META-REFERENTIAL GAMES TO LEARN COMPOSI-TIONAL LEARNING BEHAVIOURS

Anonymous authors

Paper under double-blind review

## ABSTRACT

Human beings use compositionality to generalise from past to novel experiences, assuming that past experiences can be decomposed into fundamental atomic components that can be recombined in novel ways. We frame this as the ability to learn to generalise compositionally, and refer to behaviours making use of this ability as compositional learning behaviours (CLBs). Learning CLBs requires the resolution of a binding problem (BP). While it is another feat of intelligence that human beings perform with ease, it is not the case for artificial agents. Thus, in order to build artificial agents able to collaborate with human beings, we develop a novel benchmark to investigate agents' abilities to exhibit CLBs by solving a domain-agnostic version of the BP. Taking inspiration from the Emergent Communication, we propose a meta-learning extension of referential games, entitled Meta-Referential Games, to support our benchmark, the Symbolic Behaviour Benchmark (S2B). Baseline results and error analysis show that the S2B is a compelling challenge that we hope will spur the research community to develop more capable artificial agents.

023 024 025

026 027

004

010 011

012

013

014

015

016

017

018

019

021

### 1 INTRODUCTION

Defining compositional behaviours (CBs) as "the ability to generalise from combinations of **trained**-028 on atomic components to novel re-combinations of those very same components", we can define 029 compositional learning behaviours (CLBs) as "the ability to generalise in an online fashion from a few combinations of never-before-seen atomic components to novel re-combinations of those very 031 same components". We employ the term online here to imply a few-shot learning context (Vinyals et al., 2016; Mishra et al., 2018) that demands that agents learn from, and then leverage some novel 033 information, both over the course of a single lifespan, or episode, in our case of few-shot meta-RL 034 (see Beck et al. (2023) for a review of meta-RL). Thus, in this paper, we investigate artificial agents' abilities for CLBs, which involve a few-shot learning aspect that is not present in CBs. For an 035 intuitive comparison, benchmarks that tests for CBs alone ask whether trained-and-now-frozen agents can generalise to novel combinations of those same **familiar** atomic components they have been 037 training on. This is different from what a benchmark testing for CLBs would ask, to wit, whether trained-and-now-frozen agents can generalise to novel combinations of **never-before-seen** atomic components, provided some warm-up exposition rounds to those novel atomic components. Such 040 a benchmark instantiates a meta-learning challenge where tested agents can only be successful if 041 they have learned to learn to generalise compositionally, whereas CB-testing benchmark, like 042 SCAN (Lake & Baroni, 2018) and gSCAN (Ruis et al., 2020), can be solved by agents that have only 043 learned to generalise compositionally. Thus, in order to prompt agents to acquire the skill of learning 044 to generalise compositionally, we rely on the instantiation of a theoretically-infinite distribution of different atomic components, and the underlying semantic structure they belong to, so that at each training episode our agents are prompted with novel atomic components and novel underlying 046 semantic structure. 047

Compositional Learning Behaviours as Symbolic Behaviours. Santoro et al. (2021) states that
 a symbolic entity does not exist in an objective sense but solely in relation to an *"interpreter who treats it as such"*, and it ensues that there exists a set of behaviours, i.e. *symbolic behaviours*, that
 are consequences of agents engaging with symbols. Thus, in order to evaluate artificial agents in
 terms of their ability to collaborate with humans, we can use the presence or absence of symbolic
 behaviours. Among the different characteristic of symbolic behaviours, this work will primarily focus
 on the receptivity and constructivity aspects. Receptivity aspects amount to the ability to receive

- new symbolic conventions in an online fashion. For instance, when a child introduces an adult to
  their toys' names, the adults are able to discriminate between those new names upon the next usage.
  Constructivity aspects amount to the ability to form new symbolic conventions in an online fashion.
  For instance, when facing novel situations that require collaborations, two human teammates can
  come up with novel referring expressions to easily discriminate between different events occurring.
  Both aspects refer to abilities that support collaboration. Thus, this paper develops a benchmark to
  evaluate agents' abilities in receptive and constructive behaviours, with a primary focus on CLBs.
- 061 Binding Problem & Meta-Learning. Following Greff et al. (2020), we refer to the binding problem 062 (BP) as the challenges in "dynamically and flexibly bind[/re-use] information that is distributed 063 throughout the [architecture]" of some artificial agents (modelled with artificial neural networks here). 064 We note that there is an inherent BP that requires solving for agents to exhibit CLBs. Indeed, over the course of a single episode (as opposed to a whole training process, in the case of CBs), agents 065 must dynamically identify/segregate the component values from the observation of multiple stimuli, 066 timestep after timestep, and then bind/(re-)use/(re-)combine this information (hopefully stored in 067 some memory component of their architecture) in order to respond correctly to novel stimuli.Solving 068 the BP instantiated in such a context, i.e. re-using previously-acquired information in ways that 069 serve the current situation, is another feat of intelligence that human beings perform with ease, on the contrary to current state-of-the-art artificial agents. Thus, our benchmark must emphasise 071 testing agents' abilities to exhibit CLBs by solving a version of the BP. Moreover, we argue for a 072 domain-agnostic BP, i.e. not grounded in a specific modality such as vision or audio, as doing so 073 would limit the external validity of the test. We aim for as few assumptions as possible to be made 074 about the nature of the BP we instantiate (Chollet, 2019). This is crucial to motivate the form of the 075 stimuli we employ, and we will further detail this in Section 3.2.
- 076 Language Grounding & Emergence. In order to test the quality of some symbolic behaviours, our 077 proposed benchmark needs to query the semantics that agents (the interpreters) may extract from their experience, and it must be able to do so in a referential fashion (e.g. being able to query to 079 what extent a given experience is referred to as, for instance, 'the sight of a red tomato'), similarly to most language grounding benchmarks. Subsequently, acknowledging that the simplest form of 081 collaboration is maybe the exchange of information, i.e. communication, via a given code, or language, we argue that the benchmark must therefore also allow agents to manipulate this code/language that they use to communicate. This property is known as the metalinguistic/reflexive function of 083 languages (Jakobson, 1960). It is mainly investigated in the current deep learning era within the field of Emergent Communication (Lazaridou & Baroni (2020), and see Brandizzi (2023) and 085 Denamganaï & Walker (2020a) for further reviews), via the use of variants of the referential games 086 (RGs) (Lewis, 1969). Thus, we take inspiration from the RG framework, where (i) the language 087 domain represents a semantic domain that can be probed and queried, and (ii) the reflexive function of 880 language is indeed addressed. Then, in order to instantiate different BPs at each episode, we propose 089 a meta-learning extension to RGs, entitled Meta-Referential Games, and use this framework to build 090 our benchmark. It results in our proposed Symbolic Behaviour Benchmark (S2B), which has the 091 potential to test for many aspects of symbolic behaviours.
- After review of the background (Section 2), we will present our contributions as follows: we propose the Symbolic Behaviour Benchmark to enables evaluation of symbolic behaviours in Section 3, presenting the Symbolic Continuous Stimulus (SCS) representation scheme which is able to instantiate a BP, on the contrary to common symbolic representations (Section 3.2), and our Meta-Referential Games framework, a meta-learning extension to RGs (Section 3.1); then we provide baseline results and error analysis in Section 4 for state-of-the-art RL agents and Large Language Models (LLMs -Brown et al. (2020), showing that our benchmark is a compelling challenge that we hope will spur the research community.
- 100 101

# 2 BACKGROUND

102 103

The first instance of an environment with a primary focus on efficient communication is the *signaling game* or *referential game* (RG) by Lewis (1969), where a speaker agent is asked to send a message to the listener agent, based on the *state/stimulus* of the world that it observed. The listener agent then acts upon the observed message by choosing one of the *actions* available to it. Both players' goals are aligned (it features *pure coordination/common interests*), with the aim of performing the

2

108 'best' action given the observed state. In the recent deep learning era, many variants of the RG 109 have appeared (Lazaridou & Baroni, 2020). Following the nomenclature proposed in Denamganaï 110 & Walker (2020b), Figure 1 illustrates in the general case a discriminative 2-players / L-signal / 111 N-round / K-distractors / descriptive / object-centric variant, where the speaker receives a stimulus 112 and communicates with the listener (up to N back-and-forth using messages of at most L tokens each), who additionally receives a set of K + 1 stimuli (potentially including a semantically-similar 113 stimulus as the speaker, referred to as an object-centric stimulus). The task is for the listener 114 to determine, via communication with the speaker, whether any of its observed stimuli match the 115 speaker's. We highlight here features of RGs that will be relevant to how S2B is built, and then provide 116 formalism used throughout the paper. The **number of communication rounds** N characterises 117 (i) whether the listener agent can send messages back to the speaker agent and (ii) how many 118 communication rounds can be expected before the listener agent is finally tasked to decide on an action. 119

The basic (discriminative) RG is stimulus-120 centric, which assumes that both agents 121 would be somehow embodied in the same 122 body, and they are tasked to discriminate 123 between given stimuli, that are the results of one single perception 'system'. On 124 the other hand, Choi et al. (2018) intro-125 duced an object-centric variant which in-126 corporates the issues that stem from the 127 difference of embodiment (which has been 128 later re-introduced under the name Concept 129 game by Mu & Goodman (2021)). The 130 agents must discriminate between objects 131 (or scenes) independently of the viewpoint



Figure 1: Illustration of a *discriminative 2-players / L-signal / N-round* variant of a *RG*.

132 from which they may experience them. In the object-centric variant, the game is more about bridging 133 the gap between each other's cognition rather than just finding a common language. The adjective 134 'object-centric' is used to qualify a stimulus that is different from another but actually present the same meaning (e.g. same object, but seen under a different viewpoint). Following the last communi-135 cation round, the listener outputs a decision  $(D_i^L \text{ in Figure 2})$  about whether any of the stimulus it 136 is observing matches the one (or a semantically similar one, in object-centric RGs) experienced by 137 the speaker, and if so its action index must represent the index of the stimulus it identifies as being 138 the same. The **descriptive** variant allows for none of the stimuli to be the same as the target one, 139 therefore the action of index 0 is required for success. The agent's ability to make the correct decision 140 over multiple RGs is referred to as RG accuracy. 141

Compositionality, Disentanglement & Systematicity. Compositionality is a phenomenon that 142 human beings are able to identify and leverage thanks to the assumption that reality can be decomposed 143 over a set of "disentangle[d,] underlying factors of variations" (Bengio, 2012), and our experience 144 is a noisy, entangled translation of this factorised reality. This assumption is critical to the field 145 of unsupervised learning of disentangled representations (Locatello et al., 2020) that aims to find 146 "manifold learning algorithms" (Bengio, 2012), such as variational autoencoders (VAEs (Kingma & 147 Welling, 2013)), with the particularity that the latent encoding space would consist of disentangled 148 latent variables (see Higgins et al. (2018) for a formal definition). As a concept, compositionality 149 has been the focus of many definition attempts. For instance, it can be defined as "the algebraic 150 capacity to understand and produce novel combinations from known components" (Loula et al. (2018) 151 referring to Montague (1970)) or as the property according to which "the meaning of a complex expression is a function of the meaning of its immediate syntactic parts and the way in which they are 152 combined" (Krifka, 2001). Although difficult to define, the community seems to agree on the fact 153 that it would enable learning agents to exhibit systematic generalisation abilities (also referred to as 154 combinatorial generalisation (Battaglia et al., 2018)). While often studied in relation to languages, it is 155 usually defined with a focus on behaviours. In this paper, we will refer to (linguistic) compositionality 156 when considering languages, and interchangeably compositional behaviours or systematicity to refer 157 to "the ability to entertain a given thought implies the ability to entertain thoughts with semantically 158 related contents" (Fodor & Pylyshyn, 1988). 159

160 Compositionality can be difficult to measure. Brighton & Kirby (2006)'s *topographic similarity* 

(**topsim**) which is acknowledged by the research community as the main quantitative metric (Lazaridou et al., 2018; Guo et al., 2019; Słowik et al., 2020; Chaabouni et al., 2020; Ren et al., 2020). 162 Recently, taking inspiration from disentanglement metrics, Chaabouni et al. (2020) proposed the 163 posdis (positional disentanglement) and bosdis (bag-of-symbols disentanglement) metrics, that 164 have been shown to be differently 'opinionated' when it comes to what kind of compositionality 165 they capture. As hinted at by Choi et al. (2018); Chaabouni et al. (2020) and Dessi et al. (2021), 166 compositionality and disentanglement appears to be two sides of the same coin, in as much as emergent languages are discrete and sequentially-constrained unsupervisedly-learned representations. 167 In Section 3.2, we bridge further compositional language emergence and unsupervised learning of 168 disentangled representations by asking what would an ideally-disentangled latent space look like? to build our proposed benchmark. 170

Richness of the Stimuli & Systematicity. Chaabouni et al. (2020) found that compositionality is not 171 172 necessary to bring about systematicity, as shown by the fact that non-compositional languages wielded by symbolic (generative) RG players were enough to support success in zero-shot compositional 173 tests (ZSCTs). They found that the emergence of a posdis-compositional language was a sufficient 174 condition for systematicity to emerge. Finally, they found a necessary condition to foster systematicity, 175 that we will refer to as richness of stimuli condition (Chaa-RSC). It was framed as (i) having a large 176 stimulus space  $|I| = i_{val}^{i_{attr}}$ , where  $i_{attr}$  is the number of attributes/factor dimensions, and  $i_{val}$ 177 is the number of possible values on each attribute/factor dimension, and (ii) making sure that it 178 is densely sampled during training, in order to guarantee that different values on different factor 179 dimensions have been experienced together. In a similar fashion, Hill et al. (2019) also propose a 180 richness of stimuli condition (Hill-RSC) that was framed as a data augmentation-like regularizer 181 caused by the egocentric viewpoint of the studied embodied agent. In effect, the diversity of viewpoint 182 allowing the embodied agent to observe over many perspectives the same and unique semantical 183 meaning allows a form of contrastive learning that promotes the agent's systematicity.

