

---

# EReLELA: Exploration in Reinforcement Learning via Emergent Language Abstractions

---

Anonymous Author(s)  
Affiliation  
Address  
email

## Abstract

1 Instruction-following from prompts in Natural Languages (NLs) is an important  
2 benchmark for Human-AI collaboration. Training Embodied AI agents for  
3 instruction-following with Reinforcement Learning (RL) poses a strong exploration  
4 challenge. Previous works have shown that NL-based state abstractions can  
5 help address the exploitation versus exploration trade-off in RL. However, NLs  
6 descriptions are not always readily available and are expensive to collect. We  
7 therefore propose to use the Emergent Communication paradigm, where artificial  
8 agents are free to learn an emergent language (EL) via referential games, to bridge  
9 this gap. ELs constitute cheap and readily-available abstractions, as they are the  
10 result of an unsupervised learning approach. In this paper, we investigate (i) how  
11 EL-based state abstractions compare to NL-based ones for RL in hard-exploration,  
12 procedurally-generated environments, and (ii) how properties of the referential  
13 games used to learn ELs impact the quality of the RL exploration and learning.  
14 Results indicate that the EL-guided agent, namely EReLELA, achieves similar  
15 performance as its NL-based counterparts without its limitations. Our work shows  
16 that Embodied RL agents can leverage unsupervised emergent abstractions to  
17 greatly improve their exploration skills in sparse reward settings, thus opening new  
18 research avenues between Embodied AI and Emergent Communication.

## 19 1 Introduction

20 Natural Languages (NLs) have some properties, such as compositionality and recursive syntax, that  
21 allow us to talk about infinite meanings while only using a finite number of words (or even letters,  
22 or phonemes...). In other words, it enables us to be as expressive as one might needs. However,  
23 it may be interesting sometimes to use language to abstract away from the details and only focus  
24 on the essence of a specific experience, or a specific sensory stimulus. Thus, even though NLs can  
25 sometimes be used with high expressiveness, they also can work as abstractions. For instance, using a  
26 unique utterance to refer to a lot of semantically-similar but (visually) different situations, such as the  
27 one presented in Figure 1 where the utterance ‘one can see a purple key and a green ball’ can refer  
28 to many of the first-person perspective of the embodied agent, irrespective of the actual perspective  
29 under which each object is seen.

30 Tam et al. [61] referred to that aspect as compacting/clustering a state/observation space, which is  
31 in effect segmenting it into a set of less-detailed but more-meaningful sub-spaces. We employ the  
32 term meaningful with respect the task that the embodied agent is possibly trained for. For instance,  
33 if the task consists of picking and placing objects, then it is meaningful for utterances to contain  
34 information about objects and places, but not so much to contain information about other agents in  
35 the environment, if any. In this paradigm, Tam et al. [61] and Mu et al. [51] provided some arguments  
36 towards the compacting/clustering assumption of NLs, as they used NLs oracle to build an abstraction

37 over a 3D and 2D environments. They relied upon state-of-the-art exploration algorithms, such as  
38 Random Network Distillation (RND - Burda et al. [9]) and Never-Give-Up (NGU - Badia et al. [1]),  
39 which can be difficult to deploy.

40 Thus, in this work, we aim to simplify the process of using  
41 languages as abstractions and address the limitation of using  
42 NLS, as they are expensive to harvest and not necessarily the  
43 most meaningful abstraction for any given task. Indeed, instead  
44 of state-of-the-art exploration algorithms, we show that simpler  
45 count-based approaches combined with language abstraction  
46 can be leveraged for hard-exploration tasks. And, in order to  
47 remove the reliance on NLS, we look at the field of Emergent  
48 Communication (EC) [41, 7] which have shown that artificial  
49 languages, that we refer to as emergent languages (ELs), can  
50 emerge through unsupervised learning algorithms, such as Ref-  
51 erential Games and variants [19], with structure and properties  
52 similar to NLS. Our experimental evidences show that ELs,  
53 acquired over an embodied agent’s observations in an online  
54 fashion and in parallel of its training, can be leveraged for hard-  
55 exploration tasks. We investigate what are the properties of  
56 NLS and ELs in terms of their abstraction building abilities  
57 by proposing a novel metric entitled Compactness Ambigu-  
58 ity Metric (CAM). Measures show that ELs abstractions are  
59 aligned but not similar to NLS in terms of the abstractions they  
60 perform, as the Emergent Communication context successfully  
61 picks up on the meaningful features of the environment. Indeed,  
62 EReLELA’s abstractions reflect colors in the *MultiRoom-N7-S4*  
63 environment which only features coloured, unlocked doors, but no distracting objects, or shapes in  
64 the *KeyCorridor-S3-R2* environment where it is important to pickup a relevant key, among other  
65 distractingly-shaped objects, and to open the locked door-shaped object.

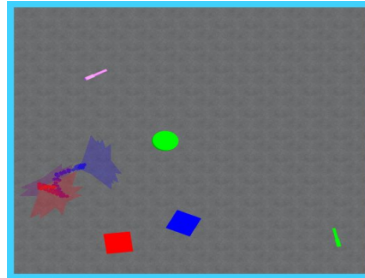


Figure 1: Top-view visualization of a wall-free 3D environment with different objects (e.g. red and blue cubes, purple and green keys, and green ball) showing the trajectory (from blue to red dots) of a randomly-walking embodied agent, with first-person perspectives highlighted at relevant timesteps using colored cones - showing the agent’s viewpoint direction when a new utterance is used to describe the first-person perspective using an oracle speaking in NL.

66 We continue by reviewing EC and RL backgrounds and notations in Section 2. After detailing our  
67 method in Section 3, we present experimental results on procedurally-generated, hard-exploration  
68 task from the MiniGrid [15] benchmarks in Section 4. Finally, we discuss in Section 5 the results  
69 presented in light of some related works and highlight possible future works.

## 70 2 Background & Notation

71 We provide details on our Reinforcement Learning (RL) settings and count-based exploration methods  
72 in Section 2.1. Then, we review Emergent Communication in Section 2.2.

### 73 2.1 Exploration vs Exploitation in Reinforcement Learning

74 An RL agent interacts with an environment in order to learn a mapping from states to actions that  
75 maximises its reward signal. Initially, both the reward signal and the dynamics of the environment,  
76 i.e. the impact that the agent actions may have on the environment, are unknown to the agent. It must  
77 explore the environment and gather information, but, all the while it is exploring, it cannot exploit the  
78 best strategy that it has found so far to maximise the currently-known reward signal. This dilemma is  
79 known as the Exploration-vs-Exploitation trade-off of RL. This dilemma is only the start of the rabbit  
80 hole, as it can even get worse. Indeed, in sparse reward environments, the reward signal is mainly  
81 zero most of the time. This context makes it very difficult for RL agents to learn anything, because RL  
82 algorithms derive feedback (i.e. gradients to update their parameters) from the reward signal that they  
83 observe from the environment. It is usually referred to as extrinsic, in order to differentiate it from an  
84 intrinsic reward signal. As the extrinsic reward is mostly zero, RL agents must exploit another signal  
85 to derive information about the currently-unknown environment. This other signal can be found in  
86 relation to the observation/state space, as RL agents can learn to seek novelty or surprise around the  
87 observation/state space and attempt to manipulate it efficiently by choosing relevant actions. Focusing  
88 on this novelty, RL agents can harvest an intrinsic reward signal, in the sense that RL agents are  
89 building it and giving it to themselves. Note that this intrinsic reward signal is very different from the

90 extrinsic reward signal, because it does not inform about the task that RL agents need to perform  
 91 in the environment. Ideally, though, it provides a graded and dense signal that the RL agent can  
 92 use to start learning anything about the environment. This is inspired by intrinsic motivation in  
 93 psychology [53]. Exploration driven by curiosity/novelty might be an important way for children  
 94 to grow and learn. Here, we focus on novelty, but the intrinsic rewards could be correlated with e.g.  
 95 impact [54], surprise [9] or familiarity of the state. The intrinsic reward signal is only a proxy for  
 96 RL agents to start to make progress into learning about the environment and eventually, hopefully  
 97 encounter some non-zero extrinsic reward signal along the way. It provides a denser reward signal  
 98 that can guide RL agents into learning internal representations about the environment’s dynamic so  
 99 that, whenever some extrinsic reward are encountered along the way, then they can efficiently bind  
 100 their previously-learned representations to those recently-encountered extrinsic rewards.

101 Formally, we study a single agent in a Markov Decision Pro-  
 102 cess (MDP) defined by the tuple  $(\mathcal{S}, \mathcal{A}, T, \mathcal{R}, \gamma)$ , referring to,  
 103 respectively, the set of states, the set of actions, the transition  
 104 function  $T : \mathcal{S} \times \mathcal{A} \rightarrow P(\mathcal{S})$  which provides the probability  
 105 distribution of the next state given a current state and action,  
 106 the reward function  $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow r$ , and the discount fac-  
 107 tor  $\gamma \in [0, 1]$ . The agent is modelled with a stochastic policy  
 108  $\pi : \mathcal{S} \rightarrow P(\mathcal{A})$  from which actions are sampled at every time step of an episode of finite time horizon  
 109  $T$ . The agent’s goal is to learn a policy which maximises its discounted expected return at time  $t$ ,  
 110 defined in equation 1. We further define  $\mathcal{R} = \lambda_{\text{ext}}\mathcal{R}^{\text{ext}} + \lambda_{\text{int}}\mathcal{R}^{\text{int}}$  as the weighted sum of the extrinsic  
 111 and intrinsic reward functions, respectively,  $\mathcal{R}^{\text{ext}}, \mathcal{R}^{\text{int}}$ , with weights  $\lambda_{\text{ext}}, \lambda_{\text{int}}$ . Indeed, while the  
 112 extrinsic reward is provided by the environment, we assume that for any tuple  $(s_t, a_t, s_{t+1})$  we can  
 113 compute an intrinsic reward.

$$R_t = \mathbb{E}_{\substack{s_{t+k+1} \sim T(s_{t+k}, a_{t+k}) \\ a_{t+k+1} \sim \pi(s_{t+k+1})}} \left[ \sum_{k=0}^T \gamma^k R(s_{t+k+1}, a_{t+k+1}) \right] \quad (1)$$

114 Stanton and Clune [58] identifies two categories of exploration strategies, to wit *across-training*,  
 115 where novelty of states, for instance, is evaluated in relation to all prior training RL episodes, and  
 116 *intra-life*, where it is evaluated solely in relation of the current RL episode. And, historically, we  
 117 can identify two types of intrinsic motivation exploration depending on how the intrinsic reward is  
 118 computed, either relying on count-based or prediction-based methods. Prediction-based methods fit  
 119 into the *across-training* category and count-based methods can actually fit in both categories but they  
 120 have mainly been instantiated in the literature as *across-training* methods after extension of *intra-life*  
 121 core mechanisms. As our proposed architecture EReLELA fit into the category of count-based  
 122 methods, we detail them further. In the context of an intrinsic reward signal correlated with surprise,  
 123 then it is necessary to quantify how much of surprise each observation/state provides. Intuitively, we  
 124 can count how many times a given observation/state has been encountered and derive from that count  
 125 our intrinsic reward. The reward would guide the RL agent to prefer rarely visited/observed states  
 126 compared to common states. This is referred to as the count-based exploration method. Count-based  
 127 exploration method were originally only applicable to tabular RL where the state space is discrete  
 128 and it is easy to compare states together. When dealing with continuous or high-dimensional state  
 129 spaces, such method is not practical. Thus, Bellemare et al. [3] proposed (and extended in Ostrovski  
 130 et al. [52]) a pseudo-count approach which was derived from increasingly more efficient density  
 131 models, and they showed success in applying it to image-based exploration environments from Atari  
 132 2600 benchmark, such as *Montezuma’s Revenge*, *Private Eye*, and *Venture*. We provide more relevant  
 133 details in Appendix B.

134 Nevertheless, hard-exploration task involving procedurally-generated environments are notoriously  
 135 difficult for count-based exploration methods. Indeed, when states are procedurally-generated, almost  
 136 all states will be showing ‘novel’ features, most times irrespectively of whether it is relevant to the  
 137 task or not. It will follow that their state (pseudo-)count will always be low and therefore the RL  
 138 agent will get feedback towards reaching all of them indefinitely, but if every state is ‘novel’ then  
 139 there is nothing to guide the agent in any specific direction that would entail to good exploration.

## 140 2.2 Emergent Communication

141 Emergent Communication is at the interface of language grounding and language emergence. While  
 142 language emergence raises the question of how to make artificial languages emerge, possibly with  
 143 similar properties to NLS, such as compositionality [2, 24, 45, 55], language grounding is concerned  
 144 with the ability to ground the meaning of (natural) language utterances into some sensory processes,

145 e.g. the visual modality. On one hand, the compositionality of ELs has been shown to further  
 146 the learnability of said languages [38, 57, 8, 45] and, on the other hand, the compositionality of  
 147 NLS promises to increase the generalisation ability of the artificial agent that would be able to  
 148 rely on them as a grounding signal, as it has been found to produce learned representations that  
 149 generalise, when measured in terms of the data-efficiency of subsequent transfer and/or curriculum  
 150 learning [27, 49, 50, 33]. Yet, emerging languages are far from being ‘natural-like’ protolanguages  
 151 [40, 10, 11], and the questions of how to constraint them to a specific semantic or a specific syntax  
 152 remain open problems. Nevertheless, some sufficient conditions can be found to further the emergence  
 153 of compositional languages and generalising learned representations [40, 43, 17, 5, 24, 39, 12, 21].

154 The backbone of the field rests on games that emphasise the functionality of languages, namely,  
 155 the ability to efficiently communicate and coordinate between agents. The first instance of such  
 156 an environment is the *Signaling Game* or *Referential Game (RG)* by Lewis [44], where a speaker  
 157 agent is asked to send a message to the listener agent, based on the *state/stimulus* of the world that it  
 158 observed. The listener agent then acts upon the observation of the message by choosing one of the  
 159 *actions* available to it in order to perform the ‘best’ *action* given the observed *state* depending on the  
 160 notion of ‘best’ *action* being defined by the interests common to both players. In RGs, typically, the  
 161 listener action is to discriminate between a target stimulus, observed by the speaker and prompting  
 162 its message generation, and some other distractor stimuli. Distractor stimuli are selected using a  
 163 distractor sampling scheme, which has been shown to impact the resulting EL [42, 43]. The listener  
 164 must discriminate correctly while relying solely on the speaker’s message. The latter defined the  
 165 discriminative variant, as opposed to the generative variant where the listener agent must reconstruct/  
 166 generate the whole target stimulus (usually played with symbolic stimuli). Visual (discriminative)  
 167 RGs have been shown to be well-suited for unsupervised representation learning, either by competing  
 168 with state-of-the-art self-supervised learning approaches on downstream classification tasks [22], or  
 169 because they have been found to further some forms of disentanglement [28, 35, 14, 46] in learned  
 170 representations [65, 18]. Such properties can enable “better up-stream performance”[63], greater  
 171 sample-efficiency, and some form of (systematic) generalization [48, 26, 59]. Thus, this paper aims  
 172 to investigate visual discriminative RGs as auxiliary tasks for RL agents.

