

Storing messages with multipartite neural cliques

Nissim Zerbib¹, Vincent Gripon², and ?³

¹*Département d'Informatique, École normale supérieure, Paris, France*

²*Département d'Électronique, Télécom Bretagne, Brest, France*

Abstract—We extend recently introduced associative memories based on clustered cliques to multipartite cliques. We propose a variant of the classic retrieving rule. We study its performance relatively to the former one for retrieving partially erased or corrupted messages. We provide both analytical and simulation results showing improvements in both networks efficiencies and resilience to damages. We compute asymptotic capacities of these networks.

Index Terms—associative memory, error correcting code, cliques, multipartite cliques, neural networks

I. INTRODUCTION

Associative memories

Recently, a new type of associative memories was proposed by Gripon and Berrou [2].

GBNN can also be used to retrieve messages from erroneous versions of them. We derive a formula for error rate in case of errors on messages.

Fault tolerance of those networks was extensively studied in [4]. We extend this to multipartite cliques and show significant improvements in resilience to faults.

II. NETWORKS OF NEURAL CLIQUES

An associative memory is able to store messages and retrieve them through a partial versions of them.

We first describe how messages can be stored in a binary clustered network of multipartite cliques. Secondly we detail how messages can be retrieved after being stored in such a network, thereby showing how such a network can be used as an associative memory.

A. Storing messages

Our goal is to store messages so as to be able to retrieve them from partial or erroneous versions of them.

We will consider a set \mathcal{M} of messages of fixed length c over a finite alphabet \mathcal{A} .

Messages of \mathcal{M} are stored in a clustered neural network, of c clusters of l units each. Each unit can take two values 0 (deactivated) and 1 (activated). The network is therefore binary. A cluster will correspond to a symbol in a message and the pattern of unit activations in a cluster will determine which symbol in the alphabet the cluster represents.

A binary constant weight code of weight a is a code where all codewords have the same weight a , i.e. the same number of 1 in them. Previous works ([1], [3]) considered the storage of binary constant weight codewords of weight 1. We study

the case of arbitrary weight. It is convenient to take \mathcal{A} as the set of all binary constant weight codewords of length l and of weight a , where a is a positive integer parameter which we will call number of activities and which satisfies $a \leq l$. We now have $\binom{l}{a} = |\mathcal{A}|$. We also set $m = |\mathcal{M}|$.

States of the network are characterized by which units are activated. A message corresponds to one, and only one, state of the network. A message y of length c can be seen as an array x of length $c \cdot l$ where for each $j \in [1, c]$, the j -th symbol of y denoted y_j is represented by $x[(j-1)l+1 \dots jl]$.¹ To sum up the k -th unit in the j -th cluster corresponds to the index $i = (j-1)l + k$. We say the network is in state x if for all $j \in [1, c]$ and $k \in [1, l]$, unit k in cluster j of the network is active if, and only if, $x[(j-1)l + k] = 1$. Taking into account this correspondence between messages y in \mathcal{M} and their clustered representation x , we will now also consider x as messages.

Storing a message y in the network amounts to consider the network in state x ; then adding all edges between activated units to the network excepting for units sharing the same cluster. Notice that edges already present in the network remain unchanged throughout the storing process. Thus we obtain a multipartite clique between activated units representing a message.

We can describe the network which is an undirected unweighted graph, by its adjacency matrix W . For a message x , xx^T is the matrix describing the graph where all edges between active units exist, and only them. Therefore we can express W in function of the following :

$\bar{W} = \max_{x \in \mathcal{M}} xx^T = \bigvee_{x \in \mathcal{M}} xx^T$ which is the matrix where all edges between activated units for each state x in the same cluster have been added. W can now be expressed as :

$$W_{ii'} = \begin{cases} 0 & \text{if } \exists j \in [1, c], (j-1)l+1 \leq i, i' \leq jl \\ \bar{W}_{ii'} & \text{otherwise} \end{cases}$$

We obtain a binary symmetric matrix W with null blocks on the diagonal.

As logical or (\vee) is associative and commutative, storing can be done online and is independent of the order in which messages are stored.

