

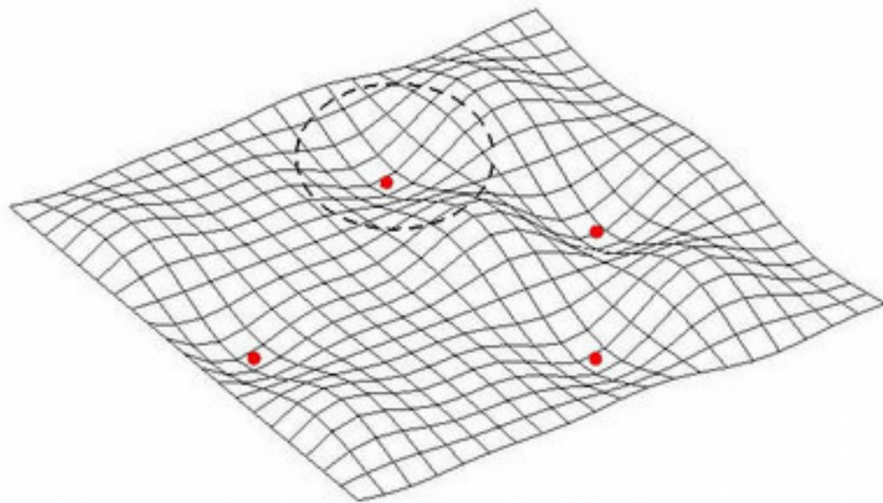
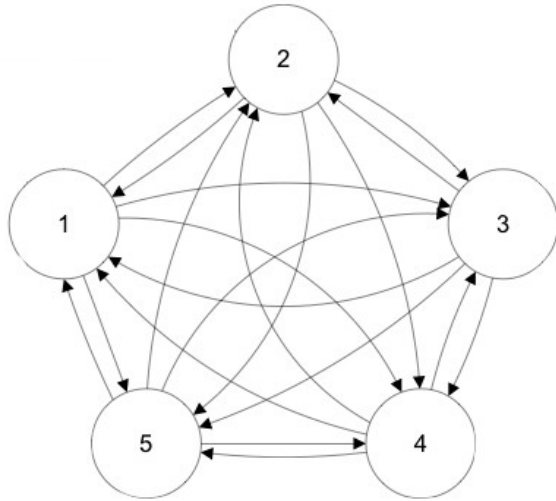
# A purely local, distributed, simple learning scheme achieves near-optimal capacity in recurrent neural networks without explicit supervision

Alireza Alemi, Carlo Baldassi, Nicolas Brunel, Riccardo Zecchina

BDA 2015



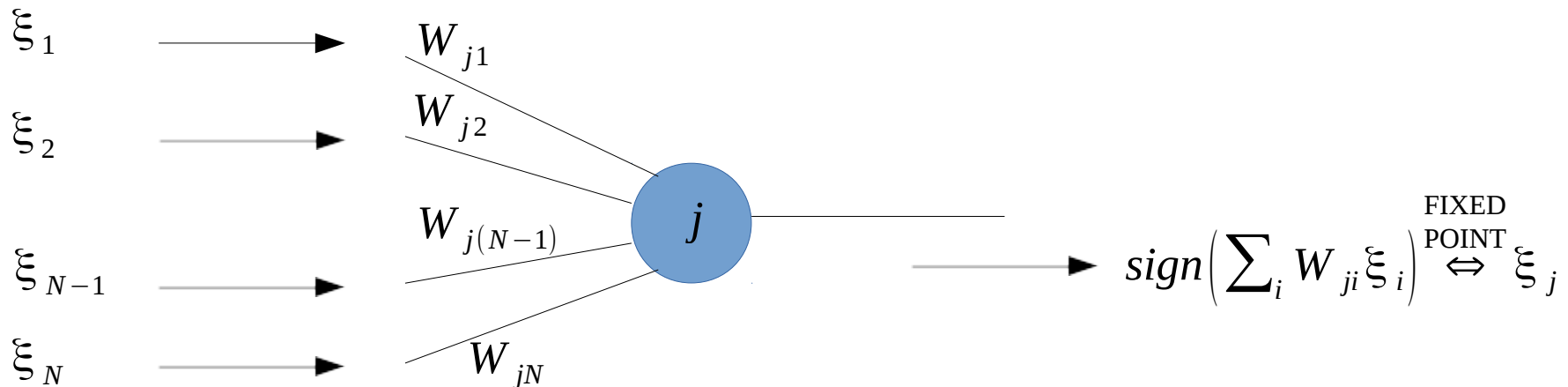
# Attractor networks



- Popular model for information storage in the brain (memorization – working memory, recognition, error-correction, ...)
- Recurrent neural networks
- Distributed model (each unit behaves independently, information is stored in the collective behaviour)
- Learning → patterns of activity are encoded as fixed points of the network dynamics
- Robustness → basins of attraction around the fixed points

# Hopfield network

- First, most popular model (1984), with many later variants
- Binary  $\pm 1$  units and patterns (perceptrons  $\rightarrow output_j = sign(\sum_i W_{ji} \xi_i)$ )
- **Hebbian learning** (“fire together  $\rightarrow$  wire together, out-of-sync  $\rightarrow$  fail to link”)
  - Rule:  $W_{ij} = \frac{1}{M} \sum_a \xi_j^a \xi_i^a$
- **Pros**: simple, local, distributed, unsupervised, some experimental support
- **Cons**: symmetric, low capacity ( $\sim 0.138N$ ), catastrophic forgetting beyond capacity



