

Computation with Synapses

吴思

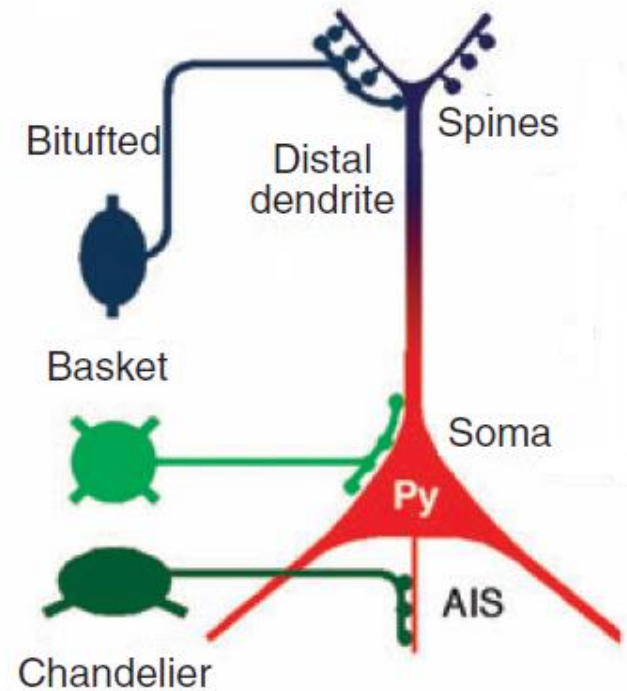
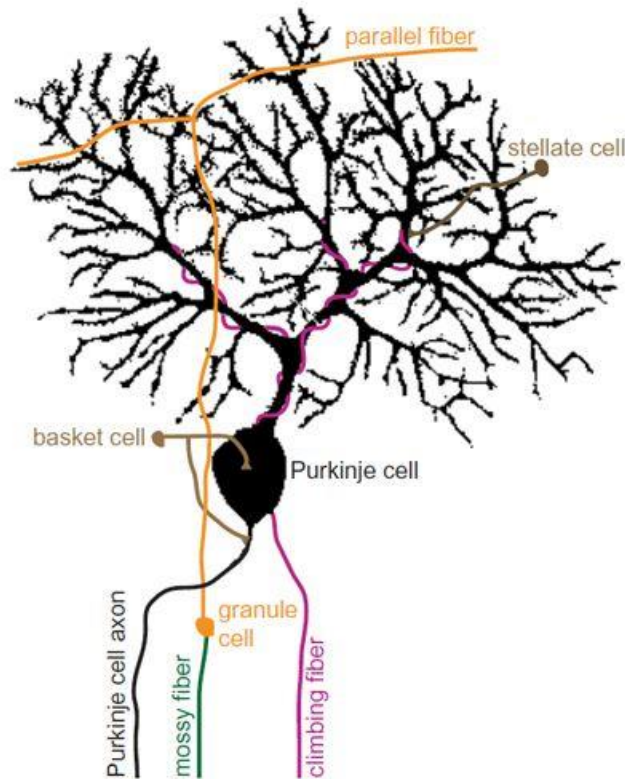
心理与认知科学学院

IDG/McGovern 脑科学研究所

北大-清华联合生命科学中心

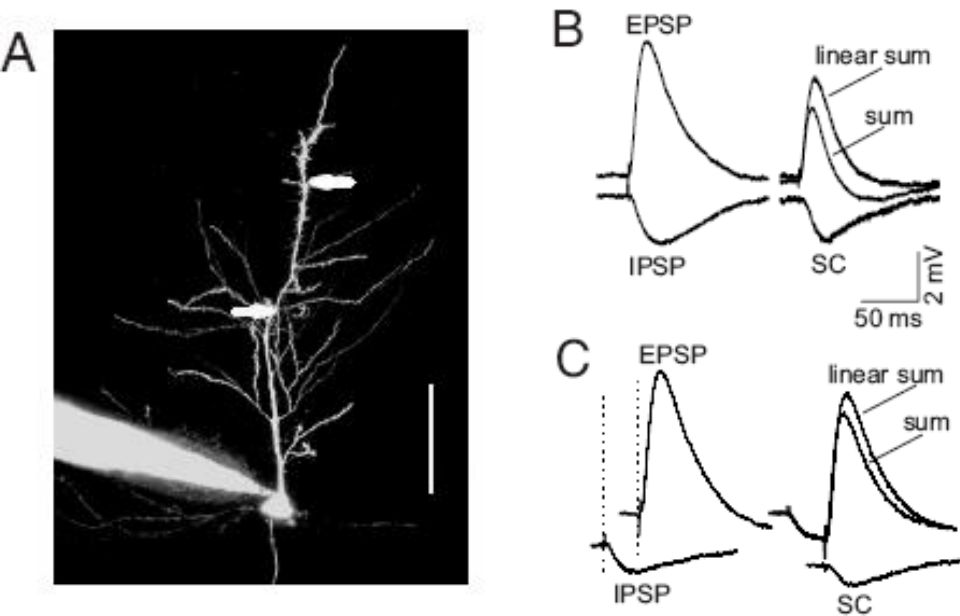
北京大学

Rich dendrites and inhibitory neurons



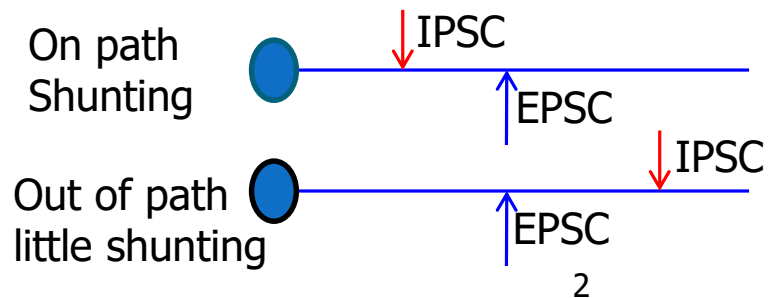
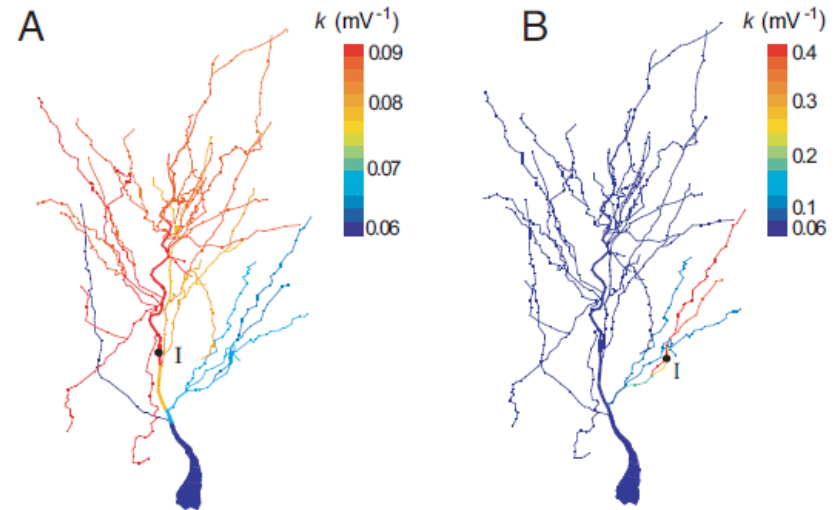
Shunting Inhibition

Measured sum of EPSP and IPSP



$$\text{Sum} = \text{EPSP} + \text{IPSP} + k * \text{EPSP} * \text{IPSP}$$

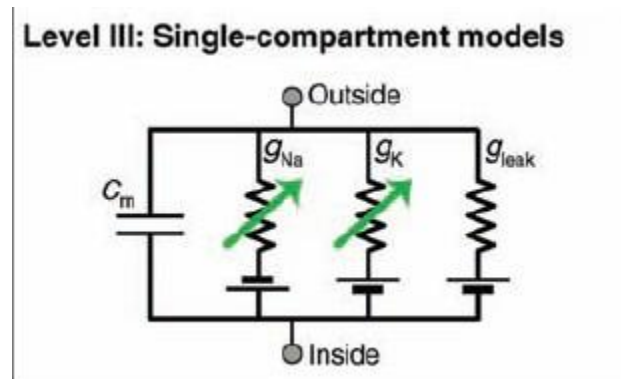
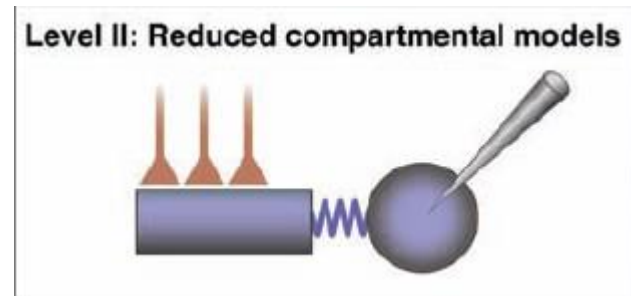
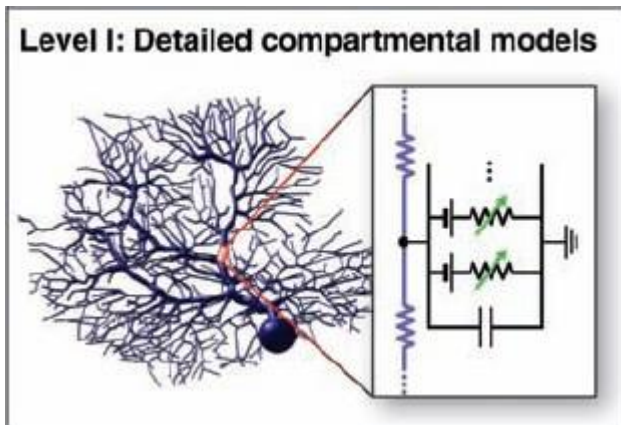
Location dependence of shunting inhibition



