

---

# Learning Attractor Dynamics for Generative Memory (Appendix)

---

Yan Wu, Greg Wayne, Karol Gregor, Timothy Lillicrap  
DeepMind  
{yanwu,gregwayne,karolg,countzero}@google.com

## A The Bayesian update rule for the Kanerva Machine

Here we reproduce the exact Bayesian update rule used in the Kanerva Machine.

$$p(\mathbf{M}|\mathbf{z}, \mathbf{w}) = \frac{p(\mathbf{z}|\mathbf{w}, \mathbf{M}) p(\mathbf{M})}{\int p(\mathbf{z}|\mathbf{w}, \mathbf{M}) p(\mathbf{M}) d\mathbf{M}} \quad (1)$$

A memory with mean  $\mathbf{R}_{t-1}$ , row covariance  $\mathbf{U}_{t-1}$  and observational noise variance  $\sigma_\xi^2$  is updated given a newly observed sample  $\mathbf{z}$  and its addressing weight  $\mathbf{w}$  by:

$$\Delta \leftarrow \mathbf{z} - \mathbf{R}_{t-1}^\top \mathbf{w} \quad (2)$$

$$\Sigma_c \leftarrow \mathbf{U}_{t-1} \mathbf{w} \quad (3)$$

$$\Sigma_z \leftarrow \mathbf{w}^\top \mathbf{U}_{t-1} \mathbf{w} + \sigma_\xi^2 \quad (4)$$

$$\mathbf{R}_t \leftarrow \mathbf{R}_{t-1} + \Sigma_c \Sigma_z^{-1} \Delta^\top \quad (5)$$

$$\mathbf{U}_t \leftarrow \mathbf{U}_{t-1} - \Sigma_c \Sigma_z^{-1} \Sigma_c^\top \quad (6)$$

Note that the new covariance  $\mathbf{U}_t$  would collapse to zero if the noise variance  $\sigma^2 \rightarrow 0$ .

The graphical model assumes an observational noise  $\mathcal{N}(\mathbf{0}, \sigma_\xi^2)$ , which results in the read out distribution given the addressing weights  $\mathbf{w}$ ,  $p(\mathbf{z}|\mathbf{R} \mathbf{w}, \sigma_\xi^2)$  (eq. 1 in matrix notation). We ignore observation noise when reading the memory and directly take  $\mathbf{z} \leftarrow \mathbf{R} \mathbf{w}$ . This simplification reduces the variance in training and is justified by the fact that the fixed observation noise does not convey any information. A strong-enough decoder will learn to remove such noise through training.

## B Parameters and Initialisation

Here we enumerate parameters of the model and their initialisations. In our experiments, the memory parameters are insensitive to initial values.

Parameters	description	Initial Value
$\mathbf{R}_0$	$K \times C$ memory prior mean matrix	$\mathcal{N}(\mathbf{0}, \mathbf{I})$
$\mathbf{U}_0 = \sigma_U^2 \mathbf{I}$	$K \times K$ memory prior covariance matrix	$\sigma_U^2 = 1.0$
$\sigma_w^2$	Addressing weight posterior variance	0.3
$\sigma_\xi^2$	Memory observation noise variance	1.0
...	Neural network weights of encoder and decoder	Glorot Initialization

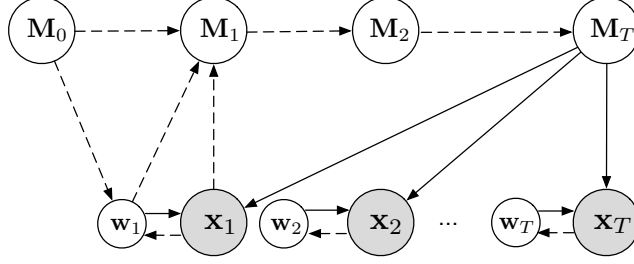


Figure 1: The probabilistic graphical model illustrating sequential online updates of memory. Dashed lines show the inference model and solid lines illustrate the generative model. For brevity, we illustrated only the inference of  $\mathbf{M}_1$  given  $\mathbf{M}_0$ ,  $\mathbf{w}_1$  and  $\mathbf{x}_1$ ; all the following steps are the same until reaching the end of the episode  $\mathbf{x}_T$ .