# Perceptron learning rule (PLR)

- On line, supervised learning rule for training individual units:
  1. Present a pattern at random:  $\xi^a$
  2. In case of error, change  $W_{ji}$  in the opposite direction, modulated by  $\xi^a$

$$\Delta W_{ji} = \eta \xi_i^a \left( \xi_j^a - \text{sign} \left( \sum_i W_{ji} \xi_i^a \right) \right)$$

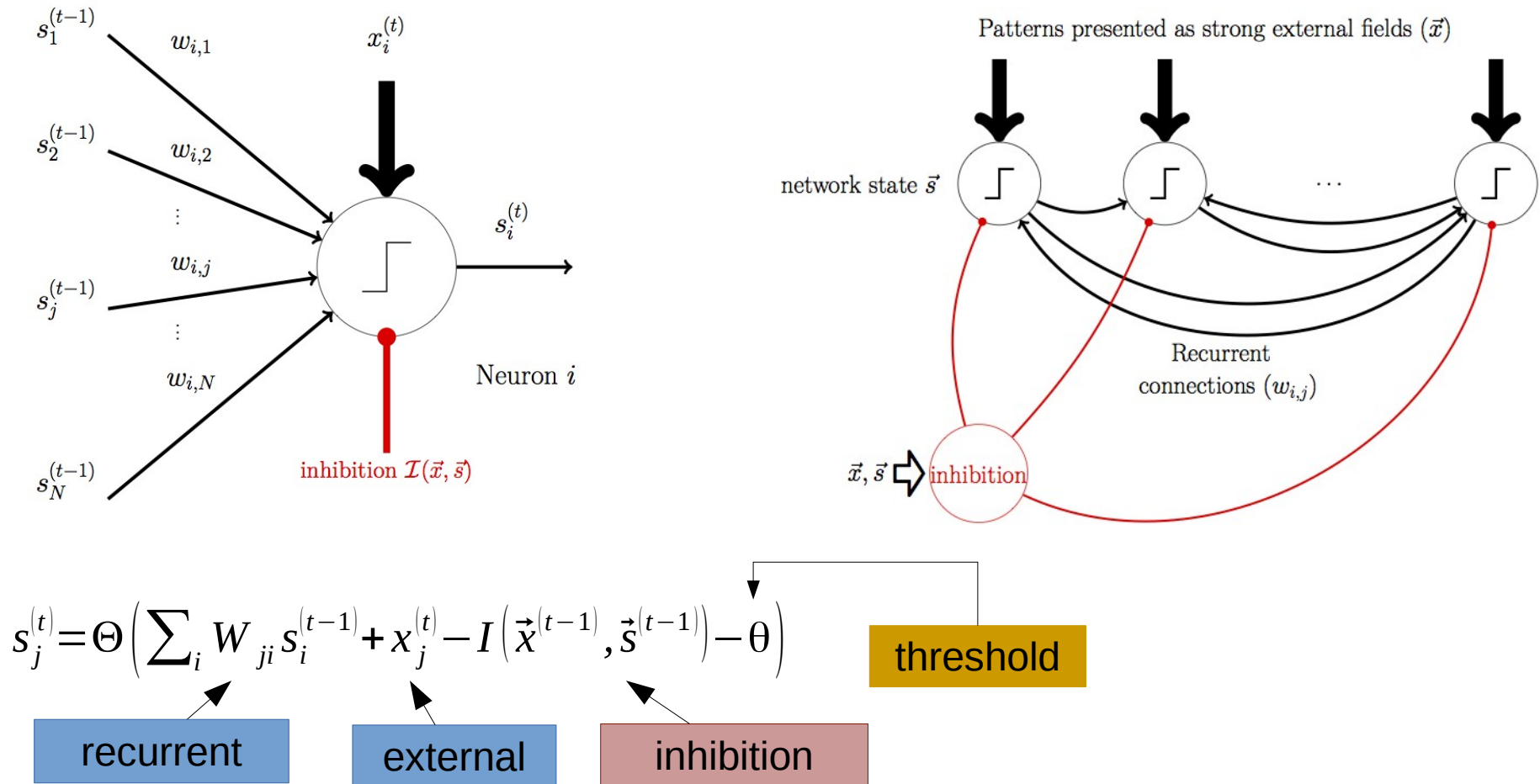
- **Pros**: able to achieve the maximal capacity ( $\sim 2N$ ) (even with **correlated patterns**), no catastrophic forgetting, allows asymmetric weights
- **Cons**: requires an explicit supervisory error signal
  - i.e. compare “output in absence of the pattern” vs “pattern itself”

# Best of both worlds?

- **Goal:** get the best of both worlds, a distributed, local, simple, unsupervised rule which achieves maximal capacity, allows asymmetric weights, has no catastrophic forgetting – in a **more realistic** setting
- **Means:** convert the PLR in an unsupervised setting, using statistical properties of the inputs
- **Main observation:** the statistic of the depolarization fields carries enough information about the error type → no need for an explicit error signal

