

## 520 A Code and Further Resources

521 We provide complementary code in the supplementary of the paper. All models were imple-  
 522 mented in Python using the PyTorch library. The total training computation required for op-  
 523 timizing all models and baselines presented in this works amounts to roughly 150 hours (wall  
 524 clock) on a single standard GPU. Visualization and additional material will be uploaded here  
 525 <https://sites.google.com/view/velap-cori/home>.

## 526 B Hyperparameters and algorithm details

527 Here we present a description of the hyperparameter used in method used during the training and  
 528 inference phases.

### 529 B.1 Model architectures

Table 2: Hyperparameters of the encoder  $\phi$

Parameter	Value
Batch-norm.	yes
Filters	[32,32,64,64]
Kernels	[4,4,4,4]
Strides	[2,2,2,2]
Activation	LeakyRelu
Dense layers	[256, 128, 32]

Table 3: Hyperparameters of dynamics model  $h$

Parameter	Value
Batch-norm.	yes
Activation	LeakyRelu
Dense layers	[128,128,128,128,32]

Table 4: Hyperparameters of action sampler  $g$  ( $\beta$ -VAE)

Parameter	Value
Batch-norm.	yes
Activation	LeakyRelu
Latent dimension	16
$\beta$ (kl-weight)	0.01
Encoder dense layers	[128,128,128, 2*16]
Decoder dense layers	[128,128,128, $d_{\text{action}}$ ,]

Table 5: Hyperparameters of policy networks  $\pi^l$  and  $\pi^g$

Parameter	Value
Batch-norm.	yes
Activation	LeakyRelu
Dense layers	[128,128,128, $d_{\text{action}}$ ]

Table 6: Hyperparameters of critic networks  $Q_k^l$  and  $Q^g$ 

Parameter	Value
Batch-norm.	yes
Activation	LeakyRelu
Dense layers	[128,128,128, 1]

## 530 B.2 Training hyperparameters

Table 7: Hyperparameters during training (model learning)

Parameter	Value
batch size	64
learning rate (all models)	0.0003
$c_0$	0.2
$c_1$ ( <i>SpiralMaze</i> )	0.001
$c_1$ ( <i>ObstacleMaze</i> )	1.0
$c_1$ (metaworld tasks)	0.01
$c_2$	0.001
$c_3$	0.001
$c_3$ (expert)	0.5
$\gamma$ (discount factor)	0.96
$d_{\mathcal{Z}}$	32
$T$ (temperature parameter)	1.0
$n_{\text{ens}}$	3

## 531 B.3 Planner and controller hyperparameters

Table 8: Hyperparameters during inference (planning)

Name	Description	Value
$n_{\text{iter}}$	Number of planner iterations	250 (500 in <i>ButtonWall</i> )
$n_{\text{sim}}$	Number of simulation steps during tree expansion	5 (10 in <i>SpiralMaze</i> , <i>ButtonWall</i> and <i>DrawerButton</i> )
$\tau_{\text{discard}}^{\text{high}}$	Q-value threshold for discarding node if too close to existing nodes in the tree	$\gamma^2$
$\tau_{\text{discard}}^{\text{low}}$	Q-value threshold for discarding node if too far from expansion node	$\gamma^{n_{\text{sim}}}$
$\tau_{\text{discard}}^{\text{std}}$	Q-value threshold for discarding node if standard deviation of ensemble prediction is too high	$1.0 - \gamma$
$\tau_{\text{neigh}}$	Q-value threshold to determine neighboring nodes	$\gamma^3$
$\tau_{\text{goal}}$	Q-value threshold to determine goal nodes	$\gamma^1$
$d_{\text{neigh}}$	Euclidean distance threshold to determine candidate neighbors	3 x upper 5-percentile of Eucl. distances between encoding of subsequent states

Table 9: Hyperparameters during inference (controller)

Parameter	Description	Value
$n_{\text{replan}}$	Planning frequency	15 (25 in <i>SpiralMaze</i> , <i>ButtonWall</i> )
$\epsilon_{\text{goal}}$	Q-threshold to determine vicinity to the goal	$\gamma^5$
$\epsilon_{\text{wp}}$	Q-threshold for switching to the next waypoint	$\gamma^3$

## 532 B.4 Training of policy and value functions

533 We use TD3-BC [48] as the base offline RL algorithm to train our local and goal policies  $\pi^l$  and  $\pi^g$ ,  
 534 respectively state-action value functions  $Q_k^l$  and  $Q^g$ . Within our planning framework  $Q_k^l$  takes an  
 535 important role as it provides us with a distance proxy. To further improve the accuracy of  $Q_k^l$ , we  
 536 use  $k$  Q networks (instead of 2 usually used in TD3). During the training update of the Q-network,  
 537 we then determine the Q-target by taking the minimum value among the predictions given by the  
 538 ensemble of Q-networks (similar to [49]). The ensemble further allows us to filter out unlikely  
 539 or out-of-distribution transitions generated during the tree expansion by thresholding the resulting  
 540 Q-values based on the minimum predicted ensemble value and the standard deviation among the  
 541 predicted values.

