

# The Enhanced Rise and Delayed Fall of Memory in a Model of Synaptic Integration: Extension to Discrete State Synapses

Terry Elliott<sup>1</sup>

Department of Electronics and Computer Science,

University of Southampton,

Highfield,

Southampton, SO17 1BJ,

United Kingdom.

**Running Title:** Filter-Based Memory Dynamics with Discrete Synapses.

May 3, 2016.

---

<sup>1</sup>Tel.: +44 (0)23 8059 6000, Fax.: +44 (0)23 8059 2783, E.-mail:  
te@ecs.soton.ac.uk.

# Abstract

Integrate-and-express models of synaptic plasticity propose that synapses may act as low-pass filters, integrating synaptic plasticity induction signals in order to discern trends before expressing synaptic plasticity. We have previously shown that synaptic filtering strongly controls destabilising fluctuations in developmental models. When applied to palimpsest memory systems that learn new memories by forgetting old ones, we have also shown that with binary-strength synapses, integrative synapses lead to an initial memory signal rise before its fall back to equilibrium. Such an initial rise is in dramatic contrast to non-integrative synapses, in which the memory signal falls monotonically. We now extend our earlier analysis of palimpsest memories with synaptic filters to consider the more general case of discrete state, multi-level synapses. We derive exact results for the memory signal dynamics and then consider various simplifying approximations. We show that multi-level synapses enhance the initial rise in the memory signal and then delay its subsequent fall by inducing a plateau-like region in the memory signal. Such dynamics significantly increase memory lifetimes, defined by a signal-to-noise ratio (SNR). We derive expressions for optimal choices of synaptic parameters (filter size, number of strength states, number of synapses) that maximise SNR memory lifetimes. However, we find that with memory lifetimes defined via mean-first-passage times, such optimality conditions do not exist, suggesting that optimality may be an artifact of SNRs.

# 1 Introduction

The Hopfield model (Hopfield, 1982) provides the general foundation for many approaches to associative memory. However, its catastrophic forgetting above a threshold memory loading renders it implausible even as a toy model of biological memory. Imposing bounds on synaptic strengths overcomes this catastrophic forgetting by turning the network into a palimpsest, storing new memories by forgetting old ones (Nadal *et al.*, 1986; Parisi, 1986). One biophysically plausible way of implementing bounds on synaptic strength is to suppose that synapses exist in only a finite set of states of synaptic strength. While some experimental evidence supports the possibility of binary-strength synapses (Petersen *et al.*, 1998; O’Connor *et al.*, 2005b), other evidence suggests the existence of ternary-strength (Montgomery & Madison, 2002, 2004) or ternary-state (O’Connor *et al.*, 2005a) synapses, while yet further evidence indicates that changes in synaptic strength may be discrete, step-like processes without necessarily addressing any possible limit on the number of states of synaptic strength (Yasuda *et al.*, 2003; Bagal *et al.*, 2005; Sobczyk & Svoboda, 2007). Many models have considered memory formation with both binary-strength synapses and more general, multi-level, discrete synapses (see, for example, Willshaw *et al.*, 1969; Tsodyks, 1990; Amit & Fusi, 1994, Fusi *et al.*, 2005, Leibold & Kempster, 2006, 2008; Rubin & Fusi, 2007; Fusi & Abbott, 2007; Barrett & van Rossum, 2008; Huang & Amit, 2010, 2011). All these related models share one feature in common: the fidelity of recall of a memory

falls monotonically in time, often exponentially fast. Much work has been devoted to extending the resulting rather short memory lifetimes in these models, but the underlying problem of monotonic memory trace decay always remains.

In previous work, we have proposed that synapses may integrate synaptic plasticity induction signals before expressing synaptic plasticity during both development (Elliott, 2008; Elliott & Lagogiannis, 2009) and memory formation (Elliott & Lagogiannis, 2012). By integrating plasticity induction signals, synapses behave as low-pass filters, suppressing high frequency noise and responding only to low frequency signals. In this way, the fluctuations in synaptic strength that destabilise both developmentally-relevant states and states of memory can be controlled (Elliott, 2011b). We applied such a filtering mechanism to memory formation in a feedforward framework with binary-strength synapses and showed that in radical contrast to non-integrative models of memory, synaptic filtering and ongoing memory storage actually facilitate an initial increase in the fidelity of recall of a stored memory (Elliott & Lagogiannis, 2012). Such a model outperforms cascade-type models (Fusi *et al.*, 2005) in most biologically-relevant regions of parameter space.

Here, we extend our earlier analysis from binary-strength to more general, discrete synapses. After discussing our general formalism in section 2, we derive exact results for the tracked memory signal in the presence of  $n$  discrete strength states in section 3. Various approximations to these results may be obtained for the large time limit or the large  $n$  limit, which facilitate the extraction of expressions for memory lifetimes, defined via a signal-to-noise ratio

(SNR). In section 4, we first explore the dynamics of the tracked memory signal for general, discrete synapses in the presence of a synaptic filter. We show that the signal rise that occurs for binary-strength synapses is present for general, discrete synapses and that, indeed, this signal rise is logarithmically enhanced as a function of  $n$ . We also show that the signal then essentially plateaus for large  $n$ , with the duration of this plateauing increasing quadratically with  $n$ . We then turn to considering memory lifetimes explicitly, and explore the dependence of memory lifetimes on the number of synapse, number of states of strength per synapse, and also on the synaptic filter size,  $\Theta$ . We find that for SNR memory lifetimes, we can trade  $n$  and  $\Theta$ , significantly reducing to biophysically realistic ranges the optimal values of  $n$  or  $\Theta$  that generate maximum SNR memory lifetimes. However, when using a mean-first-passage time (MFPT) definition of memory lifetimes (Elliott, 2014), we do not see maxima in memory lifetimes, so optimality conditions do not exist in this case. Finally, in section 5, we discuss our results and the issues raised by them.

## 2 General Formalism

We provide an outline of our general formalism here. Further details may be found elsewhere (Elliott & Lagogiannis, 2012; Elliott, 2014). Table 1 provides a summary of the main parameters and quantities used throughout.

Parameter or quantity	Description
$n$	Number of states of synaptic strength per synapse.
$N$	Number of synapses.
$\Theta$	Filter size.
$s_A$	Strength of strength state $A$ , $A \in \{1, \dots, n\}$ .
$\mathbf{S}$	Vector of strengths $s_A$ , $\mathbf{S}^T = (s_1, \dots, s_n)$ .
$S_i(t)$	Strength of synapse $i$ at time $t$ .
$h(t)$	Tracked memory signal.
$\mu(t), \sigma(t)$	Mean and standard deviation of $h(t)$ .
$\mathcal{B}, \mathcal{B}_I$	Equilibrium distribution of filter states, vector and components.
$\mathcal{A}$	Joint distribution of filter and strength states in equilibrium.
$\mathbb{M}^\pm$	Matrices implementing potentiation and depression steps on the joint distribution of filter and strength states.
$\mathbb{F}^\pm$	Matrices incrementing or decrementing filter state without threshold processes.
$\mathbb{D}^\pm$	Matrices implementing filter threshold processes.
$G_J^\pm(t)$	Densities for first escapes through filter thresholds from filter state $J$ in time $t$ .
$H_J(t)$	Probability of not having reached either filter threshold from filter state $J$ in time $t$ .
$f_{I J}(t)$	Probability of a transition from filter state $J$ to filter state $I$ in time $t$ .

$p_{I J}^{A B}(t)$	Probability of a transition from filter state $J$ and strength state $B$ to filter state $I$ and strength state $A$ in time $t$ .
$\mathbb{P}_{I J}(t)$	Matrix with elements $p_{I J}^{A B}(t)$ for given filter states $I$ and $J$ .
$\mathbb{G}_J(t)$	Matrix of escape densities from filter states $J$ for transitions to adjacent strength states.
$\mathbb{W}(t)$	Auxiliary but key matrix of filter escape densities determining transitions between strength states, defined by $\mathbb{W}(t) = \delta(t)\mathbb{I} - \mathbb{G}_0(t)$ .
$\mathbb{C}$	The matrix $\frac{1}{2}\mathbb{C}$ is a stochastic matrix implementing a symmetric random walk between two reflecting boundaries.
$\lambda_m, \mathbf{e}^m$	Eigenvalues and eigenvectors of $\frac{1}{2}\mathbb{C}$ , $m = 0, \dots, n - 1$ .
$\mathbf{v}$	$\mathbf{v}^T = (-1, 0, \dots, 0, +1)$ is a vector that captures the change in the equilibrium distribution of strengths at the storage of the tracked memory.

**Table 1.** Summary of main parameters and quantities used throughout.

## 2.1 Perceptron Formulation

We consider the possibility of  $n$  states of synaptic strength, with  $n \geq 2$ , and examine the dependence of memory lifetimes on this parameter. We index these strength states by letters such as  $A$  and  $B$ , and we define strength state  $A$  to correspond to strength

$$s_A = -1 + 2 \frac{A - 1}{n - 1}, \quad (2.1)$$

with  $A = 1, \dots, n$ , so that  $s_1 = -1$  and  $s_n = +1$ . We have scaled the strengths, regardless of  $n$ , into the interval  $[-1, +1]$  in order to facilitate comparison of

results for different  $n$ . We shall discuss the biological relevance of this scaling in the Discussion. For simplicity and mathematical tractability, we consider a single perceptron. The perceptron has  $N$  synapses with strengths  $S_i(t)$ ,  $i = 1, \dots, N$ , where  $t$  denotes time, with  $S_i(t) \in \{s_A \mid A = 1, \dots, n\}$ . As standard, the perceptron is assumed to have binary-valued inputs  $x_i \in \{-1, +1\}$  through these  $N$  synapses. The activation upon presentation of input vector  $\mathbf{x}$  is then

$$h_{\mathbf{x}}(t) = \frac{1}{N} \sum_{i=1}^N x_i S_i(t). \quad (2.2)$$

For our purposes here, we are interested only in this activation and not in any thresholding of the perceptron’s activation that generates the perceptron’s binary-valued output.

The perceptron is required to store “memories”  $\xi^\alpha$ ,  $\alpha = 0, 1, 2, \dots$ . In a discrete time formalism, memory  $\xi^\alpha$  is stored at time  $t = \alpha$ . From a biological perspective, however, a discrete time formalism for memory storage is not particularly realistic. Furthermore, we have previously shown that driving memory storage as a discrete time process eliminates covariance terms that have a detrimental impact on memory dynamics (Elliott & Lagogiannis, 2012). Using a continuous time process to drive memory storage is biologically more realistic and allows a full consideration of the resulting impact of covariance terms on memory dynamics. We therefore employ a continuous time formalism to drive memory storage. The simplest continuous time process to consider is the Poisson process. Memories are therefore stored as a Poisson process of



rate  $r$ , which we may without loss of generality take as  $r = 1$  Hz, since  $r$  may be restored in formulae by the replacement  $t \rightarrow r t$ . Despite using a Poisson process, memory  $\xi^0$  is nevertheless always stored at  $t = 0$  s; in fact, we consider it to be stored at time  $t = 0^-$  s so that the time immediately after the storage of  $\xi^0$  can be referred to simply as  $t = 0$  s. We need not specify a target perceptron output associated with memory  $\xi^\alpha$  because for an isolated perceptron we can always without loss of generality consider instead the storage of  $-\xi^\alpha$  rather than  $+\xi^\alpha$ . We then always consider the target output for any memory to be  $+1$ , so that the corresponding perceptron activation is above firing threshold. With this convention,  $\xi_i^\alpha$  is the plasticity induction signal to synapse  $i$  upon storage of memory  $\alpha$ :  $\xi_i^\alpha = +1$  requires the synapse to potentiate (strengthen) while  $\xi_i^\alpha = -1$  requires it to depress (weaken). We discuss the implementation of synaptic plasticity in response to these induction signals below. As usual, the memories are assumed to be random and uncorrelated, both across synapses and between different memories, so that  $\xi_i^\alpha = \pm 1$  with probability  $1/2$  independent of  $i$  and  $\alpha$ . We note that we do not consider the possibility of a sparse coding framework here. We discuss sparse coding in the Discussion.

Although the use of strengths  $s_A$  in the range  $[-1, +1]$  and binary-valued inputs  $x_i \in \{-1, +1\}$  may appear biologically problematic, we can always translate these ranges so that they become non-negative under an associated change in the perceptron's firing threshold. These issues are discussed elsewhere (Elliott & Lagogiannis, 2012; Elliott, 2014).

We are interested in the fidelity of recall of the first memory  $\xi^0$  by the

perceptron in the face of the ongoing storage of the subsequent memories  $\xi^\alpha$ ,  $\alpha > 0$ . The perceptron's activation upon re-presentation of  $\xi^0$  at some later time  $t$  is just  $h_{\xi^0}(t)$ , and if this activation is above firing threshold, then the memory is still stored by the perceptron. We refer to  $\xi^0$  as the tracked memory, to  $h_{\xi^0}(t)$  as the tracked memory signal or just the memory signal, and we write  $h(t) = h_{\xi^0}(t)$  for convenience. Memory lifetimes may be defined in many different ways (for examples, see Tsodyks, 1990; Leibold & Kempter, 2006; Huang & Amit, 2010; Elliott, 2014). Here, we mostly employ the SNR definition (Tsodyks, 1990; Amit & Fusi, 1994). If  $\mu(t)$  is the mean memory signal and  $\sigma(t)$  the standard deviation in the memory signal, then the SNR is defined as  $\mu(t)/\sigma(t)$ . Memory lifetime is then defined as the time  $\tau_{\text{snr}}$  at which  $\mu(t)/\sigma(t)$  falls below some defined point, which is typically taken to be unity, so that  $\tau_{\text{snr}}$  is the solution of  $\mu(\tau_{\text{snr}})/\sigma(\tau_{\text{snr}}) = 1$ . For simplicity, we will mostly use this definition here. Specifically, we will assume that a perceptron's firing threshold can always be chosen so that the mean memory signal does not become inaccessible by ever dropping below the perceptron's firing threshold. For most models, this requirement amounts to choosing a firing threshold of zero, because  $\mu(t)$  asymptotes to zero at large times. Without this assumption, memory lifetimes are severely and disastrously shortened, and become independent of  $N$  as  $N$  increases (Elliott, 2014). We will also consider for comparison a definition of memory lifetimes based on MFPTs (Elliott, 2014). In this formulation, memory lifetime is defined as the average time at which the stochastic memory signal  $h(t)$  first falls below firing threshold. Such a defini-

tion provides a more natural definition of memory lifetime, but is analytically much more difficult to study for non-trivial models of synaptic plasticity.

## 2.2 Filter-Based Synaptic Plasticity

Upon the storage of memory  $\xi^\alpha$ , the component  $\xi_i^\alpha$  is the induction signal to synapse  $i$ , indicating whether the synapse should potentiate ( $\xi_i^\alpha = +1$ ) or depress ( $\xi_i^\alpha = -1$ ). We have proposed that synaptic plasticity induction signals should be integrated by synapses before synaptic plasticity is expressed (Elliott, 2008), generating what we have termed “integrate-and-express” models of synaptic plasticity (Elliott & Lagogiannis, 2009) in analogy with integrate-and-fire models of neuronal firing. Specifically, we have proposed that a synapse may implement a discrete low-pass filter that attenuates high-frequency noise while passing a low-frequency signal (Elliott, 2011a; Elliott & Lagogiannis, 2012). The synapse essentially decides whether or not to express synaptic plasticity depending on whether or not a synaptic filter mechanism reaches upper or lower filter thresholds, for potentiation or depression, respectively.

Fig. 1 represents this synaptic filter as a continuous time Markov process. Because potentiating and depressing induction signals are equiprobable, we need only consider a symmetric filter with equal upper and lower thresholds,  $\pm\Theta$ . Filter states are indexed by letters such as  $I$  and  $J$ , with  $I \in \{-(\Theta - 1), \dots, +(\Theta - 1)\}$ , and are represented in the figure by the circles enclosing the filter states. Rightward transitions represent potentiating induction signals that cause the filter state to increment by one; conversely, leftward transitions

represent depressing induction signals. The rates of these signals, indicated on the transitions between filter states in the figure, are just the rate of memory storage ( $r$ , which we set to unity) multiplied by the probabilities that  $\xi_i^\alpha = \pm 1$  (which are  $1/2$  in both cases). If the filter is in state  $+(\Theta - 1)$  and receives a potentiating induction signal, then the filter reaches threshold, is returned to the  $I = 0$  filter state, and a potentiation step is expressed (indicated by  $\uparrow$ ), so that the synapse's strength increases from  $s_A$  to  $s_{A+1}$ . If  $s_A = +1$ , or  $A = n$ , then of course a potentiation step cannot be expressed since the synapse is already saturated at its upper strength limit; in this case, the strength remains at  $s_n = +1$ . Similarly, if the filter is in state  $-(\Theta - 1)$  and receives a depressing induction signal, it is returned to state  $I = 0$  and a depression step is expressed (indicated by  $\downarrow$ ), so that the synapse's strength decreases from  $s_A$  to  $s_{A-1}$ . If  $s_A = -1$ , or  $A = 1$ , then a depression step cannot be expressed as the synapse is already saturated at its lower strength limit; in this case, the strength remains at  $s_1 = -1$ . The synapse thus performs a random walk on its allowed strength states  $A = 1, \dots, n$  in the presence of reflecting boundaries at  $A = 0$  and  $A = n + 1$ , implementing saturation of strength. The random walk between these strength states is driven by the underlying filter threshold events.

Let  $\mathbb{F}^+$  be the  $(2\Theta - 1) \times (2\Theta - 1)$  matrix that increments the filter state by one unit but without taking the  $I = +(\Theta - 1)$  filter state back to  $I = 0$ , so without the filter upper threshold process.  $\mathbb{F}^+$  has entries of unity on its lower diagonal and zeros elsewhere. Let  $\mathbb{D}^+$  be the matrix that takes the  $I = +(\Theta - 1)$  filter state back to  $I = 0$ . This matrix has zeros everywhere except for its entry

of unity at the  $0, +(\Theta - 1)$  position in filter indices [or position  $\Theta, (2\Theta - 1)$  with conventional indexing of matrix entries]. For  $\Theta = 3$ , for example, we have

$$\mathbb{F}^+ = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{pmatrix}, \quad \mathbb{D}^+ = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}.$$

Let the matrices  $\mathbb{F}^-$  and  $\mathbb{D}^-$  be the corresponding matrices for a decrement in filter state.  $\mathbb{F}^- = (\mathbb{F}^+)^T$ , where the superscript T denotes the transpose, and the matrix  $\mathbb{D}^-$  is zero everywhere except for unity at the  $0, -(\Theta - 1)$  position in filter indices [or position  $\Theta, 1$  with conventional indexing]. For  $\Theta = 3$ , for example,

$$\mathbb{F}^- = \begin{pmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}, \quad \mathbb{D}^- = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}.$$

As there are  $n$  strength states and  $2\Theta - 1$  filter states, the joint probability distribution of the strength and filter states of a synapse is represented by a  $(2\Theta - 1)n$ -dimensional vector. We order the entries of such vectors so that the  $A^{\text{th}}$  batch of  $2\Theta - 1$  entries corresponds to the distribution of filter states when the synapse is in strength state  $A$ . Let  $\mathbb{M}^+$  denote the  $(2\Theta - 1)n \times (2\Theta - 1)n$  matrix that implements a potentiation step and  $\mathbb{M}^-$  the corresponding matrix

implementing a depression step. For example, with  $n = 4$ ,  $\mathbb{M}^+$  and  $\mathbb{M}^-$  are given schematically by

$$\mathbb{M}^+ = \begin{pmatrix} \mathbb{F}^+ & & & \\ \mathbb{D}^+ & \mathbb{F}^+ & & \\ & \mathbb{D}^+ & \mathbb{F}^+ & \\ & & \mathbb{D}^+ & \mathbb{F}^+ + \mathbb{D}^+ \end{pmatrix}, \quad (2.3)$$

and

$$\mathbb{M}^- = \begin{pmatrix} \mathbb{F}^- + \mathbb{D}^- & \mathbb{D}^- & & \\ & \mathbb{F}^- & \mathbb{D}^- & \\ & & \mathbb{F}^- & \mathbb{D}^- \\ & & & \mathbb{F}^- \end{pmatrix}, \quad (2.4)$$

where all entries are zero unless explicitly specified. The appearance of the submatrices  $\mathbb{F}^+ + \mathbb{D}^+$  or  $\mathbb{F}^- + \mathbb{D}^-$  in the relevant sub-blocks in  $\mathbb{M}^+$  or  $\mathbb{M}^-$  reflects the fact that a saturated synapse cannot potentiate or depress, respectively, any further, but that its filter state is nevertheless returned to zero when the appropriate threshold is reached.

The matrices  $\mathbb{M}^+$  and  $\mathbb{M}^-$  implement potentiation and depression steps on synapses. The matrix superposition  $\mathbb{M} = \frac{1}{2}(\mathbb{M}^+ + \mathbb{M}^-)$  implements a plasticity operation on a synapse that is potentiation with probability  $\frac{1}{2}$  and depression with probability  $\frac{1}{2}$ . For  $n = 4$  we schematically have

$$\mathbb{M} = \frac{1}{2} \begin{pmatrix} \mathbb{F}^+ + \mathbb{F}^- + \mathbb{D}^- & \mathbb{D}^- & & \\ \mathbb{D}^+ & \mathbb{F}^+ + \mathbb{F}^- & \mathbb{D}^- & \\ & \mathbb{D}^+ & \mathbb{F}^+ + \mathbb{F}^- & \mathbb{D}^- \\ & & \mathbb{D}^+ & \mathbb{F}^+ + \mathbb{F}^- + \mathbb{D}^+ \end{pmatrix}. \quad (2.5)$$

The equilibrium joint probability distribution of strength and filter states for a synapse is determined by the eigenvector of  $\mathbb{M}$  with unit eigenvalue. In equilibrium, for  $n = 2$  synapses we previously found that the probability distribution  $\mathcal{B}$  of filter states has components

$$\mathcal{B}_I = \frac{1}{\Theta^2} (\Theta - |I|), \quad (2.6)$$

regardless of the strength state (Elliott & Lagogiannis, 2012). This distribution just corresponds to the (suitably normalised) eigenvector of the matrix  $\frac{1}{2}(\mathbb{M}^+ + \mathbb{M}^- + \mathbb{D}^+ + \mathbb{D}^-)$  with unit eigenvalue. Because the distribution  $\mathcal{B}$  is symmetric about its central,  $I = 0$  component, we have that  $\mathbb{D}^+ \mathcal{B} \equiv \mathbb{D}^- \mathcal{B}$ . The vector  $\mathcal{B}$  is therefore also an eigenvector of the two matrices  $\frac{1}{2}(\mathbb{M}^+ + \mathbb{M}^- + 2\mathbb{D}^+)$  and  $\frac{1}{2}(\mathbb{M}^+ + \mathbb{M}^- + 2\mathbb{D}^-)$  with unit eigenvalue. If we consider the probability distribution

$$\mathcal{A}^T = \frac{1}{n} \left( \underbrace{\mathcal{B}^T | \dots | \mathcal{B}^T}_n \right), \quad (2.7)$$

with  $\mathcal{A}$  being a  $(2\Theta - 1)n$ -dimensional vector with  $\mathcal{B}$  occurring once for each strength state  $A$ ,  $A = 1, \dots, n$ , then it is clear from the block structure of the matrix  $\mathbb{M}$  that  $\mathcal{A}$  is an eigenvector of  $\mathbb{M}$  with unit eigenvalue. The vector  $\mathcal{A}$  is therefore the equilibrium joint distribution of filter and strength states for general  $n$ . All filter states are therefore distributed according to  $\mathcal{B}$  in equilibrium, regardless of the value of  $n$ , and we see that the strength states are themselves uniformly distributed with probability  $1/n$  in equilibrium because of the common  $1/n$  scaling factor in  $\mathcal{A}$ .

