

Algorithm -

Init - policy (π), critic(V), feature_extractor (G), action predictor (I) , forward model (F), Replay Memory (M)

For episodes in 1 to E:

For $t = 1 \dots T$:

$g^t = G(s^t)$

Sample action (a^t) = $\pi(g^t)$; $U(0,1) \geq \epsilon$
Random ; otherwise

$r^t, s^{t+1} = \text{env.step}(a^t)$

Store transition ($s^t, a^t, r^t, G(s^{t+1})$) in M

=====

Sample transition ($s^t, a^t, r^t, G(s^{t+1})$) from M

$g^t = G(s^t)$ # Recompute
 $g^{t+1} = G(s^{t+1})$ # From M

ICM Update

$a^{\text{pred}} = I(g^t, g^{t+1})$

$L_I = \text{cross_entropy}(a^{\text{pred}}, a^t)$

$g^{\text{pred}} = F(g^t, a^t)$

$L_F = \text{MSE}(g^{\text{pred}}, g^{t+1}) = r^{\text{intrinsic}}$

Update F and I . Note $G(\cdot)$ is updated at all time steps

Update A2C

$R = r^{\text{intrinsic}} + r^t + \gamma V(g^{t+1})$

$a^i = \text{sample}(\text{dist} = \pi(g^t))$

$\Theta(\pi) = \Theta(\pi) + \nabla(\pi) \log(\pi(a^i | g^t)) [R - V(g^t)]$

$\Theta(V) = \Theta(V) + \nabla(V) [R - V(g^t)]^2$

If Refresh == True:

$a^{\text{new}} = \text{one_hot}(a^i)$

$g^{\text{new}} = F(g^t, a^{\text{new}})$

Update Transition ($s^t, a^{\text{new}}, r^t, g^{\text{new}}$)
