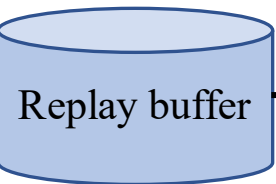
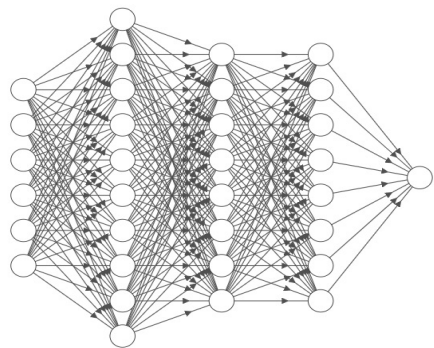


$$M_{\mathcal{B}} = (\mathcal{S}, \mathcal{S}', \mathcal{A}, P, R, \gamma)$$



Batch data used
as prior knowledge



BRL Model:

Finding the optimized policy that maximizes returns

$$\pi^*(s) = \underset{a}{\operatorname{argmax}} \text{ s.t. } (s,a) \in \mathcal{B} Q_{\mathcal{B}}^{\pi}(s, a)$$

TC_t
Thermal comfort index

$$a_t = \{a_i + \xi_{\phi}(s, a, \Phi)\}_{i=0}^n$$

Selected action decided by a VAE model,
a perturbation model ξ_{ϕ} , two Q-networks,
and a set of hyperparameters

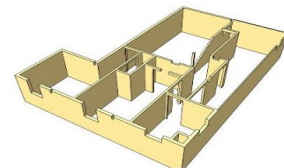
s_t

RL states feedback to calculate reward



Regression Model :
Predict thermal comfort
with current thermal states

PMV_t
PMV features



Real multi-zone
building environment