

Sampling-based probabilistic inference through neural and synaptic dynamics

Wolfgang Maass for

Robert Legenstein

Institute for Theoretical Computer Science

Graz University of Technology, Austria

I will address two time scales of sampling in neural networks:

Part I: Neural Sampling

Time scale of spikes up to short term dynamics of synapses and neurons

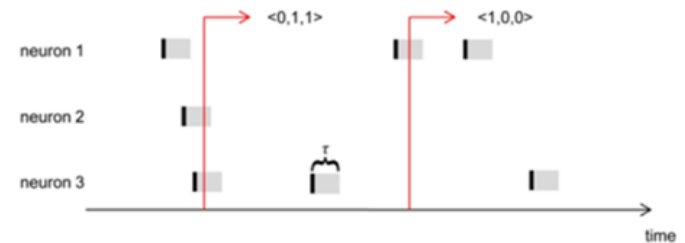
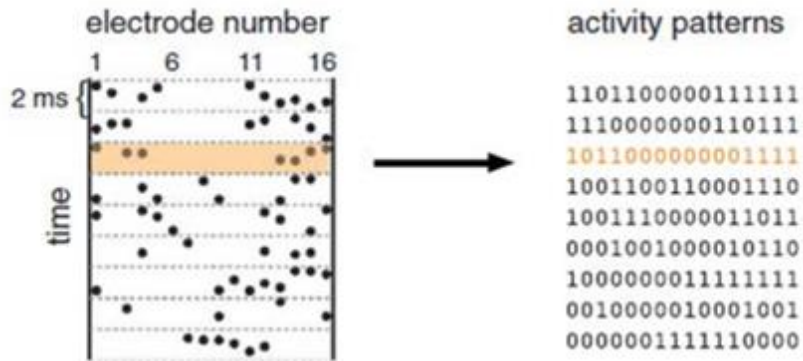
Part II: Synaptic Sampling

Time scale of synaptic plasticity and rewiring

Question: On which of these two time scales can (and does) the brain implement probabilistic inference through sampling ?

Part I: Neural Sampling

Inspiration from data: Berkes et al. recorded the distribution of „network states“ with 16 electrodes in area V1 of ferrets:



Their interpretation:

This distribution of network states becomes during development an internal model for visual inputs

P. Berkes, G. Orban, M. Lengyel, and J. Fiser. Spontaneous cortical activity reveals hallmarks of an optimal internal model of the environment. *Science*, 331:83-87, 2011

Such stochastic internal models in brain networks have in fact been proposed by many researchers

- K. Friston, “A theory of cortical responses, 2005
- E. Vul and H. Pashler, “Measuring the crowd within”, 2008
- S. Denison, E.B. Bonawitz, A. Gopnik, T.L. Griffiths, “Preschoolers sample from probability distributions”, 2010,

Some of this work suggests furthermore, that the brain carries out probabilistic inference from distributions p through MCMC sampling from p .

.

Assume for example that $p_C = p(z_1, \dots, z_K)$ is the stationary distribution over K binary random variables, each represented by one neuron, of a brain network C .

If concrete values \mathbf{e} („evidence“) are plugged in for some of these variables then the posterior marginal $p(z_1|\mathbf{e}) = \sum_{v_2, \dots, v_m} p(z_1, v_2, \dots, v_m|\mathbf{e})$ can be estimated by

observing the firing rate of the neuron v_1 that is associated with the binary random variable z_1 .

These studies suggest the analysis of stationary distributions of neural networks

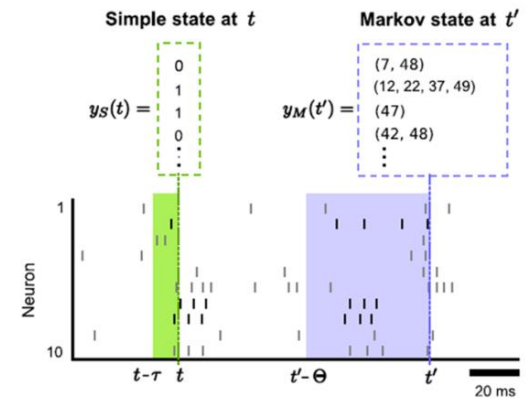
Stumbling block for theory:

The Markov chains that are defined by networks of spiking neurons are **nonreversible**, both because of their inherent dynamics (a spike causes temporally extended changes in other neurons), and because of non-symmetric synaptic connections.

