

---

# Sparse Modern Hopfield Networks

---

**André F. T. Martins**

Instituto de Telecomunicações  
Instituto Superior Técnico (Lisbon ELLIS Unit) & Unbabel  
Lisbon, Portugal  
andre.t.martins@tecnico.ulisboa.pt

**Vlad Niculae**

Informatics Institute  
University of Amsterdam  
Amsterdam, The Netherlands  
v.niculae@uva.nl

**Daniel McNamee**

Neuroscience Programme  
Champalimaud Research  
Lisbon, Portugal  
daniel.mcnamee@research.fchampalimaud.org

## Abstract

Ramsauer et al. (2021) recently pointed out a connection between modern Hopfield networks and attention heads in transformers. In this paper, we extend their framework to a broader family of energy functions which can be written as a difference of a quadratic regularizer and a Fenchel-Young loss (Blondel et al., 2020), parametrized by a generalized negentropy function  $\Omega$ . By working with Tsallis negentropies, the resulting update rules become end-to-end differentiable *sparse* transformations, establishing a new link to adaptively sparse transformers (Correia et al., 2019) and allowing for exact convergence to single memory patterns. Experiments on simulated data show a higher tendency to avoid metastable states.

## 1 Introduction

Hopfield networks are a kind of biologically-plausible neural network exhibiting associative memory capabilities (Hopfield, 1982). Their attractor dynamics makes them suitable for modeling the retrieval of episodic memories in humans and animals (Tyulmankov et al., 2021; Whittington et al., 2021). The limited storage capabilities of classical (quadratic energy) Hopfield networks were recently overcome through new energy functions and continuous state patterns (Krotov and Hopfield, 2016; Demircigil et al., 2017; Ramsauer et al., 2021), leading to exponential storage capacities, and sparking renewed interest in modern Hopfield networks. In particular, Ramsauer et al. (2021) revealed striking connections to transformer attention, via an update rule linked to the convex-concave procedure (CCCP; Yuille and Rangarajan 2003). However, this model has the often-undesirable tendency to converge to large metastable states (mixing many input patterns) instead of retrieving a single pattern.

There is a strong neurobiological motivation to seek new Hopfield energies capable of sparse selection of patterns. Sparse neural activity patterns are observed in electrophysiological recordings from many brain areas across a variety of animal species and forms a core principle of cortical computation due to their efficient coding properties (Simoncelli and Olshausen, 2001; Palm, 2013). With respect to memory formation circuits, the sparse firing of neurons in the dentate gyrus (DG), a major input pathway to the CA3 and CA1 subregions in the hippocampus, underpins its theorized role in pattern separation during memory storage (Yassa and Stark, 2011; Severa et al., 2017). Indeed, evidence suggests that the sparsified activity profiles of DG neurons aids in minimizing interference between competing memory patterns just prior to pattern completion via autoassociative dynamics downstream (Leutgeb et al., 2007; Neunuebel and Knierim, 2014).

In this paper, inspired by the framework of regularized prediction maps and Fenchel-Young losses (Blondel et al., 2020), we extend Ramsauer et al.’s (2021) energy function to a wider family induced by generalized entropies. The minimization of these energy functions leads to update rules which include as particular cases **sparsemax** (Martins and Astudillo, 2016) and the  $\alpha$ -**entmax transformations** used by adaptively sparse transformers (Peters et al., 2019; Correia et al., 2019). Unlike Ramsauer et al.’s (2021) Hopfield layers, our proposed update rules can lead to **sparse** convex combinations of input patterns and can have *exact* convergence to a fixed point in a small number of steps, maintaining end-to-end differentiability. Experiments on simulated data show frequent convergence to a *single* input pattern or small metastable states.

## 2 Modern Hopfield Networks

Let  $\mathbf{X} \in \mathbb{R}^{N \times D}$  be a matrix whose rows hold a set of examples  $\mathbf{x}_1, \dots, \mathbf{x}_N$  (“memory patterns”), where each  $\mathbf{x}_i \in \mathbb{R}^D$ , and let  $\mathbf{q}_0 \in \mathbb{R}^D$  be a query vector (or “state pattern”). The Hopfield network iteratively updates  $\mathbf{q}_t \mapsto \mathbf{q}_{t+1}$  for  $t \in \{0, 1, \dots\}$  according to a certain rule and, under certain conditions, these dynamical trajectories converge to a fixed point attractor state  $\mathbf{q}^*$  which either corresponds to one of the memorized examples, or to a mixture thereof. This update rule correspond to the minimization of an energy function, which for classic Hopfield networks (Hopfield, 1982) takes the form  $E(\mathbf{q}) = -\frac{1}{2} \mathbf{q}^\top \mathbf{W} \mathbf{q}$ , where  $\mathbf{W} = \mathbf{X}^\top \mathbf{X} \in \mathbb{R}^{D \times D}$  are the Hopfield network parameters, which when  $D \ll N$  can be seen as a “compressed memory”. In the classic Hopfield network, the state vector  $\mathbf{q}$  is further constrained to be  $\{\pm 1\}$ -valued, and the update rule is  $\mathbf{q}_{t+1} = \text{sign}(\mathbf{W} \mathbf{q}_t)$ . A limitation of this classical network is that it has only  $O(D)$  memory storage capacity.

Recent work sidestepped this limitation through alternative energy functions (Krotov and Hopfield, 2016; Demircigil et al., 2017), leading to the development of a class of models known as “modern Hopfield networks” with superlinear (often exponential) memory capacity. In Ramsauer et al. (2021), the state vector  $\mathbf{q} \in \mathbb{R}^D$  is continuous and unconstrained and the following energy is used:

$$E(\mathbf{q}) = -\text{lse}(\beta, \mathbf{X} \mathbf{q}) + \frac{1}{2} \mathbf{q}^\top \mathbf{q} + \beta^{-1} \log N + \frac{1}{2} M^2, \quad (1)$$

where  $M = \max_i \|\mathbf{x}_i\|$  and  $\text{lse}(\beta, \boldsymbol{\theta}) = \beta^{-1} \log \sum_{i=1}^N \exp(\beta \theta_i)$ .