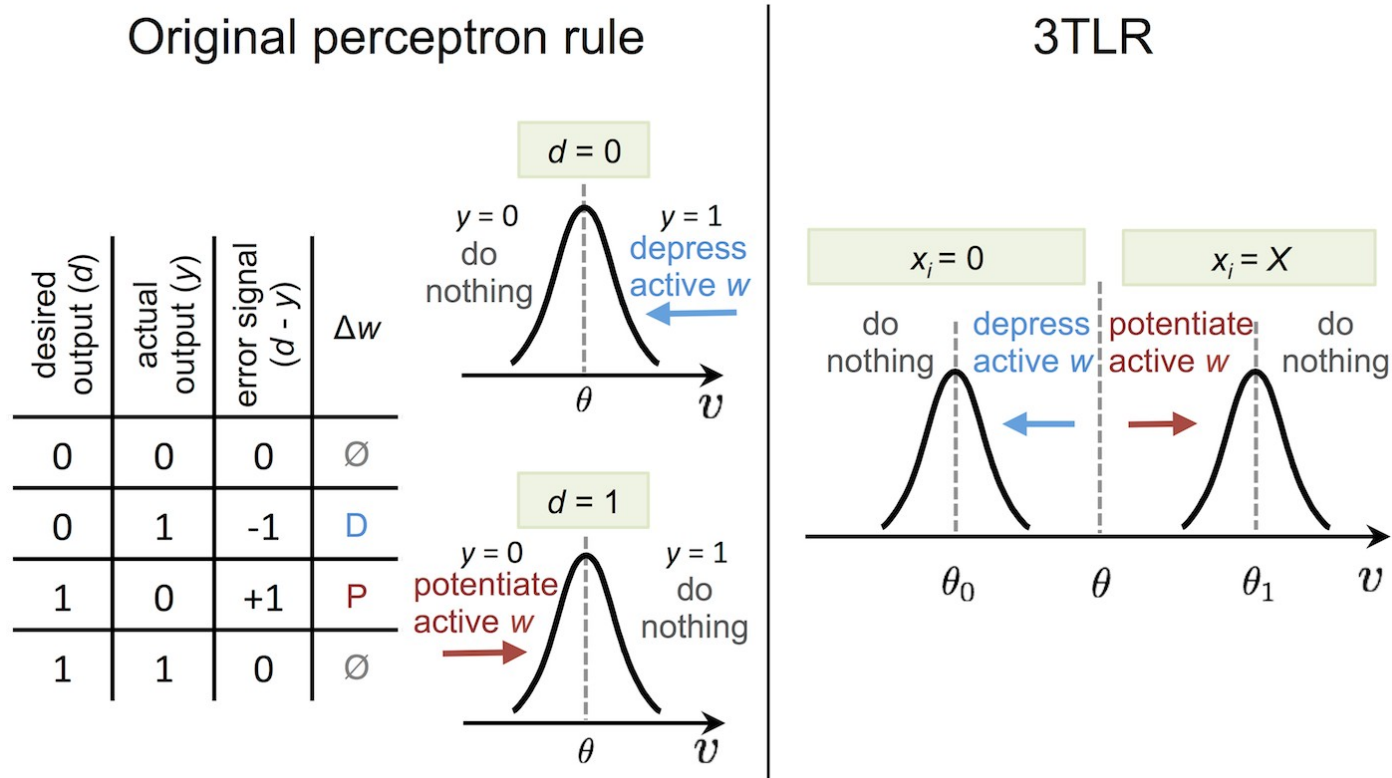
# Our network model

- Network model:** excitatory population, state-dependent inhibitory feedback for stabilization, patterns presented via external inputs