It is against the background of this equilibrium distribution  $\mathcal{A}$  that the definite memory  $\xi^0$  is stored at time  $t = 0^-$  s. Synapse  $i$  will have distribution  $\mathbb{M}^\pm \mathcal{A}$  at time  $t = 0$  s after the storage of  $\xi^0$  depending on the definite sign of  $\xi_i^0 = \pm 1$ . For all subsequent memories  $\xi^\alpha$ ,  $\alpha > 0$ , however, the relevant matrix operator is  $\mathbb{M}$ , which averages over both possible induction signals for later memories at any given synapse, allowing us to average over all possible subsequent memories as we are not interested in any particular realisation of subsequent memories. With symmetric filters (equal upper and lower thresholds,  $\pm\Theta$ ), the two distributions  $\mathbb{M}^\pm \mathcal{A}$  are exact mirror images of each other. For  $n = 2$ , we previously showed that synapses experiencing an initial potentiating induction signal ( $\xi_i^0 = +1$ ) and those experiencing an initial depressing induction signal ( $\xi_i^0 = -1$ ) contribute identically to the tracked memory signal  $h(t)$  because the roles of weak (strength  $-1$ ) and strong (strength  $+1$ ) synapses are reversed in their contributions to  $h(t)$ , depending on the signs of the induction signals (Elliott & Lagogiannis, 2012). Specifically, if we define  $\tilde{S}_i(t) = \xi_i^0 S_i(t)$ , so that

$$h(t) = \frac{1}{N} \sum_{i=1}^N \tilde{S}_i(t), \quad (2.8)$$

then the various  $\tilde{S}_i(0)$  are all identically distributed random variables, regardless of  $i$ . Furthermore, since all synapses subsequently experience only a superposition of induction signals via the same matrix operator  $\mathbb{M}$  in order to average over all possible later memories  $\xi^\alpha$ ,  $\alpha > 0$ , if  $\tilde{S}_i(t)$  are identically distributed at  $t = 0$  s, then they remain identically distributed for all time. It is



then a statement in elementary probability that

$$\mu(t) = \mathbf{E}[\tilde{S}(t)], \quad (2.9)$$

$$\sigma(t)^2 = \frac{1}{N} \text{Var}[\tilde{S}(t)] + \left(1 - \frac{1}{N}\right) \text{Cov}(t), \quad (2.10)$$

where  $\mathbf{E}[\tilde{S}(t)]$  and  $\text{Var}[\tilde{S}(t)]$  denote the mean and variance, respectively, of any one of the  $\tilde{S}_i(t)$  variables, and  $\text{Cov}(t)$  denotes the covariance between any two of them. This equivalence and resulting simplification arises because, for  $n = 2$  strength states, the two strengths are treated symmetrically. It is therefore clear that for general  $n$ , provided that the various strengths are symmetrically distributed around zero (or, in general, around their mean value), so that if for some strength state  $A$  there exists a strength state  $B$  such that  $s_A = -s_B$ , as is true for Eq. (2.1), then the same arguments go through. In fact, these arguments also go through for any model of synaptic plasticity in which potentiation and depression processes are treated completely symmetrically (Elliott, 2014), and not just for filter-based mechanisms of synaptic plasticity as considered here and elsewhere.

This equivalence between the two distributions  $\mathbb{M}^\pm \mathcal{A}$  in terms of their contributions to  $h(t)$  therefore means that we need only consider the joint probability distribution of the filter and tilded-strength (rather than strength) states, and we can therefore restrict without loss of generality to considering only, say, the distribution of  $\mathbb{M}^+ \mathcal{A}$  at time  $t = 0$  s. Defining  $\Delta = \mathbb{D}^+ \mathcal{B}$  (or equivalently  $\Delta = \mathbb{D}^- \mathcal{B}$ ), which is a  $(2\Theta - 1)$ -dimensional vector with

components  $\Delta_I = \Theta^{-2} \delta_{I,0}$ , where  $\delta_{I,J}$  is the Kronecker delta symbol, we then have that

$$(\mathbb{M}^+ \mathcal{A})^T = \frac{1}{n} \left( (\mathbb{F}^+ \mathcal{B})^T \left| \underbrace{(\mathbb{F}^+ \mathcal{B})^T + \Delta^T \right| \cdots \left| (\mathbb{F}^+ \mathcal{B})^T + \Delta^T \right| (\mathbb{F}^+ \mathcal{B})^T + 2\Delta^T \right) \quad (2.11)$$

for the probability distribution of states immediately after the storage of memory  $\xi^0$ . This follows directly from the block structure of  $\mathbb{M}^+$  expressed in terms of  $\mathbb{F}^+$  and  $\mathbb{D}^+$ . The contributions from  $\Delta$  to the  $n - 2$  intermediate strength states with  $2 \leq A \leq n - 1$  arise from the upper filter threshold processes occurring in states  $1 \leq A \leq n - 2$ . There is no such contribution to the  $A = 1$  state because there is no lower,  $A = 0$  state. There are two such contributions to the  $A = n$  state because one arises from the upper threshold process from the  $A = n - 1$  state and the other from the  $A = n$  state's own upper threshold process, which cannot induce an increment in strength because of saturation. Summing over the filter states for each strength state in Eq. (2.11), we see that

$$\left. \begin{aligned} \frac{1}{n} \sum_J (\mathbb{F}^+ \mathcal{B} + 0 \Delta)_J &= \frac{1}{n} \left(1 - \frac{1}{\Theta^2}\right), \\ \frac{1}{n} \sum_J (\mathbb{F}^+ \mathcal{B} + 1 \Delta)_J &= \frac{1}{n}, \\ \frac{1}{n} \sum_J (\mathbb{F}^+ \mathcal{B} + 2 \Delta)_J &= \frac{1}{n} \left(1 + \frac{1}{\Theta^2}\right), \end{aligned} \right\} \quad (2.12)$$

which give the probabilities of the various strength states at time  $t = 0$  s. The intermediate (tilded-)strength states  $2 \leq A \leq n - 1$  therefore continue to have probability  $1/n$  immediately after the storage of  $\xi^0$ , while the probability of state  $A = 1$  decreases to  $\frac{1}{n} \left(1 - \frac{1}{\Theta^2}\right)$  and that of  $A = n$  increases to  $\frac{1}{n} \left(1 + \frac{1}{\Theta^2}\right)$ .

The initial mean memory signal  $\mu(0)$  is therefore

$$\mu(0) = \frac{2}{n} \frac{1}{\Theta^2}. \quad (2.13)$$

Finally, we note that the matrix operator  $\mathbb{M}$  acts on the state  $\mathbb{M}^+ \mathcal{A}$  in such a way that the probabilities of the intermediate states  $2 \leq A \leq n-1$  remain unchanged, at  $1/n$ , and this remains true for all time. Only the probabilities of the  $A=1$  and  $A=n$  states change over time. These always differ equally but oppositely from  $1/n$ , because the total probability of all strength states must sum to unity.<sup>2</sup> It is therefore straightforward to see that

$$\mathbf{E}[\tilde{S}(t)^2] = \frac{1}{3} \frac{n+1}{n-1}, \quad (2.14)$$

independent of  $t$ , because  $\frac{1}{n} \sum_{A=1}^n s_A^2 = \frac{1}{3} \frac{n+1}{n-1}$  and any deviations of the  $A=1$  and  $A=n$  states from probability  $1/n$  cancel out in  $\mathbf{E}[\tilde{S}(t)^2]$  because  $s_1^2 = s_n^2 = 1$ .

### 3 Analytical Results

With our general formalism established, we may now derive an analytical expression for the mean memory signal  $\mu(t)$ . Analytical expressions for  $\sigma(t)$  are very much harder to derive because of the covariance term (Elliott & Lagogian-

---

<sup>2</sup>The maximum possible initial memory signal is thus  $2/n$ , achieved for a trivial,  $\Theta = 1$  filter.

nis, 2012). Even for  $n = 2$  the generating matrix defining the Markov process lacks a complete set of eigenvectors (that is, the matrix is defective), so tensor product methods cannot be used to extract higher-order statistics (Elliott & Lagogiannis, 2012). This difficulty remains for  $n > 2$ . For our purposes here, it is sufficient to approximate the variance  $\sigma(t)^2$  using either  $\text{Var}[\tilde{S}(t)]/N$  or even more simply  $\mathbf{E}[\tilde{S}^2(t)]/N$ , but where necessary we employ numerical matrix methods to compute the full result.

### 3.1 Laplace Transform of $\mu(t)$

Let  $f_{I|J}(t)$  be the probability of a transition from filter state  $J$  to filter state  $I$  in time  $t$  without filter thresholds being reached. An expression for  $f_{I|J}(t)$  was derived in Elliott & Lagogiannis (2012) although its explicit form is not required. Let  $G_J^\pm(t)$  be the densities for first escapes through the upper and lower filter thresholds, respectively, at time  $t$ . We have that  $G_J^\pm(t) = \frac{r}{2} f_{\pm(\Theta-1)|J}(t)$ , explicitly including the rate factor  $r$ . For a symmetric filter, we have that  $G_{-J}^\pm(t) = G_{+J}^\mp(t)$ . Let  $H_J(t)$  be the probability of not having reached either filter threshold, starting from filter state  $J$ , in time  $t$ . We have that  $H_J(t) = \sum_{I=-(\Theta-1)}^{+(\Theta-1)} f_{I|J}(t)$ , but also

$$H_J(t) = 1 - \int_0^t dt_1 [G_J^+(t_1) + G_J^-(t_1)]. \quad (3.1)$$

$H_J(t)$  is the probability that a random walk on filter states, starting from state  $J$ , has yet to reach threshold (or absorbing boundaries) in time  $t$ . The

sum  $G_J^+(t) + G_J^-(t)$  is the probability density for reaching either threshold, so the integral in Eq. (3.1) gives the integrated probability density or just the probability of having reached either threshold in time  $t$ . The probability of not having reached either threshold then follows directly. Simple expressions may be derived for the Laplace transforms of  $G_J^\pm(t)$  (Elliott, 2011a). If  $\widehat{G}_J^\pm(s)$  denote these Laplace transforms with transformed variable  $s$ , then

$$\widehat{G}_J^\pm(s) = \frac{[\Phi_+(s)]^{\Theta \pm J} - [\Phi_-(s)]^{\Theta \pm J}}{[\Phi_+(s)]^{2\Theta} - [\Phi_-(s)]^{2\Theta}}, \quad (3.2)$$

where  $\Phi_\pm(s)$  are the two solutions of  $\Phi^2 - 2(1 + s/r)\Phi + 1 = 0$ , where the rate factor  $r$  is retained for generality, so that  $\Phi_+(s)\Phi_-(s) = 1$  and  $\Phi_+(s) + \Phi_-(s) = 2(1 + s/r)$ .

Let  $p_{I|J}^{A|B}(t)$  be the probability of a transition from filter state  $J$  and strength state  $B$  to filter state  $I$  and strength state  $A$  in time  $t$ . Then, generalising the argument in Elliott & Lagogiannis (2012), we may write down the system of

renewal equations for  $p_{I|J}^{A|B}(t)$  in terms of  $f_{I|J}(t)$  and  $G_J^\pm(t)$ :

$$\begin{aligned}
p_{I|J}^{A|1}(t) &= f_{I|J}(t) \delta^{A,1} \\
&\quad + \int_0^t dt_1 \left[ p_{I|0}^{A|1}(t-t_1) G_J^-(t_1) + p_{I|0}^{A|2}(t-t_1) G_J^+(t_1) \right], \\
p_{I|J}^{A|B}(t) &= f_{I|J}(t) \delta^{A,B} \\
&\quad + \int_0^t dt_1 \left[ p_{I|0}^{A|B-1}(t-t_1) G_J^-(t_1) + p_{I|0}^{A|B+1}(t-t_1) G_J^+(t_1) \right], \\
p_{I|J}^{A|n}(t) &= f_{I|J}(t) \delta^{A,n} \\
&\quad + \int_0^t dt_1 \left[ p_{I|0}^{A|n-1}(t-t_1) G_J^-(t_1) + p_{I|0}^{A|n}(t-t_1) G_J^+(t_1) \right],
\end{aligned} \tag{3.3}$$

where the middle equation applies only for  $2 \leq B \leq n-1$  so that it is the general case rather than the boundary cases at  $B=1$  or  $B=n$ . The inhomogeneous term on the right hand sides (RHSs) arise only when  $A=B$  and thus allow the possibility of changes in filter state without any threshold processes arising. The homogeneous terms consider threshold processes leading to changes in strength state (with or without saturation or reflecting boundary dynamics), followed by further state transition processes after the first threshold process. Let  $\mathbb{P}_{I|J}(t)$  be a matrix with components  $p_{I|J}^{A|B}(t)$ , so a matrix of strength change probabilities in time  $t$  for given initial and final filter states,



$\delta(t)\mathbb{I} - \mathbb{G}_0(t)$  where  $\delta(t)$  is the Dirac delta function, or its Laplace transform,

$$\widehat{\mathbb{W}}(s) = \mathbb{I} - \widehat{\mathbb{G}}_0(s), \quad (3.6)$$

we have  $\widehat{\mathbb{P}}_{I|0}(s) = \widehat{f}_{I|0}(s) \widehat{\mathbb{W}}^{-1}(s)$ , so that

$$\widehat{\mathbb{P}}_{I|J} = \widehat{f}_{I|J} \mathbb{I} + \widehat{f}_{I|0} \widehat{\mathbb{W}}^{-1} \widehat{\mathbb{G}}_J, \quad (3.7)$$

where we drop the Laplace argument  $s$  for notational simplicity. The second term in Eq. (3.5) has transformed into  $\widehat{f}_{I|0} \widehat{\mathbb{W}}^{-1} \widehat{\mathbb{G}}_J$ . We note that, formally,  $\widehat{\mathbb{W}}^{-1} = \sum_{m=0}^{\infty} [\widehat{\mathbb{G}}_0]^m$ . A general term in this second term on the RHS of Eq. (3.7) is therefore of the form  $\widehat{f}_{I|0} [\widehat{\mathbb{G}}_0]^m \widehat{\mathbb{G}}_J$ . This product of matrices in Laplace-transform space represents a first change in strength starting from initial filter state  $J$ , followed by precisely  $m$  changes in strength associated with transitions from the zero filter state back to the zero filter state via filter threshold processes, and finally a change in filter state from the zero state to state  $I$  without any change in strength. Eq. (3.7) therefore decomposes the overall transition matrix  $\mathbb{P}_{I|J}$  into elementary, fundamental steps, and sums over all possibilities.

To compute  $\mu(t)$ , we must sum  $p_{I|J}^{A|B}(t)$  over the initial and final states with suitable weighting factors. The final filter state  $I$  is irrelevant to  $\mu(t)$  so we directly sum over  $I$ . The final strength state  $A$  must be weighted by  $s_A$ . Defining the vector  $\mathbf{S}^T = (s_1, \dots, s_n)$  as the vector of strengths  $s_A$ , we then



have

$$\sum_I \mathbf{S}^T \widehat{\mathbb{P}}_{I|J} = \widehat{H}_J \mathbf{S}^T + \widehat{H}_0 \mathbf{S}^T \widehat{\mathbb{W}}^{-1} \widehat{\mathbb{G}}_J. \quad (3.8)$$

To sum over the initial state, we must sum over the components of the vector  $\mathbb{M}^+ \mathbf{A}$  because this gives the state of the system immediately after the storage of the tracked memory  $\boldsymbol{\xi}^0$  whose mean tracked memory signal  $\mu(t)$  we wish to determine. Let  $\boldsymbol{\sigma}_J$  be the  $n$ -dimensional vector of probabilities of a synapse's strength at time  $t = 0$  s for each particular filter state  $J$ . Then we have that  $\widehat{\mu}(s) \equiv \sum_{I,J} \mathbf{S}^T \widehat{\mathbb{P}}_{I|J}(s) \boldsymbol{\sigma}_J$ , so that

$$\widehat{\mu} = \sum_J \widehat{H}_J (\mathbf{S}^T \boldsymbol{\sigma}_J) + \widehat{H}_0 \mathbf{S}^T \widehat{\mathbb{W}}^{-1} \left( \sum_J \widehat{\mathbb{G}}_J \boldsymbol{\sigma}_J \right). \quad (3.9)$$

Structurally, this equation for the mean memory signal should be compared to the equation for  $\mathbb{P}_{I|J}$  in Eq. (3.7). The first term on the RHS describes the contribution to  $\mu$  that arises from the initial storage of the tracked memory but without any subsequent changes in strength. This contribution decays away monotonically, controlled by  $H_J$ , because the probability of no changes in strength drops to zero as time increases. The second term on the RHS describes the contributions arising from subsequent changes in strength after the storage of the tracked memory. Again we may expand  $\widehat{\mathbb{W}}^{-1}$  as before and observe contributions from definite numbers of strength changes.

The vectors  $\boldsymbol{\sigma}_J$  can be read off from Eq. (2.11). Noting that the matrix  $\mathbb{F}^+$  is just a shift operator on filter states, so that  $(\mathbb{F}^+ \boldsymbol{\mathcal{B}})_J = \Theta^{-2}(\Theta - |J - 1|)$ , we

can write

$$\boldsymbol{\sigma}_J = \begin{cases} \frac{1}{n} \frac{1}{\Theta^2} (\Theta - |J - 1|) \mathbf{n} & \text{for } J \neq 0 \\ \frac{1}{n} \frac{1}{\Theta^2} (\Theta \mathbf{n} + \mathbf{v}) & \text{for } J = 0 \end{cases} \quad (3.10)$$

where  $\mathbf{v}^\text{T} = (-1, \overbrace{0, \dots, 0}^{n-2}, +1)$  and  $\mathbf{n}^\text{T} = (1, 1, \dots, 1, 1)$ , both being  $n$ -dimensional vectors. The  $J = 0$  form accounts for the various  $\mathbf{\Delta}$  contributions arising in Eq. (2.11). Note that  $\mathbf{S}^\text{T} \mathbf{n} \equiv 0$ . We then have that  $\mathbf{S}^\text{T} \boldsymbol{\sigma}_J = \frac{2}{n} \frac{1}{\Theta^2} \delta_{J,0}$ , so that the first term on the RHS of Eq. (3.9) becomes  $\frac{2}{n} \frac{1}{\Theta^2} \widehat{H}_0$ . Evaluating  $\widehat{\mathbb{G}}_J \boldsymbol{\sigma}_J$ , we readily find that

$$\widehat{\mathbb{G}}_J \boldsymbol{\sigma}_J = \begin{cases} \frac{1}{n} \frac{1}{\Theta^2} (\Theta - |J - 1|) \left[ (\widehat{G}_J^+ + \widehat{G}_J^-) \mathbf{n} + (\widehat{G}_J^+ - \widehat{G}_J^-) \mathbf{v} \right] & \text{for } J \neq 0 \\ \frac{1}{n} \frac{1}{\Theta^2} \widehat{G}_0 \left[ 2\Theta \mathbf{n} + \mathbf{u} \right] & \text{for } J = 0 \end{cases} \quad (3.11)$$

where  $\mathbf{u}^\text{T} = (-1, -1, \overbrace{0, \dots, 0}^{n-4}, +1, +1)$ .<sup>3</sup> To compute  $\mathbf{S}^\text{T} \widehat{\mathbb{W}}^{-1} \widehat{\mathbb{G}}_J \boldsymbol{\sigma}_J$ , we first write  $\widehat{\mathbb{W}} = \mathbb{I} - \widehat{G}_0 \mathbb{C}$ , setting  $G_0(t) = G_0^\pm(t)$  as  $G_0^+(t) = G_0^-(t)$  for a symmetric filter, where we define the matrix  $\mathbb{C}$  by

$$\mathbb{C} = \begin{pmatrix} 1 & 1 & & & & \\ & 1 & 0 & 1 & & \\ & & \ddots & \ddots & \ddots & \\ & & & 1 & 0 & 1 \\ & & & & 1 & 1 \end{pmatrix}. \quad (3.12)$$

---

<sup>3</sup>This form of  $\mathbf{u}$  is valid for  $n \geq 4$ . For  $n = 3$ ,  $\mathbf{u}^\text{T} = (-1, 0, +1)$  and for  $n = 2$ ,  $\mathbf{u}^\text{T} = (0, 0)$ . We may calculate with the general form for  $n \geq 4$  and then confirm that our final results are in fact also valid for  $n = 2$  and  $n = 3$ .

We may then formally write  $\widehat{\mathbb{W}}^{-1} = \sum_{m=0}^{\infty} (\widehat{G}_0)^m \mathbb{C}^m$ . We observe that if some  $n$ -dimensional vector  $\mathbf{w}$  is antisymmetric about its centre, so that  $w_A = -w_{n+1-A}$  for any component  $A$ , then so is the vector  $\mathbb{C}\mathbf{w}$ . Of course,  $\mathbf{S}$  is such a vector. Thus, the vector  $\mathbf{S}^T \widehat{\mathbb{W}}^{-1}$  is antisymmetric about its centre and so  $\mathbf{S}^T \widehat{\mathbb{W}}^{-1} \mathbf{n} \equiv 0$ . Hence, all the  $\mathbf{n}$  terms in  $\widehat{\mathbb{G}}_J \boldsymbol{\sigma}_J$  in Eq. (3.11) are killed by  $\mathbf{S}^T \widehat{\mathbb{W}}^{-1}$  and only the  $\mathbf{u}$  and  $\mathbf{v}$  terms survive. For the terms involving  $\mathbf{v}$ , we require

$$\sum_{J \neq 0} (\Theta - |J - 1|) (\widehat{G}_J^+ - \widehat{G}_J^-) = 2 \sum_{J > 0} (\widehat{G}_J^+ - \widehat{G}_J^-). \quad (3.13)$$

For the term involving  $\mathbf{u}$ , we observe that  $\mathbf{u} = \mathbb{C}\mathbf{v}$  and that

$$\widehat{\mathbb{W}}^{-1} \mathbb{C} = (\widehat{G}_0)^{-1} (\widehat{\mathbb{W}}^{-1} - \mathbb{I}). \quad (3.14)$$

Putting all this together, we may finally write

$$\widehat{\mu} = \frac{1}{n} \frac{1}{\Theta^2} \widehat{H}_0 \left[ 1 + 2 \sum_{J > 0} (\widehat{G}_J^+ - \widehat{G}_J^-) \right] \mathbf{S}^T \widehat{\mathbb{W}}^{-1} \mathbf{v}. \quad (3.15)$$

For  $n = 2$ , it is easy to see that  $\mathbf{S}^T \widehat{\mathbb{W}}^{-1} \mathbf{v} = 2$  [this is true for any density  $G_0(t)$ ], so that we obtain

$$\widehat{\mu}_2 = \frac{1}{\Theta^2} \widehat{H}_0 \left[ 1 + 2 \sum_{J > 0} (\widehat{G}_J^+ - \widehat{G}_J^-) \right], \quad (3.16)$$

which is precisely the expression that we obtained before (Elliott & Lagorian-

nis, 2012). We may then write

$$\widehat{\mu}(s) = \frac{1}{n} \widehat{\mu}_2(s) \mathbf{S}^T \widehat{\mathbb{W}}^{-1}(s) \mathbf{v}. \quad (3.17)$$

We will interpret the terms in this equation for  $\widehat{\mu}(s)$  momentarily.

