

# ALAN: Autonomously Exploring Robotic Agents in the Real World

**Anonymous Author(s)**

Affiliation

Address

email

1       **Abstract:** In order to build robotic agents that can autonomously operate in the  
2       real world, it is crucial to explore the environment. While it is possible to build  
3       agents that can learn without supervision, current methods struggle to scale to the  
4       real world. Thus, we propose, ALAN, an autonomously exploring robotic agent,  
5       that can perform many tasks in the real world with little training and interaction  
6       time. Our approach builds a shared model of the world, across all possible tasks,  
7       and continuously explores. We propose a novel intrinsic motivation reward that  
8       leverages both self-supervised agent-centric prediction as well environment-centric  
9       priors for visual change. We show our approach on a real world play kitchen setting,  
10      performing multiple manipulation tasks and showing much better exploratory  
11      performance than state-of-the-art methods. Videos can be found at [https://  
12      robo-explorer.github.io/](https://robo-explorer.github.io/)

13      **Keywords:** Exploration, World Models, Learning Control

## 14      1 Introduction

15      To gain true autonomy, robots will need to perform a diverse range of tasks in the real world. Due  
16      to the challenges of dealing with uncertainty, deep learning has emerged as a promising approach  
17      [1, 2] for robotics. A critical challenge for scaling learning based approaches to more complex  
18      settings is the task specification problem. Prior works require heavy reward engineering or human  
19      demonstrations, which is not feasible for performing large numbers of tasks. This also requires  
20      knowledge of the environment, which might be hard to obtain for every domain. Instead, if robots  
21      can collect their own data using task-agnostic objectives, then they could explore their environments  
22      and learn interesting and useful skills.

23      In the absence of explicit task definitions, the agent should have an efficient way to use all its collected  
24      experience for learning. World models [3, 4] provide a means of learning an effective low dimensional  
25      representation of observed images from diverse domains, thus enabling reuse of data for learning  
26      useful behaviors. Furthermore, if there are certain states where prediction for the world model is  
27      difficult, then it likely needs more data for the corresponding part of the environment. This gives  
28      rise to a natural intrinsic objective of maximizing model uncertainty [5, 6]. While this does lead to  
29      the discovery of interesting behavior, there has been difficulty in scaling such approaches to real  
30      world settings since collecting samples is very expensive. We ask if there is a different task-agnostic  
31      objective that can enable robots to *more efficiently explore*?

32      In order to address the above question, we present, ALAN, an efficient autonomous real robotic  
33      explorer. Our key insight is that interesting behavior for robots in the manipulation setting mostly  
34      involve interactions with objects, which cause changes in the visual features of the observations.  
35      Thus, seeking to maximize the change in these visual features can be a useful objective for robots  
36      to optimize. We note that this is an ‘environment-centric’ reward in contrast to maximizing model  
37      uncertainty which is more ‘agent-centric’, since it is concerned with controlling predictions in its own  
38      mental model. Furthermore, the environment change is a task-agnostic objective, and will generally  
39      guide the robot towards more contact-rich behavior, which is essential for learning manipulation  
40      tasks. Using the world model, the agent can learn to predict what the environment change will be in

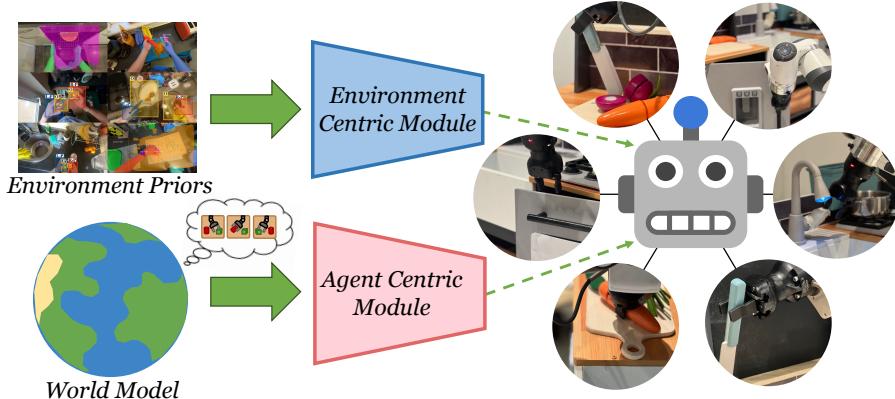


Figure 1: We propose Autonomous Learning Agents (ALAN) that can enable robots to collect rich data from their environment efficiently. The agent autonomously identifies workspace of interest, and learns a joint world model on all data. ALAN follows two intrinsically defined objectives - taking actions that cause difficulty for prediction (the agent centric module), and actions driven by environmental priors (environment centric module).

41 the low-dimensional representation space, enabling these agents to explicitly optimize for this metric  
 42 when planning, leading to trajectories that interact with objects.

