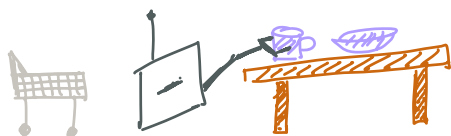
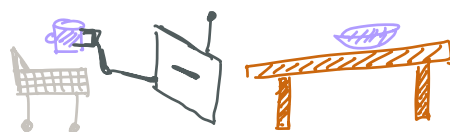


①



②



③



④



FORMULATE AS A MARKOV DECISION PROCESS (MDP)

$\langle S, A, R, T \rangle$



S:

A:

R:

$T(s' | s, a)$ .

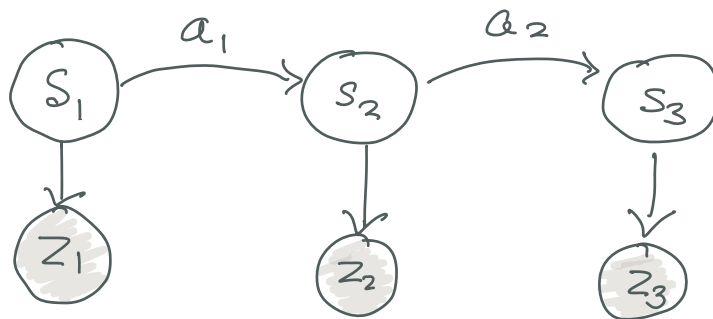
PERCEPTION

# How DO WE ESTIMATE STATE FROM OBSERVATION?



RGBD IMAGE  
(OBSERVATION  $z_t$ )

WE DON'T  
SEE STATE,  
WE SEE OBSERVATIONS



## STATE

- POSE OF OBJECTS
- CONFIGURATION OF ROBOT

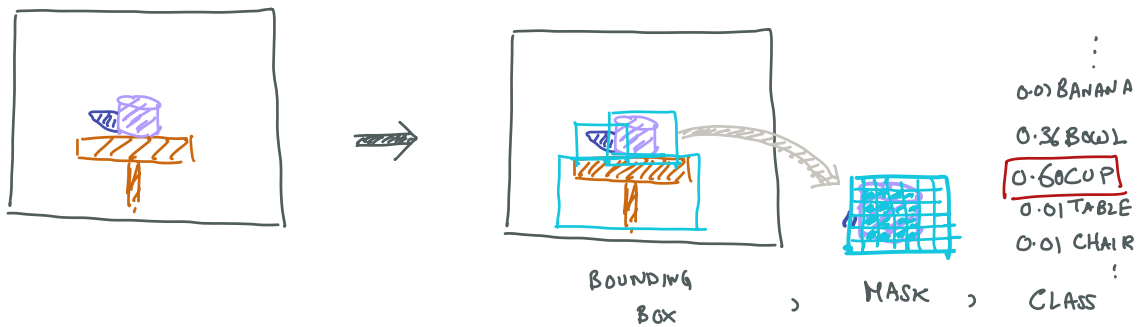
## OBSERVATION

- RGBD IMAGE
- ENCODERS ON JOINTS
- FORCE-TORQUE SENSOR
- IMU

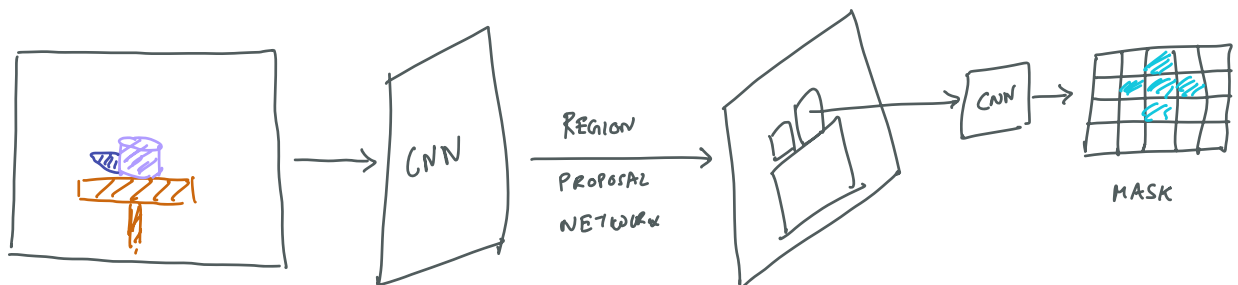
How DO WE ESTIMATE <sup>3D</sup> POSE OF OBJECTS  
FROM RGBD IMAGE ?

