

Figure 6: Domains used for our experiments in Section 4. Enumerating left to right from the top-left: first screen of Montezuma’s Revenge, MiniGrid FourRooms, Ant Medium Maze, Robosuite Door, Robosuite Lever, and Robosuite Slide.

517 A Task Descriptions

518 **Four Rooms.** This task, which is part of the MINIGRID suite [Chevalier-Boisvert et al., 2018], is
 519 an adaptation of the four rooms problem presented in the paper introducing the Options framework
 520 [Sutton et al., 1999]. Observations are 84×84 images; the underlying state-space contains approxi-
 521 mately 19×19 grid locations and 4 possible orientations of the player. We defined 4 options whose
 522 termination conditions (subgoals) were to navigate to the center of the 4 rooms. \mathcal{S}_0 was the set of all
 523 empty grid locations.

524 **Montezuma’s Revenge.** As is standard in ALE [Bellemare et al., 2013], observations are 84×84
 525 images, action space is a set of 18 discrete actions [Machado et al., 2018]. We defined start states
 526 scattered across the first room where the player was on the ground (not jumping or falling), was not
 527 too close to the skull and did not already have the key in its possession. We also defined 5 options
 528 whose termination conditions are reaching the left door, right door, bottom-right of the first screen,
 529 bottom-left of the first screen and getting the key.

530 **Robosuite manipulation tasks.** Three constrained manipulation tasks were used to study the task-
 531 oriented grasping performance of our initiation set learning algorithms: opening a door, flipping a
 532 lever, and sliding a knob. The door task was originally implemented in ROBOSUITE [Zhu et al., 2020];
 533 the others were implemented and described in the work of Rosen et al. [2022]. All three are 1-DoF
 534 articulated objects which require making and sustaining contact to manipulate. The 52-dimensional
 535 observation space consists of the robot’s proprioceptive state (joint position, joint velocity, end-
 536 effector position, gripper state, tactile data) as well as the object state (object pose, joint position,
 537 handle position). The action space employed is operational space control with variable impedance
 538 [Martín-Martín et al., 2019]: the agent controls the 6-DoF change in position and orientation of the
 539 end-effector, the 6-DoF change in stiffness, and 1-DoF gripper state. Episodes have a maximum
 540 length of 250 steps. In each task, \mathcal{S}_0 was a set of arm configurations establishing contact with the
 541 object; see Section B.2.

542 **Ant Medium Maze.** The goal location is small region around $(20, 20)$. A state is considered to
 543 satisfy a goal if the two have a euclidean distance of 0.5 units or less $R(s, g) = \|s - g\|_2 < 0.5$.
 544 The agent is evaluated by rolling out the learned policy once every 10 episodes; during evaluation,

545 the agent starts from a small region around $(0, 0)$, during training it starts at a location randomly
 546 sampled from the open locations in the maze. The task reward function is -1 for non-goal transitions
 547 and a terminating reward of 0 for reaching the goal. Episodes last a maximum of 1000 steps. The
 548 training and evaluation protocol is identical to Bagaria et al. [2021a], except for the fact that we learn
 549 initiation sets over the full state.

550 B Details about Learning Algorithm

551 B.1 Accuracy Experiment

552 Algorithm 1 is the pseudocode used for the experiments described in Section 4.1. Every episode,
 553 every option is executed from every start state in \mathcal{S}_0 . The result of that execution is recorded as
 554 ground-truth $Y_{s,o}(t)$ and stored to later compute the size of the true initiation set $|Y| = \frac{\sum_{s,o,t} Y_{s,o}(t)}{|Y_{s,o}(t)|}$.
 555 If the start was predicted to be inside the initiation set by the learning algorithm, then the trajectory
 556 generated by rolling out the option policy is used to update the policy and the initiation learner (e.g,
 557 IVF, classifier). We report the agreement between the predicted initiations and the ground-truth as an
 558 accuracy measurement for that state-option pair. The fraction of start states in \mathcal{S}_0 that lead to success
 559 is reported as a measurement of the size of the true initiation set.