### 173 3 Method

174 In this section, following the acknowledgement of a gap in terms of evaluating the abstractions  
 175 that different languages perform over different state/observation space, we start by introducing in  
 176 Section 3.1 our Compactness Ambiguity Metric (CAM) that attempts to fill in that gap. Then, in  
 177 Section 3.2, we present the EReLELA architecture that leverages EL abstractions in an *intra-life*  
 178 count-based exploration scheme for RL agents.

#### 179 3.1 Compactness Ambiguity Metric

180 In order to measure qualities related to the kind of abstraction that a language performs over stimuli,  
 181 we propose to rely on the temporal aspects of embodied agent’s trajectories in a given environment.  
 182 We build over the following intuition, represented in Figure 2: we consider two possible languages  
 183 grounded into the first-person viewpoint of an embodied agent situated in a 3D environment populated  
 184 with objects of different shapes and colors. On one hand, we have the Blue language, which is only  
 185 concerned about blue objects and its utterances only describe that they are of color blue when they  
 186 are, while, on the other hand, we have the Color language, which is describing the color of all  
 187 visible objects. Inherently, those two languages expose different semantics about the world, and  
 188 therefore they perform different abstractions. We aim to build a metric that captures how different the  
 189 semantics they expose are. To do so, we propose to arrange their respective utterances when prompted  
 190 with the very same agent’s trajectories into different timespan-focused buckets towards building  
 191 an histogram. These timespan-focused buckets reflect  $\delta(u)$  the number of consecutive timesteps  
 192  $(t_k)_{k \in [k_{\text{start}}, k_{\text{start}} + \delta(u)]}$  for which a specific utterance  $u$  would be uttered by a speaker of each language  
 193 when prompted with the stimuli in those timesteps. We will refer to these as compactness counts. For  
 194 instance the Blue language’s utterance ‘I see a blue object’ at the beginning of the trajectory occupies  
 195 twice as more consecutive timesteps as the same utterance coming from a Color language speaker (or,  
 196 its compactness count in the Blue language is twice its compactness count in the Color language).  
 197 Therefore, in the case of the Blue language, this utterance would increment the medium-length bucket,  
 198 while it would increment the short-length bucket in the case of Color language histogram. It ensues

199 that the histograms of timespan-focused buckets captures semantics exposed by each language, and  
 200 we will therefore refer to the resulting histogram as the histogram of semantic-clustering timespans.  
 201 As the toy example highlights, the histograms of semantic-clustering timespans will differ from one  
 202 language to another depending on the semantics each language expose or, in other words, depending  
 203 on the abstractions they perform. This is the first intuition on which the Compactness Ambiguity  
 204 metric is built.

205 Formally, we define  $\mathcal{L}$  as the set of all possible languages over vocabulary  $V$  with maximum sentence  
 206 length  $L$ , such that for any language  $l \in \mathcal{L}$  we denote  
 207  $\text{Sp}_l : \mathcal{S} \rightarrow l$  as a speaker agent or oracle that maps  
 208 any state/observation  $s \in \mathcal{S}$  to a caption or utterance  
 209  $u \in l$ . Thus, we can now consider  $N$  buckets whose  
 210 related timespans  $(T_i)_{i \in [1, N]}$  are sampled relative to  
 211 the maximal length  $T$  of a trajectory in the given environment, and the histogram of semantic-clustering  
 212 timespans that they induce.  
 213  
 214

215 Then, the other intuition on which the metric is built  
 216 is made evident by considering the expressivity or, its  
 217 inverse, the ambiguity, of a given language  $l$ , defined  
 218 as  $\mathcal{E}_l = \frac{\#\text{unique utterances}}{\#\text{unique stimuli}}$  with  $\#$  the set cardinality  
 219 operator. Dealing with stimuli being states/observations of a (randomly walking) embodied agent,  
 220 gathered into a dataset  $\mathcal{D}$ , the number of unique stimuli cannot be estimated reliably when dealing  
 221 with complex, continuous stimuli. Thus, the best we can rely on is a measure of relative expressivity  
 222 over a dataset, that we define as  $\mathcal{RE}_l(\mathcal{D}) = \frac{\#\text{unique utterances}}{\#\text{stimuli}} = \frac{\#\text{Sp}_l(\mathcal{D})}{|\mathcal{D}|}$ , with  $|\cdot|$  being the size  
 223 operator over collections (differing from sets in the sense that they allow duplicates). In those terms,  
 224 the relative expressivity is maximised if and only if (i)  $\#\mathcal{D} = |\mathcal{D}|$ , and (ii)  $\text{Sp}_l$  is a bijection over  
 225  $\mathcal{D}$ . On the other hand, considering that a language  $l$  performs an abstraction over  $\mathcal{D}$  is tantamount  
 226 to some stimuli  $(s, s') \in \mathcal{D}^2$  sharing the same utterance  $u = \text{Sp}_l(s) = \text{Sp}_l(s')$ , i.e. consisting of  
 227 a hash collision, meaning that the mapping  $\text{Sp}_l$  from  $\mathcal{D}$  to  $l$  would not be injective (and therefore  
 228 not bijective). Incidentally, the relative expressivity  $\mathcal{RE}_l(\mathcal{D})$  cannot be maximised, leading to the  
 229 language  $l$  being ambiguous over  $\mathcal{D}$ . In this consideration, we can see that the ambiguity of a  
 230 language (over a given dataset) can be impacted by either the extent to which an abstraction is  
 231 performed (meaning that most colliding states/observations are of consecutive timesteps) or the  
 232 extent to which the dataset is redundant (meaning  $\#\mathcal{D} \ll |\mathcal{D}|$ ). Therefore it is important that our  
 233 proposed Compactness Ambiguity Metric is built to focus on sources of ambiguities that are the  
 234 result of consecutive-timesteps states colliding, more than sources of ambiguities that are the result  
 235 of redundancy in the given dataset.

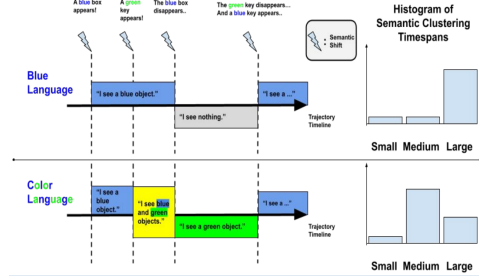


Figure 2: Toy example illustration of how different languages expose different semantics over the same observed trajectory of stimuli, and that the discrepancy in exposed semantics can be captured by an histogram of semantic-clustering timespans.

$$236 \quad \forall i \in [1, N], T_i = 1 + \lceil \lambda_i \cdot \mathcal{RA}_l(\mathcal{D}) \rceil \quad (2)$$

$$237 \quad \forall i \in [1, N], T'_i = 1 + \lceil \lambda_i \cdot T \rceil \quad (3)$$

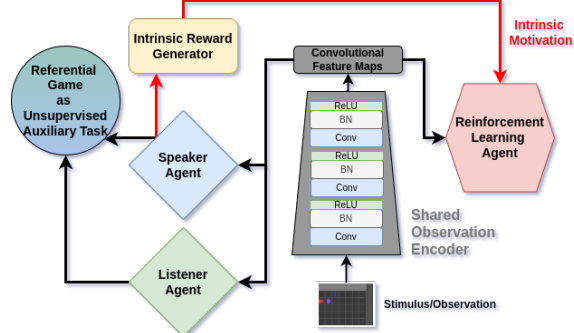
$$238 \quad \forall i \in [1, N], CA(\mathcal{D})_{T_i} = \sum_{u \in l} \frac{\#\delta_{\mathcal{D}}^{\geq T_i}(u)}{\#\delta_{\mathcal{D}}(u)} \quad (4)$$

240  
 241  $\mathcal{D}$ , we define the buckets' related timespans in relation to the relative ambiguity  $\mathcal{RA}_l(\mathcal{D}) = \frac{1}{\mathcal{RE}_l(\mathcal{D})} =$   
 242  $\frac{|\mathcal{D}|}{\#\text{Sp}_l(\mathcal{D})}$ , as shown in equation 2 with  $\lambda_i \in [0, 1]$  s.t.  $\forall(j, k), j < k \implies \lambda_j < \lambda_k$ , and  $\lceil \cdot \rceil$  being  
 243 the ceiling operator. This is in lieu of defining them in relation to the maximal length  $T$  of an agent's  
 244 trajectory in the environment, as shown in equation 3. More specifically, let us first acknowledge  
 245 decomposition of relative ambiguity over two independent quantities, one for each of its sources  
 246 being either abstraction or redundancy, such that  $\mathcal{RA}_l = \mathcal{RA}_l^{\text{redundancy}} + \mathcal{RA}_l^{\text{abstract}}$ . Then note that  
 247 the relative ambiguity is equal to the mean number of consecutive timesteps, or compactness count,  
 248 for which a given utterance would be used when the unique utterances are uniformly distributed  
 249 over the dataset  $\mathcal{D}$ . Thus, in the metric, we propose to absorb variations of relative ambiguity due to  
 250 redundancy by changing the metric's bucket setup, from Equation 3 to Equation 2. Doing so, it is true  
 251 that the metric's bucket setup will also vary when the abstraction-induced relative ambiguity varies,  
 252 we remark that the metric would not build invariance to this source of relative ambiguity since it is  
 253 taken into accounts when sorting out the different unique utterances into their relevant bucket, based  
 254

255 on the maximal number of consecutive timesteps in which they occur, as shown in equation 4 with  
 256  $\delta_{\mathcal{D}} : l \rightarrow 2^{\mathbb{N}}$  is the compactness count function that associates each utterances  $u \in l$  to its related set  
 257 of compactness counts over dataset  $\mathcal{D}$ , i.e. the set that contains numbers of consecutive timesteps  
 258 for which  $u \in l$  was uttered by  $\text{Sp}_l$ , each time it was uttered without being uttered in the previous  
 259 timestep. For instance, if we consider  $u \in l$  such that  $\text{Sp}_l^{-1}(u) = \{s_{t_1}, s_{t_1+1}, s_{t_1+2}, s_{t_2}\}$ , with  
 260  $(t_1, t_2) \in [0, T]^2$  such that  $t_2 > t_1 + 3$ , then  $\delta_{\mathcal{D}}(u) = \{3, 1\}$  because  $u$  occurred 2 non-consecutive  
 261 times over  $\mathcal{D}$  and those occurrences lasted for, respectively, 3 and 1 consecutive timesteps, i.e. for  
 262 compactness counts of 3 and 1. The superscript  $\geq T_i$  in  $\delta_{\mathcal{D}}^{\geq T_i}$  implies filtering of the output set based  
 263 on compactness counts being greater or equal to  $T_i$ . We provide in appendix C an analysis of the  
 264 sensitivity of our proposed metric, and in appendix E.1 experimental results that ascertain the internal  
 265 validity of our proposed metric, we consider a 3D room environment of MiniWorld [15], filled with 5  
 266 different, randomly-placed objects, as shown in a top-view perspective in Figure 1.

### 267 3.2 EReLELA Architecture

268 This section details the EReLELA architecture, which stands for Exploration in Reinforcement Learning  
 269 via Emergent Language Abstractions. As a count-based exploration method,  
 270 we present here its *intra-life* core mechanism, where intrinsic reward signals are derived from novelty at  
 271 the level of language utterances describing the current observation/state.  
 272 It relies on a hashing-like function (cf. Appendix B), which takes the form of the speaker agent of a referential  
 273 game (RG), to turn continuous and high-dimensional observations/states into discrete, variable-length sequences of tokens.  
 274 EReLELA is built around an RL agent augmented with an unsupervised auxiliary task, a (discriminative, here, or generative) RG, following  
 275 the UNREAL architecture from Jaderberg et al. [31], as shown in Figure 3.



276 Figure 3: EReLELA architecture consisting of a stimulus/observation encoder shared between an RL agent and the speaker and listener agents of a RG, framed as an unsupervised auxiliary task [31]. The language utterances outputted by the RG speaker agent are used in a count-based exploration method to generate intrinsic rewards for the RL agent.

286 We train the RG agents in a descriptive, discriminative RG with  $K = 256$  distractors, every  $T_{RG} =$   
 287  $32768$  gathered RL observations, on a dataset  $\mathcal{D}_{RG}$  consisting of the most recent  $|\mathcal{D}_{RG}| = 8192$   
 288 observations, among which 2048 are held-out for validation/testing-purpose, over a maximum of  
 289  $N_{RG-epoch} = 32$  epochs or until they reach a validation/testing RG accuracy greater than a given  
 290 threshold  $acc_{RG-thresh} = 90\%$ . Our preliminary experiments in Appendices D.1 and D.2 show,  
 291 respectively, that increasing the RG accuracy threshold  $acc_{RG-thresh}$  increases the sample-efficiency  
 292 of the EL-guided RL agent, and that the number of distractors  $K \in [15, 128, 256]$  is critical (even  
 293 more so than the distractor sampling scheme - which we set to be uniform unless specified otherwise),  
 294 and that it correlates positively with the performance of the RL agent. More specific details about  
 295 the RG and its agents' architectures can be found in Appendices F and G and our open-source  
 296 implementation<sup>1</sup>.

## 297 4 Experiments

298 **Agents** Our RL agent is optimized using the R2D2 algorithm from [34] with the Adam optimizer  
 299 Kingma and Ba [36]. Importantly, as it aims to maximise the weighted sum of the extrinsic and  
 300 intrinsic reward functions following equation 1, throughout this paper, we use  $\lambda_{int} = 0.1$  and  
 301  $\lambda_{ext} = 10.0$  in order to make sure that the agent pursues the external goal once the exploration of  
 302 the environment has highlighted it. Further details about the RL agent can be found in Appendix F.  
 303 For our RG agents, we consider optimization using either the Impatient-Only or the LazImpa loss  
 304 function from Rita et al. [56], but the latter is adapted to the context of a Straight-Through Gumbel-  
 305 Softmax (STGS) communication channel [25, 21], as detailed in Appendix G.1, and we refer to  
 306 it as STGS-LazImpa. Indeed, the LazImpa loss function has been shown to induce Zipf's Law of

<sup>1</sup>HIDDEN\_FOR\_REVIEW\_PURPOSE



307 Abbreviation (ZLA) in the ELs. Thus, we can investigate in the following experiments how does  
 308 **structural** similarity between NLs and ELs affect the kind of abstractions they perform, as well as  
 309 the resulting RL agent. Further details about the RG in EReLELA can be found in Appendix G.

310 **Environments.** After having considered in our preliminary experiments (cf. Appendix E.4) the 2D  
 311 environment *MultiRoom-N7-S4*, we propose below experiments in the more challenging *KeyCorridor-*  
 312 *S3-R2* environment from MiniGrid [15]. Indeed, it involves complex object manipulations, such as  
 313 (distractors) object pickup/drop and door unlocking, which requires first picking up the relevantly-  
 314 colored key object.

315 **Natural Language Oracles.** Our implementation of a NL oracle is simply describing the visible  
 316 objects in terms of their colour and shape attributes, from left to right on the agent’s perspective,  
 317 whilst also taking into account object occlusions. For instance, around the end of the trajectory  
 318 presented in Figure 1, the green key would be occluded by the blue cube, therefore the NL oracle  
 319 would provide the description ‘blue cube red cube’ alone. We also implement colour-specific and  
 320 shape-specific language oracles, which consists of filtering out from the NL oracle’s utterance the  
 321 information that each of those language abstract away, i.e. removing any shape-related word in the  
 322 case of the colour-specific language, and vice-versa.

