2022)

*To my beloved Grandmother Rinoula. My second mother, feisty and loving, who was over the moon when she learnt I was going to be a man of science. I carry your memory with me every day. This one is for you Yiayia. Love.*

# Chapter 1

# Introduction

## 1.1 Motivation

Progressive deciphering of the human brain and its constituent processing methods has inspired phenomenal advances in Artificial Intelligence (AI) in recent years [1]–[3]. The capabilities of deep learning algorithms and artificial neural networks (ANNs) have already transformed modern societies. ANNs employ finely tuned neuronal layers which perform statistical learning over input data for a given state space. This process yields great generalisation performance **over single or highly correlated learning tasks** such as image classification, speech recognition and natural language processing [4], [5]. Owing to this, AI is set to be further adopted in executable tasks [6] with strong attention on embedding it on edge-systems [7], [8].

Yet, it is evident that current deep learning algorithms bring strong inefficiencies towards realising artificial general intelligence (AGI) even in a centralised (cloud-based) environment. For context, using a comparable number of trainable synapses (the biological learning unit), approximately $10^{13}$-$10^{14}$, biological neural networks in the cerebral cortex facilitate general cognition [9], [10], a capacity much larger than their engineered counterparts. ANNs cannot achieve this capacity by learning sequentially since incoming synaptic modifications impose destructive interference on the networks' state, known as *catastrophic forgetting* [11]. In other words, ANNs must forget old memories in favour of new ones. The adoption of AI in more learning tasks of increasing complexity thus arrives at the cost of ever-increasing computational and memory resources [7], [12], [13]. Artificial memory is further challenged by the limits of Moore's law [14] prohibiting the scaling of ANNs in form factors small enough for AGI on the edge.

To address this challenge, strong focus has been put on loading AI's operations on complementary processing methods. Graphics processing units (GPUs), tensor processing units (TPUs), field-programmable gate arrays (FPGAs) and other specialised CMOS-based circuitry [15]–[18] allow more efficient, *accelerated* training and deployment of

FIGURE 1.1: Motivation of this thesis. A solution to the discrepancy of learning capacity between biological and artificial synapses and catastrophic forgetting by emulating biological consolidation in volatile memristors.

software-based ANNs. Concurrently, more specified neuromorphic technologies aim to directly emulate brain and learning functions in hardware, completely bypassing the need for software implementations of ANNs [19]–[22]. The appeal of neuromorphic systems first rises from attributes such as learning speed and energy efficiency [23]. To date though, neither implementation has demonstrated a solution to catastrophic forgetting. Since ANNs need to increase exponentially in size to yield linearly higher accuracy [24], it is natural to look for different approaches in increasing learning capacity.

One such approach seeks to increase the re-usability of ANNs allowing the same synapses to be used for **multiple and uncorrelated tasks**; a property naturally observed. Biological synapses are able to consolidate multiple memories which can be revealed at different timescales - much like a palimpsest [25]. Synapses can remember long-term plasticity events, namely potentiation (LTP) and depression (LTD) while expressing altered states in the short-term [26]. This temporal partition enables the brain to use the same resources for multiple computation processes and offers a practical partition to the issue of catastrophic forgetting.

Palimpsest storage is realised biologically via the bidirectional interaction of hidden biochemical processes affecting the manifestation of synaptic efficacy at different timescales [25], after each memory modification. These processes are characterised by their own degrees of plasticity (i.e. learning rates) and lifetimes (i.e. 'forgetting time constants'). These properties make memristive devices natural candidates for hardware palimpsest consolidation, which have already showcased their potential in synaptic emulation [20], [21], [27]–[29]. In particular, frequently overlooked volatile resistive RAM (RRAM) families are also governed by hidden electrochemical processes affecting their analogue state, much akin to biological synapses. **Thus, the primary motivation of this thesis has been the harnessing of RRAM's intrinsic volatile properties to consolidate multiple uncorrelated memories at different timescales and at palimpsest fashion** (see

FIGURE 1.2: Project objectives and thesis organisation.

Figure 1.1 for a high-level description). If appropriately tuned bidirectionally, the relaxation dynamics of RRAM can replicate the hidden biochemical states that protect memories from synaptic modifications. **Such memristive synapses can address catastrophic forgetting and pave the way for AI hardware that can learn on the Edge.**

## 1.2 Research Objectives

The research objectives of this project are as follows:

1. Assess the candidacy of volatile RRAM technologies for emulating synaptic memory consolidation.

2. Characterise the intrinsic time-dependent signatures of RRAM volatility.

3. Accurately model the intrinsic time-dependent signatures of RRAM volatility.

4. Demonstrate in hardware how palimpsest memory consolidation can naturally arise via volatile RRAM.

5. Extend the technology's principal operations towards further neuromorphic applications.

To achieve the objectives of this project a number of key steps have been taken (see Figure 1.2 for an objective outline and the subsequent thesis organisation). **Firstly**, a thorough review has been conducted to verify the suitability of TiO$_x$-based RRAM technologies for palimpsest memory consolidation. This review has covered relevant neuroscience theories of synaptic operations and memory consolidation, a background

foundation of memristive devices as synaptic emulators as well as key hardware and software learning implementations. RRAM can efficiently replicate plasticity events either in the form of LTP or LTD and plasticity rates can explicitly be modulated via appropriate stimulation tuning. Moreover, it has been concluded that RRAM volatility could provide a reference to the hidden relaxation properties of biological synapses. However, no volatility framework that could be utilised to express the necessary timescales bidirectionally has been found to exist.

Thus, **secondly**, this project has set to develop a complete framework for both characterising and modelling RRAM volatility bidirectionally. This framework aims to decipher the relationship between relaxation timescales and plasticity stimulation to support appropriate interfacing with memristive synapses. This study has shown that RRAM volatility is governed by two distinct timescales; a fast volatile regime and a slow non-volatile residue after stimulation. **Thirdly**, necessary 'bridging' algorithms have been designed to manipulate memristive volatility both in hardware and in simulations and replicate basic synaptic functions. The volatile and non-volatile timescales have been exploited to realise a fast learning, plastic short-term memory and a slower but more rigid long-term compartment in a single device.

**Fourthly**, experiments have been designed to demonstrate automatic palimpsest consolidation in volatile RRAM synapses using this dual memory capacity, accompanied by explanatory data analysis. These experiments have practically demonstrated that volatile synapses can consolidate memories in the long-term capacity and automatically protect them against multiple short-term memory modifications. This proof of concept in hardware has satisfied the main motivation of this thesis. **Finally**, owing to the ubiquity of intrinsic time referencing in biology, this project has aimed at applying volatile learning properties in further neuromorphic applications. This has been carried through in the general context of brain-computer interfacing (BCI). There, single devices have shown great potential in encoding high neuronal spiking activity without a need for pre-processing spiking data.

## 1.3   Thesis Organisation

This thesis is organised as follows. Chapter 2 surveys the literature to give an overview of memory and learning in biology and in computational models, as well as of recent engineered implementations. This includes a description of biological synapses, their main plasticity mechanisms and mathematical descriptions of memory consolidation to form the theoretical foundation of this work. Moreover, an introduction of RRAM technologies and their core properties is made, followed by a review of memristive technologies' involvement in artificial synaptic learning and in neuronal interfacing. Chapter 3 presents novel methodologies for characterising and modelling memristive

volatility under various ranges of operating conditions. The experimental results presented in this chapter link to the theory behind memory consolidation to form the foundation of the main application designed in this thesis. In Chapter 4, memristive volatility parameters are translated into synaptic stimulation protocols to demonstrate palimpsest memory consolidation in hardware. The results are both experimental and simulated and provide a thorough analysis of volatile memristor suitability for learning applications. Chapter 5 presents a brief proof of concept study on how the same learning properties of volatile memristors can be utilised to perform unsupervised temporal integration of neuronal spikes for high activity detection. Finally, Chapter 6 summarises this thesis and provides the author's recommendations on future extensions of this work.

# Chapter 2

# Memory and Learning: Biology and Engineering

## 2.1 Introduction

The ability to learn, recall and consolidate memories is a fundamental aspect of biological and artificial intelligence. Memory consolidation and learning occur biologically via the synapse, an 'adjustable connection' between neurons which regulates the transmission of signals in the brain. Naturally, this has drawn extensive scientific attention both towards understanding synaptic mechanisms as well as engineering them. This chapter introduces the synaptic mechanisms required for learning, both biologically and computationally, and gives an overview of how these mechanisms can be reverse-engineered for artificial intelligence and further neuromorphic applications.



FIGURE 2.1: Chapter's outline and objectives.

FIGURE 2.2: Mechanisms of memory consolidation. Information is initially stored in a fast and labile short-term memory. Reinforcement of memory enables it to be consolidated into a slower but stable long-term memory. A cleanup mechanism is required to ensure forgetting of information deemed unimportant to the system.

## 2.2　Memory Consolidation Conceptually

Efficient handling of information in any learning system requires both mechanisms for retaining important memories and equally importantly, forgetting less significant ones. The human brain deals with this challenge by storing memories in distinct short and long-term time phases [30] [31]. While the information in the short-term memory is easily recorded, it is also easily forgotten. Contrarily, information in the long-term memory is recorded with more difficulty but is also more stable [30]. This mismatch of stability can act as a cleanup filter in any learning system, whereby information that has not been forgotten in the short-term can progressively be stored more securely. The process of transferring information from the labile short-term phase to the more stable long-term memory is referred to as *memory consolidation* [32].

At a higher level, consolidation of memory can be perceived as a controlled process. For it to function, the system needs the two types of mentioned memory systems and four mechanisms; signal transfer, positive reinforcement, a threshold metric and cleanup. A heuristic approach is the following. Consolidation occurs via the transfer mechanism which moves information to the system's long-term memory. The transfer may only be triggered if a threshold condition has been met. In turn, memories can meet this condition if they are reinforced in the short-term. All memories that fail to reach the transfer threshold are discarded by the cleanup mechanism. Crucially, both the memory stages and the mechanisms are orthogonal modules of the consolidation model and in theory, they can be implemented in numerous ways both in vivo and in silico. A higher level depiction of the model is shown in Fig. 2.2.

Importantly, consolidation can occur both at a systems level and at a memory-unit level. That is, the distinction between short and long-term memory can either be spatial or structural within the storage system respectively. The spatial consolidation paradigm

FIGURE 2.3: Neurotransmitter release during synaptic operation [37]. The electrical signal from the presynaptic neuron is converted into chemical information through neurotransmitters and back to electrical form in the postsynaptic neuron.

is employed in conventional computing where fast but volatile RAM devices consolidate data to slower but stable hard drives. Evidence for systems consolidation in the brain involves the transfer of information from the hippocampus (*short*) to the neocortex (*long*) [33] [32]. At the unit level, consolidation can occur through changes in protein concentrations in the synapse [30] [34].

## 2.3 The Synapse: An Overview

Information processing occurs in the human brain via signal transfer between individual neurons. This signal has the form of short spikes of electrical current (Fig. 2.3 (1)) propagating from the outward end of a transmitting neuron (*axon*) to the inward end of a receiving neuron (*dendrite*) [9] [35]. The junction between these two elements is termed *synapse* and it is widely considered the principal unit of memory and computation within the brain [9] [36]. In reality, this form of neural communication is not the only one that has been observed [9]. This work, however, only considers the synapse as a connection between a transmitting (*presynaptic*) and a receiving (*postsynaptic*) neuron.

Conceptually, it is very important to understand a reduced model of how a synapse operates, to appreciate how it can function both as the memory and the computational unit of the brain. Generally, synapses can be categorised into two major types; excitatory and inhibitory. Excitatory synapses activate the postsynaptic neuron when triggered and inhibitory ones stop the electrical signal flow [9]. While information is both transmitted and received by neurons in electrical form, the synapse operates with slower, biochemical processes and makes use of concentrations of specialised molecules

called *neurotransmitters* [9] [38]. Such molecules are contained and released within *synaptic vesicles*, "packets" of quantised amounts of molecules.

The electrical potential induced by a presynaptic neuron triggers a change in potential difference within the synapse. This in turn releases quantised amounts of neurotransmitters (Fig. 2.3 (3,4)) which carry the signal to the postsynaptic neuron (Fig. 2.3 (5,6)). This signal can either be of excitatory (e.g. in the case of Glutamate-mediated synapses) or inhibitory (e.g. *GABA*) nature [39] [9]. Specifically, the effect of either excitatory or inhibitory events is determined by the number of vesicles released. The more excitatory neurotransmitters are released, the more likely the postsynaptic neuron to fire. Conversely, the more inhibitory neurotransmitters are released, the less likely is for the neuron to fire. The release of synaptic vesicles depends on the level of presynaptic activity. However, the number of (available) vesicles formed within the synapse can be variable and their synthesis depends on the history of synaptic activity [9]. The degree of likelihood that a synapse releases a specific amount of neurotransmitters and thus excites or inhibits the postsynaptic neuron is called *synaptic efficacy*. Efficacy is the equivalent of the weight of a connection in an artificial neural network (ANN). In turn, the mechanism by which efficacy changes is called *synaptic plasticity* and its two sides are *potentiation*, i.e. increase in efficacy and *depression*, i.e. decrease in efficacy.

### 2.3.1   Plasticity Mechanisms

Synaptic plasticity refers to mechanisms by which synaptic efficacy is modulated. It is explicitly activity-dependent, meaning that it is reactive to signalling neuronal inputs and is widely regarded as a primary facilitator of biological memory [40]. Plasticity can be categorised using various metrics. Most notably, it can be distinguished by the lifetime of the efficacy changes it induces (short- and long-term plasticity), as well as by the underlying mechanisms facilitating it. Additionally, higher levels of plasticity can be identified and in fact, are considered to be key mechanisms of memory consolidation. This section will give a conceptual overview of these characteristics since they will be a recurring topic in the next chapters of this thesis.

First, synaptic efficacy is known to be modulated for short-lived timeframes by what is termed short-term plasticity. The mechanism has been recorded in numerous organisms of varying complexity and is believed to have lasting ranges from milliseconds to minutes [41]. Typically, short-term plasticity is expressed relative to paired-pulse stimulations, whereby an input pulse can affect a synapse's response to a closely spaced subsequent pulse. Short-term plasticity can facilitate both an increase (facilitation) and a decrease (depression) in efficacy. The plasticity direction is normally determined by the temporal interval between pulse sequences, with small interpulse periods (< 20ms) typically (but not always) leading to short-term depression and larger (20-500ms) to

FIGURE 2.4: Schematic illustration of LTP and LTD, shown as the increase and decrease of the field excitatory postsynaptic potential (fESP) in the hippocampus over the timeframe of 1 hour [40].

short-term facilitation [41], [42]. Short-term plasticity is believed to provide biological brains with filtering capabilities. Specifically, since high-frequency stimulation can cause short-term depression, high-efficacy synapses can respond to sparse signals with high accuracy but may dampen denser input signals; effectively acting as a low-pass filter. Equivalent but opposite relationships can turn low-efficacy synapses into high-pass filters [43].

While short-term plasticity and its associated filtering capabilities can provide some interesting properties in neural networks, experiences perceived by the brain can also induce long-term changes in synaptic efficacy. These changes occur concurrently at ensembles of synaptic circuits and have been reported to persist for hours to days in mammalian brains [44], [45]. Such actions result in the permanent or semi-permanent rewiring of biological neural networks, which is the foundation of all learning. Long-term plasticity takes form both via potentiation (LTP) and depression (LTD), which are umbrella terms for various mechanisms governing efficacy modulation in multiple brain regions and for multiple purposes [40]. An illustration of both mechanisms is shown in Fig. 2.4. Particularly, the co-existence of both LTP and LTD within individual synapses and the corresponding bidirectional nature of efficacy [46], is a key justification for perceiving synaptic memory as an analogue weighted value - a perception which is a pillar of modern deep learning algorithms [1].

One of the most widely studied mechanisms for inducing long-term plasticity is spike-timing-dependent plasticity (STDP). This mechanism is of great interest since it is also the main learning rule of neuromorphic applications [23]. STDP is aligned with the long-standing Hebbian learning rule, which proposes the existence of a mechanism for increasing the signal transmission efficiency between neurons that tend to communicate frequently [47]. This is commonly summarised by the popular heuristic 'neurons that fire together, wire together', which effectively suggests that synaptic plasticity must act as a mechanism for learning in neural networks.

FIGURE 2.5: Types of STDP in biological neurons. The x-axis measures the temporal difference between the presynaptic and postsynaptic firing $t_{pre}$ - $t_{post}$ in milliseconds (ms). LTP and LTD are colour-coded in blue and red respectively.

In its basis, STDP between a presynaptic and a postsynaptic neuron is determined by the temporal relationship of their activity. STDP is determined by the magnitude of the time difference between spiking events occurring at two neurons, $t_{pre}$ - $t_{post}$, as well as the sign of this value. Most pronounced plasticity changes occur at time differences of 25-50ms, although many different directions have been observed (see Fig. 2.5 for common examples as described in [48], [49]). Most neuromorphic applications follow the learning rule shown in Fig. 2.5a [23], [29]. If the firing of a postsynaptic neuron succeeds that of the presynaptic neuron then the corresponding synapse undergoes LTP; conversely, if the postsynaptic firing precedes the presynaptic then the synapse undergoes LTD.

Strong association between LTP/LTD and memory can also be drawn due to their similar underpinning properties [50]. First, plasticity can be reinforced via repetition. Second, the change in a synapse's efficacy is correlated to the strength of an input signal. Moreover, the mechanisms by which plasticity is facilitated are mostly intrinsic to individual synapses, meaning that they do not affect their neighbours and give rise to the incredible degree of freedom found in learning brains.

## 2.3.2   Synaptic Metaplasticity

It is evident that synaptic plasticity is a key facilitator of biological memory. Changes in synaptic efficacy 'rewire' neural networks and regulate their function, thereby altering their output to specific inputs and thus enabling learning. Nevertheless, plasticity itself cannot account for memory consolidation and/or protection. After all, reinforcement can aid plasticity but it does so bidirectionally. In other words, modifications in efficacy, caused by ongoing input signals between neurons inevitably cause the degradation of older memories [25]. Intuitively, this would mean that brains should constantly keep forgetting memories in favour of new ones.

The mechanism that is regarded to be protecting memory degradation in the synapse is metaplasticity. Metaplasticity can be thought of as a higher-level attribute of synaptic plasticity. It describes the mechanism by which a synapse can undergo changes to its degree of plasticity [51]. Metaplasticity thus does not refer to the absolute changes of

synaptic efficacy but to the freedom by which efficacy can change. For instance, it has been observed that synapses in the hippocampus may be subject to stimulation that can favour LTD and limit LTP in the future [52], [53].

Synaptic metaplasticity can solve the challenge of protecting consolidated memories in the brain and thus has been a big inspiration for this project. While there are still many unknowns as to how plasticity is regulated molecularly, biological evidence for its existence has been observed [33], [54]. Moreover, computational neuroscience has made significant progress in approximating these effects. The next section will cover some key computational aspects of memory consolidation and protection, which have been utilised in the design of artificial synapses throughout this work.

## 2.4 Computational Models

### 2.4.1 Defining Memory Capacity

As it is outlined in Chapter 1 and will later be discussed in detail in Chapter 4, the primary aim of this work is to realise palimpsest memory consolidation in a hardware implementation. Since this involves the coexistence of multiple memory states within a system, it is imperative that a clear metric is defined to distinguish between palimpsest and non-palimpsest memory consolidation. To find such a metric, this work considers the possible different states that a memory unit can concurrently hold *without forgetting old information*. Importantly, the co-storage of multiple states could theoretically occur both concurrently (being able to read/recall multiple states from a memory unit at the same time) or could be partitioned in time. The latter would involve overwriting a memory state for some time while retaining the ability to reverse back to a previous state after that time has elapsed. Such an ability is believed to be present in biological synapses, has already been formulated mathematically (see Section 2.4.4 of this chapter) and is the primary inspiration for this work.

Accordingly, this work defines the *capacity* of a memory system as the total number of states that can be arranged in a temporal palimpsest fashion. For example, a typical SSD memory unit can only store 1 state since it cannot be retrieved when overwritten. A theoretical system which could store an N-th signal while retaining the ability to recall the previous N-1 signals would have capacity N. This definition should be distinguished from conventional definitions of memory capacity in the context of information theory (i.e. bits as units of signal entropy [55]) which is largely used to evaluate artificial memories [56], [57]. For example, if a theoretical memory unit can take 4 distinct states with equal probability but can only do so irreversibly with respect to previous states, then its capacity within the context of this work will be equal to 1 state and not 2 bits. Therefore, since no practical demonstration of palimpsest memory consolidation has been

shown in hardware, we can deduce that all known memory systems have a capacity of 1 state and thus the main aim of this work can be rephrased to simply demonstrate a memory capacity greater than 1.

### 2.4.2   Memory Storage

Memories are retained in the brain through changes in the efficacies of relevant synapses. To better understand this notion certain assumptions should be made. Firstly, every experience that can be remembered, can also be translated into some input stimulus on some neural ensemble. Secondly, this input has the form of presynaptic spikes directed towards some postsynaptic neurons. This activity in turn induces plasticity changes in the set of synapses connecting the neurons activated while the experience is recorded. Memory traces of a given experience can then be remembered as long as the efficacies of those synapses remain close to the values prescribed by the experience [25].

Specifically, traces of a stored memory within an ensemble of synapses lie in the modifications in their synaptic efficacies. The memory signal is strongest immediately after modification occurs. The ongoing plasticity events induced in the memory system add noise to the memory trace which is assumed to decay over time. At any point, the memory signal can be computed as the "overlap between the state of the synaptic ensemble and the pattern of modifications originally imposed by the memory being remembered" [58] (see Chapter 4).

The brain comprises a finite number of synaptic connections, the efficacy of which is limited by bounded values. Since memories are retained by ensembles of synapses, those bounds also limit the possible combinations of efficacies, i.e. memories that can be stored. Moreover, the efficacy of a given synapse undergoes ongoing plasticity changes which constantly add noise to the traces of older memories. Yet the human brain exhibits the ability to recall information for years and still keep storing new memories. Computational metaplasticity models offer a perspective on how this is achievable and an inspiration for engineering equivalent mechanisms.

### 2.4.3   Cascade Metaplasticity Model

Using a metaplastic metric, a synapse can vary from being completely *rigid* (i.e. not liable to plasticity changes) to *plastic* (i.e. very prone to plasticity changes). From a memory consolidation perspective, designing a metaplastic synapse ensures the short and long-term memory types at a structural level, as mentioned in Section 2.2.

Computationally, metaplasticity can be visualised using the cascade model of a binary synapse, developed by Fusi et. al. [59]. The synapse can obtain two efficacy states,

FIGURE 2.6: Cascade model of metaplastic synapse. Binary efficacy values consist of *weak* and *strong*. Each *potentiation* event induces a plasticity change towards plastic $strong_1$ state with transition probability $q_i$ if a synapse is in $weak_i$ state. Importantly, plasticity transition $q_i$ becomes exponentially less likely for higher values of $i$. Alternatively, it consolidates memories into more rigid $strong_i$ metaplastic states with transition probability $p_i^+$, if the synapse is in state $strong_{i-1}$. Each *depression* event acts respectively on the synapse for metaplastic transition probabilities $p_i^-$.

potentiated and depressed. Metaplasticity works as a hidden variable, where continuous plasticity events of equal polarity push the synapse deeper into states *i* with lower efficacy transition probabilities, $q_i$. The model is illustrated in Fig. 2.6 [59]. Referring to Section 2.2, the *reinforcement* mechanism of the system is the repetition of plasticity events and the *transfer* mechanism is the transition to different metaplastic states. Information is consolidated by being "entrenched" into more stable states, while information that is not reinforced can easily be overwritten.

The cascade model provides a straightforward explanation of how metaplasticity aids memory consolidation. However, it is not robust enough to support consolidation in general forms. This has been recently demonstrated in the case of binarised neural networks that employ weights with cascade metaplastic properties [60]. There, it is shown that such networks can only generalise over multiple tasks/memories which are highly correlated. This can be intuitively understood since cascade synapses can only hold one binary state which is catastrophically forgotten as soon as it is changed. Cascade metaplasticity thus favours a form of static consolidation by a binary approach towards interfering memories: either incoming memories are completely ignored if the synapse is rigid or are written by permanently overwriting the existing state. While this paradigm has been proven useful in expanding the storage of correlated memories which share the same efficacy states [60], it still challenges online AI in real learning scenarios where the nature of incoming memories is unknown a priori. What is needed is the flexibility to learn uncorrelated memories fast, and concurrently protect older consolidated counterparts. Importantly, this should be done automatically and in an unsupervised manner.

FIGURE 2.7: Synaptic model. (**a**) Chain plasticity model consisting of $u_k$ biochemical processes, where $k \in [1, m]$, for $m$ different processes. Each variable interacts with its neighbours, while $u_1$ is connected to $u_2$ and the input and $u_m$ communicates with $u_{m-1}$ and a leak term (equivalent to $u_{m+1} = 0$). Parameters $g_{k,k+1}$ are the strengths of the interactions between two variables. Combined with $C_k$, they determine the timescale on which each process operates. (**b**) Intuitive illustration of hidden variables. The chain model behaves like an ensemble of communicating beakers. Each variable $u_k$ measures the deviation of each variable from the equilibrium, as shown with $u_3$. Variables $g_{k,k+1}$ represent the connection weight between two beakers and $C_k$ reflects the area/volume of each beaker.

### 2.4.4   Palimpsest Consolidation Model

An explanation for the brain's incredible learning capacity would be that a specific synapse is used for the storage of multiple memories. Memories then, are not stored statically in a *one next to the other* fashion but are more likely to be stored in a palimpsest fashion, i.e. *one on top of the other* [61]. In such a scenario, however, plasticity events triggered by multiple memories stored in a given synapse, deviate its efficacy from values that would suffice for recalling previously stored memories.

The challenge of storing and retaining memories is encompassed by the *plasticity / rigidity dilemma* [11]. Plastic synapses are good at storing new information but their liability to further plasticity events makes them forget as easily. Rigid synapses are good at retaining old information but they fail at storing new memories. As proposed by Fusi et.al. [25], for consolidation to occur at a synaptic level, metaplasticity has to be multidimensional and a synapse must contain both plastic and rigid components. The authors argue that memories are retained by multiple biochemical processes, operating at different timescales and interacting with each other bidirectionally. Visualisation of such interaction is provided in Figure 2.7 [25].

This synaptic model consists of a visible weight $w(t)$ and multiple hidden variables operating at progressively slower timescales. Computationally, this is described by

a chain model. The first element of the chain is the synaptic weight and is also the most plastic component. Its value is very sensitive to plasticity events, resulting in high initial memory strengths. The remaining chain elements represent the hidden variables affecting synaptic metaplasticity and equilibrate around the average values of their neighbours (Fig. 2.7a).

The chain model can be expressed by a chain of communicating liquid beakers (Fig. 2.7b). The yellow beaker represents the synaptic weight and the most plastic component. Potentiation events equate to adding liquid to this beaker and depression events act oppositely. Progressive plasticity events tend to perturb the equilibrium levels of deeper hidden beakers. Thus, the information of persistent plasticity patterns becomes "entrenched" within the synapse and can dictate the trend of the synaptic weight for longer, via the continuous propagation of liquid in between the beakers.

For instance, if many potentiation events occur, slower variables deviate higher than their equilibrium level. Then, possible depression events can affect the synaptic efficacy within a fast timescale but such information becomes overwritten by the bidirectional communication in the long run. In such cases, the synapse becomes rigidly potentiated. The same unit is able to recall a depression event for a short, vulnerable timescale and a long-term potentiation over a stable long-term period. Over an ensemble of $N$ synapses, memories are consolidated structurally, in a *palimpsest* fashion. Importantly, this configuration allows both the automatic consolidation of memories, reinforced only by the frequency of storage and the protection of old memories in the long term. This is a significant advantage which can allow hardware AI to undertake multiple learning tasks without catastrophically forgetting previously learnt memories.

This work has focused on creating equivalent synaptic models on hardware by engineering the appropriate timescales required for palimpsest consolidation. Such a task can only be undertaken by utilising emerging nanotechnologies which have equivalently complex memory capabilities. Novel *memristive* technologies have thus been considered, and in particular resistive random-access memories or RRAM, the learning capabilities of which are discussed in the following section.

