# PGM in Modern AI Approaches

**Eric Walzthöny**

**Shanshan Zhang**

**2022-05-24**

# Contents



Transfer Learning

Deep Reinforcement Learning

Transformer

Deep Generative Models

Meta Learning

# Transformers

**Multi-Head Self-Attention**

$$A_h = \text{Softmax}(\alpha Q_h K_h^\top) V_h$$ **~1/2**

**Layer-norm and** residual connection
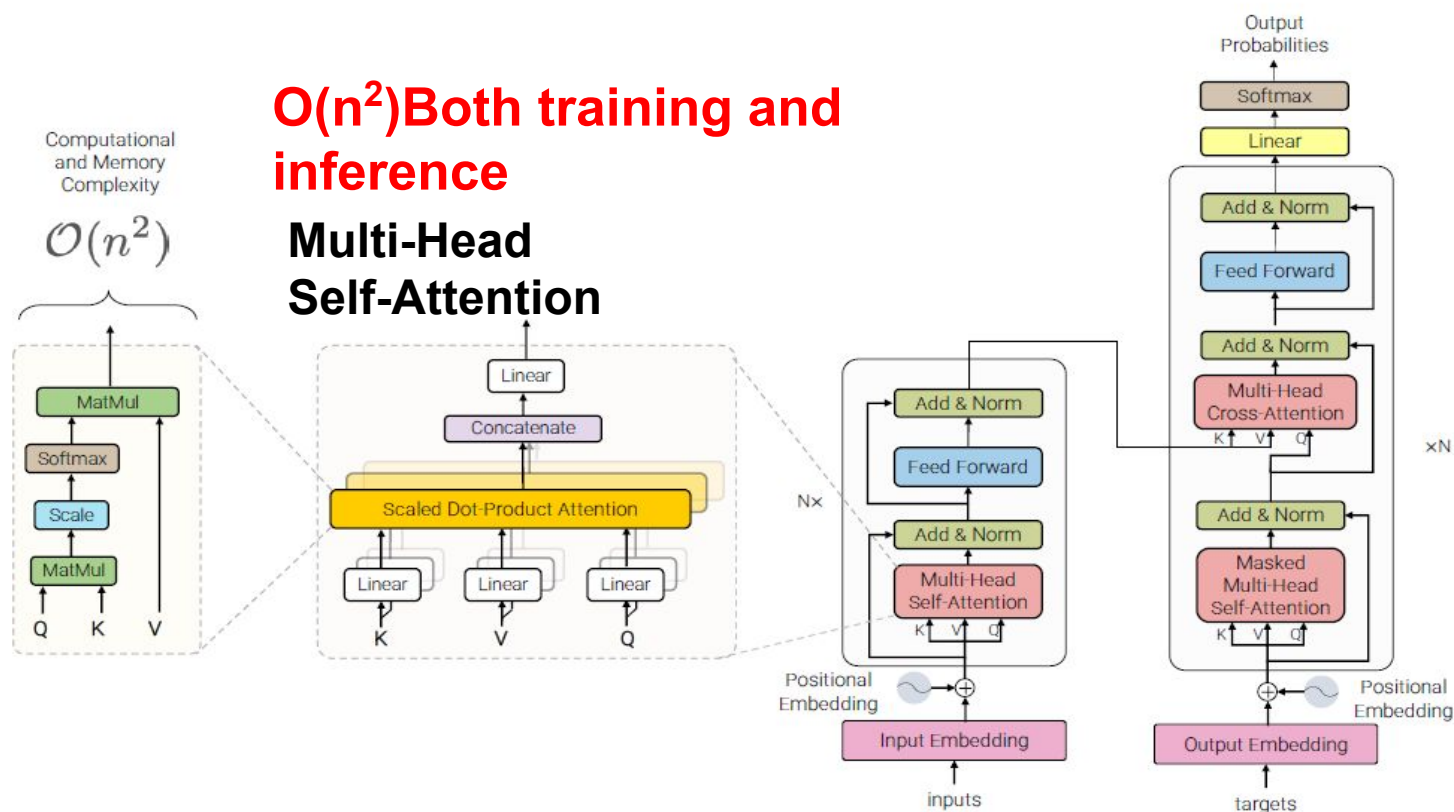
$$X = \text{LayerNorm}(F_S(X)) + X$$

**Position-wise Feed-forward layers** **~1/2**

$$F_2(ReLU(F_1(X_A)))$$

**Layer-norm and** residual connection

$$X_B = \text{LayerNorm}(X_A)) + X_A$$

Computational and Memory Complexity

$$\mathcal{O}(n^2)$$

**O(n²)Both training and inference**

**Multi-Head Self-Attention**
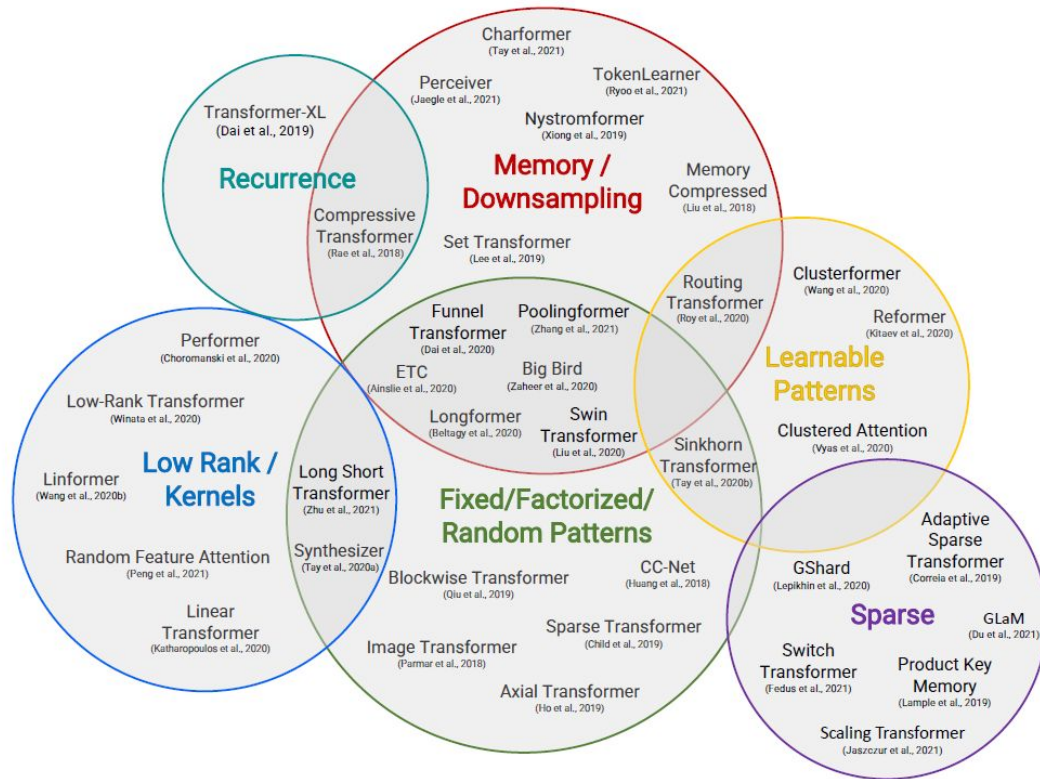