323 **Hypotheses.** We seek to validate the following hypotheses. Firstly, we consider whether NL  
 324 abstractions can help for hard-exploration in RL with a simple count-based approach (**H1**), and refer  
 325 to the relevant agent using NL abstractions to compute intrinsic rewards as NLA. We carry on with  
 326 the hypothesis that ELs can be used similarly (**H2**), and we investigate to what extent do ELs compare  
 327 to NLs in terms of abstraction. We would expect ELs to perform more meaningful abstractions than  
 328 NLs (**H3**), in the sense that their abstractions would be more aligned with the relevant features of a  
 329 given environment.

330 **Evaluation.** We employ 3 random seeds for each agent. We evaluate (H1) and (H2) using both the  
 331 success rate and the manipulation count, in the hard-exploration task of *KeyCorridor-S3-R2*. The  
 332 manipulation count is a per-episode counter incremented each time an object is successfully picked  
 333 up or dropped by the RL agent over the course of each episode. In order to evaluate both (H3.1)  
 334 and (H3.2), we use the CAM to measure the kind of abstractions performed by ELs, and compare  
 335 those measures with those of the oracles’ languages that we previously studied. We report the CAM  
 336 distances between ELs and the NL, Color language, and Shape language oracles, which is computed  
 337 as an euclidean distance in  $\mathbb{R}^6$  by considering the  $N = 6$  CAM scores for each timespans/thresholds  
 338 as vectors in this space. As we remarked that an agent’s skillfulness at the task would induce very  
 339 different trajectories (e.g. in *MultiRoom-N7-S4*, staying in the first room and only ever seeing the  
 340 first door, for an unskillfull agent, as opposed to visiting multiple rooms and observing multiple  
 341 colored-doors, for a skillfull agent), we compute the oracle languages CAM scores on the exact same  
 342 trajectories than used to compute each EL’s CAM scores.

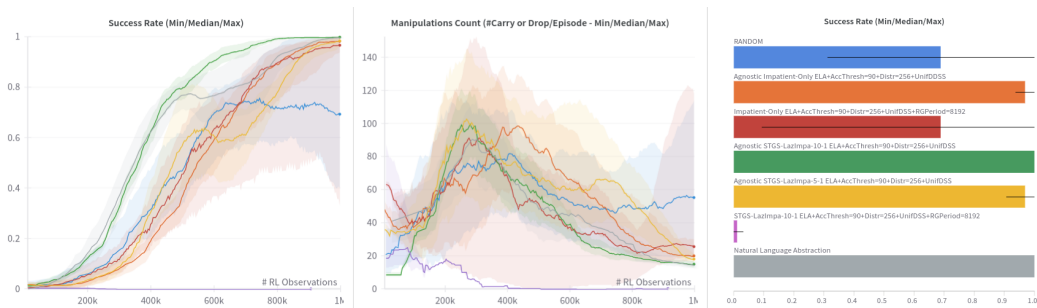


Figure 4: Success rate learning curve (left), computed as running averages over 1024 episodes each time (i.e. 32 in parallel, as there are 32 actors, over 32 running average steps), and barplot (right), along with per-episode manipulation count (middle) in *KeyCorridor-S3-R2* from MiniGrid [15], for different agents: (i) the *Natural Language Abstraction* agent (NLA) refers to using the NL oracle to compute intrinsic reward, (ii) the *STGS-LazImpa- $\beta_1$ - $\beta_2$  EReLELA* agents with  $\beta_1 = 5$  (agnostic only) or  $\beta_1 = 10$  (shared and agnostic), and  $\beta_2 = 1$ , (iii) the *Impatient-Only EReLELA* agents (shared and agnostic), and (iv) the *RANDOM* agent referring to an ablated version of EReLELA without RG training.

343 **4.1 EReLELA learns Systematic Navigational & Manipulative Exploration Skills from**  
 344 **Scratch**

345 We present in Figure 4 both the success rate of the different agents (as line plot through learning -left-,  
 346 or barplot at the end of learning -right-), and the per-episode manipulation count (middle). From  
 347 the fact that both the NLA and EReLELA agent performance converges higher or close to 80% of  
 348 success rate (except the **STGS-LazImpa-10-1**), we validate hypotheses (H1) and (H2), meaning that  
 349 it is possible to learn systematic exploration skills from both NL or EL abstractions with a simple  
 350 count-based exploration method, in 2D environments (cf. further evidence in Appendix D.1 with the  
 351 *MultiRoom-S7-R4* environment). This result puts into perspective the directions of previous literature  
 352 designing complex exploration algorithms [9, 1].

353 The sample-efficiency is better for NLA than it is for most EL-based agents, except the **Agnostic**  
 354 **STGS-LazImpa-10-1 agent**, possibly because of the fact that ELs are learned online in parallel of the  
 355 RL training, as opposed to the case of NLA which makes use of a ready-to-use oracle. Concerning  
 356 the most-sample-efficient **Agnostic STGS-LazImpa-10-1 agent**, we interpret its success to be the  
 357 result of benefiting from both a language structure ascribing to the ZLA and a performed abstraction  
 358 that is more optimal than NL oracle’s ones, because it is learned from the stimuli themselves.

359 Among the different Agnostic EReLELA agents, the final performance are not statistically-  
 360 significantly distinguishable, meaning that learning systematic exploration skills with EReLELA can  
 361 be done with some robustness to the anecdotal differences in qualities of the different ELs. On the  
 362 other hand, the shared/non-agnostic EReLELA agents’s performance are statistically-significantly  
 363 distinguishable from each other and from their agnostic versions, achieving lower performance or  
 364 even failing to learn anything in the case of the **STGS-LazImpa-10-1 EReLELA agent**. We interpret  
 365 these results as being caused by some kind of interference between the RG training and the RL  
 366 training, preventing any valuable representations from being learned in the shared observation encoder  
 367 (cf. Figure 3), thus warranting the need for future works to investigate whether a synergy can be  
 368 achieved.

369 Finally, acknowledging the **RANDOM agent**, which is the ablated version of EReLELA without  
 370 RG training, enabling still a median performance around 70% of success rate, we recall the Random  
 371 Network Distillation approach from Burda et al. [9], for they both share a randomly initialised  
 372 networked from which feedback is harvested to guide an RL agent. Thus, even more so in a 2D  
 373 environment, this ablated version is not to be confused with a lower-bound baseline but rather an  
 374 interesting ablation that enables us to show the impact of the RG training, increasing the sample-  
 375 efficiency and final performance of the resulting RL agent.

376 **4.2 EReLELA learns Meaningful Abstractions**

377 Regarding hypothesis (H3), we show in Figure 5 the CAM distances between the different agent’s  
 378 ELs and the natural, colour-specific, and shape-specific languages. We recall that in the *KeyCorridor-*  
 379 *S3-R2* environment, the most important feature is object shape as the agent must pickup a key from

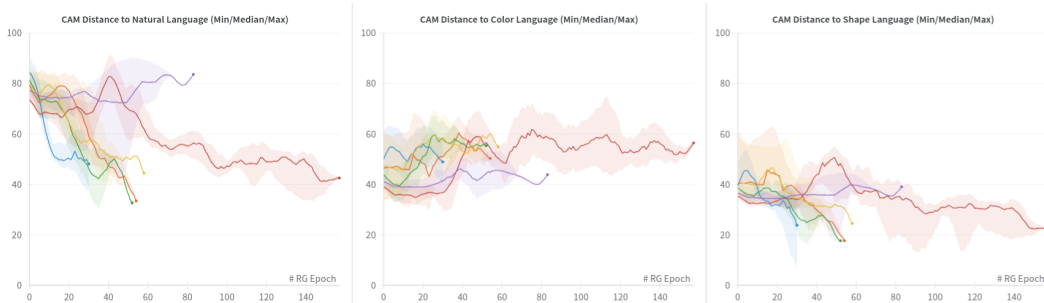


Figure 5: CAM distances to NL (left), Color language (middle), and Shape language (right), for ELs brought about in *KeyCorridor-S3-R2* from MiniGrid [15], with different agents: (i) the **STGS-LazImpa- $\beta_1$ - $\beta_2$  EReLELA** agents with  $\beta_1 = 5$  (**agnostic only**) or  $\beta_1 = 10$  (**shared** and **agnostic**), and  $\beta_2 = 1$ , (ii) the **Impatient-Only EReLELA** agents (**shared** and **agnostic**), and (iii) the **RANDOM agent** referring to an ablated version of EReLELA without RG training.



380 all other distractor objects and then use it to unlock the locked door. Thus, as we observe that  
381 most ELs’ abstractions are closer to the shape-specific language than the others, we conclude that  
382 EReLELA learns meaningful abstractions, thus validating hypothesis (H3) (cf. Appendix E.3 for  
383 further evidence in the context of *MultiRoom-N7-S4*). Further, we remark that the failing **STGS-**  
384 **LazImpa-10-1 EReLELA agent** is indeed failing because its EL’s abstractions are not highlighting  
385 shape features. When considering the shared/non-agnostic agents only, we can see that they require  
386 many more RG training epochs, meaning that they reach the accuracy threshold less often than their  
387 agnostic counterparts. We take this as further evidence for our interpretation that there might be  
388 interference between the RL objective and the RG objective.

389 We note that abstractions from ELs brought about in the contexts of the *Agnostic STGS-LazImpa*  
390 *agents* and the *Agnostic Impatient-Only agents* are the closest to that of the shape-specific language  
391 ones, and their evolution throughout learning are similar. Yet, the *Agnostic STGS-LazImpa agents*  
392 achieves statistically-significantly better sample-efficiency (cf. Figure 7). We interpret this as being  
393 caused by the ZLA structure of the ELs in the context of the *Agnostic STGS-LazImpa agents*, thus  
394 showing that NL-like structure is impacting the kind of abstractions being performed in ways that are  
395 yet to be unveiled by future works.

396 **Limitations.** With regards to the external validity of EReLELA, we acknowledge that the current  
397 work only addresses a 2D environment and therefore, despite being procedurally-generated, it presents  
398 less challenges to count-based exploration methods than in the context of 3D procedurally-generated  
399 environments. Although we provide some results in Appendix E.3 showing that EReLELA is able  
400 to learn meaningful abstractions in a 3D environment, we leave it to future work to ascertain the  
401 external validity of EReLELA by testing it in a procedurally-generated 3D environment that pose  
402 purely-navigational or navigational and manipulative exploration challenges.

## 403 5 Discussion

404 We investigated the compacting/clustering hypothesis for ELs, questioning how do NLS and ELs  
405 compare in terms of the abstractions they perform over state/observation spaces. To answer this  
406 question, we proposed a novel metric entitled Compactness Ambiguity Metric (CAM), for which we  
407 analysed the sensitivity and performed internal validation.

408 We then leveraged this metric to show that ELs abstractions are more meaningful than NLS ones,  
409 as the Emergent Communication context successfully picks up on the meaningful features of the  
410 environment.

411 Then, we have proposed the **Exploration in Reinforcement Learning via Emergent Languages**  
412 **Abstractions (EReLELA)** agent, which leverages ELs abstractions to generate intrinsic motivation  
413 rewards for an RL agent to learn systematic exploration skills. Our experimental evidences showed  
414 the performance of EReLELA in procedurally-generated, hard-exploration 2D environments from  
415 MiniGrid [15].

416 Moreover, in the parallel optimization of the RG players, we evidenced how the STGS-LazImpa loss  
417 function, which induces EL to abide by ZLA like most NLS, impacts the kind of abstraction being  
418 performed compared to baseline Impatient-Only loss function, and yields better sample-efficiency for  
419 the RL agent training.

420 Future work ought to investigate different loss functions and distractor sampling schemes, especially  
421 if playing discriminative RGs like here, as we expect, for instance, that sampling distractors more  
422 contrastively, e.g. like in Choi et al. [17], may induce the emergence of more complete, and therefore  
423 more meaningful ELs. By complete, we mean that the ELs would still be abstracting away details but  
424 also capturing more information about the underlying structure of the stimuli space, e.g. capturing  
425 both colour- and shape-related information of visible objects. In this light, we would also expect  
426 generative RGs to propose a possibly different picture that is worth investigating.

427 While we leave it to subsequent work to investigate the external validity of EReLELA and whether  
428 it transfers similarly well to 3D environments, our results open the door to a new application  
429 of the principles of Emergent Communication and ELs towards influencing/shaping the learned  
430 representations and behaviours of Embodied AI agents trained with RL.

431 **References**

- 432 [1] A. P. Badia, P. Sprechmann, A. Vitvitskyi, D. Guo, B. Piot, S. Kapturowski, O. Tieleman,  
433 M. Arjovsky, A. Pritzel, A. Bolt, et al. Never give up: Learning directed exploration strategies.  
434 In *International Conference on Learning Representations*, 2019.
- 435 [2] M. Baroni. Linguistic generalization and compositionality in modern artificial neural networks.  
436 mar 2019. URL <http://arxiv.org/abs/1904.00157>.
- 437 [3] M. Bellemare, S. Srinivasan, G. Ostrovski, T. Schaul, D. Saxton, and R. Munos. Unifying  
438 count-based exploration and intrinsic motivation. *Advances in neural information processing*  
439 *systems*, 29, 2016.
- 440 [4] L. Biewald. Experiment tracking with weights and biases, 2020. URL <https://www.wandb.com/>.  
441 Software available from wandb.com.
- 442 [5] B. Bogin, M. Geva, and J. Berant. Emergence of Communication in an Interactive World with  
443 Consistent Speakers. sep 2018. URL <http://arxiv.org/abs/1809.00549>.
- 444 [6] D. Bouchacourt and M. Baroni. How agents see things: On visual representations in an emergent  
445 language game. aug 2018. URL <http://arxiv.org/abs/1808.10696>.
- 446 [7] N. Brandizzi. Towards more human-like AI communication: A review of emergent communica-  
447 tion research. Aug. 2023.
- 448 [8] H. Brighton. Compositional syntax from cultural transmission. *MIT Press, Artificial*, 2002.  
449 URL <https://www.mitpressjournals.org/doi/abs/10.1162/106454602753694756>.
- 450 [9] Y. Burda, H. Edwards, D. Pathak, A. Storkey, T. Darrell, and A. A. Efros. Large-Scale Study of  
451 Curiosity-Driven Learning. aug 2018. URL <http://arxiv.org/abs/1808.04355>.
- 452 [10] R. Chaabouni, E. Kharitonov, E. Dupoux, and M. Baroni. Anti-efficient encoding in emergent  
453 communication. *NeurIPS*, may 2019. URL <http://arxiv.org/abs/1905.12561>.
- 454 [11] R. Chaabouni, E. Kharitonov, A. Lazaric, E. Dupoux, and M. Baroni. Word-order biases in deep-  
455 agent emergent communication. may 2019. URL <http://arxiv.org/abs/1905.12330>.
- 456 [12] R. Chaabouni, E. Kharitonov, D. Bouchacourt, E. Dupoux, and M. Baroni. Compositionality  
457 and Generalization in Emergent Languages. apr 2020. URL [http://arxiv.org/abs/2004.](http://arxiv.org/abs/2004.09124)  
458 09124.
- 459 [13] M. S. Charikar. Similarity estimation techniques from rounding algorithms. In *Proceedings of*  
460 *the thirty-fourth annual ACM symposium on Theory of computing*, pages 380–388, 2002.
- 461 [14] R. T. Q. Chen, X. Li, R. Grosse, and D. Duvenaud. Isolating sources of disentanglement in  
462 VAEs, 2018.
- 463 [15] M. Chevalier-Boisvert, B. Dai, M. Towers, R. de Lazcano, L. Willems, S. Lahlou, S. Pal, P. S.  
464 Castro, and J. Terry. Minigrid & miniworld: Modular & customizable reinforcement learning  
465 environments for goal-oriented tasks. *CoRR*, abs/2306.13831, 2023.
- 466 [16] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and  
467 Y. Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine  
468 translation. *arXiv preprint arXiv:1406.1078*, 2014.
- 469 [17] E. Choi, A. Lazaridou, and N. de Freitas. Compositional Obverter Communication Learning  
470 From Raw Visual Input. apr 2018. URL <http://arxiv.org/abs/1804.02341>.
- 471 [18] K. Denamganāi, S. Missaoui, and J. A. Walker. Visual referential games further the emergence  
472 of disentangled representations. *arXiv preprint arXiv:2304.14511*, 2023.
- 473 [19] K. Denamganāi and J. A. Walker. Referentialgym: A nomenclature and framework for language  
474 emergence & grounding in (visual) referential games. *4th NeurIPS Workshop on Emergent*  
475 *Communication*, 2020.