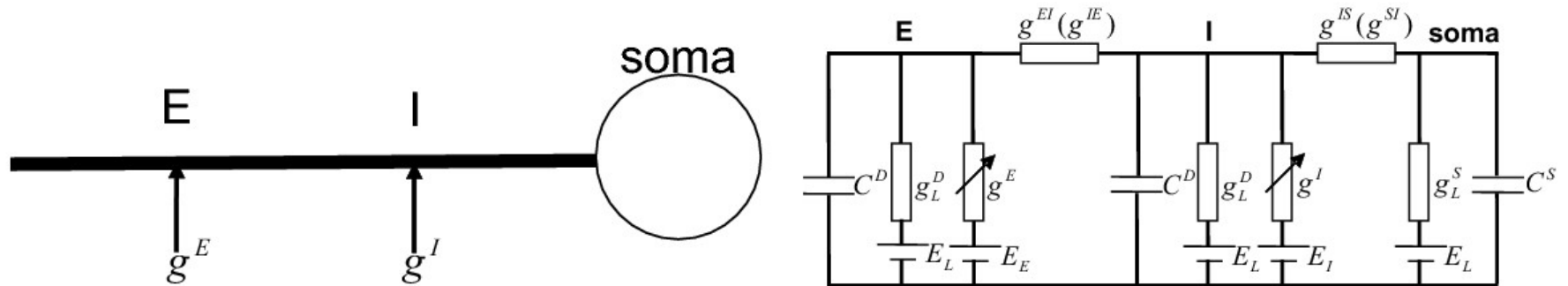
Levels of simplification

Simple models:

1. Capture the fundamental features of a systems related to brain function;
2. Allow us for analytical treatments.



A Neuron Model with Soma and a Dendrite



Dynamics of the conductance-based model with passive channels:

$$C^S \dot{v}^S = -g_L^S (v^S - E_L) - g^{IS} (v^S - v^I),$$

$$C^D \dot{v}^I = -g_L^D (v^I - E_L) - g^{SI} (v^I - v^S) - g^{EI} (v^I - v^E) - g^I (v^I - E_I),$$

$$C^D \dot{v}^E = -g_L^D (v^E - E_L) - g^{IE} (v^E - v^I) - g^E (v^E - E_E),$$

A Simplified Model

Fast and slow dynamics

$$\tau_s = C^S / (g_L^S + g^{IS}) \qquad C^S \gg C^D,$$

$$\tau_I = C^D / (g_L^D + g^{SI} + g^{EI}) \qquad g_L^S \gg g^{IS}, \quad g_L^D \gg g^{SI} \text{ and } g_L^D \gg g^{IE}$$

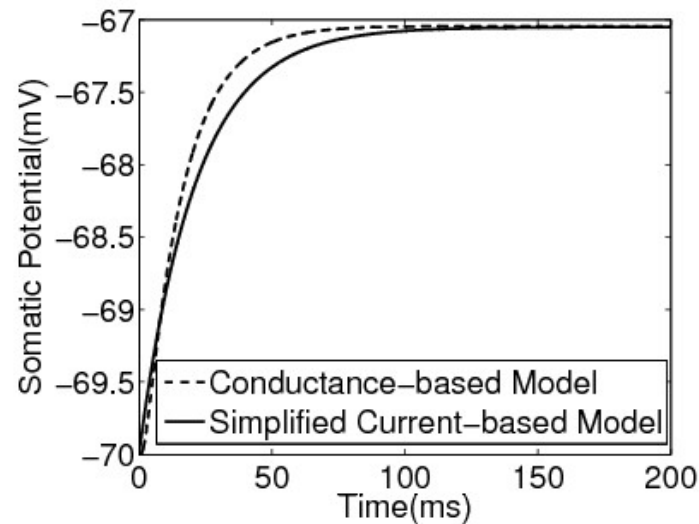
$$\tau_E = C^D / (g_L^D + g^{IE})$$

$$\tau_s \dot{v}^S = -(v^S - E_L) + f_1(g^E) + f_2(g^I) + k f_1(g^E) f_2(g^I),$$

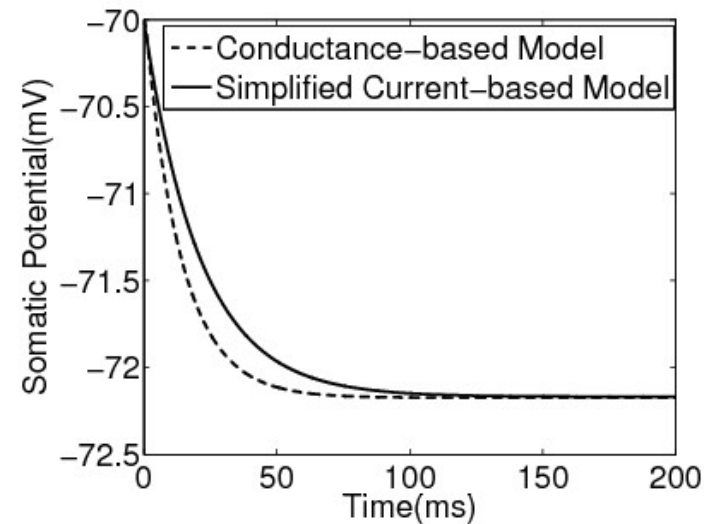
where,

$$\begin{aligned} f_1(g^E) &= \frac{\alpha^2 g^{SI} g^{IE} g^E (E_E - E_L)}{(g_L^S g_L^D + g_L^S g^{SI} + \alpha g_L^S g^{IE} + \alpha g_L^D g^{SI} + \alpha^2 g^{IE} g^{SI})(g_L^D + g^E + g^{IE})}, \\ f_2(g^I) &= \frac{\alpha g^{SI} g^I (E_I - E_L)}{(g_L^S g_L^D + g_L^S g^{SI} + \alpha g_L^D g^{SI}) + g^I (g_L^S + \alpha g^{SI}) + g^{IE} (\alpha g_L^S + \alpha^2 g^{SI})}, \\ k &= \frac{g_L^S + \alpha g^{SI}}{\alpha g^{SI} (E_L - E_I)}. \end{aligned}$$

The Goodness of Simplification



(a)



(b)

Figure 2: The conductance-based multi-compartment model vs. the simplified current-based one. (a) The excitatory input dominates, $g^E = 40nS$ and $g^I = 5nS$; (b) The inhibitory input dominates, $g^E = 5nS$ and $g^I = 40nS$. The other parameters are: $C^S = 740pF$, $C^D = 50pF$, $g_L^S = 30nS$, $g_L^D = 20nS$, $g^{SI} = 5nS$, $g^{IE} = 1nS$, $\alpha = 5$; $E_L = -70mV$, $E_E = 10mV$ and $E_I = -80mV$.

Theoretical justification of experimental results

1. The shunting strength have a larger value for inhibitory input located at the distal side of a dendrite than at the proximal side.

$$k = \frac{g_L^S + \alpha g^{SI}}{\alpha g^{SI}(E_L - E_I)}.$$

g^{SI} is smaller for inhibitory input located at the distal side of a dendrite.

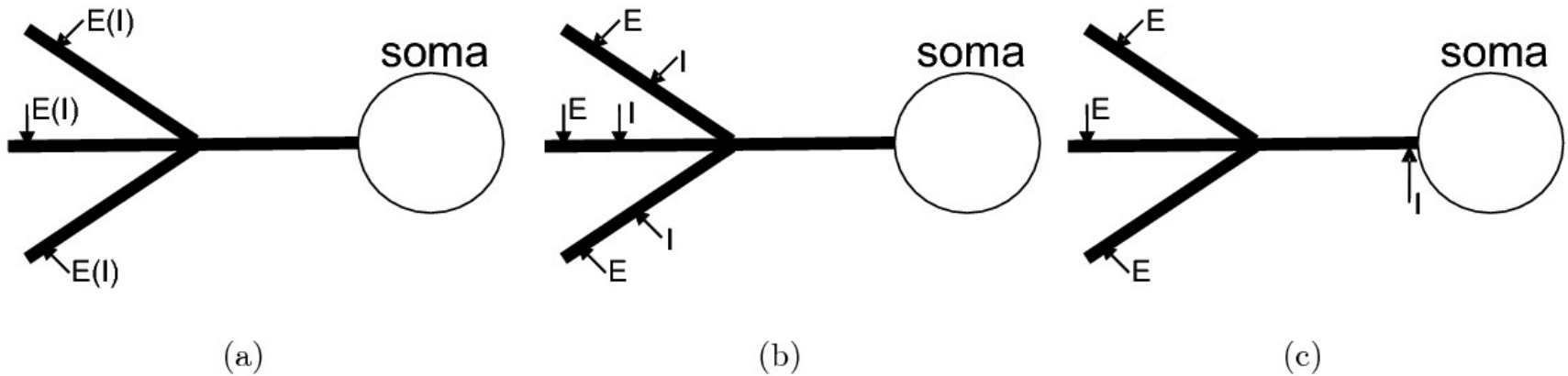
2. The shunting strength is much smaller for the out-of-path configuration.