It is striking that the form for  $\widehat{\mu}$  in Eq. (3.15) factorises into a form involving  $\widehat{\mu}_2$  and the new contribution  $\mathbf{S}^T \widehat{\mathbb{W}}^{-1} \mathbf{v}$ . Immediately after the storage of the tracked memory  $\boldsymbol{\xi}^0$  at time  $t = 0^-$  s, synapses may be in different filter states, with probability distribution governed by Eq. (2.11). Consider, however, a scenario in which all synapses at time  $t = 0$  s are in filter state  $I = 0$ . Over time they escape through the upper or lower filter thresholds and are returned to the  $I = 0$  filter state, with associated steps in synaptic strength where possible. In this scenario, synapses therefore perform a non-Markovian but renewing random walk on the strength states  $A = 1, \dots, n$  with waiting times between transitions governed by the densities  $G_0^\pm(t)$  [which for symmetric filters are equal,  $G_0^\pm(t) = G_0(t)$ ]. The matrix  $\mathbb{C}$ , or more correctly the matrix  $\frac{1}{2}\mathbb{C}$  is precisely a stochastic matrix that implements an unbiased (that is, symmetric) random walk on  $n$  discrete states between two reflecting boundaries. In section 4 of Elliott (2010), we derived general results for renewal processes on bounded intervals in the presence of non-exponential waiting times between transitions. Adapting those results to our notation here, if  $\mathbb{P}(t)$  is the matrix of transition probabilities with elements  $p^{AB}(t)$  for transitions from strength state B to strength state A in time  $t$  (dropping filter indices because they are ir-

relevant to the argument here), then we would have for the scenario considered here that [cf. Eq. (4.8) in Elliott (2010)]

$$\widehat{\mathbb{P}}(s) [\mathbb{I} - \widehat{G}_0(s) \mathbb{C}] = \widehat{H}_0(s) \mathbb{I}, \quad (3.18)$$

so that

$$\widehat{\mathbb{P}}(s) = \widehat{H}_0(s) \widehat{\mathbb{W}}^{-1}(s). \quad (3.19)$$

The interpretation of this equation for  $\widehat{\mathbb{P}}(s)$  can be seen clearly by expanding  $\widehat{\mathbb{W}}^{-1}(s)$  in powers of the stochastic matrix  $\frac{1}{2} \mathbb{C}$ ,

$$\widehat{\mathbb{P}}(s) = \widehat{H}_0(s) \sum_{m=0}^{\infty} [2 \widehat{G}_0(s) \times \frac{1}{2} \mathbb{C}]^m, \quad (3.20)$$

and then taking the inverse Laplace transform using the convolution theorem,

$$\begin{aligned} \mathbb{P}(t) &= H_0(t) \mathbb{I} + \int_0^t dt_1 H_0(t - t_1) [2 G_0(t_1) \times \frac{1}{2} \mathbb{C}] \\ &+ \int_0^t dt_1 \int_0^{t-t_1} dt_2 H_0(t - t_1 - t_2) [2 G_0(t_2) \times \frac{1}{2} \mathbb{C}] [2 G_0(t_1) \times \frac{1}{2} \mathbb{C}] \\ &+ \dots \end{aligned} \quad (3.21)$$

The third term on the RHS of this equation, for example, represents a change in strength occurring at time  $0 \leq t_1 \leq t$  with probability density  $2 G_0(t_1)$  followed by a second change in strength at later time  $0 \leq t_1 + t_2 \leq t$  with probability density  $2 G_0(t_2)$  and then no subsequent changes in strength, indicated by the presence of the waiting time function  $H_0(t - t_1 - t_2)$  over the

remaining time interval. The changes in strength are signalled by the presence of the stochastic matrix  $\frac{1}{2} \mathbb{C}$ , which implements a single step in strength space subject to possible saturation. These changes in strength are governed by filter transitions from the zero state through either filter threshold back to the zero state, so are governed by the probability density  $2G_0(t)$ . Similarly, the general term  $\widehat{H}_0(s) [2\widehat{G}_0(s) \times \frac{1}{2} \mathbb{C}]^m$  in  $\widehat{\mathbb{P}}(s)$  corresponds to the occurrence of precisely  $m$  filter threshold escape processes and therefore  $m$  possible strength changes giving rise to the  $m$  occurrences of  $2\widehat{G}_0(s) \times \frac{1}{2} \mathbb{C}$ , followed by no filter threshold escape processes, giving rise to the  $\widehat{H}_0(s)$  waiting time factor. To compute  $\mu(t)$  in this scenario, suppose that the initial probability distribution of strength states at  $t = 0$  s is governed by that generated by the storage of  $\boldsymbol{\xi}^0$ , despite considering all filter states to be  $I = 0$  at time  $t = 0$  s. This distribution, from Eq. (2.12), is just

$$\mathbf{a} = \frac{1}{n} \left( \mathbf{n} + \frac{1}{\Theta^2} \mathbf{v} \right). \quad (3.22)$$

In this scenario, the mean memory signal is just  $\mu(t) = \mathbf{S}^T \mathbb{P}(t) \mathbf{a}$ , or

$$\widehat{\mu} = \frac{1}{n} \frac{1}{\Theta^2} \widehat{H}_0 \mathbf{S}^T \widehat{\mathbb{W}}^{-1} \mathbf{v}. \quad (3.23)$$

This result is obtained essentially by integrating out synapses' internal filter states and instead dealing directly with transitions in synaptic strength (cf. Elliott, 2010). Comparing this result to Eq. (3.15), we see that we have all the

terms except for that in square brackets. Most of the terms in the full form of  $\widehat{\mu}$  therefore arise from the unbiased, renewing random walk in bounded strength space governed by the non-exponential waiting times  $G_0^\pm(t) = G_0(t)$ . The remaining term, in square brackets, must arise because of the preparation of the initial filter states, and specifically because the initial filter states are not, in general,  $I = 0$  immediately after the storage of  $\boldsymbol{\xi}^0$ . In Appendix A, we provide an alternative derivation of Eq. (3.15) using the argument in this paragraph, but taking into account the correct distribution of initial filter states at time  $t = 0$  s.

With these considerations in hand, we may now interpret the two factors,  $\widehat{\mu}_2(s)$  and  $\mathbf{S}^T \widehat{\mathbb{W}}^{-1}(s) \mathbf{v}$ , appearing in Eq. (3.17) for  $\widehat{\mu}(s)$ . Up to overall parameters and the function  $\widehat{H}_0(s)$ , the factor  $\widehat{\mu}_2(s)$  arises purely from the change in the distribution of filter states induced by the storage of the tracked memory  $\boldsymbol{\xi}^0$  and is independent of the number of states of synaptic strength,  $n$ . This change in the distribution of filter states is therefore already captured entirely by the particular case of binary synapses,  $n = 2$ , studied before (Elliott & Lagogiannis, 2012). The contribution to the tracked memory signal from multiple strength states and therefore general  $n$  is essentially confined to the second factor,  $\mathbf{S}^T \widehat{\mathbb{W}}^{-1}(s) \mathbf{v}$ . It is important to note, however, that the filter dynamics continue to exert an influence in this factor through the probability density for escape through either filter threshold from the zero state,  $2G_0(t)$ . As we have seen, the inverse matrix  $\widehat{\mathbb{W}}^{-1}(s)$  essentially generates a random walk on strength states between reflecting barriers by summing over all possi-

ble processes. The vector  $\mathbf{S}$  in  $\mathbf{S}^T \widehat{\mathbb{W}}^{-1}(s) \mathbf{v}$  weights the final states by their strengths, and the vector  $\mathbf{v}$  is proportional to the shift in the distribution of strength states induced by the storage of the tracked memory  $\boldsymbol{\xi}^0$ , since only the  $A = 1$  and  $A = n$  strength states change their probabilities. It is indeed striking, and even remarkable, that  $\widehat{\mu}(s)$  factorises in this way. However, this factorisation of course reflects the separation of the underlying dynamics into an initial phase governed by the change in filter distributions induced by the storage of  $\boldsymbol{\xi}^0$ , and a later phase governed by transitions in synaptic strength driven by a renewing but non-Markovian random walk regulated by the probability density  $2G_0(t)$ . This probability density arises because once the state of a filter immediately after the storage of  $\boldsymbol{\xi}^0$  has been forgotten due to a filter threshold process leading to a return of the filter to the zero state, later processes naturally decompose into dynamics starting from the zero filter state and returning to the zero filter state via a filter threshold process, leading (possibly) to a change in strength. We must now determine the contribution from these later processes by explicitly computing  $\mathbf{S}^T \widehat{\mathbb{W}}^{-1}(s) \mathbf{v}$ .

### 3.2 Computation of $\mathbf{S}^T \widehat{\mathbb{W}}^{-1}(s) \mathbf{v}$

In order to obtain the full form for  $\widehat{\mu}(s)$ , we must now compute  $\mathbf{S}^T \widehat{\mathbb{W}}^{-1}(s) \mathbf{v}$ . We may do this using two different methods, both of which lead to different but useful approximations, which we shall examine below.



### 3.2.1 Method 1: Eigenanalysis

The first method is a direct computation of  $\widehat{\mathbb{W}}^{-1}$ . As the matrix  $\mathbb{C}$  is tridiagonal, its eigenstructure may be computed explicitly. We give the details in Appendix B. The stochastic matrix  $\frac{1}{2}\mathbb{C}$  has eigenvalues  $\lambda_m$  given by

$$\lambda_m = \cos \frac{m\pi}{n}, \quad m = 0, \dots, n-1, \quad (3.24)$$

and corresponding normalised eigenvectors  $\mathbf{e}^m$  with components  $e_A^m$  given by

$$e_A^m = \begin{cases} \sqrt{\frac{1}{n}} & \text{for } m = 0 \\ \sqrt{\frac{2}{n}} \cos \left[ \frac{m\pi}{2n}(2A-1) \right] & \text{for } m = 1, \dots, n-1 \end{cases}. \quad (3.25)$$

The matrix inverse  $\widehat{\mathbb{W}}^{-1} = \left[ \mathbb{I} - (2\widehat{G}_0)(\frac{1}{2}\mathbb{C}) \right]^{-1}$  is therefore given by

$$\widehat{\mathbb{W}}^{-1} = \sum_{m=0}^{n-1} \frac{1}{1 - 2\widehat{G}_0 \lambda_m} \mathbf{e}^m (\mathbf{e}^m)^\top. \quad (3.26)$$

We find that

$$\mathbf{S}^\top \mathbf{e}^m = \begin{cases} 0 & \text{for } m = 0 \\ -\frac{[1 - (-1)^m]}{\sqrt{2n(n-1)}} \frac{\cos \frac{m\pi}{2n}}{\sin^2 \frac{m\pi}{2n}} & \text{for } m = 1, \dots, n-1 \end{cases}, \quad (3.27)$$

and

$$(\mathbf{e}^m)^\top \mathbf{v} = -\sqrt{\frac{2}{n}} [1 - (-1)^m] \cos \frac{m\pi}{2n}, \quad (3.28)$$

where strictly this latter equation is valid only for  $m = 1, \dots, n - 1$  but it also generates the required result (of zero) for  $m = 0$ . Finally, then, we have that

$$\mathbf{S}^T \widehat{\mathbb{W}}^{-1} \mathbf{v} = \frac{2}{n(n-1)} \sum_{m=1}^{n-1} \frac{[1 - (-1)^m]}{1 - 2\widehat{G}_0 \cos \frac{m\pi}{n}} \cot^2 \frac{m\pi}{2n}. \quad (3.29)$$

### 3.2.2 Method 2: Generating Function

The first method above for computing  $\mathbf{S}^T \widehat{\mathbb{W}}^{-1} \mathbf{v}$  involves a direct eigenanalysis of the stochastic matrix  $\frac{1}{2}\mathbb{C}$ . We may instead use a generating function approach that implicitly determines all the matrix powers  $(\frac{1}{2}\mathbb{C})^m$  without a direct eigenanalysis and in so doing obtain a very different representation of  $\mathbf{S}^T \widehat{\mathbb{W}}^{-1} \mathbf{v}$ .

The stochastic matrix  $\frac{1}{2}\mathbb{C}$  generates a random walk between reflecting boundaries, so consider the probability distribution  $\mathbf{p}(m)$  after  $m$  discrete time steps, starting from some initial distribution  $\mathbf{p}(0)$ . Then  $\mathbf{p}(m) = (\frac{1}{2}\mathbb{C})^m \mathbf{p}(0)$ . Let the components of  $\mathbf{p}(m)$  be  $p_A(m)$ ,  $A = 1, \dots, n$ , the probability of being in state  $A$  after  $m$  time steps. Define a generating function for  $p_A(m)$  over both states  $A$  and time steps  $m$  by writing

$$F(w, z) = \sum_{m=0}^{\infty} w^m \sum_{A=1}^n z^A p_A(m) = \sum_{A=1}^n z^A \sum_{m=0}^{\infty} w^m p_A(m). \quad (3.30)$$

If we set  $w = 2\widehat{G}_0$ , then  $F(2\widehat{G}_0, z)$  essentially determines  $\widehat{\mathbb{W}}^{-1}$ , because  $p_A(m) = [(\frac{1}{2}\mathbb{C})^m \mathbf{p}(0)]_A$ . By writing out the set of equations for  $f_A(m+1)$  in terms of  $f_A(m)$  and taking into account the reflecting boundary conditions, a

standard calculation (see, for example, van Kampen, 1992) produces

$$F(w, z) = \frac{w z (1 - z) \mathcal{P}_1(w) - w z^{n+1} (1 - z) \mathcal{P}_n(w) - 2 z f(z)}{w z^2 - 2 z + w}, \quad (3.31)$$

where  $\mathcal{P}_A(w) = \sum_{m=0}^{\infty} w^m p_A(m)$  and  $f(z) = \sum_{A=1}^n z^A p_A(0)$ . The function  $f(z)$  is the generating function for the initial state at  $m = 0$ . Taking the initial state to be  $B$  with probability unity, we have  $f(z) = z^B$ .  $F(w, z)$  involves the two unknown functions  $\mathcal{P}_1(w)$  and  $\mathcal{P}_n(w)$ , which arise from the reflecting boundary conditions. However, these functions can be uniquely determined via analyticity arguments in the complex  $z$  plane (Cox & Miller, 1965). Specifically,  $F(w, z)$  is by construction a polynomial of degree  $n$  in  $z$  and so cannot contain any singularities in  $z$ . Yet the denominator in Eq. (3.31) has zeros at the two roots, call them  $z_{\pm}(w)$ , of the equation  $z^2 - (2/w)z + 1 = 0$ . The numerator must therefore also have zeros at these two locations. Hence, we may deduce that

$$\mathcal{P}_1(w) = \frac{2}{w} \frac{z_+^n (1 - z_+) f(z_-) - z_-^n (1 - z_-) f(z_+)}{(1 - z_+) (1 - z_-) (z_+^n - z_-^n)}, \quad (3.32)$$

$$\mathcal{P}_n(w) = \frac{2}{w} \frac{(1 - z_+) f(z_-) - (1 - z_-) f(z_+)}{(1 - z_+) (1 - z_-) (z_+^n - z_-^n)}, \quad (3.33)$$

although we do not, in fact, need to know  $\mathcal{P}_n(w)$ .

By writing  $z^2 - (2/w)z + 1 = (1 - z_- z)(1 - z_+ z)$  in the denominator of  $F(w, z)$ , it is straightforward to obtain an expression for the coefficient of

the  $z^A$ ,  $1 \leq A \leq n$ , term in  $F(w, z)$ .<sup>4</sup> With  $f(z) = z^B$ , this coefficient is just  $\sum_{m=0}^{\infty} w^m [(\frac{1}{2}\mathbb{C})^m]_{AB}$ . Because we want to compute  $\mathbf{S}^T \widehat{\mathbb{W}}^{-1} \mathbf{v}$  with  $\mathbf{v}^T = (-1, 0, \dots, 0, +1)$ , it suffices to set  $B = 1$ . Setting  $B = n$  must by symmetry produce the same final result, up to a sign. The required coefficient is then

$$\sum_{m=0}^{\infty} w^m [(\frac{1}{2}\mathbb{C})^m]_{A,1} = \mathcal{P}_1(w) \frac{z_+^A - z_-^A}{z_+ - z_-} - \left[ \mathcal{P}_1(w) + \frac{2}{w} \right] \frac{z_+^{A-1} - z_-^{A-1}}{z_+ - z_-}, \quad (3.34)$$

with  $\mathcal{P}_1(w)$  given by Eq. (3.32), and we need

$$\begin{aligned} & \mathbf{S}^T \widehat{\mathbb{W}}^{-1} \mathbf{v} \\ &= -2 \sum_{A=1}^n s_A \sum_{m=0}^{\infty} (2\widehat{G}_0)^m [(\frac{1}{2}\mathbb{C})^m]_{A,1} \\ &= 2 \sum_{A=1}^n s_A \left\{ \left[ \mathcal{P}_1(2\widehat{G}_0) + (\widehat{G}_0)^{-1} \right] \frac{[z_+(2\widehat{G}_0)]^{A-1} - [z_-(2\widehat{G}_0)]^{A-1}}{z_+(2\widehat{G}_0) - z_-(2\widehat{G}_0)} \right. \\ & \quad \left. - \mathcal{P}_1(2\widehat{G}_0) \frac{[z_+(2\widehat{G}_0)]^A - [z_-(2\widehat{G}_0)]^A}{z_+(2\widehat{G}_0) - z_-(2\widehat{G}_0)} \right\}. \quad (3.35) \end{aligned}$$

After a great deal of tedious but routine algebra, we eventually obtain the

---

<sup>4</sup>Writing, say,  $(1 - z_- z)^{-1} = \sum_{i=0}^{\infty} (z_- z)^i$ , appears to remove the denominator and thus undermine the argument concerning its zeros. However, with this rewriting, the radii of convergence of the two resulting series must be considered, and these two radii are precisely  $|z_{\pm}(w)|$ .

remarkably simple form,

$$\mathbf{S}^T \widehat{\mathbb{W}}^{-1} \mathbf{v} = \frac{2}{n-1} \frac{1}{1-2\widehat{G}_0} \left\{ n - \sqrt{\frac{1+2\widehat{G}_0}{1-2\widehat{G}_0}} \frac{[z_+(2\widehat{G}_0)]^n - 1}{[z_+(2\widehat{G}_0)]^n + 1} \right\}. \quad (3.36)$$

We note that  $z_+(2\widehat{G}_0) \equiv \Phi_+^\ominus$  for a definite choice of sign conventions.

### 3.3 Extraction of $\mu(t)$

We now invert the Laplace transform  $\widehat{\mu}(s)$  to obtain  $\mu(t)$  for the two forms of  $\mathbf{S}^T \widehat{\mathbb{W}}^{-1}(s) \mathbf{v}$  derived above. The two expressions for  $\mathbf{S}^T \widehat{\mathbb{W}}^{-1} \mathbf{v}$  in Eqs. (3.29) and (3.36) are very different in structure, but each is useful for extracting different approximations.

#### 3.3.1 $\mu(t)$ from Eqs. (3.15) and (3.29)

Inserting  $\widehat{G}_0(s) = [\Phi_+(s)]^\ominus / \{[\Phi_+(s)]^{2\ominus} + 1\}$  into Eq. (3.29) and expanding in  $1/s$ , we find that the leading order behaviour is  $\mathbf{S}^T \widehat{\mathbb{W}}^{-1}(s) \mathbf{v} \sim 2 + \mathcal{O}(1/s)$ . For  $n = 2$ ,  $\mathbf{S}^T \widehat{\mathbb{W}}^{-1}(s) \mathbf{v} \equiv 2$ . It is therefore convenient to define

$$\begin{aligned} \widehat{\chi}_n(s) &= \frac{1}{2} \mathbf{S}^T \widehat{\mathbb{W}}^{-1}(s) \mathbf{v} - 1 \\ &= \frac{1}{n(n-1)} \sum_{m=1}^{n-1} [1 - (-1)^m] \frac{2\widehat{G}_0 \cos \frac{m\pi}{n}}{1 - 2\widehat{G}_0 \cos \frac{m\pi}{n}} \cot^2 \frac{m\pi}{2n}, \end{aligned} \quad (3.37)$$

so that the Dirac delta function  $\delta(t)$  present in  $\mathbf{S}^T \mathbb{W}^{-1}(t) \mathbf{v}$  is isolated and made explicit. By the convolution theorem, we may directly invert Eq. (3.17),

writing  $\mu(t)$  in the form

$$\mu_s(t) := \frac{n}{2} \mu(t) = \mu_2(t) + (\mu_2 * \chi_n)(t), \quad (3.38)$$

which defines the scaled form of  $\mu(t)$  to be  $\mu_s(t)$ , removing an overall factor of  $2/n$ , and where  $(f * g)(t)$  denotes the Laplace convolution,

$$(f * g)(t) = \int_0^t d\tau f(t - \tau)g(\tau). \quad (3.39)$$

Using  $\mu_s(t)$  allows us to directly compare  $\mu(t)$  for different values of  $n$ . The form of Eq. (3.38) allows us to see that  $\mu_s(t)$  deviates from and indeed exceeds  $\mu_2(t)$  by a lagged response that is determined by  $\chi_n(t)$  via the convolution  $(\mu_2 * \chi_n)(t)$ . We of course have that  $\chi_2(t) \equiv 0$ , so that  $\mu_s(t) \equiv \mu_2(t)$  for  $n = 2$ . The response is lagged because  $\chi_n(t)$  rises from zero at  $t = 0$  s. To obtain an explicit formula for  $\mu_s(t)$ , we first reproduce here the result for  $\mu_2(t)$  derived elsewhere (Elliott & Lagogiannis, 2012):

$$\begin{aligned} \mu_2(t) = & \frac{1}{\Theta^3} \sum_{l=0}^{\Theta-1} \cot^2 \frac{(2l+1)\pi}{4\Theta} \exp \left\{ -t \left[ 1 - \cos \frac{(2l+1)\pi}{2\Theta} \right] \right\} \\ & - \frac{4}{\Theta^3} \sum_{l=0}^{\lfloor \frac{\Theta-1}{2} \rfloor} \cot^2 \frac{(2l+1)\pi}{2\Theta} \exp \left\{ -t \left[ 1 - \cos \frac{(2l+1)\pi}{\Theta} \right] \right\}, \quad (3.40) \end{aligned}$$

where  $[x]$  is the floor function. We must also determine  $\chi_n(t)$  explicitly from Eq. (3.37) by computing the inverse Laplace transform. We have that

$$\frac{\widehat{G}_0(s)}{1 - 2\widehat{G}_0(s)\cos\frac{m\pi}{n}} = \frac{[\Phi_+(s)]^\Theta}{[\Phi_+(s)]^{2\Theta} - 2\cos\frac{m\pi}{n}[\Phi_+(s)]^\Theta + 1}, \quad (3.41)$$

so there are poles in  $s$  at the roots determined by the solutions of  $[\Phi_+(s)]^\Theta = \exp(\pm im\pi/n)$ . Because  $\Phi_+(s)\Phi_-(s) = 1$  and  $\Phi_+(s) + \Phi_-(s) = 2(1+s)$ , if  $\Phi_+(s) = \omega$  and  $\Phi_+(s) = \omega^*$  are a root and its complex conjugate in  $\Phi_+$ , then

$$\frac{\Phi_+(s)}{[\Phi_+(s)]^2 - (\omega + \omega^*)\Phi_+(s) + 1} = \frac{1}{2} \frac{1}{s + 1 - \frac{1}{2}(\omega + \omega^*)}, \quad (3.42)$$

so a root and its complex conjugate in  $\Phi_+$  combine to create a simple pole in  $s$  at  $s = \frac{1}{2}(\omega + \omega^*) - 1$ . With this observation, the inverse Laplace transform is routine, and we obtain

$$\mathcal{L}^{-1}\left[\frac{\widehat{G}_0(s)}{1 - 2\widehat{G}_0(s)\cos\frac{m\pi}{n}}; t\right] = \frac{1}{2\Theta} \sum_{l=0}^{\Theta-1} \frac{\sin\frac{(2l+m/n)\pi}{\Theta}}{\sin\frac{m\pi}{n}} \exp\left\{-t\left[1 - \cos\frac{(2l+m/n)\pi}{\Theta}\right]\right\}, \quad (3.43)$$

whence,

$$\begin{aligned} \chi_n(t) &= \frac{1}{\Theta n(n-1)} \sum_{m=1}^{n-1} [1 - (-1)^m] \cot^2\frac{m\pi}{2n} \cot\frac{m\pi}{n} \\ &\quad \times \sum_{l=0}^{\Theta-1} \sin\frac{(2l+m/n)\pi}{\Theta} \exp\left\{-t\left[1 - \cos\frac{(2l+m/n)\pi}{\Theta}\right]\right\}. \end{aligned} \quad (3.44)$$

Plugging Eqs. (3.40) and (3.44) into Eq. (3.38) and explicitly computing the convolution, we obtain a messy expression for  $\mu_s(t)$ .