## 2.5 RRAM Technologies

### 2.5.1 An Overview

Since their inception [62] and first realisation [63], memristive devices show great potential in neuromorphic computing [64] [65]. While several types of memristors exist, this project focuses on Resistive Random Access Memory (RRAM) devices. These are two terminal (Metal-Oxide-Metal) devices which operate as tunable resistors whose resistive states depend on the history of applied bias. Their activation process usually

also requires an invasive form of initial biasing called *electroforming* [66]. RRAM offers a significant leap in storage per area with memristors being able to store more than 6 bits of information in a single device [67].

Device resolution and history dependence make RRAM a natural integrator element, a perfect candidate for neural implementations on hardware, for example, the Tempotron [68] [69]. Moreover, memristors exhibit a direct resemblance to synaptic efficacy properties such as long-term potentiation (LTP) and depression (LTD) [21] [70]. Lately, their tunability has been examined in numerous neuromorphic applications including the encoding of neuronal spikes [71] [72], unsupervised learning [20] and in-memory computing [73] [65]. Memristors have also been used to realise metaplastic phenomena [74] [27], where plasticity effects vary depending on timing conditions. However, no prior work has associated RRAM technology with consolidation in a *palimpsest* fashion.

### 2.5.2   RRAM Switching Mechanisms

The change in a memristor's resistive state $R$ is induced via the application of appropriately invasive biasing voltages $V$. Assuming some device-specific threshold value $V_{th}$, higher *programming* amplitudes induce changes in $R$, while sub-threshold *read* amplitudes can be used to read a device's state at any time. Thus, by employing sequences of programming events accordingly, one can manipulate resistive changes $\Delta R$ within some resistive range, as illustrated in Fig. 2.8. There, a chosen RRAM device responds bidirectionally to programming events of opposite polarities. Specifically, $R$ (shown in grey data points) increases with positive voltages and decreases when negative amplitudes are applied. The apparent accumulation of $\Delta R$ with successive programming events gives rise to the intrinsic integration properties of RRAM. Moreover, the device's behaviour shows dependence on the degree of invasiveness by which programming events are applied; a common pattern that will also be prevalent in this work.

This dependence is represented with more clarity in Fig. 2.9 by [76]. (a) shows how the relative change in resistance $\frac{\Delta R}{R}$ is affected by increasing programming amplitudes in both polarities. The transition between sub-threshold voltages (boxed in red dashed lines) to invasive programming amplitudes is shown as $\frac{\Delta R}{R}$ deviates from 0. Concurrently, (b) shows a characteristic I-V relationship during $\Delta R$ which is reflected by the hysteresis in the I-V curve. This relationship is of significant importance since quantifying the parameters by which RRAM switching can be manipulated is an essential first step towards the design and implementation of memristor-based neuromorphic applications.

However, while memristive properties are promising candidates for post von-Neumann neuromorphic architectures, no unified theory yet exists to encompass their operation,

FIGURE 2.8: A typical representation of RRAM $\Delta R$, presented in [75]. Consecutive trains of positive programming pulses (shown in blue in the bottom half) cause $R$ to progressively increase (grey data points in the top half), while negative programming pulses decrease $R$. Importantly, the boundaries of $\Delta R$ saturation widen as the absolute programming amplitude increases. The switching behaviour is accurately modelled, as illustrated via the overlapping red line and grey data points.



FIGURE 2.9: RRAM switching characterisation as presented in [76]. (a) The dependence of relative $\Delta R$ on the programming pulse amplitude. As shown, the data can be categorised into two distinct regions, the sub-threshold region (boxed in red) and the invasive bias region, where device $R$ is changing. (b) I-V dependence for the same RRAM device. The apparent hysteresis loop is an indication of a change in the resistive state due to invasive biasing.

even though progress has been made [76] [77] [78]. Evidence suggests that RRAM operation may depend upon a variety of factors, such as the oxide film and the electrode materials [79]. Resistive switching may occur either through electrochemical changes in their active electrodes [80] [81], via redox reactions which induce conductivity changes of the Metal-Oxide film of the device [82], or via thermal ionic diffusion due to current induced changes [81]. Lastly, switching may occur via the alteration of the electrode-oxide interfacial potential barrier [79].

Concurrently, work has been conducted to model the behaviour of RRAM devices under bias. Results suggest that one can simulate with great accuracy resistive switching for known biasing regimes [75]. This model describes the rate of change of device resistance as:

$$\frac{dR}{dt} = s(v) \times f(R, v) \tag{2.1}$$

where $s(v)$ is the voltage sensitivity function of the device, depending only on the pulsing scheme and $f(R, v)$ is the device window function which imposes soft switching boundaries as a function of resistance $R$ and voltage $v$.

This is also shown in Fig. 2.8 where this model (red line) is used to fit raw switching data (grey). Nevertheless, this work only assumes ideal non-volatile devices and cannot be used to simulate metaplastic timescales.

### 2.5.3   RRAM Volatility

Within the context of this work, RRAM volatility is defined as the *overall monotonic* drift of device resistive state, *after* and not *during* device programming protocols. Volatility is thus defined as a strictly transient phenomenon, observed under non-invasive and sub-switching-threshold reading voltage amplitudes, that is observed only after appropriate invasive stimulation protocols. This is to be distinguished by the commonly observed stochastic fluctuations in RRAM resistive state when successive state readings are performed. Such fluctuations are present both during passive device readings [83] and during both positive and negative amplitude stimulation protocols [84][77]. Such fluctuations do not show any dependence on the device state history and in this work are considered to be and referred to as noise (see Chapter 3 - 4 as well as Appendix C). On the contrary, RRAM volatility studied in this work follows clear monotonic trends which tend to average out over multiple continuous noisy fluctuations and over longer periods of time.

Volatile behaviour has been reported in RRAM studies in the past [85][68][86]. The passive drift of resistive state over time has been linked with memory consolidation as a potential cleanup mechanism [27] (see Section 2.2). Work by Cortese et.al. suggests that

volatility is a result of ionic diffusion occurring after the removal of the applied electric field [85]. This gradually increases the electrode-oxide interfacial barrier and thus the device's resistive state. This is supported by Liu et.al who report similar diffusion phenomena [87].

It has already become evident that memristive volatility can be designed to a desired effect [88]–[90]. However, in order to fully realise the potential of the volatile RRAM families a systematic approach is needed for characterising and modelling volatility. Existing volatility models suggest that volatile device families can be switched in a non-volatile way only if a minimum energy level has been input during stimulation, while they behave in a volatile manner otherwise [68].. This is explained by separating volatile from non-volatile switching via quantised energy levels needed to induce stable changes in a device's interfacial barrier (the proposed mechanism for manipulation of R). In this scenario, a device's state is always drawn towards a local R attractor. A SPICE model has also been described and has served as inspiration for work conducted in this project. Moreover, routines have been developed to distinguish between volatile and non-volatile regions of memristive operation [86]. Furthermore, attempts to model volatility via equivalent circuits have been presented by [91].

Literature on volatile RRAM switching has focused on predicting the phenomenon within non-volatile technologies. The work that has so far been conducted in this project offers a novel, data-driven volatility model. This can be used in the development of synaptic circuits, with multiple timescale variables, extending the ideas proposed mentioned in [25] and discussed in 2.4.4.

## 2.6 Hardware Synapses

The advantages in learning that can be brought to AI hardware by direct emulation of the discussed synaptic properties have already attracted the interest of the field. This is particularly true within the broader context of memristive technologies, where multiple studies have demonstrated the technologies' intrinsic resemblance to synaptic plasticity.

Memristive synapses have first demonstrated basic plasticity rules in hardware in as LTP and LTD. Based on the analogue switching regime of non-volatile memristors, multiple implementations have shown that efficacy changes become more pronounced as potentiation/depression cycles are applied successively. These studies are also based on phase change memory (PCM) memristive technologies which undergo changes in conductance as bidirectional stimulation is applied to them to reflect potentiation and depression. Emulation of plasticity in this paradigm has been achieved both in standalone devices [29], [92], [93] or by integrating memristors into more complex systems [94], [95]. While these are significant milestones towards designing artificial synapses,

plasticity alone is not sufficient for learning and protecting memories from incoming interference; a fundamental functionality for autonomous learning systems.

To address this issue, both RRAM and PCM memristor studies have worked towards emulating synaptic metaplasticity [51], [59] directly via manipulating the learning rate of the artificial synapses. This has been shown in integrated CMOS-based synapses in the context of spiking neural networks (SNNs) [95], [96]. Metaplasticity has been a very popular research target since it can address the protection of consolidated memories from ongoing interference. Similarly, non-volatile RRAM has also exhibited tunable switching/learning rate via the direct manipulation of stimulation bias [74], [97]–[100]. While these studies have successfully shown tunable switching rates, they do not consider memory consolidation holistically. First, metaplasticity is manipulated by explicit changes in external stimulation. This implies that the need for changes in the switching rates and thus consolidation is known a priori. This is not realistic for autonomous AI where agents are expected to learn in an online fashion and in unpredictable scenarios. Moreover, all these studies consider metaplasticity bidirectionally and exclusively on either LTP or LTD scenarios. For metaplasticity to have full effect, it is essential that it also addresses catastrophic forgetting bidirectionally, allowing a synapse to protect consolidated memories from all forms of interference at the same time [59].

Finally, several studies have focused on memory protection directly, albeit in the context of passive memory lifetime. This has been a feature exclusive to volatile RRAM, which has shown a transition from short- to long-term memory lifetimes [21], [27], [101], [102]. The transition is induced by the frequency of stimulation of a particular memory and can be directly associated with the reinforcement mechanism for consolidation that is discussed in 2.2. However, the transition in all those studies is both irreversible and unidirectional. Specifically, long-term memory has only been shown via LTP, where a binary potentiated state is written for longer time periods due to successive stimulation. While this solidifies the candidature of volatile RRAM as an efficient synaptic emulator, these studies fail to consider the protection of long-term memories under interference caused by opposing synaptic modifications. Thus, the issue of catastrophic forgetting remains unaddressed.

All mentioned studies have either focused on expressing synaptic plasticity or emulating metaplasticity in a static way. Long-term memory has been discussed both in terms of reduced learning rates or longer state lifetimes. A summary of the mentioned technologies is presented in Table 2.1 in chronological order. Consequently, there still remains a clear gap in addressing memory consolidation **and** reversible protection from catastrophic forgetting in an autonomous way. This niche can be addressed by manipulating multiple consolidation timescales as discussed in 2.4.4. Filling this gap is the main objective of this work and is thoroughly discussed in the following chapters of this thesis.

| Example | Technology | Type | Speed | Metaplasticity | Consolidation | LTP/LTD | Lifetime | Timescales | Capacity |
|---|---|---|---|---|---|---|---|---|---|
| Chang et.al. [27] | WO$_x$ RRAM | Single device | 1ms | No | Yes | LTP | < minute | 1 | 1 |
| Ohno et.al. [102] | Ag$_2$S synapse | Single device | 500ms | No | Yes | LTP | ≈ 20s | 1 | 1 |
| Ambrogio et.al. [94] | GST PCM | Multi device | 250ns | No | No | LTP/LTD | n/a | n/a | 1 |
| Berdan et.al. [21] | TiO$_x$ RRAM | Single device | > s | No | Yes | LTP | 10s | 1 | 1 |
| Tan et.al. [101] | WO$_x$ RRAM | Single device | 10μs | Explicitly tunable | Yes | LTP | minutes | 1 | 1 |
| Boybat et.al. [29] | GST PCM | Single device | < μs | No | No | LTP/LTD | n/a | n/a | 1 |
| Burr et.al. [93] | GST PCM | Single device | n/a | No | No | LTP/LTD | n/a | n/a | 1 |
| Cheng et.al. [98] | YSZ RRAM | Single device | 0.75-1.5ms | History dependent | No | LTP/LTD | n/a | n/a | 1 |
| La Barbera et.al. [92] | GST PCM | Single device | 5-300ns | No | No | LTP/LTD | n/a | n/a | 1 |
| Lee et.al. [100] | KN memristor | Single device | 100μs | Explicitly tunable | No | LTP | n/a | n/a | 1 |
| Liu et.al. [99] | Graphene memristor | Single device | 100ns | Explicitly tunable | No | LTP/LTD | n/a | n/a | 1 |
| Wu et.al. [97] | HfO$_x$ RRAM | Single device | 1μs | Explicitly tunable | No | LTP/LTD | n/a | n/a | 1 |
| Brivio et.al. [96] | HfO$_x$ RRAM | Multi device | 10μs | No | No | LTP/LTD | n/a | n/a | 1 |
| Demirag et.al. [95] | GST PCM | Multi device | 100ns | No | Yes | LTP | 30s | 1 | 1 |

TABLE 2.1: Overview of existing demonstrations of artificial synapses and their corresponding learning properties.

## 2.7    Memristors as Neural Interfaces

The last objective of this work has been to extend the application spectrum of RRAM technologies beyond strict memory consolidation while maintaining focus on their learning capabilities. Thus, a last remark should be made regarding the functionality of memristors in further neuromorphic fields and in brain-computer interfacing (BCI) specifically.

Deciphering interfacing with the processes of the human brain brings the immense potential to both next-generation medicines [103], [104] and human-AI integration [105], [106]. Applications harnessing BCI have managed to bridge communication with tetraplegic patients [107], [108], facilitate the control of neuroprosthetic arms [109] and even decode speech directly from recorded brain activity [110]. Neuronal spikes, the messaging signals between neurons, have traditionally been recorded from brain cells using bespoke probing technologies [111], [112]. However, most implementations have faced bottlenecks with respect to efficiency, parallelism and real time processing since raw neuronal data typically suffer from very low signal-to-noise ratios (SNR) [71]. For brain activity to be deciphered, interfacing systems need to efficiently detect and sort neuronal spikes from input reading channels, and consequently decode the resulting spiking data meaningful information.

FIGURE 2.10: Neuronal spike detection using MIS. (**a**) A recorded trace of raw neuronal activity. The input signal is modulated such that low SNR spikes exceed the switching threshold of the RRAM sensor. (**b**) Corresponding switching history of the sensor, where resistance changes correlate to the occurrence of spike data. (**c**) The actual timeline of detected spikes is obtained via control testing software. Each spike is illustrated with a black vertical line - a total of 81 spikes are recorded. (**d**) Green-shaded regions mark instances where at least one change in R has been recorded by the MIS. A total of 74 changes have been recorded.

TiO$_x$ based RRAM technologies have already provided a solution to the spike detection and sorting problems. This has been achieved with the concept of the memristive integrating sensor (MIS) platform. MIS exploits RRAM's intrinsic switching characteristics (see Fig. 2.8) and the biasing threshold dynamics (see Fig. 2.9) to filter noisy spiking distinguish low SNR neuronal spikes and encode them into iterative state changes in resistance [71]. A high level illustration is presented in Fig. 2.10, taken from the original publication [71]. Appropriately modulated neuronal signals can produce voltage spikes that cross the switching threshold of the sensor, while noise levels remain below the threshold at all times. As such, spiking events can be encoded in analogue resistance changes as read in the sensor, completely bypassing the challenges imposed by the low SNR of the raw spike data.

MIS addresses encoding neuronal spikes in real-time and at low powers, providing a significant advantage in the field of BCI. However, spike detection alone does not carry information about the general state of a biological network. Subsequent analysis of neuronal data is therefore needed to encode the temporal firing characteristics and distinguish between regions of low and high spiking activity. While this task can

be carried out via external processing, it would be very beneficial to perform activity detection using bespoke hardware solutions that could be integrated with MIS. This frequency filtering bears a strong resemblance to the low/high pass filtering that can be expressed via short-term plasticity (as discussed in 2.3.1). Since the main objective of this thesis has been the manipulation of RRAM volatility to emulate plasticity dynamics that support memory consolidation, the work can naturally be extended towards spike activity detection. This objective has been examined briefly towards the end of this work and is presented in Chapter 5.

## 2.8 Summary

Advances in our understanding of synaptic mechanisms and their role in biological memory and learning in the last century have enabled great progress in AI and BCI. Concurrently, emerging RRAM technologies have shown the potential into emulating neuromorphic functions that are unattainable by conventional computing. This thesis aims to combine the two fields, first by creating a new framework for the previously overlooked volatile RRAM domain and then by utilising volatility to emulate time-dependent mechanisms of memory consolidation and neural interfacing.

# Chapter 3

# Bidirectional Volatility in TiO$_2$ RRAM

## 3.1 Introduction

The translation of consolidation timescales to artificial synaptic properties and vice versa can be achieved by utilising the intrinsic time dynamics of volatile RRAM. To achieve this, a thorough understanding of RRAM volatility is needed and requires both the identification of its key characteristics and the factors which affect them. This knowledge can allow a controllable mapping of plasticity functions and forgetting timescales to RRAM stimulation and time-dependent state transitions which can then be manipulated for palimpsest operations. This chapter focuses on addressing these



FIGURE 3.1: Chapter's outline and objectives.

challenges by presenting a thorough, data-driven approach for characterising and modelling RRAM volatility. The approach that has been developed in this thesis is technology- and application-agnostic and focuses on understanding how volatility unfolds after stimulation and importantly, how it can be manipulated to desired specifications. Thus, the results presented in this chapter can be decoupled from the neuromorphic perspective and although some key points are mentioned, a detailed mapping from volatility to plasticity functions is discussed in chapter 4.

This chapter is organised as follows. Section 3.2 introduces the scope within which RRAM volatility is discussed, including definitions of the key operational parameters that characterise RRAM transient relaxation. Section 3.3 presents a methodology for characterising volatility bidirectionally and identifies the main stimulation parameters that dictate its characteristics. In section 3.5, this methodology is used to briefly analyse how these characteristics can be affected by operating device temperature to identify the optimal conditions for realising the main objectives of this thesis. Then, section 3.4 uses the characterisation methodology to define a data-driven mathematical model describing volatility under various stimulation regimes. This model not only provides further insight into volatility itself but also enables the simulation of relevant behaviours for more extended application case studies in chapter 4. Finally, section 3.6 summarises this chapter.

## 3.2   Defining Volatility in RRAM

Within the context of this thesis, *volatility* will refer to the passive, overall monotonic transient change in a memristor's state R, resulting from an explicitly defined stimulation protocol (*i.e.* a set of programming pulses) and unfolding within a predefined retention time window T. This is to be distinguished by the ubiquitously observed RRAM noise [77], [83], [84] that is discussed in Chapter 2. Volatility is illustrated in Fig. 3.2 where a device under test (DUT) experiences volatile changes in R following invasive stimulation. This example is a typical representation of volatility in TiO$_2$ RRAM and serves as a baseline for defining the key volatility characteristics that will be discussed and manipulated in this thesis.

A typical cycle of RRAM volatility is shown in Fig. 3.2a. Here, volatile changes are induced using a set of invasive/stimulation pulses with amplitude $V_P$ (a write phase, shown in red) while DUT's state R is measured using non-invasive read pulses (a read phase, shown in blue) with some amplitude $V_R$ . For clarity purposes, a pair of consecutive write and read phases is defined as a retention cycle. Prior to simulation, DUT is recorded at some state $R_{pre}$ . Following a set of N stimulation pulses, passive volatility is observed (see inset). Starting from state $R_{start}$ after stimulation, and within the time window T, DUT passively converges towards an end state $R_{end}$ .

FIGURE 3.2: Experimental retention data showing typical volatile behaviour in TiO$_2$ based memristors in (**a**), following stimulation shown in (**b**). Starting from a resistive state pre-stimulation R$_{pre}$ , the device is programmed (*red*) using N pulses at voltage V$_P$ . Volatility then unfolds over time when the resistive state is read using non-invasive signals at V$_R$ (*blue*). This behaviour occurs bidirectionally, depending on the polarity of the programming stimulus. Retention data are shown more analytically in (**c**). The device is inspected for a fixed time window T. Starting from a state R$_{start}$ at t=0, it converges to R$_{end}$ at t=T. This relaxation is characterised by a time constant τ.

Any change in the resistive state ΔR = R(t < T) - R$_{pre}$ is considered volatile within the defined window. However, ΔR = R$_{end}$ - R$_{pre}$ is considered to be a non-volatile residue at time T after stimulation. When consecutive retention cycles are considered, the R$_{end}$ of one cycle is also the R$_{pre}$ of the next one. It is very important to note that thus far in this project, R(t) is only attempted to be understood within T and results are not extrapolated outside this time window. From an application perspective, this means that a designed memory system will only be observed within time T from the latest stimulation event and subsequently, refresh mechanisms can be used to preserve a desired state of a synapse. Alternatively, more work is needed towards understanding memristive drift in much larger timescales.

## 3.3   Bidirectional Volatility: Characterisation

Using the defined quantities above, this thesis has put a strong emphasis on characterising RRAM volatility, so that it can be reproduced in a controlled manner within the context of synaptic plasticity. This includes identifying the key parameters affecting volatility as well as mapping the relationship between the two counterparts, such that RRAM stimulation can emulate balanced plasticity events.

### 3.3.1   Characterisation Methodology

#### 3.3.1.1   Stimulation Parameters and RRAM Relaxation

To identify the factors that may dictate R(t), individual programming phases have been considered. Each programming phase is comprised of 4 externally controlled parameters: *(1)* a set of N identical pulses of *(2)*, fixed invasive pulse amplitude $V_P$ , *(3)* fixed duration width and *(4)* interpulse time, used to induce changes in R. Assuming that R(t) is a function of the stimulation that is exerted on DUT as well as its current state, and by considering conditions parameters (3-4) to be constant, then volatility can be expressed as:

$$R(t) = f(N, V_P, R_{pre}) \tag{3.1}$$

It follows that to decipher volatility characteristics in a data-driven approach, R(t) has to be sampled under a wide $\{N, V_P, R_{pre}\}$ state space. While $V_P$ and N can independently be chosen, controlling $R_{pre}$ presents a significant challenge since its value depends on the history retention cycles which produce the unknown phenomena that are yet to be characterised. To address this challenge, retention cycles have been applied in a staircase fashion (Fig. 3.3a). Successive sets of cycles form a staircase with a finite number of steps, each of which employs progressively more invasive programming phases. The end of one staircase is followed by a new one which includes one additional, more invasive cycle. Thus, the $R_{pre}$ -$V_P$ and $R_{pre}$ -N can be effectively swept in order to examine the role of $R_{pre}$ in R(t). The effective sampled area is shown in Fig. 3.3a.

Through the sampled areas and the resulting dependencies, one can properly understand how to utilise the input programming parameters to tailor memory events for specific needs. For instance, it is essential to know how can new memory events be mapped towards *write* events of minimum to maximum strength. By deciphering the factors that govern R(t) can be done through appropriate choices of N and $V_P$ . Here, the following key attributes are considered: the retention limit points $R_{start}$ and $R_{end}$ as well as the relaxation time constant $\tau$.

FIGURE 3.3: Data sampling methodology for volatility characterisation. (**a**) Successive cycles are used to sweep either the number of programming pulses in a cycle (N) or the amplitude of the programming pulses ($V_P$) in both polarities, two factors affecting relaxation dynamics. The invasiveness of each stimulation cycle progresses in a staircase fashion. Stemming from the nature of memristive volatility, the application of multiple retention cycles induces a drift in $R_{pre}$. Consequently, this protocol samples either the $R_{pre}$-$V_P$ or $R_{pre}$-N planes as illustrated in (**b**).

The latter is extracted by fitting retention data with the stretched exponential function, which has already been linked to memristive volatility in [27]. Stretched exponential relaxation is commonly observed in disordered systems due to dispersive transport arising by large amounts of randomly distributed trapping mechanisms [113]–[115]. Therefore, DUT relaxation is mathematically expressed by:

$$R(t) = \alpha \exp^{-\left(\frac{t}{\tau}\right)^{\beta}} + \gamma \tag{3.2}$$

Here, R(t) represents the time-dependent change in R over time t. $\alpha$ reflects the relative change in R at t=0s compared to the projected saturation point of R at t→ ∞. The relaxation time constant is determined by $\tau$, while $\beta \in$ (0,1] is an exponential stretch

FIGURE 3.4: Volatility repeatability test. The figure illustrates the deviation of $R_{end}$ from the DUT's initial state $R_{initial}$ , as successive retention cycles employ increasing absolute $V_P$ magnitudes in both polarities. Each test terminates either when the R(t) has become indistinguishable from $R_{initial}$ which is determined by a *t-test* (see text) or when a total of T = 1 hr has elapsed. The datapoint colour map indicates the time elapsed before each retention test terminates. The datapoints from cycles which have failed the t-test are highlighted in red.

factor. At $\beta$= 1, R(t) follows a normal exponential decay while as $\beta \rightarrow 0$, R(t) is temporally stretched. Finally, $\gamma$ is equal to the theoretical saturated R value as t $\rightarrow \infty$.

### 3.3.1.2   Volatility Proof of Concept

The first step towards evaluating the memristive technology in question has been the evaluation of the universality of the observed volatile phenomena, intra-DUT as well as across multiple devices. This is crucial if we want to exploit a technology that can be reliable in memory applications. To begin with, DUT is put under test in order to observe whether its volatility characteristics are repeatable.

Starting from the device's initial state $R_{initial}$, increasingly more invasive retention cycles have been applied in phases and repeatability has been examined with respect to how effectively, the original state is reinstated after stimulation has been exerted. Retention *blocks* of 3 identical retention cycles have been employed for every value of programming pulse invasiveness —$V_P$ —, before this increases by 1V. Each programming cycle consists of a single pulse of 100 μs. In all retention tests, T is variable with each retention terminating either when a t-test between R(t) and $R_{initial}$ is satisfied (a t-value = 2 is chosen) or after a total of T = 1 hr.

The results of this repeatability test are shown in Fig. 3.4. The figure plots each retention cycle independently and shows the percentage (%) difference between $R_{end}$ and $R_{initial}$ against the corresponding $V_P$ (positive cycles in circle and negative shown in square shapes). Each datapoint is colour-coded with respect to the time elapsed before each cycle is terminated.  The cycles where the t-test has failed are also highlighted, indicating no convergence within 1 hour.

In all positive cycle runs below 7V, the test terminates within the first few retention minutes, thereby passing the t-test. Indeed, all $R_{end}$ datapoints are well below a 1% difference from $R_{initial}$ . Programming pulses of 7V have passed the t-test 2/3 times, while all cycles above 7V fail the test and terminate after 1 hour has passed.  In contrast, all negative retention cycles pass the t-test, and all $R_{end}$ values lie within a 0.5% deviation from $R_{initial}$ .

The asymmetry that is apparent in opposite programming polarities is a known effect of RRAM which utilises electrodes of different materials [116].  Based on the results obtained from the repeatability test, a programming pulse amplitude of 7V has been chosen to be the reference $V_P$ value throughout this characterisation study, since it provides a more pronounced volatility effect than 6V and a better opportunity to examine the potential accumulation of non-volatile residues stemming from multiple retention cycles.  This accumulation can be translated in a memory reinforcement mechanism, whereby stronger and repeated memory events can change the state of the synapse in a more rigid way.

Concurrently, summary results on cycle-to-cycle and device-to-device variability for volatility in both directions are shown in Fig. 3.5. The data indicate that relaxation dynamics are similar across devices which is an encouraging step towards large-scale integration of the technology.  Of course, such an endeavour requires independent study which is outside the scope of this thesis.  Moreover, the fact that volatility exists bidirectionally is on its own a significant technological advancement that hasn't yet been observed in memristive families. Specifically, clearly defined relaxation trends can be explicitly directed towards upper or lower resistive states, given opposite displacements in R caused by positive or negative programming amplitudes respectively. The dependency of volatility on the direction of previous programming events and its clear overall trend, distinguish this phenomenon from stochastic fluctuations in R (i.e. noise) that are frequently observed in RRAM families (see also Chapter 2 and Appendix C). There is therefore an opportunity for exploiting this characteristic from an application perspective as bidirectional volatility can enable the realisation of plastic potentiation/depression synaptic changes.