$$k = \frac{g_L^S + \alpha g^{SI}}{\alpha g^{SI}(E_L - E_I)} \quad \text{vs.} \quad k = \frac{g_L^S + \alpha g^{SI}}{\alpha g^{SI}(E_L - E_E)}$$

On the path

Out of path

Extensions to multiple dendrites

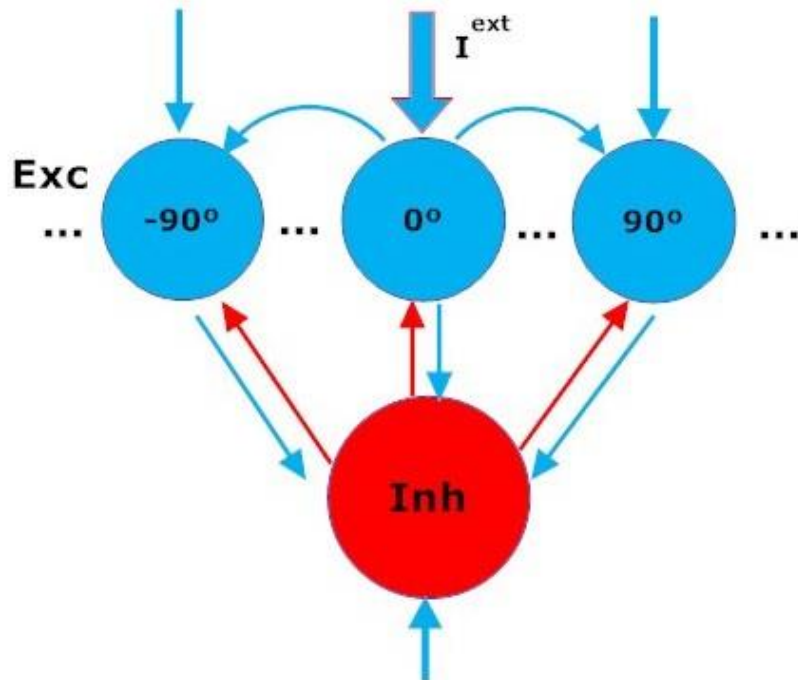


$$(a). \tau_s \dot{v}^S = -v^S + E_L + \sum_i f(g_i^{syn})$$

$$(b). \tau_s \dot{v}^S = -v^S + E_L + \sum_i [f_1(g_i^E) + f_2(g_i^I) + k f_1(g_i^E) f_2(g_i^I)]$$

$$(c). \tau_s \dot{v}^S = -v^S + E_L + \sum_i f_1(g_i^E) + f_2(g^I) + k f_2(g^I) \sum_i f_1(g_i^E)$$

Global Shunting Inhibition Realizes Divisive Normalization



Global Shunting Inhibition

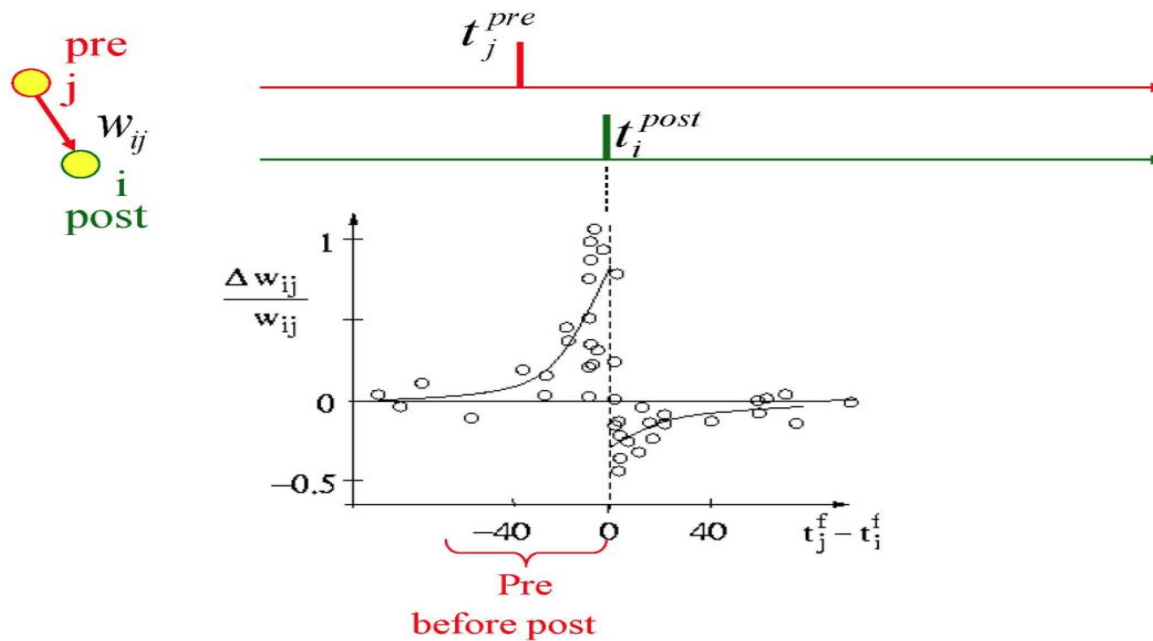
$$\text{Inh: } R^I \propto \sum_i R_i^E$$

$$\text{Exc: } R_i^E \propto I^{ext} + R^I - k I^{ext} R^I$$

$$\square \frac{I^{ext}}{1 + k \sum_i R_i^E}$$

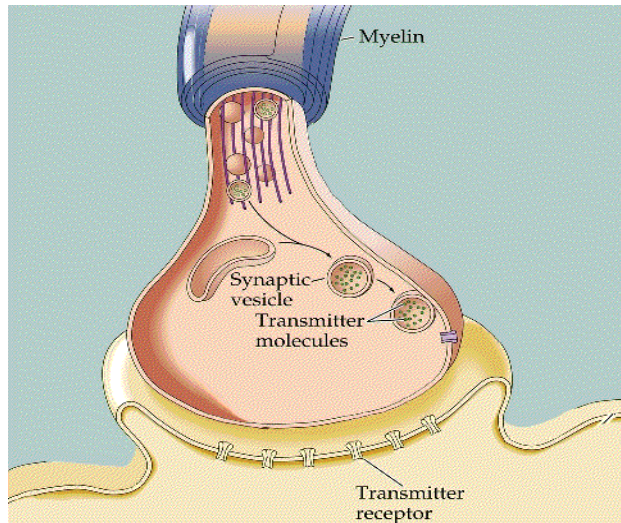
Divisive Normalization

Spike-time-dependent-plasticity (STDP)

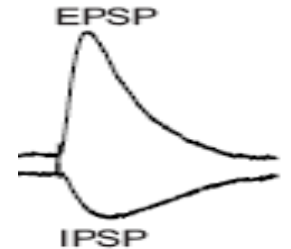


Bi & Poo, J. Neurosci. 1998

The Synapse



- Chemical Synapse: Action potential triggers the release of neurotransmitter
- Neurotransmitters diffuse across the synaptic cleft
- Neurotransmitter-gated channels open, generating postsynaptic potential
- Dependent on the sign of PSP, synapses are clarified as **excitatory** and **inhibitory** ones.



Short-Term Plasticity (STP)

- Synaptic efficacy varies temporally in short-time scales depending on the input history

- Short-term depression (STD)

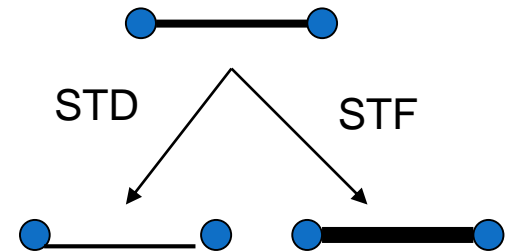
 - Neural firing depletes neurotransmitters

- Short-term facilitation (STF)

 - Neural firing elevates calcium level,

 - increasing the releasing probability

 - of neurotransmitters



Steven et al., Neuron, 1995
Markram et al, Nature, 1996

Mathematical formulation of STP

$$\begin{aligned}\frac{du}{dt} &= \frac{-u}{\tau_f} + U(1 - u^-)\delta(t - t_{sp}), \\ \frac{dx}{dt} &= \frac{1 - x}{\tau_d} - u^+x^-\delta(t - t_{sp}), \\ \frac{dI}{dt} &= -\frac{I}{\tau_s} + Au^+x^-\delta(t - t_{sp})\end{aligned}$$

- The STD effect is modeled by a normalized variable x ($0 \leq x \leq 1$), denoting the fraction of resources that remain available after neurotransmitter depletion.
- The STF effect is modeled by a utilization parameter u ($0 \leq u \leq 1$), representing the fraction of available resources ready for use (release probability).
- A denotes the response amplitude that would be produced by total release of all the neurotransmitter ($u=x=1$)
- Following a spike t_{sp} ,
 - (i) u increases due to spike-induced calcium influx to the presynaptic terminal,
 - after which (ii) a fraction u of available resources is consumed to produce the post-synaptic current.
 - Between spikes, u decays back to zero with time constant τ_f and x recovers to 1 with time constant τ_d

