
Embodied AI - Final Report

Lukas Rüttgers (2023403372), Simon Wilmes (2023403367)

Tsinghua University

1tgslks23@mails.tsinghua.edu.cn

1 1 Introduction

2 The ManiSkill2 challenge contains a broad spectrum of objects and tasks to tackle the problem of
3 generalization in reinforcement learning (RL). Generalization here covers both the generalization to
4 unseen tasks as well as the generalization to unseen objects. Agents must therefore obtain robust ab-
5 stractions of physical and visual properties of the objects, but also the task structure, to successfully
6 generalize their skills in unseen environments.
7 After presenting the two winner algorithms from the last years, we also provide a review on the
8 research on task generalization, since we feel that we especially lacked knowledge in this field. We
9 present the research in this field and come up with own research questions and ideas which could
10 improve the agent's performance in the ManiSkill setting. Finally, we also provide the videos and
11 learning curve figures of the baselines provided by the ManiSkill authorship.

12 2 Prior work

13 2.1 Winning Algorithms of the Last Years

14 In the original ManiSkill challenges, multiple teams presented ways to achieve the necessary gen-
15 eralization in view of the numerous unseen tasks and objects. In the following, we give a brief
16 summary of the winning algorithms for the Maniskill 2021 challenge and the Maniskill2 challenge,
17 which was conducted one year later.

18 2.1.1 Maniskill 2021 Winner Algorithm

19 The winner algorithm of 2021 basically relies on three major changes to the baseline algorithms
20 which were already provided along the challenge by Mu et al. (2021). While the task-specific
21 rewards yield the potential to accurately guide the robot to the desired goal configuration, the authors
22 argue that adopting to task-specific rewards would contradict the spirit of the challenge to generalize
23 to a broad spectrum of unseen tasks. For this reason, they decide to not use the provided dense
24 rewards. To avoid a sparse reward distribution, they adopt the Generative Adversarial Imitation
25 Learning (GAIL) framework by Ho and Ermon (2016), where the learned policy acts as the generator
26 G and the discriminator D tries to tell the expert trajectories and the trajectories from G apart. Most
27 importantly, D provides dense rewards that scale with the success of G to imitate the expert and
28 hence tackle the reward sparsity issue.

29 However, not only the rewards, but also the expert trajectories follow a separate distribution for
30 each task. The underlying mathematical models of the policy network naturally produce a uniform
31 distribution of trajectories, G will hence struggle to account for the non-uniform distribution of
32 expert trajectories. For this reason, they impose a handicap on D by throttling its model complexity.
33 Nonetheless, the embedding architecture of D remains the same over the entire evolution. As ob-
34 servation, D receives the robot state, the action and a point cloud with labeled segmentations of all
35 parts that are involved in the specific task. The point cloud is masked for each part separately and
36 forwarded through a dedicated PointNet. To account for the point that belong to neither of the parts,
37 another two PointNets embed the rest of the point cloud and the entire point cloud respectively.
38 Moreover, a Multilayer-Perceptron (MLP) extracts features out of the robot state and the applied
39 action and together with a bias form the last two components of the embedding, which is finally
40 concatenated.
41 This embedding is then forwarded through both a simple and a strong model, where the simple
42 model merely selects the maximum as the final feature and the strong model leverages the encoding
43 part of the transformer architecture proposed by Vaswani et al. (2017) to identify relations between
44 the different part embeddings via self-attention. The final output of D is simply a linear interpolation
45 between this two models which tends to the strong model over time.
46 Since this still does not account for the distribution offset between the non-uniform expert trajec-
47 tories and the trajectories from G , the authors fill the expert buffer with successful trajectories of G .
48 To ensure that a relatively easy task does not lead to flooding the expert buffer with trajectories from
49 this task, a fixed amount of space is reserved for each task. The resulting buffer structure is coined
50 as *Category-Level Instance Buffer* (CLIB).

51 2.1.2 ManiSkill2 2022 Winner Algorithm

52 Gao et al. (2023) sought out pragmatic solutions to these problems. Their core idea is to allow
53 a more stable learning process by splitting the training process into two parts: a general learning
54 phase, followed by a fine-tuning phase. The reasoning behind this split was that the authors noticed
55 that after an initial surge in the success probability after a certain number of training steps, the
56 algorithm would start to overfit and thereby reduce its success performance. To overcome this issue,
57 they stored the best-performing network and used that network as the starting point for the second
58 round of fine-tuning, in which they reduced the batch size and the number of samples in each step.
59 The authors claim that this allows their algorithm to focus on more minor details.
60 The general setup of their algorithm is first to use the PointNet framework by Qi et al. (2017) to
61 segment the point cloud captured by the robot into features that can directly be passed as input into a
62 traditional network. They employed Behaviour Cloning as their learning strategy for imitation learn-
63 ing, while for Reinforcement Learning, they used Proximal Policy Optimization (PPO) (Schulman
64 et al. (2017)).
65 Even though their optimization ideas are not fundamentally different from more standard ideas for
66 this task, by making these small adjustments, this setup achieved first place in all three tracks,
67 thereby showing the advantage of their approach.

68 2.2 Competence progress motivation for intrinsic rewards

69 Our ability to acquire a thorough understanding of the dynamics inherent in our environment is
70 largely attributable to our intrinsic motivation that drives us to explore beyond our comfort zone.

71 Among several psychological theories of intrinsic motivation, Hunt (1965) famously coined intrinsic
72 motivation as the strive for *optimal incongruity* between humans' expectation and their actual
73 perception. In face of the sparsity of extrinsic rewards in long-horizon reinforcement learning (RL)
74 tasks, Oudeyer and Kaplan (2007) derived intrinsic reward formulations from these psychological
75 theories and associated the notion of incongruity with *learning progress motivation*, that generates
76 high rewards for decreasing epistemic uncertainty in an agent's world model. Schmidhuber (1991)
77 proposed to measure epistemic uncertainty as the prediction error of a forward dynamics predictor
78 that is learned during exploration. Until today, most of the research on intrinsic rewards follows
79 this paradigm. While Lakshminarayanan et al. (2017), Burda et al. (2018) and Pathak et al. (2019)
80 estimate uncertainty as the disagreement between an ensemble of predictors, Sekar et al. (2020)
81 imagine rollouts inside a latent space to estimate the expected novelty of a future state.

82 However, there is also another influential psychological theory around intrinsic motivation, which
83 describes it as craving for competence (Deci and Ryan (1985)) or challenge (Csikszentmihalyi
84 (2008, cop. 1990)). It claims that humans deliberately expose themselves to self-determined chal-
85 lenges that exhibit an adequate difficulty to efficiently foster their current aptitudes. In their transfer
86 of this conceptual idea to RL, Oudeyer and Kaplan (2007) merely discussed how to express the
87 difficulty of a challenge and culminated in a coarse distance-based metric, which they coined as
88 *competence progress*. Stout and Barto (2010) concretized this in a limited setting, as they still as-
89 sumed a finite state and action space and a given fixed set of challenges, where each challenge was
90 defined by a given value function.

91 2.3 Project Motivation

92 Therefore, the computational formulation of challenge-based motivation is incomplete: how does
93 the agent discover new tasks in a self-supervised manner? How can he efficiently determine his
94 competence progress in high-dimensional infinite domains? And how does he integrate and keep
95 track of the abundant number of possible skills in one policy? This requires lifting the goal from a
96 mere parameter to part of the policy's input. But to efficiently represent such a policy, the agent must
97 learn analogies between tasks and find patterns by which to compress their representation, similar
98 to how Schmidhuber (2010) described the compression of novel patterns in observations.

99 Although these questions have received little attention in the last decade of research, we believe
100 that competence-based intrinsic rewards can push agents to discover new skills and learn affinities
101 between them, which in turn is an essential building block towards learning generalizable and robust
102 behaviour.

103 After analysing and discussing possible computational architectures for intrinsic motivation through
104 competence progress, we want to derive efficient implementations and finally verify their contribu-
105 tion to generalization to unseen tasks on the ManiSkill2 benchmark (Gu et al. (2023)).

106 3 Computational Architectures for Competence Progress Motivation

107 Firstly, we will construct an ideal, general architecture that could capture the dynamics of compe-
108 tence progress motivation. Because completely implementing this ideal formulation on our own
109 requires engineering efforts that go beyond the scope of this project, we secondly derive an alter-
110 native architecture that builds upon prior work and enables us to realise the core components of
111 our idea more efficiently. Throughout the upcoming sections in the report, we present ideas on an

112 intuitive level, and refer the interested reader to more extensive explanations and thoughts in the
113 appendix.

114 **3.1 Ideal Two-Party Formulation**

115 Let us consider the dynamics of competence progress motivation as a symbiosis between two parties,
116 the *challenger* C , and the *player* P . The goal of both parties is to maximize the competence of P .
117 To that end, C is supposed to design a challenge-based curriculum for P , while P merely tries to
118 perform as good as possible in these challenges.

119 **3.1.1 Challenge description and criteria**

120 To create an effective learning incentive, C should dynamically adapt the challenges to the current
121 learning state of P . On the one hand, P will not profit from challenges that are already well-known
122 and easy to him. On the other hand, challenges that are too complex will overwhelm P . Both cases
123 will not lead to significant progress in this challenge. As originally modeled by Oudeyer and Kaplan
124 (2007), P could therefore repeatedly attempt the challenge and measure his competence progress
125 as the increase in reward over the last trials. The larger this increase is, the higher P 's progress
126 might be. P then reports this progress to C , who in turn can reflect on the adequacy of the current
127 challenges, better understand the strengths and weaknesses of P , and thereupon improve her future
128 challenge suggestions. Such a feedback is also indispensable to track the current competence state
129 of P . The interaction between C and P hence resembles the interaction between a student and a
130 teacher who actively adjusts the student's curriculum to his progress. Besides adequate difficulty,
131 challenges should be meaningful with regard to the actual task that P is trained for.

132 Before we continue, we should first provide a formal definition of a challenge. The most naïve
133 approach would regard a challenge as a state that ought to be reached. Because the ideal state itself
134 is mostly inaccessible, we can instead represent the goal state by its associated observation from a
135 fixed frame. To abstract from this detail, we refer to both as a state $q \in \mathcal{Q}$ in the following.

136 However, a challenge description may not only comprise *where* P shall get to. The challenge may
137 also lie in *how* P is commanded to get there. To express more complicated challenges of such kind,
138 we could include multiple landmarks as intermediate sub-goals that P should traverse up to some
139 error on his way to the final goal state. For simplicity, we will however stick with the restricted
140 definition of a challenge that only includes the observation of one final goal state. Besides that, it
141 is crucial to also specify the initial state distribution of a challenge as we elaborate on in Appendix
142 A.1.

143 **3.1.2 Challenge discovery and evaluation**

144 In order to serve as a useful teacher, C should have advantageous knowledge in comparison to P ,
145 otherwise it is unlikely that C can come up with challenges that are both meaningful to the task
146 and enriching for P 's competence progress. This information asymmetry is the main reason why to
147 divide these roles into two different logical units.