### STEP 1

TRAIN AN INSTANCE SEGMENTATOR ( MASK-RCNN )  
TO PROPOSE : ( BOUNDING BOX , MASK , OBJECT CLASS )

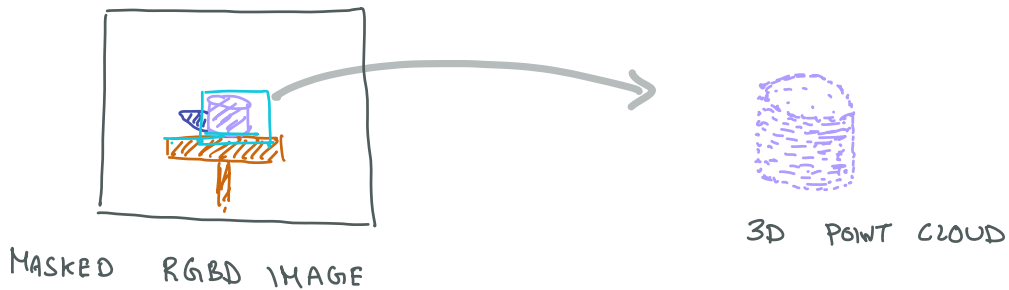


How?



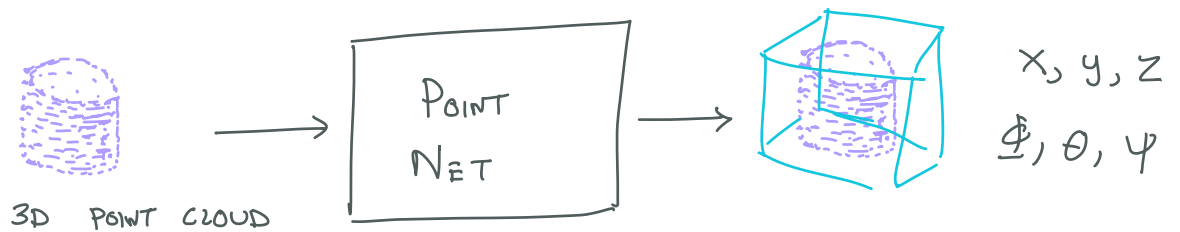
STEP 2

EXTRACT 3D POINT CLOUD FROM 2D IMAGE



STEP 3

TRAIN POINT NET TO PREDICT ACCURATE 6DOF POSE FROM POINT CLO



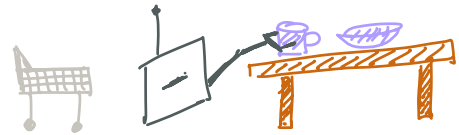
LET'S TAKE STOCK : WE KNOW STATE, WE KNOW ACTIONS

$\langle S, A, R, T \rangle$

GOAL:

LEARN A POLICY  $\pi$  THAT MAPS STATE TO ACTIONS.

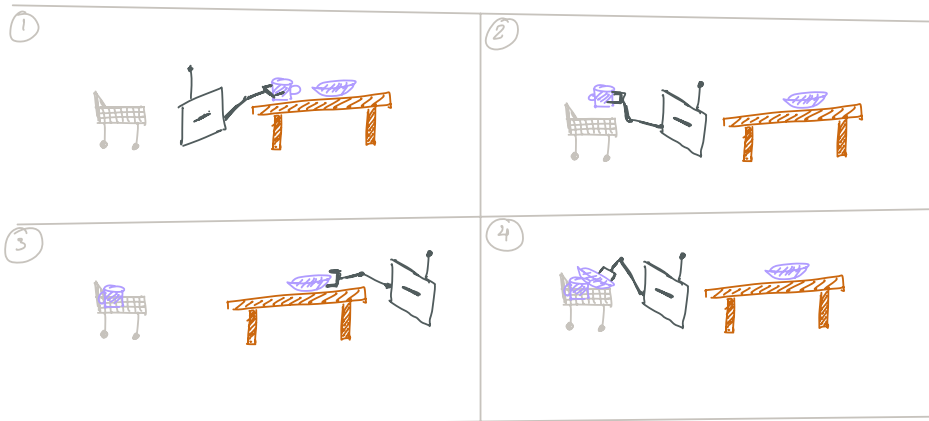
$\pi: S \rightarrow A$



WHAT IS THE SIMPLEST APPROACH ?