Dynamics of STP

$$\frac{du}{dt} = \frac{-u}{\tau_f} + U(1-u^-)\delta(t-t_{sp}),$$

$$\frac{dx}{dt} = \frac{1-x}{\tau_d} - u^+x^-\delta(t-t_{sp}),$$

$$\frac{dI}{dt} = -\frac{I}{\tau_s} + Au^+x^-\delta(t-t_{sp})$$

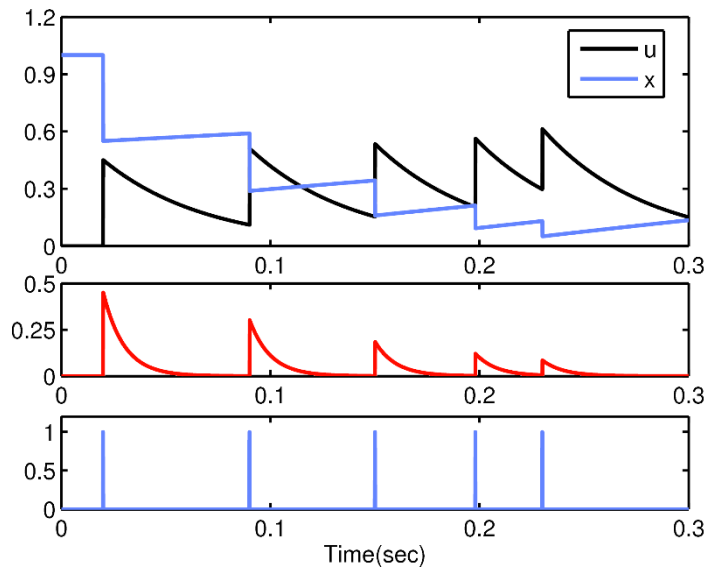
Upon arriving a spike

$$u^+ = u^- + U(1-u^-)$$

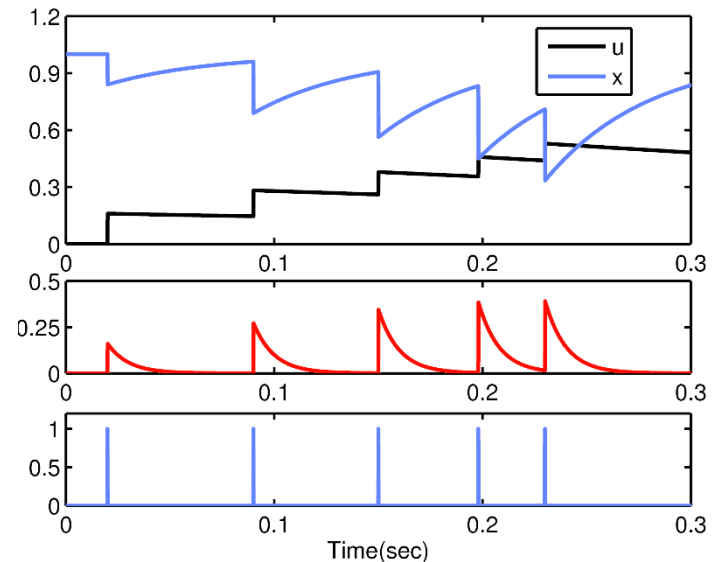
$$x^+ = x^- - u^+x^-$$

$$I^+ = I^- + Au^+x^-$$

STD-dominant



STF-dominant



Why STP

- An ubiquitous phenomenon in cortexes
- Large diversity in different cortical areas
 - STD-dominating in sensory cortexes
 - STF-dominating in prefrontal cortex
- Time scale $\sim 10\text{ms}$ - 1000ms , between fast neural signaling and slow learning
- Neural substrates for processing temporal information in the relevant scales?

Functional Roles of STP

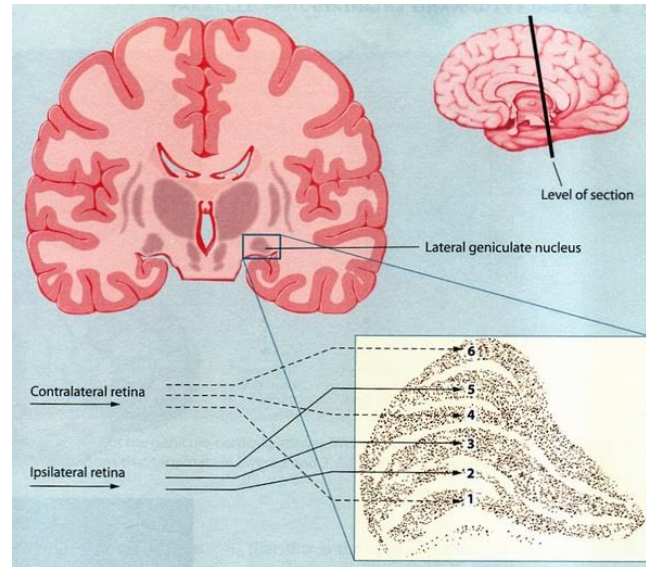
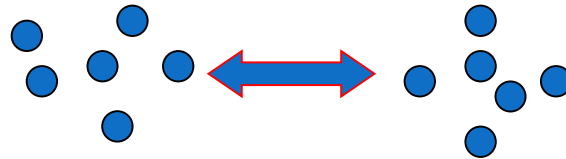
Reviews

- Wu & Tsodyks, Scholarpedia 2013
- Wu, Wong, & Tsodyks,
Neural Information Processing with Dynamical Synapses,
Frontiers in Computational Neuroscience, e-book, 2013

STP for feature binding

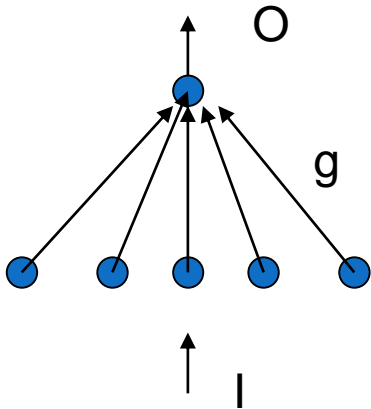
■ STF: dynamical link for feature binding

(von der Malsburg)



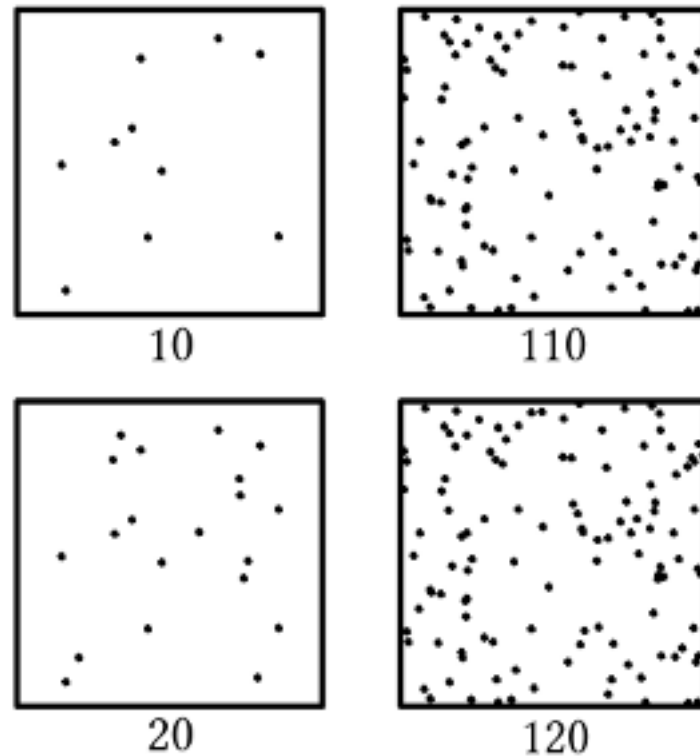
STP for Weber's law

- STD: gain control, satisfying the Weber's law (Abbott et al.)