General theoretical result

Theorem: Virtually any network C of spiking neurons with noise has a unique stationary distribution p_C of network states, to which it converges exponentially fast from any initial state.

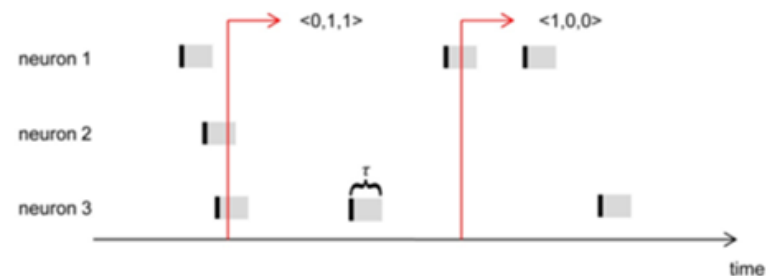
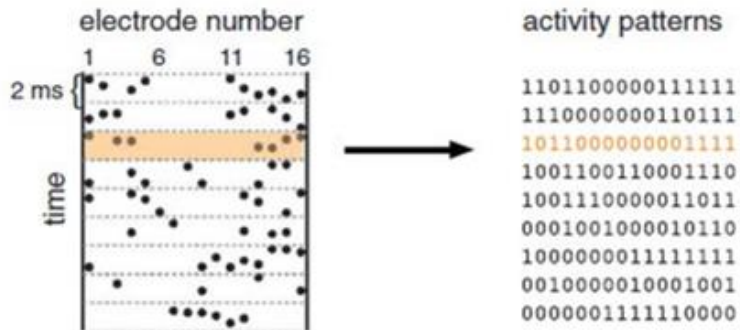
Two relevant notions
of network state:



- A biological neural network can only be viewed as a MC if one moves to temporally extended notions of network state
- These are MCs with **continuous time and continuous state spaces**
- This Theorem is one of few that also hold for data-based nonlinear models of neurons and synapses.

What types of probability distributions can arise as stationary distribution of a network of spiking neurons?

I am considering here the same convention as Berkes et al. for relating spikes to bits:



We cannot use the common detailed balance criterion for verifying that a particular distribution p is the stationary distribution, but the Neural Computability Condition (NCC) is a replacement

Theorem: The **Neural Computability Condition (NCC)** provides a sufficient condition for representing p through a network of spiking neurons:

For each RV z_k there is some neuron ν_k whose membrane potential at time t is

$$u_k(t) = \log \left(\frac{p(z_k = 1 | z_{\setminus k}(t))}{p(z_k = 0 | z_{\setminus k}(t))} \right)$$

Note:

This result holds rigorously only for an idealized type of spiking neuron with firing probability

$$\rho_k(t) = \frac{1}{\tau} \exp(u_k(t))$$

and step functions as PSPs

But numerical simulations suggest that the error is not large for biologically more realistic models.

Neural Computability Condition (NCC)

For p with only **2nd-order dependencies**, i.e.: **Boltzmann distributions**

$$p(z) = \frac{1}{Z} \exp\left(\sum_{ij} \frac{1}{2} W_{ij} z_i z_j + \sum_i b_i z_i\right)$$

A network of spiking neurons with symmetric weights automatically satisfies the NCC.

For p with **arbitrary higher-order dependencies**:

Such p can be represented directly via networks of spiking neurons with **asymmetric** connections.

In fact, such networks can **learn** to approximate any discrete distribution p from examples (drawn from p).

L. Büsing, J. Bill, B. Nessler, and W. Maass. PLOS Comp. Biol. 2011

Neural sampling is structurally different from Gibbs sampling, also for this case !