43 Our contributions can be summarized by the following. We present ALAN, an efficient real world  
 44 exploration algorithm with a novel intrinsic motivation reward that extracts better signal from the  
 45 environment by learning to predict which transitions cause large visual change. Our approach shows  
 46 strong exploratory performance and goal reaching on several manipulation tasks in a real world  
 47 toy kitchen, that outperform SOTA baselines. We also show that the data is useful for reaching  
 48 human-specified goals by training distilling the collected experience into a goal-conditioned policy.

## 49 2 Related Work

50 **Exploration** In reinforcement learning (RL), exploration has been studied in various different  
 51 contexts ranging from highly tabular to complex continuous spaces. For discrete settings, the  
 52 concept of exploration can be analytically defined as state visitation counts [7] or as probability  
 53 distributions over visited states [8, 9]. For complex states, such as images, previous works have used  
 54 neural networks to approximate state counts [10]. A common way to describe an intrinsic reward  
 55 for exploration is to use either the error [11] or uncertainty [12, 13] in prediction about how the  
 56 environment and agent would interact. Pathak et al. [14] propose a differentiable intrinsic reward  
 57 which measures the disagreement between an ensemble of models via the variance of the set. Sekar  
 58 et al. [5] leverage a disagreement intrinsic reward, but explore in an imagined space via a world model  
 59 [15, 4, 16]. These methods also do not scale to real world settings. Prior work has attempted to have  
 60 the agent provide its own supervision in the real world [17, 18, 19], however, these approaches only  
 61 trained for specific tasks and cannot share data between each one. Most of these approaches [20, 2]  
 62 suffer from sample efficiency. Our approach leverages both the agent’s own supervision, is able to  
 63 share data across manipulation tasks and leverages an initialization from human behavior priors.

64 **Self-Supervised Goal Reaching** Collecting good data and learning skills is important, however, an  
 65 agent needs to be able to distill this experience into a policy that is able to reach some form of goals.  
 66 Traditionally, prior work has learnt a policy that directly practices on multiple tasks, and conditions  
 67 on some provided goal. This goal can come from supervision [21, 22] or can come from the agent’s  
 68 own experience [23]. Training an agent in the real world requires operating from raw sensory data,  
 69 and being able to achieve goals specified by images only. Prior work has used a mix of visual  
 70 representation such as contrastive [24] or generative models [25, 17]. While visual representations  
 71 can provide some signal, using simple distance metric does not enable exploration of the environment.  
 72 Pong et al. [25] and Zhang et al. [26] employ goal proposal mechanism trained on previously seen  
 73 states. Proposing goals close to previously achieved ones tends to keep the agent stuck performing  
 74 the same types of task, and does not allow for learning of complex behaviors. Mendonca et al. [6]  
 75 propose an exploration approach that is able to leverage an imagined world model to learn to achieve

76 sampled goals. One of our key insights is to utilize a world model and practice goal reaching in  
 77 imagination, similarly to [6]. However, it is difficult to know *what* to focus on while exploring. For  
 78 this, we use offline priors from human behavior data from the internet.

### 79 3 Background

80 **Model-Based RL and Planning** Model Based RL operates in a Markov Decision Process (MDP).  
 81 An MDP is defined by a set of states  $\mathcal{S}$ , actions  $\mathcal{A}$ , transition probabilities between states conditioned  
 82 on actions,  $\mathcal{T}(s_{t+1}|s_t, a_t)$ , a initial state distribution  $\mathcal{S}_0$ , a reward function  $\mathcal{R}(s_t, a_t)$ . Specifically, the  
 83 goal of a model based RL algorithm, is to learn a function  $f_\phi(s_{t+1}|s_t, a_t)$  which best approximates  
 84 the the true transition dynamics  $\mathcal{T}$ . While planning, the Cross-Entropy Method (CEM) is used to find  
 85 the best set of actions  $a_{1:T}$ , which produce the highest reward under the trained dynamics model  $f_\phi$ .

86 **Sample Efficient Model-Based RL** In order to do as efficient model-driven learning as possible, it  
 87 is desirable to build a general *dynamics model* of the world. When dealing with high-dimensional  
 88 inputs such as images, it is often beneficial to train the dynamics of the system in a lower dimensional  
 89 space, to avoid overfitting to visual artefacts. Thus, we employ the Recurrent State-Space Model  
 90 (RSSM) from Hafner et al. [4]. Given an image input at timestep  $t$ :  $x_t$ , we obtain a low dimensional  
 91 embedding  $e_t = E_\phi(x_t)$ , where  $E$  is the ResNet-based [27] encoder from Vahdat and Kautz [28].  
 92 RSSMs model the state of the system as  $s_t$  (which contains both  $e_t$  and a hidden recurrent state  $h_t$ )  
 93 using a forward dynamics model, predicting  $s_t = f_\phi(s_{t-1}, a_{t-1})$ .