148 While *Teacher-Student Architectures* like Rapid Motor Adaptation proposed by Kumar et al. (2021)
149 provide the teacher with privileged state information, this privileged state information might not
150 always be available. We therefore undertake a more general formulation where the challenger C
151 itself has to acquire knowledge and talent in the same way as P , but might have an edge in experience
152 and information. That is, while P 's view is limited to the challenge-oriented rewards, C has access to
153 the extrinsic reward function of the actual environment and can explore it by herself in the usual way.

154 Regions where C firstly perceives high extrinsic rewards, or secondly acquires new knowledge about
155 the environment dynamics are both likely to be valuable for P too. To measure gain in knowledge,
156 we can avail to the broad battery of computational formulations for *learning progress motivation*
157 such as the epistemic uncertainty estimators mentioned in Section 2.2.

158 By this interaction with the true environment, C can assess which regions might be particularly
159 attractive for P 's curriculum. She can do so by integrating the aforementioned measurements into a
160 *challenge quality score* $S : \mathcal{Q} \rightarrow \mathbb{R}$ that estimates the quality of setting a state $q \in \mathcal{Q}$ as the goal of
161 a challenge. Most importantly, this function is not static but evolves over time, because the learning
162 progress on which the score is partially based is also subject to change.

163 C may now greedily choose the state $q^* := \arg \max_{q \in \tau} S(q)$ with the highest score as candidate for
164 the next challenge. More sophisticated approaches might also consider the scores of the entire asso-
165 ciated trajectory $\tau \ni q^*$ or at least its predecessors and successors as the local context window – we
166 borrow terminology from Natural Language Processing here – to integrate more global information.

167 3.1.3 Challenge update and multiple challenges

168 For now, let us assume C can pose only one challenge at each point in time. C stores the current
169 challenge as the observation o_g of a certain goal state g in a shared memory that both C and P have
170 access to. We will refer to this shared memory as the *challenge set* Ω .

171 After a challenge g was selected by C and undertaken by P , the challenge quality score $S(c)$ should
172 be updated based upon the feedback from P . Specifically, the learning progress upon which C orig-
173 inally based the *estimated* competence progress of P is now step-wisely replaced by the *empirical*
174 competence progress that P reports to C . This could be accounted for by a bilinear superposition
175 $\alpha \cdot \hat{x} + (1 - \alpha) \cdot x$ where \hat{x} is the competence progress estimate and x the empirically reported
176 competence progress, and $\alpha \in [0, 1]$ a parameter that starts at 1 and continues to decay to 0 for each
177 challenge trial P undergoes. The extrinsic rewards are static and do not need to be replaced in the
178 score term.

179 This sets the stage to discuss under what conditions C updates the current challenges. Intuitively,
180 this should be based upon their respective scores. If the current challenge g still proves more valuable
181 in terms of $S(c)$ than the current challenge candidate q^* , then q^* may be stored in a buffer and taken
182 into consideration in the upcoming update phases. However, as time passed by, the challenge quality
183 score of a candidate $S(q)$ might become deprecated as it does not reflect the *current* competence and
184 knowledge of P any more. For this reason, such candidate states should not be kept too long in such
185 a buffer. Alternatively, their score $S(q^*)$ could be discounted over time by a factor λ .

186 The overall reward and information flow is again depicted in Figure 1. As it illustrates, Ω might
187 in general comprise multiple challenges that introduces new degrees of freedom such as ordering
188 challenges in Ω by their value. On the other side, it also necessitates a selection strategy on P 's
189 side. If this selection strategy is in turn controlled by C , then this is equivalent to the previous
190 single-challenge setting.

191 3.1.4 Pseudo-rewards for selected challenges

192 By the nature of RL, C must accompany a challenge g with a proper reward function that lures P
193 to g . With regard to the maxim "*the way [and not the goal itself] is the goal*", we have to carefully
194 consider *what* we actually want to teach P on the way to the goal when defining this reward function.



Figure 1: Schematic overview of the reward and information flow between C and P .

195 Besides the *extrinsic* reward from the environment, the proximity to the goal state in some distance
 196 metric would provide a dense reward signal. What's more, if there are some particular regions that
 197 C wants P to explore more thoroughly on the quest to g , then C can also integrate this bias into the
 198 reward function. In Section 4.4.3 on our implementation, we specify how we eventually compose
 199 our reward function for challenges. In the following, it suffices to assume a reward function that
 200 properly integrates all of the aforementioned aspects.

201 3.1.5 Considerations for implementation

202 In general, it might be tedious to infer P 's current state of competence only from his competence
 203 progress feedback. To simplify the search for challenges of adequate difficulty, we provide C ac-
 204 cess to the current policy network of P . When taking actions during challenge exploration, C can
 205 now decide between taking an action based upon P 's current policy or following her own – more
 206 proficient – exploration network. The more C will avail to P 's policy network, the easier and well-
 207 known the traversed regions and challenges will become for P , hence a sound amount of exploration
 208 is indispensable for adequate challenge difficulty.

209 Another difficulty lies in the active curriculum that C has to design to teach P to finally achieve
 210 competence at the environment-specific task. Over the course of challenges, our changing artificial
 211 rewards might confuse P 's policy updates and in the worst case lead to deliberate unlearning of
 212 useful knowledge it has achieved in prior challenges. Even without this pitfall, C needs to ensure
 213 that through the mastery of her challenges, P 's policy eventually converges to a policy that fulfils
 214 the environment task. For that reason, we might at some point pose the environment goal state itself
 215 as the final challenge or integrate it frequently into Ω throughout the curriculum to ensure the P 's
 216 policy is sufficiently prepared for transitioning to the final goal of the environment.

217 3.2 Off-Policy Alternative

218 While our original formulation could have worked, we also thought about a different architecture.
 219 Our main concern with the original idea was the high complexity as it consisted of multiple net-
 220 works each with their own task, for which we sometimes even did not know if their specific goal
 221 was realistic. A more efficient and stable way to realise our core ideas was to build upon an already
 222 existing work and modify it to suit our needs. In our case, this was Hindsight Experience Replay
 223 (HER) by Andrychowicz et al. (2017) as their idea contains similarities that we thought were worth

224 exploiting. For one, the idea of training the model not only on trajectories generated with the environment reward but also on trajectories with an artificial reward is an essential component of both
225 training processes. In contrast to HER where the reward usually signifies a single goal state that
226 ought to be reached, our idea rather formulates the reward and challenge goal with the objective
227 of maximizing the learning experience *on the way* to the mastery of the challenge. Although this
228 necessitates to adjust our formulation from an on-policy to an off-policy learning strategy, it actually
229 introduces more freedom in the way the trajectories are generated. In particular, we do not have to
230 account for the initial state distribution problem explained in Appendix A.1.
231

232 3.2.1 Structure and challenge discovery behaviour

233 First of all, let us present the structure of our idea in detail. We stick with our original two-party
234 formulation in the sense that we regard P and C as two agents. On the one hand, C is trained
235 to accurately estimate the epistemic uncertainty of a transition. At the same time, C generates
236 trajectories by alternating between its own exploration policy and P 's policy network.

237 Afterwards, we fill two buffers with these trajectories. The first buffer comprises the complete
238 trajectories with their *extrinsic* reward. We refer to this buffer as the *original* buffer D . The second
239 buffer is dedicated to be filled with our artificial challenges and hence referred to as the *imaginary*
240 buffer I . To that end, C recorded her uncertainty estimates for each transition. She now uses
241 these estimates along with the extrinsic rewards to compute a challenge quality score of each state
242 q . For each trajectory, she truncates the trajectory after the state $q^* := \arg \max_{q \in \tau} S(q)$ with the
243 highest score. Afterwards, she assesses each truncated trajectory based on their scores and other
244 statistics and includes only those trajectories into the buffer I that exhibit sufficient quality. For
245 these sufficiently good trajectories, it replaces the extrinsic reward by a linear combination of an
246 exploration reward, the extrinsic reward, and a distance reward that is supposed to increase for
247 transitions that lead towards the end state of the truncated trajectory q^* . All of these three steps – the
248 computation of the challenge quality score, the assessment of trajectories, and the computation of
249 the new reward – are fundamental elements of our idea and crucial to the success of our algorithm.
250 For this reason, we will thoroughly elaborate on their implementation later in Section 4.

251 These two buffers will then be used to train P in an off-policy fashion. Regarding the training
252 algorithm of P , we will stick with the original proposal in the HER paper to use Soft Actor-Critic
253 (SAC) (Lakshminarayanan et al. (2017)). The concept of *challenge* has now acquired a more indirect
254 meaning. Instead of directly attempting the challenge, P is challenged to integrate the additional
255 information from the imaginary buffer I into the update of its policy. If C filled I wisely, then the
256 transitions in I will point toward the weak spots of P 's policy and provoke him to improve in these
257 aspects.

258 3.2.2 Measuring uncertainty

259 So far, we have not specified how C acquires good estimates of the uncertainty of a transition. As an-
260 nounced in Section 3.1.2, we avail to prior work and use the Dreamerv2 World Model (Hafner et al.
261 (2019)) and its Plan2Explore extension (Sekar et al. (2020)). While the idea of long-horizon plan-
262 ning in learned latent representations of the high-dimensional observation space demonstrated im-
263 pressive results in exploration-heavy environments with complicated tasks, Plan2Explore has been
264 designed to build upon this strong world model to find regions of high novelty, which is exactly what
265 we expect the exploration agent to do. Furthermore, the world model also proves useful for the com-
266 putation of our artificial rewards, as we explain later in Section 4.4.3. Reusing these prior models
267 enables us to focus on our core contributions, which already exhibit enough freedom to determine

268 success or total failure of our idea. Algorithm 1 summarizes our training process. More specific
implementation details are subject of the following section.

Algorithm 1 CPM Dreamer: Training process

Input: Main Agent $A := (C, P)$ with world model WM, image encoder ENC,
original buffer D , imaginary buffer I , number of trajectories generated per episode T ,
challenge quality score S , artificial reward formula R , challenge filter $F : \tau \mapsto \{0, 1\}$.

```

1: for all episodes do
2:   for  $T$  times do
3:     Generate one trajectory  $\tau$  by exploiting  $P$  and exploring with  $C$   $\varepsilon$ -greedily.
4:     For each state  $s_t \in \tau$ ,  $C$  estimates its uncertainty as  $u_t$ .
5:     Store  $\tau$  in  $D$ .
6:     for all transitions  $(s_t, a_t, r_t, s_{t+1}) \in \tau$  do
7:        $C$  estimates  $u_t \leftarrow$  epistemic uncertainty in  $s_t$  according to WM.
8:        $C$  embeds  $e_t \leftarrow \text{ENC}(s_t)$ .
9:        $C$  evaluates the score  $S(s_t, r_t, u_t)$ .
10:    end for
11:     $C$  obtains  $g := \arg \max_t S(s_t, r_t, u_t)$  and truncates  $\tau' := \tau_{[1:g]}$ .
12:    if  $F(\tau') = 0$  then
13:      continue with next trajectory.
14:    end if
15:    for all transitions  $(s_t, a_t, r_t, s_{t+1}) \in \tau'$  do
16:       $C$  calculates the distance improvement  $d_t \leftarrow \|e_g - e_t\| - \|e_g - e_{t+1}\|$ .
17:       $C$  replaces  $r_t$  by new reward  $r'_t := R(s_t, r_t, u_t, d_t)$ .
18:    end for
19:    Store  $\tau'$  in imaginary buffer  $I$ 
20:  end for
21:  Train WM, ENC and  $C$  off-policy with  $D$ .
22:  Train  $P$  off-policy with both  $D$  and  $I$ .
23: end for
```

Figure 2: Pseudocode for training our Competence Progress Motivation (CPM) Dreamer model. In our case, P is an archetypical SAC agent and the challenger C is the Dreamerv2 model with the Plan2Explore extension that already includes a world model and an image encoder.

