

# Attention via $\log \sum \exp$ energy

Alexander Tschantz

January 27, 2025

## 1 General Framework

**Setup.** We consider a single set of nodes  $\mathbf{v} = \{\mathbf{v}_i : i \in \{1, 2, \dots, N\}\}$ , where each node  $\mathbf{v}_i \in \mathbb{R}^{d_i}$ . The relationships between these nodes are defined by a set of  $M$  energy functions  $\{E_m : m \in \{1, 2, \dots, M\}\}$ . Each energy function  $E_m$  defines a subset of nodes acting as *children*  $C_m \subseteq \{1, 2, \dots, \text{normalization}\}$  and a subset acting as *parents*  $P_m \subseteq \{1, 2, \dots, N\}$ , which may overlap.

**Energy.** Each energy function  $E_m$  defines a similarity function:

$$\text{sim}(\mathbf{v}_c, \mathbf{v}_p) \quad : \quad \mathbb{R}^{d_c} \times \mathbb{R}^p \rightarrow \mathbb{R}, \quad (1)$$

which produces a scalar similarity between a child  $\mathbf{v}_c$  and a parent  $\mathbf{v}_p$ . Let  $\{\mathbf{v}_c\} = \{\mathbf{v}_c : c \in C_m\}$  and  $\{\mathbf{v}_p\} = \{\mathbf{v}_p : p \in P_m\}$ , the energy for  $E_m$  is defined as:

$$E_m(\{\mathbf{v}_c\}, \{\mathbf{v}_p\}) = - \sum_{c \in C_m} \ln \left( \sum_{p \in P_m} \exp(\text{sim}(\mathbf{v}_c, \mathbf{v}_p)) \right). \quad (2)$$

The global energy sums over all energy functions:

$$E(\{\mathbf{v}\}) = \sum_{m=1}^M E_m(\{\mathbf{v}_c\}, \{\mathbf{v}_p\}). \quad (3)$$

**Gradient Updates.** For a single node  $\mathbf{v}_a$ , the gradient of the global energy  $E$  w.r.t.  $\mathbf{v}_a$  decomposes into two terms. Let  $\mathcal{M}_c(a) = \{m : a \in C_m\}$  denote the energy functions where  $\mathbf{v}_a$  acts as a *child*, and  $\mathcal{M}_p(a) = \{m : a \in P_m\}$  the energy functions where  $\mathbf{v}_a$  acts as a *parent*. Then:

$$\begin{aligned} -\frac{\partial E}{\partial \mathbf{v}_a} &= \underbrace{\sum_{m \in \mathcal{M}_c(a)} \sum_{p \in P_m} \text{softmax}_p(\text{sim}(\mathbf{v}_a, \mathbf{v}_p)) \frac{\partial}{\partial \mathbf{v}_a} \text{sim}(\mathbf{v}_a, \mathbf{v}_p)}_{\mathbf{v}_a \text{ acting as a child}} \\ &+ \underbrace{\sum_{m \in \mathcal{M}_p(a)} \sum_{c \in C_m} \text{softmax}_a(\text{sim}(\mathbf{v}_c, \mathbf{v}_a)) \frac{\partial}{\partial \mathbf{v}_a} \text{sim}(\mathbf{v}_c, \mathbf{v}_a)}_{\mathbf{v}_a \text{ acting as a parent}}. \end{aligned} \quad (4)$$

The first term captures contributions from  $\mathbf{v}_a$  being explained by its parents, while the second term captures contributions from  $\mathbf{v}_a$  explaining its children.

## 2 Gaussian Mixture Models

**Setup.** We have  $N$  data points (children)  $\mathbf{x}_i \in \mathbb{R}^d$ ,  $i \in C = \{1, \dots, N\}$ , and  $K$  mixture components (parents), each with mean  $\boldsymbol{\mu}_k \in \mathbb{R}^d$  and covariance  $\boldsymbol{\Sigma}_k$ ,  $k \in P = \{1, \dots, K\}$ . Let  $\pi_k$  be the mixing proportion.

**Similarity function.** We define

$$\text{sim}(\mathbf{x}_i, \boldsymbol{\mu}_k) = \ln \pi_k - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k).$$

**Energy.**

$$E^{\text{GMM}}(\{\mathbf{x}_i\}, \{\boldsymbol{\mu}_k\}) = - \sum_{i=1}^N \ln \left( \sum_{k=1}^K \exp(\text{sim}(\mathbf{x}_i, \boldsymbol{\mu}_k)) \right). \quad (5)$$

**Gradients.** If we differentiate w.r.t.  $\boldsymbol{\mu}_k$ , then

$$-\frac{\partial E^{\text{GMM}}}{\partial \boldsymbol{\mu}_k} = \sum_{i=1}^N \text{softmax}_k(\text{sim}(\mathbf{x}_i, \boldsymbol{\mu}_k)) \boldsymbol{\Sigma}_k^{-1}(\mathbf{x}_i - \boldsymbol{\mu}_k).$$

Setting this gradient to zero yields the usual GMM M-step:

$$\boldsymbol{\mu}_k = \frac{\sum_{i=1}^N \text{softmax}_k(\text{sim}(\mathbf{x}_i, \boldsymbol{\mu}_k)) \mathbf{x}_i}{\sum_{i=1}^N \text{softmax}_k(\text{sim}(\mathbf{x}_i, \boldsymbol{\mu}_k))}.$$

### 3 Hopfield Networks

**Setup.** We have a set of *children* data vectors  $\mathbf{x}_i \in \mathbb{R}^d$ ,  $i \in C = \{1, \dots, N\}$ , and a set of *parent* memory vectors  $\mathbf{m}_\mu \in \mathbb{R}^d$ ,  $\mu \in P = \{1, \dots, K\}$ .