We do not reproduce in full this expression here because below we will obtain a much simpler expression for  $\mu_s(t)$  via Eq. (3.36). However, the advantage of the convolution representation is that we may obtain an approximation to  $\mu_s(t)$  that is valid both at small times and at large times, but not at intermediate times. Because  $\mu_s(t)$  is at small times very close to  $\mu_2(t)$  due to the lagged response from  $\chi_n(t)$ , any good enough approximation to  $\chi_n(t)$  will maintain the behaviour  $\mu_s(t) \approx \mu_2(t)$  for small  $t$ . Furthermore, if we approximate  $\chi_n(t)$  with a form that improves and becomes asymptotically exact at large  $t$ , then  $\mu_s(t)$  under this approximation will become asymptotically exact. Replacing  $\chi_n(t)$  by just its slowest decaying mode [ $m = 1$  and  $l = 0$  in Eq. (3.44)] will achieve this approximation, so writing

$$\chi_n(t) \approx \frac{2}{\Theta n(n-1)} \cot^2 \frac{\pi}{2n} \cot \frac{\pi}{n} \sin \frac{\pi}{\Theta n} \exp \left[ -t \left( 1 - \cos \frac{\pi}{\Theta n} \right) \right]. \quad (3.45)$$

A different approximation can be achieved by retaining the full sum over  $l$  in Eq. (3.44) and retaining only the  $m = 1$  contribution from the sum over  $m$ , but this approximation is necessarily more complicated than retaining just the slowest decaying,  $m = 1$  and  $l = 0$  mode in Eq. (3.44).

In Fig. 2 we illustrate the result for  $\chi_n(t)$  in Eq. (3.44) by plotting  $\chi_n(t)$  as a function of  $t$  for various choices of  $n$  and  $\Theta$ . In Fig. 2A, for varying  $n$  and fixed  $\Theta$ , we see explicitly the lagged onset of  $\chi_n(t)$ , rising from zero, reaching



a peak, and then falling back to zero. As  $n$  increases, the peak increases in amplitude, is displaced somewhat later in time, and also broadens, becoming plateau-like. In Fig. 2B we fix  $n$  and vary  $\Theta$ . As  $\Theta$  increases,  $\chi_n(t)$  is overall and significantly (note the log scale) scaled down, and its peak is increasingly displaced towards later times. Comparing  $\chi_n(t)$  to the probability density  $2G_0(t)$  for escape through either filter threshold starting from the zero filter state, plotted in Fig. 2C, we see that  $\chi_n(t)$  tracks  $2G_0(t)$  quite closely for smaller times, with the overall scale of  $\chi_n(t)$  being set by  $2G_0(t)$ . In fact, if we expand the form for  $\widehat{\chi}_n(s)$  in Eq. (3.36) in terms of  $2\widehat{G}_0(s)$ , we find that the leading order behaviour is governed by

$$\widehat{\chi}_n(s) \sim \frac{n-2}{n-1} \times 2\widehat{G}_0(s), \quad (3.46)$$

so that  $\chi_n(t) \sim \frac{n-2}{n-1} \times 2G_0(t)$  for small times. The escape density  $2G_0(t)$  therefore explicitly sets the scale for  $\chi_n(t)$ , up to the overall factor of  $\frac{n-2}{n-1}$ . As  $n$  increases, this factor approaches unity, and indeed we see all the curves for  $\chi_n(t)$  for increasing  $n$  in Fig. 2A converging for smaller times.

### 3.3.2 $\mu(t)$ from Eqs. (3.15) and (3.36)

We now obtain  $\mu(t)$  using the form for  $\mathbf{S}^T \widehat{\mathbb{W}}^{-1}(s) \mathbf{v}$  in Eq. (3.36). Rather than exploiting the convolution structure explicit in the product of two Laplace transforms, we instead reduce Eq. (3.15) to its simplest form before evaluating the inverse transform. We know that  $\widehat{\mu} = \frac{1}{n} \widehat{\mu}_2 \mathbf{S}^T \widehat{\mathbb{W}}^{-1} \mathbf{v}$  and from Eq. (A.50)

in Elliott & Lagogiannis (2012), we have that

$$\widehat{\mu}_2 = \frac{1}{\Theta^2} \frac{2\Phi_+(\Phi_++1)}{(\Phi_+^{2\Theta}+1)(\Phi_+^\Theta+1)} \left( \frac{\Phi_+^\Theta-1}{\Phi_+-1} \right)^3. \quad (3.47)$$

Writing  $\mathbf{S}^\top \widehat{\mathbb{W}}^{-1}(s)\mathbf{v}$  in Eq. (3.36) out in terms of  $\Phi_+(s)$ , we obtain

$$\mathbf{S}^\top \widehat{\mathbb{W}}^{-1}(s)\mathbf{v} = \frac{2n}{n-1} \frac{\Phi_+^{2\Theta}+1}{(\Phi_+^\Theta-1)^2} - \frac{2}{n-1} \frac{(\Phi_+^{2\Theta}+1)(\Phi_+^\Theta+1)}{(\Phi_+^\Theta-1)^3} \frac{\Phi_+^{\Theta n}-1}{\Phi_+^{\Theta n}+1}, \quad (3.48)$$

so we have

$$\widehat{\mu}_s = \frac{1}{\Theta^2} \frac{1}{n-1} \frac{2\Phi_+(\Phi_++1)}{(\Phi_+-1)^3} \left[ n \frac{\Phi_+^\Theta-1}{\Phi_+^\Theta+1} - \frac{\Phi_+^{\Theta n}-1}{\Phi_+^{\Theta n}+1} \right]. \quad (3.49)$$

The inverse Laplace transform is routine, and we obtain

$$\begin{aligned} \mu_s(t) = \frac{2}{\Theta^3(n-1)} & \left( \frac{1}{n} \sum_{l=0}^{\lfloor \frac{\Theta n-1}{2} \rfloor} \cot^2 \frac{(2l+1)\pi}{2\Theta n} \exp \left\{ -t \left[ 1 - \cos \frac{(2l+1)\pi}{\Theta n} \right] \right\} \right. \\ & \left. - n \sum_{l=0}^{\lfloor \frac{\Theta-1}{2} \rfloor} \cot^2 \frac{(2l+1)\pi}{2\Theta} \exp \left\{ -t \left[ 1 - \cos \frac{(2l+1)\pi}{\Theta} \right] \right\} \right). \end{aligned} \quad (3.50)$$

This general  $n$  form is striking in its similarity to the  $n = 2$  form in Eq. (3.40).

For  $n = 2$ ,  $\mu_s(t)$  reduces identically to  $\mu_2(t)$  because  $\lfloor \Theta - \frac{1}{2} \rfloor \equiv \Theta - 1$ . By again taking the slowest decaying mode from each of the two sums in Eq. (3.50), we

obtain an extremely simple approximation to  $\mu_s(t)$ ,

$$\mu_s(t) \approx \frac{2}{\Theta^3(n-1)} \left\{ \frac{1}{n} \cot^2 \frac{\pi}{2\Theta n} \exp \left[ -t \left( 1 - \cos \frac{\pi}{\Theta n} \right) \right] - n \cot^2 \frac{\pi}{2\Theta} \exp \left[ -t \left( 1 - \cos \frac{\pi}{\Theta} \right) \right] \right\}, \quad (3.51)$$

or, taking only the slowest decaying term, the even simpler

$$\mu_s(t) \approx \frac{2}{\Theta^3 n(n-1)} \cot^2 \frac{\pi}{2\Theta n} \exp \left[ -t \left( 1 - \cos \frac{\pi}{\Theta n} \right) \right]. \quad (3.52)$$

This latter form is especially useful for determining SNR memory lifetimes because it gives an extremely simple expression for  $\mu_s(t)$  at large times.

Finally, we may extract some large  $n$  limits by considering the limiting behaviour of Eq. (3.36). Throwing away terms that are exponentially suppressed in  $n$ , we have

$$\mathbf{S}^T \widehat{\mathbb{W}}^{-1} \mathbf{v} \sim \frac{2}{n-1} \frac{1}{1-2\widehat{G}_0} \left[ n - \sqrt{\frac{1+2\widehat{G}_0}{1-2\widehat{G}_0}} \right] \quad (3.53)$$

as an  $\mathcal{O}(1/n)$  approximation, and as an  $\mathcal{O}(1)$  approximation we have

$$\mathbf{S}^T \widehat{\mathbb{W}}^{-1} \mathbf{v} \sim \frac{2}{1-2\widehat{G}_0}. \quad (3.54)$$

The two corresponding forms of  $\mu_s(t)$  are then

$$\begin{aligned} \mu_s(t) \sim & \frac{n}{n-1} \left( \frac{1}{\Theta} - \frac{2}{\Theta^3} \sum_{l=0}^{\lfloor \frac{\Theta-1}{2} \rfloor} \cot^2 \frac{(2l+1)\pi}{2\Theta} \exp \left\{ -t \left[ 1 - \cos \frac{(2l+1)\pi}{2\Theta} \right] \right\} \right) \\ & - \frac{1}{n-1} \frac{1}{\Theta^2} [(1+2t)I_0(t) + 2tI_1(t)] \exp(-t), \end{aligned} \quad (3.55)$$

and

$$\mu_s(t) \sim \frac{1}{\Theta} - \frac{2}{\Theta^3} \sum_{l=0}^{\lfloor \frac{\Theta-1}{2} \rfloor} \cot^2 \frac{(2l+1)\pi}{2\Theta} \exp \left\{ -t \left[ 1 - \cos \frac{(2l+1)\pi}{2\Theta} \right] \right\} \quad (3.56)$$

$$\rightarrow \frac{1}{\Theta} \text{ as } t \rightarrow \infty, \quad (3.57)$$

where  $I_0(t)$  and  $I_1(t)$  are modified Bessel functions of the first kind. Note the remarkable behaviour that  $\mu_s(t)$  asymptotes to a constant,  $1/\Theta$ , for large  $t$  in the formal,  $n \rightarrow \infty$  limit. We may compare this  $\mathcal{O}(1)$  approximation to the large  $n$  form of the two-decay approximation in Eq. (3.51),

$$\mu_s(t) \sim \frac{8}{\pi^2} \frac{1}{\Theta} - \frac{2}{\Theta^3} \cot^2 \frac{\pi}{2\Theta} \exp \left[ -t \left( 1 - \cos \frac{\pi}{\Theta} \right) \right]. \quad (3.58)$$

This large  $n$  form of the two-decay approximation therefore underestimates the asymptotic behaviour  $\mu_s(t) \sim 1/\Theta$  by around 19%, since  $8/\pi^2 \approx 0.81$ .

### 3.4 Results for a Stochastic Updater Synapse

To facilitate comparison, we also consider a simple, stochastic updater synapse (Tsodyks, 1990). Such a synapse expresses potentiation or depression steps

with a fixed probability  $p$  upon receipt of a plasticity induction signal. Previously, we compared such a synapse to a filter-based synapse, but only for  $n = 2$  (Elliott & Lagogiannis, 2012).

The transition matrix for single, one-step changes in synaptic strength for a stochastic updater synapse is simply  $\mathbb{J} = (1 - p)\mathbb{I} + \frac{1}{2}p\mathbb{C}$ , so the transition matrix for changes in synaptic strength in time  $t$  is

$$\begin{aligned}\mathbb{P}(t) &= \exp[t(\mathbb{J} - \mathbb{I})] \\ &= \sum_{m=0}^{n-1} \mathbf{e}^m (\mathbf{e}^m)^\top \exp[-pt(1 - \cos \frac{m\pi}{n})],\end{aligned}\quad (3.59)$$

where the second line follows immediately from the eigenstructure of  $\frac{1}{2}\mathbb{C}$  considered above in section 3.2.1. The mean memory signal is as usual  $\mu(t) = \mathbf{S}^\top \mathbb{P}(t) \mathbf{a}$ , where  $\mathbf{a}$  is the probability distribution of strength states immediately after the storage of memory  $\boldsymbol{\xi}^0$ . This distribution is just [cf. Eq. (3.22)]

$$\mathbf{a} = \frac{1}{n} (\mathbf{n} + p\mathbf{v}). \quad (3.60)$$

We then obtain

$$\mu(t) = \frac{2p}{n^2(n-1)} \sum_{m=1}^{n-1} [1 - (-1)^m] \cot^2 \frac{m\pi}{2n} \exp[-pt(1 - \cos \frac{m\pi}{n})]. \quad (3.61)$$

For  $n = 2$ ,  $\mu(t)$  reduces to  $\mu(t) = p \exp(-pt)$ , as it should (Elliott & Lagogiannis, 2012).

Because a stochastic updater synapse has no internal states, the argument

leading to Eq. (3.23) applies with the replacements  $1/\Theta^2 \rightarrow p$  and  $G_0(t) \rightarrow \frac{1}{2} p \exp(-pt)$ , where  $G_0(t)$  is just the density for the expression of either a potentiation or a depression step. We may then use the large  $n$  approximations to  $\mathbf{S}^T \widehat{\mathbb{W}}^{-1} \mathbf{v}$  in Eqs. (3.53) and (3.54) to obtain large  $n$  forms for  $\mu(t)$  for a stochastic updater. We obtain

$$\mu_s(t) \sim \frac{n}{n-1} p - \frac{1}{n-1} p [(1 + 2pt)I_0(pt) + 2pt I_1(pt)] \exp(-pt), \quad (3.62)$$

and

$$\mu_s(t) \sim p, \quad (3.63)$$

respectively. We note that  $\mu_s(t) \sim \mu_s(0)$ , in the formal limit,  $n \rightarrow \infty$ , so that the scaled mean memory signal remains fixed for all time in this limit.

For a stochastic updater, we may also readily compute  $\mathbf{E}[\tilde{S}_1(t)\tilde{S}_2(t)]$  and hence  $\text{Cov}(t)$ . To determine  $\mathbf{E}[\tilde{S}_1(t)\tilde{S}_2(t)]$  we must work in the tensor product space defined by any pair of synapses. The transition matrix over time  $t$  for joint changes in strength between any pair of synapses is  $\exp[t(\mathbb{J} \otimes \mathbb{J} - \mathbb{I} \otimes \mathbb{I})]$ , so

$$\mathbf{E}[\tilde{S}_1(t)\tilde{S}_2(t)] = \mathbf{S}^T \otimes \mathbf{S}^T \exp[t(\mathbb{J} \otimes \mathbb{J} - \mathbb{I} \otimes \mathbb{I})] \mathbf{a} \otimes \mathbf{a}. \quad (3.64)$$

Letting the eigenvalues of  $\mathbb{J}$  be  $\Lambda_m$ , where  $\Lambda_m = 1 - p(1 - \cos \frac{m\pi}{n})$ , for  $m = 0, \dots, n-1$ , then  $\mathbf{e}^{m_1} \otimes \mathbf{e}^{m_2}$  is an eigenvector of  $\mathbb{J} \otimes \mathbb{J} - \mathbb{I} \otimes \mathbb{I}$  with eigenvalue

$\Lambda_{m_1}\Lambda_{m_2} - 1$ . We then have

$$\begin{aligned} \mathbf{E}[\tilde{S}_1(t)\tilde{S}_2(t)] &= \frac{4p^2}{n^4(n-1)^2} \sum_{m_1, m_2=1}^{n-1} [1 - (-1)^{m_1}] [1 - (-1)^{m_2}] \cot^2 \frac{m_1\pi}{2n} \cot^2 \frac{m_2\pi}{2n} \\ &\quad \times \exp [t (\Lambda_{m_1}\Lambda_{m_2} - 1)]. \end{aligned} \quad (3.65)$$

For  $n = 2$ , this reduces to  $\mathbf{E}[\tilde{S}_1(t)\tilde{S}_2(t)] = p^2 \exp [-tp(2-p)]$ , as before (Elliott & Lagogiannis, 2012). Eq. (3.65) does not factorise over the two sums and so reduce to  $\mu(t)^2$  because in general  $(\Lambda_{m_1}\Lambda_{m_2} - 1) \neq (\Lambda_{m_1} - 1) + (\Lambda_{m_2} - 1)$ .

## 4 Memory Performance

Having derived exact and approximate results for  $\mu(t)$  above, we may now employ these results to examine the memory signal dynamics and their impact on memory lifetimes.

### 4.1 Memory Signal Dynamics

We first consider the dynamics of the mean memory signal  $\mu(t)$  and its scaled version,  $\mu_s(t)$ . In Fig. 3 we verify our analytical results by comparing  $\mu(t)$  to simulation results for various choices of parameters. Full details of our simulation protocols may be found elsewhere (Elliott & Lagogiannis, 2012; Elliott, 2014). Figs. 3A and 3B compare results for  $\Theta = 4$  and  $\Theta = 8$ , respectively, each for four different choices of  $n$ , as indicated. The agreement between analytical

and simulation results is exact.

From Figs. 3A and 3B we see that  $\mu(t)$  for  $n > 2$  continues to exhibit the very striking initial increase in the memory signal that we first reported for binary,  $n = 2$  synapses in the presence of a synaptic filter that integrates synaptic plasticity induction signals before expressing synaptic plasticity (Elliott & Lagogiannis, 2012). Although  $\mu(t)$  is progressively scaled down as  $n$  increases, this signal rise is sustained somewhat longer as  $n$  increases, and then its subsequent fall back to zero is attenuated and drastically slowed down so that, at least for the parameters used in these figures,  $\mu(t)$  for larger  $n$  will eventually exceed  $\mu(t)$  for smaller  $n$ , despite this overall scaling down. By using the scaled mean memory signal  $\mu_s(t) = \frac{n}{2}\mu(t)$ , we may more easily visually compare the memory signal dynamics for different values of  $n$ . Thus, in Figs. 3C and 3D, the shift in the peak signal location as  $n$  increases is clearer. We also see that in terms of the scaled signal, increasing  $n$  increases not only the time at which the signal reaches its peak value but also the absolute value of the peak signal. For small times,  $\mu_s(t)$  is essentially independent of  $n$ , with  $\mu_s(t) \approx \mu_2(t)$ . The lagged influence of larger values of  $n$  then becomes apparent as  $t$  increases, with  $\mu_s(t)$  taking longer to reach higher signal peaks, and then taking longer to relax back to equilibrium as  $t \rightarrow \infty$ . These dynamics for  $\mu(t)$  or  $\mu_s(t)$  reflect the dynamics for  $\chi_n(t)$  discussed in Fig. 2.

For  $n = 2$  synaptic strength states, the signal rise occurs because the storage of the initial memory  $\xi^0$  at  $t = 0^-$  s induces a systematic bias in the distribution of synaptic filter states (Elliott & Lagogiannis, 2012). Those synapses



experiencing an initial potentiating induction signal become more likely subsequently to remain strong or to potentiate if not strong, and correspondingly those synapses experiencing an initial depressing induction signal become subsequently biased to remain weak or become weak. The signal rise for  $n = 2$  occurs while this bias in filter states exists: the biased filter states induce progressively greater asymmetry in the distribution of synaptic strength states, which translates directly into the rising memory signal. Once the filter states have re-equilibrated and lost their bias, the memory signal stops rising and then starts to fall as the distribution of synaptic strength states re-equilibrates (Elliott & Lagogiannis, 2012). For  $n > 2$ , essentially identical dynamics occur. However, for  $n > 2$  synaptic strength states, the re-equilibration of filter states takes somewhat longer and the subsequent re-equilibration of synaptic strength states then takes considerably longer. We will quantify these statements using the approximate forms for  $\mu(t)$  derived above.

Before this quantification, we compare in Fig. 4 the various approximations to  $\mu_s(t)$  derived above to its exact form. In Fig. 4A, we examine the convolution form of  $\mu_s(t)$  in Eq. (3.38) together with the approximation to  $\chi_n(t)$  in Eq. (3.45). By construction, this approximate form for  $\mu_s(t)$  agrees with the exact form for both small times and large times, with discrepancies only at intermediate times. Because this approximation still involves a convolution over  $\mu_2(t)$ , using it to obtain the location of the signal peak would require yet further approximations. In Figs. 4B and 4C we consider the one- and two-decay forms of  $\mu_s(t)$  in Eq. (3.52) and (3.51), respectively. Both forms necessarily

agree with  $\mu_s(t)$  for large times, while the two-decay form at least captures the key dynamics of the initial signal rise and its subsequent fall, even if the approximation somewhat underestimates the amplitude of the signal peak. The one-decay form is particularly simple, and because it matches the exact form for large times, it is especially useful for determining memory lifetimes. The two-decay form is useful for estimating the location and amplitude of the signal peak. Finally, in Fig. 4D we examine the large  $n$ ,  $\mathcal{O}(1/n)$  form for  $\mu_s(t)$  in Eq. (3.55). The agreement with the exact result for all but large times, even for the relatively small value of  $n = 10$  used here, is remarkably good. Comparing Eq. (3.55) with the exact form for  $\mu_s(t)$  in Eq. (3.50), we see that the approximate form retains the second sum in Eq. (3.50), which does not depend on  $n$  (except through overall multipliers) and essentially replaces the first,  $n$ -dependent sum by two Bessel functions. As the second sum is also present in  $\mu_2(t)$  and since  $\mu_s(t) \approx \mu_2(t)$  for small  $t$ , the first sum in Eq. (3.50) must essentially encode the  $n$ -dependent signal dynamics that lead to changes in the location and amplitude of the signal peak and its subsequent, slower decay. Since the replacement of the first sum by the two Bessel functions constitutes an  $\mathcal{O}(1/n)$  approximation that becomes increasingly good as  $n$  increases, the striking agreement in Fig. 4D should therefore perhaps not be too surprising. From a purely numerical perspective, then, this large  $n$  approximation provides an extremely efficient approximation to Eq. (3.50) as we do not need to compute the first sum with its nearly  $n \Theta/2$  terms. However, because it does not capture the large time dynamics correctly, this approximation is not useful

for estimating memory lifetimes.