B. Retrieving messages

Messages are retrieved by iterative application of a correcting rule. We adapt the rule first developed in [2] to multipartite cliques.

¹This is a clustered representation. x is separated in c clusters of l units.

We make a small modification the classic retrieving rule to fully exploit the local regularity of the code (constant weight of a). We keep activated all units that achieve a score equal or greater than the w -th score. (where w is a positive integer parameter). $w = 1$ amounts to taking the maximum and so is exactly the same as the classic retrieving rule.

Intuitively the choice of $w = a$ is optimal, as each cluster should eventually contain exactly a active units to at least represent a symbol in \mathcal{A} .

We consider an altered version \tilde{x} of a message x which we are trying to retrieve. We set $\tilde{x}^0 = \tilde{x}$. The network is in state \tilde{x} at the beginning of the retrieving. The state of the network at the beginning of the step t will be denoted by \tilde{x}_t . The message will be successfully retrieved if at the end of the algorithm, the network is in state x .

At each step t , we compute the scores, that is to mean the number of active units each unit is connected to, plus a memory effect γ accounting for units begin already active. The array of scores at step t is : $s_t = \gamma\tilde{x}_t + W\tilde{x}_t$

To detail the final step of the retrieving rule, we introduce a function select such as for an array v , $\text{select}(w, v)$ is equal to the w -th greatest element of v (counting repetitions).²

$$\tilde{x}_{t+1}[(j-1)c+k] = \begin{cases} 1 & \text{if } s_t[(j-1)c+k] \geq \text{select}(w, s_t[(j-1)l+1 \dots jl]) \\ 0 & \text{otherwise} \end{cases}$$

We can for example fix a maximum number of steps to the algorithm or stop at one point if $\tilde{x}_{t+1} = \tilde{x}_t$. Increasing the number of iterations can of course improve the retrieval rate.

We now study the performance of the new retrieving rule on distinct problems of retrieval : first partially erased messages, then erroneous messages.

III. RETRIEVAL PERFORMANCE

Messages will be considered as independent and identically distributed according to the uniform distribution. We compute theoretical error rates after one iteration of the retrieving rule. We compare those results to empirical error rates obtained through simulations. We also compare efficiencies of distinct networks for distinct values of the parameter a (weight of the constant weight code) in the sense of information theory and show improvements over the case where $a = 1$.

A. Retrieving partially erased messages

One can see a partial version of a message x as the result of erasures of symbols of x . Some symbols of x are replaced by the blank character \top . Erasures of a symbol in a message corresponds to an "erased" cluster in the network. An erased cluster is a cluster of units which are all inactive, that is to mean set to 0.

² $\text{select}(3, [3, 2, 1, 2, 0]) = 2$

1) *Analytical result:* We first compute the theoretical density of the network, that is to mean the probability for a random edge (i, i') to be present in the network (i.e. $W_{ii'} = 1$). We then express the theoretical error rate after one iteration of the retrieving rule in function of the density. All this computations (on error rates) are done assuming existences of edges between random units are independent, which is false but is a convenient approximation that gives results matching very closely results from simulations.

As the probability for i (and i') to be active in its cluster for one message x is $\frac{a}{l}$ providing messages are uniformly distributed, it follows that the probability for the edge (i, i') of not being added in the network for x is $1 - (a/l)^2$. Since messages are independent, the density d , which is the probability for an edge to be in the network, can be expressed as :

$$d = 1 - \left(1 - \left(\frac{a}{l}\right)^2\right)^m \quad (1)$$

, recalling that m is the number of messages in \mathcal{M} .

Suppose we are trying to retrieve a message where c_e clusters have been erased. Thanks to the memory effect γ , activated units in a non-erased cluster stay activated and are the only ones being so in their cluster (as the others can't achieve a higher score thanks to not being already activated). For units in erased clusters, the maximum attainable score is $a(c - c_e)$. Units in erased clusters corresponding to the original message achieve this score. Since scores follow a binomial law, the probability for a random unit to achieve the maximum score in such a cluster, that is to mean $a(c - c_e)$ is $d^{a(c - c_e)}$.

Since unit activations are independent (as messages are themselves independent), the probability for none of the vertices of one erased cluster, excepting the correct ones, to achieve the maximum is $(1 - d^{a(c - c_e)})^{l - a}$.