94 This system is learned via variational inference, leveraging the posterior  $q_\phi(s_t|s_{t-1}, a_{t-1}, e_t)$  and an  
 95 image posterior  $p_\phi(x_t|s_t)$ . Throughout the paper we refer to the RSSM as a world model, and we  
 96 train it using the ELBO loss [29, 30].

**Intrinsic Motivation** When learning a dynamics model of the world,  $f(s_{t+1}|s_t, a_t)$ , it is possible  
 to use the quality of the model as an intrinsic reward. For instance, Pathak et al. [11] use model  
 prediction error as reward

$$r_t = \|f(s_{t+1}|s_t, a_t) - s_{t+1}\|$$

However, this formulation is dependent on environment dynamics, and thus needs a policy-gradient  
 driven approach to optimizer it. Thus, Pathak et al. [14] propose to minimize the *disagreement*  
 between an ensemble of dynamics model  $f_{\phi^{(k)}}$  for  $k = 1, \dots, M$ . The disagreement reward can be  
 described as

$$\mathbb{E}_{s_t, a_t, s_{t+1} \sim \rho(s)} [\text{Var}_k(f_{\phi^{(k)}})]$$

### 97 4 Autonomous Real World Robot Learning

98 Intelligent agents should be able to perform diverse tasks in  
 99 complex, real world environments. There are three major challenges to this. (1) There is a large space of possible interactions,  
 100 especially in continuous control. (2) It is difficult to obtain any  
 101 reward signal without human supervision. (3) There is a large  
 102 cost for collecting experience with real hardware.

103 To this end, we propose ALAN, an autonomous robot learning  
 104 algorithm that is able to efficiently explore in the real world,  
 105 and achieve various goals, which can be self-supervised or hu-  
 106 man provided. Firstly, we use offline visual data to reduce the  
 107 search space for the robot. ALAN identifies the locations of  
 108 potential interesting and complex interactions for the robot. Sec-  
 109 ondly, ALAN defines a novel intrinsic exploration objective  
 110 for the agent to direct its behavior. This novel objective has an  
 111 environment-centric component and an agent-centric component.  
 112 We explain how these insights enable efficient autonomous learning below.



Figure 2: Visualizations of the object detections we use [31]. We sample masks for each workspace.

114 **4.1 Priors from offline data**

115 In any given real world environment, there are many interesting aspects to manipulate. However,  
 116 there are even more parts of the environment which are akin to free space, and do not provide much  
 117 insight to the robot. Exploration methods often get stuck in “free space” with the agent moving itself  
 118 only. One such way to avoid this is to leverage visual priors from offline data, helping understand  
 119 *what* to explore. One instantiation of this is using object-detectors. Recent models [31] are robust  
 120 and can identify objects even in cluttered scenes. We can use these to initialize the robot close to the  
 121 object. Consider the scene image at time  $t$ ,  $I_t$ . We process the image using a visual object detector  
 122  $O_\phi$  and obtain  $k$  masks  $M_1, \dots, M_k$ . At training time, we sample  $M_i$ , obtain the image center of  
 123 the mask  $(u_i, v_i)$ , and convert using an approximate transform  $(R, t)$  obtained from calibration. We  
 124 obtain a trajectory  $\tau$  from the central reset pose  $p_0$  to  $p_i$ , which is the pose of the centre of the mask  
 125 in the world coordinate frame. We use an off-the-shelf RRT\* [32] planner to reach these locations.

126 **4.2 World Model**

127 Traditional RL methods struggle to train in the real world while interacting with many different  
 128 tasks due to sample efficiency. To overcome this, we enable ALAN to leverage *all* the data seen for  
 129 every task. Naively combining all the data the robot sees will not work, as it will override previous  
 130 experience. Thus, we leverage the RSSM model [4, 16, 6, 5] discussed above. Since our inputs  
 131 are high dimensional, we augment the RSSM to learn the agent learns a compact representations  
 132 of its environment by learning to predict observed data trajectories using the RSSM model ([4]),  
 133 which has recurrence to model long-range dependencies and also accounts for the stochasticity of  
 134 the environment. More concretely, this *world model*,  $W_\theta$  can be described using the following  
 135 components.

Image Encoder	$h_t = [e_t^1, e_t^2, e_t^3], e_t^i = \phi_i(x_t)$	Dynamics Prior	$p_\theta(s_{t+1} s_t, a_t)$
Image Decoder	$f_\theta(x_t^i s_t), i = 1, 2, 3$	Dynamics Posterior	$q_\theta(s_{t+1} s_t, a_t, h_{t+1})$
Embed Decoder	$g_\theta(e_t^i s_t), i = 1, 2, 3$	Change Predictor	$r_\theta(c_t^i e_t^i, e_{t-1}^i)$

(1)

136 Our Image Encoder and Decoder employ the visual encoding approach from Vahdat and Kautz [33],  
 137 while the rest of our world model follows a similar structure to [4]. The use of the world model  
 138 enables our approach to not only work for combining data from many tasks but it also allows for  
 139 efficient exploration, as we can use the model think about where to in fact explore. We call this  
 140 *prospective agent-centric exploration*.

