
Extending spike sorting for simultaneous large-scale extra-cellular electrical stimulation and recording

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Simultaneous electrical stimulation and recording of the neural tissue in micro-
2 electrode arrays (MEAs) has been proposed as a prominent technology for achiev-
3 ing new advances in neuroscience. However, to fully exploit its capabilities in
4 large-scale contexts some technical problems have to be overcome. One of such
5 problems has to do with the corruptions in recordings induced by the stimulus (arti-
6 facts), which limits heavily the identifiability of neural signals. Here we provide
7 a principled solution to this problem: by modeling the stimulation artifact as a
8 Gaussian Process we are able to faithfully represent its structure, and to come up
9 with an algorithm that quickly estimates and subtracts this artifact. The effectiveness
10 of our method is demonstrated in both real and simulated large-scale data (512
11 electrodes).

12 1 Introduction

13 There is a consensus that many new developments in systems neuroscience and neural engineering
14 (e.g. prosthetic devices) will rely in our ability to control neural activity through the exogenous
15 stimulation of the neural tissue. In this line, the feedback control of the neural tissue has been
16 proposed as a particularly fruitful research avenue. However, such closed-loop interactions require a
17 proper read-out of elicited neural activity in response to stimuli, calling for a methodological shift, as
18 many current methods and technologies for this read-out were conceived for the passive observation
19 of neurons, but break down in regimes dominated by stimulation.

20 We focus on the simultaneous extra-cellular electrical stimulation and recording, where a proper
21 computational infrastructure for the analysis of neural signals in the large-scale case is still lacking,
22 as spike sorting—the method that allows identification of neurons from their electrophysiological
23 spatio-temporal fingerprints—cannot handle the corruptions induced by electrical stimuli, since these
24 corruptions—the stimulation artifacts—can be of much greater magnitude than and overlap with the
25 actual neural activity.

26 Although many approaches have been proposed to tackle this problem, they all fall short in important
27 aspects: they are often based on very restrictive assumptions and don't exploit the inherently spatial
28 setup imposed by the MEAs. In consequence, often there is no other solution than throwing away an
29 important proportion of recordings, which leads to biased research results that not necessarily have to
30 do with the regimes where the most interesting neuronal dynamics do occur.

31 In this work we develop a modern large-scale principled framework for the analysis of neural data in
32 this regime corrupted by artifacts, based on a comprehensive account of the variability of the artifact
33 in the spatio-temporal and stimulus dimensions. Specifically, we characterize this highly structured
34 artifact and model it in the Gaussian Process framework, where we leverage recent advances in
35 machine learning to enable a fast and scalable implementation. Then, we come up with a spike

36 sorting algorithm and demonstrate its effectiveness by comparison to human-curated ground truth
 37 from the primate retina. Although some features of our method are context-dependent, we discuss
 38 extensions to other scenarios, stressing the generality of our approach.

39 2 Method

40 In this section we develop a method for identifying neural activity in response to electrical stimulation.
 41 At the highest level, it is based on the following generative model

$$Y = A + s + \epsilon, \quad (1)$$

42 where Y represent the observed traces, A is the stimulation artifact, s is the neural activity and ϵ is a
 43 noise term. We divide the exposition of our method in four parts. First, we describe the nature of the
 44 data that is available for analysis (2.1). Second, we describe the structure of the stimulation artifacts
 45 (2.2). Third, we propose a Gaussian Process model for these stimulation artifacts (2.3). Finally, we
 46 come up with an algorithm that produces an estimate of A and s (2.4).

47 2.1 Setup

48 Data is made up by: (1) recordings in response to stimulation, including stimulation covariates and
 49 (2) electrical images (EI); the spatio-temporal fingerprints of targeted neurons.

50 Regarding (1), we assume voltage traces are available over a set of E electrodes (the array) and a
 51 discrete time window of length T (representing a multiple of the sampling rate) in response to J
 52 different stimuli on a single stimulating electrode (fixed). Although stimuli can depend on many
 53 features (e.g. specific structure of the current pulses that are passed through through stimulating
 54 electrode) we assume strength or amplitude of stimulus can be finally summarized in terms of the
 55 quantities a_j , assumed increasing. Also, for each stimulus, we assume a number of n_j repetitions are
 56 available.