184 185

# **3** Symbolic Behaviour Benchmark

186 187 188

We present a version of the S2B that focuses on evaluating receptive and constructive behaviour traits<sup>1</sup>. This evaluation relies on a single task built around 2-players multi-agent RL (MARL) episodes. Each episode consists of a series of RGs (cf. lines 11 and 17 in Alg. 5 calling Alg. 3). We denote one such episode as a meta-RG and detail it in Section 3.1. Each RG in a meta-RG follows the formalism laid out in Section 2. The only difference is that both players speak simultaneously, rather than turn-by-turn. Consequently, they observe their partner's message upon the next RL step. Each RG consists of N + 2 RL steps, where N is the *number of communication rounds* available to the players (cf. Section 2).

At each RL step, each player observe a stimulus and a message coming from their partner. Throughout the first N + 1 steps, these stimuli remain constant. They are *object-centric*, and may be similar or different from one player to the other. We recall that the common goal of RG players is to figure out whether they observe similar stimuli or not. Stimuli are represented using the the Symbolic Continuous Stimulus (SCS) representation, which we detail in Section 3.2.

From those observations, each player acts from different actions spaces. Each player's action space is dependent on their role in the RG, as the speaker or the listener. In the first N steps, both players' action space only allow them to send messages to their partner. Then, at step N + 1, the listener player must decide whether they observe a similar stimulus than the speaker. Thus, the listener's action space only allows the decision-related actions. And, the speaker's action space only allows a *no-operation* (NO-OP) action. In practice, the environment provides the players with masks that identify valid actions. If the players still choose invalid actions, the environment simply ignores them.

Finally, after the listener has provided their decision, step N + 2 provides feedback to the listener player. This feedback consist of two elements. First, the environment reward is non-null. And, secondly, the listener's stimulus is the same that the speaker was observing throughout the current RG (cf. line 12 and 18 in Alg. 5). Note, that it is the exact stimulus rather than an object-centric sample. In Figure 4, we present SCS-represented stimuli, observed by a speaker over the course of a typical episode.

214 215

<sup>&</sup>lt;sup>1</sup>HIDDEN\_FOR\_REVIEW\_PURPOSE



224 225

226

230

Figure 2: Left: Sampling of the necessary components to create the i-th RG  $(RG_i)$  of a meta-RG. The target stimulus (red) and the object-centric target stimulus (purple) are both sampled from the Target Distribution  $TD_i$ , a set of O different stimuli representing the same latent semantic 227 meaning. The latter set and a set of K distractor stimuli (orange) are both sampled from a dataset of 228 SCS-represented stimuli (Dataset), which is instantiated from the current episode's symbolic space, 229 whose semantic structure is sampled out of the meta-distribution of available semantic structure over  $N_{dim}$ -dimensioned symbolic spaces. Right: Illustration of the resulting meta-RG with a focus 231 on the i-th RG  $RG_i$ . The speaker agent receives at each step the target stimulus  $s_0^i$  and distractor 232 stimuli  $(s_k^i)_{k \in [1:K]}$ , while the listener agent receives an object-centric version of the target stimulus 233  $s_0^{\prime i}$  or a distractor stimulus (randomly sampled), and other distractor stimuli  $(s_k^i)_{k \in [1:K]}$ , with the 234 exception of the Listener Feedback step where the listener agent receives feedback in the form 235 of the exact target stimulus  $s_0^i$ . The Listener Feedback step takes place after the listener agent has 236 provided a decision  $D_i^L$  about whether the target meaning is observed or not and in which stimuli is 237 it instantiated, guided by the vocabulary-permutated message  $M_i^S$  from the speaker agent.

238 239 240

#### 3.1 META-REFERENTIAL GAMES

241 A Meta-RG consists of a series of RGs, played over specific stimuli. These stimuli are sampled from 242 a given  $N_{\text{dim}}$ -dimensioned symbolic space. To be successful, players must learn to communicate 243 about stimuli using initially-ungrounded symbols. They must ground symbols into the symbolic 244 space they observe. But each new episode brings about a new communication channel and a new 245 symbolic space. The symbolic space is renewed by changing its semantic structure, which we detail 246 further in Section 3.2. And the communication channel has its symbols being randomly permutated, 247 as detailed below. Players of a Meta-RG must learn to adapt to new stimuli and new communication 248 channels. This renewal mechanism, in a meta-learning way, defines a distribution of adaptation tasks. 249 Each meta-RG evaluates the ability of the players to adapt.

250 The series of RGs that a Meta-RG consists of can be separated over two phases: a supporting 251 and querying/ZSCT phase. The supporting phase helps players learn the new semantics of the 252 symbolic space. Players must also agree on communication conventions by grounding the newly 253 randomized symbols. They prepare for the querying/ZSCT phase, during which ZSCT-purposed RGs 254 are played. They rely on target stimuli consisting of new combinations of the supporting stimuli's 255 atomic components. Focusing on novel combinations of never-before-seen components, prior to the current episode, is the main advantage of Meta-RGs over RGs. Indeed, this focus allows evaluation 256 of players' ability to learn to generalise compositionally. 257

- 258 Players that mastered the skill of learning to generalise compositionally will exhibit high RG 259 accuracy during the querying/ZSCT phase. But those who lacks the (meta-)learning/adaptation part 260 will not. Algorithms 4 and 5 (in Appendix A) contrast how a common RG differ from a meta-RG. 261 Indeed, the supporting phase of a Meta-RG does not involve updating the parameters/weights of the learning agents. This is in keeping with the meta-learning framework of the few-shot learning 262 kind. Its instantiation in the context of Meta-RG becomes striking when comparing the positions and 263 dependencies of lines 21 in Alg. 5 and 6 in Alg. 4. 264
- 265 The supporting phase lasts until all possible atomic component values on each latent dimension have 266 been shown for at least S shots (cf. lines 3 - 7 in Alg. 5). Thus, it will amount to at least S different 267 target stimuli being shown. That is to say at least S supporting-phase RGs. And, at most, there will actually be less supporting-phase RGs than the number of possible supporting-purposed stimuli in 268 the current episode's symbolic space/dataset. This is due to the focus on familiarising players with 269 atomic component values rather than stimuli themselves. On the other hand, the querying/ZSCT phase

- will present all the testing-purposed stimuli. We emphasise again that all RGs from both phases are played without the agents' parameters changing in-between. Indeed, learning CLBs involve agents adapting in an online/few-shot learning setting. The semantic structure of the symbolic space is randomly sampled at the beginning of each episode (cf. lines 2 - 3 in Alg. 5). The reward function proposed to both agents is null except on the N + 2-th step. It will be +1 if the listener decided correctly or, during the querying phase only, -2 if incorrect (cf. line 21 in Alg. 5).
- 276 Stimulus Representation. Meta-RG players must be able to deal with stimuli from  $N_{dim}$ -277 dimensioned symbolic spaces of varying semantic structures. Thus, it is necessary that stimulus 278 representation shapes remain constant from one semantic structure to another. This is not the case 279 for the common one-/multi-hot-encoded representation scheme. Thus, we propose the Symbolic 280 Continuous Stimulus (SCS) representation scheme, detailed in Section 3.2. Thanks to its shape *invariance property*, once a number of latent/factor dimension  $N_{dim}$  is choosen, it allows generation 281 of many different semantically structured symbolic spaces while maintaining a consistent stimulus 282 representation shape. Figure 2 highlights the structure of an episode, and its reliance on differently 283 semantically structured  $N_{dim}$ -dimensioned symbolic spaces. 284
- 285 Vocabulary Permutation. We remark that only changing the semantic structure of the symbolic 286 space is not sufficient to force MARL agents to adapt in each episode. Indeed, they can learn to cheat 287 by relying on an episode-invariant (and therefore independent of the instantiated semantic structure) emergent language (EL). This cheating EL consists of encoding the continuous values of the SCS 288 representation like an analog-to-digital converter would. It would map a fine-enough partition of 289 the SCS range onto a fixed vocabulary in a bijective fashion (see Appendix D for more details). 290 Therefore, to prevent the MARL agents from relying on such a cheating EL, we employ a vocabulary 291 permutation scheme (Cope & Schoots, 2021) that samples at the beginning of each episode a random 292 permutation of the vocabulary symbols (cf. line 1 in Alg. 2). This approach is bears some similarity 293 with the Other-Play algorithm from Hu et al. (2020). 294
- **Richness of the Stimulus.** We further bridge the gap between Hill-RSC and Chaa-RSC by allowing 295 the number of object-centric samples O and the number of shots S to be parameterized in the 296 benchmark. S represents the minimal number of times any given atomic component value may be 297 observed throughout the course of an episode. Intuitively, throughout their lifespan, an embodied 298 observer may only observe a given component a limited number of times (e.g. considering the value 299 'blue', on the latent/factor dimension 'color', being observed once within a 'blue car' stimulus, and 300 another time within a 'blue cup' stimulus). These parameters allow experimenters to account for both 301 the Chaa-RSC's sampling density of the different stimulus components and Hill-RSC's diversity of 302 viewpoints.
- 303 304 305

### 3.2 Symbolic Continuous Stimulus representation

306 Building about successes of the field of unsuper-307 vised learning of disentangled representations (Hig-308 gins et al., 2018), to the question what would an 309 ideally-disentangled latent space look like?, we pro-310 pose the Symbolic Continuous Stimulus (SCS) rep-311 resentation and provide numerical evidence of it in 312 Appendix E.2. It is continuous and relying on Gaus-313 sian kernels, and it has the particularity of enabling 314 the representation of stimuli sampled from differently 315 semantically structured symbolic spaces while maintaining the same representation shape (later referred 316 as the *shape invariance property*), as opposed to the 317 one-/multi-hot encoded (OHE/MHE) vector represen-318 tation commonly used when dealing with symbolic 319 spaces. 320

321 While the SCS representation is inspired by vectors 322 sampled from VAE's latent spaces, this representation

	d(0)=4 / d(1)=2 / d(2)=3		d(0)=3 / d(1)=3 / d(2)=3	
Example Latent Stimuli	( <mark>4 2</mark> 3)	(3 <b>1 2</b> )	(3 <b>1</b> 2)	(2 1 1)
MHE Sample Representation	( 000 <mark>1</mark> 0100 0010 )	( 0010 <b>1</b> 000 0 <b>1</b> 00 )	( 001 100 010 )	( 0 <mark>1</mark> 0 100 100 )
SCS Sample Representation	(0.82 0.55 0.74	0.25 -0.62 0.04	0.95 -0.90 0.10	-0.10 -0.85 -0.74

Figure 3: OHE/MHE and SCS representations of example latent stimuli for two differently-structured symbolic spaces with  $N_{dim} = 3$ , i.e. on the left for d(0) = 4, d(1) = 2, d(2) = 3, and on the right for d(0) = 3, d(1) = 3, d(2) = 3. Note the shape invariance property of the SCS representation, as its shape remains unchanged by the change in semantic structure of the symbolic space, on the contrary to the OHE/MHE representations.