## C Sequential Variational Inference for Memory

The log-likelihood for any  $\mathbf{x}_{\leq T}$  can be decompose as a sum of a variational lower-bound and KL-divergences as:

$$\begin{aligned} \ln p(\mathbf{x}_{\leq T}) &= \ln \frac{p(\mathbf{x}_{\leq T}, \mathbf{w}_{\leq T}, \mathbf{M})}{p(\mathbf{w}_{\leq T}, \mathbf{M} | \mathbf{x}_{\leq T})} \\ &= \left\langle \ln \frac{p(\mathbf{x}_{\leq T} | \mathbf{w}_{\leq T}, \mathbf{M}) p(\mathbf{w}_{\leq T}) p(\mathbf{M}) q(\mathbf{w}_{\leq T}) q(\mathbf{M})}{p(\mathbf{w}_{\leq T} | \mathbf{x}_{\leq T}, \mathbf{M}) p(\mathbf{M} | \mathbf{x}_{\leq T}) q(\mathbf{w}_{\leq T}) q(\mathbf{M})} \right\rangle_{q(\mathbf{M}) q(\mathbf{w}_{\leq T})} \quad (7) \\ &= \mathcal{L}_T + \sum_{t=1}^T \langle \text{D}_{\text{KL}}(q(\mathbf{w}_t) \| p(\mathbf{w}_t | \mathbf{x}_t, \mathbf{M})) \rangle_{q(\mathbf{M})} + \text{D}_{\text{KL}}(q(\mathbf{M}) \| p(\mathbf{M} | \mathbf{x}_{\leq T})) \end{aligned}$$

$$\mathcal{L}_T = \sum_{t=1}^T \left( \langle \ln p(\mathbf{x}_t | \mathbf{w}_t, \mathbf{M}) \rangle_{q(\mathbf{w}_t) q(\mathbf{M})} - \text{D}_{\text{KL}}(q(\mathbf{w}_t) \| p(\mathbf{w}_t)) \right) - \text{D}_{\text{KL}}(q(\mathbf{M}) \| p(\mathbf{M})) \quad (8)$$

However, as we noted in the main text, it is hard to maximise  $\mathcal{L}_T$  directly, since we can not compute  $q(\mathbf{M}) \approx p(\mathbf{M} | \mathbf{x}_{\leq T})$  directly.

To derive a sequential update rule of the memory to compute  $q(\mathbf{M}) \approx p(\mathbf{M} | \mathbf{x}_{\leq T})$ , we consider updating the memory for step  $t$  of an episode. This assumes memory from the previous update  $q(\mathbf{M}_{t-1}; \mathbf{R}_{t-1}, \mathbf{U}_{t-1}) \approx p(\mathbf{M}_{t-1} | \mathbf{x}_{\leq t-1})$  is given, so that we can decompose  $\ln p(\mathbf{x}_{\leq t})$  conditioned on  $\mathbf{M}_{t-1}$ :

$$\begin{aligned} \ln p(\mathbf{x}_{\leq t}) &= \underbrace{\langle \ln p(\mathbf{x}_{\leq t} | \mathbf{M}_{t-1}) \rangle_{q(\mathbf{M}_{t-1})} - \text{D}_{\text{KL}}(q(\mathbf{M}_{t-1}) \| p(\mathbf{M}_{t-1}))}_{\mathcal{L}_t} \\ &\quad + \text{D}_{\text{KL}}(q(\mathbf{M}_{t-1}) \| p(\mathbf{M}_{t-1} | \mathbf{x}_{\leq t})) \end{aligned} \quad (9)$$

where we have a likelihood lower-bound  $\mathcal{L}_t$ , which has the same form as  $\mathcal{L}_T$ . This lower-bound is tight when  $q(\mathbf{M}_{t-1}) = p(\mathbf{M}_{t-1} | \mathbf{x}_{\leq t})$ . This suggests, ideally, that the memory  $q(\mathbf{M}_{t-1})$  at step  $t-1$  needs to be *predictive* of the next observation  $\mathbf{x}_t$ , in addition to accumulating information from the  $\mathbf{x}_{\leq t-1}$  that are already used in computing  $q(\mathbf{M}_{t-1})$ .

As illustrated in Figure 1, we assume a deterministic transition  $q(\mathbf{M}_t | \mathbf{M}_{t-1}) = \delta(\mathbf{M}_t - \mathbf{M}_{t-1})$ , so the prior of  $q(\mathbf{M}_t)$  simplifies to  $\int q(\mathbf{M}_t | \mathbf{M}_{t-1}) q(\mathbf{M}_{t-1}) d\mathbf{M}_{t-1} = \int \delta(\mathbf{M}_t - \mathbf{M}_{t-1}) q(\mathbf{M}_{t-1}) d\mathbf{M}_{t-1} = q(\mathbf{M}_t)$ . We can then *recursively* expand the likelihood term in eq. 9, similar to eq. 7 and eq. 8 (omitting the expectation over  $q(\mathbf{M}_{t-1})$ ):