Ramsauer et al. (2021) have revealed an interesting relation between the updates in this modern Hopfield network and the attention layers in transformers. Namely, the minimization of the energy (1) using the concave-convex procedure (CCCP; Yuille and Rangarajan 2003) leads to the update rule

$$\mathbf{q}_{t+1} = \mathbf{X}^\top \text{softmax}(\beta \mathbf{X} \mathbf{q}_t). \quad (2)$$

When  $\beta = \frac{1}{\sqrt{D}}$ , each update matches the computation performed in the attention layer of a transformer with a single attention head and with identity projection matrices. This triggered interest in developing variants of Hopfield layers which can be used as drop-in replacements for multi-head attention layers.

While Ramsauer et al. (2021) have derived useful theoretical properties of these networks (including their exponential storage capacity under some assumptions), the use of softmax in the update rule (2) prevents exact convergence and may lead to undesirable metastable states in some situations, as illustrated in toy experiments in §5. We explore in this paper a more general energy function which uses sparse transformations as an attempt to overcome these drawbacks.

## 3 $\Omega$ -Regularized Prediction Maps and Fenchel-Young Losses

Our contribution is rooted in the concept of Fenchel-Young losses, which we next review. Let  $\Delta_{N-1} := \{\mathbf{p} \in \mathbb{R}^N : \mathbf{p} \geq \mathbf{0}, \mathbf{1}^\top \mathbf{p} = 1\}$  denote the probability simplex, whose elements are probability vectors of length  $N$ . Given a convex function  $\Omega : \Delta_{N-1} \rightarrow \mathbb{R}$ , the  $\Omega$ -**regularized prediction map** ( $\Omega$ -RPM; Blondel et al. 2020),  $\hat{\mathbf{p}}_\Omega : \mathbb{R}^N \rightarrow \Delta_{N-1}$ , is:

$$\hat{\mathbf{p}}_\Omega(\boldsymbol{\theta}) = \arg \max_{\mathbf{p} \in \Delta_{N-1}} \boldsymbol{\theta}^\top \mathbf{p} - \Omega(\mathbf{p}). \quad (3)$$

One example of an  $\Omega$ -RPM is the **softmax** transformation, obtained when  $\Omega(\mathbf{p}) = \sum_{i=1}^N p_i \log p_i$  is the Shannon negentropy. Another example is the **sparsemax** transformation, obtained when  $\Omega(\mathbf{p}) =$

$\frac{1}{2}\|\mathbf{p}\|^2$  (Martins and Astudillo, 2016). The sparsemax corresponds to the Euclidean projection onto the probability simplex. Softmax and sparsemax are both particular cases of  $\alpha$ -**entmax transformations** (Peters et al., 2019), parametrized by a scalar  $\alpha \geq 0$  (called the entropic index), which correspond to the following choice of regularizer, called the **Tsallis  $\alpha$ -negentropy** (Tsallis, 1988):

$$\Omega_\alpha(\mathbf{p}) = (-1 + \|\mathbf{p}\|_\alpha^\alpha) / \alpha(\alpha - 1). \quad (4)$$

Note that, when  $\alpha \rightarrow 1$ ,  $\Omega_\alpha$  becomes Shannon’s negentropy and the corresponding  $\Omega$ -RPM is the softmax, and when  $\alpha = 2$ , it becomes the  $\ell_2$ -norm (up to a constant) and we recover the sparsemax.

Let  $\Omega^*$  be the convex conjugate of  $\Omega$ ,  $\Omega^*(\boldsymbol{\theta}) = \max_{\mathbf{p} \in \Delta_{N-1}} \boldsymbol{\theta}^\top \mathbf{p} - \Omega(\mathbf{p})$ . The  $\Omega$ -RPM in (3) equals the gradient map of  $\Omega^*$ ,  $\hat{\mathbf{p}}_\Omega(\boldsymbol{\theta}) = \nabla \Omega^*(\boldsymbol{\theta})$ . Note also that we have  $\Omega^*(\boldsymbol{\theta}) = \boldsymbol{\theta}^\top \hat{\mathbf{p}}_\Omega(\boldsymbol{\theta}) - \Omega(\hat{\mathbf{p}}_\Omega(\boldsymbol{\theta}))$ . The **Fenchel-Young loss** induced by  $\Omega$  (Blondel et al., 2020) is the function defined as

$$L_\Omega(\boldsymbol{\theta}, \mathbf{p}) = \Omega(\mathbf{p}) + \Omega^*(\boldsymbol{\theta}) - \boldsymbol{\theta}^\top \mathbf{p}. \quad (5)$$

When  $\Omega$  is the Shannon negentropy,  $\Omega^*(\boldsymbol{\theta}) = \text{lse}(1, \boldsymbol{\theta})$ , and  $L_\Omega$  is the **cross-entropy loss**, up to a constant (Blondel et al., 2020, §3.2). Intuitively, Fenchel-Young losses quantify how “compatible” a score vector  $\boldsymbol{\theta} \in \mathbb{R}^N$  (e.g., logits) is to a desired probability vector  $\mathbf{p} \in \Delta_{N-1}$ . Additionally, (Blondel et al., 2020, Prop. 2):

1. Fenchel-Young losses are non-negative,  $L_\Omega(\boldsymbol{\theta}, \mathbf{p}) \geq 0$ , with equality iff  $\mathbf{p} = \hat{\mathbf{p}}_\Omega(\boldsymbol{\theta})$ .
2. Fenchel-Young losses are convex on  $\boldsymbol{\theta}$  and their gradient is  $\nabla_{\boldsymbol{\theta}} L_\Omega(\boldsymbol{\theta}, \mathbf{p}) = -\mathbf{p} + \hat{\mathbf{p}}_\Omega(\boldsymbol{\theta})$ .