Attention is all you need https://arxiv.org/pdf/1706.03762.pdf

# Different Normalization

# Efficient Transformers
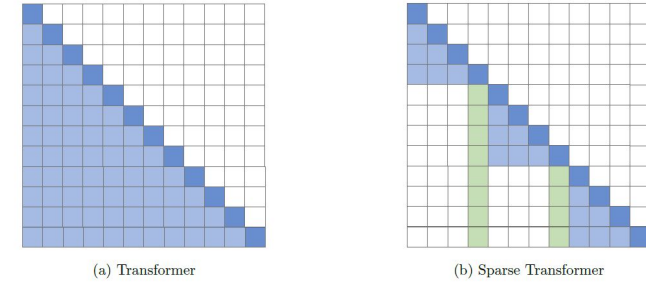


## Sparse Transformer



Figure 4: Illustration of patterns of the attention matrix for dense self-attention in Transformers and sparse fixed attention in Sparse Transformers. Blue in the right diagram represents the local self-attention while green represents the strided component of the sparse attention.
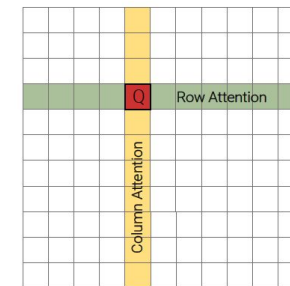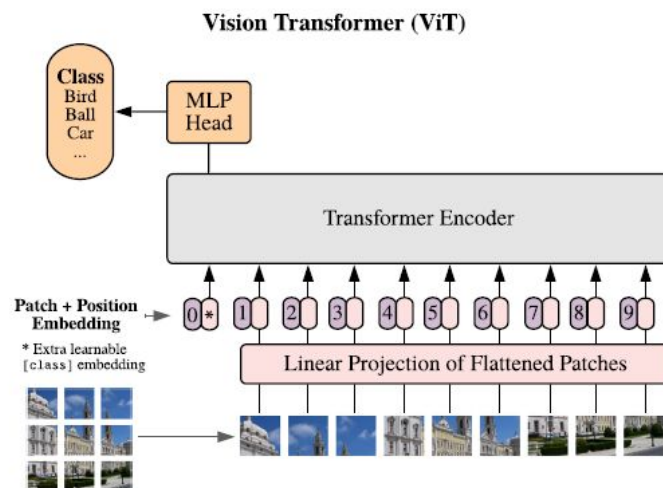
## Axical Transformer



Figure 5: Attention span in Axial Transformer on a two-dimensional input.
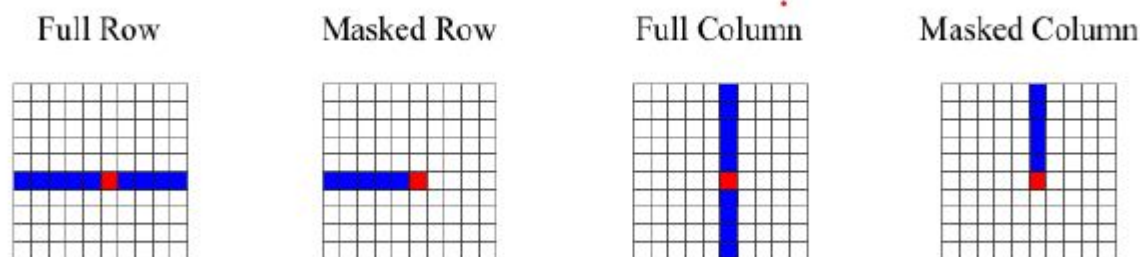
Efficient Transformers: https://arxiv.org/pdf/2009.06732.pdf

# Sparse vs. axial attention

**Vision Transformer (ViT)**

**Axial attention over axis $k$ can be implemented by transposing all axes except $k$ to the batch axis, calling standard attention as a subroutine**

$$O(N^{(d-1)/d})$$

$$O(n\sqrt[p]{n})$$



(a) Transformer      (b) Sparse Transformer (strided)      (c) Sparse Transformer (fixed)

Formally, $A_i^{(1)} = \{t, t+1, ..., i\}$ for $t = \max(0, i - l)$ and $A_i^{(2)} = \{j : (i - j) \bmod l = 0\}$. This pattern can be visualized in Figure 3(b).

Formally, $A_i^{(1)} = \{j : (\lfloor j/l \rfloor = \lfloor i/l \rfloor)\}$, where the brackets denote the floor operation, and $A_i^{(2)} = \{j : j \bmod l \in \{t, t+1, ..., l\}\}$, where $t = l - c$ and $c$ is a hyperparameter.