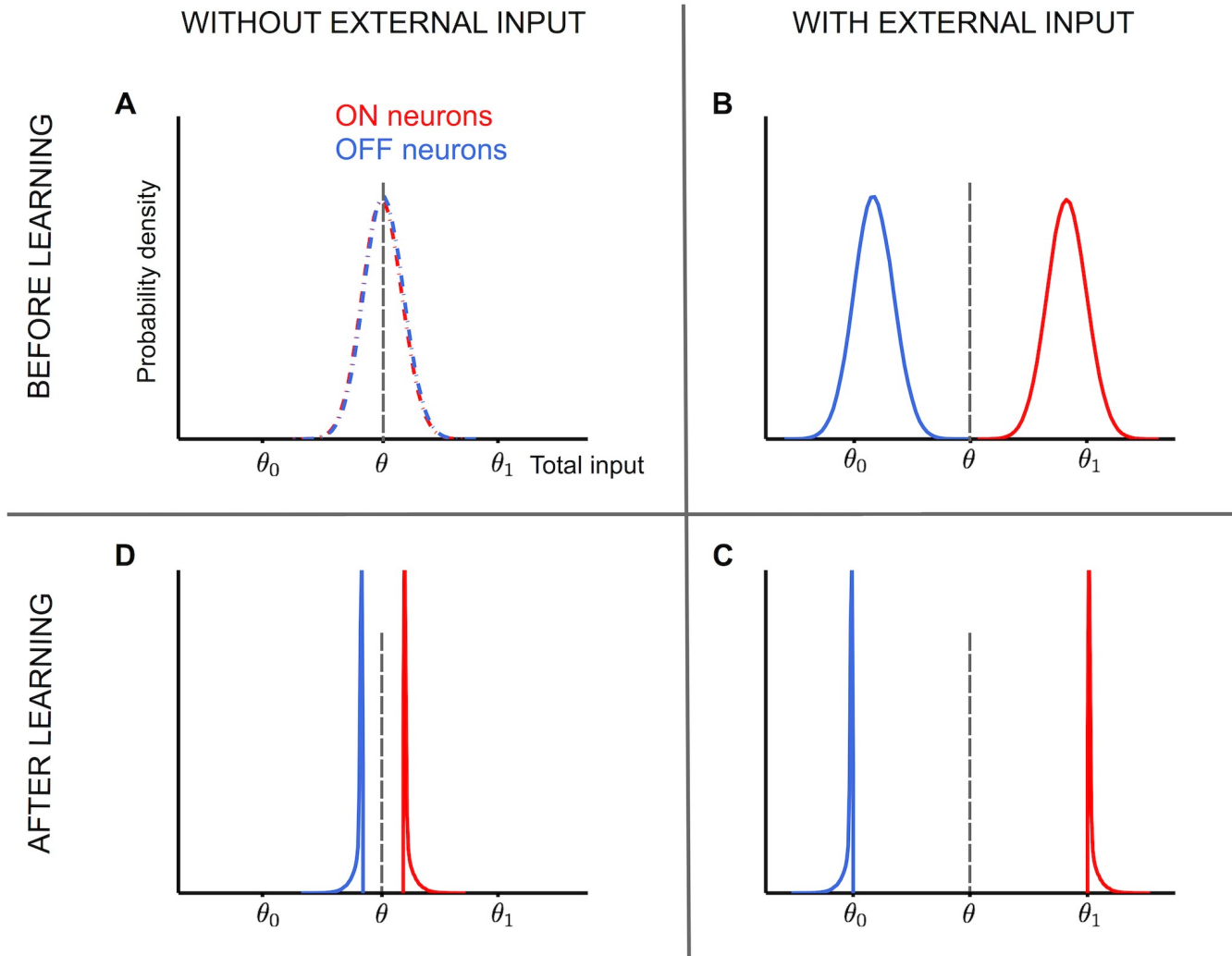


# A three-threshold learning rule (3TLR)

- Converting the PLR into an **unsupervised** rule: 3TLR
- **Crucial observation**: depolarizations  $\sum_i W_{ji} s_i$  are distributed according to a Gaussian of width  $O(\sqrt{N})$ ; external inputs  $x_j$  make them bimodal



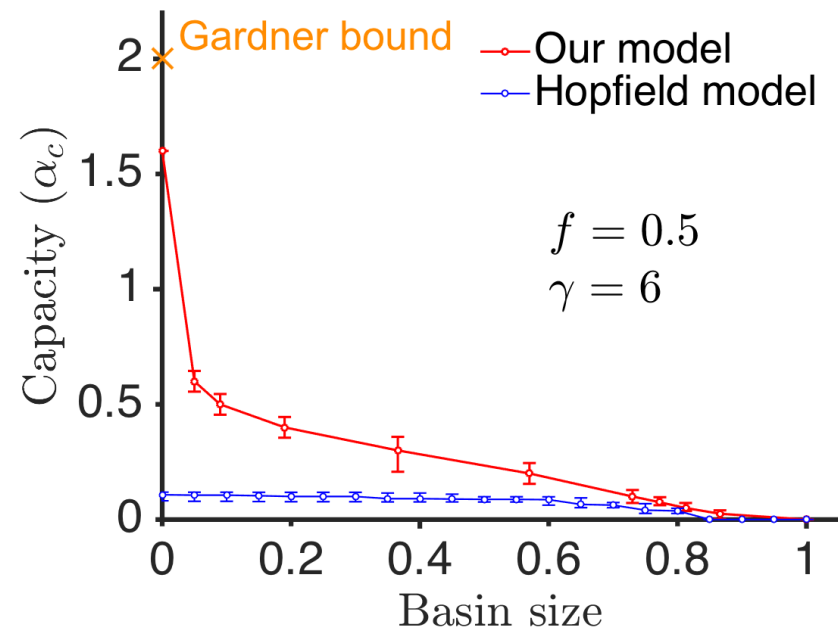
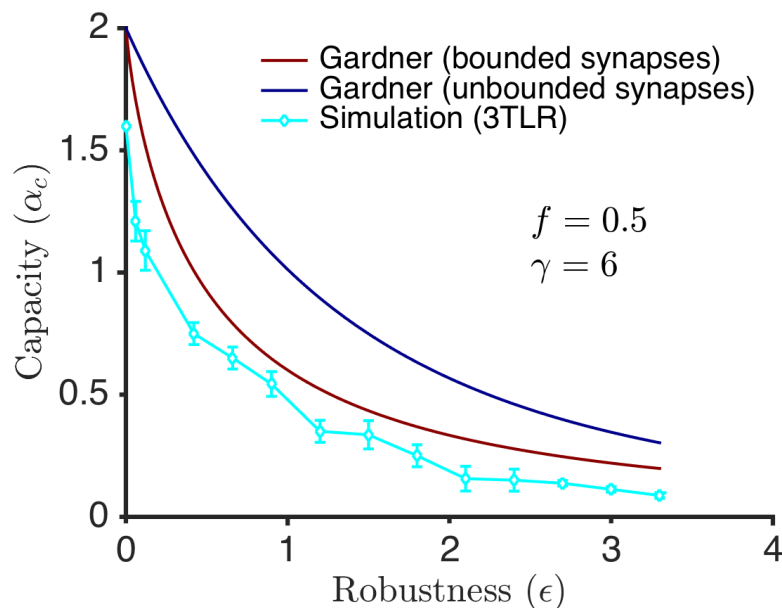
# 3TLR in action