269

270 4 Implementation Details

271 As already mentioned before, we realised P as a standard SAC model and availed to the SAC
272 implementation in the ManiSkill2-Learn repository as reference.

273 Regarding the controller, it proved most effective to use the pd_ee_delta_pose controller, as we
274 reason more thoroughly in Appendix B.3.

275 Because it provides the richest learning signal, we decided to receive dense and not sparse rewards.

276 As reported on GitHub, multi-threaded rendering issues prevented the parallel rollout of the agent’s
277 policy in the environment during training. For this reason, one single agent instance was executed
278 in the environment during environment, which slowed down the training process.

279 4.1 Dreamer and Plan2Explore

280 On the other hand, we implemented the challenger C as the Dreamerv2 World Model with the
281 Plan2Explore extension. Unfortunately, there were a few conflicts between these models. Firstly,

282 Dreamer expects RGB inputs while our implementation of SAC relies on a PointNet. For this reason,
283 we adjusted the observation wrappers to receive both RGB and Point Cloud observations from the
284 environment.

285 For the Point Cloud input, we sample additional 50 points near the goal observation if the
286 goal state is visualized in the environment, such as is in PickCube-v0, PickSingleYCB-v0, and
287 PickSingleECAD-v0. This was recommended by the authors and proved helpful in performance.
288 The RGB images are captured from both the end-effector camera and the external camera, and were
289 scaled to 32×32 pixels.

290 Secondly, Dreamer and Plan2Explore are both implemented in TensorFlow and hence incom-
291 patible with PyTorch’s automatic differentiation package autograd that the reference imple-
292 mentation of SAC was built upon. For that reason, we decided to reimplement Dreamer and
293 Plan2Explore in PyTorch. Although there was already a reference reimplementation of both mod-
294 ules at <https://github.com/jsikyoon/dreamer-torch/tree/main>, it turned out to contain sev-
295 eral data-processing errors and unfinished code passages – especially in the Plan2Explore extension
296 – that needed to be fixed in the first place.

297 Furthermore, we also decided to deviate from the reference implementation in a few places and
298 adjusted the numerous hyperparameters to the needs of our tasks. This includes but is not limited to
299 the depth and kernels of the CNN encoder and decoder, the hidden state of the GRU, the size of the
300 latent state, the layer size and depth of the actor, value, and reward head, and their learning rates.
301 We extensively elaborate on the structure of Dreamer, our deviations, and our considerations for
302 fine-tuning in Appendix B.1. For similar insights in our implementation of Plan2Explore, where we
303 most importantly propose a different metric to rate the uncertainty of a state, we refer the interested
304 reader to Appendix B.2.

305 4.2 Challenge Exploration Decay

306 When generating trajectories for the replay buffers I and D , we started with a high initial challenge
307 exploration probability ε and applied a decay that steadily enlarges the impact of P ’s policy network
308 on the generated trajectories. The decay model should firstly ensure that ε remains high enough in
309 the initial phase where P ’s policy network still exhibits much randomness. Secondly, it should have
310 a gentle decay in the long-term to maintain a small amount of challenge exploration. We present the
311 impacts of several decay models and the empirical advantages of our final choice in Appendix ??.

312 4.3 HER and SAC

313 The SAC agent was trained on both the original and imaginary replay buffer. In each update step,
314 one batch was randomly sampled from the original buffer with a probability of 75% and from the
315 imaginary buffer with a probability of 25%. The agent was then trained with this batch and this ran-
316 dom process was repeated six times. In this way, we could ensure a sufficiently accurate perception
317 of the true rewards in the environment while still providing imaginary challenges enough impact.

318 Moreover, we chose the size of the original buffer to contain 50 full trajectories, while the imaginary
319 buffer could only contain 10 trajectories. This is because the extrinsic rewards stay permanently
320 valid, while the imaginary rewards builds on time-variant uncertainty and knowledge and their con-
321 sistency decays with time.

322 **4.4 Challenge Choice**

323 In the following, we will explain how we implemented the challenge quality metric that rates each
324 state in the generated trajectory. Secondly, we present how we filter trajectories and states to obtain
325 our final challenge trajectories for the imaginary buffer I . Thirdly, we derive the artificial reward
326 terms that replace the extrinsic rewards in these challenges.

327 **4.4.1 Challenge Quality Metric**

328 Promising regions for challenges should not only be novel to P , but also exhibit some relevance to
329 the task. We hence decided to compute our challenge quality metric as a linear combination of the
330 extrinsic reward r_t and the uncertainty estimate u_t as

$$S(s_t, r_t, u_t) = c \cdot r_t + u_t, \text{ where } c \cdot r_t \approx 0.75 \cdot u_t.$$

331 Furthermore, u_t is computed as a linear combination of the KL loss between the stochastic prior and
332 posterior of Dreamer's latent state KL-LOSS and the disagreement within the predictor ensemble of
333 the Plan2Explore extension DISAG.

$$u_t = c_1 \cdot \text{KL-LOSS} + \text{DISAG}, \text{ where } c_1 \cdot \text{KLLOSS} \approx 0.1 \cdot \text{DISAG}.$$

334 Because all values are subject to significant changes throughout the training process, the coefficients
335 c and c_1 must be dynamically adjusted. Since this requirement appears useful for many hyperpa-
336 rameters and coefficients, we implemented a general mechanism to adjust hyperparameters if they
337 deviate from a target ratio too strongly. This mechanism is throughly introduced in Appendix B.7.

338 **4.4.2 Challenge Filtering with Self-Balancing Strictness**

339 Not necessarily all trajectories serve as useful challenges for I . For that reason, a balanced filtering
340 strategy $F : \tau \mapsto \{0, 1\}$ that discards trajectories with too feeble attributes is crucial to maintain
341 a high challenge quality. Intuitively, F validates three criteria. If τ does not fulfil any of them,
342 it is discarded. Firstly, the exploration reward u_t should not be negligibly small compared to r_t .
343 Otherwise, the challenge is not novel enough. Secondly, the reward r_t should not be abysmal low
344 to ensure that a challenge is also related to the actual task. Here, "low" must be regarded relative
345 to the previous reward history of P . Thirdly, there should be some development on the way to the
346 challenge goal g . By definition, g has the highest quality score on the entire trajectory. But if this
347 improvement is relatively small compared to the average scores, then there will unlikely be any
348 advantage of reaching g .

349 We thoroughly explain how we eventually realise these three criteria in a stable way in Appendix
350 B.5.2.

351 **4.4.3 Progress-Sensitive Artificial Rewards**

352 Similar considerations as for the goal quality scores had to be made when defining how the reward
353 term for artificial challenge trajectories is constituted.

354 Briefly, our reward term r'_t consists out of four terms. Firstly, the extrinsic reward r_t , secondly, the
355 uncertainty term u_t , thirdly, the improvement in the *embedding distance* to the goal state $d_{emb,t}$, and
356 finally, the improvement in the *step distance* to the goal state $d_{step,t}$.

357 The embedding distance improvement $d_{emb,t}$ is computed as

$$d_{emb,t} := \|e_g - e_t\|_1 - \|e_g - e_{t+1}\|_1,$$

358 where $\|\cdot\|_1$ denotes the Minkowski distance with $p = 1$.

359 Moreover, the step distance $d_{step,t}$ embedding is just a reward that increases proportionally towards
360 the goal state to further lure the SAC agent to explore this state in future trajectories.

361 To preserve consistency with the extrinsic reward of the same transitions in the original buffer, we
362 scale the other three terms u_t , $d_{emb,t}$, and $d_{step,t}$ to meet certain target ratios to the extrinsic reward
363 r_t in average.

364 Specifically, we ensure that the extrinsic reward is approximately 1.5 times higher than the explo-
365 ration reward u_t , while the mean of the distance terms is targeted to be 0.3 times smaller than u_t .

366 To meet the aforementioned target ratios, we again avail to our dynamic hyperparameter adjustment
367 mechanism that we thoroughly introduce in Appendix B.7.

368 To remain in a similar range with the extrinsic rewards in the original buffer, the entire resulting
369 linear combination is then scaled by 0.5.

370 We thoroughly explain this design decision of the artificial rewards in Appendix B.6.

371 4.5 Dynamic Hyperparameter Targeting

372 Throughout the learning progress of the world model and the agent, most values like rewards, uncer-
373 tainty, or embeddings will fluctuate strongly. Static hyperparameters can not account for this strong
374 deviations in magnitude. By dynamically adjusting the hyperparameters to always meet a certain tar-
375 get ratio, we can ensure their balanced interplay independent of environment-specific value ranges.

376 For this reason, we implemented a general dynamic hyperparameter adjustment mechanism that
377 we applied to almost all of the aforementioned scales for the artificial reward terms, the challenge
378 quality metric, or the uncertainty term.

379 In all of these cases, we formulate a target ratio r^* that the hyperparameter h over some comparison
380 statistic g should satisfy. If this ratio deviated too strongly, that is $h < \omega_{low} \cdot r \cdot g$ or $h > \omega_{high} \cdot r \cdot g$
381 for some static tolerance bounds $\omega_{low}, \omega_{high}$, the parameter h was increased or decreased by a static
382 scale $\lambda > 1$. Here, the tolerances and the adjustment factor were chosen depending how strictly we
383 wanted to satisfy this target ratio and how high we expected its variance to be.

384 We thoroughly explain our implementation in Appendix B.7.

385 5 Experiments

386 5.1 Baselines

387 As requested, we implemented the baselines as seen in
388 examples/tutorials/reinforcement-learning/sb3_ppo_liftcube_rgbd.py. We ran the
389 code for all 5 environments: "PickCube-v0", "PickSingleEGAD-v0", "PickSingleYCB-v0",
390 "StackCube-v0", "TurnFaucet-v0". Because of a rendering error in the cabinet mesh, we did
391 not conduct experiments on "OpenCabinetDoor-v1". Although the baseline was able to achieve
392 sound results in the LiftCube-v0 environment during our mid term report, it fails in all of the five
393 environments, which are more challenging. The authors reported that these baselines performed

394 better in an earlier version 0.3.0 of ManiSkill, but after update their rendering to 0.4.0, only the
DAPG baseline seemed to preserve good performance according to their GitHub page.