For Tsallis negentropies  $\Omega_\alpha$  with  $\alpha > 1$ , a margin property holds (Blondel et al., 2020, Prop. 7):

$$\forall i \in [N], \quad L_{\Omega_\alpha}(\boldsymbol{\theta}, \mathbf{e}_i) = 0 \iff \hat{\mathbf{p}}_{\Omega_\alpha}(\boldsymbol{\theta}) = \mathbf{e}_i \iff \theta_i - \max_{j \neq i} \theta_j \geq (\alpha - 1)^{-1}. \quad (6)$$

We will use these properties in the next section to define and analyze sparse Hopfield networks.

## 4 Sparse Hopfield Networks

We now use  $\Omega$ -RPMs and Fenchel-Young losses to define a new class of energy functions associated to modern Hopfield networks. We assume that the regularizer  $\Omega$  is a **generalized negentropy**, i.e., null when  $\mathbf{p}$  is peaked, strictly convex, and permutation-invariant (see Appendix A). These conditions imply that  $\Omega \leq 0$  and that  $\Omega$  is minimized when  $\mathbf{p} = \mathbf{1}/N$  is the uniform distribution (Blondel et al., 2020, Prop. 4). The Tsallis negentropy (4) satisfies these properties for  $\alpha \geq 1$ .

We define the **Hopfield-Fenchel-Young energy** as

$$E(\mathbf{q}) = \underbrace{-\beta^{-1} L_\Omega(\beta \mathbf{X} \mathbf{q}; \mathbf{1}/N)}_{E_{\text{concave}}(\mathbf{q})} + \underbrace{\frac{1}{2} \|\mathbf{q} - \boldsymbol{\mu}_X\|^2 + \frac{1}{2} (M^2 - \|\boldsymbol{\mu}_X\|^2)}_{E_{\text{convex}}(\mathbf{q})}, \quad (7)$$

where  $\boldsymbol{\mu}_X := \mathbf{X}^\top \mathbf{1}/N \in \mathbb{R}^D$  is the empirical mean of the patterns. This energy extends that of (1), which is recovered when  $\Omega$  is Shannon’s negentropy, in which case  $E_{\text{concave}}(\mathbf{q}) = -\text{lse}(\beta, \mathbf{X} \mathbf{q}) + \beta^{-1} \log N + \mathbf{q}^\top \boldsymbol{\mu}_X$ . The concavity of  $E_{\text{concave}}$  holds from the convexity of Fenchel-Young losses on its first argument and from the fact that composition of a convex function with an affine map is convex. The convexity of  $E_{\text{convex}}$  comes from the fact that it is a quadratic function.<sup>1</sup>

There are two terms competing when minimizing the energy function (7) with respect to  $\mathbf{q}$ :

- Minimizing  $E_{\text{concave}}$  is equivalent to *maximizing*  $L_\Omega(\beta \mathbf{X} \mathbf{q}; \mathbf{1}/N)$ , which pushes for state patterns  $\mathbf{q}$  as far from possible from a uniform average and close to a single memory pattern.
- Minimizing  $E_{\text{convex}}$  serves as a regularization, encouraging the state pattern to stay close to  $\boldsymbol{\mu}_X$ .

The next result, proved in Appendix B, establishes bounds and derives the Hopfield update rule for our energy function, generalizing Ramsauer et al. (2021, Lemma A.1 and Theorem A.1).

<sup>1</sup>Up to constants, for  $\Omega_1$  this is the same convex-concave decomposition of Ramsauer et al. (2021)

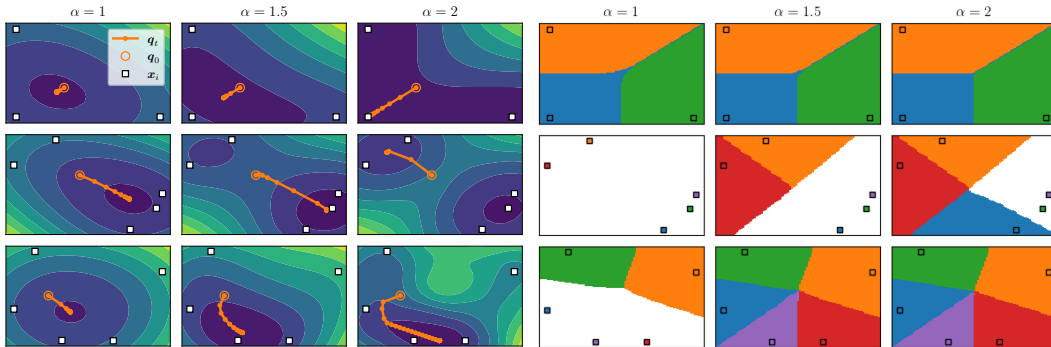


Figure 1: Left: contours of the energy function and optimization trajectory of the CCCP iteration ( $\beta = 1$ ). Right: attraction basins associated with each pattern. (White sections do not converge to a single pattern but to a metastable state;  $\beta = 10$ ; for  $\alpha = 1$  we allow a tolerance of  $\epsilon = .01$ ).

**Proposition 1.** *Let the query  $\mathbf{q}$  be in the convex hull of the rows of  $\mathbf{X}$ , i.e.,  $\mathbf{q} = \mathbf{X}^\top \mathbf{p}$  for some  $\mathbf{p} \in \Delta_{N-1}$ . Then, the energy (7) satisfies  $0 \leq E(\mathbf{q}) \leq \min \{2M^2, -\beta^{-1}\Omega(\mathbf{1}/N) + \frac{1}{2}M^2\}$ . Furthermore, minimizing (7) by CCCP (Yuille and Rangarajan, 2003) leads to the updates:*

$$\mathbf{q}_{t+1} = \mathbf{X}^\top \hat{\mathbf{p}}_\Omega(\beta \mathbf{X} \mathbf{q}_t). \quad (8)$$

