

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

Alexander Tschantz

January 30, 2025

# Overview

$\log \sum \exp$  framework

$\log \sum \exp$  examples

$\log \sum \exp$  graphical models

$\log \sum \exp$  spatio-temporal model

Renormalised mixture models

## $\log \sum \exp$ energy

We consider energy functions that map a set of parents  $\{\mathbf{v}_p \in \mathbb{R}^{d_p} : p \in P\}$  to a set of children  $\{\mathbf{v}_c \in \mathbb{R}^{d_c} : c \in C\}$ .

Each energy function defines a similarity function (with parameters  $\theta$ ) that measures agreement between a parent and child vector:

$$\text{sim}_{\theta}(\mathbf{v}_c, \mathbf{v}_p) : \mathbb{R}^{d_c} \times \mathbb{R}^{d_p} \rightarrow \mathbb{R}.$$

The  $\log \sum \exp$  energy function is then given by:

$$E(\{\mathbf{v}_c\}, \{\mathbf{v}_p\}, \theta) = - \sum_{c \in C} \ln \left( \sum_{p \in P} \exp(\text{sim}_{\theta}(\mathbf{v}_c, \mathbf{v}_p)) \right).$$

## $\log \sum \exp$ derivatives

We define *attention* as:

$$\alpha_{c,p} = \text{softmax}_p(\text{sim}_\theta(\mathbf{v}_c, \mathbf{v}_p)) = \frac{\exp(\text{sim}_\theta(\mathbf{v}_c, \mathbf{v}_p))}{\sum_{p' \in P} \exp(\text{sim}_\theta(\mathbf{v}_c, \mathbf{v}_{p'}))}.$$

The derivatives with are then given by:

$$\begin{aligned} -\frac{\partial E}{\partial \mathbf{v}_c} &= \sum_{p \in P} \alpha_{c,p} \frac{\partial \text{sim}_\theta(\mathbf{v}_c, \mathbf{v}_p)}{\partial \mathbf{v}_c}, \\ -\frac{\partial E}{\partial \mathbf{v}_p} &= \sum_{c \in C} \alpha_{c,p} \frac{\partial \text{sim}_\theta(\mathbf{v}_c, \mathbf{v}_p)}{\partial \mathbf{v}_p}, \\ -\frac{\partial E}{\partial \theta} &= \sum_{c \in C} \sum_{p \in P} \alpha_{c,p} \frac{\partial \text{sim}_\theta(\mathbf{v}_c, \mathbf{v}_p)}{\partial \theta}. \end{aligned}$$

## $\log \sum \exp$ derivatives

The gradient of the energy function can be interpreted as an *expected value* over a discrete distribution  $\alpha_{c,p} = P(p \mid c)$ :

$$-\frac{\partial E}{\partial \mathbf{v}_c} = \mathbb{E}_{p \sim P(p|c)} \left[ \frac{\partial \text{sim}_\theta(\mathbf{v}_c, \mathbf{v}_p)}{\partial \mathbf{v}_c} \right] \quad (1)$$

While the child gradient is an exact expectation under  $P(p \mid c)$ , the parent gradient does not correspond exactly to an expectation under  $P(c \mid p)$ ; it is instead proportional to it through Bayes rule.

## $\log \sum \exp$ graphical models

We consider a single set of  $N$  nodes  $\mathbf{v} = \{\mathbf{v}_i : i \in \{1, 2, \dots, N\}\}$ , where each node  $\mathbf{v}_i \in \mathbb{R}^{d_i}$ , and  $M$  energy functions  $\{E_m : m \in \{1, 2, \dots, M\}\}$ .

Each  $E_m$  has a similarity function  $\text{sim}_m(\cdot)$  with parameters  $\theta_m$  and defines a subset of nodes acting as *children*  $C_m \subseteq \{1, 2, \dots, N\}$  and a subset acting as *parents*  $P_m \subseteq \{1, 2, \dots, N\}$ , which may overlap.