(Jonke, Habenschuss, Maass, Arxiv 2015)

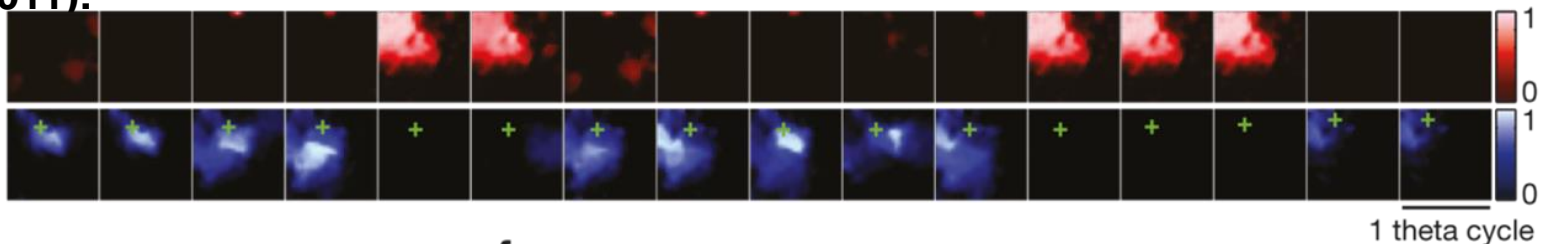
D. Pecevski, L. Büsing, W. Maass, PLOS Comp. Biol., 2011

D. Pecevski, W. Maass, 2015 (under review)

Some open questions about neural sampling

- What is the speed (rhythm) of neural sampling in the brain?

Apparently fastest known sampling-like dynamics (Jezek et al, Nature 2011):



Blue and red colors indicate strong correlation with the place map of the blue or red maze („context“), Hardly any theta cycles with mixed „contexts“ were found

- To what extent does the brain use results of neural sampling for decision making etc?
- Can stationary distributions be learnt from examples (that are generated by some external distribution p^*) by networks of spiking neurons?
- How are behaviourally relevant random variables encoded?

Part II: Synaptic Sampling

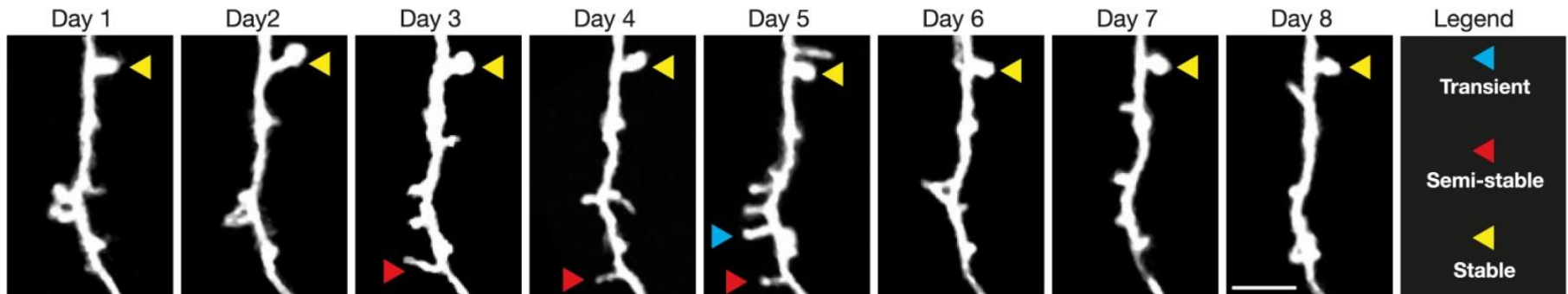
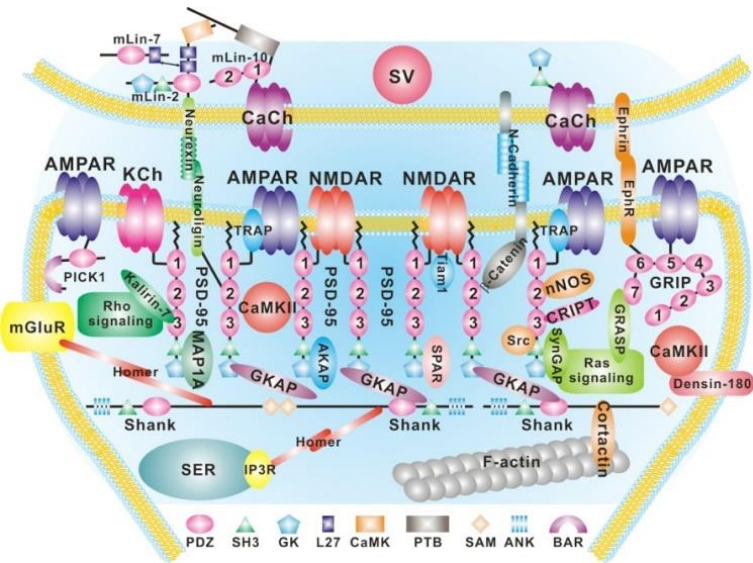
Many biological network parameters are fluctuating (more or less) all the time:

A postsynaptic density consists of over 1000 different types of proteins, many in small numbers.

Since these molecules have a lifetime of only weeks or months, their number is subject to permanent stochastic fluctuations.

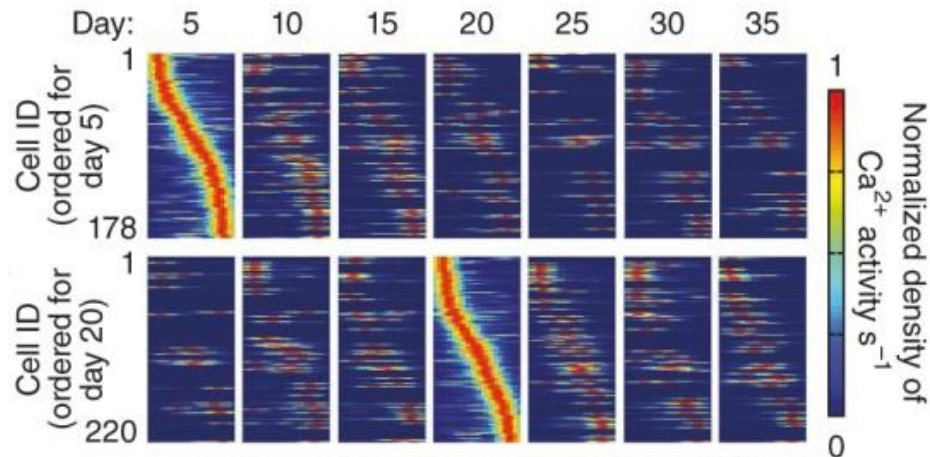
Receptors etc. are subject to Brownian motion within the membrane.

Furthermore axons sprout and dendritic spines come and go on a time scale of days (even in adult cortex, perhaps even in the absence of neural activity)



Data from Svoboda Lab

Longterm recordings show that *neural codes drift* on the time-scale of weeks and months



Ziv, Y., Burns, L. D., Cocker, E. D., Hamel, E. O., Ghosh, K. K., Kitch, L. J., ... & Schnitzer, M. J.. Long-term dynamics of CA1 hippocampal place codes. *Nature Neuroscience*, 2013