**Similarity function.**

$$\text{sim}(\mathbf{x}_i, \mathbf{m}_\mu) = \mathbf{x}_i^\top \mathbf{m}_\mu.$$

**Energy.**

$$E^{\text{Hopfield}}(\{\mathbf{x}_i\}, \{\mathbf{m}_\mu\}) = - \sum_{i=1}^N \ln \left( \sum_{\mu=1}^K \exp(\mathbf{x}_i^\top \mathbf{m}_\mu) \right). \quad (6)$$

**Gradients.**

$$-\frac{\partial E^{\text{Hopfield}}}{\partial \mathbf{x}_i} = \sum_{\mu=1}^K \text{softmax}_\mu(\mathbf{x}_i^\top \mathbf{m}_\mu) \mathbf{m}_\mu. \quad (7)$$

$$-\frac{\partial E^{\text{Hopfield}}}{\partial \mathbf{m}_\mu} = \sum_{i=1}^N \text{softmax}_\mu(\mathbf{x}_i^\top \mathbf{m}_\mu) \mathbf{x}_i. \quad (8)$$

### 4 Slot Attention

**Setup.** Let  $\mathbf{x}_j \in \mathbb{R}^d$ ,  $j \in C = \{1, \dots, N\}$  be the children (tokens), and  $\boldsymbol{\mu}_i \in \mathbb{R}^d$ ,  $i \in P = \{1, \dots, S\}$  be the parents (slots). We typically apply linear transforms  $\mathbf{W}_K, \mathbf{W}_Q \in \mathbb{R}^{d \times d}$  to form

$$\text{sim}(\mathbf{x}_j, \boldsymbol{\mu}_i) = (\mathbf{W}_K \mathbf{x}_j)^\top (\mathbf{W}_Q \boldsymbol{\mu}_i).$$

**Energy.**

$$E^{\text{Slot}}(\{\mathbf{x}_j\}, \{\boldsymbol{\mu}_i\}) = - \sum_{j=1}^N \ln \left( \sum_{i=1}^S \exp(\text{sim}(\mathbf{x}_j, \boldsymbol{\mu}_i)) \right). \quad (9)$$

**Gradients.**

$$-\frac{\partial E^{\text{Slot}}}{\partial \boldsymbol{\mu}_i} = \sum_{j=1}^N \text{softmax}_i(\text{sim}(\mathbf{x}_j, \boldsymbol{\mu}_i)) \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{x}_j. \quad (10)$$

$$-\frac{\partial E^{\text{Slot}}}{\partial \mathbf{x}_j} = \sum_{i=1}^S \text{softmax}_i(\text{sim}(\mathbf{x}_j, \boldsymbol{\mu}_i)) \mathbf{W}_K^\top \mathbf{W}_Q \boldsymbol{\mu}_i. \quad (11)$$

## 5 Self-Attention

**Setup.** In self-attention, every node can act as both a child (query) and a parent (key). Concretely, let us have  $N$  tokens  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ . We form

$$\mathbf{q}_i = \mathbf{W}^Q \mathbf{x}_i, \quad \mathbf{k}_i = \mathbf{W}^K \mathbf{x}_i,$$

for  $i = 1, \dots, N$ . Thus the set  $C = \{1, \dots, N\}$  and  $P = \{1, \dots, N\}$  coincide, with

$$\text{sim}(\mathbf{x}_c, \mathbf{x}_p) = (\mathbf{W}^Q \mathbf{x}_c)^\top (\mathbf{W}^K \mathbf{x}_p).$$

**Energy.**

$$E^{\text{SA}}(\{\mathbf{x}_i\}) = - \sum_{c=1}^N \ln \left( \sum_{p=1}^N \exp \left( (\mathbf{W}^Q \mathbf{x}_c)^\top (\mathbf{W}^K \mathbf{x}_p) \right) \right). \quad (12)$$

**Gradients.** Since each  $\mathbf{x}_i$  is *both* a child and a parent, its gradient is a sum of two terms (the child side and the parent side). Writing it out explicitly:

$$\begin{aligned} -\frac{\partial E^{\text{SA}}}{\partial \mathbf{x}_i} &= \underbrace{\sum_{p=1}^N \text{softmax}_p \left( (\mathbf{W}^Q \mathbf{x}_i)^\top (\mathbf{W}^K \mathbf{x}_p) \right) \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{x}_p}_{\text{child } i \text{ being explained by parents } p} \\ &\quad + \underbrace{\sum_{c=1}^N \text{softmax}_i \left( (\mathbf{W}^Q \mathbf{x}_c)^\top (\mathbf{W}^K \mathbf{x}_i) \right) \mathbf{W}_K^\top \mathbf{W}_Q \mathbf{x}_c}_{\text{parent } i \text{ explaining children } c}. \end{aligned} \quad (13)$$

## 6 Spatiotemporal Attention

**Setup.** We consider a hierarchy of  $L$  layers, each containing  $K_l$  latent variables of dimension  $D_l$ . Concretely, let  $\mathbf{x}_{k,t}^{(l)} \in \mathbb{R}^{D_l}$  denote the  $k$ -th variable (or “slot”) in layer  $l$  at time  $t$ . The lowest layer ( $l = 1$ ) has  $K_1 = N$  observed variables (e.g., pixels) and dimension  $D_1$  (e.g.,  $[x, y, r, g, b]$ ), while higher layers have separate dimensions  $D_l$  and numbers of slots  $K_l$ . Our goal is to define an energy that couples these variables *vertically* (across layers), *concurrently* (within the same layer and time), and *temporally* (across time).

**Similarity Functions.** We introduce three types of similarity, each with its own projection matrices. For inter-layer connections (linking layers  $l - 1$  and  $l$ ), we define:

$$\text{sim}_v^{(l)}(\mathbf{x}, \mathbf{x}') = (\mathbf{W}_v^{K,(l)} \mathbf{x})^\top (\mathbf{W}_v^{Q,(l)} \mathbf{x}'), \quad (14)$$

For intra-layer (slots within the same layer and time):

$$\text{sim}_c^{(l)}(\mathbf{x}, \mathbf{x}') = (\mathbf{W}_c^{Q,(l)} \mathbf{x})^\top (\mathbf{W}_c^{K,(l)} \mathbf{x}'), \quad (15)$$

For temporal connections (the same slot across different times):

$$\text{sim}_t^{(l)}(\mathbf{x}, \mathbf{x}') = (\mathbf{W}_t^{Q,(l)} \mathbf{x})^\top (\mathbf{W}_t^{K,(l)} \mathbf{x}'). \quad (16)$$

Here,  $\mathbf{W}_\bullet^{Q,(l)}$  and  $\mathbf{W}_\bullet^{K,(l)}$  are learnable projection matrices for layer  $l$ . The intra-layer and temporal parameters parallel the key-query mechanism in self-attention, while the inter-layer parameters are analogous to the inverted attention mechanism used in slot attention.