57 Regarding (2), we assume EIs are available for all of the N neurons under study: each of these
 58 templates is represented as a $E \times T'$ matrix containing an estimate of the voltage deflections produced
 59 by a spike over the array in a length T' time window, and aligning the onset of a spike to an arbitrary
 60 value. EIs can be obtained in a separate experiment in the absence of electrical stimulation, using
 61 standard large-scale spike sorting methods. Figure 4 contains an example of many EIs obtained
 62 during such kind of experiment.

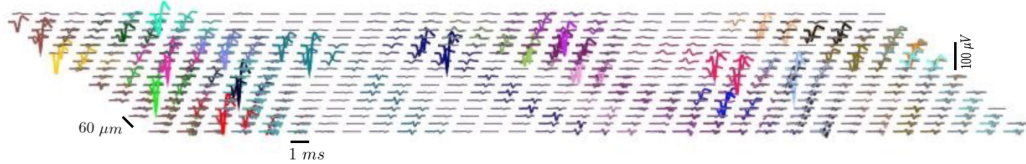


Figure 1: Overlapping EIs of 24 neurons (different colors) over a MEA, aligned to onset of spiking at $t = 0.5ms$. Each trace represent the time course of voltage at a certain electrode (for each neuron, traces are plotted only in the electrodes that contribute the most to the EI).

63 2.2 Stimulation Artifacts

64 Electrical stimulation experiments where neural responses are ablated (e.g. using the neurotoxin
 65 TTX) provide qualitative insights about the structure of the stimulation artifact $A(e, t, j, i)$; that is,
 66 how it varies as a function of all the relevant covariates: space (represented by electrode, e), time
 67 (t), amplitude of stimulus a_j and stimulus repetition i . The very first distinction to make is between
 68 stimulating and non-stimulating electrodes: magnitude of artifacts are much stronger in the former
 69 than in the latter, which is specially problematic for the recording of neural activity in response to
 70 somatic stimulation. However, because the modeling is more context-specific for this electrode we
 71 defer its treatment for the appendix and focus here on the non-stimulating electrodes, where the
 72 artifact has the follow properties: i) in the time component, its magnitude peaks following the onset
 73 of stimulus, and then stabilizes smoothly ii) in the spatial component, artifact decays smoothly with

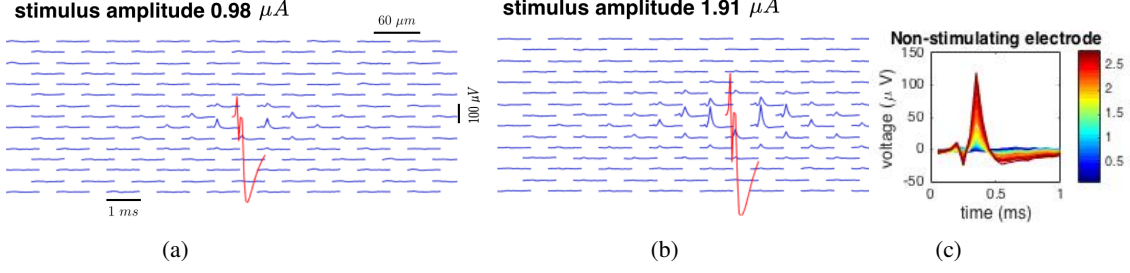


Figure 2: Illustration of main features of artifact in a TTX experiment. (a) and (b): artifacts on a neighborhood of the stimulating electrode (red). Each trace represent the time course of voltage at a certain electrode. (c) Artifact, as a function of amplitude of stimulus (different colors, in μA) for one of the nearest neighbors of the stimulating electrode.

distance to the stimulating electrode iii) in the stimulus component, magnitude of artifact increases smoothly with increasing stimulus strength iv) repetitions of stimulation lead to the same artifact up to random fluctuations. Then, omitting such trial-dependent random fluctuations, we can understand artifact as a stack of movies $A(e, t, j)$, one for each amplitude of stimulation a_j . Figure 2 depicts the main features of this artifact and illustrates why electrical stimulation thwarts the identification of neural activity.

We finish our description of the artifact with a useful metaphor that will help us think in the structure of the artifact and the problem in more concrete terms: one can think of the artifact as the wave that results from throwing a food pellet to a pond, where bigger pellets (stronger stimulus) leads to waves of larger magnitude. Our aim is to observe the fish (neurons) that come to eat the food, but this is difficult because the perturbations make the water blurry.