542 Our models  $\pi^l$  and  $Q_k^l$  describe goal-reaching navigation policy and state-action value functions  
 543 which require a set of goal-conditioned reaching experiences for training. Since our original dataset  
 544  $\mathcal{D}$  might not provide such data, we augment it using data augmentation via hindsight relabeling. In  
 545 particular, we create a new dataset  $\mathcal{D}'$  creating transitions  $(z_t, a_t, r_t, z_{t+1}, z^g, \gamma) \in \mathcal{D}'$  based on the  
 546 existing transitions in  $\mathcal{D}$  by relabeling the values of  $r_t, \gamma$  ( $\gamma$  also indicates terminal condition, i.e.  
 547  $\gamma = 0$ ) and adding a new goal state  $z^g$ . In this regards, we apply a combination of three different  
 548 relabeling strategies (a) set goal  $z^g$  to be next state of the original transitions and set  $\gamma = 0$  and  
 549  $r_t = 1$  (b) sample  $z^g$  from the set of future states within the same trajectory and set  $\gamma = 0$ ;  $r_t = 0$   
 550 (c) sample  $z^g$  from another trajectory in the data and  $\gamma = 0$  and  $r_t = 0$ .

## 551 B.5 Additional details about planning method

552 **Neighbor computation** To determine if a newly sampled node  $z_{\text{new}}$  is novel, we check its simi-  
 553 larity to existing nodes in the tree by evaluating the state-action value function. Yet, evaluating the  
 554 values network for all nodes in the tree results in an enormous computations overhead. Yet, we can  
 555 significantly reduce this computation by first determining a set of candidate neighbors around  $z_{\text{new}}$   
 556 using the Euclidean metric and a distance threshold  $d_{\text{neigh}}$ . In practice, we found it useful to define  
 557  $d_{\text{neigh}}$  based on the statistics of Euclidean distances between subsequent states in the dataset (see  
 558 App. B.3).

559 **Batch processing** The method in Alg. 1 describes an iterative schema for which at every expan-  
 560 sion step one new node is generated and evaluated. Yet, some steps can be computed in parallel on  
 561 a GPU in order to speed up the planning time. For a practical implementation, we therefore suggest  
 562 to parallelize the tree expansion by sampling multiple expansion nodes at once and generating new  
 563 nodes by passing batches through the neural network dynamics model. Similarly, we can compute  
 564 state-action values in batches instead of assessing one new nodes at a time. For discussions about  
 565 highly-parallelized implementations of classical RRT-like planners, we refer to [58, 59].

## 566 B.6 Additional details about MPC controller

567 The below Alg. 2 presents the pseudocode for our MPC controller. The function  
 568 `update_waypoint( $z_{\text{curr}}, g^*$ )` determines the next waypoint  $z_{\text{wp}}$  which we seeks to achieve using  
 569 our local policy. In particular, we estimate the value between the current state and waypoint and  
 570 switch to the next element in  $g^*$  if the predicted value surpasses a threshold  $\epsilon_{\text{wp}}$ , i.e  $Q_{\text{min}}^{\text{curr}, \text{wp}} > \epsilon_{\text{wp}}$ .  
 571 Once the  $z_{\text{curr}}$  gets close to the goal, we disable planning and determine actions based on  $\pi^g$ . To  
 572 determine vicinity to the goal, we check if the prediction of the global value function  $Q^g$  exceeds a  
 573 threshold  $\epsilon_{\text{goal}}$ .