Scores in clusters being also independent, the probability for none of the vertices in any erased cluster, excepting the correct ones, to achieve the maximum in their respective clusters is $(1 - d^{a(c - c_e)})^{c_e(l - a)}$.

Whence error rate in retrieving messages is :

$$P_{err} = 1 - \left(1 - d^{a(c - c_e)}\right)^{c_e(l - a)} \quad (2)$$

In order to compare distinct networks using distinct alphabets (due to variations of the weight parameter a), we define the efficiency of the network relatively to the hardness of the problem (how much it is efficient compared to a perfect associative memory). A perfect associative memory would use exactly the same number of bits than the original messages and retrieve messages by maximum likelihood (ambiguities causing errors). Efficiency is the amount of information retrieved divided by the hardness of the problem.

The network is a binary symmetric matrix with null diagonal and can therefore be stored in memory as $\frac{c(c-1)l^2}{2}$ bits. Since messages are not stored in an ordered fashion, the entropy of the set of messages is $\log_2\left(\frac{\binom{l}{a}^c}{m!}\right) =_{m \rightarrow +\infty} m(c \log_2 \binom{l}{a} - \log_2(m) + \frac{1}{\log_2} + o(1))$ (thanks to Stirling's formula). The ratio

of information effectively stored in the network is $P_{retrieve} \times \frac{2m(c \log_2(\frac{l}{a}) - \log_2(m) + \frac{1}{\log 2})}{c(c-1)l^2}$.

Ambiguities would occur in a perfect associative memory if at least two messages were identical over non-erased letters. The probability it doesn't happen is $(1 - (\frac{l}{a})^{c-c_e})^{m-1}$.

$$\eta_m = P_{retrieve} \frac{2m(c \log_2(\frac{l}{a}) - \log_2(m) + \frac{1}{\log 2})}{c(c-1)l^2} \times (1 - (\frac{1}{a})^{c-c_e})^{-(m-1)}$$

The maximum ratio of information that can be stored in the network is therefore $\eta = \max_x \eta_x$, which we will define as the efficiency of the network.

2) *Simulations*: Theoretical and empirical densities match very closely, see Figure 1.

Simulations agree with the analytical results for error rate, albeit we can observe a small shift, see Figure 2. For one iteration, as long as the parameter w is between 1 and a , there is no effect on retrieving. However for 4 iterations, the choice of $w = a$ allows a shift of about 5000 messages stored on top of those stored for $w = 1$. The variant rule can thus be considered a significant improvement over the classic one for multipartite clique networks (recall that $w = 1$ amounts to the classic retrieving rule).

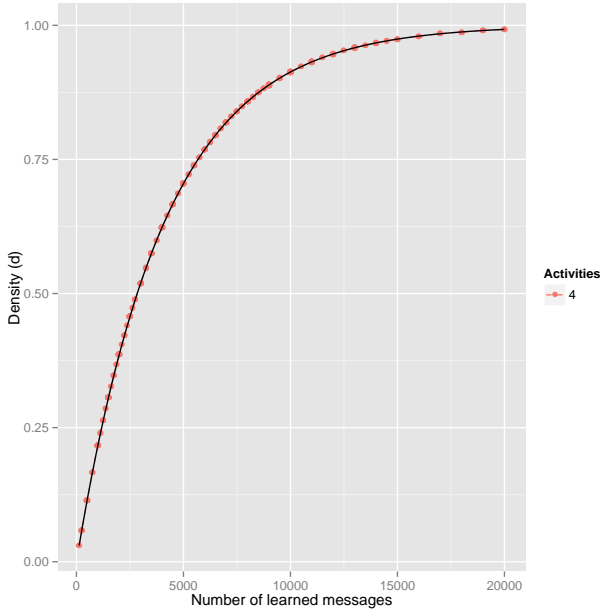


Figure 1. Theoretical and empirical densities for a network of 8 clusters, 256 units per cluster and 4 activated units per cluster

Thanks to the efficiency defined in the previous section, we can easily compare distinct choices of parameters a (we now take $w = a$ since it is the most efficient choice in practice) over distinct number of erased clusters. Figure 3 shows improvements in efficiencies over parameter $a = 1$. The more difficult the problem is, the more redundancies introduced by multipartite cliques are useful.