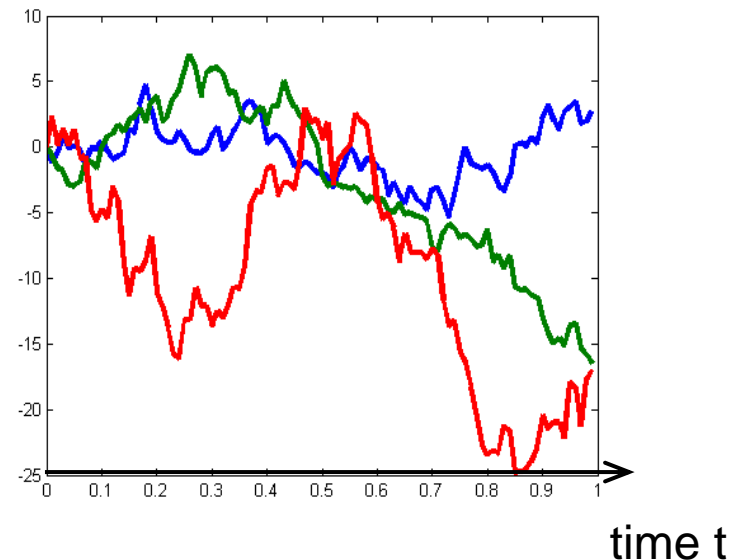
See also:

Rokni, U., Richardson, A. G., Bizzi, E., & Seung. Motor learning with unstable neural representations. *Neuron*, 2007

and forthcoming new data.

Mathematical framework for capturing these phenomena: „Synaptic Sampling“

- We model the evolution of network parameters through Stochastic Differential Equations (SDEs): $d\theta_i = \underbrace{\left(b \frac{\partial}{\partial \theta_i} \log p^*(\boldsymbol{\theta})\right)}_{\text{drift}} dt + \underbrace{\sqrt{2Tb} \cdot d\mathcal{W}_i}_{\text{diffusion}}$
- The diffusion term $d\mathcal{W}_i$ in the SDE denotes an infinitesimal step of a random walk („Brownian motion), whose temporal evolution from time \mathbf{s} to time \mathbf{t} satisfies $\mathcal{W}_i^t - \mathcal{W}_i^s \sim \text{NORMAL}(\mathbf{0}, t - s)$.



Mathematical framework for capturing these phenomena: „Synaptic Sampling“

Integration of this SDE yields infinitely many different solutions for the evolution of the parameters.

But the **evolution of their probability density** is given by a deterministic PDE (*Fokker-Planck equation*):

$$\frac{\partial}{\partial t} p_{FP}(\boldsymbol{\theta}, t) = \sum_i -\frac{\partial}{\partial \theta_i} \left(\left(b \frac{\partial}{\partial \theta_i} \log p^*(\theta_i | \mathbf{x}, \boldsymbol{\theta}_{\setminus i}) \right) p_{FP}(\boldsymbol{\theta}, t) \right) + \frac{\partial^2}{\partial \theta_i^2} (Tb p_{FP}(\boldsymbol{\theta}, t))$$

By setting the left-hand side to 0, this FP-equation makes the resulting **stationary distribution** $\frac{1}{Z} p^*(\boldsymbol{\theta})^{\frac{1}{T}}$ for the vector $\boldsymbol{\theta}$ of all network parameters θ_i explicit.

Implication: One can **program into stochastic plasticity rules**

$$d\theta_i = \left(b \frac{\partial}{\partial \theta_i} \log p^*(\boldsymbol{\theta}) \right) dt + \sqrt{2Tb} \cdot d\mathcal{W}_i$$

desired target distributions $\frac{1}{Z} p^*(\boldsymbol{\theta})^{\frac{1}{T}}$ of the parameters.

Hence synaptic sampling can implement sampling from a posterior distribution of network parameters

synaptic sampling
with prior $p_S(\boldsymbol{\theta})$

unsupervised learning (generative models)

$$p^*(\boldsymbol{\theta}|\mathbf{x}) \propto p_S(\boldsymbol{\theta}) p_{\mathcal{N}}(\mathbf{x}|\boldsymbol{\theta})$$

where

- \mathbf{x} are repeatedly occurring network inputs
- $p_{\mathcal{N}}(\mathbf{x}|\boldsymbol{\theta})$ is the generative model provided by a neural network \mathcal{N} with parameters $\boldsymbol{\theta}$