is not learned and is not aimed to help the agent performing its task. It is solely meant to make it possible to define a distribution over infinitely many semantic/symbolic spaces, while instantiating

324 a BP for the agent to resolve. Indeed, contrary to OHE/MHE representation, observation of one 325 stimulus is not sufficient to derive the nature of the underlying semantic space that the current episode 326 instantiates (cf. Figure 3 for comparison). Rather, it is only via a kernel density estimation on multiple 327 samples (over multiple timesteps) that the semantic space's nature can be inferred, thus requiring the 328 agent to segregated and (re)combine information that is distributed over multiple observations. In other words, the benchmark instantiates a domain-agnostic BP. We provide in Appendix E.1 some numerical evidence to the fact that the SCS representation differentiates itself from the OHE/MHE 330 representation because it instantiates a BP. Deriving the SCS representation from an idealised VAE's 331 latent encoding of stimuli of any domain makes it a domain-agnostic representation, which is an 332 advantage compared to previous benchmark because domain-specific information can therefore not 333 be leveraged to solve the benchmark. 334

In details, the semantic structure of an N<sub>dim</sub>-335 dimensioned symbolic space is the tuple 336  $(d(i))_{i \in [1;N_{dim}]}$  where  $N_{dim}$  is the number of 337 latent/factor dimensions, d(i) is the **number of possi-**338 ble symbolic values for each latent/factor dimension 339 *i*. Stimuli in the SCS representation are vectors 340 sampled from the continuous space  $[-1, +1]^{N_{dim}}$ . In 341 comparison, stimuli in the OHE/MHE representation 342 are vectors from the discrete space  $\{0,1\}^{d_{OHE}}$ 343 where  $d_{OHE} = \sum_{i=1}^{N_{dim}} d(i)$  depends on the d(i)'s. 344 Note that SCS-represented stimuli have a shape that 345 does not depend on the d(i)'s values, this is the 346 shape invariance property of the SCS representation 347 (see Figure 3 for illustration).

In the SCS representation, the d(i)'s do not shape the 349 stimuli but only the semantic structure, i.e. represen-350 tation and semantics are disentangled from each other. 351 The d(i)'s shape the semantic by enforcing, for each 352 factor dimension *i*, a partitionaing of the [-1, +1]353 range into d(i) value sections. Each partition corre-354 sponds to one of the d(i) symbolic values available on the *i*-th factor dimension. Having explained how 355 to build the SCS representation sampling space, we 356 now describe how to sample stimuli from it. It starts 357 with instantiating a specific latent meaning/symbol, 358 embodied by latent values l(i) on each factor dimen-359 sion i, such that  $l(i) \in [1; d(i)]$ . Then, the i-th entry 360 of the stimulus is populated with a sample from a 361 corresponding Gaussian distribution over the l(i)-362 th partition of the [-1, +1 range. It is denoted as  $g_{l(i)} \sim \mathcal{N}(\mu_{l(i)}, \sigma_{l(i)})$ , where  $\mu_{l(i)}$  is the mean of 364 the Gaussian distribution, uniformly sampled to fall 365 within the range of the l(i)-th partition, and  $\sigma_{l(i)}$  is the standard deviation of the Gaussian distribution, uniformly sampled over the range  $\left[\frac{2}{12d(i)}, \frac{2}{6d(i)}\right]$ .  $\mu_{l(i)}$ 366 367 and  $\sigma_{l(i)}$  are sampled in order to guarantee (i) that 368 the scale of the Gaussian distribution is large enough, 369 but (ii) not larger than the size of the partition section 370 it should fit in. Figure 4 shows an example of such 371 instantiation of the different Gaussian distributions 372 over each factor dimensions' [-1, +1] range. 373



Figure 4: Visualisation of the SCSrepresented stimuli (column) observed by the speaker agent at each RG over the course of one meta-RG, with  $N_{dim} = 3$  and d(0) = 5, d(1) = 5, d(2) = 3. The supporting phase lasted for 19 RGs. For each factor dimension  $i \in [0; 2]$ , we present on the right side of each plot the kernel density estimations of the Gaussian kernels  $\mathcal{N}(\mu_{l(i)}, \sigma_{l(i)})$  of each latent value available on that factor dimension  $l(i) \in [1; d(i)]$ . Colours of dots, used to represent the sampled value  $g_{l(i)}$ , imply the latent value l(i)'s Gaussian kernel from which said continuous value was sampled. For each factor dimension, there is no overlap between the different latent values' Gaussian kernels.

# 374 375

376

348

### 4 EXPERIMENTS

377 **Agent Architecture.** The architectures of the RL agents that we consider are detailed in Appendix C. Optimization is performed via an R2D2 algorithm(Kapturowski et al., 2018) augmented with both the

Value Decomposition Network (Sunehag et al., 2017) and the Simplified Action Decoder approach (Hu & Foerster, 2019). As preliminary results showed poor performance, we follow Hill et al. (2020) and add an auxiliary reconstruction task to promote agents learning to use their core memory module. It consists of a mean squared-error between the stimuli observed at a given time step and a prediction conditioned on the current state of the core memory module after processing the current stimuli.

4.1 LEARNING CLBS IS OUT-OF-REACH TO STATE-OF-THE-ART MARL

386 Playing a meta-RG, the speaker aims at 387 each episode to make emerge a new lan-388 guage (constructivity) and the listener aims to acquire it (receptivity) as fast as possible, 389 before the querying-phase of the episode 390 comes around. Critically, we assume that 391 both agents must perform in accordance 392 with the principles of CLBs as it is the 393 only resolution approach. Indeed, there 394 is no success without a generalizing and 395 easy-to-learn EL, or, in other words, a (lin-396 guistically) compositional EL (Brighton & 397 Kirby, 2001; Brighton, 2002). Thus, we 398 investigate whether agents are able to coordinate to learn to perform CLBs from 399

Table 1: Meta-RG ZSCT and Ease-of-Acquisition (EoA) ZSCT accuracies and linguistic compositionality measures ( $\% \pm$  s.t.d.) for the multi-agent context after a sampling budget of 500k. The last column shows linguist results when evaluating the Posdis-Speaker (PS).

	Sh	Shots		
Metric	S = 1	S=2		
$Acc_{ZSCT}$ $\uparrow$	$53.6 \pm 4.7$	$51.6 \pm 2.2$	N/A	
$Acc_{\rm EoA}$ $\uparrow$	$50.6\pm8.8$	$50.6\pm5.8$	N/A	
topsim ↑	$29.6 \pm 16.8$	$21.3 \pm 16.6$	$96.7\pm0$	
posdis $\uparrow$	$23.7\pm20.8$	$13.8 \pm 12.8$	$92.0\pm0$	
bosdis †	$25.6\pm22.9$	$19.1 \pm 17.5$	$11.6 \pm 0$	

scratch, which is tantamount to learning receptivity and constructivity aspects of CLBs in parallel.

401 Evaluation & Results. We report the performance and compositionality of the behaviours in the multi-402 agent context in Table 1, on 3 random seeds of an LSTM-based model in the task with  $N_{dim} = 3$ , 403  $V_{min} = 2, V_{max} = 5, O = 4$ , and S = 1 or 2. As we assume no success without emergence of a 404 (linguistically) compositional language, we measure the linguistic compositionality profile of the emerging languages by, firstly, freezing the speaker agent's internal state (i.e. LSTM's hidden and 405 cell states) at the end of an episode and query what would be its subsequent utterances for all stimuli 406 in the latest episode's dataset (see Figure 2), and then compute the different compositionality metrics 407 on this collection of utterances. We compare the compositionality profile of the ELs to that of a 408 compositional language, in the sense of the posdis compositionality metric (Chaabouni et al., 2020) 409 (see Figure 7 and Table 6 in Appendix C.2). This language is produced by a fixed, rule-based agent 410 that we will refer to as the Posdis-Speaker (PS). Similarly, after the latest episode ends and the speaker 411 agent's internal state is frozen, we evaluate the EoA of the emerging languages by training a new, 412 non-meta/common listener agent for 512 epochs on the latest episode's dataset with the frozen 413 speaker agent using a descriptive-only/object-centric common RG and report its ZSCT accuracy (see 414 Algorithm 3). Table 1 shows Acczscr being around chance-level (50%), thus the meta-RL agents fail 415 to coordinate together, despite the simplicity of the setting, meaning that learning CLBs from scratch 416 is currently out-of-reach to state-of-the-art MARL agents, and therefore show the importance of our benchmark. As the linguistic compositionality measures are very low compared to the PS agent, and 417 since the chance-leveled  $Acc_{EoA}$  implies that the emerging languages are not easy to learn, it leads us 418 to think that the poor MARL performance is due to the lack of compositional language emergence. 419

420 421

422

383 384

385

### 4.2 SINGLE-AGENT LISTENER-FOCUSED RL CONTEXT

Seeing that the multi-agent benchmark is out of reach to state-of-the-art cooperative MARL agents, we investigate a simplification along two axises. Firstly, we simplify to a single-agent RL problem by instantiating a fixed, rule-based agent as the speaker, which should remove any issues related to agents learning in parallel to coordinate. Secondly, we use the Posdis-Speaker agent, which should remove any issues related to the emergence of assumed-necessary compositional languages, which corresponds to the constructivity aspects of CLBs. These simplifications allow us to focus our investigation on the receptivity aspects of CLBs, which relates to the ability from the listener agent to acquire and leverage a newly-encountered compositional language at each episode.

- 430 431
- 4.2.1 Symbol-Manipulation Induction Biases are Valuable

432 Firstly, in the simplest setting of O = 1 and S =433 1 we have the size that sum hal maximulation hi

1, we hypothesise that symbol-manipulation bi ases, such as efficient memory-addressing mech-

anism (e.g. attention) and greater algorithm-learning abilities (e.g. explicit memory), should

Table 2: Meta-RG ZSCT accuracies ( $\% \pm$  s.t.d.).

	LSTM	ESBN	DCEM
$Acc_{ZSCT} \uparrow$	$86.0\pm0.1$	$89.4\pm2.8$	$81.9\pm0.6$

improve performance, and propose to test the Emergent Symbol Binding Network (ESBN) (Webb
et al., 2020), the Dual-Coding Episodic Memory (DCEM) (Hill et al., 2020) and compare to baseline
LSTM (Hochreiter & Schmidhuber, 1997).

440 Evaluation & Results. We report in Table 2 the final ZSCT accuracies in the setting of  $N_{dim} = 3$ , 441  $V_{min} = 2, V_{max} = 3$ , with a sampling budget of 10M observations and 3 random seeds per 442 architecture. LSTM performing better than DCEM is presumably due to the difficulty of the latter in learning to use its complex memory scheme (preliminary experiments involving a Differentiable 443 Neural Computer (DNC - Graves et al. (2016)), on which the DCEM is built, show it struggling 444 to learn to use its memory compared to LSTM - cf Appendix E.3). On the other hand, we interpret 445 the best performance of the ESBN as being due to it being built over the LSTM, thus allowing its 446 complex memory scheme to be bypassed until it becomes useful. We validate our hypothesis but 447 carry on experimenting with the simpler LSTM model in order to facilitate analysis. 448

449

450

4.3 RECEPTIVITY ASPECTS OF CLBS CAN BE LEARNED SUB-OPTIMALLY

451 **Hypotheses.** The SCS representation instanti-452 ates a BP even when O = 1 (cf. Appendix E.1), 453 and we suppose that when O increases the BP's 454 complexity increases. Thus, it would stand to 455 reason to expect performance to decrease when 456 O increases (Hyp. 1). On the other hand, we would expect that increasing S would provide 457 the learning agent with a denser sampling (in 458 order to fulfill Chaa-RSC (ii)), and thus perfor-459

Table 3: Meta-RG ZSCT accuracies ( $\% \pm s.t.d.$ ).

	Shots		
Samples	S = 1	S=2	S = 4
O = 1	$62.2\pm3.7$	$73.5\pm2.4$	$75.0\pm2.3$
O = 4	$62.8 \pm 0.8$	$62.6 \pm 1.7$	$60.2 \pm 2.2$
O = 16	$64.9 \pm 1.7$	$62.0 \pm 2.0$	$61.8 \pm 2.1$

mance is expected to increase as S increases (Hyp. 2). Indeed, increasing S amounts to giving more opportunities for the agents to estimate each Gaussian, thus relaxing the instantiated BP's complexity.