**Energy** At each layer  $l$ , the total energy is split into three terms:

$$\begin{aligned}
E_v^{(l)} &= - \sum_{t=1}^T \sum_{c=1}^{K_{l-1}} \underbrace{\ln \left( \sum_{k=1}^{K_l} \exp(\text{sim}_v^{(l)}(\mathbf{x}_{c,t}^{(l-1)}, \mathbf{x}_{k,t}^{(l)})) \right)}_{\text{Inter-layer energy}}, \\
E_c^{(l)} &= - \sum_{t=1}^T \sum_{k=1}^{K_l} \underbrace{\ln \left( \sum_{\substack{k'=1 \\ k' \neq k}}^{K_l} \exp(\text{sim}_c^{(l)}(\mathbf{x}_{k,t}^{(l)}, \mathbf{x}_{k',t}^{(l)})) \right)}_{\text{Intra-layer energy}}, \\
E_t^{(l)} &= - \sum_{k=1}^{K_l} \sum_{t=2}^T \underbrace{\ln \left( \sum_{t' < t} \exp(\text{sim}_t^{(l)}(\mathbf{x}_{k,t}^{(l)}, \mathbf{x}_{k,t'}^{(l)})) \right)}_{\text{Temporal energy}}.
\end{aligned} \tag{17}$$

Summing these over all layers  $l \in \{1, \dots, L\}$  defines the total energy  $E$ . These terms correspond to inter-layer connections ( $E_v^{(l)}$ ) and match the energy function for slot attention, intra-layer connections ( $E_c^{(l)}$ ) which match the energy function for self attention, and temporal connections ( $E_t^{(l)}$ ), which match the energy function for *causal* self attention (as each variable is only explained by past variables).

**Gradients** Consider a single variable  $\mathbf{x}_{k,t}^{(l)}$ . Its gradient with respect to the energy decomposes into four parts:

$$-\frac{\partial E}{\partial \mathbf{x}_{k,t}^{(l)}} = \underbrace{-\frac{\partial E_v^{(l)}}{\partial \mathbf{x}_{k,t}^{(l)}}}_{\text{bottom-up}} + \underbrace{-\frac{\partial E_v^{(l+1)}}{\partial \mathbf{x}_{k,t}^{(l)}}}_{\text{top-down}} + \underbrace{-\frac{\partial E_c^{(l)}}{\partial \mathbf{x}_{k,t}^{(l)}}}_{\text{intra-layer}} + \underbrace{-\frac{\partial E_t^{(l)}}{\partial \mathbf{x}_{k,t}^{(l)}}}_{\text{temporal}}. \tag{18}$$

**Bottom-up:** We define a similarity matrix  $\mathbf{A}_{c,k}$ , representing the interactions between  $\mathbf{x}_{k,t}^{(l)}$  in layer  $l$  and  $\mathbf{x}_{c,t}^{(l-1)}$  in the layer below.

$$\mathbf{A}_{c,k} = \text{sim}_v^{(l)}(\mathbf{x}_{c,t}^{(l-1)}, \mathbf{x}_{k,t}^{(l)}), \tag{19}$$

The gradient with respect to  $\mathbf{x}_{k,t}^{(l)}$  is:

$$-\frac{\partial E_v^{(l)}}{\partial \mathbf{x}_{k,t}^{(l)}} = \sum_{c=1}^{K_{l-1}} \text{softmax}_k(\mathbf{A}_{c,k}) \mathbf{W}_v^{Q,(l)\top} \mathbf{W}_v^{K,(l)} \mathbf{x}_{c,t}^{(l-1)}. \tag{20}$$

This term aggregates contributions from the children in layer  $l-1$ , weighted by the attention  $\text{softmax}_k(\mathbf{A}_{c,k})$ , and is analogous to Slot Attention or Gaussian mixture models.

**Top-down:** We define a similarity matrix  $\mathbf{A}_{k,p}$ , representing the interactions between  $\mathbf{x}_{k,t}^{(l)}$  in layer  $l$  and  $\mathbf{x}_{p,t}^{(l+1)}$  in the layer above.

$$\mathbf{A}_{k,p} = \text{sim}_v^{(l+1)}(\mathbf{x}_{k,t}^{(l)}, \mathbf{x}_{p,t}^{(l+1)}), \tag{21}$$

The gradient with respect to  $\mathbf{x}_{k,t}^{(l)}$  is:

$$-\frac{\partial E_v^{(l+1)}}{\partial \mathbf{x}_{k,t}^{(l)}} = \sum_{p=1}^{K_{l+1}} \text{softmax}_p(\mathbf{A}_{k,p}) \mathbf{W}_v^{K,(l+1)\top} \mathbf{W}_v^{Q,(l+1)} \mathbf{x}_{p,t}^{(l+1)}. \tag{22}$$

This term captures the influence of  $\mathbf{x}_{k,t}^{(l)}$  being treated as a child, weighted by the attention  $\text{softmax}_p(\mathbf{A}_{k,p})$ , and reflects the top-down influence from layer  $l+1$ , and is analogous to Hopfield attention or the parent term in self-attention.