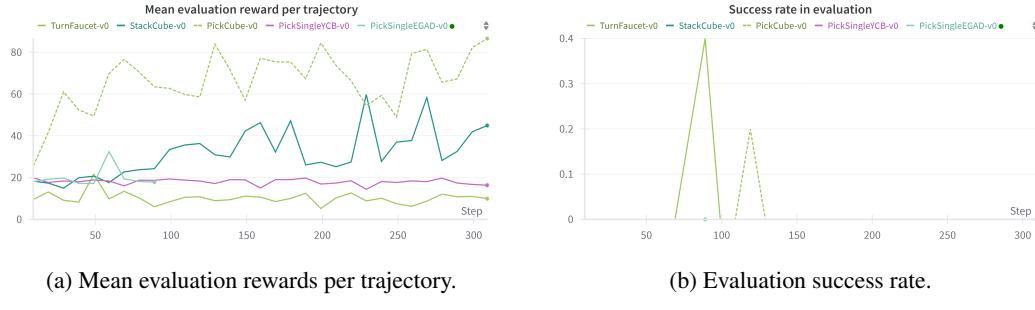


Figure 3: Only in the PickCube-v0 environment, the PPO SB3 baseline succeeds seldom during evaluation. In all other environments, it never succeeds.

395

396 5.2 CPM Dreamer

397 To showcase the potential of our idea, we compare two handicapped versions of our model with a
398 strong Dreamer model, a PPO that relies on PointCloud observations, and the PPO SB3 baseline that
399 receives RGBD inputs on the PickCube-v0 task. Despite these handicaps, our models achieve the
400 most frequent success rates already during an early training stage. Although the PPO PointCloud
401 model dominates all other models with regard to rewards, it does not succeed once over a far longer
402 period. Although our handicapped CPM Dreamers already achieve some success, it occurs fairly
403 seldom, which we account to the exaggerated exploration and weakness of the models respectively.
404 The success and evaluation rates are displayed in Figure 4.

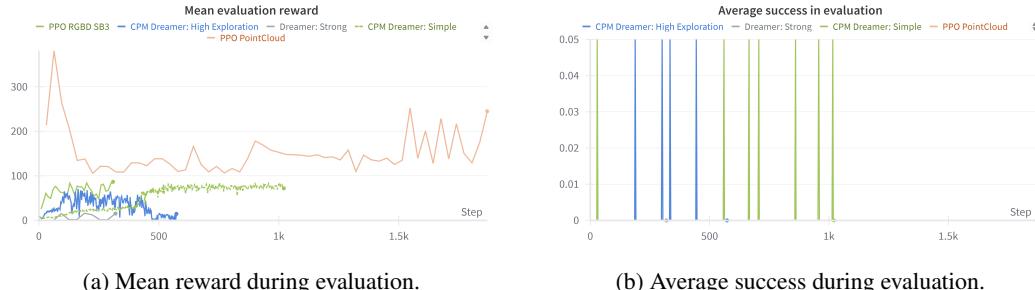


Figure 4: Five models are evaluated on the PickCube-v0 task. Among them are two handicapped versions of our CPM Dreamer model. While the CPM Dreamer: High Exploration model has a far too small exploration decay and therefore does not explore without letting the player P contribute to the trajectories and hence foster his task-specific knowledge, the CPM Dreamer: Simple

405 As the training procedures for stronger versions of our models unfortunately crashed at least after
406 a few hours because the GPU memory was quite occupied in the last days, we can not demonstrate
407 yet how large the advantage of our final, fine-tuned model over the baselines actually is.

408 But as these weak versions already exhibit an advantage in evaluation regarding the eventual success,
409 we believe that our idea might contain at least some components such as the dynamic hyperparameter
410 adjustment, the self-balancing challenge filtering and evaluation that prove useful in providing a
411 good learning experience to our main SAC agent.

412 More experiments in the Appendix C underscore the ability of our model to learn reasonable latent
413 representations of the environment and adjust to new ranges of time-variant hyperparameters.

414 6 Conclusion

415 However, there are still possible improvements of the model. More research could be conducted on
416 more sophisticated distance terms in the goal evaluation reward, or the self-balancing policy in goal
417 filtering.

418 References

- 419 M. Andrychowicz, F. Wolski, A. Ray, J. Schneider, R. Fong, P. Welinder, B. McGrew, J. Tobin,
420 O. Pieter Abbeel, and W. Zaremba. Hindsight experience replay. *Advances in neural information*
421 *processing systems*, 30, 2017.
- 422 Y. Burda, H. Edwards, A. Storkey, and O. Klimov. Exploration by random network distillation.
423 *arXiv preprint arXiv:1810.12894*, 2018.
- 424 M. Csikszentmihalyi. *Flow: The psychology of optimal experience*. HarperPerennial ModernClassics,
425 New York [etc.], 1st harper perennial modern classics ed. edition, 2008, cop. 1990. ISBN
426 0060162538.
- 427 E. L. Deci and R. M. Ryan. *Intrinsic motivation and self-determination in human behavior*. Springer Science+Business Media, New York, 1985. ISBN 978-1-4899-2273-1. doi:
428 10.1007/978-1-4899-2271-7.
- 430 F. Gao, X. Li, J. Yu, and F. Shaung. A two-stage fine-tuning strategy for generalizable manipulation
431 skill of embodied ai, 2023.
- 432 J. Gu, F. Xiang, X. Li, Z. Ling, X. Liu, T. Mu, Y. Tang, S. Tao, X. Wei, Y. Yao, X. Yuan, P. Xie,
433 Z. Huang, R. Chen, and H. Su. Maniskill2: A unified benchmark for generalizable manipulation
434 skills. In *International Conference on Learning Representations*, 2023.
- 435 D. Hafner, T. Lillicrap, J. Ba, and M. Norouzi. Dream to control: Learning behaviors by latent
436 imagination. *arXiv preprint arXiv:1912.01603*, 2019.
- 437 J. Ho and S. Ermon. Generative adversarial imitation learning. *CoRR*, abs/1606.03476, 2016. URL
438 <http://arxiv.org/abs/1606.03476>.
- 439 J. M. Hunt. Intrinsic motivation and its role in psychological development. *Nebraska Symposium*
440 *on Motivation*, 13:189–282, 1965.
- 441 A. Kumar, Z. Fu, D. Pathak, and J. Malik. Rma: Rapid motor adaptation for legged robots. *arXiv*
442 *preprint arXiv:2107.04034*, 2021.
- 443 B. Lakshminarayanan, A. Pritzel, and C. Blundell. Simple and scalable predictive uncertainty esti-
444 mation using deep ensembles. *Advances in neural information processing systems*, 30, 2017.
- 445 T. Mu, Z. Ling, F. Xiang, D. Yang, X. Li, S. Tao, Z. Huang, Z. Jia, and H. Su. Maniskill:
446 Generalizable manipulation skill benchmark with large-scale demonstrations. *arXiv preprint*
447 *arXiv:2107.14483*, 2021.

- 448 P.-Y. Oudeyer and F. Kaplan. What is intrinsic motivation? a typology of computational approaches.
 449 *Frontiers in neurorobotics*, 1:6, 2007. doi: 10.3389/neuro.12.006.2007.
- 450 D. Pathak, D. Gandhi, and A. Gupta. Self-supervised exploration via disagreement. In *International
 451 conference on machine learning*, pages 5062–5071. PMLR, 2019.
- 452 C. R. Qi, H. Su, K. Mo, and L. J. Guibas. Pointnet: Deep learning on point sets for 3d classification
 453 and segmentation, 2017.
- 454 J. Schmidhuber. A possibility for implementing curiosity and boredom in model-building neural
 455 controllers. In *Proc. of the international conference on simulation of adaptive behavior: From
 456 animals to animats*, pages 222–227, 1991.
- 457 J. Schmidhuber. Formal theory of creativity, fun, and intrinsic motivation (1990–2010). *IEEE
 458 Transactions on Autonomous Mental Development*, 2(3):230–247, 2010. ISSN 1943-0604. doi:
 459 10.1109/TAMD.2010.2056368.
- 460 J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization
 461 algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 462 R. Sekar, O. Rybkin, K. Daniilidis, P. Abbeel, D. Hafner, and D. Pathak. Planning to explore via self-
 463 supervised world models. In *International Conference on Machine Learning*, pages 8583–8592.
 464 PMLR, 2020.
- 465 A. Stout and A. G. Barto. Competence progress intrinsic motivation. In *2010 IEEE 9th International
 466 Conference on Development and Learning (ICDL 2010)*, pages 257–262, Piscataway, NJ, 2010.
 467 IEEE. ISBN 978-1-4244-6900-0. doi: 10.1109/DEVLRN.2010.5578835.
- 468 A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polo-
 469 sukhin. Attention is all you need. *CoRR*, abs/1706.03762, 2017. URL <http://arxiv.org/abs/1706.03762>.

471 A Architecture Details

472 A.1 Including Start States into the Challenge Description

- 473 Usually, a precise formulation of a challenge should also at least implicitly include a starting region
 474 or more formally a distribution of initial states. If we would challenge P to climb a hill along a
 475 difficult path but allow him to spawn in all regions around the hill, the content of the challenge might
 476 eventually strongly shift depending on where P spawns. Over and above, we can only guarantee that
 477 our reward function has been designed for the specific path we wanted to guide P over. Perhaps, it
 478 will even provide misleading signals if P approaches g from other paths. If P therefore commits to a
 479 challenge, we must ensure to adjust the initial state distribution to the current challenge specification.
- 480 Unfortunately, not all simulation environments allow this manipulation of the initial state, as was
 481 the case for ManiSkill2. This provided us another reason to formulate a restricted alternative for the
 482 scope of this project.

483 **B Implementation Details**

484 **B.1 Dreamer**

485 To begin with, we explain the Dreamer pipeline based on the parameters in the original implemen-
486 tation. Afterwards, we will discuss and reason our own adjustments.

487 The original Dreamer implementation received images of input size 64×64 . To the end of planning
488 in a latent space of the high-dimensional image observations, the fundamental objective of Dreamer
489 is to obtain compact and rich low-dimensional representations of the observations that facilitate the
490 understanding of the environment dynamics with minimal loss in information.

491 Firstly, a four-layer CNN processes the raw observation into an embedding vector of a size of 2^5 ·
492 depth, where depth=32 initially. From this more compact representation, a forward dynamics
493 model obtains a hidden latent state, which comprises both a deterministic part and a stochastic part,
494 where the deterministic part occupies 200 dimensions and the stochastic part 50 dimensions.

495 This forward dynamics model comprises a one layer GRU cell of hidden state size 200 and is
496 wrapped by two feed-forward fully connected layers. In case of the prior latent state, these net-
497 works map the stochastic part of the previous latent state to the hidden state size and vice-versa. For
498 the posterior stochastic latent state that is also conditioned on the embedding of the actual observa-
499 tion, the GRU is directly bypassed and replaced by another fully-connected feed-forward layer that
500 maps the embedding and deterministic part of the prior latent state to the stochastic posterior of the
501 latent state vector.

502 To ensure that this latent state vector preserves as much information from the environment as pos-
503 sible, the World Model architecture culminates in two heads, which predict the image observation
504 and the perceived reward associated with the observation respectively. While the reward head is im-
505 plemented as a two-layer MLP of hidden layer dimension 400, the head that predicts the raw image
506 observation is a Decoder CNN that is supposed to complement the original dimension. The decon-
507 volution layers are preceded by a fully-connected layer to again bridge the gap between different
508 dimensionalities in the latent space.