# VAE in NN perspective

# Deep Generative MODELS

## Explicit probabilistic models

Provide an explicit parametric specification of the distribution of $x$

l Tractable likelihood function $p\#(x)$
l E.g., Deep generative model parameterized with NNs **(e.g., VAEs)**
l

$$P_\emptyset(x|z) = N(x; \mu_\emptyset, \sigma)$$

$$P(z) = N(x; 0, I)$$
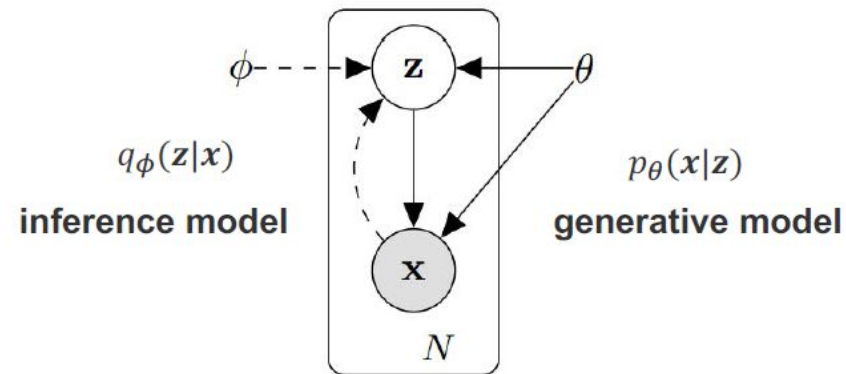


$\phi \dashrightarrow$ **z** $\leftarrow \theta$

$q_\phi(z|x)$

**inference model**

$p_\theta(x|z)$

**generative model**

**x**

$N$

Figure courtesy: Kingma & Welling, 2014

VAE: https://arxiv.org/pdf/1312.6114.pdf

# Variational Inference in VAE

Maximize log likelihood

$$\log[Pr(\mathbf{x}|\boldsymbol{\phi})] = \log\left[\int Pr(\mathbf{x},\mathbf{h}|\boldsymbol{\phi})d\mathbf{h}\right]$$

Jensen's inequality

$$\text{ELBO}[\boldsymbol{\theta},\boldsymbol{\phi}] = \int q(\mathbf{h}|\boldsymbol{\theta})\log\left[\frac{Pr(\mathbf{x},\mathbf{h}|\boldsymbol{\phi})}{q(\mathbf{h}|\boldsymbol{\theta})}\right]d\mathbf{h}.$$

$$= \int q(\mathbf{h}|\mathbf{x},\boldsymbol{\theta})\log[Pr(\mathbf{x}|\mathbf{h},\boldsymbol{\phi})]d\mathbf{h} - D_{KL}[q(\mathbf{h}|\mathbf{x},\boldsymbol{\theta}),Pr(\mathbf{h})]$$

Minimize Reconstruction Error

Minimize

$$\sum_{i=1}^{3}(exp(\sigma_i) - (1+\sigma_i) + (m_i)^2)$$



Variational distribution should be similar to prior

Data example should have high probability

$q(\mathbf{h}|\mathbf{x},\theta)$

$\text{ELBO}[\theta,\phi]$

Loss function

$\mathbf{x}$ → $g[\mathbf{x},\theta]$ → $\mu$, $\Sigma$ → $\mathbf{h}^* = \mu + \Sigma^{1/2}\epsilon^*$ → $\mathbf{f}[\mathbf{h}^*,\phi]$ → $Pr(\mathbf{x}|\mathbf{h}^*,\phi)$

$\text{Norm}_\epsilon[\mathbf{0},\mathbf{I}]$ → $\epsilon^*$

Sample

# DEEP GENERATIVE MODEL



## Implicit probabilistic models – **GANs**
- Defines a stochastic process to simulate data $x = G_\theta(z)$
- Define an implicit distribution over x: $P_{g_\theta}(x)$
- Do not require tractable likelihood function

$$\min_G \max_D V(D, G) = E_{x \sim p_{data}(x)}[\log D(x)]$$
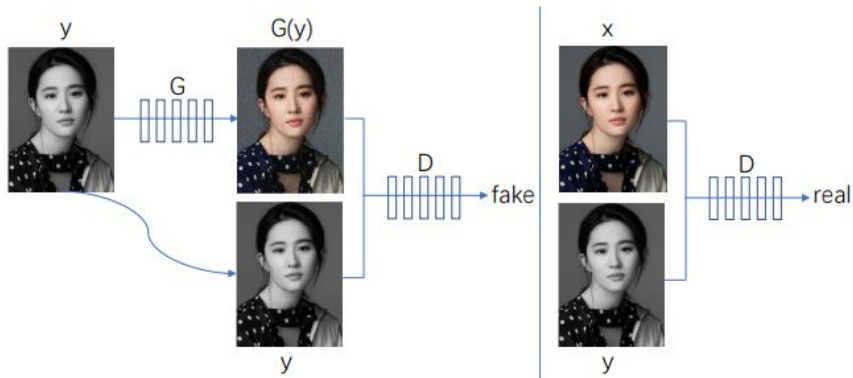$$+ E_{z \sim p_z(z)}[\log(1 - D(G(z)))].$$