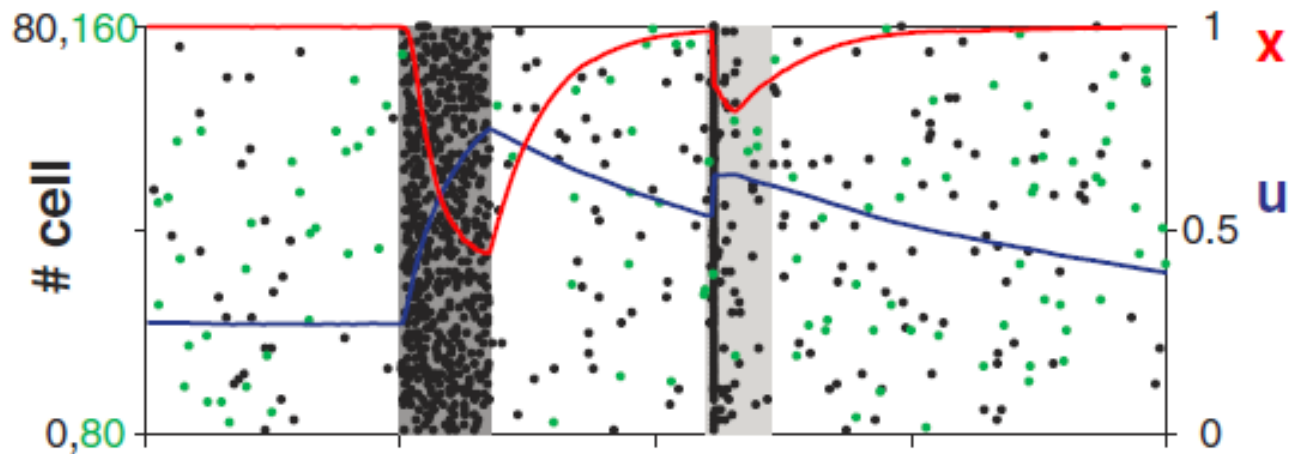


$$p = k \ln \frac{S}{S_0}$$

$$dp = k \frac{dS}{S}$$



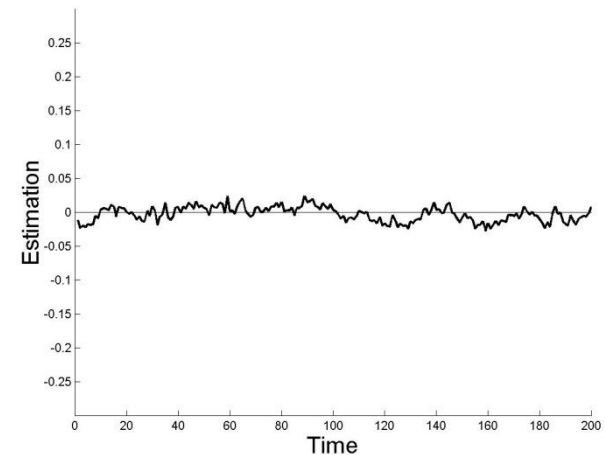
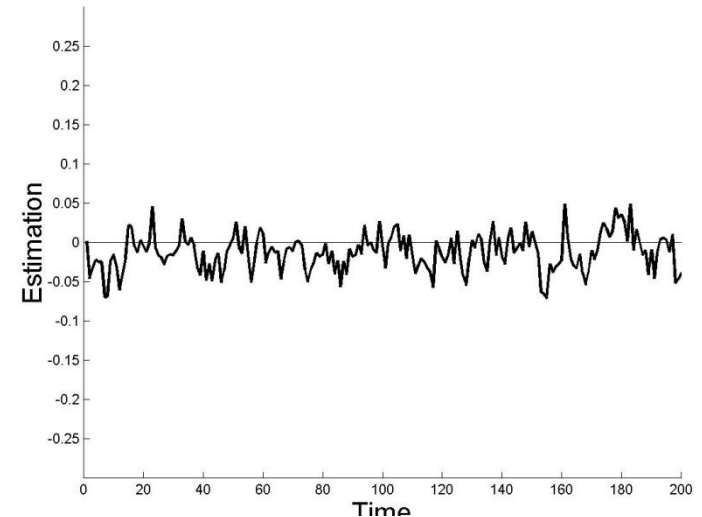
STP for population spike



Transient population response followed by silence

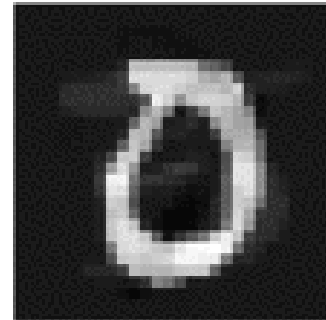
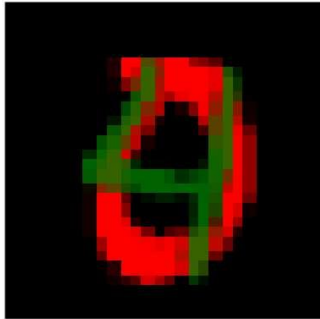
STF implements temporal Bayesian decoding

$$J(x, x', t) = J_0(x, x') + W(x, x', t)$$
$$\tau \frac{dW(x, x', t)}{dt} = -W(x, x', t) + \lambda r(x)r(x')$$

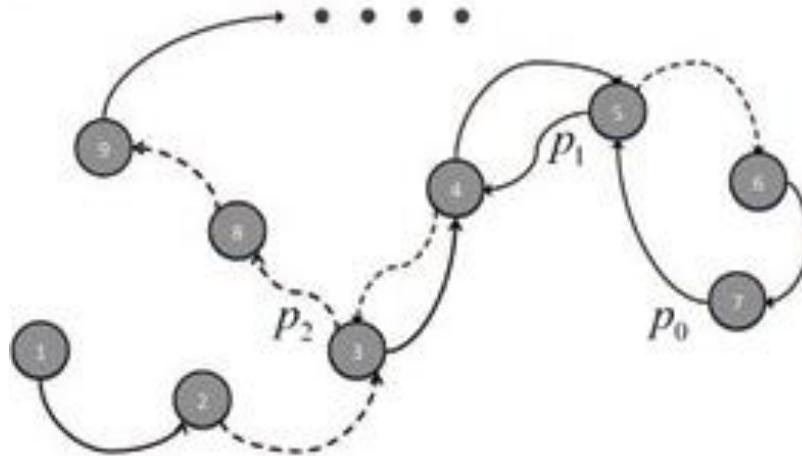


STP for state switching

STD induces switching between attractors

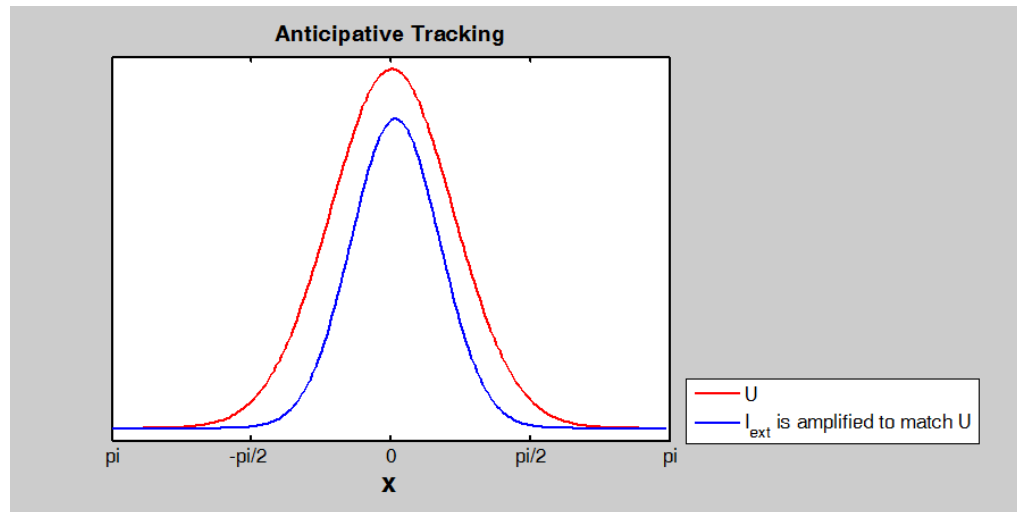


STP for memory search



- Neuronal connections encode similarity
- STD triggers state hopping

STP implements anticipative tracking



Short-term Memory with Graded Lifetime

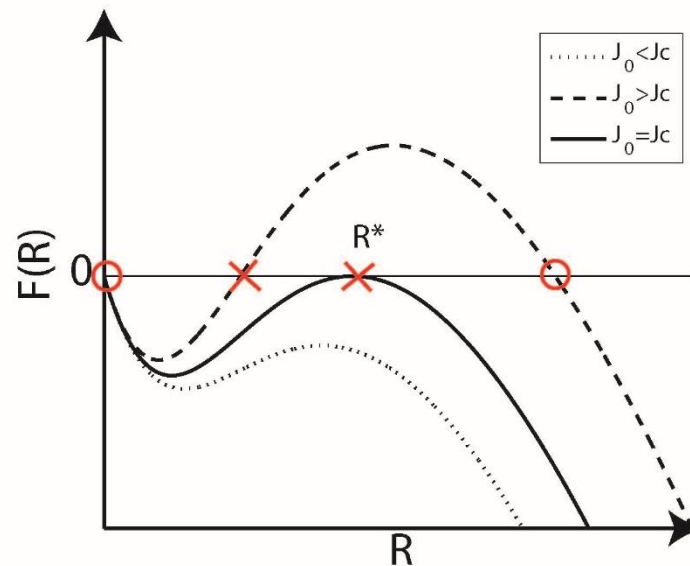
$$\frac{du}{dt} = \frac{U - u_{\mu}}{\tau_f} + U(1 - u)R$$