**Intra-layer:** We define two similarity matrices,  $\mathbf{A}_{k,k'}$  and  $\mathbf{A}_{k',k}$ , representing the bidirectional interactions between  $\mathbf{x}_{k,t}^{(l)}$  and other slots  $\mathbf{x}_{k',t}^{(l)}$ :

$$\begin{aligned}\mathbf{A}_{k,k'} &= \text{sim}_c^{(l)}(\mathbf{x}_{k,t}^{(l)}, \mathbf{x}_{k',t}^{(l)}), \\ \mathbf{A}_{k',k} &= \text{sim}_c^{(l)}(\mathbf{x}_{k',t}^{(l)}, \mathbf{x}_{k,t}^{(l)}),\end{aligned}\tag{23}$$

where

$$\text{sim}_c^{(l)}(\mathbf{x}_{k,t}^{(l)}, \mathbf{x}_{k',t}^{(l)}) = (\mathbf{W}_c^{Q,(l)} \mathbf{x}_{k,t}^{(l)})^\top (\mathbf{W}_c^{K,(l)} \mathbf{x}_{k',t}^{(l)}).$$

The gradient with respect to  $\mathbf{x}_{k,t}^{(l)}$  is:

$$\begin{aligned}-\frac{\partial E_c^{(l)}}{\partial \mathbf{x}_{k,t}^{(l)}} &= \underbrace{\sum_{k' \neq k} \text{softmax}_k(\mathbf{A}_{k',k}) \mathbf{W}_c^{K,(l)\top} \mathbf{W}_c^{Q,(l)} \mathbf{x}_{k',t}^{(l)}}_{\mathbf{x}_{k,t}^{(l)} \text{ acting as parent (explaining others)}} \\ &\quad + \underbrace{\sum_{k' \neq k} \text{softmax}_{k'}(\mathbf{A}_{k,k'}) \mathbf{W}_c^{Q,(l)\top} \mathbf{W}_c^{K,(l)} \mathbf{x}_{k',t}^{(l)}}_{\mathbf{x}_{k,t}^{(l)} \text{ acting as child (being explained by others)}}.\end{aligned}\tag{24}$$

The first term captures the influence of  $\mathbf{x}_{k,t}^{(l)}$  as a parent, aggregating the contributions from other slots  $\mathbf{x}_{k',t}^{(l)}$  it explains, weighted by  $\text{softmax}_{k'}(\mathbf{A}_{k,k'})$ . The second term reflects  $\mathbf{x}_{k,t}^{(l)}$ 's role as a child, being explained by other slots  $\mathbf{x}_{k',t}^{(l)}$ , weighted by  $\text{softmax}_k(\mathbf{A}_{k',k})$ .

**Temporal:**

**Setup.** We define a similarity matrix  $\mathbf{A}_{t,t'}$  representing the interaction between a slot  $\mathbf{x}_{k,t}^{(l)}$  at time  $t$  and its past representation  $\mathbf{x}_{k,t'}^{(l)}$  at time  $t'$ , for  $t' < t$ :

$$\mathbf{A}_{t,t'} = \text{sim}_t^{(l)}(\mathbf{x}_{k,t}^{(l)}, \mathbf{x}_{k,t'}^{(l)}) = (\mathbf{W}_t^{Q,(l)} \mathbf{x}_{k,t}^{(l)})^\top (\mathbf{W}_t^{K,(l)} \mathbf{x}_{k,t'}^{(l)}).\tag{25}$$

**Gradient.** Assuming causal attention (i.e. no future times), the gradient w.r.t.  $\mathbf{x}_{k,t}^{(l)}$  is:

$$-\frac{\partial E_t^{(l)}}{\partial \mathbf{x}_{k,t}^{(l)}} = \sum_{t' < t} \text{softmax}_{t'}(\mathbf{A}_{t,t'}) \mathbf{W}_t^{Q,(l)\top} \mathbf{W}_t^{K,(l)} \mathbf{x}_{k,t'}^{(l)}.\tag{26}$$

Thus  $\mathbf{x}_{k,t}^{(l)}$  is influenced only by its own past  $\mathbf{x}_{k,t'}^{(l)}$  ( $t' < t$ ), as in causal self-attention.

**Remark (Full Version).** We can generalize to let each latent  $\mathbf{s}_{k,t}$  attend *all* data features and *all* slot latents at *any* time. This merges the three energies into a single large log-sum-exp, imposing direct competition among data, concurrency, and temporal parents, but increases complexity to  $\mathcal{O}((N+K)T)^2$ .

## 7 Dynamic Attention Mixtures

**Setup.** We consider a single variable  $\mathbf{x} \in \mathbb{R}^D$  at the current time step, which is influenced by  $K$  other entities, indexed by  $k \in \{1, \dots, K\}$ . These entities may represent: (i) previous states of the same variable across time (temporal parents), or (ii) other variables or slots at the same time step (concurrent parents). The influence of these  $K$  entities determines the value of  $\mathbf{y} \in \mathbb{R}^D$ , where  $\mathbf{x} \rightarrow \mathbf{y}$  evolves through  $K$  modes or parameterized transformations. Switching between these modes is governed by an attention mechanism over the  $K$  entities.

**Energy Function.** The energy function combines self-attention over parents and compatibility with dynamical modes. It is defined as:

$$E(\mathbf{x}, \mathbf{y}) = -\ln \sum_{k=1}^K \exp(\text{sim}_{\text{parents}}(\mathbf{x}, \mathbf{x}_k) + \text{sim}_{\text{dynamics}}(\mathbf{x}, \mathbf{y}; V_k)), \quad (27)$$

where

$$\text{sim}_{\text{parents}}(\mathbf{x}, \mathbf{x}_k) = (\mathbf{W}^Q \mathbf{x})^\top (\mathbf{W}^K \mathbf{x}_k), \quad (28)$$

$$\text{sim}_{\text{dynamics}}(\mathbf{x}, \mathbf{y}; V_k) = -\frac{1}{2} \|\mathbf{y} - V_k \mathbf{x}\|^2. \quad (29)$$

The term  $\text{sim}_{\text{parents}}(\mathbf{x}, \mathbf{x}_k)$  measures the similarity between the current state  $\mathbf{x}$  and the  $k$ -th parent through learnable query and key matrices  $\mathbf{W}^Q$  and  $\mathbf{W}^K$ . The term  $\text{sim}_{\text{dynamics}}(\mathbf{x}, \mathbf{y}; V_k)$  evaluates the compatibility of the  $k$ -th mode  $V_k \in \mathbb{R}^{D \times D}$  in predicting  $\mathbf{y}$  from  $\mathbf{x}$ .