141 **4.3 Prospective Agent-Centric Exploration**

142 When interacting with the world, the robot must carefully decide how to spend its resources. Knowing  
 143 in advance about what interactions might provide some benefit to the robot is very important. The  
 144 robot can directly leverage the world model,  $W_\theta$  for this. Given some initial observation  $x_1$ , the  
 145 robot could sample actions  $a_1, \dots, a_t$  that it thinks might be interesting. Concretely, the robot defines  
 146 an *intrinsic* motivation score,  $r_t^i$ , as the robot does not have any other external reward. This is  
 147 derived from function  $I(\cdot)$  for any transition that it images  $z_t, a_t, z_{t+1}$ . We sample many such action  
 148 sequences,  $\mathcal{A}_1, \dots, \mathcal{A}_M$  and run them through the intrinsic function.  $W_\theta(x_1, \mathcal{A}_m)$  describes the output  
 149 of the world model  $z_1, z_2, \dots, z_T$ , all states in imagination. The goal for the agent is to maximize:

$$\mathcal{A}^* = \operatorname{argmax}_m I(W_\theta(x_0, \mathcal{A}_m)) \quad (2)$$

150 An important aspect to our exploration approach is to define the agent-centric exploration metric as  
 151 the *disagreement* between a set of predictive models  $f_{\phi^{(k)}}$ , as explained in the previous section:

$$\mathbb{E}_{z \sim \rho(z)} [\operatorname{Var}_k(f_{\phi^{(k)}}(z_{t+1}|z_t, a_t))] \quad (3)$$

---

**Algorithm 1** Autonomous Learning Agent (ALAN)

---

```
1: Run visual processing (eg: object detectors) to get workspaces  $\mathcal{T}_1, \mathcal{T}_1 \dots \mathcal{T}_n$ .  
2: Collect random trajectories for each workspace  $\mathcal{T}_i$   
3: Train a joint World Model  $\mathcal{W}$  on all the aggregated data  
4: while exploring do  
5:   Randomly select workspace  $\mathcal{T}_i$  to explore.  
6:   Initialize workspace specific model  $\mathcal{W}_i = \mathcal{W}$  and workspace specific data  $\mathcal{D}_i$   
7:   while number of new trajectories is less than N do  
8:     Train  $\mathcal{W}_i$  on  $\mathcal{D}_i$ .  
9:     Run CEM with model  $\mathcal{W}_i$  to maximize intrinsic objective 5  
10:    end while  
11:   Add  $\mathcal{D}_i$  to  $\mathcal{D}$   
12:   Train joint model  $\mathcal{W}$  on aggregated data  $\mathcal{D}$   
13: end while
```

---

152 Here,  $z_t$  is sampled from the RSSM ( $W_\theta$ ) described previously. ALAN iterates over these intrinsic  
153 rewards for action sequences  $\mathcal{A}$  in imagination, in a similar fashion to the CEM algorithm. We  
154 use this intrinsic disagreement objective as it allows the robot to explore skills. While this will  
155 eventually learn meaningful skills, the robot will be curious about spurious actions. To avoid this,  
156 it is important for the exploration to be grounded in the *environment*.

157 **4.4 Environment-Centric Exploration**

158 Grounding the robot with how the environment  
159 behaves is important to guide exploration. No  
160 matter the skills of the robot, it should always be  
161 exploring some important aspects of the environ-  
162 ment. We denote this as Environment-Centric  
163 exploration. This often involves interactions  
164 with objects which cause changes in the envi-  
165 ronment [34, 35]. Given a scene observation at  
166  $x_t$ , we can compute this using visual features  $\Psi$ ,  
167 for instance from a pretrained visual network.  
168 Specifically we use a semantic segmentation net-  
169 work [36]. We define a function  $f_c$  that measures the *change* in the environment.

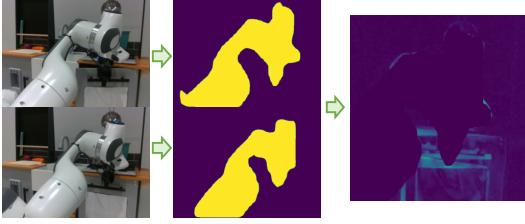


Figure 3: We present a visualization of our environment centric metric, described in (4), right shows the change metric between the two images.