509 Availing to this world model, one task-specific agent and one exploration agent are trained to op-
510 timize their value functions and actions inside the latent space learned by the world model. The
511 exploration agent is an implementation of Plan2Explore, which we will focus on later. The task-
512 specific agent first imagines a trajectory of length $T = 15$ inside the latent space, starting from an
513 initial latent state representation provided by the Dreamer unit. It then uses the reward head from
514 Dreamer to generate rewards for these latent transitions and update its value and policy network
515 with them. The stability of these parameter updates strongly relies on a world model that already
516 captured meaningful representations.

517 **B.2 Plan2Explore**

518 Because the purpose of Dreamer in our setting is to explore the environment and hopefully find
519 unseen regions that yield enriching challenges and experiences to the main agent, the exploration
520 agent – albeit small – is the most important final unit. Plan2Explore avails to the world model
521 learned by dreamer to estimate the expected novelty of states that yet lie ahead multiple steps.
522 Mathematically, it models uncertainty as the disagreement between a predictor ensemble that shall
523 predict a certain target value from the current latent state and action. This target value can either be
524 the next latent state, embedding, or more specific parts of the latent state. While the original design

525 choice favoured the embedding as target value, we first decided to target the entire next latent state.
526 Another implementation focuses only on predicting the stochastic part of the next latent state, but
527 this seemed incomplete to us since the transition dynamics for the deterministic latent state matter
528 as well in the estimation of uncertainty.

529 The predictor heads were simple MLPs whose output was interpreted as a distribution. Uncertainty
530 was consequently measured as the average standard deviation between the mode of this distribution
531 across all dimensions of the target vector. This renders aleatoric uncertainty irrelevant to the final
532 uncertainty measure since the predictors should in expectation all converge to the similar mean for
533 states that innately exhibit variance.

534 For the following argument, it is worth reporting that the final average standard deviation term took
535 values around 0.1 and slowly decreased to 0. In the original implementation, this average standard
536 deviation was projected by the log function to obtain the final uncertainty scalar. Effectively, this
537 implies that a state where the predictor ensemble achieves a disagreement measure of 0.001 receives
538 a negative term with a far larger magnitude than a state where the disagreement is "only" about
539 0.01. However, both of these values are extremely low and represent states that the agent already
540 sufficiently explored. Our metric does not need to be sensitive in this value region. More importantly,
541 the metric should filter the most *unknown* states and be more sensitive between higher values close
542 to 1.0. For this reason, we decided to replace the log term by an exp term and beforehand scaled
543 the disagreement by a value that dynamically adjusts itself to the overall progress in uncertainty to
544 remain in the optimal value range where exp provides the desired exponential increase for higher
545 disagreement. Specifically, we consider a distribution of disagreement where the highest quartile
546 has a value around 1.0 as good.

547 Because it is reasonable to expect large deviations in the disagreement value range throughout the
548 training process, we decided to rescale the disagreement values with a scale factor that dynamically
549 adjusts over time. To that end, we formulate a tolerance interval that the disagreement values should
550 in average not exceed or undercut. Otherwise, the scale factor is multiplied or divided by a small
551 factor of around 1.5. We considered both the mean and the max of the current batch of disagreement
552 values and formulated these tolerance intervals over both statistics. This dynamic, loose adjustment
553 is later introduced more generally, since we followed the same principle to control the coefficients
554 in our reward term for artificial goals. To be fair, one could also formulate such a target value and
555 tolerance intervals for the log metric and ensure that the best values still remain positive. However,
556 deviations from this target value would incur far more strong oscillations in the exploration terms. In
557 experiments, we still observed situations where almost all exploration terms were negative because
558 the average scaled predictor disagreement slightly decreased below 1.0. This eventually discour-
559 aged exploration of the main agent since the uncertainty term was an essential component for the
560 assessment of trajectories as possible imaginary goals.

561 We believe that the exp term provides a more stable behaviour for increasingly small disagreements
562 as it maps them all to values slightly greater than 1.0 and is quite insensitive to small deviations
563 in these small values. We could hence safely subtract 1.0 from the final exp term to lower the
564 uncertainty metrics for these undesirable, sufficiently explored states down to values quite close to
565 0, while unseen states that achieved higher rewards still obtained values greater than 1.0.

566 B.3 Controllers

567 The default joint_pd_delta_pos controller tracks only the position of the robot's joints, but not its
568 rotations. On the other hand, the pd_ee_delta_pose controller also includes the rotation matrix and

569 focuses on the end-effector frame. This proves useful for many of the pick-and-place tasks that we
570 deal with in this report. For this reason, we have observed a tremendous performance improvement
571 when switching from `joint_pd_delta_pos` to `pd_ee_delta_pose`. This was also reported by the
572 authors. The `pd_ee_target_delta_pose` controller that applies the ideal transformation matrix that
573 the current action would lead to does not achieve a similarly good performance and was hence not
574 selected for our final model.

575 Unfortunately, the authors realised all controllers as PD controllers and not the more performant PID
576 controller. Moreover, the stiffness and damping coefficients in the PD formula could not be adjusted.
577 Since the PD controllers exhibited strongly different performance, it is likely that the coefficients in
578 the PD controller error term could profit from more fine-tuning.

579 **B.4 Decay Models for Challenge Exploration**

580 **B.5 Challenge Choice**

581 **B.5.1 Challenge Quality Metric**

582 After having collected new trajectories through the interplay of our task-related main agent P and
583 the exploration-oriented Dreamer C , we inspect the new trajectories to find state regions where we
584 expect enriching experiences for P . Mainly two factors contribute to this *challenge quality* of a state
585 s_t .

586 Firstly, the task-specific, extrinsic reward r_t is a good indicator of how useful a certain state might
587 prove for eventually mastering the environment task. Secondly, the epistemic uncertainty measured
588 by Dreamer's Plan2Explore extension signifies novel regions that could aid P to abstract from its
589 limited experience and eventually push forward to unseen regions in the optimization landscape. At
590 the same time, both terms should not dominate the quality score of a state too much.

591 On the one hand, the agent might lose sight of the actual task if we considered only the epistemic
592 uncertainty. Plan2Explore was originally intended to obtain a task-agnostic, broad distribution of
593 states, which is only of limited benefit to our agent since some of them will yield little benefit for our
594 current task. This trade-off between exploration and goal orientation reminds of the RRT planning
595 algorithm, where we sample the states not purely uniform, but bias towards sampling the goal state
596 with a certain probability to ensure that the final tree is expanding towards the task-specific goal state.
597 On the other hand, too little exploration might render the Dreamer agent superfluous and narrow
598 down the behaviour of the agent to its selective experience, which reflects the usual combination of
599 exploration and exploitation in the earliest reinforcement learning algorithms.

600 For this reason, we decided to target a ratio of the extrinsic reward over the exploration reward of
601 $c^* = 0.75$ to slightly favour exploration and hence target

$$S(s_t, r_t, u_t) = c \cdot r_t + u_t, \text{ where } c \cdot r_t \approx c^* \cdot u_t, \quad c^* = 0.75.$$

602 We do not need another coefficient to scale the entire term because we only compare these terms
603 *relatively* to each other. Since the Dreamer agent already computes the KL divergence between the
604 prior and posterior belief states conditioned on the embedding of the current image, we included this
605 value with a ratio of $c_1^* = 0.1$ to the ensemble predictor disagreement into the uncertainty term.

$$u_t = c_1 \cdot \text{KL-LOSS} + \text{DISAG}, \text{ where } c_1 \cdot \text{KLLOSS} \approx c_1^* \cdot \text{DISAG}, \quad c_1^* = 0.1.$$

606 Because the latent state was also split into a deterministic part and stochastic distribution, we felt
 607 that it would have the same ability as the predictor disagreement DISAG to mitigate the influence
 608 of aleatoric uncertainty by steadily converging to the same mode in the distribution that minimizes
 609 the mean squared error. But we still remained conservative and chose the impact of this KL term
 610 to be 10 times smaller than the disagreement. Just as mentioned before when discussing the scale
 611 of the disagreement value, we ensured that the coefficient of the KL value and the coefficient of the
 612 extrinsic reward remained in a tolerance interval around the aforementioned target values. Again,
 613 we will discuss our implementation of such a dynamic scale adjustment in more detail in Appendix
 614 [B.7](#).

615 **B.5.2 Challenge Filtering with Self-Balancing Strictness**

616 After each state was assigned a goal quality score S , we next had to pick goal states out of the batch
 617 of trajectories whose metric seemed promising to consider this state as an artificial goal. To avoid
 618 possible error sources, we refrained from overcomplicating the selection strategy at first and defined
 619 the state $g := \arg \max_t S(s_t, r_t, u_t)$ with the highest goal quality metric in each trajectory as a
 620 *challenge candidate*.

621 We then ensured that these candidates met five criteria $F_i : \tau \mapsto \{0, 1\}$, $i = 1, \dots, 5$ that related
 622 either to the sole candidate state or the entire trajectory from which the candidate stemmed. While
 623 the last two criteria were hard criteria that each challenge must meet, the challenge candidate must
 624 only satisfy the *majority* of the first three ones.

625 To validate the criteria for the entire trajectory, we computed its the mean exploration reward \bar{u} and
 626 the mean extrinsic reward \bar{r} .

627 Firstly, the exploration benefit of the entire trajectory should be sufficiently high, else the goal state
 628 was mainly selected because of its high extrinsic reward. In such a case, we believe that the trajectory
 629 will already enrich the agent in the original replay buffer. If \bar{u} is negligibly small compared to \bar{r} , we
 630 hence rejected the goal candidate. Here, negligible refers to falling short of \bar{r} by a tolerance factor
 631 ω_1 that was initialized with 2.0. To reasonably compare the magnitudes of \bar{r} and \bar{u} , we avail to the
 632 fact that after some iterations, the coefficient c in the goal quality term S of the last section [B.5.1](#) will
 633 have been sufficiently adjusted to correctly track the target ratio c^* between r_t and u_t . Expressed
 634 formally, it should hold that $c \cdot r_t \approx c^* \cdot u_t$, therefore we have

$$F_1(\tau, g) := \begin{cases} 0, & \omega_1 \cdot c \cdot \bar{r} > c^* \cdot \bar{u} \\ 1, & \omega_1 \cdot c \cdot \bar{r} \leq c^* \cdot \bar{u} \end{cases}.$$

635 However, this initial value $\omega_1 = 2.0$ is heuristically set as many other hyperparameters are and
 636 experimental evaluation proved that this value mostly achieved a good balance, as a static value it
 637 could not account for the inevitable fluctuations in \bar{r} and \bar{u} . We hence implemented a mechanism
 638 that we call *self-balancing strictness*. The intuition behind this mechanism is the observation that
 639 humans also adjust their strictness in judgment contextually. We don't blame a small child for not
 640 walking stably and crawling over the ground. But as it grows, so grow our expectations on which
 641 we base our judgment. Another more related example are relative grading principles which are
 642 practice in some schools or universities to relate the performance judgment to the performance of
 643 an entire population. In our case, the increasing certainty about the environment and the increasing
 644 reward can similarly be regarded as such a developmental progress to which we want to adjust our
 645 judgment.