- 476 [20] K. Denamganai and J. A. Walker. Referentialgym: A framework for language emergence &  
477 grounding in (visual) referential games. *4th NeurIPS Workshop on Emergent Communication*,  
478 2020.
- 479 [21] K. Denamganai and J. A. Walker. On (emergent) systematic generalisation and compositionality  
480 in visual referential games with straight-through gumbel-softmax estimator. *4th NeurIPS*  
481 *Workshop on Emergent Communication*, 2020.
- 482 [22] R. Dessi, E. Kharitonov, and M. Baroni. Interpretable agent communication from scratch (with  
483 a generic visual processor emerging on the side). May 2021.
- 484 [23] T. Eccles, Y. Bachrach, G. Lever, A. Lazaridou, and T. Graepel. Biases for emergent communi-  
485 cation in multi-agent reinforcement learning. Dec. 2019.
- 486 [24] S. Guo, Y. Ren, S. Havrylov, S. Frank, I. Titov, and K. Smith. The emergence of compositional  
487 languages for numeric concepts through iterated learning in neural agents. *arXiv preprint*  
488 *arXiv:1910.05291*, 2019.
- 489 [25] S. Havrylov and I. Titov. Emergence of Language with Multi-agent Games: Learning to  
490 Communicate with Sequences of Symbols. may 2017. URL [http://arxiv.org/abs/1705.](http://arxiv.org/abs/1705.11192)  
491 11192.
- 492 [26] I. Higgins, A. Pal, A. Rusu, L. Matthey, C. Burgess, A. Pritzel, M. Botvinick, C. Blundell,  
493 and A. Lerchner. DARLA: Improving Zero-Shot Transfer in Reinforcement Learning. URL  
494 <https://arxiv.org/pdf/1707.08475.pdf>.
- 495 [27] I. Higgins, N. Sonnerat, L. Matthey, A. Pal, C. P. Burgess, M. Botvinick, D. Hassabis, and  
496 A. Lerchner. SCAN: Learning Abstract Hierarchical Compositional Visual Concepts. jul 2017.  
497 URL <http://arxiv.org/abs/1707.03389>.
- 498 [28] I. Higgins, D. Amos, D. Pfau, S. Racaniere, L. Matthey, D. Rezende, and A. Lerchner. Towards  
499 a Definition of Disentangled Representations. dec 2018. URL [http://arxiv.org/abs/1812.](http://arxiv.org/abs/1812.02230)  
500 02230.
- 501 [29] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):  
502 1735–1780, 1997.
- 503 [30] D. Horgan, J. Quan, D. Budden, G. Barth-Maron, M. Hessel, H. Van Hasselt, and D. Silver.  
504 Distributed prioritized experience replay. *arXiv preprint arXiv:1803.00933*, 2018.
- 505 [31] M. Jaderberg, V. Mnih, W. M. Czarnecki, T. Schaul, J. Z. Leibo, D. Silver, and K. Kavukcuoglu.  
506 Reinforcement learning with unsupervised auxiliary tasks. In *International Conference on*  
507 *Learning Representations*, 2016.
- 508 [32] N. Jaques, A. Lazaridou, E. Hughes, C. Gulcehre, P. A. Ortega, D. Strouse, J. Z. Leibo, and  
509 N. De Freitas. Social influence as intrinsic motivation for multi-agent deep reinforcement  
510 learning. *arXiv preprint arXiv:1810.08647*, 2018.
- 511 [33] Y. Jiang, S. Gu, K. Murphy, and C. Finn. Language as an Abstraction for Hierarchical Deep  
512 Reinforcement Learning. jun 2019. URL <http://arxiv.org/abs/1906.07343>.
- 513 [34] S. Kapturowski, G. Ostrovski, J. Quan, R. Munos, and W. Dabney. Recurrent experience replay  
514 in distributed reinforcement learning. In *International conference on learning representations*,  
515 2018.
- 516 [35] H. Kim and A. Mnih. Disentangling by factorising. *arXiv preprint arXiv:1802.05983*, 2018.
- 517 [36] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint*  
518 *arXiv:1412.6980*, 2014.
- 519 [37] D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*,  
520 2013.
- 521 [38] S. Kirby. Learning, bottlenecks and the evolution of recursive syntax. 2002.

- 522 [39] T. Korbak, J. Zubek, Ł. Kuciński, P. Miłoś, and J. Rączaszek-Leonardi. Developmentally  
523 motivated emergence of compositional communication via template transfer. oct 2019. URL  
524 <http://arxiv.org/abs/1910.06079>.
- 525 [40] S. Kottur, J. M. F. Moura, S. Lee, and D. Batra. Natural Language Does Not Emerge 'Naturally'  
526 in Multi-Agent Dialog. jun 2017. URL <http://arxiv.org/abs/1706.08502>.
- 527 [41] A. Lazaridou and M. Baroni. Emergent Multi-Agent communication in the deep learning era.  
528 June 2020.
- 529 [42] A. Lazaridou, A. Peysakhovich, and M. Baroni. Multi-Agent Cooperation and the Emergence  
530 of (Natural) Language. dec 2016. URL <http://arxiv.org/abs/1612.07182>.
- 531 [43] A. Lazaridou, K. M. Hermann, K. Tuyls, and S. Clark. Emergence of Linguistic Communication  
532 from Referential Games with Symbolic and Pixel Input. apr 2018. URL <http://arxiv.org/abs/1804.03984>.
- 534 [44] D. Lewis. Convention: A philosophical study. 1969.
- 535 [45] F. Li and M. Bowling. Ease-of-Teaching and Language Structure from Emergent Communica-  
536 tion. jun 2019. URL <http://arxiv.org/abs/1906.02403>.
- 537 [46] F. Locatello, S. Bauer, M. Lucic, G. Rätsch, S. Gelly, B. Schölkopf, and O. Bachem. A sober  
538 look at the unsupervised learning of disentangled representations and their evaluation. Oct.  
539 2020.
- 540 [47] R. Lowe, J. Foerster, Y.-L. Boureau, J. Pineau, and Y. Dauphin. On the Pitfalls of Measuring  
541 Emergent Communication. mar 2019. URL <http://arxiv.org/abs/1903.05168>.
- 542 [48] M. L. Montero, C. J. Ludwig, R. P. Costa, G. Malhotra, and J. Bowers. The role of disentangl-  
543 e-ment in generalisation. In *International Conference on Learning Representations*, 2021. URL  
544 <https://openreview.net/forum?id=qbH974jKUVy>.
- 545 [49] I. Mordatch and P. Abbeel. Emergence of Grounded Compositional Language in Multi-Agent  
546 Populations. URL <https://arxiv.org/pdf/1703.04908.pdf>.
- 547 [50] K. Moritz Hermann, F. Hill, S. Green, F. Wang, R. Faulkner, H. Soyer, D. Szepesvari, W. M.  
548 Czarnecki, M. Jaderberg, D. Teplyashin, M. Wainwright, C. Apps, D. Hassabis, P. Blunsom,  
549 and D. London. Grounded Language Learning in a Simulated 3D World. URL <https://arxiv.org/pdf/1706.06551.pdf>.
- 551 [51] J. Mu, V. Zhong, R. Raileanu, M. Jiang, N. Goodman, T. Rocktäschel, and E. Grefenstette.  
552 Improving intrinsic exploration with language abstractions. *Advances in Neural Information*  
553 *Processing Systems*, 35:33947–33960, 2022.
- 554 [52] G. Ostrovski, M. G. Bellemare, A. Oord, and R. Munos. Count-based exploration with neural  
555 density models. In *International conference on machine learning*, pages 2721–2730. PMLR,  
556 2017.
- 557 [53] P.-Y. Oudeyer and F. Kaplan. How can we define intrinsic motivation? In *the 8th International*  
558 *Conference on Epigenetic Robotics: Modeling Cognitive Development in Robotic Systems*. Lund  
559 University Cognitive Studies, Lund: LUCS, Brighton, 2008.
- 560 [54] R. Raileanu and T. Rocktäschel. Ride: Rewarding impact-driven exploration for procedurally-  
561 generated environments. In *International Conference on Learning Representations*, 2019.
- 562 [55] Y. Ren, S. Guo, M. Labeau, S. B. Cohen, and S. Kirby. Compositional Languages Emerge in a  
563 Neural Iterated Learning Model. feb 2020. URL <http://arxiv.org/abs/2002.01365>.
- 564 [56] M. Rita, R. Chaabouni, and E. Dupoux. "lazimpa": Lazy and impatient neural agents learn to  
565 communicate efficiently. *arXiv preprint arXiv:2010.01878*, 2020.
- 566 [57] K. Smith, S. Kirby, H. B. A. Life, and U. 2003. Iterated learning: A framework for the emergence  
567 of language. *Artificial Life*, 9(4):371–389, 2003. URL <https://www.mitpressjournals.org/doi/abs/10.1162/106454603322694825>.

- 569 [58] C. Stanton and J. Clune. Deep curiosity search: Intra-life exploration can improve performance  
570 on challenging deep reinforcement learning problems. *arXiv preprint arXiv:1806.00553*, 2018.
- 571 [59] X. Steenbrugge, S. Leroux, T. Verbelen, and B. Dhoedt. Improving generalization for abstract  
572 reasoning tasks using disentangled feature representations. Nov. 2018.
- 573 [60] U. Strauss, P. Grzybek, and G. Altmann. *Word length and word frequency*. Springer, 2007.
- 574 [61] A. Tam, N. Rabinowitz, A. Lampinen, N. A. Roy, S. Chan, D. Strouse, J. Wang, A. Banino,  
575 and F. Hill. Semantic exploration from language abstractions and pretrained representations.  
576 *Advances in Neural Information Processing Systems*, 35:25377–25389, 2022.
- 577 [62] H. Tang, R. Houthoofd, D. Foote, A. Stooke, X. Chen, Y. Duan, J. Schulman, F. De Turck, and  
578 P. Abbeel. Exploration: A study of count-based exploration for deep reinforcement learning.  
579 arxiv e-prints, page. *arXiv preprint arXiv:1611.04717*, 2016.
- 580 [63] S. van Steenkiste, F. Locatello, J. Schmidhuber, and O. Bachem. Are disentangled representa-  
581 tions helpful for abstract visual reasoning? May 2019.
- 582 [64] Z. Wang, T. Schaul, M. Hessel, H. Hasselt, M. Lanctot, and N. Freitas. Dueling network  
583 architectures for deep reinforcement learning. In *International conference on machine learning*,  
584 pages 1995–2003. PMLR, 2016.
- 585 [65] Z. Xu, M. Niethammer, and C. Raffel. Compositional generalization in unsupervised com-  
586 positional representation learning: A study on disentanglement and emergent language. Oct.  
587 2022.
- 588 [66] G. K. Zipf. *Human behavior and the principle of least effort: An introduction to human ecology*.  
589 Ravenio Books, 2016.

## 590 **NeurIPS Paper Checklist**

### 591 **1. Claims**

592 Question: Do the main claims made in the abstract and introduction accurately reflect the  
593 paper's contributions and scope?

594 Answer: [\[Yes\]](#)

595 Justification: Contribution/Claim # 1, i.e. a comparison between emergent and natural lan-  
596 guages with respect to the kind of abstractions they perform, is substantiated in Section E.1,  
597 where we verify the internal validity of the metric we propose for quantitative compari-  
598 son, and Section E.2 where measures using our proposed metrics on different natural or  
599 emergent languages are presented and discussed. Contribution/Claim # 2, i.e. simple count-  
600 based exploration methods guided by natural or emergent language abstractions are helpful  
601 for exploration in reinforcement learning over hard-exploration, procedurally-generated  
602 environments, is substantiated in Section E.3.

603 Guidelines:

- 604 • The answer NA means that the abstract and introduction do not include the claims  
605 made in the paper.
- 606 • The abstract and/or introduction should clearly state the claims made, including the  
607 contributions made in the paper and important assumptions and limitations. A No or  
608 NA answer to this question will not be perceived well by the reviewers.
- 609 • The claims made should match theoretical and experimental results, and reflect how  
610 much the results can be expected to generalize to other settings.
- 611 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
612 are not attained by the paper.

### 613 **2. Limitations**

614 Question: Does the paper discuss the limitations of the work performed by the authors?

615 Answer: [\[Yes\]](#)

616 Justification: We discuss limitations at the end of Section 4.

617 Guidelines:

- 618 • The answer NA means that the paper has no limitation while the answer No means that  
619 the paper has limitations, but those are not discussed in the paper.
- 620 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 621 • The paper should point out any strong assumptions and how robust the results are to  
622 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
623 model well-specification, asymptotic approximations only holding locally). The authors  
624 should reflect on how these assumptions might be violated in practice and what the  
625 implications would be.
- 626 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
627 only tested on a few datasets or with a few runs. In general, empirical results often  
628 depend on implicit assumptions, which should be articulated.
- 629 • The authors should reflect on the factors that influence the performance of the approach.  
630 For example, a facial recognition algorithm may perform poorly when image resolution  
631 is low or images are taken in low lighting. Or a speech-to-text system might not be  
632 used reliably to provide closed captions for online lectures because it fails to handle  
633 technical jargon.
- 634 • The authors should discuss the computational efficiency of the proposed algorithms  
635 and how they scale with dataset size.
- 636 • If applicable, the authors should discuss possible limitations of their approach to  
637 address problems of privacy and fairness.
- 638 • While the authors might fear that complete honesty about limitations might be used by  
639 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
640 limitations that aren't acknowledged in the paper. The authors should use their best  
641 judgment and recognize that individual actions in favor of transparency play an impor-  
642 tant role in developing norms that preserve the integrity of the community. Reviewers  
643 will be specifically instructed to not penalize honesty concerning limitations.



