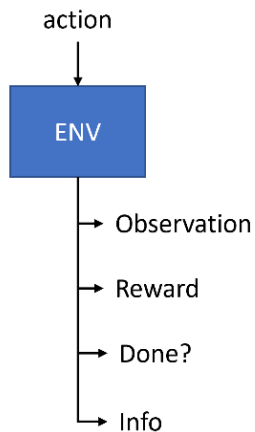


Supplementary Material



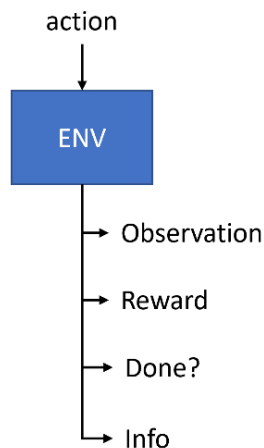
- Action
 - Defines step direction in which the robot moves the block
 - Is either left, right, forward, backward, up or down (6 directions)
 - Each direction axis has a different movement length in the world space

Direction	Left	Right	Forward	Backward	Down	Up
Length	-0.04	0.04	-0.02	0.02	-0.07	0.07

- Observation/state
 - Consist of world space coordinates and 5 laser observations
 - 3 cartesian coordinates (x, y, z) that are **normalized**
 - Output are float values between 0.0 and 1.0

State	X	Y	Z	L1	L2	L3	L4	L5
Val	[0, 1]	[0, 1]	[0, 1]	[0, 1]	[0, 1]	[0, 1]	[0, 1]	[0, 1]

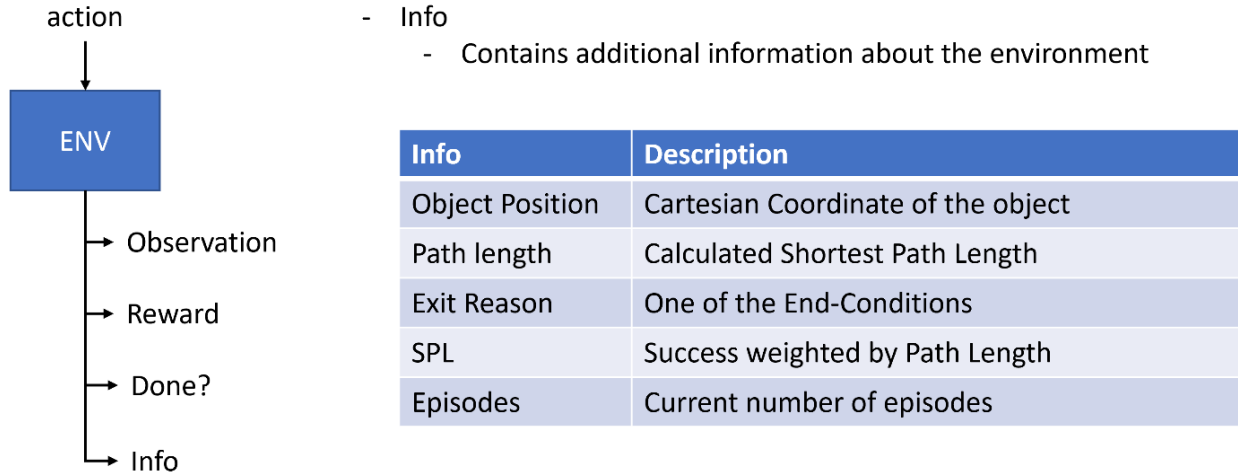
Supplementary Figure 1: Detailed definition of the simulation environment for use in OpenAI Gym – definition of the continuous action space and observation (or state).



- Reward
 - A value corresponding to the reward of each step
 - Every step that doesn't guide to the goal has negative reward
- Done
 - Boolean: If True, an end-condition occurred and the episode ends
 - End-Condition are Collisions/Out-of-bounds, reaching goal, and when maximum of 160 steps per episode is reached (over max steps)