$$\frac{dx_\mu}{dt} = \frac{1 - x_\mu}{\tau_d} - uxR$$

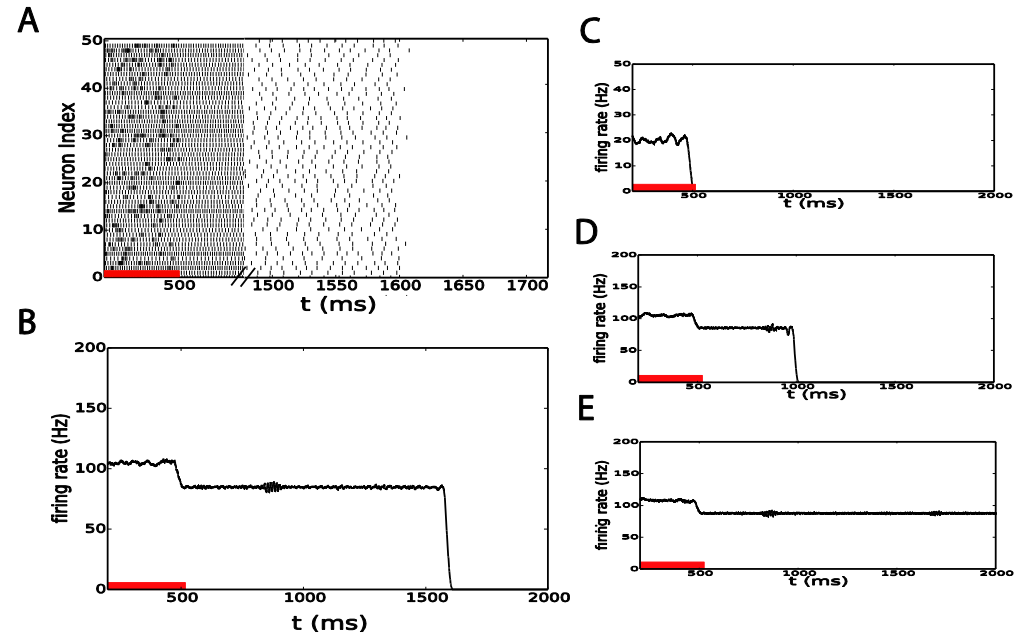
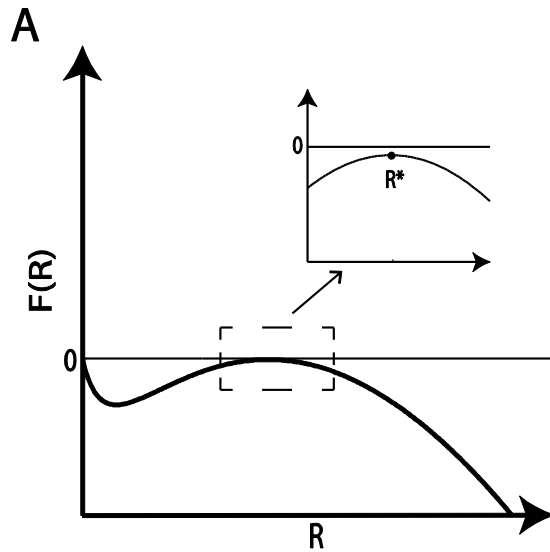
$$\tau_s \frac{dh}{dt} = -h + J_0 u x R + I$$

$$R = \max(\beta h, 0)$$

$$\tau_s \frac{dR}{dt} = -R + \frac{J_0 \beta \tau_f U R^2}{1 + \tau_f U R + \tau_d \tau_f U R^2} \equiv F(R)$$



Short-term Memory with Graded Lifetime



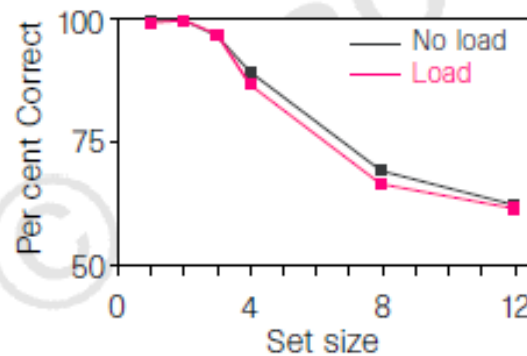
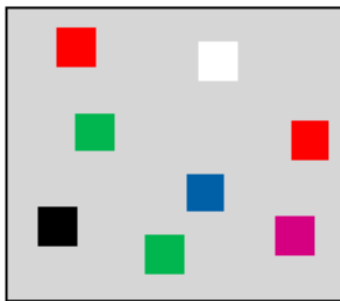
- The network is marginally unstable
- Persistent activity is closed naturally

- The lifetime of persistent activity varies with the combination of STF & STD
- Diversely distributed STP in cortical areas hold different memory times

Interplay between STF & STD determines the lifetime of neural activity

Working Memory

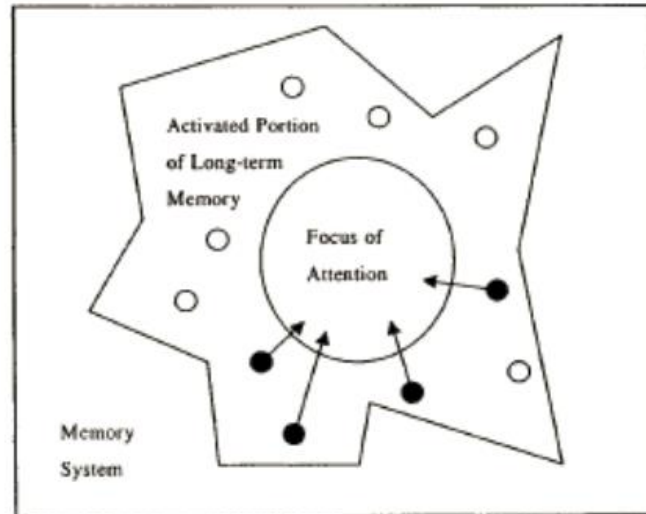
- ✓ **Working memory** refers to short-term storage and manipulation of information for cognitive functions .
--- Miller et al, 1960



Steven J. Luck & Edward K. Vogel,
Nature, 1997

- ✓ The sample array consisted of 1-12 colored squares and was presented for 100ms
 - ✓ This was followed by a 900ms blanked delay interval
 - ✓ A 2000ms presentation of the test array, which was either identical to the sample or differed in the color of one of the squares
 - ✓ Whether the two arrays were identical or different in terms of a single feature.
-
- ✓ **Working memory capacity** is extremely limited, ranging between 3 and 6 items for most healthy human participants.
--- Cowan et al, 2001

Focus of Attention



Cowan et al, 2001

- ✓ The brain possesses a specialized buffer, or 'focus of attention', where memory items can be temporarily placed for short periods of time and removed when needed.
- ✓ The working memory capacity corresponds to the size of the buffer.
- ✓ **The neural implementation** of the focus of attention and its size, as well as the way memory items are placed and removed from it, **remains unresolved**.

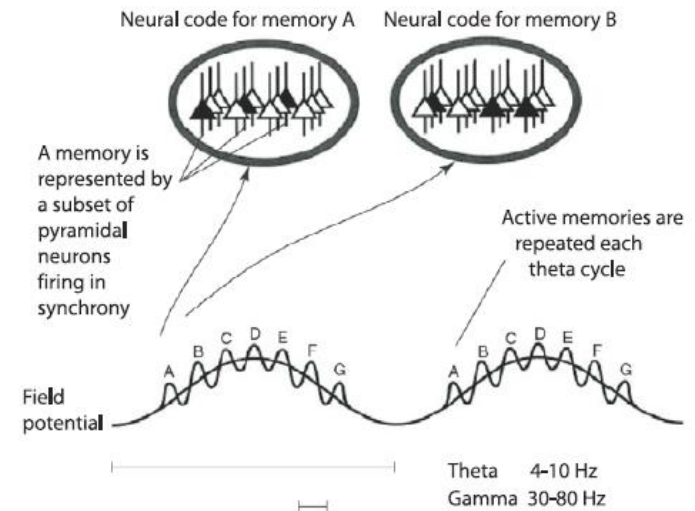
Hypotheses

✓ Persistent neural firing

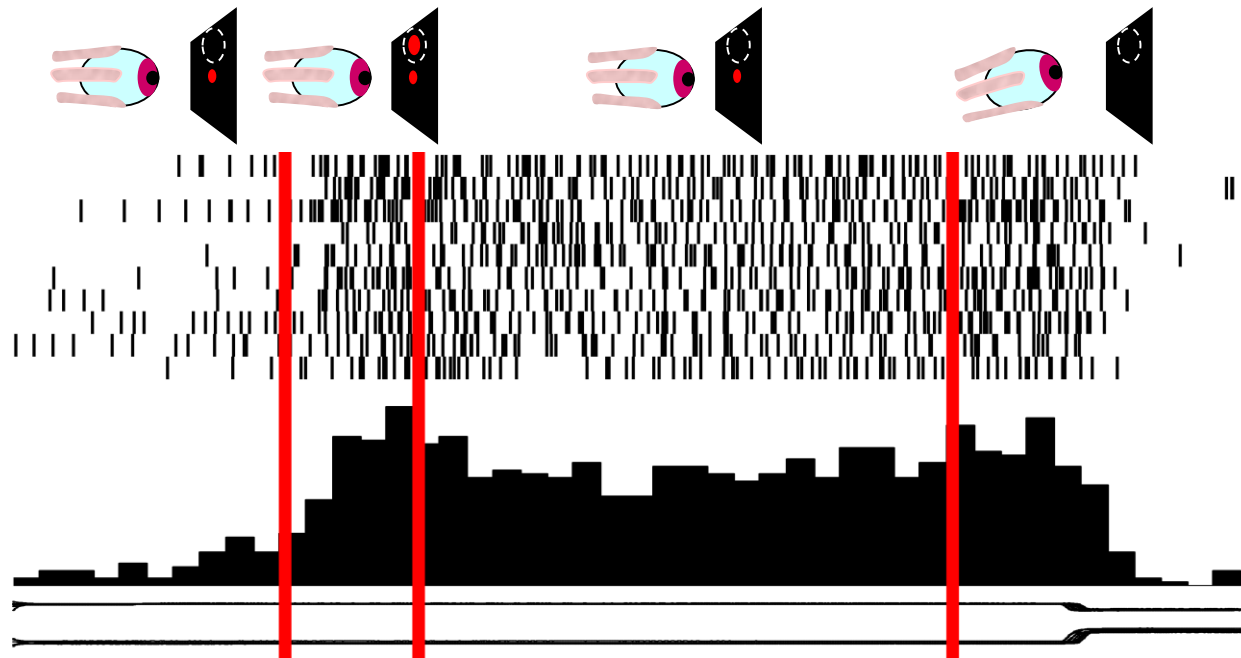
- Working memory is mediated by persistent activity of neurons encoding the corresponding items in long-term memory.
- The maximal number of items simultaneously active depends on the characteristics of the network in a complex way.
- No fundamental upper limit on working memory in this model.

