

The paper uses a model that takes as inputs $(\vec{o}_t, r_t, d_t, \vec{h}_t)$. I.e., the same architecture as the DQN network. However, they only input a single observation \vec{o}_t , no other variables.

Outside of training, the algorithm keeps a rolling average for the $\mu_{\vec{z}}$ and $\Sigma_{\vec{z}}$. Not sure why, these statistics are only used to compute the loss.