*In particular, when  $\Omega = \Omega_\alpha$  (the Tsallis  $\alpha$ -entropy (4)), the  $\Omega$ -RPM is the  $\alpha$ -entmax transformation, corresponding to the adaptively sparse transformer of Correia et al. (2019).*

We now show that, with  $\Omega = \Omega_\alpha$  and for  $\alpha > 1$  (the sparse case), the memory patterns can be stationary points of the energy (7). This result is stronger than that of Ramsauer et al. (2021) for their energy (which is ours for  $\alpha = 1$ ), according to which memory patterns are only  $\epsilon$ -close to stationary points, where a small  $\epsilon = O(\exp(-\beta))$  requires a low temperature (large  $\beta$ ). The proof (Appendix C) relies on the margin property stated in (6). Following Ramsauer et al. (2021, Def. 2), we define the separation of pattern  $\mathbf{x}_i$  from data as  $\Delta_i = \mathbf{x}_i^\top \mathbf{x}_i - \max_{j \neq i} \mathbf{x}_i^\top \mathbf{x}_j$ .

**Proposition 2.** *Assume  $\Omega = \Omega_\alpha$  with  $\alpha > 1$ , and let  $\mathbf{x}_i$  be a memory pattern outside the convex hull of the other memory patterns. Then,  $\mathbf{x}_i$  is a stationary point of the energy (7) iff  $\Delta_i \geq \frac{1}{(\alpha-1)\beta}$ . In addition, if the initial query  $\mathbf{q}_0$  satisfies  $\mathbf{q}_0^\top (\mathbf{x}_i - \mathbf{x}_j) \geq \frac{1}{(\alpha-1)\beta}$  for all  $j \neq i$ , then the update rule (8) converges to  $\mathbf{x}_i$  exactly in one iteration. Moreover, if the patterns are normalized and  $\Delta_i \geq \frac{1}{(\alpha-1)\beta} + 2\epsilon$ , then any  $\mathbf{q}_0$   $\epsilon$ -close to  $\mathbf{x}_i$  ( $\|\mathbf{q}_0 - \mathbf{x}_i\| \leq \epsilon$ ) will converge to  $\mathbf{x}_i$  in one iteration.*

We next validate these results on simulated data.

## 5 Experiments

We repeatedly generate random patterns on a unit sphere and assess the frequency with which the update rule (8) leads to metastable states of different sizes (Table 1). Figure 1 shows optimization trajectories for several queries and pattern configurations, along with the basins of attraction for the three methods (a larger  $\beta$  is needed to allow the  $\alpha = 1$  model to get  $\epsilon$ -close to a single pattern). We use  $\alpha \in \{1, 1.5, 2\}$  since for those cases the  $\Omega_\alpha$ -RPM admits an exact algorithm (Peters et al., 2019).<sup>2</sup> Additional plots are shown in Appendix D. Overall, we observe that, with  $\alpha = 1$  (which corresponds to Ramsauer et al. (2021)) large metastable states abound, whereas for larger  $\alpha$  many updates converge exactly to a single pattern.

Table 1: Distribution of metastable state cardinalities (in %), for uniform patterns on a unit sphere ( $N = 10, D = 5$ ), and a uniform query on the unit ball. (Estimated using 1000 random trials;  $\beta = 4$ ; for  $\alpha = 1$  we threshold at 0.01).

$\ \mathbf{p}\ _0$	$\alpha = 1$	$\alpha = 1.5$	$\alpha = 2$
1	0.0	23.9	72.3
2	0.0	44.4	26.7
3	0.5	25.0	1.0
4	1.9	5.9	0.0
5	9.6	0.8	0.0
6	20.0	0.0	0.0
7	23.5	0.0	0.0
8	25.7	0.0	0.0
9	15.4	0.0	0.0
10	3.4	0.0	0.0

<sup>2</sup>In practice any  $\alpha > 1$  can be used with bisection, but we did not want approximations to affect our results.

## 6 Conclusions and Related Work

We proposed new Hopfield energies, linked to Fenchel-Young losses, which generalize the framework of Ramsauer et al. (2021). For suitable choices of generalized entropy, linked to Tsallis entropies, the resulting update equations mimic the sparse attention mechanism of adaptively sparse transformers (Correia et al., 2019). Our proposed models have interesting properties, including *exact* convergence to single patterns (not just to a nearby region), while maintaining end-to-end differentiability. We provide theoretical conditions for convergence in one iteration, along with toy experiments highlighting the usefulness of the new energies.

Concurrently to our work, Hu et al. (2023) recently proposed a model for sparse modern Hopfield networks along with a memory retrieval error bound provably tighter than the dense analog of Ramsauer et al. (2021). Their energy can be seen as a particular case of our Hopfield-Fenchel-Young energy, specifically the  $\alpha = 2$  case (sparsemax). Our work differs in that we consider the more general scenario where  $\alpha > 1$  ( $\alpha$ -entmax), we make a connection to Fenchel-Young losses, and we use the margin property of  $\alpha$ -entmax to prove exact convergence to single patterns under the conditions of Proposition 2. Future work will examine the suitability of the proposed approach to real-world problems, by exploring the use of learnable sparse Hopfield layers in concrete applications.

## Acknowledgments

This work was supported by the European Research Council (DECOLLAGE, ERC-2022-CoG 101088763), by the Portuguese Recovery and Resilience Plan through project C645008882-00000055 (Center for Responsible AI), by Fundação para a Ciência e Tecnologia through contract UIDB/50008/2020, and by the Dutch Research Council (NWO) via VI.Veni.212.228.