FIGURE 3.5: Retention tests over multiple devices in order to investigate volatility variability across them. Each device is subject to 5 consecutive retention cycles employing N = 500 pulses at ± 7V, as illustrated in (**c**). (**a**) DUT is examined in 3 consecutive cycles and exhibits identical relaxation characteristics. A downward trend in R can be explained by the accumulation of non-volatile residues in each retention. (**b**) The protocol is applied to 3 devices and all candidates show equivalent responses. Importantly, (**a**) and (**b**) both show that our devices exhibit volatile characteristics in both directions.

### 3.3.1.3   Characterisation Protocol and Sampling Space

To obtain a clear dataset on RRAM relaxation, a relatively large T value has been chosen. Specifically, volatility has been recorded for a total time T of 2 minutes. This has been done to investigate whether it is feasible to a) induce strong volatile changes fast enough and b) retain these changes for a long time period. This decision balances the assumption that learning systems in online scenarios may be subject to constant stimulation, with the for RRAM synapses to store short-term memories without a need for frequent reinforcement. Such conditions construct an ideal scenario for addressing the plasticity-rigidity dilemma mentioned in [25].

Volatility dependencies on $V_P$ and N have been examined independently. When $V_P$ is swept, N is kept constant value of 10 pulses per cycle and vice versa, $|V_P| = 7V$, when N is variable. Detailed descriptions of the distinct *Voltage Dependence* and *Pulse Dependence* protocols are presented in Alg. 1. The characterisation parameters used throughout this work are summarised in Table 3.1.

---

**Algorithm 1:** Characterisation Protocol

---

**Voltage dependence:** Sweep $V_P$ , N = 10;

$\mathbf{V} = [1; 0.5; 9V]$;

**for** $i \leftarrow 0$ **to** *size of V* **do**

    **for** $V_P \leftarrow V[0]$ **to** *V[i]* **do**

        Apply N pulses with $V_P$ , pw, $\Delta$t parameters;

        **while** $t \leq T$ **do**

            Read with $V_R$ ;

**Pulse dependence** : Sweep N, $|V_P| = 7V$;

$\mathbf{N} = [1] \cup [10; 10; 100] \cup [200; 100; 1000]$;

**for** $i \leftarrow 0$ **to** *size of N* **do**

    **for** $N \leftarrow N[0]$ **to** *N[i]* **do**

        Apply N pulses with $V_P$ , pw, $\Delta$t parameters **while** $t \leq T$ **do**

            Read with $V_R$ ;

---

TABLE 3.1: Volatility characterisation protocol: model parameters.

| Parameter | Description | Value | Unit |
|:---:|:---:|:---:|:---:|
| $V_R$ | Reading voltage | 0.2 | V |
| **$V_P$ Sweep** | | | |
| $V_R$ | Programming voltage | 1 - 9 | V |
| $V_{step}$ | Parametric step in V | 0.5 | V |
| N | Pulses per retention cycle | 10 | N/A |
| **N Sweep** | | | |
| $V_P$ | Programming voltage | 7 | V |
| N | Pulses per retention cycle | 1 - 1000 | N/A |
| $N_{step}$ | Step in $N \in [10, 100)$ | 10 | N/A |
| $N_{step}$ | Step in $N \in [100, 1000]$ | 100 | N/A |
| pw | Programming pulse width | 100 | µs |
| $\Delta$t | Program interpulse time | 1 | ms |
| T | Retention time window | 2 | minutes |

FIGURE 3.6: The effect of V$_P$ amplitude on DUT switching, expressed through the relative change $DeltaR_{start} \equiv \frac{R_{start} - R_{pre}}{R_{pre}} \times 100\%$ (shown in the z-axis), on the R$_{pre}$ - V$_P$ plane (x- and y-axes respectively). The results are shown separately for positive stimulation in (a) and negative in (b). The projections on the x-y plane reflect the effective sampled area. Datapoints are colour-coded with respect to their z-axis values. As shown via the x-z and y-z projections, $\Delta R_{start}$ depends mostly on the value of V$_P$ and not on R$_{pre}$. The mean $\Delta R_{start}$ values for each V$_P$ are shown in red.

### 3.3.2   Effects of Stimulation Regime on Volatile Characteristics

#### 3.3.2.1   Programming Pulse Amplitude

Results have been extracted from the Voltage Dependence characterisation protocol and are presented in Fig. 3.6. The aim of this part of the experiment is to identify whether the strength of incoming memory events can be normalised in between some minimum and maximum switching values when written on a synapse by appropriately regulating V$_P$ .

The relative difference between R$_{start}$ and R$_{pre}$ for each retention cycle, i.e. $\Delta R_{start} \equiv \frac{R_{start} - R_{pre}}{R_{pre}} \times 100\%$, is shown separately for positive and negative $|V_P| \in [1, 9V]$ in Fig. 3.6a-b respectively. Each retention programming phase employs N = 10 identical pulses. The datapoints are presented against V$_P$ and R$_{pre}$ in a 3D fashion. $\Delta R_{start}$ increases monotonically in magnitude with increasing biasing amplitudes, although switching is opposite, according to the polarity, as expected from the results presented in Fig. 3.5.

Importantly, $\Delta R_{start}$ only seems to be affected by the magnitude of V$_P$ and no direct relationship with R$_{pre}$ is apparent. Positive retention cycles begin to have a significant impact on $\Delta R_{start}$ for V$_P$ > 3V where its value increases linearly with voltage until it saturates close to V$_P$ = 8V. These results are in line with previous works suggesting

the dependence of RRAM switching on pulsing amplitude [75], although the exponential relationship between switching sensitivity and $V_P$ is not shown here. Similar behaviours are observed for negative values of $V_P$ although the effect is less pronounced here (a maximum value of ~ 5.5% is noted for positive $V_P$ and a value of ~ 3% for negative).

The asymmetry that is apparent in opposite programming polarities is a known effect of RRAM which utilises electrodes of different materials [116]. This is an indication that if a consolidation application requires symmetrical potentiation and depression events then the technology will need to be operated at different voltages.

### 3.3.2.2 Programming Pulse Number

By now, one can draw a straightforward analogy to link biasing voltage $V_P$ and the strength of the presynaptic potential that is exerted on biological synapses during neural memory consolidation. Intuitively, this can be loosely framed as a *time-invariant* strength by which information is written in a memory system, much like the discrepancy between normalised input values in an ANN. The term time-invariant is used because sole manipulation of $V_P$ can encode the weight of input information but it does not necessarily reflect on the timeframe during which such information is being written for in a real-world scenario. This gap is extremely important as one would expect that information which is present to a learning agent for a longer time would be further consolidated than the information available only for a brief period, should the writing strength, i.e. $V_P$, remain constant. The inclusion of spatiotemporal degrees of artificial memories could effectively be modulated by the number of programming pulses used to change synaptic efficacies. To that extent, it is imperative that the relationship between N and RRAM volatility is further understood.

The second part of the characterisation protocol relates the switching and relaxation characteristics of DUT to $R_{pre}$ and N. Experimental data referring to positive biasing are shown in Fig. 3.7a-c while opposite polarity induced data are shown in Fig. 3.7d-f.

Initially, $\Delta R_{start}$ has been plotted in a similar fashion to Fig. 3.6. In both polarity scenarios, change in R shows a clear dependence on the value of N. For positive biasing, the mean values of $\Delta R_{start}$ for every value of N (shown in red on the y-z projection) appear to initially decay exponentially with the number of pulses (especially for $N \leq 100$) where it increases in magnitude from −1% to −8%. This value tends towards saturation at a value close to 12% for larger numbers $N \geq 500$.

On the other hand, DUT switching shows little dependence on the value of $R_{pre}$ as reflected by the scattering of the $R_{pre}$ - $\Delta R_{start}$ axis projection. However, it could be speculated that the saturation in $\Delta R_{start}$ is a combination of both the diminishing effects of N and the fact that DUT approaches its resistive boundary. In such a case, it is

FIGURE 3.7: The effect of N on DUT switching for both polarities. $\Delta R_{start}$ is represented similarly to Fig. 3.6 but in this case is plotted on top of the $R_{pre}$ - N plane. **Positive stimulation:** As per (a), the mean values of $\Delta R_{start}$ (shown in red) exponentially decay with increasing numbers of pulses. No clear dependence is shown between $\Delta R_{start}$ and $R_{pre}$. The final retention cycle (N = 1000) is illustrated in (b). The stretched exponential model fits excellently the raw retention data. Moreover, the (%) non-volatile residue, i.e. $\Delta R_{end}$ is plotted against N in (c). The red points indicate the mean values for each number of pulses. $\Delta R_{end}$ initially increases with N but saturates quickly for larger numbers of pulses. **Negative stimulation:** DUT shows a similar but opposite behaviour. $\Delta R_{start}$ datapoints for each retention cycle are shown in (d). The retention data for N = 1000 have been fitted in (e). Non-volatile relative residues are shown in (f).

more accurate to say that $R_{pre}$ doesn't affect $\Delta R_{start}$ only as long as the device remains sufficiently far from that boundary.

By comparing Fig. 3.6a to 3.7a it is clear that within the chosen $V_P$ - N state space, it is more effective to induce pronounced changes in R by employing larger values of N (the largest value in 3.6a is about 50% smaller than in 3.7a). Hence, at this point, it is easier to illustrate clear volatile behaviour in DUT, as shown in Fig. 3.7b, where raw retention data from the last retention cycle (N = 1000) have been fitted using Eq. 3.2. The data reveal that DUT volatility operates in two timescales; a rapid change in R which is prevalent for t < 10s, followed by a slower transient relaxation until t = T. Interestingly, the rate of change in R(t) seems to have decreased significantly within the chosen retention window, although no signs of absolute saturation have been observed.

The (%) non-volatile residue $\Delta R_{end}$ is plotted against N in Fig. 3.7c, with the mean value for each N value shown in red. The data show discrepancies between $R_{end}$ and $R_{pre}$ are less than 1%. Judging from such a figure alone, one could assume that a window of 2 minutes is appropriate for fully volatile operations. Nevertheless, as reflected by the

drift in $R_{pre}$ in 3.7a, even this small percentage difference can accumulate to produce significant changes in R when consecutive programming cycles are applied on DUT. This accumulated drift in R can be harnessed as an intrinsic consolidation mechanism within RRAM-based synapses.

Similar but opposite results have been obtained via the negative stimulation characterisation protocol. The obtained $\Delta R_{start}$ surface is plotted in Fig. 3.7d and depicts the same saturating pattern with the increase of N. The mean values (in red) show an exponential decay for $N \leq 100$ followed by a saturating increase which reaches close to 7.5%. This is explained by the lower switching sensitivity of DUT at negative values of $V_P$ .

The retention data obtained for N = 1000 are again plotted against time in Fig. 3.7e. The fitted model (shown in orange) still shows an initial exponential change in R but now signals a slower change over time than that compared in 3.7b. Again, asymmetries are observed between DUT operations at opposite polarities, which will have to be accounted for in future memory applications. The non-volatile residues are presented in 3.7f.

### 3.3.2.3 Relaxation Time Constant

It has already been shown that the volatile behaviour observed in the device technology under test not only unfolds in multiple timescales (remember the rapid change in R followed by a slower transient trajectory in Fig. 3.7b,e) but is also related to relaxation direction (remember the slower recovery shown for negative polarity in the same figures). Further understanding the dependencies of the relaxation time constant $\tau$ can introduce additional flexibility in the operation of RRAM-based synapses, as varying $\tau$ could be used to emulate the multiple consolidation timescales that are discussed in 2.4.4.

Extracted $\tau$ values for each retention fit have been plotted against $R_{start}$ and N in Fig. 3.8a,c, in order to examine whether DUT relaxation depends a) on the device's state when stimulation ends and/or b) on the level of invasiveness employed in each programming phase. As shown by the $R_{start}$ - $\tau$ and N - $\tau$ projections, volatility time constants seem to only depend on N. In the case of positive polarity, a linear relationship between the two values is observed, with $\tau$ ranging from a few milliseconds (ms) to 3s.

Not all retention cycles have been fitted successfully due to the lesser switching effects induced by lower values of N; specifically, only 35% of retention cycles using N = 1 have been fitted successfully, while all other cycles were unsuccessful. An example of a failed retention fit for N = 1 is shown in 3.8b. Due to the low input $E_{total}$, the initial switching of DUT is not sufficient to produce clear relaxation phenomena. The raw

FIGURE 3.8: Dependence of relaxation time constants $\tau$ on $R_{start}$ and N as determined by stretched exponential fits on retention data for positive (**a**) and negative (**c**) stimulus polarities. The datapoints are plotted in a similar fashion to the 3D plots in Fig. 3.6 - 3.7 and are positioned independently to achieve the greatest visibility. $R_{start}$ is plotted in the x-axes of (**a**) and (**c**), while N is plotted on the y-axes. The mean values of $\tau$ for each value of N are shown in red on the x-z projections. **Positive stimulation:** An example of a failed retention fit for N = 1 is shown in (**b**). **Negative stimulation:** An example of a failed retention fit for N = 1 is shown in (**d**) The stretched exponential model fits failure rate for N ≤ 100 is shown in (**e**).

data (purple points) depict an almost instant saturation of R after $R_{start}$ . Thus the fitted model (green) fails to encapsulate the stretched exponential characteristics of volatility.

Stemming from the reduced switching sensitivity of DUT to negative programming cycles, the resulting retention data are less pronounced, something which has propagated down to the extracted $\tau$ values. Specifically, a significant portion of the retention cycles that employ small N values has been fitted unsuccessfully because the resulting volatility is either non-existent or unfolds in faster timescales than those that could be captured. All $\tau$ values larger than 120s (maximum retention window) have been omitted. The relaxation constant $\tau$ appears to be at least 1 order of magnitude larger than the results obtained in 3.8a. It is yet unclear why this is observed, but one possible explanation could be that the natural tendency of DUT to reach its high resistive boundary allows for more rapid relaxation in that direction, i.e. following positive stimulation. Again, an example of an unclear retention dataset is shown in Fig. 3.8d while the success rates for all retention cycles with N ≤ 100 is shown in 3.8e.

### 3.3.3 Characterisation Discussion

The introduction and application of this volatility characterisation protocol have allowed for the identification of the key factors that dictate DUT relaxation over time. Importantly, the findings that have thus far been presented can be linked to the principal goals of this project, i.e. consolidation in multiple timescales.

To begin with, it has been shown for the first time that RRAM volatility can unfold bidirectionally (see Fig. 3.5). Effectively, an artificial synapse can deviate away from a more stable state, i.e. $R_{pre}$ in either direction. In memory terminology, one can think of $R_{pre}$ as a more rigid consolidated state and volatile $\Delta R_{start}$ a plastic change in efficacy which stores new information only in the short term. Moreover, it has been suggested that $V_P$ can be manipulated to modulate the strength of a new memory event. The number of pulses N further dictates both the switching and relaxation behaviours of DUT and it is a possible option for encoding time-related memory events (for instance, higher N values if the information is available for longer time periods to the learning agent). It is interesting to note that within the chosen retention window T = 2 minutes, a small but effective accumulation of non-volatile residues has been observed to push DUT away from its $R_{initial}$. With reference to the consolidation mechanisms that are discussed in 2.2, successive programming cycles can be utilised as a memory reinforcement mechanism, while the gradual drift away from $R_{initial}$ can act as a transfer mechanism from short-term to long-term memory. The next steps in this project lie in identifying how the mentioned long-term change in a synapse's state can coexist with additional volatile changes that effectively store information in a time-dependent palimpsest fashion.

Data from Fig. 3.8 already indicate that the timescales of plastic memory storage can be altered by the level of invasiveness by which a new memory event is stored in the synapse. However, there is still a lack of a systematic approach towards modelling the dependencies of RRAM volatility within a wide state space. Indeed, the next steps in designing the desired technology lie in accurately simulating RRAM volatility within a device's operating resistive range as well as a range of biasing conditions. This is essential since the identification of mathematical dependencies can firstly enable a component analysis on these devices and hence, secondly lay the groundwork for new application-oriented experimental work. Thus, the next section of this report deals with modelling volatility within a predefined range of the R(t) parameters that are defined in Eq. 3.1.

## 3.4    Bidirectional Volatility: Modelling

### 3.4.1    Modelling State Space and Methodology

The volatility characterisation protocol that has been introduced in 3.3 has been re-employed in order to gather sufficient data for modelling RRAM relaxation characteristics. Specifically, the complete range of programming pulses used is N $\in [1, 1000]$ and absolute programming amplitudes $V_P = \{3, 5, 6, 7, 9\}$V are used for both polarities. The device under test (DUT) that has been characterised in 3.3 has also been used in this section along with two new, randomly chosen devices; for simplicity, these will be referred to as DUT1, DUT2 and DUT3.

The only difference in this modelling work is the chosen value for T which will be equal to 60s instead of 120s. This choice has been made for two reasons. Firstly, it allows for a 50% reduction in the experimental run-time which as is, has still required more than 40 hours of consecutive measurements for all 3 devices. Moreover, as it has been observed in 3.3, volatility unfolds in a combination of exponential and fast changes along with slower and saturating transient phenomena. Hence, by T = 1 min most of the characteristic change in R has already occurred. At the same time, this T value may also result in a higher accumulation of non-volatile residues which would be crucial to quantify.

Retention data have been modelled using Eq. 3.2, which is shown once more below:

$$R(t) = \alpha \exp^{-(\frac{t}{\tau})^\beta} + \gamma$$

The dependencies of each parameter of R(t) are analysed separately. For simplicity, only results from positive programming cycles are presented, although the corresponding dependencies hold true in both polarities and the effectiveness of our model for bidirectional volatility is shown at the end of this section. For illustrative purposes, only device 1, DUT1 is analysed explicitly and will be referred to as DUT. This work considers individual fittings unsuccessful either if the resulting relaxation time constant is greater than T or if the corresponding $R^2$ value is less than 0.1.

### 3.4.2    Effects of Stimulation Regime on Modelling Parameters

#### 3.4.2.1    Switching Component

The maximum volatile switching that is observed immediately after a device has been stimulated, as it has been expressed by $\Delta R_{start}$ in 3.3 can be indirectly reflected from the

FIGURE 3.9: Switching component $\alpha$, as extracted from Eq. 3.2 for positive stimulation. All data values are shown in 3D space against $R_{pre}$ and N in (a). For a clearer observation of the parameter's dependency on N and $V_P$, $\alpha$ is plotted only against the former in (b). For each distinct programming amplitude, the change of $\alpha$ can be modelled using a combination of two linear components (inset equation in (b)) with very little error as per (c). Each component of $\alpha(N)$ then follows a loose linear relationship with increasing $V_P$ which is depicted in (d-f).

$\alpha(N, V_P)$ parameters of the stretched exponential function, since from Eq. 3.2, R(t=0) = $\alpha$. The extracted $\alpha$ values from the successful retention fits can be examined in Fig. 3.9.

Specifically, Fig. 3.9a sees the parameters unfold at different areas of the $R_{pre}$ - N plane for several $V_P$ values. As has been the case with $\Delta R_{start}$, it is vague whether $R_{pre}$ directly affects switching, as shown in 3.9a. The projections on the N - $\alpha$ plane are plotted in Fig. 3.9b. The effect of increasing $V_P$ is apparent, especially for larger N values, where it leads to larger switching that may reach up to 1 order of magnitude in difference (3V vs. 9V).

For each unique programming amplitude, $\alpha$ appears to follow a steep linear path for N ≤ 100, followed by a more saturated linear change for larger pulse numbers. The relationship resembles that of exponential decay, as it was also seen in Fig. 3.7, but it is more effectively described as a combination of two linear components. Eq. 3.3 fits $\alpha$ using the same number of parameters and also boasts a low mean error, especially in low N values. This error is shown in Fig. 3.9c.

$$\alpha(N) = \frac{\alpha_\alpha \cdot N \cdot (\alpha_\beta \cdot N + \alpha_\gamma)}{\alpha_\alpha \cdot N + (\alpha_\beta \cdot N + \alpha_\gamma)} \tag{3.3}$$

An increased error spread for smaller N can be explained by the fact that the corresponding programming cycles are not invasive enough to produce clear relaxation phenomena. In fact, retention cycles induced either via smaller N or $V_P$ values are expected to present the greatest challenge when it comes to modelling volatility signatures.

After modelling α(N) using Eq. 3.3, each individual parameter $\alpha_{\alpha,\beta,\gamma}$ is plotted against $V_P$ in subfigures 3.9d-f for all 3 devices under test. It is shown that all subjects follow the same patterns. The negative value of $\alpha_a$ decreases linearly with $V_P$ , indicating that the effect of smaller N values becomes steeper at higher voltages. Contrarily, $\alpha_\beta$ increases with $V_P$ which reflects the decrease in the steepness of Eq. 3.3 when larger numbers of pulses are employed. This could occur because more invasive programming cycles push the DUT closer to its resistive boundary and limit its capacity to switch further. Lastly, $\alpha_\gamma$ decreases with increasing voltage which accounts for an overall vertical offset between each curve in 3.9b.

Every parameter $\alpha_x, x \in [\alpha, \beta, \gamma]$ of Eq. 3.3 can be modelled for $V_P$ as follows:

$$\alpha_x(V) = \alpha_{x0} \cdot V + \alpha_{x1} \tag{3.4}$$

The results suggest that switching is mostly a function of N and $V_P$. However, $V_P$ related discrepancies are only clear given sufficiently large N values, i.e. at least a few tens of pulses. This suggests that the devices' switching depends heavily on the time for which stimulation is present. The main effect of $V_P$ is then the fact that it amplifies switching and allows DUT to be pushed to further resistive states, as it has been observed in non-volatile cases [75].

#### 3.4.2.2   Relaxation Time Constant

When the devices are left to relax post-stimulation, the relaxation time constant $\tau$ is a direct indicator of $\frac{dR}{dt}$. As it has been seen in Fig. 3.8, the more invasive the stimulation, the longer it takes for DUT to come close to its corresponding $R_{pre}$. From a memory consolidation perspective, stronger memory events should induce efficacy changes that are preserved for longer timeframes. The data presented in Fig. 3.10 show this variability can be quantified in the technology under test.

The time constant of the stretched exponential function is plotted against N in Fig. 3.10a, for all $V_P$ sets. Contrary to intuition, $\tau$ does not appear to change monotonically with N. Indeed, datapoints that correspond to the smallest N values accumulate at the

FIGURE 3.10: A linear dependence between the relaxation time constant $\tau$ and N is observed for N > 100 in (a) for positive $V_P$ values. The anomalies in $\tau$ values for N $\leq$ 100 are considered to be the results of unsuccessful fittings of Eq. 3.2. Individual parameters of $\tau(N)$ are subsequently plotted against $V_P$ where an increase with programming amplitude is also observed.

largest $\tau$ values (around 3s) and decrease until N $\approx$ 100. At that point, $\tau$ exhibits a linear overall increase with the number of pulses. Moreover, as $V_P$ increases an overall upward shift is shown in this trend which is accompanied by a minor increase in the linear slope. It is assumed that $\tau$ values for small N are artefacts of unsuccessful fittings on the stretched exponential, caused by insufficient DUT stimulation as discussed extensively in 3.3. The linear relationship that is observed between $\tau$ and N > 100 can be modelled using Eq. 3.5.

$$\tau(N) = \tau_\alpha \cdot N + \tau_\beta \tag{3.5}$$

Subfigures 3.10b-c show how the linear components of $\tau(N)$ change with $V_P$ for all DUTs. The slope component $\tau_\alpha$ increases with $V_P$ however the values themselves are very small ($\times 10^{-2}$). The upward shift in the $\tau(N)$ function is reflected from the linear increase of $\tau_\beta$ with $V_P$. Again, the results are consistent across all 3 devices. It can be concluded that more invasive programming stimuli induce volatile changes that last longer. This effect holds both when it is only manifested through N, and when $V_P$ is increased since the latter amplifies the effect of the former.

So far, it has been demonstrated that both DUT switching and volatility time constants can be altered by appropriately manipulating the biasing parameters used for stimulation. Importantly, the fact that relaxation characteristics can be - to some extent - engineered to desire by $V_P$ or N independently allows more flexibility in the design and optimisation of RRAM-based synapses.

FIGURE 3.11: The stretch factor of R(t), β, undergoes a rapid exponential decay with increasing N values in (**a**). Moreover, this trend shifts to lower overall values as $V_P$ increases. Again, a high error margin in the fittings is apparent for small N values due to the noise produced by less pronounced volatility retention cycles (see inset). The parameters of β(N) are plotted against $V_P$ in subfigures (**b-d**). While the datapoints mostly agree on a trend for $V_P$  3V, the results are less clear for the lowest amplitude, which again could be explained by the unreliable volatile results in programming cycles of low invasiveness.

### 3.4.2.3   Stretch Factor

One of the advantages of the stretched exponential function when it comes to modelling R(t) is that its factor $\beta \in (0, 1]$ allows relaxation to be stretched over the time domain. Once again, β is evaluated in the same way as in previous sections, and it is plotted against N in Fig. 3.11a.

Here, it can be observed that programming phases which employ the smallest N values result in less stretched decays, or even fully exponential ($\beta = 1$), irrespective of $V_P$, as this is apparent both for 3V and 9V. However, the stretch factor decays rapidly as N increases until it saturates around $N \approx 200$ for all amplitudes. Moreover, the trend for β is almost identical for all amplitudes except for $V_P = 3V$, where the saturation point is higher. This decaying relationship has been modelled using Eq. 3.6.

$$\beta(N) = \beta_\alpha \cdot e^{-\frac{N}{\beta_\tau}} + \beta_\gamma \tag{3.6}$$

The inset in Fig. 3.11a shows the fitting error of the Eq. 3.6 which, for the same reasons as with the other stretch exponential parameters is higher in the lesser invasive retention cycles. Individual parameters from the fitting function have been plotted again against $V_P$ for all DUTs in Fig. 3.11b-d. They can be expressed mathematically in terms of $V_P$ as a 1st order polynomial (Eq. 3.4). The decrease of $\beta_\alpha$ indicates that higher amplitudes result in more stretched relaxations even for lower N values. Furthermore,

FIGURE 3.12: The offset between each retention cycle's $R_{pre}$ and the DUT's projected saturation point $\gamma$ is hinted to depend on the index of successive cycles, or loops. However, such a trend is most clear in the case of $V_P = 9V$ and can be modelled with a logarithmic function (see subfigure (**a**)). A vague tendency for a negative average offset is shown for initial loops, which however crosses to positive as more retention cycles are employed. The positive offset can explain that it becomes progressively more difficult to keep DUT further closer to its resistive boundary. (**b-c**) verify the relationship for all devices under test.

$\beta_\tau$ remains constant for all amplitudes close to N = 50, signalling a very fast saturation, which is also evident in subfigure a. Note that there is increased variation in the case of $V_P = 3V$ between different devices. However, while discrepancy exists, all cases have values below N = 200 which are in line with the saturation observed in subfigure a. Finally, the decrease in $\beta_{\gamma_0}$ describes how more invasive amplitudes lead to more stretched versions of R(t).

#### 3.4.2.4 Saturation Offset

This modelling approach does not guarantee any extrapolation of R(t) outside the time window T that retention has been observed. However, the approximation of $R_{end}$ largely depends on the projected saturation point of the stretched exponential function, $\gamma$ (note that R(t → ∞) →$\gamma$ ). Interestingly, when initially observed, $\gamma$ seemed to depend solely on the value of $R_{pre}$, with the former changing linearly with the latter and varying within a ± error margin. However, it has later been shown that the sign of the difference between $R_{pre}$ and $\gamma$ plays a significant part in the evolution of R(t) in each retention cycle. Specifically, whether the projected offset = $\gamma$ - $R_{pre}$ , is positive or negative determines the evolution of R(t) and the accumulation of non-volatile residues in successive retention cycles.

It has been found that the offset only loosely depends on the cycle index within each experiment as shown in Fig. 3.12a. Whether DUT is projected to saturate with a non-volatile residue or reach or even surpass its original value seems to depend on the history of previous stimulations. Generally, within a set of consecutive retention cycles, early stimulation events drive DUT away from its initial state and a non-volatile residue is left within T. However, as progressive stimulations push DUT further towards its resistive boundary, DUT exhibits increasingly strong tendencies to return to its original initial state. This implies that the sign of $\gamma$ - $R_{pre}$ may change as the simulation progresses. This phenomenon requires further, dedicated study. Effectively, this suggests that the accumulation of the non-volatile component "lags" the accumulated change in successive $R_{pre}$ values. In other words, as we continue to apply successive stimulation blocks the final state of the retention state increasingly tends to move towards an $R_{end}$ value that is higher than the $R_{pre}$ of the corresponding run.