# Simulation results

- Many-fold increase in capacity w.r.t. Hopfield network, even though the theoretical capacity is lower
- Works in sparse regime, with correlated patterns etc.
- Depends on the external inputs being strong enough (although...)
- After learning, many synapses are off (“silent”) → sparsification



# Further comments, future directions

- Parameter tuning → unsupervised pre-training phase (learn the general statistics of the inputs)
- 3TLR can be framed within the BCM theory (with additional specifications)
- Experimental evidence?
- The transformation could be applied to (essentially) any supervised rule (e.g. discrete synapses)
- Future direction: more realistic neurons (firing-rate, integrate and fire, HH) and general scenario

# Thanks!

see details in (to be published soon):

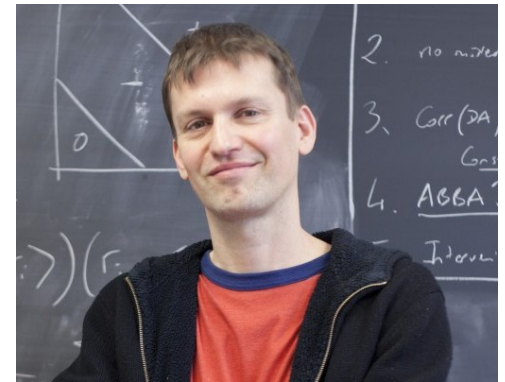
*A Three-Threshold Learning Rule Approaches the Maximal Capacity of Recurrent Neural Networks*, A. Alemi, C. Baldassi, N. Brunel and R. Zecchina, *Plos. Comp. Biol.* 2015