644 **3. Theory Assumptions and Proofs**

645 Question: For each theoretical result, does the paper provide the full set of assumptions and  
646 a complete (and correct) proof?

647 Answer: [Yes]

648 Justification: Our only theoretical results is found in Appendix C with the full set of  
649 assumptions and a complete and correct proof.

650 Guidelines:

- 651 • The answer NA means that the paper does not include theoretical results.
- 652 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
653 referenced.
- 654 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 655 • The proofs can either appear in the main paper or the supplemental material, but if  
656 they appear in the supplemental material, the authors are encouraged to provide a short  
657 proof sketch to provide intuition.
- 658 • Inversely, any informal proof provided in the core of the paper should be complemented  
659 by formal proofs provided in appendix or supplemental material.
- 660 • Theorems and Lemmas that the proof relies upon should be properly referenced.

661 **4. Experimental Result Reproducibility**

662 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
663 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
664 of the paper (regardless of whether the code and data are provided or not)?

665 Answer: [Yes]

666 Justification: All the information needed to reproduce the main experimental results and  
667 appendices experimental results are discussed both in Sections 3 or 4 for critical (and new)  
668 hyperparameters, and in Appendices G and F for hyperparameters introduced in previous  
669 works.

670 Guidelines:

- 671 • The answer NA means that the paper does not include experiments.
- 672 • If the paper includes experiments, a No answer to this question will not be perceived  
673 well by the reviewers: Making the paper reproducible is important, regardless of  
674 whether the code and data are provided or not.
- 675 • If the contribution is a dataset and/or model, the authors should describe the steps taken  
676 to make their results reproducible or verifiable.
- 677 • Depending on the contribution, reproducibility can be accomplished in various ways.  
678 For example, if the contribution is a novel architecture, describing the architecture fully  
679 might suffice, or if the contribution is a specific model and empirical evaluation, it may  
680 be necessary to either make it possible for others to replicate the model with the same  
681 dataset, or provide access to the model. In general, releasing code and data is often  
682 one good way to accomplish this, but reproducibility can also be provided via detailed  
683 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
684 of a large language model), releasing of a model checkpoint, or other means that are  
685 appropriate to the research performed.
- 686 • While NeurIPS does not require releasing code, the conference does require all submis-  
687 sions to provide some reasonable avenue for reproducibility, which may depend on the  
688 nature of the contribution. For example
  - 689 (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
690 to reproduce that algorithm.
  - 691 (b) If the contribution is primarily a new model architecture, the paper should describe  
692 the architecture clearly and fully.
  - 693 (c) If the contribution is a new model (e.g., a large language model), then there should  
694 either be a way to access this model for reproducing the results or a way to reproduce  
695 the model (e.g., with an open-source dataset or instructions for how to construct  
696 the dataset).

697 (d) We recognize that reproducibility may be tricky in some cases, in which case  
698 authors are welcome to describe the particular way they provide for reproducibility.  
699 In the case of closed-source models, it may be that access to the model is limited in  
700 some way (e.g., to registered users), but it should be possible for other researchers  
701 to have some path to reproducing or verifying the results.

## 702 5. Open access to data and code

703 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
704 tions to faithfully reproduce the main experimental results, as described in supplemental  
705 material?

706 Answer: [Yes]

707 Justification: The open-access code contains a README.md file with sufficient instructions  
708 to faithfully reproduce the main experimental results.

709 Guidelines:

- 710 • The answer NA means that paper does not include experiments requiring code.
- 711 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
712 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 713 • While we encourage the release of code and data, we understand that this might not be  
714 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
715 including code, unless this is central to the contribution (e.g., for a new open-source  
716 benchmark).
- 717 • The instructions should contain the exact command and environment needed to run to  
718 reproduce the results. See the NeurIPS code and data submission guidelines ([https://  
719 nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 720 • The authors should provide instructions on data access and preparation, including how  
721 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 722 • The authors should provide scripts to reproduce all experimental results for the new  
723 proposed method and baselines. If only a subset of experiments are reproducible, they  
724 should state which ones are omitted from the script and why.
- 725 • At submission time, to preserve anonymity, the authors should release anonymized  
726 versions (if applicable).
- 727 • Providing as much information as possible in supplemental material (appended to the  
728 paper) is recommended, but including URLs to data and code is permitted.

## 729 6. Experimental Setting/Details

730 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
731 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
732 results?

733 Answer: [Yes]

734 Justification: All the information needed to reproduce the main experimental results and  
735 appendices experimental results are discussed both in Sections 3 or 4 for critical (and newly-  
736 introduced) hyperparameters, and in Appendices G and F for hyperparameters introduced  
737 in previous works.

738 Guidelines:

- 739 • The answer NA means that the paper does not include experiments.
- 740 • The experimental setting should be presented in the core of the paper to a level of detail  
741 that is necessary to appreciate the results and make sense of them.
- 742 • The full details can be provided either with the code, in appendix, or as supplemental  
743 material.

## 744 7. Experiment Statistical Significance

745 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
746 information about the statistical significance of the experiments?

747 Answer: [Yes]

748 Justification: All plots (barplots or line plots) contains in the title the type of information  
749 about the statistical significance of the experiments (i.e. min/median/max, meaning that the  
750 shaded area reflect the min and max values of the distribution while the bar or line reflects  
751 the median of the distribution).

752 Guidelines:

- 753 • The answer NA means that the paper does not include experiments.
- 754 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
755 dence intervals, or statistical significance tests, at least for the experiments that support  
756 the main claims of the paper.
- 757 • The factors of variability that the error bars are capturing should be clearly stated (for  
758 example, train/test split, initialization, random drawing of some parameter, or overall  
759 run with given experimental conditions).
- 760 • The method for calculating the error bars should be explained (closed form formula,  
761 call to a library function, bootstrap, etc.)
- 762 • The assumptions made should be given (e.g., Normally distributed errors).
- 763 • It should be clear whether the error bar is the standard deviation or the standard error  
764 of the mean.
- 765 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
766 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
767 of Normality of errors is not verified.
- 768 • For asymmetric distributions, the authors should be careful not to show in tables or  
769 figures symmetric error bars that would yield results that are out of range (e.g. negative  
770 error rates).
- 771 • If error bars are reported in tables or plots, The authors should explain in the text how  
772 they were calculated and reference the corresponding figures or tables in the text.

## 773 8. Experiments Compute Resources

774 Question: For each experiment, does the paper provide sufficient information on the com-  
775 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
776 the experiments?

777 Answer: [Yes]

778 Justification: Section F contains sufficient information on the computer resources needed to  
779 reproduce the experiments.

780 Guidelines:

- 781 • The answer NA means that the paper does not include experiments.
- 782 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
783 or cloud provider, including relevant memory and storage.
- 784 • The paper should provide the amount of compute required for each of the individual  
785 experimental runs as well as estimate the total compute.
- 786 • The paper should disclose whether the full research project required more compute  
787 than the experiments reported in the paper (e.g., preliminary or failed experiments that  
788 didn't make it into the paper).

## 789 9. Code Of Ethics

790 Question: Does the research conducted in the paper conform, in every respect, with the  
791 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

792 Answer: [Yes]

793 Justification: The research conducted in the paper conform in every respect with the NeurIPS  
794 Code of Ethics.

795 Guidelines:

- 796 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 797 • If the authors answer No, they should explain the special circumstances that require a  
798 deviation from the Code of Ethics.

- 799 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
800 eration due to laws or regulations in their jurisdiction).

## 801 10. Broader Impacts

802 Question: Does the paper discuss both potential positive societal impacts and negative  
803 societal impacts of the work performed?

804 Answer: [Yes]

805 Justification: The paper contains a Broader Impact discussion in Appendix A.

806 Guidelines:

- 807 • The answer NA means that there is no societal impact of the work performed.
- 808 • If the authors answer NA or No, they should explain why their work has no societal  
809 impact or why the paper does not address societal impact.
- 810 • Examples of negative societal impacts include potential malicious or unintended uses  
811 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
812 (e.g., deployment of technologies that could make decisions that unfairly impact specific  
813 groups), privacy considerations, and security considerations.
- 814 • The conference expects that many papers will be foundational research and not tied  
815 to particular applications, let alone deployments. However, if there is a direct path to  
816 any negative applications, the authors should point it out. For example, it is legitimate  
817 to point out that an improvement in the quality of generative models could be used to  
818 generate deepfakes for disinformation. On the other hand, it is not needed to point out  
819 that a generic algorithm for optimizing neural networks could enable people to train  
820 models that generate Deepfakes faster.
- 821 • The authors should consider possible harms that could arise when the technology is  
822 being used as intended and functioning correctly, harms that could arise when the  
823 technology is being used as intended but gives incorrect results, and harms following  
824 from (intentional or unintentional) misuse of the technology.
- 825 • If there are negative societal impacts, the authors could also discuss possible mitigation  
826 strategies (e.g., gated release of models, providing defenses in addition to attacks,  
827 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
828 feedback over time, improving the efficiency and accessibility of ML).

## 829 11. Safeguards

830 Question: Does the paper describe safeguards that have been put in place for responsible  
831 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
832 image generators, or scraped datasets)?

833 Answer: [NA]

834 Justification: The paper does release data or models that have any risk for misuses.

835 Guidelines:

- 836 • The answer NA means that the paper poses no such risks.
- 837 • Released models that have a high risk for misuse or dual-use should be released with  
838 necessary safeguards to allow for controlled use of the model, for example by requiring  
839 that users adhere to usage guidelines or restrictions to access the model or implementing  
840 safety filters.
- 841 • Datasets that have been scraped from the Internet could pose safety risks. The authors  
842 should describe how they avoided releasing unsafe images.
- 843 • We recognize that providing effective safeguards is challenging, and many papers do  
844 not require this, but we encourage authors to take this into account and make a best  
845 faith effort.

## 846 12. Licenses for existing assets

847 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
848 the paper, properly credited and are the license and terms of use explicitly mentioned and  
849 properly respected?

850 Answer: [NA]

851 Justification: Apart from the environments from MiniGrid [15], the paper does not use  
852 existing assets.

853 Guidelines:

- 854 • The answer NA means that the paper does not use existing assets.
- 855 • The authors should cite the original paper that produced the code package or dataset.
- 856 • The authors should state which version of the asset is used and, if possible, include a  
857 URL.
- 858 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 859 • For scraped data from a particular source (e.g., website), the copyright and terms of  
860 service of that source should be provided.
- 861 • If assets are released, the license, copyright information, and terms of use in the  
862 package should be provided. For popular datasets, `paperswithcode.com/datasets`  
863 has curated licenses for some datasets. Their licensing guide can help determine the  
864 license of a dataset.
- 865 • For existing datasets that are re-packaged, both the original license and the license of  
866 the derived asset (if it has changed) should be provided.
- 867 • If this information is not available online, the authors are encouraged to reach out to  
868 the asset's creators.

### 869 13. New Assets

870 Question: Are new assets introduced in the paper well documented and is the documentation  
871 provided alongside the assets?

872 Answer: [NA]

873 Justification: The paper does not release new assets.

874 Guidelines:

- 875 • The answer NA means that the paper does not release new assets.
- 876 • Researchers should communicate the details of the dataset/code/model as part of their  
877 submissions via structured templates. This includes details about training, license,  
878 limitations, etc.
- 879 • The paper should discuss whether and how consent was obtained from people whose  
880 asset is used.
- 881 • At submission time, remember to anonymize your assets (if applicable). You can either  
882 create an anonymized URL or include an anonymized zip file.

### 883 14. Crowdsourcing and Research with Human Subjects

884 Question: For crowdsourcing experiments and research with human subjects, does the paper  
885 include the full text of instructions given to participants and screenshots, if applicable, as  
886 well as details about compensation (if any)?

887 Answer: [NA]

888 Justification: The paper does not involve experiments with human subjects nor crowdsourc-  
889 ing.

890 Guidelines:

- 891 • The answer NA means that the paper does not involve crowdsourcing nor research with  
892 human subjects.
- 893 • Including this information in the supplemental material is fine, but if the main contribu-  
894 tion of the paper involves human subjects, then as much detail as possible should be  
895 included in the main paper.
- 896 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
897 or other labor should be paid at least the minimum wage in the country of the data  
898 collector.

### 899 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 900 Subjects

901 Question: Does the paper describe potential risks incurred by study participants, whether  
902 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
903 approvals (or an equivalent approval/review based on the requirements of your country or  
904 institution) were obtained?

905 Answer: [NA]

906 Justification: The paper does not involve crowdsourcing nor research with human subjects.

907 Guidelines:

- 908 • The answer NA means that the paper does not involve crowdsourcing nor research with  
909 human subjects.
- 910 • Depending on the country in which research is conducted, IRB approval (or equivalent)  
911 may be required for any human subjects research. If you obtained IRB approval, you  
912 should clearly state this in the paper.
- 913 • We recognize that the procedures for this may vary significantly between institutions  
914 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
915 guidelines for their institution.
- 916 • For initial submissions, do not include any information that would break anonymity (if  
917 applicable), such as the institution conducting the review.



918 **A Broader impact**

919 No technology is safe from being used for malicious purposes, which equally applies to our research.  
920 However, we view many of the ethical concerns surrounding research to be mitigated in the present  
921 case. These include data-related concerns such as fair use or issues surrounding use of human subjects,  
922 given that our data consists solely of simulations.

923 With regards to the ethical aspects related to its inclusion in the field of Artificial Intelligence, we argue  
924 that our work aims to have positive outcomes on the development of human-machine interfaces since  
925 we investigate, among other things, alignment of emergent languages with natural-like languages.

926 The current state of our work does not allow extrapolation towards negative outcomes. We believe  
927 that this work is of benefit to the research community of reinforcement learning, language emergence  
928 and grounding, in their current state.

929 **B Further details on Count-Based Exploration**

930 Another approach to counting states from continuous and/or high-dimensional state spaces is by  
931 relying on hashing functions, so that states become tractable. Indeed, Tang et al. [62] have shown that  
932 a generalisation of classical counting techniques through hashing can provide an appropriate signal  
933 for exploration in continuous and/or high-dimensional environments where informed exploration is  
934 required. In effect, they proposed to discretise the state space  $\mathcal{S}$  with a hash function  $\phi : \mathcal{S} \rightarrow \mathbb{Z}^k$ ,  
935 with  $k \in \mathbb{N} \setminus \{0\}$ , to derive an exploration bonus of the form  $r^+(s) = \frac{\beta}{\sqrt{n(\phi(s))}}$  where  $\beta \in \mathbb{R}^+$  is a  
936 bonus coefficient and  $n(\cdot)$  is a count initialised at zero for the whole range of  $\phi$  and updated at each  
937 step  $t$  of the RL loop by increasing by 1 the count  $n(\phi(s_t))$  related to the current observation/state  
938  $s_t$ . Performance is dependent on the hash function  $\phi$ , and especially in terms of granularity of the  
939 discretisation it induces. Indeed, it would be desirable that the ‘similar’ states result in hashing  
940 collisions while the ‘distant’ states would not. To this end, they propose to use locality-sensitive  
941 hashing (LSH) such as SimHash [13], resulting in the following:

$$\phi(s) = \text{sgn}(Ag(s)) \in \{-1, 1\}^k, \tag{5}$$

942 where  $\text{sgn}$  is the sign function,  $A \in \mathbb{R}^{k \times D}$  is a matrix with each entry drawn i.i.d. from a standard  
943 Gaussian distribution, and  $g : S \rightarrow \mathbb{R}^D$  is an optional preprocessing function. Note that increasing  
944  $k$  leads to higher granularity and therefore decreases the number of hashing collisions. Tang et al.  
945 [62] reports great results on the Atari 2600 benchmarks, both with and without a learnable  $g$  that is  
946 modelled as the encoder of an autoencoder (AE).