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

The full energy is the sum over all energy functions:

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

# $\log \sum$ exp expectation maximisation

## Expectation Step:

- ▶ Perform gradient descent on each  $\mathbf{v}_i$  to find the optimal latent states:

$$\mathbf{v}_i^* = \arg \min_{\mathbf{v}_i} E(\{\mathbf{v}_i\}, \{\theta_m\}).$$

## Maximisation Step:

- ▶ Given the updated  $\mathbf{v}_i^*$ , take a gradient step on each  $\theta_m$ :

$$\theta_m = \theta_m - \eta \frac{\partial E}{\partial \theta_m} \Big|_{\mathbf{v}_i^*}.$$

## $\log \sum \exp$ framework

```
class Energy:
    def similarity(self, *args):
        raise NotImplementedError

    def __call__(self, *args):
        # energy function
        sim_matrix = self.similarity(*args)
        return -sum(logsumexp(sim_matrix, axis=1))

dx = grad(energy)(x, y)
```



# Gaussian Mixture Model

We define a set of child nodes  $\{\mathbf{x}_i\}_{i=1}^N$  and parent means  $\{\boldsymbol{\mu}_k\}_{k=1}^K$ . The covariance matrices  $\{\boldsymbol{\Sigma}_k\}_{k=1}^K$  are treated as parameters  $\theta_k = \boldsymbol{\Sigma}_k$  (note that we could have also included means in the similarity parameters).

The similarity function is given by:

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

The energy function is:

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

The similarity function can be interpreted as the log probability of  $\mathbf{x}_i$  under the conditional Gaussian distribution  $\mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ , up to a normalization constant.

# Gaussian Mixture Model

The gradient for  $\mu_k$  is given by:

$$-\frac{\partial \mathbf{E}^{\text{GMM}}}{\partial \mu_k} = \sum_{i=1}^N \underbrace{\text{softmax}_k(\text{sim}(\mathbf{x}_i, \mu_k))}_{\alpha_{i,k}} \Sigma_k^{-1}(\mathbf{x}_i - \mu_k).$$

Solving for zero gives the fixed point update for  $\mu_k$ :

$$\mu_k = \frac{\sum_{i=1}^N \alpha_{i,k} \mathbf{x}_i}{\sum_{i=1}^N \alpha_{i,k}}.$$

# Hopfield Attention

We define a set of child nodes  $\{\mathbf{x}_i\}_{i=1}^N$ , each  $\mathbf{x}_i \in \mathbb{R}^d$ , and a set of parent memory vectors  $\{\mathbf{m}_\mu\}_{\mu=1}^K$ , each  $\mathbf{m}_\mu \in \mathbb{R}^d$ .

The energy function is given by:

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

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

The gradients are:

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

$$-\frac{\partial E}{\partial \mathbf{m}_\mu} = \sum_{i=1}^N \alpha_{i,\mu} \mathbf{x}_i, \quad -\frac{\partial E}{\partial \mathbf{x}_i} = \sum_{\mu=1}^K \alpha_{i,\mu} \mathbf{m}_\mu,$$

# Slot Attention

We define a set of child (token) nodes  $\{\mathbf{x}_j\}_{j=1}^N$  and a set of parent (slot) nodes  $\{\boldsymbol{\mu}_i\}_{i=1}^S$ . The parameters  $\theta$  consist of two projection matrices  $\mathbf{W}_K$  and  $\mathbf{W}_Q$ .

The energy function is:

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

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

The gradients are:

$$\alpha_{j,i} = \text{softmax}_i(\text{sim}(\mathbf{x}_j, \boldsymbol{\mu}_i)).$$

$$-\frac{\partial E}{\partial \boldsymbol{\mu}_i} = \sum_{j=1}^N \alpha_{j,i} \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{x}_j, \quad -\frac{\partial E}{\partial \mathbf{x}_j} = \sum_{i=1}^S \alpha_{j,i} \mathbf{W}_K^\top \mathbf{W}_Q \boldsymbol{\mu}_i.$$