Alireza Alemi

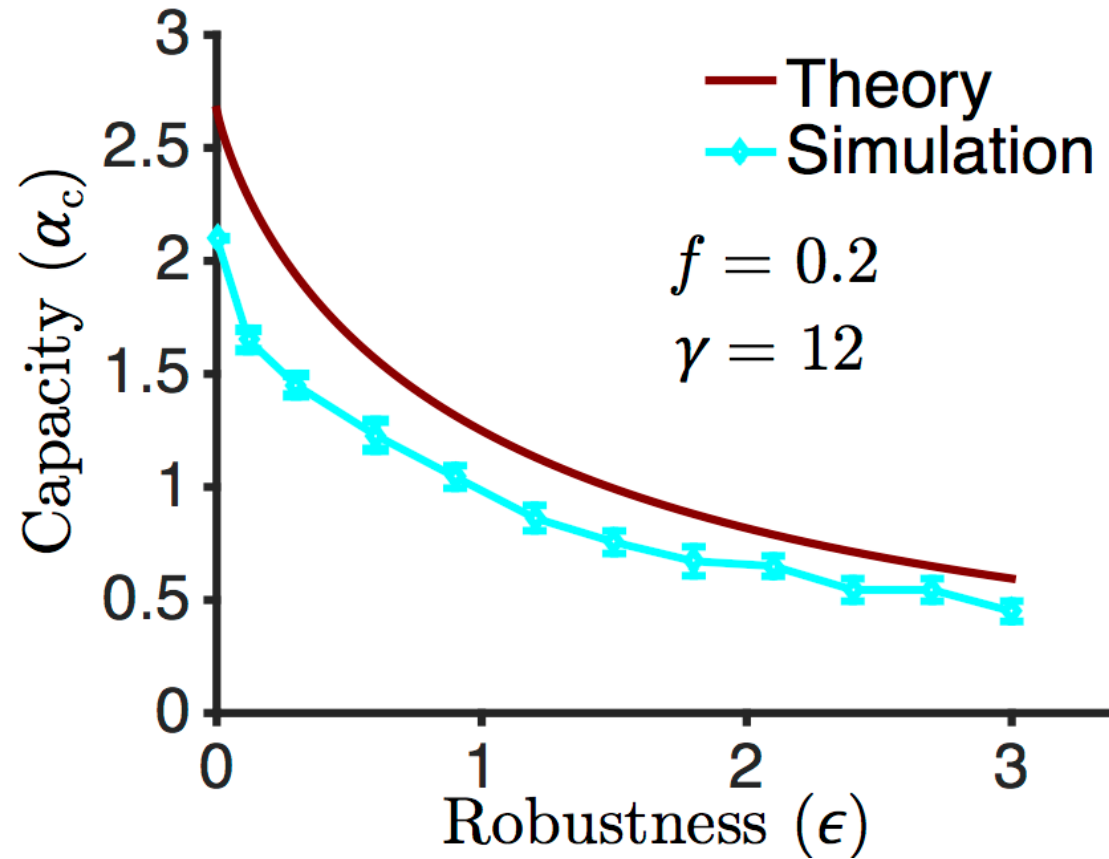


Riccardo Zecchina

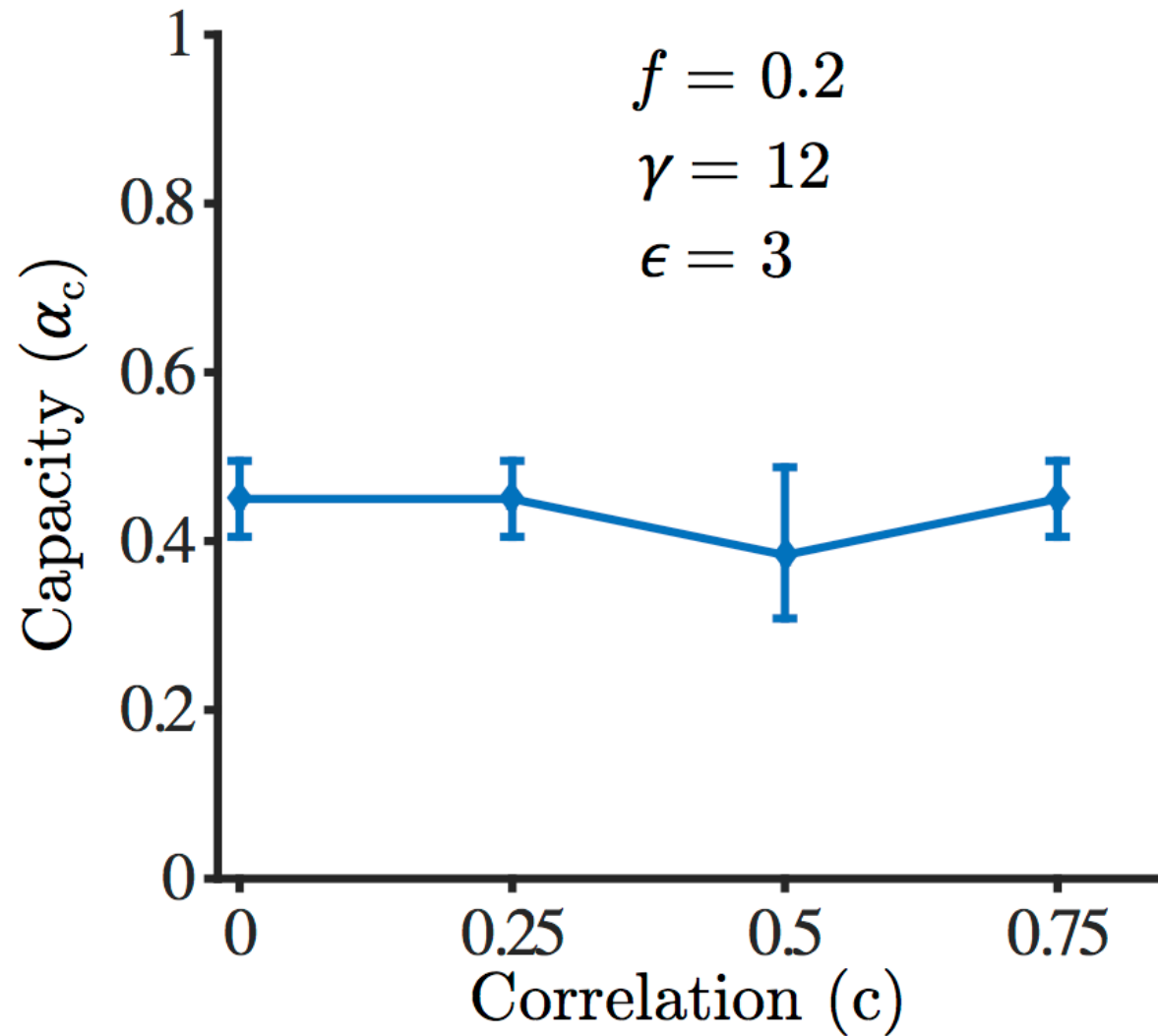


Nicolas Brunel

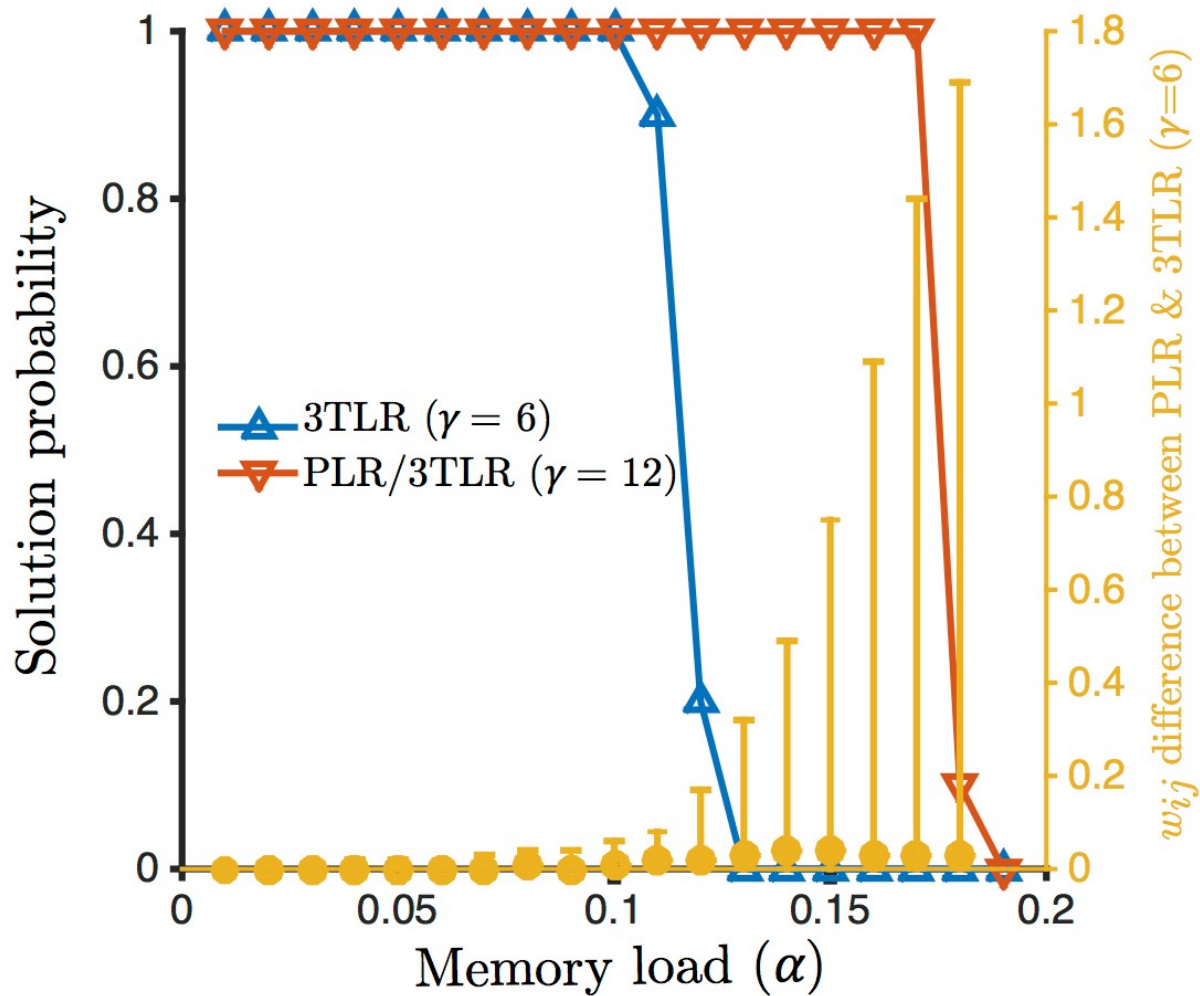