$$\begin{aligned} \ln p(\mathbf{x}_{\leq t} | \mathbf{M}_{t-1}) &= \mathcal{L}'_t + \sum_{t'=1}^t \text{D}_{\text{KL}}(q(\mathbf{w}_{t'}) \| p(\mathbf{w}_{t'} | \mathbf{x}_{t'}, \mathbf{M}_{t-1})) \\ &\quad + \text{D}_{\text{KL}}(q(\mathbf{M}_t | \mathbf{w}_t) \| p(\mathbf{M}_t | \mathbf{x}_{\leq t}, \mathbf{w}_{\leq t})) \end{aligned} \quad (10)$$

$$\begin{aligned}\mathcal{L}'_t = & \sum_{t'=1}^t \left( \langle \ln p(\mathbf{x}_{t'} | \mathbf{w}_{t'}, \mathbf{M}_t) \rangle_{q(\mathbf{w}_{t'}), q(\mathbf{M}_t | \mathbf{w}_t)} - D_{\text{KL}}(q(\mathbf{w}_{t'}) \| p(\mathbf{w}_{t'})) \right) \\ & - \langle D_{\text{KL}}(q(\mathbf{M}_t | \mathbf{w}_t) \| q(\mathbf{M}_{t-1})) \rangle_{q(\mathbf{w}_t)}\end{aligned}\quad (11)$$

The above  $\mathcal{L}'_t$  is easier to maximise, since it only depends on  $\mathbf{x}_t$ , and on  $\mathbf{M}_{t-1}$ , which we assume we know. We can minimise the gap between  $\mathcal{L}'_t$  and  $\ln p(\mathbf{x}_{\leq t} | \mathbf{M}_{t-1})$  by minimising  $D_{\text{KL}}(q(\mathbf{M}_t | \mathbf{w}_t) \| p(\mathbf{M}_t | \mathbf{x}_{\leq t}, \mathbf{w}_{\leq t}))$  using the Bayes' update rule (Appendix A), and minimising  $\sum_{t'=1}^t D_{\text{KL}}(q(\mathbf{w}_{t'}) \| p(\mathbf{w}_{t'} | \mathbf{x}_{t'}, \mathbf{M}_{t-1}))$  using dynamic addressing (Section 3.1).

We can tighten  $\mathcal{L}_t$  by allowing further updating iterations as shown by the optional step in Algorithm 1. This is likely to be a tighter lower-bound, since generally the KL-divergence  $D_{\text{KL}}(q(\mathbf{M}_t) \| p(\mathbf{M}_{t-1} | \mathbf{x}_{\leq t})) < D_{\text{KL}}(q(\mathbf{M}_{t-1}) \| p(\mathbf{M}_{t-1} | \mathbf{x}_{\leq t-1}))$  after incorporating information from  $\mathbf{x}_t$ . This process can be repeated until this KL-divergence is tightened to its minimum.

From the above equations, we have the inequality

$$\ln p(\mathbf{x}_{\leq t}) \geq \mathcal{L}_t \geq \underbrace{\mathcal{L}'_t - D_{\text{KL}}(q(\mathbf{M}_{t-1}) \| p(\mathbf{M}_{t-1}))}_{\mathcal{B}_t} \quad (12)$$

Therefore, we can maximise  $\ln p(\mathbf{x}_{\leq t})$  by maximising the lower-bound  $\mathcal{B}_t$ . Naively, in eq. 10, all the  $t$  terms in  $\sum_{t'=1}^t D_{\text{KL}}(q(\mathbf{w}_{t'}) \| p(\mathbf{w}_{t'} | \mathbf{x}_{t'}, \mathbf{M}_{t-1}))$  need to be minimised at step  $t$ . This would result in  $\mathcal{O}(T^2)$  cost in both inferring  $q(\mathbf{w}_{\leq T})$  and  $q(\mathbf{M}_T)$ . To reduce the computational cost, we keep previously  $q(\mathbf{w}_{\leq t-1})$ , and only infer  $q(\mathbf{w}_t)$ , resulting in only  $\mathcal{O}(T)$  cost. The trade-off is a looser lower-bound  $\mathcal{L}'_t$  and therefore a looser  $\mathcal{B}_t$ , since some of the  $D_{\text{KL}}(q(\mathbf{w}_{t'}) \| p(\mathbf{w}_{t'} | \mathbf{x}_{t'}, \mathbf{M}_{t-1}))$  for  $t' < t$  may not be minimised.