170 Here, the first term is a heuristic image change function, and  $\Psi(x)$  are the semantic segmentation  
171 features. This can be used as an exploration bonus to the intrinsic disagreement approach described.  
172 However, this change can only be measured on collected trajectories, which means that we cannot  
173 imagine potentially environment-changing action sequences. In order to reconcile this, we distill a  
174 network  $r_\theta$  fits world model image embeddings  $e_t$  of images  $x_t$  to  $f_c(x_t, x_{t-1}) = c_t$ . This change  
175 predictor network's loss can be described as the  $\mathcal{L}_c = \|f_c(x_t, x_{t-1}) - r_\theta(c_t | e_t, e_{t-1})\|_2$ .  
176 This gives rise to our intrinsic objective used for classifying action sequences  $a_{1:T} = \mathcal{A}, I_\theta(\cdot)$ :

$$I_\theta(\mathcal{A}) = \sum_{t=1}^T \mathbb{E}_{z \sim \rho(z_t)} [\text{Var}_k(f_{\phi^{(k)}}(z_{t+1} | z_t, a_t))] + r_\theta(c_t | e_t, e_{t-1}) \quad (5)$$

177 **4.5 Continual Exploration and Goal Reaching**

178 We summarize our approach, ALAN, in this section. The first part of the approach, involves collecting  
179 exploration experience using our world model  $W_\theta$ . We assume that this has been trained on a little  
180 bit of data so that the predictions are grounded in the environment. We use the sampling approach  
181 described in Equation 2 with ranking function that uses our agent and environment centric metric



Figure 4: We explore on 6 real world settings as pictured above (left): fridge, tap, knife, stovepot, vegetables and cabinet. On the right, we present a full view of all the worksapces [37].

182 (Equation 5). The robot collects exploration data on each region of the environment,  $\mathcal{T}_i$  and shares  
183 the data between them.

184 An important aspect to our approach is how to train on one region,  $\mathcal{T}_i$  without confusing it with  
185 another region of the workspace  $\mathcal{T}_j$ . ALAN is able to continuously train, using data from previous  
186 experiences. Specifically, we propose always tuning the world model,  $W_\theta$  on all the data, but when  
187 the robot explores the workspace  $\mathcal{T}_i$ , we propose only finetuning on data from that region. This allows  
188 for the model to focus on a task, and yet aggregate experiences.

189 **Achieving goals** Given the contact-rich data collected by the exploration controllers, the question  
190 is how can we efficiently and effectively use this data to learn manipulation skills? Since there is  
191 no provided objective, the agent will have to set its own goals from the collected experience, and  
192 practice reaching them. It is possible for the agent to sample goals from previously seen exploration  
193 data. Since the agent sees "interesting" data, any possible state can be a goal. Concretely, given some  
194 human sampled goal images,  $x_g$ , we train a goal conditioned policy  $\pi_{GC}$  that can learn from the  
195 exploratory data. The policy is trained on world model embeddings  $z_t$ , and a randomly sampled goal  
196 from the same trajectory  $z_g$ . The policy predicts the action  $a_t$  taken at timestep  $t$ .

$$\pi^* = \operatorname{argmin}_{i,j} \sum ||\pi_{GC}(z_i|z_j) - a_j||_2 \quad (6)$$

197 Since our method has seen interesting trajectories, it is more likely to see diverse goals, and thus  
198 when a human provided goal  $z_{gh}$  is given, more likely to reach it. We find that ALAN provides just  
199 this, as the exploration of our method is meaningful and leads to a lot of rich data.

## 200 5 Experimental Setup

201 In our experiments, we seek to address the following questions : 1) Does our system enable au-  
202 tonomous exploration and discovery of useful manipulation skills in a complex real world environment  
203 ? 2) How does the quality of this data compare to prior approached? 3) Is it possible to use this data  
204 to reach human specified goals?

205 **Real World Setup** We tested our system on a Franka Panda 7dof robot, and on a real-world  
206 kitchen play-set (taken from the bridge dataset [37]), which has many diverse objects and possible  
207 manipulation tasks, comprising a very large search space (112cm X 92cm X 110cm). Specifically, the  
208 6 object regions that are returned by our human prior, the object detector from Zhou et al. [31]. These  
209 regions were chosen based on high confidence from the model (Figure 2). The following locations  
210 are sampled: the right cabinet, burner knobs, stove pot, tap, knife, vegetables on the right counter-top,  
211 and the left cabinet. We sample different meaningful goal images for each of these regions. For  
212 example, this involves lifting the knife up in the air, removing the pot, opening or closing the cabinet.