The actual factors that form this dependency are not yet fully understood. Moreover, the only strong indications that this trend exists are most prevalent for $V_P$ = 9V. There, the observed $R_{end}$ in early cycles is up to 8 kOhms less than the corresponding $R_{pre}$. However, this offset rapidly decays to 0 and later increases with a decreasing rate as cycles progress. This relationship has been modelled using a logarithmic function as shown in Eq. 3.7 and is shown in Fig. 3.12a. The logarithmic relationship has been chosen because it encapsulates both the early offset convergence to 0 as well as its slowing increase when larger than 0 in later stages.

$$offset = \gamma - R_{pre} = \gamma_\alpha \cdot \ln{(cycle)} + \gamma_\beta \tag{3.7}$$

The fitting parameters of the logarithmic function are plotted against $V_P$ for DUT1-3 in Fig. 3.12b-c. The fact that $\gamma_\alpha$ increases exponentially with $V_P$ verifies that the relationship is mostly present in invasive voltage amplitudes. Parameter $\gamma_\alpha$ has been modelled using Eq. 3.8. The fact that larger $V_P$ values can cause the accumulation of larger non-volatile residues can be reflected from the linear decrease of $\gamma_\beta$ as $V_P$ increases. The parameter has been modelled using Eq. 3.4.

$$\gamma_\alpha(V) = \gamma_{\alpha 0} \cdot e^V + \gamma_{\alpha 1} \tag{3.8}$$

While the logarithmic function models the observed behaviour to an extent, it should be noted that there is still more ground to cover in truly understanding volatility saturation. Specifically, the observed interplay between the offset and retention cycles could be a contaminated reflection of a hidden global relaxation. If that would be the case then the offset would not depend on the retention cycle number per-se, but rather on a timing metric that deals with how long the device has been left to relax. And if that is true, then similar logarithmic relationships between offset and cycle numbers could

FIGURE 3.13: The performance of the model is verified via simulation for both positive and negative programming cycles. Raw data from successive retention cycles that employ increasing values of N are shown for different amplitudes. In both positive (a-b) and negative (c-d) polarity cases, the retention data are fitted adequately by the model which is shown in black.

occur from different initial resistive states. This however is a speculation that requires further independent studying.

### 3.4.2.5 Model Simulation

Thus far, results have only been shown for modelling the positive bias induced volatility. However, volatility in the opposite polarity can be studied using the same procedure. All the necessary parameters used to model volatility in DUT1 are included in Appendix A.

The extracted model parameters have been used to simulate DUT1 volatility bidirectionally and at different voltage amplitudes. Raw data from consecutive retention cycles employing N = 400, 500, 600, 700 pulses at $V_P = \pm$ 3, 7, 9V are shown in Fig. 3.13a-b/c-d for positive and negative polarities respectively. Relatively large values of N have been chosen to ensure that the relaxation presented unfolds clearly enough. Moreover, the two limit values of $V_P$ have been chosen along with an intermediate 7V that is optimum for visibility purposes. The simulated retention curves sit on top of the retention data and can be seen in bold black lines.

In the case of $V_P$ = 3V, an increase in N results in more pronounced switching but the time constants for all cycles remain fast and limited by the small programming amplitude, as expected from Fig. 3.10. Contrarily, the simulated R(t) for $V_P$ = 9V reflects both

more prevalent volatile switching as well as relaxation over larger time constants. In addition, an inspection of simulations for negative stimuli reveals both a lower sensitivity of RRAM to increasing programming amplitude as well as significantly slower transient responses, characteristics in line with the volatility characterisation results presented in 3.3.

### 3.4.3   Modelling Discussion

The results that have been obtained in this modelling study allow for better quantification of the key factors dictating the devices' relaxation characteristics. This is a significant step towards understanding volatile RRAM families alone, but it also adds further insight into the development of metaplastic synaptic circuits.

Due to the quantification of DUT bipolar switching, biasing parameters can now be adjusted to manipulate potentiation and depression events so that they produce symmetric changes in synaptic efficacy. At the same time, the asymmetric relaxation time constants that occur for the two polarities pose a significant challenge in engineering identical behaviour for both types of write events. Specifically, in order to produce equal plasticity changes in both directions a tradeoff has to be paid in much longer volatility timescales after depression events. Vice versa, should the application focus on identical timescales, then switching will need to be unbalanced. This tradeoff has to be examined thoroughly in order to understand what is best for the purposes of consolidation. Lastly, unexpected mechanisms have been revealed concerning the non-volatile residues that result from consecutive retention cycles. The logarithmic relationship between DUT's saturation offset and the sequencing of unipolar programming phases could be utilised in order to emulate plasticity changes that are induced in a progressively harder way, much like the saturating boundary approach of computational synaptic models [25].

Nevertheless, the capabilities of this modelling method are crucially limited by certain factors that should not be overlooked. Firstly, there is no guarantee for the extrapolation performance of the stretched exponential function, either in the time domain or in the state-input space. Indeed, reliable extrapolation outside the observed retention window is still unknown. Moreover, if the volatile operation needs to be optimised for minimised energy consumption, then observation of the relaxation dynamics, caused by less invasive stimuli, alone remains a major challenge that has to be overcome.

## 3.5   Thermal Effects on RRAM Volatility

The importance of developing the described characterisation methodology is evident as the routine enables a deeper understanding of volatile operations. Since the methodology is agnostic to operational parameters, volatility dependencies under an expanded range of input conditions can be examined. As such, this project has also been able to consider how environmental factors may affect RRAM volatility and specifically, how initial volatile responses depend on a device's operating temperature.

This brief case study focuses on the same data-driven approach that is described in 3.3 and aims to explain how the relative magnitude of the initial volatile change $\Delta R_{start}$ and the relaxation time constant $\tau$ are affected by temperature. It is important to note that all the application demonstrations presented in this thesis are based on normal room-temperature environments. However, understanding the discrepancies induced by temperature is an essential future-proofing step since modern integrated electronic systems are inevitably subject to operational heat cycles.

### 3.5.1   Thermal Volatility State Space and Methodology

The experimental variables have been altered from 3.3 to focus more on temperature rather than stimulation intensity. However, volatility has again been characterised under different stimulation amplitudes since it has already been since that increasing $V_P$ can induce significant shifts in the baseline behaviour of volatile RRAM. The resulting experimental methodology is as follows.

Consecutive retention cycles have again been employed to uncover volatility under continuous stimulation. Here, each programming phase employs N = 500 identical pulses of 100μs width, evenly spaced 1ms apart. This choice for N has been made since it allows for sufficient volatile switching without being necessarily making N itself the primary effect driver (see Fig. 3.7). Each cycle uses progressively more invasive programming amplitudes in the range $V_P = \pm [1.5, 5.0]$V, while volatility is observed for a total window T = 1 minute. Each read phase uses non-invasive pulses at 0.2V. This experiment has been repeated at 4 different operating temperatures $\{22, 40, 55, 75\}°$C on 3 different on-wafer devices (Au/$TiO_2$/Pt with 20nm active oxide thickness). Temperature levels have been controlled using a thermal chuck in full contact with the test wafer, while all other experimental parameters have been kept constant. Temperature levels have been swept randomly and between each configuration, each DUT has been allowed to return to room temperature ($\approx 22°$C) to minimise potential history biases.

The experiment has been conducted for a minimum of 3 times on each DUT. As mentioned, relative volatile changes are considered with respect to $R_{pre}$ for each retention

FIGURE 3.14: (**a**) 3D surface showing the effects of operating temperature and programming amplitude on the relative volatile switching % $\Delta R_{start}$ for positive stimulation. The two parameters are isolated in projections for (**b**) switching-temperature relationship and (**c**) switching-amplitude relationship. Respective results for negative stimulation are shown in (**d-f**). Resistance changes are normalised with respect to $R_{pre}$.

cycle. This ensures normalisation across different temperatures which induce baseline changes in R irrespective of any further stimulation. Hence, this study considers $\Delta R_{start}(\%) = \frac{\Delta R_{start} - R_{pre}}{R_{pre}} \times 100\%$. Additionally, Eq. 3.2 has again been utilised to extract $\tau$ values for each cycle. The mean result values across all experiments are presented in the following sections.

### 3.5.2   Thermal Volatility Results and Discussion

By observing the surface plot at Fig. 3.14a and its projections against $V_P$ (c), the relative volatile change induced increases linearly with $V_P$, which is consistent with the results at Fig. 3.6. Interestingly, temperature (see subplot b) affects switching more subtly, by tuning its variance across increasing $V_P$. More specifically, increasing temperatures limit the range of relative switching by increasing the effect of lower amplitudes and decreasing the effect of larger ones.

The pattern is similar for the relationship between $V_P$ and $\Delta R_{start}$ (%) for negative stimulation but opposite with respect to temperature (see Fig. 3.14d-f). Here, temperature

FIGURE 3.15: (**a**) 3D surface showing the effects of operating temperature and programming amplitude on volatility relaxation time constant $\tau$ for positive stimulation. The two parameters are isolated in projections for (**b**) $\tau$ - temperature relationship and (**c**) $\tau$ - amplitude relationship. Respective results for negative stimulation are shown in (**d-f**). Resistance changes are normalised with respect to $R_{pre}$.

enhances switching at more invasive $V_P$ values but interestingly, it seems to play little part for lower absolute amplitudes.

The extracted $\tau$ values have been extracted from all retention cycles and are presented in a similar manner in Fig. 3.15. The results in positive stimulation polarity agree in principle with how volatility has already been characterised with respect to $V_P$. Specifically, more invasive stimulation leads to relatively longer time constants, stemming from more pronounced volatile responses. The results in (b) show little evidence to support a relationship between $\tau$ and temperature although there is some hint of similarity to the observed effects on $\Delta R_{start}$ (%), with increasing temperature effectively limiting the range of $\tau$ at more invasive stimulation.

However, these relationships are less clear in the case of negative stimulation (Fig. 3.15d-f). There, a saturating increase in $\tau$ is observed as $|V_P|$ increases, which is however highly affected by the operating temperature of the DUT. As temperature rises, the dependence of $\tau$ on $V_P$ becomes minimised but also drops significantly in absolute values. This is counterintuitive when considering the inverse effect of $\Delta R_{start}$ (%) and temperature that is seen in Fig. 3.14.

By assessing the results from Fig. 3.14 and Fig. 3.15, a clear difference is noted for thermal effects on volatility towards opposite directions. This is first evident by noticing how temperature affects $\Delta R_{\text{start}}$ (%) differently in positive and negative stimulations. At positive bias, lower temperature enhances volatile switching while at negative it clearly inhibits volatility. Conversely, higher temperature seems to limit $\Delta R_{\text{start}}$ (%) at positive regimes and enhance it at negative programming cycles. This discrepancy suggests again the possible existence of two distinct mechanisms that govern volatility in each stimulation polarity, something that can also be argued by inspecting the differences in switching sensitivity shown in Fig. 3.6. Nonetheless, the value ranges for $\Delta R_{\text{start}}$ (%) in both polarities remain similar which would allow for consistency in volatile switching application across RRAM heat cycles with the appropriate fine-tuning of other stimulation parameters such as $V_P$ and/or N.

The eventuality is further suggested by the complete dissimilarity in how volatility's time constant and temperature are interconnected across positive and negative stimulation. $\tau$ cannot be equilibrated across polarities via temperature, with values at positive stimulation being one order lower than in negative. In general, temperature-induced effects also appear to be enhanced at more invasive biasing regimes. Contrarily, at lower voltages, see $\pm 2$V, the results show little variability stemming from less pronounced volatile responses in the first place. Given that this thesis will later consider translating biasing regimes of variable intensity into synaptic plasticity protocols, it has been concluded that temperature will not be utilised as a volatility manipulator to achieve more stable results. Thus, all application-oriented experiments and results shown in this thesis will assume room operating temperature $\approx 22°$C.

## 3.6   Summary

Memristive volatility has traditionally been perceived as a drawback or an obstacle, rather than a feature. Consequently, while it has been studied before, its true potential has been overlooked. This chapter has presented an extensive data-driven study, designed to uncover the governing characteristics of RRAM volatility, which importantly, has been demonstrated bidirectionally for the first time. The results shown here set a framework for characterising and modelling volatility, under various stimulation conditions, observation timeframes and operating temperatures. Naturally, this methodology is application-agnostic and can be useful for any technology design with the need for intrinsic time referencing.

In the context of this thesis, bidirectional volatility has shown great potential in emulating time-dependent plasticity dynamics in a twofold way. With regard to the palimpsest memory consolidation, the identification of two volatile timescales is a promising pathway towards the emulation of plastic and rigid memories in a single artificial synapse.

Plasticity characteristics such as volatile switching and forgetting rate can be fine-tuned with stimulation parameters so that we have balanced LTP and LTD. The non-volatile residue can be indirectly changed in the long term with progressive stimulation towards any direction. Moreover, volatility can be used to emulate the high- and low-pass frequency filters that are associated with short-term plasticity. These mechanisms show strong potential in the domain of neuronal activity detection since volatile RRAM could be fine-tuned to filter low or high spiking frequencies.

It should be noted that all the work that is presented in the next chapters of this thesis has been conducted at operating room temperatures, stemming from the complex relationship between volatility and temperature.

# Chapter 4

# Palimpsest Memory Consolidation

## 4.1 Introduction

This chapter presents a proof of concept demonstration of palimpsest consolidation in hardware synapses. The rapid growth of AI and particularly deep learning, have increasingly demanding needs for computational resources, without however providing the same learning capacity that is met in biological systems [9], [10]. Palimpsest operation allows multiple memories to be written sequentially at different timescales and without causing catastrophic deterioration of previously stored counterparts [25]. This translates into more generalised learning and therefore provides a solution to the mentioned capacity discrepancy. Here, it is shown how emerging volatile RRAM technologies can emulate the processes that drive biological palimpsest consolidation under a variety of learning scenarios. It is found that this technology



FIGURE 4.1: Chapter's outline and objectives.

| RRAM | Synapse | Memory |
|------|---------|--------|
| Resistance (R) | Weight ($w_i$) | Memory state |
| Stimulation ($\pm V_P$) | Plasticity | Memory write |
| Switching ($\Delta R_{start}$) | Plastic $\Delta w_i$ | Plastic consolidation timescale |
| Non-volatile residue ($\Delta R_{end}$) | Rigid $\Delta w_i$ | Rigid consolidation timescale |
| Volatility time constants ($\tau$) | Weight decay | Memory degradation |

TABLE 4.1: Translation from device to synapse to system level. The pillar concepts of RRAM operation are presented with their corresponding representations in the synapse which in turn dictate the overall performance of the memory system.

aids significantly the learning flexibility of memory networks which showcases how RRAM synapses can expand the capabilities of AI hardware towards more generalised applications.

**This chapter is organised as follows**: Section 4.2 details a framework for mapping the volatile properties of single RRAM devices onto synaptic plasticity operations. Section 4.3 presents the fundamental demonstration of palimpsest consolidation in hardware. Section 4.4 extends the scale of a memory network comprised of RRAM synapses and evaluates its performance under continuous learning of random and uncorrelated memory signals. Section 4.5 adapts the functionality of the studied memory networks in the context of a visual working memory for binary images. Section 4.6 analyses the key findings of this work and section 4.7 summarises and concludes this chapter.

## 4.2   System Framework

The first step towards a practical demonstration of how memristive synapses can facilitate memory consolidation has been the layout of an appropriate operational framework. Specifically, the thus-far agnostic volatile RRAM properties need to be mapped to equivalent neuromorphic functions to allow the devices to emulate biological synapses. This mapping is outlined in Table 4.1.

As thoroughly discussed in Chapter 3, the technology studied in this project is characterised by 4 key controllable parameters. First, the visible state R of the memristive synapses can be recorded at any time. This state equates to the weight/efficacy $w_i$ of a synapse *i*, which in turn reflects the state of memory at a system level. By appropriate stimulation of invasive positive or negative programming pulses, the synapse can undergo plasticity changes of potentiation and depression respectively. Plasticity is thus the main facilitator of memory writing in the studied system. The main novelty of this approach is the exploitation of all volatile parameters after stimulation. The initial volatile switching $\Delta R_{start}$ can effectively be used as a plastic mechanism for fast writing of a memory state which however decays over time. Concurrently, the state of $\Delta R_{end}$ reflects a rigid accumulation of $\Delta w_i$ which cannot easily be overwritten. While

the plastic component can flexibly be used to overwrite an existing memory state without a serious penalty over time, it is the progressive shift in the rigid component that results in memory consolidation. This is perfectly analogous to biology where memories are only protected in the long term given sufficient stimulation. Finally, since the volatile time constant $\tau$ dictates how quickly R transitions closer to $R_{end}$, it also dictates how quickly temporary memory overwrites decay in the memory system. Using this framework, one can finally obtain a true outline for palimpsest consolidation in RRAM.

Having laid this foundation, this thesis can now introduce the base operations that entail a memristive synapse. It is important to note that while R can be read in an analogue way, this demonstration takes a reductionist approach and considers binary synaptic operation only. This decision can be justified for two reasons. First, it is the simplest mode of operation and it is imperative to verify that palimpsest consolidation can be implemented in this fashion before further extending this investigation. Second, there is sufficient information to suggest diminishing returns in memory performance if memory resolution is extended beyond binary states (for details see Chapter 2).

To emulate the binary synaptic network, single volatile devices have been used following the mapping in Table 4.1, where R is read in an analogue fashion and binary states are determined with reference to a predefined threshold state $R_{thres}$. Specifically, plasticity events are implemented using Eq. 4.1 while the binary weight is determined using Eq. 4.2.

$$\text{plasticity event} = \begin{cases} \text{potentiation} & \text{if } V > 0 \\ \text{depression} & \text{otherwise} \end{cases} \tag{4.1}$$

$$w = \begin{cases} 1 & \text{if } R < R_{thres} \\ 0 & \text{otherwise} \end{cases} \tag{4.2}$$

Moreover, it has been essential for the analysis of the following application case studies to quantify the noise levels carried from the noise that is inherent in volatile R(t) over to binarised states. This is especially important for simulations of larger memory networks. To do so, the transient changes in the analogue synaptic states have again been modelled using the stretched exponential relationship for volatility, introduced in Eq. 3.2. Using this model's fittings on real device data, noise has been calculated as the percentage (%) difference between the ideal and real data (see Eq. 4.3).

$$(\%)\Delta R = \text{noise} = \frac{R_{ideal} - R}{R} \times 100\% \tag{4.3}$$

As it will be discussed later on, RRAM noise has been found to follow a Gaussian (normal) distribution. The distribution's mean value $\mu$ and standard deviation $\sigma$ can

then be extracted such that noise can be included stochastically in the analysis of the following experiments using Eq. 4.4.

$$p(\text{noise}) = \frac{1}{\sqrt{2\pi\sigma^2}}e^{-\frac{(\text{noise}-\mu)^2}{2\sigma^2}} \tag{4.4}$$

The rules defined in this section apply to the rest of this chapter. The next sections present the results obtained from various application studies and assess the capacity of volatile RRAM synapses to support palimpsest memory consolidation.

## 4.3   Consolidation Demonstration in Memristive Synapses

### 4.3.1   Memristive Synapse Set-up

Before conducting consolidation experiments, it has been essential to configure the corresponding stimulation and reading parameters, such that the memristive synapses will perform optimally with respect to case-specific goals. For reference, this thesis has had access to $Pt/TiO_2/Au$ volatile RRAM devices. Pt and Au refer to the $20\times20\mu m^2$ bottom and top electrodes respectively and the thickness of the active oxide layer is $25nm^2$. Initially, simple I-V measurements have been made on the device samples to identify the basic operation regime (see Fig. 4.2). Based on these results it has been decided that the passive read voltage $V_R$ will be set to 0.5V. This is appropriately tuned so that the synaptic reads are non-invasive but also lie above the measurement noise threshold.

Next, a simple parametric sweep study has been conducted to examine the dependence of device volatility on stimulation amplitude $V_P$. Specifically, this study has considered single pulses with $V_P$ values in the range $\pm[5,8]V$ and a 100$\mu s$ pulse width. The results are shown in Fig. 4.3 and include the analogue device history during stimulation and the biasing protocols in use. The resistance values caused by stimulation (equivalent $R_{start}$ values from Chapter 3) are shown in red, while the datapoints collected during relaxation are shown in blue. A choice of $V_P = \pm7V$ at 100$\mu s$ width has been made so that large enough volatile efficacy changes are induced and thus binary state transitions can be observed in the short-term from potentiated states to depressed and from depressed to potentiated. Finally, extensive stimulation experiments have been conducted using this basing regime to ensure the overall endurance and reliability of this technology. These results are shown in Appendix B.

These results support $TiO_2$ RRAM as a prime candidate for palimpsest consolidation due to the devices' intrinsic controllable and bidirectional volatility. Thus, palimpsest phenomena can finally be demonstrated in hardware on a single memristive synapse

FIGURE 4.2: I-V measurements on a $TiO_2$ memristive synapse. The absolute values of reading current are shown in a log scale. The read voltage value for this experiment is chosen at 0.5V.

and are illustrated in Fig. 4.4. In this study, the aim has been to consolidate by repetition a binary state $S_1$ via long-term potentiation (LTP) and then overwrite it with the competing state $S_0$ in the short term. The stimulation profile used is this described in 4.3.1, with write events being spread at 30s apart.

The resulting efficacy and stimulation histories are shown in Fig. 4.4b-c. There, 3 unique consolidation stages are distinguished. These stages are directly relatable to the beaker theory by [25] and the analogy is illustrated in Fig. 4.4d. Stage 1: Potentiation events induce volatile changes which drive R below a defined binary threshold $R_{thres}$ in a fast decaying manner. $S_1$ is thus expressed in the synapse but only in the short-term. The non-volatile residue remains above $R_{thres}$ and thus $S_0$ is reinstated as the synapse's long-term memory state. Stage 2: Further reinforcement of $S_1$ occurs via successive potentiation events. This accumulates changes in the hidden non-volatile residue which now lies below $R_{thres}$ and hence, $S_1$ is consolidated in the long-term. Stage 3: Depression plasticity events cause R to increase but rapidly decay. R crosses $R_{thres}$ in the opposite direction which translates to the expression of $S_0$ in the short-term, before $S_1$ is reinstated again, owing to its previous consolidation.

In this scenario, the memristive synapse can protect a hidden consolidated memory ($S_1$) while expressing another opposite memory ($S_0$) temporarily atop it. This is an effective doubling of memory capacity, distributed across the two volatility timescales and also a first demonstration of palimpsest consolidation in hardware. Some interesting properties of this technology should be mentioned. First, potentiation events during Stage

FIGURE 4.3: Examination of induced volatility by single pulse stimulation of varying under potentiation (**a-b**) and negative polarity (**c-d**), with $|V| \in [5-8]$V. The state recorded following stimulation and the passive relaxation datapoints are shown in red and blue respectively in the resistance plots. Stimulation levels at $\pm$ 7V have been chosen for the hardware demonstration of palimpsest consolidation due to their strong corresponding volatile phenomena.

2 have diminishing effects in consolidating $S_1$ further. This is expected since RRAM switching is known to operate within soft resistive bounds [75]. This means that successive stimulation eventually fails to induce a stronger consolidation of one state since the synapse's analogue state is not entrenched further away from $R_{thres}$. This bounded synaptic efficacy is known to aid the capacity of memory networks since it regulates asymmetric stimulation [25], and here it is intrinsic at the device level. Moreover, it is seen in Stage 3 that opposing memory signals inevitably regress the synapse's analogue state towards $R_{thres}$. The synapse necessarily trades some rigidity in the slow timescale to benefit from increased plasticity in the short term. Finally, the artificial synapse shows asymmetric behaviour in potentiation and depression events. However, it has been used in this study because its high ratio between volatile and non-volatile plasticity changes showcases clearly palimpsest properties.

FIGURE 4.4: Demonstration of a palimpsest memristive synapse. (**a**) Real images of volatile RRAM. (**b**) Typical operation of a binary volatile synapse. Consecutive plasticity events (potentiation) produce pronounced changes in synaptic efficacy which are short-lived (fast timescale - short-term memory) and accumulating, smaller and more stable changes (slow timescale - long-term memory). The applied stimulation pattern is shown in (**c**). (**d**) Schematic association of the three consolidation stages to the beaker theory. The evolution of the first beaker's liquid level (synaptic efficacy) is determined over time by the state of the hidden second beaker.

## 4.3.2 Consolidation of Fully Destructive Memories

The next logical step in assessing the synapses' capabilities is to integrate them into a basic memory network and consolidate competing signals of a larger size. This study examines the network's performance in the worst-case scenario of two antipodal, fully competing memories. Signals $M_1$ ($\triangledown$ = [101100]) and $M_2$ ($\blacktriangledown$ = [010011]) have been presented in memory for 15 and 3 consecutive events respectively. This is is the same stimulation pattern shown in Fig. 4.4; again, memory signals are distributed 30s apart from each other. Each bit $b_n$, $n \in \{0,1,...5\}$ of memories $M_{1,2}$ is written on an individual synapse $w_n$.

This has been chosen a posteriori to fit best with the history of all synapses' states. The raw data from all synapses are shown in Fig. 4.5. As seen, not all device samples show exhibit identical behaviour. The binary threshold value has been chosen as $R_{thres}$ = 10.6MΩ, such that it allows all synapses to overwrite their consolidated states in the short term. It should be noted that while independent studies on device uniformity would be essential for future adoptions of this technology, the following demonstration is successful in providing a conceptual proof of concept and de-risking.

FIGURE 4.5: Analogue state history of hardware memristive synapses. Each subplot shows individual bits [0-5] of antipodal vectors $M_{1,2}$. Each stimulation event is followed by a 30s retention period. Real resistance data are plotted along with ideal fitted data. The binary state threshold (dashed line) is tuned so that all synapses can overwrite $M_1$ after observation of $M_2$.

FIGURE 4.6: Noise distribution of hardware synapses shown individually for each bit $\in [0,5]$, storing signals $M_{1,2}$. Noise has been sampled as the difference between the real and ideal resistance data from Fig. 4.5.

By utilising Eq. 4.3 the noise distributions for the studied synapses can be extracted (see Fig. 4.6). This will later be used in generalising the network's behaviour.

It has been essential to demonstrate consolidation in this zero-sum game scenario because any generalisation beyond this would imply a general increase in performance. Since $M_1$ is fully destructive to $M_2$ and vice versa, writing one signal requires a complete overwrite of the other. Random and uncorrelated memory signals would be on average 50% and thus only the other 50% would be subject to interference, leading to overall less memory degradation. Correlated memories would enjoy higher overlaps and so on.

The experiment is illustrated more holistically in Fig. 4.7a. Each $b_n$ of signals $M_{1,2}$, is fed in a synapse with binary weight $w_n$. Potentiated and depressed bits are shown in green and red respectively. The stimulation timeline of each memory is also shown. The analogue resistance values (same as in Fig. 4.5) give insight analogous to the consolidation stages presented in Fig. 4.4. Specifically, with each time $M_1$ is presented, individual synapses are entrenched further away from $R_{thres}$ and consolidate $M_1$ more stably. Consequently, initial presentation of $M_2$ is unsuccessful in overwriting $M_1$. This stems due to the observed asymmetry in RRAM switching; synapses $w_{1,4,5}$ which have undergone LTD cannot cross $R_{thres}$ during the first two presentations of $M_2$. However, the analogue state of all synapses retreats closer to $R_{thres}$ in all cases of interference.

FIGURE 4.7: Palimpsest coexistence of short- and long-term memories in memristive synapses. (**a**) The network is comprised of 6 with analogue weights $w_n$, $n \in \{0, 1, ..., 5\}$ and is subject to plasticity instructions for individual bits $b_n$ of antipodal signals $M_{1,2}$. Orange and purple represent potentiation and depression respectively. (**b**) The evolution of networks overlaps to signal $M_1$. 100% corresponds to perfect storage of $M_1$ and 0% implies perfect storage of $M_2$. (**c**) Probability of partial recollection of $M_1$ (overlap of at least x%), where $x \in \{50\%, 80\%, 100\%\}$).