Once  $q(\mathbf{M}_t | \mathbf{w}_t)$  is computed, we can compute the marginal for the next step as:

$$\begin{aligned}q(\mathbf{M}_t) &= \int q(\mathbf{M}_t | \mathbf{w}_t) q(\mathbf{w}_t) d\mathbf{w}_t \\ &\approx q(\mathbf{M}_t | \mathbf{w}_t) \Big|_{\mathbf{w}^* = \arg\max_{\mathbf{w}} p(\mathbf{w} | \mathbf{x}_q, \mathbf{M})}\end{aligned}\quad (13)$$

Memory updating is nonlinear, so the integral is not analytically tractable. A simple approximation is to use the mode of  $q(\mathbf{w}_t)$ , which is the mean  $\mu_{\mathbf{w}}$ . At this point, we carry forward the approximation of  $q(\mathbf{M}_{t-1})$  to  $q(\mathbf{M}_t)$ , which can be used for the  $t + 1$  step update. This procedure can start from  $t = 0$  and continue until  $t = T$ ; we thus obtain an approximate memory posterior  $q(\mathbf{M}_T)$  by maximising the lower-bound  $\mathcal{B}_T$ . Thus, the sequential update of memory, as summarised in Algorithm 1, maximises a lower-bound of the episode log-likelihood.

## D The Lease-Square Problem in Inference and Prediction

This sections shows that solving the same least squares problem is involved in both of the following problems:

1. minimising the KL-divergence between  $\mathbf{w}$  during inference (section 3.1),  $\mu_{\mathbf{w}}^* = \arg\min_{\mu_{\mathbf{w}}} D_{\text{KL}}(q(\mathbf{w}) \| p(\mathbf{w} | \mathbf{x}, \mathbf{M}))$
2. approximating the predictive distribution  $q(\hat{\mathbf{x}} | \mathbf{x}_q, \mathbf{M})$  (Section 3.3),  $\mathbf{w}^* = \arg\max_{\mathbf{w}} p(\mathbf{w} | \mathbf{x}_q, \mathbf{M})$

We first re-write the KL-divergence using its definition:

$$\begin{aligned}D_{\text{KL}}(q(\mathbf{w}) \| p(\mathbf{w} | \mathbf{x}, \mathbf{M})) &= \int q(\mathbf{w}) \ln \frac{q(\mathbf{w})}{p(\mathbf{w} | \mathbf{x}, \mathbf{M})} d\mathbf{w} \\ &= -H[q(\mathbf{w})] - \langle \ln p(\mathbf{w} | \mathbf{x}, \mathbf{M}) \rangle_{q(\mathbf{w})}\end{aligned}\quad (14)$$

where the first entropy term is a constant that depends on the fixed variance  $\sigma_{\mathbf{w}}^2$ . Therefore, minimising this KL-divergence is equivalent to maximising  $\langle \ln p(\mathbf{w} | \mathbf{x}, \mathbf{M}) \rangle_{q(\mathbf{w})}$ .

This posterior distribution over  $\mathbf{w}$  can be expanded using Bayes' rule:

$$\begin{aligned}\ln p(\mathbf{w}|\mathbf{x}, \mathbf{M}) &= \ln \frac{p(\mathbf{x}|\mathbf{w}, \mathbf{M}) p(\mathbf{w})}{p(\mathbf{x}|\mathbf{M})} \\ &= \ln p(\mathbf{x}|\mathbf{w}, \mathbf{M}) + \ln p(\mathbf{w}) + \dots \\ &\approx -\frac{\|e(\mathbf{x}) - \mathbf{M}^\top \mathbf{w}\|^2}{2\sigma_\xi^2(\mathbf{x})} - \frac{1}{2} \|\mathbf{w}\|^2 + \dots\end{aligned}\quad (15)$$

We omitted terms that do not depend on  $\mathbf{w}$ , including various normalising constants. In addition, the last line used the encoding projection  $e(\mathbf{x}) \rightarrow \mathbf{z}$  to transform the distribution over  $\mathbf{x}$  to that over  $\mathbf{z}$ . When  $e(\mathbf{x})$  is invertible, the Jacobian factor  $\frac{\mathbf{z}}{\partial \mathbf{x}} = \frac{\partial e(\mathbf{x})}{\partial \mathbf{x}}$  resulting from the distribution transform is well-defined and can be omitted since it does not depend on  $\mathbf{w}$ . However, the assumption of bijection is unlikely to be strictly satisfied by the neural network encoder/decoder pair, so the relation is approximate.