Algorithm 1 Accuracy/Size Experiment Procedure

Inputs: Option termination conditions $\beta_o, \forall o \in \mathcal{O}$, start states \mathcal{S}_0 , number of episodes `n_episodes`.
Outputs: Accuracy table A and Initiation Size table S ; both map state-option pairs to a list of booleans.

```

Initialize goal-conditioned policy  $\pi_\theta : \mathcal{S} \times \mathcal{G} \rightarrow a \in \mathcal{A}$ .
Initialize Initiation Value Function (IVF)  $\mathcal{V}_\phi : \mathcal{S} \times \mathcal{G} \rightarrow \mathbb{R}$ .
Initialize binary classifier for each option  $\mathcal{I}_o(s; \psi) \rightarrow \{0, 1\}$ .
Initialize replay buffers for Rainbow  $B_R$  and IVF  $B_I$ .
Initialize buffers to store positive and negative examples for each option’s initiation classifier.
Initialize tables  $A$  and  $S$  as mapping each state-option pair to an empty list.
for episode  $\in \text{range}(\text{n\_episodes})$  do
  for start state  $s_0 \in \mathcal{S}_0$  do
    for option  $o \in \mathcal{O}$  do
      Reset simulator to  $s_0$ .
      Record option  $o$ ’s initiation decision  $X = \mathcal{I}_o(s_0; \psi)$ .
      Rollout option policy  $\pi_o(s_0, g \sim \beta_o)$  to get trajectory  $\tau$  and next state  $s'$ .
      Record whether the option policy reached the goal  $Y = s' \in \beta_o$ .
      Record accuracy  $A[s_0][o].\text{append}(\mathbb{1}(X = Y))$ .
      Update ground-truth size table  $S[s_0][o].\text{append}(Y)$ .
      if predicted initiation  $X = 1$  then
        Add trajectory  $\tau$  to policy’s replay  $B_R$ .
        Relabel trajectory  $\tau$  with initiation cumulant  $c_0 : \mathcal{S} \rightarrow \{0, 1\}$ .
        Add relabeled trajectory to IVF’s replay  $B_I$ .
        Add trajectory  $\tau$  to  $o$ ’s positive/negative example buffer.
      end if
    end for
  end for
  Sample minibatch and update  $\pi_o$  using Rainbow.
  Sample minibatch and update  $\mathcal{V}_\phi$  using TD(0).
  for option  $o \in \mathcal{O}$  do
    Compute weights  $w(s)$  for all training examples using Equation 1.
    Update  $o$ ’s initiation classifier by minimizing weighted cross-entropy loss.
  end for
end for

```

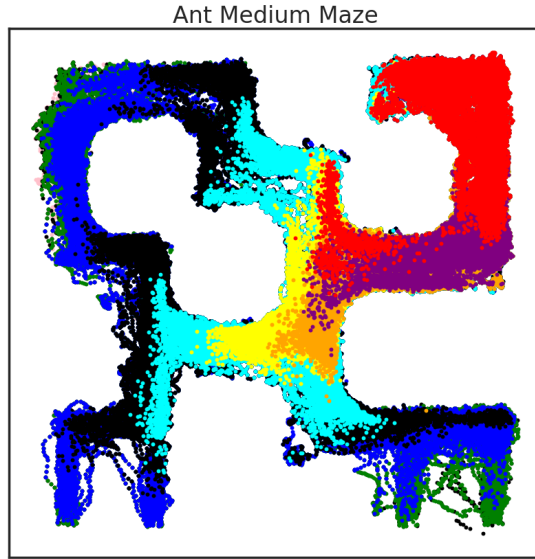


Figure 7: Initiation sets learned by deep skill chaining using the weighted classification approach in ANT MEDIUM MAZE. The task involves navigating the ant from the bottom-left to the top-right. Each color in the scatter plot denotes the initiation set of a different option; although the plot shows the location of the ant in the maze, the initiation set is learned using the full 30-dimensional state.

560 B.2 Robot Manipulation

561 The task-specific grasping problem is typically phrased as identifying grasp poses $g \in \mathbb{SE}(3)$ that
 562 afford task success. In practice, the difficulty of this problem is compounded by the fact that,
 563 for redundant manipulators, each grasp pose g yields an infinite number of corresponding arm
 564 configurations (solutions to the inverse kinematics problem). Explicitly, this relation is governed by
 565 the manipulator’s forward kinematics $f : \mathcal{C} \rightarrow \mathbb{SE}(3)$ which maps (typically) 7-DoF configurations of
 566 the arm $q \in \mathcal{C}$ to poses in Cartesian space. In practice, only a subset of these configurations for a given
 567 grasp pose enable successful manipulation [Schiavi et al., 2022]. As a result, we task the initiation
 568 set learning algorithm with choosing start states directly from the space of arm configurations \mathcal{C} .

569 We generate collision-free grasp poses on each object using off-the-shelf grasp generation method
 570 GPG [Ten Pas et al., 2017] and corresponding arm poses using IKFLOW [Ames et al., 2022]. We
 571 chose to generate 50 grasp poses with 5 random inverse kinematics solutions each yielding a total of
 572 250 starting configurations for each task.

573 **Reward function.** The reward functions are implemented as progress toward 1-DoF object joint
 574 position goals. The agent receives reward when the current joint position exceeds its previous
 575 maximum in a given episode.