B. Retrieving corrupted messages

1) *Analytical result*: We provide a formula for error rate in case of messages with corrupted clusters.

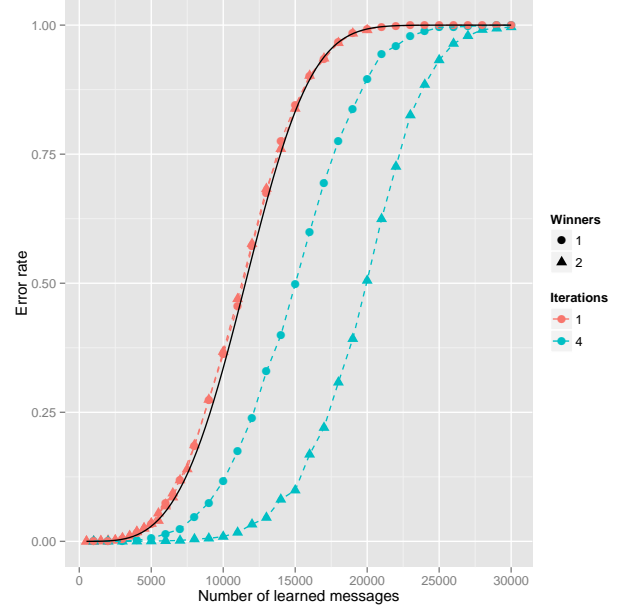


Figure 2. Evolution of error rate for increasing number of stored messages, 2 activities per cluster, 4 clusters, 512 units per cluster, two erasures, each point is the mean of 5 networks with 1000 sampled messages, analytical result (for one iteration) in continuous black

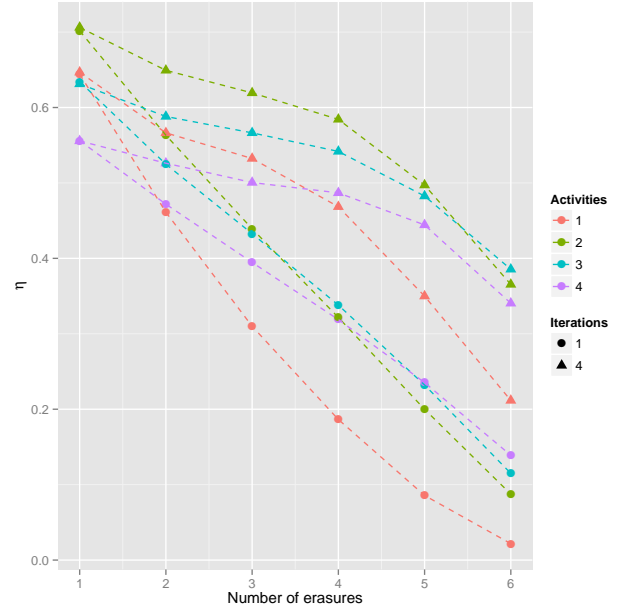


Figure 3. Comparison between different number of activities per cluster : 8 clusters, 256 units per cluster, 1000 sampled messages, pools of 5

Messages can be corrupted, our models for errors is : given c_e the number of erroneous clusters, we draw at random c_e clusters and in each of these draw at random, uniformly, a new letter in the alphabet that will replace the old one (it could be the same).

Knowing the density of the network, which follows formula (1), we can compute the laws of probabilities for scores achieved by different type of units in different clusters. We

can divide units in 5 groups. Subscript g will indicate a unit being part of the message (it should be activated), subscript b the contrary (the unit should not be activated). Subscript a means the unit is already activated. Subscript c corresponds to a correct cluster, e to an erroneous cluster. Correct units in uncorrupted clusters achieve at least a score of $a(c-c_e-1)+\gamma$. Being connected to a given activated unit in a corrupted cluster occurs with probability d , since errors are identically distributed. Moreover correct units see their scores shift, as do all activated units respectively thanks to being part of the original message and the memory effect γ . We can express those probability laws as : For example an activated good unit in a correct cluster starts with a base score of $a(c-c_e-1)+\gamma$. The probability it has to be connected to a random activated unit in an erroneous cluster is d . Thereby its law is :

$$P(n_{agc} = a(c-c_e-1)+\gamma+x) = \binom{a c_e}{x} d^x (1-d)^{a c_e - x}$$

Correct units in corrupted clusters achieve at least a score of $a(c-c_e)$. $P(n_{ge} = a(c-c_e)+x) = \binom{a(c-c_e-1)}{x} d^x (1-d)^{a(c-c_e-1)-x}$

