

Weekly Review

13/04/21

- **Tasks**
- Updated rewards to include positive values ✓
- Understand intuition behind SAC ⚠ *In progress*
- Conduct experiments and evaluate results using GUI and plots ⚠ *In progress*
- **Problems**
- Results of experiment still not improved
- **To-Do Items for Next Week**
- Compare organization of current repo with RLPYT original repository
- Update states, values, rewards till you get
- Define more reward functions, states, actions for our use case
- **To-Do Later**
- Explore usage of intermediate testing on simulation before sim-to-real transfer
- Define use-case (for different type of towels (colour, texture, etc.) / one type)

RL Problem for obtaining one flat seam

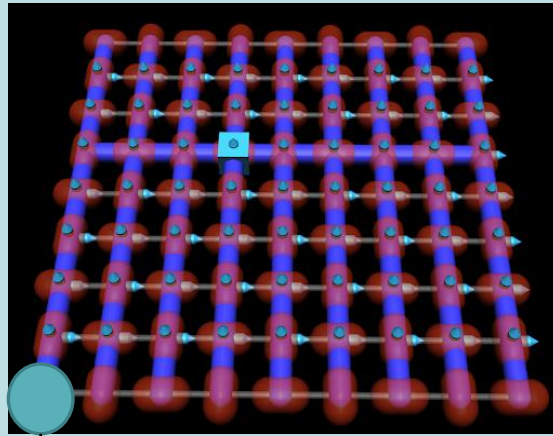
13/04/21

Goal

Obtain one flat seam

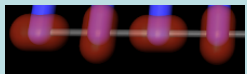
Given

Cloth in mujoco represented by 64 particles in 8×8 grid



Corner 1

Observations



$[x, y, z]$ positions of 4 points adjacent to corner 1

Actions

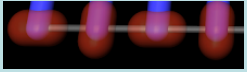
Random $[x, y]$ movement of corner 1

RL Problem for obtaining one flat seam

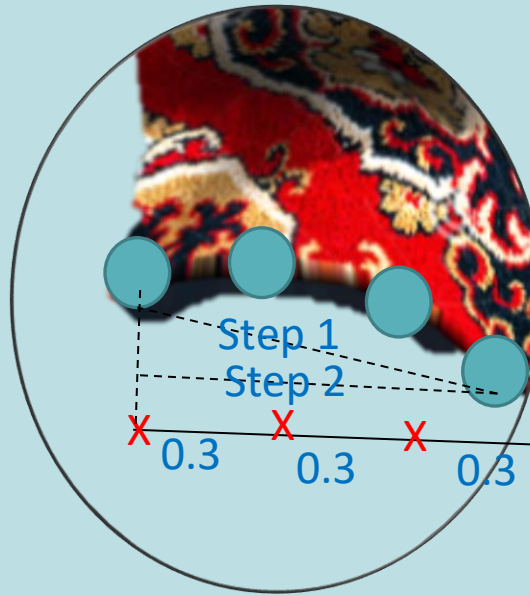
13/04/21

Goal

Obtain one flat seam -> Corner particle + 3 adjacent particles in a straight line