462 Evaluation & Results. We report in table 3 ZSCT accuracies on LSTM-based models (6 random 463 seeds per settings) with  $N_{dim} = 3$  and  $V_{min} = 2$ ,  $V_{max} = 5$ . The chance threshold is 50%. When S = 1, increasing O is surprisingly correlated with non-significant increases in performance/sys-464 tematicity. On the other hand, when S > 1, accuracy distributions stay similar or decrease while O 465 increases. Thus, overall, Hyp. 1 tends to be validated. Regarding Hyp. 2, when O = 1, increasing 466 S (and with it the density of the sampling of the input space, i.e. Chaa-RSC (ii)) correlates with 467 increases in systematicity. Thus, despite the difference of settings between common RG, in Chaabouni 468 et al. (2020), and meta-RG here, we retrieve a similar result that Chaa-RSC promotes systematicity. 469 On the other hand, our results show a statistically significant distinction between BPs of complexity 470 associated with O > 1 and those associated with O = 1. Indeed, when O > 1, our results contradict 471 Hyp.2 since accuracy distributions remain the same or decrease when S increases. Acknowledging 472 the LSTMs' notorious difficulty with integrating/binding information from past to present inputs 473 over long dependencies, we explain these results based on the fact that increasing S also increases the length of each RL episode, thus the 'algorithm' learned by LSTM-based agents might fail to 474 adequately estimate Gaussian kernel densities associated with each component value. 475

476 477

478

### 4.4 LLMs perform below chance-level on S2B as Listener

Hypotheses. With LLMs being trained on curated human conversations, we investigate whether they have acquired skills in terms of receptivity aspects of CLBs. Thus, we test them without fine-tuning in the same listener-focused RL setting as in Section 4.3, with the same set of hypotheses. Our prompts are presented in Appendix B, but note that in this experiment we only make use of the listener prompt, as the speaker agent is played by our PS rule-based agent.

**Evaluation & Results.** We report in table 4 ZSCT accuracies on Mixtral-8x7B-Instruct-v0.1 (Mixtral) and OpenAI GPT-4o-mini (GPT) models, with  $N_{dim} = 3$  and  $V_{min} = 2$ ,  $V_{max} = 5$ . Over all contexts (O = 1 or 4 and S = 1 or 2), we observe poor performance (below chance level at 50%) from both

486 tested LLMs, thus providing very clear signal that the benchmark is challenging, even in this simpler 487 listener-focused RL setting, with low (and therefore simpler) values of  $(N_d im, V_m in, V_m ax)$ . 488

In more details, firstly, we ob-489 serve statistically non-significant 490 increases of performance when 491 S is increased from 1 to 2, for 492 both tested LLMs, which com-493 fort Chaa-RSC but the lack of 494 statistical significance prevent us 495 from drawing any conclusion. Then, as O is increased from 496 1 to 4, we observe a surprising 497 increase of performance, espe- M 498 cially marked for GPT-40-mini, G 499 500

Table 4: Meta-RG ZSCT accuracies (mean  $\% \pm$  s.t.d.) for Mixtral-8x7B-Instruct-v0.1 (Mixtral) and OpenAI GPT-4o-mini (GPT), with  $N_{\text{dim}} = 3$ ,  $V_{min} = 2$ ,  $V_{\text{max}} = 5$ , O = 1 and S = 1 or 2. Evaluation is performed over 5 seeds per setting and 64 Meta-RG episodes per seed, using the HuggingFace Text-Generation Inference API or OpenAI API, with Structured Outputs.

	O = 1		O = 4	
	S = 1	S=2	S = 1	S = 2
lixtral	$45.6\pm10.5$	$49.3 \pm 13.7$	$48.6 \pm 17.6$	$49.9\pm7.9$
PT	$33.1 \pm 11.4$	$36.8 \pm 12.7$	$39.8 \pm 11.6$	$42.9\pm3.8$
n ciani	ficent			

albeit far from being statistically non-significant.

We assume that these below-chance level results are, at least in part, caused by the well-studied limitation of LLMs when dealing with numerical values (Lee et al., 2023; Shen et al., 2023). We leave it to future works to investigate whether enhanced embeddings of numerical data (e.g. McLeish et al. (2024); Schwartz et al. (2024)) enable greater performance.

5 DISCUSSION

501

502

504 505 506

507 508

509 Compositional Behaviours vs CLBs. The learning of compositional behaviours (CBs) is one of the central study in language grounding with benchmarks like SCAN (Lake & Baroni, 2018) and 510 gSCAN (Ruis et al., 2020), as well as in the subfield of Emergent Communication (see Brandizzi 511 (2023); Boldt & Mortensen (2023) for reviews), but none investigates nor allow testing for CLBs. 512 Thus, our benchmark aims to fill in this gap. At a high level, SCAN and gSCAN are build to 513 evaluate compositional behaviours (CBs) only, because the (supervised-learning) agents tackling 514 those benchmarks are only facing a single, fixed underlying semantic structure. Thus, they get to 515 learn what atomic components compose this single, fixed underlying semantic structure over multiple 516 'lives' (over multiple update iterations of the agent's parameter weights). For an intuitive comparison, 517 SCAN and gSCAN benchmarks ask whether trained-and-now-frozen agents can generalise to novel 518 combinations of **those same familiar atomic components they have been training on**. This is 519 different from our proposed S2B, which intuitively asks whether trained-and-now-frozen agents 520 can generalise to novel combinations of **never-before-seen atomic components**, provided some warm-up exposition rounds to those atomic components. Warm-up exposition rounds do not 521 involve update of the agent's parameters. Thus, S2B instantiates meta-learning challenges where the 522 agents can only be successful if they have learned to learn to generalise compositionally, whereas 523 SCAN and gSCAN can be solved by agents that have only learned to generalise compositionally. 524 Indeed, S2B instantiates a theoretically-infinite distribution of different atomic components by way of 525 controlling the underlying semantic structure they belong to. Contrary to SCAN and gSCAN, agents 526 tackling the S2B are facing an infinity of always-changing underlying semantic structure. S2B is a 527 meta-learning benchmark, whereas SCAN and gSCAN are not, in principle. 528

That being said, SCAN and gSCAN can be updated to become meta-learning benchmarks, which is 529 what Lake (2019) and Lake & Baroni (2023) did to train their agents. Indeed, without making the 530 nuance, Lake (2019) and Lake & Baroni (2023) actually use CLBs as a training paradigm, where a 531 meta-learning extension of the sequence-to-sequence learning setting (i.e. CLB training) is shown to 532 enable human-like systematic CBs at test-time. Contrary to our work, they evaluate AI's abilities 533 towards SCAN-specific CBs after SCAN-specific CLBs training. We propose to go further because 534 we propose to train for CLBs and test for CLBs too, because we argue that CBs are actually not useful in open-ended contexts. We refer to open-ended contexts as real-world situations where agents would 536 encounter more diverse semantic structure than just the single one that they have been trained on. We 537 argue that CLBs are very useful in open-ended contexts. Moreover, even though Lake (2019) and Lake & Baroni (2023) extended SCAN to be used as a meta-learning benchmark, we argue that there 538 is a missing element, which is the instantiation of a novel binding problem (BP) in each task. Indeed, their extension made use of an OHE/MHE representation which we have argued in Section 3.2 and

Appendix E.1 that it does not instantiate a BP, contrary to the SCS representation which we use in our proposed S2B. Given their demonstration of the potential of CLBs, we leverage our proposed Meta-RG framework to propose a domain-agnostic CLB-focused benchmark for evaluation of CLBs abilities themselves, in order to address novel research questions around CLBs.

Symbolic Behaviours & Binding Problem. Following Santoro et al. (2021)'s definition of symbolic behaviours, our benchmark is the first specifically-principled benchmark to evaluate systematically artificial agents's abilities towards any symbolic behaviours. Similarly, while most challenging benchmark instantiates a version of the BP, as described by Greff et al. (2020), there is currently no principled benchmark that specifically investigates whether BP can be solved by artificial agents. Thus, not only does our benchmark fill that other gap, but it also instantiate a domain-agnostic version of the BP, which is critical in order to ascertain the external validity of conclusions that may be drawn from it. Indeed, domain-agnosticity guards us against confounders that could make the task solvable without fully solving the BP, e.g. by gaming some domain-specific aspects (Chollet, 2019). 

Limitations. Our experiments only evaluated state-of-the-art RL models and LLMs in the simplest configuration of our benchmark. We leave it to future works to investigate more complex configurations and evaluate other classes of models, such as neuro-symbolic models (Yu et al., 2023) or LLMs augmented with prompting methods that further reasoning abilities (Wei et al., 2022) and/or abilities to experiment (Lampinen et al., 2022) and correct themselves (Shinn et al., 2023).

In summary, we have proposed a novel benchmark to investigate artificial agents abilities at learning CLBs, by casting the problem of learning CLBs as a meta-reinforcement learning problem. It uses our proposed extension to RGs, entitled Meta-Referential Games, which contains an instantiation of a domain-agnostic BP. We provided baseline results for both the multi-agent tasks and the single-agent listener-focused tasks of learning CLBs in the context of our proposed benchmark. Our analysis of the behaviours in the multi-agent context highlighted the complexity for the speaker agent to invent a compositional language. But, when the language is already compositional, then a learning listener is able to acquire it and coordinate, albeit sub-optimally, with a rule-based speaker, in some of the simplest settings of our benchmark. Symbol-manipulation induction biases were found to be valuable, but, overall, our results show that our proposed benchmark is currently out of reach for current state-of-the-art artificial agents, as further exemplified by LLMs performing even below chance-level when instantiated as a listener and paired with our rule-based speaker. Thus, we hope our benchmark will spur the research community towards developing more capable artificial agents. 

# 594 REFERENCES

596 597 598	Peter W Battaglia, Jessica B Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, Caglar Gulcehre, Francis Song, Andrew Ballard, Justin Gilmer, George Dahl, Ashish Vaswani, Kelsey Allen, Charles Nash, Victoria Langston, Chris Dyer, Nicolas Heess, Daan Wierstra, Pushmeet Kohli,
599 600 601	Matt Botvinick, Oriol Vinyals, Yujia Li, and Razvan Pascanu. Relational inductive biases, deep learning, and graph networks. 2018. URL https://arxiv.org/pdf/1806.01261.pdf.
602 603	Jacob Beck, Risto Vuorio, Evan Zheran Liu, Zheng Xiong, Luisa Zintgraf, Chelsea Finn, and Shimon Whiteson. A survey of meta-reinforcement learning. <i>arXiv preprint arXiv:2301.08028</i> , 2023.
604 605	Yoshua Bengio. Deep learning of representations for unsupervised and transfer learning. <i>Conf. Proc. IEEE Eng. Med. Biol. Soc.</i> , 27:17–37, 2012.
607 608	Brendon Boldt and David R Mortensen. A review of the applications of deep learning-based emergent communication. <i>Transactions on Machine Learning Research</i> , 2023.
609 610 611	Nicolo' Brandizzi. Towards more human-like AI communication: A review of emergent communica- tion research. August 2023.
612 613 614	Henry Brighton. Compositional syntax from cultural transmission. <i>MIT Press</i> , Artificial, 2002. URL https://www.mitpressjournals.org/doi/abs/10.1162/106454602753694756.
616 617 618	Henry Brighton and Simon Kirby. The survival of the smallest: Stability conditions for the cultural evolution of compositional language. In <i>European Conference on Artificial Life</i> , pp. 592–601. Springer, 2001.
619 620 621 622	Henry Brighton and Simon Kirby. Understanding Linguistic Evolution by Visualizing the Emergence of Topographic Mappings. Artificial Life, 12(2):229–242, jan 2006. ISSN 1064-5462. doi: 10. 1162/artl.2006.12.2.229. URL http://www.mitpressjournals.org/doi/10.1162/ artl.2006.12.2.229.
623 624 625 626	Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. <i>Advances in neural information processing systems</i> , 33:1877–1901, 2020.
627 628 629	Rahma Chaabouni, Eugene Kharitonov, Diane Bouchacourt, Emmanuel Dupoux, and Marco Baroni. Compositionality and Generalization in Emergent Languages. apr 2020. URL http://arxiv. org/abs/2004.09124.
630 631 632	Ricky T Q Chen, Xuechen Li, Roger Grosse, and David Duvenaud. Isolating sources of disentangle- ment in VAEs, 2018.
633 634	Edward Choi, Angeliki Lazaridou, and Nando de Freitas. Compositional Obverter Communication Learning From Raw Visual Input. apr 2018. URL http://arxiv.org/abs/1804.02341.
636	François Chollet. On the Measure of Intelligence. Technical report, 2019.
637 638 639	Dylan Cope and Nandi Schoots. Learning to communicate with strangers via channel randomisation methods. <i>arXiv preprint arXiv:2104.09557</i> , 2021.
640 641 642	Kevin Denamganaï and James A. Walker. Referentialgym: A nomenclature and framework for language emergence & grounding in (visual) referential games. <i>4th NeurIPS Workshop on Emergent Communication</i> , 2020a.
643 644 645	Kevin Denamganaï and James Alfred Walker. Referentialgym: A framework for language emergence & grounding in (visual) referential games. <i>4th NeurIPS Workshop on Emergent Communication</i> , 2020b.
647	Roberto Dessi, Eugene Kharitonov, and Marco Baroni. Interpretable agent communication from scratch (with a generic visual processor emerging on the side). May 2021.