## References

- Blondel, M., Martins, A. F., and Niculae, V. (2020). Learning with Fenchel-Young losses. *The Journal of Machine Learning Research*, 21(1):1314–1382.
- Correia, G. M., Niculae, V., and Martins, A. F. (2019). Adaptively sparse transformers. In *Proceedings of EMNLP-IJCNLP*.
- Demircigil, M., Heusel, J., Löwe, M., Uppang, S., and Vermet, F. (2017). On a model of associative memory with huge storage capacity. *Journal of Statistical Physics*, 168:288–299.
- Hopfield, J. J. (1982). Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the National Academy of Sciences*, 79(8):2554–2558.
- Hu, J. Y.-C., Yang, D., Wu, D., Xu, C., Chen, B.-Y., and Liu, H. (2023). On sparse modern hopfield model. *arXiv preprint arXiv:2309.12673*.
- Krotov, D. and Hopfield, J. J. (2016). Dense associative memory for pattern recognition. *Advances in neural information processing systems*, 29.
- Leutgeb, J., Leutgeb, S., Moser, M., and Moser, E. (2007). Pattern separation in the dentate gyrus and CA3 of the hippocampus. *Science*, 315(5814):961–966.
- Martins, A. and Astudillo, R. (2016). From softmax to sparsemax: A sparse model of attention and multi-label classification. In *Proceedings of ICML*.
- Neunuebel, J. and Knierim, J. (2014). CA3 retrieves coherent representations from degraded input: direct evidence for CA3 pattern completion and dentate gyrus pattern separation. *Neuron*, 81(2):416–427.
- Palm, G. (2013). Neural associative memories and sparse coding. *Neural Netw*, 37:165–171.
- Peters, B., Niculae, V., and Martins, A. F. (2019). Sparse sequence-to-sequence models. In *Proceedings of ACL*.

- Ramsauer, H., Schäfl, B., Lehner, J., Seidl, P., Widrich, M., Gruber, L., Holzleitner, M., Adler, T., Kreil, D., Kopp, M. K., Klambauer, G., Brandstetter, J., and Hochreiter, S. (2021). Hopfield networks is all you need. In *Proceedings of ICLR*.
- Severa, W., Parekh, O., James, C., and Aimone, J. (2017). A combinatorial model for dentate gyrus sparse coding. *Neural Computation*, 29(1):94–117.
- Simoncelli, E. P. and Olshausen, B. A. (2001). Natural image statistics and neural representation. *Annual Review of Neuroscience*, 24(1):1193–1216.
- Tsallis, C. (1988). Possible generalization of boltzmann-gibbs statistics. *Journal of Statistical Physics*, 52:479–487.
- Tyulmankov, D., Fang, C., Vadaparty, A., and Yang, G. R. (2021). Biological learning in key-value memory networks. *Advances in Neural Information Processing Systems*.
- Whittington, J. C., Warren, J., and Behrens, T. E. (2021). Relating transformers to models and neural representations of the hippocampal formation. In *Proceedings of ICLR*.
- Yassa, M. and Stark, C. (2011). Pattern separation in the hippocampus. *Trends Neurosci*, 34(10):515–525.
- Yuille, A. L. and Rangarajan, A. (2003). The concave-convex procedure. *Neural computation*, 15(4):915–936.

## A Generalized negentropies

**Definition 1.** (Blondel et al., 2020, §4.1). A function  $\Omega : \Delta_{N-1} \rightarrow \mathbb{R}$  is a generalized negentropy iff it satisfies the following properties:

1. *Zero negentropy:*  $\Omega(\mathbf{p}) = 0$  if  $\mathbf{p}$  is a one-hot vector (delta distribution), i.e.,  $\mathbf{p} = \mathbf{e}_i$  for any  $i \in \{1, \dots, N\}$ .
2. *Strict convexity:*  $\Omega((1-\alpha)\mathbf{p} + \alpha\mathbf{p}') < (1-\alpha)\Omega(\mathbf{p}) + \alpha\Omega(\mathbf{p}')$ .
3. *Permutation invariance:*  $\Omega(\mathbf{P}\mathbf{p}) = \Omega(\mathbf{p})$  for any permutation matrix  $\mathbf{P}$  (i.e., square matrices with a single 1 in each row and each column, zero elsewhere).

## B Proof of Proposition 1

We start by proving that  $E(\mathbf{q}) \geq 0$ . We show first that for any  $\Omega$  satisfying conditions 1–3 above, we have

$$L_\Omega(\boldsymbol{\theta}; \mathbf{1}/N) \leq \max_i \theta_i - \mathbf{1}^\top \boldsymbol{\theta}/N. \quad (9)$$

From the definition of  $\Omega^*$  and the fact that  $\Omega(\mathbf{p}) \geq \Omega(\mathbf{1}/N)$  for any  $\mathbf{p} \in \Delta_{N-1}$ , we have that, for any  $\boldsymbol{\theta}$ ,  $\Omega^*(\boldsymbol{\theta}) = \max_{\mathbf{p} \in \Delta_{N-1}} \boldsymbol{\theta}^\top \mathbf{p} - \Omega(\mathbf{p}) \leq \max_{\mathbf{p} \in \Delta_{N-1}} \boldsymbol{\theta}^\top \mathbf{p} - \Omega(\mathbf{1}/N) \leq \max_i \theta_i - \Omega(\mathbf{1}/N)$ , which leads to (9).

Let now  $k = \arg \max_i \mathbf{q}^\top \mathbf{x}_i$ , i.e.,  $\mathbf{x}_k$  is the pattern most similar to the query  $\mathbf{q}$ . Therefore, we have

$$\begin{aligned} E(\mathbf{q}) &= -\beta^{-1} L_\Omega(\beta \mathbf{X} \mathbf{q}; \mathbf{1}/N) + \frac{1}{2} \|\mathbf{q} - \boldsymbol{\mu}_X\|^2 + \frac{1}{2} (M^2 - \|\boldsymbol{\mu}_X\|^2) \\ &\geq -\beta^{-1} (\beta \max_i \mathbf{q}^\top \mathbf{x}_i - \beta \mathbf{1}^\top \mathbf{X} \mathbf{q}/N) + \frac{1}{2} \|\mathbf{q} - \boldsymbol{\mu}_X\|^2 + \frac{1}{2} (M^2 - \|\boldsymbol{\mu}_X\|^2) \\ &= -\mathbf{q}^\top \mathbf{x}_k + \mathbf{q}^\top \boldsymbol{\mu}_X + \frac{1}{2} \|\mathbf{q} - \boldsymbol{\mu}_X\|^2 + \frac{1}{2} (M^2 - \|\boldsymbol{\mu}_X\|^2) \\ &= -\mathbf{q}^\top \mathbf{x}_k + \frac{1}{2} \|\mathbf{q}\|^2 + \frac{1}{2} M^2 \\ &\geq -\mathbf{q}^\top \mathbf{x}_k + \frac{1}{2} \|\mathbf{q}\|^2 + \frac{1}{2} \|\mathbf{x}_k\|^2 \\ &= \frac{1}{2} \|\mathbf{x}_k - \mathbf{q}\|^2 \geq 0. \end{aligned}$$