- Generate data from a deterministic equation given parameters and random noise $z \sim N(0, I)$

- Intractable to evaluate likelihood

**Mode Collapse?!**

**Motivating entropy regularization!!**

Adding the entropy of G($z$)) to the objective



Goodfellow et al., 2014
https://arxiv.org/pdf/1406.2661.pdf

# GAN Variants

## Pixel2Pixel



### 3.1.1.1 Original minimax game:
The objective function of GANs [3] is

$$\min_{G} \max_{D} V(D,G) = E_{x \sim p_{data}(x)} [\log D(x)] + E_{z \sim p_z(z)} [\log (1 - D(G(z)))]. \qquad (1)$$

$\log D(x)$ is the cross-entropy between $[1 \quad 0]^T$ and $[D(x) \quad 1 - D(x)]^T$. Similarly, $\log(1 - D(G(z)))$ is the cross-entropy between $[0 \quad 1]^T$ and $[D(G(z)) \quad 1 - D(G(z))]^T$. For fixed $G$, the optimal discriminator $D$ is given by [3]:

$$D_G^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)}. \qquad (2)$$

The minmax game in (1) can be reformulated as:

$$C(G) = \max_{D} V(D,G)$$
$$= E_{x \sim p_{data}} [\log D_G^*(x)]$$
$$\quad + E_{z \sim p_z} [\log(1 - D_G^*(G(z)))]$$
$$= E_{x \sim p_{data}} [\log D_G^*(x)] + E_{x \sim p_g} [\log(1 - D_G^*(x))] \qquad (3)$$
$$= E_{x \sim p_{data}} \left[ \log \frac{p_{data}(x)}{\frac{1}{2}(p_{data}(x) + p_g(x))} \right]$$
$$\quad + E_{x \sim p_g} \left[ \frac{p_g(x)}{\frac{1}{2}(p_{data}(x) + p_g(x))} \right] - 2 \log 2.$$

The definition of KullbackLeibler (KL) divergence and Jensen-Shannon (JS) divergence between two probabilistic distributions $p(x)$ and $q(x)$ are defined as

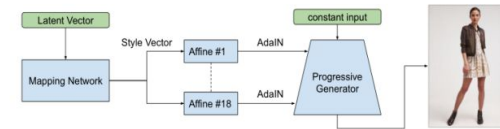$$KL(p \| q) = \int p(x) \log \frac{p(x)}{q(x)} dx, \qquad (4)$$

$$JS(p \| q) = \frac{1}{2} KL(p \| \frac{p+q}{2}) + \frac{1}{2} KL(q \| \frac{p+q}{2}). \qquad (5)$$
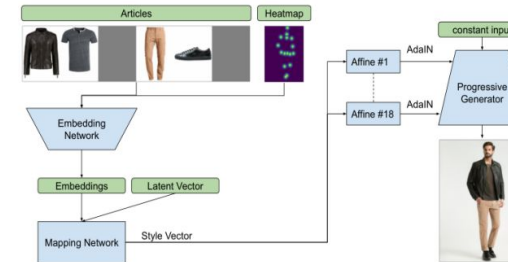
Therefore, (3) is equal to

$$C(G) = KL(p_{data} \| \frac{p_{data} + p_g}{2}) + KL(p_g \| \frac{p_{data} + p_g}{2})$$
$$\quad - 2 \log 2 \qquad (6)$$
$$= 2JS(p_{data} \| p_g) - 2 \log 2.$$