# Self-Attention

We define a set of nodes  $\{\mathbf{x}_i\}_{i=1}^N$ , where each  $\mathbf{x}_i$  serves as both a child (query) and a parent (key). The parameters  $\theta$  consist of projection matrices  $\mathbf{W}^Q$  and  $\mathbf{W}^K$ , forming:

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

The energy function is:

$$\text{sim}(\mathbf{x}_c, \mathbf{x}_p) = \mathbf{q}_c^\top \mathbf{k}_p.$$
$$E^{\text{SA}} = - \sum_{c=1}^N \ln \left( \sum_{p=1}^N \exp(\mathbf{q}_c^\top \mathbf{k}_p) \right).$$

The gradients are:

$$\alpha_{c,p} = \text{softmax}_p(\mathbf{q}_c^\top \mathbf{k}_p).$$
$$-\frac{\partial E}{\partial \mathbf{x}_i} = \underbrace{\sum_{p=1}^N \alpha_{i,p} \mathbf{W}^{Q\top} \mathbf{W}^K \mathbf{x}_p}_{\text{Child side}} + \underbrace{\sum_{c=1}^N \alpha_{c,i} \mathbf{W}^{K\top} \mathbf{W}^Q \mathbf{x}_c}_{\text{Parent side}}.$$

# Linear Mixture Model

We define a set of child nodes  $\{\mathbf{x}_i\}_{i=1}^N$  and a set of parent nodes  $\{\mathbf{z}_k\}_{k=1}^K$ . The parameters  $\theta = \{(\mathbf{A}_k, \mathbf{b}_k, \Sigma_k)\}_{k=1}^K$  define a linear mapping:

$$\mathbf{x}_i = \mathbf{A}_k \mathbf{z}_k + \mathbf{b}_k + \epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \Sigma_k).$$

The energy function is:

$$\text{sim}_k(\mathbf{x}_i, \mathbf{z}_k) = -\frac{1}{2}(\mathbf{x}_i - \mathbf{A}_k \mathbf{z}_k - \mathbf{b}_k)^\top \Sigma_k^{-1}(\mathbf{x}_i - \mathbf{A}_k \mathbf{z}_k - \mathbf{b}_k).$$

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

Note this is the energy function for the switching linear dynamical system (SLDS) model, and our simple latent attention (SLA) model.

# Linear Mixture Model

The attention weights are:

$$\alpha_{i,k} = \text{softmax}_k(\text{sim}_k(\mathbf{x}_i, \mathbf{z}_k)).$$

The gradient w.r.t.  $\mathbf{x}_i$  is:

$$-\frac{\partial E}{\partial \mathbf{x}_i} = \sum_{k=1}^K \alpha_{i,k} \Sigma_k^{-1} (\mathbf{x}_i - \mathbf{A}_k \mathbf{z}_k - \mathbf{b}_k).$$

The fixed-point update for  $\mathbf{x}_i$  is:

$$\mathbf{x}_i^* = \sum_{k=1}^K \alpha_{i,k} (\mathbf{A}_k \mathbf{z}_k + \mathbf{b}_k).$$

# Non-Linear Mixture Model

We define a set of child nodes  $\{\mathbf{x}_i\}_{i=1}^N$  and a set of parent nodes  $\{\mathbf{z}_k\}_{k=1}^K$ . The parameters  $\theta = \{\mathbf{A}_k, \mathbf{b}_k, \sigma\}$  define a non-linear mapping:

$$f(\mathbf{z}_k) = \sigma(\mathbf{A}_k \mathbf{z}_k + \mathbf{b}_k),$$

where  $\sigma(\cdot)$  is a non-linearity.