648 649 650	Jerry A Fodor and Zenon W Pylyshyn. Connectionism and cognitive architecture: A critical analysis. <i>Cognition</i> , 28(1-2):3–71, 1988.
651 652 653	<ul> <li>Alex Graves, Greg Wayne, Malcolm Reynolds, Tim Harley, Ivo Danihelka, Agnieszka Grabska-Barwińska, Sergio Gómez Colmenarejo, Edward Grefenstette, Tiago Ramalho, John Agapiou, et al. Hybrid computing using a neural network with dynamic external memory. <i>Nature</i>, 538(7626): 471–476, 2016.</li> </ul>
655 656	Klaus Greff, Sjoerd van Steenkiste, and Jürgen Schmidhuber. On the binding problem in artificial neural networks. <i>arXiv preprint arXiv:2012.05208</i> , 2020.
658 659 660	Shangmin Guo, Yi Ren, Serhii Havrylov, Stella Frank, Ivan Titov, and Kenny Smith. The emergence of compositional languages for numeric concepts through iterated learning in neural agents. <i>arXiv</i> preprint arXiv:1910.05291, 2019.
661 662 663	Irina Higgins, David Amos, David Pfau, Sebastien Racaniere, Loic Matthey, Danilo Rezende, and Alexander Lerchner. Towards a Definition of Disentangled Representations. dec 2018. URL http://arxiv.org/abs/1812.02230.
664 665 666	Felix Hill, Andrew Lampinen, Rosalia Schneider, Stephen Clark, Matthew Botvinick, James L McClelland, and Adam Santoro. Environmental drivers of systematicity and generalization in a situated agent. October 2019.
667 668 669	Felix Hill, Olivier Tieleman, Tamara von Glehn, Nathaniel Wong, Hamza Merzic, and Stephen Clark DeepMind. Grounded language learning fast and slow. Technical report, 2020.
670 671 672	Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. <i>Neural computation</i> , 9(8): 1735–1780, 1997.
673 674 675	Dan Horgan, John Quan, David Budden, Gabriel Barth-Maron, Matteo Hessel, Hado Van Hasselt, and David Silver. Distributed prioritized experience replay. <i>arXiv preprint arXiv:1803.00933</i> , 2018.
676 677 678	Hengyuan Hu and Jakob N Foerster. Simplified action decoder for deep multi-agent reinforcement learning. In <i>International Conference on Learning Representations</i> , 2019.
679 680	Hengyuan Hu, Adam Lerer, Alex Peysakhovich, and Jakob Foerster. "other-play" for zero-shot coordination. In <i>International Conference on Machine Learning</i> , pp. 4399–4410. PMLR, 2020.
681 682	Roman Jakobson. Linguistics and poetics. In <i>Style in language</i> , pp. 350–377. MA: MIT Press, 1960.
683 684 685	Steven Kapturowski, Georg Ostrovski, John Quan, Remi Munos, and Will Dabney. Recurrent experience replay in distributed reinforcement learning. In <i>International conference on learning representations</i> , 2018.
686 687 688	Hyunjik Kim and Andriy Mnih. Disentangling by factorising. <i>arXiv preprint arXiv:1802.05983</i> , 2018.
689 690	D P Kingma and M Welling. Auto-encoding variational bayes. <i>arXiv preprint arXiv:1312.6114</i> , 2013.
692 693	Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. <i>arXiv preprint arXiv:1412.6980</i> , 2014.
694 695	Manfred Krifka. Compositionality. <i>The MIT encyclopedia of the cognitive sciences</i> , pp. 152–153, 2001.
697 698	Brenden M Lake. Compositional generalization through meta sequence-to-sequence learning. <i>Advances in neural information processing systems</i> , 32, 2019.
699 700 701	Brenden M. Lake and Marco Baroni. Generalization without systematicity: On the compositional skills of sequence-to-sequence recurrent networks. <i>35th International Conference on Machine Learning, ICML 2018</i> , 7:4487–4499, oct 2018. URL http://arxiv.org/abs/1711.00350.

702 703	Brenden M Lake and Marco Baroni. Human-like systematic generalization through a meta-learning neural network. <i>Nature</i> , pp. 1–7, 2023.
704 705 706 707 708	Andrew K Lampinen, Nicholas Roy, Ishita Dasgupta, Stephanie CY Chan, Allison Tam, James Mcclelland, Chen Yan, Adam Santoro, Neil C Rabinowitz, Jane Wang, et al. Tell me why! explanations support learning relational and causal structure. In <i>International Conference on Machine Learning</i> , pp. 11868–11890. PMLR, 2022.
709 710 711	Angeliki Lazaridou and Marco Baroni. Emergent Multi-Agent communication in the deep learning era. June 2020.
712 713 714	Angeliki Lazaridou, Karl Moritz Hermann, Karl Tuyls, and Stephen Clark. Emergence of Linguistic Communication from Referential Games with Symbolic and Pixel Input. apr 2018. URL http://arxiv.org/abs/1804.03984.
715 716 717	Nayoung Lee, Kartik Sreenivasan, Jason D Lee, Kangwook Lee, and Dimitris Papailiopoulos. Teaching arithmetic to small transformers. <i>arXiv preprint arXiv:2307.03381</i> , 2023.
718	David Lewis. Convention: A philosophical study. 1969.
719 720 721 722	Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Rätsch, Sylvain Gelly, Bernhard Schölkopf, and Olivier Bachem. A sober look at the unsupervised learning of disentangled representations and their evaluation. October 2020.
723 724 725	João Loula, Marco Baroni, and Brenden M. Lake. Rearranging the Familiar: Testing Compositional Generalization in Recurrent Networks. jul 2018. URL http://arxiv.org/abs/1807. 07545.
726 727 728 729	Sean McLeish, Arpit Bansal, Alex Stein, Neel Jain, John Kirchenbauer, Brian R Bartoldson, Bhavya Kailkhura, Abhinav Bhatele, Jonas Geiping, Avi Schwarzschild, et al. Transformers can do arithmetic with the right embeddings. <i>arXiv preprint arXiv:2405.17399</i> , 2024.
730 731	Nikhil Mishra, Mostafa Rohaninejad, Xi Chen, and Pieter Abbeel. A simple neural attentive meta- learner. In <i>International Conference on Learning Representations</i> , 2018.
732 733	Richard Montague. Universal grammar. Theoria, 36(3):373–398, 1970.
734 735	Jesse Mu and Noah Goodman. Emergent communication of generalizations. Advances in Neural Information Processing Systems, 34:17994–18007, 2021.
736 737 738 739	Yi Ren, Shangmin Guo, Matthieu Labeau, Shay B. Cohen, and Simon Kirby. Compositional Languages Emerge in a Neural Iterated Learning Model. feb 2020. URL http://arxiv.org/abs/2002.01365.
740 741	Karl Ridgeway and Michael C Mozer. Learning deep disentangled embeddings with the F-Statistic loss, 2018.
743 744	Laura Ruis, Jacob Andreas, Marco Baroni, Diane Bouchacourt, and Brenden M Lake. A benchmark for systematic generalization in grounded language understanding. March 2020.
745 746 747	Adam Santoro, Andrew Lampinen, Kory Mathewson, Timothy Lillicrap, and David Raposo. Symbolic behaviour in artificial intelligence. <i>arXiv preprint arXiv:2102.03406</i> , 2021.
748 749 750	Eli Schwartz, Leshem Choshen, Joseph Shtok, Sivan Doveh, Leonid Karlinsky, and Assaf Arbelle. Numerologic: Number encoding for enhanced llms' numerical reasoning. <i>arXiv preprint arXiv:2404.00459</i> , 2024.
751 752 753	Ruoqi Shen, Sébastien Bubeck, Ronen Eldan, Yin Tat Lee, Yuanzhi Li, and Yi Zhang. Positional description matters for transformers arithmetic. <i>arXiv preprint arXiv:2311.14737</i> , 2023.
754 755	Noah Shinn, Federico Cassano, Beck Labash, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning.(2023). <i>arXiv preprint cs.AI/2303.11366</i> , 2023.

756 757 758 750	Agnieszka Słowik, Abhinav Gupta, William L. Hamilton, Mateja Jamnik, Sean B. Holden, and Christopher Pal. Exploring Structural Inductive Biases in Emergent Communication. feb 2020. URL http://arxiv.org/abs/2002.01335.
759 760 761 762	Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z Leibo, Karl Tuyls, et al. Value-decomposition networks for cooperative multi-agent learning. <i>arXiv preprint arXiv:1706.05296</i> , 2017.
763 764	Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. volume 29, 2016.
765 766 767 768	Ziyu Wang, Tom Schaul, Matteo Hessel, Hado Hasselt, Marc Lanctot, and Nando Freitas. Dueling network architectures for deep reinforcement learning. In <i>International conference on machine learning</i> , pp. 1995–2003. PMLR, 2016.
769 770	Taylor Whittington Webb, Ishan Sinha, and Jonathan Cohen. Emergent symbols through binding in external memory. In <i>International Conference on Learning Representations</i> , 2020.
771 772 773 774	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837, 2022.
775 776	Dongran Yu, Bo Yang, Dayou Liu, Hui Wang, and Shirui Pan. A survey on neural-symbolic learning systems. <i>Neural Networks</i> , 2023.
777	
778	
779	
780	
781	
782	
783	
704	
700	
787	
788	
789	
790	
791	
792	
793	
794	
795	
796	
797	
798	
799	
800	
800	
002	
80/	
805	
806	
807	
808	
809	

810 811 812 813 814 ON ALGORITHMIC DETAILS OF META-REFERENTIAL GAMES А 815 816 In this section, we detail algorithmically how Meta-Referential Games differ from common RGs. We 817 start by presenting in Algorithm 4 an overview of the common RGs, taking place inside a common 818 supervised learning loop, where we highlight the following: 819 820 (i) preparation of the data on which the referential game is played (highlighted in green), 821 822 (ii) elements pertaining to the supervised learning loop (highlighted in red). 823 824 Helper functions are detailed in Algorithm 1, 2 and 3. Next, we can now show in greater and 825 contrastive details the Meta-Referential Game algorithm in Algorithm 5, where we highlight the 826 following: 827 828 (i) preparation of the data on which the referential game is played (highlighted in green), 829 830 (ii) elements pertaining to the **meta-learning loop** (highlighted in blue). 831 832 (iii) elements pertaining to setup of a Meta-Referential Game (highlighted in red). 833 834 835 Algorithm 1: Helper function : DataPrep 836 Given 837 • a target stimuli  $s_0$ , 838 • a dataset of stimuli Dataset, 839 840 • O: Number of Object-Centric samples in each Target Distribution over stimuli  $TD(\cdot)$ . 841 • K : Number of distractor stimuli to provide to the listener agent. 842 • FullObs : Boolean defining whether the speaker agent has full (or partial) observation. 843 844 • DescrRatio : Descriptive ratio in the range [0, 1] defining how often the listener agent is observing the same semantic as the speaker agent. 845 846  $s'_0, D^{Target} \leftarrow s_0, 0;$ 847 if random(0,1) > DescrRatio then /\* Exclude target stimulus from listener's observation: 848 \*/  $s'_0 \sim \text{Dataset} - TD(s_0);$ 849  $D^{Target} \leftarrow K+1;$ 850 end 851 else if O > 1 then 852 Sample an Object-Centric distractor  $s'_0 \sim TD(s_0)$ ; 853 end 854 Sample K distractor stimuli from Dataset  $-TD(s_0)$ :  $(s_i)_{i \in [1,K]} \sim \text{Dataset} - TD(s_0)$ ; 855  $Obs_{\text{Speaker}} \leftarrow \{s_0\}$ ; if *FullObs* then 856  $Obs_{\text{Speaker}} \leftarrow \{s_0\} \cup \{s_i | \forall i \in [1, K]\};$ 857 end 858  $Obs_{\text{Listener}} \leftarrow \{s'_0\} \cup \{s_i | \forall i \in [1, K]\};$ /\* Shuffle listener observations and update index of target 859 decision: \*/ 860  $Obs_{Listener}, D^{Target} \leftarrow Shuffle(Obs_{Listener}, D^{Target});$ 861 **Output** :  $Obs_{Speaker}, Obs_{Listener}, D^{Target};$ 862