The zero value of energy is attained when  $\mathbf{X} = \mathbf{1}\mathbf{q}^\top$  (all patterns are equal to the query), in which case  $\boldsymbol{\mu}_X = \mathbf{q}$ ,  $M = \|\mathbf{q}\| = \|\boldsymbol{\mu}_X\|$ , and we get  $E_{\text{convex}}(\mathbf{q}) = E_{\text{concave}}(\mathbf{q}) = 0$ .

Now we prove the two upper bounds. For that, note that, for any  $\mathbf{p} \in \Delta_{N-1}$ , we have  $0 \leq L_\Omega(\boldsymbol{\theta}, \mathbf{p}) = L_\Omega(\boldsymbol{\theta}, \mathbf{1}/N) - \Omega(\mathbf{1}/N) + (\mathbf{p} - \mathbf{1}/N)^\top \boldsymbol{\theta} \leq L_\Omega(\boldsymbol{\theta}, \mathbf{1}/N) - \Omega(\mathbf{1}/N) - (\mathbf{p} - \mathbf{1}/N)^\top \boldsymbol{\theta}$ , due to the assumptions 1–3 which ensure  $\Omega$  is non-positive. That is,  $L_\Omega(\boldsymbol{\theta}, \mathbf{1}/N) \geq \Omega(\mathbf{1}/N) + (\mathbf{p} - \mathbf{1}/N)^\top \boldsymbol{\theta}$ . Therefore, with  $\mathbf{q} = \mathbf{X}^\top \mathbf{p}$ , we get

$$E_{\text{concave}}(\mathbf{q}) \leq -\beta^{-1} \Omega(\mathbf{1}/N) - \mathbf{p}^\top \mathbf{X} \mathbf{q} + \mathbf{q}^\top \boldsymbol{\mu}_X = -\beta^{-1} \Omega(\mathbf{1}/N) - \|\mathbf{q}\|^2 + \mathbf{q}^\top \boldsymbol{\mu}_X,$$

$$\text{and } E(\mathbf{q}) = E_{\text{concave}}(\mathbf{q}) + E_{\text{convex}}(\mathbf{q}) \leq -\beta^{-1} \Omega(\mathbf{1}/N) - \|\mathbf{q}\|^2 + \mathbf{q}^\top \boldsymbol{\mu}_X + \frac{1}{2} \|\mathbf{q} - \boldsymbol{\mu}_X\|^2 + \frac{1}{2} (M^2 - \|\boldsymbol{\mu}_X\|^2) = -\beta^{-1} \Omega(\mathbf{1}/N) - \frac{1}{2} \|\mathbf{q}\|^2 + \frac{1}{2} M^2 \leq -\beta^{-1} \Omega(\mathbf{1}/N) + \frac{1}{2} M^2.$$

To show the second upper bound, use the fact that  $E_{\text{concave}}(\mathbf{q}) \leq 0$ , which leads to  $E(\mathbf{q}) \leq E_{\text{convex}}(\mathbf{q}) = \frac{1}{2} \|\mathbf{q} - \boldsymbol{\mu}_X\|^2 + \frac{1}{2} (M^2 - \|\boldsymbol{\mu}_X\|^2) = \frac{1}{2} \|\mathbf{q}\|^2 - \mathbf{q}^\top \boldsymbol{\mu}_X + \frac{1}{2} M^2$ . Note that  $\|\mathbf{q}\| = \|\mathbf{X}^\top \mathbf{p}\| \leq \sum_i p_i \|\mathbf{x}_i\| \leq M$  and that, from the Cauchy-Schwarz inequality, we have  $-\mathbf{q}^\top \boldsymbol{\mu}_X \leq \|\boldsymbol{\mu}_X\| \|\mathbf{q}\| \leq M^2$ . Therefore, we obtain  $E(\mathbf{q}) \leq \frac{1}{2} \|\mathbf{q}\|^2 - \mathbf{q}^\top \boldsymbol{\mu}_X + \frac{1}{2} M^2 \leq \frac{1}{2} M^2 + M^2 + \frac{1}{2} M^2 = 2M^2$ .

We now turn to the update rule. The CCCP algorithm works as follows: at the  $t^{\text{th}}$  iteration, it linearizes the concave function  $E_{\text{concave}}$  by using a first-order Taylor approximation around  $\mathbf{q}_t$ ,

$$E_{\text{concave}}(\mathbf{q}) \approx \tilde{E}_{\text{concave}}(\mathbf{q}) := E_{\text{concave}}(\mathbf{q}_t) + \left( \frac{\partial E_{\text{concave}}(\mathbf{q}_t)}{\partial \mathbf{q}} \right)^\top (\mathbf{q} - \mathbf{q}_t).$$