947 **C Sensitivity Analysis of the Compactness Ambiguity Metric**

948 Based on derivative-based local sensitivity analysis, we propose an intuitive proof of our claim that  
 949 defining timespans in relation to the relative ambiguity reduces the sensibility to variations induced  
 950 by redundancy-based ambiguity in the resulting metric, compared to defining timespans in relation to  
 951 the the maximal length  $T$  of an agent's trajectory in the environment. To do so, we assume:

- 952 (i) that there exists two differentiable function  $f_i, f'_i$  such that for all  $i \in [1, N]$ , we have  
 953  $CA(\mathcal{D})_{T_i} = f_i(\mathcal{D}, \mathcal{R}\mathcal{A}_l^{\text{redundancy}}, \mathcal{R}\mathcal{A}_l^{\text{abstract}})$  when  $T_i$  is defined according to Equation 2,  
 954 and respectively with  $f'_i$  when using  $T'_i$  from Equation 3, and  
 955 (ii) that their partial derivatives with respect to  $T_i$  or  $T'_i$  are negative. Indeed,  $T_i$  and  $T'_i$  are  
 956 involved into filtering operations reducing the value of the numerator in Equation 4, therefore  
 957 any increase of their values would result in decreasing the overall metric output, which  
 958 implies that their partial derivatives with  $f_i$  and  $f'_i$  must be negative.

959 With those assumptions, we show that  $f_i$ 's sensitivity to redundancy-induced ambiguity  $\mathcal{R}\mathcal{A}_l^{\text{redundancy}}$   
 960 is less than that of  $f'_i$ :

*Proof.*

$$\begin{aligned} \frac{\partial f_i}{\partial \mathcal{R}\mathcal{A}_l^{\text{redundancy}}} &= \frac{\partial f_i}{\partial CC_{\mathcal{D}}} \cdot \frac{\partial CC_{\mathcal{D}}}{\partial \mathcal{R}\mathcal{A}_l^{\text{redundancy}}} + \frac{\partial f_i}{\partial T_i} \cdot \frac{\partial T_i}{\partial \mathcal{R}\mathcal{A}_l^{\text{redundancy}}} && \text{(from Assump. (i) about } f_i) \\ \iff \frac{\partial f_i}{\partial \mathcal{R}\mathcal{A}_l^{\text{redundancy}}} &= \frac{\partial f'_i}{\partial \mathcal{R}\mathcal{A}_l^{\text{redundancy}}} + \frac{\partial f_i}{\partial T_i} \cdot \frac{\partial T_i}{\partial \mathcal{R}\mathcal{A}_l^{\text{redundancy}}} && \text{(from Assump. (i) about } f'_i) \\ \iff \frac{\partial f_i}{\partial \mathcal{R}\mathcal{A}_l^{\text{redundancy}}} &= \frac{\partial f'_i}{\partial \mathcal{R}\mathcal{A}_l^{\text{redundancy}}} + \frac{\partial f_i}{\partial T_i} \cdot \lambda_i \\ \implies \left| \frac{\partial f_i}{\partial \mathcal{R}\mathcal{A}_l^{\text{redundancy}}} \right| &\leq \left| \frac{\partial f'_i}{\partial \mathcal{R}\mathcal{A}_l^{\text{redundancy}}} \right| && \text{(since } \frac{\partial f_i}{\partial T_i} \cdot \lambda_i \leq 0 \text{ from Assump. (ii))} \end{aligned}$$

961

□

## 962 D Preliminary Experiments

### 963 D.1 Impact of Referential Game Accuracy

964 In this experiments, we investigate whether the RG accuracy impacts the RL agent training, in the  
965 context of the *MultiRoom-N7-S4* environment from *MiniGrid* [15], with an RL sampling budget of  
966  $1M$  observations.

967 **Hypothesis.** We seek to validate the following hypotheses, **(PH1)** : the sample-efficiency of the  
968 RL agent is dependant on the quality of the RG players, as parameterised by the  $acc_{RG-thresh}$   
969 hyperparameter.

970 **Evaluation.** We report both the success rate and the coverage count in the hard-exploration task of  
971 *MultiRoom-N7-S4*. To compute the coverage count, we overlay a grid of tiles over the environment’s  
972 possible locations/cells of the agents and we count the number of different tiles visited by the RL  
973 agent over the course of each episode. We use 3 random seeds for each agent. In order to evaluate the  
974 impact of the RG accuracy strictly in terms of the kind of abstractions that are being performed by the  
975 resulting EL, we use the *Impatient-Only* loss function (removing the impact of the hyperparameter of  
976 the scheduling function  $\alpha(\cdot)$  from the *Lazy* term of the *STGS-LazImpa* loss function), and we employ  
977 an **agnostic** version of our proposed EReLELA agent, i.e. **without sharing the observation encoder**  
978 **between the RG players and the RL agent**. We present results for two different RG accuracy  
979 threshold  $acc_{RG-thresh} = 60\%$  (**green**) or  $acc_{RG-thresh} = 80\%$  (**red**), and compare against, as an  
980 upper bound the *Natural Language Abstraction* agent (**blue**), which refers to using the NL oracle to  
981 compute intrinsic reward, and, as a lower bound an ablated version of EReLELA without RG training  
982 (**orange**).

983 **Results.** We present results in Figure 6. We observe statistically significant differences between  
984 the performances (in terms of success rate, cf. Figure 6(left)) of the two EReLELA agents with  
985  $acc_{RG-thresh} = 60\%$  or  $acc_{RG-thresh} = 80\%$ , thus validating hypothesis (PH1). We observe that  
986 higher RG accuracy threshold lead to higher sample-efficiency.

987 As a sanity check, we plot the results of the ablated EReLELA agent without RG training, and we were  
988 expecting it to perform poorer than any other agent since the quality of its RG players is the lowest, at  
989 chance level. Yet, we observe that it performs on par with the best  $acc_{RG-thresh} = 80\%$ -EReLELA  
990 agent. While puzzling, we propose a possible explanation in the observation that the test-time relative  
991 expressivity of the ablated agent is higher than that of the least-performing,  $acc_{RG-thresh} = 60\%$ -  
992 EReLELA agent, and on par with that of the best-performing,  $acc_{RG-thresh} = 80\%$ -EReLELA  
993 agent, at the beginning of the RL agent training process. Thus, we interpret this as follows: the  
994 randomly-initialised ablated agent’s EL is possibly performing an abstraction over the observation

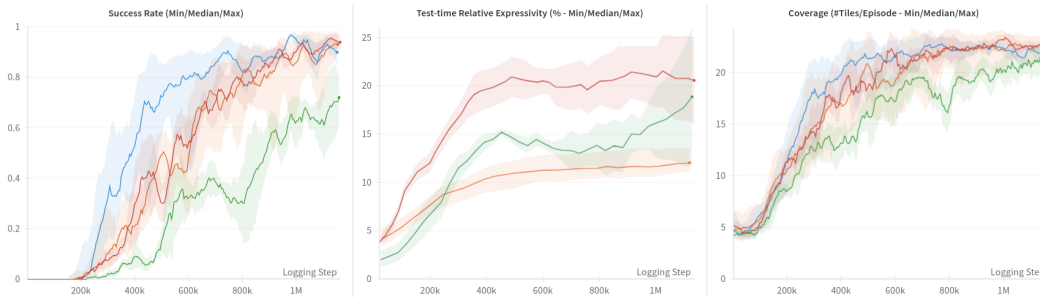


Figure 6: Success rate (left), test-time relative expressivity (middle), and per-episode coverage count (right) in *MultiRoom-N7-S4* from *MiniGrid* [15], computed as running averages over 256 episodes each time (i.e. 32 in parallel, as there are 32 actors, over 8 running average steps), for different agents: (i) the *Natural Language Abstraction* agent (**blue**) refers to using the NL oracle to compute intrinsic reward, the *Agnostic Impatient-Only EReLELA* agent refers to our proposed architecture **without sharing the observation encoder between the RG players and the RL agent**, using the *Impatient-Only* loss function to optimize the RG players, with an RG accuracy threshold  $acc_{RG-thresh} = 60\%$  (ii - **green**) or  $acc_{RG-thresh} = 80\%$  (iii - **red**), and (iv) an ablated version without RG training (**orange**).

995 space that is good-enough for the RL agent to start learning exploration skills, the same way the  
 996 random network in the context of the RND agent from Burda et al. [9] probably does, and increasing  
 997 the quality of the RG players may only be a sufficient condition to increasing the sample-efficiency  
 998 of the EL-guided RL agent.

## 999 D.2 Impact of Referential Game Distractors

1000 In this experiments, we investigate whether the RG’s number of distractors  $K$  and distractor sampling  
 1001 scheme impacts the RL agent training, in the context of the *KeyCorridor-S3-R2* environment from  
 1002 *MiniGrid* [15], with an RL sampling budget of  $1M$  observations.

1003 **Hypothesis.** We seek to validate the following hypotheses, **(PH2)** : the sample-efficiency of the RL  
 1004 agent is dependant on the number of distractors  $K$  and the distractor sampling scheme.

1005 **Evaluation.** We report the success rate in the hard-exploration task of *KeyCorridor-S3-R2*. We  
 1006 use 3 random seeds for each agent. Like previously, we use the *Impatient-Only* loss function (to  
 1007 remove the impact of the hyperparameter of the scheduling function  $\alpha(\cdot)$  from the *Lazy* term of  
 1008 the *STGS-LazImpa* loss function), and we employ an **agnostic** version of our proposed EReLELA  
 1009 agent, i.e. **without sharing the observation encoder between the RG players and the RL agent**.  
 1010 We present results for three different number of distractors  $K \in [15, 128, 256]$  and two different  
 1011 sampling scheme between *UnifDSS* corresponding to uniformly sampling distractors over the whole  
 1012 training dataset, or *Sim50DSS* corresponding to sampling distractors 50% of the time from the same  
 1013 RL episode than the current target stimulus is from and, the rest of the time following *UnifDSS*.  
 1014 Following results in Appendix D.1, we set the RG accuracy threshold  $acc_{RG-thresh} \in [80\%, 90\%]$ .

1015 **Results.** We present results in Figure 7. We observe statistically significant differences between the  
 1016 performances of the different EReLELA agents, thus validating hypothesis (PH2). Our results show  
 1017 that (i) the number of distractors  $K$  is the most impactful parameter and it correlates positively with  
 1018 the resulting performance, irrespective of the distractor sampling scheme used, and, indeed, (ii) while  
 1019 the *Sim50DSS* seems to provide better performance than *UnifDSS* for low numbers of distractors  
 1020  $K = 15$ , although not statistically-significantly, the table is turned when considering high number of  
 1021 distractors  $K = 256$  where the *UnifDSS* yields statistically significantly better performance than the  
 1022 *Sim50DSS*.

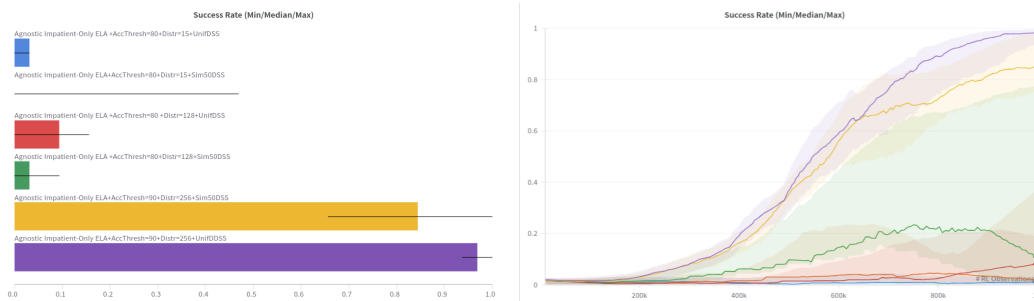


Figure 7: Final success rate barplot (left) and success rate throughout learning (right) in *KeyCorridor-S3-R2* from *MiniGrid* [15], computed as running averages over 1024 episodes each time (i.e. 32 in parallel, as there are 32 actors, over 32 running average steps), for the *Agnostic Impatient-Only EReLELA* agent, which refers to our proposed architecture **without sharing the observation encoder between the RG players and the RL agent**, using the *Impatient-Only* loss function to optimize the RG players, with different number of distractors  $K$  and distractors sampling schemes: with RG accuracy threshold  $acc_{RG-thresh} = 80\%$ , (i)  $K = 15$  and *UnifDSS* or *Sim50DSS*, (ii)  $K = 128$  and *UnifDSS* or *Sim50DSS*, or with RG accuracy threshold  $acc_{RG-thresh} = 90\%$ , (iii)  $K = 256$  and *UnifDSS* or *Sim50DSS*.

1023 **E Further Experiments**

1024 **E.1 Experiment #1: CAM Metric Internal Validity**

1025 **Environment.** We consider a 3D room environment of MiniWorld [15], where the agent’s observation  
 1026 is egocentric, as a first-person viewpoint. The room is filled with 5 different, randomly-placed objects,  
 1027 with different shapes (among ball, box or key) and colours (among). The dimensions simulate a 12  
 1028 by 5 meters room, like shown in a top-view perspective in Figure 1.

1029 **Hypothesis.** In this experiments, we seek to validate two hypotheses, **(H1.1)** : the Compactness  
 1030 Ambiguity Metric captures something that is related to the kind of abstraction a language performs,  
 1031 and **(H1.2)** : the Compactness Ambiguity Metric allows a graduated comparison of different kind  
 1032 of abstractions being performed, meaning that it allows discrimination between different kind of  
 1033 abstractions.

1034 **Evaluation.** In order to compute the metric, we use 5 seeds to gather random walk trajectories in our  
 1035 environment, for each language. In order to evaluate (H1.1), we propose to measure a language that  
 1036 is built to present no meaningful abstractions and we expect the measure to be close to null. We build  
 1037 a language that performs no meaningful abstraction from the natural language oracles by shuffling  
 1038 its utterances over the set of agent trajectories that are used to compute the metric, meaning that  
 1039 the mapping between temporally-sensitive stimuli and linguistic utterances is rendered completely  
 1040 random.

1041 Then, in order to evaluate (H1.2), we show experimental evidences that the metric allows qualitative  
 1042 discrimination between the different languages built above from the natural language oracles, which  
 1043 are build to perform different kind of abstractions.

1044 **Results.** We present results of the metric with  $N = 6$  timespans in Figure 8, for  $\lambda_0 = 0.0306125$ ,  
 1045  $\lambda_1 = 0.06125$ ,  $\lambda_2 = 0.125$ ,  $\lambda_3 = 0.25$ ,  $\lambda_4 = 0.5$  and  $\lambda_5 = 0.75$ . As the shuffled (natural) language  
 1046 measure is almost null on all timespans/thresholds, we validate hypothesis (H1.1).