reinforcement learning

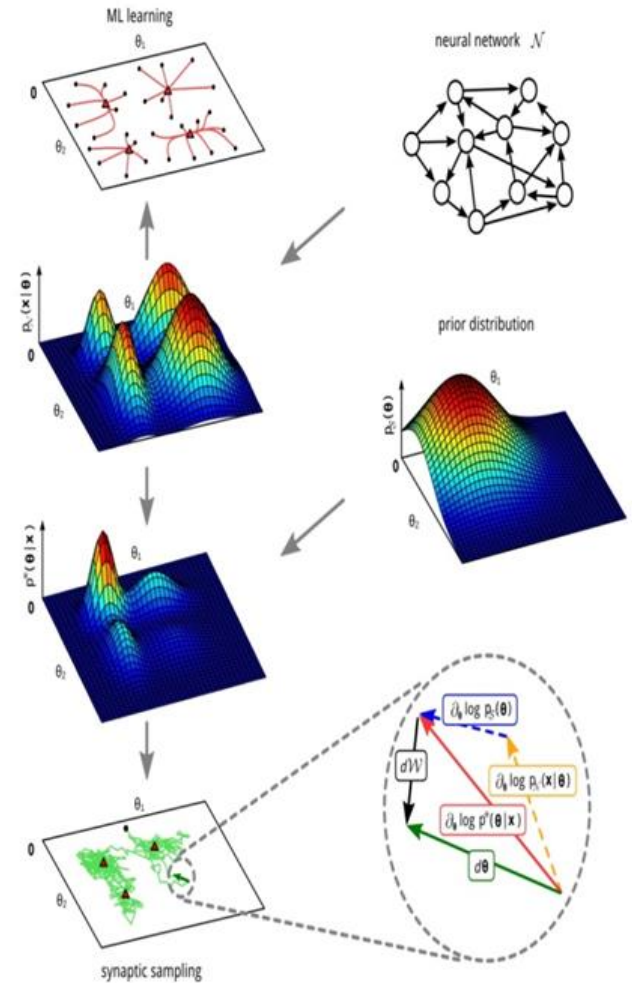
$$p^*(\boldsymbol{\theta}) \propto p_S(\boldsymbol{\theta}) \cdot p_{\mathcal{N}}(R = \max|\boldsymbol{\theta})$$

where R signals reward

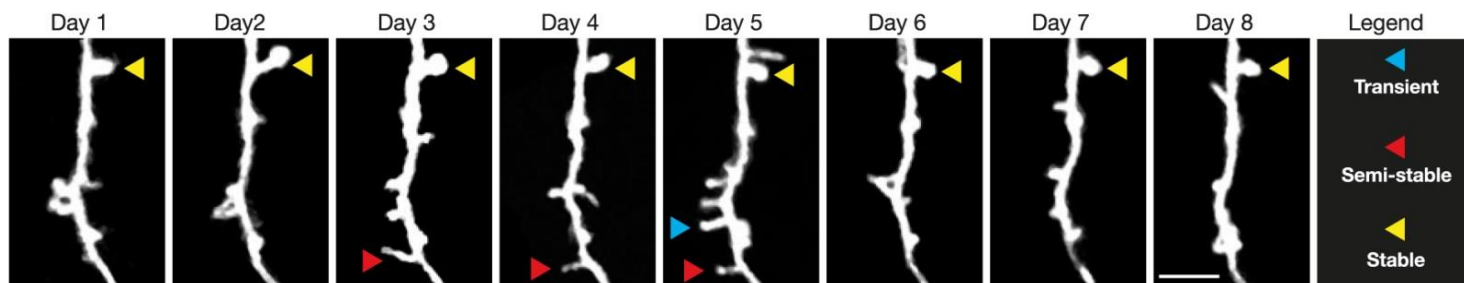
This integrates policy gradient RL with probabilistic inference.