Reward

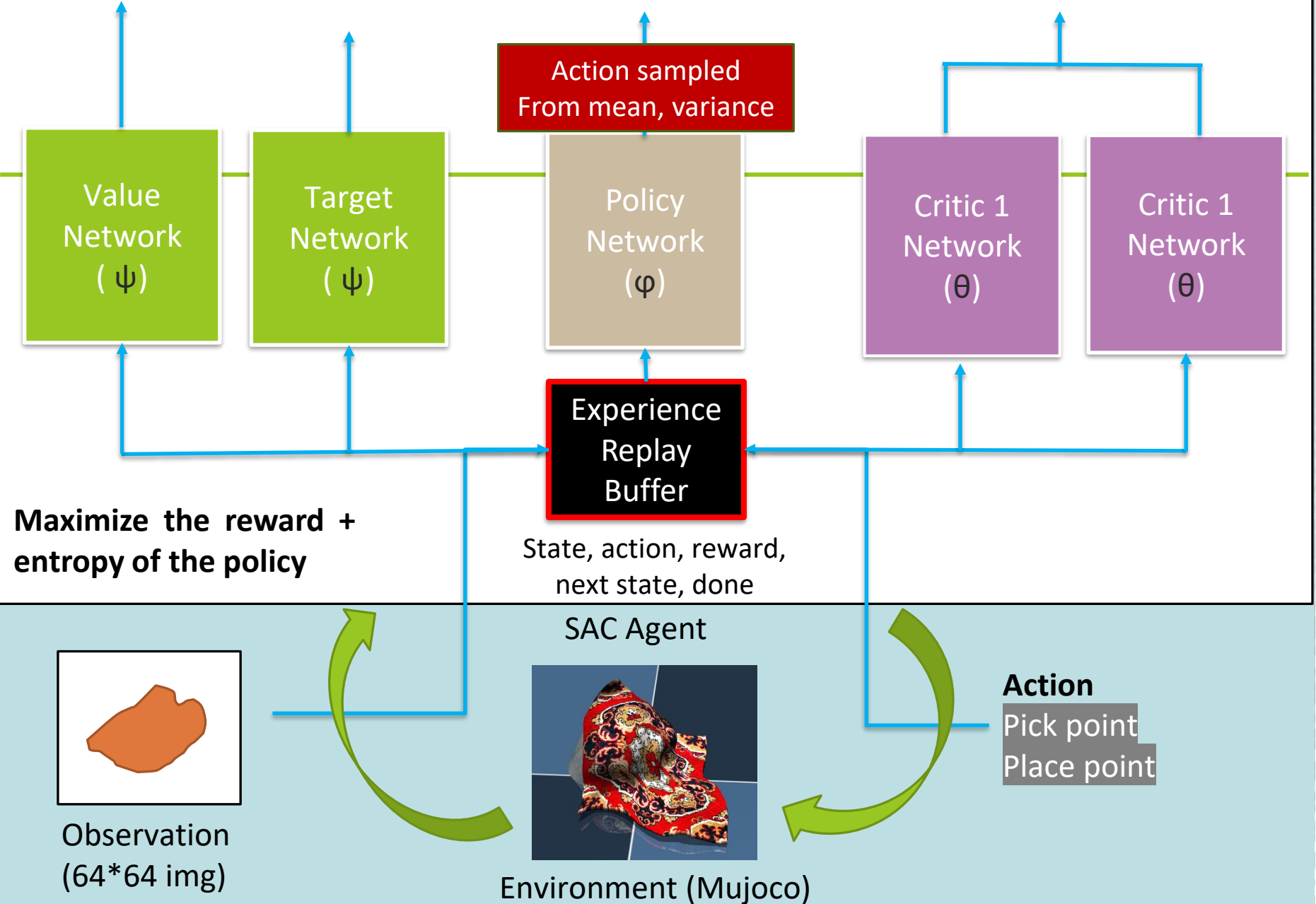


1. Join 1st point and last point
2. Project on x,y plane
3. Reward is proportional to :

*(1 - 1 * (x,y,z) distance from the ideal line) * 10*

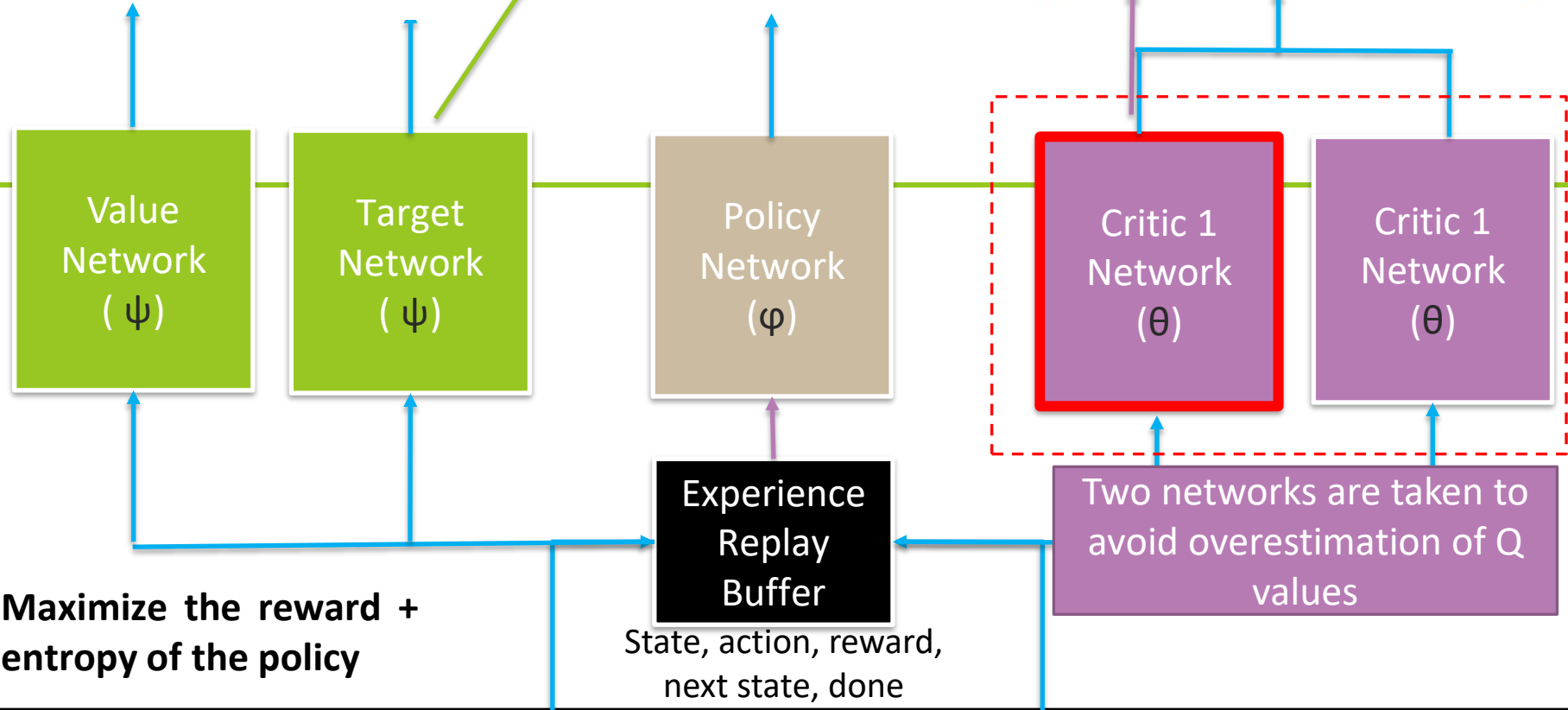
Updated reward

STEP 1 : Sample actions from value network and store state, action, reward, next state, done in Experience Replay Buffer



STEP 2 : We train the Quality Network by minimizing the following error

$$\hat{Q}(s_t, a_t) = r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim p} [V_{\tilde{\psi}}(s_{t+1})]$$
$$J_Q(\theta) = \mathbb{E}_{(s_t, a_t) \sim \mathcal{D}} \left[\frac{1}{2} \left(Q_{\theta}(s_t, a_t) - \hat{Q}(s_t, a_t) \right)^2 \right]$$



Maximize the reward + entropy of the policy

State, action, reward, next state, done

Two networks are taken to avoid overestimation of Q values



Observation (64*64 img)