Algor	ithm 2: Helper function : MetaRGDatasetPreparation
Given	
	• V : Vocabulary (finite set of tokens available),
	• $N_{\text{dim}}$ : Number of attribute/factor dimensions in the symbolic spaces,
	• $V_{min}$ : Minimum number of possible values on each attribute/factor dimensions in the symbolic spaces,
	• $V_{max}$ : Maximum number of possible values on each attribute/factor dimensions in the symbolic spaces,
Initial	ise random permutation of vocabulary: $V' \leftarrow RandomPerm(V)$
Sampl	le semantic structure: $(d(i))_{i \in [1, N_{\text{dim}}]} \sim \mathcal{U}(V_{min}; V_{max})^{\mathcal{W}_{\text{dim}}};$
Gener Split (	ate symbolic space/dataset $D((d(i))_{i \in [1, N_{dim}]});$
reada	ability);
Outpu	at : $V', D((d(i))_{i \in [1, N_{dim}]}), D^{\text{support}}, D^{\text{query}};$
Algor	ithm 3. Helper function - PlayPG
Algor	ithm 3: Helper function : PlayRG
Algor Given	ithm 3: Helper function : PlayRG • Speaker and Listener agents.
Algor Given	ithm 3: Helper function : PlayRG         :         • Speaker and Listener agents,         • Set of speaker observations Observations (baser)
Algor Given	ithm 3: Helper function : PlayRG         :         • Speaker and Listener agents,         • Set of speaker observations Obs <sub>Speaker</sub> ,         • Sat of listener observations Obs_
Algor Given	ithm 3: Helper function : PlayRG         :         • Speaker and Listener agents,         • Set of speaker observations Obs <sub>Speaker</sub> ,         • Set of listener observations Obs <sub>Listener</sub> ,         • Number of speaker observations Obs <sub>Listener</sub> ,
Algor Given	ithm 3: Helper function : PlayRG         :         Speaker and Listener agents,         • Set of speaker observations Obs <sub>Speaker</sub> ,         • Set of listener observations Obs <sub>Listener</sub> ,         • N : Number of communication rounds to play,
Algor Given	<ul> <li>ithm 3: Helper function : PlayRG</li> <li>Speaker and Listener agents,</li> <li>Set of speaker observations Obs<sub>Speaker</sub>,</li> <li>Set of listener observations Obs<sub>Listener</sub>,</li> <li>N : Number of communication rounds to play,</li> <li>L : Maximum length of each message,</li> </ul>
Algor Given	<ul> <li>ithm 3: Helper function : PlayRG</li> <li>Speaker and Listener agents,</li> <li>Set of speaker observations Obs<sub>Speaker</sub>,</li> <li>Set of listener observations Obs<sub>Listener</sub>,</li> <li>Set of listener observations Obs<sub>Listener</sub>,</li> <li>N : Number of communication rounds to play,</li> <li>L : Maximum length of each message,</li> <li>V : Vocabulary (finite set of tokens available),</li> </ul>
Algor Given	ithm 3: Helper function : PlayRG         i         • Speaker and Listener agents,         • Set of speaker observations $Obs_{Speaker}$ ,         • Set of listener observations $Obs_{Listener}$ ,         • N : Number of communication rounds to play,         • L : Maximum length of each message,         • V : Vocabulary (finite set of tokens available),         ute message $M^S = \text{Speaker}(Obs_{\text{Speaker}} \emptyset);$
Algor Given Comp Initial	ithm 3: Helper function : PlayRG : • Speaker and Listener agents, • Set of speaker observations $Obs_{\text{Speaker}}$ , • Set of listener observations $Obs_{\text{Listener}}$ , • N : Number of communication rounds to play, • L : Maximum length of each message, • V : Vocabulary (finite set of tokens available), ute message $M^S = \text{Speaker}(Obs_{\text{Speaker}} \emptyset)$ ; ise Communication Channel History: CommH $\leftarrow [M^S]$ ;
Algor Given Comp Initial for ro	ithm 3: Helper function : PlayRG : • Speaker and Listener agents, • Set of speaker observations $Obs_{\text{Speaker}}$ , • Set of listener observations $Obs_{\text{Listener}}$ , • N : Number of communication rounds to play, • L : Maximum length of each message, • V : Vocabulary (finite set of tokens available), ute message $M^S = \text{Speaker}(Obs_{\text{Speaker}} \emptyset)$ ; ise Communication Channel History: CommH $\leftarrow [M^S]$ ; und = 0, N  do
Algor Given Comp Initial for ro.	ithm 3: Helper function : PlayRG • Speaker and Listener agents, • Set of speaker observations $Obs_{\text{Speaker}}$ , • Set of listener observations $Obs_{\text{Listener}}$ , • N : Number of communication rounds to play, • L : Maximum length of each message, • V : Vocabulary (finite set of tokens available), ute message $M^S = \text{Speaker}(Obs_{\text{Speaker}} \emptyset)$ ; ise Communication Channel History: CommH $\leftarrow [M^S]$ ; und = 0, N  do ompute Listener's reply $M_{\text{round}}^L$ = Listener( $Obs_{\text{Listener}} \text{CommH}$ );
Algor Given Comp Initial for ro.	ithm 3: Helper function : PlayRG • Speaker and Listener agents, • Set of speaker observations $Obs_{\text{Speaker}}$ , • Set of listener observations $Obs_{\text{Listener}}$ , • N : Number of communication rounds to play, • L : Maximum length of each message, • V : Vocabulary (finite set of tokens available), ute message $M^S = \text{Speaker}(Obs_{\text{Speaker}} \emptyset)$ ; ise Communication Channel History: CommH $\leftarrow [M^S]$ ; und = 0, N  do ompute Listener's reply $M_{\text{round}}^L$ , = Listener( $Obs_{\text{Listener}} \text{CommH}$ ); ommH $\leftarrow \text{CommH} + [M_{\text{round}}^L]$ ;
Algor Given Initial for ro.	ithm 3: Helper function : PlayRG • Speaker and Listener agents, • Set of speaker observations $Obs_{\text{Speaker}}$ , • Set of listener observations $Obs_{\text{Listener}}$ , • N : Number of communication rounds to play, • L : Maximum length of each message, • V : Vocabulary (finite set of tokens available), ute message $M^S = \text{Speaker}(Obs_{\text{Speaker}} \emptyset)$ ; ise Communication Channel History: CommH $\leftarrow [M^S]$ ; und = 0, N  do ompute Listener's reply $M_{\text{round}}^L = \text{Listener}(Obs_{\text{Listener}} \text{CommH})$ ; ommH $\leftarrow \text{CommH} + [M_{\text{round}}^L]$ ; $mmH \leftarrow \text{CommH} + [M_{\text{round}}^L]$ ;
Algor Given Comp Initial for ro. Co Co Co end	ithm 3: Helper function : PlayRG : • Speaker and Listener agents, • Set of speaker observations $Obs_{\text{Speaker}}$ , • Set of listener observations $Obs_{\text{Listener}}$ , • N : Number of communication rounds to play, • L : Maximum length of each message, • V : Vocabulary (finite set of tokens available), ute message $M^S = \text{Speaker}(Obs_{\text{Speaker}} \emptyset)$ ; ise Communication Channel History: CommH $\leftarrow [M^S]$ ; und = 0, N  do ompute Listener's reply $M_{\text{round}}^L$ , = Listener( $Obs_{\text{Listener}} \text{CommH}$ ); ommH $\leftarrow \text{CommH} + [M_{\text{round}}^L]$ ; ompute Speaker's reply $M_{\text{round}}^S$ = Speaker( $Obs_{\text{Speaker}} \text{CommH}$ ); ommH $\leftarrow \text{CommH} + [M_{\text{round}}^L]$ ;
Algor Given Comp Initial for ro. Co Co end Comp	ithm 3: Helper function : PlayRG : • Speaker and Listener agents, • Set of speaker observations $Obs_{Speaker}$ , • Set of listener observations $Obs_{Listener}$ , • N : Number of communication rounds to play, • L : Maximum length of each message, • V : Vocabulary (finite set of tokens available), ute message $M^S = \text{Speaker}(Obs_{\text{Speaker}} \emptyset)$ ; ise Communication Channel History: CommH $\leftarrow [M^S]$ ; und = 0, N  do sommt $\leftarrow \text{CommH} + [M^L_{round}]$ ; $\text{pompute Speaker's reply } M^L_{round} = \text{Speaker}(Obs_{\text{Speaker}} \text{CommH})$ ; $\text{pommH} \leftarrow \text{CommH} + [M^S_{round}]$ ; $\text{ute listener decision } D^L = \text{Listener}(Obs_{mult}) = (\text{CommH})$ ;

```
918
        Algorithm 4: Common Referential Game inside a Common Supervised Learning Loop
919
        Given
920
               • a dataset of stimuli Dataset,
921
               • a set of hyperparameters defining the RG:
922
923
                    - O: Number of Object-Centric samples in each Target Distribution over stimuli TD(\cdot).
924
                    - N: Number of communication rounds to play.
925
                    - L : Maximum length of each message.
926
                    - V : Vocabulary (finite set of tokens available).
927
                    - K: Number of distractor stimuli to provide to the listener agent.
928
                    - FullObs : Boolean defining whether the speaker agent has full (or partial) observation.
929
                    - Descriptive ratio in the range [0, 1] defining how often the listener agent
930
                       is observing the same semantic as the speaker agent.
931
                    - \mathcal{L}: Loss function to use in the agents update.
932
        Initialize:
933
934
               • Speaker(\cdot) and Listener(\cdot) agents.
935
        Systematically split Dataset into training and testing dataset, D<sup>train</sup> and D<sup>test</sup>;
936
        for epoch = 1, N_{epoch} do
937
            for target stimulus s_0 \in D^{train} do
938
                /* Preparation of observations and target decision:
                                                                                                       */
939
                Obs_{Speaker}, Obs_{Listener}, D^{Target} \leftarrow DataPrep(Dataset, s_0, O, K, FullObs, DescrRatio)
940
                /* Play Referential Game:
                                                                                                       */
941
                D^{L}, _ = PlayRG(Speaker, Listener, Obs_{Speaker}, Obs_{Listener}, N, L, V);
942
                /* Supervised Learning Parameters Update on Training
943
                     Stimulus Only:
                                                                                                       */
                Update both speaker and listener agents' parameters using the loss \mathcal{L}(D^{Target}, D^L);
944
945
            end
946
            Initialise ZSCT accuracy: Acc_{ZSCT} \leftarrow 0;
            for target stimulus s_0 \in D^{test} do
947
                /* Preparation of observations and target decision:
                                                                                                       */
948
                Obs_{Speaker}, Obs_{Listener}, D^{Target} \leftarrow DataPrep(Dataset, s_0, O, K, FullObs, DescrRatio)
949
                /* Play Referential Game:
                                                                                                       */
950
                D^{L}, _ = PlayRG(Speaker, Listener, Obs_{Speaker}, Obs_{Listener}, N, L, V);
951
                /* Update ZSCT Accuracy:
                                                                                                       */
952
                Acc_{ZSCT} \leftarrow Update(Acc_{ZSCT}, D^{Target}, D^L);
953
            end
954
        end
955
956
957
958
959
960
961
962
963
964
965
966
967
```