1047 We observe that we can qualitatively discriminate between each evaluated language’s measures since  
 1048 the histograms are statistically different. Moreover, language abstractions scores are inversely corre-  
 1049 lated with the amount of information being abstracted away, i.e. attribute-value-specific languages’  
 1050 abstraction score lower than colour/shape-specific languages abstraction, which score lower than  
 1051 natural language abstractions. Thus, we can see that the metric is graduated and that the graduation  
 1052 follows the amount of abstraction being performed by each language. This allows us to validate  
 1053 hypothesis (H1.2).



Figure 8: Interval validity measures of Compactness Ambiguity Metric for  $N = 6$  timespans/thresholds, with  $\lambda_0 = 0.0306125$ ,  $\lambda_1 = 0.06125$ ,  $\lambda_2 = 0.125$ ,  $\lambda_3 = 0.25$ ,  $\lambda_4 = 0.5$  and  $\lambda_5 = 0.75$ , for different languages built to perform different kind of abstraction. We can qualitatively discriminate between each languages, and validate that the shuffled (natural) language’s meaningless abstraction scores almost null.

1054 **E.2 Experiment #2: Qualities of Emergent Languages Abstractions in 3D environment**

1055 In this experiment, we investigate what kind of abstractions do ELs perform over a 3D environment,  
 1056 in comparison to some natural languages abstractions, as detailed at the beginning of Section 4. For  
 1057 further precision, we also implement attribute-value-specific language oracles with the same filtering  
 1058 approach. For instance, for the green value on the colour attribute, we would obtain a green-only  
 1059 language oracle whose utterances could be ‘EoS’ if no visible object is green, or ‘green green’ if there  
 1060 are two green objects visible in the agent’s observation. We consider the same 3D room environment  
 1061 of MiniWorld [15] as in Section E.1, i.e. the agent’s observation is egocentric, as a first-person  
 1062 viewpoint and the room is filled with 5 different, randomly-placed objects, with different shapes  
 1063 (among ball, box or key) and colours (among). The dimensions simulate a 12 by 5 meters room, like  
 1064 shown in a top-view perspective in Figure 1.

1065 **Hypothesis.** We seek to validate the following hypotheses, **(H2.1)** : ELs build meaningful abstractions,  
 1066 and **(H2.2)** : ELs brought about using the STGS-LazImpa loss function (type II) perform more  
 1067 meaningful abstractions than Impatient-Only baseline (type I).

1068 **Evaluation.** In order to make the CAM measures, we use 5 seeds to gather random walk trajectories  
 1069 in our environment, for each language. In order to evaluate both (H2.1) and (H2.2), we use the CAM  
 1070 to measure the kind of abstractions performed by ELs brought about in the two different EReLELA  
 1071 settings, with Impatient-Only or STGS-LazImpa losses, and compare those measures with those of  
 1072 the oracles’ languages that we previously studied.

1073 **Results.** We present results of the metric with  $N = 6$  timespans in Figure 9. We observe statistically  
 1074 significant differences between ELs of type I and II, with type I’s abstraction being similar to a Blue-  
 1075 specific language’s abstraction (timespans 0 – 4) or a Ball-specific language’s abstraction (timespans  
 1076 1 – 3), and type II’s abstraction not really resembling any of the oracle languages’ abstractions, but  
 1077 still being meaningful with scores increasing along with the length of the considered timespans. Thus,  
 1078 we validate hypothesis (H2.1), but cannot conclude on hypothesis (H2.2), unless we consider that  
 1079 CAM scores related to longer timespans are more meaningful, for instance.

1080 **E.3 Experiment #3: Learning Purely-Navigational Systematic Exploration Skills from**  
 1081 **Scratch**

1082 In the following, we present an experiment in the *MultiRoom-N7-S4* environment from *MiniGrid* [15],  
 1083 which is possibly less challenging than *KeyCorridor-S3-R2*, presented in the Section 4, for it does  
 1084 not involve as many complex object manipulation (e.g. only open/close doors, no unlocking of  
 1085 doors – which requires the corresponding key to be firstly picked up – nor pickup/drop keys or  
 1086 other objects as distractors), but still poses a **purely-navigational** hard-exploration challenge. We  
 1087 report results on the **agnostic** version of our proposed EReLELA architecture, that is to say **without**  
 1088 **sharing the observation encoder between both RG players and the RL agent**, in order to guard  
 1089 ourselves against the impact of possible confounders found in multi-task optimization, such as possible



Figure 9: Measures of Compactness Ambiguity Metric for  $N = 6$  timespans/thresholds, with  $\lambda_0 = 0.0306125$ ,  $\lambda_1 = 0.06125$ ,  $\lambda_2 = 0.125$ ,  $\lambda_3 = 0.25$ ,  $\lambda_4 = 0.5$  and  $\lambda_5 = 0.75$ , comparing ELs (Type I and II) with different oracles’ languages built to perform different kind of abstraction.



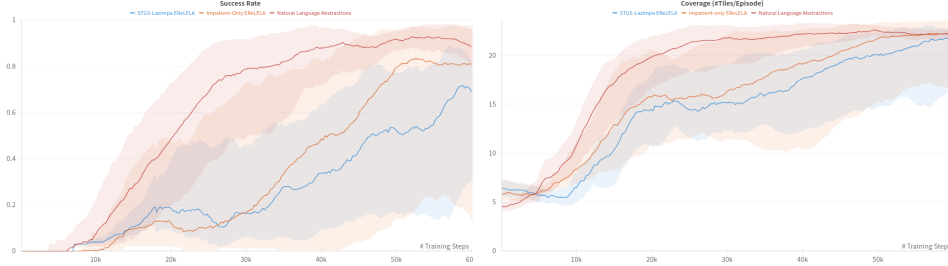


Figure 10: Success rate (left) and per-episode coverage count (right) in *MultiRoom-N7-S4* from MiniGrid [15], computed as running averages over 1024 episodes each time (i.e. 32 in parallel, as there are 32 actors, over 32 running average steps), for different agents: (i) the *Natural Language Abstraction* agent (NLA) refers to using the NL oracle to compute intrinsic reward, (ii) the *STGS-LazImpa EReLELA* agent refers to our proposed architecture, EReLELA, using the STGS-LazImpa loss function to optimize the RG players, and (iii) the *Impatient-Only EReLELA* agent refers to the same architecture without the lazy-speaker loss to optimize the RG players.

1090 interference between the RL-objective-induced gradients and the RG-training-induced gradients. We  
 1091 use an RG accuracy threshold  $acc_{RG-thresh} = 65\%$  and a number of training distractors  $K = 3$   
 1092 (like at testing/validation time).

1093 **Hypotheses.** We consider whether NL abstractions can help for a purely-navigational hard-  
 1094 exploration task in RL with a count-based approach (**H3.0**), and refer to the relevant agent using  
 1095 NL abstractions to compute intrinsic rewards as NLA. Then, we make the hypothesis that ELs can  
 1096 be used similarly (**H3.1**), and we investigate to what extent do ELs compare to NLA in terms of  
 1097 abstraction performed, in this purely-navigational task. In the case of (H3.1) being verified, we would  
 1098 expect ELs to perform similar abstractions as NLA (**H3.2**).

1099 **Evaluation.** We evaluate (H3.0) and (H3.1) using both the success rate and the coverage count. To  
 1100 compute the coverage count, we overlay a grid of tiles over the environment’s possible locations/cells  
 1101 of the agents and we count the number of different tiles visited by the RL agent over the course of  
 1102 each episode. To evaluate (H3.2), we compute the CAM scores of both the ELs and the oracles’  
 1103 natural, color-specific, and shape-specific languages. As we remarked that an agent’s skillfulness at  
 1104 the task would induce very different trajectories (e.g. in *MultiRoom-N7-S4*, staying in the first room  
 1105 and only ever seeing the first door, for an unskillfull agent, as opposed to visiting multiple rooms  
 1106 and observing multiple colored-doors, for a skillfull agent), we compute the oracle languages CAM  
 1107 scores on the exact same trajectories than used to compute each EL’s CAM scores.

1108 **Results.** We present in Figure 10(left) the success rate of the different agents, and the per-episode  
 1109 coverage count in Figure 10(right). From the fact that both the NLA and EReLELA agent performance  
 1110 converges higher or close to 80% of success rate, we validate hypotheses (H0) and (H3.1), in the  
 1111 context of the *MultiRoom-N7-S4* environment. We remark that the sample-efficiency is slightly better  
 1112 for NLA than it is for EL-based agents, possibly because of the fact that ELs are learned online  
 1113 in parallel of the RL training, as opposed to the case of NLA which makes use of a ready-to-use  
 1114 oracle. Among the two EReLELA agents, the learning curves are not statistically-significantly

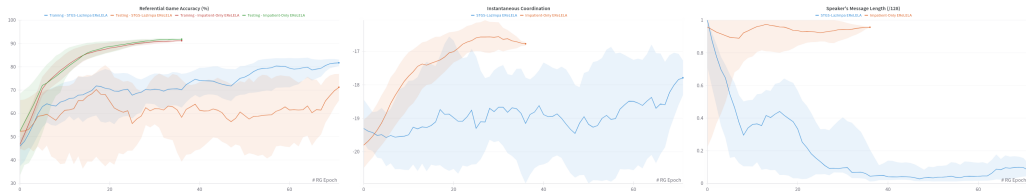


Figure 11: Performance and qualities of the ELs brought about in the context of both (i) the *STGS-LazImpa EReLELA* agent, and (ii) the *Impatient-Only EReLELA* agent, with respect to both the training- and validation/testing-time RG accuracy (left), the validation/test-time Instantaneous Coordination [32, 47, 23](middle), and the validation/testing-time length of the speaker’s messages (as a ratio over the max sentence length  $L = 128$  - right).

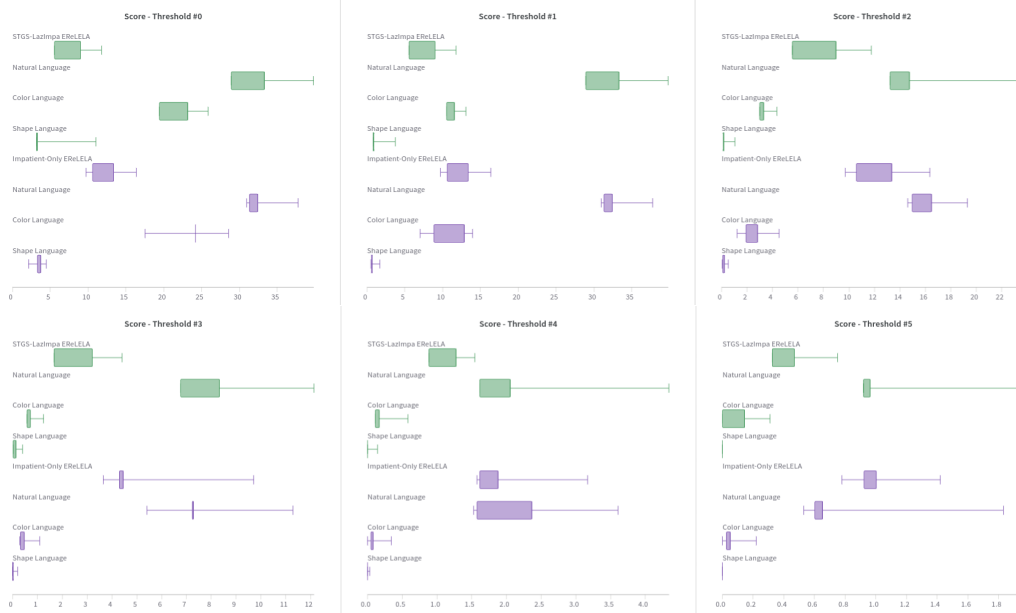


Figure 12: Comparison of Compactness Ambiguity Metric scores for  $N = 6$  timespans/thresholds, with  $\lambda_0 = 0.0306125$ ,  $\lambda_1 = 0.06125$ ,  $\lambda_2 = 0.125$ ,  $\lambda_3 = 0.25$ ,  $\lambda_4 = 0.5$  and  $\lambda_5 = 0.75$ , between the abstractions performed by ELs brought about in the context of both (i) the *STGS-LazImpa EReLELA* agent (in green, first rows) and (ii) the *Impatient-Only EReLELA* agent (in purple, bottom rows), and the abstractions performed by the natural, colour-specific, and shape-specific languages, computed on the very same agent trajectories.

1115 distinguishable, meaning that learning systematic exploration skills with EReLELA can be done with  
 1116 some robustness to the anecdotal differences in qualities of the different ELs due to using different  
 1117 optimization losses. Indeed, we also report in Figure 11 both the training- and validation/testing-time  
 1118 RG accuracies (on the left), the validation/testing-time Instantaneous Coordination (in the middle –  
 1119 Jaques et al. [32], Lowe et al. [47], Eccles et al. [23]), and the validation/testing-time length of the RG  
 1120 speaker’s messages (on the right), showing that the ELs brought about in the two different contexts  
 1121 perform differently in terms of their RG objective and have different qualities, but these discrepancies  
 1122 do not seem to impact the RL agents learning equally well from the different abstractions they  
 1123 perform (as evidenced in the next paragraph).

1124 Next, with regards to hypothesis (H3.2), we investigate whether the two contexts bring about ELs  
 1125 that perform different abstractions, and how do these relate to the abstractions performed by natural,  
 1126 colour-specific, and shape-specific languages, by showing in Figure 12 their CAM scores. We  
 1127 observe that both contexts result in ELs performing abstractions similar or better than colour-specific  
 1128 languages, which is to be expected as (door) colours are the most salient features of the environment.  
 1129 Indeed, the only two shapes or objects visible are ‘wall’ and ‘door’, whereas there are more than  
 1130 7 different colours of interest. In the context of the Impatient-Only EReLELA agent, the EL’s  
 1131 abstractions are scoring very similarly to NL abstractions, as we consider longer timespans (from  
 1132 timespans #2 to #5). We could hypothesise that without the lazy-ness constraint the speaker agent  
 1133 may be given enough capacity to compress/express information pertaining to the location of visible  
 1134 objects, as this information is the only one that is captured by the NL oracle but not captured by the  
 1135 shape- and colour-specific languages.

#### 1136 E.4 Experiment #4: Quantifying RL Agents’ Learning Progress?

1137 In the context of RGs, the speed at which a language emerges (in terms of sampled observations, or  
 1138 number of games played) may possibly remain constant, when the data and the player architectures  
 1139 are fixed. Thus, when the data changes, the rate of language emergence may change too. Incidentally,  
 1140 we are entitled to ponder whether some properties of the data, which here are RL trajectories, would  
 1141 influence the rate of language emergence and how?

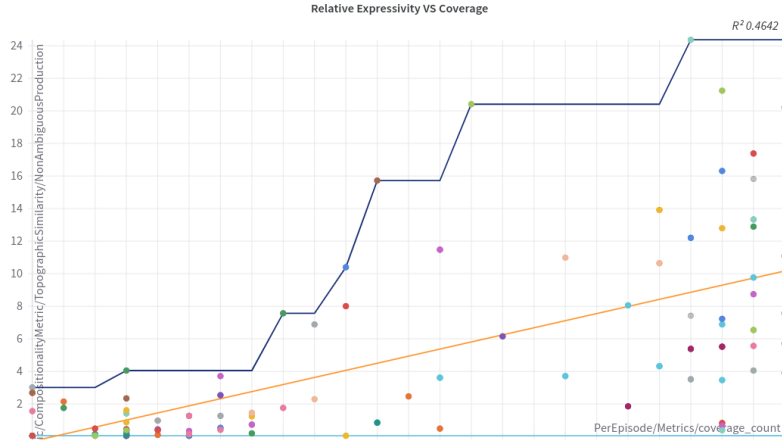


Figure 13: Relative expressivity of the EL as a function of the per-episode coverage of the RL agent, at the end of training, over multiple runs with different hyperparameters during a W&B Sweep [4].