The similarity function is:

$$\text{sim}_k(\mathbf{x}_i, \mathbf{z}_k) = -\frac{1}{2} \|\mathbf{x}_i - f(\mathbf{z}_k)\|^2.$$

The energy function is:

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



# Non-Linear Mixture Model: Gradients

The attention weights are:

$$\alpha_{i,k} = \text{softmax}_k(\text{sim}_k(\mathbf{x}_i, \mathbf{z}_k)).$$

The gradient w.r.t.  $\mathbf{x}_i$  is:

$$-\frac{\partial E}{\partial \mathbf{x}_i} = \sum_{k=1}^K \alpha_{i,k} (\mathbf{x}_i - f(\mathbf{z}_k)).$$

The gradient w.r.t.  $\mathbf{z}_k$  is:

$$-\frac{\partial E}{\partial \mathbf{z}_k} = \sum_{i=1}^N \alpha_{i,k} \frac{\partial f(\mathbf{z}_k)}{\partial \mathbf{z}_k} (\mathbf{x}_i - f(\mathbf{z}_k)).$$

Note these are these are predictive coding updates, weighted by the attention  $\alpha_{i,k}$ .

## Kronecker $\log \sum \exp$

We define a single set of child nodes  $\{\mathbf{x}_i\}_{i=1}^N$  and *two* sets of parents:

$$\{\mathbf{z}_k\}_{k=1}^K, \quad \{\mathbf{u}_\ell\}_{\ell=1}^L.$$

We combine these factors *multiplicatively* via a *Kronecker* structure, giving a single  $\log \sum \exp$  term over all pairs  $(k, \ell)$ .

**Similarity:** For each child  $\mathbf{x}_i$  and each parent pair  $(k, \ell)$ , define

$$\text{sim}(\mathbf{x}_i; \mathbf{z}_k, \mathbf{u}_\ell) = \text{sim}_z(\mathbf{x}_i, \mathbf{z}_k) + \text{sim}_u(\mathbf{x}_i, \mathbf{u}_\ell),$$

$$E^{\text{Kron}} = - \sum_{i=1}^N \ln \left( \sum_{k=1}^K \sum_{\ell=1}^L \exp(\text{sim}(\mathbf{x}_i; \mathbf{z}_k, \mathbf{u}_\ell)) \right).$$

## Kronecker log $\sum$ exp

The resulting attention is now a *joint* softmax over all parent pairs:

$$\begin{aligned}\alpha_{i,k,\ell} &= \text{softmax}_{k,\ell}(\text{sim}(\mathbf{x}_i; \mathbf{z}_k, \mathbf{u}_\ell)) \\ &= \frac{\exp(\text{sim}(\mathbf{x}_i; \mathbf{z}_k, \mathbf{u}_\ell))}{\sum_{k'=1}^K \sum_{\ell'=1}^L \exp(\text{sim}(\mathbf{x}_i; \mathbf{z}_{k'}, \mathbf{u}_{\ell'}))}.\end{aligned}$$

**Gradients:** For a child  $\mathbf{x}_i$ ,

$$-\frac{\partial E}{\partial \mathbf{x}_i} = \sum_{k=1}^K \sum_{\ell=1}^L \alpha_{i,k,\ell} \frac{\partial}{\partial \mathbf{x}_i} \text{sim}(\mathbf{x}_i; \mathbf{z}_k, \mathbf{u}_\ell).$$

Similarly, for a parent  $\mathbf{z}_k$ :

$$-\frac{\partial E}{\partial \mathbf{z}_k} = \sum_{i=1}^N \sum_{\ell=1}^L \alpha_{i,k,\ell} \frac{\partial}{\partial \mathbf{z}_k} \text{sim}_z(\mathbf{x}_i, \mathbf{z}_k),$$

# Kronecker log $\sum \exp$ (Discrete Case)

**Setup:** We define a single discrete child variable  $\mathbf{x} \in \{1, \dots, D\}$  and two discrete parent variables:

$$\mathbf{z} \in \{1, \dots, K\}, \quad \mathbf{u} \in \{1, \dots, L\}.$$

The parameters  $\boldsymbol{\theta} \in \mathbb{R}^{D \times K \times L}$  represent a joint categorical distribution.