Taking the expectation of the above quadratic equation over the Gaussian distribution  $q(\mathbf{w})$  results in the same quadratic form:

$$\begin{aligned}\langle \ln p(\mathbf{w}|\mathbf{x}, \mathbf{M}) \rangle_{q(\mathbf{w})} &\approx \left\langle -\frac{\|e(\mathbf{x}) - \mathbf{M}^\top \mathbf{w}\|^2}{2\sigma_\xi^2} - \frac{1}{2} \|\mathbf{w}\|^2 \right\rangle_{q(\mathbf{w})} + \dots \\ &= -\frac{\|e(\mathbf{x}) - \mathbf{M}^\top \mu_{\mathbf{w}}\|^2}{2\sigma_\xi^2} - \frac{1}{2} \|\mu_{\mathbf{w}}\|^2 + \frac{\sigma_{\mathbf{w}}^2}{2\sigma_\xi^2} \text{Tr}(\mathbf{M} \mathbf{M}^\top) + \sigma_{\mathbf{w}}^2 C + \dots\end{aligned}\quad (16)$$

where the last two terms do not depend on  $\mu_{\mathbf{w}}$ . Therefore, both inference and prediction involve solving the same least-squares problem.

## E Proof of Attractor Dynamics

Here we show that in a well trained model, a pattern  $\mathbf{x}^*$  in the memory is *asymptotically stable* under the dynamics, so that a state near  $\mathbf{x}^*$  will converge to it. By ‘‘a well trained model’’, we assume that pattern  $\mathbf{x}^*$  is a local maximum of the ELBO (eq. 4 in the main text)

$$\mathcal{L}(\mathbf{x}^*) = \langle \ln p(\mathbf{x}^*|\mathbf{w}^*, \mathbf{M}) \rangle_{q(\mathbf{M})} - \text{D}_{\text{KL}}(q(\mathbf{w}^*)\|p(\mathbf{w})) - \text{D}_{\text{KL}}(q(\mathbf{M})\|p(\mathbf{M}|\mathbf{x}_{\leq T})) \quad (17)$$

When  $\mathcal{L}(\mathbf{x}^*)$  is at maximum, the energy we defined in eq. 14 (copied below) would be at its local minimum. This follows since the negative energy is just the first 2 terms of  $\mathcal{L}$  without the KL-divergence between  $\mathbf{M}$ , which is a constant when the memory is fixed.

$$\mathcal{E}(\mathbf{x}, \mathbf{w}) = -\langle \ln p(\mathbf{x}|\mathbf{w}, \mathbf{M}) \rangle_{q(\mathbf{M})} + \text{D}_{\text{KL}}(q_t(\mathbf{w})\|p(\mathbf{w})) \quad (18)$$

Section 3.5 of the main text shows  $\mathcal{E}(\mathbf{x}_n, \mathbf{w}_n) \leq \mathcal{E}(\mathbf{x}_{n-1}, \mathbf{w}_{n-1})$  under the predictive dynamics. Therefore, we can construct a Lyapunov function candidate as:

$$V(\mathbf{x}, \mathbf{w}) = \mathcal{E}(\mathbf{x}, \mathbf{w}) - \mathcal{E}(\mathbf{w}^*, \mathbf{w}^*) \quad (19)$$

which satisfies:

$$V(\mathbf{x}^*, \mathbf{w}^*) = 0 \quad (20)$$

$$V(\mathbf{x}, \mathbf{w}) > 0 \quad \forall (\mathbf{x}, \mathbf{w}) \neq (\mathbf{x}^*, \mathbf{w}^*) \quad (21)$$

$$V(\mathbf{x}_n, \mathbf{w}_n) < V(\mathbf{x}_{n-1}, \mathbf{w}_{n-1}) \quad \forall (\mathbf{x}_n, \mathbf{w}_n) \neq (\mathbf{x}^*, \mathbf{w}^*) \quad (22)$$

Therefore, according to Lyapunov Stability theory, state  $\mathbf{x}^*$  is asymptotically stable and serves as a point attractor in the system.

## F Other Practical Considerations

For readers interested in applying the DKM, a few variants of Algorithm 1 may be worth considering. First, instead of using  $\mathcal{L}_T$  in eq.8 (eq. 3 in the main text) as the objective, an alternative objective is  $\mathcal{B}_T$ . Although this lower-bound tends to be less tight than  $\mathcal{L}_T$ , it is also cheaper to compute. It may be particularly useful in online settings, since we only need to run through an episode once to compute  $q(\mathbf{M}_T)$ . This bounds can be further tightened by: 1. Using the optional step in Algorithm 1. We recommend starting with 2 or 3 steps. 2. Minimising a few other intermediate  $\mathcal{B}_t$ 's, for  $0 < t < T$ . This may be helpful in the case of long episodes wherein gradient propagation through the entire episode is infeasible.