1142 **Hypothesis.** We hypothesise that as the RL agent gets more skillful, the expressivity of the emergent  
 1143 language increases (**H4.1**). Indeed, at each RG training epoch, the size of the dataset is fixed, and as  
 1144 the stimuli gets more diverse when the RL agent gets more skillful at exploring, the RG training will  
 1145 prompt the EL to increase its expressivity.

1146 **Evaluation.** To verify our hypothesis, we propose to measure the skillfulness of the RL agent in  
 1147 terms of exploration using the per-episode coverage count metric, and we measure the expressivity of  
 1148 the EL via the test-time (Relative) Expressivity after each RG training epoch.

1149 **Results.** We present results in Figure 13, that show the (relative) expressivity of the ELs does exhibit  
 1150 variations throughout the learning process of the RL agent. And, if we perform a regression analysis  
 1151 with each runs in terms of the per-episode coverage count of the RL agent on the x-axis and the  
 1152 expressivity of the ELs on the y-axis, we obtain a high coefficient of determination between the two  
 1153 metrics,  $R^2 = 0.4642$ . Thus, we conclude that the (relative) expressivity of the ELs in EReLELA can  
 1154 provide a way to quantify the progress of the RL agent, at least when it comes to exploration skills.

1155 **Limitations.** Exploration skills translates directly into diversity of the stimuli being observed, and  
 1156 therefore it prompts any RG players to increase the expressivity of their communication protocol,  
 1157 but it remains to be seen whether this effect is valid in any environment. For instance, it is unclear  
 1158 whether a skillful player in any other video game would induce the same effect on the diversity of  
 1159 the stimuli encountered. Thus, it is worth investigating whether this correlation holds for other genre  
 1160 of environments and skills, which we leave to future works.

1161 **F Agent Architecture**

1162 The ERELELA architecture is made up of three differentiable agents, the language-conditioned RL  
 1163 agent and the two RG agents (speaker and listener). Each agent contains at least a visual/observation  
 1164 encoder module that can be shared between agents. Both RG agents contain a language module that is  
 1165 not shared. The *listener* agent additionally incorporates a third decision module that combines the  
 1166 outputs of the other two modules. The RL agent similarly incorporates a third decision module with  
 1167 the addition that this third module contains a recurrent network, acting as core memory module for  
 1168 the agent. Using the Straight-Through Gumbel-Softmax (STGS) approach in the communication  
 1169 channel of the RG, the *speaker* agent is prompted to produce the output string of symbols with a  
 1170 *Start-of-Sentence* symbol and the visual module’s output as an initial hidden state while the *listener*  
 1171 agent consumes the string of symbols with the null vector as the initial hidden state. In the following  
 1172 subsections, we detail each module architecture in depth.

1173 **Visual Module.** The visual module  $f(\cdot)$  consists of the *Shared Observation Encoder*, which can be  
 1174 shared between all the different agents. The former consists of three blocks of convolutional layers  
 1175 of sizes 8, 4, 3 with strides 4, 3, 1, each followed by a 2D batch normalization layer and a ReLU  
 1176 non-linear activation function. The two first convolutional layers have 32 filters, whilst the last one  
 1177 has 64. The bias parameters of the convolutional layers are not used, as it is common when using  
 1178 batch normalisation layers. Inputs are stimuli consisting of RGB frames of the environment resized  
 1179 to  $64 \times 64$ .

1180 **Language Module.** The language module  $g(\cdot)$  consists of some learned Embedding followed by  
 1181 either a one-layer GRU network [16] in the case of the RL agent, or a one-layer LSTM network [29]  
 1182 in the case of the RG agents. In the context of the *listener* agent, the input message  $m = (m_i)_{i \in [1, L]}$   
 1183 (produced by the *speaker* agent) is represented as a string of one-hot encoded vectors of dimension  
 1184  $|V|$  and embedded in an embedding space of dimension 64 via a learned Embedding. The output  
 1185 of the *listener* agent’s language module,  $g^l(\cdot)$ , is the last hidden state of the RNN layer,  $h_L^l =$   
 1186  $g^l(m_L, h_{L-1}^l)$ . In the context of the *speaker* agent’s language module  $g^s(\cdot)$ , the output is the  
 1187 message  $m = (m_i)_{i \in [1, L]}$  consisting of one-hot encoded vectors of dimension  $|V|$ , which are sampled  
 1188 using the STGS approach from a categorical distribution  $Cat(p_i)$  where  $p_i = \text{Softmax}(\nu(h_i^s))$ ,  
 1189 provided  $\nu$  is an affine transformation and  $h_i^s = g^s(m_{i-1}, h_{i-1}^s)$ .  $h_0^s = f(s_t)$  is the output of the  
 1190 visual module, given the target stimulus  $s_t$ .

1191 **Decision Module.** From the RL agent to the RG’s listener agent, the decision module are very  
 1192 different since their outputs are either, respectively, in the action space  $\mathcal{A}$  or the space of distributions  
 1193 over  $K + 1$  stimuli (i.e. discriminating between distractors and target stimuli). For the RL agent, the  
 1194 decision module takes as input a concatenated vector comprising the output of visual module, after  
 1195 it has been processed by a 3-layer fully-connected network with 256, 128 and 64 hidden units with  
 1196 ReLU non-linear activation functions, and some other information relevant to the RL context (e.g.  
 1197 previous reward and previous action selected, following the recipe in Kapturowski et al. [34]). The  
 1198 resulting concatenated vector is then fed to the core memory module, a one-layer LSTM network [29]  
 1199 with 1024 hidden units, which feeds into the advantage and value heads of a 1-layer dueling network  
 1200 [64].

1201 In the case of the RG’s listener agent, similarly to Havrylov and Titov [25], the decision module  
 1202 builds a probability distribution over a set of  $K + 1$  stimuli/images  $(s_0, \dots, s_K)$ , consisting of  $K$   
 1203 distractor stimuli and the target stimulus, provided in a random order, given a message  $m$  using the  
 1204 scalar product:

$$p((d_i)_{i \in [0, K]} | (s_i)_{i \in [0, K]}; m) = \text{Softmax}\left((h_L^l \cdot f(s_i)^T)_{i \in [0, K]}\right). \quad (6)$$

1205 Regarding optimization of the RL agent, table 1 highlights the hyperparameters used for the off-policy  
 1206 RL algorithm, R2D2[34]. More details can be found, for reproducibility purposes, in our open-source  
 1207 implementation at HIDDEN-FOR-REVIEW-PURPOSES.

1208 Each run can be done on less than 2Gb of VRAM, and the amount of training time for a run, with e.g.  
 1209 one NVIDIA GTX1080 Ti, is between 24 and 48 hours depending on the architecture (e.g. shared or  
 1210 agnostic).

Table 1: Hyper-parameter values relevant to R2D2 in the EReLELA architecture presented. All missing parameters follow the ones in Ape-X [30].

R2D2	
Number of actors	32
Actor update interval	1 env. step
Sequence unroll length	20
Sequence length overlap	10
Sequence burn-in length	10
N-steps return	3
Replay buffer size	$1 \times 10^4$ obs.
Priority exponent	0.9
Importance sampling exponent	0.6
Discount $\gamma$	0.98
Minibatch size	64
Optimizer	Adam [36]
Learning rate	$6.25 \times 10^{-5}$
Adam $\epsilon$	$10^{-12}$
Target network update interval	2500
Value function rescaling	updates
	None

## 1211 G On the Referential Game in EReLELA

1212 We follow the nomenclature proposed in Denamganai and Walker [20] and focus on a *descrip-*  
 1213 *tive object-centric (partially-observable) 2-players/L = 10-signal/N = 0-round/K-distractor* RG  
 1214 variant.

1215 The descriptiveness implies that the target stimulus may not be passed to the listener agent, but  
 1216 instead replaced with a descriptive distractor. In effect, the listener agent’s decision module therefore  
 1217 outputs a  $K + 2$ -logit distribution where the  $K + 2$ -th logit represents the meaning/prediction that a  
 1218 descriptive distractor has been introduced and none of the  $K + 1$  stimuli is the target stimulus that  
 1219 the speaker agent was ‘talking’ about. The addition is made following Denamganai et al. [18] as a  
 1220 learnable logit value,  $logit_{no-target}$ , it is an extra parameter of the model. In this case the decision  
 1221 module output is no longer as specified in Equation 6, but rather as follows:

$$p((d_i)_{i \in [0, K+1]} | (s_i)_{i \in [0, K]}; m) = \text{Softmax}\left((h_L^l \cdot f(s_i)^T)_{i \in [0, K]} \cup \{logit_{no-target}\}\right). \quad (7)$$

1222 The descriptiveness is ideal but not necessary in order to employ the listener agent as a predicate  
 1223 function for the hindsight experience replay scheme. Thus, in the main results of the paper, we  
 1224 present the version without descriptiveness.

1225 The object-centrism is achieved via application of data augmentation schemes before feeding stimuli  
 1226 to any RG agent, following Dessi et al. [22] but using Gaussian Blur transformation alone, as it was  
 1227 found sufficient in practice.

1228 We optimize the RG agents with either the Impatient-Only STGS loss and the STGS-LazImpa loss.

1229 In the remainder of this section, we detail the STGS-LazImpa loss that we employed to optimize the  
 1230 referential game agents.

### 1231 G.1 STGS-LazImpa Loss

1232 Emergent languages rarely bears the core properties of natural languages [40, 6, 43, 12], such as  
 1233 Zipf’s law of Abbreviation (ZLA). In the context of natural languages, this is an empirical law which  
 1234 states that the more frequent a word is, the shorter it tends to be [66, 60]. Rita et al. [56] proposed  
 1235 LazImpa in order to make emergent languages follow ZLA.

1236 To do so, Lazimpa adds to the speaker and listener agents some constraints to make the speaker  
 1237 lazy and the listener impatient. Thus, denoting those constraints as  $\mathcal{L}_{STGS-lazy}$  and  $\mathcal{L}_{impatient}$ , we  
 1238 obtain the STGS-LazImpa loss as follows:

$$\mathcal{L}_{STGS-LazImpa}(m, (s_i)_{i \in [0, K]}) = \mathcal{L}_{STGS-lazy}(m) + \mathcal{L}_{impatient}(m, (s_i)_{i \in [0, K]}). \quad (8)$$

1239 In the following, we detail those two constraints.

1240 **Lazy Speaker.** The Lazy Speaker agent has the same architecture as common speakers. The  
 1241 ‘Laziness’ is originally implemented as a cost on the length of the message  $m$  directly applied to the  
 1242 loss, of the following form:

$$\mathcal{L}_{lazy}(m) = \alpha(acc) \cdot |m| \quad (9)$$

1243 where  $acc$  represents the current accuracy estimates of the referential games being played, and  $\alpha$   
 1244 is a scheduling function as follows:  $\alpha : accuracy \in [0, 1] \mapsto \frac{accuracy^{\beta_1}}{\beta_2}$ , with  $(\beta_1, \beta_2) = (45, 10)$ .  
 1245 It is aimed to adaptively penalize depending on the message length. Since the laziness loss is  
 1246 not differentiable, they ought to employ a REINFORCE-based algorithm for the purpose of credit  
 1247 assignement of the speaker agent.

1248 In this work, we use the STGS communication channel, which has been shown to be more sample-  
 1249 efficient than REINFORCE-based algorithms [25], but it requires the loss functions to be differen-  
 1250 tiable. Therefore, we modify the laziness loss by taking inspiration from the variational autoencoders  
 1251 (VAE) literature [37].

1252 The length of the speaker’s message is controlled by the appearance of the EoS token, wherever  
 1253 it appears during the message generation process that is where the message is complete and its  
 1254 length is fixed. Symbols of the message at each position are sampled from a distribution over all  
 1255 the tokens in the vocabulary that the listener agent outputs. Let  $(W_l)$  be this distribution over all  
 1256 tokens  $w \in V$  at position  $l \in [1, L]$ , such that  $\forall l \in [1, L], m_l \sim (W_l)$ . We devise the laziness loss  
 1257 as a Kullbach-Leibler divergence  $D_{KL}(\cdot|\cdot)$  between these distribution and the distribution  $(W_{EoS})$   
 1258 which attributes all its weight on the EoS token. Thus, we dissuade the listener agent from outputting  
 1259 distributions over tokens that deviate too much from the EoS-focused distribution  $(W_{EoS})$ , at each  
 1260 position  $l$  with varying coefficients  $\beta(l)$ . The coefficient function  $\beta : [1, L] \rightarrow \mathbb{R}$  must be monotonically  
 1261 increasing. We obtain our STGS-lazyness loss as follows:

$$\mathcal{L}_{STGS-lazy}(m) = \alpha(acc) \cdot \sum_{l \in [1, L]} \beta(l) D_{KL}((W_{EoS})|(W_l)) \quad (10)$$

1262 **Impatient Listener.** Our implementation of the Impatient Listener agent follows the original work  
 1263 of Rita et al. [56]: it is designed to guess the target stimulus as soon as possible, rather than solely  
 1264 upon reading the EoS token at the end of the speaker’s message  $m$ . Thus, following Equation 6, the  
 1265 Impatient Listener agent outputs a probability distribution over a set of  $K + 1$  stimuli  $(s_0, \dots, s_K)$  for  
 1266 all sub-parts/prefixes of the message  $m = (m_1, \dots, m_l)_{l \in [1, L]} = (m_{\leq l})_{l \in [1, L]}$ :

$$\forall l \in [1, L], p((\mathbf{d}_i^{\leq l})_{i \in [0, K]} | (s_i)_{i \in [0, K]}; \mathbf{m}^{\leq l}) = \text{Softmax}\left(\left(\mathbf{h}_{\leq l} \cdot f(s_i)^T\right)_{i \in [0, K]}\right), \quad (11)$$

1267 where  $\mathbf{h}_{\leq l}$  is the hidden state/output of the recurrent network in the language module after consuming  
 1268 tokens of the message from position 1 to position  $l$  included.

1269 Thus, we obtain a sequence of  $L$  probability distributions, which can each be contrasted, using the  
 1270 loss of the user’s choice, against the target distribution  $(D_{target})$  attributing all its weights on the  
 1271 decision  $d_{target}$  where the target stimulus was presented to the listener agent. Here, we employ  
 1272 Havrylov and Titov [25]’s Hinge loss. Denoting it as  $\mathbb{L}(\cdot)$ , we obtain the impatient loss as follows:

$$\mathcal{L}_{impatient/\mathbb{L}}(m, (s_i)_{i \in [0, K]}) = \frac{1}{L} \sum_{l \in [1, L]} \mathbb{L}((d_{i \in [0, K]}^{\leq l}, (D_{target})). \quad (12)$$