**Attention Weights.** The attention weights, which determine the influence of each parent and its associated dynamics mode  $V_k$ , are defined as:

$$\alpha_k = \frac{\exp(\text{sim}_{\text{parents}}(\mathbf{x}, \mathbf{x}_k) + \text{sim}_{\text{dynamics}}(\mathbf{x}, \mathbf{y}; V_k))}{\sum_{j=1}^K \exp(\text{sim}_{\text{parents}}(\mathbf{x}, \mathbf{x}_j) + \text{sim}_{\text{dynamics}}(\mathbf{x}, \mathbf{y}; V_j))}. \quad (30)$$

**Gradient with respect to  $\mathbf{y}$ .** The derivative of the energy function with respect to  $\mathbf{y}$  is:

$$\frac{\partial E}{\partial \mathbf{y}} = -\sum_{k=1}^K \alpha_k \frac{\partial}{\partial \mathbf{y}} \text{sim}_{\text{dynamics}}(\mathbf{x}, \mathbf{y}; V_k). \quad (31)$$

$$\frac{\partial}{\partial \mathbf{y}} \text{sim}_{\text{dynamics}}(\mathbf{x}, \mathbf{y}; V_k) = \mathbf{y} - V_k \mathbf{x}, \quad (32)$$

$$\frac{\partial E}{\partial \mathbf{y}} = -\sum_{k=1}^K \alpha_k (\mathbf{y} - V_k \mathbf{x}). \quad (33)$$

**Gradient with respect to  $\mathbf{x}$ .** The derivative of the energy function with respect to  $\mathbf{x}$  is:

$$\frac{\partial E}{\partial \mathbf{x}} = -\sum_{k=1}^K \alpha_k \left( \frac{\partial}{\partial \mathbf{x}} \text{sim}_{\text{parents}}(\mathbf{x}, \mathbf{x}_k) + \frac{\partial}{\partial \mathbf{x}} \text{sim}_{\text{dynamics}}(\mathbf{x}, \mathbf{y}; V_k) \right). \quad (34)$$

$$\frac{\partial}{\partial \mathbf{x}} \text{sim}_{\text{parents}}(\mathbf{x}, \mathbf{x}_k) = W_Q W_K^\top \mathbf{x}_k, \quad (35)$$

$$\frac{\partial}{\partial \mathbf{x}} \text{sim}_{\text{dynamics}}(\mathbf{x}, \mathbf{y}; V_k) = V_k^\top (\mathbf{y} - V_k \mathbf{x}). \quad (36)$$

$$\frac{\partial E}{\partial \mathbf{x}} = -\sum_{k=1}^K \alpha_k (W_Q W_K^\top \mathbf{x}_k + V_k^\top (\mathbf{y} - V_k \mathbf{x})). \quad (37)$$

## 8 Predictive coding

**Setup.** We again have child vectors  $\{\mathbf{x}_i\}_{i=1}^N$  and a set of parent  $\boldsymbol{\mu}_k \in \mathbb{R}^d, k = 1, \dots, K$ . However, rather than a direct difference  $\mathbf{x}_i - \boldsymbol{\mu}_k$ , let us assume the *model* maps  $\boldsymbol{\mu}_k$  through some non-linear function  $f_\phi(\cdot)$  before comparing to  $\mathbf{x}_i$ . For instance,  $f_\phi$  could be a neural network.

**Similarity function.** Define

$$\text{sim}(\mathbf{x}_i, \boldsymbol{\mu}_k) = -\frac{1}{2} \|\mathbf{x}_i - f_\phi(\boldsymbol{\mu}_k)\|^2.$$

**Energy.**

$$E^{\text{PC}}(\{\mathbf{x}_i\}, \{\boldsymbol{\mu}_k\}) = - \sum_{i=1}^N \ln \left( \sum_{k=1}^K \exp \left( -\frac{1}{2} \|\mathbf{x}_i - f_\phi(\boldsymbol{\mu}_k)\|^2 \right) \right). \quad (38)$$

**Gradients.**

$$\begin{aligned} \alpha_{i,k} &= \text{softmax}_k \left( -\frac{1}{2} \|\mathbf{x}_i - f_\phi(\boldsymbol{\mu}_k)\|^2 \right). \\ -\frac{\partial E^{\text{PC}}}{\partial \mathbf{x}_i} &= \sum_{k=1}^K \alpha_{i,k} (\mathbf{x}_i - f_\phi(\boldsymbol{\mu}_k)). \\ -\frac{\partial E^{\text{PC}}}{\partial \boldsymbol{\mu}_k} &= \sum_{i=1}^N \alpha_{i,k} \underbrace{\frac{\partial f_\phi(\boldsymbol{\mu}_k)}{\partial \boldsymbol{\mu}_k}}_{\text{Jacobian of } f_\phi} (\mathbf{x}_i - f_\phi(\boldsymbol{\mu}_k)). \end{aligned}$$

Here,  $\frac{\partial f_\phi(\boldsymbol{\mu}_k)}{\partial \boldsymbol{\mu}_k}$  is the  $d \times d$  Jacobian (or more general shape if  $\boldsymbol{\mu}_k$  and  $f_\phi(\boldsymbol{\mu}_k)$  differ in dimension).

## 9 Cross Attention

**Setup.** We have a set of child vectors (queries)  $\mathbf{Q} \in \mathbb{R}^{d \times N_Q}$  and a set of parent vectors (keys)  $\mathbf{K} \in \mathbb{R}^{d \times N_K}$ . Let

$$C = \{1, \dots, N_Q\}, \quad P = \{1, \dots, N_K\},$$

so  $\mathbf{v}_c = \mathbf{q}_c$  is the  $c$ -th query, and  $\mathbf{v}_p = \mathbf{k}_p$  is the  $p$ -th key. Suppose we have learnable weight matrices  $\mathbf{W}^Q, \mathbf{W}^K \in \mathbb{R}^{d \times d}$ . Then

$$\mathbf{q}_c = \mathbf{W}^Q \mathbf{x}_c^Q, \quad \mathbf{k}_p = \mathbf{W}^K \mathbf{x}_p^K,$$

where  $\mathbf{x}_c^Q$  is the raw  $c$ -th query token and  $\mathbf{x}_p^K$  the raw  $p$ -th key token.