Functional implications of synaptic sampling for learning

1. Better generalization (predicted by MacKay, 1992)
2. Structural plasticity can easily be integrated with synaptic plasticity in a principled manner
3. Automatic self-repair capabilities of the network (without requiring a clever supervisor that switches learning back on after a perturbation)

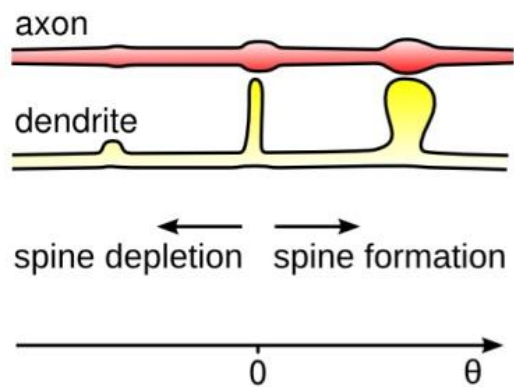


Spine dynamics and synaptic plasticity can easily be integrated into a SDE for a parameter that regulates both



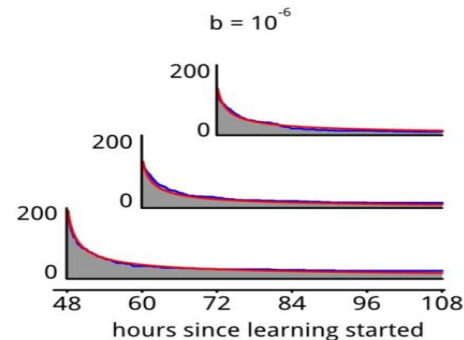
Ansatz: A single parameter θ_i controls the spine volume and – once a synaptic connection has been formed – the weight of this synaptic connection.

parametrization of spine motility through the parameter θ



Not only **STDP**, but also **experimentally observed power-law survival curves for synaptic connections** are reproduced by this combined rule:

Experimental data from (Löwenstein, Kuras, Ruml, J. of Neuroscience, in press)

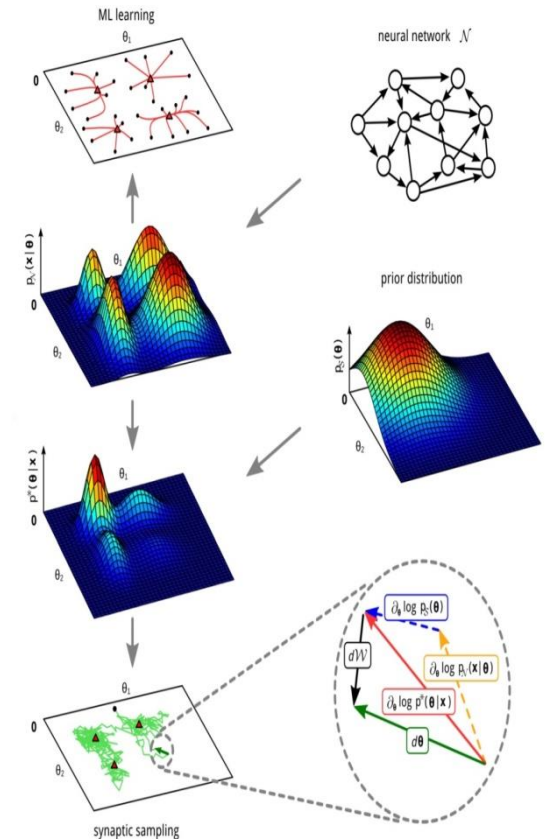


- some slides with unpublished results are deleted

Summary of part II: Synaptic sampling

One arrives at a new understanding of network learning

- The parameter vector permanently wanders around (with varying speed) in some very high-dimensional space of parameter values θ
- It spends most of its time in a low-dimensional sub-manifold where both the prior and the performance (likelihood of inputs, or reward probability) are high
- Hence changes in the input distribution or lesions are no big deal, since the parameter vector θ does not have a „permanent home“ anyway
- Priors enable the network to combine experience dependent learning with structural rules in a theoretically optimal way (**Bayesian inference**)



Credits:



Robert Legenstein



David Kappel



Stefan Habenschuss



Dejan Pecevski



Zeno Jonke



Lars Buesing