Activated units benefits of the memory effect γ :

$$P(n_{abe} = \gamma+x) = \binom{a(c-1)}{x} d^x (1-d)^{a(c-1)-x}$$

Finally $P(n_{be} = x) = P(n_{bc} = x) = \binom{a(c-1)}{x} d^x (1-d)^{a(c-1)-x}$

A message is retrieved after one iteration, if and only if correct units achieve strictly the best scores.

The probability for all good units of a correct cluster (resp. of an erroneous cluster) to achieve a score equal or greater than x , with at least one having a score of x is $P_{ac}(x) = \sum_{k=1}^a \binom{a}{k} P(n_{agc} = x)^k P(n_{agc} > x)^{a-k}$ (resp. $P_e(x) = \sum_{k=1}^a \binom{a}{k} P(n_{ge} = x)^k P(n_{ge} > x)^{a-k}$)

$$P_{retrieve} = \left[\sum_{x=1}^{a(c-1)+\gamma} P_{ac}(x) P(n_{bc} < x)^{l-a} \right]^{c-c_e} \times \left[\sum_{x=1}^{a(c-1)+\gamma} P_e(x) P(n_{be} < x)^{l-2a} P(n_{abe} < x)^a \right]^{c_e} \quad (3)$$

with the approximation that there is few activated units that should be kept active in an erroneous cluster.

$$P_{err} = 1 - P_{retrieve}$$

2) *Simulations*: See Figure 4 compares the two rules. "a-winners take all" better than "winner takes all" with errors instead of erasures, even for one iteration.

$$\eta_m = P_{retrieve} \frac{2m(c \log_2 \binom{l}{a} - \log_2(m) + \frac{1}{\log 2})}{c(c-1)l^2} \times \left[\sum_{k=c_e+1}^c \left(1 - \frac{1}{\binom{l}{a}}\right)^k \left(\frac{1}{\binom{l}{a}}\right)^{c-k} \right]^{-(m-1)}$$

IV. RESILIENCE

A. Analytical result

As information on a message is shared across edges with multiple activated units per cluster, such networks are more resilient to damages (for one message information is carried by $a^2 \binom{c}{2}$ edges. v_c will indicate a correct unit that should be activated, v the other type of units.

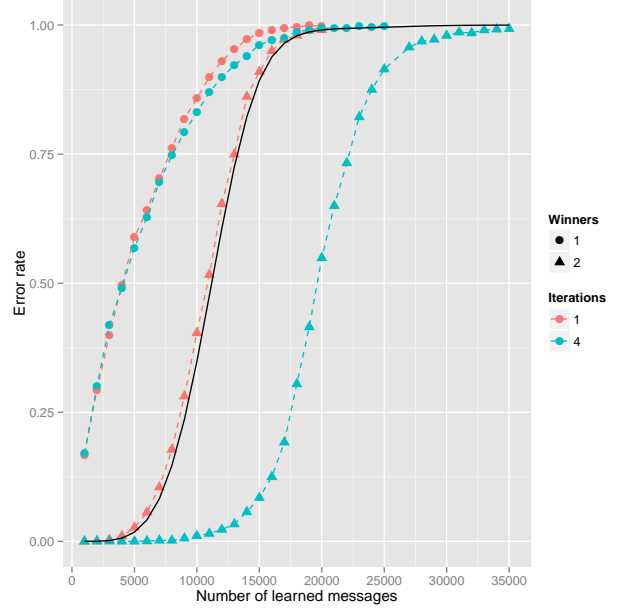


Figure 4. 2 activities per cluster, 4 clusters, 512 units per cluster, one erroneous cluster, each point is the mean of 5 networks with 1000 sampled messages, analytical result in continuous black

Our deviation model is the same as [4]³, that is random binary noise on edges. We generalise results of [4] to multi-partite cliques.