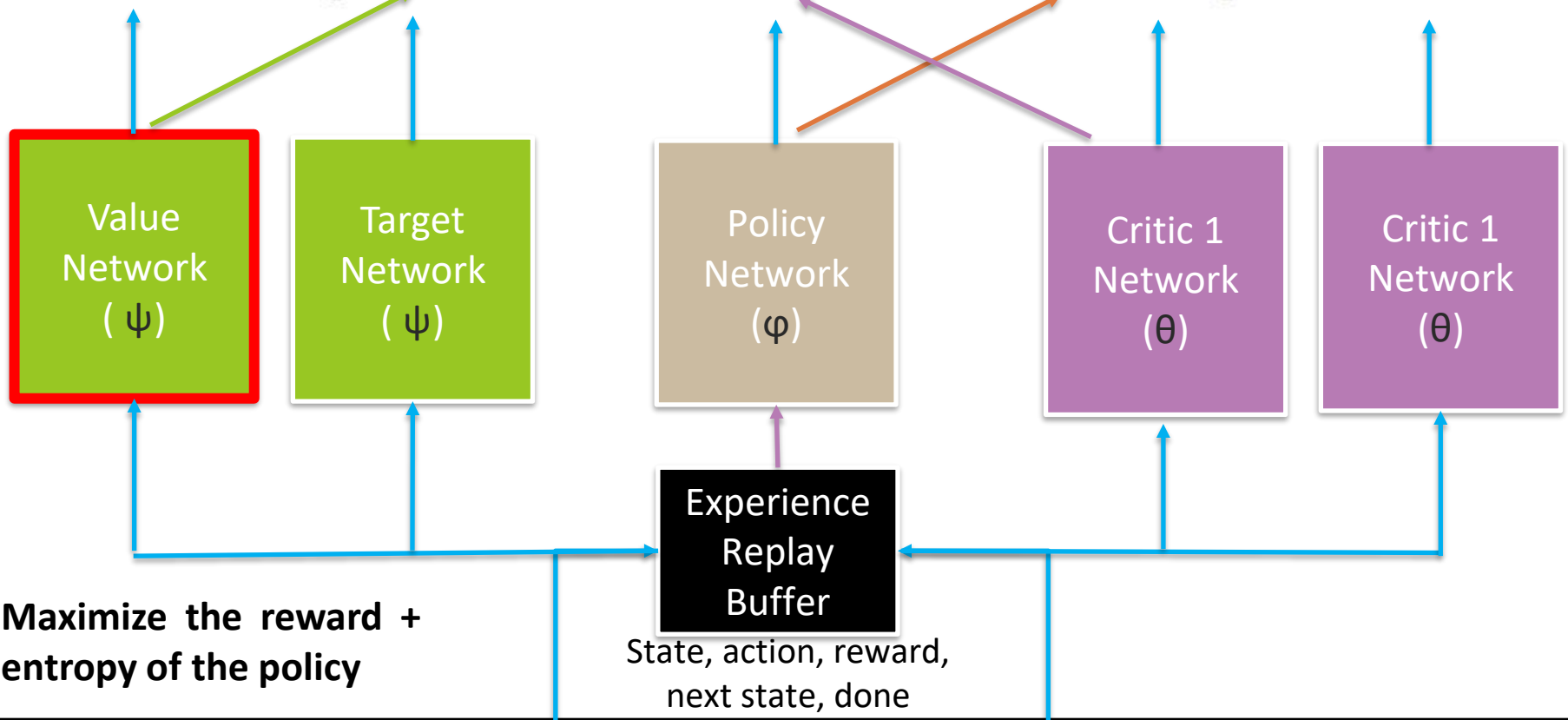


Environment (Mujoco)