2.3 A scalable GP model for stimulation artifact

Here we propose a generative model for the artifact; namely, we model observed artifacts as draws from a GP, $A \sim GP(\mu, K)$. We choose the GP framework since it has shown to provide a parsimonious framework for the estimation of non-linear functions in the context of uncertainty, and because it leads to a straightforward implementation of all the usual operations that are relevant for our purposes (e.g. extrapolation) in terms of some conditional gaussian distributions. Now we address how to make a sensible choice of the parameters (μ, K) , that is, which parametric families are to be chosen and how to perform statistical inference in there.

Regarding μ , we follow the standard in the applied statistics community: μ is a centering parameter and all the non-random aspects of data should be captured by it. In our case this component is given by what we call the switching artifact, a waveform $A_0 = A_0(e, t)$ that is present regardless of the amplitude of stimulation. As a result, we can estimate $\hat{\mu}$ by taking the mean of recordings at the lowest amplitude of stimulation (see appendix for details) leading to a zero-mean Gaussian process.

The choice of K is more involved: it is well-known that the bottleneck in kernel methods lies in the inversion of K , which has complexity $O(d^3)$. In our context, $d = T \times E \times J \sim 10^6$, which renders the above inversion intractable unless further kernel structure is exploited. Here we implement a recent insight in the machine learning community: imposing a Kronecker product structure in K leads to tractable and scalable inferences. In detail, we assume the following decomposition

$$K = \rho K_t \otimes K_e \otimes K_s + \sigma^2 I_d. \quad (2)$$

where K_t , K_e and K_s are the kernels that account for variations in the time, space and stimulus dimensions of data, respectively. The dimensionless quantity ρ is used to control the overall magnitude of variability and the term $\sigma^2 I_d$ is included in acknowledgement that what is finally observed is a noise-corrupted version of the actual artifact.

Now this structured Kernel has been stated it remains to specify parametric families for the elementary kernels K_t, K_e, K_s . As shown below, we construct them from the Matérn family, using extra parameters to account for the non-stationarities described in 2.3.

110 2.3.1 A non-stationary family of kernels

111 Consider the Matérn(3/2) Kernel, the continuous version of an AR(2) process. It's (stationary)
 112 Covariance is given by

$$K_\lambda(x_1, x_2) = K_\lambda(\delta = |x_1 - x_2|) = \left(1 + \sqrt{3}\delta\lambda\right) \exp\left(-\sqrt{3}\delta\lambda\right). \quad (3)$$

113 Here $\lambda > 0$ represents the (inverse) lengthscale and determines how fast correlations decay with
 114 distance. We induce non-stationarities by introducing the family of unnormalized gamma densities
 115 $d_{\alpha,\beta}(\cdot)$:

$$d_{\alpha,\beta}(x) = \exp(-x\beta)x^\alpha. \quad (4)$$

116 By an appropriate choice of the pair $(\alpha, \beta) > 0$ we aim to expressively represent non-stationary
 117 'bumps' in variability. Then, to actually come up with a non-stationarity family we consider the
 118 process $Z \equiv Z(x) = d_{\alpha,\beta}(x)Y(x)$ where $Y \sim GP(0, K_\lambda)$. Then, Z is a *bona fide* GP¹ with the
 119 following covariance matrix ($D_{\alpha,\beta}$ is a diagonal matrix with entries $d_{\alpha,\beta}(\cdot)$):

$$K(\lambda, \alpha, \beta) = D_{\alpha,\beta}K_\lambda D_{\alpha,\beta}. \quad (5)$$

120 We choose all the kernels K_t, K_e, K_s as in (5): for the time and stimulus Kernel we use time and
 121 stimuli as covariates (δ in equation (3) and x in (4)), respectively. However, the case of the spatial
 122 Kernel is different: while δ represents distance between recording electrodes, now x represents
 123 distance between stimulating and recording electrodes, as it is this distance the one that covaries with
 124 artifact magnitude.