# Sparse case



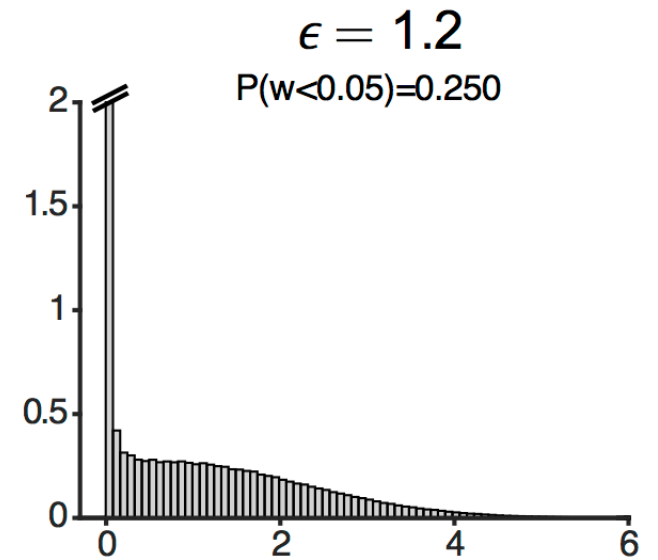
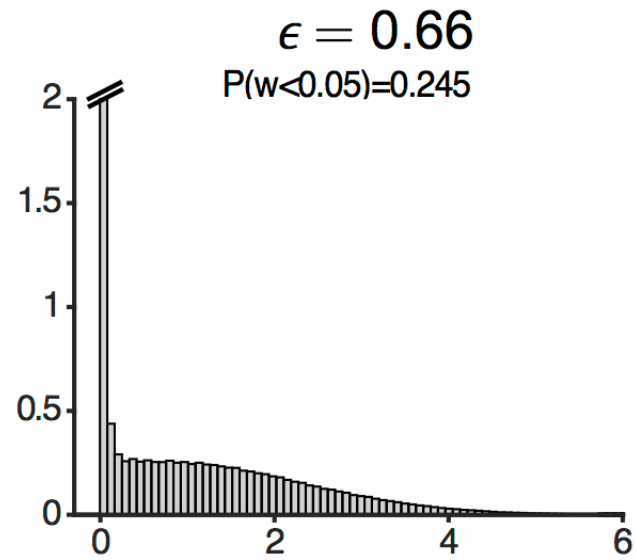
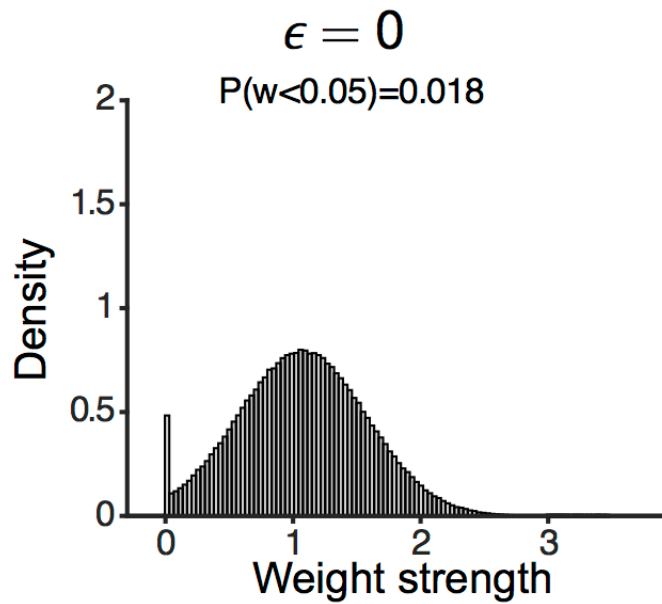
# Correlated patterns