We assume $\gamma = \infty$ for one iteration as it prevents activated neurons to be deactivated. With probability $1-\psi$, one correct vertex is still connected to a given activated unit in an intact cluster. Thus $P(n_{vc} = x) = \binom{a(c_k-1)}{x} (1-\psi)^x \psi^{a(c-c_e-1)-x}$

The probability for an edge to be present in the network after damages is : $P_+ = \psi(1-d) + (1-\psi)d$. So for other units $P(n_v = x) = \binom{a(c-c_e)}{x} P_+^x (1-P_+)^{a(c-c_e)-x}$.

The probability for all correct units to achieve a score equal or greater than x , with at least one achieving a score of x is $P_c(x) = \sum_{k=1}^a \binom{a}{k} P(n_{vc} = x)^k P(n_{vc} > x)^{a-k}$

It follows that $P(\text{no other vertex activated in one cluster}) = \sum_{x=1}^{a(c-c_e)} P_c(x) P(n_v < x)^{l-a}$.

$$P_{retrieve} = \left(\sum_{x=1}^{a(c-c_e)} P_c(x) P(n_v < x)^{l-a} \right)^{c_e} \quad (4)$$

B. Simulations

With random noise on edges, hardness of the problem must take into account the capacity of the binary symmetric channel (the closer ψ is of 0.5, the harder the problem is). Our new efficiency is written as : $\eta_{m,\psi} = P_{retrieve} m \frac{2(c \log_2 \binom{l}{a} - \log_2(m) + \frac{1}{\log 2})}{c(c-1)l^2} \times (1 - \binom{l}{a}^{c-c_e})^{-(m-1)} \times \frac{1}{1 + \psi \log_2(\psi) + (1-\psi) \log_2(1-\psi)}$
 $\eta_\psi = \max_x \eta_{x,\psi}$

³Compared to this article, we do not draw at random a letter in case of ambiguity, which hurts performance a bit.

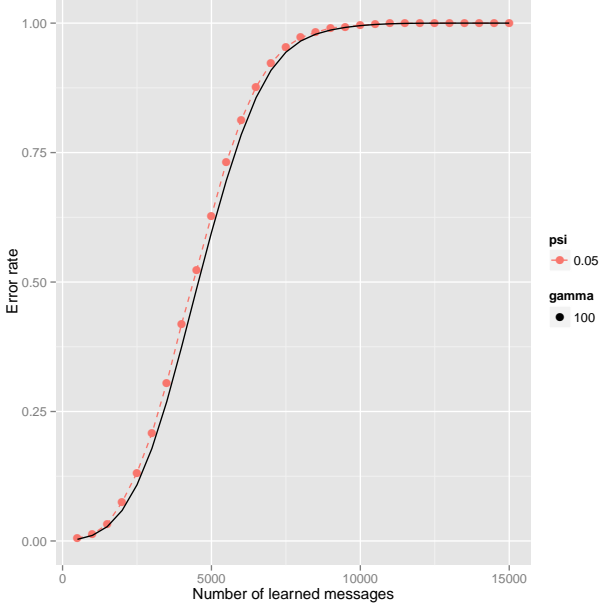


Figure 5. $\psi = 0.05$, 2 activities per cluster, 8 clusters, 256 units per cluster, 4 erasures, each point is the mean of 10 networks with 1000 sampled messages, analytical result in continuous black