125 2.3.2 Kernel Learning

126 From equations (2), (3) and (5) our GP for the artifact is completely specified in terms of the
 127 parameters $\theta = (\rho, \alpha, \lambda, \beta)$ (one triplet (α, λ, β) for each of the space, time and stimulus kernels)
 128 and σ^2 . In other word, we encode all the peculiarities — lengthscales, orders of magnitudes, size of
 129 non-stationarities — of observed artifacts in different experiments through these parameters. Let's
 130 then call $K^\theta = \rho K_t \otimes K_e \otimes K_s$ and $K^{(\theta, \sigma^2)} = K^\theta + \sigma^2 I_d$. Our parameter learning problem is
 131 simply the maximization of the following function:

$$\max_{\theta} \log p(\tilde{A}|\theta, \sigma^2) = -\frac{1}{2} \tilde{A}^t \left(K^{(\theta, \sigma^2)}\right)^{-1} \tilde{A} - \frac{1}{2} \log \left|K^{(\theta, \sigma^2)}\right|. \quad (6)$$

132 Notice first that to avoid numerical instabilities we not include σ^2 in the optimization; alternatively
 133 one can come up with simple estimates of the background noise by using stimulus-free data. Also,
 134 here \tilde{A} represents some guess of the artifact: naturally the artifact itself is not revealed (that would
 135 solve the problem) but as we are only concerned with parameters that inform about the global and
 136 overall properties of the artifact, using a guess instead would not lead to big changes in θ . In practice,
 137 we use data corrupted by neural activity, that is $\tilde{A} \equiv Y$. This is a sensible decision as neural activity
 138 is sparse and small compared to artifact, and does indeed lead to results that are very close to guesses
 139 \tilde{A} from TTX experiments (where available). Figure 3 shows an example of learned Kernels K_t, K_e
 140 and K_s using this approach.

141 2.4 Algorithm

142 We start by fully specifying the different aspects in equation (1), where we aim to produce estimates
 143 \hat{A} and \hat{s} . We adopt the convention that whenever indexes are specified we refer to the subvector of
 144 the huge vector (either Y, s, A)² that created by fixing that component (e.g. $A(t, e)$ is the artifact at
 145 time t and electrode e moving freely along different stimulus amplitudes). With these convention, to
 146 clearly state trial-wise dependencies, we re-express equation (1) as:

$$Y(t, e, j, i) = A(t, e, j) + s(t, e, j, i) + \epsilon(t, e, j, i) \quad t \leq T, e \leq E, j \leq J, i \leq n_j. \quad (7)$$

¹Indeed, the above covariance can be thought as the Hadamard or pointwise product between K_λ and the degenerate Kernel dd^t with d being the diagonal of $D_{\alpha,\beta}$.

²The same convention also applies for submatrices of K^θ or $K^{(\theta, \sigma^2)}$.

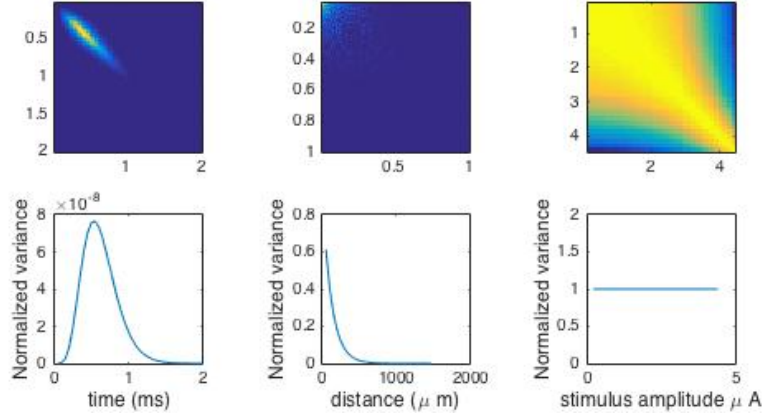


Figure 3: Example of kernels K_t , K_e and K_s learned from real data with $\tilde{A} = Y$

Noise is assumed uncorrelated gaussian with variance σ^2 , that is, $\epsilon \sim \mathcal{N}(0, \sigma^2 I_{d'})$ with $d' = T \times E \times J \times \sum_{i=1}^J n_j$. This assumption is rather restrictive and we assume it here for computational ease, but defer the reader to the supplementary material for a sketch of a formulation that takes into account correlated noise.

In agreement with 2.3 we take $A \sim GP(0, K^\theta)$. This means we conceptualize A as a noise-free artifact and then $A + \epsilon$ can be thought as its noise corrupted version, with covariance $K^{(\theta, \sigma^2)}$.