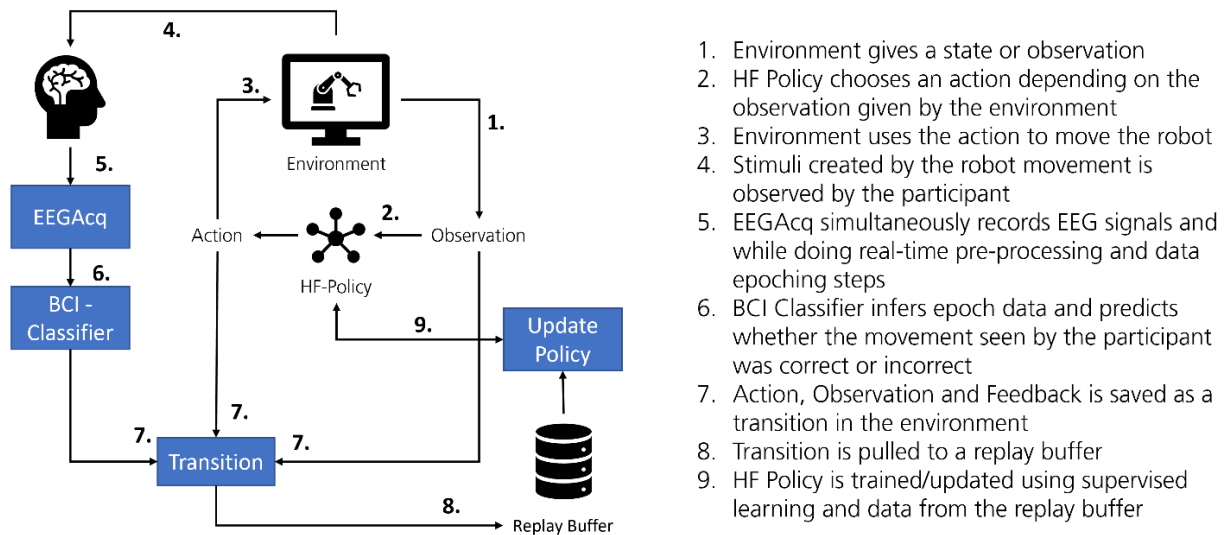
End-Condition	Reward
Step	-0.05
Out of bounds / unreachable position	-8
Plane Collision	-8
Wall Collision	-8
Self Collision	-8
Over Max Steps	-7
Goal Reached	12

Supplementary Figure 2: Detailed definition of the simulation environment for use in OpenAI Gym – definition of reward and done condition (end of a learning episode).



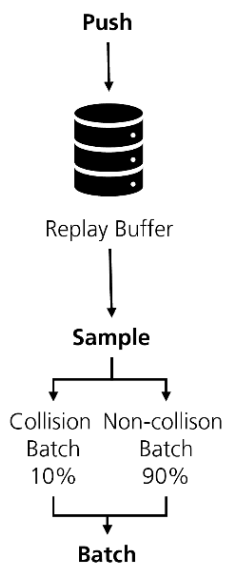
Supplementary Figure 3: Detailed definition of the simulation environment for use in OpenAI Gym – definition of the info.

Pipeline to train a Human Feedback (HF) Policy



Supplementary Figure 4: Detailed description of the real-time pipeline for supervised learning of a human feedback policy function (fully connected neural network).

Optimized Replay Buffer



- Push
 - Pushes Transition (state, action, feedback from BCI) into replay buffer
 - Pushes also information **if Transition was a collision**
- Sample
 - Samples Transitions from replay buffer into a **batch (batch-size=32)**
 - Batch consists of collision and non-collision batches
 - collision batch consists of Transitions where collisions happened
 - Whole batch consist of a maximum of 10% collision batches
 - => Optimized Batch Sampling
 - Makes sure that collision will always be trained
 - Collision-samples depend not on User-Feedback therefore the feedback-labels of those samples are 100% accurate

Supplementary Figure 5: Optimized replay buffer for training the human feedback policy via a fully connected neural network.

```
HFPolicyNet(  
  (linear): Sequential(  
    (0): Linear(in_features=8, out_features=32, bias=True)  
    (1): ReLU()  
    (2): Linear(in_features=32, out_features=6, bias=True)  
    (3): Softmax(dim=1)  
  )  
)
```

Supplementary Figure 6: Fully connected neural network for the human feedback policy.

```
Actor(  
  (net): Sequential(  
    (0): BatchNorm1d(8, eps=1e-05, momentum=0.1, affine=True,  
track_running_stats=True)  
    (1): Linear(in_features=8, out_features=64, bias=True)  
    (2): ReLU()  
    (3): BatchNorm1d(64, eps=1e-05, momentum=0.1, affine=True,  
track_running_stats=True)  
    (4): Linear(in_features=64, out_features=64, bias=True)  
    (5): ReLU()  
    (6): BatchNorm1d(64, eps=1e-05, momentum=0.1, affine=True,  
track_running_stats=True)  
    (7): Linear(in_features=64, out_features=6, bias=True)  
  )  
)  
  
Critic(  
  (net): Sequential(  
    (0): Linear(in_features=14, out_features=64, bias=True)  
    (1): ReLU()  
    (2): Linear(in_features=64, out_features=64, bias=True)  
    (3): ReLU()  
    (4): Linear(in_features=64, out_features=1, bias=True)  
  )  
)
```

Supplementary Figure 7: Deep deterministic policy gradient network architecture consisting of actor-critic neural networks.