- 968 969
- 970 971

Algorithm 5: Meta-Referential Game inside a Meta-Learning Loop	
Given :	
• $N_{episode}$ , $N_{dim}$ : Number of episodes, and number of attribute/factor dimensions,	
• $S$ : Minimum number of Shots over which each possible value on each attribute/fac	ctor
dimension ought to be observed by the agents (as part of a target stimulus).	
• V	/factor
dimensions in the symbolic spaces.	incetor
• $TSS(\mathcal{D}   S, S)$ : Target stimulus sampling function which samples from dataset $\mathcal{D}$	aiven a
set of previously sampled stimuli S while maximising the likelihood that each pos	sible
value on each attribute/factor dimension are sampled at least S times.	JIOIC
• a set of hyperparameters defining the RG:	
O N where Colling Contributions is the Transformer Distribution of the	
- O: Number of Object-Centric samples in each Target Distribution over stimuli	$TD(\cdot).$
-N: Number of communication rounds to play.	
-L: Maximum length of each message.	
-V: Vocabulary (finite set of tokens available).	
- $K$ : Number of distractor stimuli to provide to the listener agent.	
- FullObs : Boolean defining whether the speaker agent has full (or partial) obse	rvation.
– DescrRatio : Descriptive ratio in the range [0, 1] defining how often the listene	r agent
is observing the same semantic as the speaker agent.	
Initialize :	
• Sneaker( $\cdot$ ) and Listener( $\cdot$ ) agents	
for the third sector ( ) agents.	
Ior $episode = 1, N_{episode}$ do	
$V' D = D^{\text{support}} D^{\text{query}} $ $M \text{ sta} D C D \text{ starset} D \text{ support} D^{\text{query}} $ $V = V = V$	~/
$V$ , $D_{episode}$ , $D_{episode}$ , $D_{episode} \leftarrow MetarGDataSetPreparation(V, N_{dim}, V_{min}, V_{max});$	
initialise set of sampled supporting sumuli: $\mathcal{S}^{\text{regat}} \leftarrow \emptyset$ ;	
Sample training purposed target stimulus $a^i = TSS(D^{\text{support}} S^{\text{support}})$	
Sample training-purposed target similar $S_0 \sim I SS(D_{\text{episode}}, S^{(n)}, S)$	
$\{S_{i}^{\text{orrel}} \leftarrow S_{i}^{\text{orrel}} \cup \{S_{0}^{\text{o}}\}; i \leftarrow i + 1;$	
Initialise RG index: $i \leftarrow 0$ :	:5,
/* Supporting Phase:	*/
for target stimulus $s_0^i \in S^{support}$ do	, i
$D_{i}^{Target} \leftarrow DataPren(D^{support}, s_{i}^{i}, O, K, FullObs, Descr Rati$	o).
$\mathcal{D}$ and $\mathcal{D}$ and $\mathcal{D}$ episode, $\mathcal{D}_i$ , $\mathcal{D}_i$ , $\mathcal{D}_i$ episode, $\mathcal{D}_i$ ,	0),
$D_i^L, CommH_i = \text{PlayRG}(\text{Speaker}, \text{Listener}, Obs_{\text{Speaker}}^i, Obs_{\text{Listener}}^i, N, L, V');$	
$=$ Listener( $Obs_{S_{rescher}}^{i} CommH_{i}$ ); /* Listener-Feedback Step	*/
end	
/* Querying/ZSCT Phase:	*/
Initialise ZSCT accuracy: $Acc_{ZSCT} \leftarrow 0$ ;	
<b>for</b> target stimulus $s_0^i \in D_{episode}^{query}$ <b>do</b>	
$Obs_{\text{Speaker}}^{i}, Obs_{\text{Listener}}^{i}, \dot{D}_{i}^{Target} \leftarrow DataPrep(D_{\text{episode}}, s_{0}^{i}, O, K, \text{FullObs}, \text{DescrRati})$	o);
	*
$D_i^L, CommH_i = PlayRG(Speaker, Listener, Obs_{Speaker}^i, Obs_{Listener}^i, N, L, V');$	
_,_=Listener( $Obs_{Speaker}^{i} CommH_{i}$ ); /* Listener-Feedback Step	*/
/* Update ZSCT Accuracy:	*/
$Acc_{\text{ZSCT}} \leftarrow \text{Update}(Acc_{\text{ZSCT}}, D_i^{Target}, D_i^L); \ i \leftarrow i+1;$	
end	
/* Meta-Learning Parameters Update on Whole Episode:	*/
(1 if $D_i^{Target} == D_i^L$	
Update both agents using rewards $R_i = \begin{cases} 0 & \text{otherwise, during supporting phase:} \end{cases}$	
-2 otherwise during merving phase	
end	

1026 1027 1028

#### В **PROMPTS FOR LANGUAGE MODELS**

1029 1030 1031

1032

1033

1034

1035

1036

1037

1039

1040

1041

1042

1043

1045

1047

1048

1049

1050

1051

1052

1053

1055

1056

1057

1058

1061

1062

1063

1064

1067

1068

1069

#### **Speaker & Listener - Detailed Prompt Examples** Speaker Prompt - RG #0 - Step 1 Listener Prompt - RG #0 - Step 1 You and your partner are playing a sequence You and your partner are playing a sequence $\hookrightarrow$ of referential games. You are the speaker $\hookrightarrow$ of referential games. You are the $\hookrightarrow$ listener. In the first phase, you will get accounted In the first phase, you will get accounted 2 $\hookrightarrow$ with the atomic components of the $\hookrightarrow$ with the atomic components of the $\hookrightarrow$ possible observations. Then, the game $\hookrightarrow$ possible observations. Then, the game $\hookrightarrow$ counter will restart, and you will be $\hookrightarrow$ counter will restart, and you will be $\hookrightarrow$ tested with new observations, combining $\hookrightarrow$ tested with new observations, combining $\hookrightarrow$ the same atomic components in novel ways. $\hookrightarrow$ the same atomic components in novel ways. At each game, each of you observes a At each game, each of you observes a 3 $\hookrightarrow$ stimulus, which represents a latent $\hookrightarrow$ stimulus, which represents a latent $\hookrightarrow$ meaning, and your common goal is to $\hookrightarrow$ meaning, and your common goal is to $\hookrightarrow$ figure out whether you are observing $\hookrightarrow$ figure out whether you are observing $\hookrightarrow$ different or similar latent meanings. You $\hookrightarrow$ different or similar latent meanings. To $\hookrightarrow$ can communicate with your partner using $\, \hookrightarrow \,$ help you do so, your partner can send you $\hookrightarrow$ the communication channel. The $\hookrightarrow$ messages using the communication channel $\hookrightarrow$ communication channel is made up of 10 $\hookrightarrow$ , which is made up of 10 symbols that can $\hookrightarrow$ symbols that you can combine together to $\hookrightarrow$ be combined together to form a sentence $\hookrightarrow$ of maximum length 5. $\hookrightarrow$ form a sentence of maximum length 5. $\hookrightarrow$ Beware that symbol 0 is grounded already. Beware that symbol 0 is grounded already. It 4 It is the end-of-message symbol. It $\hookrightarrow$ is the end-of-message symbol. It means 1046 $\hookrightarrow$ means that any symbol that comes after it $\hookrightarrow$ that any symbol that comes after it will $\rightarrow$ will be ignored and regularised into $\hookrightarrow$ be ignored and regularised into symbol 0. $\hookrightarrow$ symbol 0. From one game to the next, you should aim to Starting game #0, this is the new stimulus: 6 → [0.346 -0.524 0.868]. be consistent so that your partner can $\rightarrow$ figure out the code that you are using to communicate and decrypt messages towards You are an expert in the matter. Given the 8 $\hookrightarrow$ information above, answer the following fulfilling your common goal. $\hookrightarrow$ question(s) to the best of your abilities $\hookrightarrow$ . Starting game #0, this is the new stimulus: Ouestion #1: Are you observing a stimulus 10 1054 You are an expert in the matter. Given the $\hookrightarrow$ representing the same latent meaning as $\hookrightarrow$ the stimulus observed by your partner? ightarrow information above, answer the following $\hookrightarrow$ question(s) to the best of your abilities Answer either 0.:'Yes' or 1.:'No'. 11 $\hookrightarrow$ 12 Question #2: What message should you send 13 Question #1: Do you think your partner $\hookrightarrow$ your partner to better coordinate with → understands your messages? $\hookrightarrow$ them towards fulfilling your common goal? Answer either 0.:'Yes' or 1.:'No'. The message is made up of 5 symbols, each of $\hookrightarrow$ which can be filled with one of the 10 Question #2: What message should you send to $\hookrightarrow$ vocab symbols. For example: [1, 9, 2, 6, $\hookrightarrow$ your partner to better coordinate $\hookrightarrow$ 5]. $\hookrightarrow$ together towards fulfilling your common This question corresponds to 5 implicit 15 $\hookrightarrow$ goal? $\hookrightarrow$ questions, one for each of the 5 symbols The message is made up of 5 symbols, each of $\hookrightarrow$ of the message. Thus, each possible $\hookrightarrow$ which can be filled with one of the 10 $\hookrightarrow$ answer id is between 0 and 9, $\hookrightarrow$ vocab symbols. For example: [0, 0, 0, 0, $\hookrightarrow$ corresponding to one of the 10 vocab $\rightarrow$ 01. $\hookrightarrow$ symbols. This question corresponds to 5 implicit $\hookrightarrow$ questions, one for each of the 5 symbols $\hookrightarrow$ of the message. Thus, each possible $\hookrightarrow$ answer id is between 0 and 9, $\hookrightarrow$ corresponding to one of the 10 vocab $\hookrightarrow$ symbols.

Figure 5: Detailed prompts for the speaker and listener agents at the start of a Meta-Referential Game 1070 (RG #0), starting with an explanation of the context and the rules (lines 1-4), followed by information 1071 about the (previous and) current game(s) (line 6 - cf. Figure 6 - RG #1 Step 1 for more details 1072 presenting previous games). The prompt is ended (from line 8) with two multi-choice questions 1073 referring to the two actions that agents may perform, i.e. providing a decision and a message that is 1074 to be communicated to the other player at the next step. Note that both of those actions are not 1075 necessary for both players, i.e. the decision is only necessary for the listener, and the message 1076 is only necessary for the speaker in the experiments presented here. Our proposed benchmark 1077 incorporates the setting where the listener agent is allowed to communicate back to the speaker, but 1078 we leave investigation of this more complex setting to subsequent works. RL actions that are not 1079 necessary are simply ignored by the benchmark RL environment.



Figure 6: Partial prompts for the speaker and listener agents, from step 2 of RG #0 to step 1 of RG #1, highlighting information changes from one step to the next. Note the effect of the vocabulary permutation scheme transforming the message send by the speaker into a vocabulary-permutated version of the message observed by the listener, mainly at step 2. Step 3 corresponds to the Listener Feedback step where the listener is presented with the exact stimulus that the speaker was observing during the current RG.

	R2D2	
Number of actors	3	32
Actor parameter update interval	1 environ	ment step
Sequence unroll length	2	20
Sequence length overlap	10	
Sequence burn-in length	10	
N-steps return	ŕ	3
Replay buffer size	$5  imes 10^4$ ol	bservations
Priority exponent	0	.9
Importance sampling exponent	0	.6
Discount $\gamma$	0.9	997
Minibatch size	3	32
Optimizer	Adam (Kingm	na & Ba, 2014)
Optimizer settings	learning rate $= 6.25$	$5 \times 10^{-5}, \epsilon = 10^{-12}$
Target network update interval	2500 updates	
Value function rescaling	None	
Core	Memory Module	
LSTM (Hochreiter & Schmidhuber, 1997)	DNC (Graves et al., 2016)	
Number of layers 2	LSTM-controller settings	2 hidden layers of size 128
Hidden layer size 256, 128	Memory settings	128 slots of size 32
Activation function ReLU	Read/write heads	2 reading ; 1 writing

Table 5: Hyper-parameters values used in R2D2, with LSTM or DNC as the core memory module.
 All missing parameters follow the ones in Ape-X (Horgan et al., 2018).