Then, it computes a new iterate by solving the convex optimization problem  $\mathbf{q}_{t+1} := \arg \min_{\mathbf{q}} E_{\text{convex}}(\mathbf{q}) + \tilde{E}_{\text{concave}}(\mathbf{q})$ , which leads to the equation  $\nabla E_{\text{convex}}(\mathbf{q}_{t+1}) = -\nabla E_{\text{concave}}(\mathbf{q}_t)$ . Using the fact that  $\nabla L_{\Omega}(\boldsymbol{\theta}, \mathbf{p}) = \hat{\mathbf{p}}_{\Omega}(\boldsymbol{\theta}) - \mathbf{p}$  and the chain rule leads to

$$\begin{aligned} \nabla E_{\text{concave}}(\mathbf{q}) &= -\beta^{-1} \nabla_{\mathbf{q}} L_{\Omega}(\beta \mathbf{X} \mathbf{q}; \mathbf{1}/N) = \mathbf{X}^{\top} (\mathbf{1}/N - \hat{\mathbf{p}}_{\Omega}(\beta \mathbf{X} \mathbf{q})) \\ &= \boldsymbol{\mu}_{\mathbf{X}} - \mathbf{X}^{\top} \hat{\mathbf{p}}_{\Omega}(\beta \mathbf{X} \mathbf{q}) \\ \nabla E_{\text{convex}}(\mathbf{q}) &= \mathbf{q} - \boldsymbol{\mu}_{\mathbf{X}}, \end{aligned} \tag{10}$$

giving the update equation (8).

## C Proof of Proposition 2

A stationary point is a solution of the equation  $-\nabla E_{\text{concave}}(\mathbf{q}) = \nabla E_{\text{convex}}(\mathbf{q})$ . Using the expression for gradients (10), this is equivalent to  $\mathbf{q} = \mathbf{X}^{\top} \hat{\mathbf{p}}_{\Omega}(\beta \mathbf{X} \mathbf{q})$ . If  $\mathbf{x}_i = \mathbf{X}^{\top} \mathbf{e}_i$  is not a convex combination of the other memory patterns,  $\mathbf{x}_i$  is a stationary point iff  $\hat{\mathbf{p}}_{\Omega}(\beta \mathbf{X} \mathbf{x}_i) = \mathbf{e}_i$ . We now use the margin property of  $\alpha$ -entmax transformations (6), according to which the latter is equivalent to  $\beta \mathbf{x}_i^{\top} \mathbf{x}_i - \max_{j \neq i} \beta \mathbf{x}_i^{\top} \mathbf{x}_j \geq \frac{1}{\alpha-1}$ . Noting that the left hand side equals  $\beta \Delta_i$  leads to the desired result.

If the initial query satisfies  $\mathbf{q}_0^{\top} (\mathbf{x}_i - \mathbf{x}_j) \geq \frac{1}{(\alpha-1)\beta}$  for all  $j \neq i$ , we have again from the margin property that  $\hat{\mathbf{p}}_{\Omega}(\beta \mathbf{X} \mathbf{q}_0) = \mathbf{e}_i$ , which combined to the previous claim ensures convergence in one step to  $\mathbf{x}_i$ .

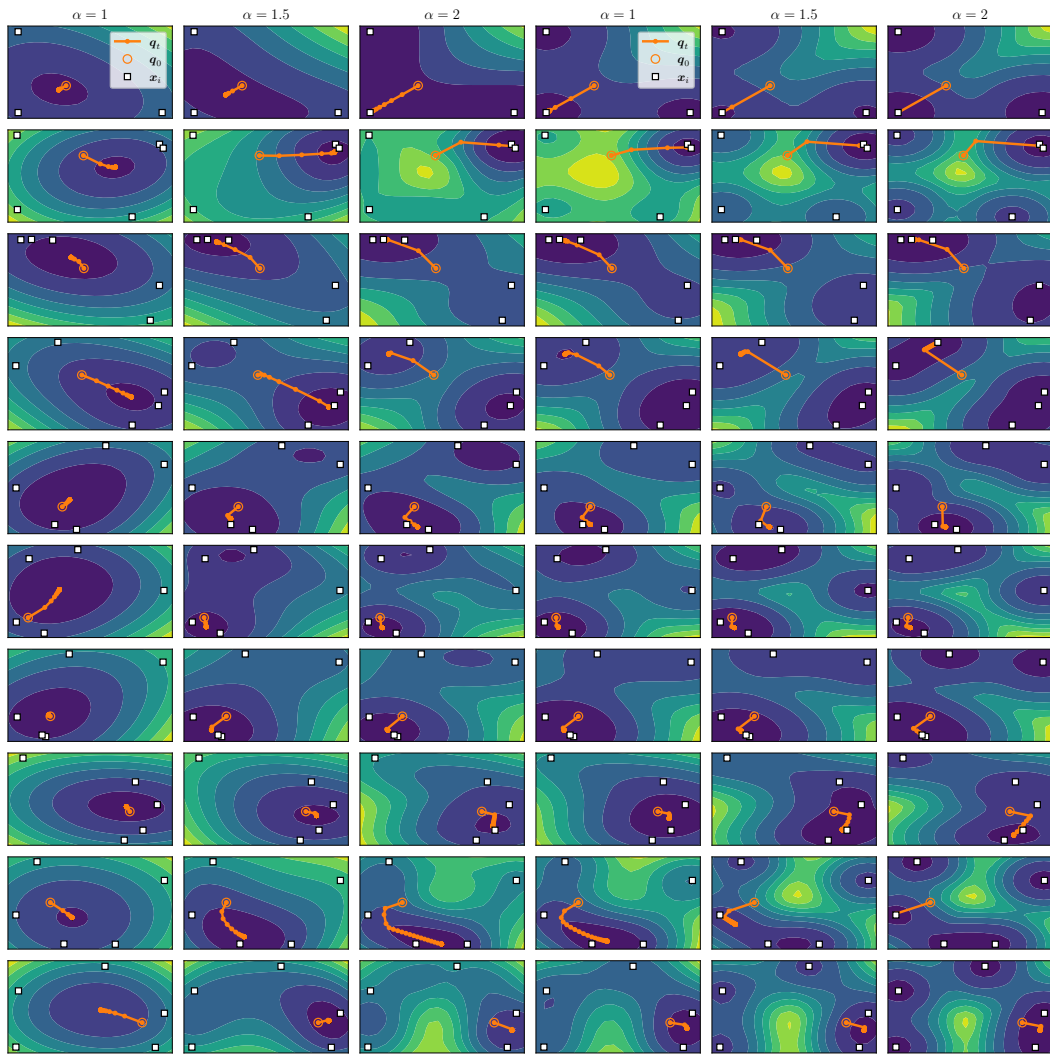
Finally, note that, if  $\mathbf{q}_0$  is  $\epsilon$ -close to  $\mathbf{x}_i$ , we have  $\mathbf{q}_0 = \mathbf{x}_i + \epsilon \mathbf{r}$  for some vector  $\mathbf{r}$  with  $\|\mathbf{r}\| = 1$ . Therefore, we have