---

**Algorithm 2** MPC controller

---

```
Given:  $s_{\text{init}}, n_{\text{replan}}, n_{\text{max steps}}, \phi, \pi^l$ 
 $z_{\text{curr}} \leftarrow \phi(s_{\text{init}})$  ▷ Map state to latent encoding
 $i \leftarrow 0$ 
while goal not achieved and  $i < n_{\text{max steps}}$  do
  if not  $i \bmod n_{\text{replan}}$  then ▷ Replan every  $n_{\text{replan}}$  steps
    Build tree  $\mathcal{T}$  rooted at  $z_{\text{curr}}$  for  $n_{\text{iter}}$  steps (Alg. 1).
    Determine  $g^* = \{z_{\text{curr}}, z_1, \dots, z_n\}$  given  $\mathcal{T}$  (Eq. 8)
  end if
   $z_{\text{wp}} \leftarrow \text{update\_waypoint}(z_{\text{curr}}, g^*)$  ▷ Update waypoint if close enough
   $a \leftarrow \pi^l(z_{\text{curr}}, z_{\text{wp}})$  ▷ Compute action given policy
  Execute  $a$  and update state
   $z_{\text{curr}} \leftarrow \phi(s_{\text{curr}})$ 
end while
```

---

## 574 C Evaluation Environments

### 575 C.1 Description of block environments

576 Similar to the evaluation environments in [14], we implement two long-horizon navigation tasks  
577 whose underlying state space is relatively low-dimensional in order to facilitate illustration and  
578 visual inspect of learned embeddings using dimensionality reduction techniques such as Isomap  
579 [60]. For both environments, the a block robot is controlled using velocity commands while it's  
580 position in constrained to a planar surface.

#### 581 C.1.1 SpiralMaze

582 To solve this task, the block agent must navigate form the outer end of the spiral-shaped corridor  
583 to the inner region ( colored in red; see Fig. 2). The maximum allowed number of episode steps is  
584 limited to 300. To generate training data, the agent is placed randomly at a collision free position  
585 in the workspace and random actions sequences are applied by subsequently adding Gaussian noise  
586 an initially sampled random uniform action. For testing, the agent's position is sampled uniformly  
587 within a small region close to the outer end of the spiral-shaped corridor.

#### 588 C.1.2 ObstacleMaze

589 In this environment, the agent must navigate towards the upper wall of the workspace (color in  
590 red; see Fig. 2). To achieve the goal the agent must take actions around two obstacles which are  
591 randomly placed around the center of the workspace at the beginning of each new episode. The  
592 maximum allows number of environment steps is set to 100. For testing, the agent is initialized to  
593 random configuration close to wall which is opposite to the goal. We used the same random data  
594 collection policy as for the *SpiralMaze* task.

### 595 C.2 Description of manipulation environments

596 We adapted and implemented several robot manipulation environments based on the Metaworld [17]  
597 robot benchmark tasks. The underlying physics simulator in this regard is Mujoco [61]. To enable  
598 visual manipulation, similar to the problems studied in [6], we enable background rendering of RGB  
599 images from a static viewpoint. The robot is controlled by commanding desired endeffector and  
600 gripper opening displacements resulting in a 4-dimensional action space. While *WindowClose* and  
601 *FaucetClose* were with small modifications adapted from the [17], we evaluate two new scenarios  
602 *ButtonWall* and *DrawerButton* which were particularly desired to study our method under extreme  
603 sparse reward conditions over a long temporal horizon which requires trajectory "stitching" to find  
604 a solution policy.

605 For data collection, we use a suboptimal policies which takes random actions (additive Gaussian  
606 noise) most of the time and with a low probability takes an action generated by a scripted expert

policy. Table. 10 provides insight about the number of samples and trajectories in the training data and as well presents the portion of successful actions (reward=1). For all manipulation tasks, we set the maximum permitted environmental steps at 150, with the exception of the "ButtonWall" scenario, where we allow up to 250 steps during the evaluation phase.

### C.2.1 WindowClose

In order to accomplish this task, the robotic arm must successfully open a window by shifting a specific handle sideways. We implement environmental variability by randomly determining the x-y location of the window object in each episode. During the data collection stage, we randomly set the positioning of the end-effector above the surface of the table. However, we restrict the sampling of expert actions to areas close to the objective (window handle). This approach is intended to guarantee that the strategy employed necessitates to "stitch" different trajectories together to reach the objective and complete the task when starting from states that are farther away. To ensure challenging planning situations during testing, we initiate the robot at a significant distance away from the target.

### C.2.2 FaucetClose

This task is similar to the *WindowClose* task, but it requires the agent to use its end-effector to close a faucet instead. In addition, we employ analogous strategies for data gathering and scenario creation as those used in the *WindowClose* environment.

### C.2.3 ButtonWall

In this scenario, the robot's end-effector is required to navigate around a wall structure before pressing a button. The wall's location is randomly set at the beginning of each episode. Furthermore, a height limitation is imposed on the end-effector to ensure that the agent takes a more extended path around the wall, as opposed to simply elevating the end-effector. The dataset was produced by either placing the agent in front of the wall, near the button, or far behind the wall. However, expert samples in the dataset only exist for scenarios when starting closer to the goal. For testing purposes, the end-effector is sampled within a restricted region behind the wall to increase the planning task's complexity.