The two-decay approximation to  $\mu_s(t)$  in Eq. (3.51) provides, as we have seen, an analytically very simple form that captures the key dynamics of both the mean memory signal rise and its subsequent fall. Using this two-decay form, we can estimate the location of the signal peak as

$$rt_{\text{peak}} \approx \frac{2}{\cos \frac{\pi}{\Theta n} - \cos \frac{\pi}{\Theta}} \log_e \left( n \frac{\cos \frac{\pi}{2\Theta}}{\cos \frac{\pi}{2\Theta n}} \right) \quad (4.1)$$

$$\sim \frac{4\Theta^2}{\pi^2} \frac{n^2}{n^2 - 1} \log_e n \sim \frac{4\Theta^2}{\pi^2} \log_e n \quad (4.2)$$

where the second and third forms follow for  $n$  and  $\Theta$  large enough. The location of the signal peak therefore grows relatively mildly as a function of  $n$ , growing only logarithmically. An estimate of the peak scaled memory signal is therefore

$$\mu_s^{\text{peak}} \approx \frac{8}{\pi^2} \frac{1}{\Theta} \left( 1 + \frac{1}{n} \right) n^{-2/(n^2-1)}, \quad (4.3)$$

for  $n$  and  $\Theta$  large enough. We note that as a function of  $n$ , this estimate for  $\mu_s^{\text{peak}}$  first rises and then falls, asymptoting to  $(8/\pi^2)(1/\Theta)$ . Fig. 3, however, indicates that  $\mu_s^{\text{peak}}$  actually increases monotonically as a function of  $n$ , asymptoting to  $1/\Theta$ . This failure of the estimate for  $\mu_s^{\text{peak}}$  to capture the exact, quantitative behaviour of  $\mu_s^{\text{peak}}$  is not an artifact of the large  $n$  approximations used to derive it, but rather is a direct consequence of the two-decay approximation to  $\mu_s(t)$ . Specifically, if we insert the exact location of the peak of the two-decay form of  $\mu_s(t)$  in Eq. (4.1) into this two-decay form for  $\mu_s(t)$ , then

it still exhibits the same overall dependence on  $n$  as the form in Eq. (4.3). In order to obtain better approximations to  $\mu_s^{\text{peak}}$ , we would have to use better approximations to  $\mu_s(t)$ , but then analytical approximations for  $\mu_s^{\text{peak}}$  would become intractable.

Although the two-decay estimate of  $\mu_s^{\text{peak}}$  is only fair, we may nevertheless use it to determine the qualitative dependence of the fall in  $\mu_s(t)$  after attaining its peak value on  $n$ . In particular, we may compute the time at which  $\mu_s(t)$  or equivalently  $\mu(t)$  drops to some fraction,  $1 - \delta$  for  $\delta$  small, of its peak value. If this time is large enough so that we may use the one-decay form for  $\mu_s(t)$ , then we can use Eq. (3.52) to solve the equation  $\mu_s(t) = (1 - \delta)\mu_s^{\text{peak}}$  for  $t = t_{(1-\delta)\text{peak}} > t_{\text{peak}}$ . We obtain

$$rt_{(1-\delta)\text{peak}} \approx \frac{2\Theta^2 n^2}{\pi^2} \delta, \quad (4.4)$$

again for  $n$  and  $\Theta$  large enough, and for  $\delta$  small enough. This time at which the mean memory signal has fallen from its peak value by some fraction  $\delta \ll 1$  therefore grows quadratically in  $n$ .

The results in Eqs. (4.2) and (4.4) quantify our earlier statement above regarding the re-equilibration of filter and strength states after the storage of memory  $\xi^0$  at time  $t = 0^-$ . Filter states re-equilibrate only marginally (logarithmically) more slowly as  $n$  increases, but strength states re-equilibrate significantly (quadratically) more slowly as  $n$  increases. Of course, the number of filter states,  $2\Theta - 1$ , does not by construction depend on the number of

strength states  $n$ , so the mild dependence of the re-equilibration of filter states on  $n$  may appear surprising. However, this dependence simply reflects the fact that it takes longer for the systematic biasing in filter states over all strength states [see Eq. (2.11)] to work out of the system.

We continue to explore these themes in Fig. 5 by examining  $\mu_s(t)$  as  $n$  increases. In Fig. 5A, we show  $\mu_s(t)$  for progressively doubling values of  $n$ , from  $n = 2^1$  to  $n = 2^{12} = 4096$ . Larger values of  $n$  are likely biologically implausible but we consider them for illustrative purposes. We see clearly the progressively but only mildly increasing signal peak location as  $n$  increases. We also see the manner in which  $\mu_s(t)$  attains the large  $n$  limit. The scaled memory signal rises to very close to  $1/\Theta$ , which is the formal, asymptotic limit in Eq. (3.57), remains very close to this value for increasingly long periods of time as  $n$  increases, and then eventually peels away from it to relax back to its equilibrium value. Fig. 5B indicates that precisely the same dynamics are present in the two-decay approximation to  $\mu_s(s)$  in Eq. (3.51). The two-decay approximation underestimates the asymptotic value of  $1/\Theta$ , instead giving  $(8/\pi^2)(1/\Theta)$ , but otherwise captures identical memory signal dynamics. In Fig. 5C we examine the dependence of the location of the signal peak on  $n$ , both for the exact form for  $\mu_s(t)$  (location determined numerically) and for the two-decay form, for which the location is given by Eq. (4.1). The onset of logarithmic behaviour for both forms as  $n$  increases is clear. We see that the two-decay form systematically overestimates the location of the peak signal, but the qualitative dependence on  $n$  is correct. Fig. 5D shows the time at which

$\mu_s(t)$  has fallen to 95% of its peak value, again both for the exact form of  $\mu_s(t)$  and the two-decay form. These times for both forms of  $\mu_s(t)$  are determined numerically from Eqs. (3.50) and (3.51). We do not use the approximate solution in Eq. (4.4) for the two-decay form, although this approximate solution agrees with the numerically-determined solution for  $n$  large enough. For both forms of  $\mu_s(t)$  we see the onset of linearity in the log-log plot, and in fact the gradient approaches precisely two, indicating quadratic growth in  $n$ , consistent with Eq. (4.4). Again, the two-decay form systematically overestimates the exact location of 95% of signal peak, but is qualitatively correct.

## 4.2 Memory Lifetimes

Having examined in the detail the dynamics of  $\mu(t)$  and  $\mu_s(t)$ , we may now consider the impact of these dynamics on memory lifetimes, gauged principally by the SNR  $\mu(t)/\sigma(t)$ .

As we have discussed elsewhere (Elliott & Lagogiannis, 2012; Elliott, 2014), the presence of covariance between any pair of synapses' strengths induced by driving memory storage as a continuous-time process considerably complicates the computation of memory lifetimes using SNRs. However, if memory lifetimes are sufficiently long that the covariance has died away by the time that the SNR reaches unity, then we may safely approximate the variance  $\sigma(t)^2$  in Eq. (2.10) by dropping the covariance term. We may usually approximate even further by replacing  $\text{Var}[S(t)]$  by  $\mathbf{E}[S(t)^2]$  because  $\mu(t)^2$  is often small enough to neglect. We validate these approximations for the more general,

$n > 2$  case in Fig. 6. First, in Fig. 6A, we compare  $\sigma(t)$  determined exactly via numerical matrix methods to results obtained in simulation. The agreement is exact, although there is inevitably a little more noise in second order statistics obtained via simulation compared to first order statistics. The variance  $\sigma(t)^2$  exhibits a double peak. Both peaks are increasingly suppressed as  $n$  increases, but the second peak is suppressed more significantly than the first peak. This enhanced suppression of the second peak is evident in the covariance between pairs of synapses' strength, shown in Fig. 6B. Such enhanced suppression for increasing  $n$ , which we expect to be associated with increasing memory lifetimes due to the sustained memory signal  $\mu(t)$ , helps to ensure that the covariance term in Eq. (2.10) plays an increasingly insignificant role in the determination of memory lifetimes. In Figs. 6C and 6D we plot the SNR  $\mu(t)/\sigma(t)$  for two sets of parameters leading, respectively, to shorter memory lifetimes and longer memory lifetimes. For both figures, we use the three different forms of  $\sigma(t)^2$  discussed above: the full form in Eq. (2.10); Eq. (2.10) without the covariance term;  $\mathbf{E}[S(t)^2]/N$ . In both figures, for small times there are large differences in these different forms of SNR, specifically between the exact form and its approximate forms. The approximate forms significantly overestimate the SNR because the large covariance significantly reduces the SNR. For large times, however, the covariance dies away. If it dies away more rapidly than the mean memory signal, then its effect becomes negligible and all three forms of SNR coincide. In Fig. 6C, for a shorter memory lifetime, this coinciding of the three forms does not occur before  $\mu(t)/\sigma(t)$  reaches unity. In this case, use of the

approximate forms of the SNR leads to an overestimate of memory lifetime by around 20%. In Fig. 6D, however, for a longer memory lifetime, this coinciding of the three forms does occur before  $\mu(t)/\sigma(t)$  reaches unity, and use of the approximate forms of SNR leads to only a 0.5% error in the computed memory lifetimes. Generally speaking, memory lifetimes are increased by increasing  $\Theta$  or  $n$  (or both), so for  $\Theta$  and  $n$  sufficiently large, we may approximate  $\sigma(t)^2$  either by dropping the covariance term in Eq. (2.10) or even more simply by writing  $\sigma(t)^2 \approx \mathbf{E}[S(t)^2]/N$ . For large times, we may use the one-decay form of  $\mu(t)$  as an approximation to the exact form of  $\mu(t)$ . We may then easily solve the SNR condition  $\mu(\tau_{\text{snr}})/\sigma(\tau_{\text{snr}}) = 1$  for  $\sigma(t)^2 \approx \mathbf{E}[\tilde{S}(t)^2]/N$ , giving

$$r\tau_{\text{snr}} = \frac{-1}{1 - \cos \frac{\pi}{\Theta n}} \log_e \left[ \frac{1}{4} \frac{1}{\sqrt{3N}} \sqrt{\frac{n+1}{n-1}} \Theta^3 n^2 (n-1) \tan^2 \frac{\pi}{2\Theta n} \right] \quad (4.5)$$

$$\sim \frac{\Theta^2 n^2}{\pi^2} \log_e \left( \frac{768}{\pi^4} \frac{N}{\Theta^2 n^2} \right), \quad (4.6)$$

with the second line following for  $n$  large enough. We note in passing that this solution for  $\tau_{\text{snr}}$  is completely unchanged even if we do not scale synaptic strengths by an overall factor of  $n$ . Such scaling of course modifies both  $\mu(t)$  and  $\sigma(t)$  identically, and any such scaling therefore drops out entirely from the ratio  $\mu(t)/\sigma(t)$  used to determine SNR memory lifetimes.

In Fig. 7 we plot the evolution of the SNR in time for various choices of  $n$  (Fig. 7A) and  $N$  (Fig. 7B) using  $\sigma(t)^2 \approx \{\mathbf{E}[S(t)^2] - \mu(t)^2\}/N$ , and the dependence of memory lifetimes  $\tau_{\text{snr}}$  on  $N$  for various choices of  $n$  (Fig. 7C). In Fig. 7A, we see that although  $\mu(t)/\sigma(t)$  is overall scaled down as  $n$  increases,



$\mu(t)/\sigma(t)$  is sustained for longer for increasing  $n$  because of the dynamics in  $\mu(t)$  discussed earlier. Up to a point, therefore, the solution to  $\mu(t)/\sigma(t) = 1$  increases as  $n$  increases, so that memory lifetimes increase with  $n$ . This increase in memory lifetimes as a function of  $n$  is, at least initially, quadratic in  $n$ . For example, we see from Fig. 7A that in moving from  $n = 2$  to  $n = 8$ , the solution to  $\mu(t)/\sigma(t) = 1$  increases around 100-fold, and similarly in moving from  $n = 4$  to  $n = 16$ . However, the overall scaling down of  $\mu(t)$  as  $n$  increases wins out, so that for intermediate values of  $n$ , memory lifetimes stop increasing with  $n$  and start decreasing (e.g., there is hardly any change in the solution to  $\mu(t)/\sigma(t) = 1$  in Fig. 7A in moving from  $n = 64$  to  $n = 128$ ), and eventually  $\mu(t)/\sigma(t)$  never exceeds unity at all. These interactions lead to an optimal choice for  $n$  for maximally enhancing memory lifetimes that we will explore below. As a function of  $N$ , memory lifetimes consistently and uniformly increase as  $N$  increases, as shown in Fig. 7B. However, the enhancement in  $\mu(t)/\sigma(t)$  as a function of  $N$  is only logarithmic in  $N$ , as seen in Eq. (4.5). Increasing  $N$  to increase memory lifetimes is therefore an extremely inefficient method, and a biologically very expensive one. Finally, Fig. 7C explicitly plots memory lifetimes as a function of  $N$  for various choices of  $n$ . We determine  $\tau_{\text{snr}}$  both from the exact form of  $\mu(t)$  using numerical methods (thick lines) and from the solution to the one-decay form in Eq. (4.5) (thin lines). For  $N$  large enough, these two solutions coincide exactly. For small  $N$ , however, the SNR will not exceed unity, so a solution to  $\mu(t)/\sigma(t) = 1$  does not exist. A critical value of  $N$  exists above which such a solution exists, so there is a bifurcation

as  $N$  increases. The one-decay solution in Eq. (4.5), however, always exists, although this solution will be negative (and thus meaningless) for  $N$  too small. We may regard the transition from negative solutions to positive positions as akin to a pseudo-bifurcation. The exact and approximate memory lifetimes therefore typically differ near the (real or pseudo-) bifurcation point, but are qualitatively very similar in the small  $N$  regime.

As discussed, there is a tension between the enhanced longevity of  $\mu(t)$  caused by increasing  $n$  and the overall scaling down of  $\mu(t)$  as  $n$  increases. This tension is clear in Eq. (4.6). Memory lifetimes grow initially quadratically in  $n$  but the logarithmic term eventually wins out, halting and then reversing the growth, so that memory lifetimes start falling at some critical value of  $n$ . This critical value of  $n$  is an optimal value, call it  $n_{\text{opt}}$ , for maximising memory lifetimes. For fixed  $N$  and  $\Theta$ , it is given by

$$n_{\text{opt}} = \sqrt{\frac{768}{\pi^4 e} \frac{\sqrt{N}}{\Theta}} \approx 1.70 \frac{\sqrt{N}}{\Theta}. \quad (4.7)$$

Of course, alternatively we may regard both  $n$  and  $\Theta$  as fixed and obtain an optimal value of  $N$ ,  $N_{\text{opt}}$ , or regard both  $n$  and  $N$  as fixed and obtain an optimal value of  $\Theta$ ,  $\Theta_{\text{opt}}$ . Inserting  $n_{\text{opt}}$  into Eq. (4.6), we obtain the maximum possible value of  $\tau_{\text{snr}}$ ,

$$r\tau_{\text{snr}}^{\text{max}} = \frac{768}{\pi^6 e} N \approx 0.29N. \quad (4.8)$$

Notice that an optimal choice of either  $n$  or  $\Theta$  has therefore replaced the feeble logarithmic growth of  $\tau_{\text{snr}}$  in  $N$  with linear growth in  $N$ . Alternatively, we may

optimally choose  $N$  and instead obtain

$$r\tau_{\text{snr}}^{\text{max}} = \frac{\Theta^2 n^2}{\pi^2} \approx 0.10 \Theta^2 n^2. \quad (4.9)$$

In this case, the attenuated and reversed growth of  $\tau_{\text{snr}}$  for large  $n$  has been removed, so that for optimally chosen  $N$ ,  $\tau_{\text{snr}}$  grows quadratically in  $n$  (or  $\Theta$ ) with no attenuation. For the specific case of  $n = 2$ , binary synapses, for  $\Theta$  large enough, Eq. (4.6) becomes

$$r\tau_{\text{snr}} \sim \frac{4\Theta^2}{\pi^2} \log_e \left( \frac{256}{\pi^4} \frac{N}{\Theta^2} \right), \quad (4.10)$$

and for optimally chosen  $\Theta$ ,

$$\Theta_{\text{opt}} = \sqrt{\frac{256}{\pi^4 e}} \sqrt{N} \approx 0.98 \sqrt{N}, \quad (4.11)$$

we then obtain,

$$r\tau_{\text{snr}}^{\text{max}} = \frac{1024}{\pi^6 e} N \approx 0.39N. \quad (4.12)$$

Comparing Eqs. (4.8) and (4.12), at first blush it may appear somewhat perverse that we obtain a 33% higher upper limit to memory lifetimes for  $n = 2$ , binary synapses than for  $n > 2$  synapses. However, these upper limits are realised only for optimal parameter choices. Comparing Eqs. (4.7) and (4.11), we see that the numerical values for optimal parameter choices may differ dramatically. For  $N = 10^4$  synapse, Eq. (4.11) tells us that for binary synapses,

maximum memory lifetimes are achieved for the rather unrealistic filter size of  $\Theta_{\text{opt}} \approx 100$  states. However, from Eq. (4.7) for non-binary synapses,  $n_{\text{opt}} \approx 170/\Theta$  or equivalently  $\Theta_{\text{opt}} \approx 170/n$ . These optimal values for  $n > 2$  synapses may therefore be brought down to more biologically plausible values, of order 10. Thus, although slightly longer maximum memory lifetimes are theoretically achievable with binary synapses compared to non-binary synapses, such a possibility is not in fact realisable with biologically plausible choices of parameters. Non-binary synapses give only slightly reduced theoretically possible maximum memory lifetimes, but importantly these limits are entirely realisable with biologically plausible choices of parameters.

Fig. 8 illustrates many of these issues for non-binary synapses. In Figs. 8A and 8B we plot SNR memory lifetimes  $\tau_{\text{snr}}$  against  $n$  for various choices of  $\Theta$  and  $N$ , respectively. Memory lifetimes are determined numerically. We explicitly see a maximum in  $\tau_{\text{snr}}$  at a particular,  $\Theta$ - and  $N$ -dependent value of  $n$ , and we may confirm that this value coincides with that given in Eq. (4.7). We note that from Eq. (4.7), doubling  $\Theta$  halves  $n_{\text{opt}}$ , while Eq. (4.6) indicates that if the product  $\Theta n$  is fixed, then  $\tau_{\text{snr}}$  is unchanged. This reciprocal relationship between  $n$  and  $\Theta$  while leaving  $\tau_{\text{snr}}$  unchanged is clear in Fig. 8A. Furthermore, Eq. (4.8) indicates the maximum possible value of  $\tau_{\text{snr}}$  as a function of  $N$ , for optimally chosen  $n$  (or  $\Theta$ ). For  $N = 10^4$  used in Fig. 8A, this maximum possible SNR memory lifetime is around  $3000/r$  s, and we see precisely this maximum value for each choice of  $\Theta$  in Fig. 8A for some particular (optimal) value of  $n$ . According to Eq. (4.7), increasing  $N$  by a factor of 10 increases  $n_{\text{opt}}$  by a

factor of  $\sqrt{10} \approx 3.2$ , and this is also confirmed in Fig. 8B. More important, in Fig. 8B we observe the maximum possible value of  $\tau_{\text{snr}}$  growing linearly rather than logarithmically with  $N$ , being maximised for optimal choices of  $n$  (or  $\Theta$ ). Again, we see directly from Fig. 8B that the maximal value of  $\tau_{\text{snr}}$  is around  $0.3N/r$ , agreeing with Eq. (4.8).

In Fig. 8C, we compare SNR and MFPT memory lifetimes as a function of  $n$  for fixed  $\Theta$  and  $N$ . SNR lifetimes are determined as usual numerically while MFPT lifetimes are determined in simulation using methods described in detail elsewhere (Elliott, 2014). There is good agreement between  $\tau_{\text{snr}}$  and  $\tau_{\text{mfpt}}$  for smaller  $n$ . However, although  $\tau_{\text{snr}}$  exhibits its characteristic profile as a function of  $n$ , with an optimal value  $n_{\text{opt}}$  maximising  $\tau_{\text{snr}}$ , strikingly  $\tau_{\text{mfpt}}$  exhibits no such maximum but instead rises monotonically with  $n$  without any indication of optimality. We may resolve this apparent inconsistency between SNR and MFPT definitions of memory lifetimes by considering the variance in the first passage lifetimes. In Fig. 8D we plot the one standard deviation region around the MFPT-defined lifetimes. As  $n$  increases, the variance in the first passage lifetime increases so that for relatively small  $n$ , the MFPT-defined lifetime becomes indistinguishable from zero at the level of one standard deviation in the first passage lifetime. To understand this radical difference between the behaviour of the MFPT and SNR memory lifetimes, we must consider the storage of the tracked memory,  $\xi^0$ . Immediately after its storage, the mean

and variance in the initial memory signal, call it  $h_0 = h(0)$ , are

$$\mu_0 = \frac{2}{n} \frac{1}{\Theta^2}, \quad (4.13)$$

$$\sigma_0^2 = \frac{1}{N} \left[ \frac{1}{3} \frac{n+1}{n-1} - \mu_0^2 \right], \quad (4.14)$$

respectively. If  $\mu_0 \leq \sigma_0$ , then according to the SNR definition of memory lifetimes, the memory  $\xi^0$  is not stored successfully at  $t = 0^-$  s and its lifetime is zero. However,  $h_0$  has a non-singular distribution, so there is a finite probability  $\text{Prob}[h_0 > \sigma_0]$  that  $\xi^0$  is in fact stored successfully, in the sense that for some particular realisation,  $h_0$  is distinguishable from zero at the one standard deviation level. MFPTs are averages over both unsuccessful (giving a first passage time of zero) and successful (giving a first passage time greater than zero) initial memory storage events. An MFPT lifetime may therefore be positive while an SNR lifetime is zero. To compute  $\text{Prob}[h_0 > \sigma_0]$ , we observe that because  $h_0$  is the sum over  $N$  identically-distributed, independent random variables, by the central limit theorem  $h_0$  will be distributed very nearly as a normal distribution,  $N(\mu_0, \sigma_0^2)$ , for  $N$  large enough. Thus, the probability of the successful storage of  $\xi^0$  is

$$\begin{aligned} \text{Prob}[h_0 > \sigma_0] &\approx \frac{1}{2} \left[ 1 + \text{erf} \left( \frac{\mu_0 - \sigma_0}{\sigma_0 \sqrt{2}} \right) \right] \\ &\sim \frac{1}{2} \left[ 1 + \text{erf} \left( -\frac{1}{\sqrt{2}} \right) \right] \approx 0.16, \end{aligned} \quad (4.15)$$

where the asymptotic form arises for  $n$  or  $\Theta$  large enough; erf is the error func-

tion. Thus, the probability that  $\xi_0$  is successfully stored remains reasonably non-zero, at around 16%, even in the large  $n$  limit. This is despite the fact that the mean memory signal scales down inversely with  $n$ . As  $n$  increases, the probability of the successful storage of  $\xi^0$  decreases and asymptotes to around 0.16. A significant fraction of attempted storage events fail and return a zero first passage time; this fraction increases, but asymptotes, with  $n$ , accounting for the increasing variance in the first passage times. Those 16% of storage events that are successful, however, will contribute increasingly significantly to the MFPT because of both the  $\Theta$ -dependent signal rise and the increasing longevity of  $\mu(t)$  and thus of successfully stored memories, as  $n$  increases. These dynamics explain why  $\tau_{\text{mfpt}}$  behaves rather differently from  $\tau_{\text{snr}}$  in Fig. 8C but also why the variance in the first passage lifetimes grows to swamp the MFPT, as seen Fig. 8D.