**Energy** We define a single similarity function that depends on all three variables:

$$\text{sim}(\mathbf{x}, \mathbf{z}, \mathbf{u}) = \ln \theta_{\mathbf{x}, \mathbf{z}, \mathbf{u}}.$$

$$E^{\text{Kron}} = - \sum_{i=1}^N \ln \left( \sum_{k=1}^K \sum_{\ell=1}^L \exp(\text{sim}(\mathbf{x}_i, \mathbf{z}_k, \mathbf{u}_\ell)) \right).$$

# Kronecker log $\sum \exp$ (Discrete Case)

**Tensor Attention:** This induces a joint softmax over both parent variables:

$$\alpha_{i,k,\ell} = \frac{\exp(\text{sim}(x_i, z_k, u_\ell))}{\sum_{k'=1}^K \sum_{\ell'=1}^L \exp(\text{sim}(x_i, z_{k'}, u_{\ell'}))}.$$

## Interpretation:

- ▶ Parent variables  $(z_k, u_\ell)$  conspire to jointly explain child variables  $x_i$ .
- ▶ The attention  $\alpha_{i,k,\ell}$  generalizes categorical mixture models to factorial structures.

## Multi log $\sum$ exp

We define a set of child nodes  $\{\mathbf{x}_i\}_{i=1}^N$  and a set of parent nodes  $\{\mathbf{z}_k\}_{k=1}^K$ . Each parent  $k$  is associated with multiple similarity terms, indexed by different factors.

**Similarity functions:** Each similarity term contributes independently to the energy function:

$$\text{sim}_1(\mathbf{x}_i, \mathbf{z}_k), \quad \text{sim}_2(\mathbf{x}_i, \mathbf{z}_k), \quad \dots, \quad \text{sim}_M(\mathbf{x}_i, \mathbf{z}_k).$$

**Energy function:**

$$E^{\text{Multi}} = - \sum_{i=1}^N \ln \left( \sum_{k=1}^K \exp \left( \sum_{m=1}^M \text{sim}_m(\mathbf{x}_i, \mathbf{z}_k) \right) \right).$$

Here, all similarity terms are summed before the softmax, leading to a *joint* mixture model over multiple factors.

## Multi log $\sum$ exp

The resulting attention weight is:

$$\alpha_{i,k} = \text{softmax}_k \left( \sum_{m=1}^M \text{sim}_m(\mathbf{x}_i, \mathbf{z}_k) \right).$$

With gradients:

$$-\frac{\partial E}{\partial \mathbf{x}_i} = \sum_{k=1}^K \alpha_{i,k} \sum_{m=1}^M \frac{\partial}{\partial \mathbf{x}_i} \text{sim}_m(\mathbf{x}_i, \mathbf{z}_k).$$

$$-\frac{\partial E}{\partial \mathbf{z}_k} = \sum_{i=1}^N \alpha_{i,k} \sum_{m=1}^M \frac{\partial}{\partial \mathbf{z}_k} \text{sim}_m(\mathbf{x}_i, \mathbf{z}_k).$$

This formulation is useful when multiple similarities contribute to the similarity between a child and parent.

# Block-slot Attention

**Setup:** We define a set of child nodes  $\{\mathbf{x}_i\}_{i=1}^N$ , a set of slot parent nodes  $\{\mathbf{z}_k\}_{k=1}^K$ , and a set of memory parent nodes  $\{\mathbf{m}_\mu\}_{\mu=1}^M$ . The parameters  $\theta = \{\mathbf{W}_K, \mathbf{W}_Q\}$  consist of projection matrices.

**Energy:** Each child  $\mathbf{x}_i$  is compared to slot parents  $\mathbf{z}_k$  and memory parents  $\mathbf{m}_\mu$ :

$$\text{sim}_z(\mathbf{x}_i, \mathbf{z}_k) = (\mathbf{W}_K \mathbf{x}_i)^\top (\mathbf{W}_Q \mathbf{z}_k),$$

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

$$E^{\text{BlockSlot}} = - \sum_{i=1}^N \ln \left( \sum_{k=1}^K \exp(\text{sim}_z(\mathbf{x}_i, \mathbf{z}_k)) \right) \\ - \sum_{i=1}^N \ln \left( \sum_{\mu=1}^M \exp(\text{sim}_m(\mathbf{x}_i, \mathbf{m}_\mu)) \right).$$



# Energy Transformer

**Setup:** We define a set of nodes  $\{\mathbf{x}_i\}_{i=1}^N$  and a set of memory nodes  $\{\mathbf{m}_\mu\}_{\mu=1}^M$ . The parameters  $\theta = \{\mathbf{W}^Q, \mathbf{W}^K\}$  consist of projection matrices for queries and keys.