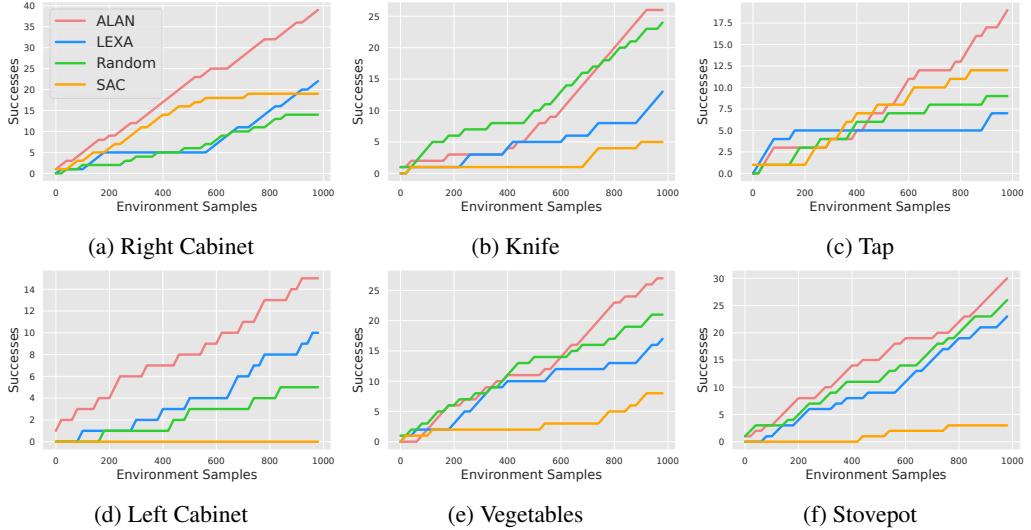


Figure 5: Coincidental success for exploration on our six workspaces. We can see that ALAN sees a lot more success in every task compared to baseline exploration approaches.

213 **Training Procedure** For each of the regions, we first collect a random dataset of 50 trajectories.  
 214 We then train a joint world model on these  $W_\theta$ . As described above, ALAN trains the joint world  
 215 model on all workspaces continually, but when exploring an area it only uses data from that region.  
 216 The process of data collection takes about 45 seconds per trajectory. The finetuning process for  $W_\theta$   
 217 takes about 20 minutes. Our method uses the NVAE [28] architecture, as well as the RSSM from  
 218 Mendonca et al. [6]. To extract the environment centric metric, we train a Mask RCNN model [36] on  
 219 200 instances of the robot. We extract features after applying this mask to make them agent agnostic.

220 **Baselines and Ablations** It is important to test the different design decisions and choices made  
 221 in designing ALAN. Firstly, we test if our novel exploration metric, which combines both agent  
 222 and environment centric exploration, works better than exploration approaches that just use an  
 223 agent centric approach, such as Mendonca et al. [6]. This method is one of our baselines, which  
 224 we call LEXA. It uses the same approach and training procedure as ALAN, however, does not use  
 225 our environment centric intrinsic reward described in Equation 4, and only uses Equation 3 as a  
 226 motivation score. Secondly, we compare whether the world model,  $W_\theta$  is useful. We call this method,  
 227 Rand, as it collects random data only, this method uses a similar goal reaching approach as Pari et al.  
 228 [38]. Finally, we test our method against state of the art RL approach, Soft Actor-Critic (SAC) [39].  
 229 This approach using Q-Learning, but with our environment centric reward (Equation 4). The goal  
 230 achieving mechanism for all is the same goal conditioned policy learning procedure described in the  
 231 section above (Equation 6). We also test the goal conditioned policy against an RL policy trained in  
 232 the world model’s imagination, similar to Mendonca et al. [6].

## 233 6 Results

234 **Exploration** We present results of the exploratory data collected for ALAN, in Figure 5. It is  
 235 important to measure the quality of the data collected. To this end, we define an exploration metric  
 236 that measures the *number* of successful interactions over time. The x-axis in Figure 5 is the number  
 237 of environment samples taken, and the y-axis is the number of successes. We can see that in all of the  
 238 tasks, our approach (red) outperforms the baseline approaches by a significant margin. Empirically,  
 239 we have observed this as well, as we have seen the robot interact a lot more the workspace in  
 240 meaningful ways. LEXA [6], which does not use our environment-centric intrinsic module tends to  
 241 perform worse when more precise control is required, for example when handling the knife or the tap.  
 242 On the other hand, these are the tasks where Random tends to perform better. We see that overall,



Figure 6: Manually specified goals used for evaluation

243 apart from ALAN, SAC [39] tends to perform better than the baseline. This shows that latching on to  
 244 our intrinsic motivation (visual change) objective does lead to success.

245 Additionally, we see that ALAN continuously improves as it sees more data. Towards the end of  
 246 the exploration process, it is almost constantly seeing successes. This shows that ALAN is in fact  
 247 getting better with more exploration data. We suspect running our approach for longer will lead to  
 248 even better data across all the tasks. This does not hold true for other methods, including LEXA,  
 249 which follows a similar data aggregation and word model scheme as us.

250 **Achieving Goals** It is important to show that  
 251 the exploratory data can be used for performing  
 252 human specified tasks. For this, we use goal  
 253 conditioned supervised learning to learn a goal  
 254 reaching policy from the collected data, and also  
 255 compare to an approach that trains a policy using  
 256 model imagination (LEXA). We try 10 different  
 257 trials for the goals shown in 6, and present the  
 258 average success rates in Table 1. We can see  
 259 that our approach significantly outperforms the  
 260 random data baseline, while being comparable  
 261 or better than SAC and LEXA. The importance  
 262 of the world model is apparent from the cabinet opening task, as the approaches which do not use  
 263 it fail to get any success. The difference between the tap goals is that goal1 requires more precise  
 264 control as the handle is not pushed all the way down. Our approach is the only one that gets high  
 265 success on this task, while remaining competitive on the other tap goal as well.