### 4.3 Comparison to Stochastic Updater

In Fig. 9, we briefly consider results for a stochastic updater synapse, comparing them to those for a filter-based synapse. We set  $p = 1/25 = 0.04$ , which is equivalent, at least in terms of the initial mean memory signal, to a  $\Theta = 5$  filter-based synapse. In Fig. 9A, we plot  $\mu_s(t)$  for progressively increasing  $n$ , similarly to Fig. 5A for a filter-based synapse. For a stochastic updater,  $\mu_s(t)$  falls monotonically, but as  $n$  increases,  $\mu_s(t)$  takes progressively longer to peel away from its initial value of  $p$ . In the formal,  $n \rightarrow \infty$  limit,  $\mu_s(t) \equiv p$  for all time. Figs. 9B and 9C show the equivalents of Figs. 7A and 7C. Because

there is no signal growth but only decay for a stochastic updater, the overall scaling down of  $\mu(t)$  as  $n$  increases has a relatively greater impact on the SNR  $\mu(t)/\sigma(t)$ , so that for smaller values of  $n$  than for a matched filter-based synapse, a stochastic updater's SNR will not exceed unity. This then directly and significantly impacts memory lifetimes, as seen in Fig. 7C, with larger values of  $N$  being required to achieve positive memory lifetimes.

We may repeat the analysis above of optimal parameter choices and maximally enhancing  $\tau_{\text{snr}}$ . The one-decay form of Eq. (3.61) is

$$\mu(t) \approx \frac{4p}{n^2(n-1)} \cot^2 \frac{\pi}{2n} \exp \left[ -pt \left( 1 - \cos \frac{\pi}{n} \right) \right], \quad (4.16)$$

which should be compared to the one-decay form in Eq. (3.52) for a filter-based synapse. Structurally these one-decay forms for a filter-based synapse and a stochastic updater synapse are virtually identical, ultimately reflecting the underlying random walk on  $n$  strength states in both cases. From this one-decay form for a stochastic updater, we obtain

$$\begin{aligned} rp\tau_{\text{snr}} &= \frac{-1}{1 - \cos \frac{\pi}{n}} \log_e \left[ \frac{1}{4} \frac{1}{\sqrt{3N}} \sqrt{\frac{n+1}{n-1}} \frac{1}{p} n^2 (n-1) \tan^2 \frac{\pi}{2n} \right] \\ &\sim \frac{n^2}{\pi^2} \log_e \left( \frac{768}{\pi^4} \frac{p^2 N}{n^2} \right). \end{aligned} \quad (4.17)$$

For optimality we require,

$$n_{\text{opt}} = p \sqrt{\frac{768}{\pi^4 e}} \sqrt{N} \approx 1.70 p \sqrt{N}, \quad (4.18)$$



giving rise to

$$r\tau_{\text{snr}}^{\text{max}} = \frac{768}{\pi^6 e} pN \approx 0.29 pN. \quad (4.19)$$

This equation is identical to Eq. (4.8) except for the presence of a factor of  $p$ . By setting  $p = 1$ , we may therefore obtain the same maximum theoretically possible SNR memory lifetimes for a stochastic updater synapse as for a filter-based synapse. Again, given the simplicity of a  $p = 1$  stochastic updater synapse compared to a filter-based synapse, this may appear perverse. However, we must again consider whether this theoretical limit is in fact biologically realisable for a stochastic updater synapse. Setting  $p = 1$  to extremise  $\tau_{\text{snr}}^{\text{max}}$  forces  $n_{\text{opt}} \approx 1.70\sqrt{N}$  and this, for  $N = 10^4$ , gives  $n_{\text{opt}} \approx 170$  states, which is implausibly large. We may trade  $n_{\text{opt}}$  and  $p$ , just as we traded  $n_{\text{opt}}$  and  $\Theta$  for a filter-based synapse, and so reduce  $n_{\text{opt}}$  by reducing  $p$ . However, we then significantly reduce  $\tau_{\text{snr}}^{\text{max}} \approx 0.29 pN$ . For example, a choice of  $p = 1/10$  to reduce  $n_{\text{opt}}$  to around 17 for  $N = 10^4$  synapses would reduce  $\tau_{\text{snr}}^{\text{max}} \approx 0.29 pN$  to  $0.03N$ . A stochastic updater synapse therefore cannot reconcile the dual requirements of biologically plausible values of  $n$  while maintaining reasonable SNR memory lifetimes. In contrast, a filter-based synapse reconciles these requirements easily and very naturally.

## 5 Discussion

We have extended our earlier analysis of associative memory in a feedforward framework with binary-strength, integrative, filter-based synapses to consider

the more general case of discrete, multi-level synapses with  $n$  states of synaptic strength. The natural, easily-generalisable structure of our filter model ensures that exact analytical results may continue to be derived for the memory signal dynamics. Good, simplifying approximations to these exact results, specifically in the large time or large  $n$  limits, may also be extracted, facilitating the derivation of good estimates of SNR memory lifetimes.

For binary synapses, we previously observed that the memory signal initially rises before it reaches a maximum and then begins to fall back to its equilibrium value (Elliott & Lagogiannis, 2012). Ongoing memory storage actually facilitates this initial memory signal rise, in radical contrast to all other related but non-integrative models, in which the memory signal falls monotonically, and often exponentially fast. With general, discrete-state synapses, we continue to observe similar memory signal dynamics to the binary-strength case. Indeed, these dynamics are further enhanced for  $n > 2$ . The location of the memory signal peak is somewhat enhanced, increasing logarithmically with  $n$ . Relative to the binary-strength case, the (scaled) memory signal also reaches higher values. The memory signal then starts to fall, but as  $n$  increases, a quasi-plateau in the memory signal emerges, so that the memory signal remains approximately constant, falling only very slowly, for a period of time that grows quadratically with  $n$ . The usual, exponential decay then takes over, with the memory signal falling to equilibrium. These overall memory signal dynamics translate into SNR memory lifetimes that are initially quadratically enhanced as  $n$  increases, but the overall scaling down of the

memory signal with  $n$  eventually reduces and then cuts off this enhancement, reducing memory lifetimes above some optimal choice of  $n$ .

Using our exact results or the various approximations to them, an explicit formula for SNR memory lifetimes was derived. These lifetimes may be optimised by optimal selections of either  $n$ ,  $\Theta$  or  $N$ , with the other two parameters regarded as constants. A maximum memory lifetime of approximately  $0.29N/r$  s is attainable with either  $n$  or  $\Theta$  set optimally (when these parameters are adequately large to make approximations valid), so that SNR memory lifetimes grow linearly rather than logarithmically with the number of synapses. For the specific case of binary,  $n = 2$  synapses, this upper limit reaches around  $0.39N/r$  s, although it is harder to realise, biologically-speaking. With an optimal choice of  $N$  instead, memory lifetimes grow quadratically in both  $n$  and  $\Theta$ .

Many authors have previously considered memory dynamics with binary-strength or discrete synapses in either recurrent or feedforward networks, using various metrics to gauge memory lifetimes or memory capacity (see, for example, Willshaw *et al.*, 1969; Tsodyks, 1990; Amit & Fusi, 1994, Fusi *et al.*, 2005, Leibold & Kempter, 2006, 2008; Rubin & Fusi, 2007; Fusi & Abbott, 2007; Barrett & van Rossum, 2008; Huang & Amit, 2010, 2011). For example, Amit & Fusi (1994) observed that the usual logarithmic dependence of memory lifetimes on  $N$  could be overcome by setting parameters optimally. Fusi & Abbott (2007) showed that under optimal choices with discrete synapses, memory lifetimes increase quadratically with  $n$ , under balanced potentiation

and depression, as here. They argue, however, that such balanced processes constitute fine-tuning. In other work, Fusi *et al.* (2005) showed that cascade-type, binary-strength synapses do not exhibit such fine-tuning problems, and that memory lifetimes can grow like  $N^{2/3}$ , although the exponent appears to depend quite sensitively on the precise details of the fits used and may be somewhat lower (our own unpublished observations). Barrett & van Rossum (2008) considered an information-theoretic approach, computing Shannon information per synapse with discrete synapses. Under optimal choices of the learning rule (that is, the entire rule is optimised rather than merely a handful of synaptic parameters), they found that discrete, bounded synapses can perform similarly to continuous, unbounded synapses, although it is harder to relate their results directly to ours because of the focus on information content.

The argument that balanced potentiation and depression represents undesirable fine-tuning (Fusi & Abbott, 2007) is based, however, on an incomplete analysis, as we have argued elsewhere (Elliott, 2010a). Such analyses are typically based on Markov models in which the rates of potentiation and depression processes are free parameters that are set, and *fixed*, by hand. Essentially, then, the postsynaptic firing rate is decoupled from the presynaptic firing rate. Yet, with synaptic and other types of plasticity, changes in synaptic strengths will feed directly back into changes in postsynaptic firing rates. In models of synaptic plasticity that are intrinsically stable, perhaps because of an inherent fixed-point structure, and that do not require artificial hard (or soft) bounds to prevent run-away learning to stabilise them, synaptic strengths will evolve to

fixed points in which, essentially by definition, potentiation and depression processes are, on average, precisely balanced (see, for example, Bienenstock *et al.*, 1982; Burkitt *et al.* 2004, Appleby & Elliott, 2006). By effectively removing the potentially stabilising mechanisms of synaptic plasticity that couple postsynaptic firing to presynaptic firing, analyses that consider fixed firing rates therefore somewhat ironically fail to consider properly the very mechanisms of synaptic plasticity (for example leading to memory storage) that are under investigation. On this view, any viable model of synaptic plasticity should reasonably be expected to dynamically regulate the postsynaptic firing rate precisely so that depression and potentiation processes are (on average) balanced. We would not expect real synapses to rely on saturation of synaptic strengths to impose stability, because this would decrease the effective dynamical range of strengths available to synapses. Even with a finite, discrete set of strength states, we would not expect synapses' strengths to be clustered at the upper or lower ends of the available range.

Perhaps a more significant argument against maximising memory lifetimes via the device of optimally setting parameter values is that it is unlikely that neurons can tune either  $n$  (the number of available strength states) or  $N$  (the number of synapses) in the manner required. Indeed, we would expect large variations in  $N$  between neurons, and perhaps even large variations in  $n$  within any given neuron's many synapses. It is therefore unclear how a real memory system could realise optimal memory lifetimes. Furthermore, as we have seen, the optimal values of, for example,  $n$  can be implausibly large, from a biological

point of view. Given that synaptic plasticity at individual synapses is governed by a relatively small available pool of large macromolecules (see, for example, Harris & Stevens, 1989; Nusser *et al.*, 1998; Bagal *et al.*, 2005; Miller *et al.*, 2005; Asrican *et al.*, 2007) it is unlikely that the number of states of strength, or indeed the number of putative filter states, available to a synapse exceeds more than a few tens at most. Yet, a stochastic updater synapse with  $p = 1$  requires  $n_{\text{opt}} \approx 170$  and a filter-based, binary-strength synapse requires  $\Theta_{\text{opt}} \approx 100$ , for  $N = 10^4$ . These numbers are likely an order of magnitude too high. Moreover, there is the tantalising possibility that such optimality conditions are purely an artifact of the SNR metric used to gauge memory lifetimes, because we do not in fact see memory lifetimes exhibiting a maximum under variation of parameters when using MFPTs: they appear to continue to grow indefinitely. Although the variance in memory lifetimes increases, this is principally due to increasing memory encoding failure. However, the failure rate saturates at somewhat under 100%. At the worst, 16% of storage attempts succeed, so a memory system could simply employ roughly a six-fold redundancy in its architecture to ensure successful memory storage even in this worst-case scenario. If there is one thing that the mammalian brain does not lack, it is neurons.

Notwithstanding the considerations in the preceding paragraph, we note that in a filter-based model with discrete synapses, SNR memory lifetimes depend on both  $n$  and  $\Theta$ , with these playing essentially identical roles. Because the key combination is the product  $\Theta n$ , we may trade one for the other.

Specifically, increasing  $\Theta$  allows us to decrease  $n$ . Although we should be cautious of naïve optimality arguments for the reasons discussed above, by considering  $n$  and  $\Theta$  to be of order 10, we can satisfy the dual constraints of biologically-plausible synaptic values and memory lifetimes that are not too far to the left or too far to the right of the peak in maximum SNR memory lifetimes (see Fig. 8A). We might refer to this as the “sweet spot” or the “Goldilocks regime”. Without the presence of the synaptic filter variable  $\Theta$ , the requirement that  $n$  be plausibly-sized would mean that memory lifetimes are too short, while the requirement that memory lifetimes are not too short would require that  $n$  is implausibly large. Introducing  $\Theta$  allows us to bring  $n$  down to reasonable values and simultaneously ensure that memory lifetimes are not significantly compromised. Essentially, we trade “external” synaptic states (states of strength) for “internal” synaptic states (filter states). If we imagine these each to be encoded by the configurations of a single molecule (or a small ensemble of identical molecules), then for a filter-based synapse, we require two (sets of) molecules with a small handful of states, while a non-filter-based synapse requires a single molecule (or a single set of molecules) with potentially of order one hundred states.

We have also argued before that the filter size  $\Theta$  may be under dynamic regulation, with larger  $\Theta$  being used to stabilise existing memories and smaller  $\Theta$  being used to promote the rapid learning of new memories (Elliott, 2011a; Elliott & Lagogiannis, 2012). Such dynamic regulation could be achieved by regulating the kinase and phosphatase activity that is known to be critical in

synaptic plasticity (Malenka *et al.*, 1989; Mulkey *et al.*, 1993; Ferrell, 1996; Lisman & Zhabotinsky, 2001; Pi & Lisman, 2008; Pagani *et al.*, 2009). Although we think it unlikely that synapses regulate their filter thresholds to achieve precisely optimal memory lifetimes (because there is too much variability between synapses and neurons), the dependence of memory lifetimes on  $\Theta$  derived above does raise the possibility that regulating  $\Theta$  could be used as a mechanism for actively changing memories' average lifetimes.

If synaptic strengths are genuinely discrete, then how many strength states does a synapse typically have? Some evidence suggests that synapses are binary (Petersen *et al.*, 1998; O'Connor *et al.*, 2005b) while other evidence suggests that they may be ternary (Montgomery & Madison, 2002, 2004). The interpretation of such evidence is highly problematic. For example, Montgomery & Madison report the existence of one strength state (or range of strength states) for “naïve” synapses that have not undergone a long term potentiation (LTP) or long term depression (LTD) protocol, and two other strength states for synapses that have undergone LTP or LTD. We have previously argued, however, that LTP and LTD protocols may be saturating, forcing synapses to the extremes of their possible strength range, while naïve, unstimulated synapses may more naturally occupy a large range of possible strengths (Elliott, 2010a). Specifically, we showed that if a synapse has access to around 10 or more states of strength, and if LTP and LTD protocols are indeed saturating, then synapses would exhibit an effective, ternary-like structure (Elliott, 2010a). In the same work, we also used spike-timing-dependent plasticity (STDP) data from sin-



gle spine-head experiments (Harvey & Svoboda, 2007) to try to estimate the number of states of strength available to a synapse. Such an estimate can be extracted from the prediction that the classic exponential curves of STDP would be cut off at shorter interspike intervals due to saturation effects with discrete synapses (Elliott, 2010a). Although the data set is far too small for definitive conclusions, it is at least intriguing that our fits to Harvey & Svoboda’s data preferred fewer states of strength rather than more. Specifically, fits are better for  $n$  around 10. Such numbers should of course be interpreted with extreme caution, but it is at least consistent that all these analyses, including the one performed here, cohere and suggest that the number of states of synaptic strength available to a synapse may be of order 10. A very recent study, based on a full reconstruction of a small volume of hippocampal neuropil, corroborates this order by finding evidence for 26 distinguishable synaptic strengths (Bartol *et al.*, 2015).

With the assumption that synaptic strengths scale with  $n$  as in Eq. (2.1), the overall mean memory signal is scaled down as  $n$  increases. While many authors choose to scale in this manner, some do not [for example, Barrett & van Rossum (2008) do not]. From the point of view of SNR memory lifetimes, this scaling is completely irrelevant, because overall scale factors cancel in such a ratio. Physiologically, however, for larger values of  $n$ , would we expect a neuron’s activation level (akin to its membrane potential) to remain low and its firing rate therefore to remain low? As we have seen, the maximum possible initial memory signal is  $2/n$ , which is a function of both the scaling down and

the uniform, equilibrium distribution of synaptic strengths. The scaled mean memory signal  $\mu_s(t) = \frac{n}{2}\mu(t)$  removes this scaling down and can reach its maximum possible value of unity. If a neuron is relatively quiescent, we would expect homeostatic plasticity mechanisms (Turrigiano & Nelson, 2004) to intervene and restore a neuron’s activity by scaling up its synaptic strengths. Therefore, although the distinction between  $\mu(t)$  and  $\mu_s(t)$  is irrelevant for SNR memory lifetimes, it is critically important from a physiological point of view in determining a neuron’s absolute response to a memory. It seems more likely that  $\mu_s(t)$  is the relevant indicator of the tracked memory signal. On this view, a synaptic filter in the presence of more general, discrete synapses actually increases the absolute memory signal peak compared to binary-strength synapses.

In discussing our filter-based, integrative model of synaptic plasticity and related but non-integrative models, we have not considered the transition from early-phase to late-phase plasticity, which is governed by protein synthesis-dependent processes (Reymann & Frey, 2007). Late-phase plasticity appears to be controlled by synaptic tagging at strongly-stimulated synapses and their subsequent capture of plasticity-related proteins (PRPs) (Frey & Morris, 1998). However, weakly-stimulated synapses can also capture these PRPs and express late-phase rather than just early-phase plasticity if stimulated sufficiently closely in time to strongly-stimulated synapses. A few models of the transition from early-phase to late-phase plasticity exist (Clopath *et al.*, 2008; Barrett *et al.*, 2009; Pappper *et al.*, 2011). Such models could exhibit a delayed

augmentation of the memory signal because of the cross-capture of PRPs by weakly-stimulated synapses. Such dynamics would be reminiscent of the signal rise in our filter-based model. However, in the former case, these dynamics are driven heterosynaptically and still non-integratively, while in the latter case, they are driven homosynaptically and via explicit, integrative mechanisms.

Finally, we note that in our analysis of memory lifetimes above, we have not considered the possibility of a sparse coding regime. Sparseness may mean either that only a small fraction of neurons is active in any given population, or that all (or many) neurons are active but only with low firing rates. Sparse coding extends memory lifetimes when memories are correlated (Tsodyks & Feigel'man, 1988). It achieves this both by reducing the rate of plasticity induction signals experienced at individual synapses (which merely dilates time) and by reducing the interference between correlated memories. The role of sparseness in more complicated models of synaptic plasticity has also been considered (Rubin & Fusi, 2007; Leibold & Kempter, 2008). Were we to extend our analysis above to consider sparse coding and correlated memories, we would not expect the conclusions of these earlier works to be fundamentally modified. That is, we would expect that sparseness serves to extend memory lifetimes even further in a filter-based framework such as that considered here. Our present focus has been on establishing our basic framework and showing that it operates as a viable model of synaptic plasticity in different contexts. Future work could extend to including a detailed analysis of the impact of sparseness on memory lifetimes in our filter-based approach to synaptic

plasticity.

## Appendix A Averaging Over Initial Filter

### States

We show that by averaging over the initial filter states immediately after the storage of  $\xi^0$ , we may modify the argument leading to Eq. (3.19) and instead directly derive Eq. (3.15), therefore showing explicitly that the term in square brackets in Eq. (3.15) arises directly from averaging over the initial filter states.

We saw in Eq. (3.20) that the general term  $\widehat{H}_0(s) [2\widehat{G}_0(s) \times \frac{1}{2}\mathbb{C}]^m$  corresponds to the occurrence of precisely  $m$  filter threshold escape processes each with total escape density  $2G_0(t)$ , followed by no filter threshold escape processes, giving rise to the  $\widehat{H}_0(s)$  factor. Specifically, the first strength change process in the scenario considered there is governed by the probability density  $2G_0(t)$  because that scenario considered all filters to be prepared initially in the zero state. In order to account for the more general initial filter state distribution at time  $t = 0$  s, it is therefore enough to modify the very first transition process and its associated waiting time, and thereafter consider the identical renewal processes that are associated with the probability density  $2G_0(t)$  and stochastic matrix  $\frac{1}{2}\mathbb{C}$  for all subsequent changes in strength. Let the matrix  $\mathbb{K}(t)$  encode the average filter threshold escape densities immediately after the storage of  $\xi^0$  and let the diagonal matrix  $\mathbb{T}(t)$  encode the waiting times for these average filter threshold processes. Then we may modify Eq. (3.20) to

account for the initial filter states by instead writing

$$\begin{aligned}
\widehat{\mathbb{P}}(s) &= \widehat{\mathbb{T}}(s) + \widehat{H}_0(s) \widehat{\mathbb{K}}(s) + \widehat{H}_0(s) [\widehat{G}_0(s) \mathbb{C}] \widehat{\mathbb{K}}(s) \\
&\quad + \widehat{H}_0(s) [\widehat{G}_0(s) \mathbb{C}]^2 \widehat{\mathbb{K}}(s) + \widehat{H}_0(s) [\widehat{G}_0(s) \mathbb{C}]^3 \widehat{\mathbb{K}}(s) + \dots \\
&= \widehat{\mathbb{T}}(s) + \widehat{H}_0(s) [\mathbb{I} - \widehat{G}_0(s) \mathbb{C}]^{-1} \widehat{\mathbb{K}}(s) \\
&\equiv \widehat{\mathbb{T}}(s) + \widehat{H}_0(s) \widehat{\mathbb{W}}^{-1}(s) \widehat{\mathbb{K}}(s). \tag{A.1}
\end{aligned}$$

To determine  $\mathbb{K}(t)$ , we define  $\rho_I^A$  to be the probability that a synapse in strength state  $A$  is in filter state  $I$  immediately after the storage of  $\boldsymbol{\xi}^0$ . We then have that the average filter threshold escape densities are just

$$K_A^\pm(t) = \sum_I \rho_I^A G_I^\pm(t), \tag{A.2}$$

and the corresponding waiting times are  $T_A(t) = 1 - \int_0^t dt_1 [K_A^+(t_1) + K_A^-(t_1)]$ .