### C.2.4 DrawerButton

In this scenario, the agent is tasked to first close a drawer using its end-effector and subsequently press a button. To train the agent, we develop a dataset by separately collecting trajectories for each subtask. This approach necessitates the use of a method capable of combining different trajectories in the data to devise a solution that achieves the overall task goal.

## C.3 Composition of training dataset

The table below presents the composition of our training datasets. Each context in this regards, refers to a new environment initialization (excl. agent) such as the position of obstacles.

Table 10: Composition of training datasets for each environment

Environment	Num. contexts	Traj. per context	Max. traj. length	Successful transitions
<i>SpiralMaze</i>	1	1000	20	0.12 %
<i>ObstacleMaze</i>	250	20	20	0.11 %
<i>WindowClose</i>	200	10	50	0.48 %
<i>ButtonWall</i>	200	10	50	0.16 %
<i>FaucetClose</i>	200	10	50	0.31 %
<i>DrawerButton</i>	150	20	50	0.16 %

## 642 D Baseline methods

643 To enable a fair comparison between different methods, we use the same underlying representa-  
644 tion/encoder  $\phi$  and dynamics model  $h$  in the evaluation of all baselines. For assessing the quality  
645 and impact of our representation learner, please review the experimental ablation study in App. E.2.

### 646 D.1 BC and BC ( $\mathcal{D}^*$ )

647 Simple behavioral cloning baselines for which we use the same network architecture as our policy  
648 networks (see Table 5) and train using a mean-squared error objective. For  $\mathcal{D}^*$  we train only on the  
649 subset of successful episodes in the dataset. The train both methods for  $3 \cdot 10^5$  iterations using a  
650 learning rate of  $3 \cdot 10^{-4}$  and batches of size 128.

### 651 D.2 TD3-BC [48]

652 This baseline resembles the underlying global policy  $\pi^g$  used in VELAP. It provides us with an  
653 intuition how well pure offline RL performs without adding any planning methods on top.

### 654 D.3 IQL [54]

655 This method presents a state of the art model-free offline RL baseline which utilizes expectile regres-  
656 sion to estimate state-conditional expectiles of the target values in order to avoid querying values of  
657 out-of-sample actions during training. The train IQL for  $3 \cdot 10^5$  iterations, a learning rate of  $3 \cdot 10^{-4}$ ,  
658 batches of size 256 and  $\tau = 0.7$  and  $\beta = 3.0$ .

### 659 D.4 MPPI

660 We implement a trajectory optimization baseline similar to the model-based planning algorithm  
661 introduced in [7]. The method in [7] presents an adaption of MPPI specifically for the online rein-  
662 forcement learning setting which seeks to optimize the expected return of sampled trajectories. To  
663 estimate the return, a learned model is used to predict the reward for each trajectory node while a  
664 learned Q-function predicts the future return beyond the specified planning horizon. Since rewards  
665 in our evaluation environments are sparse, a learned model of the environment reward carries guid-  
666 ance for the trajectories optimization as most states have 0 reward. Therefore, we adapt the objective  
667 in [7] and instead use the accumulated sum of state-action values as the optimization criterion. This  
668 type of scoring function in model-based RL has recently been discussed in [62]. To implement this  
669 baseline, we utilize the Q-function of TD3-BC. For all environments, we use 1000 samples per iteration,  
670 a planning horizon of 50, elite size 64 and 5 iterations. Replanning is done every 5 environment  
671 steps.

### 672 D.5 MBOP [35]

673 MBOP presents an adaptation of MPPI which was particularly designed for the offline RL setting.  
674 It generates new candidate trajectories by adding small amount of Gaussian noise to the actions  
675 predicted by a behavioral-cloned policy. To evaluate the quality of the rollouts it uses a truncated  
676 value function trained on the offline data. Due to the sparse nature of rewards in our experiments, we  
677 found that both the behavioral-cloned policies and the truncated value function were insufficient to  
678 generate farsighted behaviors that solve our tasks. To accommodate for the long planning horizons,  
679 we instead sample action for a TD3-BC policy and use the corresponding Q-functions to score  
680 candidate trajectory during the optimization update. For all environments, we use 1000 samples per  
681 iteration, a planning horizon of 50, elite size 64, 5 iterations and a  $\beta$  parameter of 0.7. Replanning  
682 is done every 5 environment steps.