✓ Nested theta-gamma Oscillation

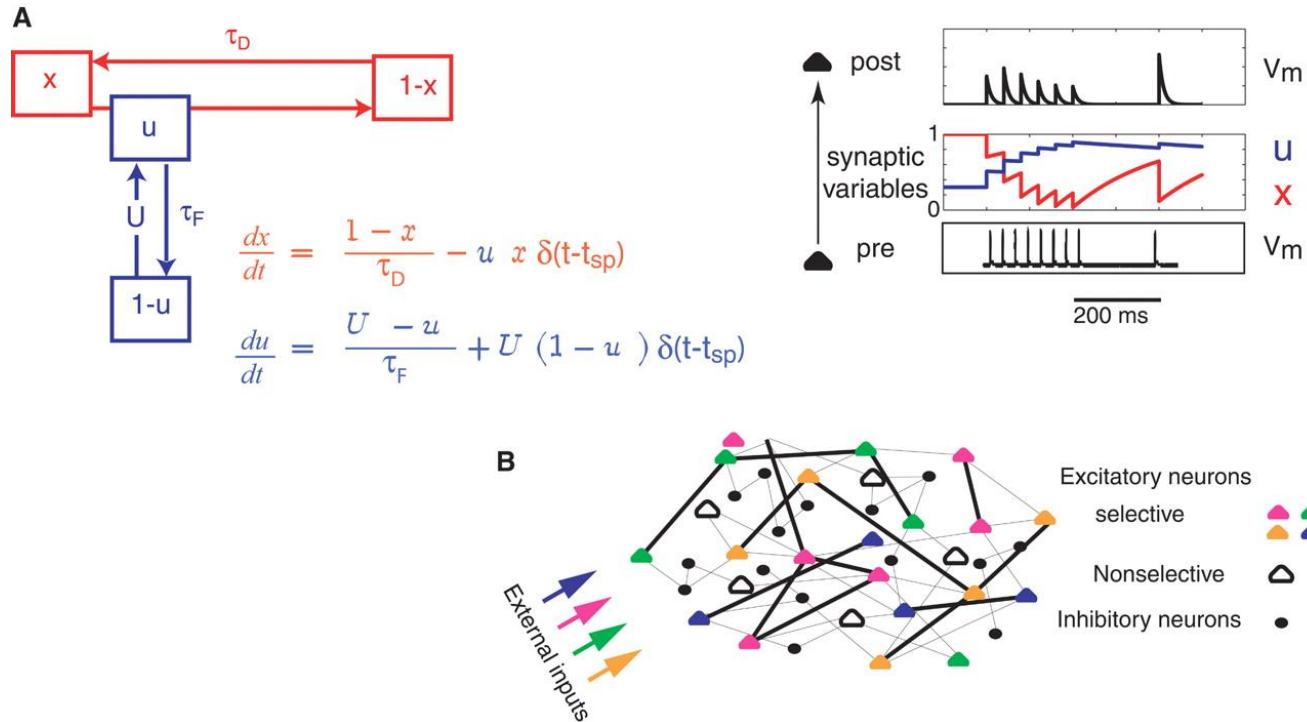
- Multiple memories are stored in the nested theta-gamma oscillatory patterns, and can be activated sequentially at different moments of time. ----- (Lisman, **Nature**, 1995)
- Each memory is stored in a different high-frequency (gamma wave) subcycle embedded in a low-frequency oscillation (theta wave).
- The WM capacity depends on the model parameters?



Working Memory with Persistent Firing



Working Memory without Persistent Firing



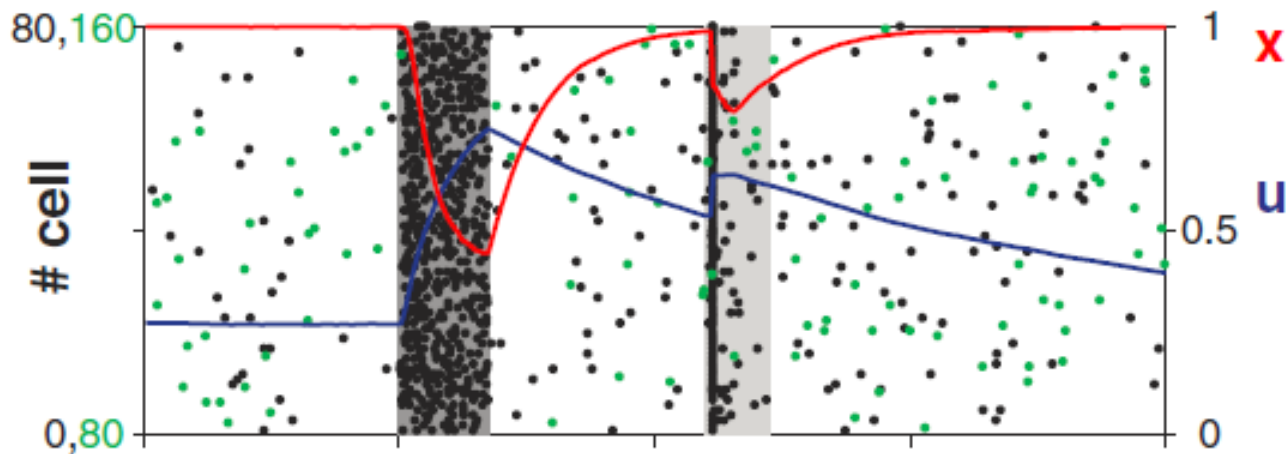
Facilitated synapses hold the memory trace of an external input

Gianluigi Mongillo et al. Science 2008

Working Memory via Short-term Plasticity

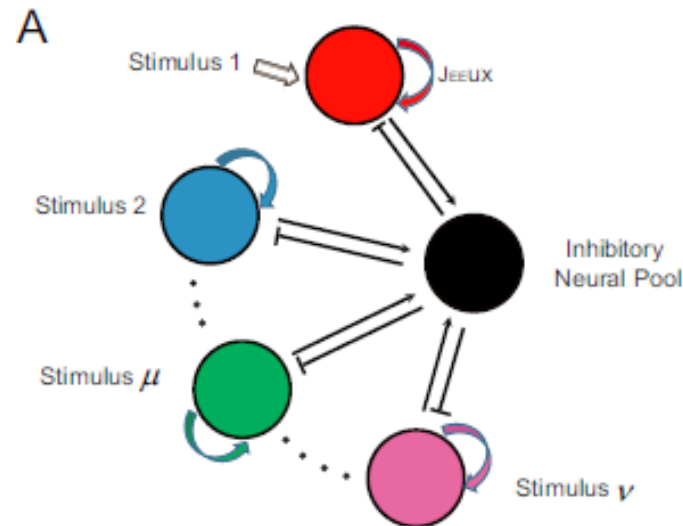
- Memory is retained in the facilitated recurrent connections between neurons

-----Mongillo, Barak and Tsodyks, **Science**, 2008



- ✓ An external stimulus is loaded into memory by activating the corresponding subpopulation of neurons, afterwards the neural population return to the spontaneous state.
- ✓ A non-specific read-out signal can selectively recall the memory item, since the corresponding neural population has larger synaptic strengths.

A Working Memory System based on STP



- ✓ Each memory item is stored at an excitatory neural cluster
- ✓ All excitatory neural clusters are connected to an inhibitory neuron pool;
- ✓ The inhibitory neuron pool feeds back to all excitatory clusters, preventing explosive responses and inducing competition among excitatory clusters.
- ✓ No overlap between excitatory neural clusters.

The Network Dynamics

$$\frac{du_{\mu}}{dt} = \frac{U - u_{\mu}}{\tau_f} + U(1 - u_{\mu})R_{\mu}$$

$$\tau_f = 1.5s, \quad \tau_d = 0.2s$$

$$\frac{dx_{\mu}}{dt} = \frac{1 - x_{\mu}}{\tau_d} - u_{\mu}x_{\mu}R_{\mu}$$

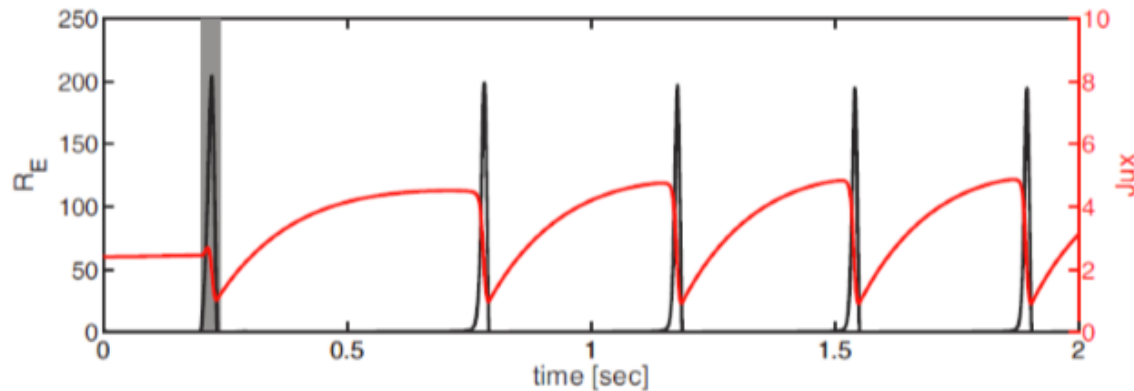
$$\tau \frac{dh_{\mu}}{dt} = -h_{\mu} + J_{EE}u_{\mu}x_{\mu}R_{\mu} - J_{EI}R_I + I_b + I_e(t)$$