576 **Parameterization.** As described in Section A, the agent receives proprioceptive and object-state
 577 observations and controls the manipulator’s end-effector pose and impedance. The learning algorithm
 578 employed is TD3 [Fujimoto et al., 2018]. Goal-conditioning is omitted in these experiments as they
 579 have a single goal and a single option.

580 B.3 Deep Skill Chaining

581 Deep skill chaining (DSC) [Konidaris and Barto, 2009, Bagaria and Konidaris, 2020] proceeds
 582 recursively backward from the goal: the termination region of the first option is the task goal (Line 2,
 583 Algorithm 3); the agent first learns an option policy and initiation set for that option. Then, it learns

584 another option whose termination region is the initiation set of the first option; this process continues
 585 until there is some option whose initiation set covers the start state. The result is a collection of
 586 options that can be sequenced to reliably reach the goal.

587 Since the initiation sets of options are the subgoals for other options, the entire process is sensitive to
 588 the way in which the initiation sets are learned: poorly estimated initiation sets can lead to subgoals
 589 that do not improve the agent’s ability to reach the task goal.

590 Details about line 12 of Algorithm 2 differ based on the method used to learn the initiation set. When
 591 using the pure GVF approach, we perform as many minibatch gradient updates as the length of the
 592 option rollout; when using weighted classification, we recompute weights using Eq 1 for all training
 593 examples and then proceed to minimize weighted cross-entropy loss (3 epochs, batch size 128).
 594 When using classification (weighted or unweighted), we boost the contribution of the minority class
 595 by the ratio of the size of the majority class to that of the minority class.

Algorithm 2 Robust DSC Rollout

Inputs. Skill Chain \mathcal{O}

Hyperparameters. Option horizon H

- 1: Initialize empty trajectory buffer \mathcal{B}
 - 2: **for** each timestep t **do**
 - 3: Select option o using policy over options $\pi_{\mathcal{O}}(s_t)$
 - 4: Sample a goal for selected option: $g \sim \beta_o$
 - 5: Execute option policy $\pi_o(\cdot|g)$ in the environment
 - 6: Add trajectory $\tau = \bigcup_{i=0}^{H-1} (s_i, o, a_i, s_{i+1}, g)$ to \mathcal{B}
 - 7: **if** final state s_H reached goal g **then**
 - 8: Add τ to o ’s list of positive examples
 - 9: **else**
 - 10: Add τ to o ’s list of negative examples
 - 11: **end if**
 - 12: Refit option o ’s initiation classifier
 - 13: Add τ to replay buffer and update π_o using TD3
 - 14: **end for**
 - 15: **return** $\mathcal{B} = \bigcup_t (s_t, o_t, a_t, s_{t+1}, g_t)$
-

Algorithm 3 Robust DSC Algorithm

Inputs. Start state s_0 , Goal region g .

- 1: Initialize global option o_G such that $\mathcal{I}_{o_G}(\cdot) = 1$
 - 2: Initialize goal option o_g such that $\beta_{o_g} = g$
 - 3: Initialize skill chain \mathcal{O} with $\{o_g\}$
 - 4: **for** each episode **do**
 - 5: transitions = ROLLOUT(\mathcal{O})
 - 6: **if** $s_0 \notin \mathcal{I}_o, \forall o \in \mathcal{O}$ **then**
 - 7: Create new option ω
 - 8: Add ω to skill chain \mathcal{O}
 - 9: **end if**
 - 10: **end for**
-

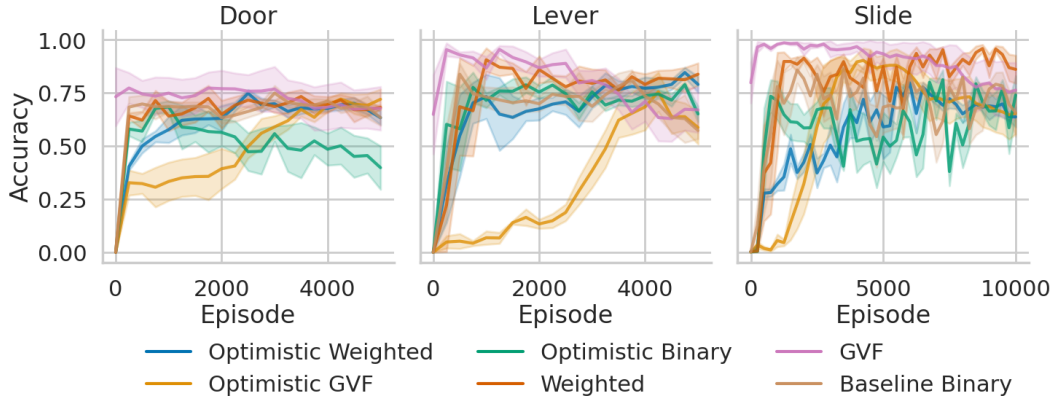
596 **Picking a goal for option execution.** Line 4 of Algorithm 2 samples a goal from the option’s
 597 termination region. To implement this sampling procedure, we consider the option’s parent ω in
 598 the chain (the parent option is the one whose initiation set is being targeted by the current option o).
 599 We enumerate the positive examples used to train \mathcal{I}_ω and pick the goal with the highest initiation
 600 probability under the current option. This process is done iteratively backward from the goal: the first
 601 goal is the task goal, the next one is the positive example closest (in terms of highest IVF value) to
 602 the task goal and so on.