**Similarity function.**

$$\text{sim}(\mathbf{q}_c, \mathbf{k}_p) = \mathbf{q}_c^\top \mathbf{k}_p.$$

**Energy.**

$$E^{\text{Cross}}(\{\mathbf{q}_c\}, \{\mathbf{k}_p\}) = - \sum_{c=1}^{N_Q} \ln \left( \sum_{p=1}^{N_K} \exp(\mathbf{q}_c^\top \mathbf{k}_p) \right). \quad (39)$$

**Gradients.**

$$-\frac{\partial E^{\text{Cross}}}{\partial \mathbf{q}_c} = \sum_{p=1}^{N_K} \text{softmax}_p(\mathbf{q}_c^\top \mathbf{k}_p) \mathbf{k}_p. \quad (40)$$

$$-\frac{\partial E^{\text{Cross}}}{\partial \mathbf{k}_p} = \sum_{c=1}^{N_Q} \text{softmax}_p(\mathbf{q}_c^\top \mathbf{k}_p) \mathbf{q}_c. \quad (41)$$

When mapping back to the raw tokens  $\mathbf{x}_c^Q$  or  $\mathbf{x}_p^K$ , chain-rule multiplies by  $\mathbf{W}^Q$  or  $\mathbf{W}^K$ , respectively.

## 10 Switching Linear Dynamical System

**Setup.** Consider a time series of observations  $\{\mathbf{x}_t\}_{t=1}^T$ , where each  $\mathbf{x}_t \in \mathbb{R}^d$ . We assume there are  $K$  distinct (linear) dynamical modes, each with parameters  $\{\mathbf{A}_k, \mathbf{b}_k, \boldsymbol{\Sigma}_k\}$ . Let  $\pi_k$  be the mixing weight of mode  $k$ . A typical *switching linear dynamical system* (SLDS) posits:

$$\mathbf{x}_{t+1} \approx \mathbf{A}_k \mathbf{x}_t + \mathbf{b}_k, \quad k \in \{1, \dots, K\},$$

with Gaussian noise  $\boldsymbol{\Sigma}_k$ . We treat  $\mathbf{x}_{t+1}$  as a *child* and the  $\{\mathbf{A}_k, \mathbf{b}_k\}$  (together with  $\mathbf{x}_t$ ) as *parents* in a mixture-of-linear-dynamics fashion.

**Similarity function.** Define, for each mode  $k$ ,

$$\text{sim}(\mathbf{x}_{t+1}, \mathbf{x}_t; \mathbf{A}_k, \mathbf{b}_k) = \ln \pi_k - \frac{1}{2} (\mathbf{x}_{t+1} - \mathbf{A}_k \mathbf{x}_t - \mathbf{b}_k)^\top \Sigma_k^{-1} (\mathbf{x}_{t+1} - \mathbf{A}_k \mathbf{x}_t - \mathbf{b}_k).$$

**Energy.** Summing over all time steps  $t = 1, \dots, T-1$ , we write

$$E^{\text{SLDS}}(\{\mathbf{x}_t\}, \{\mathbf{A}_k, \mathbf{b}_k\}) = - \sum_{t=1}^{T-1} \ln \left( \sum_{k=1}^K \exp(\text{sim}(\mathbf{x}_{t+1}, \mathbf{x}_t; \mathbf{A}_k, \mathbf{b}_k)) \right). \quad (42)$$

**Gradients.** Let

$$\alpha_{t,k} = \text{softmax}_k(\text{sim}(\mathbf{x}_{t+1}, \mathbf{x}_t; \mathbf{A}_k, \mathbf{b}_k)),$$

i.e. the normalized exponent for mode  $k$ .

$$\begin{aligned} -\frac{\partial E^{\text{SLDS}}}{\partial \mathbf{x}_{t+1}} &= \sum_{k=1}^K \alpha_{t,k} \Sigma_k^{-1} (\mathbf{x}_{t+1} - \mathbf{A}_k \mathbf{x}_t - \mathbf{b}_k). \\ -\frac{\partial E^{\text{SLDS}}}{\partial \mathbf{x}_t} &= \sum_{k=1}^K \alpha_{t,k} (-\mathbf{A}_k^\top \Sigma_k^{-1}) (\mathbf{x}_{t+1} - \mathbf{A}_k \mathbf{x}_t - \mathbf{b}_k) \quad (t = 1, \dots, T-1). \\ -\frac{\partial E^{\text{SLDS}}}{\partial \mathbf{A}_k} &= \sum_{t=1}^{T-1} \alpha_{t,k} \Sigma_k^{-1} (\mathbf{x}_{t+1} - \mathbf{A}_k \mathbf{x}_t - \mathbf{b}_k) \mathbf{x}_t^\top. \\ -\frac{\partial E^{\text{SLDS}}}{\partial \mathbf{b}_k} &= \sum_{t=1}^{T-1} \alpha_{t,k} \Sigma_k^{-1} (\mathbf{x}_{t+1} - \mathbf{A}_k \mathbf{x}_t - \mathbf{b}_k). \end{aligned}$$

## 11 Layer Normalization

**Setup.** We consider a batch of input vectors  $\{\mathbf{x}_i\}_{i=1}^B$ , where each vector  $\mathbf{x}_i \in \mathbb{R}^D$ . Each vector is normalized by subtracting the mean and dividing by the standard deviation, with learnable scaling and bias parameters  $\gamma, \delta \in \mathbb{R}^D$ . For numerical stability, a small constant  $\epsilon > 0$  is added to the variance.