646 Even in the static case where the assessed population does not experience any change in performance,
 647 such a dynamic judgment adjustment also proves useful for another problem. Namely, we
 648 want to avoid that we reject all goals, which will effectively leave the replay buffer of imaginary
 649 goals empty. On the other side, accepting each goal candidate will not filter good trajectories from
 650 worse trajectories, of which there certainly will be many.

651 We hence employ a separate tolerance ω_i for each filter F_i . Every time a goal candidate g was
 652 rejected because $F_i(\tau, g) = 0$, its tolerance value ω_i is adjusted by a slight factor $\gamma_F \approx 1$ such that
 653 F_i becomes a slightly more loose criterion in the future. The opposite adjustment is taken when a
 654 goal was accepted by this criteria. Because this causes an exponential growth or decay for extreme
 655 initial biases towards rejection or acceptance, the value will quickly stabilise itself and moreover
 656 dynamically adjust to the quality of the candidate states.

657 Symmetrically, one could check whether \bar{r} is negligibly small compared to \bar{u} . However, we refrain
 658 from rejecting candidate states based on this criterion. The reason is that the buffer for the imaginary
 659 goals is mainly dedicated to fostering exploration and should not punish states because of their nov-
 660 elty. However, we already stressed our concerns for diverging too exhaustively into task-irrelevant
 661 areas and hence formulated another criteria related to the above.

662 Namely, we ensure that the extrinsic reward of a goal trajectory is not too low in comparison to the
 663 total average reward previously achieved by the agent in a certain time window.

664 This first requires measuring this total average reward. We decided to extend our agent by a scalar
 665 variable that we call the *mean reward memory* r_{mem} . Every time the agent was called with new tra-
 666 jectories, the agent computed the mean reward over all these trajectories as r_{new} . Then we updated
 667 the reward memory in a bilinear fashion as

$$r_{mem} = \gamma_r \cdot r_{mem} + (1 - \gamma) \cdot r_{new}.$$

668 Here, γ_r is a reward memory decay, which we set to 0.99. This means that rewards that were
 669 perceived 100 steps ago become roughly three times less important than the current rewards. We
 670 were worried that smaller values might entail a too short-sighted memory that is sensitive to short-
 671 term oscillations in the perceived reward and oriented ourselves on the discount factor γ that was
 672 also set to 0.99.

673 With this mean reward memory in hand, we formulated the second criteria in an analogous fashion
 674 to the first criteria, but this time compared the average reward of the entire trajectory with the mean
 675 reward in the memory.

$$F_2(\tau, g) := \begin{cases} 0, & \omega_2 \cdot \bar{r} < r_{mem} \\ 1, & \omega_2 \cdot \bar{r} \geq r_{mem} \end{cases}.$$

676 We ensured that F_2 tolerates short-term deviations in the overall reward of an agent which we ex-
 677 pected to occur by nature of the commonly unstable optimization procedures for hard optimization
 678 landscapes. After experimental evaluations, we culminated in a bound with a comparably larger
 679 value $\omega_2 = 8.0$ than ω_1 .

680 From now on, we focus on the challenge quality score of the candidate goal state g . Our third
 681 criterion ensures that there should be some development on the way to the challenge goal g . By def-
 682 inition, g has the highest quality score on the entire trajectory. But if this improvement is relatively
 683 small compared to the average scores on τ , then there will unlikely be any advantage of reaching
 684 g and lifting g to a “goal” state is indeed questionable. We hence verified whether the goal quality
 685 score of the candidate state was at least $\omega_3 := 2.0$ times higher than the average goal quality score.

686 Of course, this value was dynamically adjusted in the same fashion as the prior values. In the ex-
 687 periments, we observed that this criteria turned out to be one of the strictest, since many trajectories
 688 did not include any "highlight". We later experimentally evaluate the influence of these filters in
 689 Appendix C.2.2 However, the majority of proposed candidates finally passed all these three criteria.
 690 Finally, the two last criteria seal some edge cases we encountered during debugging. The fourth
 691 criterion F_4 rejects candidates with negative scores, which happened in early tests where the log
 692 implementation of the predictor ensemble disagreement was not yet replaced by our exp metric.
 693 With the exp metric at hand, negative rewards should be impossible, but we left it as an assertion.
 694 The fifth criterion F_5 rejected candidates that were the very first states in a trajectory. Because we
 695 decided to truncate trajectories after their goal state, this resulted in erratic behaviour, which we
 696 avoided this way.
 697 We further have to note that we removed the first five steps of each trajectory, because the initial
 698 prior latent state of each trajectory consisted only of 0s. Because the initial state in the environment
 699 is mostly random, it is impossible to correctly predict the next latent state even if the same action
 700 is given. Therefore, the predictor ensembles always had their highest disagreement in the first state
 701 of each trajectory, which is a correct behaviour, but undesired for our case. After three steps, the
 702 predictions were fully grounded on the actual observations and hence pertinent for our challenge
 703 discovery.
 704 If a candidate (τ, g) passed all five checks, the trajectory τ was truncated after this state. To integrate
 705 them as imaginary goal experiences into the buffer, we yet have to define the artificial reward term,
 706 which we discuss in the following section.

707 B.6 Progress-sensitive Artificial Rewards

708 Similar considerations as for the challenge quality metric in Appendix B.5.1 had to be made when
 709 defining how the reward term for artificial goals should be constituted.
 710 If we neglect the dense extrinsic reward term r_t too strongly, then the reward experience in the
 711 imaginary trajectories will directly disagree with the reward signal in the original trajectories and
 712 hence unnecessarily disturb the agent. On the contrary, we want to push the agent to explore unseen
 713 regions, which again requires a significant weight of the exploration term u_t already used above for
 714 the goal quality score.
 715 Moreover, both of these terms do not actively raise the reward towards the specific imaginary goal
 716 in the trajectory. To enlarge the emphasis on the goal state g to eventually curl the agent to the goal
 717 region, we included two more terms that reflect a *distance metric* to g .
 718 We refer to the first term as the *embedding distance* d_{emb} , which is literally the *improvement* in
 719 the embedding distance to the goal embedding, where we make use of Dreamer's image encoder.
 720 Because the trajectory is a result of exploration, we cannot expect it to lead straight to the goal. In
 721 the lecture on offline RL, Professor Xu introduced a similar problem, where the agent has to learn
 722 to stitch good local trajectories together to efficiently approach the goal state without unnecessary
 723 detours.
 724 By this assumption, we would punish detours in the trajectory and encourage the agent to directly
 725 approach the goal. However, it is foolish to directly take the embedding, because the initial states
 726 that are far away but directly lead to the goal still have a large distance while actions that detour
 727 from the goal state will still receive good scores if the state they originate from is close to the goal

728 state. What matters more is the *direction* in which the agent proceeds. For this reason, the *change* in
 729 the embedding distance seems a more useful metric to assess how an action roughly contributes to
 730 reaching the goal. In this way, actions that moved away from the goal – or at least from its embedding
 731 –, received a slightly negative reward. Therefore, the embedding distance d_{emb} is computed as

$$d_{emb,t} := \|e_g - e_t\|_1 - \|e_g - e_{t+1}\|_1,$$

732 where $\|\cdot\|_1$ denotes the Minkowski distance with $p = 1$.

733 Finally, we also included the simple step distance d_{step} to the goal as last heuristic for our arti-
 734 ficial reward term. When s_t is the goal state in a trajectory, then reaching s_i will yield a reward
 735 proportional to i .

$$d_{step,i} := \frac{i}{t}, \quad \text{where } t \text{ is the index of the goal state.}$$

736 Our intention behind this reward was to complement potential inaccuracies in the embedding dis-
 737 tance d_{emb} , which could provide brittle signals especially at the beginning where the encoder net-
 738 work of Dreamer still has not acquired a good embedding.

739 To mitigate the dissonance between the imaginary rewards and the actual rewards, we targeted the
 740 extrinsic reward term to be slightly more important than the exploration reward term, namely by
 741 a factor of 1.5. This is converse to their relation in the goal quality score where the exploration
 742 term was slightly more important, where it is not necessary to care for accordance with the extrinsic
 743 rewards.

744 In this way, we could still lay an emphasis on trajectories that promise a high exploration reward
 745 in the goal evaluation stage without evoking a severe dissonance in the experienced reward that the
 746 main agent receives to learn. However, the reoccurrence of this trajectory in the imaginary goal
 747 buffer will hone the agent’s knowledge of this possibly unfamiliar area.

748 On the other hand, the two distance terms were both targeted to be approximately 0.3 times less
 749 important than the exploration term. This way, the sum of purely fictitious rewards and the extrinsic
 750 reward coarsely balance each other out and still leave a lot of freedom for the artificial reward
 751 component to lure the agent into regions that the extrinsic reward might not highlight.

752 Finally, the scaling so far was relative to the exploration reward and does not ensure that the magni-
 753 tude of rewards are approximately equal to the true original rewards on this trajectory. Practically,
 754 we reformulated the scales to be relative to the *extrinsic* reward term. Because of this change, the
 755 artificial reward term is approximately two times higher than the true extrinsic reward term if all
 756 scales meet their aforementioned target ratios. To slightly emphasize trajectories in the imaginary
 757 buffer, we only scaled the entire reward term by 0.5, leaving a slight possible advantage that fully
 758 depends on the exploration and distance terms.

759 Hence, the final artificial reward term writes as

$$\begin{aligned} r'_t &:= 0.5 \cdot (r_t + c_{expl} \cdot u_t + c_{emb} \cdot d_{emb,t} + c_{step} \cdot d_{step,t}), \quad \text{such that} \\ \bar{r} &\approx c_{expl}^* \cdot c_{expl} \cdot \bar{u}, \quad c_{expl}^* = 1.5, \\ c_{emb}^* \cdot \bar{u} &\approx c_{emb} \cdot \bar{d}_{emb}, \quad c_{emb}^* = 0.3, \\ c_{step}^* \cdot \bar{u} &\approx c_{step} \cdot \bar{d}_{step}, \quad c_{step}^* = 0.3, \end{aligned}$$

760 where \bar{x} again represents the mean of x over the entire trajectory.

761 **B.7 Dynamic Hyperparameter Targeting**

762 In the last paragraphs, we have often mentioned that we dynamically adjust hyperparameters to
763 meet a certain ratio between the scaled value and another comparison value. In view of the highly
764 dynamic learning progress, we are concerned that static hyperparameters will cause scales to drift
765 off if new achievements or novel challenges affect the agents learning experience.

766 To that end, we implemented a mechanism to dynamically adjust the scale coefficient c of a hyper-
767 parameter h to meet a certain target ratio r^* with regard to another comparison value g up to a fixed
768 tolerance ω . Although the target ratio r^* of hyperparameters is also another static hyperparameter,
769 it resides on a more abstract, coarse layer. Fixing these strategic hyperparameters to a fixed value
770 is less problematic than fixing more concrete, operational hyperparameters to a specific value. For
771 each hyperparameter h that was subject to this dynamic adjustment, we defined an adjustment factor
772 λ , which was usually chosen to be between 1.1 and 1.5 depending on how large we expected or
773 desired h or its ratio to g to fluctuate.