$$\begin{aligned} \mathbf{q}_0^{\top} (\mathbf{x}_i - \mathbf{x}_j) &= (\mathbf{x}_i + \epsilon \mathbf{r})^{\top} (\mathbf{x}_i - \mathbf{x}_j) \\ &\geq \Delta_i + \epsilon \mathbf{r}^{\top} (\mathbf{x}_i - \mathbf{x}_j) \\ &\geq \Delta_i - \epsilon \underbrace{\|\mathbf{r}\|}_{=1} \|\mathbf{x}_i - \mathbf{x}_j\|, \end{aligned} \tag{11}$$

where we invoked the Cauchy-Schwarz inequality in the last step. Since the patterns are normalized, we have from the triangle inequality that  $\|\mathbf{x}_i - \mathbf{x}_j\| \leq \|\mathbf{x}_i\| + \|\mathbf{x}_j\| = 2$ ; using the assumption that  $\Delta_i \geq \frac{1}{(\alpha-1)\beta} + 2\epsilon$ , we obtain  $\mathbf{q}_0^{\top} (\mathbf{x}_i - \mathbf{x}_j) \geq \frac{1}{(\alpha-i)\beta}$ , which from the previous points ensures convergence to  $\mathbf{x}_i$  in one iteration.

## D Additional Plots

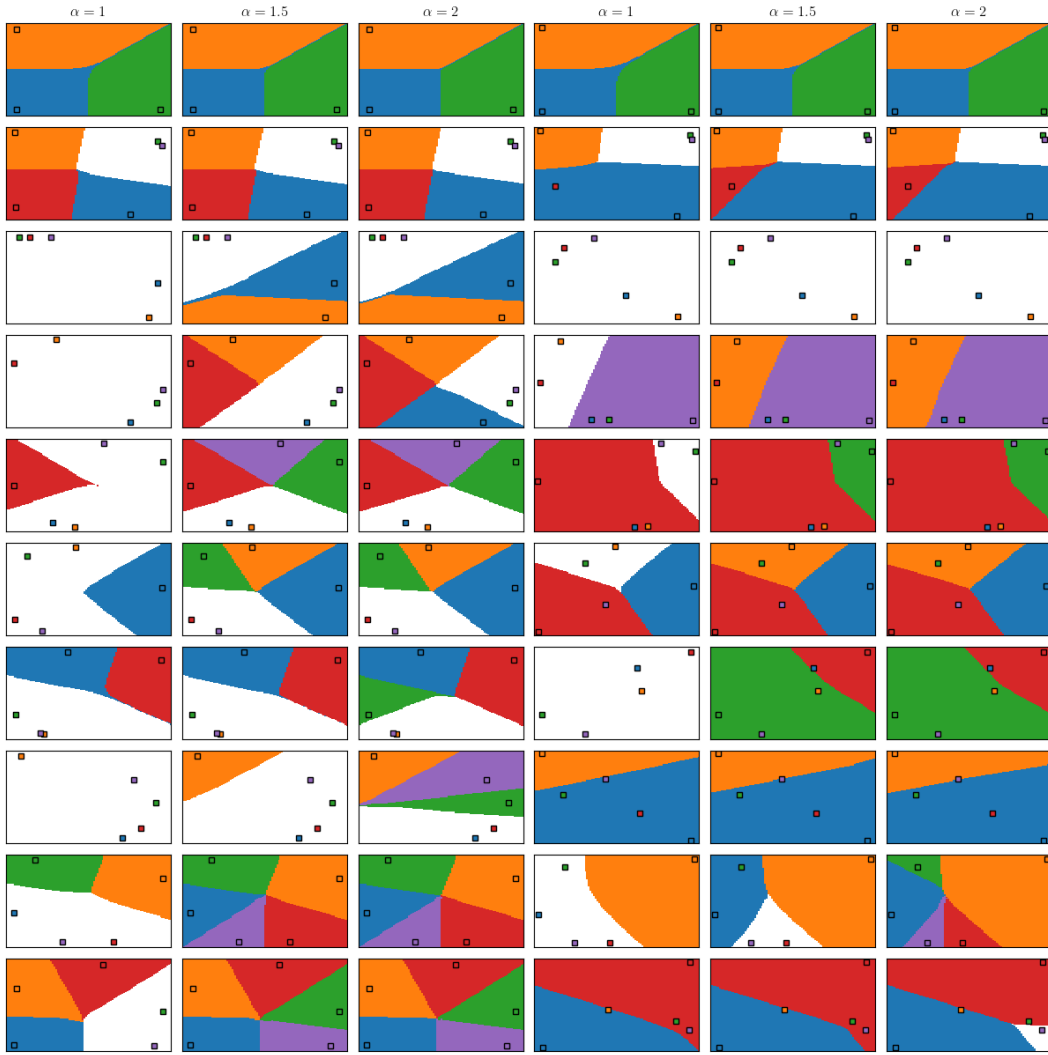




(a)  $\beta = 1$

(b)  $\beta = 4$

Figure 2: Additional energy contour and CCCP optimization trajectory plots.



(a)  $\beta = 10$ , normalized patterns

(b)  $\beta = 4$ , non-normalized patterns

Figure 3: Additional attraction basin plots. A tolerance of  $\epsilon = 0.01$  is allowed when  $\alpha = 1$ .