Eventually, $M_2$ temporarily overwrites $M_1$ in the short-term, before the latter is reinstated as the long-term state.

The aggregate state of the network is shown in Fig. 4.7c, where the overlap between its state and $M_1$ is plotted against time. An x% $M_1$ overlap implies a (100-x)% overlap with $M_2$. The plot only shows the signal overlap history from the start of the experiment to before the third presentation of $M_1$ at t = 60s and from the first presentation of $M_2$ at t = 450s until the end of the experiment. The intermediate timeframe carries no significant information since $M_1$ is deeply consolidated at a constant 100% overlap. The quantised nature of the result can be attributed to the small network size which reduces the signal resolution. Progressive interference caused by $M_2$ makes reinstation of $M_1$ slower and noisier but not impossible. The increased noise in the overlap can be explained due to the intrinsic RRAM noise which makes a bigger difference as synapses are closer to $R_{thres}$.

Memory degradation is examined further by simulating this stimulation history using the ideal data fittings from Fig. 4.5 and the noise distribution form Fig. 4.6. A total of 200 simulations have been run, where the ideal analogue synapse state is contaminated randomly with noise using Eq. 4.4. The results are aggregated in Fig. 4.7c, showing the probability of recalling $M_1$ either fully, or partially (at least 80% or 50% overlap), following interference caused by $M_2$. Again, the reinstation of $M_1$ becomes progressively

more difficult. However, partial recall is still easily obtainable (at least 80% overlap is observed 80% of the time, while 50% is obtained almost ubiquitously). This relative metric is crucial because it suggests that memory performance can ultimately be assessed via the application's target recall accuracy. For instance, where most AI algorithms thrive on retaining aggregated representations of similar input signals [1], partial memory recall is already an acceptable result.

## 4.4    Memory System Operation

### 4.4.1    System Configuration

The previous section of this thesis considers memory consolidation where there is full competition between individual signals. It has been important to demonstrate this on hardware as a baseline proof of concept in the worst-case scenario. However, the capabilities of this technology can be generalised into the context of random and uncorrelated memories, where destructive interference is less severe and thus degradation of long-term memories is slower. To obtain a higher resolution in examining the network's behaviour, it has been decided to design a virtual network consisting of 100 identical synapses. Simulating this configuration has been essential in this stage since state of art RRAM interfacing capabilities are currently prohibiting a large-scale construction in hardware.

Importantly, while the setup of memristive synapses in the previous section has prioritised large ratios between non-volatile switching and non-volatile residues, an extended study cannot disregard the asymmetry between different plasticity directions. Any asymmetry between potentiation and depression could attract the non-volatile residue in the long run which would prohibit simulating in the range of hundreds of memories. For this reason, symmetrical responses for LTP and LTD have been sought after, in favour of symmetrical stimulation, high switching ratios or fast write speeds. Consequently, the constructed synaptic model used in the following simulations has been based on a different operation regime, which will now be discussed, and relies heavily on the methodologies discussed in Chapter 3.

To identify symmetric response along both plasticity directions, volatility has once again been examined on a device level. The results are shown in Fig. 4.8 where volatility has been induced bidirectionally on two RRAM devices, with alternate polarity. Here, each write event consists of 500 identical pulses, each 500µs each and at 10µs apart. RRAM relaxation is now measured only for 10s. Potentiation occurs at 1.4V and depression at -2.6V. Resistance has been read at 1.0V. Here, device response is at equilibrium and the extracted RRAM model can be extrapolated towards more extensive stimulation histories without suffering from 1 state attraction.

FIGURE 4.8: Retention study using two RRAM devices, serving as a reference for the simulated synaptic network. Individual retention cycles last 10s under both stimulation and potentiation events. The asymmetric stimulation profile ensures symmetric volatile responses where no plasticity direction is dominant enough to form a long-term attractor.

| Positive Stimulation | | Negative Stimulation | |
|---|---|---|---|
| Parameter | Value | Parameter | Value |
| Model Parameters | | | |
| $\alpha_p$ | -2130731.9 | $\alpha_n$ | 1166065.3 |
| $\tau_p$ | 0.728 | $\tau_n$ | 1.241 |
| $\beta_p$ | 0.537 | $\beta_n$ | 0.641 |
| $\gamma_p$ | -514203.7 | $\gamma_n$ | 512316.6 |
| Noise Parameters | | | |
| $\mu_p$ | 0.003 | $\mu_n$ | 0.004 |
| $\sigma_p$ | 0.557 | $\sigma_n$ | 0.646 |

TABLE 4.2: Extracted operation parameters for simulated memristive synapses.

Ideal retention data and noise levels have been extracted using Eq. 3.2 and Eq. 4.3 respectively on the data from Fig. 4.8. Using these parameters, a network of 100 identical parameters has been constructed. The binary threshold $R_{thres}$ has now been chosen at 11MΩ. The parameters used to construct these synapses are shown in Table 4.2. The detailed model parameters and noise distributions can be found in Appendix C. It should be noted that while the time constant $\tau$ for depression events is almost twice the value of its potentiation counterpart, the difference is negligible with respect to the overall 10s time window between memory events.

FIGURE 4.9: Memory network performance under continuous stimulation. (**a**) The temporal evolution of long-term memory (LTM) along with 3 randomly chosen short-term $ST_{1-3}$ memories, expressed via signal overlap (%). (**b**) History of all random STMs after LTM consolidation. The time where $STM_x$ overlap is higher and smaller than LTM is marked by red and blue respectively. (**c**) The distribution of all STM total lifetimes in the left y-axis along with the corresponding PDF and CDF in the right y-axis. (**d**) Average expected STM lifetime after LTM has been consolidated with various strength levels $s$. The green line ($s = 2$) corresponds to the study shown in **a-c**.

## 4.4.2   Random Memory Stream

The constructed memory network, consisting of 100 identical synapses is subject to an ongoing stream of 500 random and uncorrelated memory signals, evenly spaced at 10s apart. The goal of this study is to long-term memory degradation on a larger scale and under more persistent interference. First, after 10 random memories have been written in the network, a long-term memory (LTM) is consolidated with an intensity of $s = 2$ repetitions. The value of s is low to prevent strong consolidation of LTM in the rigid timescale. Thereafter, random short-term memories (STMs) are written with $s = 1$ and interfere with LTM. The results are shown in Fig. 4.9a, where the overlap between the network state and LTM is plotted against incoming memories. This is also repeated for 3 randomly chosen $STM_{1,2,3}$. Spikes indicate the time when a signal is written in memory and thus when it is most strongly expressed. The signal overlaps are shown against the 50% noise floor (a random signal will be on average 50% similar to LTM). Each time an STM is written in the network, the relative strength of LTM is surpassed before it is reinstated as the strongest memory

Owing to the low consolidation intensity and the interference caused by the random memories prior to LTM, the memory is never expressed perfectly in the network. Stochastic stimulation can cause local entrenchment in parts of the synaptic network which acts as an obstacle for LTM. Moreover, the overlap of all monitored signals falls over time. This reflects the overall degradation of all memories when the network is subject to continuous stimulation. It suggests that in scenarios where strong overlap is required, repetition will be essential; something that is intuitive to learning. Lastly, all observed signals retain an overlap above 50% for a long time, even if the time during which they are dominant in the network is relatively short. Even in one-shot scenarios, this network shows very high capacity for familiarity recalls and can distinguish between memories that have been presented to it before or not.

As has been previously discussed, the minimum signal overlaps that can constitute a successful memory recall are relative with respect to application specifications. Hence, it is more insightful to focus on the relative strength of LTM compared to the observed STM signals. For each $STM_x$ signal written after LTM, Fig. 4.9b observes which memory is more dominant (has higher % overlap with the network's state) during the total 10s prior to a new memory event. Blue and red regions imply that the LTM is stronger and weaker than $STM_x$ respectively. The first 100 STM signals only overwrite LTM temporarily, with the latter being swiftly restored about 2s after interference. However, the older LTM becomes in terms of the network's stimulation history the more challenging its restoration becomes. Indeed, the consolidated pattern degrades, recent signals become more potent and restoration deteriorates.

To quantify this further, STM lifetime is defined as the total time period when $STM_x$ is stronger than LTM, for a given 10s period. The distribution of STM lifetime is shown in Fig. 4.9c. Specifically, the histogram of LTM lifetimes is shown on the left y-axis, while the distribution's probability (PDF) and cumulative density functions (CDF) are plotted in the right y-axis (information on calculating these metrics can be found in Appendix C). Three distinct distribution regions can be identified. First, about 50% of the signals have a lifetime in the range between 0 and 2 seconds. These are mainly signals written at the beginning of the experiment where LTM is stored more rigidly. Then, the region in the [2, 10) seconds range is sparsely populated. This stems from the fact that after 2s, the volatile component of RRAM switching has relaxed (see $\tau$ values in Table 4.2). This means that there is a smaller probability for a synapse's state to cross $R_{thres}$. Finally, the last bin at 10s aggregates all instances where LTM fails to be reinstated altogether. This dynamic is also evident in the PDF curve in blue. Overall, LTM seems to be reinstated as the dominant signal within 4s of interference in about 80% of the instances (see CDF in red). The system thus exhibits great capacity to protect LTM, albeit with some inevitable deterioration in strength.

Nevertheless, these results arise by consolidating LTM at the minimum repetition intensity, i.e. s=2. By increasing s, the system is able to retain LTM signals more rigidly,

even by completely blocking incoming STMs. Fig. 4.9d shows the expected $STM_x$ lifetimes for a range of consolidation intensities s = $\{2, 3, 4, 5, 6\}$. Expected lifetime values have been calculated by computing the mean lifetime value using a sliding window on the last 150 occurrences. For datapoints before the $150^{th}$ STM signal, the complete history is considered. With increasing s, the short-term plasticity of the network and the ability to overwrite LTM deteriorate significantly. The lifetimes of the first STM signals decrease by about 1 order of magnitude as s increases by 1. Interestingly, at s=5, the first few tenths, and at s=6, the first hundreds of STMs completely fail to surpass LTM. At these intensities, the synapses' rigid state (RRAM non-volatile residue) has shifted away from $R_{thres}$ at distances greater than the plastic volatile changes induced by stimulation.

Consequently, this technology can also show metaplastic properties [59] where the plastic efficacy component is completely overruled. Here, the importance of the plasticity-rigidity dilemma in learning capacity is highlighted. Weak consolidation enables the synapse to act flexibly in the short term but at the expense of faster LTM deterioration. The tradeoff for strong consolidation of LTM is rigid metaplasticity, whereby the synapse is unable to store palimpsest memories altogether.

## 4.5 Visual Working Memory

Thus far, the memristive synaptic networks have been tested in the context of fully catastrophic interference (4.3.2) and random, uncorrelated memory streams (4.4.2). In this section, the simulated synapses that have been introduced in 4.4.2 are used to examine a network's learning capacity with statistically correlated signals, which is expected to enhance performance. This is done in the context of a vision network with short-term attention complementing its long-term memory. This configuration takes inspiration from how biological visual working memory is believed to function [117].

This time, the network is comprised of $100 \times 100$ identical synapses and is subject to a stream of incoming binary images. Specifically, one image is consolidated in the long-term memory (LTM) with an intensity s=3, followed by a one-shot presentation of two short-term images $ST_{1,2}$. In addition to evaluating the network's capacity to protect LTM, its ability to consolidate LTM's statistically significant information is also examined. This is done by only presenting LTM implicitly, via 3 randomly contaminated/noisy variations of NLTM. Memory correlation is thus observed at two levels. All memory signals are at least 70% similar, while NLTM variations are at least 80% similar.

The corresponding stimulation timeline and the binary image memories are shown in Fig. 4.10a-b. Time T=0s is referenced to the instance of $ST_1$'s presentation to the network. Each image pixel corresponds to a binary state written in an individual synapse.

FIGURE 4.10: Consolidation of binary images in a palimpsest network. (**a**) Stimulation timeline, including three noisy variations of long-term memory (NLTM / LTM) and two images stored in the short-term, $ST_{1,2}$. (**b**) The original image signals, including a randomly chosen NLTM variation. (**c**) Temporal snapshots of the network's binary states. The first image depicts the network's state before $ST_1$ is first observed. Time T = 0s starts when $ST_1$ is written in memory. Both instances when $ST_1$ and $ST_2$ are written are noted. This timeline shows the gradual degradation of both STs followed by the reinstation of LTM. (**d**) The overlap between the network's state and LTM over the timeline outlined in **a**. The experiment has been repeated for different NLTM noise levels. Effective reconstruction of LTM is observed even if the noiseless memory is never presented in the network, with LTM being stored more accurately than the individual NLTM overlaps, even after interference caused by $ST_{1,2}$.

The showcased NLTM image has been chosen randomly. Following the described stimulation timeline, temporal snapshots of the network's state in the scenario where LTM contamination levels are at 10% are shown in Fig. 4.10c. The first snapshot before T=0s is the last observation of the network before $ST_1$ is written in memory. At T=0s, $ST_1$ successfully overwrites the network's state; however, this transition is short-lived and LTM is reinstated. Subsequent snapshots show the transition back to LTM, which by T=9s has been recovered; importantly at a higher resolution than before $ST_1$'s presentation. $ST_2$ is written in the network at T=10s. Its decay visibly occurs in two stages. First, the pixels that are common in LTM and $ST_1$ are reinstated faster due to a stronger implicit consolidation intensity. However, the overall recovery of LTM is not aided since the transition towards it is slower. This is also comparable at Fig. 4.10d where LTM reinstation after $ST_2$ is slower than after $ST_1$; another manifestation of memory degradation due to interference. Eventually, LTM is still recovered at a greater resolution than prior to $ST_{1,2}$.

The network's capacity to detect statistical significance in the NLTM signals without supervision is qualitatively verified by inspecting the first and the last memory snapshots. The end state is in fact a much cleaner representation of LTM than a given NLTM signal. This occurs because noise is spontaneous and uncorrelated and thus it is not carried through successive NLTM write events. Random information which is only presented sparsely is consolidated weakly, if at all (see Fig. 4.9). However, the core information of LTM still has a high chance of being observed at a higher intensity. Effectively, the network has an intrinsic cleanup filter derived from its metaplastic properties.

These unsupervised denoising properties are more thoroughly examined in Fig. 4.10d. The experiment has been repeated for different NLTM noise levels, namely 0%, 2%, 5%, 10%, 20% and LTM signal overlaps have been plotted against time (the 10% overlap corresponds to the results discussed above). The overlap shown at the first NLTM occurrence reflects the corresponding noise level. In all cases, after writing all 3 NLTM signals, the resulting LTM overlap increases significantly and entirely spontaneously. In the worst case 20% contamination scenario, the LTM overlap increases by almost 10%. While these levels are subsequently contaminated by $ST_{1,2}$, the final recordings still show that LTM resolution has increased without any explicit instructions.

## 4.6 Discussion

The results presented in this chapter have covered a range of operational scenarios that demonstrate the concept of palimpsest memory consolidation in artificial synapses. The proposed synaptic model operates in a binary fashion, by quantising the analogue RRAM state R with respect to $R_{thres}$. Palimpsest capabilities are enabled via the trivial

| Example | Technology | Type | Speed | Metaplasticity | Consolidation | LTP/LTD | Lifetime | Timescales | Capacity |
|---|---|---|---|---|---|---|---|---|---|
| Chang et.al. [27] | WO$_x$ RRAM | Single device | 1ms | No | Yes | LTP | < minute | 1 | 1 |
| Ohno et.al. [102] | Ag$_2$S synapse | Single device | 500ms | No | Yes | LTP | ≈ 20s | 1 | 1 |
| Ambrogio et.al. [94] | GST PCM | Multi device | 250ns | No | No | LTP/LTD | n/a | n/a | 1 |
| Berdan et.al. [21] | TiO$_x$ RRAM | Single device | > s | No | Yes | LTP | 10s | 1 | 1 |
| Tan et.al. [101] | WO$_x$ RRAM | Single device | 10μs | Explicitly tunable | Yes | LTP | minutes | 1 | 1 |
| Boybat et.al. [29] | GST PCM | Single device | < μs | No | No | LTP/LTD | n/a | n/a | 1 |
| Burr et.al. [93] | GST PCM | Single device | n/a | No | No | LTP/LTD | n/a | n/a | 1 |
| Cheng et.al. [98] | YSZ RRAM | Single device | 0.75-1.5ms | History dependent | No | LTP/LTD | n/a | n/a | 1 |
| La Barbera et.al. [92] | GST PCM | Single device | 5-300ns | No | No | LTP/LTD | n/a | n/a | 1 |
| Lee et.al. [100] | KN memristor | Single device | 100μs | Explicitly tunable | No | LTP | n/a | n/a | 1 |
| Liu et.al. [99] | Graphene memristor | Single device | 100ns | Explicitly tunable | No | LTP/LTD | n/a | n/a | 1 |
| Wu et.al. [97] | HfO$_x$ RRAM | Single device | 1μs | Explicitly tunable | No | LTP/LTD | n/a | n/a | 1 |
| Brivio et.al. [96] | HfO$_x$ RRAM | Multi device | 10μs | No | No | LTP/LTD | n/a | n/a | 1 |
| Demirag et.al. [95] | GST PCM | Multi device | 100ns | No | Yes | LTP | 30s | 1 | 1 |
| **This work** | **TiO$_2$ RRAM** | **Single device** | **100μs - 250ms** | **Yes (implicit in device)** | **Yes** | **LTP/LTD** | **10s-30s** | **2** | **2** |

TABLE 4.3: Comparison of RRAM palimpsest synapses and previous implementations.

manipulation of volatile state changes (a plastic timescale) and the more challenging accumulation of non-volatile residues (a rigid timescale). This work presents for the first time how memories can be stored in a palimpsest fashion and in multiple timescales in hardware. The key differentiator between the results of this thesis and previous hardware implementations of synaptic memory consolidation is the bidirectional nature of memristive volatility which has been fully exploited for higher capacity and automatic long-term memory protection. A comparison of these palimpsest RRAM synapses and other technologies reviewed in Chapter 2 is outlined in Table 4.3.

Binary synaptic models are known to exhibit very good learning capabilities both when examined in the context of ANNs [60], and when examined more abstractly [118]. The synaptic resolution could in theory be increased to more distinct states, although it would require more sophisticated manipulation of the two consolidation timescales, which has not been covered in this thesis. This manipulation could rise from more complex stimulation frameworks. For instance, plasticity events of different intensities could be incorporated within a single learning experiment. However, it is unclear how the increase in stimulation profiles' resolution could facilitate an increase in learning resolution.

As discussed in Chapter 2, volatility has already been utilised for unidirectional transitions from short- to long-term memory states, which are expressed via longer binary memory lifetimes [27], [101], [102]. The differentiating factor in the results presented

here is that consolidated memories are also protected against incoming streams of interfering signals, which hasn't been examined until now. Crucially, this protection occurs automatically, without explicit external cues; this flexibility can provide autonomy for learning agents to infer things worth learning in an unsupervised manner (see clean-up filter mechanisms discussed in Fig. 4.10). Moreover, RRAM synapses enjoy a doubled memory capacity, enabled via bidirectional volatility and its two distinct timescales.

In all experiments that have been conducted, memory retrieval occurs on inevitably contaminated versions of original signals. The aggregation of multiple learnt representations into averaged memory signals is ubiquitous in neuro-inspired systems and hence partial retrieval does not present a significant problem per-se. The difference in the context of palimpsest consolidation is the additive deterioration of LTM signals caused by STMs. This though is another manifestation of the plasticity-rigidity dilemma and in fact, there may not be a decisive answer with respect to memory protection in absolute terms. There may not be an absolute winner between palimpsest memories and full metaplastic properties; some applications may benefit from short-term flexibility and memory overwrites while others may require heavier metaplastic bias and memory protection. Fortunately, with appropriate tuning (stimulation profile and device level) both can rise naturally in learning environments and will be dictated by the importance of the learning stimulus, as expressed by its presentation intensity.

Palimpsest synaptic networks can be evaluated through several prisms. First, the capacity to recall multiple memories concurrently is tied to the number of consolidation timescales that are inherent within the synapses. Consequently, these synapses can only utilise two memory slots, the short- and long-term memories. This means that non-consolidated STM signals have to constantly compete with each other before interfering with the LTM. That being said, in cases where recall accuracy does not need to be high, multiple STMs can be remembered above the random noise floor without affecting the LTM, as shown in Fig. 4.9a. Moreover, palimpsest capabilities thrive in regimes where random and uncorrelated memories are learnt and can even support consolidation when fully destructive interference is considered. This is a significant advantage over more conventionally metaplastic synaptic networks (see cascade metaplasticity [59]) which are able to increase learning capacity only for highly correlated memory signals [60].

An additional advantage of this technology rises from the flexibility palimpsest consolidation gives to online memory operations. This is most apparent in Fig. 4.10 and loosely resembles how neurological real estate is used efficiently in visual working memory systems [117], [119]–[121]. Judging by the way the simulated network handles multiple image inputs, there are things beyond absolute capacity expansion that can benefit learning. The advantage comes in the flexibility to utilise short-term memories without suffering from long-term memory degradation. These networks can thus

be used for multiple purposes without explicitly requiring a linearly expanding memory space. This interplay between memories and the inability to recall everything at once may not be a significant penalty since to some extent, learning tasks are inevitably performed sequentially.

This technology also shows a very high capacity for familiarity calls. Albeit with low recall accuracy, the familiarity filter is able to recognise signals that have been previously seen in memory and is not limited to a given LTM and STM signal. In fact, as seen in Fig. 4.9a, more than 50 STM signals can pass the familiarity threshold simultaneously. Familiarity can enable short-term attention mechanisms where a learning agent can infer significance on previously observed signals - a capability that has proven very advantageous to modern machine learning algorithms [122].

Last but not least, the designed networks are able to perform unsupervised memory reconstruction for consolidated signals, a property linked with the CA3 area in the hippocampus [123]. The dual temporal capacity that is prevalent in the technology is similar to molecular bi-stable switching which has been shown to support biological plasticity [124]. This switching, expressed by the accumulation of non-volatile residues during LTP/LTD can be a potential emulator of the calcium/calmodulin-dependent protein kinase II (CaMKII), which is believed to be a major facilitator of memory in the molecular scale [125], [126].

Further expanding the capacity for consolidated memories as well as enhancing the initial signal overlap of all memory signals requires independent research on manipulating volatility timescales at the RRAM level. Crucially, such research should inspect the candidacy of different material technologies beyond $TiO_2$. The absolute increase in memory capacity would inevitably require the identification and control of linearly increasing volatility timescales - something that to this date has not been observed in the literature. Studies at the device level would also play a significant role in engineering operational parameters such as memory write speeds, volatility time constants, energy efficiency etc. For instance, Chapter 3 shows an extensive range of volatility data under different stimulation protocols and operating at larger retention windows (up to 2 minutes). This hints that palimpsest operations can be fine-tuned to meet case-specific needs. To that end, this is a future research path that is of utmost importance for the development and potential commercialisation of this technology.

## 4.7   Summary

This chapter presented **a**) a cohesive framework for translating device-level RRAM volatility into synaptic plasticity functions, **b**) experimental results which demonstrate how palimpsest memory consolidation can be facilitated by memristive synapses and **c**) evaluated the technology's learning capabilities in a range of operational scenarios.

Plasticity functions and learning rates were emulated by the appropriate tuning of bipolar DUT stimulation. Consolidation timescales unfolded naturally and related directly to passive DUT volatility. Memristive synapses were operated at a binary configuration and their intrinsic volatile dynamics emulated two consolidation timescales; a plastic and a rigid one. This dual memory capacity facilitates the consolidation of a memory in the rigid/long-term timescale (LTM), expressed via the accumulation of non-volatile residues in specific binary states, and its temporal overwrite by multiple signals in the short-term (STM). This overwrite does not however induce loss of LTM which is essential in expanding the learning capabilities of AI hardware. Palimpsest capabilities rise automatically in memristive synapses and LTM protection is independently inferred by the intensity by which a memory is written in the network. Moreover, by appropriate tuning of the plasticity/stimulation profiles, signals can be written in or close to a one-shot fashion. Future research towards optimising and precise control of the consolidation timescales would be most fruitful. Additional properties such as noise filtering and high-capacity familiarity recall also occur naturally due to the dynamics of palimpsest operation. These characteristics can provide immense flexibility in AI agents since learning attention can temporarily be shifted to multiple tasks without suffering the tradeoff of catastrophic forgetting of previous memories.

# Chapter 5

# Detection of High Neuronal Activity

## 5.1 Introduction

The underlying operating principles of RRAM volatility (discussed in Chapter 3) and the corresponding learning timescales (demonstrated in Chapter 4) can be extended to neuromorphic applications beyond memory consolidation. Such a natural extension can be directed towards the area of brain-computer interfacing (BCI), where the ability to automatically distinguish between phases of low and high spiking activity in recorded brain regions is of great importance [103], [109], [110]. The state of spiking activity in independent regions can typically carry high-level information on the cognitive functions undertaken by a subject in real-time. High activity can be recorded by accumulating information (or integrating) regarding the firing history of an input. This



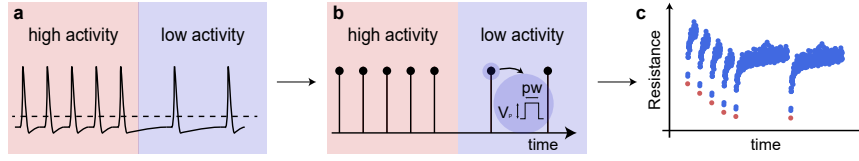FIGURE 5.1: Chapter's outline and objectives.

FIGURE 5.2: Conceptual schematic of neuronal activity detection experiment. (**a**) Schematic representation of real neuronal spikes recorded from a single neuron. The stimulation timeline can be categorised into two regions of high and low activity. (**b**) Translation of neuronal spikes into typical stimulation pulses. The programming history preserves the temporal characteristics of the raw recordings and can be used directly to induce accumulating volatile changes in a RRAM sensor. (**c**) Recorded RRAM volatility, resulting from the stimulation history in b. A single RRAM sensor is subject to said stimulation history which results in an accumulation of non-volatile residues over time. The sensor's state immediately following a programming pulse is shown in red. Resistance is red passively in between pulses and is shown in blue. Denser or sparser stimulation results in larger or lesser accumulation of non-volatile residues respectively. This discrepancy can act as a clean-up filter which can *learn* high activity regions and *forget* them during low activity timeframes.

can typically be examined on the activity of a single firing neuron (temporal integration) or cumulatively across multiple neuronal channels (spatiotemporal integration). Interestingly, unidirectionally induced volatility and the corresponding accumulation of non-volatile residues in RRAM devices can intrinsically carry information regarding the frequency of neuronal activity in a given input, using the same learning principles discussed in 4.5.

This brief chapter presents results from a preliminary study in linking neural activity detection and RRAM's integration properties. These results have been supported by the University of Padova which has provided a real dataset of neural activity recordings from living animals. It is organised as follows:

## 5.2   Aim and Methodology

The aim of this study has been to identify whether volatile RRAM devices can encode the frequency of neuronal spiking activity in their state. Importantly, activity detection needs to occur with the minimum amount of pre-processing, meaning that spike batching cannot be allowed. For this proof of concept, activity states have been quantised to *low* and *high*. This approach builds directly on the same application principles derived for palimpsest consolidation in memristive synapses (Chapter 4) and is thus a natural extension of this thesis' main objectives. In particular, the RRAM sensor that is examined relies entirely on the same technology and operates at a reduced unidirectional capacity.