## 266 7 Discussion and Limitations

267 We present ALAN, an autonomously exploring agent that can aggregate data and experience across  
 268 multiple tasks. Our approach hinges on two major components, world-model based agent-centric  
 269 intrinsic motivation as well as environment-centric visual change. This reward in the absence of true  
 270 task rewards helps our agent achieve success in the real world. We also show that for a subset of  
 271 tasks, it is possible to distill this experience into a goal-conditioned reaching policy.

272 **Limitations** While enabling efficient exploration, our approach is centered on environment driven  
 273 priors. Firstly, if the object detections are wrong or there are inaccuracies in the motion planning or  
 274 calibration, this will lead to a bad initialization for exploration. This can be prevented by training  
 275 a more robust policy with a larger workspace. Secondly, our approach hinges on the quality of the  
 276 world model. When in the same kitchen, this model is able to aggregate data, however, this may  
 277 not hold when there is too much variation between scenes or environments. Finally, our approach is  
 278 limited by the quality of the visual change. If the visual change detects spurious correlations (such as  
 279 lighting, artefacts etc), then it will focus on the wrong things. We saw to some extent, for example,  
 280 the change metric focusing on shadows or lighting conditions at some points in time. In the future,  
 281 we hope to address these limitations of robustness, diversity in the environments and reliability of the  
 282 visual change metric.

283 **References**

- 284 [1] S. Levine, C. Finn, T. Darrell, and P. Abbeel. End-to-end training of deep visuomotor policies.  
285 *JMLR*, 2016.
- 286 [2] D. Kalashnikov, A. Irpan, P. Pastor, J. Ibarz, A. Herzog, E. Jang, D. Quillen, E. Holly, M. Kalakr-  
287 ishnan, V. Vanhoucke, et al. Qt-opt: Scalable deep reinforcement learning for vision-based  
288 robotic manipulation. *arXiv preprint arXiv:1806.10293*, 2018.
- 289 [3] D. Ha and J. Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018.
- 290 [4] D. Hafner, T. Lillicrap, J. Ba, and M. Norouzi. Dream to control: Learning behaviors by latent  
291 imagination. *arXiv preprint arXiv:1912.01603*, 2019.
- 292 [5] R. Sekar, O. Rybkin, K. Daniilidis, P. Abbeel, D. Hafner, and D. Pathak. Planning to explore  
293 via self-supervised world models. *ICML*, 2020.
- 294 [6] R. Mendonca, O. Rybkin, K. Daniilidis, D. Hafner, and D. Pathak. Discovering and achieving  
295 goals via world models. *NeurIPS*, 2021.
- 296 [7] A. Strehl and M. Littman. An analysis of model-based interval estimation for markov decision  
297 processes. *Journal of Computer and System Sciences*, 2008.
- 298 [8] M. O. Duff and A. Barto. *Optimal Learning: Computational procedures for Bayes-adaptive*  
299 *Markov decision processes*. PhD thesis, University of Massachusetts at Amherst, 2002.
- 300 [9] P. Poupart, N. Vlassis, J. Hoey, and K. Regan. An analytic solution to discrete bayesian  
301 reinforcement learning. In *ICML*, 2006.
- 302 [10] M. Bellemare, S. Srinivasan, G. Ostrovski, T. Schaul, D. Saxton, and R. Munos. Unifying  
303 count-based exploration and intrinsic motivation. In *Advances in Neural Information Processing*  
304 *Systems*, pages 1471–1479, 2016.
- 305 [11] D. Pathak, P. Agrawal, A. A. Efros, and T. Darrell. Curiosity-driven exploration by self-  
306 supervised prediction. In *ICML*, 2017.
- 307 [12] K. Lowrey, A. Rajeswaran, S. Kakade, E. Todorov, and I. Mordatch. Plan online, learn offline:  
308 Efficient learning and exploration via model-based control. *arXiv preprint arXiv:1811.01848*,  
309 2018.
- 310 [13] I. Osband, J. Aslanides, and A. Cassirer. Randomized prior functions for deep reinforcement  
311 learning. In *Advances in Neural Information Processing Systems*, pages 8617–8629, 2018.
- 312 [14] D. Pathak, D. Gandhi, and A. Gupta. Self-supervised exploration via disagreement. In  
313 *International Conference on Machine Learning*, pages 5062–5071, 2019.
- 314 [15] D. Hafner, T. Lillicrap, I. Fischer, R. Villegas, D. Ha, H. Lee, and J. Davidson. Learning latent  
315 dynamics for planning from pixels. *arXiv preprint arXiv:1811.04551*, 2018.
- 316 [16] D. Hafner, T. Lillicrap, M. Norouzi, and J. Ba. Mastering atari with discrete world models.  
317 *arXiv preprint arXiv:2010.02193*, 2020.
- 318 [17] A. V. Nair, V. Pong, M. Dalal, S. Bahl, S. Lin, and S. Levine. Visual reinforcement learning  
319 with imagined goals. In *NeurIPS*, pages 9191–9200, 2018.
- 320 [18] D. Pathak, P. Mahmoudieh, G. Luo, P. Agrawal, D. Chen, Y. Shentu, E. Shelhamer, J. Malik,  
321 A. A. Efros, and T. Darrell. Zero-shot visual imitation. In *Proceedings of the IEEE conference*  
322 *on computer vision and pattern recognition workshops*, pages 2050–2053, 2018.
- 323 [19] L. Pinto and A. Gupta. Supersizing self-supervision: Learning to grasp from 50k tries and 700  
324 robot hours. *ICRA*, 2016.