774 In each update step of the model, we also check if any dynamic hyperparameter h violated their
775 respective tolerances and adjusted them by λ if necessary. To estimate how well the current h sat-
776 isfies the ratio r^* in a certain time window, we had to obtain a statistic that provided contextually
777 relevant information on the distribution of this value. In most cases, the average value over trajec-
778 tories or batches of trajectories was taken as \bar{h} . In the optimal case, $\bar{h} = r^* \cdot \bar{g}$, that is h should
779 in average meet the static target ratio to the other value g . We again stress that g is non-static but
780 also time-variant. If this target was missed too strongly, h was updated by λ . This allowed the scale
781 to quickly respond to severe violations of the target ratio with exponential growth or decay. In this
782 way, we could ensure that some values are kept in a certain range and neither explode nor diminish.
783 In the experiments in Appendix ??, we demonstrate that this mechanism was indispensable to tackle
784 strong natural fluctuations of external values such as reward, uncertainty, or embedding distance.
785 The following Algorithm 2 specifies this dynamic hyperparameter adjustment.

Algorithm 2 CPM Dreamer: Training process

Input: Hyperparameter statistic \bar{h} , adjustment $\lambda > 1.0$, comparison value statistic \bar{g} ,

target ratio r^* , lower tolerance ω_{low} , upper tolerance ω_{high} .

```
1: if  $\bar{h} < \omega_{low} \cdot r^* \cdot \bar{g}$  then
2:    $\bar{h} = \bar{h} * \lambda$ 
3: end if
4: if  $\bar{h} > \omega_{high} \cdot r^* \cdot \bar{g}$  then
5:    $\bar{h} = \bar{h} / \lambda$ 
6: end if
```

786 **C Experiments**

787 The following experiments were done on the PickCube-v0 environment. They either compare
788 handicapped versions of our CPM (Competence Progress Motivation) Dreamer with a strong, self-
789 standing Dreamer model and PPO models, or experimentally verify the capabilities of our CPM
790 Dreamer.

791 **C.1 Challenge discovery**

792 Figure 5 compares the experienced reward and success during training. We stress that in this early
 793 stage, the exploration rate of our CPM Dreamer models are still high and are plotted in Figure 8a.
 794 For this reason, the generated trajectories are mostly oriented at exploring unseen states with the
 795 Plan2Explore extension and not oriented towards succeeding in the task. Despite that, our models
 796 still achieves seldom success. While the PPO PointCloud model achieves far higher rewards and
 797 also slightly more frequent successes, it does not achieve any success during evaluation, while our
 798 model also sometimes achieves success there (Fig. 4b).

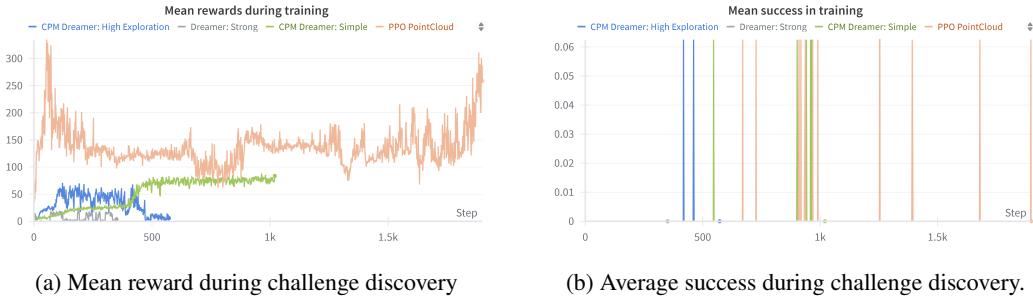


Figure 5: During training, the agent $A = (C, P)$ uses both the Plan2Explore extension of Dreamer C and the SAC agent P to act ε -greedily, where ε is initialized with a high value and decays over time. The resulting behaviour is therefore mainly oriented at exploring unseen states and *not* at trying to master the task. Despite of that, A sometimes even succeeds in the task. During training, the PPO PointCloud model achieves higher success rates and rewards than our CPM Dreamer models, the PPO PointCloud does not achieve any success during evaluation (Fig. 4).

799 **C.2 Challenge choice**

800 **C.2.1 Challenge quality metric**

801 The extrinsic reward and the uncertainty that both constitute the challenge quality metric are subject
 802 to high fluctuation (Figures 6a, 6b). Accordingly, the coefficient for the extrinsic reward quickly
 803 adapts to these deviations and in this way manages to balance the influence of both terms to the final
 804 challenge quality metric S (Fig. 6c).

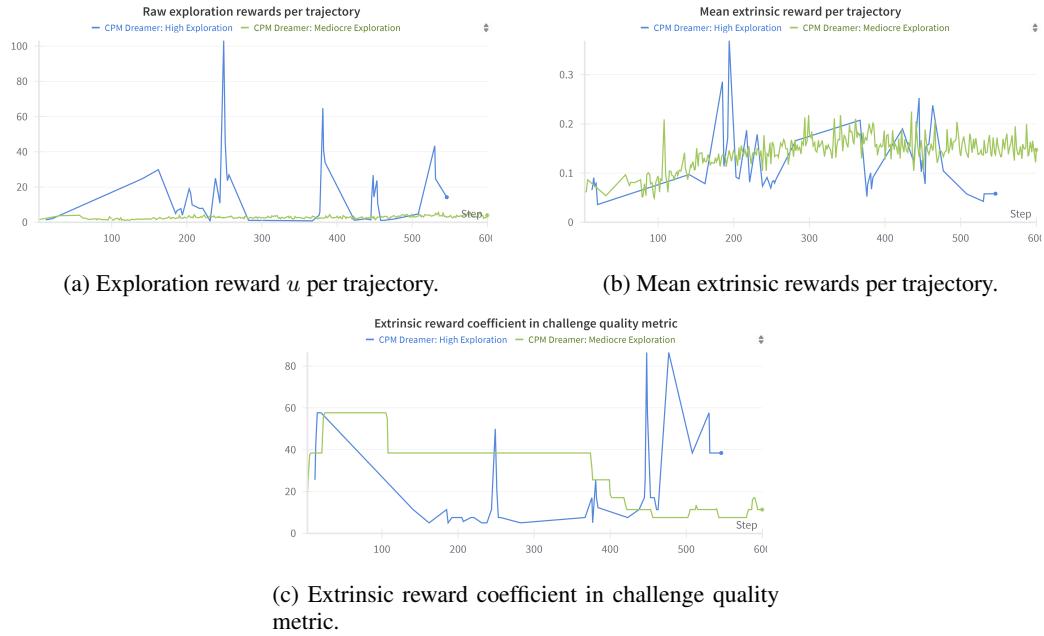


Figure 6: Overview of the exploration and extrinsic rewards and the coefficient that dynamically adjusts to their fluctuating magnitude to satisfy the target ratio.

805 C.2.2 Challenge filters

806 As Figure 7 illustrates for all three tolerances that they dynamically adjust their strictness to the
807 amount of accepted challenge trajectory.

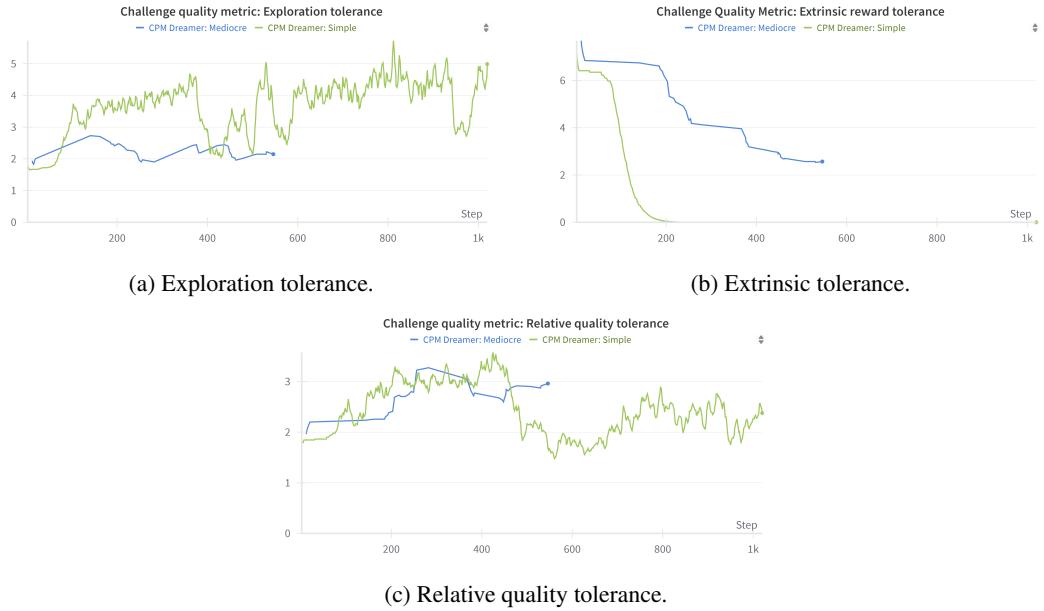


Figure 7: Our self-balancing strictness in challenge filtering dynamically adjusts to the prior acceptance and rejection behaviour.

808 **C.3 Artificial rewards**

809 Similar to Figure 6, this experiment again underscores the ability of our dynamic hyperparameter
 810 adjustment mechanism to account for strong fluctuations in time-variant values. In this case, the
 811 exploration reward experiences severe peaks (Fig. 8b). The coefficients for the extrinsic reward (Fig.
 812 8c) and the step distance (Fig. 8d) both mimic this increase to similarly amplify the influence of their
 813 terms, but also quickly decrease after the exploration rewards go down again. This demonstrates that
 814 our reward term remains robust against different value scales.