A schematic representation of this study is shown in Fig. 5.2. First, a dataset of neuronal spiking recordings with accurate time reference is required to dictate the sensor's stimulation. This dataset has been kindly provided by Prof. Stefano Vassaneli's research group at the University of Padova. For the results presented here, recordings from a single firing neuron are considered. Fig. 5.2a illustrates a schematic of a neuron's output which can be distinguished into two activity regions: high and low activity shown in red and blue respectively. The time signatures of individual spikes have then been translated into programming pulses, exerted on the RRAM sensor (see Fig. 5.2b). These pulses are unidirectional and have a predefined amplitude and pulse width. The sensor's state is continuously monitored using non-invasive pulses, as defined in Chapter 3. Here, each spike event has been translated to a single pulse of 100μs and 3V amplitude. Resistance has been measured passively at 1.75V to minimise reading noise. Importantly, the time interval between spiking events/programming pulses is not fixed as in previous studies of this thesis. Instead, variable time windows between pulses should reflect the temporal characteristics of the examined firing neuron. These cannot be known a priori; it is precisely this variability in stimulation frequency which can be exploited to encode quantised activity states in the RRAM.

A typical sensor's volatile response is shown in Fig. 5.2c. The device's state immediately after stimulation is shown in red, while retention data are shown in blue. The dataset follows the stimulation profile shown in Fig. 5.2b. The difference in the stimulation frequency directly affects the accumulation of non-volatile residues in the sensor's analogue sensor. Specifically, denser stimulation inevitably lowers overall resistance while sparse stimulation allows the device to relax to higher analogue states. It is that anticipated that by appropriately defining a) the sensor stimulation profile and b) an activity detection threshold, a single volatile RRAM device can accurately encode the temporal characteristics of neuronal spikes.

## 5.3 Results and Discussion

Before examining the RRAM sensor's capabilities directly, an analysis of the input spiking dataset has first been carried out. The dataset used in the experiment is shown in Fig. 5.3.

An input history of total of 20 seconds has been chosen (see Fig. 5.3b), particularly because it contains both regions of relatively high spiking frequency as well as counterparts of no spiking activity. This is a good test case for assessing the sensor's sensitivity. Corresponding spike firing rates averaged over different bins or observation time windows can then be seen in Fig. 5.3a. These windows are $bin = \{0.1, 0.15, 0.25\}$s. The rates have been calculated as a moving average over the specified bin values; for spikes occurring earlier than the bin size, the thus far total number of spikes has been considered.
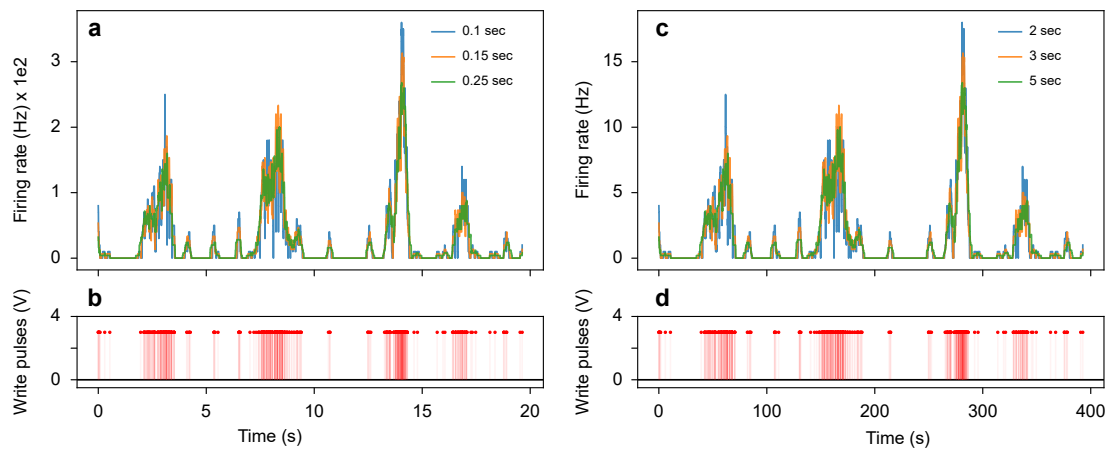
FIGURE 5.3: Temporal analysis of real spiking data. (**a-b**) The firing rate of a recorded neuron's spikes and the equivalent stimulation profile are shown for a total of 20s. The rate is calculated as a moving average over different bins/observation windows (see legend). Larger bins reduce the sensitivity or peaks of the firing rate which otherwise follows the same trends that are observed in b. (**c-d**) The same stimulation profile has been stretched in time by a factor of ×20. The chosen observation windows have also been scaled accordingly.

The results show 4 regions of relatively high activity with peaks above 100kHz for all t values. It is important to note that all peaks in Fig. 5.3a lag slightly behind the high activity centres in Fig. 5.3b since by default the activity frequency can only be inferred after the phenomenon has been observed.

A key limitation in this experiment however has been the mismatch between the spiking frequency and the temporal resolution with which RRAM devices could be stimulated and passively read. Specifically, the mean interpulse time between two spikes is $t_{mean} \approx 0.4s$ while the median interpulse is $t_{median} \approx 2ms$. To make matters worse, the minimum interpulse between successive spike events is $t_{min} \approx 0.04ms$. For this level of granularity to be accurately expressed in the RRAM sensor, a) write speeds need to be negligible with respect to $t_{min}$ and b) read speeds need to be less than or equal to $t_{min}$. These constraints are essential to bypass the need for batching multiple spikes and achieve seamless detection.

Unfortunately, they have lied beyond the current capabilities of the available interfacing instrumentation. To mitigate this, the experiment has run on a scaled timeframe such the resolution of spiking history can be compatible with the available interfacing capabilities. Fig. 5.3c-d shows the equivalent spike dataset after a scale factor of ×20 has been applied. The bin sizes and the corresponding firing rates are also proportionally scaled.

The scaled spiking dataset has then been applied to a RRAM sensor. The results of the experiment are shown in Fig. 5.4, with the resulting stimulation history shown in 5.4c. The analogue resistive state of the sensor is illustrated in Fig. 5.4a. Datapoints recorded at stimulation are coloured in red while retention points are coded in blue. The device
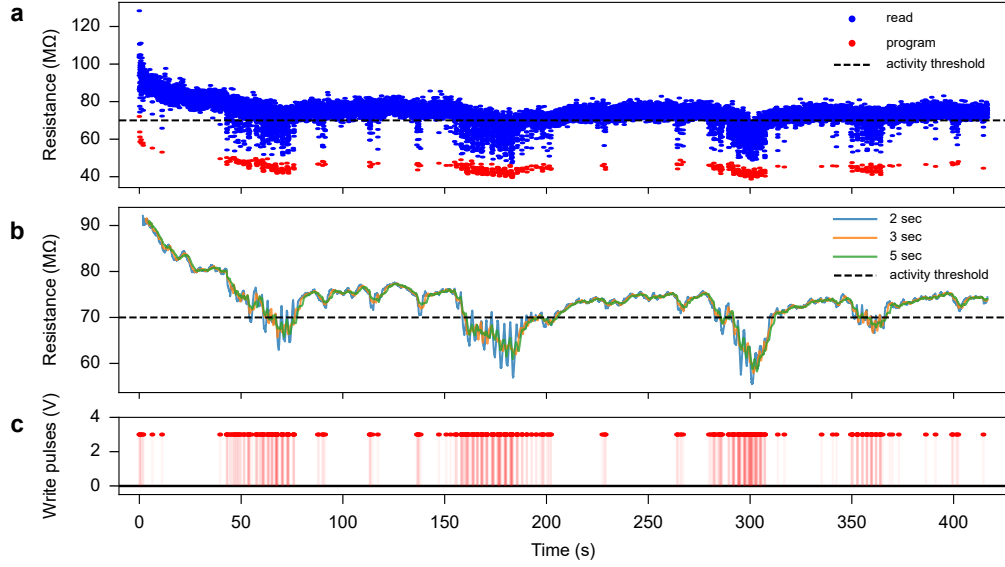
FIGURE 5.4: Neuronal activity detection - RRAM sensor's stimulation and state history.  (**a**) Raw data of the RRAM sensor's analogue state following the stimulation pattern of the spiking history in Fig. 5.3d.  The data obtained after stimulation are shown in red, while passive retention data are shown in blue.  The activity detection threshold is shown in a dashed black line. (**b**) The sensor's averaged state is calculated from the raw data using the same observation windows used in Fig. 5.3c.  Again, the activity threshold is shown. (**c**) The actual stimulation regime that is applied on the sensor.

exhibits typical volatile responses, as those discussed in the previous chapters of this thesis.  Importantly, a very high ratio between the switching datapoints and the relaxed state at equilibrium is observed (approximately a 50% decrease in R).  Moreover, the sensor naturally reaches an equilibrium close to R = 70MΩ.  For this reason, the activity detection threshold has also been chosen at the same value and is included with a black dashed line.

However, it is evident that the device's immediate stimulation (red datapoints) is almost identical irrespective of the spiking frequency and thus the detection threshold is crossed even due to minimal activity.  The raw resistive state is thus insufficient for encoding high neuronal activity.  The issue can be addressed by revisiting the averaging method used on the spike timestamp dataset.  A moving average on the analogue resistance is shown in Fig. 5.4b for using the same bin sizes as Fig. 5.3c.  The technique can effectively filter out spurious volatile changes that are caused by sparse neuronal activity.  The dips in the averaged resistive state follow the regions of high activity and are effectively inverse representations of the firing rate depicted in Fig. 5.3c.

This approach has the advantage of bypassing the need for pre-processing the spiking dataset.  It is thus making it compatible with the MIS platform which is capable of encoding individual spikes [71] (see 2.7), although bespoke linking hardware for translating detected spikes into programmable pulses is still missing from a fully integrated

application. This is a significant advantage since it paves the path of activity detection in real-time; a key requirement for many future interfacing applications. Activity has not yet been able to be decoded directly on the sensor's state, which still has to be averaged over time for binary detection. This indicates an additional memory and post-processing cost which are not necessarily prohibitively taxing. This has been caused by the high switching sensitivity of the sensor, whereby all spiking events cause a maximum displacement in R. However, the sensor could be highly accurate if the aim was to sense all spiking activity irrespective of frequency, especially by adjusting the detection threshold at around 50M$\Omega$. It would also be equivalent to calculating a moving average with a bin size comparable to individual pulse widths. Once again, the need for a thorough study of memristive volatility at the material level can be made. In an ideal scenario, volatile switching and relaxation time constants can be manipulated in order to modulate to control the apparent strength of individual programming events, which equates to the effective timeframe over which the sensor averages spiking events before detection.

It should also be noted that current device write and read speeds form a significant obstacle towards performing the experiment closer to real-time, i.e. using a stretch factor ×1 on the spiking dataset. Right now, it is highly challenging to bias RRAM devices at timeframes comparable to $t_{min}$. To accommodate this, bespoke hardware is needed such that stimulation can occur at such frequencies. Also, devices operating at lower baseline resistances would reduce reading times and overall energy required. Last but not least, the capabilities of the sensor have only been examined temporally and deal with a single input neuron channel. In many practices, especially in non-invasive activity measurements, spiking data are expected to source spatially from multiple neurons. This added complexity will most probably burden the sensor with a greater need for temporal resolution. Thus, it is imperative that RRAM volatility is optimised at the device level before more sophisticated implementations are realised.

## 5.4   Summary

Brain-computer interfacing is a rapidly developing field with immense potential towards next-generation medicine and integrating biological and artificial intelligence. However, interpreting complex brain signals is both a computational challenge and a hazard due to the invasive nature of the specialised neural recording. Thus, whether this recording is wireless or indeed invasive, signals are typically noisy and unstructured. This brief chapter has built upon existing methods of denoising low SNR neuronal signals to present a proof of concept for efficiently detecting high neuronal activity, by filtering the frequency of input spikes.

Activity detection has been encoded into the state of a volatile RRAM sensor, which emulates short-term plasticity mechanisms of high-pass signal filtering. Volatile changes that are induced by individual spikes cause the sensor to reach a binary state of high activity, which is then forgotten during periods of no stimulation. The experiments have been conducted on a stretched stimulation timeline to mitigate systematic mismatch between the temporal resolution of the spiking data and the resolution of the RRAM interfacing apparatus. Nevertheless, the relative frequency signatures of the spikes have been preserved. Activity detection can be performed accurately by averaging the state history of the sensor. This extra denoising step has been essential since the sensor has been too sensitive to stimulation. Nevertheless, by appropriately tuning device volatility, such that resistance changes do not saturate immediately and follow a cumulative increase (which has been demonstrated both in Chapters 3 and 4), temporal averaging can become an intrinsic device property. This was not achieved in this work due to the timing limitations of the project.

# Chapter 6

# Conclusions and Future Work

## 6.1 Conclusion

This thesis has successfully merged the theory behind time-dependent mechanisms that govern synaptic plasticity and the intrinsic properties of volatile RRAM technologies, to emulate various learning properties in artificial synapses. It provides evidence for the first demonstration of palimpsest memory consolidation in hardware, thereby setting a promising pathway for more dynamic and generalised learning in embedded AI.

To do so, this thesis has firstly developed a comprehensive data-driven approach for characterising and modelling volatility in $TiO_2$ RRAM devices. This has been done with the aim of demonstrating how, if understood properly, previously overlooked volatile phenomena can be interpreted as key operation mechanisms for novel application paradigms. It was thus shown for the first time that volatility can be manipulated bidirectionally, its characteristics can be controlled by appropriately tuned stimulation regimes and it manifests itself at various timeframes (here ranging from a few seconds up to 2 minutes). The presented characterisation protocol was framed with flexibility in mind, meaning that it can be adjusted to application-specific needs, ranging far from memory consolidation. Modelling was conducted for each of the main volatility characteristics individually, in order to uncover their underlying relationships with device stimulation parameters. These steps have been essential in translating the mentioned plasticity mechanisms into tangible properties of RRAM. However, the protocol was also naturally extended to investigate the relationships between volatility and operating temperatures. This was also a manifestation of the protocol's usability outside the main scopes of this thesis, further validating its importance in future RRAM-related research.

Gaining a new understanding of memristive volatility and being able to induce it bidi-rectionally in a controllable manner was the missing link between theoretical and prac-tical implementations of palimpsest memory consolidation. RRAM stimulation was distinguished at two levels; a very plastic, large change in a resistive state characterised by a fast decay and a slower, non-volatile residue which falls close to the state be-fore stimulation, accumulates over time and is more cumbersome to overwrite. These two variables were translated to two distinct consolidation timescales, a plastic and a rigid, emulating short- and long-term memory compartments respectively. Operated in a binary efficacy fashion, the flexibility of the plastic timescale allows for opposing memory states to be written in the synapses without causing the permanent loss of the rigid memory state which is attracted via the non-volatile residue. While the address-ing of catastrophic forgetting amongst highly correlated memories has previously been demonstrated, palimpsest functionality extends this ability towards memory signals that are random in nature. It allows for even fully destructive interference to over-write long-term memories in the short term. This technology was also demonstrated as a visual working memory where volatility intrinsically emulated short-term plasticity mechanisms which perform unsupervised denoising during consolidation.

The existence of only two distinct volatility timescales also limited the absolute tem-poral capacity of the RRAM synapse to an equal number of memory compartments; a short-term memory followed by the reinstation of long-term memory. Whether this ca-pacity can be expanded, remains unanswered and requires further investigation into manipulating volatility at the material level. Nevertheless, memory networks con-structed with this technology exhibited a very high capacity for familiarity filtering, with up to hundreds of random signals scoring higher recall overlaps than noise. This hints that while short-term signals decay easily, they may still have an advantage in reinstation when compared to previously unseen signals. This is very important since it could enable memory networks to recall large amounts of information with mini-mum effort, thus allowing learning agents to shift between learning tasks effortlessly. Moreover, the binary resolution of the synapses was only imposed artificially; RRAM synapses practically operate at an analogue resistive regime. Controlling palimpsest operations in an analogue fashion would be substantially more challenging but could yield significant improvements in the dynamic range of these memory networks. The mentioned attributes of this technology are unique to volatile RRAM, owing to the technology's rich intrinsic time dynamics. These advantages that have been outlined thus cannot be implemented with standard CMOS (see Section 4.6 for a thorougher analysis of said advantages). Nevertheless, due to several practical inefficiencies that still hinder RRAM, CMOS can still overcome the lack of palimpsest consolidation due to its seer advantage in metrics such as speed and efficiency. The need for addressing these RRAM inefficiencies is discussed in the next section of this chapter.

Lastly, the short-term plasticity mechanisms which gave rise to the denoising properties of the RRAM synapses were also utilised as a proof of concept for high neuronal activity detection. A volatile RRAM sensor was used in a binary fashion, to encode the frequency of prerecorded streams of neuronal spikes. Owing to the sensor's relaxation timescales, regions of high stimulation frequency caused the crossing of a predefined detection threshold, while timeframes of low stimulation led to the state's forgetting. Nevertheless, due to the lack of optimisation, the sensor activity was not directly encoded in the RRAM state but rather on a moving average of its history. Nevertheless, other results presented in this thesis suggest that with appropriate tuning activity detection could be directly read from the sensor's state. This was not able to be conducted in this thesis due to the time limitations of this project. Moreover, this proof of concept dealt with spiking data, stretched in the time domain, while preserving the relative temporal characteristics. This was done to address the mismatch between the average resolution of individual spikes and the available RRAM interfacing speeds.

To conclude, this thesis has completed an interlinked journey from neuroscience theories of synaptic plasticity, through investigations on RRAM volatility at the device level and finally to the design and demonstration of applications for palimpsest memory consolidation for AI hardware and neuronal activity detection for BCI. The author of this thesis has found the volatile operation regime of TiO$_2$ memristors to bring great potential in the design and development of neuromorphic applications. The initial interest in RRAM technologies has been in their aptitude for high-capacity memory storage but in reality, many neuronal functions are intrinsically time-dependent and show repeatable characteristics of decaying. Palimpsest consolidation has long been a goal of neuromorphic engineering due to its elegant solution to catastrophic forgetting. When developed appropriately, it has the potential to expand the capabilities of online learning agents, without the need for external computation. It is only natural that the same principles have proven beneficial in the domain of neural interfacing since both rely on synaptic plasticity principles.

Still, science often proves to be a multi-headed beast and the results presented here only point to further investigation to increase the readiness level of this technology. To that extent, the author recommends the following future investigation areas to extend this research area.

## 6.2 Author's Recommendations for Future Work

1. **Process development of volatile RRAM technologies.** Bidirectional volatility has proven to be an igniting feature for many neuromorphic applications. Nevertheless, the development of volatility on the device level is far from a mature

process. It would be most interesting to investigate how volatility can be engineered and manipulated. A concentrated effort should be made to link repeatable volatile characteristics with device parameters such as electrode materials, active oxide (potentially beyond $TiO_2$) thickness, device area, etc. In the context of memory consolidation, fine-tuning volatility translates to engineering the synapses' sensitivity to consolidation, apparent metaplasticity and forgetting timescales. With regard to neural activity detection, volatile characteristics dictate the RRAM sensor's frequency sensitivity, which can eliminate the need for historical averaging.

**Performance metrics against conventional CMOS.** As briefly mentioned in the previous point, this work has effectively de-risked this technology in the neuromorphic contexts of memory consolidation and neuronal activity detection. However, it still falls behind standard CMOS in several key performance metrics which still makes its commercial adoption challenging. This section is complementary to the previous point and provides more quantitative goals for the future of this work. First, even though it has been established that standard CMOS (e.g. DRAM, SSDs) cannot support palimpsest consolidation, it has a clear advantage in device/area metrics. Specifically, with transistors sizing a few nanometers (nm) and RRAM ranging at the micrometre (µm) scale, the advantages of palimpsest capabilities are still outweighed by the density gap between the two technologies. Therefore, it is imperative that volatile RRAM is scaled down to CMOS competitive sizes, whilst preserving its intrinsic properties. In a similar vein, it is essential that RRAM write/read speeds improve by significant margins. Throughout this study, RRAM programming pulses have been in the microseconds (µm) range with individual pulses being mostly separated by 1 millisecond (ms). These parameters result in Megahertz (MHz) operating frequencies which again cannot compete with CMOS at the Gigahertz (GHz) domain. Further experiments need to verify the limit frequencies at which the technology can operate. Operating speeds can also be increased via the size reduction of individual RRAM cells, hence this effort is complementary to the previous comment. Last but not least, while the prospect of palimpsest consolidation promises further flexibility for autonomous AI agents, this is hampered by the technology's current power hunger. This work has interfaced with RRAM devices at voltage amplitudes up to 9V (exerted for long time periods as previously mentioned) in order to induce adequately observable volatile phenomena. The power that is currently required to exploit the time dynamics of volatile RRAM is simply too much to be sustained in embedded AI hardware and not competitive against CMOS alternatives. To conclude, what RRAM currently lacks in performance metrics against CMOS could in part be addressed by reducing the technology's size. That being said, it is equally important to conduct research in material science for alternative RRAM configurations which could improve performance in all metrics. This work is an

absolutely essential step towards commercialising palimpsest memory consolidation, with the capabilities presented in this work. The following suggestions refer to further improving the conceptual and operational capabilities of volatile RRAM.

2. **Extension of consolidation timescales beyond 2 (plastic and rigid).** Potentially extending volatility timescales beyond the fast and slow relaxation domains are tightly linked to the previous route of further research. It is currently unknown whether more than two distinct timescales can be engineered yet doing so would increase the capacity of the synapses and thus the flexibility of corresponding memory networks.

3. **Design of large-scale palimpsest networks in hardware.** The practical demonstration of consolidation in memristive synapses has mostly been on the device level. The largest purely hardware-based network consisted of 6 individual synapses, while larger networks have been simulated using the volatility modelling methods developed in this thesis. This has been done due to the development and interfacing limitations accompanying the studied RRAM technology. The next logical step in verifying this application would be to develop memristive arrays of identically operating devices (this is tied to optimising volatile devices processing) as well as bespoke hardware for interfacing them concurrently.

4. **Extension of synaptic operation in the analogue domain.** The palimpsest synapses studied in this thesis have operated strictly in a binary fashion, i.e. with two states arbitrarily quantised from the analogue resistance of the RRAM devices. Binary synapses have in the past proven to be sufficient for efficient learning in ANNs but it would still be worth investigating architectures utilising synapses of higher resolution. The implementation would face its own challenges, mainly due to the volatile nature of the synapses, which would very easily corrupt 'perfect' analogue states. A practical solution could be a higher resolution quantisation of analogue resistance into more than two states. This would also benefit from the identification of further consolidation timescales, stemming from more complex volatile relaxations.

5. **Investigation on decayed memory reinstation.** While short-term memories during the random stimulation study of palimpsest consolidation have decayed quickly following further interference, the networks have shown an incredible capacity to store hundreds of signals above the noise level. This hints that palimpsest networks have the ability to express previously stored signals more efficiently, even if they have practically been forgotten. Assuming that memory in biological organisms also has the ability to reach a certain neuronal state, this detail

may have significant consequences in the design of artificial intelligence. Further study should be conducted to verify whether forgotten memories have advantages in reinstation. If so, limited capacity can be mitigated with improved reconstructive abilities.

6. **Optimisation of RRAM volatility for neuronal activity detection.** Due to the time limitations of this project, the RRAM sensor used for the corresponding study was not optimised for the experiment. Sensors with tuned switching sensitivity and forgetting time constants need to be used to verify whether activity detection can be directly stored in the state of the device.

7. **Improvement of RRAM interfacing speeds.** The study for neuronal activity detection was conducted on an artificially scaled spiking data timeline. The reason behind this step was the fact that RRAM writes and read times were significantly larger than the average interpulse between real spike data. Hence, to prevent the need to spike batching and to allow true real-time detection, the development of appropriate interfacing equipment is essential.

# Appendix A

# Volatility Model Parameters

The following table includes the parameters extracted by fitting volatility data in section 3.4 using Eq. 3.2.

TABLE A.1: Volatility extracted model parameters.

| Parameter | Positive | Negative |
|:---:|:---:|:---:|
| *alpha* | | |
| $\alpha_{\alpha 0}$ | -248.14 | -167.74 |
| $\alpha_{\alpha 1}$ | 105.37 | 264.22 |
| $\alpha_{\beta 0}$ | 1.80 | -3.48 |
| $\alpha_{\beta 1}$ | -14.26 | -1.36 |
| $\alpha_{\gamma 0}$ | $-1.11 \cdot 10^4$ | $-4.62 \cdot 10^3$ |
| $\alpha_{\gamma 1}$ | $-1.15 \cdot 10^4$ | $1.82 \cdot 10^4$ |
| *time constant* | | |
| $\tau_{\alpha 0}$ | $2.73 \cdot 10^{-4}$ | $-1.5 \cdot 10^3$ |
| $\tau_{\alpha 1}$ | $4.48 \cdot 10^{-6}$ | $3.47 \cdot 10^{-3}$ |
| $\tau_{\beta 0}$ | $2.87 \cdot 10^{-2}$ | -0.26 |
| $\tau_{\beta 1}$ | 0.125 | 7.98 |
| *beta* | | |
| $\beta_{\alpha 0}$ | -1.14 | $1.52 \cdot 10^{-2}$ |
| $\beta_{\alpha 1}$ | 8.48 | 0.49 |
| $\beta_{\tau 0}$ | 0.55 | 1.57 |
| $\beta_{\tau 1}$ | 20.88 | 20.71 |
| $\beta_{\gamma 0}$ | -0.02 | $6.24 \cdot 10^{-3}$ |
| $\beta_{\gamma 1}$ | 0.44 | 0.365 |
| *offset* | | |
| $\gamma_{\alpha 0}$ | 0.104 | 239.69∗ |
| $\gamma_{\alpha 1}$ | 0 | -390.03∗ |
| $\gamma_{\beta 0}$ | -647.76 | -464.71 |
| $\gamma_{\beta 1}$ | $3.04 \cdot 10^3$ | 1238.58 |

∗$\gamma_\alpha$ uses Eq. 3.4 for negative polarity.