# 1159 1160

## C AGENT ARCHITECTURE & TRAINING

1161 1162

The baseline RL agents that we consider use a 3-layer fully-connected network with 512, 256, and 1163 finally 128 hidden units, with ReLU activations, with the stimulus being fed as input. The output 1164 is then concatenated with the message coming from the other agent in a OHE/MHE representation, 1165 mainly, as well as all other information necessary for the agent to identify the current step, i.e. the 1166 previous reward value (either +1 and 0 during the training phase or +1 and -2 during testing phase), 1167 its previous action in one-hot encoding, an OHE/MHE-represented index of the communication 1168 round (out of N possible values), an OHE/MHE-represented index of the agent's role (speaker or listener) in the current game, an OHE/MHE-represented index of the current phase (either 'training' 1169 or 'testing'), an OHE/MHE representation of the previous RG's result (either success or failure), the 1170 previous RG's reward, and an OHE/MHE mask over the action space, clarifying which actions are 1171 available to the agent in the current step. The resulting concatenated vector is processed by another 1172 3-layer fully-connected network with 512, 256, and 256 hidden units, and ReLU activations, and then 1173 fed to the core memory module, which is here a 2-layers LSTM (Hochreiter & Schmidhuber, 1997) 1174 with 256 and 128 hidden units, which feeds into the advantage and value heads of a 1-layer dueling 1175 network (Wang et al., 2016).

Table 5 highlights the hyperparameters used for the learning agent architecture and the learning algorithm, R2D2(Kapturowski et al., 2018). More details can be found, for reproducibility purposes, in our open-source implementation at HIDDEN\_FOR\_REVIEW\_PURPOSE.

Training was performed for each run on 1 NVIDIA GTX1080 Ti, and the average amount of training time for a run is 18 hours for LSTM-based models, 40 hours for ESBN-based models, and 52 hours for DCEM-based models.

1183

1184 C.1 ESBN & DCEM 1185

The ESBN-based and DCEM-based models that we consider have the same architectures and parameters than in their respective original work from Webb et al. (2020) and Hill et al. (2020), with the exception of the stimuli encoding networks, which are similar to the LSTM-based model.

# 1188 C.2 RULE-BASED SPEAKER AGENT

1190 The rule-based speaker agents used in the single-agent task, where only the listener agent is a 1191 learning agent, speaks a compositional language in the sense of the posdis metric (Chaabouni et al., 1192 2020), as presented in Table 6 for  $N_{dim} = 3$ , a maximum sentence length of L = 4, and vocabulary 1193 size  $|V| \ge max_i d(i) = 5$ , assuming a semantical space such that  $\forall i \in [1,3], d(i) = 5$ .



Figure 7: Visualisation on each column of the messages sent by the posdis-compositional rulebased speaker agent over the course of the episode presented in Figure 4. Colours are encoding the information of the token index, as a visual cue.

## 1215 D CHEATING LANGUAGE

12161217The agents can develop a cheating language, cheating in the<br/>sense that it could be episode/task-invariant (and thus semantic<br/>structure invariant). This emerging cheating language would<br/>encode the continuous values of the SCS representation like an<br/>analog-to-digital converter would, by mapping a fine-enough<br/>partition of the [-1, +1] range onto the vocabulary in a bijective<br/>fashion.

For instance, for a vocabulary size ||V|| = 10, each symbol can be unequivocally mapped onto  $\frac{2}{10}$ -th increments over [-1, +1], and, by communicating  $N_{dim}$  symbols (assuming  $N_{dim} \leq L$ ), the speaker agents can communicate to the listener the (digitized) continuous value on each dimension *i* of the SCS-

Table 6: Examples of the latent stimulus to language utterance mapping of the posdis-compositional rule-based speaker agent. Note that token 0 is the EoS token.

Latent Dims			Comp. Language
#1	#2	#3	Tokens
0	1	2	1, 2, 3, 0
1	3	4	2, 4, 5, 0
2	5	0	3, 6, 1, 0
3	1	2	4, 2, 3, 0
4	3	4	5, 4, 5, 0

represented stimulus. If  $max_j d(j) \le ||V||$  then the cheating language is expressive-enough for the speaker agent to digitize all possible stimulus without solving the binding problem, i.e. without inferring the semantic structure. Similarly, it is expressive-enough for the listener agent to convert the spoken utterances to continuous/analog-like values over the [-1, +1] range, thus enabling the listener agent to skirt the binding problem when trying to discriminate the target stimulus from the different stimuli it observes.

1234

1213 1214

### E FURTHER EXPERIMENTS:

1236

1237 1238

E.1 ON THE BP INSTANTIATED BY THE SCS REPRESENTATION

Hypothesis. The SCS representation differs from the OHE/MHE one primarily in terms of the binding problem (Greff et al., 2020) that the former instantiates while the latter does not. Indeed, the semantic structure can only be inferred after observing multiple SCS-represented stimuli. We hypothesised that it is via the *dynamic binding of information* extracted from each observations that

an estimation of a density distribution over each dimension *i*'s [-1, +1] range can be performed. And, estimating such density distribution is tantamount to estimating the number of likely gaussian distributions that partition each [-1, +1] range.

1245 **Evaluation.** Towards highlighting that there is a binding problem taking place, we show results of 1246 baseline RL agents (similar to main experiments in Section 4) evaluated on a simple single-agent 1247 recall task. The Recall task structure borrows from few-shot learning tasks as it presents over 2 shots 1248 all the stimuli of the instantiated symbolic space (not to be confused with the case for Meta-RG 1249 where all the latent/factor dimensions' values are being presented over S shots – Meta-RGs do not 1250 necessarily sample the whole instantiated symbolic space at each episode, but the Recall task does). 1251 Each shot consists of a series of recall games, one for each stimulus that can be sampled from an  $N_{dim} = 3$ -dimensioned symbolic space. The semantic structure  $(d(i))_{i \in [1; N_{dim}]}$  of the symbolic 1252 space is randomly sampled at the beginning of each episode, i.e.  $d(i) \sim \mathcal{U}(2; 5)$ , where  $\mathcal{U}(2; 5)$  is the 1253 uniform discrete distribution over the integers in [2; 5], and the number of object-centric samples is 1254 O = 1, in order to remove any confounder from object-centrism. 1255

1256 Each recall game consists of two steps: in the first step, a stimulus is presented to the RL agent, and 1257 only a *no-operation* (NO-OP) action is made available, while, on the second step, the agent is asked 1258 to infer/recall the **discrete** l(i) **latent value** (as opposed to the representation of it that it observed, 1259 either in the SCS or OHE/MHE form) that the previously-presented stimulus had instantiated, on a given *i*-th dimension, where value *i* for the current game is uniformly sampled from  $\mathcal{U}(1; N_{dim})$ 1260 at the beginning of each game. The value of i is communicated to the agent via the observation 1261 on this second step of different stimulus that in the first step: it is a zeroed out stimulus with the 1262 exception of a 1 on the *i*-th dimension on which the inference/recall must be performed when using 1263 SCS representation, or over all the OHE/MHE dimensions that can encode a value for the *i*-th latent 1264 factor/attribute when using the OHE/MHE representation. On the second step, the agent's available 1265 action space now consists of discrete actions over the range  $[1; max_i d(j)]$ , where  $max_i d(j)$  is a 1266 hyperparameter of the task representing the maximum number of latent values for any latent/factor 1267 dimension. In our experiments,  $max_i d(j) = 5$ . While the agent is rewarded at each game for 1268 recalling correctly, we only focus on the performance over the games of the second shot, i.e. on the 1269 games where the agent has theoretically received enough information to infer the density distribution 1270 over each dimension i's [-1, +1] range. Indeed, observing the whole symbolic space once (on the first shot) is sufficient (albeit not necessary, specifically in the case of the OHE/MHE representation). 1271

1272 **Results.** Figure 8 details the recall accuracy over all the 1273 games of the second shot of our baseline RL agent through-1274 out learning. There is a large gap of asymptotic perfor-1275 mance depending on whether the Recall task is evaluated 1276 using OHE/MHE or SCS representations. We attribute the poor performance in the SCS context to the instantia-1277 tion of a BP. We note again that during those experiments 1278 the number of object-centric samples was kept at O = 1, 1279 thus emphasising that the BP is solely depending on the 1280 use of the SCS representation and does not require object-1281 centrism. 1282



Figure 8: 5-ways 2-shots accuracies on the Recall task with different stimulus representation (OHE:blue ; SCS; orange).

- 1283 1284 E.2 ON THE IDEALLY-DISENTANGLED-NESS OF THE 1285 SCS REPRESENTATION
- 1286 In this section, we verify our hypothesis that the SCS representation yields ideally-disentangled 1287 stimuli. We report on the FactorVAE Score Kim & Mnih (2018), the Mutual Information Gap 1288 (MIG) Chen et al. (2018), and the Modularity Score Ridgeway & Mozer (2018) as they have been 1289 shown to be part of the metrics that correlate the least among each other (Locatello et al., 2020), 1290 thus representing different desiderata/definitions for disentanglement. We report on the  $N_{dim} = 3$ -1291 dimensioned symbolic spaces with  $\forall j, d(j) = 5$  and O = 5. The measurements are of 100.0%, 94.8, and 98.9% for, respectivily, the FactorVAE Score, the MIG, and the Modularity Score, thus validating our design hypothesis about the SCS representation. We remark that the MIG and Modularity 1293 Score are sensitive to the number of object-centric samples O, which can be seen decreasing the 1294 measurements as low as 64.4% and 66.6% for O = 1. The FactorVAE Score is not affected, possibly 1295 due to its reliance on a deterministic classifier.

# 1296 E.3 AUXILIARY RECONSTRUCTION LOSS

In the following, we investigate and compare the performance when using an LSTM (Hochreiter & Schmidhuber, 1997) or a Differentiable Neural Computer (DNC) (Graves et al., 2016) as core memory module, with or without the auxiliary reconstruction loss inspired from Hill et al. (2020).

In the case of the LSTM, the prediction network of the reconstruction loss takes as input the LSTM hidden states, while in the case of the DNC, the input is the memory. Figure 9b shows the stimulus reconstruction accuracies for both architectures, highlighting a greater data-efficiency (and resulting asymptotic performance in the current observation budget) of the LSTM-based architecture, compared to the DNC-based one.

1306 Figure 9a shows the 4-ways (3 distractors descriptive meta-RGs) ZSCT accuracies of the different 1307 agents throughout learning. The ZSCT accuracy is the accuracy over querying-/testing-purpose 1308 stimuli only, after the agent has observed for two consecutive times (i.e. S = 2) the supportive 1309 training-purpose stimuli for the current episode. The DNC-based architecture has difficulty learning 1310 how to use its memory, even with the use of the auxiliary reconstruction loss, and therefore it utterly fails to reach better-than-chance ZSCT accuracies. On the otherhand, the LSTM-based architecture is 1311 fairly successful on the auxiliary reconstruction task, but it is not sufficient for training on the main 1312 task to really take-off. As expected from the fact that the benchmark instantiates a binding problem 1313 that requires relational responding, our results hint at the fact that the ability to use memory towards 1314 deriving valuable relations between stimuli seen at different time-steps is primordial. Indeed, only the 1315 agent that has the ability to use its memory element towards recalling stimuli starts to perform at a 1316 better-than-chance level. Thus, the auxiliary reconstruction loss is an important element to drive some 1317 success on the task, but it is also clearly not sufficient, and the rather poor results that we achieved 1318 using these baseline agents indicates that new inductive biases must be investigated to be able to 1319 solve the problem posed in our proposed benchmark.

1320 1321

1322 1323

1332

# F BROADER IMPACT

No technology is safe from being used for malicious purposes, which equally applies to our research. 1324 However, aiming to develop artificial agents that relies on the same symbolic behaviours and the same 1325 social assumptions (e.g. using CLBs) than human beings is aiming to reduce misunderstanding be-1326 tween human and machines. Thus, the current work is targeting benevolent applications. Subsequent 1327 works around the benchmark that we propose are prompted to focus on emerging protocols in general 1328 (not just posdis-compositional languages), while still aiming to provide a better understanding of 1329 artificial agent's symbolic behaviour biases and differences, especially when compared to human 1330 beings, thus aiming to guard against possible misunderstandings and misaligned behaviours. The 1331 current state of this work does not allow discussion of potential negative societal impact.



Figure 9: (a): 4-ways (3 distractors) zero-shot compositional test accuracies of different architectures. 5 seeds for architectures with DNC and LSTM, and 2 seeds for runs with DNC+Rec and LSTM+Rec, where the auxiliary reconstruction loss is used. (b): Stimulus reconstruction accuracies for the architectures augmented with the auxiliary reconstruction task. Accuracies are computed on binary values corresponding to each stimulus' latent dimension's reconstructed value being close enough to the ground truth value, with a threshold of 0.05 on each dimension, which correspond to a deviation tolerance of 2.5% since the range in which SCS stimuli are instantiated is [-1, 1].