- 325 [20] H. Zhu, J. Yu, A. Gupta, D. Shah, K. Hartikainen, A. Singh, V. Kumar, and S. Levine. The  
326 ingredients of real-world robotic reinforcement learning. *arXiv preprint arXiv:2004.12570*,  
327 2020.
- 328 [21] L. P. Kaelbling. Learning to achieve goals. In *IJCAI*, pages 1094–1099. Citeseer, 1993.
- 329 [22] T. Schaul, D. Horgan, K. Gregor, and D. Silver. Universal value function approximators. In  
330 *International conference on machine learning*, pages 1312–1320. PMLR, 2015.
- 331 [23] M. Andrychowicz, F. Wolski, A. Ray, J. Schneider, R. Fong, P. Welinder, B. McGrew, J. Tobin,  
332 P. Abbeel, and W. Zaremba. Hindsight experience replay. *arXiv preprint arXiv:1707.01495*,  
333 2017.
- 334 [24] D. Warde-Farley, T. Van de Wiele, T. Kulkarni, C. Ionescu, S. Hansen, and V. Mnih. Unsuper-  
335 vised control through non-parametric discriminative rewards. *arXiv preprint arXiv:1811.11359*,  
336 2018.
- 337 [25] V. H. Pong, M. Dalal, S. Lin, A. Nair, S. Bahl, and S. Levine. Skew-fit: State-covering  
338 self-supervised reinforcement learning. *arXiv preprint arXiv:1903.03698*, 2019.
- 339 [26] Y. Zhang, P. Abbeel, and L. Pinto. Automatic curriculum learning through value disagreement.  
340 *Advances in Neural Information Processing Systems*, 33, 2020.
- 341 [27] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *CoRR*,  
342 abs/1512.03385, 2015. URL <http://arxiv.org/abs/1512.03385>.
- 343 [28] A. Vahdat and J. Kautz. Nvae: A deep hierarchical variational autoencoder. *Advances in Neural*  
344 *Information Processing Systems*, 33:19667–19679, 2020.
- 345 [29] D. J. Rezende, S. Mohamed, and D. Wierstra. Stochastic backpropagation and approximate  
346 inference in deep generative models. *arXiv preprint arXiv:1401.4082*, 2014.
- 347 [30] D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*,  
348 2013.
- 349 [31] X. Zhou, R. Girdhar, A. Joulin, P. Krähenbühl, and I. Misra. Detecting twenty-thousand classes  
350 using image-level supervision. *arXiv preprint arXiv:2201.02605*, 2022.
- 351 [32] S. Karaman, M. R. Walter, A. Perez, E. Frazzoli, and S. Teller. Anytime motion planning using  
352 the rrt. In *2011 IEEE International Conference on Robotics and Automation*, pages 1478–1483.  
353 IEEE, 2011.
- 354 [33] A. Vahdat and J. Kautz. Nvae: A deep hierarchical variational autoencoder. *Advances in Neural*  
355 *Information Processing Systems*, 33:19667–19679, 2020.
- 356 [34] S. Parisi, V. Dean, D. Pathak, and A. Gupta. Interesting object, curious agent: Learning task-  
357 agnostic exploration. *Advances in Neural Information Processing Systems*, 34:20516–20530,  
358 2021.
- 359 [35] S. Bahl, A. Gupta, and D. Pathak. Human-to-robot imitation in the wild. In *RSS*, 2022.
- 360 [36] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask r-cnn. In *ICCV*, 2017.
- 361 [37] F. Ebert, Y. Yang, K. Schmeckpeper, B. Bucher, G. Georgakis, K. Daniilidis, C. Finn, and  
362 S. Levine. Bridge data: Boosting generalization of robotic skills with cross-domain datasets.  
363 *arXiv preprint arXiv:2109.13396*, 2021.
- 364 [38] J. Pari, N. Muhammad, S. P. Arunachalam, L. Pinto, et al. The surprising effectiveness of  
365 representation learning for visual imitation. *arXiv preprint arXiv:2112.01511*, 2021.
- 366 [39] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy  
367 deep reinforcement learning with a stochastic actor. *arXiv preprint arXiv:1801.01290*, 2018.