We represent neural activity s as the superimposed activity of the N neurons involved, each represented by an appropriate convolution of the EI. In detail, we first consider binary vectors b^n that indicate our constraint on spike occurrence and timing: for each trial and neuron there can be at most one spike, which can only occur on a pre-defined discrete set of times (a strict subset of the entire recording window to prevent aliasing). This is the case in reality, as the recording window is comparable with the refractory period of neurons. Then, the actual activities s^n are obtained through products with the Toeplitz matrices M^n that contain copies of the EIs with spike onset aligned at different times. All of this is summarized in equation (8):³

$$s = \sum_{n=1}^N s^n, s^n = M^n b^n. \quad (8)$$

Now we turn to the parameter estimation problem, based on posterior inference for

$$p(A, s | Y, \theta, \sigma^2) \propto p(Y | s, A, \sigma^2) p(A | \theta). \quad (9)$$

The non-convexity of the set over which the binary vectors b_n are defined makes this problem difficult as many local optima do exist in practice and, as a result, for global optimization there may not be a better alternative than to look at a huge number of possible cases. We circumvent this cumbersome global optimization by taking a principled greedy approach, with two main characteristics i) joint optimization over A and s is addressed with alternate ascent and ii) data is divided in batches corresponding to the same stimulus amplitude, and the analysis for the $j + 1$ -th batch starts only after definite estimates $\hat{s}_{1, \dots, j}$ and $\hat{A}_{[j]}$ ⁴ have already been produced. Moreover, this latter estimate of the artifact is used to produce an initial estimate for A_{j+1} . Intuitively, we borrow strength from lower amplitude of stimulation conditions to counteract the worsened effect of artifacts and increased responses of neurons due to high amplitudes of stimulation.

Now we describe in detail these features, which finally lead to the algorithm that is sketched at the end of this section.

³Note the slight abuse of notation: while the above definitions make sense trial-wise, when considering the entire dataset, actually one deals with many matrices M^n concatenated vertically.

⁴ $[j]$ denotes the set $\{1, \dots, j\}$

2.4.1 Coordinate Ascent

Given the batch Y_j and an initial artifact estimate A_j^0 (see 2.4.2) we alternate between neural activity estimation given a current artifact estimate, and artifact estimation \hat{A}_j given the current estimate of neural activity. This alternate optimization stops when changes in \hat{s}_j are sufficiently small, or nonexistent.

Matching pursuit for spike maximization:

Given the current artifact estimate \hat{A}_j we maximize the conditional for neural activity, which turns out to be the following sparse regression problem:

$$\min_{b_j^n \in S, n=1, \dots, N} \|(Y_j - \hat{A}_j) - \sum_{n=1}^n M^n b_j^n\|^2. \quad (10)$$

The set S embodies our constraints on spike occurrence and timing. Intuitively, we pursue to place spikes so that they will match with the residuals $(Y_j - \hat{A}_j)$ the best; that is, will lead to the smallest sum of squares. For computational ease, we proceed greedily: by concatenating all the M^n and b^n we can find the neuron and spike time that decreases the above residuals the most, and after extracting the corresponding (b^n, M^n) from (b, M) we can continue iteratively until no spike placement of the remaining neurons leads to further decrease in the sum of squares.

Bayesian regression for spike maximization:

Given the current estimate of neural activity \hat{s}_j we maximize the conditional of the artifact, that is:

$$\max_{A_j} p(A_j | Y_j, \hat{s}_j, \theta, \sigma^2), \quad (11)$$

which — by well known properties of the gaussian distribution — leads to the posterior mean estimator (the overline indicates mean across the n_j stimulation repetitions):

$$\hat{A}_j = E(A_j | Y_j, \hat{s}_j, \theta, \sigma^2) = K_{j,j}^{\theta} \left(K_{j,j}^{(\theta, \sigma^2)} \right)^{-1} (\bar{Y}_j - \bar{\hat{s}}_j). \quad (12)$$

2.4.2 Iteration over batches and artifact extrapolation

The procedure described in 2.4.1 is repeated in a for loop that iterates through the batches corresponding to different stimulus strengths, from the lowest to the highest. Also, when doing $j \rightarrow j+1$ an initial estimate for the artifact A_{j+1}^0 is generated by extrapolating from the current, faithful, estimate of the artifact up to the j -th batch. This extrapolation is easily implemented as the mean of the noise-free posterior distribution in this GP setup, that is:

$$A_{j+1}^0 = E(A_{j+1} | \hat{A}_{[j]} \theta, \sigma^2) = K_{(j+1),[j]}^{\theta} \left(K_{([j],[j])}^{(\theta, \sigma^2)} \right)^{-1} \hat{A}_{[j]}. \quad (13)$$

Importantly, in practice this initial estimate ends up being extremely useful: without a proper estimate of the artifact, coordinate ascent will often lead to poor optima, but the accurate estimates provided by this extrapolation prevent that from being the case.