603 **Learned initiation sets.** Figure 7 shows the initiation sets learned by DSC when using the weighted
 604 classification approach. Starting from the bottom-left of the maze, the agent successively targets the

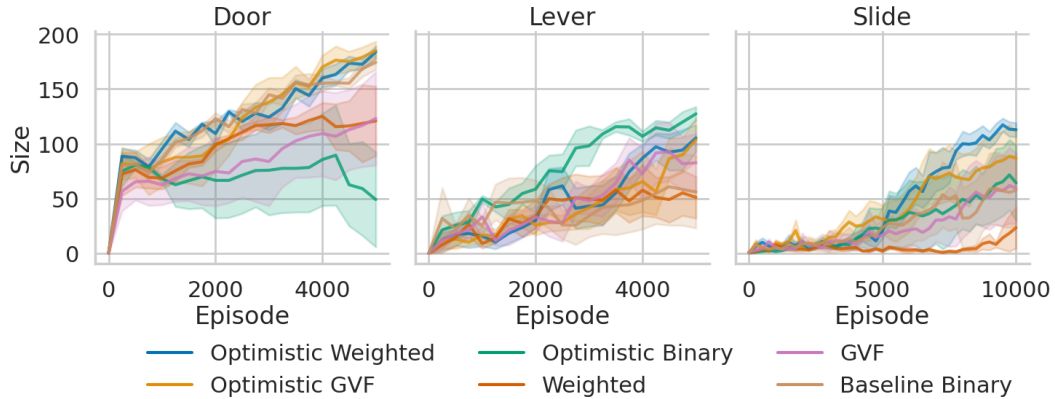
605 next option’s initiation set, until it reaches the task goal at the top-right of the maze. The plot was
 606 generated by querying the learned initiation classifiers on the states in the agent’s replay buffer at the
 607 end of training; only the (x, y) coordinates of those states are visualized.

608 C Additional Manipulation Experiments

609 Initiation set accuracy and true size are computed during training by performing an analogous
 610 procedure to Algorithm 1. Periodically, the manipulator was reset to each candidate start state and
 611 the initiation prediction was compared with the outcome of a policy rollout. Initiation set accuracy is
 612 visualized in Figure 8a. The methods generally converge to similar accuracy. True initiation set size
 613 is plotted in Figure 8b; size increases with optimism and correlates with success rates.



(a) Initiation set accuracy for manipulation domains.



(b) Initiation set size for manipulation domains (out of 250 start states).

Figure 8: (a) Accuracy of the learned initiation sets in the robot manipulation domains. (b) The size of the “true” initiation sets measured by performing Monte Carlo rollouts of the option policy.

614 D Hyperparameters

615 Rainbow was used for policy learning in Section 4.1, TD3 was used in the other experiments. Their
 616 hyperparameters (Tables 1 and 4) were not tuned and are either identical to the original paper
 617 implementation or borrowed from Bagaria et al. [2021a]. The bonus scale c (described in Sec 3.3)
 618 was tuned over the set $\{0.05, 0.1, 0.25, 0.5, 1.0\}$, the best performing hyperparameters are listed in
 619 Table 2.

Parameter	Value
Replay buffer size	10^6
Critic Learning rate	$3 \cdot 10^{-4}$
Actor Learning rate	$3 \cdot 10^{-4}$
Optimizer	Adam
Target Update Rate	$5 \cdot 10^{-3}$
Batch size	256
Iterations per time step	1
Discount Factor	0.99
Output Normalization	False

Table 1: TD3 Hyperparameters for Robosuite and DSC Experiments

Method	Bonus Scale
Optimistic Binary	0.1
Optimistic GVF	0.5
Optimistic Weighted	0.5

Table 2: Exploration Hyperparameters for Robosuite Experiments.

Parameter	Value
Replay buffer size	$3 \cdot 10^5$
Replay start size	1024
Learning rate	10^{-4}

Table 3: Rainbow Hyperparameters for Accuracy Experiments

Parameter	Value
Learning rate	10^{-4}
Optimizer	Adam
Replay buffer size	10^5
Batch size	32
Threshold	0.5
Target network update rate	$5 \cdot 10^{-3}$

Table 4: IVF Hyperparameters