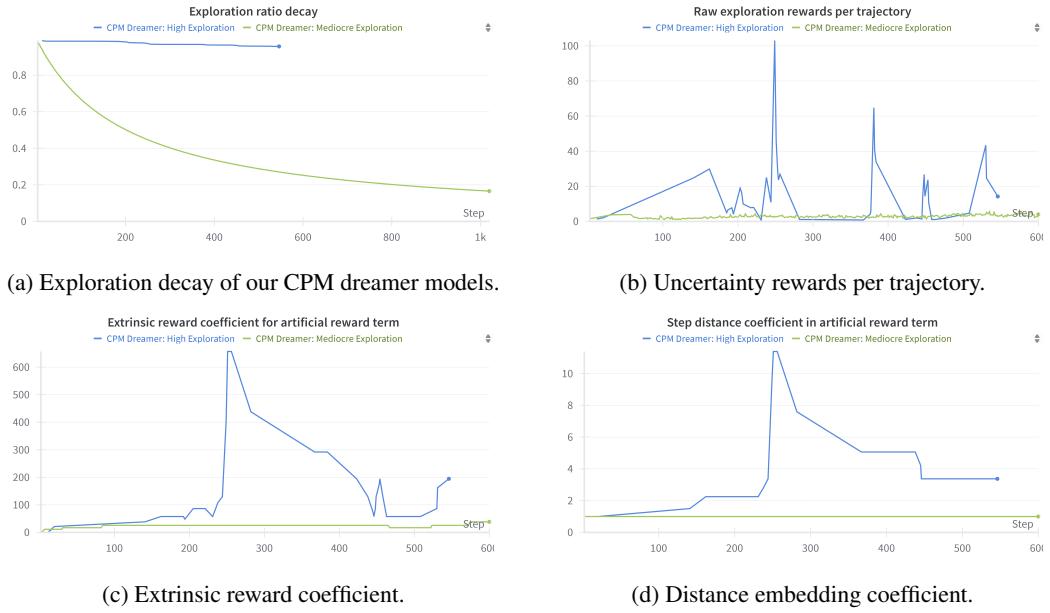


Figure 8: Overview of our dynamic coefficients development in view of a highly fluctuating uncertainty during challenge discovery. Owing to the dynamic adjustment, the other terms can adjust to the fluctuating level of the exploration reward term (uncertainty). We showcase these adjustments for two models with highly different exploration decays.

815 **C.4 Dreamer**

816 **C.4.1 Plan2Explore**

817 The following Figure 9 is supposed to visualize how the pre-exponentiation scale (Fig. 9b) adjusts
 818 itself to the current range of the predictor ensemble loss (Fig. 9a). This scale hence adjusts the range
 819 to ensure that the most unfamiliar states are significantly highlighted in comparison to more known
 820 ones. As expected, the strong Dreamer has the largest prediction capabilities and hence the lowest
 821 possible loss. The pre-exponentiation scales and hence lowered to sensitize the filter. Because we
 822 applied a loose upper threshold, this scale remains stably low even when the predictor ensemble face
 823 a new unknown region which severely aggravates their loss.

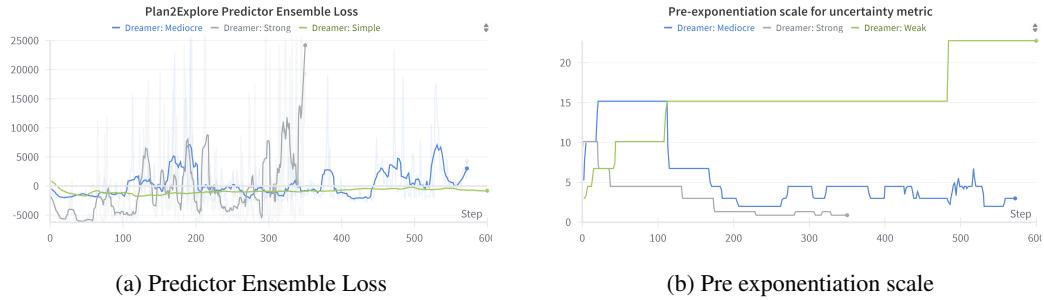


Figure 9: Juxtaposition of Plan2Explore’s predictor ensemble loss with the associated adjustment of the dynamic scale coefficient before the uncertainty terms u_t are fed into \exp .

824 C.4.2 Model Loss

- 825 The Dreamer model is indeed capable of learning meaningful latent representations of the observa-
 826 tions, as the following experiment results summarized in Figure 10 demonstrate.
 827 On the one hand, the image head loss converges at roughly 5.7 and achieves a low absolute aver-
 828 age pixel prediction error of less than 0.1. It is important to note that the pixel values have been
 829 normalized to $[-0.5, 0.5]$. Interestingly, the loss for the simple Dreamer and the strong Dreamer is
 830 similar, although the strong Dreamer achieves far better reconstruction results (Fig. 10b). We can
 831 not explain these results, as we would have expected the loss to be more lower for the strong model.
 832 Similarly, the reward model reduces its absolute and relative prediction error (Fig. 10f) proportional
 833 to the reward loss (Fig. 10e). To much exploration seems to prohibit the model from truly grasping
 834 the nature of the rewards, as it explores regions that are too unrelated to the task.

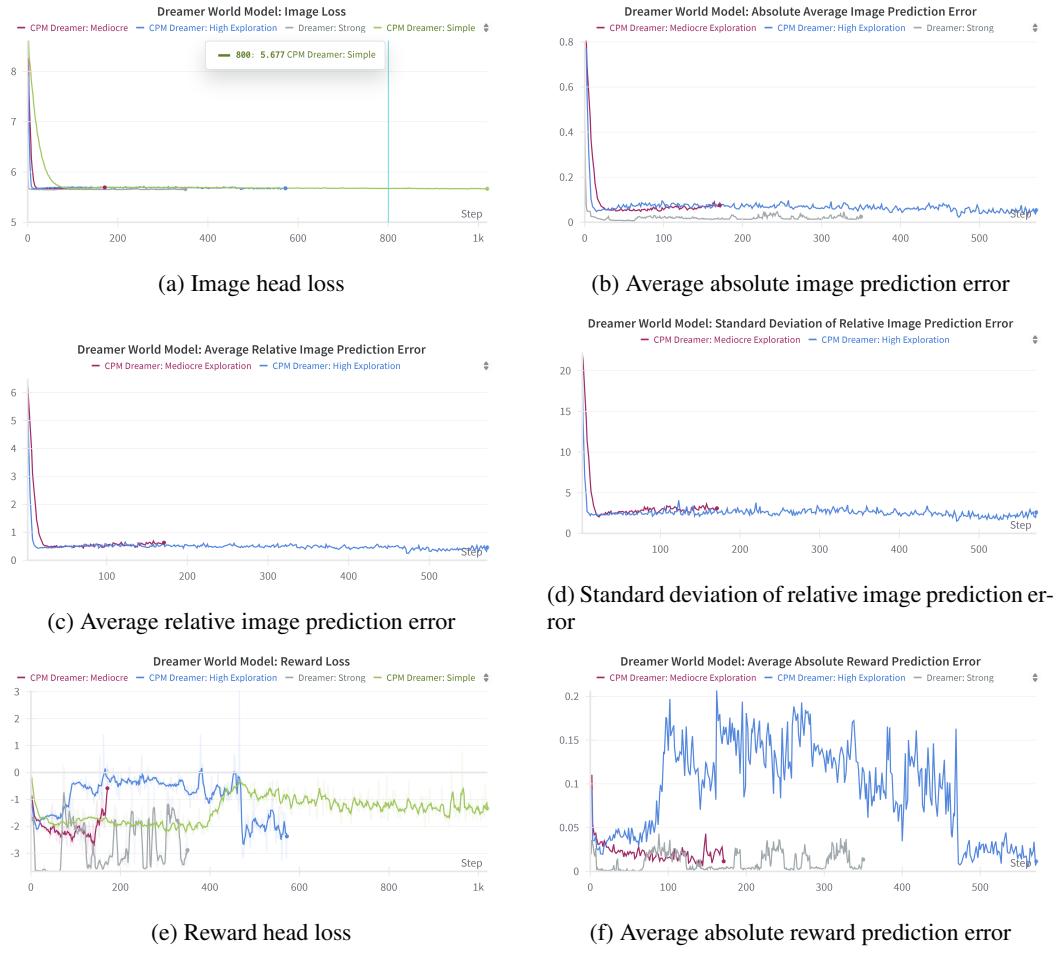


Figure 10: Dreamer’s model loss reasonably coincides with accurate prediction of both the image observation and the reward. The image decoder and reward head are trained to predict the image observation and the perceived reward from the hidden latent state of the world model.

835 C.5 SAC Agent

836 Finally, we juxtapose the loss and Q values of the SAC agent for both the original and the imaginary
 837 reward in Figure 11. Owing to the similar scale of the new reward term, the Q values and loss values
 838 for both buffers are roughly in the same range for both buffers and do not exhibit large dissonance.

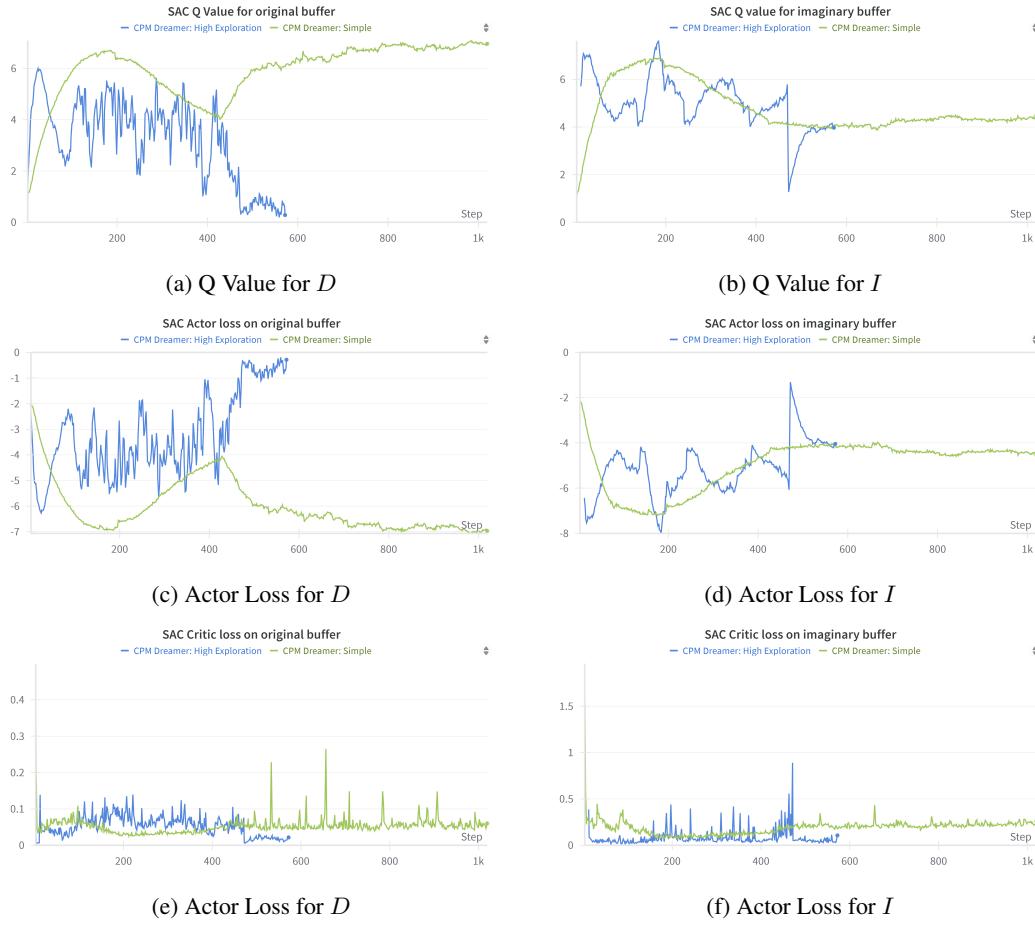


Figure 11: Juxtaposition of the SAC agent’s Q values, actor, and critic loss for both the original buffer D and the imaginary buffer I .