## 683 D.6 IRIS [36]

684 IRIS presents an offline RL methods that was desired particularly for sparse reward environments.  
685 It uses a hierarchical decomposition of the planning agent into a low-level and a high-level policy.  
686 A conditional VAE model is trained on offline trajectory data and predicts a set of suitable subgoals  
687 states that are  $n$ -step ahead of the current states. The high-level policy is essentially represented  
688 by a Q-function which chooses the highest values subgoal states among the set of generated  
689 candidate states. The low-level policy is then used to navigate towards the predicted subgoal. In  
690 all experiments we implement IRIS by training an conditional VAE to predict subgoal states at a  
691 horizon of 5 and sample a candidate set of size 256. To implement both the low-level and high-level  
692 policy, we use TD3-BC as the base RL algorithm.

## 693 D.7 IRIS (multi-step)

694 We evaluate an extension of IRIS in which we use the state prediction model to generate multi-step  
695 rollouts of suitable subgoals. This strategy increases the exploration horizon and allows to choose  
696 the best subgoal from a larger and potentially more diverse set of states. It exploration strategy can  
697 also be seen as random shooting of coarse subgoal sequences. In all experiments, we generate 256  
698 different trajectories using rollouts of length 5 and a conditional generative model to predict states  
699 for a horizon of 5. In our evaluation, we found this method to sometimes perform worse than IRIS.  
700 We attribute this to the fact that the global policy doesn't align with the capabilities of the local one,  
701 which occasionally results in the selection of subgoal states that might not be attainable.

## 702 E Supplementary Experiments and Analysis

### 703 E.1 Physical hardware experiments

704 For the real-world validation of our method, we collected 200 episodes of data for the sponge ( $\sim$   
705 15000 samples) and 150 episodes of data for the rope manipulation ( $\sim$  15000 samples) tasks. Train-  
706 ing data was generated by operating the robot through a gamepad and took less than 1 hour per task.  
707 The collected dataset consist largely of suboptimal trajectories. Successful transition (positive re-  
708 ward + episode termination) were labeled during data. We form states by stacking three subsequent  
709 images taken by a static camera. The results of a comparison against BC, BC ( $\mathcal{D}^*$ ) and IRIS are  
710 presented in Table 11.

Table 11: Results of physical robot experiments (successful episodes)

Environment	BC	BC $\mathcal{D}^*$	IRIS	VELAP
Sponge	5/20	6/20	6/20	14/20
Rope	0/20	0/20	2/20	8/20

711 **E.2 Ablating the impact of the learned representation**

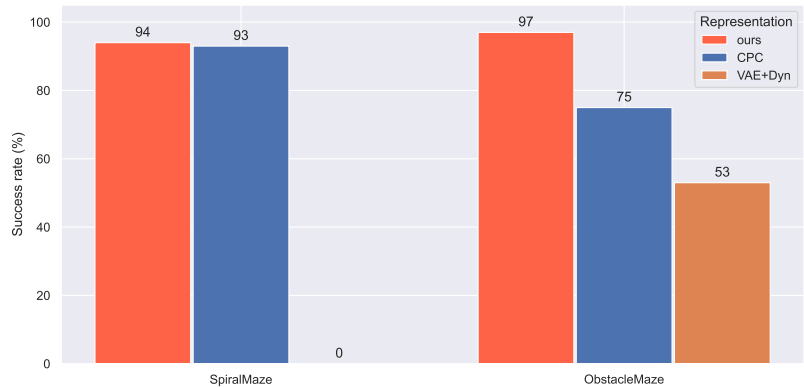


Figure 6: Impact of type of representation on the performance of our planner.

712 **E.3 Influence of the dynamics loss**

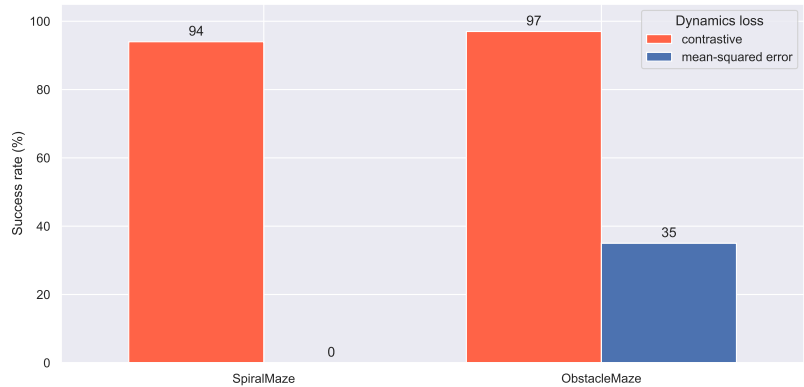


Figure 7: Impact of type of representation on the performance of our planner.