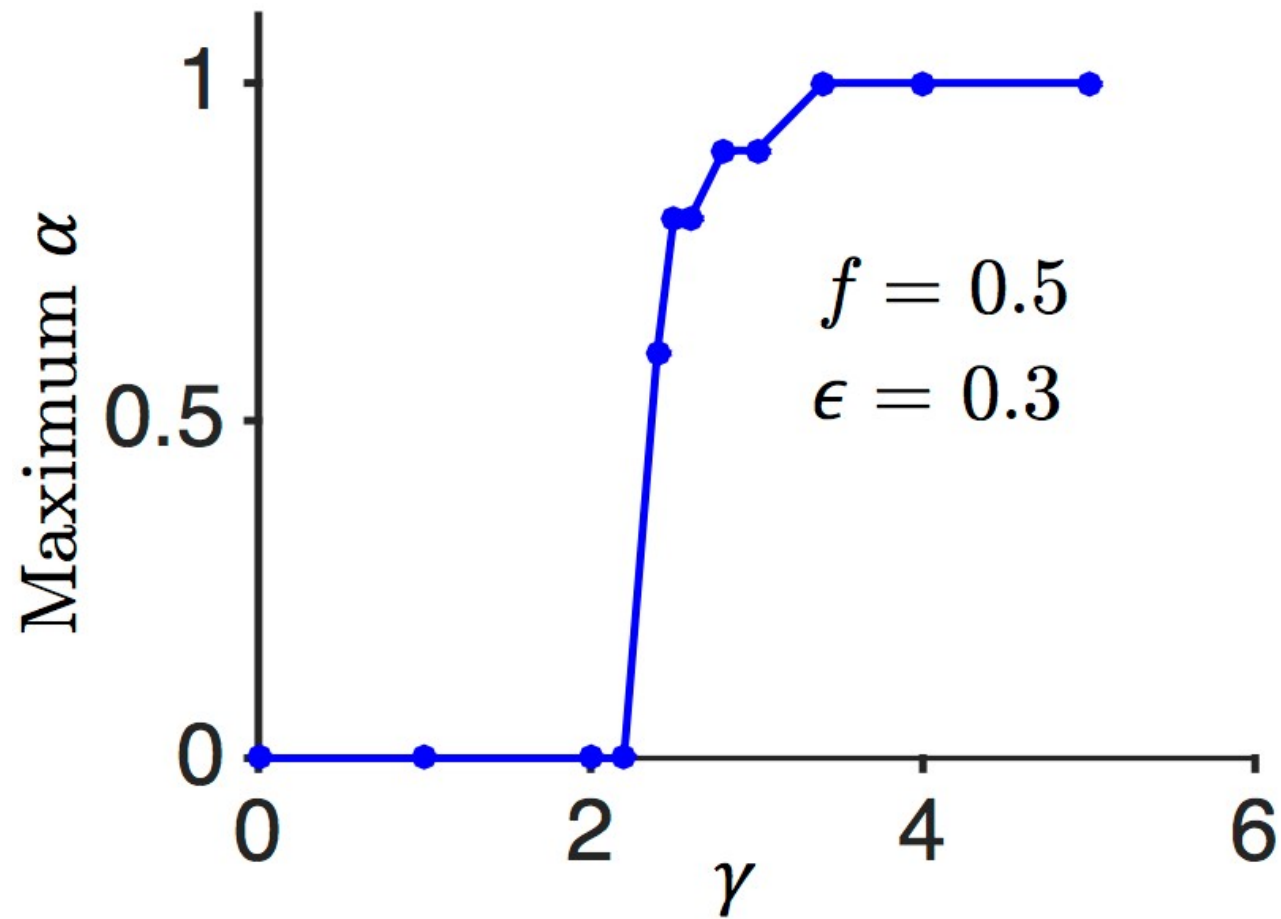
# 3TLR vs PRL



# Distribution of synaptic weights



# External field strength





# Weights symmetry