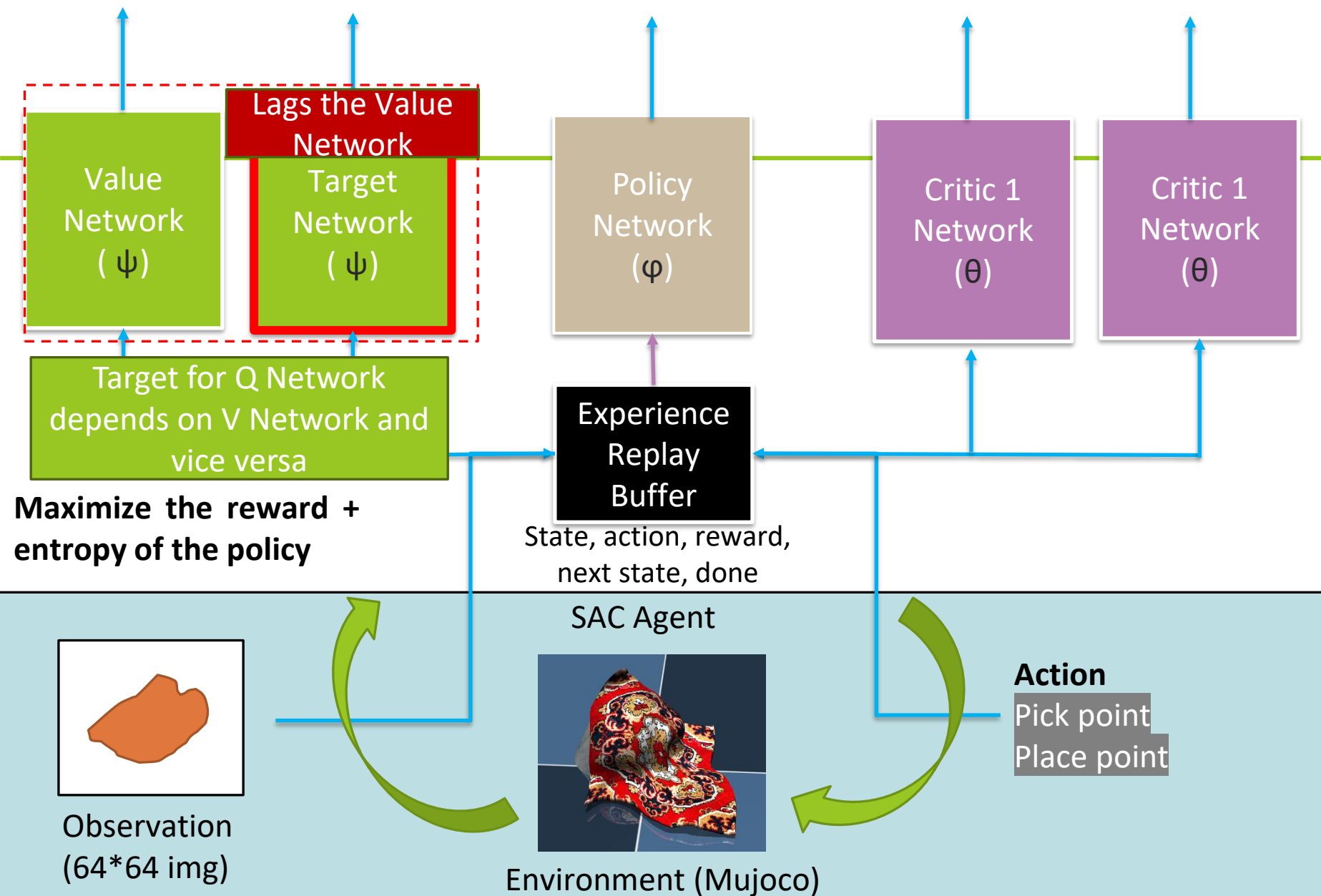
Action
Pick point
Place point

STEP 3 : We train the Value Network by minimizing the following

$$J_V(\psi) = \mathbb{E}_{s_t \sim \mathcal{D}} \left[\frac{1}{2} \left(V_{\psi}(s_t) - \mathbb{E}_{a_t \sim \pi_{\phi}} [Q_{\theta}(s_t, a_t) - \log \pi_{\phi}(a_t | s_t)] \right)^2 \right]$$



STEP 4 : We update the Target Network



Plan update

13/04/21

V3	13/04/2021	<p>Planned</p> <p>Phase 3 : Implementation : 52 days (mid Feb- early Apr)</p> <ul style="list-style-type: none"> a) Setting up the Reinforcement Learning Platform and Simulation environment : 13 days b) Prepare a custom implementation taking existing states, actions, rewards : 9 days c) Redefine actions and rewards for our use case : 15 days d) Test the pipeline and iterate : 15 days <p>Phase 4 (Additional): Sim-to-real transfer : 28 days (early Apr- early May)</p> <p><i>Master's Thesis Registration</i></p> <ul style="list-style-type: none"> d) Perform domain randomization : 7 days e) Transfer to real robot : 21 days <p>Phase 5 (Additional): Enhancement : 31 days (May)</p> <ul style="list-style-type: none"> a) Exploration of alternate sim-to-real transfer approaches and Exploration of acceleration strategies : 10 days b) Incorporation of a proven acceleration approach : 15 days 	Require more time to get valid results with SAC Algorithm
----	------------	--	---

Plan update

13/04/21

Update

Phase 3 : Implementation : 52 +22 days (mid Feb-end Apr)

- e) Setting up the Reinforcement Learning Platform and Simulation environment : 13 days
- f) Prepare a custom implementation taking existing states, actions, rewards : 9 days
- g) Redefine actions and rewards for our use case : 15 days
- h) Test the pipeline and iterate : 15 days
(+22 days for g) and h))

Phase 4 (Additional): Sim-to-real transfer : 31 days

(May)

Master's Thesis Registration

- f) Perform domain randomization : 7 days
- g) Transfer to real robot : 21 days

THANK YOU