**Energy.** The energy for layer normalization is given by:

$$E^{\text{LN}}(\{\mathbf{x}_i\}) = \sum_{i=1}^B \left[ \gamma \sqrt{\frac{1}{D} \sum_{j=1}^D (x_{ij} - \bar{\mathbf{x}}_i)^2 + \epsilon} + \sum_{j=1}^D \delta_j x_{ij} \right],$$

where:

$$\bar{\mathbf{x}}_i = \frac{1}{D} \sum_{j=1}^D x_{ij}.$$

**Derivative.** The normalized outputs are obtained as the derivative of the energy with respect to the inputs:

$$\begin{aligned} \frac{\partial E^{\text{LN}}}{\partial x_{ij}} &= \gamma_j \frac{x_{ij} - \bar{\mathbf{x}}_i}{\sqrt{\frac{1}{D} \sum_{k=1}^D (x_{ik} - \bar{\mathbf{x}}_i)^2 + \epsilon}} + \delta_j. \\ \bar{\mathbf{x}}_i &= \frac{1}{D} \sum_{k=1}^D x_{ik}. \end{aligned}$$



## 12 Coordinate ascent

Consider an energy of the form

$$E(\{\mathbf{x}_j\}, \{\boldsymbol{\mu}_i\}; \theta) = - \sum_{j=1}^N \ln \left( \sum_{i=1}^S \exp(\text{sim}_\theta(\mathbf{x}_j, \boldsymbol{\mu}_i)) \right),$$

with  $\theta$  denoting the parameters of the similarity function. In an EM procedure, we alternate:

- **E-step:** Update latent variables using fixed  $\theta$ .
- **M-step:** Update parameters  $\theta$  in closed form given fixed latent assignments.

For *slot attention*, we have:

$$\begin{aligned} \theta &= \{\mathbf{W}_K, \mathbf{W}_Q\}, \\ \text{sim}_\theta(\mathbf{x}_j, \boldsymbol{\mu}_i) &= (\mathbf{W}_K \mathbf{x}_j)^\top (\mathbf{W}_Q \boldsymbol{\mu}_i). \\ \mathbf{A} \text{ with entries } A_{ji} &= \text{softmax}_i \left( (\mathbf{W}_K \mathbf{x}_j)^\top (\mathbf{W}_Q \boldsymbol{\mu}_i) \right). \end{aligned}$$

### E-step: Update Slots

Fix  $\mathbf{W}_K, \mathbf{W}_Q$ . The gradient for each slot  $\boldsymbol{\mu}_i$  is

$$-\frac{\partial E^{\text{Slot}}}{\partial \boldsymbol{\mu}_i} = \sum_{j=1}^N A_{ji} \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{x}_j.$$

Use this gradient to iteratively update  $\{\boldsymbol{\mu}_i\}$  until convergence.

### M-step: Update Parameters

Given fixed slots  $\{\boldsymbol{\mu}_i\}$  and attention matrix  $\mathbf{A}$ , we aim to update the parameters  $\mathbf{W}_K, \mathbf{W}_Q$ . This is motivated by setting the gradient of the energy with respect to these parameters to zero:

$$-\frac{\partial E^{\text{Slot}}}{\partial \mathbf{W}_K} = 0, \quad -\frac{\partial E^{\text{Slot}}}{\partial \mathbf{W}_Q} = 0.$$

Under the fixed assignments provided by  $\mathbf{A}$  and slots  $\{\boldsymbol{\mu}_i\}$ , these conditions are equivalent to solving a weighted least-squares problem. Specifically, we consider minimizing the objective

$$\min_{\mathbf{W}_K, \mathbf{W}_Q} \sum_{j=1}^N \sum_{i=1}^S A_{ji} \left\| \mathbf{W}_K \mathbf{x}_j - \mathbf{W}_Q \boldsymbol{\mu}_i \right\|^2,$$

since the stationary point of this quadratic form corresponds to zero gradients with respect to  $\mathbf{W}_K, \mathbf{W}_Q$ .

**Update  $\mathbf{W}_Q$ :** Differentiate w.r.t.  $\mathbf{W}_Q$ , set to zero:

$$\sum_{j,i} A_{ji} \left( \mathbf{W}_K \mathbf{x}_j - \mathbf{W}_Q \boldsymbol{\mu}_i \right) \boldsymbol{\mu}_i^\top = 0,$$

yielding

$$\mathbf{W}_Q = \left( \sum_{j,i} A_{ji} \mathbf{W}_K \mathbf{x}_j \boldsymbol{\mu}_i^\top \right) \left( \sum_{j,i} A_{ji} \boldsymbol{\mu}_i \boldsymbol{\mu}_i^\top \right)^{-1}.$$

**Update  $\mathbf{W}_K$ :** Similarly, differentiate w.r.t.  $\mathbf{W}_K$ , set to zero:

$$\sum_{j,i} A_{ji} \left( \mathbf{W}_K \mathbf{x}_j - \mathbf{W}_Q \boldsymbol{\mu}_i \right) \mathbf{x}_j^\top = 0,$$

yielding

$$\mathbf{W}_K = \left( \sum_{j,i} A_{ji} \mathbf{W}_Q \boldsymbol{\mu}_i \mathbf{x}_j^\top \right) \left( \sum_{j,i} A_{ji} \mathbf{x}_j \mathbf{x}_j^\top \right)^{-1}.$$

**Iterate:** Alternate between the E-step (updating  $\{\boldsymbol{\mu}_i\}$ ) and the M-step (updating  $\mathbf{W}_K, \mathbf{W}_Q$ ) until convergence.

## A Appendix: Derivation of Gradient Updates

Let us consider a generic term:

$$-\ln\left(\sum_{m=1}^M \exp(f_m(\mathbf{x}))\right),$$

where  $\mathbf{x} \in \mathbb{R}^d$  is some variable, and each  $f_m$  is a scalar function. We compute its gradient:

$$\begin{aligned} \frac{\partial}{\partial \mathbf{x}} \left[ -\ln\left(\sum_{m=1}^M \exp(f_m(\mathbf{x}))\right) \right] &= -\frac{1}{\sum_{m'} \exp(f_{m'}(\mathbf{x}))} \sum_{m=1}^M \exp(f_m(\mathbf{x})) \frac{\partial f_m(\mathbf{x})}{\partial \mathbf{x}} \\ &= -\sum_{m=1}^M \left[ \frac{\exp(f_m(\mathbf{x}))}{\sum_{m'} \exp(f_{m'}(\mathbf{x}))} \right] \frac{\partial f_m(\mathbf{x})}{\partial \mathbf{x}}. \end{aligned}$$

Defining  $\text{softmax}_m(f(\mathbf{x})) = \frac{\exp(f_m(\mathbf{x}))}{\sum_{m'} \exp(f_{m'}(\mathbf{x}))}$ , this is

$$-\sum_{m=1}^M \text{softmax}_m(f(\mathbf{x})) \frac{\partial f_m(\mathbf{x})}{\partial \mathbf{x}},$$

which matches the softmax-weighted gradient structure.

