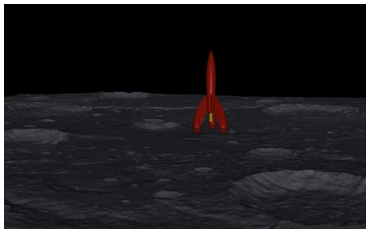




Agent

T
Action = $\begin{bmatrix} \theta_p \\ \theta_y \end{bmatrix}$

Rewards x_t
Observation = $\begin{bmatrix} \dot{x}_t \\ m \end{bmatrix}$
 x_{target}



Environment