Algorithm: Spike sorting with electrical stimulation artifacts.

Input: traces $Y = (Y_j)_{j=1, \dots, J}$, in response to J stimulus. EIs of N neurons.

Output: estimates of artifact and neural activity (spikes) for each neuron.

Initialization: estimate σ^2 and θ as in (6). (Can re-use these parameter for data from different stimulating electrodes)

Then: repeat for $j=1, \dots, J$

Construct A_j^0 from $[j-1]$ using equation (13) ($A_1^0 \equiv 0$)

repeat until little or no changes in binary vectors b_j^n

- Estimate \hat{s}_j by greedily finding spikes as in (10) until no spike of no neuron leads to increase in likelihood.
 - Estimate \hat{A}_j using (12) to increase in likelihood.
-

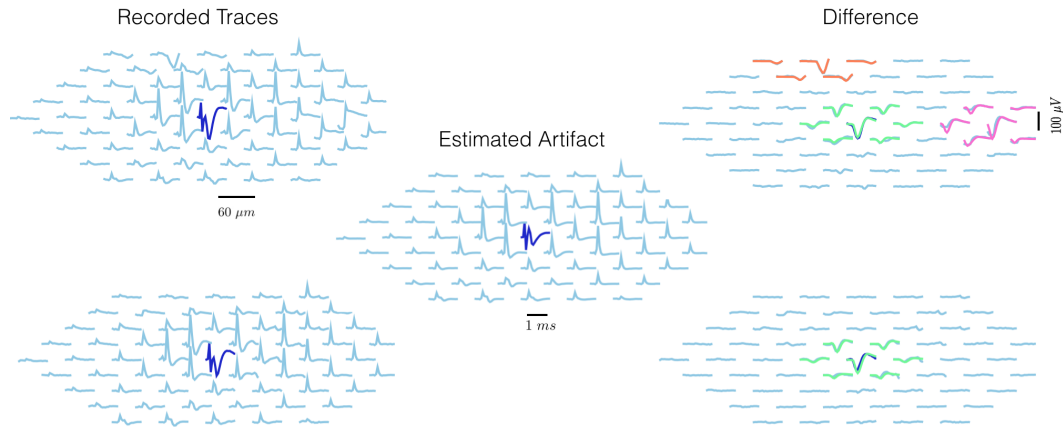


Figure 4

3 Results

4 Appendix

4.1 Stimulating Electrode

Second, regarding the stimulating electrode: i) generally speaking, magnitude of artifact is much greater than in the non-stimulating electrodes ii) variations in time also peak right after the onset of stimulation but oscillations are wilder and iii) artifact magnitude increases with stimulus strength smoothness is lost: due to stimulation hardware idiosyncrasies, there might be some breakpoint stimulus strengths such that artifact shape and magnitude will ostensibly change after such strengths are surpassed.

Acknowledgments

Use unnumbered third level headings for the acknowledgments. All acknowledgments go at the end of the paper. Do not include acknowledgments in the anonymized submission, only in the final paper.

References

References follow the acknowledgments. Use unnumbered first-level heading for the references. Any choice of citation style is acceptable as long as you are consistent. It is permissible to reduce the font size to small (9 point) when listing the references. **Remember that you can use a ninth page as long as it contains only cited references.**

- [1] Alexander, J.A. & Mozer, M.C. (1995) Template-based algorithms for connectionist rule extraction. In G. Tesauero, D.S. Touretzky and T.K. Leen (eds.), *Advances in Neural Information Processing Systems 7*, pp. 609–616. Cambridge, MA: MIT Press.
- [2] Bower, J.M. & Beeman, D. (1995) *The Book of GENESIS: Exploring Realistic Neural Models with the GEneral NEural Simulation System*. New York: TELOS/Springer-Verlag.
- [3] Hasselmo, M.E., Schnell, E. & Barkai, E. (1995) Dynamics of learning and recall at excitatory recurrent synapses and cholinergic modulation in rat hippocampal region CA3. *Journal of Neuroscience* **15**(7):5249-5262.