# BEHAVIOR CLONING

STEP 1: HUMAN PROVIDES DEMONSTRATIONS



$$D = \left\{ \begin{array}{l} (s_1^*, a_1^*, s_2^*, a_2^*, \dots, s_T^*, a_T^*) \\ (s_1^*, a_1^*, s_2^*, a_2^*, \dots, s_T^*, a_T^*) \\ (s_1^*, a_1^*, s_2^*, a_2^*, \dots, s_T^*, a_T^*) \end{array} \right\}$$

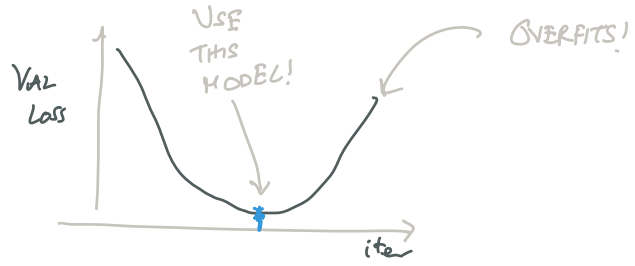
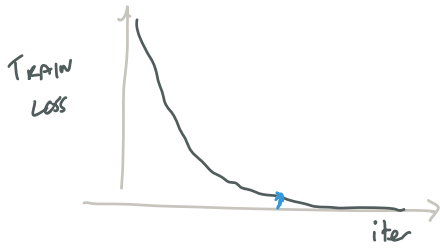
STEP 2: TRAIN A POLICY TO MAP STATES TO ACTIONS

$$\arg \min_{\pi} E_{s^*, a^* \sim \pi^*} \ell(a^*, \pi(s^*))$$

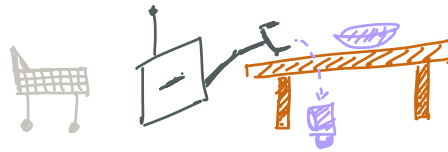
HOW DO WE IMPLEMENT THIS IN PRACTICE?

WHAT IS THE LOSS FOR  
DISCRETE ACTIONS  
CONTINUOUS ACTIONS:

# STEP 3: CHECK VALIDATION LOSS



How CAN BC FAIL?





# DAGGER

INSIGHT: ASK THE HUMAN EXPERT FOR "CORRECTIONS"  
ON STATES THE ROBOT VISITS.

$$\operatorname{argmin}_{\pi} E_{s \sim \pi} \ell(\pi^*(s), \pi(s))$$



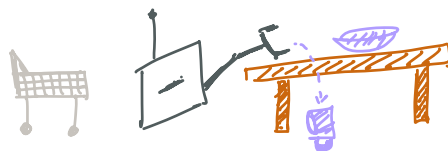
## DAGGER ALGORITHM

$\pi_0 \leftarrow$  INITIALIZE POLICY  
WITH BEHAVIOR CLONING.

$D \leftarrow \{\}$  INITIALIZE  
EMPTY DATA  
BUFFER

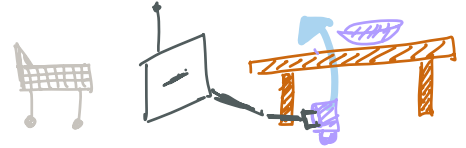
FOR  $i = 1 \dots N$

| ROLLOUT  $\pi_i$   
( $s_1, a_1, s_2, a_2, \dots$ )