From Eqs. (2.11) and (3.10) we deduce that

$$K_A^\pm(t) = \begin{cases} \frac{1}{\Theta^2 - 1} [L^\pm(t) + 0 G_0(t)] & \text{for } A = 1 \\ \frac{1}{\Theta^2 + 0} [L^\pm(t) + 1 G_0(t)] & \text{for } 2 \leq A \leq n - 1 \\ \frac{1}{\Theta^2 + 1} [L^\pm(t) + 2 G_0(t)] & \text{for } A = n \end{cases}, \tag{A.3}$$

where we have defined  $L^\pm(t) = \sum_I (\Theta - |I - 1|) G_I^\pm(t)$ . If  $\Theta = 1$ , then  $K_1^\pm(t)$



The mean memory signal  $\mathbf{S}^T \widehat{\mathbb{P}}(s) \mathbf{a}$  then reduces exactly to Eq. (3.15), showing that the term in square brackets arises precisely from averaging over the initial filter states immediately after the storage of  $\boldsymbol{\xi}^0$  while the other terms arise from the non-Markovian random walk in synaptic strength with waiting times governed by  $G_0(t)$ .

## Appendix B Eigenstructure of $\mathbb{C}$

We define  $\mathbb{C}_n$  to be the  $n \times n$  matrix  $\mathbb{C}$  in Eq. (3.12), and the auxiliary matrix  $\mathbb{E}_n$  to be identical to  $\mathbb{C}_n$  except that the 1, 1 entry of  $\mathbb{E}_n$  is zero instead of unity.

We then have, schematically,

$$\mathbb{C}_{n+1} = \left( \begin{array}{c|ccc} 1 & 1 & 0 & \dots \\ \hline 1 & & & \\ 0 & & \mathbb{E}_n & \\ \vdots & & & \end{array} \right) \quad \text{and} \quad \mathbb{E}_{n+1} = \left( \begin{array}{c|ccc} 0 & 1 & 0 & \dots \\ \hline 1 & & & \\ 0 & & \mathbb{E}_n & \\ \vdots & & & \end{array} \right). \quad (\text{B.1})$$

We define  $\mathbb{E}_1 = 1$  so that  $\mathbb{C}_2$  and  $\mathbb{E}_2$  are correct. We define  $x_n = \det(\mathbb{C}_n - \Lambda \mathbb{I}_n)$  and  $y_n = \det(\mathbb{E}_n - \Lambda \mathbb{I}_n)$ , where  $\mathbb{I}_n$  is the  $n \times n$  identity matrix. With these definitions, from the tridiagonal structure of  $\mathbb{C}_n$  and  $\mathbb{E}_n$  we then have the recurrence relations

$$x_n = (1 - \Lambda) y_{n-1} - y_{n-2}, \quad (\text{B.2})$$

$$y_n = -\Lambda y_{n-1} - y_{n-2}, \quad (\text{B.3})$$

where  $y_1 = 1 - \Lambda$  and we take  $y_0 \equiv 1$  to ensure that  $x_2$  and  $y_2$  are correct and hence all subsequent values. The solution of Eq. (B.3) is

$$y_n = \frac{\Psi_+^n (1 + \Psi_+) - \Psi_-^n (1 + \Psi_-)}{\Psi_+ - \Psi_-}, \quad (\text{B.4})$$

where  $\Psi_{\pm} = \frac{1}{2}(-\Lambda \pm \sqrt{\Lambda^2 - 4})$ , or  $\Psi_+ + \Psi_- = -\Lambda$  and  $\Psi_+ \Psi_- = 1$ . Inserting this solution for  $y_n$  into Eq. (B.2), we find that

$$x_n = (2 - \Lambda) \frac{\Psi_+^n - \Psi_-^n}{\Psi_+ - \Psi_-}. \quad (\text{B.5})$$

The RHS always simplifies to a polynomial of degree  $n$  in  $\Lambda$ , as it should. The  $n$  eigenvalues  $\Lambda_m$ ,  $m = 0, \dots, n - 1$ , of  $\mathbb{C}_n$  are therefore  $\Lambda_0 = 2$  and  $\Lambda = -\Psi_+ - 1/\Psi_+$  where the  $\Psi_+$  are determined from the  $n - 1$  solutions of  $\Psi_+^{2n} = 1$  in the lower half of the complex plane, or  $\Psi_+ = -e^{im\pi/n}$ ,  $m = 1, \dots, n - 1$ . (The solutions in the upper half of the complex plane merely enumerate the eigenvalues in reverse order, or amount to the re-definition  $m \rightarrow n - m$ .) We therefore have  $\Lambda_m = 2 \cos \frac{m\pi}{n}$ , for  $m = 1, \dots, n - 1$ , although  $m = 0$  also correctly reproduces  $\Lambda_0 = 2$ . Hence, we may write  $\Lambda_m = 2 \cos \frac{m\pi}{n}$  for  $m = 0, \dots, n - 1$ . The eigenvalues of  $\frac{1}{2}\mathbb{C}$  are then just  $\lambda_m = \frac{1}{2}\Lambda_m = \cos \frac{m\pi}{n}$ , as stated in Eq. (3.24).

We now determine the normalised eigenvector  $\mathbf{e}^m$  of  $\mathbb{C}$  corresponding to eigenvalue  $\Lambda_m$ . The vector  $\mathbf{e}^m$  has components  $e_A^m$ ,  $A = 1, \dots, n$ . Ignoring for the moment the two boundary cases, these components must satisfy the



recurrence relation

$$e_{A-1}^m - \Lambda_m e_A^m + e_{A+1}^m = 0. \quad (\text{B.6})$$

Comparing this equation to the recurrence relation for  $y_n$ ,

$$y_{n-1} + \Lambda_m y_n + y_{n+1} = 0,$$

the components  $e_A^m$  satisfy  $e_A^m \propto (-1)^{A-1} y_{A-1}$ , or

$$\mathbf{e}^m \propto (+y_0, -y_1, +y_2, -y_3, \dots, (-1)^{n-1} y_{n-1})^T. \quad (\text{B.7})$$

It is important to note that the sequence of determinants in this expression is computed using Eq. (B.4) but with  $n$  in  $\Psi_{\pm} = -e^{\pm im\pi/n}$  held fixed at the dimensionality of  $\mathbf{e}^m$ . Using Eq. (B.4), we then find that

$$e_A^m \propto \sin \frac{Am\pi}{n} - \sin \frac{(A-1)m\pi}{n}, \quad (\text{B.8})$$

for  $m > 0$ , and  $e_A^0 \propto 1$  for the particular case  $m = 0$ . We may check explicitly that the two boundary equations,

$$\begin{aligned} 0 &= (1 - \Lambda_m) e_1^m + e_2^m, \\ 0 &= e_{n-1}^m + (1 - \Lambda_m) e_n^m, \end{aligned}$$

are also satisfied. Normalising and simplifying, we finally obtain the results for  $e_A^m$  stated in Eq. (3.25).

## References

- Amit, D.J., & Fusi, S. (1994). Learning in neural networks with material synapses. *Neural Comput.*, **6**, 957–982.
- Appleby, P.A., & Elliott, T. (2006). Stable competitive dynamics emerge from multispikes interactions in a stochastic model of spike-timing-dependent plasticity. *Neural Comput.*, **18**, 2414–2464.
- Asrican, B., Lisman, J., & Otmakhov, N. (2007). Synaptic strength of individual spines correlates with bound  $\text{Ca}^{2+}$ -calmodulin-dependent kinase II. *J. Neurosci.*, **27**, 14007–14011.
- Bagal, A.A., Kao, J.P.Y., Tang, C.-M., & Thompson, S.M. (2005). Long-term potentiation of exogenous glutamate responses at single dendritic spines. *Proc. Natl. Acad. Sci. U.S.A.*, **102**, 14434–14439.
- Barrett, A., Billings, G., Morris, R., & van Rossum, M. (2009). State based model of long-term potentiation and synaptic tagging and capture. *PLoS Comput. Biol.*, **5**, e1000259.
- Barrett, A.B., & van Rossum, M.C.W. (2008). Optimal learning rules for discrete synapses. *PLoS Comput. Biol.*, **4**, e1000230.
- Bartol, T.M., Bromer, C., Kinney, J., Chirillo, M.A., Bourne, J.N., Harris, K.M., & Sejnowski, T.J. (2015). Nanoconnectomic upper bound on the variability of synaptic plasticity. *eLife*, **4**, e10778.

- Bienenstock, E.L., Cooper, L.N., & Munro, P.W. (1982). Theory for the development of neuron selectivity: Orientation specificity and binocular interaction in visual cortex. *J. Neurosci.*, **2**, 32–48.
- Burkitt, A.N., Meffin, H., & Grayden, D.B. (2004). Spike-timing-dependent plasticity: The relationship to rate-based learning for models with weight dynamics determined by a stable fixed point. *Neural Comput.*, **16**, 885–940.
- Clopath, C., Ziegler, L., Vasilaki, E., Büsing, L., & Gerstner, W. (2008). Tag-trigger-consolidation: A model of early and late long-term-potential and depression. *PLoS Comput. Biol.*, **4**, e1000258.
- Cox, D.R., & Miller, H.D. (1965). *The Theory of Stochastic Processes*. London: Methuen.
- Elliott, T. (2008). Temporal dynamics of rate-based plasticity rules in a stochastic model of spike-timing-dependent plasticity. *Neural Comput.*, **20**, 2253–2307.
- Elliott, T. (2010a). Discrete states of synaptic strength in a stochastic model of spike-timing-dependent plasticity. *Neural Comput.*, **22**, 244–272.
- Elliott, T. (2010b). A non-Markovian random walk underlies a stochastic model of spike-timing-dependent plasticity. *Neural Comput.*, **22**, 1180–1230.
- Elliott, T. (2011a). The mean time to express synaptic plasticity in stochas-

- tic, integrate-and-express models of synaptic plasticity induction. *Neural Comput.*, **23**, 124–159.
- Elliott, T. (2011b). Stability against fluctuations: Scaling, bifurcations and spontaneous symmetry breaking in stochastic models of synaptic plasticity. *Neural Comput.*, **23**, 674–734.
- Elliott, T. (2014). Memory nearly on a spring: A mean first passage time approach to memory lifetimes. *Neural Comput.*, **26**, 1873–1923.
- Elliott, T., & Lagogiannis, K. (2009). Taming fluctuations in a stochastic model of spike-timing-dependent plasticity. *Neural Comput.*, **21**, 3363–3407.
- Elliott, T., & Lagogiannis, K. (2012). The rise and fall of memory in a model of synaptic integration. *Neural Comput.*, **24**, 2604–2654.
- Ferrell, J.E. (1996). Tripping the switch fantastic: how a protein kinase cascade can convert graded inputs into switch-like outputs. *Trends Biochem. Sci.*, **21**, 460–466.
- Frey, U., & Morris, R.G.M. (1998). Synaptic tagging: Implications for the late maintenance of hippocampal long-term potentiation. *Trends Neurosci.*, **21**, 181–188.
- Fusi, S., & Abbott, L.F. (2007). Limits on the memory storage capacity of bounded synapses. *Nature Neurosci.*, **10**, 485–493.

- Fusi, S., Drew, P.J., & Abbott, L.F. (2005). Cascade models of synaptically stored memories. *Neuron*, **45**, 599–611.
- Harris, K.M., & Stevens, J.K. (1989). Dendritic spines of CA1 pyramidal cells in the rat hippocampus: serial electron microscopy with reference to their biophysical characteristics. *J. Neurosci.*, **9**, 2982–2887.
- Harvey, C.D., & Svoboda, K. (2007). Locally dynamic synaptic learning rules in pyramidal neuron dendrites. *Nature*, **450**, 1195–2002.
- Hopfield, J.J. (1982). Neural networks and physical systems with emergent collective computational abilities. *Proc. Natl. Acad. Sci. U.S.A.*, **79**, 2554–2558.
- Huang, Y., & Amit, Y. (2010). Precise capacity analysis in binary networks with multiple coding level inputs. *Neural Comput.*, **22**, 660–688.
- Huang, Y., & Amit, Y. (2011). Capacity analysis in multi-state synaptic models: A retrieval probability perspective. *J. Comput. Neurosci.*, **30**, 699–720.
- Leibold, C., & Kempster, R. (2006). Memory capacity for sequences in a recurrent network with biological constraints. *Neural Comput.*, **18**, 904–941.
- Leibold, C., & Kempster, R. (2008). Sparseness constrains the prolongation of memory lifetime via synaptic metaplasticity. *Cereb. Cortex*, **18**, 67–77.
- Lisman, J., & Zhabotinsky, A.M. (2001). A model of synaptic memory:

- A CaMKII/PP1 switch that potentiates transmission by organizing an AMPA receptor anchoring assembly. *Neuron*, **31**, 191–201.
- Malenka, R.C., Kauer, J.A., Perkel, D.J., Mauk, M.D., Kelly, P.T., Nicoll, R.A., & Waxham, M.N. (1989). An essential role for postsynaptic calmodulin and protein kinase activity in long-term potentiation. *Nature*, **340**, 554–557.
- Miller, P., Zhabotinsky, A.M., Lisman, J.E., & Wang, X.-J. (2005). The stability of a stochastic CaMKII switch: Dependence on the number of enzyme molecules and protein turnover. *PLoS Biol.*, **3**, 0705.
- Montgomery, J.M., & Madison, D.V. (2002). State-dependent heterogeneity in synaptic depression between pyramidal cell pairs. *Neuron*, **33**, 765–777.
- Montgomery, J.M., & Madison, D.V. (2004). Discrete synaptic states define a major mechanism of synapse plasticity. *Trends Neurosci.*, **27**, 744–750.
- Mulkey, R.M., Herron, C.E., & Malenka, R.C. (1993). An essential role for protein phosphatases in hippocampal long-term depression. *Science*, **261**, 1104–1107.
- Nadal, J.P., Toulouse, G., Changeux, J.P., & Dehaene, S. (1986). Networks of formal neurons and memory palimpsests. *Europhys. Lett.*, **1**, 535–542.
- Nusser, Z., Lujan, R., Laube, G., Roberts, J.D., Molnar, E., & Somogyi, P. (1998). Cell type and pathway dependence of synaptic AMPA receptor number and variability in the hippocampus. *Neuron*, **21**, 545–559.

- O'Connor, D.H., Wittenberg, G.M., & Wang, S.S.-H. (2005a). Dissection of bidirectional synaptic plasticity into saturable unidirectional process. *J. Neurophysiol.*, **94**, 1565–1573.
- O'Connor, D.H., Wittenberg, G.M., & Wang, S.S.-H. (2005b). Graded bidirectional synaptic plasticity is composed of switch-like unitary events. *Proc. Natl. Acad. Sci. U.S.A.*, **102**, 9679–9684.
- Pagani, M.R., Oishi, K., Gelb, B.D., & Zhong, Y. (2009). The phosphatase SHP2 regulates the spacing effect for long-term memory induction. *Cell*, **139**, 186–198.
- Päpper, M., Kempter, R., & Leibold, C. (2011). Synaptic tagging, evaluation of memories, and the distal reward problem. *Learn. and Mem.*, **18**, 58–70.
- Parisi, G. (1986). A memory which forgets. *J. Phys. A: Math. Gen.*, **19**, L617–L620.
- Petersen, C.C.H., Malenka, R.C., Nicoll, R.A., & Hopfield, J.J. (1998). All-or-none potentiation at CA3-CA1 synapses. *Proc. Natl. Acad. Sci. U.S.A.*, **95**, 4732–4737.
- Pi, H.J., & Lisman, J.E. (2008). Coupled phosphatase and kinase switches produce the tristability required for long-term potentiation and long-term depression. *J. Neurosci.*, **28**, 13132–13138.
- Reymann, K., & Frey, J. (2007). The late maintenance of hippocampal LTP:

- Requirements, phases, ‘synaptic tagging’, ‘late-associativity’ and implications. *Neuropharm.*, **52**, 24–40.
- Rubin, D.D.B.D., & Fusi, S. (2007). Long memory lifetimes require complex synapses and limited sparseness. *Front. Comput. Neurosci.*, **1**, 7.
- Sobczyk, A., & Svoboda, K. (2007). Activity-dependent plasticity of the NMDA-receptor fractional  $\text{Ca}^{2+}$  current. *Neuron*, **53**, 17–24.
- Tsodyks, M.V. (1990). Associative memory in neural networks with binary synapses. *Mod. Phys. Lett. B*, **4**, 713–716.
- Tsodyks, M.V., & Feigel’man, M.V. (1988). The enhanced storage capacity in neural networks with low activity levels. *Europhys. Letts.*, **6**, 101–105.
- Turrigiano, G.G., & Nelson, S.B. (2004). Homeostatic plasticity in the developing nervous system. *Nature Rev. Neurosci.*, **5**, 97–107.
- van Kampen, N.G. (1992). *Stochastic Processes in Physics and Chemistry*. Amsterdam: Elsevier.
- Willshaw, D.J., Duneman, O.P., & Longuet-Higgins, H. (1969). Nonholographic associative memory. *Nature*, **222**, 960–962.
- Yasuda, R., Sabatini, B.L., & Svoboda, K. (2003). Plasticity of calcium channels in dendritic spines. *Nature Neurosci.*, **6**, 948–955.



## Figure Captions

**Figure 1:** A filter-based mechanism for the integration of synaptic plasticity induction signals leading to the expression of synaptic plasticity at filter thresholds. Synaptic filter states are represented by the circled numbers,  $-(\Theta - 1), \dots, +(\Theta - 1)$ . Plasticity induction signals occur at Poisson rate  $r$ , with potentiation signals (arrows  $\uparrow$  and  $\uparrow\uparrow$ ) and depression signals (arrows  $\downarrow$  and  $\downarrow\downarrow$ ) being equiprobable, with probability  $\frac{1}{2}$ . Potentiating induction signals acting on filter states  $-(\Theta - 1), \dots, +(\Theta - 2)$  lead only to increments in filter state (indicated by  $\uparrow$ ), while a potentiating induction signal acting on filter state  $+(\Theta - 1)$  causes the filter to reach its upper threshold, leading to the expression of a potentiation step if possible ( $\uparrow\uparrow$ ) and resetting the filter state to zero. Similarly, depressing induction signals acting on states  $-(\Theta - 2), \dots, +(\Theta - 1)$  decrement the filter state ( $\downarrow$ ) while a depressing induction signal acting on state  $-(\Theta - 1)$  leads to the expression of a depression step if possible ( $\downarrow\downarrow$ ) and resetting the filter state to zero.

**Figure 2:** The function  $\chi_n(t)$  as a function of time, for different choices of  $n$  and  $\Theta$ . (A)  $\chi_n(t)$  for  $n = 2^2, 2^3, 2^4, 2^5$  and  $2^6$  (moving right to left in the graph, with smaller  $n$  corresponding to overall smaller  $\chi_n(t)$ ) for the particular choice,  $\Theta = 4$ , as indicated. (B)  $\chi_n(t)$  for  $\Theta = 2, 4, 6, 8$  and  $10$  (moving top to bottom in the graph, with smaller  $\Theta$  corresponding to larger maxima for  $\chi_n(t)$ ) for the particular choice,  $n = 8$ , as indicated. (C) For comparison, we

show the probability density for escape through either filter threshold starting from the zero filter state,  $2G_0(t)$ , for the same values of  $\Theta$  used in (B), with again smaller values of  $\Theta$  corresponding to larger maxima for  $2G_0(t)$ . The initial, small time profile of  $\chi_n(t)$  in  $B$  follows very closely that for  $2G_0(t)$ .

**Figure 3:** Validation of analytical results for the mean memory signal. (A) Comparison between analytical and simulation results for  $\Theta = 4$  for different choices of  $n$ , as indicated. The agreement is exact. Simulations are averaged over  $10^4$  trials, and we have used  $N = 10^4$  synapses for better self-averaging within trials, although the analytical result for  $\mu(t)$  is independent of  $N$ . (B) Same as A, except that  $\Theta = 8$ . (C) and (D) show the scaled mean memory signal  $\mu_s(t) = \frac{n}{2}\mu(t)$ . The scaling allows us to see much more clearly that increasing  $n$  causes the memory signal to rise for longer and decay at least initially more slowly.

**Figure 4:** Comparison between exact and approximate results for  $\mu_s(t)$ . (A) Use of the approximate form for  $\chi(t)$  given in Eq. (3.45) in the convolution form for  $\mu_s(t)$  in Eq. (3.38). This approximation preserves  $\mu_s(t) \approx \mu_2(t)$  at small times and asymptotes to the exact form of  $\mu_s(t)$  at large times, deviating from  $\mu_s(t)$  only at intermediate times. (B) Use of the one-decay form in Eq. (3.52) in which only the slowest decaying term in Eq. (3.50) is retained. This approximation, being asymptotically exact at large times, is useful for determining memory lifetimes. (C) Use of the two-decay form in Eq. (3.51),

which retains the slowest decaying term from each of the separate sums in Eq. (3.51). This approximation is useful for estimating the location of the peak in the mean memory signal and its subsequent decay dynamics. (D) Use of the large  $n$  form in Eq. (3.55). This approximation for  $\mu_s(t)$  is essentially exact for all except large times. In all cases we have set  $\Theta = 6$  and  $n = 10$ , with even this relatively small value of  $n$  being sufficiently large for the large  $n$  limiting forms of  $\mu_s(t)$  to provide good approximations.

**Figure 5:** Dynamics of scaled mean memory signal for large  $n$ . (A)  $\mu_s(t)$  for  $n = 2^1, 2^2, 2^3, 2^4, 2^5, 2^6, 2^7, 2^8, 2^9, 2^{10}, 2^{11}, 2^{12}$  (solid lines, moving left to right in the figure), and for the formal  $\mathcal{O}(1)$  limit,  $n \rightarrow \infty$  in Eq. (3.56) (dotted line, which can only be distinguished from the solid lines at large times). (B) Comparison between the exact form of  $\mu_s(t)$  (solid lines) and the two-decay form in Eq. (3.51) (dashed lines) for  $n = 2^2, 2^5, 2^8, 2^{11}$  (moving left to right in the figure). Also shown as the dotted lines are the large  $n$  forms in Eqs. (3.56) and (3.58). (C) The location of the peak in the mean memory signal from both the exact form and the two-decay form of  $\mu_s(t)$ , as a function of  $n$ . The location increases only logarithmically in  $n$  for both forms, although the two-decay approximation systematically overestimates the peak location. (D) The time at which the mean memory signal falls to 95% of its peak value. For  $n$  sufficiently large (which is only  $n \sim 10$  or  $n \sim 20$ ), this time grows quadratically in  $n$ . Again, the two-decay approximation overestimates the location of 95% of the signal peak.