**Energy:** Each node  $\mathbf{x}_i$  attends to itself (self-attention) and to memory nodes  $\mathbf{m}_\mu$  (Hopfield attention):

$$\text{sim}_{\text{self}}(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{W}^Q \mathbf{x}_i)^\top (\mathbf{W}^K \mathbf{x}_j),$$

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

$$E^{\text{ET}} = - \sum_{i=1}^N \ln \left( \sum_{j=1}^N \exp(\text{sim}_{\text{self}}(\mathbf{x}_i, \mathbf{x}_j)) \right) \\ - \sum_{i=1}^N \ln \left( \sum_{\mu=1}^M \exp(\text{sim}_{\text{memory}}(\mathbf{x}_i, \mathbf{m}_\mu)) \right).$$

# Atari Model

**Setup:** We define a sequence of child nodes  $\{\mathbf{x}_t\}_{t=1}^T$ , a sequence of latent parent nodes  $\{\mathbf{z}_t\}_{t=1}^T$ , and a set of mode parameters  $\{\mathbf{A}_k, \mathbf{b}_k, \Sigma_k\}_{k=1}^K$ . The parameters  $\theta = \{\mathbf{A}_k, \mathbf{b}_k, \Sigma_k\}$  define the mappings.

**Similarity Functions:** The child  $\mathbf{x}_t$  is assigned to a mixture component  $\mathbf{z}_t$  using:

$$\text{sim}_{\text{GMM}}(\mathbf{x}_t, \mathbf{z}_t) = -\frac{1}{2}(\mathbf{x}_t - \mathbf{z}_t)^\top \Sigma_k^{-1}(\mathbf{x}_t - \mathbf{z}_t).$$

The latent child  $\mathbf{z}_t$  is explained by a linear mixture model:

$$\text{sim}_{\text{LM}}(\mathbf{z}_t, \mathbf{z}_{t-1}) = -\frac{1}{2}(\mathbf{z}_t - \mathbf{A}_k \mathbf{z}_{t-1} - \mathbf{b}_k)^\top \Sigma_k^{-1}(\mathbf{z}_t - \mathbf{A}_k \mathbf{z}_{t-1} - \mathbf{b}_k).$$

# Atari Model

## Energy:

$$E^{\text{Atari}} = - \sum_{t=1}^T \ln \left( \sum_{k=1}^K \exp(\text{sim}_{\text{GMM}}(\mathbf{x}_t, \mathbf{z}_t)) \right) \\ - \sum_{t=1}^{T-1} \ln \left( \sum_{k=1}^K \exp(\text{sim}_{\text{LM}}(\mathbf{z}_t, \mathbf{z}_{t-1})) \right).$$

The gradient for the latent parent  $\mathbf{z}_t$  combines both responsibilities:

$$-\frac{\partial E}{\partial \mathbf{z}_t} = \sum_{k=1}^K \alpha_{t,k} \Sigma_k^{-1}(\mathbf{x}_t - \mathbf{z}_t) + \sum_{k=1}^K \beta_{t,k} \Sigma_k^{-1}(\mathbf{z}_t - \mathbf{A}_k \mathbf{z}_{t-1} - \mathbf{b}_k).$$

# Spatiotemporal Attention

**Hierarchy:** We define a hierarchy of  $L$  layers. Each layer  $l$  contains  $K_l$  latent variables, where each variable evolves over  $T$  time steps. Concretely, let:

$$\mathbf{x}_{k,t}^{(l)} \in \mathbb{R}^{D_l}$$

denote the  $k$ -th variable (or "slot") in layer  $l$  at time  $t$ , with  $D_l$ -dimensional representations. The number of variables  $K_l$  may differ across layers.