QUERY HUMAN  $\pi^*$  FOR  
CORRECT ACTIONS

$(s_1, \pi^*(s_1), s_2, \pi^*(s_2), \dots)$



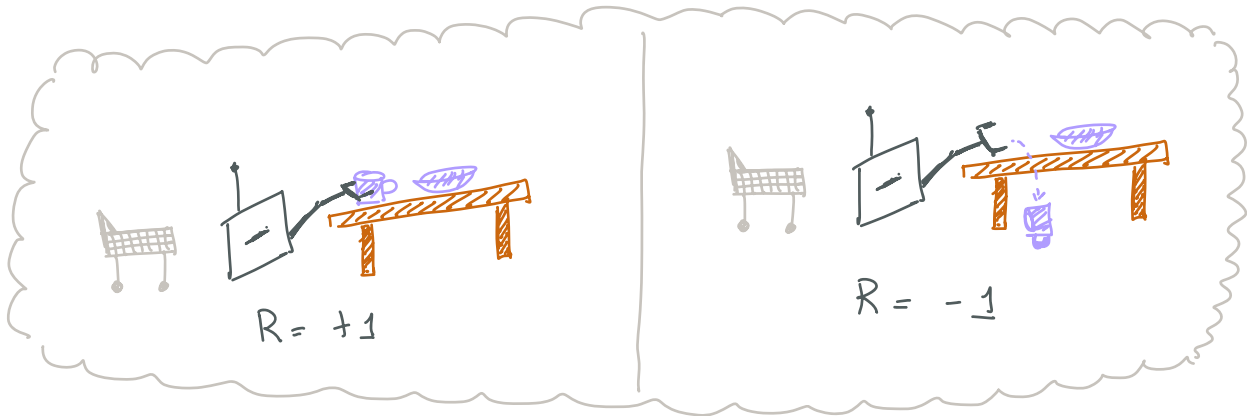
$D \leftarrow D \cup \{ (s_1, \pi^*(s_1), s_2, \pi^*(s_2), \dots) \}$

$\pi_i \leftarrow \text{TRAIN}(D)$

PRACTICAL ISSUES WITH DAGGER?

ASK HUMAN TO DESIGN REWARD FUNCTION.

$\langle S, A, R, T \rangle$

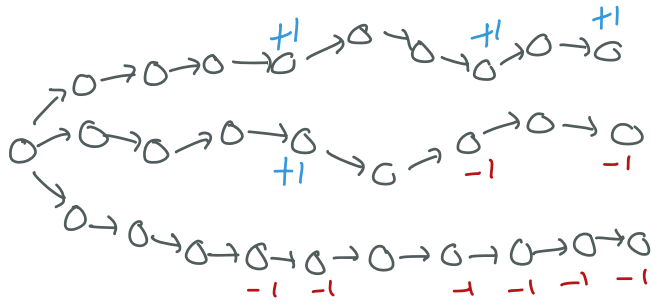


HOW CAN WE TRAIN A POLICY  $\pi$   
BY INTERACTING WITH WORLD  $T(s'|s, a)$ ?

## POLICY GRADIENT

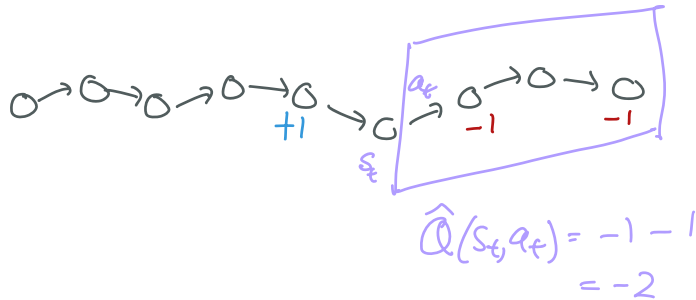
IDEA: ROLLOUT POLICIES, ESTIMATE VALUE, UPDATE.

$$\operatorname{argmax}_{\pi_{\theta}} E_{s_t \sim \pi_{\theta}} \left[ \sum_{t=1}^T r(s_t, a_t) \right]$$



## POLICY GRADIENT THEOREM

$$\theta' \leftarrow \theta + \eta \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \tilde{Q}(s_t, a_t) \right]$$



# REINFORCE

---

$\pi_0 \leftarrow$  INITIALIZE POLICY

FOR  $i = 1 \dots N$

ROLLOUT  $\pi_i$ :

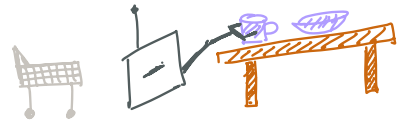
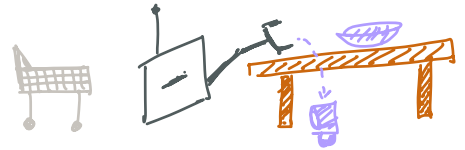
$(s_1, a_1, r_1, s_2, a_2, r_2, \dots)$

COMPUTE  $\hat{Q}$ : (REWARD-TO-GO)

$(s_1, a_1, \hat{Q}_1, s_2, a_2, \hat{Q}_2, \dots)$

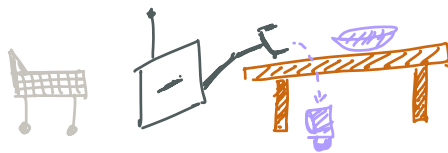
UPDATE POLICY  $\pi_\theta$ :

$$\theta' \leftarrow \theta + \eta \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \hat{Q}(s_t, a_t) \right]$$



PROBLEMS WITH REINFORCE ?