## B Appendix

Here, we collect the explicit partial derivatives of sim for each model discussed.

### Gaussian Mixture Models

$$\begin{aligned} \text{sim}(\mathbf{x}_i, \boldsymbol{\mu}_k) &= \ln \pi_k - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k). \\ \frac{\partial}{\partial \boldsymbol{\mu}_k} \text{sim}(\mathbf{x}_i, \boldsymbol{\mu}_k) &= \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k). \end{aligned}$$

### Cross Attention

$$\begin{aligned} \text{sim}(\mathbf{q}_c, \mathbf{k}_p) &= \mathbf{q}_c^\top \mathbf{k}_p. \\ \frac{\partial}{\partial \mathbf{q}_c} \text{sim}(\mathbf{q}_c, \mathbf{k}_p) &= \mathbf{k}_p, \quad \frac{\partial}{\partial \mathbf{k}_p} \text{sim}(\mathbf{q}_c, \mathbf{k}_p) = \mathbf{q}_c. \end{aligned}$$

### Hopfield Networks

$$\begin{aligned} \text{sim}(\mathbf{x}_i, \mathbf{m}_\mu) &= \mathbf{x}_i^\top \mathbf{m}_\mu. \\ \frac{\partial}{\partial \mathbf{x}_i} \text{sim}(\mathbf{x}_i, \mathbf{m}_\mu) &= \mathbf{m}_\mu, \quad \frac{\partial}{\partial \mathbf{m}_\mu} \text{sim}(\mathbf{x}_i, \mathbf{m}_\mu) = \mathbf{x}_i. \end{aligned}$$

### Slot Attention

$$\begin{aligned} \text{sim}(\mathbf{x}_j, \boldsymbol{\mu}_i) &= (\mathbf{W}_K \mathbf{x}_j)^\top (\mathbf{W}_Q \boldsymbol{\mu}_i). \\ \frac{\partial}{\partial \mathbf{x}_j} \text{sim}(\mathbf{x}_j, \boldsymbol{\mu}_i) &= \mathbf{W}_K^\top \mathbf{W}_Q \boldsymbol{\mu}_i, \quad \frac{\partial}{\partial \boldsymbol{\mu}_i} \text{sim}(\mathbf{x}_j, \boldsymbol{\mu}_i) = \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{x}_j. \end{aligned}$$

### Self-Attention

$$\begin{aligned} \text{sim}(\mathbf{x}_c, \mathbf{x}_p) &= (\mathbf{W}^Q \mathbf{x}_c)^\top (\mathbf{W}^K \mathbf{x}_p). \\ \frac{\partial}{\partial \mathbf{x}_c} \text{sim}(\mathbf{x}_c, \mathbf{x}_p) &= \mathbf{W}^{Q\top} \mathbf{W}^K \mathbf{x}_p, \quad \frac{\partial}{\partial \mathbf{x}_p} \text{sim}(\mathbf{x}_c, \mathbf{x}_p) = \mathbf{W}^{K\top} \mathbf{W}^Q \mathbf{x}_c. \end{aligned}$$

## Switching Linear Dynamical System

$$\begin{aligned}
\text{sim}(\mathbf{x}_{t+1}, \mathbf{x}_t; \mathbf{A}_k, \mathbf{b}_k) &= \ln \pi_k - \frac{1}{2} (\mathbf{x}_{t+1} - \mathbf{A}_k \mathbf{x}_t - \mathbf{b}_k)^\top \Sigma_k^{-1} (\mathbf{x}_{t+1} - \mathbf{A}_k \mathbf{x}_t - \mathbf{b}_k). \\
\frac{\partial}{\partial \mathbf{x}_{t+1}} \text{sim}(\mathbf{x}_{t+1}, \mathbf{x}_t; \mathbf{A}_k, \mathbf{b}_k) &= \Sigma_k^{-1} (\mathbf{x}_{t+1} - \mathbf{A}_k \mathbf{x}_t - \mathbf{b}_k), \\
\frac{\partial}{\partial \mathbf{x}_t} \text{sim}(\mathbf{x}_{t+1}, \mathbf{x}_t; \mathbf{A}_k, \mathbf{b}_k) &= -\mathbf{A}_k^\top \Sigma_k^{-1} (\mathbf{x}_{t+1} - \mathbf{A}_k \mathbf{x}_t - \mathbf{b}_k), \\
\frac{\partial}{\partial \mathbf{A}_k} \text{sim}(\mathbf{x}_{t+1}, \mathbf{x}_t; \mathbf{A}_k, \mathbf{b}_k) &= \Sigma_k^{-1} (\mathbf{x}_{t+1} - \mathbf{A}_k \mathbf{x}_t - \mathbf{b}_k) \mathbf{x}_t^\top, \\
\frac{\partial}{\partial \mathbf{b}_k} \text{sim}(\mathbf{x}_{t+1}, \mathbf{x}_t; \mathbf{A}_k, \mathbf{b}_k) &= \Sigma_k^{-1} (\mathbf{x}_{t+1} - \mathbf{A}_k \mathbf{x}_t - \mathbf{b}_k).
\end{aligned}$$

## Predictive Coding

$$\begin{aligned}
\text{sim}(\mathbf{x}_i, \boldsymbol{\mu}_k) &= -\frac{1}{2} \|\mathbf{x}_i - f_\phi(\boldsymbol{\mu}_k)\|^2. \\
\frac{\partial}{\partial \mathbf{x}_i} \text{sim}(\mathbf{x}_i, \boldsymbol{\mu}_k) &= \mathbf{x}_i - f_\phi(\boldsymbol{\mu}_k), \\
\frac{\partial}{\partial \boldsymbol{\mu}_k} \text{sim}(\mathbf{x}_i, \boldsymbol{\mu}_k) &= \frac{\partial f_\phi(\boldsymbol{\mu}_k)}{\partial \boldsymbol{\mu}_k} (\mathbf{x}_i - f_\phi(\boldsymbol{\mu}_k)).
\end{aligned}$$