## Structure:

- ▶ **Inter-layer (vertical):** Each variable  $\mathbf{x}_{k,t}^{(l)}$  is influenced by variables from the layer below ( $l - 1$ ) at the same time  $t$ .
- ▶ **Intra-layer (concurrent):** Variables within the same layer  $l$  interact at each time step  $t$ .
- ▶ **Temporal:** Each variable  $\mathbf{x}_{k,t}^{(l)}$  is influenced by its own past states  $\mathbf{x}_{k,t'}^{(l)}$  for  $t' < t$ .

# Spatiotemporal Attention

**Inter-layer (vertical):**

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

**Intra-layer (concurrent):**

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

**Temporal (past states):**

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

# Spatiotemporal Attention

**Inter-layer Energy (Slot attention):**

$$E_v^{(l)} = - \sum_{t=1}^T \sum_{c=1}^{K_{l-1}} \ln \left( \sum_{k=1}^{K_l} \exp(\text{sim}_v^{(l)}(\mathbf{x}_{c,t}^{(l-1)}, \mathbf{x}_{k,t}^{(l)})) \right).$$

**Intra-layer Energy (Self attention):**

$$E_c^{(l)} = - \sum_{t=1}^T \sum_{k=1}^{K_l} \ln \left( \sum_{\substack{k'=1 \\ k' \neq k}}^{K_l} \exp(\text{sim}_c^{(l)}(\mathbf{x}_{k,t}^{(l)}, \mathbf{x}_{k',t}^{(l)})) \right).$$

**Temporal Energy (Casual self attention):**

$$E_t^{(l)} = - \sum_{k=1}^{K_l} \sum_{t=2}^T \ln \left( \sum_{t' < t} \exp(\text{sim}_t^{(l)}(\mathbf{x}_{k,t}^{(l)}, \mathbf{x}_{k,t'}^{(l)})) \right).$$

# Spatiotemporal Attention

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

## Bottom-up gradient

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

## Top-down gradient

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

# Spatiotemporal Attention

## Intra-layer gradient

$$\begin{aligned} -\frac{\partial E_c^{(l)}}{\partial \mathbf{x}_{k,t}^{(l)}} &= \sum_{k' \neq k} \text{softmax}_k(\mathbf{A}_{k',k}) \mathbf{W}_c^{K,(l)\top} \mathbf{W}_c^{Q,(l)} \mathbf{x}_{k',t}^{(l)} \\ &\quad + \sum_{k' \neq k} \text{softmax}_{k'}(\mathbf{A}_{k,k'}) \mathbf{W}_c^{Q,(l)\top} \mathbf{W}_c^{K,(l)} \mathbf{x}_{k',t}^{(l)}. \end{aligned}$$

## Temporal gradient

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

**Remark:** By merging inter-layer, intra-layer, and temporal connections into a single large log-sum-exp term, we increase flexibility but at the cost of  $\mathcal{O}((N+K)T)^2$  complexity.



# Renormalised mixture models

We define a hierarchy of  $L$  layers, each containing  $K_l$  latent variables of dimension  $D_l$ , where  $K_l$  decreases with  $l$  ( $K_1 > K_2 > \dots > K_L$ ).

Each variable  $\mathbf{x}_{k,t}^{(l)}$  at layer  $l$  receives input only from a local subset of variables in layer  $l - 1$ , forming a *receptive field*:

$$\mathcal{R}_k^{(l)} \subseteq \{1, \dots, K_{l-1}\}.$$

In this setting, intra-layer and temporal energy remain the same.

