| Hyper-parameter | Ours | DDPG |
|---|---|---|
| Critic Learning Rate | $10^{-3}$ | $10^{-3}$ |
| Critic Regularization | None | $10^{-2} \cdot ||\theta||^2$ |
| Actor Learning Rate | $10^{-3}$ | $10^{-4}$ |
| Actor Regularization | None | None |
| Optimizer | Adam | Adam |
| Target Update Rate ($\tau$) | $5 \cdot 10^{-3}$ | $10^{-3}$ |
| Batch Size | 100 | 64 |
| Iterations per time step | 1 | 1 |
| Discount Factor | 0.99 | 0.99 |
| Reward Scaling | 1.0 | 1.0 |
| Normalized Observations | False | True |
| Gradient Clipping | False | False |
| Exploration Policy | $\mathcal{N}(0, 0.1)$ | OU, $\theta = 0.15, \mu = 0, \sigma = 0.2$ |