Thus, the objective function of GANs is related to both KL divergence and JS divergence.
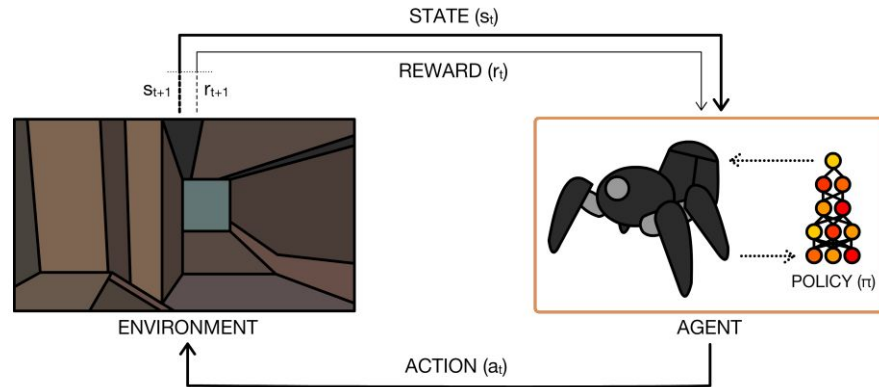
## Conditional



(a) Unconditional



GAN Survey (including 400+ variants):
https://arxiv.org/pdf/2001.06937.pdf

# Deep Reinforcement learning



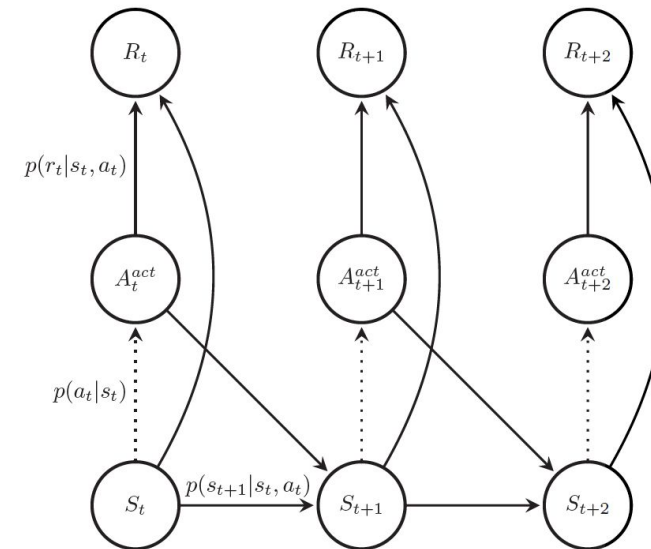There are four key ingredients for a RL system:

**Environment**: Physical world in which the agent operates

**State**: Current situation of the agent

**Reward**: Feedback from the environment

**Policy**: Method to map agent's state to actions

**Directed Acyclic Graph For Markov Decision Process**



DRL+PGM survey:  https://arxiv.org/pdf/1906.10025.pdf
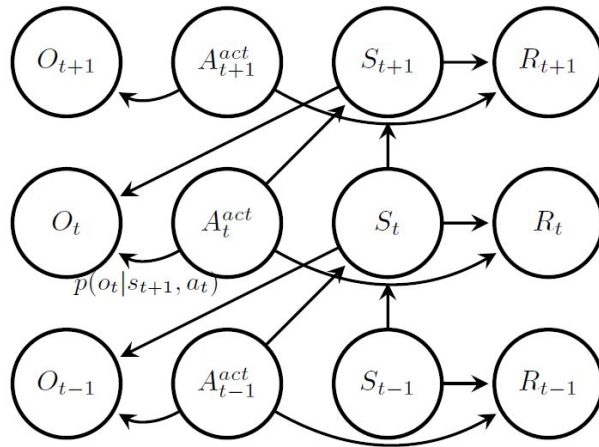
# Partially Observable MDP