# Appendix B
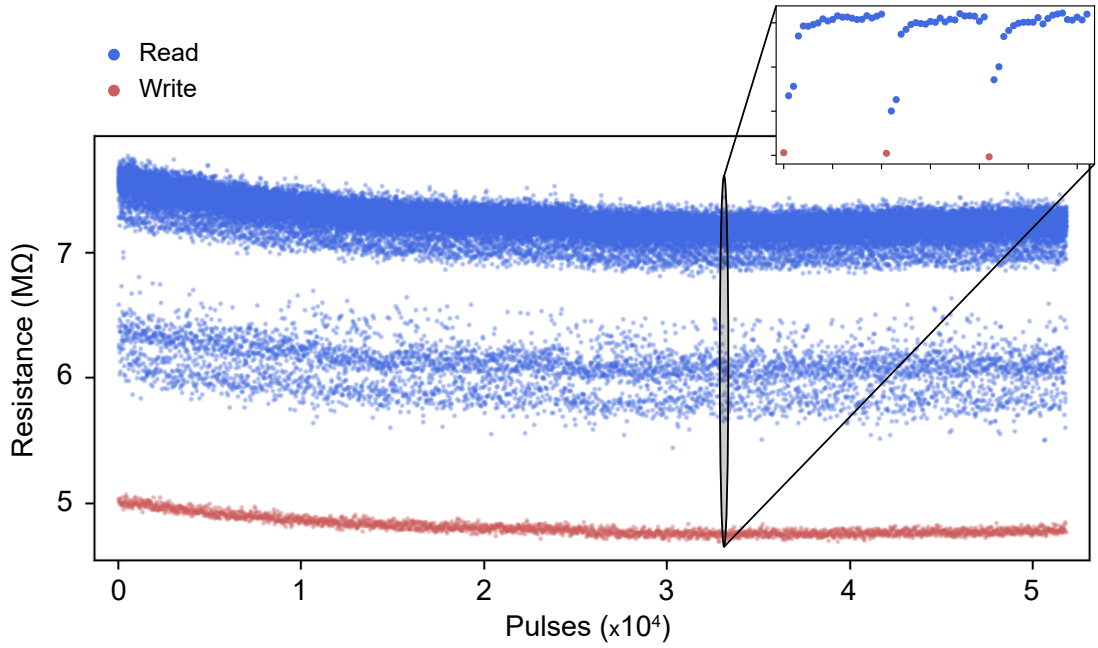
# Memristive Synapse Eligibility Study



FIGURE B.1: Observation of typical volatile behaviour in $TiO_2$ RRAM using positive polarity stimulation under a long time period. A device is subject to 2500 consecutive retention cycles. The resistance points measured at stimulation are shown in red, while the relaxation datapoints are shown in blue. Each retention cycle lasts 2 seconds. The stimulation parameters are those discussed in 4.3.1. The inset shows individual retention cycles in higher resolution. This studies shows the technology's eligibility for long-term stimulation.
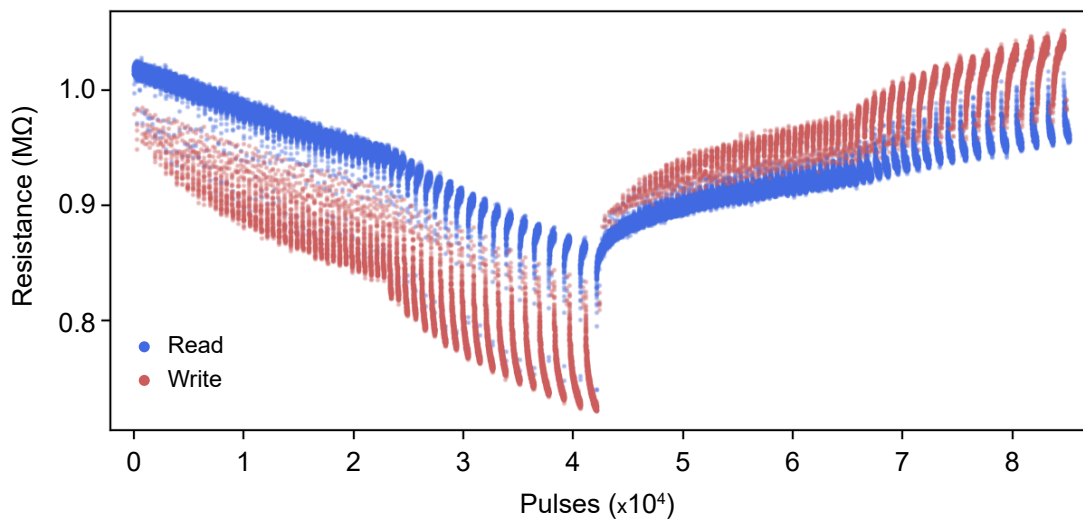
FIGURE B.2: Observation of bidirectional volatility in $TiO_2$ RRAM under continuous stimulation. The stimulation parameters are those discussed in 4.3.1 but each cycle employs a variable number N of identical pulses. N ranges from 1 to 1000 (1, then [5, 95] with a step of 5, then [100, 1000] with a step of 100) with a control retention cycle using 5 cycles in between each step. The total experiment run time concludes just short of 8.5 hours. This showcases the reversibility of the device's state under bidirectional stimulation, which ensures the overwrite of binary states and allows palimpsest consolidation in the memristive synapse.

# Appendix C

# Simulated Memory Network Set-up

## C.1 STM Lifetime Statistics

The histogram of STM lifetime occurrences has been derived by quantising the total 10s observation window using a bin of size 0.1s. Assuming $m$ corresponding bins, the total number of occurences for each bin are denoted by $o_m$ and form a resulting occurrence vector $O = \{o_0, o_1, ..., o_{m-1}\}$.

The PDF of the distribution can be obtained by dividing $O$ by the total number of STM signals, as shown below:

$$PDF = \frac{\mathbf{O}}{\#STM\ signals} \times 100\%$$

Moreover, the CDF is the cumulative sum of the PDF for each bin $i \in [0, m-1]$, such that:

$$CDF_k = \sum_{k=0}^{i} PDF_i$$

## C.2 Network Model

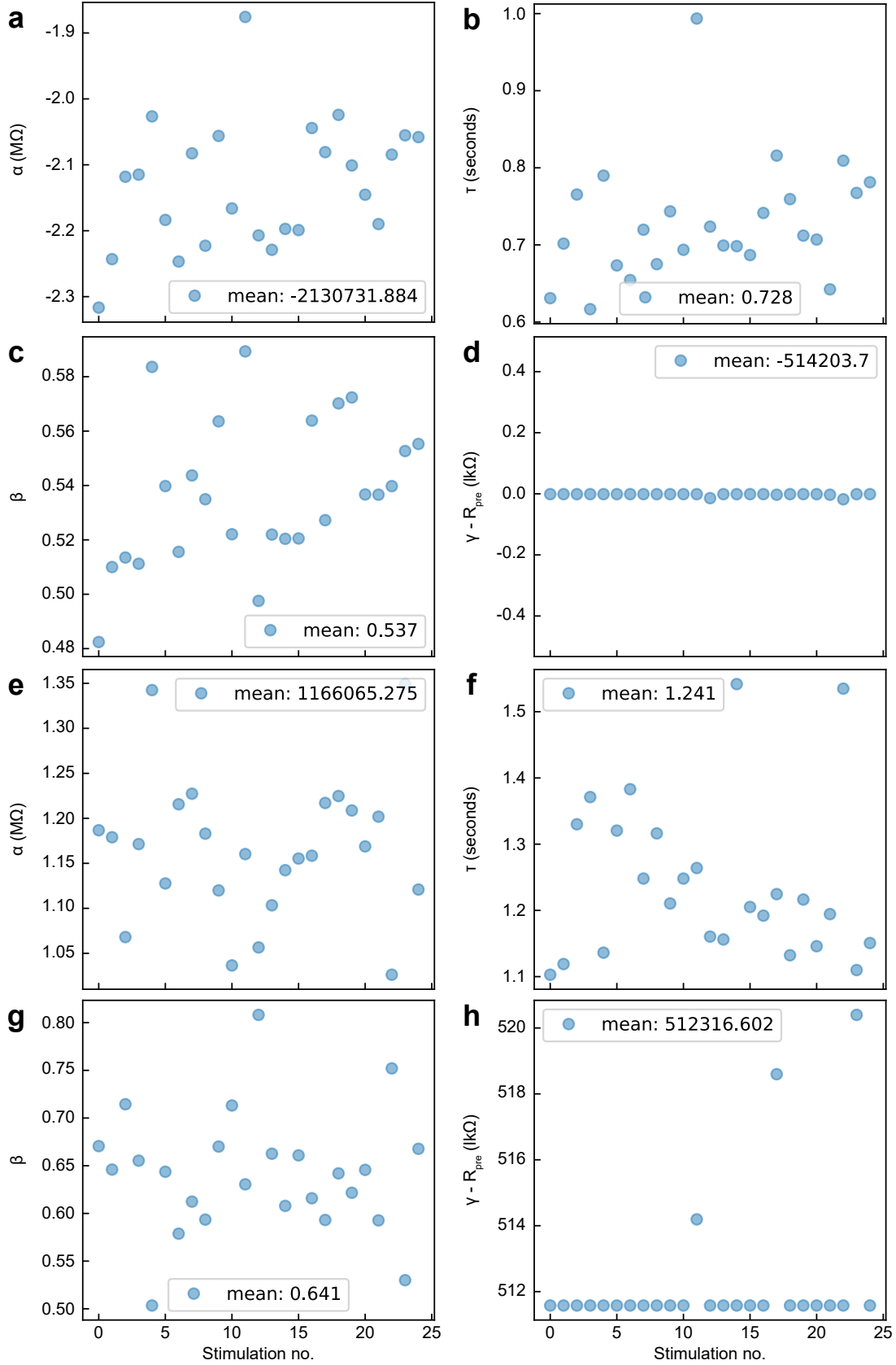$$R(t) = \alpha e^{-\left(\frac{t}{\tau}\right)^{\beta}} + \gamma$$



FIGURE C.1: Volatile R(t) parameters for Eq. 3.2. The results refer to the two device retention study shown in Fig. 4.8. Each subplot shows the average parameter value across $DUT_{1,2}$ for each retention cycle. The mean value across all runs is shown in the subplot legend. (**a-d**) Parameters for positive stimulation/potentiation. (**e-h**) Parameters for negative stimulation/depression.
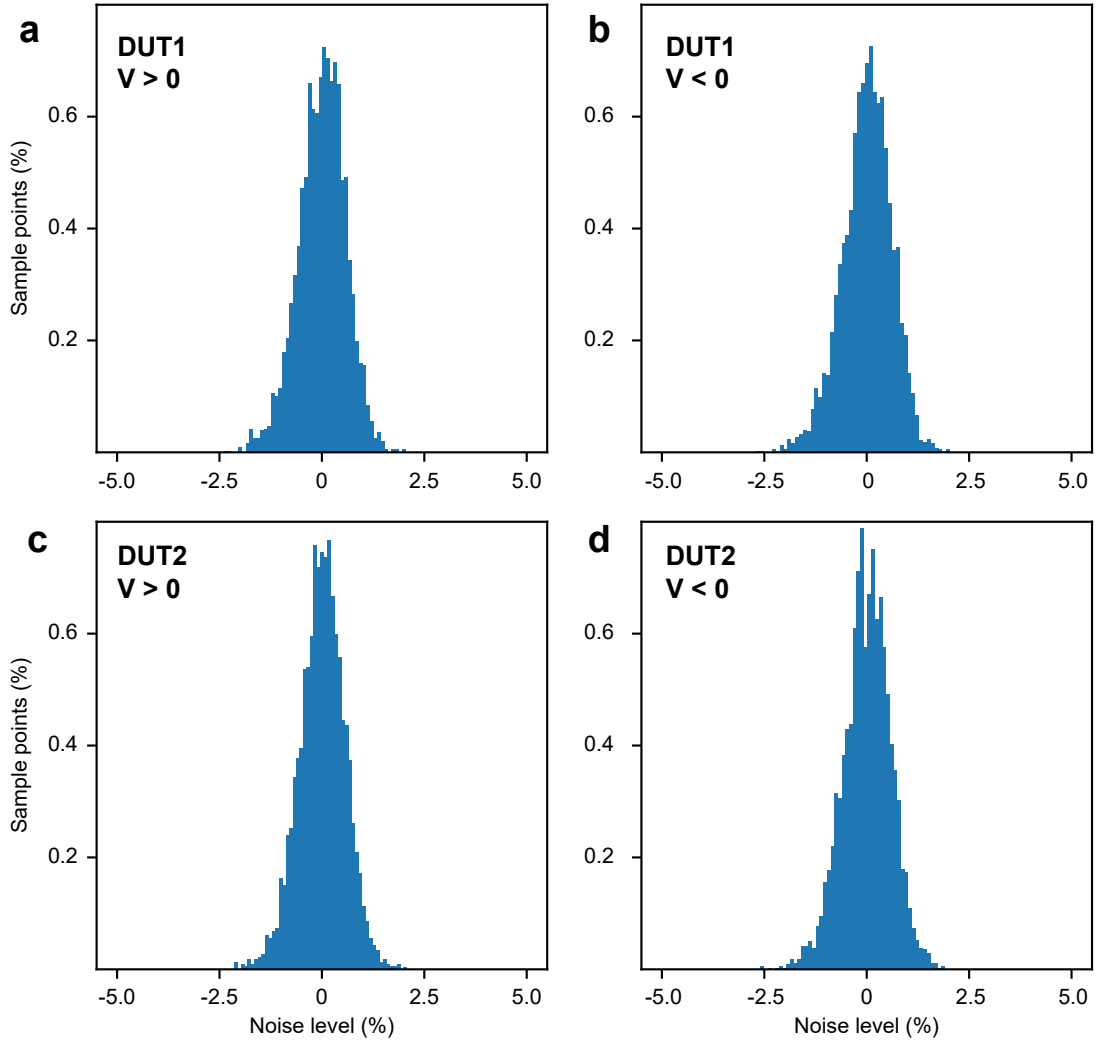
FIGURE C.2: Noise distribution from RRAM devices used to construct the simulated memory network (see Fig. 4.8). Noise data are sampled as the difference between the raw resistances and their ideal fittings.

# Appendix D

# Contributions

1. C. Giotis, A. Serb, S. Stathopoulos, L. Michalas, A. Khiat and T. Prodromakis, "Bidirectional Volatile Signatures of Metal–Oxide Memristors—Part I: Characterization," *IEEE Transactions on Electron Devices*, vol. 67, no. 11, pp. 5158-5165, Nov. 2020, doi: 10.1109/TED.2020.3014854.

2. C. Giotis, A. Serb, S. Stathopoulos and T. Prodromakis, "Bidirectional Volatile Signatures of Metal-Oxide Memristors—Part II: Modelling," *IEEE Transactions on Electron Devices*, vol. 67, no. 11, pp. 5166-5173, Nov. 2020, doi: 10.1109/TED.2020.3022343.

3. T. Abbey, C. Giotis, A. Serb, S. Stathopoulos and T. Prodromakis, "Thermal Effects on Initial Volatile Response and Relaxation Dynamics of Resistive RAM Devices" *IEEE Electron Device Letters*, vol. 43, no. 3, pp. 386-389, March 2022, doi: 10.1109/LED.2022.3145620.

4. C. Giotis, A. Serb, V. Manouras, S. Stathopoulos, and T. Prodromakis, "Palimpsest Memories Stored in Memristive Synapses," *Science Advances*, vol. 8, no. 25, Jun. 2021, doi: 10.1126/sciadv.abn7920.

# Bibliography

[1] Y. Lecun, Y. Bengio, and G. Hinton, *Deep learning*. Nature Publishing Group, May 2015, vol. 521, pp. 436–444. DOI: 10.1038/nature14539. [Online]. Available: http://www.nature.com/articles/nature14539.

[2] D. Silver, A. Huang, C. J. Maddison, *et al.*, "Mastering the game of Go with deep neural networks and tree search.," *Nature*, vol. 529, no. 7587, pp. 484–9, Jan. 2016, ISSN: 1476-4687. DOI: 10.1038/nature16961. [Online]. Available: http://www.ncbi.nlm.nih.gov/pubmed/26819042.

[3] D. Silver, J. Schrittwieser, K. Simonyan, *et al.*, "Mastering the game of Go without human knowledge.," *Nature*, vol. 550, no. 7676, pp. 354–359, 2017, ISSN: 1476-4687. DOI: 10.1038/nature24270. [Online]. Available: https://www.nature.com/articles/nature24270.pdf.

[4] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," *Communications of the ACM*, vol. 60, no. 6, pp. 84–90, 2017, ISSN: 15577317. DOI: 10.1145/3065386. [Online]. Available: http://code.google.com/p/cuda-convnet/.

[5] J. Devlin, M. W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *NAACL HLT 2019 - 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies - Proceedings of the Conference*, vol. 1, Association for Computational Linguistics (ACL), Oct. 2019, pp. 4171–4186, ISBN: 9781950737130. DOI: 10.48550/arxiv.1810.04805. [Online]. Available: https://arxiv.org/abs/1810.04805v2.

[6] J. McKendrick, "AI Adoption Skyrocketed Over the Last 18 Months," *Harvard Business Review*, pp. 1–7, 2021. [Online]. Available: https://hbr.org/2021/09/ai-adoption-skyrocketed-over-the-last-18-months.

[7] X. Xu, Y. Ding, S. X. Hu, *et al.*, "Scaling for edge inference of deep neural networks," *Nature Electronics*, vol. 1, no. 4, pp. 216–222, Apr. 2018, ISSN: 2520-1131. DOI: 10.1038/s41928-018-0059-3. [Online]. Available: https://doi.org/10.1038/s41928-018-0059-3http://www.nature.com/articles/s41928-018-0059-3.

[8]     N. Dey, A. E. Hassanien, C. Bhatt, A. S. Ashour, and S. Satapathy, *Internet of Things and Big Data Analytics Toward Next-Generation Intelligence*. Springer International Publishing, Apr. 2018, ISBN: 978-3-319-60434-3. DOI: 10.1007/978-3-319-60435-0.

[9]     G. M. Shepherd, *The Synaptic Organization of the Brain*, G. M. Shepherd, Ed. Oxford University Press, Jan. 2004, ISBN: 9780195159561. DOI: 10.1093/acprof:oso/9780195159561.001.1. [Online]. Available: https://oxford.universitypressscholarship.com/view/10.1093/acprof:oso/9780195159561.001.1/acprof-9780195159561.

[10]    C. Koch, *Biophysics of Computation: Information Processing in Single Neurons*. Oxford University Press, Nov. 1998, ISBN: 9780195104912. DOI: 10.1093/oso/9780195104912.001.0001. [Online]. Available: https://oxford.universitypressscholarship.com/view/10.1093/oso/9780195104912.001.0001/isbn-9780195104912.

[11]    M. McCloskey and N. J. Cohen, "Catastrophic Interference in Connectionist Networks: The Sequential Learning Problem," *Psychology of Learning and Motivation - Advances in Research and Theory*, vol. 24, no. C, pp. 109–165, Jan. 1989, ISSN: 00797421. DOI: 10.1016/S0079-7421(08)60536-8. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0079742108605368.

[12]    V. Sze, Y. H. Chen, T. J. Yang, and J. S. Emer, "Efficient Processing of Deep Neural Networks: A Tutorial and Survey," *Proceedings of the IEEE*, vol. 105, no. 12, pp. 2295–2329, 2017, ISSN: 15582256. DOI: 10.1109/JPROC.2017.2761740. [Online]. Available: http://www.ieee.org/publications_standards/publications/rights/index.html.

[13]    A. Canziani, A. Paszke, and E. Culurciello, "An Analysis of Deep Neural Network Models for Practical Applications," 2016. [Online]. Available: http://arxiv.org/abs/1605.07678.

[14]    C. Toumey, "Less is Moore," *Nature Nanotechnology*, vol. 11, no. 1, pp. 2–3, 2016, ISSN: 17483395. DOI: 10.1038/nnano.2015.318. [Online]. Available: http://go.nature.com/k5H3wJ.

[15]    Y. H. Chen, T. Krishna, J. S. Emer, and V. Sze, "Eyeriss: An Energy-Efficient Reconfigurable Accelerator for Deep Convolutional Neural Networks," *IEEE Journal of Solid-State Circuits*, vol. 52, no. 1, pp. 127–138, Jan. 2017, ISSN: 00189200. DOI: 10.1109/JSSC.2016.2616357.

[16]    T. Luo, S. Liu, L. Li, *et al.*, "DaDianNao: A neural network supercomputer," *IEEE Transactions on Computers*, vol. 66, no. 1, pp. 73–88, 2017, ISSN: 00189340. DOI: 10.1109/TC.2016.2574353. [Online]. Available: http://www.ieee.org/publications_standards/publications/rights/index.html.

[17] N. P. Jouppi, C. Young, N. Patil, *et al.*, "In-datacenter performance analysis of a tensor processing unit," in *Proceedings - International Symposium on Computer Architecture*, vol. Part F1286, Institute of Electrical and Electronics Engineers Inc., Apr. 2017, pp. 1–12, ISBN: 9781450348928. DOI: 10.1145/3079856.3080246. [Online]. Available: https://arxiv.org/abs/1704.04760v1.

[18] C. Zhang, P. Li, G. Sun, Y. Guan, B. Xiao, and J. Cong, "Optimizing FPGA-based accelerator design for deep convolutional neural networks," in *FPGA 2015 - 2015 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, Association for Computing Machinery, Inc, Feb. 2015, pp. 161–170, ISBN: 9781450333153. DOI: 10.1145/2684746.2689060.

[19] C. Mead, *Analog VLSI and Neural Systems*. 1989, ISBN: 0201059924. [Online]. Available: https://dl.acm.org/citation.cfm?id=64998.

[20] A Serb, J Bill, A Khiat, R Berdan, R Legenstein, and T Prodromakis, "Unsupervised learning in probabilistic neural networks with multi-state metal-oxide memristive synapses," *Nat Commun*, vol. 7, p. 12 611, 2016. DOI: 10.1038/ncomms12611. [Online]. Available: https://www.ncbi.nlm.nih.gov/pubmed/27681181.

[21] R Berdan, E Vasilaki, A Khiat, G Indiveri, A Serb, and T Prodromakis, "Emulating short-term synaptic dynamics with memristive devices," *Sci Rep*, vol. 6, p. 18 639, 2016. DOI: 10.1038/srep18639. [Online]. Available: https://www.ncbi.nlm.nih.gov/pubmed/26725838.

[22] S. H. Jo, T. Chang, I. Ebong, B. B. Bhadviya, P. Mazumder, and W. Lu, "Nanoscale memristor device as synapse in neuromorphic systems," *Nano Letters*, vol. 10, no. 4, pp. 1297–1301, Apr. 2010, ISSN: 15306984. DOI: 10.1021/nl904092h. [Online]. Available: https://pubs.acs.org/doi/full/10.1021/nl904092h.

[23] G. Indiveri and S. C. Liu, *Memory and Information Processing in Neuromorphic Systems*, 2015. DOI: 10.1109/JPROC.2015.2444094. [Online]. Available: https://ieeexplore.ieee.org/document/7159144.

[24] S. Liang and R Srikant, "Why Deep Neural Networks for Function Approximation?," Oct. 2016. DOI: 10.48550/arxiv.1610.04161. [Online]. Available: https://arxiv.org/abs/1610.04161v2.

[25] M. K. Benna and S Fusi, "Computational principles of synaptic memory consolidation," *Nat Neurosci*, vol. 19, no. 12, pp. 1697–1706, 2016. DOI: 10.1038/nn.4401. [Online]. Available: https://www.ncbi.nlm.nih.gov/pubmed/27694992.

[26] W. C. Abraham, "Metaplasticity: Tuning synapses and networks for plasticity," *Nature Reviews Neuroscience*, vol. 9, no. 5, pp. 387–399, 2008, ISSN: 1471003X. DOI: 10.1038/nrn2356. [Online]. Available: www.nature.com/reviews/neuro.

[27] T. Chang, S. H. Jo, and W. Lu, "Short-term memory to long-term memory transition in a nanoscale memristor," *ACS Nano*, vol. 5, no. 9, pp. 7669–7676, 2011, ISSN: 19360851. DOI: 10.1021/nn202983n. [Online]. Available: www.acsnano.org.

[28] J. H. Yoon, Z. Wang, K. M. Kim, *et al.*, "An artificial nociceptor based on a diffusive memristor," *Nature Communications*, vol. 9, no. 1, 2018, ISSN: 20411723. DOI: 10.1038/s41467-017-02572-3. [Online]. Available: www.nature.com/naturecommunications.

[29] I. Boybat, M. Le Gallo, S. R. Nandakumar, *et al.*, "Neuromorphic computing with multi-memristive synapses," *Nature Communications*, vol. 9, no. 1, pp. 1–12, Dec. 2018, ISSN: 20411723. DOI: 10.1038/s41467-018-04933-y. [Online]. Available: www.nature.com/naturecommunications.

[30] D. J. F. de Quervain, B Roozendaal, R. M. Nitsch, J. L. McGaugh, and C Hock, "Acute cortisone adminisration impairs retrieval of long-term declarative memory in humans," *Nature America*, 2000.

[31] A. Roxin and S. Fusi, "Efficient Partitioning of Memory Systems and Its Importance for Memory Consolidation," *PLoS Computational Biology*, vol. 9, no. 7, p. 1003146, 2013, ISSN: 1553734X. DOI: 10.1371/journal.pcbi.1003146. [Online]. Available: www.ploscompbiol.org.

[32] L. R. Squire, L. Genzel, J. T. Wixted, and R. G. Morris, "Memory consolidation," *Cold Spring Harbor Perspectives in Biology*, vol. 7, no. 8, 2015, ISSN: 19430264. DOI: 10.1101/cshperspect.a021766. [Online]. Available: http://cshperspectives.cshlp.org/.

[33] T Shallice, P Fletcher, C. D. Frith, P Grasby, R. S. Frackowiak, and R. J. Dolan, "Brain regions associated with acquisition and retrieval of verbal episodic memory," *Nature*, vol. 368, no. 6472, pp. 633–635, 1994. DOI: 10.1038/368633a0. [Online]. Available: https://www.ncbi.nlm.nih.gov/pubmed/8145849.

[34] C. M. Alberini, "Mechanisms of memory stabilization: Are consolidation and reconsolidation similar or distinct processes?" *Trends in Neurosciences*, vol. 28, no. 1, pp. 51–56, Jan. 2005, ISSN: 01662236. DOI: 10.1016/j.tins.2004.11.001. [Online]. Available: https://linkinghub.elsevier.com/retrieve/pii/S0166223604003558.

[35] C. Koch, T. Poggio, and V. Torre, "Retinal ganglion cells: a functional interpretation of dendritic morphology.," *Philosophical transactions of the Royal Society of London. Series B, Biological sciences*, vol. 298, no. 1090, pp. 227–263, Jul. 1982, ISSN: 09628436. DOI: 10.1098/rstb.1982.0084. [Online]. Available: http://rstb.royalsocietypublishing.org/cgi/doi/10.1098/rstb.1982.0084.

[36] G. M. Shepherd and R. K. Brayton, "Logic operations are properties of computer-simulated interactions between excitable dendritic spines," *Neuroscience*, vol. 21, no. 1, pp. 151–165, Apr. 1987, ISSN: 03064522. DOI: 10.1016/0306-4522(87)90329-0. [Online]. Available: https://www.sciencedirect.com/science/article/abs/pii/0306452287903290?via%3Dihub.

[37] *Neurotransmitter release*. [Online]. Available: https://www.britannica.com/science/neurotransmitter-release/media/1/410790/66782.

[38] E. R. Kandel, "The molecular biology of memory: cAMP, PKA, CRE, CREB-1, CREB-2, and CPEB.," *Molecular brain*, vol. 5, p. 14, 2012, ISSN: 1756-6606. DOI: 10.1186/1756-6606-5-14. [Online]. Available: http://www.molecularbrain.com/content/5/1/14http://www.ncbi.nlm.nih.gov/pubmed/22583753%0Ahttp://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC3514210.

[39] T. M. Newpher and M. D. Ehlers, "Glutamate Receptor Dynamics in Dendritic Microdomains," *Neuron*, vol. 58, no. 4, pp. 472–497, May 2008, ISSN: 08966273. DOI: 10.1016/j.neuron.2008.04.030. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S089662730800408X?via%3Dihub.

[40] A. Citri and R. C. Malenka, "Synaptic plasticity: Multiple forms, functions, and mechanisms," *Neuropsychopharmacology*, vol. 33, no. 1, pp. 18–41, 2008, ISSN: 0893133X. DOI: 10.1038/sj.npp.1301559. [Online]. Available: www.neuropsychopharmacology.org.

[41] R. S. Zucker and W. G. Regehr, "Short-term synaptic plasticity," *Annual Review of Physiology*, vol. 64, pp. 355–405, 2002, ISSN: 00664278. DOI: 10.1146/annurev.physiol.64.092501.114547. [Online]. Available: https://pubmed.ncbi.nlm.nih.gov/11826273/.

[42] B. Katz and R. Miledi, "The role of calcium in neuromuscular facilitation," *The Journal of Physiology*, vol. 195, no. 2, pp. 481–492, Mar. 1968, ISSN: 14697793. DOI: 10.1113/jphysiol.1968.sp008469. [Online]. Available: https://pubmed.ncbi.nlm.nih.gov/4296699/.

[43] L. F. Abbott and W. G. Regehr, "Synaptic computation," *Nature*, 2004, ISSN: 00280836. DOI: https://doi.org/10.1038/nature03010. [Online]. Available: www.nature.com/nature.

[44] T. V. Bliss and T. Lømo, "Long-lasting potentiation of synaptic transmission in the dentate area of the anaesthetized rabbit following stimulation of the perforant path," *The Journal of Physiology*, vol. 232, no. 2, pp. 331–356, Jul. 1973, ISSN: 14697793. DOI: 10.1113/jphysiol.1973.sp010273. [Online]. Available: https://pubmed.ncbi.nlm.nih.gov/4727084/.

[45] T. V. Bliss and A. R. Gardner-Medwin, "Long-lasting potentiation of synaptic transmission in the dentate area of the unanaesthetized rabbit following stimulation of the perforant path," *The Journal of Physiology*, vol. 232, no. 2, pp. 357–374, Jul. 1973, ISSN: 14697793. DOI: 10.1113/jphysiol.1973.sp010274. [Online]. Available: https://pubmed.ncbi.nlm.nih.gov/4727085/.

[46] R. M. Mulkey and R. C. Malenka, "Mechanisms underlying induction of homosynaptic long-term depression in area CA1 of the hippocampus," *Neuron*, vol. 9, no. 5, pp. 967–975, 1992, ISSN: 08966273. DOI: 10.1016/0896-6273(92)90248-C. [Online]. Available: https://pubmed.ncbi.nlm.nih.gov/1419003/.