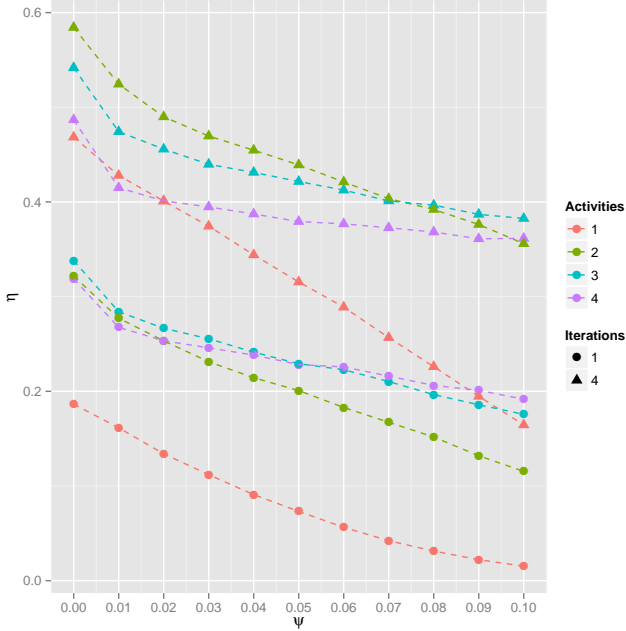


Figure 6. Comparison of evolutions of efficiency with increasing damages between different number of activities per cluster : 8 clusters, 256 units per cluster

We compute values represented on figure 6 by scanning on the number of messages, so as to find the maximum efficiency.

After a transition, efficiencies seem to follow a linear decreasing low according to ψ .

Absolute values of slopes diminish for increasing number of activities per cluster. We can conclude usage of multipartite cliques improve the resilience of the network.

Figure 6 seems to indicate a linear relationship (after a transition for multiple activities) between $\eta_\psi - \eta_0$ and ψ .

V. MAXIMUM THEORETICAL EFFICIENCY

It was shown in [5] that maximal asymptotic capacity of Willshaw networks is $\log 2$. We adapt this result to multipartite clique clustered networks.

The entropy of the messages set can be approximated as $H(\mathcal{M}) \approx mc \log_2 \binom{l}{a}$. From $d = 1 - \left(1 - \left(\frac{a}{l}\right)^2\right)^m$, we can deduce $m \sim \frac{\log_2(1-d)}{\log_2(1-(\frac{a}{l})^2)} \sim -\frac{l^2 \log_2(1-d) \log 2}{a^2}$. We set $\log_2 \binom{l}{a} = kc$ (otherwise density approaches 0 or 1, and efficiency 0). It follows that $\eta \sim \frac{mc \log_2 \binom{l}{a}}{l^2 \binom{c}{2}} \sim \frac{kc^2 \log_2(1-d) \log 2}{a^2 \binom{c}{2}}$.
 $\eta \sim -2ka^{-2} \log 2 \times \log_2(1-d)$.

We can bound the probability $P_{\exists e}$ for a message being "wrongly present" in the network, that is to mean adding it will not add any edge to the network albeit it is not in \mathcal{M} . By union-bound $P_{\exists e} \leq \binom{l}{a}^c d^{a^2 \binom{c}{2}} = 2^{c \log_2 \binom{l}{a} + a^2 \binom{c}{2} \log_2(d)} \approx 2^{c^2(k + \frac{a^2}{2} \log_2(d))}$.

Fixing $k^* = -\frac{a^2}{2} \log_2(d)$ gives us $P_{\exists e}$ and $\eta \sim \log_2 d \log_2(1-d) \log 2$ and therefore

$$\eta_{max} = \log 2$$

Multipartite cliques don't improve theoretical asymptotic efficiency (for not accepting messages that weren't previously seen by the network).

VI. CONCLUSION

REFERENCES

- [1] Vincent Gripon and Claude Berrou. A simple and efficient way to store many messages using neural cliques. In *Proceedings of IEEE Symposium on Computational Intelligence, Cognitive Algorithms, Mind, and Brain*, pages 54–58, Paris, France, April 2011.
- [2] Vincent Gripon and Claude Berrou. Sparse neural networks with large learning diversity. *IEEE Transactions on Neural Networks*, 22(7):1087–1096, July 2011.
- [3] Vincent Gripon and Claude Berrou. Nearly-optimal associative memories based on distributed constant weight codes. In *Proceedings of Information Theory and Applications Workshop*, pages 269–273, San Diego, CA, USA, February 2012.
- [4] François Leduc-Primeau, Vincent Gripon, Michael Rabbat, and Warren Gross. Cluster-based associative memories built from unreliable storage. In *ICASSP*, May 2014. To appear.
- [5] Günther Palm. On associative memories. 1980.