WHY IS MODEL-FREE RL CHALLENGING  
FOR ROBOTICS?



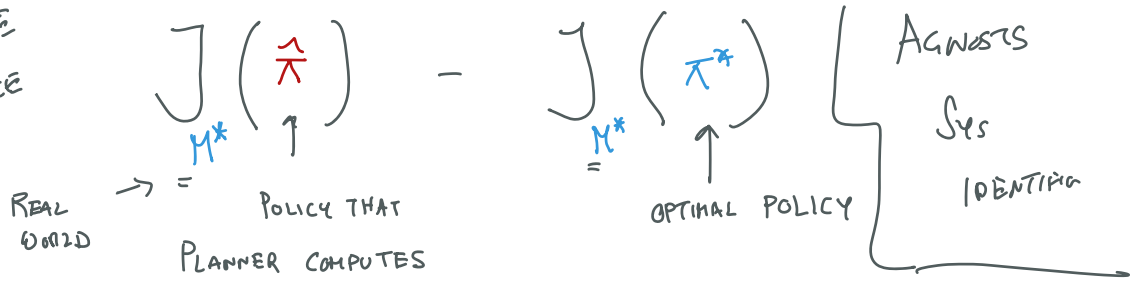
WHAT IF WE LEARN  $T(s'|s,a)$  ?

$\langle S, A, R, T \rangle$

WHAT HAPPENS IF WE PLAN WITH A MODEL  
LEARNED ONLY FROM HUMAN DEMONSTRATION?



PERFORMANCE  
DIFFERENCE



$$\leq \text{MODEL LOSS ON STATES LEARNER VISITS} + \text{MODEL LOSS ON STATES EXPERT VISITS} + \text{SUBOPTIMALITY OF YOUR PLANNER}$$



# DAGGER FOR MODEL BASED RL

GIVEN: DATA FROM AN EXPERT POLICY  $\pi^*$

$$D_{\text{exp}} = \left\{ (s_i, a_i, s'_i) \right\}_{i=1}^N$$

INIT:

INITIALIZE A MODEL  $\hat{M}_1(s, a) \rightarrow s'$

$$L(M) = \frac{1}{|D_{\text{exp}}|} \sum_{s, a, s' \in D_{\text{exp}}} \| \hat{M}(s, a) - s' \|^2 \quad \text{Optimize loss}$$

$D_{\text{learner}} = \{ \}$  # INITIALIZE LEARNER DATA  
 FOR  $i = 1 \dots N$  (# OUTER LOOP  
 #  $N$  IS ROUNDS OF  
 DATA COLLECT)

# CALL A PLANNER ON MY MODEL  $M_i$

$$\hat{\pi}_i = \arg \min_{\pi \in \Pi} J(\pi)_{M_i}$$

# COLLECT DATA BY RUNNING POLICY  $\hat{\pi}_i$  IN REAL WORLD.

$$D_{\text{learner}}^i = \left\{ s_i, a_i, s'_i \right\}_{i=1}^N$$

# AGGREGATE LEARNER DATA

$$D_{\text{learner}} = D_{\text{learner}} \cup D'_{\text{learner}}$$

# TRAIN MODEL.

$$L_1(\hat{M}) = \frac{1}{|D_{\text{train}}|} \sum_{\substack{(s,a,s') \\ \in D_{\text{train}}}} \|\hat{M}(s,a) - s'\|^2$$

$$L_2(\hat{M}) = \frac{1}{|D_{\text{exp}}|} \sum_{\substack{(s,a,s') \in \\ D_{\text{exp}}}} \|\hat{M}(s,a) - s'\|^2$$

# GET NEW MODEL

$$\hat{M}_{\text{fit}} = \text{OPTIMIZE} (L_1(\hat{M}) + L_2(\hat{M}))$$

POLICY

$$s \rightarrow \begin{matrix} \circ & \circ \\ \circ & \circ \end{matrix} \rightarrow a.$$

MODEL

$$\begin{matrix} s \\ a \end{matrix} \rightarrow \begin{matrix} \circ & \circ \\ \circ & \circ \end{matrix} \rightarrow s'.$$

D