[47] D. O. Hebb, *The organization of behavior; a neuropsychological theory.* Oxford, England: Wiley, 1949, pp. xix, 335–xix, 335.

[48] L. F. Abbott and S. B. Nelson, "Synaptic plasticity: Taming the beast," *Nature Neuroscience*, vol. 3, no. 11s, pp. 1178–1183, 2000, ISSN: 15461726. DOI: 10.1038/81453. [Online]. Available: https://www.nature.com/articles/nn1100_1178.

[49] N. Caporale and Y. Dan, "Spike timing-dependent plasticity: A Hebbian learning rule," *Annual Review of Neuroscience*, vol. 31, pp. 25–46, 2008, ISSN: 0147006X. DOI: 10.1146/annurev.neuro.31.060407.125639. [Online]. Available: www.annualreviews.org.

[50] R. A. Nicoll, J. A. Kauer, and R. C. Malenka, "The current excitement in long term potentiation," *Neuron*, vol. 1, no. 2, pp. 97–103, 1988, ISSN: 08966273. DOI: 10.1016/0896-6273(88)90193-6. [Online]. Available: https://pubmed.ncbi.nlm.nih.gov/2856092/.

[51] W. C. Abraham and M. F. Bear, "Metaplasticity: The plasticity of synaptic plasticity," *Trends in Neurosciences*, vol. 19, no. 4, pp. 126–130, Apr. 1996, ISSN: 01662236. DOI: 10.1016/S0166-2236(96)80018-X. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S016622369680018X?via%3Dihub.

[52] Y. Y. Huang, A. Colino, D. K. Selig, and R. C. Malenka, "The influence of prior synaptic activity on the induction of long-term potentiation," *Science*, vol. 255, no. 5045, pp. 730–733, 1992, ISSN: 00368075. DOI: 10.1126/science.1346729. [Online]. Available: https://pubmed.ncbi.nlm.nih.gov/1346729/.

[53] H. Wang and J. J. Wagner, "Priming-induced shift in synaptic plasticity in the rat hippocampus," *Journal of Neurophysiology*, vol. 82, no. 4, pp. 2024–2028, 1999, ISSN: 00223077. DOI: 10.1152/jn.1999.82.4.2024. [Online]. Available: https://pubmed.ncbi.nlm.nih.gov/10515995/.

[54] R. Urbanczik and W. Senn, "Learning by the Dendritic Prediction of Somatic Spiking," *Neuron*, vol. 81, no. 3, pp. 521–528, Feb. 2014, ISSN: 08966273. DOI: 10.1016/j.neuron.2013.11.030. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0896627313011276?via%3Dihub.

[55]   C. E. Shannon, "A Mathematical Theory of Communication," *Bell System Technical Journal*, vol. 27, no. 3, pp. 379–423, 1948, ISSN: 15387305. DOI: `10.1002/j.1538-7305.1948.tb01338.x`.

[56]   H. Luo, T. Shahroodi, H. Hassan, *et al.*, "CLR-DRAM: A Low-Cost DRAM Architecture Enabling Dynamic Capacity-Latency Trade-Off," in *Proceedings - International Symposium on Computer Architecture*, vol. 2020-May, Institute of Electrical and Electronics Engineers Inc., May 2020, pp. 666–679, ISBN: 9781728146614. DOI: `10.1109/ISCA45697.2020.00061`.

[57]   D. Lee, Y. Kim, V. Seshadri, J. Liu, L. Subramanian, and O. Mutlu, "Tiered-latency DRAM: A low latency and low cost DRAM architecture," in *Proceedings - International Symposium on High-Performance Computer Architecture*, 2013, pp. 615–626, ISBN: 9781467355858. DOI: `10.1109/HPCA.2013.6522354`.

[58]   S Fusi, "Computational models of long term plasticity and memory," 2017.

[59]   S Fusi, P. J. Drew, and L. F. Abbott, "Cascade models of synaptically stored memories," *Neuron*, vol. 45, no. 4, pp. 599–611, 2005. DOI: `10.1016/j.neuron.2005.02.001`. [Online]. Available: `https://www.ncbi.nlm.nih.gov/pubmed/15721245`.

[60]   A. Laborieux, M. Ernoult, T. Hirtzlin, and D. Querlioz, "Synaptic metaplasticity in binarized neural networks," *Nature Communications*, vol. 12, no. 1, p. 2549, Dec. 2021, ISSN: 2041-1723. DOI: `10.1038/s41467-021-22768-y`. [Online]. Available: `https://doi.org/10.1038/s41467-021-22768-yhttp://www.nature.com/articles/s41467-021-22768-y`.

[61]   C. Savin, P. Dayan, and M. Lengyel, "Two is better than one : distinct roles for familiarity and recollection in retrieving palimpsest memories," *Advances in Neural Information Processing Systems 24*, pp. 1305–1313, 2011. [Online]. Available: `https://papers.nips.cc/paper/4364-two-is-better-than-one-distinct-roles-for-familiarity-and-recollection-in-retrieving-palimpsest-memories`.

[62]   L. O. Chua, "Memristor—The Missing Circuit Element," *IEEE Transactions on Circuit Theory*, vol. 18, no. 5, pp. 507–519, 1971, ISSN: 00189324. DOI: `10.1109/TCT.1971.1083337`. [Online]. Available: `http://ieeexplore.ieee.org/document/1083337/`.

[63]   D. B. Strukov, G. S. Snider, D. R. Stewart, and R. S. Williams, "The missing memristor found.," *Nature*, 2008, ISSN: 1476-4687. DOI: `10.1038/nature06932`.

[64]   R. B. Jacobs-Gedrim, S. Agarwal, R. S. Goeke, *et al.*, "Analog high resistance bilayer RRAM device for hardware acceleration of neuromorphic computation," *Citation: Journal of Applied Physics*, vol. 124, no. 20, p. 202 101, Nov. 2018, ISSN: 10897550. DOI: `10.1063/1.5042432`. [Online]. Available: `http://aip.scitation.org/doi/10.1063/1.5042432http://aip.scitation.org/toc/jap/124/20`.

[65]   A. Serb, A. Khiat, and T. Prodromakis, "Seamlessly fused digital-analogue re-configurable computing using memristors," *Nature Communications*, vol. 9, no. 1, p. 2170, Dec. 2018, ISSN: 20411723. DOI: 10.1038/s41467-018-04624-8. [Online]. Available: http://www.nature.com/articles/s41467-018-04624-8.

[66]   H. Jiang and Q. Xia, "Effect of voltage polarity and amplitude on electroforming of TiO 2 based memristive devices," *Nanoscale*, vol. 5, no. 8, pp. 3257–3261, 2013, ISSN: 20403364. DOI: 10.1039/c3nr00622k. [Online]. Available: www.rsc.org/nanoscale.

[67]   S Stathopoulos, A Khiat, M Trapatseli, *et al.*, "Multibit memory operation of metal-oxide bi-layer memristors," *Sci Rep*, vol. 7, no. 1, p. 17 532, 2017. DOI: 10.1038/s41598-017-17785-1. [Online]. Available: https://www.ncbi.nlm.nih.gov/pubmed/29235524.

[68]   R. Berdan, C. Lim, A. Khiat, C. Papavassiliou, and T. Prodromakis, "A memristor SPICE model accounting for volatile characteristics of practical ReRAM," *IEEE Electron Device Letters*, vol. 35, no. 1, pp. 135–137, Jan. 2014, ISSN: 07413106. DOI: 10.1109/LED.2013.2291158. [Online]. Available: http://ieeexplore.ieee.org/document/6680642/.

[69]   R. Gütig and H. Sompolinsky, "The tempotron: a neuron that learns spike timing–based decisions," *Nature Neuroscience*, vol. 9, no. 3, pp. 420–428, 2006. DOI: 10.1038/nn1643. [Online]. Available: https://dx.doi.org/10.1038/nn1643.

[70]   S. L. Wei, E Vasilaki, A Khiat, I Salaoru, R Berdan, and T Prodromakis, "Emulating long-term synaptic dynamics with memristive devices,"

[71]   I. Gupta, A. Serb, A. Khiat, R. Zeitler, S. Vassanelli, and T. Prodromakis, "Real-time encoding and compression of neuronal spikes by metal-oxide memristors," *Nature Communications*, vol. 7, no. 1, p. 12 805, Dec. 2016, ISSN: 20411723. DOI: 10.1038/ncomms12805. [Online]. Available: http://www.nature.com/articles/ncomms12805.

[72]   I Gupta, A Serb, A Khiat, R Zeitler, S Vassanelli, and T Prodromakis, "Sub 100 nW Volatile Nano-Metal-Oxide Memristor as Synaptic-Like Encoder of Neuronal Spikes," *IEEE Trans Biomed Circuits Syst*, vol. 12, no. 2, pp. 351–359, 2018. DOI: 10.1109/TBCAS.2018.2797939. [Online]. Available: https://www.ncbi.nlm.nih.gov/pubmed/29570062.

[73]   A Serb, A Khiat, and T Prodromakis, *Practical demonstration of a memristive fuse*, arXiv, 2016.

[74]   X. Zhu, C. Du, Y. Jeong, and W. D. Lu, "Emulation of synaptic metaplasticity in memristors," *Nanoscale*, vol. 9, no. 1, pp. 45–51, 2017, ISSN: 20403372. DOI: 10.1039/c6nr08024c. [Online]. Available: www.rsc.org/nanoscale.

[75] I. Messaris, A. Serb, S. Stathopoulos, A. Khiat, S. Nikolaidis, and T. Prodromakis, "A Data-Driven Verilog-A ReRAM Model," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 37, no. 12, pp. 3151–3162, 2018. DOI: 10.1109/tcad.2018.2791468.

[76] L Michalas, S. Stathopoulos, A Khiat, and T. Prodromakis, "An electrical characterisation methodology for identifying the switching mechanism in TiO2 memristive stacks," *Scientific Reports*, vol. 9, no. 1, 2019, ISSN: 20452322. DOI: 10.1038/s41598-019-44607-3. [Online]. Available: https://doi.org/10.1038/s41598-019-44607-3.

[77] S. Stathopoulos, A. Serb, A. Khiat, M. Ogorzalek, and T. Prodromakis, "A Memristive Switching Uncertainty Model," *IEEE Transactions on Electron Devices*, vol. 66, no. 7, pp. 2946–2953, Jul. 2019, ISSN: 15579646. DOI: 10.1109/TED.2019.2918102.

[78] S Stathopoulos, L Michalas, A Khiat, A Serb, and T Prodromakis, "A techgnology agnostic RRAM characterisation methodology protocol,"

[79] I. Valov, "Interfacial interactions and their impact on redox-based resistive switching memories (ReRAMs)," *Semiconductor Science and Technology*, vol. 32, no. 9, p. 093006, Sep. 2017, ISSN: 13616641. DOI: 10.1088/1361-6641/aa78cd. [Online]. Available: http://stacks.iop.org/0268-1242/32/i=9/a=093006?key=crossref.4e3144ced4fa96dd0b94c9ce0aad6af8.

[80] I. H. Inoue and A. Sawa, "Resistive switchings in transition-metal oxides," in *Functional Metal Oxides: New Science and Novel Applications*, 6, vol. 11, Jun. 2013, pp. 443–463, ISBN: 9783527654864. DOI: 10.1002/9783527654864.ch16. [Online]. Available: https://linkinghub.elsevier.com/retrieve/pii/S1369702108701196.

[81] J. J. Yang, M. D. Pickett, X. Li, D. A. Ohlberg, D. R. Stewart, and R. S. Williams, "Memristive switching mechanism for metal/oxide/metal nanodevices," *Nature Nanotechnology*, vol. 3, no. 7, pp. 429–433, 2008, ISSN: 17483395. DOI: 10.1038/nnano.2008.160. [Online]. Available: www.nature.com/naturenanotechnology.

[82] C. Wang, H. Wu, B. Gao, T. Zhang, Y. Yang, and H. Qian, "Conduction mechanisms, dynamics and stability in ReRAMs," *Microelectronic Engineering*, vol. 187-188, pp. 121–133, Feb. 2018, ISSN: 01679317. DOI: 10.1016/j.mee.2017.11.003. [Online]. Available: https://linkinghub.elsevier.com/retrieve/pii/S0167931717303672.

[83] A. Fantini, G. Gorine, R. Degraeve, *et al.*, "Intrinsic program instability in HfO2 RRAM and consequences on program algorithms," in *Technical Digest - International Electron Devices Meeting, IEDM*, Institute of Electrical and Electronics Engineers Inc., Feb. 2015, pp. 1–7, ISBN: 9781467398930. DOI: 10.1109/IEDM.2015.7409648.

[84]  D. M. Nminibapiel, D. Veksler, P. R. Shrestha, *et al.*, "Impact of RRAM read fluctuations on the program-verify approach," *IEEE Electron Device Letters*, vol. 38, no. 6, pp. 736–739, Jun. 2017, ISSN: 07413106. DOI: 10.1109/LED.2017.2696002.

[85]  S. Cortese, M. Trapatseli, A. Khiat, and T. Prodromakis, "On the origin of resistive switching volatility in Ni/TiO2/Ni stacks," *Journal of Applied Physics*, vol. 120, no. 6, p. 65 104, 2016, ISSN: 10897550. DOI: 10.1063/1.4960690. [Online]. Available: https://doi.org/10.1063/1.4960690.

[86]  I. Gupta, A. Serb, R. Berdan, A. Khiat, and T. Prodromakis, "Volatility Characterization for RRAM Devices," *IEEE Electron Device Letters*, vol. 38, no. 1, pp. 28–31, 2017. DOI: 10.1109/led.2016.2631631.

[87]  T. Liu, M. Verma, Y. Kang, and M. Orlowski, "Volatile resistive switching in Cu/TaO x/$\delta$-Cu/Pt devices," *Applied Physics Letters*, vol. 101, no. 7, p. 73 510, 2012, ISSN: 00036951. DOI: 10.1063/1.4746276. [Online]. Available: https://doi.org/10.1063/1.4746276.

[88]  E. Covi, D. Ielmini, Y. H. Lin, *et al.*, "A volatile RRAM synapse for neuromorphic computing," in *2019 26th IEEE International Conference on Electronics, Circuits and Systems, ICECS 2019*, Institute of Electrical and Electronics Engineers Inc., Nov. 2019, pp. 903–906, ISBN: 9781728109961. DOI: 10.1109/ICECS46596.2019.8965044.

[89]  W. Wang, M. Laudato, E. Ambrosi, *et al.*, "Volatile Resistive Switching Memory Based on Ag Ion Drift/Diffusion Part I: Numerical Modeling," *IEEE Transactions on Electron Devices*, vol. 66, no. 9, pp. 3795–3801, Sep. 2019, ISSN: 15579646. DOI: 10.1109/TED.2019.2928890.

[90]  W. Wang, M. Laudato, E. Ambrosi, *et al.*, "Volatile Resistive Switching Memory Based on Ag Ion Drift/Diffusion - Part II: Compact Modeling," *IEEE Transactions on Electron Devices*, vol. 66, no. 9, pp. 3802–3808, Sep. 2019, ISSN: 15579646. DOI: 10.1109/TED.2019.2928888.

[91]  A. Ascoli, S. Slesazeck, H. Mähne, R. Tetzlaff, and T. Mikolajick, "Nonlinear Dynamics of a Locally-Active Memristor," in *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 62, 2015, pp. 1165–1174. DOI: 10.1109/TCSI.2015.2413152.

[92]  S. La Barbera, D. R. Ly, G. Navarro, *et al.*, "Narrow Heater Bottom Electrode-Based Phase Change Memory as a Bidirectional Artificial Synapse," *Advanced Electronic Materials*, vol. 4, no. 9, 2018, ISSN: 2199160X. DOI: 10.1002/aelm.201800223. [Online]. Available: https://doi.org/10.1002/aelm.201800223..

[93] G. W. Burr, R. M. Shelby, S. Sidler, *et al.*, "Experimental Demonstration and Tolerancing of a Large-Scale Neural Network (165 000 Synapses) Using Phase-Change Memory as the Synaptic Weight Element," *IEEE Transactions on Electron Devices*, vol. 62, no. 11, pp. 3498–3507, 2015, ISSN: 00189383. DOI: 10.1109/TED.2015.2439635.

[94] S. Ambrogio, N. Ciocchini, M. Laudato, *et al.*, "Unsupervised learning by spike timing dependent plasticity in phase change memory (PCM) synapses," *Frontiers in Neuroscience*, vol. 10, no. MAR, p. 56, 2016, ISSN: 1662453X. DOI: 10.3389/fnins.2016.00056. [Online]. Available: www.frontiersin.org.

[95] Y. Demirag, F. Moro, T. Dalgaty, *et al.*, "PCM-trace: Scalable synaptic eligibility traces with resistivity drift of phase-change materials," *Proceedings - IEEE International Symposium on Circuits and Systems*, vol. 2021-May, pp. 1–6, 2021, ISSN: 02714310. DOI: 10.1109/ISCAS51556.2021.9401446.

[96] S Brivio, D Conti, M. V. Nair, *et al.*, "Extended memory lifetime in spiking neural networks employing memristive synapses with nonlinear conductance dynamics," *Nanotechnology*, vol. 30, no. 1, 2019, ISSN: 13616528. DOI: 10.1088/1361-6528/aae81c. [Online]. Available: https://doi.org/10.1088/1361-6528/aae81c.

[97] Q. Wu, H. Wang, Q. Luo, *et al.*, "Full imitation of synaptic metaplasticity based on memristor devices," *Nanoscale*, vol. 10, no. 13, pp. 5875–5881, Apr. 2018, ISSN: 20403372. DOI: 10.1039/c8nr00222c.

[98] C. Cheng, Y. Li, T. Zhang, *et al.*, "Bipolar to unipolar mode transition and imitation of metaplasticity in oxide based memristors with enhanced ionic conductivity," *Journal of Applied Physics*, vol. 124, no. 15, p. 152 103, Oct. 2018, ISSN: 10897550. DOI: 10.1063/1.5037962. [Online]. Available: http://aip.scitation.org/doi/10.1063/1.5037962.

[99] B. Liu, Z. Liu, I. S. Chiu, *et al.*, "Programmable Synaptic Metaplasticity and below Femtojoule Spiking Energy Realized in Graphene-Based Neuromorphic Memristor," *ACS Applied Materials and Interfaces*, vol. 10, no. 24, pp. 20 237–20 243, 2018, ISSN: 19448252. DOI: 10.1021/acsami.8b04685. [Online]. Available: https://pubs.acs.org/doi/10.1021/acsami.8b04685.

[100] T. H. Lee, H. G. Hwang, J. U. Woo, D. H. Kim, T. W. Kim, and S. Nahm, "Synaptic Plasticity and Metaplasticity of Biological Synapse Realized in a KNbO3 Memristor for Application to Artificial Synapse," *ACS Applied Materials and Interfaces*, vol. 10, no. 30, pp. 25 673–25 682, 2018, ISSN: 19448252. DOI: 10.1021/acsami.8b04550. [Online]. Available: www.acsami.org.

[101] Z. H. Tan, R. Yang, K. Terabe, X. B. Yin, X. D. Zhang, and X. Guo, "Synaptic Metaplasticity Realized in Oxide Memristive Devices," *Advanced Materials*, vol. 28, no. 2, pp. 377–384, Jan. 2016, ISSN: 15214095. DOI: 10.1002/adma.201503575. [Online]. Available: http://doi.wiley.com/10.1002/adma.201503575.

[102]  T. Ohno, T. Hasegawa, T. Tsuruoka, K. Terabe, J. K. Gimzewski, and M. Aono, "Short-term plasticity and long-term potentiation mimicked in single inorganic synapses," *Nature Materials*, vol. 10, no. 8, pp. 591–595, Jun. 2011, ISSN: 14764660. DOI: 10.1038/nmat3054. [Online]. Available: www.nature.com/naturematerials.

[103]  J. A. Frank, M. J. Antonini, and P. Anikeeva, "Next-generation interfaces for studying neural function," *Nature Biotechnology*, vol. 37, no. 9, pp. 1013–1023, Aug. 2019, ISSN: 15461696. DOI: 10.1038/s41587-019-0198-8. [Online]. Available: https://www.nature.com/articles/s41587-019-0198-8.

[104]  J. Rivnay, H. Wang, L. Fenno, K. Deisseroth, and G. G. Malliaras, "Next-generation probes, particles, and proteins for neural interfacing," *Science Advances*, vol. 3, no. 6, Jun. 2017, ISSN: 23752548. DOI: 10.1126/sciadv.1601649. [Online]. Available: https://www.science.org/doi/abs/10.1126/sciadv.1601649.

[105]  M. A. Nicolelis, "Actions from thoughts," *Nature*, vol. 409, no. 6818, pp. 403–407, Jan. 2001, ISSN: 00280836. DOI: 10.1038/35053191. [Online]. Available: https://www.nature.com/articles/35053191.

[106]  E. Musk and Neuralink, "An integrated brain-machine interface platform with thousands of channels," *Journal of Medical Internet Research*, vol. 21, no. 10, p. 16 194, 2019, ISSN: 14388871. DOI: 10.2196/16194. [Online]. Available: http://jmir.org/2019/10/e16321/Commentin:http://jmir.org/2019/10/e16339/Commentin:http://jmir.org/2019/10/e16356/Commentin:http://jmir.org/2019/11/e16344/.

[107]  L. R. Hochberg, M. D. Serruya, G. M. Friehs, *et al.*, "Neuronal ensemble control of prosthetic devices by a human with tetraplegia," *Nature*, vol. 442, no. 7099, pp. 164–171, Jul. 2006, ISSN: 00280836. DOI: 10.1038/nature04970. [Online]. Available: https://www.nature.com/articles/nature04970.

[108]  W. Wang, J. L. Collinger, A. D. Degenhart, *et al.*, "An Electrocorticographic Brain Interface in an Individual with Tetraplegia," *PLoS ONE*, vol. 8, no. 2, Feb. 2013, ISSN: 19326203. DOI: 10.1371/journal.pone.0055344. [Online]. Available: https://pubmed.ncbi.nlm.nih.gov/23405137/.

[109]  L. R. Hochberg, D. Bacher, B. Jarosiewicz, *et al.*, "Reach and grasp by people with tetraplegia using a neurally controlled robotic arm," *Nature*, vol. 485, no. 7398, pp. 372–375, May 2012, ISSN: 00280836. DOI: 10.1038/nature11076. [Online]. Available: https://www.nature.com/articles/nature11076.

[110]  G. K. Anumanchipalli, J. Chartier, and E. F. Chang, "Speech synthesis from neural decoding of spoken sentences," *Nature*, vol. 568, no. 7753, pp. 493–498, Apr. 2019, ISSN: 14764687. DOI: 10.1038/s41586-019-1119-1. [Online]. Available: https://www.nature.com/articles/s41586-019-1119-1.

[111] U. Frey, U. Egert, F. Heer, S. Hafizovic, and A. Hierlemann, "Microelectronic system for high-resolution mapping of extracellular electric fields applied to brain slices," *Biosensors and Bioelectronics*, vol. 24, no. 7, pp. 2191–2198, Mar. 2009, ISSN: 09565663. DOI: 10.1016/j.bios.2008.11.028.

[112] T. J. Blanche, M. A. Spacek, J. F. Hetke, and N. V. Swindale, "Polytrodes: High-density silicon electrode arrays for large-scale multiunit recording," *Journal of Neurophysiology*, vol. 93, no. 5, pp. 2987–3000, May 2005, ISSN: 00223077. DOI: 10.1152/jn.01023.2004. [Online]. Available: https://pubmed.ncbi.nlm.nih.gov/15548620/.

[113] J. C. Phillips and J. M. Vandenberg, "Subensembles and Kohlrausch relaxation in electronic and molecular glasses," *Journal of Physics Condensed Matter*, vol. 9, no. 18, 1997, ISSN: 09538984. DOI: 10.1088/0953-8984/9/18/001.

[114] J Kakalios, R. A. Street, and W. B. Jackson, "Stretched-exponential relaxation arising from dispersive diffusion of hydrogen in amorphous silicon," *Physical Review Letters*, vol. 59, no. 9, pp. 1037–1040, 1987, ISSN: 00319007. DOI: 10.1103/PhysRevLett.59.1037.

[115] B Sturman, E Podivilov, and M Gorkunov, "Origin of stretched exponential relaxation for hopping-transport models," *Physical Review Letters*, vol. 91, no. 17, 2003, ISSN: 10797114. DOI: 10.1103/PhysRevLett.91.176602.

[116] K. M. Kim, J. Zhang, C. Graves, *et al.*, "Low-Power, Self-Rectifying, and Forming-Free Memristor with an Asymmetric Programing Voltage for a High-Density Crossbar Application," *Nano Letters*, vol. 16, no. 11, pp. 6724–6732, Nov. 2016, ISSN: 15306992. DOI: 10.1021/acs.nanolett.6b01781. [Online]. Available: https://pubs.acs.org/doi/10.1021/acs.nanolett.6b01781.

[117] G. Mongillo, O. Barak, and M. Tsodyks, "Synaptic Theory of Working Memory," *Science*, vol. 319, no. 5869, pp. 1543–1546, Mar. 2008, ISSN: 10959203. DOI: 10.1126/science.1150769. [Online]. Available: https://pubmed.ncbi.nlm.nih.gov/18339943/.

[118] M. K. Benna and S Fusi, "Efficient online learning with low-precision synaptic variables," English, *2017 Fifty-First Asilomar Conference on Signals, Systems, and Computers*, pp. 1610–1614, 2017.

[119] L. Matthey, P. M. Bays, and P. Dayan, "A Probabilistic Palimpsest Model of Visual Short-term Memory," *PLoS Computational Biology*, vol. 11, no. 1, p. 1004003, Jan. 2015, ISSN: 15537358. DOI: 10.1371/journal.pcbi.1004003. [Online]. Available: http://www.cns.nyu.edu/malab/.

[120] S. J. Luck and E. K. Vogel, "The capacity of visual working memory for features and conjunctions," *Nature*, vol. 390, no. 6657, pp. 279–284, Nov. 1997, ISSN: 00280836. DOI: 10.1038/36846. [Online]. Available: https://www.nature.com/articles/36846.

[121] W. Zhang and S. J. Luck, "Discrete fixed-resolution representations in visual working memory," *Nature*, vol. 453, no. 7192, pp. 233–235, May 2008, ISSN: 14764687. DOI: 10.1038/nature06860. [Online]. Available: www.nature.com/nature..

[122] A. Vaswani, N. Shazeer, N. Parmar, *et al.*, "Attention is all you need," in *Advances in Neural Information Processing Systems*, vol. 2017-Decem, 2017, pp. 5999–6009.

[123] C Savin, P Dayan, and M Lengyel, "Optimal recall from bounded metaplastic synapses: predicting functional adaptations in hippocampal area CA3," *PLoS Comput Biol*, vol. 10, no. 2, e1003489, 2014. DOI: 10.1371/journal.pcbi.1003489. [Online]. Available: https://www.ncbi.nlm.nih.gov/pubmed/24586137.

[124] U. S. Bhalla, "Molecular computation in neurons: A modeling perspective," *Current Opinion in Neurobiology*, vol. 25, pp. 31–37, Apr. 2014, ISSN: 09594388. DOI: 10.1016/j.conb.2013.11.006.

[125] J. Lisman, H. Schulman, and H. Cline, "The molecular basis of CaMKII function in synaptic and behavioural memory," *Nature Reviews Neuroscience*, vol. 3, no. 3, pp. 175–190, 2002, ISSN: 14710048. DOI: 10.1038/nrn753. [Online]. Available: https://www.nature.com/articles/nrn753.

[126] P. Miller, A. M. Zhabotinsky, J. E. Lisman, and X. J. Wang, "The stability of a stochastic CaMKII switch: Dependence on the number of enzyme molecules and protein turnover," *PLoS Biology*, vol. 3, no. 4, pp. 0705–0717, 2005, ISSN: 15449173. DOI: 10.1371/journal.pbio.0030107. [Online]. Available: www.plosbiology.org.