## Renormalised Mixture Model

Each variable  $\mathbf{x}_{k,t}^{(l)}$  at layer  $l$  is explained only by a local receptive field  $\mathcal{R}_k^{(l)}$  in layer  $l - 1$ :

$$E_v^{(l)} = - \sum_{t=1}^T \sum_{k=1}^{K_l} \ln \left( \sum_{c \in \mathcal{R}_k^{(l)}} \exp(\text{sim}_v^{(l)}(\mathbf{x}_{c,t}^{(l-1)}, \mathbf{x}_{k,t}^{(l)})) \right).$$

Define:

$$\mathbf{A}_{c,k} = \text{sim}_v^{(l)}(\mathbf{x}_{c,t}^{(l-1)}, \mathbf{x}_{k,t}^{(l)}), \quad \mathbf{A}_{k,p} = \text{sim}_v^{(l+1)}(\mathbf{x}_{k,t}^{(l)}, \mathbf{x}_{p,t}^{(l+1)}).$$

Then the bottom-up and top-down gradients are:

$$-\frac{\partial E_v^{(l)}}{\partial \mathbf{x}_{k,t}^{(l)}} = \sum_{c \in \mathcal{R}_k^{(l)}} \text{softmax}_c(\mathbf{A}_{c,k}) \frac{\partial}{\partial \mathbf{x}_{k,t}^{(l)}} \text{sim}_v^{(l)}(\mathbf{x}_{c,t}^{(l-1)}, \mathbf{x}_{k,t}^{(l)}),$$

$$-\frac{\partial E_v^{(l+1)}}{\partial \mathbf{x}_{k,t}^{(l)}} = \sum_{p \in \mathcal{R}_k^{(l+1)}} \text{softmax}_p(\mathbf{A}_{k,p}) \frac{\partial}{\partial \mathbf{x}_{k,t}^{(l)}} \text{sim}_v^{(l+1)}(\mathbf{x}_{k,t}^{(l)}, \mathbf{x}_{p,t}^{(l+1)}).$$

## Bayesian model expansion

If a new data point  $\mathbf{x}_{N+1}$  is not well explained by existing parent components, we introduce a new parent to explain it.

We evaluate the energy contribution of the new data point:

$$E_{N+1} = -\ln \left( \sum_{k=1}^K \exp(\text{sim}(\mathbf{x}_{N+1}, \mathbf{z}_k)) \right).$$

If  $E_{N+1} > \tau$  (for some threshold  $\tau$ ), then  $\mathbf{x}_{N+1}$  is not sufficiently explained, and we add a new parent component  $\mathbf{z}_{K+1}$ , where the new parent is initialized based on  $\mathbf{x}_{N+1}$ .

# Coordinate Ascent

## Energy Function:

We consider an energy of the form

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

where  $\theta$  denotes the parameters of the similarity function.

## EM Procedure:

- ▶ **E-step:** Update latent variables  $\{\boldsymbol{\mu}_i\}$  given fixed  $\theta$ .
- ▶ **M-step:** Update parameters  $\theta$  in closed form given current slots.

# Coordinate Ascent

## Slot Attention Example:

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

Define the attention matrix  $\mathbf{A}$  with entries

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

**Goal:** Perform an **E-step** (update slots  $\boldsymbol{\mu}_i$ ) and **M-step** (update  $\mathbf{W}_K, \mathbf{W}_Q$ ).

# Coordinate Ascent

## M-step: Update Parameters

Given fixed slots  $\{\mu_i\}$  and attention  $\mathbf{A}$ , we update  $\mathbf{W}_K, \mathbf{W}_Q$ .  
We consider minimizing

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

Setting the gradient to zero yields a weighted least-squares problem in  $\mathbf{W}_K$  and  $\mathbf{W}_Q$ .

## Key idea:

- ▶  $\mathbf{A}$  is fixed (like responsibilities in EM).
- ▶ Solve for  $\mathbf{W}_K, \mathbf{W}_Q$  in closed form.

# Coordinate Ascent

## Closed-Form Updates

Differentiate and set to zero:

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

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

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

# Open questions

- ▶ Hierarchical POMDPs and hierarchical mixture models differ only in their similarity functions.
- ▶ Crucially, hierarchical mixture models induce a discrete state space (in terms of their attention matrices).
- ▶ This is seen in SLDS, where we model the dynamics of the induced discrete state space (which mode is active at a given time) as a HMM.
- ▶ Is there a way to combine the two models, with continuous mixture models on the bottom and discrete mixture models on the top?