Fig. 4. Probabilistic Graphical Model for POMDP

*The distributions on other edges are omitted in last slide*

Partially Observable Markov Decision process with its DAG representation shows that the agent could only observe the state partially by observing $O_t$ through a non invertible function of the next state $S_{t+1}$ and the action $a_t$, as indicated the Figure by $P(O_t| S_{t+1}, a_t)$
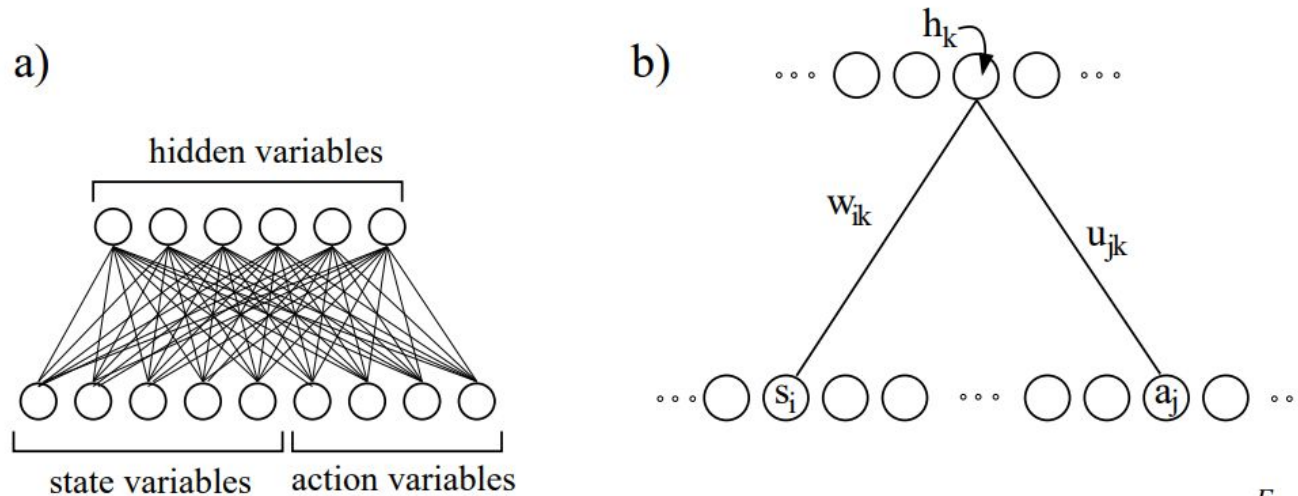
For POMDP, belief state $b_t$ $\quad \sum_{\mathcal{S}} b(S_t) = 1$

$$b_{t+1}(s_{t+1})$$
$$=p(s_{t+1} \mid o_t, a_t, b_t)$$
$$=p(o_t \mid s_{t+1}, a_t)\frac{\sum_{s_t} p(s_{t+1} \mid s_t, a_t)p(s_t \mid a_t, b_t)}{p(o_t \mid a_t, b_t)}$$

# RBM in DRL

$$p(a|s) = 1/Z(s)e^{-F(s,a)/T} = 1/Z(s)e^{Q(s,a)/T}$$

Approximate the value function of an MDP with the negative free energy of the restricted Boltzmann machine.

a)

hidden variables



state variables    action variables

*The state and action variables will be assumed to be discrete, and will be represented by the visible binary variables of the restricted Boltzmann machine.*

b)



Undirected graph defines a joint probability distribution over state and action pairs through hidden state

$$P(\mathbf{v}, \mathbf{h}) = \frac{\exp(-E(\mathbf{v}, \mathbf{h}))}{\sum_{\widehat{\mathbf{v}}, \widehat{\mathbf{h}}} \exp(-E(\widehat{\mathbf{v}}, \widehat{\mathbf{h}}))},$$

**Function Approximation Q(S, a)**

$$\exp(-F(\mathbf{v})) = \sum_{\mathbf{h}} \exp(-E(\mathbf{v}, \mathbf{h})),$$

$$P(\mathbf{v}) = \frac{\exp(-F(\mathbf{v}))}{\sum_{\widehat{\mathbf{v}}} \exp(-F(\widehat{\mathbf{v}}))}.$$

**Temporal Difference Learning**

$$E_{TD}(\mathbf{s}^t, \mathbf{a}^t) = [r^t + \gamma Q(\mathbf{s}^{t+1}, \mathbf{a}^{t+1})] - Q(\mathbf{s}^t, \mathbf{a}^t).$$

$$\frac{\partial F(\mathbf{v})}{\partial w_{ik}} = -v_i \langle h_k \rangle_{P(h_k|\mathbf{v})}$$

**The negative free energy to approximate the state action value function Q(s, a).**