**Figure 6:** Determination of memory lifetimes based on various approximations to the signal-to-noise ratio. (A) Validation of numerical matrix methods by comparing the numerically-computed variance  $\sigma(t)^2$  in the tracked memory signal to that obtained from simulations. The agreement is exact. Simulations are averaged over  $10^4$  trials. (B) The dependence of the covariance between two synapses' strengths on  $n$  over time. The second peak in covariance is increasingly suppressed relative to the first peak as  $n$  increases. (C) The SNR  $\mu(t)/\sigma(t)$  for three different forms of  $\sigma(t)$ , as indicated: the full, exact form in Eq. (2.10); the full form but without the covariance term, so just  $\sigma(t)^2 \approx \{\mathbf{E}[S(t)^2] - \mu(t)^2\}/N$ ; the simplest form  $\sigma(t)^2 \approx \mathbf{E}[S(t)^2]/N$ . When the memory lifetime is relatively small, approximations to  $\sigma(t)$  overestimate  $r\tau_{\text{snr}}$ . For this choice of  $\Theta$ ,  $n$  and  $N$ , the approximate forms overestimate  $r\tau_{\text{snr}}$  by around 20%. (D) As C, but for different parameters, leading to longer memory lifetimes. In this case the overestimate is merely around 0.5%.

**Figure 7:** Parameter-dependence of SNRs and memory lifetimes, using  $\sigma(t)^2 \approx \{\mathbf{E}[S(t)^2] - \mu(t)^2\}/N$ . (A) SNRs for different values of  $n$  as a function of time. Although increasing  $n$  suppresses  $\mu(t)$  by  $1/n$ , the mean memory signal near peak is sustained roughly speaking  $n^2$ -fold longer, offsetting this  $1/n$  suppression. (B) SNRs for different values of  $N$  as a function of time. Increasing  $N$  also increases the SNR, but memory lifetimes increase only logarithmically with  $N$ . (C) Dependence of SNR memory lifetimes on  $N$  for different values of

$n$ . Thick lines show results from the exact form for  $\mu(t)$ , while thin lines show for comparison results from the one-decay form for  $\mu(t)$  leading to Eq. (4.5). For  $N$  too small, the SNR never exceeds unity, so a bifurcation exists in  $N$ , defining a critical value of  $N$  above which memories are successfully stored. The equation for  $\tau_{\text{snr}}$  in Eq. (4.5) takes no account of this bifurcation.

**Figure 8:** Optimal values of  $n$  or  $N$  establish maximal memory lifetimes.

(A) Dependence of SNR memory lifetimes on  $n$  for various choices of  $\Theta$ , for  $N = 10^4$ . For each choice of  $\Theta$  there exists an optimal value of  $n$  that maximises memory lifetimes. (B) Dependence of SNR memory lifetimes on  $n$  for various choices of  $N$ , for  $\Theta = 5$ . Again, an optimal value of  $n$  maximises memory lifetimes. (C) Comparison between SNR and MFPT memory lifetimes. Although the SNR memory lifetime exhibits a maximum for a particular value of  $n$ , the MFPT memory lifetime does not. (D) For the same parameters used in C, we show the one standard deviation region around the MFPT, determined from the variance in the first passage times. For  $n$  large enough, the MFPT is indistinguishable from zero at the one standard deviation level.

**Figure 9:** Memory performance of a stochastic updater synapse with  $p =$

$1/25$ . (A)  $\mu_s(t)$  for  $n = 2^1, 2^2, 2^3, 2^4, 2^5, 2^6, 2^7, 2^8, 2^9, 2^{10}, 2^{11}, 2^{12}$  (moving left to right in the figure). (B) SNRs for different values of  $n$  as a function of time, using  $\sigma(t)^2 \approx \{\mathbf{E}[S(t)^2] - \mu(t)^2\}/N$ . (C) Dependence of SNR memory lifetimes on  $N$  for different values of  $n$ . Lifetimes are determined numerically

rather than from a one-decay approximation to Eq. (3.61).

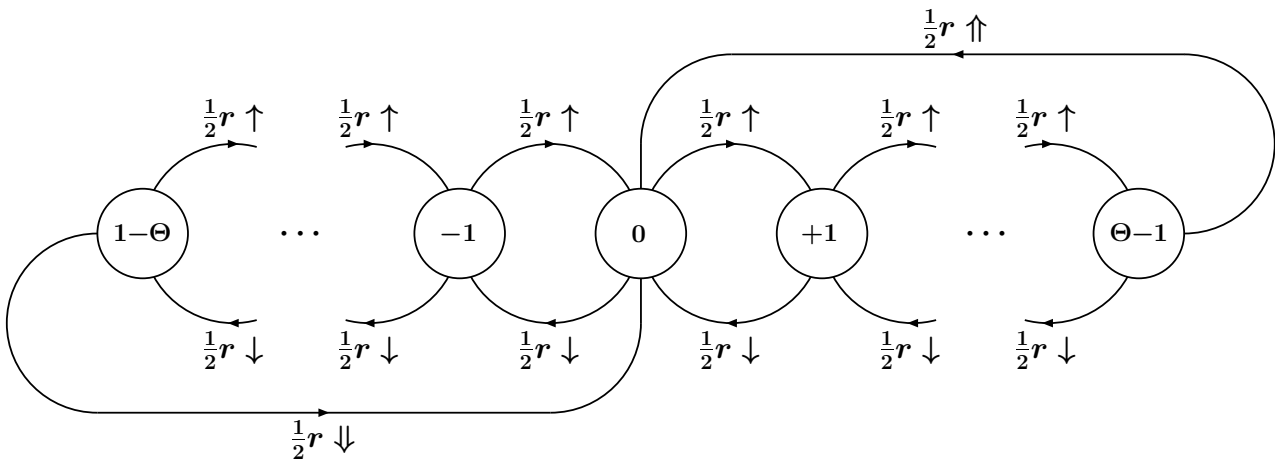
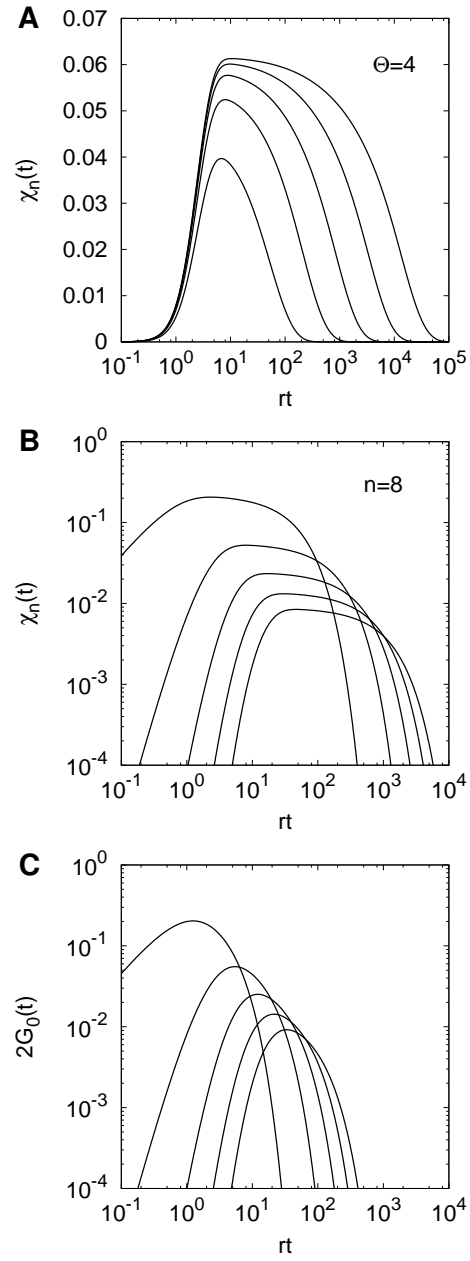


Figure 1



**Figure 2**



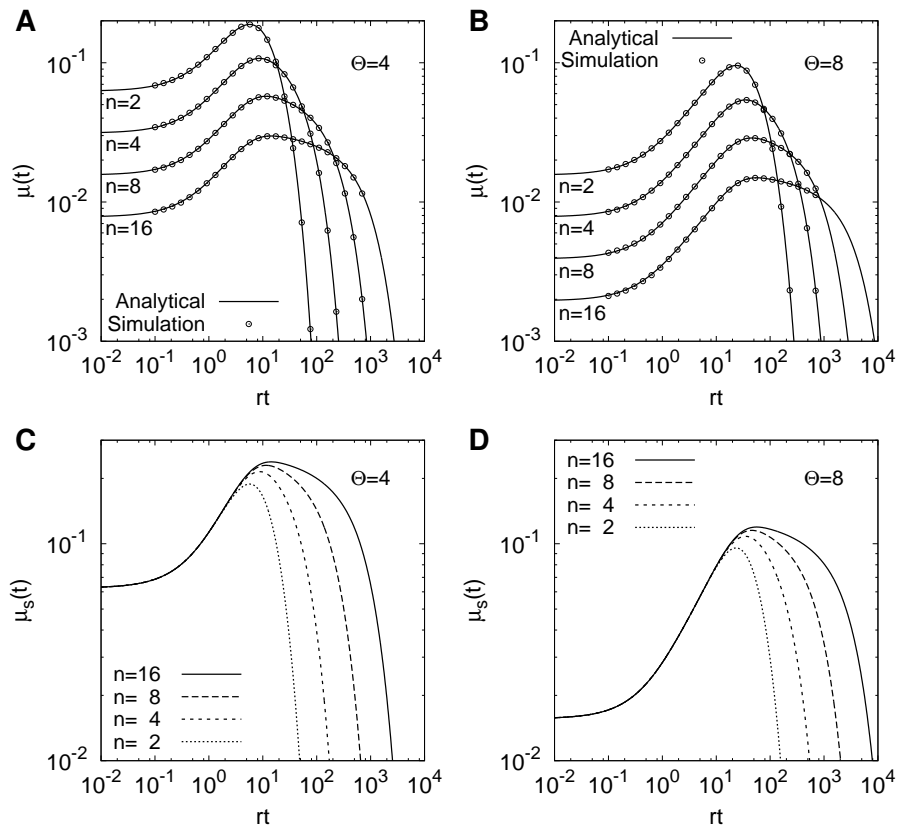


Figure 3

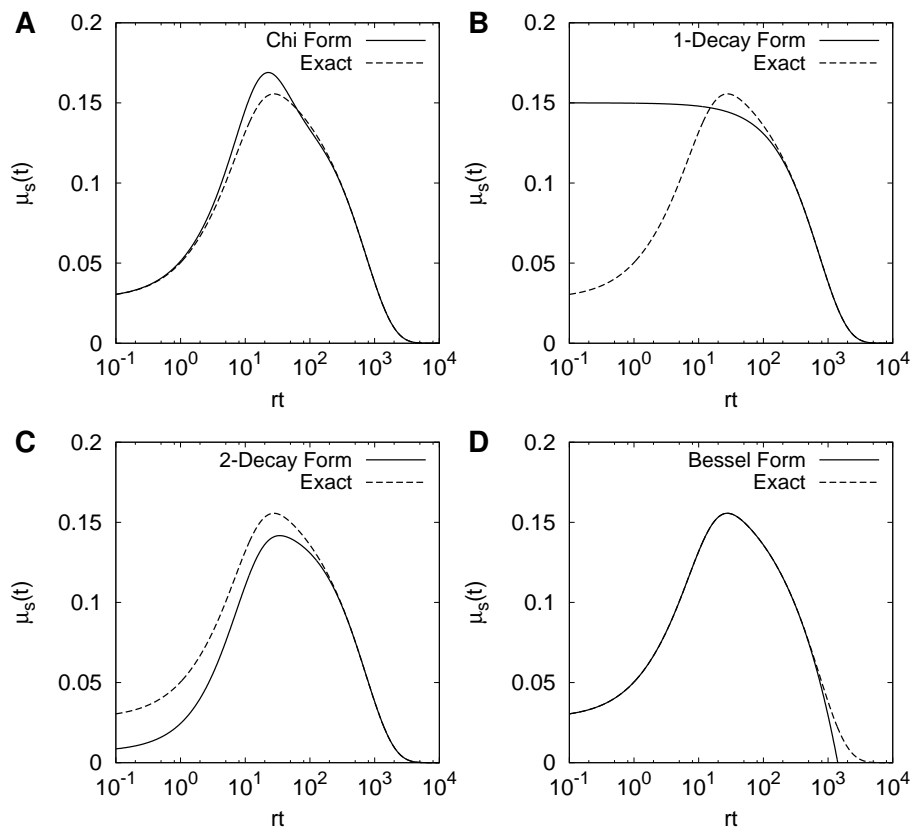


Figure 4

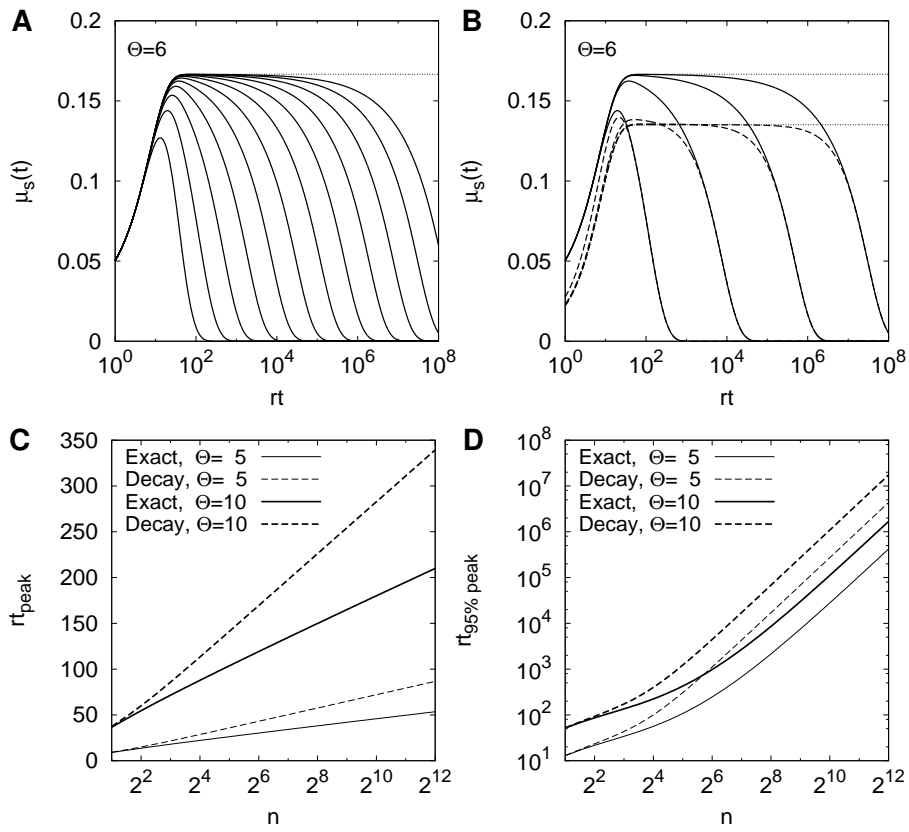


Figure 5

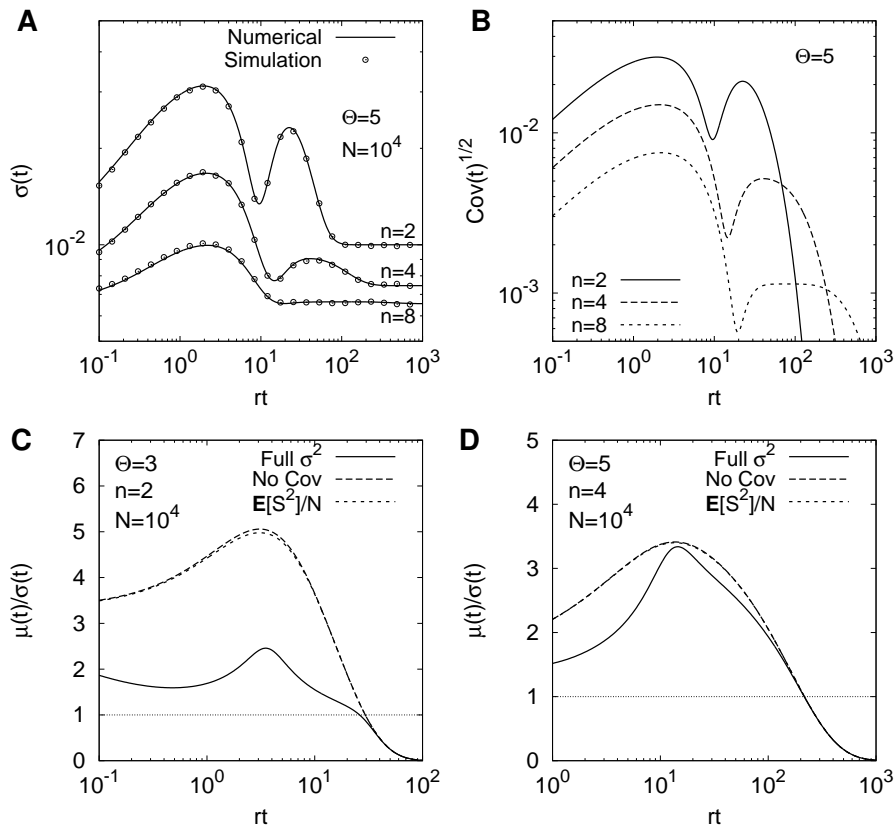


Figure 6

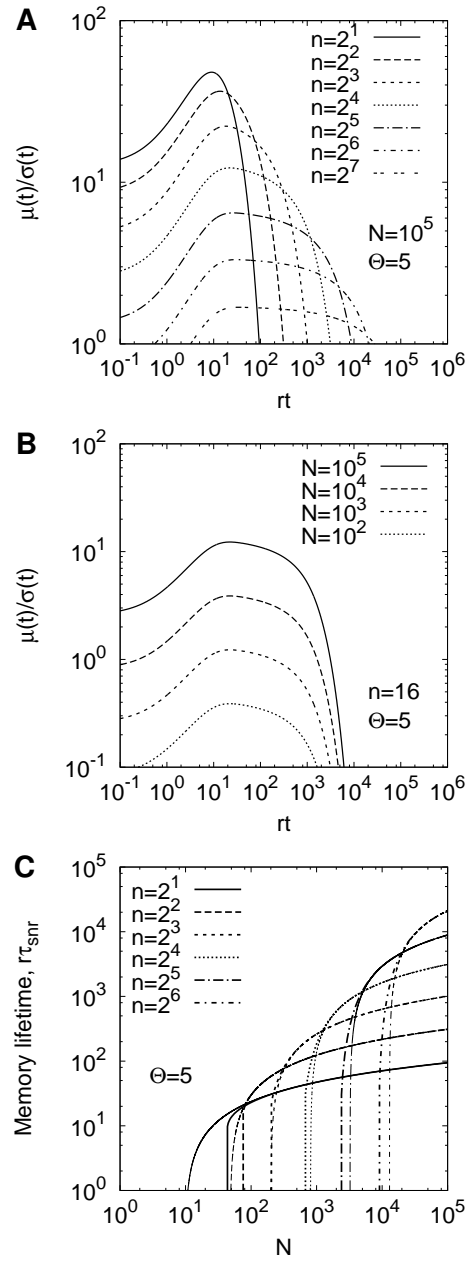


Figure 7

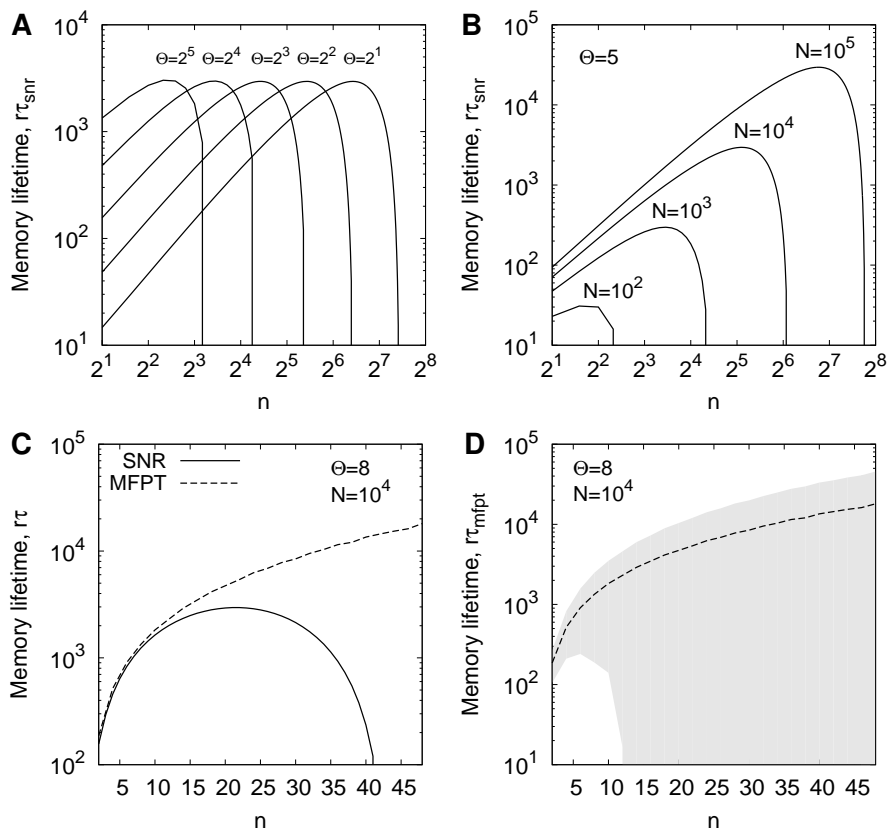


Figure 8

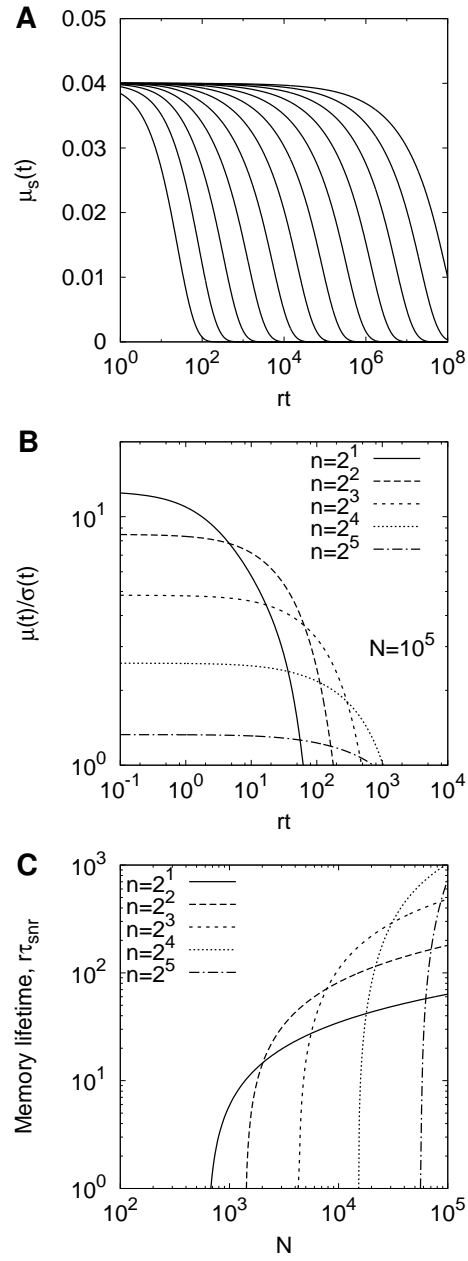


Figure 9