$$\tau \frac{dh_I}{dt} = -h_I + J_{IE} \sum_{\nu} R_{\nu}$$

$$R = g(h) = \alpha \ln(1 + \exp(h/\alpha))$$

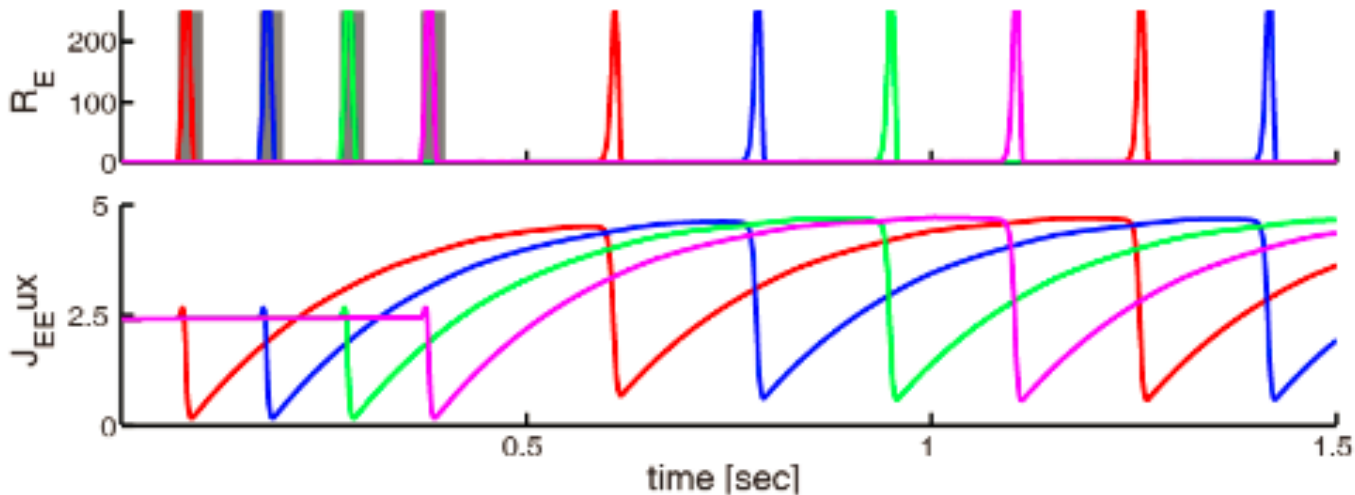
I_b refers to the arousal signal when the neural system is engaged in a WM task.

Loading & Retrieving a Memory Item



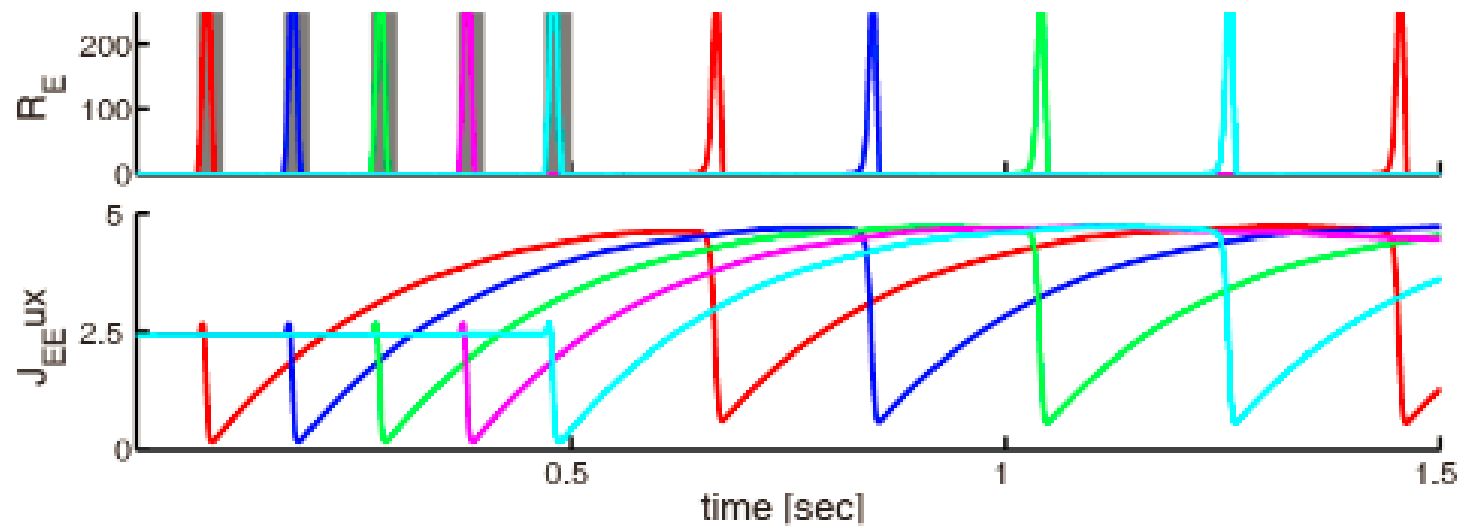
- ◇ **Dynamics:** with a strong enough I_b , a neural cluster falls into a limit-circle state, i.e., the cluster generates PS periodically.
- ◇ **Loading:** due to STP, a neural cluster generates PS, in response to an external transient input, corresponding to a memory item is loaded.
- ◇ **Retrieving:** with a strong enough arousal level (large I_b), once a memory item is loaded, the cluster will generate PS periodically, corresponding to the retrieval of the loaded memory.

Storing & Recalling Multiple Items Sequentially



- ◇ Loading multiple memory items one-by-one, with one item at one excitatory cluster
- ◇ Neural clusters recall the memorized items in the same order as how they are loaded.

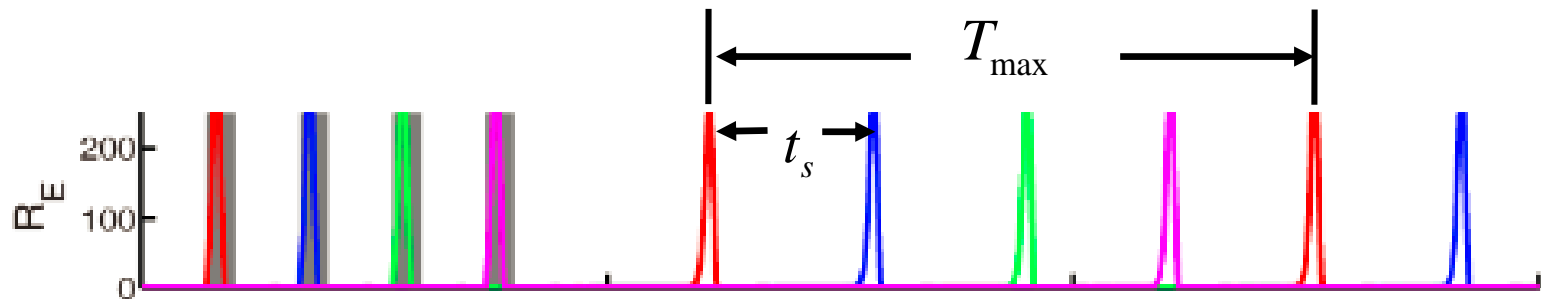
Capacity of Working Memory



The Rise of Capacity: Theoretical Analysis

The maximum number of items can be maintained in WM is determined by the ratio of two factors:

- ✓ The maximum period of T_{\max} of the limit cycle of the network, i.e., the maximum time between subsequent reactivation of each cluster.
- ✓ The temporal separation between two consecutive PSs, referred to as t_s



The capacity of WM is given by the maximum number of PSs that can be accommodated in a single period of the limit cycle,

$$N_C \approx T_{\max} / t_s$$

Capacity of Working Memory

$$N_C \approx \frac{\tau_d}{\tau} \frac{\ln \frac{\tau_f/\tau_d}{1-U}}{\ln \frac{|h_0|}{I_b - I_{\text{crit}}}} + C$$

- ✓ The WM capacity scales with the ratio of two time constants, one characterizing the synaptic depression and the other one – synaptic current decay time.
- ✓ The WM capacity is controlled by the background excitation that should be above the critical level below which no items can be maintained in WM.

References

1. Zhang D, Li Y, Rasch MJ and Wu S (2013) Nonlinear multiplicative dendritic integration in neuron and network models. *Front. Comput. Neurosci.* 7:56. doi: 10.3389/fncom.2013.00056
2. Danke Zhang, Xichun Zhang, Malte Rasch, Si Wu (2013). Divisive Normalization by Shunting Inhibition in Neural Networks. Internal Joint Conference on Artificial Intelligence-Workshop on Intelligence Science (IJCAI-WIS2013), Beijing 2013.
3. Si Wu & Misha Tsodyks, Short-term synaptic plasticity. Scholarpedia 2013
4. Wu, S., Wong, K. Y. M., & Tsodyks, M. (2013). Neural information processing with dynamical synapses. *Frontiers in Computational Neuroscience*, 7(188), 1. doi:10.3389/fncom.2013.00188.
5. T. Haga and T. Fukai. (2018) Recurrent networks model for learning goal-directed sequences through reverse replay. *eLife* 7:e34171.
6. Yuanyuan Mi, Mikhail Katkov, Misha Tsodyks, Synaptic Correlates of Working Memory Capacity, *Neuron*, 93, 323–330, 2017
7. J. Ba et al. Using fast weights to attend to the recent past. NIPS 2016.